

# **SAS<sup>®</sup> OPTGRAPH**

## **Procedure 14.1**

### **Graph Algorithms and Network Analysis**

The correct bibliographic citation for this manual is as follows: SAS Institute Inc. 2015. *SAS® OPTGRAPH Procedure 14.1: Graph Algorithms and Network Analysis*. Cary, NC: SAS Institute Inc.

**SAS® OPTGRAPH Procedure 14.1: Graph Algorithms and Network Analysis**

Copyright © 2015, SAS Institute Inc., Cary, NC, USA

All Rights Reserved. Produced in the United States of America.

**For a hard-copy book:** No part of this publication may be reproduced, stored in a retrieval system, or transmitted, in any form or by any means, electronic, mechanical, photocopying, or otherwise, without the prior written permission of the publisher, SAS Institute Inc.

**For a web download or e-book:** Your use of this publication shall be governed by the terms established by the vendor at the time you acquire this publication.

The scanning, uploading, and distribution of this book via the Internet or any other means without the permission of the publisher is illegal and punishable by law. Please purchase only authorized electronic editions and do not participate in or encourage electronic piracy of copyrighted materials. Your support of others' rights is appreciated.

**U.S. Government License Rights; Restricted Rights:** The Software and its documentation is commercial computer software developed at private expense and is provided with RESTRICTED RIGHTS to the United States Government. Use, duplication, or disclosure of the Software by the United States Government is subject to the license terms of this Agreement pursuant to, as applicable, FAR 12.212, DFAR 227.7202-1(a), DFAR 227.7202-3(a), and DFAR 227.7202-4, and, to the extent required under U.S. federal law, the minimum restricted rights as set out in FAR 52.227-19 (DEC 2007). If FAR 52.227-19 is applicable, this provision serves as notice under clause (c) thereof and no other notice is required to be affixed to the Software or documentation. The Government's rights in Software and documentation shall be only those set forth in this Agreement.

SAS Institute Inc., SAS Campus Drive, Cary, NC 27513-2414

July 2015

SAS® and all other SAS Institute Inc. product or service names are registered trademarks or trademarks of SAS Institute Inc. in the USA and other countries. ® indicates USA registration.

Other brand and product names are trademarks of their respective companies.

# Contents

---

Chapter 1. The OPTGRAPH Procedure . . . . .	1
---	---

<b>Index</b>	<b>237</b>
--------------	------------





# Credits

---

## Documentation

Writing	Matthew Galati, Yi Liao
Editing	Anne Baxter
Documentation Support	Tim Arnold, Melanie Gratton, Daniel Underwood
Technical Review	Mustafa Kabul, Charles B. Kelly, Minghui Liu, Michelle Opp, Rob Pratt

---

## Software

PROC OPTGRAPH	Matthew Galati, Yi Liao
---------------	-------------------------

---

## Support Groups

Software Testing	Mustafa Kabul, Charles B. Kelly, Minghui Liu, Rob Pratt
Technical Support	Tonya Chapman



# Chapter 1

## The OPTGRAPH Procedure

### Contents

---

Overview: OPTGRAPH Procedure . . . . .	<b>3</b>
Getting Started: OPTGRAPH Procedure . . . . .	<b>4</b>
Road Network Shortest Path . . . . .	4
Authority in U.S. Supreme Court Precedence . . . . .	7
Syntax: OPTGRAPH Procedure . . . . .	<b>10</b>
Functional Summary . . . . .	10
PROC OPTGRAPH Statement . . . . .	17
BICONCOMP Statement . . . . .	20
CENTRALITY Statement . . . . .	20
CLIQUE Statement . . . . .	25
COMMUNITY Statement . . . . .	26
CONCOMP Statement . . . . .	28
CORE Statement . . . . .	29
CYCLE Statement . . . . .	30
DATA_LINKS_VAR Statement . . . . .	32
DATA_MATRIX_VAR Statement . . . . .	33
DATA_NODES_VAR Statement . . . . .	33
EIGENVECTOR Statement . . . . .	33
LINEAR_ASSIGNMENT Statement . . . . .	34
MINCOSTFLOW Statement . . . . .	35
MINCUT Statement . . . . .	36
MINSPANTREE Statement . . . . .	37
PERFORMANCE Statement . . . . .	38
REACH Statement . . . . .	39
SHORTPATH Statement . . . . .	40
SUMMARY Statement . . . . .	42
TRANSITIVE_CLOSURE Statement . . . . .	44
TSP Statement . . . . .	44
Details: OPTGRAPH Procedure . . . . .	<b>49</b>
Graph Input Data . . . . .	49
Matrix Input Data . . . . .	57
Data Input Order . . . . .	58
Parallel Processing . . . . .	59
Numeric Limitations . . . . .	59
Size Limitations . . . . .	61
Common Notation and Assumptions . . . . .	61

Biconnected Components and Articulation Points . . . . .	61
Centrality . . . . .	65
Clique . . . . .	86
Community Detection . . . . .	89
Connected Components . . . . .	97
Core Decomposition . . . . .	102
Cycle . . . . .	107
Eigenvector Problem . . . . .	112
Linear Assignment (Matching) . . . . .	115
Minimum-Cost Network Flow . . . . .	116
Minimum Cut . . . . .	124
Minimum Spanning Tree . . . . .	128
Reach (Ego) Network . . . . .	130
Shortest Path . . . . .	142
Summary . . . . .	153
Transitive Closure . . . . .	162
Traveling Salesman Problem . . . . .	164
Macro Variables . . . . .	172
ODS Table Names . . . . .	184
Examples: OPTGRAPH Procedure . . . . .	<b>187</b>
Example 1.1: Articulation Points in a Terrorist Network . . . . .	187
Example 1.2: Influence Centrality for Project Groups in a Research Department . . . .	189
Example 1.3: Betweenness and Closeness Centrality for Computer Network Topology	193
Example 1.4: Betweenness and Closeness Centrality for Project Groups in a Research Department . . . . .	196
Example 1.5: Eigenvector Centrality for Word Sense Disambiguation . . . . .	199
Example 1.6: Centrality Metrics for Project Groups in a Research Department . . . .	202
Example 1.7: Community Detection on Zachary's Karate Club Data . . . . .	205
Example 1.8: Recursive Community Detection on Zachary's Karate Club Data . . . .	209
Example 1.9: Cycle Detection for Kidney Donor Exchange . . . . .	211
Example 1.10: Linear Assignment Problem for Minimizing Swim Times . . . . .	216
Example 1.11: Linear Assignment Problem, Sparse Format versus Dense Format . . . .	219
Example 1.12: Minimum Spanning Tree for Computer Network Topology . . . . .	222
Example 1.13: Transitive Closure for Identification of Circular Dependencies in a Bug Tracking System . . . . .	223
Example 1.14: Reach Networks for Computation of Market Coverage of a Terrorist Network . . . . .	226
Example 1.15: Traveling Salesman Tour through US Capital Cities . . . . .	229
References . . . . .	<b>234</b>

---

## Overview: OPTGRAPH Procedure

The OPTGRAPH procedure includes a number of graph theory, combinatorial optimization, and network analysis algorithms. The algorithm classes are listed in [Table 1.1](#).

**Table 1.1** Algorithm Classes in PROC OPTGRAPH

Algorithm Class	PROC OPTGRAPH Statement
Biconnected components	BICONCOMP
Centrality metrics	CENTRALITY
Maximal cliques	CLIQUE
Community detection	COMMUNITY
Connected components	CONCOMP
Core decomposition	CORE
Cycle detection	CYCLE
Eigenvector problem	EIGENVECTOR
Weighted matching	LINEAR_ASSIGNMENT
Minimum-cost network flow	MINCOSTFLOW
Minimum cut	MINCUT
Minimum spanning tree	MINSPANTREE
Reach networks	REACH
Shortest path	SHORTPATH
Graph summary	SUMMARY
Transitive closure	TRANSITIVE_CLOSURE
Traveling salesman	TSP

You can use the OPTGRAPH procedure to analyze relationships between entities. These relationships are typically defined by using a *graph*. A graph  $G = (N, A)$  is defined over a set  $N$  of nodes and a set  $A$  of arcs. A *node* is an abstract representation of some entity (or object), and an *arc* defines some relationship (or connection) between two nodes. The terms *node* and *vertex* are often interchanged in describing an entity. The term *arc* is often interchanged with the term *edge* or *link* when describing a connection.

You can check the SAS log for the version number being used in any invocation of PROC OPTGRAPH. The following statements check the version:

```
proc optgraph;
run;
```

Then the log displays the version number as shown in [Figure 1.1](#).

**Figure 1.1** Version Number Displayed in Log

---

```
NOTE: -----
NOTE: Running OPTGRAPH version 14.1.
NOTE: -----
NOTE: The OPTGRAPH procedure is executing in single-machine mode.
NOTE: -----
```

---

## Getting Started: OPTGRAPH Procedure

Since graphs are abstract objects, their analyses have applications in many different fields of study, including social sciences, linguistics, biology, transportation, marketing, and so on. This document shows a few potential applications through simple examples.

This section shows two introductory examples for getting started with the OPTGRAPH procedure. For more detail about the input formats expected and the various algorithms available, see the sections “[Details: OPTGRAPH Procedure](#)” on page 49 and “[Examples: OPTGRAPH Procedure](#)” on page 187.

## Road Network Shortest Path

Consider the following road network between a SAS employee’s home in Raleigh, NC, and the SAS headquarters in Cary, NC.

In this road network (graph), the links are the roads and the nodes are intersections between roads. For each road, you assign a *link attribute* in the variable `time_to_travel` to describe the number of minutes that it takes to drive from one node to another. The following data were collected using Google Maps (Google 2011), which gives an approximate number of minutes to traverse between two points, based on the length of the road and the typical speed during normal traffic patterns:

```
data LinkSetInRoadNC10am;
  input start_inter $1-20 end_inter $20-40 miles miles_per_hour;
  datalines;
614CapitalBlvd      Capital/WadeAve      0.6  25
614CapitalBlvd      Capital/US70W        0.6  25
614CapitalBlvd      Capital/US440W       3.0  45
Capital/WadeAve      WadeAve/RaleighExpy 3.0  40
Capital/US70W        US70W/US440W        3.2  60
US70W/US440W        US440W/RaleighExpy  2.7  60
Capital/US440W       US440W/RaleighExpy  6.7  60
US440W/RaleighExpy  RaleighExpy/US40W   3.0  60
WadeAve/RaleighExpy RaleighExpy/US40W   3.0  60
RaleighExpy/US40W   US40W/HarrisonAve   1.3  55
US40W/HarrisonAve   SASCampusDrive      0.5  25
;

data LinkSetInRoadNC10am;
  set LinkSetInRoadNC10am;
  time_to_travel = miles * 1/miles_per_hour * 60;
run;
```

Using PROC OPTGRAPH, you want to find the route that yields the shortest path between home (614 Capital Blvd) and the SAS headquarters (SAS Campus Drive). This can be done with the SHORTPATH statement as follows:

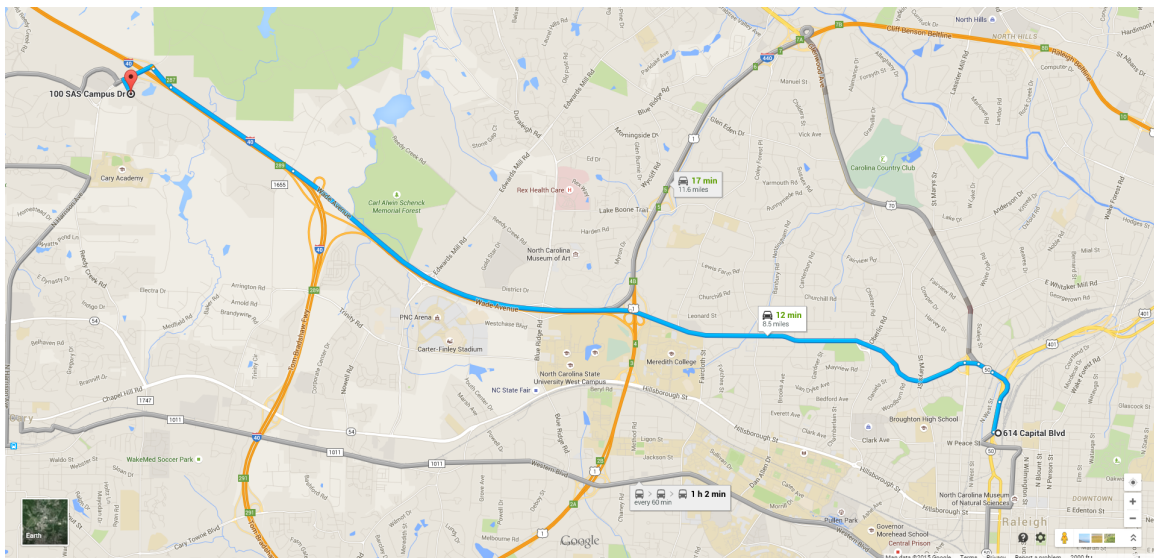
```
proc optgraph
  data_links      = LinkSetInRoadNC10am;
  data_links_var
    from          = start_inter
    to            = end_inter
    weight        = time_to_travel;
  shortpath
    out_paths     = ShortPath
    source        = "614CapitalBlvd"
    sink          = "SASCampusDrive";
run;
```

For more information about shortest path algorithms in PROC OPTGRAPH, see the section “[Shortest Path](#)” on page 142. [Figure 1.2](#) displays the output data set ShortPath, which shows the best route to take to minimize travel time at 10:00 a.m. This route is also shown in Google Maps in [Figure 1.3](#).

**Figure 1.2** Shortest Path for Road Network at 10:00 A.M.

order	start_inter	end_inter	time_to_travel
1	614CapitalBlvd	Capital/WadeAve	1.4400
2	Capital/WadeAve	WadeAve/RaleighExpy	4.5000
3	WadeAve/RaleighExpy	RaleighExpy/US40W	3.0000
4	RaleighExpy/US40W	US40W/HarrisonAve	1.4182
5	US40W/HarrisonAve	SASCampusDrive	1.2000
			<b>11.5582</b>

**Figure 1.3** Shortest Path for Road Network at 10:00 A.M. in Google Maps



Now suppose that it is rush hour (5:00 p.m.) and the time to traverse the roads has changed because of traffic patterns. You want to find the route that is the shortest path for going home from SAS headquarters under different speed assumptions due to traffic. The following data set lists approximate travel times and speeds for driving in the opposite direction:

```
data LinkSetInRoadNC5pm;
  input start_inter $1-20 end_inter $20-40 miles miles_per_hour;
  datalines;
614CapitalBlvd      Capital/WadeAve      0.6  25
614CapitalBlvd      Capital/US70W      0.6  25
614CapitalBlvd      Capital/US440W      3.0  45
Capital/WadeAve      WadeAve/RaleighExpy  3.0  25 /*high traffic*/
Capital/US70W        US70W/US440W      3.2  60
US70W/US440W        US440W/RaleighExpy  2.7  60
Capital/US440W       US440W/RaleighExpy  6.7  60
US440W/RaleighExpy  RaleighExpy/US40W    3.0  60
WadeAve/RaleighExpy RaleighExpy/US40W    3.0  60
RaleighExpy/US40W   US40W/HarrisonAve    1.3  55
US40W/HarrisonAve   SASCampusDrive      0.5  25
;

data LinkSetInRoadNC5pm;
  set LinkSetInRoadNC5pm;
  time_to_travel = miles * 1/miles_per_hour * 60;
run;
```

The following statements are similar to the first PROC OPTGRAPH run, except that they use the data set LinkSetInRoadNC5pm and the SOURCE and SINK option values are reversed:

```
proc optgraph
  data_links    = LinkSetInRoadNC5pm;
  data_links_var
    from        = start_inter
    to          = end_inter
    weight      = time_to_travel;
  shortpath
    out_paths   = ShortPath
    source      = "SASCampusDrive"
    sink        = "614CapitalBlvd";
run;
```

Now, the output data set ShortPath, shown in [Figure 1.4](#), shows the best route for going home. Because the traffic on Wade Avenue is usually heavy at this time of day, the route home is different from the route to work.

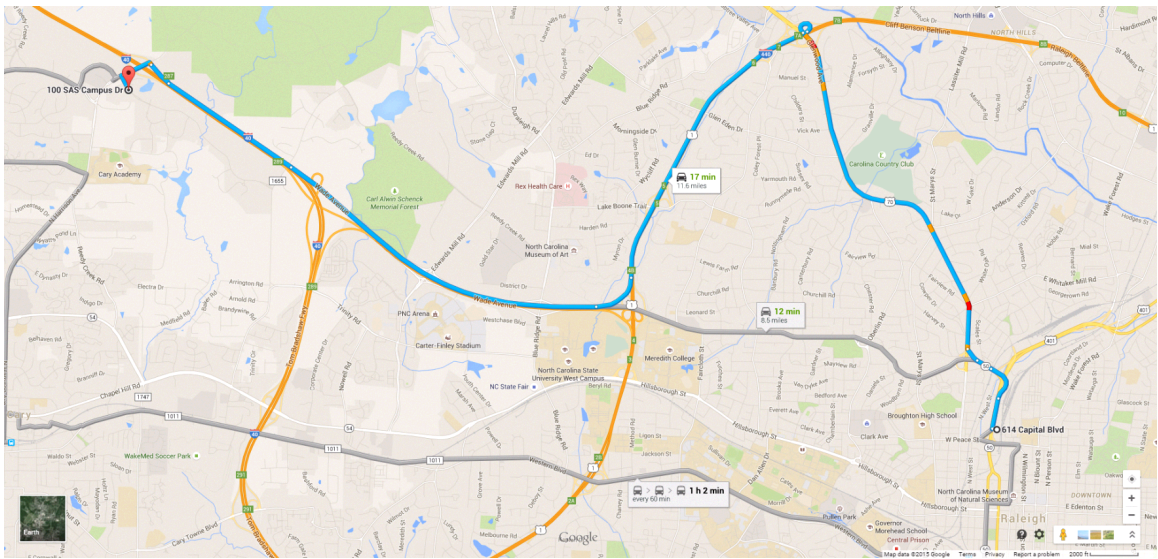


Figure 1.4 Shortest Path for Road Network at 5:00 P.M.

order	start_inter	end_inter	time_to_travel
1	SASCampusDrive	US40W/HarrisonAve	1.2000
2	US40W/HarrisonAve	RaleighExpy/US40W	1.4182
3	RaleighExpy/US40W	US440W/RaleighExpy	3.0000
4	US440W/RaleighExpy	US70W/US440W	2.7000
5	US70W/US440W	Capital/US70W	3.2000
6	Capital/US70W	614CapitalBlvd	1.4400
			12.9582

This new route is shown in Google Maps in Figure 1.5.

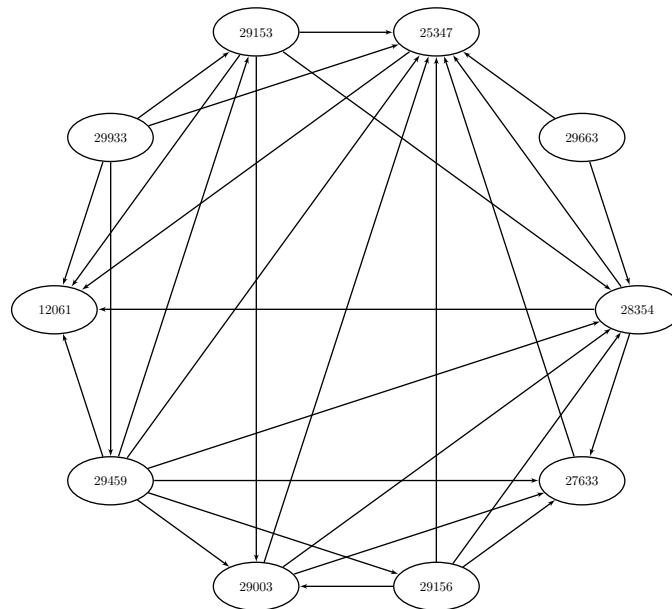
Figure 1.5 Shortest Path for Road Network at 5:00 P.M. in Google Maps



Authority in U.S. Supreme Court Precedence

This example looks at the use of precedents in court cases. Consider the judge’s problem of identifying precedent court cases that are most relevant and important to the current case. This application of network analysis was published in Fowler and Joen 2008. Because of norms inherited from 19th century English law, judges are encouraged to follow precedent in order to take advantage of “accumulated experience of many judges responding to the arguments and evidence of many lawyers” (Landes and Posner 1976). In network analysis, one way to define the importance of a previous case is to look at the network of citations used in related cases. That is, if a particular case *A* cited case *B* to help support its argument, then a link exists from *A* to *B* in the citation network.

Given such a citation network, you can then use a metric known as *authority score* to rank the importance of these cases. This metric is explained in more detail in the section “Hub and Authority Scoring” on page 76. Figure 1.6 shows a small representative subset of the citation network for landmark abortion decisions from the example in Fowler and Joen 2008.

**Figure 1.6** Citation Network for Some U.S. Supreme Court Cases

The data set Cases stores a mapping between case name and the case identifier:

```
data Cases;
  length case_id 8 case_name $80;
  input case_id 1-5 case_name $ 7-80;
  datalines;
12061 Jacobson v. Massachusetts, 197 U.S. 11 (1905)
25347 Roe vs. Wade, 410 U.S. 113 (1973)
27633 Akron vs. Akron Cntr for Repro-Health, 462 U.S. 416 (1983)
28354 Thornburgh vs. American College, 476 U.S. 747 (1986)
29003 Webster vs. Repro-Health Services, 492 U.S. 490 (1989)
29153 Cruzan v. Director, MO Dept of Health, 497 U.S. 261 (1990)
29155 Georgia v. South Carolina, 497 U.S. 376 (1990)
29156 Hodgson v. Minnesota, 497 U.S. 417 (1990)
29459 Planned Parenthood of SE PA vs. Casey, 505 U.S. 833 (1992)
29663 Madsen v. Women's Health Ctr., 512 U.S. 753 (1994)
29933 Wash. v. Glucksberg, 521 U.S. 702 (1997)
;
```

The data set LinkSetInCourt provides the citation network between case identifiers:

```
data LinkSetInCourt;
  input from_case to_case @@;
  datalines;
27633 25347 28354 25347 28354 27633 29003 25347 29003 27633
29003 28354 29459 25347 29459 27633 29459 28354 29459 29003
25347 12061 28354 12061 29459 12061 29933 25347 29933 29459
29933 12061 29933 29153 29663 25347 29663 28354 29153 12061
29153 28354 29153 29003 29153 25347 29459 29153 29156 27633
29156 28354 29156 29003 29156 25347 29459 29156
;
```

You can calculate the authority scores of each case by using the CENTRALITY statement with the AUTH= option, as follows:

```
proc optgraph
  direction      = directed
  data_links     = LinkSetInCourt
  out_nodes      = NodeSetOut;
  data_links_var
    from         = from_case
    to           = to_case;
  centrality
    auth         = unweight;
run;
```

The output data set NodeSetOut contains the authority score for each case (node). Then, the following DATA step combines the case names with the case identifiers and sorts on the score:

```
data NodeSetOut(drop=rc case_id);
  if _n_=1 then do;
    declare hash h(dataset:'cases');
    h.definekey('case_id');
    h.definedata('case_name');
    h.definedone();
  end;
  set NodeSetOut;
  length case_id 8 case_name $80;
  rc=h.find(key:node);
run;

proc sort data=NodeSetOut;
  by descending centr_auth_unwt;
run;
```

As expected, *Roe vs. Wade* (1973) has the highest authority ranking since it is most often cited by other cases.

**Figure 1.7** Authority Ranking of Landmark U.S. Supreme Court Cases

node	centr_auth_unwt	case_name
25347	1.00000	Roe vs. Wade, 410 U.S. 113 (1973)
28354	0.72262	Thornburgh vs. American College, 476 U.S. 747 (1986)
12061	0.61717	Jacobson v. Massachusetts, 197 U.S. 11 (1905)
27633	0.59831	Akron vs. Akron Cntr for Repro-Health, 462 U.S. 416 (1983)
29003	0.50930	Webster vs. Repro-Health Services, 492 U.S. 490 (1989)
29153	0.31742	Cruzan v. Director, MO Dept of Health, 497 U.S. 261 (1990)
29156	0.20968	Hodgson v. Minnesota, 497 U.S. 417 (1990)
29459	0.10775	Planned Parenthood of SE PA vs. Casey, 505 U.S. 833 (1992)
29933	0.00000	Wash. v. Glucksberg, 521 U.S. 702 (1997)
29663	0.00000	Madsen v. Women's Health Ctr., 512 U.S. 753 (1994)

In such a small example, it is somewhat easy to see which cases have the most influence by looking at the directed graph of citations. As discussed in Fowler and Joen 2008, the real advantage of such an analysis can be seen when examining all the citations for all 30,288 cases available in their data.

---

## Syntax: OPTGRAPH Procedure

**PROC OPTGRAPH** *options* ;

*Data Input Statements:*

**DATA\_LINKS\_VAR** < *options* > ;  
**DATA\_MATRIX\_VAR** < *column1, column2, ...* > ;  
**DATA\_NODES\_VAR** < *options* > ;

*Algorithm Statements:*

**BICONCOMP** < *option* > ;  
**CENTRALITY** < *options* > ;  
**CLIQUE** < *options* > ;  
**COMMUNITY** < *options* > ;  
**CONCOMP** < *options* > ;  
**CORE** < *options* > ;  
**CYCLE** < *options* > ;  
**EIGENVECTOR** < *options* > ;  
**LINEAR\_ASSIGNMENT** < *options* > ;  
**MINCOSTFLOW** < *options* > ;  
**MINCUT** < *options* > ;  
**MINSPANTREE** < *options* > ;  
**REACH** < *options* > ;  
**SHORTPATH** < *options* > ;  
**SUMMARY** < *options* > ;  
**TRANSITIVE\_CLOSURE** < *options* > ;  
**TSP** < *options* > ;

*Performance Statement:*

**PERFORMANCE** < *options* > ;

PROC OPTGRAPH statements are divided into four main categories: the **PROC** statement, the **data input** statements, the **algorithm** statements, and the **PERFORMANCE** statement. The PROC statement invokes the procedure and sets option values that are used across multiple algorithms. The data input statements control the names of the variables that PROC OPTGRAPH expects in the data input. The algorithm statements determine which algorithms are run and set options for each individual algorithm. The PERFORMANCE statement specifies performance options for multithreaded computing.

The section “**Functional Summary**” on page 10 provides a quick reference for each of the options for each statement. Each statement is then described in more detail in its own section; the PROC OPTGRAPH statement is described first, and sections that describe all other statements are presented in alphabetical order.

---

## Functional Summary

Table 1.2 summarizes the statements and options available with PROC OPTGRAPH.

**Table 1.2** Functional Summary

Description	Option
<b>PROC OPTGRAPH Options</b>	
<b>Input</b>	
Specifies the link data set	DATA_LINKS=
Specifies the matrix data set	DATA_MATRIX=
Specifies the node data set	DATA_NODES=
Specifies the node subset data set	DATA_NODES_SUB=
<b>Output</b>	
Specifies the link output data set	OUT_LINKS=
Specifies the node output data set	OUT_NODES=
<b>Options</b>	
Specifies the subgraph filter level	FILTER_SUBGRAPH=
Specifies the graph direction	GRAPH_DIRECTION=
Specifies the internal graph format	GRAPH_INTERNAL_FORMAT=
Includes self links	INCLUDE_SELFLINK
Specifies the overall log level	LOGLEVEL=
Specifies whether time units are in CPU time or real time	TIMETYPE=
<b>Data Input Statements</b>	
<b>DATA_LINKS_VAR Options</b>	
Specifies the data set variable name for the <i>from</i> nodes	FROM=
Specifies the data set variable name for the link flow lower bounds	LOWER=
Specifies the data set variable name for the <i>to</i> nodes	TO=
Specifies the data set variable name for the link flow upper bounds	UPPER=
Specifies the data set variable name for the link weights	WEIGHT=
<b>DATA_MATRIX_VAR</b>	
Specifies the data set variable names for the matrix	
<b>DATA_NODES_VAR Options</b>	
Specifies the data set variable name for cluster identifiers	CLUSTER=
Specifies the data set variable name for the nodes	NODE=
Specifies the data set variable name for node weights	WEIGHT=
Specifies the data set variable name for auxiliary node weights	WEIGHT2=
<b>Algorithm Statements</b>	
<b>BICONCOMP Option</b>	
Specifies the log level for biconnected components	LOGLEVEL=
<b>CENTRALITY Options</b>	
Calculates authority centrality and specifies the type to process	AUTH=
Calculates betweenness centrality and specifies the type to process	BETWEEN=
Specifies whether to normalize the betweenness calculation	BETWEEN_NORM=
Decomposes the calculations for centrality by cluster (or subgraph)	BY_CLUSTER
Calculates closeness centrality and specifies the type to process	CLOSE=
Specifies the accounting method for no paths in closeness	CLOSE_NOPATH=
Calculates the node clustering coefficients	CLUSTERING_COEF
Calculates degree centrality and specifies the type to process	DEGREE=

**Table 1.2** (continued)

Description	Option
Calculates eigenvector centrality and specifies the type to process	EIGEN=
Specifies the algorithm to use for eigenvector calculation	EIGEN_ALGORITHM=
Specifies the maximum number of iterations for eigenvector calculation	EIGEN_MAXITER=
Calculates hub centrality and specifies the type to process	HUB=
Calculates influence centrality and specifies the type to process	INFLUENCE=
Specifies the iteration log frequency (nodes)	LOGFREQNODE=
Specifies the iteration log frequency (seconds)	LOGFREQTIME=
Specifies the log level for centrality	LOGLEVEL=
Specifies the subgraph node size to run separately	SUBSIZESWITCH=
Specifies the data set variable to use for weight2 in centrality	WEIGHT2=
<b>CLIQUE Options</b>	
Specifies the log level for clique calculations	LOGLEVEL=
Specifies the maximum number of cliques to return during clique calculations	MAXCLIQUES=
Specifies the maximum amount of time to spend calculating cliques	MAXTIME=
Specifies the output data set for cliques	OUT=
<b>COMMUNITY Options</b>	
Specifies the community detection algorithm	ALGORITHM=
Specifies the percentage of small-weight links to be removed	LINK_REMOVAL_RATIO=
Specifies the log level for community detection	LOGLEVEL=
Specifies the maximum number of iterations for community detection	MAXITER=
Specifies the output data set for inter-community links	OUT_COMM_LINKS=
Specifies the output data set for community summary table	OUT_COMMUNITY=
Specifies the output data set for community level summary table	OUT_LEVEL=
Specifies the output data set for community overlap table	OUT_OVERLAP=
Specifies the random factor in the parallel label propagation algorithm	RANDOM_FACTOR=
Specifies the random seed for the parallel label propagation algorithm	RANDOM_SEED=
Applies the recursive option to break large communities	RECURSIVE
Specifies the resolution list for community detection	RESOLUTION_LIST=
Specifies the modularity tolerance value for community detection	TOLERANCE=
<b>CONCOMP Options</b>	
Specifies the algorithm to use for connected components	ALGORITHM=
Specifies the log level for connected components	LOGLEVEL=
<b>CORE Options</b>	
Specifies the type of core to process	LINKS=
Specifies the log level for the core algorithm	LOGLEVEL=
Specifies the maximum amount of time to spend in the core algorithm	MAXTIME=
<b>CYCLE Options</b>	
Specifies the log level for the cycle algorithm	LOGLEVEL=
Specifies the maximum number of cycles to return during cycle calculations	MAXCYCLES=
Specifies the maximum length for the cycles found	MAXLENGTH=
Specifies the maximum link weight for the cycles found	MAXLINKWEIGHT=



**Table 1.2** (continued)

Description	Option
Specifies the maximum node weight for the cycles found	MAXNODEWEIGHT=
Specifies the maximum amount of time to spend calculating cycles	MAXTIME=
Specifies the minimum length for the cycles found	MINLENGTH=
Specifies the minimum link weight for the cycles found	MINLINKWEIGHT=
Specifies the minimum node weight for the cycles found	MINNODEWEIGHT=
Specifies the mode for the cycle calculations	MODE=
Specifies the output data set for cycles	OUT=
<b>EIGENVECTOR Options</b>	
Specifies the algebraic type of eigenvalues to calculate	EIGENVALUES=
Specifies the log level for eigenvector calculations	LOGLEVEL=
Specifies the maximum number of iterations for eigenvector calculation	MAXITER=
Specifies the number of eigenvectors to calculate	NEIGEN=
Specifies the output data set for eigenvectors	OUT=
<b>LINEAR_ASSIGNMENT Options</b>	
Specifies the data set variable names for the linear assignment identifiers	ID=( )
Specifies the log level for the linear assignment algorithm	LOGLEVEL=
Specifies the output data set for linear assignment	OUT=
Specifies the data set variable names for costs (or weights)	WEIGHT=( )
<b>MINCOSTFLOW Options</b>	
Specifies the iteration log frequency	LOGFREQ=
Specifies the log level for the minimum-cost network flow algorithm	LOGLEVEL=
Specifies the maximum amount of time to spend calculating the optimal flow	MAXTIME=
<b>MINCUT Options</b>	
Specifies the log level for the minimum-cut algorithm	LOGLEVEL=
Specifies the maximum number of cuts to return	MAXNUMCUTS=
Specifies the maximum weight of the cuts to return	MAXWEIGHT=
Specifies the output data set for minimum cut	OUT=
<b>MINSPANTREE Options</b>	
Specifies the log level for the minimum spanning tree algorithm	LOGLEVEL=
Specifies the output data set for minimum spanning tree	OUT=
<b>REACH Options</b>	
Decomposes the calculations for reach by cluster (or subgraph)	BY_CLUSTER
Calculates the directed reach counts	DIGRAPH
Treats each node as a source in reach calculations	EACH_SOURCE
Ignores the source node in reach counts	IGNORE_SELF
Specifies the maximum number of links to allow in the reach calculations	MAXREACH=
Specifies the iteration log frequency (seconds)	LOGFREQTIME=
Specifies the log level for reach calculations	LOGLEVEL=
Specifies the output data set for reach counts	OUT_COUNTS=
Specifies the output data set for reach counts (limit=1)	OUT_COUNTS1=

Table 1.2 (continued)

Description	Option
Specifies the output data set for reach counts (limit=2)	OUT_COUNTS2=
Specifies the output data set for reach links	OUT_LINKS=
Specifies the output data set for reach nodes	OUT_NODES=
<b>SHORTPATH Options</b>	
Specifies the iteration log frequency (nodes)	LOGFREQ=
Specifies the log level for shortest paths	LOGLEVEL=
Specifies the output data set for shortest paths	OUT_PATHS=
Specifies the output data set for shortest path summaries	OUT_WEIGHTS=
Specifies the type of output for shortest paths results	PATHS=
Specifies the sink node for shortest paths calculations	SINK=
Specifies the source node for shortest paths calculations	SOURCE=
Specifies whether to use weights in calculating shortest paths	USEWEIGHT=
Specifies the data set variable name for the auxiliary link weights	WEIGHT2=
<b>SUMMARY Options</b>	
Calculates information about biconnected components	BICONCOMP
Decomposes the calculations for summary by cluster (or subgraph)	BY_CLUSTER
Calculates information about connected components	CONCOMP
Calculates the approximate diameter and chooses the weight type	DIAMETER_APPROX=
Specifies the iteration log frequency (nodes)	LOGFREQNODE=
Specifies the iteration log frequency (seconds)	LOGFREQTIME=
Specifies the log level for summary calculations	LOGLEVEL=
Specifies the output data set for summary results	OUT=
Calculates information about shortest paths and chooses the weight type	SHORTPATH=
Specifies the subgraph node size to run separately	SUBSIZESWITCH=
<b>TRANSITIVE_CLOSURE Options</b>	
Specifies the log level for transitive closure	LOGLEVEL=
Specifies the output data set for transitive closure results	OUT=
<b>TSP Options</b>	
Specifies the stopping criterion based on the absolute objective gap	ABSOBJGAP=
Specifies the level of conflict search	CONFLICTSEARCH=
Specifies the cutoff value for branch-and-bound node removal	CUTOFF=
Specifies the overall cut strategy level	CUTSTRATEGY=
Emphasizes feasibility or optimality	EMPHASIS=
Specifies the initial and primal heuristics level	HEURISTICS=
Specifies the frequency of printing the branch-and-bound node log	LOGFREQ=
Specifies the log level for the traveling salesman algorithm	LOGLEVEL=
Specifies the maximum number of branch-and-bound nodes to be processed	MAXNODES=
Specifies the maximum number of solutions to be found	MAXSOLS=
Specifies the maximum amount of time to spend in the algorithm	MAXTIME=
Specifies whether to use a mixed integer linear programming solver	MILP=
Specifies the branch-and-bound node selection strategy	NODESEL=
Specifies the output data set for traveling salesman problem	OUT=



**Table 1.2** (continued)

Description	Option
Specifies the probing level	PROBE=
Specifies the stopping criterion based on the relative objective gap	RELOBJGAP=
Specifies the number of simplex iterations to be performed on each variable in the strong branching strategy	STRONGITER=
Specifies the number of candidates for the strong branching strategy	STRONGLLEN=
Specifies the stopping criterion based on the target objective value	TARGET=
Specifies the rule for selecting branching variable	VARSEL=

For more information about the options available for the PERFORMANCE statement, see the section “PERFORMANCE Statement” on page 38.

Table 1.3 lists the valid input formats, GRAPH\_DIRECTION= values, and GRAPH\_INTERNAL\_FORMAT= values for each statement in the OPTGRAPH procedure.

**Table 1.3** Supported Input Formats and Graph Types by Statement

Statement	Input Format		DIRECTION		INTERNAL_FORMAT	
	Graph	Matrix	UNDIRECTED	DIRECTED	THIN	FULL
BICONCOMP	X		X			X
CENTRALITY						
AUTH=, HUB=	X			X		X
EIGEN=	X		X			X
BETWEEN=, CLOSE=, CLUSTERING_COEF, DEGREE=, INFLUENCE=,	X		X	X		X
CENTRALITY / BY_CLUSTER						
AUTH=, HUB=	X			X	X	X
EIGEN=	X		X	X	X	X
BETWEEN=, CLOSE=, CLUSTERING_COEF, DEGREE=, INFLUENCE=,	X		X	X	X	X
CLIQUE	X		X			X
COMMUNITY						
ALGORITHM=						
LOUVAIN, LABEL_PROP	X		X		X	X
PARALLEL_LABEL_PROP	X		X	X	X	X
CONCOMP						
ALGORITHM=						
DFS	X		X	X		X
UNION_FIND	X		X		X	X
CORE	X		X	X		X

**Table 1.3** (continued)

Statement	Input Format		DIRECTION		INTERNAL_FORMAT	
	Graph	Matrix	UNDIRECTED	DIRECTED	THIN	FULL
CYCLE	X		X	X		X
EIGENVECTOR	X	X	X			X
LINEAR_ASSIGNMENT	X	X		X		X
MINCOSTFLOW	X			X	X	X
MINCUT	X		X			X
MINSPANTREE	X		X		X	X
REACH	X		X	X		X
REACH / BY_CLUSTER	X		X	X	X	X
SHORTPATH	X		X	X		X
SUMMARY	X		X	X		X
SUMMARY / BY_CLUSTER	X		X	X	X	X
TRANSITIVE_CLOSURE	X		X	X		X
TSP	X		X	X		X

Table 1.4 indicates for each algorithm statement in the OPTGRAPH procedure which output data set options you can specify and whether the algorithm populates the data sets specified in the **OUT\_NODES=** and **OUT\_LINKS=** options in the PROC OPTGRAPH statement.

**Table 1.4** Output Options by Statement

Statement	OUT_NODES	OUT_LINKS	Algorithm Statement Options
BICONCOMP	X	X	
CENTRALITY AUTH=, CLOSE=, CLUSTERING_COEF, DEGREE=, EIGEN=, HUB=, INFLUENCE=	X		
BETWEEN=	X	X	
CLIQUE			OUT=
COMMUNITY ALGORITHM= LOUVAIN, LABEL_PROP, PARALLEL_LABEL_PROP	X		OUT_COMM_LINKS=, OUT_COMMUNITY=, OUT_LEVEL=, OUT_OVERLAP=
CONCOMP	X		
CORE	X		
CYCLE			OUT=
EIGENVECTOR			OUT=
LINEAR_ASSIGNMENT			OUT=
MINCOSTFLOW		X	
MINCUT	X		OUT=

**Table 1.4** (continued)

Statement	OUT_NODES	OUT_LINKS	Algorithm Statement Options
MINSPANTREE			OUT=
REACH			OUT_COUNTS=, OUT_LINKS=, OUT_NODES=
BY_CLUSTER BY_CLUSTER and EACH_SOURCE			OUT_COUNTS=, OUT_NODES= OUT_COUNTS1=, OUT_COUNTS2=
SHORTPATH			OUT_PATHS=, OUT_WEIGHTS=
SUMMARY	X		OUT=
TRANSITIVE_CLOSURE			OUT=
TSP	X		OUT=

## PROC OPTGRAPH Statement

**PROC OPTGRAPH** < options > ;

The PROC OPTGRAPH statement invokes the OPTGRAPH procedure. You can specify the following *options* to define the input and output data sets, the log levels, and various other processing controls:

**DATA\_LINKS**=SAS-data-set

**LINKS**=SAS-data-set

specifies the input data set that contains the graph link information, where the links are defined as a list.

See the section “[Link Input Data](#)” on page 50 for more information.

**DATA\_MATRIX**=SAS-data-set

**MATRIX**=SAS-data-set

specifies the input data set that contains the matrix to be processed. This is a generic matrix (as opposed to an adjacency matrix, which defines an underlying graph).

See the section “[Matrix Input Data](#)” on page 57 for more information.

**DATA\_NODES**=SAS-data-set

**NODES**=SAS-data-set

specifies the input data set that contains the graph node information.

See the section “[Node Input Data](#)” on page 53 for more information.

**DATA\_NODES\_SUB**=SAS-data-set

**NODES\_SUB**=SAS-data-set

specifies the input data set that contains the graph node subset information.

See the section “[Node Subset Input Data](#)” on page 54 for more information.

**FILTER\_SUBGRAPH=number**

specifies the minimum number of nodes allowed in a subgraph when processing is decomposed by cluster. When the BY\_CLUSTER option is also specified in another statement, any subgraph whose number of nodes is less than or equal to *number* is skipped. The default setting is 0, so nothing is filtered by default.

See the section “[Graph Input Data](#)” on page 49 for more information.

**GRAPH\_DIRECTION=DIRECTED | UNDIRECTED****DIRECTION=DIRECTED | UNDIRECTED**

specifies whether the input graph should be considered directed or undirected.

**Table 1.5** Values for the GRAPH\_DIRECTION= Option

Option Value	Description
DIRECTED	Specifies the graph as directed. In a directed graph, each link $(i, j)$ has a direction that defines how something (for example, information) might flow over that link. In link $(i, j)$ , information flows from node $i$ to node $j$ ( $i \rightarrow j$ ). The node $i$ is called the <i>source (tail)</i> node, and $j$ is called the <i>sink (head)</i> node.
UNDIRECTED	Specifies the graph as undirected. In an undirected graph, each link $\{i, j\}$ has no direction and information can flow in either direction. That is, $\{i, j\} = \{j, i\}$ . This is the default.

By default, GRAPH\_DIRECTION=UNDIRECTED. See the section “[Graph Input Data](#)” on page 49 for more information.

**GRAPH\_INTERNAL\_FORMAT=FULL | THIN****INTERNAL\_FORMAT=FULL | THIN**

requests the internal graph format for the algorithms to use.

**Table 1.6** Values for the GRAPH\_INTERNAL\_FORMAT= Option

Option Value	Description
FULL	Stores the graph in standard (full) format. This is the default.
THIN	Stores the graph in thin format. This option can improve performance in some cases both by reducing memory and by simplifying the construction of the internal data structures. The thin format causes PROC OPTGRAPH to skip the removal of duplicate links when it reads in the graph. So this option should be used with caution. For some algorithms, the thin format is not allowed and this option is ignored. Setting GRAPH_INTERNAL_FORMAT=THIN can often be helpful when you do calculations that are decomposed by subgraph.

See the section “[Graph Input Data](#)” on page 49 for more information.

**INCLUDE\_SELFLINK**

includes self links—for example,  $(i, i)$ —when an input graph is read. By default, when PROC OPTGRAPH reads the [DATA\\_LINKS=](#) data set, it removes all self links.

**LOGLEVEL=***number* | *string*

controls the amount of information that is displayed in the SAS log. Each algorithm has its own specific log level. This setting sets the log level for all algorithms except those for which you specify the LOGLEVEL= option in the algorithm statement. [Table 1.7](#) describes the valid values for this option.

**Table 1.7** Values for LOGLEVEL= Option

<i>number</i>	<i>string</i>	Description
0	NONE	Turns off all procedure-related messages in the SAS log
1	BASIC	Displays a basic summary of the input, output, and algorithmic processing
2	MODERATE	Displays a summary of the input, output, and algorithmic processing
3	AGGRESSIVE	Displays a detailed summary of the input, output, and algorithmic processing

By default, LOGLEVEL=BASIC.

**OUT\_LINKS=***SAS-data-set*

specifies the output data set to contain the graph link information along with any results from the various algorithms that calculate metrics on links.

See the various algorithm sections for examples of the content of this output data set.

**OUT\_NODES=***SAS-data-set*

specifies the output data set to contain the graph node information along with any results from the various algorithms that calculate metrics on nodes.

See the various algorithm sections for examples of the content of this output data set.

**STANDARDIZED\_LABELS**

specifies that the input graph data is in a standardized format described in section “[Standardized Labels](#)” on page 55.

**TIMETYPE=***number* | *string*

specifies whether CPU time or real time is used for the MAXTIME= option for each applicable algorithm. [Table 1.8](#) describes the valid values of the TIMETYPE= option.

**Table 1.8** Values for TIMETYPE= Option

<i>number</i>	<i>string</i>	Description
0	CPU	Specifies units of CPU time
1	REAL	Specifies units of real time

By default, TIMETYPE=CPU.

## BICONCOMP Statement

**BICONCOMP** < option > ;

The BICONCOMP statement requests that PROC OPTGRAPH find biconnected components and articulation points of an undirected input graph.

See the section “[Biconnected Components and Articulation Points](#)” on page 61 for more information.

You can specify the following *option* in the BICONCOMP statement.

**LOGLEVEL**=*number* | *string*

controls the amount of information that is displayed in the SAS log. [Table 1.9](#) describes the valid values for this option.

**Table 1.9** Values for LOGLEVEL= Option

<i>number</i>	<i>string</i>	Description
0	NONE	Turns off all algorithm-related messages in the SAS log
1	BASIC	Displays a basic summary of the algorithmic processing
2	MODERATE	Displays a summary of the algorithmic processing
3	AGGRESSIVE	Displays a detailed summary of the algorithmic processing

The default is the value that is specified in the **LOGLEVEL=** option in the PROC OPTGRAPH statement (or BASIC if that option is not specified).

## CENTRALITY Statement

**CENTRALITY** < options > ;

The CENTRALITY statement enables you to select which centrality metrics to calculate for the given input graph. It also enables you to specify options for particular metrics. The resulting metrics are included in the node output data set (specified in the OUT\_NODES= option) or the link output data set (specified in the OUT\_LINKS= option).

The centrality metrics are described in the section “[Centrality](#)” on page 65.

You can specify the following *options* in the CENTRALITY statement.

**AUTH=WEIGHT** | **UNWEIGHT** | **BOTH**

specifies which type of authority centrality to calculate.

**Table 1.10** Values for the AUTH= Option

Option Value	Description
WEIGHT	Calculates authority centrality based on the weighted graph.
UNWEIGHT	Calculates authority centrality based on the unweighted graph.
BOTH	Calculates authority centrality based on both weighted and unweighted graphs.

If the input graph does not contain weights, then WEIGHT and UNWEIGHT both give the same results (using 1.0 for each link weight). This centrality metric can be used only for directed graphs. The authority centrality metric is described in the section “[Hub and Authority Scoring](#)” on page 76.

#### **BETWEEN=WEIGHT | UNWEIGHT | BOTH**

specifies which type of betweenness centrality to calculate.

**Table 1.11** Values for the BETWEEN= Option

Option Value	Description
WEIGHT	Calculates betweenness centrality based on the weighted graph.
UNWEIGHT	Calculates betweenness centrality based on the unweighted graph.
BOTH	Calculates betweenness centrality based on both weighted and unweighted graphs.

If the input graph does not contain weights, then WEIGHT and UNWEIGHT both give the same results (using 1.0 for each link weight). If the OUT\_NODES= option is specified in the PROC OPTGRAPH statement, the node betweenness metric is produced. If the OUT\_LINKS= option is specified, the link betweenness metric is produced. The betweenness centrality metric is described in the section “[Betweenness Centrality](#)” on page 72.

#### **BETWEEN\_NORM=YES | NO**

specifies whether to normalize the betweenness centrality metrics.

**Table 1.12** Values for the BETWEEN\_NORM= Option

Option Value	Description
YES	Normalizes the betweenness metrics. This is the default.
NO	Does not normalize the betweenness metrics.

The normalization factor for betweenness centrality is described in the section “[Betweenness Centrality](#)” on page 72.

#### **BY\_CLUSTER**

decomposes the calculations by cluster (or subgraph). If this option is specified, PROC OPTGRAPH looks for a definition of the clusters in the input data set specified by the DATA\_NODES= option in the PROC OPTGRAPH statement. The use of the BY\_CLUSTER option is described in the section “[Processing by Cluster](#)” on page 78.

#### **CLOSE=WEIGHT | UNWEIGHT | BOTH**

specifies which type of closeness centrality to calculate.

**Table 1.13** Values for the CLOSE= Option

Option Value	Description
WEIGHT	Calculates closeness centrality based on the weighted graph.
UNWEIGHT	Calculates closeness centrality based on the unweighted graph.
BOTH	Calculates closeness centrality based on both weighted and unweighted graphs.

If the input graph does not contain weights, then WEIGHT and UNWEIGHT both give the same results (using 1.0 for each link weight). The closeness centrality metric is described in the section “[Closeness Centrality](#)” on page 69.

**CLOSE\_NOPATH=NNODES | DIAMETER | ZERO | HARMONIC**

specifies a method for accounting for a shortest path distance between two nodes when a path does not exist (disconnected nodes).

**Table 1.14** Values for the CLOSE\_NOPATH= Option

Option Value	Description
NNODES	Uses the number of nodes as a shortest path distance between disconnected nodes. This option cannot be used when calculating weighted closeness centrality.
DIAMETER	Uses the graph diameter (plus one) as a shortest path distance between disconnected nodes. This is the default.
ZERO	Uses zero as a shortest path distance between disconnected nodes.
HARMONIC	Uses the harmonic formula for calculating closeness centrality.

For each option, there is a slight variation in the formula for the closeness centrality metric. These differences are described in the section “[Closeness Centrality](#)” on page 69.

**CLUSTERING\_COEF**

calculates the node clustering coefficient. The cluster coefficient is described in the section “[Clustering Coefficient](#)” on page 67.

**DEGREE=IN | OUT | BOTH**

specifies which type of degree centrality to calculate for the input graph.

**Table 1.15** Values for the DEGREE= Option

Option Value	Description
IN	Calculates degree based on in-links.
OUT	Calculates degree based on out-links.
BOTH	Calculates degree based on in-links and out-links.

For an undirected graph, the option values IN and BOTH are ignored, because there is only one notion of degree, which corresponds to the degree of out-links. The degree centrality metric is described in the section “[Degree Centrality](#)” on page 65.

**EIGEN=WEIGHT | UNWEIGHT | BOTH**

specifies which type of eigenvector centrality to calculate.

**Table 1.16** Values for the EIGEN= Option

Option Value	Description
WEIGHT	Calculates eigenvector centrality based on the weighted graph.
UNWEIGHT	Calculates eigenvector centrality based on the unweighted graph.



**Table 1.16** (continued)

Option Value	Description
BOTH	Calculates eigenvector centrality based on both weighted and unweighted graphs.

If the input graph does not contain weights, then WEIGHT and UNWEIGHT both give the same results (using 1.0 for each link weight). This centrality metric can be used only for undirected graphs. The eigenvector centrality metric is described in the section “[Eigenvector Centrality](#)” on page 74.

#### **EIGEN\_ALGORITHM=AUTOMATIC | JACOBI\_DAVIDSON | POWER**

specifies the algorithm to use in calculating centrality metrics that require solving eigensystems (EIGEN, HUB, and AUTH).

**Table 1.17** Values for the EIGEN\_ALGORITHM= Option

Option Value	Description
AUTOMATIC	Requests that PROC OPTGRAPH automatically determine the eigensolver to use. This is the default.
JACOBI_DAVIDSON (JD)	Uses a variant of the Jacobi-Davidson algorithm for solving eigensystems (Sleijpen and van der Vorst 2000). This is used as the default for the eigenvector metric on undirected graphs and the hub and authority metrics.
POWER	Uses the power method to calculate eigenvectors. This is used as the default for the eigenvector metric on directed graphs.

#### **EIGEN\_MAXITER=number**

specifies the maximum number of iterations to use for eigenvector calculations to limit the amount of computation time spent when convergence is slow. By default, EIGEN\_MAXITER=10,000.

#### **HUB=WEIGHT | UNWEIGHT | BOTH**

specifies which type of hub centrality to calculate.

**Table 1.18** Values for the HUB= Option

Option Value	Description
WEIGHT	Calculates hub centrality based on the weighted graph.
UNWEIGHT	Calculates hub centrality based on the unweighted graph.
BOTH	Calculates hub centrality based on both weighted and unweighted graphs.

If the input graph does not contain weights, then WEIGHT and UNWEIGHT both give the same results (using 1.0 for each link weight). This centrality metric can be used only for directed graphs. The hub centrality metric is described in the section “[Hub and Authority Scoring](#)” on page 76.

#### **INFLUENCE=WEIGHT | UNWEIGHT | BOTH**

specifies which type of influence centrality to calculate.

**Table 1.19** Values for the INFLUENCE= Option

Option Value	Description
WEIGHT	Calculates influence centrality based on the weighted graph.
UNWEIGHT	Calculates influence centrality based on the unweighted graph.
BOTH	Calculates influence centrality based on both weighted and unweighted graphs.

If the input graph does not contain weights, then WEIGHT and UNWEIGHT both give the same results (using 1.0 for each link weight). The influence centrality metric is described in the section “Influence Centrality” on page 66.

**LOGFREQNODE=number**

controls the frequency for displaying iteration logs for some of the centrality metrics. For computationally intensive algorithms such as betweenness and closeness centrality, this option displays progress every *number* nodes. If you also specify the BY\_CLUSTER option in this statement or a value greater than 1 for the NTHREADS= option in the PERFORMANCE statement, this option is ignored and the display frequency is determined by using the LOGFREQTIME= option instead. The value of *number* can be any integer greater than or equal to 1; the default is determined automatically based on the size of the graph. Setting this value too low can hurt performance on large-scale graphs.

**LOGFREQTIME=number**

controls the frequency for displaying iteration logs for some of the centrality metrics. For computationally intensive algorithms such as betweenness and closeness centrality, this option displays progress every *number* seconds. If you specify a value greater than 1 for the NTHREADS= option in the PERFORMANCE statement, PROC OPTGRAPH displays the number of nodes that have completed. If you specify the BY\_CLUSTER option, PROC OPTGRAPH displays the number of subgraphs that have completed. The value of *number* can be any integer greater than or equal to 1; the default is 5. Setting this value too low can hurt performance on large-scale graphs.

**LOGLEVEL=number | string**

controls the amount of information that is displayed in the SAS log. Table 1.20 describes the valid values for this option.

**Table 1.20** Values for LOGLEVEL= Option

<i>number</i>	<i>string</i>	Description
0	NONE	Turns off all algorithm-related messages in the SAS log
1	BASIC	Displays a basic summary of the algorithmic processing
2	MODERATE	Displays a summary of the algorithmic processing including a progress log using the interval that is specified in the LOGFREQNODE= or LOGFREQTIME= option
3	AGGRESSIVE	Displays a detailed summary of the algorithmic processing including a progress log using the interval that is specified in the LOGFREQNODE= or LOGFREQTIME= option

The default is the value that is specified in the LOGLEVEL= option in the PROC OPTGRAPH statement (or BASIC if that option is not specified).

**SUBSIZESWITCH=***number*

specifies the size of the subgraphs (number of nodes) to run separately when you also specify the BY\_CLUSTER option in this statement and a value greater than 1 for the NTHREADS= option in the [PERFORMANCE](#) statement. When PROC OPTGRAPH processes summary by subgraphs, it uses thread logic to simultaneously process *n* subgraphs, where *n* is the number of threads specified in the NTHREADS= option in the PERFORMANCE statement. Subgraphs that have more nodes than *number* are processed sequentially, enabling the threading to be done at the centrality metric level. The default is 10,000.

**WEIGHT2=***column*

specifies the data set variable name for a second link weight. The value of the variable named in *column* must be numeric. The use of this option is described in more detail in the section “[Weight Interpretation](#)” on page 78.

---

## CLIQUE Statement

**CLIQUE** < *options* > ;

The CLIQUE statement invokes an algorithm that finds maximal cliques on the input graph. Maximal cliques are described in the section “[Clique](#)” on page 86.

You can specify the following *options* in the CLIQUE statement:

**LOGLEVEL=***number* | *string*

controls the amount of information that is displayed in the SAS log. [Table 1.21](#) describes the valid values for this option.

**Table 1.21** Values for LOGLEVEL= Option

<i>number</i>	<i>string</i>	Description
0	NONE	Turns off all algorithm-related messages in the SAS log
1	BASIC	Displays a basic summary of the algorithmic processing
2	MODERATE	Displays a summary of the algorithmic processing
3	AGGRESSIVE	Displays a detailed summary of the algorithmic processing

The default is the value that is specified in the [LOGLEVEL=](#) option in the PROC OPTGRAPH statement (or BASIC if that option is not specified).

**MAXCLIQUES=***number*

specifies the maximum number of cliques to return during clique calculations. The default is the positive number that has the largest absolute value that can be represented in your operating environment.

**MAXTIME=***number*

specifies the maximum amount of time to spend calculating cliques. The type of time (either CPU time or real time) is determined by the value of the [TIMETYPE=](#) option. The value of *number* can be any positive number; the default value is the positive number that has the largest absolute value that can be represented in your operating environment.

**OUT=SAS-data-set**

specifies the output data set to contain the maximal cliques.

## COMMUNITY Statement

**COMMUNITY** < options > ;

The COMMUNITY statement invokes an algorithm that detects communities of the input graph. Community detection is described in the section “Community Detection” on page 89.

You can specify the following *options* in the COMMUNITY statement:

**ALGORITHM=LOUVAIN | LABEL\_PROP | PARALLEL\_LABEL\_PROP**

specifies whether to use the Louvain algorithm (LOUVAIN), the label propagation algorithm (LABEL\_PROP), or the parallel label propagation algorithm (PARALLEL\_LABEL\_PROP). The Louvain algorithm is the default.

For more information about this option, see the sections “Community Detection” on page 89 and “Parallel Community Detection” on page 91.

**LINK\_REMOVAL\_RATIO=number**

defines the percentage of small-weight links to be removed around each node neighborhood. A link is usually removed if its weight is relatively smaller than the weights of neighboring links. Suppose that node *A* links to node *B* and to node *C*, link  $A \rightarrow B$  has weight of 100, and link  $A \rightarrow C$  has weight of 1. When nodes are grouped into communities, link  $A \rightarrow B$  is much more important than link  $A \rightarrow C$  because it contributes much more to the overall modularity value. Therefore, link  $A \rightarrow C$  can be dropped from the network if dropping it does not disconnect node *C* from the network. If the LINK\_REMOVAL\_RATIO= option is specified, then the links that are incident to each node are examined. If the weight of any link is less than  $(number/100) * max\_link\_weight$ , where *max\_link\_weight* is the maximum link weight among all links incident to this node, it is removed provided that its removal does not disconnect any node from the network. This option can often dramatically improve the running time of large graphs. The valid range is between 0 and 100. The default value is 10.

**LOGLEVEL=number | string**

controls the amount of information that is displayed in the SAS log. Table 1.22 describes the valid values for this option.

**Table 1.22** Values for LOGLEVEL= Option

<i>number</i>	<i>string</i>	Description
0	NONE	Turns off all algorithm-related messages in the SAS log
1	BASIC	Displays a basic summary of the algorithmic processing
2	MODERATE	Displays a summary of the algorithmic processing
3	AGGRESSIVE	Displays a detailed summary of the algorithmic processing

The default is the value that you specify in the LOGLEVEL= option in the PROC OPTGRAPH statement (or BASIC if that option is not specified).

**MAXITER=number**

specifies the maximum number of iterations allowed in the algorithm. The default is 20 when **ALGORITHM=LOUVAIN** and 100 when **ALGORITHM=LABEL\_PROP** or **ALGORITHM=PARALLEL\_LABEL\_PROP**.

**OUT\_COMM\_LINKS=SAS-data-set**

specifies the output data set that describes the links between communities.

**OUT\_COMMUNITY=SAS-data-set**

specifies the output data set that contains the number of nodes in each community.

**OUT\_LEVEL=SAS-data-set**

specifies the output data set that contains community information at different resolution levels.

**OUT\_OVERLAP=SAS-data-set**

specifies the output data set that describes the intensity of each node.

**RANDOM\_FACTOR=number**

specifies the random factor for the parallel label propagation algorithm. Specify a *number* between 0 and 1. At each iteration,  $number \times 100\%$  of the nodes are randomly selected to skip the label propagation step. The default is 0.15, which means that 15% of nodes skip the label propagation step at each iteration.

**RANDOM\_SEED=number**

specifies the random seed for the parallel label propagation algorithm. At each iteration, some nodes are randomly selected to skip the label propagation step, based on the value that you specify in the **RANDOM\_FACTOR=** option. To choose a different set of random samples, specify a *number* in the **RANDOM\_SEED=** option. By default, **RANDOM\_SEED=1234**.

**RECURSIVE (options)**

requests that the algorithm recursively break down large communities into smaller ones until the specified conditions are satisfied. This option starts with the keyword **RECURSIVE** followed by any combination of three suboptions enclosed in parentheses—for example, **RECURSIVE (MAX\_COMM\_SIZE=500)** or **RECURSIVE (MAX\_COMM\_SIZE=1000 MAX\_DIAMETER=3 RELATION=AND)**.

**Table 1.23** RECURSIVE options

option	Description
<b>MAX_COMM_SIZE=</b>	Specifies the maximum number of nodes to be contained in any community.
<b>MAX_DIAMETER=</b>	Specifies the maximum number of links on the shortest paths between any pair of nodes in any community.
<b>RELATION=</b>	Specifies the relationship between the values of <b>MAX_COMM_SIZE=</b> and <b>MAX_DIAMETER=</b> options. If <b>RELATION=AND</b> , then recursive splitting continues until both <b>MAX_COMM_SIZE</b> and <b>MAX_DIAMETER</b> conditions are satisfied. If <b>RELATION=OR</b> , then recursive splitting continues until either the <b>MAX_COMM_SIZE</b> or the <b>MAX_DIAMETER</b> condition is satisfied. The valid values are <b>AND</b> and <b>OR</b> . The default is <b>OR</b> .

The MAX\_DIAMETER= option is ignored when you specify **ALGORITHM=PARALLEL\_LABEL\_PROP**.

**RESOLUTION\_LIST=***num\_list*

specifies a list of resolution values that are separated by spaces (for example, 4.3 2.1 1.0 0.6 0.2). The OPTGRAPH procedure interprets the RESOLUTION\_LIST= option differently depending on the value of the **ALGORITHM=** option:

- When **ALGORITHM=LOUVAIN**, specifying multiple resolution values enables you to see how communities are merged at various resolution levels. A larger parameter value indicates a higher resolution. For example, resolution 4.3 produces more communities than resolution 0.2. By default, RESOLUTION\_LIST=1.0. When you also specify the **RECURSIVE** option, the first value in the resolution list is used and the other values are ignored.
- When **ALGORITHM=LABEL\_PROP**, PROC OPTGRAPH ignores the RESOLUTION\_LIST= option. It uses the default value of 1.0.
- When **ALGORITHM=PARALLEL\_LABEL\_PROP**, specifying multiple resolution values requests that the OPTGRAPH procedure perform community detection multiple times, each time with a different resolution value. By default, RESOLUTION\_LIST=0.001. In this case, the RESOLUTION\_LIST= option is fully compatible with the **RECURSIVE** option.

For more information about the use of the RESOLUTION\_LIST= option, see the section “[Large Community](#)” on page 92.

**TOLERANCE=***number*

**MODULARITY=***number*

specifies the tolerance value for when to stop iterations. When you specify **ALGORITHM=LOUVAIN**, the algorithm stops iterations when the percentage modularity gain between two consecutive iterations falls within the specified tolerance value. When you specify **ALGORITHM=LABEL\_PROP** or **ALGORITHM=PARALLEL\_LABEL\_PROP**, the algorithm stops iterations when the percentage of label changes for all nodes in the graph falls within the tolerance specified by *number*. The valid range is strictly between 0 and 1. By default, TOLERANCE=0.01.

---

## CONCOMP Statement

**CONCOMP** < *options* > ;

The CONCOMP statement invokes an algorithm that finds the connected components of the input graph. Connected components are described in the section “[Connected Components](#)” on page 97.

You can specify the following *options* in the CONCOMP statement:

**ALGORITHM=DFS | UNION\_FIND**

specifies the algorithm to use for calculating connected components.

**Table 1.24** Values for the ALGORITHM= Option

Option Value	Description
DFS	Uses the depth-first search algorithm for connected components. You cannot specify this value when you specify GRAPH_INTERNAL_FORMAT=THIN in the PROC OPTGRAPH statement.
UNION_FIND	Uses the union-find algorithm for connected components. You can specify this value with either the THIN or FULL value for the GRAPH_INTERNAL_FORMAT= option in the PROC OPTGRAPH statement. This value can be faster than DFS when used with GRAPH_INTERNAL_FORMAT=THIN. However, you can use it only with undirected graphs.

By default, ALGORITHM=UNION\_FIND for undirected graphs, and ALGORITHM=DFS for directed graphs.

**LOGLEVEL=***number* | *string*

controls the amount of information that is displayed in the SAS log. [Table 1.25](#) describes the valid values for this option.

**Table 1.25** Values for LOGLEVEL= Option

<i>number</i>	<i>string</i>	Description
0	NONE	Turns off all algorithm-related messages in the SAS log
1	BASIC	Displays a basic summary of the algorithmic processing
2	MODERATE	Displays a summary of the algorithmic processing
3	AGGRESSIVE	Displays a detailed summary of the algorithmic processing

The default is the value that is specified in the **LOGLEVEL=** option in the PROC OPTGRAPH statement (or BASIC if that option is not specified).

## CORE Statement

**CORE** < *options* > ;

The CORE statement invokes an algorithm that finds the core decomposition of the input graph. Core decompositions are described in the section “[Core Decomposition](#)” on page 102.

You can specify the following *options* in the CORE statement:

**LINKS=**IN | OUT | BOTH

specifies which type of cores to calculate for a directed graph. You can choose to calculate the cores based on in-links (IN), out-links (OUT), or both (BOTH). For an undirected graph, core applies only to out-links.

**Table 1.26** Values for the LINKS= Option

Option Value	Description
IN	Calculates core based on in-links.
OUT	Calculates core based on out-links. This is the default.
BOTH	Calculates core based on in-links and out-links.

**LOGLEVEL=***number* | *string*

controls the amount of information that is displayed in the SAS log. [Table 1.27](#) describes the valid values for this option.

**Table 1.27** Values for LOGLEVEL= Option

<i>number</i>	<i>string</i>	Description
0	NONE	Turns off all algorithm-related messages in the SAS log
1	BASIC	Displays a basic summary of the algorithmic processing
2	MODERATE	Displays a summary of the algorithmic processing
3	AGGRESSIVE	Displays a detailed summary of the algorithmic processing

The default is the value that is specified in the **LOGLEVEL=** option in the PROC OPTGRAPH statement (or BASIC if that option is not specified).

**MAXTIME=***number*

specifies the maximum amount of time to spend in the core decomposition algorithm. The type of time (either CPU time or real time) is determined by the value of the **TIMETYPE=** option. The value of *number* can be any positive number; the default value is the positive number that has the largest absolute value that can be represented in your operating environment.

## CYCLE Statement

**CYCLE** < *options* > ;

The CYCLE statement invokes an algorithm that finds the cycles (or the existence of a cycle) in the input graph. Cycles are described in the section “[Cycle](#)” on page 107.

You can specify the following *options* in the CYCLE statement:

**LOGLEVEL=***number* | *string*

controls the amount of information that is displayed in the SAS log. [Table 1.28](#) describes the valid values for this option.

**Table 1.28** Values for LOGLEVEL= Option

<i>number</i>	<i>string</i>	Description
0	NONE	Turns off all algorithm-related messages in the SAS log
1	BASIC	Displays a basic summary of the algorithmic processing
2	MODERATE	Displays a summary of the algorithmic processing
3	AGGRESSIVE	Displays a detailed summary of the algorithmic processing



The default is the value that is specified in the **LOGLEVEL=** option in the PROC OPTGRAPH statement (or BASIC if that option is not specified).

**MAXCYCLES=***number*

specifies the maximum number of cycles to return. The default is the positive number that has the largest absolute value representable in your operating environment. This option works only when you also specify **MODE=ALL\_CYCLES**.

**MAXLENGTH=***number*

specifies the maximum number of links to allow in a cycle. Any cycle whose length is greater than *number* is removed from the results. The default is the positive number that has the largest absolute value that can be represented in your operating environment. By default, nothing is removed from the results. This option works only when you also specify **MODE=ALL\_CYCLES**.

**MAXLINKWEIGHT=***number*

specifies the maximum sum of link weights to allow in a cycle. Any cycle whose sum of link weights is greater than *number* is removed from the results. The default is the positive number that has the largest absolute value that can be represented in your operating environment. By default, nothing is filtered. This option works only when you also specify **MODE=ALL\_CYCLES**.

**MAXNODEWEIGHT=***number*

specifies the maximum sum of node weights to allow in a cycle. Any cycle whose sum of node weights is greater than *number* is removed from the results. The default is the positive number that has the largest absolute value that can be represented in your operating environment. By default, nothing is filtered. This option works only when you also specify **MODE=ALL\_CYCLES**.

**MAXTIME=***number*

specifies the maximum amount of time to spend finding cycles. The type of time (either CPU time or real time) is determined by the value of the **TIMETYPE=** option. The value of *number* can be any positive number; the default value is the positive number that has the largest absolute value that can be represented in your operating environment. This option works only when you also specify **MODE=ALL\_CYCLES**.

**MINLENGTH=***number*

specifies the minimum number of links to allow in a cycle. Any cycle that has fewer links than *number* is removed from the results. By default **MINLENGTH=1** and nothing is filtered. This option works only when you also specify **MODE=ALL\_CYCLES**.

**MINLINKWEIGHT=***number*

specifies the minimum sum of link weights to allow in a cycle. Any cycle whose sum of link weights is less than *number* is removed from the results. The default is the negative number that has the largest absolute value that can be represented in your operating environment. By default, nothing is filtered. This option works only when you also specify **MODE=ALL\_CYCLES**.

**MINNODEWEIGHT=***number*

specifies the minimum sum of node weights to allow in a cycle. Any cycle whose sum of node weights is less than *number* is removed from the results. The default is the negative number that has the largest absolute value that can be represented in your operating environment. By default, nothing is filtered. This option works only when you also specify **MODE=ALL\_CYCLES**.

**MODE=ALL\_CYCLES | FIRST\_CYCLE**

specifies the mode for processing cycles.

**Table 1.29** Values for the MODE= Option

Option Value	Description
ALL_CYCLES	Returns all (unique, elementary) cycles found.
FIRST_CYCLE	Returns the first cycle found. This is the default.

**OUT=SAS-data-set**

specifies the output data set to contain the cycles found.

---

## DATA\_LINKS\_VAR Statement

**DATA\_LINKS\_VAR** < options > ;

**LINKS\_VAR** < options > ;

The DATA\_LINKS\_VAR statement enables you to explicitly define the data set variable names for PROC OPTGRAPH to use when it reads the data set that is specified in the [DATA\\_LINKS=](#) option in the PROC OPTGRAPH statement. The format of the links input data set is defined in the section “[Link Input Data](#)” on page 50.

You can specify the following *options* in the DATA\_LINKS\_VAR statement:

**FROM=column**

specifies the data set variable name for the *from* nodes. The value of the *column* variable can be numeric or character.

**LOWER=column**

specifies the data set variable name for the link flow lower bounds. The value of the *column* variable must be numeric.

**TO=column**

specifies the data set variable name for the *to* nodes. The value of the *column* variable can be numeric or character.

**UPPER=column**

specifies the data set variable name for the link flow upper bounds. The value of the *column* variable must be numeric.

**WEIGHT=column**

specifies the data set variable name for the link weights. The value of the *column* variable must be numeric.

---

## DATA\_MATRIX\_VAR Statement

**DATA\_MATRIX\_VAR** <column1,column2,...> ;

**MATRIX\_VAR** <column1,column2,...> ;

The DATA\_MATRIX\_VAR statement enables you to explicitly define the data set variable names for PROC OPTGRAPH to use when it reads the data set that is specified in the DATA\_MATRIX= option in the PROC OPTGRAPH statement. The format of the matrix input data set is defined in the section “[Matrix Input Data](#)” on page 57. The value of each *column* variable must be numeric.

---

## DATA\_NODES\_VAR Statement

**DATA\_NODES\_VAR** < options > ;

**NODES\_VAR** < options > ;

The DATA\_NODES\_VAR statement enables you to explicitly define the data set variable names for PROC OPTGRAPH to use when it reads the data set that is specified in the DATA\_NODES= option in the PROC OPTGRAPH statement. The format of the node input data set is defined in the section “[Node Input Data](#)” on page 53.

You can specify the following *options* in the DATA\_NODES\_VAR statement:

**CLUSTER=column**

specifies the data set variable name for clusters identifiers. The value of the *column* variable must be numeric.

**NODE=column**

specifies the data set variable name for the nodes. The value of the *column* variable can be numeric or character.

**WEIGHT=column**

specifies the data set variable name for node weights. The value of the *column* variable must be numeric.

**WEIGHT2=column**

specifies the data set variable name for auxiliary node weights. The value of the *column* variable must be numeric.

---

## EIGENVECTOR Statement

**EIGENVECTOR** < options > ;

**EIGEN** < options > ;

The EIGENVECTOR statement invokes a variant of the Jacobi-Davidson algorithm (Sleijpen and van der Vorst 2000) that finds eigenvectors (and eigenvalues) for symmetric matrices. The matrix is typically defined in the input data set that is specified in the DATA\_MATRIX= option in the PROC OPTGRAPH statement.

The matrix can also be input as a graph by using the `DATA_LINKS=` option in the PROC OPTGRAPH statement. Internally, the graph is converted into a (sparse) adjacency matrix.

Eigenvectors and eigenvalues are described in the section “[Eigenvector Problem](#)” on page 112.

You can specify the following *options* in the EIGENVECTOR statement:

#### **EIGENVALUES=LA | SA**

specifies the type of eigenvector to calculate. [Table 1.30](#) describes the valid values for this option.

**Table 1.30** Values for the EIGENVALUES= Option

Option Value	Description
LA	Calculates the $n$ largest algebraic eigenvalues (and their corresponding eigenvectors), where $n$ is the value of the NEIGEN= option. This is the default.
SA	Calculates the $n$ smallest algebraic eigenvalues (and their corresponding eigenvectors), where $n$ is the value of the NEIGEN= option.

#### **LOGLEVEL=number | string**

controls the amount of information that is displayed in the SAS log. [Table 1.31](#) describes the valid values for this option.

**Table 1.31** Values for LOGLEVEL= Option

<i>number</i>	<i>string</i>	Description
0	NONE	Turns off all algorithm-related messages in the SAS log
1	BASIC	Displays a basic summary of the algorithmic processing
2	MODERATE	Displays a summary of the algorithmic processing
3	AGGRESSIVE	Displays a detailed summary of the algorithmic processing

The default is the value that is specified in the `LOGLEVEL=` option in the PROC OPTGRAPH statement (or BASIC if that option is not specified).

#### **MAXITER=number**

specifies the maximum number of matrix-vector multiplications used in the Jacobi-Davidson algorithm to calculate eigenvectors. By default, MAXITER=10,000.

#### **NEIGEN=number**

specifies the number of eigenvalues (and their corresponding eigenvectors) to generate. This value must be less than or equal to the dimension of the matrix. By default, NEIGEN=1.

#### **OUT=SAS-data-set**

specifies the output data set to contain the eigenvectors (and eigenvalues) found.

## **LINEAR\_ASSIGNMENT Statement**

**LINEAR\_ASSIGNMENT** < options > ;

**LAP** < options > ;

The **LINEAR\_ASSIGNMENT** statement invokes an algorithm that solves the minimal-cost linear assignment problem. In graph terms, this problem is also known as the minimum link-weighted matching problem on a bipartite graph. The input data (the cost matrix) is typically defined in the input data set that is specified in the **DATA\_MATRIX=** option in the **PROC OPTGRAPH** statement. The data can also be defined as a directed graph by specifying the **DATA\_LINKS=** option in the **PROC OPTGRAPH** statement, where the costs are defined as link weights. Internally, the graph is treated as a bipartite graph in which the *from* nodes define one part and the *to* nodes define the other part.

The linear assignment problem is described in the section “[Linear Assignment \(Matching\)](#)” on page 115.

You can specify the following *options* in the **LINEAR\_ASSIGNMENT** statement:

**ID=**(*<column1,column2,...>*)

specifies the data set variable names that identify the matrix rows (*from* nodes). The information in these columns is carried to the output data set that is specified in the **OUT=** option. The value of each *column* variable can be numeric or character.

**LOGLEVEL=***number* | *string*

controls the amount of information that is displayed in the SAS log. [Table 1.32](#) describes the valid values for this option.

**Table 1.32** Values for LOGLEVEL= Option

<i>number</i>	<i>string</i>	Description
0	NONE	Turns off all algorithm-related messages in the SAS log
1	BASIC	Displays a basic summary of the algorithmic processing
2	MODERATE	Displays a summary of the algorithmic processing
3	AGGRESSIVE	Displays a detailed summary of the algorithmic processing

The default is the value that is specified in the **LOGLEVEL=** option in the **PROC OPTGRAPH** statement (or BASIC if that option is not specified).

**OUT=***SAS-data-set*

specifies the output data set to contain the solution to the linear assignment problem.

**WEIGHT=**(*<column1,column2,...>*)

specifies the data set variable names for the cost matrix. The value of each *column* variable must be numeric. If this option is not specified, the matrix is assumed to be defined by all of the numeric variables in the data set (excluding those specified in the **ID=** option).

---

## MINCOSTFLOW Statement

**MINCOSTFLOW** *< options >* ;

**MCF** *< options >* ;

The **MINCOSTFLOW** statement invokes an algorithm that solves the minimum-cost network flow problem on an input graph.

The minimum-cost network flow problem is described in the section “[Minimum-Cost Network Flow](#)” on page 116.

You can specify the following *options* in the MINCOSTFLOW statement:

**LOGFREQ=number**

controls the frequency for displaying iteration logs for minimum-cost network flow calculations that use the network simplex algorithm. For graphs that contain one component, this option displays progress every *number* simplex iterations, and the default is 10,000. For graphs that contain multiple components, when you also specify LOGLEVEL=MODERATE, this option displays progress after processing every *number* components, and the default is based on the number of components. When you also specify LOGLEVEL=AGGRESSIVE, the simplex iteration log for each component is displayed with frequency *number*.

The value of *number* can be any integer greater than or equal to 1. Setting this value too low can hurt performance on large-scale graphs.

**LOGLEVEL=number | string**

controls the amount of information that is displayed in the SAS log. [Table 1.33](#) describes the valid values for this option.

**Table 1.33** Values for LOGLEVEL= Option

<i>number</i>	<i>string</i>	Description
0	NONE	Turns off all algorithm-related messages in the SAS log
1	BASIC	Displays a basic summary of the algorithmic processing
2	MODERATE	Displays a summary of the algorithmic processing including a progress log using the interval that is specified in the LOGFREQ= option
3	AGGRESSIVE	Displays a detailed summary of the algorithmic processing including a progress log using the interval that is specified in the LOGFREQ= option

The default is the value that is specified in the LOGLEVEL= option in the PROC OPTGRAPH statement (or BASIC if that option is not specified).

**MAXTIME=number**

specifies the maximum amount of time to spend calculating minimum-cost network flows. The type of time (either CPU time or real time) is determined by the value of the TIMETYPE= option. The value of *number* can be any positive number; the default value is the positive number that has the largest absolute value that can be represented in your operating environment.

## MINCUT Statement

**MINCUT** < options > ;

The MINCUT statement invokes an algorithm that finds the minimum link-weighted cut of an input graph.

The minimum-cut problem is described in the section “[Minimum Cut](#)” on page 124.

You can specify the following *options* in the MINCUT statement:

**LOGLEVEL=***number* | *string*

controls the amount of information that is displayed in the SAS log. [Table 1.34](#) describes the valid values for this option.

**Table 1.34** Values for LOGLEVEL= Option

<i>number</i>	<i>string</i>	Description
0	NONE	Turns off all algorithm-related messages in the SAS log
1	BASIC	Displays a basic summary of the algorithmic processing
2	MODERATE	Displays a summary of the algorithmic processing
3	AGGRESSIVE	Displays a detailed summary of the algorithmic processing

The default is the value that is specified in the **LOGLEVEL=** option in the PROC OPTGRAPH statement (or BASIC if that option is not specified).

**MAXNUMCUTS=***number*

specifies the maximum number of cuts to return from the algorithm. The minimal cut and any others found during the search, up to *number*, are returned. By default, MAXNUMCUTS=1.

**MAXWEIGHT=***number*

specifies the maximum weight of the cuts to return from the algorithm. Only cuts that have weight less than or equal to *number* are returned. The default is the positive number that has the largest absolute value that can be represented in your operating environment.

**OUT=***SAS-data-set*

specifies the output data set to contain the solution to the minimum-cut problem.

## MINSPANTREE Statement

**MINSPANTREE** < *options* > ;

**MST** < *options* > ;

The MINSPANTREE statement invokes an algorithm that solves the minimum link-weighted spanning tree problem on an input graph.

The minimum spanning tree problem is described in the section “[Minimum Spanning Tree](#)” on page 128.

You can specify the following *options* in the MINSPANTREE statement:

**LOGLEVEL=***number* | *string*

controls the amount of information that is displayed in the SAS log. [Table 1.35](#) describes the valid values for this option.

**Table 1.35** Values for LOGLEVEL= Option

<i>number</i>	<i>string</i>	Description
0	NONE	Turns off all algorithm-related messages in the SAS log
1	BASIC	Displays a basic summary of the algorithmic processing
2	MODERATE	Displays a summary of the algorithmic processing
3	AGGRESSIVE	Displays a detailed summary of the algorithmic processing

The default is the value that is specified in the **LOGLEVEL=** option in the PROC OPTGRAPH statement (or BASIC if that option is not specified).

**OUT=SAS-data-set**

specifies the output data set to contain the solution to the minimum link-weighted spanning tree problem.

---

## PERFORMANCE Statement

**PERFORMANCE** < *performance-options* > ;

The PERFORMANCE statement specifies performance options for multithreaded computing and requests detailed results about the performance characteristics of the OPTGRAPH procedure.

The PERFORMANCE statement enables you to control the number of threads used and the output of the ODS table that reports procedure timing. When you specify the PERFORMANCE statement, the PerformanceInfo ODS table is produced. This table lists performance characteristics such as execution mode and number of threads.

You can specify the following *performance-options* in the PERFORMANCE statement:

### DETAILS

requests that PROC OPTGRAPH produce the Timing ODS table, which shows a breakdown of the time used in each step of the procedure.

### NTHREADS=*number* | CPUCOUNT

specifies the number of threads that PROC OPTGRAPH can use. This option overrides the SAS system option THREADS | NOTHREADS. The value of *number* can be any integer between 1 and 256, inclusive. The default value is CPUCOUNT, which sets the thread count to the number determined by the SAS system option CPUCOUNT=.

Setting this option to a number greater than the number of available cores might result in reduced performance. Specifying a high *number* does not guarantee shorter solution time; the actual change in solution time depends on the computing hardware and the scalability of the underlying algorithms in the PROC OPTGRAPH. In some circumstances, the OPTGRAPH procedure might use fewer threads than the specified *number* because the procedure's internal algorithms have determined that a smaller number is preferable.

For example, the following call to PROC OPTGRAPH uses eight threads to read the data input in parallel:

```
proc optgraph
  data_links      = LinkSetIn
  graph_direction = directed
  out_nodes       = NodeSetOut;
  performance
    nthreads      = 8;
run;
```



## REACH Statement

**REACH** < options > ;

The REACH statement invokes an algorithm that calculates the reach (ego) network on an input graph.

The reach network is described in the section “[Reach \(Ego\) Network](#)” on page 130.

You can specify the following *options* in the REACH statement:

### BY\_CLUSTER

decomposes the calculations by cluster (subgraph). If this option is specified, PROC OPTGRAPH looks for a definition of the clusters in the input data set specified in the DATA\_NODES= option in the PROC OPTGRAPH statement. If BY\_CLUSTER is specified, the reach network links output (specified in the OUT\_LINKS= option) cannot be generated.

### DIGRAPH

calculates the directed reach counts when computing the reach networks and includes the directed counts in the resulting output data set that is specified in the OUT\_COUNTS= option. This option is ignored unless you specify MAXREACH=1 in the REACH statement.

### EACH\_SOURCE

treats each node as a source and calculates a reach network from each one.

### IGNORE\_SELF

ignores the source nodes in the reach network node counts.

### MAXREACH=*number*

specifies the maximum number of links to allow from each source node in a reach network. By default, MAXREACH=1.

### LOGFREQTIME=*number*

displays iteration logs for the reach algorithm every *number* seconds. When PROC OPTGRAPH runs the reach algorithm, it displays the number of source networks that have completed. When you also specify the BY\_CLUSTER option in the REACH statement, PROC OPTGRAPH displays the number of subgraphs that have completed. The value of *number* can be any integer greater than or equal to 1; the default is 5. Setting this value too low can hurt performance on large-scale graphs.

### LOGLEVEL=*number*

controls the amount of information that is displayed in the SAS log. [Table 1.36](#) describes the valid values for this option.

**Table 1.36** Values for LOGLEVEL= Option

<i>number</i>	<i>string</i>	Description
0	NONE	Turns off all algorithm-related messages in the SAS log
1	BASIC	Displays a basic summary of the algorithmic processing
2	MODERATE	Displays a summary of the algorithmic processing
3	AGGRESSIVE	Displays a detailed summary of the algorithmic processing

The default is the value that is specified in the **LOGLEVEL=** option in the PROC OPTGRAPH statement (or BASIC if that option is not specified).

**OUT\_COUNTS=SAS-data-set**

specifies the output data set to contain the node counts in each reach network.

**OUT\_COUNTS1=SAS-data-set**

specifies the output data set to contain the node counts in each reach network for the special case of calculating only counts that have limit 1 and 2. This data set holds the counts with MAXREACH=1. This option works only when the EACH\_SOURCE and BY\_CLUSTER options are specified.

**OUT\_COUNTS2=SAS-data-set**

specifies the output data set to contain the node counts in each reach network for the special case of calculating only counts that have limit 1 and 2. This data set holds the counts with MAXREACH=2. This option works only when the EACH\_SOURCE and BY\_CLUSTER options are specified.

**OUT\_LINKS=SAS-data-set**

specifies the output data set to contain the links in each reach network.

**OUT\_NODES=SAS-data-set**

specifies the output data set to contain the nodes in each reach network.

---

## SHORTPATH Statement

**SHORTPATH** < options > ;

The SHORTPATH statement invokes an algorithm that calculates shortest paths between sets of nodes on the input graph.

The shortest path algorithm is described in the section “[Shortest Path](#)” on page 142.

You can specify the following *options* in the SHORTPATH statement:

**LOGFREQ=number**

displays iteration logs for shortest path calculations every *number* nodes. The value of *number* can be any integer greater than or equal to 1. The default is determined automatically based on the size of the graph. Setting this value too low can hurt performance on large-scale graphs.

**LOGLEVEL=number**

controls the amount of information that is displayed in the SAS log. [Table 1.37](#) describes the valid values for this option.

**Table 1.37** Values for LOGLEVEL= Option

<i>number</i>	<i>string</i>	<b>Description</b>
0	NONE	Turns off all algorithm-related messages in the SAS log
1	BASIC	Displays a basic summary of the algorithmic processing
2	MODERATE	Displays a summary of the algorithmic processing
3	AGGRESSIVE	Displays a detailed summary of the algorithmic processing

The default is the value that is specified in the **LOGLEVEL=** option in the PROC OPTGRAPH statement (or BASIC if that option is not specified).

**OUT\_PATHS=SAS-data-set**

**OUT=SAS-data-set**

specifies the output data set to contain the shortest paths.

**OUT\_WEIGHTS=SAS-data-set**

specifies the output data set to contain the shortest path summaries.

**PATHS=ALL | LONGEST | SHORTEST**

specifies the type of output to produce in the output data set that is specified in the OUT\_PATHS= option.

**Table 1.38** Values for the PATHS= Option

Option Value	Description
ALL	Outputs shortest paths for all pairs of source-sinks.
LONGEST	Outputs shortest paths for the source-sink pair with the longest (finite) length. If other source-sink pairs (up to 100) have equally long length, they are also output.
SHORTEST	Outputs shortest paths for the source-sink pair with the shortest length. If other source-sink pairs (up to 100) have equally short length, they are also output.

By default, PATHS=ALL.

**SINK=sink-node**

specifies the sink node for shortest paths calculations. This setting overrides the use of the variable sink in the data set that is specified in the DATA\_NODES\_SUB= option in the PROC OPTGRAPH statement.

**SOURCE=source-node**

specifies the source node for shortest paths calculations. This setting overrides the use of the variable source in the data set that is specified in the DATA\_NODES\_SUB= option in the PROC OPTGRAPH statement.

**USEWEIGHT=YES | NO**

specifies whether to use link weights (if they exist) in calculating shortest paths.

**Table 1.39** Values for the WEIGHT= Option

Option Value	Description
YES	Uses weights (if they exist) in shortest path calculations. This is the default.
NO	Does not use weights in shortest path calculations.

**WEIGHT2=column**

specifies the data set variable name for the auxiliary link weights. The value of the *column* variable must be numeric.

---

## SUMMARY Statement

**SUMMARY** < options > ;

The SUMMARY statement invokes an algorithm that calculates various summary metrics on an input graph.

The summary metrics are described in the section “[Summary](#)” on page 153.

You can specify the following *options* in the SUMMARY statement:

**BICONCOMP**

specifies whether to calculate information about biconnected components. The graph must be undirected.

**BY\_CLUSTER**

specifies whether to decompose the calculations by cluster (or subgraph). If this option is specified, PROC OPTGRAPH looks for a definition of the clusters in the input data set specified in the DATA\_NODES= option.

**CONCOMP**

specifies whether to calculate information about connected components.

**DIAMETER\_APPROX=WEIGHT | UNWEIGHT | BOTH**

specifies whether to calculate information about the approximate diameter and what type of calculations to perform. Use this option when calculating the exact diameter (by calculating all shortest paths) is too expensive.

**Table 1.40** Values for the DIAMETER\_APPROX= Option

Option Value	Description
WEIGHT	Calculates approximate diameter based on the weighted graph.
UNWEIGHT	Calculates approximate diameter based on the unweighted graph.
BOTH	Calculates approximate diameter based on both weighted and unweighted graphs.

If the input graph does not contain weights, then WEIGHT and UNWEIGHT both give the same results (using 1.0 for each link weight). This option works only for undirected graphs.

**LOGFREQNODE=number**

controls the frequency for displaying iteration logs for some of the summary metrics. For computationally intensive summary metrics such as shortest path, this option displays progress every *number* nodes. If you also specify the BY\_CLUSTER option in this statement or a value greater than 1 for the NTHREADS= option in the [PERFORMANCE](#) statement, this option is ignored and the display frequency is determined by using the [LOGFREQTIME=](#) option instead. The value of *number* can be

any integer greater than or equal to 1. The default is determined automatically based on the size of the graph. Setting this value too low can hurt performance on large-scale graphs.

**LOGFREQTIME=number**

controls the frequency for displaying iteration logs for some of the summary metrics. For computationally intensive summary metrics such as shortest path, this option displays progress every *number* seconds. When you specify a value greater than 1 for the NTHREADS= option in the **PERFORMANCE** statement, PROC OPTGRAPH displays the number of nodes that have completed. When you specify the BY\_CLUSTER option, PROC OPTGRAPH displays the number of subgraphs that have completed. The value of *number* can be any integer greater than or equal to 1; the default is 5. Setting this value too low can hurt performance on large-scale graphs.

**LOGLEVEL=number**

controls the amount of information that is displayed in the SAS log. Table 1.41 describes the valid values for this option.

**Table 1.41** Values for LOGLEVEL= Option

<i>number</i>	<i>string</i>	Description
0	NONE	Turns off all algorithm-related messages in the SAS log
1	BASIC	Displays a basic summary of the algorithmic processing
2	MODERATE	Displays a summary of the algorithmic processing
3	AGGRESSIVE	Displays a detailed summary of the algorithmic processing

The default is the value that is specified in the **LOGLEVEL=** option in the PROC OPTGRAPH statement (or BASIC if that option is not specified).

**OUT=SAS-data-set**

specifies the output data set to contain the summary results.

**SHORTPATH=WEIGHT | UNWEIGHT | BOTH**

specifies whether to calculate information about shortest paths and what type of calculations to perform.

**Table 1.42** Values for the SHORTPATH= Option

Option Value	Description
WEIGHT	Calculates shortest paths based on the weighted graph.
UNWEIGHT	Calculates shortest paths based on the unweighted graph.
BOTH	Calculates shortest paths based on both weighted and unweighted graphs.

If the input graph does not contain weights, then WEIGHT and UNWEIGHT both give the same results (using 1.0 for each link weight).

**SUBSIZESWITCH=number**

specifies the size of the subgraphs (number of nodes) to run separately when you also specify the BY\_CLUSTER option in this statement and a value greater than 1 for the NTHREADS= option in the **PERFORMANCE** statement. When PROC OPTGRAPH processes summary by subgraphs, it uses thread logic to simultaneously process *n* subgraphs, where *n* is the number of threads specified in the NTHREADS= option in the PERFORMANCE statement. Subgraphs that have more nodes than

*number* are processed sequentially, enabling the threading to be done at the summary metric level. The default is 10,000.

---

## TRANSITIVE\_CLOSURE Statement

**TRANSITIVE\_CLOSURE** < *options* > ;

**TRANSCL** < *options* > ;

The TRANSITIVE\_CLOSURE statement invokes an algorithm that calculates the transitive closure of an input graph.

Transitive closure is described in the section “[Transitive Closure](#)” on page 162.

You can specify the following *options* in the TRANSITIVE\_CLOSURE statement:

**LOGLEVEL=***number*

controls the amount of information that is displayed in the SAS log. [Table 1.43](#) describes the valid values for this option.

**Table 1.43** Values for LOGLEVEL= Option

<i>number</i>	<i>string</i>	Description
0	NONE	Turns off all algorithm-related messages in the SAS log
1	BASIC	Displays a basic summary of the algorithmic processing
2	MODERATE	Displays a summary of the algorithmic processing
3	AGGRESSIVE	Displays a detailed summary of the algorithmic processing

The default is the value that is specified in the **LOGLEVEL=** option in the PROC OPTGRAPH statement (or BASIC if that option is not specified).

**OUT=***SAS-data-set*

specifies the output data set to contain the transitive closure results.

---

## TSP Statement

**TSP** < *options* > ;

The TSP statement invokes an algorithm that solves the traveling salesman problem.

The traveling salesman problem is described in the section “[Traveling Salesman Problem](#)” on page 164. The algorithm that is used to solve this problem is built around the same method as is used in PROC OPTMILP: a branch-and-cut algorithm. Many of the following options are the same as those described for the OPTMILP procedure in the *SAS/OR User’s Guide: Mathematical Programming*.

You can specify the following *options*:

**ABSOBJGAP=number**

specifies a stopping criterion. When the absolute difference between the best integer objective and the objective of the best remaining branch-and-bound node becomes less than the value of *number*, the solver stops. The value of *number* can be any nonnegative number; the default value is 1E-6.

**CONFLICTSEARCH=number | string**

specifies the level of conflict search that PROC OPTGRAPH performs. The solver performs a conflict search to find clauses that result from infeasible subproblems that arise in the search tree. [Table 1.44](#) describes the valid values for this option.

**Table 1.44** Values for CONFLICTSEARCH= Option

<i>number</i>	<i>string</i>	Description
-1	AUTOMATIC	Performs a conflict search based on a strategy that is determined by PROC OPTGRAPH
0	NONE	Disables conflict search
1	MODERATE	Performs a moderate conflict search
2	AGGRESSIVE	Performs an aggressive conflict search

By default, CONFLICTSEARCH=AUTOMATIC.

**CUTOFF=number**

cuts off any branch-and-bound nodes in a minimization problem that has an objective value that is greater than *number*. The value of *number* can be any number; the default value is the positive number that has the largest absolute value that can be represented in your operating environment.

**CUTSTRATEGY=number | string**

specifies the level of mixed integer linear programming cutting planes to be generated by PROC OPTGRAPH. TSP-specific cutting planes are always generated. [Table 1.45](#) describes the valid values for this option.

**Table 1.45** Values for CUTSTRATEGY= Option

<i>number</i>	<i>string</i>	Description
-1	AUTOMATIC	Generates cutting planes based on a strategy determined by the mixed integer linear programming solver
0	NONE	Disables generation of mixed integer linear programming cutting planes (some TSP-specific cutting planes are still active for validity)
1	MODERATE	Uses a moderate cut strategy
2	AGGRESSIVE	Uses an aggressive cut strategy

By default, CUTSTRATEGY=NONE.

**EMPHASIS=number | string**

specifies a search emphasis *option*. [Table 1.46](#) describes the valid values for this option.

**Table 1.46** Values for EMPHASIS= Option

<i>number</i>	<i>string</i>	Description
0	BALANCE	Performs a balanced search
1	OPTIMAL	Emphasizes optimality over feasibility
2	FEASIBLE	Emphasizes feasibility over optimality

By default, EMPHASIS=BALANCE.

**HEURISTICS=***number* | *string*

controls the level of initial and primal heuristics that PROC OPTGRAPH applies. This level determines how frequently PROC OPTGRAPH applies primal heuristics during the branch-and-bound tree search. It also affects the maximum number of iterations that are allowed in iterative heuristics. Some computationally expensive heuristics might be disabled by the solver at less aggressive levels. Table 1.47 lists the valid values for this option.

**Table 1.47** Values for HEURISTICS= Option

<i>number</i>	<i>string</i>	Description
-1	AUTOMATIC	Applies the default level of heuristics
0	NONE	Disables all initial and primal heuristics
1	BASIC	Applies basic initial and primal heuristics at low frequency
2	MODERATE	Applies most initial and primal heuristics at moderate frequency
3	AGGRESSIVE	Applies all initial primal heuristics at high frequency

By default, HEURISTICS=AUTOMATIC.

**LOGFREQ=***number*

specifies how often to print information in the branch-and-bound node log. The value of *number* can be any nonnegative integer up to the largest four-byte signed integer, which is  $2^{31} - 1$ . The default value is 100. If *number* is set to 0, then the node log is disabled. If *number* is positive, then an entry is made in the node log at the first node, at the last node, and at intervals that are controlled by the value of *number*. An entry is also made each time a better integer solution is found.

**LOGLEVEL=***number* | *string*

controls the amount of information displayed in the SAS log by the solver, from a short description of presolve information and summary to details at each branch-and-bound node. Table 1.48 describes the valid values for this option.

**Table 1.48** Values for LOGLEVEL= Option

<i>number</i>	<i>string</i>	Description
0	NONE	Turns off all solver-related messages in the SAS log
1	BASIC	Displays a solver summary after stopping
2	MODERATE	Prints a solver summary and a node log by using the interval that is specified in the LOGFREQ= option
3	AGGRESSIVE	Prints a detailed solver summary and a node log by using the interval that is specified in the LOGFREQ= option



The default value is MODERATE.

**MAXNODES=number**

specifies the maximum number of branch-and-bound nodes to be processed. The value of *number* can be any nonnegative integer up to the largest four-byte signed integer, which is  $2^{31} - 1$ . The default value is  $2^{31} - 1$ .

**MAXSOLS=number**

specifies a stopping criterion. If *number* solutions have been found, then the procedure stops. The value of *number* can be any positive integer up to the largest four-byte signed integer, which is  $2^{31} - 1$ . The default value is  $2^{31} - 1$ .

**MAXTIME=number**

specifies the maximum amount of time to spend solving the traveling salesman problem. The type of time (either CPU time or real time) is determined by the value of the **TIMETYPE=** option. The value of *number* can be any positive number; the default value is the positive number that has the largest absolute value that can be represented in your operating environment.

**MILP=number | string**

specifies whether to use a mixed integer linear programming (MILP) solver for solving the traveling salesman problem. The MILP solver attempts to find the overall best TSP tour by using a branch-and-bound based algorithm. This algorithm can be expensive for large-scale problems. If MILP=OFF, then PROC OPTGRAPH uses its initial heuristics to find a feasible, but not necessarily optimal, tour as quickly as possible. [Table 1.49](#) describes the valid values for this option.

**Table 1.49** Values for MILP= Option

<i>number</i>	<i>string</i>	Description
1	ON	Uses a mixed integer linear programming solver
0	OFF	Does not use a mixed integer linear programming solver

By default, MILP=ON.

**NODESEL=number | string**

specifies the branch-and-bound node selection strategy option. [Table 1.50](#) describes the valid values for this option.

**Table 1.50** Values for NODESEL= Option

<i>number</i>	<i>string</i>	Description
-1	AUTOMATIC	Uses automatic node selection
0	BESTBOUND	Chooses the node that has the best relaxed objective (best-bound-first strategy)
1	BESTESTIMATE	Chooses the node that has the best estimate of the integer objective value (best-estimate-first strategy)
2	DEPTH	Chooses the most recently created node (depth-first strategy)

By default, NODESEL=AUTOMATIC. For more information about node selection, see Chapter 13, “The OPTMILP Procedure” (*SAS/OR User’s Guide: Mathematical Programming*).

**OUT=SAS-data-set**

specifies the output data set to contain the solution to the traveling salesman problem.

**PROBE=number | string**

specifies a probing *option*. [Table 1.51](#) describes the valid values for this option.

**Table 1.51** Values for PROBE= Option

<i>number</i>	<i>string</i>	Description
–1	AUTOMATIC	Uses an automatic probing strategy
0	NONE	Disables probing
1	MODERATE	Uses the probing moderately
2	AGGRESSIVE	Uses the probing aggressively

By default, PROBE=NONE.

**RELOBJGAP=number**

specifies a stopping criterion that is based on the best integer objective (BestInteger) and the objective of the best remaining node (BestBound). The relative objective gap is equal to

$$| \text{BestInteger} - \text{BestBound} | / (1\text{E} - 10 + | \text{BestBound} |)$$

When this value becomes less than the specified gap size *number*, the solver stops. The value of *number* can be any nonnegative number. By default, RELOBJGAP=1E–4.

**STRONGITER=number | AUTOMATIC**

specifies the number of simplex iterations that PROC OPTGRAPH performs for each variable in the candidate list when it uses the strong branching variable selection strategy. The value of *number* can be any positive integer up to the largest four-byte signed integer, which is  $2^{31} - 1$ . If you specify the keyword AUTOMATIC or the value –1, PROC OPTGRAPH uses the default value, which it calculates automatically.

**STRONGLEN=number | AUTOMATIC**

specifies the number of candidates that PROC OPTGRAPH considers when it uses the strong branching variable selection strategy. The value of *number* can be any positive integer up to the largest four-byte signed integer, which is  $2^{31} - 1$ . If you specify the keyword AUTOMATIC or the value –1, PROC OPTGRAPH uses the default value, which it calculates automatically.

**TARGET=number**

specifies a stopping criterion for minimization problems. If the best integer objective is better than or equal to *number*, the solver stops. The value of *number* can be any number; the default is the negative number that has the largest absolute value that can be represented in your operating environment.

**VARSEL=number | string**

specifies the rule for selecting the branching variable. [Table 1.52](#) describes the valid values for this option.

**Table 1.52** Values for VARSEL= Option

number	string	Description
-1	AUTOMATIC	Uses automatic branching variable selection
0	MAXINFEAS	Chooses the variable that has maximum infeasibility
1	MININFEAS	Chooses the variable that has minimum infeasibility
2	PSEUDO	Chooses a branching variable based on pseudocost
3	STRONG	Uses the strong branching variable selection strategy

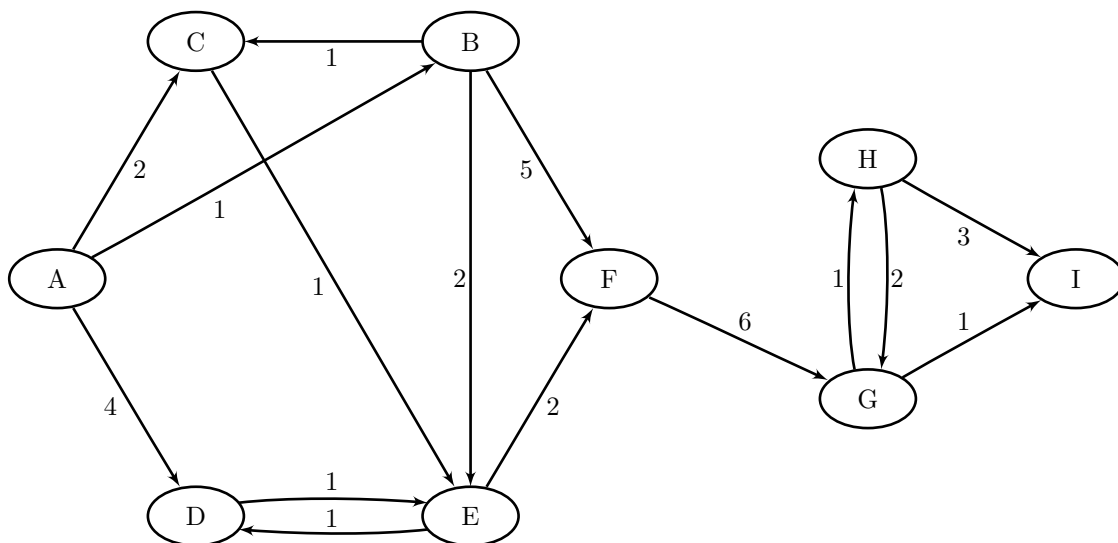
By default, VARSEL=AUTOMATIC. For more information about variable selection, see Chapter 13, “The OPTMILP Procedure” (*SAS/OR User’s Guide: Mathematical Programming*).

## Details: OPTGRAPH Procedure

### Graph Input Data

This section describes how to input a graph for analysis by PROC OPTGRAPH. Let  $G = (N, A)$  define a graph with a set  $N$  of nodes and a set  $A$  of links.

Consider the directed graph shown in Figure 1.8.

**Figure 1.8** A Simple Directed Graph

Notice that each node and link has associated attributes: a node label and a link weight.

## Link Input Data

The `DATA_LINKS=` option in the `PROC OPTGRAPH` statement defines the data set that contains the list of links in the graph. A link is represented as a pair of nodes, which are defined by using either numeric or character labels. The links data set is expected to contain some combination of the following possible variables:

- `from`: the *from* node (this variable can be numeric or character)
- `to`: the *to* node (this variable can be numeric or character)
- `weight`: the link weight (this variable must be numeric)
- `weight2`: the auxiliary link weight (this variable must be numeric)
- `lower`: the link flow lower bound (this variable must be numeric)
- `upper`: the link flow upper bound (this variable must be numeric)

As described in the `GRAPH_DIRECTION=` option, if the graph is undirected, the *from* and *to* labels are interchangeable. If the weights are not given for algorithms that call for link weights, they are all assumed to be 1.

The data set variable names can have any values that you want. If you use nonstandard names, you must identify the variables by using the `DATA_LINKS_VAR` statement, as described in the section “`DATA_LINKS_VAR` Statement” on page 32.

For example, the following two data sets identify the same graph:

```
data LinkSetInA;
  input from $ to $ weight;
  datalines;
A B 1
A C 2
A D 4
;

data LinkSetInB;
  input source_node $ sink_node $ value;
  datalines;
A B 1
A C 2
A D 4
;
```

These data sets can be presented to `PROC OPTGRAPH` by using the following equivalent statements:

```
proc optgraph
  data_links = LinkSetInA;
run;

proc optgraph
  data_links = LinkSetInB;
  data_links_var
    from = source_node
    to = sink_node
    weight = value;
run;
```

The directed graph *G* shown in Figure 1.8 can be represented by the following links data set LinkSetIn:

```
data LinkSetIn;
  input from $ to $ weight @@;
  datalines;
A B 1 A C 2 A D 4 B C 1 B E 2
B F 5 C E 1 D E 1 E D 1 E F 2
F G 6 G H 1 G I 1 H G 2 H I 3
;
```

The following statements read in this graph, declare it as a directed graph, and output the resulting links and nodes data sets. These statements do not run any algorithms, so the resulting output contains only the input graph.

```
proc optgraph
  graph_direction = directed
  data_links = LinkSetIn
  out_nodes = NodeSetOut
  out_links = LinkSetOut;
run;
```

The data set NodeSetOut, shown in Figure 1.9, now contains the nodes that are read from the input link data set. The variable node shows the label associated with each node.

Figure 1.9 Node Data Set of a Simple Directed Graph

node
A
B
C
D
E
F
G
H
I

The data set LinkSetOut, shown in Figure 1.10, contains the links that were read from the input link data set. The variables from and to show the associated node labels.

**Figure 1.10** Link Data Set of a Simple Directed Graph

Obs	from	to	weight
1	A	B	1
2	A	C	2
3	A	D	4
4	B	C	1
5	B	E	2
6	B	F	5
7	C	E	1
8	D	E	1
9	E	D	1
10	E	F	2
11	F	G	6
12	G	H	1
13	G	I	1
14	H	G	2
15	H	I	3

If you define this graph as undirected, then reciprocal links (for example,  $D \rightarrow E$  and  $D \leftarrow E$ ) are treated as the same link, and duplicates are removed. PROC OPTGRAPH takes the first occurrence of the link and ignores the others. By default, GRAPH\_DIRECTION=UNDIRECTED, so you can just remove this option to declare the graph as undirected.

```
proc optgraph
  data_links = LinkSetIn
  out_nodes  = NodeSetOut
  out_links  = LinkSetOut;
run;
```

The progress of the procedure is shown in Figure 1.11. The log now shows the links (and their observation identifiers) that were declared as duplicates and removed.

**Figure 1.11** PROC OPTGRAPH Log: Link Data Set of a Simple Undirected Graph

```
NOTE: -----
NOTE: Running OPTGRAPH version 14.1.
NOTE: -----
NOTE: The OPTGRAPH procedure is executing in single-machine mode.
NOTE: -----
NOTE: Data input used 0.00 (cpu: 0.00) seconds.
WARNING: Link (E,D) in observation 9 of the DATA_LINKS= data set is a duplicate and is ignored.
WARNING: Link (H,G) in observation 14 of the DATA_LINKS= data set is a duplicate and is ignored.
NOTE: The number of nodes in the input graph is 9.
NOTE: The number of links in the input graph is 13.
NOTE: -----
NOTE: Data output used 0.00 (cpu: 0.00) seconds.
NOTE: -----
NOTE: The data set WORK.NODESETOUT has 9 observations and 1 variables.
NOTE: The data set WORK.LINKSETOUT has 13 observations and 3 variables.
NOTE: -----
```

The data set NodeSetOut is equivalent to the one shown in [Figure 1.9](#). However, the new links data set LinkSetOut shown in [Figure 1.12](#) contains two fewer links than before, because duplicates are removed.

**Figure 1.12** Link Data Set of a Simple Undirected Graph

Obs	from	to	weight
1	A	B	1
2	A	C	2
3	A	D	4
4	B	C	1
5	B	E	2
6	B	F	5
7	C	E	1
8	D	E	1
9	E	F	2
10	F	G	6
11	G	H	1
12	G	I	1
13	H	I	3

Certain algorithms can perform more efficiently when you specify `GRAPH_INTERNAL_FORMAT=THIN` in the `PROC OPTGRAPH` statement. However, when you specify this option, `PROC OPTGRAPH` does not remove duplicate links. Instead, you should use appropriate `DATA` steps to clean your data before calling `PROC OPTGRAPH`.

## Node Input Data

The `DATA_NODES=` option in the `PROC OPTGRAPH` statement defines the data set that contains the list of nodes in the graph. This data set is used to define clusters (subgraphs) or to assign node weights.

The nodes data set is expected to contain some combination of the following possible variables:

- `node`: the node label (this variable can be numeric or character)
- `cluster`: the node cluster identifier (this variable must be numeric)
- `weight`: the node weight (this variable must be numeric)
- `weight2`: the auxiliary node weight (this variable must be numeric)

The variable `cluster` is used to define clusters (subgraphs) for decomposing the input graph into subgraphs for processing. This is useful for the algorithms specified in the `CENTRALITY`, `REACH`, and `SUMMARY` statements. The use of the variable `cluster` is explained in more detail in the section “[Processing by Cluster](#)” on page 78.

You can specify any values that you want for the data set variable names. If you use nonstandard names, you must identify the variables by using the `DATA_NODES_VAR` statement, as described in the section “[DATA\\_NODES\\_VAR Statement](#)” on page 33.

The data set that is specified in the `DATA_LINKS=` option defines the set of nodes that are incident to some link. If the graph contains a node that has no links (called a *singleton node*), then this node must be defined

in the DATA\_NODES data set. The following is an example of a graph with three links but four nodes, including a singleton node D:

```
data NodeSetIn;
    input label $ @@;
    datalines;
A B C D
;

data LinkSetInS;
    input from $ to $ weight;
    datalines;
A B 1
A C 2
B C 1
;
```

If you specify duplicate entries in the node data set, PROC OPTGRAPH takes the first occurrence of the node and ignores the others. A warning is printed to the log.

### Node Subset Input Data

For some algorithms, you might want to process only a subset of the nodes that appear in the input graph. You can accomplish this by using the DATA\_NODES\_SUB= option in the PROC OPTGRAPH statement. You can use the node subset data set in conjunction with the SHORTPATH or REACH statement. (See the sections “[Shortest Path](#)” on page 142 and “[Reach \(Ego\) Network](#)” on page 130, respectively.) The node subset data set is expected to contain some combination of the following variables:

- node: the node label (this variable can be numeric or character)
- source: whether to process this node as a source node in shortest path algorithms (this variable must be numeric)
- sink: whether to process this node as a sink node in shortest path algorithms (this variable must be numeric)
- reach: for the reach algorithm, the index of the source subgraph for processing (this variable must be numeric)

[Table 1.53](#) shows how PROC OPTGRAPH processes nodes for each algorithm type. The missing indicator (.) can also be used in place of 0 to designate that a node is not to be processed.

**Table 1.53** Determining How to Process a Node

Algorithm Type	Variable Designations	Example Shown In:
Shortest path	A value of 0 for the source variable designates that the node is not to be processed as a source; a value of 1 designates that the node is to be processed as a source. The same values can be used for the sink variable to designate whether the node is to be processed as a sink.	The section “ <a href="#">Shortest Path</a> ” on page 142



**Table 1.53** (continued)

Algorithm Type	Variable Designations	Example Shown In:
Reach	A value of 0 for the reach variable designates that the node is not to be processed. A value greater than 0 defines a marker for the source subgraph to which this node belongs. All nodes with the same marker are processed together as source nodes.	The section “ <a href="#">Reach (Ego) Network</a> ” on page 130

A representative example of a node subset data set that might be used with the graph in [Figure 1.8](#) is as follows:

```
data NodeSubSetIn;
    input node $ reach source sink;
    datalines;
A 1 1 .
F 2 . 1
E 2 1 .
;
```

The data set NodeSubSetIn indicates that you want to process the following:

- the reach network from the subgraph defined by node A
- the reach network from the subgraph defined by nodes F and E
- the shortest paths for the source-sink pairs in  $\{A, E\} \times \{F\}$

## Standardized Labels

For large-scale graphs, the processing stage that reads the nodes and links into memory can be time-consuming. Under the following assumptions, you can use the STANDARDIZED\_LABELS option in the PROC OPTGRAPH statement to greatly speed up this stage:

1. The link data set variables from and to are numeric type.
2. The node and node subset data set variable node is numeric type.
3. The node labels start from 0 and are consecutive nonnegative integers.

Consider the following links data set that uses numeric labels:

```
data LinkSetIn;
    input from to weight;
    datalines;
0 1 1
3 0 2
1 5 1
;
```

Using default settings, the following statements echo back link and node data sets that contain three links and four nodes, respectively:

```
proc optgraph
  data_links = LinkSetIn
  out_nodes  = NodeSetOut
  out_links  = LinkSetOut;
run;
```

The log is shown in [Figure 1.13](#).

**Figure 1.13** PROC OPTGRAPH Log: A Simple Undirected Graph

---

```
NOTE: -----
NOTE: Running OPTGRAPH version 14.1.
NOTE: -----
NOTE: The OPTGRAPH procedure is executing in single-machine mode.
NOTE: -----
NOTE: Data input used 0.00 (cpu: 0.00) seconds.
NOTE: The number of nodes in the input graph is 4.
NOTE: The number of links in the input graph is 3.
NOTE: -----
NOTE: Data output used 0.00 (cpu: 0.00) seconds.
NOTE: -----
NOTE: The data set WORK.NODESETOUT has 4 observations and 1 variables.
NOTE: The data set WORK.LINKSETOUT has 3 observations and 3 variables.
```

---

The data set NodeSetOut, shown in [Figure 1.14](#), contains the unique numeric node labels, {0, 1, 3, 5}.

**Figure 1.14** Node Data Set of a Simple Directed Graph

Obs	node
1	0
2	1
3	3
4	5

Using standardized labels, the same input data set defines a graph that has six (not four) nodes:

```
proc optgraph
  standardized_labels
  data_links = LinkSetIn
  out_nodes  = NodeSetOut
  out_links  = LinkSetOut;
run;
```

The log that results from using standardized labels is shown in [Figure 1.15](#).

**Figure 1.15** PROC OPTGRAPH Log: A Simple Undirected Graph Using Standardized Labels

---

```

NOTE: -----
NOTE: Running OPTGRAPH version 14.1.
NOTE: -----
NOTE: The OPTGRAPH procedure is executing in single-machine mode.
NOTE: -----
NOTE: Data input used 0.00 (cpu: 0.00) seconds.
NOTE: The number of nodes in the input graph is 6.
NOTE: The number of links in the input graph is 3.
NOTE: The number of singleton nodes in the input graph is 2.
NOTE: -----
NOTE: Data output used 0.00 (cpu: 0.00) seconds.
NOTE: -----
NOTE: The data set WORK.NODESETOUT has 6 observations and 1 variables.
NOTE: The data set WORK.LINKSETOUT has 3 observations and 3 variables.

```

---

The data set NodeSetOut, shown in [Figure 1.16](#), now contains all node labels from 0 to 5, based on the assumptions when you use the STANDARDIZED\_LABELS option.

**Figure 1.16** Node Data Set of a Simple Directed Graph

Obs	node
1	0
2	1
3	2
4	3
5	4
6	5

When you use standardized labels, the DATA\_NODES= input order (which can be arbitrary) is not preserved in the OUT\_NODES= output data set. Instead, the order is ascending, starting from zero.

---

## Matrix Input Data

This section describes the matrix input format that you can use with some of the algorithms in PROC OPTGRAPH. The DATA\_MATRIX= option in the PROC OPTGRAPH statement defines the data set that contains the matrix values. You can specify any values that you want for the data set variable names (the columns) by using the DATA\_MATRIX\_VAR statement, as described in the section “[DATA\\_MATRIX\\_VAR Statement](#)” on page 33. If you do not specify a DATA\_MATRIX\_VAR statement, then PROC OPTGRAPH assumes that all numeric variables in the data set are to be used in defining the matrix.

The following statements find the principal eigenvector of the square symmetric matrix that is defined in the data set Matrix:

```

data Matrix;
    input coll-col5;
    datalines;
1 0 2 6 1
0 2 3 0 1
2 3 1 0 2
6 0 0 0 0
1 1 2 0 0
;

proc optgraph
    data_matrix    = Matrix;
    eigenvector
        eigenvalues = LA
        nEigen      = 1
        out          = EigenVector;
run;

```

The following statements solve the linear assignment problem for the cost matrix that is defined in the data set CostMatrix:

```

data CostMatrix;
    input back breast fly free;
    datalines;
35.1 36.7 28.3 36.1
34.6 32.6 26.9 26.2
31.3 33.9 27.1 31.2
28.6 34.1 29.1 30.3
32.9 32.2 26.6 24.0
27.8 32.5 27.8 27.0
26.3 27.6 23.5 22.4
29.0 24.0 27.9 25.4
27.2 33.8 25.2 24.1
27.0 29.2 23.0 21.9
;

proc optgraph
    data_matrix = CostMatrix;
    data_matrix_var
        back--free;
    linear_assignment
        out      = LinearAssign;
run;

```

---

## Data Input Order

Many algorithms are sensitive to the order in which PROC OPTGRAPH reads the data. If the order of the nodes or links is changed, either by you or by some parameter setting, the final result might change. In some cases, this difference is simply a permutation of identifiers (for example, connected components). In other cases, when you use local optimization (for example, community detection) or discrete branching decisions (for example, the traveling salesman problem), the final result might be a local (or alternative) solution. Two parameters that could have such an effect are the STANDARDIZED\_LABELS and NTHREADS= options. Both of these options can change the internal order of the nodes and links.

---

## Parallel Processing

PROC OPTGRAPH can take advantage of multicore chip technology by processing the graph input data in parallel. To enable PROC OPTGRAPH to process in parallel, you can use the NTHREADS= option in the PERFORMANCE statement to specify the number of threads to use.

In addition, a number of the algorithms in PROC OPTGRAPH can also take advantage of multiple cores. There are two ways in which PROC OPTGRAPH can decompose the computational work in order to take advantage of parallel processing: by node and by subgraph.

To process the nodes of the graph individually, set the NTHREADS= option to some value greater than 1. You can do this for the centrality metrics closeness (see the section “[Closeness Centrality](#)” on page 69) and betweenness (see the section “[Betweenness Centrality](#)” on page 72). An example of this is shown in “[Example 1.4: Betweenness and Closeness Centrality for Project Groups in a Research Department](#)” on page 196.

To process the subgraphs of the original graph individually, set the NTHREADS= option to some value greater than 1, and designate the clusters in each node by using the cluster variable in the nodes data set, as described in the section “[Node Input Data](#)” on page 53. You can do this for centrality metrics, reach networks, and summary statistics. (See the sections “[Centrality](#)” on page 65, “[Reach \(Ego\) Network](#)” on page 130, and “[Summary](#)” on page 153, respectively.) A common use for this feature is to first decompose the original graph into communities or components. (See the sections “[Community Detection](#)” on page 89 and “[Connected Components](#)” on page 97, respectively.) Then, from these results, define the clusters in the node data set and run the analysis of each subgraph individually and in parallel. PROC OPTGRAPH takes care of all the accounting with the associated decomposition and returns results in terms of the original graph. An example of this process for centrality is shown in the section “[Processing by Cluster](#)” on page 78.

You can improve the performance of the OPTGRAPH procedure by running it in distributed computing mode. For more information about the high-performance features of the OPTGRAPH procedure, see *SAS OPTGRAPH Procedure: High-Performance Features*.

**NOTE:** Distributed computing mode requires SAS High-Performance Analytics software.

---

## Numeric Limitations

Extremely large or extremely small numerical values might cause computational difficulties for some of the algorithms in PROC OPTGRAPH. For this reason, each algorithm restricts the magnitude of the data values to a particular threshold number. If the user data values exceed this threshold, PROC OPTGRAPH issues an error message. The value of the threshold limit is different for each algorithm and depends on the operating environment. The threshold limits are listed in [Table 1.54](#), where  $M$  is defined as the largest absolute value representable in your operating environment.

**Table 1.54** Threshold Limits by Statement

Statement	Matrix	Graph Links				Graph Nodes	
		weight	weight2	lower	upper	weight	weight2
CENTRALITY AUTH=, EIGEN=, HUB= BETWEEN=, CLOSE= INFLUENCE=		1e20 $\sqrt{M}$ $\sqrt{M}$	$\sqrt{M}$			$\sqrt{M}$	
COMMUNITY		$\sqrt{M}$					
CYCLE		$\sqrt{M}$				$\sqrt{M}$	
EIGENVECTOR	1e20	1e20					
LINEAR_ASSIGNMENT	$\sqrt{M}$	$\sqrt{M}$					
MINCOSTFLOW		1e15		1e15	1e15	1e15	1e15
MINCUT		$\sqrt{M}$					
MINSPANTREE		$\sqrt{M}$					
REACH						$\sqrt{M}$	
SHORTPATH		$\sqrt{M}$	$\sqrt{M}$				
SUMMARY DIAMETER_APPROX=, SHORTPATH=		$\sqrt{M}$					
TSP		1e20					

To obtain these limits, use the SAS function constant. For example, the following DATA step assigns  $\sqrt{M}$  to a variable x and prints that value to the log:

```
data _null_;
  x = constant('SQRTBIG');
  put x=;
run;
```

## Missing Values

For most of the algorithms in PROC OPTGRAPH, there is no valid interpretation for a missing value. If the user data contain a missing value, PROC OPTGRAPH issues an error message. One exception is for the minimum-cost network flow solver when you are setting the link or node bounds. In this case, a missing value is interpreted as the default bound value, as described in the section “[Minimum-Cost Network Flow](#)” on page 116. Another exception is the linear assignment problem when you are using the matrix input format. A missing value in this case defines an invalid assignment between a row and a column of the matrix. An example of this is shown in the section “[Linear Assignment \(Matching\)](#)” on page 115.

## Negative Link Weights

For certain algorithms in PROC OPTGRAPH, a negative link weight is not allowed. The following algorithms issue an error message if a negative link weight is provided:

- CENTRALITY
  - AUTH=, BETWEEN=, CLOSE=, EIGEN=, HUB=

- COMMUNITY
- MINCUT

## Zero Link Weights

For the community detection algorithm, a zero-valued link weight is not allowed. If a zero-valued link weight is provided, the community detection algorithm issues an error message.

## Size Limitations

PROC OPTGRAPH can handle any graph whose number of nodes and links are less than or equal to 2,147,483,647 (the maximum that can be represented by a 32-bit integer). This maximum also applies to 64-bit systems. For graphs of two billion nodes (or links), memory limitations also become a limiting factor. For example, see the discussion of memory requirements for the community detection algorithm in the section “[Memory Requirement](#)” on page 91

If the data from your problem require a graph with more than two billion nodes (or links), there is typically a heuristic way to break the network into smaller networks based on problem-specific attributes. Then, using DATA steps, you can process each of the smaller networks iteratively through repeated calls to PROC OPTGRAPH. By using DATA steps, you can also often work around memory limitations, because the full graph exists only on the disk and never resides in memory.

## Common Notation and Assumptions

This section introduces some common notation and assumptions used throughout the chapter.

A *complete graph*, denoted  $K(N)$ , is a graph in which every pair of nodes in  $N$  is connected by a link. The number of links in  $K(N)$  is described in [Table 1.55](#).

**Table 1.55** Formulas for Number of Links in  $K(N)$

Graph Direction	Default	INCLUDE_SELFLINK
Directed	$ N ^2 -  N $	$ N ^2$
Undirected	$\frac{ N ^2 -  N }{2}$	$\frac{ N ^2 +  N }{2}$

## Biconnected Components and Articulation Points

A *biconnected component* of a graph  $G = (N, A)$  is a connected subgraph that cannot be broken into disconnected pieces by deleting any single node (and its incident links). An *articulation point* is a node of a graph whose removal would cause an increase in the number of connected components. Articulation points can be important when you analyze any graph that represents a communications network. Consider an articulation point  $i \in N$  which, if removed, disconnects the graph into two components  $C^1$  and  $C^2$ . All paths in  $G$  between some nodes in  $C^1$  and some nodes in  $C^2$  must pass through node  $i$ . In this sense, articulation

points are critical to communication. Examples of where articulation points are important are airline hubs, electric circuits, network wires, protein bonds, traffic routers, and numerous other industrial applications.

In PROC OPTGRAPH, you can find biconnected components and articulation points of an input graph by invoking the BICONCOMP statement. This algorithm works only with undirected graphs.

The results of the biconnected components algorithm are written to the output links data set that is specified in the OUT\_LINKS= option in the PROC OPTGRAPH statement. For each link in the links data set, the variable biconcomp identifies its component. The component identifiers are numbered sequentially starting from 1. The results of the articulation points are written to the output nodes data set that is specified in the OUT\_NODES= option in the PROC OPTGRAPH statement. For each node in the nodes data set, the variable artpoint is either 1 (if the node is an articulation point) or 0 (otherwise).

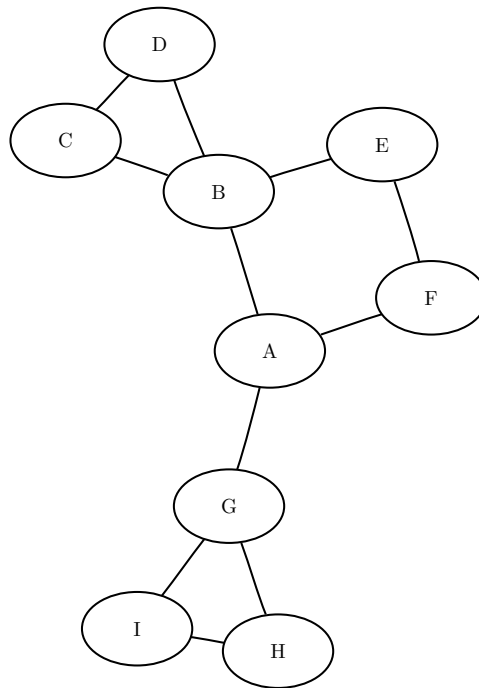
The biconnected components algorithm reports status information in a macro variable called \_OPTGRAPH\_BICONCOMP\_. For more information about this macro variable see the section “[Macro Variable \\_OPTGRAPH\\_BICONCOMP\\_](#)” on page 176.

The algorithm that PROC OPTGRAPH uses to compute biconnected components is a variant of depth-first search (Tarjan 1972). This algorithm runs in time  $O(|N| + |A|)$  and therefore should scale to very large graphs.

### Biconnected Components of a Simple Undirected Graph

This section illustrates the use of the biconnected components algorithm on the simple undirected graph  $G$  that is shown in [Figure 1.17](#).

**Figure 1.17** A Simple Undirected Graph  $G$



The undirected graph  $G$  can be represented by the following links data set LinkSetInBiCC:



```

data LinkSetInBiCC;
    input from $ to $ @@;
    datalines;
A B  A F  A G  B C  B D
B E  C D  E F  G I  G H
H I
;

```

The following statements calculate the biconnected components and articulation points and output the results in the data sets LinkSetOut and NodeSetOut:

```

proc optgraph
    data_links = LinkSetInBiCC
    out_links   = LinkSetOut
    out_nodes   = NodeSetOut;
    biconcomp;
run;

```

The data set LinkSetOut now contains the biconnected components of the input graph, as shown in [Figure 1.18](#).

**Figure 1.18** Biconnected Components of a Simple Undirected Graph

from	to	biconcomp
A	B	2
A	F	2
A	G	4
B	C	1
B	D	1
B	E	2
C	D	1
E	F	2
G	I	3
G	H	3
H	I	3

In addition, the data set NodeSetOut contains the articulation points of the input graph, as shown in [Figure 1.19](#).

**Figure 1.19** Articulation Points of a Simple Undirected Graph

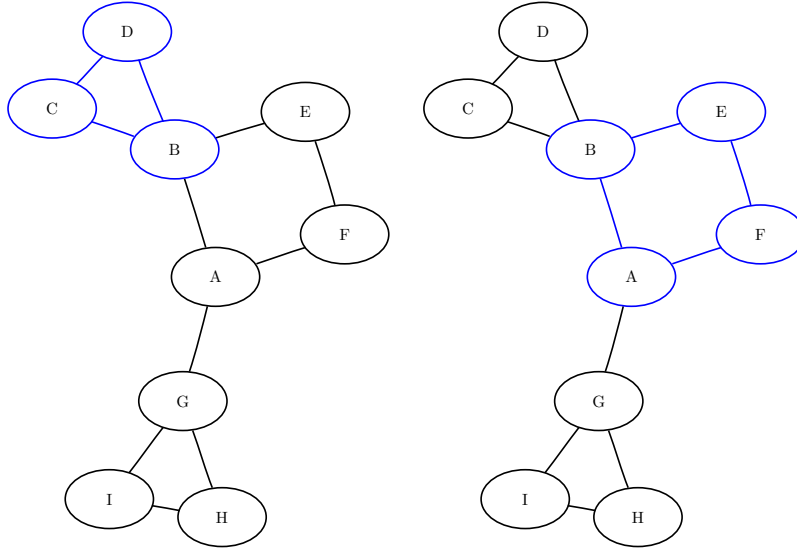
node	artpoint
A	1
B	1
F	0
G	1
C	0
D	0
E	0
I	0
H	0

The biconnected components are shown graphically in Figure 1.20 and Figure 1.21.

**Figure 1.20** Biconnected Components  $C^1$  and  $C^2$

$$C^1 = \{B, C, D\}$$

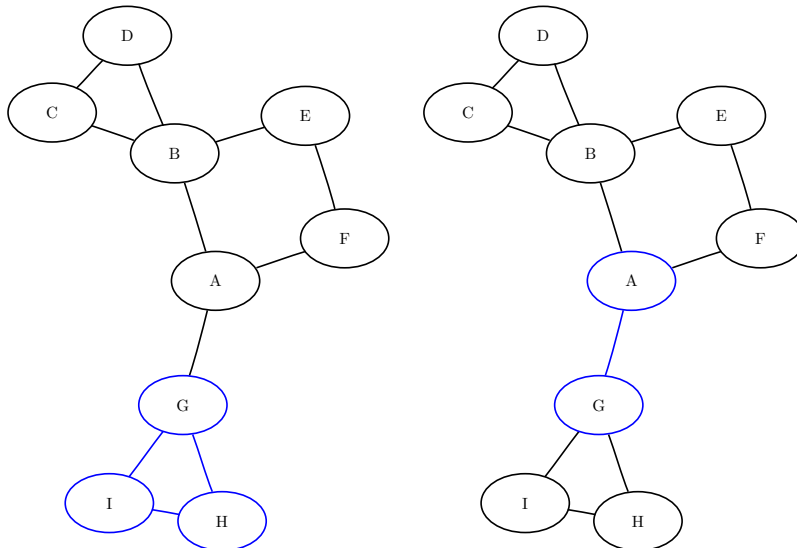
$$C^2 = \{A, B, E, F\}$$



**Figure 1.21** Biconnected Components  $C^3$  and  $C^4$

$$C^3 = \{G, H, I\}$$

$$C^4 = \{A, G\}$$



For a more detailed example, see “Example 1.1: Articulation Points in a Terrorist Network” on page 187.

## Centrality

In general terms, the *centrality* of a node or link in a graph gives some indication of its relative importance within a graph. In the field of network analysis, many different types of centrality metrics are used to better understand levels of prominence. For a good review of centrality metrics, see Newman 2010.

You can use the CENTRALITY statement in PROC OPTGRAPH to calculate several of these metrics. The options for this statement are described in the section “CENTRALITY Statement” on page 20.

The CENTRALITY statement reports status information in a macro variable called \_OPTGRAPH\_CENTRALITY\_. For more information about this macro variable, see the section “Macro Variable \_OPTGRAPH\_CENTRALITY\_” on page 177.

The following sections describe each of the possible centrality metrics that can be calculated in PROC OPTGRAPH.

### Degree Centrality

The *degree* of a node  $v$  in an undirected graph is the number of links that are incident to node  $v$ . The *out-degree* of a node in a directed graph is the number of out-links incident to that node; the *in-degree* is the number of in-links incident. The term *degree* and *out-degree* are interchangeable for an undirected graph. *Degree centrality* is simply the (in- or out-) degree of a node and can be interpreted as some form of relative importance to a network. For example, in a network where nodes are people and you are tracking the flow of a virus, the degree centrality gives some idea of the magnitude of the risk of spreading the virus. People with a higher out-degree can lead to a quicker and more widespread transmission. In a friendship network, in-degree often indicates popularity.

Degree centrality is calculated according to the value specified for the DEGREE= option in the CENTRALITY statement. The results are provided in the node output data set that is specified in the OUT\_NODES= option in the PROC OPTGRAPH statement.

The algorithm used by PROC OPTGRAPH to compute degree centrality is a simple lookup, runs in time  $O(|N|)$ , and therefore should scale to very large graphs.

As a simple example, consider again the directed graph in Figure 1.8 with data set LinkSetIn defined in the section “Link Input Data” on page 50. The following statements calculate the degree centrality for both in- and out-degree:

```
proc optgraph
  graph_direction = directed
  data_links      = LinkSetIn
  out_nodes       = NodeSetOut;
  centrality
    degree        = both;
run;
```

The node data set NodeSetOut now contains the degree centrality of the input graph. For a directed graph, the data set provides the in-degree (variable centr\_degree\_in), the out-degree (variable centr\_degree\_out), and the degree that is the sum of in- and out-degrees (variable centr\_degree). This data set is shown in Figure 1.22.

**Figure 1.22** Degree Centrality of a Simple Directed Graph

node	centr_degree_in	centr_degree_out	centr_degree
A	0	3	3
B	1	3	4
C	2	1	3
D	2	1	3
E	3	2	5
F	2	1	3
G	2	2	4
H	1	2	3
I	2	0	2

## Influence Centrality

*Influence centrality* is a generalization of degree centrality that considers the link and node weights of adjacent nodes ( $C_1$ ) in addition to the link weights of nodes that are adjacent to adjacent nodes ( $C_2$ ). The metric  $C_1$  is referred to as *first-order influence centrality*, and the metric  $C_2$  is referred to as *second-order influence centrality*.

Let  $w_{uv}$  define the link weight for link  $(u, v)$ , and let  $w_u$  define the node weight for node  $u$ . Let  $\delta_u$  represent the list of nodes connected to node  $u$ ; this list is called the adjacency list. For directed graphs, the adjacency list corresponds to the nodes in the out-links. The general formula for influence centrality is

$$C_1(u) = \frac{\sum_{v \in \delta_u} w_{uv}}{\sum_{v \in N} w_v}$$

$$C_2(u) = \sum_{v \in \delta_u} C_1(v)$$

As the name suggests, this metric gives some indication of potential influence, performance, or ability to transfer knowledge.

Influence centrality is calculated according to the value of the **INFLUENCE=** option in the **CENTRALITY** statement. The results are provided in the node output data set that is specified in the **OUT\_NODES=** option in the **PROC OPTGRAPH** statement.

The algorithm used by **PROC OPTGRAPH** to compute influence centrality is a simple traversal, runs in time  $O(|A|)$ , and therefore should scale to very large graphs.

Consider again the directed graph in [Figure 1.8](#). Ignore the weights and just calculate the  $C_1$  and  $C_2$  metrics based on connections (that is, consider all link and node weights as 1). The following statements calculate the unweighted influence centrality:

```
proc optgraph
  graph_direction = directed
  data_links      = LinkSetIn
  out_nodes       = NodeSetOut;
  centrality
    influence     = unweight;
run;
```

The node data set NodeSetOut now contains the unweighted influence centrality of the input graph, including the  $C_1$  variable centr\_influence1\_unwt and the  $C_2$  variable centr\_influence2\_unwt. This data set is shown in Figure 1.23.

**Figure 1.23** Influence Centrality of a Simple Directed Graph

node	centr_influence1_unwt	centr_influence2_unwt
A	0.33333	0.55556
B	0.33333	0.44444
C	0.11111	0.22222
D	0.11111	0.22222
E	0.22222	0.22222
F	0.11111	0.22222
G	0.22222	0.22222
H	0.22222	0.22222
I	0.00000	0.00000

For a more detailed example, see “[Example 1.2: Influence Centrality for Project Groups in a Research Department](#)” on page 189.

## Clustering Coefficient

The *clustering coefficient* for a node is the number of links between the nodes within its neighborhood divided by the number of links that could possibly exist between them (their induced complete graph).

Let  $N_i$  represent the list of nodes that are connected to node  $i$  (excluding itself). The formula for the clustering coefficient is

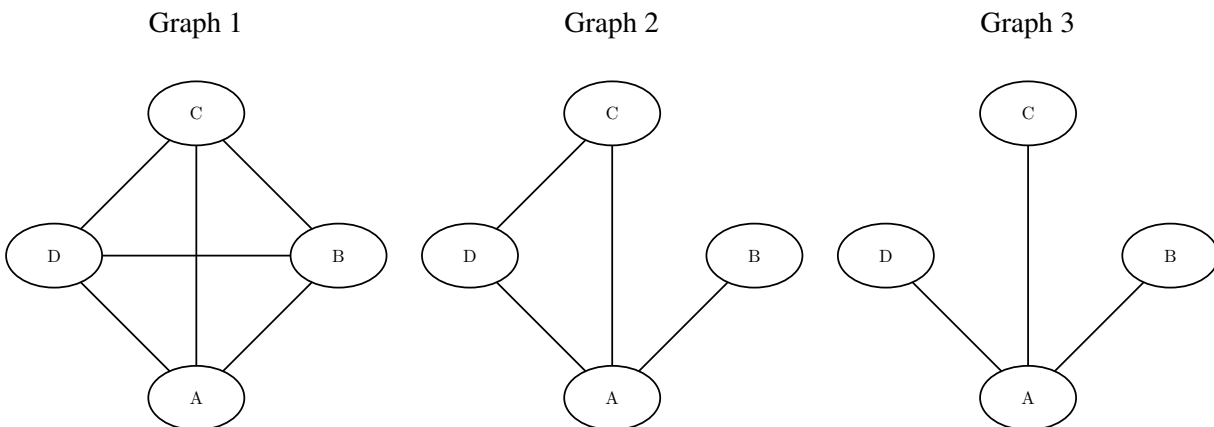
$$C(i) = \frac{|\{(u, v) \in A : u, v \in N_i\}|}{|K(N_i)|}$$

For a particular node  $i$ , the clustering coefficient determines how close the subgraph induced by its neighbor set  $N_i$  is to being a clique (complete subgraph). In social networks, a high clustering coefficient can help predict relationships that might not be known, confirmed, or realized yet. The fact that person  $i$  knows person  $j$  and person  $j$  knows person  $k$  does not guarantee that person  $i$  knows person  $k$ , but it is much more likely that person  $i$  knows person  $k$  than that person  $i$  knows some random person.

The clustering coefficient is calculated when you specify the **CLUSTERING\_COEF** option in the CENTRALITY statement. The results are provided in the node output data set that you specify in the OUT\_NODES= option in the PROC OPTGRAPH statement.

The algorithm that PROC OPTGRAPH uses to compute the clustering coefficient runs in time  $O(|N|^3)$ . Therefore, this algorithm is not expected to scale to very large graphs.

Consider the three undirected graphs on four nodes shown in Figure 1.24.

**Figure 1.24** Three Undirected Graphs

Define the three link data sets as follows:

```
data LinkSetInCC1;
  input from $ to $ @@;
  datalines;
A B  A C  A D
B C  B D  C D
;

data LinkSetInCC2;
  input from $ to $ @@;
  datalines;
A B  A C  A D
C D
;

data LinkSetInCC3;
  input from $ to $ @@;
  datalines;
A B  A C  A D
;
```

The following statements use three calls to PROC OPTGRAPH to calculate the clustering coefficients for each graph:

```
proc optgraph
  data_links = LinkSetInCC1
  out_nodes  = NodeSetOut1;
  centrality
    clustering_coef;
run;

proc optgraph
  data_links = LinkSetInCC2
  out_nodes  = NodeSetOut2;
  centrality
    clustering_coef;
run;
```

```

proc optgraph
  data_links = LinkSetInCC3
  out_nodes  = NodeSetOut3;
  centrality
    clustering_coef;
run;

```

The node data sets provide the clustering coefficients for each graph (variable `centr_cluster`), as shown in Figure 1.25 through Figure 1.27.

**Figure 1.25** Clustering Coefficient of a Simple Undirected Graph 1

node	centr_cluster
A	1
B	1
C	1
D	1

**Figure 1.26** Clustering Coefficient of a Simple Undirected Graph 2

node	centr_cluster
A	0.33333
B	0.00000
C	1.00000
D	1.00000

**Figure 1.27** Clustering Coefficient of a Simple Undirected Graph 3

node	centr_cluster
A	0
B	0
C	0
D	0

## Closeness Centrality

*Closeness centrality* is the reciprocal of the average of the shortest path (geodesic) distances from a particular node to all other nodes. Closeness can be thought of as a measure of how long it takes information to spread from a particular node to other nodes in the network.

Define  $d_{uv}$  to be the shortest path distance from node  $u$  to node  $v$ .

### **Closeness Centrality for an Undirected Graph**

For an undirected graph,  $R(u) = \{v \in N : d_{uv} < \infty\}$  is the set of *reachable nodes* from node  $u$ . The set of *unreachable nodes* from node  $u$  is  $N \setminus R(u) = \{v \in N : d_{uv} = \infty\}$ . The `CLOSE_NOPATH=` option specifies how to handle unreachable nodes.

For the special case in which all nodes are unreachable from node  $u$ , the closeness centrality is defined as 0. Otherwise, closeness centrality is calculated as

$$C_c(u) = s(u) \left( \frac{n(u)}{\sum_{v \in R(u)} d_{uv} + |N \setminus R(u)| p} \right)$$

where  $p$  defines a penalty parameter for unreachable nodes,  $n(u)$  defines the number of nodes that are considered in calculating the average, and  $s(u)$  is a scaling factor, as shown in Table 1.56.

**Table 1.56** Formulas for CLOSE\_NOPATH= Option for Undirected Graphs

CLOSE_NOPATH=	$p$	$n(u)$	$s(u)$
DIAMETER	$\max_{(i,j) \in A} \{d_{ij} : d_{ij} < \infty\} + 1$	$ N  - 1$	1
NNODES	$ N $	$ N  - 1$	1
ZERO	0	$ R(u)  - 1$	$\frac{ R(u) -1}{ N -1}$

### Closeness Centrality for a Directed Graph

For a directed graph,  $R^{\text{out}}(u) = \{v \in N : d_{uv} < \infty\}$  is the set of *reachable nodes* from node  $u$ , whereas  $R^{\text{in}}(u) = \{v \in N : d_{vu} < \infty\}$  is the set of nodes from which a finite path exists to node  $u$ . The set of *unreachable nodes* from node  $u$  is  $N \setminus R^{\text{out}}(u) = \{v \in N : d_{uv} = \infty\}$ , whereas the set of nodes from which a finite path to node  $u$  does not exist is  $N \setminus R^{\text{in}}(u) = \{v \in N : d_{vu} = \infty\}$ .

For the special case in which all nodes are unreachable from node  $u$ , the out-closeness centrality is defined as 0. Otherwise, out-closeness centrality is calculated as

$$C_c^{\text{out}}(u) = s^{\text{out}}(u) \left( \frac{n^{\text{out}}(u)}{\sum_{v \in R^{\text{out}}(u)} d_{uv} + |N \setminus R^{\text{out}}(u)| p} \right)$$

where  $n^{\text{out}}(u)$  defines the number of nodes that are considered in calculating the average and  $s^{\text{out}}(u)$  is a scaling factor, as shown in Table 1.57.

For the special case in which node  $u$  is unreachable from all other nodes, the in-closeness centrality is defined as 0. Otherwise, in-closeness centrality is calculated as

$$C_c^{\text{in}}(u) = s^{\text{in}}(u) \left( \frac{n^{\text{in}}(u)}{\sum_{v \in R^{\text{in}}(u)} d_{vu} + |N \setminus R^{\text{in}}(u)| p} \right)$$

where  $n^{\text{in}}(u)$  defines the number of nodes that are considered in calculating the average and  $s^{\text{in}}(u)$  is a scaling factor, as shown in Table 1.57.

**Table 1.57** Formulas for CLOSE\_NOPATH= Option for Directed Graphs

CLOSE_NOPATH=	$n^{\text{out}}(u)$	$s^{\text{out}}(u)$	$n^{\text{in}}(u)$	$s^{\text{in}}(u)$
DIAMETER	$ N  - 1$	1	$ N  - 1$	1
NNODES	$ N  - 1$	1	$ N  - 1$	1
ZERO	$ R^{\text{out}}(u)  - 1$	$\frac{ R^{\text{out}}(u) -1}{ N -1}$	$ R^{\text{in}}(u)  - 1$	$\frac{ R^{\text{in}}(u) -1}{ N -1}$



The overall closeness centrality for directed graphs is calculated as

$$C_c(u) = \frac{C_c^{\text{out}}(u) + C_c^{\text{in}}(u)}{2}$$

### Harmonic Centrality

*Harmonic centrality*, as described in Rochat (2009), is a variant of closeness centrality that attempts to simplify the treatment of unreachable nodes by calculating the average of the reciprocal of the shortest path distances from a particular node to all other nodes. The formula for harmonic centrality is

$$C_h(u) = \frac{1}{|N| - 1} \sum_{v \in N \setminus \{u\}} \frac{1}{d_{uv}}$$

To enable the calculation of harmonic centrality, use the CLOSE\_NOPATH=HARMONIC option.

Closeness centrality is calculated according to the value of the **CLOSE=** option in the CENTRALITY statement. The results are provided in the node output data set that you specify in the OUT\_NODES= option in the PROC OPTGRAPH statement. If CLOSE=WEIGHT (or BOTH), then the shortest paths are calculated with respect to the weighted graph. Because the metric uses shortest paths to determine closeness, the weight and the closeness metric are inversely related. In general, the lower the weight, the higher the contribution to the closeness metric.

The algorithm that PROC OPTGRAPH uses to compute closeness centrality relies on calculating shortest paths for all source-sink pairs and runs in time  $O(|N| \times (|N| \log |N| + |A|))$ . Therefore, this algorithm is not expected to scale to very large graphs. Because the shortest path calculations can be computed independently (for each source node), you can speed up the algorithm by specifying the NTHREADS= option in the PERFORMANCE statement.

Consider again the directed graph in [Figure 1.8](#) with the data set LinkSetIn, which is defined in the section “[Link Input Data](#)” on page 50. The following statements calculate the closeness centrality for both the weighted and unweighted graphs:

```
proc optgraph
  graph_direction = directed
  data_links      = LinkSetIn
  out_nodes       = NodeSetOut;
  centrality
    close         = both;
run;
```

The node data set NodeSetOut now contains the weighted and unweighted directed closeness centrality of the input graph. The data set provides the unweighted closeness (the centr\_close\_unwt variable), in-closeness (the centr\_close\_in\_unwt variable), and out-closeness (the centr\_close\_out\_unwt variable). It also provides the weighted variants centr\_close\_wt, centr\_close\_in\_wt, and centr\_close\_out\_wt. This data set is shown in [Figure 1.28](#).

**Figure 1.28** Closeness Centrality of a Simple Directed Graph

node	centr_close_wt	centr_close_in_wt	centr_close_out_wt	centr_close_unwt	centr_close_in_unwt	centr_close_out_unwt
A	0.08000	0.00000	0.16000	0.22222	0.00000	0.44444
B	0.11621	0.08696	0.14545	0.33333	0.22222	0.44444
C	0.11496	0.09877	0.13115	0.27885	0.25000	0.30769
D	0.12007	0.12903	0.11111	0.29178	0.30769	0.27586
E	0.12662	0.13559	0.11765	0.32000	0.32000	0.32000
F	0.11849	0.14286	0.09412	0.30725	0.34783	0.26667
G	0.10882	0.11765	0.10000	0.32500	0.40000	0.25000
H	0.10299	0.10959	0.09639	0.27885	0.30769	0.25000
I	0.06349	0.12698	0.00000	0.18182	0.36364	0.00000

## Betweenness Centrality

*Betweenness centrality* counts the number of times a particular node (or link) occurs on shortest paths between other nodes. Betweenness can be thought of as a measure of the control that a node (or link) has over the communication flow among the rest of the network. In this sense, the nodes (or links) that have high betweenness are the *gatekeepers* of information, because of their relative location in the network.

The formula for node betweenness centrality is

$$C_b(u) = \sum_{\substack{s \neq u \neq t \in N \\ s \neq t}} \frac{\sigma_{st}(u)}{\sigma_{st}}$$

where  $\sigma_{st}$  is the number of shortest paths from  $s$  to  $t$  and  $\sigma_{st}(u)$  is the number of shortest paths from  $s$  to  $t$  that pass through node  $u$ .

The formula for link betweenness centrality is

$$C_b(u, v) = \sum_{\substack{s, t \in N \\ s \neq t}} \frac{\sigma_{st}(u, v)}{\sigma_{st}}$$

where  $\sigma_{st}(u, v)$  is the number of shortest paths from  $s$  to  $t$  that pass through link  $(u, v)$ .

By default, this metric is normalized by dividing through by two times the number of pairs of nodes, not including  $u$ , which is  $(|N|-1)(|N|-2)$ . You can disable this normalization by using the **BETWEEN\_NORM=** option.

For directed graphs, because the paths are directed, only the *out-betweenness* is computed. To get the *in-betweenness*, you must reverse all the directions of the graph and run the procedure again. This can be accomplished by simply using the DATA\_LINKS\_VAR statement to reverse the interpretation of *from* and *to*.

Betweenness centrality is calculated according to the value of the **BETWEEN=** option in the CENTRALITY statement. The node betweenness results are provided in the node output data set that is specified in the OUT\_NODES= option in the PROC OPTGRAPH statement. The link betweenness results are provided in the link output data set that is specified in the OUT\_LINKS= option in the PROC OPTGRAPH statement. Like closeness, if BETWEEN=WEIGHT (or BOTH), then the calculation of shortest paths is done using the weighted graph. Because the metric uses shortest paths to determine betweenness, the weight and the

betweenness metric are inversely related. In general the lower the weight, the higher the contribution to the betweenness metric.

The algorithm used by PROC OPTGRAPH to compute betweenness centrality relies on calculating shortest paths for all source-sink pairs and runs in time  $O(|N| \times (|N| \log |N| + |A|))$ . Therefore, it is not expected to scale to very large graphs. Similar to closeness centrality, because shortest path computations can be calculated independently (for each source node), the algorithm can be sped up by using the NTHREADS= option in the PERFORMANCE statement. When closeness and betweenness centrality are run together, PROC OPTGRAPH calculates both metrics in one run.

Consider again the directed graph in Figure 1.8 with data set LinkSetIn defined in the section “Link Input Data” on page 50. The following statements calculate the betweenness centrality for both the weighted and unweighted graphs:

```
proc optgraph
  graph_direction = directed
  data_links      = LinkSetIn
  out_links       = LinkSetOut
  out_nodes       = NodeSetOut;
  centrality
    between       = both;
run;
```

The node data set NodeSetOut now contains the weighted (variable centr\_between\_wt) and unweighted (variable centr\_between\_unwt) node betweenness centrality of the input graph. This data set is shown in Figure 1.29.

**Figure 1.29** Node Betweenness Centrality of a Simple Directed Graph

node	centr_between_wt	centr_between_unwt
A	0.00000	0.00000
B	0.07738	0.07738
C	0.12202	0.00595
D	0.00000	0.00595
E	0.33482	0.17857
F	0.26786	0.26786
G	0.22321	0.21429
H	0.00000	0.00000
I	0.00000	0.00000

In addition, the link data set LinkSetOut contains the weighted (variable centr\_between\_wt) and unweighted (variable centr\_between\_unwt) link betweenness centrality of the input graph. This data set is shown in Figure 1.30.

**Figure 1.30** Link Betweenness Centrality of a Simple Directed Graph

from	to	weight	centr_between_wt	centr_between_unwt
A	B	1	0.09524	0.09524
A	C	2	0.04315	0.02381
A	D	4	0.00446	0.02381
B	C	1	0.11458	0.01786
B	E	2	0.08780	0.04167
B	F	5	0.00000	0.14286
C	E	1	0.22917	0.11310
D	E	1	0.08929	0.09524
E	D	1	0.06696	0.05357
E	F	2	0.35714	0.21429
F	G	6	0.32143	0.32143
G	H	1	0.12500	0.12500
G	I	1	0.13393	0.12500
H	G	2	0.02679	0.01786
H	I	3	0.00893	0.01786

For more detailed examples, see “[Example 1.3: Betweenness and Closeness Centrality for Computer Network Topology](#)” on page 193 and “[Example 1.4: Betweenness and Closeness Centrality for Project Groups in a Research Department](#)” on page 196.

## Eigenvector Centrality

*Eigenvector centrality* is an extension of degree centrality in which *centrality points* are awarded for each neighbor. However, not all neighbors are equally important. Intuitively, a connection to an important node should contribute more to the centrality score than a connection to a less important node. This is the basic idea behind eigenvector centrality. Eigenvector centrality of a node is defined to be proportional to the sum of the scores of all nodes that are connected to it. Mathematically, it is

$$x_i = \frac{1}{\lambda} \sum_{j \in \delta_i} x_j = \frac{1}{\lambda} \sum_{j \in N} A_{ij} x_j$$

where  $x_i$  is the eigenvector centrality of node  $i$ ,  $\lambda$  is a constant,  $\delta_i$  is the set of nodes that connect to node  $i$ , and  $A_{ij}$  is the weight of the link from node  $i$  to node  $j$ .

Eigenvector centrality can be written as an eigenvector equation in matrix form as

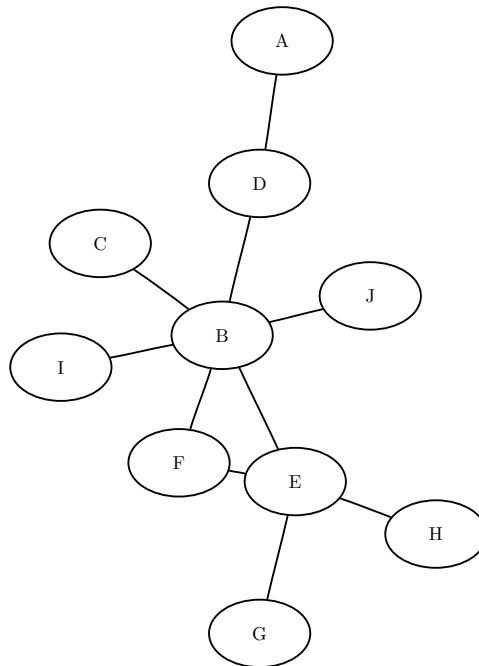
$$Ax = \lambda x$$

As the preceding equation shows,  $x$  is the eigenvector and  $\lambda$  is the eigenvalue. Because  $x$  should be positive, only the principal eigenvector that corresponds to the largest eigenvalue is of interest.

Eigenvector centrality is calculated according to the value that you specify in the **EIGEN=** option in the CENTRALITY statement. The results are provided in the node output data set that you specify in the OUT\_NODES= option in the PROC OPTGRAPH statement.

The following example illustrates the use of eigenvector centrality on the undirected graph  $G$  shown in Figure 1.31.

**Figure 1.31** Eigenvector Centrality Example of a Simple Undirected Graph



The graph can be represented by the following links data set LinkSetIn:

```

data LinkSetIn;
  input from $ to $ @@;
  datalines;
A D B C B D B E B F
B I B J E F E G E H
;

```

The following statements compute the eigenvector centrality:

```

proc optgraph
  data_links      = LinkSetIn
  out_nodes       = NodeSetOut;
  centrality
    eigen         = unweight;
run;

```

The data set NodeSetOut now contains the eigenvector centrality of each node. It is shown in Figure 1.32.

**Figure 1.32** Eigenvector Centrality Output

node	centr_eigen_unwt
B	1.00000
E	0.75919
F	0.61981
D	0.40226
I	0.35233
J	0.35233
C	0.35233
G	0.26749
H	0.26749
A	0.14173

Even though nodes F and D both have the same degree of 2, node F has a higher eigenvector centrality than node D. This is because node F links to two important nodes (B and E), whereas node D links to one important node (B) and one unimportant node (A).

For a more detailed example, see “[Example 1.5: Eigenvector Centrality for Word Sense Disambiguation](#)” on page 199.

## Hub and Authority Scoring

*Hub and authority centrality* was originally developed by Kleinberg (1998) to rank the importance of web pages. Certain web pages are important in the sense that they point to many important pages (called *hubs*). On the other hand, some web pages are important because they are pointed to by many important pages (called *authorities*). In other words, a good hub node is one that points to many good authorities, and a good authority node is one that is pointed to by many good hub nodes. This idea can be applied to many other types of graphs besides web pages. For example, it can be applied to a citation network for journal articles. A review article that cites many good authority papers has a high hub score, whereas a paper that is referenced by many other papers has a high authority score. The section “[Authority in U.S. Supreme Court Precedence](#)” on page 7 shows a similar example.

The authority centrality of a node is proportional to the sum of the hub centrality of nodes that point to it. Similarly, the hub centrality of a node is proportional to the sum of the authorities of nodes that it points to. That is,

$$x_i = \alpha \sum_{j \in N} A_{ij} y_j$$

$$y_i = \beta \sum_{j \in N} A_{ji} x_j$$

where  $x_i$  is the authority centrality of node  $i$ ,  $y_i$  is the hub centrality of node  $i$ ,  $A_{ij}$  is the weight of the link from node  $i$  to node  $j$ , and  $\alpha$  and  $\beta$  are constants.

The definition can be written in matrix form as follows:

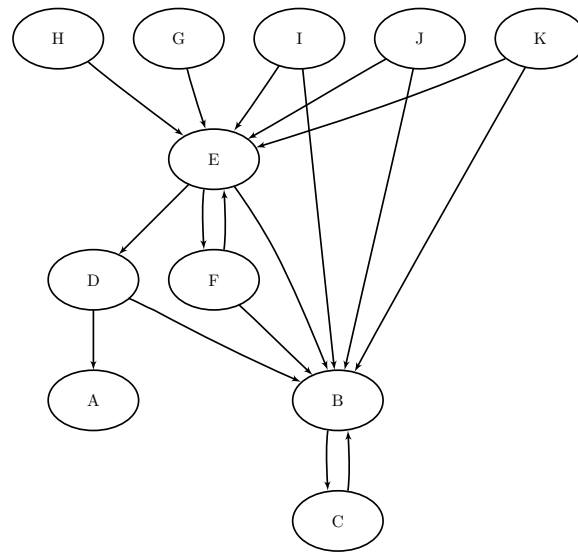
$$AA^T x = \lambda x$$

$$A^T A y = \lambda y$$

Thus, the authority and hub centralities are the principal eigenvectors of  $A^T A$  and  $AA^T$ , respectively. To solve this eigenvector problem, PROC OPTGRAPH provides two algorithms: the Jacobi-Davidson algorithm and the power method. You use the EIGEN\_ALGORITHM= option in the CENTRALITY statement to specify which algorithm to use. JACOBI\_DAVIDSON, which is the default, is a state-of-the-art package for solving large-scale eigenvalue problems (Sleijpen and van der Vorst 2000). The power method is one of the standard algorithms for solving eigenvalue problems, but it converges slowly for certain problems.

The following example illustrates the use of hub and authority scoring on the directed graph  $G$  shown in Figure 1.33. Each node represents a web page. If web page  $i$  has a hyperlink that points to web page  $j$ , then there is a directed link from  $i$  to  $j$ .

**Figure 1.33** Hub and Authority Centrality Example of a Simple Directed Graph



The graph can be represented by the following links data set LinkSetIn:

```

data LinkSetIn;
    input from $ to $ @@;
    datalines;
B C  C B  D A  D B  E B
E D  E F  F B  F E  G E
H E  I E  I B  J E  J B
K B  K E
;

```

The following statements compute hub and authority centrality:

```

proc optgraph
    graph_direction    = directed
    data_links         = LinkSetIn
    out_nodes          = NodeSetOut;
    centrality
        hub            = unweight
        auth           = unweight;
run;

```

The data set NodeSetOut now contains the hub and authority scores of each node. It is shown in [Figure 1.34](#).

**Figure 1.34** Hub and Authority Centrality Output

node	centr_hub_unwt	centr_auth_unwt
B	0.00000	1.00000
C	0.54135	0.00000
D	0.59703	0.11466
A	0.00000	0.10287
E	0.66549	0.84725
F	1.00000	0.11466
G	0.45865	0.00000
H	0.45865	0.00000
I	1.00000	0.00000
J	1.00000	0.00000
K	1.00000	0.00000

The output shows that nodes B and E have high authority scores because they have many incoming links. Nodes F, I, J, K have high hub scores because they all point to good authority nodes B and E.

## Weight Interpretation

In certain situations, you might want to calculate various centrality metrics on the same weighted graph. As described above, closeness and betweenness centrality have inverse relationships with the link weights, because these metrics are calculated using shortest paths. So the lower the weight, the higher the contribution to the centrality metric. All of the other metrics are direct relationships. That is, the higher the weight, the higher the contribution to the centrality metric.

To calculate these metrics in one invocation of PROC OPTGRAPH, you can use the **WEIGHT2=** option. The variable defined by this option is used as link weights for closeness and betweenness calculations whereas all other metrics use the standard weight variable.

For a more detailed example, see “[Example 1.6: Centrality Metrics for Project Groups in a Research Department](#)” on page 202, which uses the **WEIGHT2=** option.

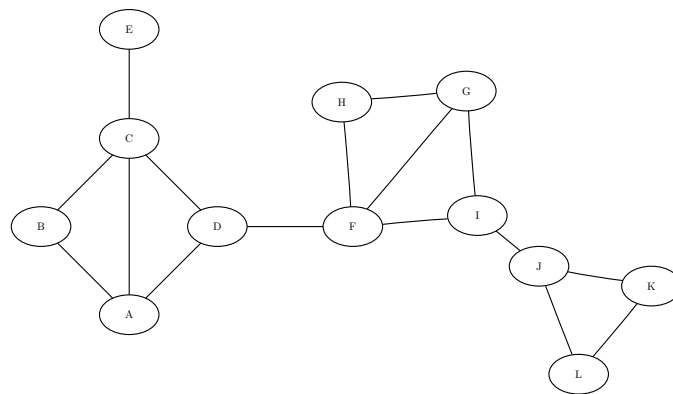
## Processing by Cluster

You can process a number of induced subgraphs of a graph with only one call to PROC OPTGRAPH by using the **BY\_CLUSTER** option in the CENTRALITY statement. This section shows an example of how to use this option.

### *Centrality by Cluster for a Simple Undirected Graph*

Consider the graph depicted in [Figure 1.35](#).

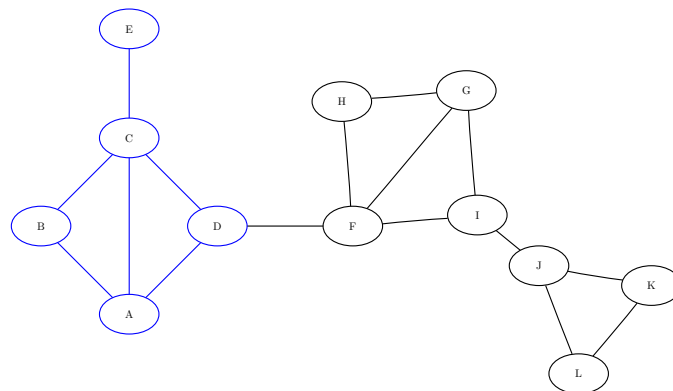


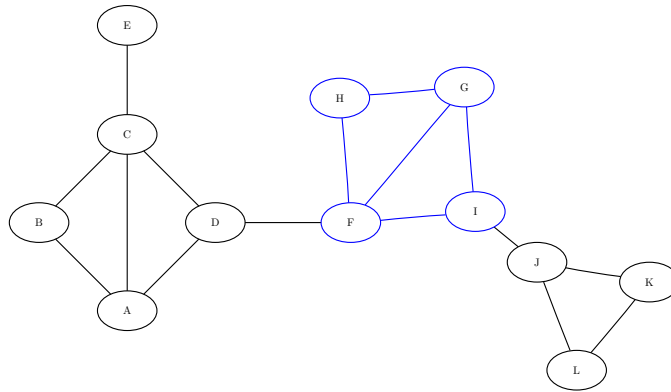
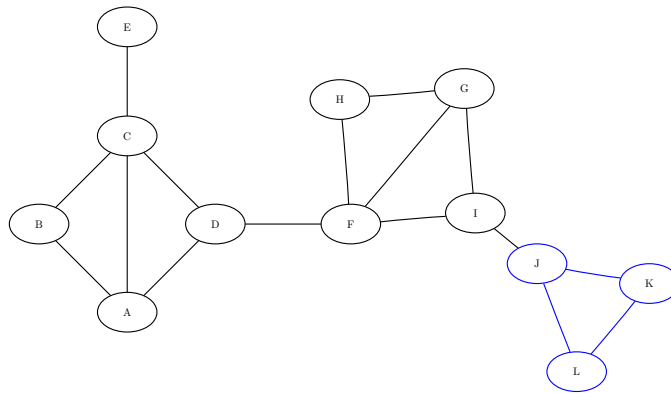
**Figure 1.35** Undirected Graph

The following statements create the data set LinkSetIn:

```
data LinkSetIn;
  input from $ to $ @@;
  datalines;
A B  A C  A D  B C  C D
C E  D F  F G  F H  F I
G H  G I  I J  J K  J L
K L
;
```

The graph seems to have three distinct parts, which are connected by just a few links. Assume that you have already partitioned the set into three sets of nodes:  $N^0 = \{A, B, C, D, E\}$ ,  $N^1 = \{F, G, H, I\}$ , and  $N^2 = \{J, K, L\}$ . The induced subgraphs on these three sets of nodes are shown in blue in Figure 1.36 through Figure 1.38. Notice that links that connect different partitions have been removed.

**Figure 1.36** Subgraph  $N^0 = \{A, B, C, D, E\}$ 

**Figure 1.37** Subgraph  $N^1 = \{F, G, H, I\}$ **Figure 1.38** Subgraph  $N^2 = \{J, K, L\}$ 

The following data sets define the three induced subgraphs:

```
data LinkSetIn0;
  input from $ to $ @@;
  datalines;
A B  A C  A D  B C  C D  C E
;
```

```
data LinkSetIn1;
  input from $ to $ @@;
  datalines;
F G  F H  F I  G H  G I
;
```

```
data LinkSetIn2;
  input from $ to $ @@;
  datalines;
J K  J L  K L
;
```

To calculate centrality metrics on the three subgraphs, you could run PROC OPTGRAPH three times, as follows:

```
proc optgraph
  data_links = LinkSetIn0
  out_nodes  = NodeSetOut0;
  centrality
    degree    = out
    influence  = unweight
    close     = unweight
    between   = unweight
    eigen     = unweight;
run;

proc optgraph
  data_links = LinkSetIn1
  out_nodes  = NodeSetOut1;
  centrality
    degree    = out
    influence  = unweight
    close     = unweight
    between   = unweight
    eigen     = unweight;
run;

proc optgraph
  data_links = LinkSetIn2
  out_nodes  = NodeSetOut2;
  centrality
    degree    = out
    influence  = unweight
    close     = unweight
    between   = unweight
    eigen     = unweight;
run;
```

This produces the results shown in [Figure 1.39](#) through [Figure 1.41](#).

**Figure 1.39** Centrality for Induced Subgraph 0

node	centr_degree_out	centr_eigen_unwt	centr_close_unwt	centr_between_unwt	centr_influence1_unwt	centr_influence2_unwt
A	3	0.89897	0.80000	0.08333	0.6	1.6
B	2	0.70711	0.66667	0.00000	0.4	1.4
C	4	1.00000	1.00000	0.58333	0.8	1.6
D	2	0.70711	0.66667	0.00000	0.4	1.4
E	1	0.37236	0.57143	0.00000	0.2	0.8

**Figure 1.40** Centrality for Induced Subgraph 1

node	centr_degree_out	centr_eigen_unwt	centr_close_unwt	centr_between_unwt	centr_influence1_unwt	centr_influence2_unwt
F	3	1.00000	1.00	0.16667	0.75	1.75
G	3	1.00000	1.00	0.16667	0.75	1.75
H	2	0.78078	0.75	0.00000	0.50	1.50
I	2	0.78078	0.75	0.00000	0.50	1.50

**Figure 1.41** Centrality for Induced Subgraph 2

node	centr_degree_out	centr_eigen_unwt	centr_close_unwt	centr_between_unwt	centr_influence1_unwt	centr_influence2_unwt
J	2	1	1	0	0.66667	1.33333
K	2	1	1	0	0.66667	1.33333
L	2	1	1	0	0.66667	1.33333

A much more efficient way to process these graphs is to define the partition by using the cluster variable in the nodes data set and using the BY\_CLUSTER option. Define the partitions of the original graph as follows:

```
data NodeSetIn;
  input node $ cluster @@;
  datalines;
A 0 B 0 C 0 D 0 E 0
F 1 G 1 H 1 I 1
J 2 K 2 L 2
;
```

Now, using one call to PROC OPTGRAPH, you can process all three induced subgraphs. In addition, because the processing of these subgraphs is completely independent, you can do the processing in parallel by using the NTHREADS= option in the PERFORMANCE statement.

```
proc optgraph
  loglevel      = moderate
  data_nodes    = NodeSetIn
  data_links    = LinkSetIn
  out_nodes     = NodeSetOut;
  performance
    nthreads    = 3;
  centrality
    by_cluster
      degree     = out
      influence  = unweight
      close      = unweight
      between    = unweight
      eigen      = unweight;
run;
%put &_OPTGRAPH_;
%put &_OPTGRAPH_CENTRALITY_;
```

Assuming that your machine has at least three cores, all three subgraphs are processed simultaneously with one call to PROC OPTGRAPH. The progress of the procedure is shown in Figure 1.42.

**Figure 1.42** PROC OPTGRAPH Log: Centrality by Cluster for a Simple Undirected Graph

---

```

NOTE: -----
NOTE: -----
NOTE: Running OPTGRAPH version 14.1.
NOTE: -----
NOTE: -----
NOTE: The OPTGRAPH procedure is executing in single-machine mode.
NOTE: -----
NOTE: -----
NOTE: Reading the nodes data set.
NOTE: There were 12 observations read from the data set WORK.NODESETIN.
NOTE: Reading the links data set.
NOTE: There were 16 observations read from the data set WORK.LINKSETIN.
NOTE: Data input used 0.01 (cpu: 0.00) seconds.
NOTE: Building the input graph storage used 0.00 (cpu: 0.00) seconds.
NOTE: The input graph storage is using 0.0 MBs (peak: 0.0 MBs) of memory.
NOTE: The number of nodes in the input graph is 12.
NOTE: The number of links in the input graph is 16.
NOTE: -----
NOTE: -----
NOTE: Processing centrality metrics by cluster.
NOTE: -----
NOTE: Using the CLUSTER variable from the DATA_NODES= data set to partition the graph.
NOTE: Links that cross subgraphs are ignored.
NOTE: -----
NOTE: Distribution of 3 subgraphs by number of nodes:
           3 subgraphs of size [    3,    10] (100.0%)
NOTE: -----
NOTE: Processing centrality metrics by cluster using 3 threads.

           Algorithm          SubGraphs  Complete      Cpu      Real
           centrality              3      100%      Time      Time
                                0.00      0.00
NOTE: Processing centrality metrics used 0.0 MBs of memory.
NOTE: Processing centrality metrics by cluster used 0.00 (cpu: 0.00) seconds.
NOTE: -----
NOTE: -----
NOTE: Creating nodes data set output.
NOTE: Data output used 0.00 (cpu: 0.00) seconds.
NOTE: -----
NOTE: -----
NOTE: The data set WORK.NODESETOUT has 12 observations and 8 variables.
STATUS=OK  CENTRALITY=OK
STATUS=OK  CPU_TIME=0.00  REAL_TIME=0.00

```

---

The results are shown in [Figure 1.43](#).

**Figure 1.43** Centrality for All Induced Subgraphs

node	cluster	centr_degree_out	centr_eigen_unwt	centr_close_unwt	centr_between_unwt	centr_influence1_unwt	centr_influence2_unwt
A	0	3	0.89897	0.80000	0.08333	0.60000	1.60000
B	0	2	0.70711	0.66667	0.00000	0.40000	1.40000
C	0	4	1.00000	1.00000	0.58333	0.80000	1.60000
D	0	2	0.70711	0.66667	0.00000	0.40000	1.40000
E	0	1	0.37236	0.57143	0.00000	0.20000	0.80000
F	1	3	1.00000	1.00000	0.16667	0.75000	1.75000
G	1	3	1.00000	1.00000	0.16667	0.75000	1.75000
H	1	2	0.78078	0.75000	0.00000	0.50000	1.50000
I	1	2	0.78078	0.75000	0.00000	0.50000	1.50000
J	2	2	1.00000	1.00000	0.00000	0.66667	1.33333
K	2	2	1.00000	1.00000	0.00000	0.66667	1.33333
L	2	2	1.00000	1.00000	0.00000	0.66667	1.33333

**Centrality by Community for a Simple Undirected Graph**

The partition defined in the data set *NodeSetIn* could have also been calculated by PROC OPTGRAPH using a method called *community detection*. This method is discussed in the section “Community Detection” on page 89. First, call the community detection method as follows:

```
proc optgraph
  data_links = LinkSetIn
  out_nodes  = Communities;
  community;
run;
```

The resulting output is a partition of the nodes of the original graph into *communities*. The data set *Communities* is shown in Figure 1.44.

**Figure 1.44** Communities for a Simple Undirected Graph

node	community_1
A	0
B	0
C	0
D	0
E	0
F	1
G	1
H	1
I	1
J	2
K	2
L	2

To calculate centrality by induced subgraph, you can simply use the communities output as the nodes data set input and use the *DATA\_NODES\_VAR* statement to define the cluster variable:

```

proc optgraph
  data_nodes    = Communities
  data_links    = LinkSetIn
  out_nodes     = NodeSetOut;
  data_nodes_var
    cluster     = community_1;
  performance
    nthreads    = 3;
  centrality
    by_cluster
      degree     = out
      influence   = unweight
      close      = unweight
      between    = unweight
      eigen      = unweight;
run;

```

This gives the same results as before, when you manually defined the partition. These results are shown in [Figure 1.43](#).

### **Centrality by Filtered Community for a Simple Undirected Graph**

In some situations, the community detection algorithm might find a large number of small communities. Those communities might not be relevant, and you might want to focus only on communities of a certain size. When you use the BY\_CLUSTER option, you can also use the `FILTER_SUBGRAPH=` option to ignore any subgraph whose number of nodes is less than or equal to a certain size. This can save on computation time, and the resulting output contains only the subgraphs of interest.

Returning to the data in the section “Centrality by Community for a Simple Undirected Graph” on page 84, you can use the filtering option as follows:

```

proc optgraph
  filter_subgraph = 3
  data_nodes      = Communities
  data_links      = LinkSetIn
  out_nodes       = NodeSetOut;
  data_nodes_var
    cluster       = community_1;
  performance
    nthreads      = 3;
  centrality
    by_cluster
      degree       = out
      influence     = unweight
      close        = unweight
      between      = unweight
      eigen        = unweight;
run;

```

The results, shown in [Figure 1.45](#), now contain only those subgraphs with node size greater than 3.

**Figure 1.45** Centrality for Some Induced Subgraphs

node	community_1	centr_degree_out	centr_eigen_unwt	centr_close_unwt	centr_between_unwt
A	0	3	0.89897	0.80000	0.08333
B	0	2	0.70711	0.66667	0.00000
C	0	4	1.00000	1.00000	0.58333
D	0	2	0.70711	0.66667	0.00000
E	0	1	0.37236	0.57143	0.00000
F	1	3	1.00000	1.00000	0.16667
G	1	3	1.00000	1.00000	0.16667
H	1	2	0.78078	0.75000	0.00000
I	1	2	0.78078	0.75000	0.00000

centr_influence1_unwt	centr_influence2_unwt
0.60	1.60
0.40	1.40
0.80	1.60
0.40	1.40
0.20	0.80
0.75	1.75
0.75	1.75
0.50	1.50
0.50	1.50

## Clique

A *clique* of a graph  $G = (N, A)$  is an induced subgraph that is a complete graph. Every node in a clique is connected to every other node in that clique. A *maximal clique* is a clique that is not a subset of the nodes of any larger clique. That is, it is a set  $C$  of nodes such that every pair of nodes in  $C$  is connected by a link and every node not in  $C$  is missing a link to at least one node in  $C$ . The number of maximal cliques in a particular graph can be very large and can grow exponentially with every node added. Finding cliques in graphs has applications in numerous industries including bioinformatics, social networks, electrical engineering, and chemistry.

You can find the maximal cliques of an input graph by invoking the CLIQUE statement. The options for this statement are described in the section “[CLIQUE Statement](#)” on page 25. The clique algorithm works only with undirected graphs.

The results of the clique algorithm are written to the output data set that is specified in the OUT= option in the CLIQUE statement. Each node of each clique is listed in the output data set along with the variable clique to identify the clique to which it belongs. A node can appear multiple times in this data set if it belongs to multiple cliques.

The clique algorithm reports status information in a macro variable called `_OPTGRAPH_CLIQU_`. For more information about this macro variable, see the section “[Macro Variable \\_OPTGRAPH\\_CLIQU\\_](#)” on page 177.

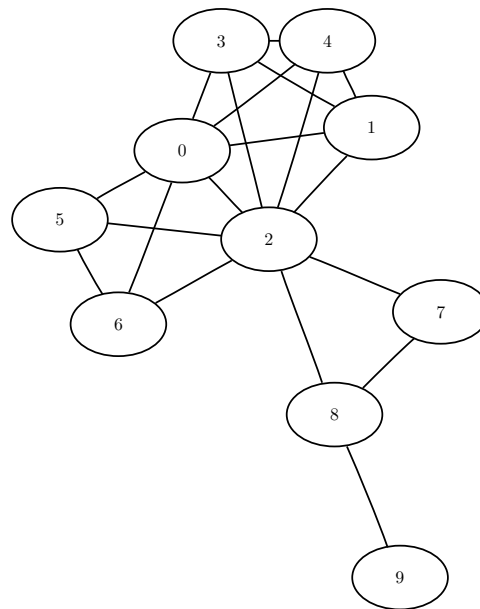
The algorithm that PROC OPTGRAPH uses to compute maximal cliques is a variant of the Bron-Kerbosch algorithm (Bron and Kerbosch 1973; Harley 2003). Enumerating all maximal cliques is NP-hard, so this algorithm typically does not scale to very large graphs.



## Maximal Cliques of a Simple Undirected Graph

This section illustrates the use of the clique algorithm on the simple undirected graph  $G$  that is shown in Figure 1.46.

**Figure 1.46** A Simple Undirected Graph  $G$



The undirected graph  $G$  can be represented by the following links data set LinkSetIn:

```
data LinkSetIn;
  input from to @@;
  datalines;
0 1 0 2 0 3 0 4 0 5
0 6 1 2 1 3 1 4 2 3
2 4 2 5 2 6 2 7 2 8
3 4 5 6 7 8 8 9
;
```

The following statements calculate the maximal cliques, output the results in the data set Cliques, and use the SQL procedure as a convenient way to create a table CliqueSizes of clique sizes:

```
proc optgraph
  data_links = LinkSetIn;
  clique
    out      = Cliques;
run;

proc sql;
  create table CliqueSizes as
  select clique, count(*) as size
  from Cliques
  group by clique
  order by size desc;
quit;
```

The data set Cliques now contains the maximal cliques of the input graph; it is shown in Figure 1.47.

**Figure 1.47** Maximal Cliques of a Simple Undirected Graph

cliq node	
1	0
1	2
1	1
1	3
1	4
2	0
2	2
2	5
2	6
3	2
3	8
3	7
4	8
4	9

In addition, the data set CliqueSizes contains the number of nodes in each clique; it is shown in Figure 1.48.

**Figure 1.48** Sizes of Maximal Cliques of a Simple Undirected Graph

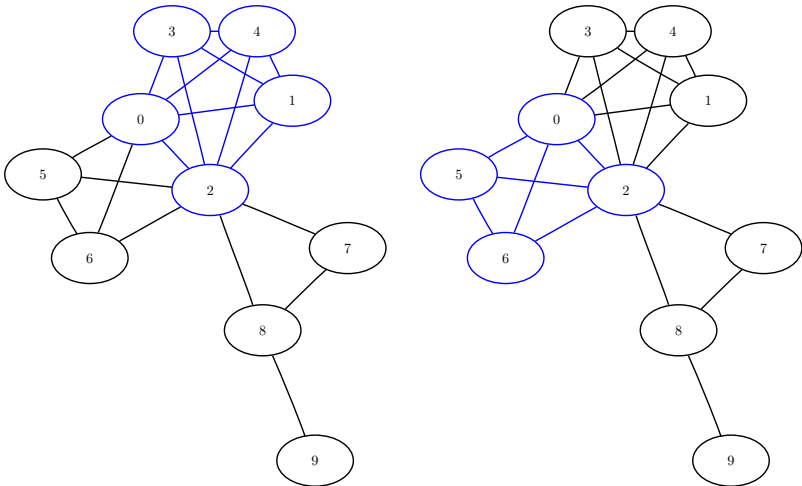
cliq size	
1	5
2	4
3	3
4	2

The maximal cliques are shown graphically in Figure 1.49 and Figure 1.50.

**Figure 1.49** Maximal Cliques  $C^1$  and  $C^2$

$C^1 = \{0, 1, 2, 3, 4\}$

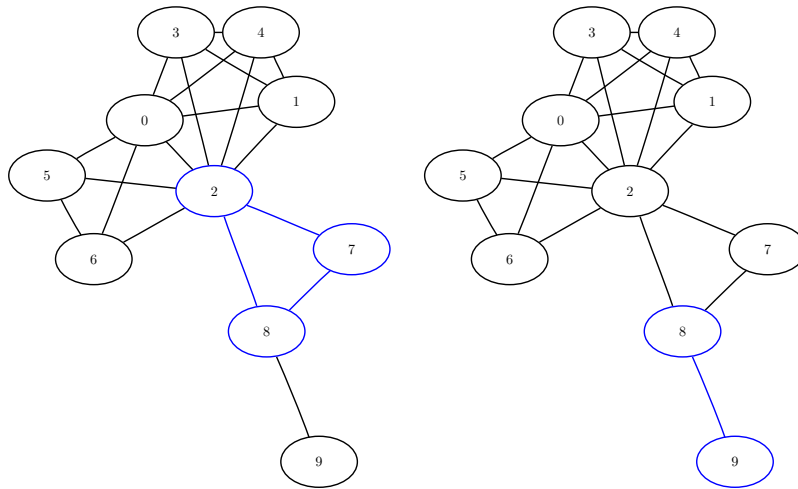
$C^2 = \{0, 2, 5, 6\}$



**Figure 1.50** Maximal Cliques  $C^3$  and  $C^4$ 

$$C^3 = \{2, 7, 8\}$$

$$C^4 = \{8, 9\}$$



## Community Detection

*Community detection* partitions a graph into communities such that the links within the community subgraphs are more densely connected than the links between communities.

In PROC OPTGRAPH, community detection can be determined by using the COMMUNITY statement. The options for this statement are described in the section “[COMMUNITY Statement](#)” on page 26.

The COMMUNITY statement reports status information in a macro variable called \_OPTGRAPH\_COMMUNITY\_. For more information about this macro variable, see the section “[Macro Variable \\_OPTGRAPH\\_COMMUNITY\\_](#)” on page 178.

When you specify `ALGORITHM=PARALLEL_LABEL_PROP` in the COMMUNITY statement, community detection supports both undirected and directed graphs. When you specify `ALGORITHM=LOUVAIN` or `ALGORITHM=LABEL_PROP` in the COMMUNITY statement, community detection is supported only for undirected graphs. For directed graphs, you need to aggregate directed links into undirected links before you call the algorithm. For example, suppose there are two directed links: a link from  $i$  to  $j$  with a link weight of 4.3, and a link from  $j$  to  $i$  with a link weight of 3.2. One common aggregation strategy is to sum the link weights together. Using this strategy, the weight of the undirected link between  $i$  and  $j$  is 7.5.

PROC OPTGRAPH implements three heuristic algorithms for finding communities: the LOUVAIN algorithm proposed in Blondel et al. (2008), the label propagation algorithm proposed in Raghavan, Albert, and Kumara (2007), and the parallel label propagation algorithm developed by SAS (patent pending).

Given a graph  $G = (N, A)$ , all three algorithms run in time  $O(k|A|)$ , where  $k$  is the average number of links per node. The Louvain algorithm aims to optimize modularity, which is one of the most popular merit functions for community detection. Modularity is a measure of the quality of a division of a graph into communities. The *modularity* of a division is defined to be the fraction of the links that fall within the communities minus the expected fraction if the links were distributed at random, assuming that you do not change the degree of each node.

Mathematically, modularity is defined as

$$Q = \frac{1}{2w} \sum_{(u,v) \in A} \left( w_{uv} - \frac{w_u w_v}{2w} \right) \Delta(c_u, c_v)$$

$$w = \sum_{(u,v) \in A} w_{uv}$$

$$w_u = \sum_{v \in \delta_u} w_{uv}$$

where  $Q$  is the modularity,  $w_{uv}$  is the link weight between node  $u$  and  $v$ ,  $\delta_u$  is the set of nodes that connects to node  $u$ ,  $w_u$  is the sum of link weights incident to node  $u$ ,  $w$  is the sum of link weights of the graph,  $c_u$  is the community to which node  $u$  belongs, and  $\Delta(c_u, c_v)$  is the Kronecker delta symbol, defined as

$$\Delta(c_u, c_v) = \begin{cases} 1 & \text{if } c_u = c_v \\ 0 & \text{otherwise} \end{cases}$$

The following is a brief description of the Louvain algorithm:

1. Initialize each node as its own community.
2. Move each node from its current community to the neighboring community that increases modularity the most. Repeat this step until modularity cannot be improved.
3. Group the nodes in each community into a supernode. Construct a new graph based on supernodes. Repeat these steps until modularity cannot be further improved or the maximum number of iterations has been reached.

The more recently proposed label propagation algorithm moves a node to a community that most of its neighbors belong to. Extensive testing by Lancichinetti and Fortunato (2009) has empirically demonstrated that the label propagation algorithm performs as well as the Louvain method in most cases.

The following is a brief description of the label propagation algorithm:

1. Initialize each node as its own community.
2. Move each node from its current community to the neighboring community that has the maximum sum of link weights to the current node; break ties randomly if necessary. Repeat this step until there are no more movements.

The parallel label propagation algorithm is an extension of the basic label propagation algorithm. During each iteration, rather than updating node labels sequentially, nodes update their labels simultaneously by using the node label information from the previous iteration. In this approach, node labels can be updated in parallel. However, simultaneous updating of this nature often leads to oscillating labels because of the bipartite subgraph structure often present in large graphs. To address this issue, at each iteration the parallel algorithm skips the labeling step at some randomly chosen nodes in order to break the bipartite structure. You can control the random samples that the algorithm takes by specifying the `RANDOM_FACTOR=` or `RANDOM_SEED=` options in the `COMMUNITY` statement.

As you can see from their descriptions, all three algorithms adopt a heuristic local optimization approach. The final result often depends on the sequence of nodes that are presented in the links input data set. Therefore, if the sequence of nodes in the links data set has been changed, the result is likely to be slightly different.

## Parallel Community Detection

Parallel community detection can be invoked by specifying `ALGORITHM=PARALLEL_LABEL_PROP` in the `COMMUNITY` statement. The computation is executed with multiple threads on a single computer. The number of threads being used can be controlled by specifying the `NTHREADS=` option in the `PERFORMANCE` statement.

The following statements demonstrate how to invoke parallel community detection using eight threads:

```
proc optgraph
  data_links      = links
  graph_direction = directed
  out_nodes       = outNodes;
  performance
    nthreads      = 8;
  community
    algorithm      = parallel_label_prop
    out_community  = outComm;
run;
```

## Memory Requirement

When you specify `GRAPH_INTERNAL_FORMAT=THIN` in the `PROC OPTGRAPH` statement and `ALGORITHM=LOUVAIN` or `ALGORITHM=LABEL_PROP` in the `COMMUNITY` statement, the memory (number of bytes) required for community detection can be estimated approximately as follows given a graph  $G = (N, A)$ :

$$(2 \times |A| + |N|) \times \text{sizeof(int)} + (3 \times |A| + |N|) \times \text{sizeof(double)}$$

When you specify `GRAPH_INTERNAL_FORMAT=THIN` and `ALGORITHM=PARALLEL_LABEL_PROP`, the memory required for community detection is approximately twice this amount.

Assume that your machine architecture is such that an integer is 4 bytes and a double is 8 bytes. Then, a graph with 100 million nodes and 650 million links would require approximately 21 gigabytes (GB) of memory when you specify `ALGORITHM=LOUVAIN` or `ALGORITHM=LABEL_PROP`:

$$(2 \times 650M + 100M) \times 4 + (3 \times 650M + 100M) \times 8 = 21GB$$

The same graph would require approximately 42 GB if you specify `ALGORITHM=PARALLEL_LABEL_PROP`.

This is only an estimate for the amount of memory that is required. `PROC OPTGRAPH` itself might require more memory to maintain the input and output data structures. In addition, other running processes might take away from the available memory.

`PROC OPTGRAPH` uses significantly more memory if `GRAPH_INTERNAL_FORMAT=FULL`. It is recommended that you use `GRAPH_INTERNAL_FORMAT=THIN` when you apply community detection on large graphs.

## Graph Direction

If you specify `ALGORITHM=PARALLEL_LABEL_PROP` in the `COMMUNITY` statement, community detection supports both undirected and directed graphs. However, you should be careful in deciding whether to model your problem as an undirected or a directed graph. For an undirected graph, the algorithm finds

communities based on the density of the subgraphs. For a directed graph, the algorithm finds communities based on the information flow along the directed links. That is, the algorithm propagates the community ID along the outgoing links of a node. Therefore, nodes are likely to be in the same community if they form circles along the outgoing links. If the directed graph lacks this circle structure, the nodes are likely to switch between communities during the computation. As a result, the algorithm does not converge well and cannot find a good community structure in the graph.

## Large Community

It has often been observed in practice that the number of nodes contained in communities (produced by community detection algorithms) usually follows a power law distribution. That is, a few communities contain a very large number of nodes, whereas most communities contain a small number of nodes. This is especially true for large graphs. PROC OPTGRAPH provides two approaches to alleviate this problem: one uses the RECURSIVE option, and the other uses the RESOLUTION\_LIST= option.

### Recursive

You can apply the RECURSIVE option to recursively break large communities into smaller ones. At the first step, PROC OPTGRAPH processes data as if no RECURSIVE option were specified. At the end of this step, it checks whether the community result satisfies the RECURSIVE option criteria. If the community result satisfies these criteria, PROC OPTGRAPH stops iterations and outputs results. Otherwise, it treats each large community as an independent graph and recursively applies community detection on top of it.

In certain cases, a community is not further split even if it does not meet the recursive criteria that you specified. One example is a star-shaped community that contains 200 nodes while MAX\_COMM\_SIZE is specified as 100; another example is a symmetric community whose diameter is 2 while MAX\_DIAMETER is specified as 1.

### Resolution List

The second way to combat the problem, provided you have specified ALGORITHM=LOUVAIN in the COMMUNITY statement, is to assign a larger value than the default value of 1 to the RESOLUTION\_LIST= option. When ALGORITHM=LOUVAIN, the value assigned to the RESOLUTION\_LIST= option can be interpreted as follows: Suppose the resolution value is  $x$ . Two communities are merged if the sum of the weights of intercommunity links is at least  $x$  times the expected value of the same sum if the graph is reconfigured randomly. Therefore, a larger resolution value produces more communities, each of which contains a smaller number of nodes. However, there is no explicit formula to detail the number of nodes in communities with respect to the resolution value. You must use trial and error to get to the expected community size. More information about resolution value is available in Ronhovde and Nussinov 2010.

If you specify ALGORITHM=LOUVAIN, you can supply multiple resolution values at one time. If you supply multiple resolution values at one time, PROC OPTGRAPH detects communities at the highest resolution level first, then merges communities at a lower resolution, and repeats the process until it reaches the lowest level. This process enables you to see how the communities are merged at different levels. Due to the local nature of this optimization algorithm, two different runs do not produce the same result if the two runs share a common resolution level. For example, the algorithm can produce different results at resolution 0.5 in two runs: one with RESOLUTION\_LIST = 1 0.7 0.5, and the other with RESOLUTION\_LIST = 1 0.5.

If you specify ALGORITHM=PARALLEL\_LABEL\_PROP in the COMMUNITY statement, the resolution value can be interpreted as the minimal density of communities in an undirected and unweighted graph. The *density* of a community is defined as the number of links inside the community divided by the total number

of possible links. A larger resolution value likely results in communities that contain fewer nodes. For more information about resolution values for label propagation, see Traag, Van Dooren, and Nesterov (2011).

If you supply multiple resolution values at one time and you specify `ALGORITHM=PARALLEL_LABEL_PROP`, the `OPTGRAPH` procedure performs community detection multiple times, each time with a different resolution value. This is equivalent to calling the `OPTGRAPH` procedure several times, each time with a different (single) resolution value specified for the `RESOLUTION_LIST=` option.

If you specify `ALGORITHM=PARALLEL_LABEL_PROP` in the `COMMUNITY` statement, the value that is specified in the `RESOLUTION_LIST=` option has a major impact on the running time of the algorithm. When a large resolution value is specified, the algorithm is likely to create many tiny communities, and nodes are likely to change communities between iterations. Therefore, the algorithm might not converge properly. On the other hand, when the resolution value is small, the algorithm might find some very large communities, such as a community that contains more than a million nodes. In this case, if you specify the `RECURSIVE` option, the algorithm spends a long time in the recursive step in order to break large communities into smaller ones.

The recommended approach is to first experiment with a set of resolution values without using the `RECURSIVE` option. At the end of the run, examine the resulting modularity values and the community size distributions. Remove the resolution values that lead to small modularity values or huge communities. Then add the `RECURSIVE` option to the `COMMUNITY` statement, if desired, and run `PROC OPTGRAPH` again.

“[Example 1.7: Community Detection on Zachary’s Karate Club Data](#)” on page 205 shows the use of the `RESOLUTION_LIST=` option in the calculation of communities.

### Large Graphs

When you are dealing with large graphs, the following practices are recommended:

- Use `GRAPH_INTERNAL_FORMAT=THIN` instead of `GRAPH_INTERNAL_FORMAT=FULL`. This enables `PROC OPTGRAPH` to store the data in memory compactly.
- Use the `LINK_REMOVAL_RATIO=` option to remove unimportant links. This practice can often dramatically improve the running time of large graphs.

### Output Data Sets

Community detection produces up to five output data sets. In these data sets, if you specify `ALGORITHM=LOUVAIN` or `ALGORITHM=LABEL_PROP` in the `COMMUNITY` statement, resolution level numbers are in decreasing order of the values that are specified in the `RESOLUTION_LIST=` option. That is, resolution level 1 corresponds to the largest value specified in the `RESOLUTION_LIST=` option, and resolution level  $K$  corresponds to the smallest value specified in the `RESOLUTION_LIST=` option. For example, if `RESOLUTION_LIST=2.5 3.1 0.6`, then resolution level 1 is at value 3.1, resolution level 2 is at value 2.5, and resolution level 3 is at value 0.6.

If you specify `ALGORITHM=PARALLEL_LABEL_PROP` in the `COMMUNITY` statement, resolution level numbers are in the same order as the values that are specified in the `RESOLUTION_LIST=` option. For example, if `RESOLUTION_LIST=0.001 0.005 0.01`, then resolution level 1 is at value 0.001, resolution level 2 is at value 0.005, and resolution level 3 is at value 0.01.

**OUT\_NODES= Data Set**

This data set describes the community identifier of each node. If multiple resolution values have been specified, the data set reports the community identifier of each node at each resolution level. The data set contains the following columns:

- node: node label
- community\_*i*: community identifier at resolution level *i*, where *i* is the resolution level number as previously described. There are *K* such columns if *K* different values are specified in the RESOLUTION\_LIST= option.

**OUT\_LEVEL= Data Set**

This data set describes the number of communities and their corresponding modularity values at various resolution levels. It contains the following columns:

- level: resolution level number
- resolution: resolution value
- communities: number of communities at the current resolution level
- modularity: modularity value at the current resolution level

**OUT\_COMMUNITY= Data Set**

This data set describes the number of nodes in each community. It contains the following columns:

- level: resolution level number
- resolution: resolution value
- community: community identifier
- nodes: number of nodes contained in the community

**OUT\_OVERLAP= Data Set**

This data set describes the intensity of each node. At the end of community detection, a node could have links that connect to multiple communities. The intensity of a node is computed as the sum of the link weights that connect to nodes in the specified community divided by the total link weights of the node. This data set is computationally expensive to produce, and it requires a large amount of disk space. Therefore, this data set is not produced if you specify multiple resolution values in the RESOLUTION\_LIST= option and ALGORITHM=PARALLEL\_LABEL\_PROP. However, if ALGORITHM=LOUVAIN, the data set is produced and will only contain results corresponding to the smallest value of the RESOLUTION\_LIST= option. The data set contains the following columns:

- node: node label
- community: community identifier
- intensity: intensity of the node that belongs to the community



**OUT\_COMM\_LINKS= Data Set**

This data set describes how communities are connected. If you specify `ALGORITHM=LOUVAIN` or `ALGORITHM=LABEL_PROP` in the `COMMUNITY` statement, this data set contains the following columns:

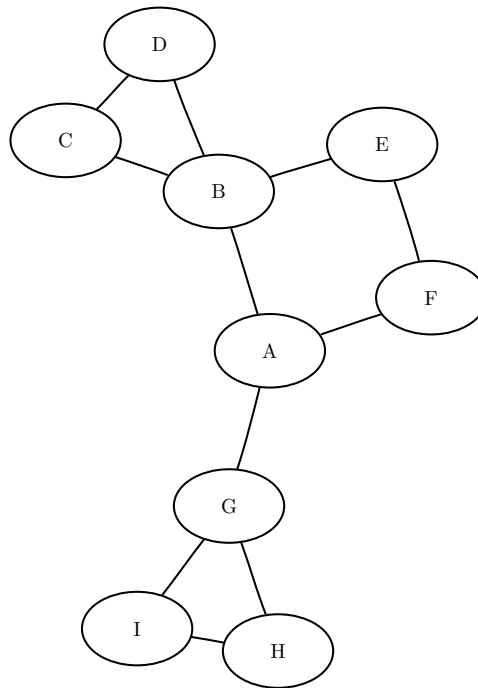
- `level`: resolution level number
- `resolution`: resolution value
- `from_community`: community identifier of the *from* community
- `to_community`: community identifier of the *to* community
- `link_weight`: sum of link weights of all links between `from_community` and `to_community`

This data set is not produced if you specify `ALGORITHM=PARALLEL_LABEL_PROP` together with multiple resolution values in the `RESOLUTION_LIST=` option.

**Community Detection on a Simple Graph**

This section illustrates the use of the community detection algorithm on the simple undirected graph  $G$  that is shown in Figure 1.51.

**Figure 1.51** A Simple Undirected Graph  $G$



The undirected graph  $G$  can be represented by the following links data set `LinkSetIn`:

```

data LinkSetIn;
  input from $ to $ @@;
  datalines;
A B  A F  A G  B C  B D

```

```

B E C D E F G I G H
H I
;

```

The following statements perform community detection and output the results in the specified data sets. The Louvain algorithm is used by default because no value is specified for the **ALGORITHM=** option.

```

proc optgraph
  data_links    = LinkSetIn
  out_nodes     = NodeSetOut;
  community
    resolution_list = 1.0 0.5
    out_level       = CommLevelOut
    out_community   = CommOut
    out_overlap     = CommOverlapOut
    out_comm_links  = CommLinksOut;
run;

```

The data set NodeSetOut contains the community identifier of each node and is shown in Figure 1.52.

**Figure 1.52** Community Detection on a Simple Graph: Nodes Output

node	community_1	community_2
A	0	0
B	1	0
F	0	0
G	2	1
C	1	0
D	1	0
E	0	0
I	2	1
H	2	1

The data set CommLevelOut contains summary information at each resolution level and is shown in Figure 1.53.

**Figure 1.53** Community Detection on a Simple Graph: Level Output

level	resolution	communities	modularity
1	1.0	3	0.39256
2	0.5	2	0.34298

The data set CommOut contains the number of nodes in each community and is shown in Figure 1.54.

**Figure 1.54** Community Detection on a Simple Graph: Community Summary

level	resolution	community	nodes
1	1.0	0	3
1	1.0	1	3
1	1.0	2	3
2	0.5	0	6
2	0.5	1	3

The data set CommOverlapOut contains community overlap information and is shown in Figure 1.55.

**Figure 1.55** Community Detection on a Simple Graph: Community Overlap

node	community	intensity
A	0	0.66667
A	1	0.33333
B	0	1.00000
F	0	1.00000
G	0	0.33333
G	1	0.66667
C	0	1.00000
D	0	1.00000
E	0	1.00000
I	1	1.00000
H	1	1.00000

The data set CommLinksOut describes how the communities are connected and is shown in Figure 1.56.

**Figure 1.56** Community Detection on a Simple Graph: Intercommunity Links

level	resolution	from_community	to_community	link_weight
1	1.0	0	1	2
1	1.0	0	2	1
2	0.5	0	1	1

## Connected Components

A *connected component* of a graph is a set of nodes that are all reachable from each other. That is, if two nodes are in the same component, then there exists a path between them. For a directed graph, there are two types of components: a *strongly connected component* has a directed path between any two nodes, and a *weakly connected component* ignores direction and requires only that a path exist between any two nodes.

In PROC OPTGRAPH, you can invoke connected components by using the CONCOMP statement. The options for this statement are described in the section “[CONCOMP Statement](#)” on page 28.

There are two main algorithms for finding connected components in an undirected graph: a depth-first search algorithm (ALGORITHM=DFS) and a union-find algorithm (ALGORITHM=UNION\_FIND). For a graph  $G = (N, A)$ , both algorithms run in time  $O(|N| + |A|)$  and can usually scale to very large graphs. The default is the union-find algorithm. For directed graphs, only the depth-first search algorithm is available.

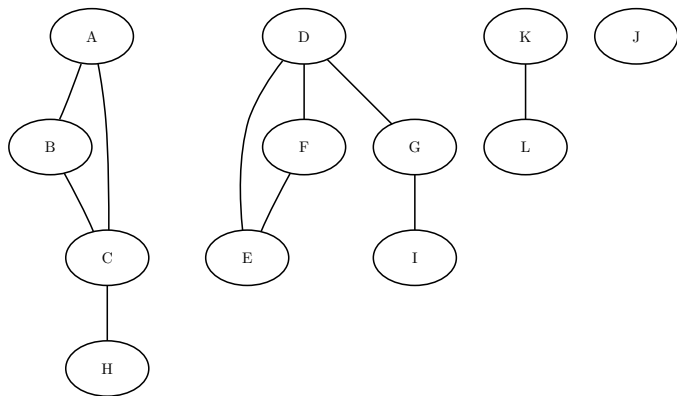
The results of the connected components algorithm are written to the output node data set that you specify in the OUT\_NODES= option in the PROC OPTGRAPH statement. For each node in the node data set, the variable concomp identifies its component. The component identifiers are numbered sequentially starting from 1.

The connected components algorithm reports status information in a macro variable called \_OPTGRAPH\_CONCOMP\_. For more information about this macro variable, see the section “[Macro Variable \\_OPTGRAPH\\_CONCOMP\\_](#)” on page 178.

Connected Components of a Simple Undirected Graph

This section illustrates the use of the connected components algorithm on the simple undirected graph *G* that is shown in Figure 1.57.

Figure 1.57 A Simple Undirected Graph *G*



The undirected graph *G* can be represented by the following links data set, LinkSetIn:

```
data LinkSetIn;
  input from $ to $ @@;
  datalines;
A B A C B C C H D E D F D G F E G I K L
;
```

The following statements calculate the connected components and output the results in the data set NodeSetOut:

```
proc optgraph
  data_links = LinkSetIn
  out_nodes = NodeSetOut;
  concomp;
run;
```

The data set NodeSetOut contains the connected components of the input graph and is shown in Figure 1.58.

Figure 1.58 Connected Components of a Simple Undirected Graph

node	concomp
A	1
B	1
C	1
H	1
D	2
E	2
F	2
G	2
I	2
K	3
L	3

Notice that the graph is defined by using only the links data set. As seen in [Figure 1.57](#), this graph also contains a singleton node labeled J, which has no associated links. By definition, this node defines its own component. But because the input graph was defined by using only the links data set, it did not show up in the results data set. To define a graph by using nodes that have no associated links, you should also define the input nodes data set. In this case, define the nodes data set `NodeSetIn` as follows:

```
data NodeSetIn;
  input node $ @@;
  datalines;
A B C D E F G H I J K L
;
```

Now, when you calculate the connected components, you define the input graph by using both the nodes input data set and the links input data set:

```
proc optgraph
  data_nodes = NodeSetIn
  data_links = LinkSetIn
  out_nodes  = NodeSetOut;
  concomp;
run;
```

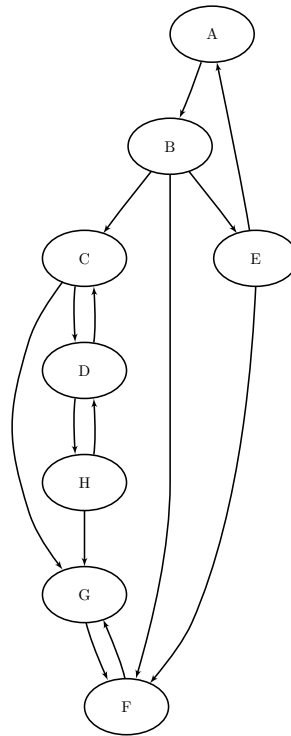
The resulting data set, `NodeSetOut`, includes the singleton node J as its own component, as shown in [Figure 1.59](#).

**Figure 1.59** Connected Components of a Simple Undirected Graph

node	concomp
A	1
B	1
C	1
D	2
E	2
F	2
G	2
H	1
I	2
J	3
K	4
L	4

## Connected Components of a Simple Directed Graph

This section illustrates the use of the connected components algorithm on the simple directed graph *G* that is shown in [Figure 1.60](#).

**Figure 1.60** A Simple Directed Graph *G*

The directed graph *G* can be represented by the following links data set, LinkSetIn:

```

data LinkSetIn;
  input from $ to $ @@;
  datalines;
A B  B C  B E  B F  C G
C D  D C  D H  E A  E F
F G  G F  H G  H D
;

```

The following statements calculate the connected components and output the results in the data set NodeSetOut:

```

proc optgraph
  graph_direction = directed
  data_links      = LinkSetIn
  out_nodes       = NodeSetOut;
  concomp;
run;

```

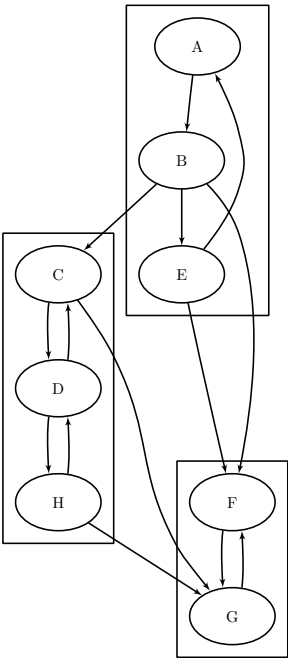
The data set NodeSetOut, shown in [Figure 1.61](#), now contains the connected components of the input graph.

**Figure 1.61** Connected Components of a Simple Directed Graph

node	concomp
A	3
B	3
C	2
E	3
F	1
G	1
D	2
H	2

The connected components are represented graphically in [Figure 1.62](#).

**Figure 1.62** Strongly Connected Components of Graph *G*



## Core Decomposition

An alternative to community detection for detecting cohesive subgroups is a method for extracting *k*-cores, known as *core decomposition*. Although this method is generally not as powerful as community detection for extracting a detailed community structure, it can give a coarse approximation of cohesive structure at a very low computational cost. Let  $G = (N, A)$  define a graph with nodes  $N$  and links  $A$ , and let  $G_S = (S, A_S)$  be an induced subgraph on nodes  $S$ . The subgraph  $G_S$  is a *k*-core if and only if for every node  $v \in S$ , the degree of  $v$  is greater than or equal to  $k$  and  $G_S$  is the maximum subgraph with this property. By definition, the cores are nested. That is, if  $G_{S_k}$  is a *k*-core of size  $k$ , then  $G_{S_{k+1}}$  is contained in  $G_{S_k}$ .

In PROC OPTGRAPH, you can invoke core decomposition by using the CORE statement. The options for this statement are described in the section “[CORE Statement](#)” on page 29.

The results of the core decomposition algorithm are given in the output node data set that is specified in the OUT\_NODES= option in the PROC OPTGRAPH statement. For each node in the node data set, the variable core\_out identifies its *core number*, the highest-order core that contains this node. The core identifiers are numbered sequentially starting from 0.

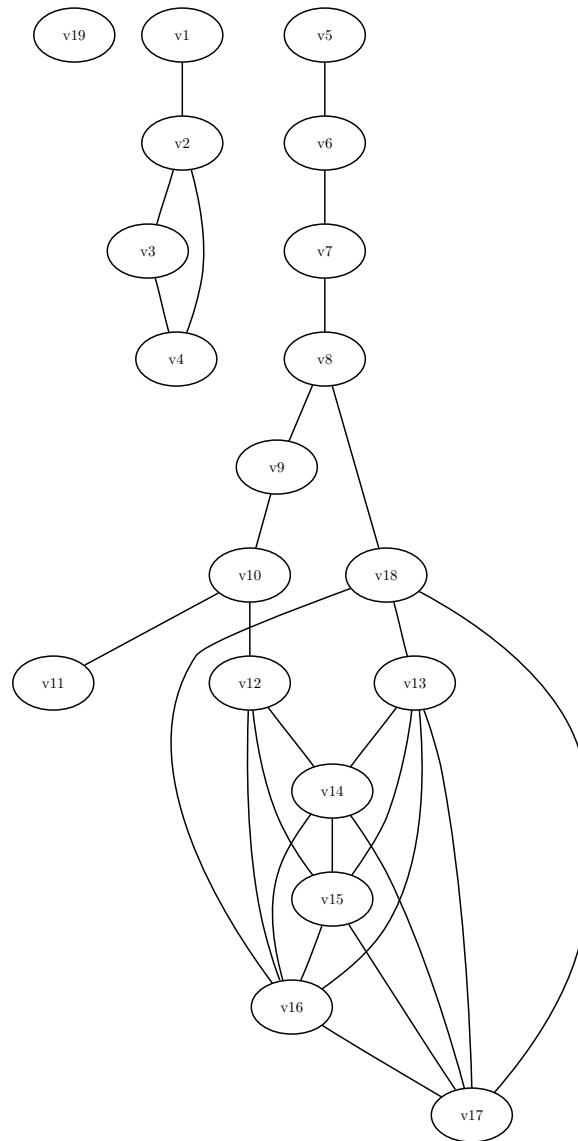
The core decomposition algorithm reports status information in a macro variable called \_OPTGRAPH\_CORE\_. For more information about this macro variable, see the section “[Macro Variable \\_OPTGRAPH\\_CORE\\_](#)” on page 179.

The algorithm used for core decomposition is based on the work presented in Batagelj and Zaversnik 2003. This algorithm runs in time  $O(|A|)$  and therefore should scale to very large graphs.

## Core Decomposition of a Simple Undirected Graph

This section illustrates the use of the core decomposition algorithm on the simple undirected graph  $G$  that is shown in [Figure 1.63](#).



**Figure 1.63** Simple Undirected Graph

The undirected graph  $G$  can be represented using the following nodes data set `NodeSetIn` and links data set `LinkSetIn`:

```
data NodeSetIn;
    input node $ @@;
    datalines;
v1    v2    v3    v4    v5
v6    v7    v8    v9    v10
v11   v12   v13   v14   v15
v16   v17   v18   v19
;

data LinkSetIn;
    input from $ to $ @@;
    datalines;
v1    v2    v5    v6    v6    v7    v7    v8    v10   v11
v2    v3    v3    v4    v2    v4    v8    v9    v9    v10
v8    v18   v10   v12   v13   v14   v13   v15   v13   v16
v13   v17   v14   v15   v14   v16   v14   v17   v15   v16
v15   v17   v16   v17   v18   v13   v18   v17   v18   v16
v12   v14   v12   v15   v12   v16
;
```

The following statements calculate the core decomposition and output the results in the data set `NodeSetOut`:

```
proc optgraph
    data_nodes = NodeSetIn
    data_links = LinkSetIn
    out_nodes  = NodeSetOut;
    core;
run;
```

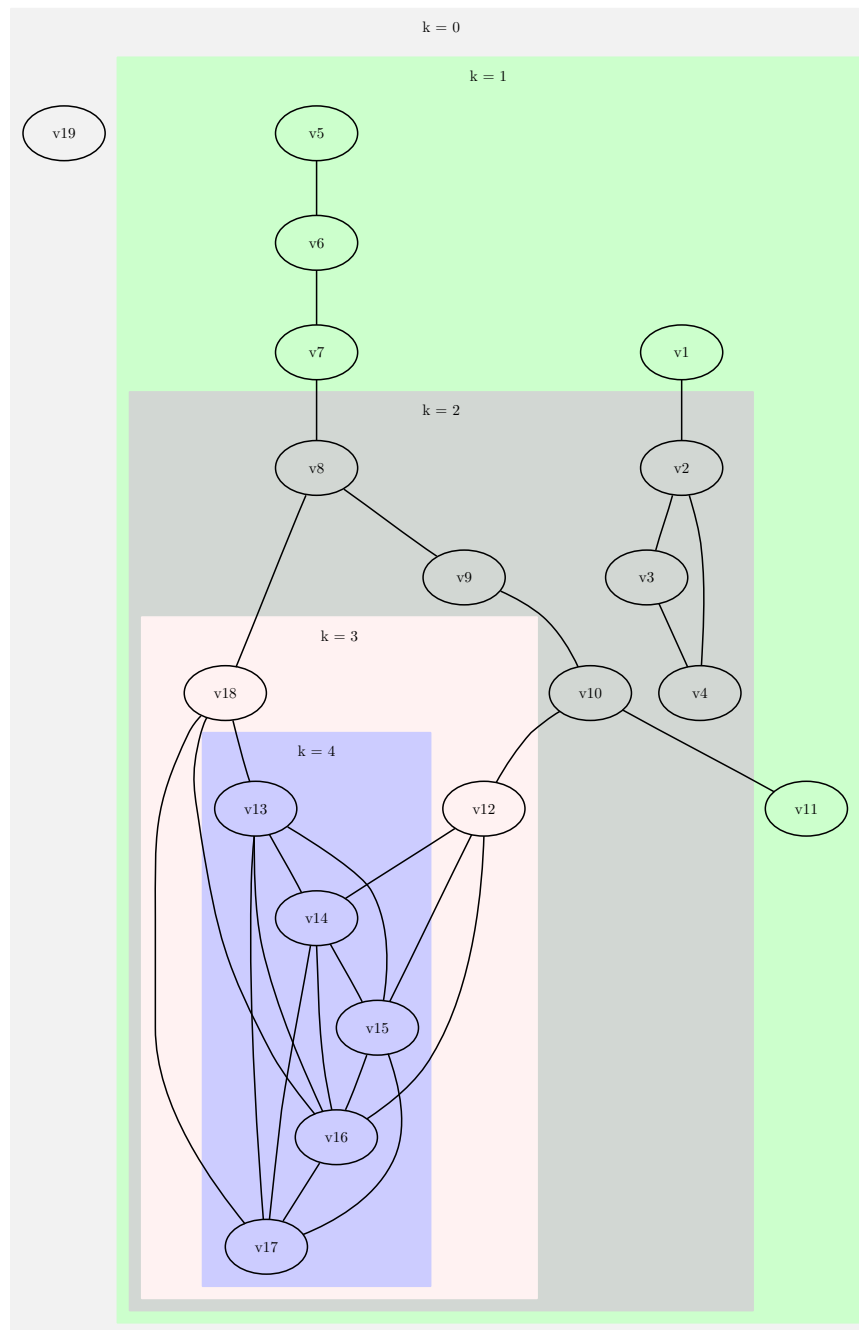
The node data set NodeSetOut contains the core number (variable core\_out) for each node and is shown in Figure 1.64.

**Figure 1.64** Core Decomposition of a Simple Undirected Graph

node	core_out
v19	0
v1	1
v5	1
v6	1
v7	1
v11	1
v2	2
v3	2
v4	2
v8	2
v9	2
v10	2
v12	3
v18	3
v13	4
v14	4
v15	4
v16	4
v17	4

Figure 1.65 shows the graph layered by its core number.

**Figure 1.65** Core Decomposition



## Cycle

A *path* in a graph is a sequence of nodes, each of which has a link to the next node in the sequence. An *elementary cycle* is a path in which the start node and end node are the same and otherwise no node appears more than once in the sequence.

In PROC OPTGRAPH, you can find (or just count) the elementary cycles of an input graph by invoking the CYCLE statement. The options for this statement are described in the section “[CYCLE Statement](#)” on page 30. To find the cycles and report them in an output data set, use the `OUT=` option. To simply count the cycles, do not use the `OUT=` option.

For undirected graphs, each link represents two directed links. For this reason, the following cycles are filtered out: trivial cycles ( $A \rightarrow B \rightarrow A$ ) and duplicate cycles that are found by traversing a cycle in both directions ( $A \rightarrow B \rightarrow C \rightarrow A$  and  $A \rightarrow C \rightarrow B \rightarrow A$ ).

The results of the cycle detection algorithm are written to the output data set that you specify in the `OUT=` option in the CYCLE statement. Each node of each cycle is listed in the `OUT=` data set along with the variable `cycle` to identify the cycle to which it belongs. The variable `order` defines the order (sequence) of the node in the cycle.

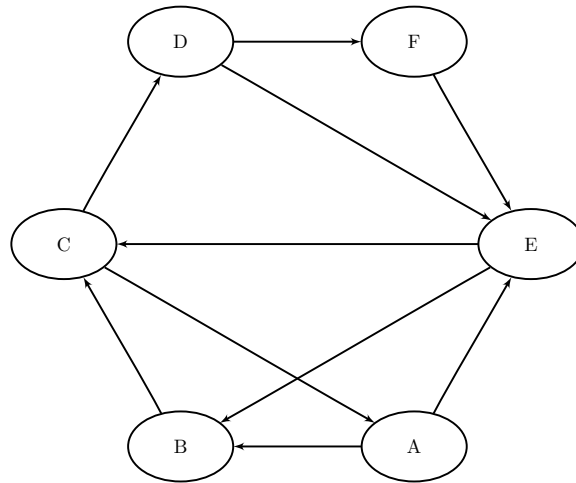
The cycle detection algorithm reports status information in a macro variable called `_OPTGRAPH_CYCLE_`. For more information about this macro variable, see the section “[Macro Variable \\_OPTGRAPH\\_CYCLE\\_](#)” on page 179.

The algorithm that PROC OPTGRAPH uses to compute all cycles is a variant of the algorithm in Johnson (1975). This algorithm runs in time  $O((|N| + |A|)(c + 1))$ , where  $c$  is the number of elementary cycles in the graph. So the algorithm should scale to large graphs that contain few cycles. However, some graphs can have a very large number of cycles, so the algorithm might not scale.

If `MODE=ALL_CYCLES` and there are many cycles, the `OUT=` data set can become very large. It might be beneficial to check the number of cycles before you try to create the `OUT=` data set. When you specify `MODE=FIRST_CYCLE`, the algorithm returns the first cycle that it finds and stops processing. This should run relatively quickly. For large-scale graphs, the `MINLINKWEIGHT=` and `MAXLINKWEIGHT=` options might increase the computation time. For more information about these options, see the section “[CYCLE Statement](#)” on page 30.

## Cycle Detection of a Simple Directed Graph

This section provides a simple example of using the cycle detection algorithm on the simple directed graph  $G$  that is shown in [Figure 1.66](#). Two other examples are “[Example 1.9: Cycle Detection for Kidney Donor Exchange](#)” on page 211, which shows the use of cycle detection for optimizing a kidney donor exchange, and “[Example 1.13: Transitive Closure for Identification of Circular Dependencies in a Bug Tracking System](#)” on page 223, which shows another application of cycle detection.

**Figure 1.66** A Simple Directed Graph  $G$ 

The directed graph  $G$  can be represented by the following links data set, LinkSetIn:

```

data LinkSetIn;
  input from $ to $ @@;
  datalines;
A B A E B C C A C D
D E D F E B E C F E
;

```

The following statements check whether the graph has a cycle:

```

proc optgraph
  graph_direction = directed
  data_links      = LinkSetIn;
  cycle
    mode          = first_cycle;
run;
%put &_OPTGRAPH_;
%put &_OPTGRAPH_CYCLE_;

```

The result is written to the log of the OPTGRAPH procedure, as shown in [Figure 1.67](#).

**Figure 1.67** PROC OPTGRAPH Log: Check the Existence of a Cycle in a Simple Directed Graph

```

NOTE: -----
NOTE: Running OPTGRAPH version 14.1.
NOTE: -----
NOTE: The OPTGRAPH procedure is executing in single-machine mode.
NOTE: -----
NOTE: Data input used 0.00 (cpu: 0.00) seconds.
NOTE: The number of nodes in the input graph is 6.
NOTE: The number of links in the input graph is 10.
NOTE: -----
NOTE: Processing cycle detection.
NOTE: The graph does have a cycle.
NOTE: Processing cycle detection used 0.00 (cpu: 0.00) seconds.
NOTE: -----
NOTE: Data output used 0.00 (cpu: 0.00) seconds.
NOTE: -----
STATUS=OK  CYCLE=OK
STATUS=OK  NUM_CYCLES=1  CPU_TIME=0.00  REAL_TIME=0.00

```

The following statements count the number of cycles in the graph:

```

proc optgraph
  graph_direction = directed
  data_links      = LinkSetIn;
  cycle
    mode          = all_cycles;
run;
%put &_OPTGRAPH_;
%put &_OPTGRAPH_CYCLE_;

```

The result is written to the log of the OPTGRAPH procedure, as shown in [Figure 1.68](#).

**Figure 1.68** PROC OPTGRAPH Log: Count the Number of Cycles in a Simple Directed Graph

```

NOTE: -----
NOTE: Running OPTGRAPH version 14.1.
NOTE: -----
NOTE: The OPTGRAPH procedure is executing in single-machine mode.
NOTE: -----
NOTE: Data input used 0.00 (cpu: 0.00) seconds.
NOTE: The number of nodes in the input graph is 6.
NOTE: The number of links in the input graph is 10.
NOTE: -----
NOTE: Processing cycle detection.
NOTE: The graph has 7 cycles.
NOTE: Processing cycle detection used 0.00 (cpu: 0.00) seconds.
NOTE: -----
NOTE: Data output used 0.00 (cpu: 0.00) seconds.
NOTE: -----
STATUS=OK  CYCLE=OK
STATUS=OK  NUM_CYCLES=7  CPU_TIME=0.00  REAL_TIME=0.00

```

The following statements return the first cycle found in the graph:

```
proc optgraph
  graph_direction = directed
  data_links      = LinkSetIn;
  cycle
    out           = Cycles
    mode          = first_cycle;
run;
```

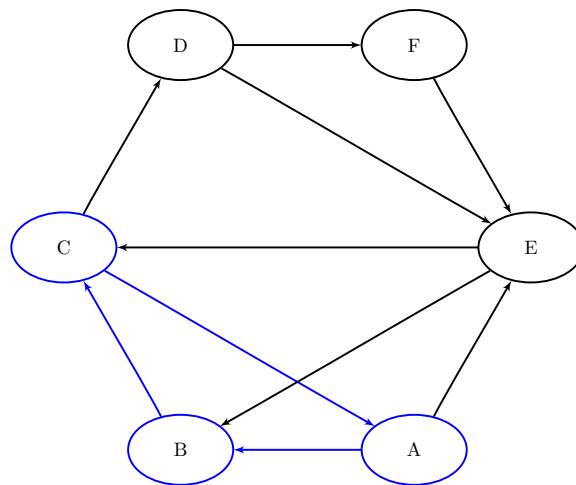
The data set Cycles now contains the first cycle found in the input graph; it is shown in Figure 1.69.

**Figure 1.69** First Cycle Found in a Simple Directed Graph

cycle	order	node
1	1	A
1	2	B
1	3	C
1	4	A

The first cycle that is found in the input graph is shown graphically in Figure 1.70.

**Figure 1.70**  $A \rightarrow B \rightarrow C \rightarrow A$



The following statements return all the cycles in the graph:

```
proc optgraph
  graph_direction = directed
  data_links      = LinkSetIn;
  cycle
    out           = Cycles
    mode          = all_cycles;
run;
```

The data set Cycles now contains all the cycles in the input graph; it is shown in Figure 1.71.

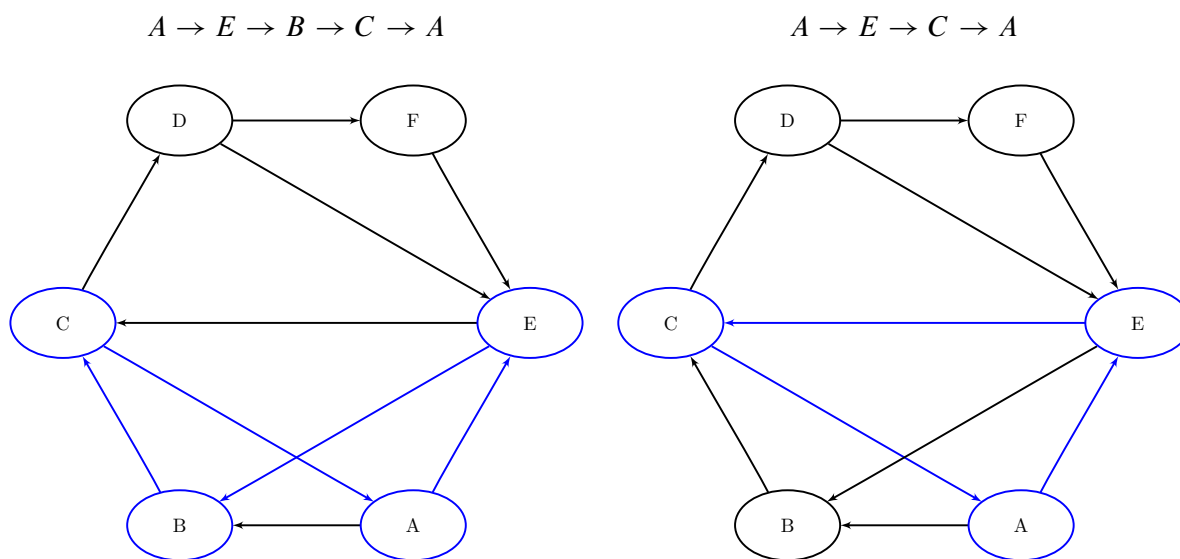


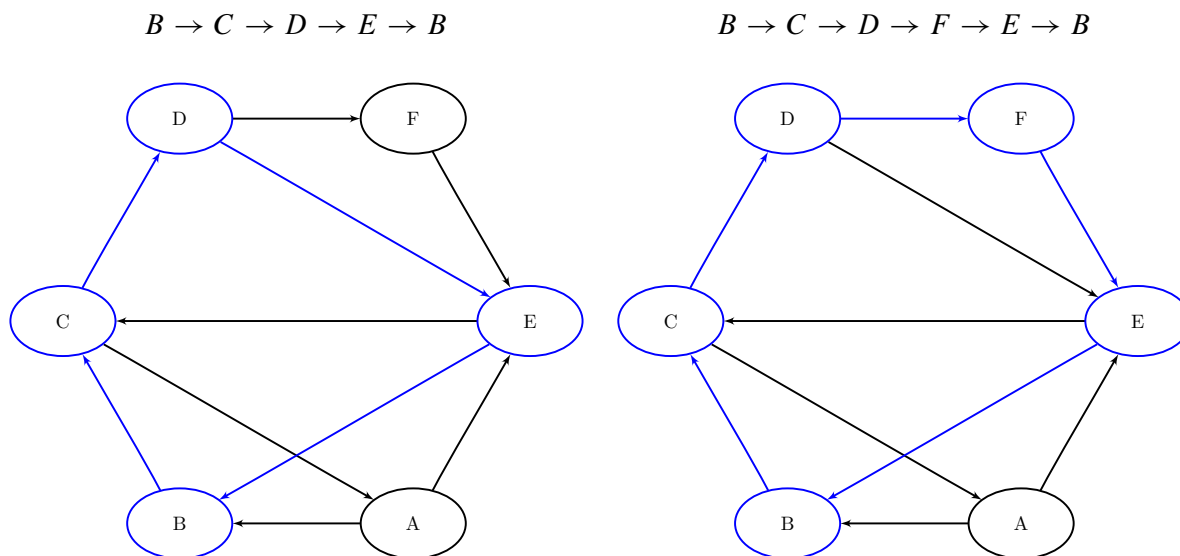
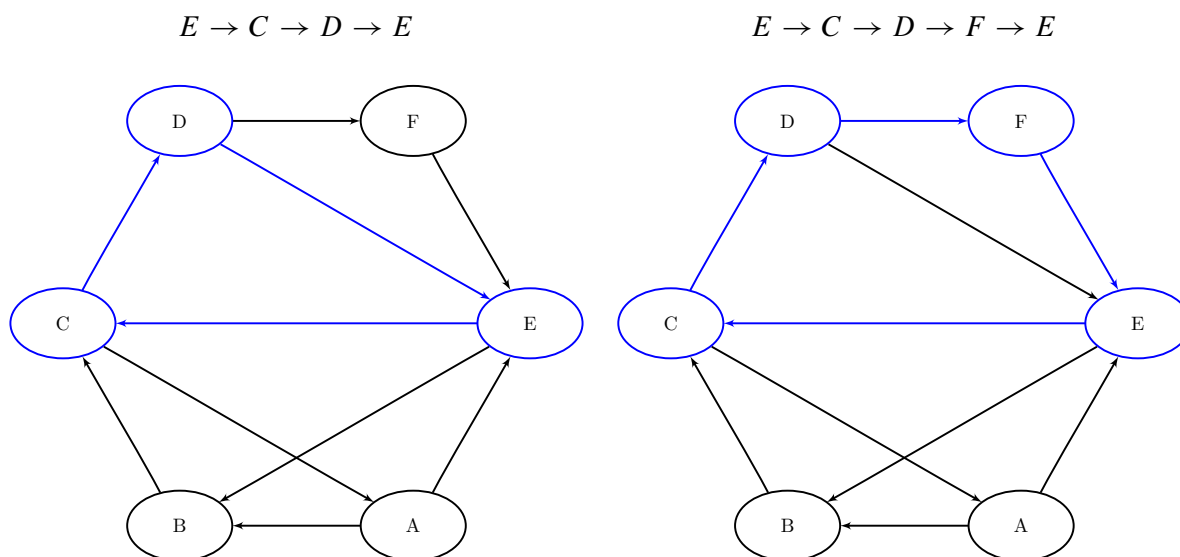
**Figure 1.71** All Cycles in a Simple Directed Graph

cycle	order	node
1	1	A
1	2	B
1	3	C
1	4	A
2	1	A
2	2	E
2	3	B
2	4	C
2	5	A
3	1	A
3	2	E
3	3	C
3	4	A
4	1	B
4	2	C
4	3	D
4	4	E

cycle	order	node
4	5	B
5	1	B
5	2	C
5	3	D
5	4	F
5	5	E
5	6	B
6	1	E
6	2	C
6	3	D
6	4	E
7	1	E
7	2	C
7	3	D
7	4	F
7	5	E

The six additional cycles are shown graphically in [Figure 1.72](#) through [Figure 1.74](#).

**Figure 1.72** Cycles

**Figure 1.73** Cycles**Figure 1.74** Cycles

## Eigenvector Problem

For a given square matrix  $A$ , the *eigenvectors* of the matrix are those nonzero vectors that remain proportional to the original vector after being multiplied by  $A$ . That is, upon multiplication, an eigenvector changes magnitude, but not direction. The corresponding amount that the vector changes in magnitude is called the *eigenvalue*. Mathematically, a nonzero vector  $v$  and scalar  $\lambda$  is an eigenvector/value pair if and only if it satisfies the equation  $Av = \lambda v$ .

In PROC OPTGRAPH, you can calculate eigenvectors of a given matrix by invoking the EIGENVECTOR statement. The options for this statement are described in the section “[EIGENVECTOR Statement](#)” on page 33.

The EIGENVECTOR statement reports status information in a macro variable called `_OPTGRAPH_EIGEN_`. For more information about this macro variable, see the section “[Macro Variable `\_OPTGRAPH\_EIGEN\_`](#)” on page 179.

Although the matrix is typically defined in the input data set specified in the `DATA_MATRIX=` option, it can also be presented in the form of a graph by using the `DATA_LINKS=` option. In this case, the graph is converted to the corresponding adjacency matrix. This conversion enables you to calculate the eigenvectors of very large matrices, since the data format for a graph is very sparse. Because of memory limitations, the matrix format is useful only for relatively small matrices. Because the matrix must be symmetric, the graph input format works only for undirected graphs.

The algorithm that PROC OPTGRAPH uses to solve the eigensystem is a variant of the Jacobi-Davidson algorithm (Sleijpen and van der Vorst 2000). This algorithm uses sparse computations for efficiency and is designed to find a small number of extremal eigenvectors. If you want to find all the eigenvectors and your matrix is relatively small, the best alternative is to use the dense solver in the IML procedure. (See the *SAS/IML User's Guide*.)

### Eigenvector Problem for a Small Matrix with Dense Input

This section shows the calculation of the principal eigenvectors of a small matrix with the following dense input:

```
data MatrixSetIn;
    input coll-col5;
    datalines;
1 0 2 6 1
0 2 3 0 1
2 3 1 0 2
6 0 0 0 0
1 1 2 0 0
;
```

The following statements calculate the two algebraically largest eigenvalues for the matrix defined in the data set MatrixSetIn:

```
proc optgraph
    data_matrix    = MatrixSetIn;
    eigenvector
        eigenvalues = LA
        nEigen      = 2
        out         = EigenMatrixOut;
run;
```

For a matrix with  $n$  columns, and `NEIGEN`= $m$  requested eigenpairs, the algebraically largest eigenvalue is given in the last observation ( $n + 1$ ) of the column `eigen_1`. The corresponding eigenvector is given in the same column in observations 1 through  $n$ . The second largest is given in column `eigen_2`, and so on, up to column `eigen_m`.

In this case, the resulting data set EigenMatrixOut (shown in [Figure 1.75](#)) gives the two largest eigenvector and eigenvalue pairs in columns `eigen_1` and `eigen_2`. The first five observations give the values of the eigenvectors (one for each column of the matrix), and the sixth observation gives the corresponding eigenvalue.

**Figure 1.75** Eigenvector Problem for a Small Matrix with Dense Input

obs	eigen_1	eigen_2
1	0.65778	0.32280
2	0.26459	-0.64125
3	0.40078	-0.49082
4	0.53174	0.40988
5	0.23227	-0.27513
.	7.42209	4.72527

### Eigenvector Problem for a Small Matrix with Sparse Input

This section shows the use of a sparse input format for the eigenvector problem. The following statements define the same matrix that is used in the section “[Eigenvector Problem for a Small Matrix with Dense Input](#)” on page 113, but they represent it sparsely in the form of graph links:

```
data LinkSetIn;
    input from to weight;
    datalines;
0 0 1
0 2 2
0 3 6
0 4 1
1 1 2
1 2 3
1 4 1
2 2 1
2 4 2
;
```

Notice that there are self links  $i \rightarrow i$ . These correspond to the diagonal entries in the matrix that is defined in the data set `MatrixSetIn`. By default, PROC OPTGRAPH ignores self links. Therefore, in the sparse format, you must use the `INCLUDE_SELFLINK` option to match the dense matrix from the section “[Eigenvector Problem for a Small Matrix with Dense Input](#)” on page 113. Now you can calculate the same eigenvectors using sparse input as follows:

```
proc optgraph
    include_selflink
    data_links      = LinkSetIn;
    eigenvector
        eigenvalues = LA
        nEigen      = 2
        out         = EigenLinksOut;
run;
```

The output is shown in [Figure 1.76](#).

**Figure 1.76** Eigenvector Problem for a Small Matrix with Sparse Input

node	eigen_1	eigen_2
0	0.65778	0.32280
2	0.40078	-0.49082
3	0.53174	0.40988
4	0.23227	-0.27513
1	0.26459	-0.64125
.	7.42209	4.72527

## Linear Assignment (Matching)

The *linear assignment problem* (LAP) is a fundamental problem in combinatorial optimization that involves assigning workers to tasks at minimal costs. In graph theoretic terms, LAP is equivalent to finding a minimum-weight matching in a weighted bipartite directed graph. In a *bipartite graph*, the nodes can be divided into two disjoint sets  $S$  (workers) and  $T$  (tasks) such that every link connects a node in  $S$  to a node in  $T$ . That is, the node sets  $S$  and  $T$  are independent. The concept of assigning workers to tasks can be generalized to the assignment of any abstract object from one group to some abstract object from a second group.

The linear assignment problem can be formulated as an integer programming optimization problem. The form of the problem depends on the sizes of the two input sets,  $S$  and  $T$ . Let  $A$  represent the set of possible assignments between sets  $S$  and  $T$ . In the bipartite graph, these assignments are the links. If  $|S| \geq |T|$ , then the following optimization problem is solved:

$$\begin{aligned}
 &\text{minimize} && \sum_{(i,j) \in A} c_{ij} x_{ij} \\
 &\text{subject to} && \sum_{(i,j) \in A} x_{ij} \leq 1 \quad i \in S \\
 &&& \sum_{(i,j) \in A} x_{ij} = 1 \quad j \in T \\
 &&& x_{ij} \in \{0, 1\} \quad (i, j) \in A
 \end{aligned}$$

This model allows for some elements of set  $S$  (workers) to go unassigned (if  $|S| > |T|$ ). However, if  $|S| < |T|$ , then the following optimization problem is solved:

$$\begin{aligned}
 &\text{minimize} && \sum_{(i,j) \in A} c_{ij} x_{ij} \\
 &\text{subject to} && \sum_{(i,j) \in A} x_{ij} = 1 \quad i \in S \\
 &&& \sum_{(i,j) \in A} x_{ij} \leq 1 \quad j \in T \\
 &&& x_{ij} \in \{0, 1\} \quad (i, j) \in A
 \end{aligned}$$

This model allows for some elements of set  $T$  (tasks) to go unassigned.

In PROC OPTGRAPH, you can invoke the linear assignment problem solver by using the `LINEAR_ASSIGNMENT` statement. The options for this statement are described in the section “[LINEAR\\_ASSIGNMENT Statement](#)” on page 34. The algorithm that the PROC OPTGRAPH uses for solving a LAP is based on augmentation of shortest paths (Jonker and Volgenant 1987). This algorithm can be applied to either matrix data input (see the section “[Matrix Input Data](#)” on page 57) or graph data input (see the section “[Graph Input Data](#)” on page 49) as long as the graph is bipartite.

The resulting assignment (or matching) is contained in the output data set that is specified in the `OUT=` option in the `LINEAR_ASSIGNMENT` statement.

The linear assignment problem solver reports status information in a macro variable called `_OPTGRAPH_LAP_`. For more information about this macro variable, see the section “[Macro Variable \\_OPTGRAPH\\_LAP\\_](#)” on page 180.

For a detailed example, see “[Example 1.10: Linear Assignment Problem for Minimizing Swim Times](#)” on page 216.

---

## Minimum-Cost Network Flow

The *minimum-cost network flow* (MCF) problem is a fundamental problem in network analysis that involves sending flow over a network at minimal cost. Let  $G = (N, A)$  be a directed graph. For each link  $(i, j) \in A$ , associate a cost per unit of flow, designated as  $c_{ij}$ . The demand (or supply) at each node  $i \in N$  is designated as  $b_i$ , where  $b_i \geq 0$  denotes a supply node and  $b_i < 0$  denotes a demand node. These values must be within  $[b_i^l, b_i^u]$ . Define decision variables  $x_{ij}$  that denote the amount of flow sent from node  $i$  to node  $j$ . The amount of flow that can be sent across each link is bounded to be within  $[l_{ij}, u_{ij}]$ . The problem can be modeled as a linear programming problem as follows:

$$\begin{aligned} &\text{minimize} && \sum_{(i,j) \in A} c_{ij} x_{ij} \\ &\text{subject to} && b_i^l \leq \sum_{(i,j) \in A} x_{ij} - \sum_{(j,i) \in A} x_{ji} \leq b_i^u \quad i \in N \\ &&& l_{ij} \leq x_{ij} \leq u_{ij} \quad (i, j) \in A \end{aligned}$$

When  $b_i = b_i^l = b_i^u$  for all nodes  $i \in N$ , the problem is called a *pure network flow problem*. For these problems, the sum of the supplies and demands must be equal to 0 to ensure that a feasible solution exists.

In PROC OPTGRAPH, you can invoke the minimum-cost network flow solver by using the `MINCOSTFLOW` statement. The options for this statement are described in the section “[MINCOSTFLOW Statement](#)” on page 35.

The minimum-cost network flow solver reports status information in a macro variable called `_OPTGRAPH_MCF_`. For more information about this macro variable, see the section “[Macro Variable \\_OPTGRAPH\\_MCF\\_](#)” on page 180.

The algorithm that PROC OPTGRAPH uses to solve the MCF problem is a variant of the primal network simplex algorithm (Ahuja, Magnanti, and Orlin 1993). Sometimes the directed graph  $G$  is disconnected. In this case, the problem is first decomposed into its weakly connected components, and then each minimum-cost flow problem is solved separately.

The input for the network is the standard graph input, which is described in the section “[Graph Input Data](#)” on page 49. The links data set, which is specified in the DATA\_LINKS= option in the PROC OPTGRAPH statement, contains the following columns:

- weight, which defines the link cost  $c_{ij}$
- lower, which defines the link lower bound  $l_{ij}$ . The default is 0.
- upper, which defines the link upper bound  $u_{ij}$ . The default is  $\infty$ .

The nodes data set, which is specified in the DATA\_NODES= option in the PROC OPTGRAPH statement, can contain the following columns:

- weight, which defines the node supply lower bound  $b_i^l$ . The default is 0.
- weight2, which defines the node supply upper bound  $b_i^u$ . The default is  $\infty$ .

To define a pure network in which the node supply must be met exactly, use the weight variable only. You do not need to specify all the node supply bounds. For any missing node, the solver uses a lower and upper bound of 0.

To explicitly define an upper bound of  $\infty$ , use the special missing value, (.I). To explicitly define a lower bound of  $-\infty$ , use the special missing value, (.M). Related to infinite bounds, the following scenarios are not supported:

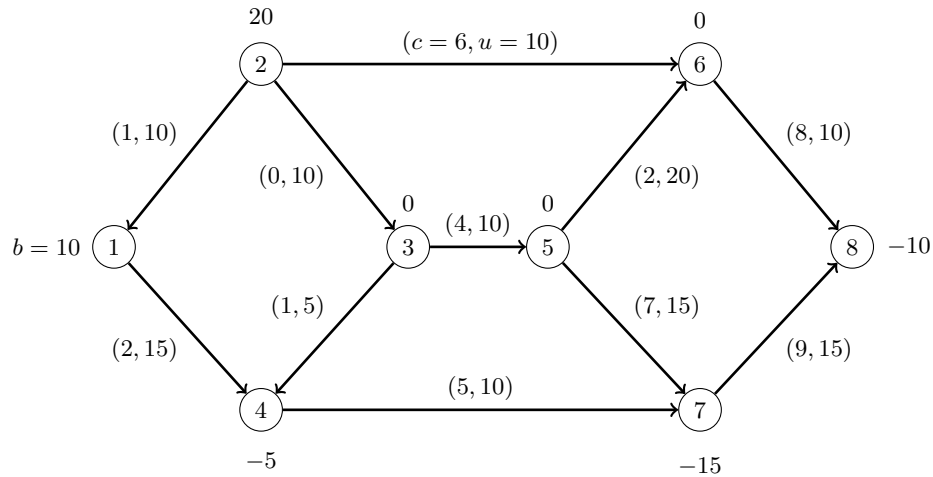
- The flow on a link must be bounded from below ( $l_{ij} = -\infty$  is not allowed).
- Flow balance constraints cannot be *free* ( $b_i^l = -\infty$  and  $b_i^u = \infty$  is not allowed).

The resulting optimal flow through the network is written to the links output data set, which is specified in the OUT\_LINKS= option in the PROC OPTGRAPH statement.

### Minimum-Cost Network Flow for a Simple Directed Graph

This example demonstrates how to use the network simplex algorithm to find a minimum-cost flow in a directed graph. Consider the directed graph in Figure 1.77, which appears in Ahuja, Magnanti, and Orlin (1993).

**Figure 1.77** Minimum-Cost Network Flow Problem: Data



The directed graph  $G$  can be represented by the following links data set LinkSetIn and nodes data set NodeSetIn:

```

data LinkSetIn;
    input from to weight upper;
    datalines;
1 4 2 15
2 1 1 10
2 3 0 10
2 6 6 10
3 4 1 5
3 5 4 10
4 7 5 10
5 6 2 20
5 7 7 15
6 8 8 10
7 8 9 15
;

data NodeSetIn;
    input node weight;
    datalines;
1 10
2 20
4 -5
7 -15
8 -10
;

```



You can use the following call to PROC OPTGRAPH to find a minimum-cost flow:

```
proc optgraph
  loglevel      = moderate
  graph_direction = directed
  data_links    = LinkSetIn
  data_nodes    = NodeSetIn
  out_links     = LinkSetOut;
  mincostflow
    logfreq     = 1;
run;
%put &_OPTGRAPH_;
%put &_OPTGRAPH_MCF_;
```

The progress of the procedure is shown in [Figure 1.78](#).

**Figure 1.78** PROC OPTGRAPH Log for Minimum-Cost Network Flow

```

NOTE: -----
NOTE: -----
NOTE: Running OPTGRAPH version 14.1.
NOTE: -----
NOTE: -----
NOTE: The OPTGRAPH procedure is executing in single-machine mode.
NOTE: -----
NOTE: -----
NOTE: Reading the nodes data set.
NOTE: There were 5 observations read from the data set WORK.NODESETIN.
NOTE: Reading the links data set.
NOTE: There were 11 observations read from the data set WORK.LINKSETIN.
NOTE: Data input used 0.00 (cpu: 0.00) seconds.
NOTE: Building the input graph storage used 0.00 (cpu: 0.00) seconds.
NOTE: The input graph storage is using 0.0 MBs (peak: 0.0 MBs) of memory.
NOTE: The number of nodes in the input graph is 8.
NOTE: The number of links in the input graph is 11.
NOTE: -----
NOTE: -----
NOTE: Processing the minimum-cost network flow problem.
NOTE: The network has 1 connected component.

```

	Primal	Primal	Dual	
Iteration	Objective	Infeasibility	Infeasibility	Time
1	0.000000E+00	2.000000E+01	8.900000E+01	0.00
2	0.000000E+00	2.000000E+01	8.900000E+01	0.00
3	5.000000E+00	1.500000E+01	8.400000E+01	0.00
4	5.000000E+00	1.500000E+01	8.300000E+01	0.00
5	5.000000E+00	1.500000E+01	8.300000E+01	0.00
6	7.500000E+01	1.500000E+01	7.900000E+01	0.00
7	1.300000E+02	1.000000E+01	7.600000E+01	0.00
8	1.300000E+02	1.000000E+01	7.600000E+01	0.00
9	1.300000E+02	1.000000E+01	7.600000E+01	0.00
10	2.700000E+02	0.000000E+00	0.000000E+00	0.00

```

NOTE: The Network Simplex solve time is 0.00 seconds.
NOTE: Objective = 270.
NOTE: Processing the minimum-cost network flow problem used 0.00 (cpu: 0.00) seconds.
NOTE: -----
NOTE: Creating links data set output.
NOTE: Data output used 0.00 (cpu: 0.00) seconds.
NOTE: -----
NOTE: -----
NOTE: The data set WORK.LINKSETOUT has 11 observations and 5 variables.
STATUS=OK   MCF=OPTIMAL
STATUS=OPTIMAL  OBJECTIVE=270  CPU_TIME=0.00  REAL_TIME=0.00

```

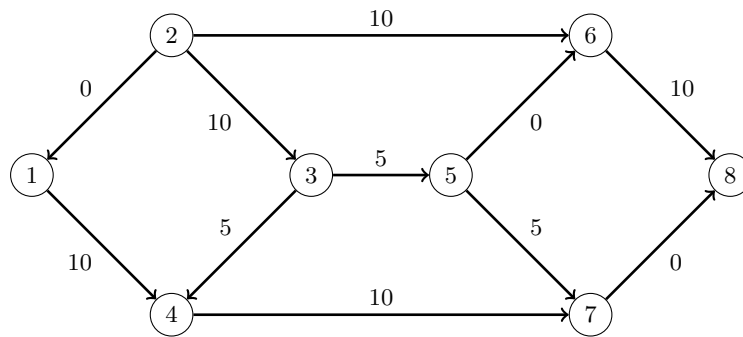
The optimal solution is displayed in [Figure 1.79](#).

**Figure 1.79** Minimum-Cost Network Flow Problem: Optimal Solution

	Obs	from	to	upper	weight	mcf_flow
	1	1	4	15	2	10
	2	2	1	10	1	0
	3	2	3	10	0	10
	4	2	6	10	6	10
	5	3	4	5	1	5
	6	3	5	10	4	5
	7	4	7	10	5	10
	8	5	6	20	2	0
	9	5	7	15	7	5
	10	6	8	10	8	10
	11	7	8	15	9	0

The optimal solution is represented graphically in [Figure 1.80](#).

**Figure 1.80** Minimum-Cost Network Flow Problem: Optimal Solution



## Minimum-Cost Network Flow with Flexible Supply and Demand

Using the same directed graph shown in [Figure 1.77](#), this example demonstrates a network that has a flexible supply and demand. Consider the following adjustments to the node bounds:

- Node 1 has an infinite supply, but it still requires at least 10 units to be sent.
- Node 4 is a throughput node that can now handle an infinite amount of demand.
- Node 8 has a flexible demand. It requires between 6 and 10 units.

You use the special missing values (.I) to represent infinity and (.M) to represent minus infinity. The adjusted node bounds can be represented by the following nodes data set:

```
data NodeSetIn;
  input node weight weight2;
  datalines;
1  10  .I
2  20  20
4  .M  -5
7 -15 -15
8 -10 -6
;
```

You can use the following call to PROC OPTGRAPH to find a minimum-cost flow:

```
proc optgraph
  loglevel          = moderate
  graph_direction   = directed
  data_links        = LinkSetIn
  data_nodes        = NodeSetIn
  out_links         = LinkSetOut;
  mincostflow
    logfreq         = 1;
run;
%put &_OPTGRAPH_;
%put &_OPTGRAPH_MCF_;
```

The progress of the procedure is shown in [Figure 1.81](#).

**Figure 1.81** PROC OPTGRAPH Log for Minimum-Cost Network Flow

```

NOTE: -----
NOTE: -----
NOTE: Running OPTGRAPH version 14.1.
NOTE: -----
NOTE: -----
NOTE: The OPTGRAPH procedure is executing in single-machine mode.
NOTE: -----
NOTE: -----
NOTE: Reading the nodes data set.
NOTE: There were 5 observations read from the data set WORK.NODESETIN.
NOTE: Reading the links data set.
NOTE: There were 11 observations read from the data set WORK.LINKSETIN.
NOTE: Data input used 0.00 (cpu: 0.00) seconds.
NOTE: Building the input graph storage used 0.00 (cpu: 0.00) seconds.
NOTE: The input graph storage is using 0.0 MBs (peak: 0.0 MBs) of memory.
NOTE: The number of nodes in the input graph is 8.
NOTE: The number of links in the input graph is 11.
NOTE: -----
NOTE: -----
NOTE: Processing the minimum-cost network flow problem.
NOTE: The network has 1 connected component.

```

	Primal	Primal	Dual	
Iteration	Objective	Infeasibility	Infeasibility	Time
1	1.000000E+01	2.000000E+01	1.730000E+02	0.00
2	1.000000E+01	2.000000E+01	1.730000E+02	0.00
3	4.500000E+01	2.000000E+01	1.740000E+02	0.00
4	4.500000E+01	2.000000E+01	1.740000E+02	0.00
5	7.500000E+01	1.500000E+01	1.660000E+02	0.00
6	7.500000E+01	1.500000E+01	1.660000E+02	0.00
7	1.590000E+02	9.000000E+00	8.700000E+01	0.00
8	1.590000E+02	9.000000E+00	1.690000E+02	0.00
9	2.140000E+02	4.000000E+00	8.700000E+01	0.00
10	2.140000E+02	4.000000E+00	8.700000E+01	0.00
11	2.260000E+02	0.000000E+00	0.000000E+00	0.00
12	2.260000E+02	0.000000E+00	0.000000E+00	0.00
13	2.260000E+02	0.000000E+00	0.000000E+00	0.00
14	2.260000E+02	0.000000E+00	0.000000E+00	0.00

```

NOTE: The Network Simplex solve time is 0.00 seconds.
NOTE: Objective = 226.
NOTE: Processing the minimum-cost network flow problem used 0.00 (cpu: 0.00) seconds.
NOTE: -----
NOTE: Creating links data set output.
NOTE: Data output used 0.00 (cpu: 0.00) seconds.
NOTE: -----
NOTE: -----
NOTE: The data set WORK.LINKSETOUT has 11 observations and 5 variables.
STATUS=OK  MCF=OPTIMAL
STATUS=OPTIMAL  OBJECTIVE=226  CPU_TIME=0.00  REAL_TIME=0.00

```

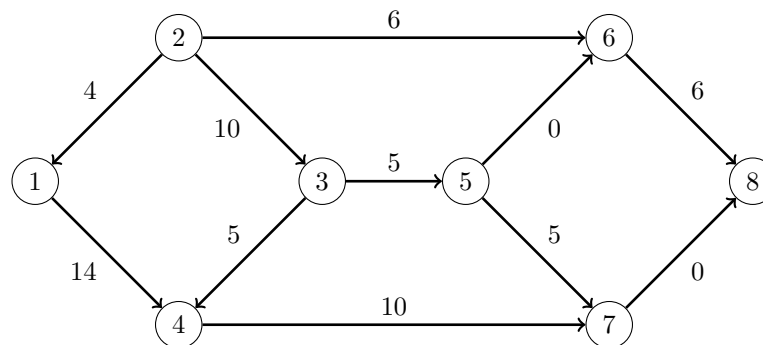
The optimal solution is displayed in Figure 1.82.

**Figure 1.82** Minimum-Cost Network Flow Problem: Optimal Solution

Obs	from	to	upper	weight	mcf_flow
1	1	4	15	2	14
2	2	1	10	1	4
3	2	3	10	0	10
4	2	6	10	6	6
5	3	4	5	1	5
6	3	5	10	4	5
7	4	7	10	5	10
8	5	6	20	2	0
9	5	7	15	7	5
10	6	8	10	8	6
11	7	8	15	9	0

The optimal solution is represented graphically in Figure 1.83.

**Figure 1.83** Minimum-Cost Network Flow Problem: Optimal Solution



## Minimum Cut

A *cut* is a partition of the nodes of a graph into two disjoint subsets. The *cut-set* is the set of links whose *from* and *to* nodes are in different subsets of the partition. A *minimum cut* of an undirected graph is a cut whose cut-set has the smallest link metric, which is measured as follows: For an unweighted graph, the link metric is the number of links in the cut-set. For a weighted graph, the link metric is the sum of the link weights in the cut-set.

In PROC OPTGRAPH, you can invoke the minimum-cut algorithm by using the MINCUT statement. The options for this statement are described in the section “[MINCUT Statement](#)” on page 36. This algorithm can be used only on undirected graphs.

If the value of the MAXNUMCUTS= option is greater than 1, then the algorithm can return more than one set of cuts. The resulting cuts can be described in terms of partitions of the nodes of the graph or the links in the cut-sets. The node partition is specified by the mincut\_*i* variable, for each cut *i*, in the data set that is specified in the OUT\_NODES= option in the PROC OPTGRAPH statement. Each node is assigned the value 0 or 1, which defines the side of the partition to which it belongs. The cut-set is defined in the output data set

that is specified in the OUT= option in the MINCUT statement. This data set lists the links and their weights for each cut.

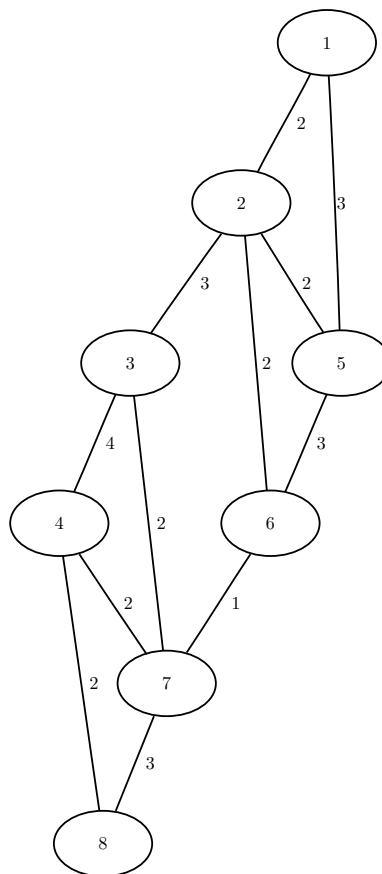
The minimum-cut algorithm reports status information in a macro variable called `_OPTGRAPH_MINCUT_`. For more information about this macro variable, see the section “[Macro Variable \\_OPTGRAPH\\_MINCUT\\_](#)” on page 181.

PROC OPTGRAPH uses the Stoer-Wagner algorithm (Stoer and Wagner 1997) to compute the minimum cuts. This algorithm runs in time  $O(|N||A| + |N|^2 \log |N|)$ .

### Minimum Cut for a Simple Undirected Graph

As a simple example, consider the weighted undirected graph in [Figure 1.84](#).

**Figure 1.84** A Simple Undirected Graph



The links data set can be represented as follows:

```

data LinkSetIn;
  input from to weight @@;
  datalines;
1 2 2 1 5 3 2 3 3 2 5 2 2 6 2
3 4 4 3 7 2 4 7 2 4 8 2 5 6 3
6 7 1 7 8 3
;

```

The following statements calculate minimum cuts in the graph and output the results in the data set MinCut:

```
proc optgraph
  loglevel      = moderate
  out_nodes     = NodeSetOut
  data_links    = LinkSetIn;
  mincut
    out         = MinCut
    maxnumcuts  = 3;
run;
%put &_OPTGRAPH_;
%put &_OPTGRAPH_MINCUT_;
```

The progress of the procedure is shown in Figure 1.85.

**Figure 1.85** PROC OPTGRAPH Log for Minimum Cut

---

```
NOTE: -----
NOTE: -----
NOTE: Running OPTGRAPH version 14.1.
NOTE: -----
NOTE: -----
NOTE: The OPTGRAPH procedure is executing in single-machine mode.
NOTE: -----
NOTE: -----
NOTE: Reading the links data set.
NOTE: There were 12 observations read from the data set WORK.LINKSETIN.
NOTE: Data input used 0.00 (cpu: 0.00) seconds.
NOTE: Building the input graph storage used 0.00 (cpu: 0.00) seconds.
NOTE: The input graph storage is using 0.0 MBs (peak: 0.0 MBs) of memory.
NOTE: The number of nodes in the input graph is 8.
NOTE: The number of links in the input graph is 12.
NOTE: -----
NOTE: -----
NOTE: Processing the minimum-cut problem.
NOTE: The minimum-cut algorithm found 3 cuts.
NOTE: The cut 1 has weight 4.
NOTE: The cut 2 has weight 5.
NOTE: The cut 3 has weight 5.
NOTE: Processing the minimum-cut problem used 0.00 (cpu: 0.00) seconds.
NOTE: -----
NOTE: -----
NOTE: Creating nodes data set output.
NOTE: Creating minimum-cut data set output.
NOTE: Data output used 0.00 (cpu: 0.00) seconds.
NOTE: -----
NOTE: -----
NOTE: The data set WORK.NODESETOUT has 8 observations and 4 variables.
NOTE: The data set WORK.MINCUT has 6 observations and 4 variables.
STATUS=OK  MINCUT=OPTIMAL
STATUS=OPTIMAL  OBJECTIVE=4  CPU_TIME=0.00  REAL_TIME=0.00
```

---



The data set NodeSetOut now contains the partition of the nodes for each cut, shown in [Figure 1.86](#).

**Figure 1.86** Minimum Cut Node Partition

node	mincut_1	mincut_2	mincut_3
1	1	1	1
2	1	1	0
5	1	1	0
3	0	1	0
6	1	1	0
4	0	1	0
7	0	1	0
8	0	0	0

The data set MinCut contains the links in the cut-sets for each cut. This data set is shown in [Figure 1.87](#), which also shows each cut separately.

**Figure 1.87** Minimum Cut Sets

mincut	from	to	weight
1	2	3	3
1	6	7	1
2	4	8	2
2	7	8	3
3	1	2	2
3	1	5	3

**mincut=1**

from	to	weight
2	3	3
6	7	1
<b>mincut</b>		<b>4</b>

**mincut=2**

from	to	weight
4	8	2
7	8	3
<b>mincut</b>		<b>5</b>

**mincut=3**

from	to	weight
1	2	2
1	5	3
<b>mincut</b>		<b>5</b>
		<b>14</b>

## Minimum Spanning Tree

A *spanning tree* of a connected undirected graph is a subgraph that is a tree that connects all the nodes together. When weights have been assigned to the links, a *minimum spanning tree* (MST) is a spanning tree whose sum of link weights is less than or equal to the sum of link weights of every other spanning tree. More generally, any undirected graph (not necessarily connected) has a *minimum spanning forest*, which is a union of minimum spanning trees of its connected components.

In PROC OPTGRAPH, you can invoke the minimum spanning tree algorithm by using the MINSPANTREE statement. The options for this statement are described in the section “[MINSPANTREE Statement](#)” on page 37. This algorithm can be used only on undirected graphs.

The resulting minimum spanning tree is contained in the output data set that is specified in the OUT= option in the MINSPANTREE statement.

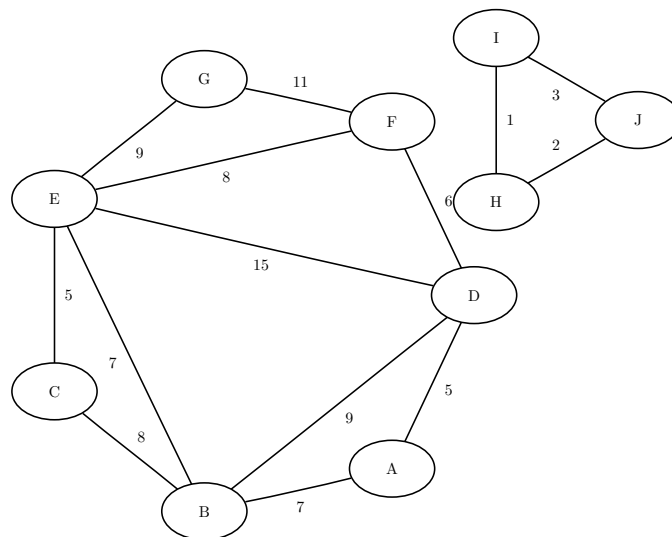
The minimum spanning tree algorithm reports status information in a macro variable called \_OPTGRAPH\_MST\_. For more information about this macro variable, see the section “[Macro Variable \\_OPTGRAPH\\_MST\\_](#)” on page 181.

PROC OPTGRAPH uses Kruskal’s algorithm (Kruskal 1956) to compute the minimum spanning tree. This algorithm runs in time  $O(|A| \log |N|)$  and therefore should scale to very large graphs.

### Minimum Spanning Tree for a Simple Undirected Graph

As a simple example, consider the weighted undirected graph in [Figure 1.88](#).

**Figure 1.88** A Simple Undirected Graph



The links data set can be represented as follows:

```

data LinkSetIn;
  input from $ to $ weight @@;
  datalines;
A B 7 A D 5 B C 8 B D 9 B E 7

```

```

C E 5 D E 15 D F 6 E F 8 E G 9
F G 11 H I 1 I J 3 H J 2
;

```

The following statements calculate a minimum spanning forest and output the results in the data set MinSpanForest:

```

proc optgraph
  data_links = LinkSetIn;
  minspantree
    out      = MinSpanForest;
run;

```

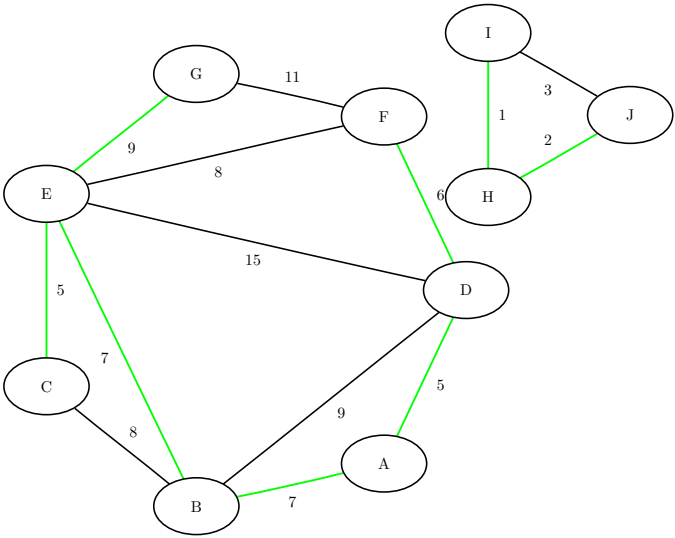
The data set MinSpanForest now contains the links that belong to a minimum spanning forest, which is shown in [Figure 1.89](#).

**Figure 1.89** Minimum Spanning Forest

from to weight		
H	I	1
H	J	2
C	E	5
A	D	5
D	F	6
A	B	7
B	E	7
E	G	9
		<b>42</b>

The minimal cost links are shown in green in [Figure 1.90](#).

**Figure 1.90** Minimum Spanning Forest



For a more detailed example, see “[Example 1.12: Minimum Spanning Tree for Computer Network Topology](#)” on page 222.

## Reach (Ego) Network

The *reach network* of a graph  $G = (N, A)$  is a graph  $G_L^R = (N_L^R, A_L^R)$  that is defined as the induced subgraph over the set of nodes  $N_L^R$  that are reachable in  $L$  steps (or hops) from a set  $S$  of nodes, called the *source nodes*. Reach networks are often referred to as *ego networks* in the context of social networks, since they focus around the neighbors of one (or more) particular individuals.

In PROC OPTGRAPH, reach networks can be calculated by using the REACH statement. The options for this statement are described in the section “[REACH Statement](#)” on page 39.

The REACH statement reports status information in a macro variable called `_OPTGRAPH_REACH_`. For more information about this macro variable, see the section “[Macro Variable \\_OPTGRAPH\\_REACH\\_](#)” on page 181.

In most cases, the set of source nodes from which to calculate reach are defined in a node subset data set, as described in the section “[Node Subset Input Data](#)” on page 54. The node subset data set can be used to define several sets of sources nodes. Each source node set is used to calculate the reach networks. The reach network identifier is given in the node subset data set’s `reach` column. When you use the `EACH_SOURCE` option, every node in the original graph’s node set  $N$  is used to find a reach network from each node separately.

## Output Data Sets

Depending on the options selected, the reach network algorithm produces output data sets as described in the following sections.

### ***OUT\_NODES= Data Set***

This data set describes the nodes in each reach network that are found from each set of source nodes. The data set contains the following columns:

- `node`: node label for each node in each reach network
- `reach`: reach network identifier (which defines the set of source nodes that was used)

### ***OUT\_LINKS= Data Set***

This data set describes the links in each reach network that are found from each set of source nodes. Output of the reach network links can sometimes be more costly computationally, relative to calculating only the nodes or counts in the reach networks. This option does not work when you use the `BY_CLUSTER` option. The data set contains the following columns:

- `from`: the *from* node label for each link in each reach network
- `to`: the *to* node label for each link in each reach network
- `reach`: reach network identifier (which defines the set of source nodes that was used)

**OUT\_COUNTS= Data Set**

This data set describes the number of nodes in each reach network for each set of sources nodes. The data set contains the following columns:

- **node**: node label for each node in the source node sets
- **reach**: reach network identifier (which defines the set of source nodes that was used)
- **count**: the number of nodes reachable using outgoing links from the source nodes
- **count\_not**: the number of nodes not reachable using outgoing links from the source nodes

If the graph is directed and you use the **DIGRAPH** option, then the OUT\_COUNTS= data set contains the following additional columns:

- **count\_in**: the number of nodes reachable using incoming links from the source node
- **count\_out**: the number of nodes reachable using outgoing links from the source node (equivalent to count)
- **count\_in\_or\_out**: the number of nodes reachable using incoming or outgoing links (but not both) from the source node
- **count\_in\_and\_out**: the number of nodes reachable using both incoming and outgoing links from the source node

If node weights are present, the OUT\_COUNTS= data set contains the following additional columns:

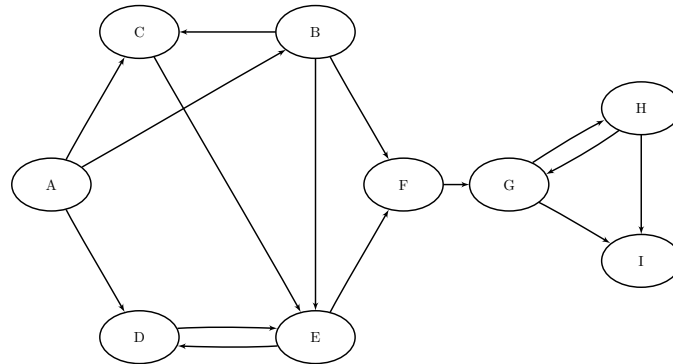
- **count\_wt**: the sum of the weights of the nodes reachable using outgoing links from the source node
- **count\_not\_wt**: the sum of the weights of the nodes not reachable from the source node
- **count\_in\_wt**: the sum of the weights of the nodes reachable using incoming links from the source node
- **count\_out\_wt**: the sum of the weights of the nodes reachable using outgoing links from the source node
- **count\_in\_or\_out\_wt**: the sum of the weights of the nodes reachable using incoming or outgoing links (but not both) from the source node
- **count\_in\_and\_out\_wt**: the sum of the weights of the nodes reachable using both incoming and outgoing links from the source node

When you want to calculate hop limits of 1 and 2 on the same graph, you can use the OUT\_COUNTS1= and OUT\_COUNTS2= options to do this in one call. This option works only when the EACH\_SOURCE and BY\_CLUSTER options are specified.

## Reach Network of a Simple Directed Graph

This section illustrates the use of the reach networks algorithm on the simple directed graph  $G$  that is shown in Figure 1.91.

**Figure 1.91** Simple Directed Graph  $G$



The directed graph  $G$  can be represented using the following links data set LinkSetIn:

```

data LinkSetIn;
  input from $ to $ @@;
  datalines;
A B  A C  A D  B C  B E
B F  C E  D E  E D  E F
F G  G H  G I  H G  H I
;

```

Consider two sets of source nodes,  $S_1 = \{A, G\}$  and  $S_2 = \{B\}$ . These can be defined separately in two node subset data sets as follows:

```

data NodeSubSetIn1;
  input node $ reach;
  datalines;
A 1
G 1
;

data NodeSubSetIn2;
  input node $ reach;
  datalines;
B 1
;

```

For the first set of source nodes, you can use the following statements to calculate the reach network with a hop limit of 1:

```

proc optgraph
  graph_direction = directed
  data_links      = LinkSetIn
  data_nodes_sub  = NodeSubSetIn1;
  reach
    out_nodes     = ReachNodes1
    out_links     = ReachLinks1
    out_counts    = ReachCounts1
    maxreach      = 1;
run;

```

The data sets ReachNodes1, ReachLinks1, and ReachCounts1 now contain the nodes, links, and counts of the reach network, respectively, that come from  $S_1$ .

**Figure 1.92** Reach Network for  $S_1 = \{A, G\}$  with Hop Limit of 1

#### ReachNodes1

reach	node
1	A
1	B
1	C
1	D
1	G
1	H
1	I

#### ReachLinks1

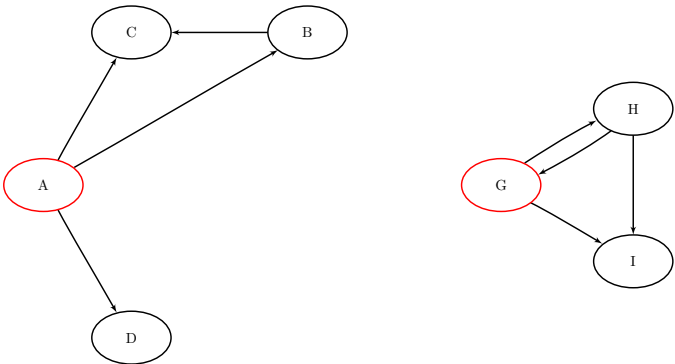
reach	from	to
1	A	B
1	A	C
1	A	D
1	B	C
1	G	H
1	H	G
1	G	I
1	H	I

#### ReachCounts1

reach	node	count	count_not
1	A	7	2
1	G	7	2

The results are displayed graphically in [Figure 1.93](#).

**Figure 1.93** Reach Network for  $S_1 = \{A, G\}$  with Hop Limit of 1



For the second set of source nodes, you can use the following statements to calculate the reach network with a hop limit of 2:

```
proc optgraph
  graph_direction = directed
  data_links      = LinkSetIn
  data_nodes_sub  = NodeSubSetIn2;
  reach
    out_nodes     = ReachNodes2
    out_links     = ReachLinks2
    out_counts    = ReachCounts2
    maxreach      = 2;
run;
```

The data sets ReachNodes2, ReachLinks2, and ReachCounts2 now contain the nodes, links, and counts of the reach network, respectively, that come from  $S_2$ .

**Figure 1.94** Reach Network for  $S_2 = \{B\}$  with Hop Limit of 2

**ReachNodes2**

reach	node
1	B
1	C
1	D
1	E
1	F
1	G



Figure 1.94 continued

ReachLinks2

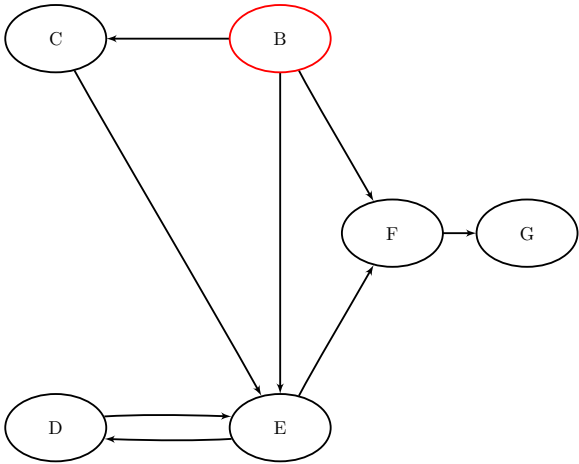
reach from to		
1	B	C
1	B	E
1	B	F
1	C	E
1	D	E
1	E	D
1	E	F
1	F	G

ReachCounts2

reach	node	count	count_not
1	B	6	3

The results are displayed graphically in Figure 1.95.

Figure 1.95 Reach Network for  $S_1 = \{B\}$  with Hop Limit of 2



Processing Multiple Reach Networks in One Pass

You can process a set of reach networks from one graph in one pass using one node subset data set. The MAXREACH= option applies to all of the reach networks requested. If the node subset data set column reach is set to 0 or missing (.), then the node is not processed. If the column reach is set to a value greater than 0, then the node is processed with other nodes by using the same marker.

Consider again the graph shown in Figure 1.91, now with source node sets  $S_1 = \{C\}$  and  $S_2 = \{A, H\}$ . These source node sets can be defined together as follows:

```

data NodeSubSetIn;
  input node $ reach;
  datalines;
A 2
C 1
H 2
;

```

You can use the following statements to process the two one-hop-limit reach networks in one pass:

```

proc optgraph
  graph_direction = directed
  data_links      = LinkSetIn
  data_nodes_sub  = NodeSubSetIn;
  reach
    out_nodes     = ReachNodes
    out_links     = ReachLinks
    out_counts    = ReachCounts
    maxreach      = 1;
run;

```

The data sets ReachNodes, ReachLinks, and ReachCounts now contain the nodes, links, and counts of the reach networks, respectively, that come from  $S_1$  and  $S_2$ .

**Figure 1.96** Reach Networks for  $S_1 = \{C\}$  and  $S_2 = \{A, H\}$  with Hop Limit of 1

#### ReachNodes

reach	node
1	C
1	E
2	A
2	B
2	C
2	D
2	G
2	H
2	I

#### ReachLinks

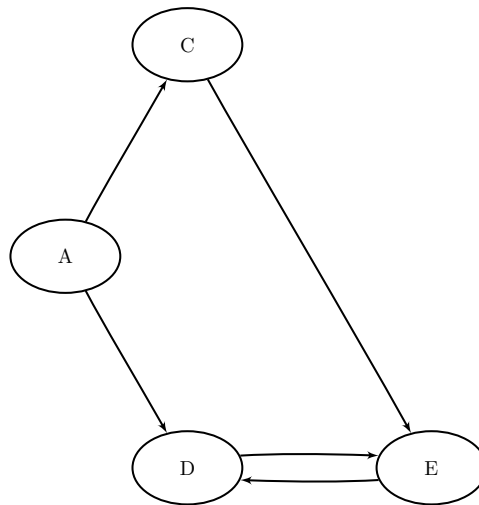
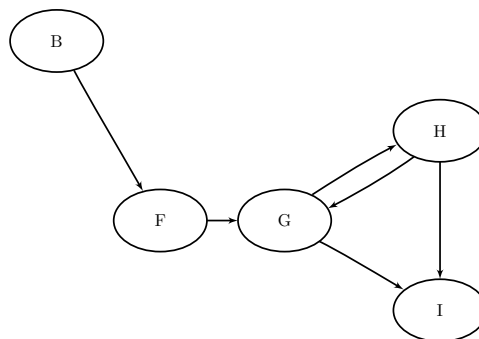
reach	from	to
1	C	E
2	A	B
2	A	C
2	A	D
2	B	C
2	G	H
2	H	G
2	G	I
2	H	I

**Figure 1.96** *continued***ReachCounts**

reach	node	count	count_not
1	C	2	7
2	A	7	2
2	H	7	2

**Processing Reach Networks by Cluster**

Similar to the usage for centrality described in the section “[Processing by Cluster](#)” on page 78, you can use the **BY\_CLUSTER** option in the REACH statement to process a number of induced subgraphs of a graph with only one call to PROC OPTGRAPH. In this section, you want to work on the subgraphs that are induced by node subsets  $N_0 = \{A, C, D, E\}$  and  $N_1 = \{B, F, G, H, I\}$  for the directed graph shown in [Figure 1.91](#). The induced subgraphs are shown graphically in [Figure 1.97](#) and [Figure 1.98](#).

**Figure 1.97** Induced Subgraph for  $N^0 = \{A, C, D, E\}$ **Figure 1.98** Induced Subgraph for  $N^1 = \{B, F, G, H, I\}$ 

Define the subgraphs in the nodes data set by using the cluster variable as follows:

```

data NodeSetIn;
  input node $ cluster @@;
  datalines;
A 0 B 1 C 0 D 0 E 0
F 1 G 1 H 1 I 1
;

```

In the node subset data set, define the source nodes set  $S = \{B, C\}$  by using the reach variable as follows:

```

data NodeSubSetIn;
  input node $ reach;
  datalines;
B 1
C 1
;

```

To process the two-hop-limit reach network for each induced subgraph, you can use the following statements:

```

proc optgraph
  graph_direction = directed
  data_links      = LinkSetIn
  data_nodes      = NodeSetIn
  data_nodes_sub  = NodeSubSetIn;
  performance
    nthreads      = 2;
  reach
    by_cluster
    out_nodes     = ReachNodes
    out_counts    = ReachCounts
    maxreach      = 2;
run;

```

Notice in this example that you can process each subgraph in parallel by using the NTHREADS= option in the **PERFORMANCE** statement.

The data sets ReachNodes and ReachCounts now contain the nodes and counts of the reach networks, respectively, that come from  $S$  for each induced subgraph.

**Figure 1.99** Reach Networks for  $S = \{B, C\}$  with Hop Limit of 2 for Induced Subgraphs

#### ReachNodes

reach	node	cluster
1	B	1
1	C	0
1	D	0
1	E	0
1	F	1
1	G	1

Figure 1.99 *continued***ReachCounts**

reach	node	cluster	count	count_not
1	C	0	3	1
1	B	1	3	2

Notice that since you are operating on the induced subgraphs (not the original graph), node B cannot reach nodes C and E because they are not in its induced subgraph.

**Processing Multiple Reach Networks in One Pass by Cluster**

You can also process several reach networks in one pass while looking over decomposed subgraphs. Consider the same original graph and subgraphs from the section “[Processing Reach Networks by Cluster](#)” on page 137. Now, suppose you want the one-hop-limit reach network where each original node is its own source node subset. Define nine source sets by using the node subset data set as follows:

```
data NodeSubSetIn;
  input node $ reach @@;
  datalines;
A 1 B 2 C 3 D 4 E 5
F 6 G 7 H 8 I 9
;
```

Then, to calculate the reach networks (including the directed graph counts) for each source node set on the induced subgraphs, use the following statements:

```
proc optgraph
  graph_direction = directed
  data_links      = LinkSetIn
  data_nodes      = NodeSetIn
  data_nodes_sub  = NodeSubSetIn;
  performance
    nthreads      = 2;
  reach
    by_cluster
    digraph
    out_nodes      = ReachNodes
    out_counts     = ReachCounts
    maxreach       = 1;
run;
```

Notice that you can do the same thing using the EACH\_SOURCE option. In this case, you do not need the subset data set.

```
proc optgraph
  graph_direction = directed
  data_links      = LinkSetIn
  data_nodes      = NodeSetIn;
  performance
    nthreads      = 2;
  reach
```

```

each_source
by_cluster
digraph
out_nodes    = ReachNodes
out_counts   = ReachCounts
maxreach     = 1;
run;

```

The resulting data sets ReachNodes and ReachCounts are displayed in [Figure 1.100](#).

**Figure 1.100** Reach Networks for Each Source for Induced Subgraphs with a Node Hop Limit of 1

### ReachNodes

reach	node	cluster
1	A	0
1	C	0
1	D	0
2	B	1
2	F	1
3	C	0
3	E	0
4	D	0
4	E	0
5	D	0
5	E	0
6	F	1
6	G	1
7	G	1
7	H	1
7	I	1
8	G	1
8	H	1
8	I	1
9	I	1

### ReachCounts

reach	node	cluster	count	count_not	count_in	count_out	count_in_or_out	count_in_and_out
1	A	0	3	1	0	3	3	0
2	B	1	2	3	0	2	2	0
3	C	0	2	2	1	2	3	0
4	D	0	2	2	2	2	2	1
5	E	0	2	2	2	2	2	1
6	F	1	2	3	1	2	3	0
7	G	1	3	2	2	3	3	1
8	H	1	3	2	1	3	2	1
9	I	1	1	4	2	1	3	0

## Processing Each Source Reach Network for Hop Limits of Both 1 and 2 in One Pass by Cluster

In this section, suppose you want to calculate the one-hop- and two-hop-limit reach counts on the same graph for each source node on a set of induced subgraphs. You can do this in one pass by using the OUT\_COUNTS1= and OUT\_COUNTS2= options, as follows:

```
proc optgraph
  graph_direction = directed
  data_links      = LinkSetIn
  data_nodes      = NodeSetIn;
  performance
    nthreads      = 2;
  reach
    each_source
    by_cluster
    out_counts1    = ReachCounts1
    out_counts2    = ReachCounts2;
run;
```

The resulting data sets ReachCounts1 and ReachCounts1 are displayed in [Figure 1.101](#).

**Figure 1.101** Reach Counts for Each Source Node for Induced Subgraphs with a Hop Limit of 1 and 2

**ReachCounts1**

reach	node	cluster	count	count_not
1	A	0	3	1
3	C	0	2	2
4	D	0	2	2
5	E	0	2	2
2	B	1	2	3
6	F	1	2	3
7	G	1	3	2
8	H	1	3	2
9	I	1	1	4

**ReachCounts2**

reach	node	cluster	count	count_not
1	A	0	4	0
3	C	0	3	1
4	D	0	2	2
5	E	0	2	2
2	B	1	3	2
6	F	1	4	1
7	G	1	3	2
8	H	1	3	2
9	I	1	1	4

For a more detailed example, see “[Example 1.14: Reach Networks for Computation of Market Coverage of a Terrorist Network](#)” on page 226.

## Shortest Path

A *shortest path* between two nodes  $u$  and  $v$  in a graph is a path that starts at  $u$  and ends at  $v$  and has the lowest total link weight. The starting node is called the *source node*, and the ending node is called the *sink node*.

In PROC OPTGRAPH, you can calculate shortest paths by using the SHORTPATH statement. The options for this statement are described in the section “[SHORTPATH Statement](#)” on page 40.

The shortest path algorithm reports status information in a macro variable called \_OPTGRAPH\_SHORTPATH\_. For more information about this macro variable, see the section “[Macro Variable \\_OPTGRAPH\\_SHORTPATH\\_](#)” on page 182.

By default, PROC OPTGRAPH finds shortest paths for all pairs. That is, it finds a shortest path for each possible combination of source and sink nodes. Alternatively, you can use the SOURCE= option to fix a particular source node and find shortest paths from the fixed source node to all possible sink nodes. Conversely, by using the SINK= option, you can fix a sink node and find shortest paths from all possible source nodes to the fixed sink node. By using both options together, you can request one particular shortest path for a specific source-sink pair. In addition, you can use the DATA\_NODES\_SUB= option to define a list of source-sink pairs to process, as described in the section “[Node Subset Input Data](#)” on page 54. The following sections show examples of these options.

Which algorithm PROC OPTGRAPH uses to find shortest paths depends on the data. The algorithm and run-time complexity for each link type are shown in [Table 1.58](#).

**Table 1.58** Algorithms for Shortest Paths

Link Type	Algorithm	Complexity (per Source Node)
Unweighted	Breadth-first search	$O( N  +  A )$
Weighted (nonnegative)	Dijkstra’s algorithm	$O( N  \log  N  +  A )$
Weighted (positive and negative allowed)	Bellman-Ford algorithm	$O( N  A )$

Details for each algorithm can be found in Ahuja, Magnanti, and Orlin (1993).

For weighted graphs, the algorithm uses the weight variable that is defined in the links data set to evaluate a path’s total weight (cost). You can also use the WEIGHT2= option in the SHORTPATH statement to define an auxiliary weight. The auxiliary weight is not used in the algorithm to evaluate a path’s total weight. It is calculated only for the sake of reporting the total auxiliary weight for each shortest path.

## Output Data Sets

The shortest path algorithm produces up to two output data sets. The output data set that you specify in the OUT\_PATHS= option contains the links of a shortest path for each source-sink pair combination. The output data set that you specify in the OUT\_WEIGHTS= option contains the total weight for the shortest path for each source-sink pair combination.

### OUT\_PATHS= Data Set

The OUT\_PATHS= data set contains the links present in the shortest path for each source-sink pair. For large graphs and a large requested number of source-sink pairs, this output data set can be extremely large. For extremely large sets, generating the output can sometimes take longer than computing the shortest paths. For



example, using the US road network data for the state of New York, the data contain a directed graph that has 264,346 nodes. Finding the shortest path for all pairs from only one source node results in 140,969,120 observations, which is a data set of size 11 GB. Finding shortest paths for all pairs from all nodes would produce an enormous output data set.

The OUT\_PATHS= data set contains the following columns:

- source: the source node label of this shortest path
- sink: the sink node label of this shortest path
- order: for this source-sink pair, the order of this link in a shortest path
- from: the *from* node label of this link in a shortest path
- to: the *to* node label of this link in a shortest path
- weight: the weight of this link in a shortest path
- weight2: the auxiliary weight of this link

#### **OUT\_WEIGHTS= Data Set**

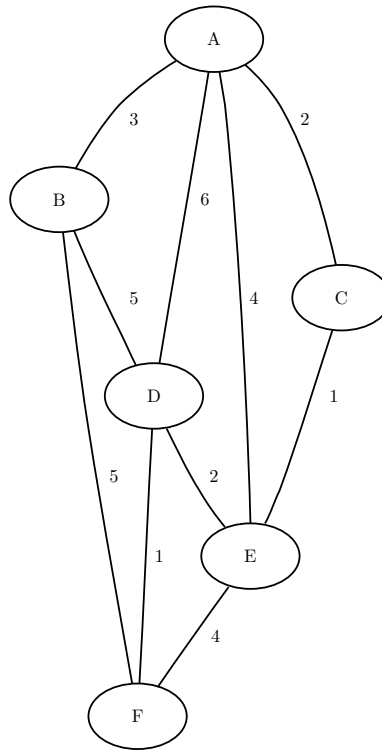
The OUT\_WEIGHTS= data set contains the total weight (and total auxiliary weight) for the shortest path for each source-sink pair.

This data set contains the following columns:

- source: the source node label of this shortest path
- sink: the sink node label of this shortest path
- path\_weight: the total weight of the shortest path for this source-sink pair
- path\_weight2: the total auxiliary weight of the shortest path for this source-sink pair

### **Shortest Paths for All Pairs**

This example illustrates the use of the shortest path algorithm for all source-sink pairs on the simple undirected graph *G* shown in [Figure 1.102](#).

**Figure 1.102** A Simple Undirected Graph  $G$ 

The undirected graph  $G$  can be represented by the following links data set, LinkSetIn:

```

data LinkSetIn;
  input from $ to $ weight @@;
  datalines;
A B 3  A C 2  A D 6  A E 4  B D 5
B F 5  C E 1  D E 2  D F 1  E F 4
;

```

The following statements calculate shortest paths for all source-sink pairs:

```

proc optgraph
  data_links      = LinkSetIn;
  shortpath
    out_weights = ShortPathW
    out_paths   = ShortPathP;
run;

```

The data set ShortPathP contains the shortest paths and is shown in [Figure 1.103](#).

**Figure 1.103** All-Pairs Shortest Paths

source	sink	order	from	to	weight
A	B	1	A	B	3
A	C	1	A	C	2
A	D	1	A	C	2
A	D	2	C	E	1
A	D	3	E	D	2
A	E	1	A	C	2
A	E	2	C	E	1
A	F	1	A	C	2
A	F	2	C	E	1
A	F	3	E	D	2
A	F	4	D	F	1
B	A	1	B	A	3
B	C	1	B	A	3
B	C	2	A	C	2
B	D	1	B	D	5
B	E	1	B	A	3
B	E	2	A	C	2
B	E	3	C	E	1
B	F	1	B	F	5
C	A	1	C	A	2
C	B	1	C	A	2
C	B	2	A	B	3
C	D	1	C	E	1
C	D	2	E	D	2
C	E	1	C	E	1
C	F	1	C	E	1
C	F	2	E	D	2
C	F	3	D	F	1

source	sink	order	from	to	weight
D	A	1	D	E	2
D	A	2	E	C	1
D	A	3	C	A	2
D	B	1	D	B	5
D	C	1	D	E	2
D	C	2	E	C	1
D	E	1	D	E	2
D	F	1	D	F	1
E	A	1	E	C	1
E	A	2	C	A	2
E	B	1	E	C	1
E	B	2	C	A	2
E	B	3	A	B	3
E	C	1	E	C	1
E	D	1	E	D	2
E	F	1	E	D	2
E	F	2	D	F	1
F	A	1	F	D	1
F	A	2	D	E	2
F	A	3	E	C	1
F	A	4	C	A	2
F	B	1	F	B	5
F	C	1	F	D	1
F	C	2	D	E	2
F	C	3	E	C	1
F	D	1	F	D	1
F	E	1	F	D	1
F	E	2	D	E	2

The data set ShortPathW contains the path weight for the shortest paths of each source-sink pair and is shown in Figure 1.104.

**Figure 1.104** All-Pairs Shortest Paths Summary

source	sink	path_weight	source	sink	path_weight
A	B	3	D	A	5
A	C	2	D	B	5
A	D	5	D	C	3
A	E	3	D	E	2
A	F	6	D	F	1
B	A	3	E	A	3
B	C	5	E	B	6
B	D	5	E	C	1
B	E	6	E	D	2
B	F	5	E	F	3
C	A	2	F	A	6
C	B	5	F	B	5
C	D	3	F	C	4
C	E	1	F	D	1
C	F	4	F	E	3

When you are interested only in the source-sink pair that has the longest shortest path, you can use the `PATHS=` option. This option affects only the output processing; it does not affect the computation. All the designated source-sink shortest paths are calculated, but only the longest ones are written to the output data set.

The following statements display only the longest shortest paths:

```
proc optgraph
  data_links = LinkSetIn;
  shortpath
    paths = longest
    out_paths = ShortPathLong;
run;
```

The data set `ShortPathLong` now contains the longest shortest paths and is shown in [Figure 1.105](#).

**Figure 1.105** Longest Shortest Paths

source	sink	order	from	to	weight
A	F	1	A	C	2
A	F	2	C	E	1
A	F	3	E	D	2
A	F	4	D	F	1
B	E	1	B	A	3
B	E	2	A	C	2
B	E	3	C	E	1
E	B	1	E	C	1
E	B	2	C	A	2
E	B	3	A	B	3
F	A	1	F	D	1
F	A	2	D	E	2
F	A	3	E	C	1
F	A	4	C	A	2

## Shortest Paths for a Subset of Source-Sink Pairs

This section illustrates the use of a node subset data set, the `DATA_NODES_SUB=` option, and the shortest path algorithm to calculate shortest paths for a subset of source-sink pairs. The data set variables `source` and `sink` are used as indicators to specify which pairs to process. The marked source nodes define a set  $S$ , and the marked sink nodes define a set  $T$ . PROC OPTGRAPH then calculates all the source-sink pairs in the crossproduct of these two sets.

For example, the following DATA step tells PROC OPTGRAPH to calculate the pairs in  $S \times T = \{A, C\} \times \{B, F\}$ :

```
data NodeSubSetIn;
    input node $ source sink;
    datalines;
A 1 0
C 1 0
B 0 1
F 0 1
;
```

The following statements calculate a shortest path for the four combinations of source-sink pairs:

```
proc optgraph
    data_nodes_sub = NodeSubSetIn
    data_links      = LinkSetIn;
    shortpath
        out_paths   = ShortPath;
run;
```

The data set ShortPath contains the shortest paths and is shown in Figure 1.106.

**Figure 1.106** Shortest Paths for a Subset of Source-Sink Pairs

source	sink	order	from	to	weight
A	B	1	A	B	3
A	F	1	A	C	2
A	F	2	C	E	1
A	F	3	E	D	2
A	F	4	D	F	1
C	B	1	C	A	2
C	B	2	A	B	3
C	F	1	C	E	1
C	F	2	E	D	2
C	F	3	D	F	1

## Shortest Paths for a Subset of Source or Sink Pairs

This section illustrates the use of the shortest path algorithm to calculate shortest paths between a subset of source (or sink) nodes and all other sink (or source) nodes.

In this case, you designate the subset of source (or sink) nodes in the node subset data set by specifying the `source` (or `sink`). By specifying only one of the variables, you indicate that you want PROC OPTGRAPH to calculate all pairs from a subset of source nodes (or to calculate all pairs to a subset of sink nodes).

For example, the following DATA step designates nodes *B* and *E* as source nodes:

```
data NodeSubSetIn;
  input node $ source;
  datalines;
B 1
E 1
;
```

You can use the same PROC OPTGRAPH call as is used in the section “Shortest Paths for a Subset of Source-Sink Pairs” on page 147 to calculate all the shortest paths from nodes *B* and *E*. The data set ShortPath contains the shortest paths and is shown in Figure 1.107.

**Figure 1.107** Shortest Paths for a Subset of Source Pairs

source	sink	order	from	to	weight
B	A	1	B	A	3
B	C	1	B	A	3
B	C	2	A	C	2
B	D	1	B	D	5
B	E	1	B	A	3
B	E	2	A	C	2
B	E	3	C	E	1
B	F	1	B	F	5
E	A	1	E	C	1
E	A	2	C	A	2
E	B	1	E	C	1
E	B	2	C	A	2
E	B	3	A	B	3
E	C	1	E	C	1
E	D	1	E	D	2
E	F	1	E	D	2
E	F	2	D	F	1

Conversely, the following DATA step designates nodes *B* and *E* as sink nodes:

```
data NodeSubSetIn;
  input node $ sink;
  datalines;
B 1
E 1
;
```

You can use the same PROC OPTGRAPH call again to calculate all the shortest paths to nodes *B* and *E*. The data set ShortPath contains the shortest paths and is shown in Figure 1.108.

**Figure 1.108** Shortest Paths for a Subset of Sink Pairs

source	sink	order	from	to	weight
A	B	1	A	B	3
A	E	1	A	C	2
A	E	2	C	E	1
B	E	1	B	A	3
B	E	2	A	C	2
B	E	3	C	E	1
C	B	1	C	A	2
C	B	2	A	B	3
C	E	1	C	E	1
D	B	1	D	B	5
D	E	1	D	E	2
E	B	1	E	C	1
E	B	2	C	A	2
E	B	3	A	B	3
F	B	1	F	B	5
F	E	1	F	D	1
F	E	2	D	E	2

### Shortest Paths for One Source-Sink Pair

This section illustrates the use of the shortest path algorithm to calculate shortest paths between one source-sink pair by using the SOURCE= and SINK= options.

The following statements calculate a shortest path between node *C* and node *F*:

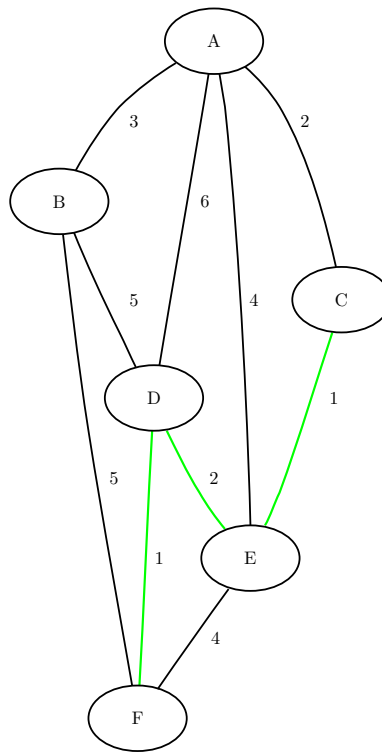
```
proc optgraph
  data_links = LinkSetIn;
  shortpath
    source = C
    sink   = F
    out_paths = ShortPath;
run;
```

The data set ShortPath contains this shortest path and is shown in [Figure 1.109](#).

**Figure 1.109** Shortest Paths for One Source-Sink Pair

source	sink	order	from	to	weight
C	F	1	C	E	1
C	F	2	E	D	2
C	F	3	D	F	1

The shortest path is shown graphically in [Figure 1.110](#).

**Figure 1.110** Shortest Path between Nodes *C* and *F*

### Shortest Paths with Auxiliary Weight Calculation

This section illustrates the use of the shortest path algorithm with auxiliary weights to calculate the shortest paths between all source-sink pairs.

Consider a links data set in which the auxiliary weight is a counter for each link:

```

data LinkSetIn;
  input from $ to $ weight count @@;
  datalines;
A B 3 1  A C 2 1  A D 6 1  A E 4 1  B D 5 1
B F 5 1  C E 1 1  D E 2 1  D F 1 1  E F 4 1
;
  
```



The following statements calculate shortest paths for all source-sink pairs:

```
proc optgraph
  data_links      = LinkSetIn;
  shortpath
    weight2       = count
    out_weights   = ShortPathW;
run;
```

The data set ShortPathW contains the total path weight for shortest paths in each source-sink pair and is shown in Figure 1.111. Because the variable count in LinkSetIn has a value of 1 for all links, the value in the output data set variable path\_weights2 contains the number of links in each shortest path.

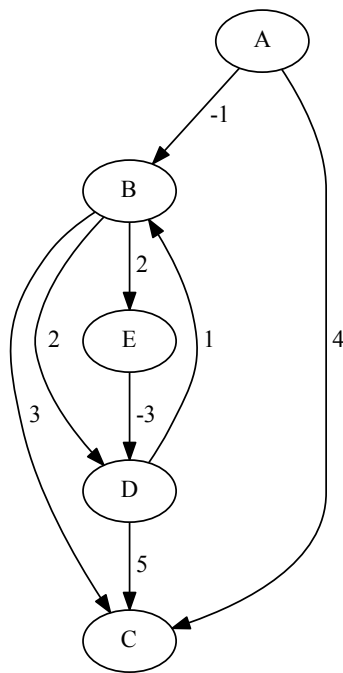
**Figure 1.111** Shortest Paths Including Auxiliary Weights in Calculation

source	sink	path_weight	path_weight2	source	sink	path_weight	path_weight2
A	B	3	1	D	A	5	3
A	C	2	1	D	B	5	1
A	D	5	3	D	C	3	2
A	E	3	2	D	E	2	1
A	F	6	4	D	F	1	1
B	A	3	1	E	A	3	2
B	C	5	2	E	B	6	3
B	D	5	1	E	C	1	1
B	E	6	3	E	D	2	1
B	F	5	1	E	F	3	2
C	A	2	1	F	A	6	4
C	B	5	2	F	B	5	1
C	D	3	2	F	C	4	3
C	E	1	1	F	D	1	1
C	F	4	3	F	E	3	2

The section “Road Network Shortest Path” on page 4 shows an example of using the shortest path algorithm to minimize travel to and from work based on traffic conditions.

### Shortest Paths with Negative Link Weights

This section illustrates the use of the shortest path algorithm on a simple directed graph  $G$  with negative link weights, shown in Figure 1.112.

**Figure 1.112** A Simple Directed Graph  $G$  with Negative Link Weights

You can represent the directed graph  $G$  by using the following links data set LinkSetIn:

```

data LinkSetIn;
  input from $ to $ weight @@;
  datalines;
A B -1  A C  4  B C  3  B D  2  B E  2
D B  1  D C  5  E D -3
;

```

The following statements calculate the shortest paths between source node  $E$  and sink node  $B$ :

```

proc optgraph
  direction      = directed
  data_links     = LinkSetIn;
  shortpath
    source       = E
    sink         = B
    out_paths    = ShortPathP;
run;

```

The data set ShortPathP contains the shortest path from node  $E$  to node  $B$  and is shown in Figure 1.113.

**Figure 1.113** Shortest Paths with Negative Link Weights

source	sink	order	from	to	weight
E	B	1	E	D	-3
E	B	2	D	B	1

Now, consider the following adjustment to the weight of link  $(B, E)$ :

```
data LinkSetIn;
  set LinkSetIn;
  if (from="B" and to="E") then
    weight=1;
run;
```

In this case, there is a negative weight cycle  $(E \rightarrow D \rightarrow B \rightarrow E)$ . The Bellman-Ford algorithm catches this and produces an error, as shown in Figure 1.114.

**Figure 1.114** PROC OPTGRAPH Log: Negative Weight Cycle

---

```
NOTE: -----
NOTE: Running OPTGRAPH version 14.1.
NOTE: -----
NOTE: The OPTGRAPH procedure is executing in single-machine mode.
NOTE: -----
NOTE: Data input used 0.00 (cpu: 0.00) seconds.
NOTE: The number of nodes in the input graph is 5.
NOTE: The number of links in the input graph is 8.
NOTE: -----
NOTE: Processing the shortest paths problem.
ERROR: The graph contains a negative weight cycle.
NOTE: Processing the shortest paths problem used 0.00 (cpu: 0.00) seconds.
NOTE: -----
NOTE: Data output used 0.00 (cpu: 0.00) seconds.
NOTE: -----
NOTE: The data set WORK.SHORTPATHP has 0 observations and 6 variables.
STATUS=OK  SHORTPATH=ERROR
```

---

## Summary

In PROC OPTGRAPH, you can calculate various summary statistics for a graph and its nodes by using the SUMMARY statement. The options for this statement are described in the section “[SUMMARY Statement](#)” on page 42.

The SUMMARY statement reports status information in a macro variable called `_OPTGRAPH_SUMMARY_`. For more information about this macro variable, see the section “[Macro Variable \\_OPTGRAPH\\_SUMMARY\\_](#)” on page 182.

## Output Data Sets

The summary statistics that are produced are broken into two categories: statistics on the entire graph and statistics on the nodes and links of the graph. The latter statistics are appended to the output nodes and links data sets that you specify in the `OUT_NODES=` and `OUT_LINKS=` option in the PROC OPTGRAPH statement. The former statistics are contained in the data set that you specify in the `OUT=` option in the SUMMARY statement.

Let  $\delta(u)$  represent the list of nodes that are connected to node  $u$  in an undirected graph. In a directed graph,  $\delta^{\text{out}}(u)$  represents the list of nodes that are connected *from* node  $u$  (out-links), and  $\delta^{\text{in}}(u)$  represents the list of nodes that are connected *to* node  $u$  (in-links).

### OUT= Data Set

By default, the summary output data set that you specify in the OUT= option in the SUMMARY statement contains the following columns:

- nodes: the number of nodes in the graph ( $|N|$ )
- links: the number of links in the graph ( $|A|$ )
- avg\_links\_per\_node: the average number of links per node
- density: the number of links in the graph divided by the number of links in a complete graph  $\left(\frac{|A|}{|K(N)|}\right)$
- self\_links\_ignored: the number of self-links ignored
- dup\_links\_ignored: the number of duplicate links ignored
- leaf\_nodes: the number of leaf nodes
  - Undirected graph:  $u \in N$  such that  $\delta(u) = 1$
  - Directed graph:  $u \in N$  such that  $\delta^{\text{out}}(u) = 0$  and  $\delta^{\text{out}}(u) + \delta^{\text{in}}(u) > 0$
- singleton\_nodes: the number of singleton nodes
  - Undirected graph:  $u \in N$  such that  $\delta(u) = 0$
  - Directed graph:  $u \in N$  such that  $\delta^{\text{out}}(u) + \delta^{\text{in}}(u) = 0$

You can produce statistics about the connectedness of the graph by using the CONCOMP and BICONCOMP options. For more information about connected components and biconnected components, see the sections “[Connected Components](#)” on page 97 and “[Biconnected Components and Articulation Points](#)” on page 61, respectively. If you use the CONCOMP and BICONCOMP options, the following columns also appear in the summary output data set for undirected graphs:

- concomp: the number of connected components in the graph
- biconcomp: the number of biconnected components in the graph
- artpoints: the number of articulation points in the graph
- isolated\_pairs: the number of isolated pairs of nodes (a connected component of size 2)
- isolated\_stars: the number of isolated stars, (a connected component,  $C$ , of size greater than 2 with):
  - one node  $i$  with  $\delta(i) = |C| - 1$  and all other nodes  $u \in C \setminus \{i\}$  with  $\delta(u) = 1$

The following columns appear for directed graphs:

- concomp: the number of strongly connected components in the graph

- `isolated_pairs`: the number of isolated pairs of nodes (a weakly connected component of size 2)
- `isolated_stars_out`: the number of isolated outward stars (a weakly connected component,  $C$ , of size greater than 2 with):
  - one node  $i$  with  $\delta^{\text{out}}(i) = |C| - 1$  and all other nodes  $u \in C \setminus \{i\}$  with  $\delta^{\text{in}}(u) = 1$
- `isolated_stars_in`: the number of isolated inward stars (a weakly connected component,  $C$ , of size greater than 2 with):
  - one node  $i$  with  $\delta^{\text{in}}(i) = |C| - 1$  and all other nodes  $u \in C \setminus \{i\}$  with  $\delta^{\text{out}}(u) = 1$

You can produce statistics about the shortest paths in the graph by using the `SHORTPATH=` option. The *diameter* of a graph is the longest shortest path distance of all possible source-sink pairs in the graph. Calculating the diameter of a graph is computationally expensive, because it involves calculating shortest paths for all pairs. For undirected graphs, an approximate method is available based on Boitmanis et al. (2006). You can invoke the algorithm by using the `DIAMETER_APPROX=` option. The exact method runs in time  $O(|N| \times (|N| \log |N| + |A|))$ ; the approximate method runs in time  $O(|A| \sqrt{|N|})$  with an additive error of  $O(\sqrt{|N|})$ . For more information about shortest paths, see the section “[Shortest Path](#)” on page 142. If you use the `SHORTPATH=` option, the following columns also appear in the summary output data set:

- `diameter_wt`: longest weighted shortest path in the graph
- `diameter_unwt`: longest unweighted shortest path in the graph
- `diameter_approx_wt`: approximate longest weighted shortest path in the graph
- `diameter_approx_unwt`: approximate longest unweighted shortest path in the graph
- `avg_shortpath_wt`: average weighted shortest path in the graph
- `avg_shortpath_unwt`: average unweighted shortest path in the graph

Depending on which other options you specify, some of these columns might not appear in the summary output data set.

### **`OUT_NODES=` Data Set**

In addition, you can produce summary statistics about the nodes of the graph. By default, the following columns are appended to the data set that you specify in the `OUT_NODES=` option in the `PROC OPTGRAPH` statement:

- `sum_in_and_out_wt`: sum of the link weights from and to the node
- `leaf_node`: 1, if the node is a leaf node; otherwise, 0
- `singleton_node`: 1, if the node is a singleton node; otherwise, 0
- `isolated_pair`: the identifier, if the node is in an isolated pair; otherwise, missing (.)
- `neighbor_leaf_nodes`: the number of leaf nodes connected to the node

In addition, the following column is appended for undirected graphs:

- `isolated_star`: the identifier, if the node is in an isolated star; otherwise, missing (.)

The following columns are appended for directed graphs:

- `isolated_star_out`: the identifier, if the node is in an isolated outward star; otherwise, missing (.)
- `isolated_star_in`: the identifier, if the node is in an isolated inward star; otherwise, missing (.)

You can produce statistics about the shortest path distances to and from nodes in the graph by using the `SHORTPATH=` option. The *eccentricity* of a node  $u$  is the longest shortest path distance of all possible shortest path distances between  $u$  and any other node. If you use the `SHORTPATH=` option, the following columns also appear in the nodes output data set:

- `eccentr_out_wt`: the longest weighted shortest path distance from the node
- `eccentr_out_unwt`: the longest unweighted shortest path distance from the node
- `eccentr_in_wt`: the longest weighted shortest path distance to the node
- `eccentr_in_unwt`: the longest unweighted shortest path distance to the node

### ***OUT\_LINKS= Data Set***

In addition, you can produce summary statistics about the links of the graph. By default, the following columns are appended to the data set that you specify in the `OUT_LINKS=` option in the `PROC OPTGRAPH` statement, for undirected graphs:

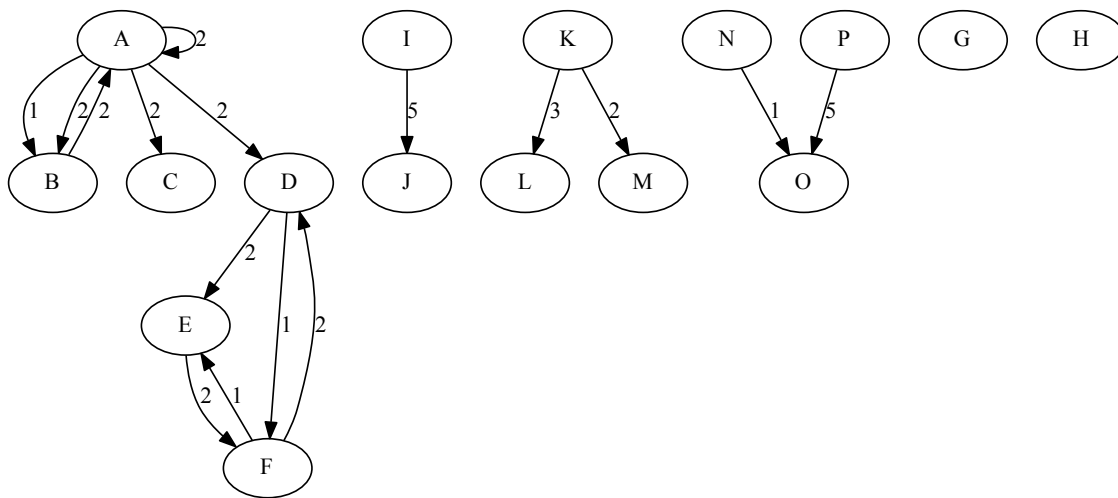
- `isolated_pair`: the identifier, if the link is in an isolated pair; otherwise, missing (.)
- `isolated_star`: the identifier, if the link is in an isolated star; otherwise, missing (.)

The following columns are appended for directed graphs:

- `isolated_star_out`: the identifier, if the link is in an isolated outward star; otherwise, missing (.)
- `isolated_star_in`: the identifier, if the link is in an isolated inward star; otherwise, missing (.)

## **Summary Statistics of a Simple Directed Graph**

This section illustrates the calculation of summary statistics on the simple directed graph  $G$  that is shown in Figure 1.115.

**Figure 1.115** A Simple Directed Graph *G*

You can represent the directed graph *G* by using the following nodes data set, NodeSetIn, and links data set, LinkSetIn:

```

data NodeSetIn;
  input node $ @@;
  datalines;
A B C D E F G H I J K L M N O P
;
data LinkSetIn;
  input from $ to $ weight @@;
  datalines;
A B 1 A C 2 A D 2 B A 2 D E 2
D F 1 E F 2 F D 2 F E 1 A A 2
A B 2 I J 5 K L 3 K M 2 N O 1
P O 1
;

```

The following statements calculate the default summary statistics and output the results in the data set Summary:

```

proc optgraph
  graph_direction = directed
  data_nodes      = NodeSetIn
  data_links      = LinkSetIn;
  summary
    out           = Summary;
run;

```

The data set Summary contains the default summary statistics of the input graph and is shown in Figure 1.116.

**Figure 1.116** Graph Summary Statistics of a Simple Directed Graph

nodes	links	avg_links_per_node	density	self_links_ignored	dup_links_ignored	leaf_nodes	singleton_nodes
16	14	0.875	0.058333	1	1	5	2

The following statements calculate the default summary statistics and information about the connectedness of the graph, and they output the results in the data set Summary:

```
proc optgraph
  graph_direction = directed
  data_nodes      = NodeSetIn
  data_links      = LinkSetIn;
  summary
    concomp
    out           = Summary;
run;
```

The data set Summary contains the summary statistics of the input graph and is shown in Figure 1.117.

**Figure 1.117** Graph Summary and Connectedness Statistics of a Simple Directed Graph

nodes	links	avg_links_per_node	density	self_links_ignored	dup_links_ignored	leaf_nodes	singleton_nodes
16	14	0.875	0.058333	1	1	5	2

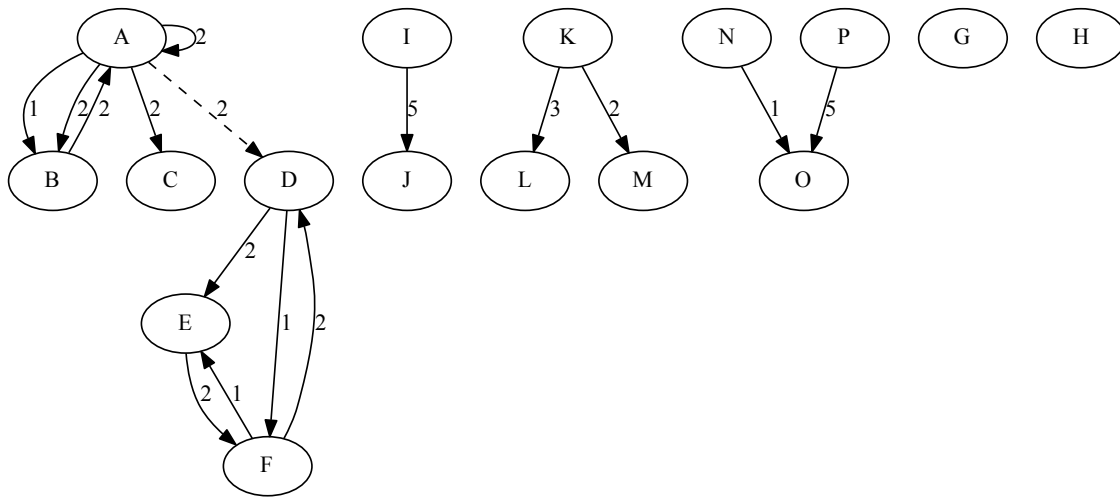
  

concomp	isolated_pairs	isolated_stars_out	isolated_stars_in
13	1	1	1

## Summary Statistics of a Simple Directed Graph by Cluster

Similar to how you can use the BY\_CLUSTER option in the CENTRALITY statement, as described in the section “Processing by Cluster” on page 78, you can process a number of induced subgraphs of a graph with only one call to PROC OPTGRAPH by using the BY\_CLUSTER option in the SUMMARY statement. In this section, you want to work on the subgraphs that are induced by node subsets  $N_0 = \{A, B, C\}$ ,  $N_1 = \{D, E, F\}$ , and  $N_3 = \{G, H, I, J, K, L, M, N, O, P\}$  for the directed graph shown in Figure 1.115. The induced subgraphs are shown graphically in Figure 1.118 (the dashed link is removed).



**Figure 1.118** Induced Subgraphs of Graph *G*

Define the subgraphs in the nodes data set by using the cluster variable as follows:

```
data NodeSetIn;
  input node $ cluster @@;
  datalines;
A 0 B 0 C 0 D 1 E 1
F 1 G 2 H 2 I 2 J 2
K 2 L 2 M 2 N 2 O 2
P 2
;
```

The following statements process the summary statistics for each induced subgraph:

```
proc optgraph
  graph_direction = directed
  data_links      = LinkSetIn
  data_nodes      = NodeSetIn
  out_links       = LinkSetOut
  out_nodes       = NodeSetOut;
  performance
    nthreads      = 2;
  summary
    by_cluster
    concomp
    out            = Summary;
run;
```

Notice in this example that you can process each subgraph in parallel by using the `NTHREADS=` option in the **PERFORMANCE** statement.

The data sets `Summary`, `NodeSetOut`, and `LinkSetOut` now contain the summary statistics for each induced subgraph; they are shown in [Figure 1.119](#).

**Figure 1.119** Summary Statistics for Induced Subgraphs of Graph *G***Summary**

cluster	nodes	links	avg_links_per_node	density	self_links_ignored	dup_links_ignored	leaf_nodes
0	3	3	1.00000	0.50000		1	1
1	3	5	1.66667	0.83333	0	0	0
2	10	5	0.50000	0.05556	0	0	4

singleton_nodes	concomp	isolated_pairs	isolated_stars_out	isolated_stars_in
0	2	0	0	0
0	1	0	0	0
2	10	1	1	1

**NodeSetOut**

node	cluster	sum_in_and_out_wt	leaf_node	singleton_node	isolated_pair	isolated_star_out	isolated_star_in	neighbor_leaf_nodes
A	0	5	0	0	.	.	.	1
B	0	3	0	0	.	.	.	0
C	0	2	1	0	.	.	.	0
D	1	5	0	0	.	.	.	0
E	1	5	0	0	.	.	.	0
F	1	6	0	0	.	.	.	0
G	2	0	0	1	.	.	.	0
H	2	0	0	1	.	.	.	0
I	2	5	0	0	1	.	.	1
J	2	5	1	0	1	.	.	0
K	2	5	0	0	.	1	.	2
L	2	3	1	0	.	1	.	0
M	2	2	1	0	.	1	.	0
N	2	1	0	0	.	.	1	1
O	2	2	1	0	.	.	1	0
P	2	1	0	0	.	.	1	1

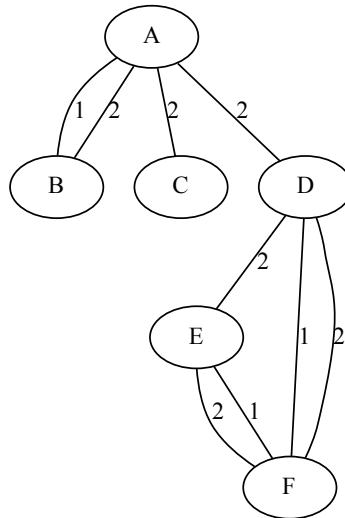
**LinkSetOut**

from	to	cluster	weight	isolated_pair	isolated_star_out	isolated_star_in
A	B	0	1	.	.	.
A	C	0	2	.	.	.
A	D	.	2	.	.	.
B	A	0	2	.	.	.
D	E	1	2	.	.	.
D	F	1	1	.	.	.
E	F	1	2	.	.	.
F	D	1	2	.	.	.
F	E	1	1	.	.	.
I	J	2	5	1	.	.
K	L	2	3	.	1	.
K	M	2	2	.	1	.
N	O	2	1	.	.	1
P	O	2	1	.	.	1

## Summary Statistics of a Simple Undirected Graph

This section illustrates the calculation of summary and shortest path statistics on the simple undirected graph  $G$  that is shown in Figure 1.120.

**Figure 1.120** A Simple Undirected Graph  $G$



You can represent the undirected graph  $G$  by using the following links data set, LinkSetIn:

```

data LinkSetIn;
  input from $ to $ weight @@;
  datalines;
A B 1 A C 2 A D 2 B A 2 D E 2
D F 1 E F 2 F D 2 F E 1
;

```

The following statements calculate the default summary statistics and information about shortest path distances of the graph, and they output the results in the data set Summary. In addition, node statistics are produced and output in the data set NodeSetOut.

```

proc optgraph
  data_links      = LinkSetIn
  out_nodes       = NodeSetOut;
  summary
    out           = Summary
    shortestpath  = weight;
run;

```

The data sets Summary and NodeSetOut now contain the summary statistics of the input graph, which are shown in Figure 1.121.

**Figure 1.121** Graph Summary and Shortest Path Statistics of a Simple Undirected Graph**Summary**

nodes	links	avg_links_per_node	density	self_links_ignored	dup_links_ignored	leaf_nodes	singleton_nodes	diameter_wt	avg_shortpath_wt
6	6	1	0.4	0	3	2	0	6	3.13333

**NodeSetOut**

node	sum_in_and_out_wt	eccentr_wt_out	leaf_node	singleton_node	neighbor_leaf_nodes
A	5	4	0	0	2
B	1	5	1	0	0
C	2	6	1	0	0
D	5	4	0	0	0
E	4	6	0	0	0
F	3	5	0	0	0

**Transitive Closure**

The *transitive closure* of a graph  $G$  is a graph  $G^T = (N, A^T)$  such that for all  $i, j \in N$  there is a link  $(i, j) \in A^T$  if and only if there exists a path from  $i$  to  $j$  in  $G$ .

The transitive closure of a graph can help to efficiently answer questions about reachability. Suppose you want to answer the question of whether you can get from node  $i$  to node  $j$  in the original graph  $G$ . Given the transitive closure  $G^T$  of  $G$ , you can simply check for the existence of link  $(i, j)$  to answer the question. Transitive closure has many applications, including speeding up the processing of structured query languages, which are often used in databases.

In PROC OPTGRAPH, you can invoke the transitive closure algorithm by using the TRANSITIVE\_CLOSURE statement. The options for this statement are described in the section “[TRANSITIVE\\_CLOSURE Statement](#)” on page 44.

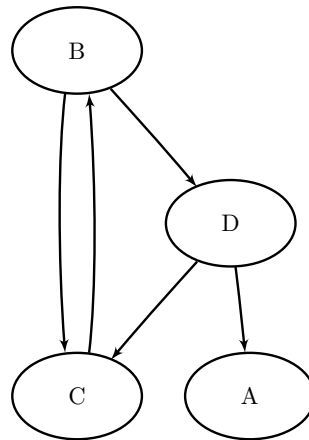
The results of the transitive closure algorithm are written to the output data set that is specified in the OUT= option in the TRANSITIVE\_CLOSURE statement. The links that define the transitive closure are listed in the output data set with variable names from and to.

The transitive closure algorithm reports status information in a macro variable called \_OPTGRAPH\_TRANSCL\_. For more information about this macro variable, see the section “[Macro Variable \\_OPTGRAPH\\_TRANSCL\\_](#)” on page 183.

The algorithm that the PROC OPTGRAPH uses to compute transitive closure is a sparse version of the Floyd-Warshall algorithm (Cormen, Leiserson, and Rivest 1990). This algorithm runs in time  $O(|N|^3)$  and therefore might not scale to very large graphs.

**Transitive Closure of a Simple Directed Graph**

This example illustrates the use of the transitive closure algorithm on the simple directed graph  $G$ , which is shown in [Figure 1.122](#).

**Figure 1.122** A Simple Directed Graph  $G$ 

The directed graph  $G$  can be represented by the following links data set LinkSetIn:

```

data LinkSetIn;
  input from $ to $ @@;
  datalines;
B C  B D  C B  D A  D C
;

```

The following statements calculate the transitive closure and output the results in the data set TransClosure:

```

proc optgraph
  graph_direction = directed
  data_links      = LinkSetIn;
  transitive_closure
    out           = TransClosure;
run;

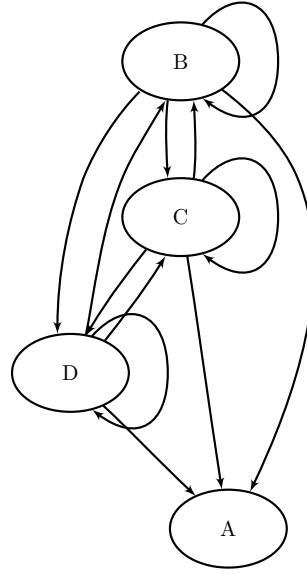
```

The data set TransClosure contains the transitive closure of  $G$  and is shown in [Figure 1.123](#).

**Figure 1.123** Transitive Closure of a Simple Directed Graph

from to	
B	C
C	B
B	D
D	C
D	A
C	C
C	D
B	B
D	B
D	D
B	A
C	A

The transitive closure of  $G$  is shown graphically in [Figure 1.124](#).

**Figure 1.124** Transitive Closure of  $G$ 

For a more detailed example, see “[Example 1.13: Transitive Closure for Identification of Circular Dependencies in a Bug Tracking System](#)” on page 223.

## Traveling Salesman Problem

The *traveling salesman problem* (TSP) finds a minimum-cost tour in a graph,  $G$ , that has a node set,  $N$ , and a link set,  $A$ . A *path* in a graph is a sequence of nodes, each of which has a link to the next node in the sequence. An *elementary cycle* is a path in which the start node and end node are the same and no node appears more than once in the sequence. A *Hamiltonian cycle* (or *tour*) is an elementary cycle that visits every node. In solving the TSP, then, the goal is to find a Hamiltonian cycle of minimum total cost, where the total cost is the sum of the costs of the links in the tour. Associated with each link  $(i, j) \in A$  are a binary variable  $x_{ij}$ , which indicates whether link  $x_{ij}$  is part of the tour, and a cost  $c_{ij}$ . Let  $\delta(S) = \{(i, j) \in A \mid i \in S, j \notin S\}$ . Then an integer linear programming formulation of the TSP (for an undirected graph  $G$ ) is as follows:

$$\begin{aligned}
 &\text{minimize} && \sum_{(i,j) \in A} c_{ij} x_{ij} \\
 &\text{subject to} && \sum_{(i,j) \in \delta(i)} x_{i,j} = 2 \quad i \in N && (\text{two\_match}) \\
 &&& \sum_{(i,j) \in \delta(S)} x_{ij} \geq 2 \quad S \subset N, 2 \leq |S| \leq |N| - 1 && (\text{subtour\_elim}) \\
 &&& x_{ij} \in \{0, 1\} && (i, j) \in A
 \end{aligned}$$

The equations (two\_match) are the *matching constraints*, which ensure that each node has degree two in the subgraph. The inequalities (subtour\_elim) are the *subtour elimination constraints* (SECs), which enforce connectivity.

For a directed graph,  $G$ , the same formulation and solution approach is used on an expanded graph  $G'$ , as described in Kumar and Li (1994). PROC OPTGRAPH takes care of the construction of the expanded graph and returns the solution in terms of the original input graph.

In practical terms, you can think of the TSP in the context of a routing problem in which each node is a city and the links are roads that connect those cities. If you know the distance between each pair of cities, the goal is to find the shortest possible route that visits each city exactly once. The TSP has applications in planning, logistics, manufacturing, genomics, and many other areas.

In PROC OPTGRAPH, you can invoke the traveling salesman problem solver by using the TSP statement. The options for this statement are described in the section “[TSP Statement](#)” on page 44.

The traveling salesman problem solver reports status information in a macro variable called `_OPTGRAPH_TSP_`. For more information about this macro variable, see the section “[Macro Variable `\_OPTGRAPH\_TSP\_`](#)” on page 183.

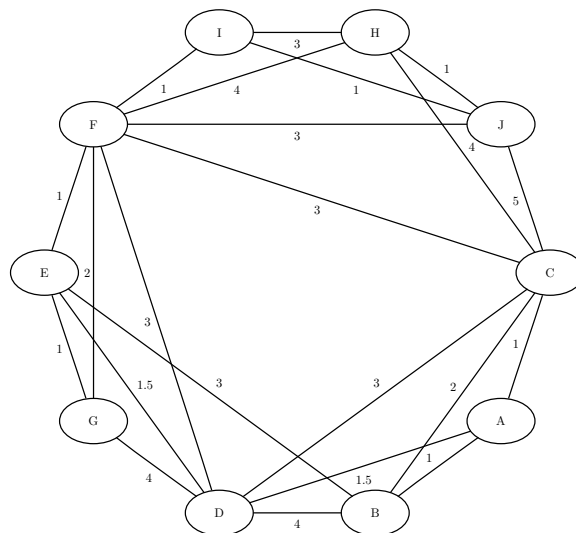
The algorithm that PROC OPTGRAPH uses for solving the TSP is based on a variant of the branch-and-cut process described in Applegate et al. (2006).

The resulting tour is represented in two ways: in the data set that you specify in the OUT\_NODES= option in the PROC OPTGRAPH statement, the tour is specified as a sequence of nodes; in the data set that you specify in the OUT= option of the TSP statement, the tour is specified as a list of links in the optimal tour.

## Traveling Salesman Problem Applied to a Simple Undirected Graph

As a simple example, consider the weighted undirected graph in Figure 1.125.

**Figure 1.125** A Simple Undirected Graph



You can represent the links data set as follows:

```
data LinkSetIn;
  input from $ to $ weight @@;
  datalines;
A B 1.0 A C 1.0 A D 1.5 B C 2.0 B D 4.0
B E 3.0 C D 3.0 C F 3.0 C H 4.0 D E 1.5
D F 3.0 D G 4.0 E F 1.0 E G 1.0 F G 2.0
F H 4.0 H I 3.0 I J 1.0 C J 5.0 F J 3.0
F I 1.0 H J 1.0
;
```

The following statements calculate an optimal traveling salesman tour and output the results in the data sets TSPTour and NodeSetOut:

```
proc optgraph
  loglevel    = moderate
  data_links  = LinkSetIn
  out_nodes   = NodeSetOut;
  tsp
    out       = TSPTour;
run;
%put &_OPTGRAPH_;
%put &_OPTGRAPH_TSP_;
```

The progress of the OPTGRAPH procedure is shown in [Figure 1.126](#).



**Figure 1.126** PROC OPTGRAPH Log: Optimal Traveling Salesman Tour of a Simple Undirected Graph

---

```

NOTE: -----
NOTE: -----
NOTE: Running OPTGRAPH version 14.1.
NOTE: -----
NOTE: -----
NOTE: The OPTGRAPH procedure is executing in single-machine mode.
NOTE: -----
NOTE: -----
NOTE: Reading the links data set.
NOTE: There were 22 observations read from the data set WORK.LINKSETIN.
NOTE: Data input used 0.00 (cpu: 0.02) seconds.
NOTE: Building the input graph storage used 0.00 (cpu: 0.00) seconds.
NOTE: The input graph storage is using 0.0 MBs (peak: 0.0 MBs) of memory.
NOTE: The number of nodes in the input graph is 10.
NOTE: The number of links in the input graph is 22.
NOTE: -----
NOTE: -----
NOTE: Processing the traveling salesman problem.
NOTE: The initial TSP heuristics found a tour with cost 16 using 0.00 (cpu: 0.00) seconds.
NOTE: The MILP presolver value NONE is applied.
NOTE: The MILP solver is called.
NOTE: The Branch and Cut algorithm is used.

```

Node	Active	Sols	BestInteger	BestBound	Gap	Time
0	1	1	16.0000000	15.5005000	3.22%	0
0	0	1	16.0000000	16.0000000	0.00%	0

```

NOTE: Objective = 16.
NOTE: Processing the traveling salesman problem used 0.01 (cpu: 0.00) seconds.
NOTE: -----
NOTE: -----
NOTE: Creating nodes data set output.
NOTE: Creating traveling salesman data set output.
NOTE: Data output used 0.00 (cpu: 0.02) seconds.
NOTE: -----
NOTE: -----
NOTE: The data set WORK.NODESETOUT has 10 observations and 2 variables.
NOTE: The data set WORK.TSPTOUR has 10 observations and 3 variables.
STATUS=OK  TSP=OPTIMAL
STATUS=OPTIMAL  OBJECTIVE=16  RELATIVE_GAP=0  ABSOLUTE_GAP=0  PRIMAL_INFEASIBILITY=0
BOUND_INFEASIBILITY=0  INTEGER_INFEASIBILITY=0  BEST_BOUND=16  NODES=1  ITERATIONS=17
CPU_TIME=0.00  REAL_TIME=0.01

```

---

The data set NodeSetOut now contains a sequence of nodes in the optimal tour and is shown in [Figure 1.127](#).

**Figure 1.127** Nodes in the Optimal Traveling Salesman Tour

node	tsp_order
A	1
B	2
C	3
H	4
J	5
I	6
F	7
G	8
E	9
D	10

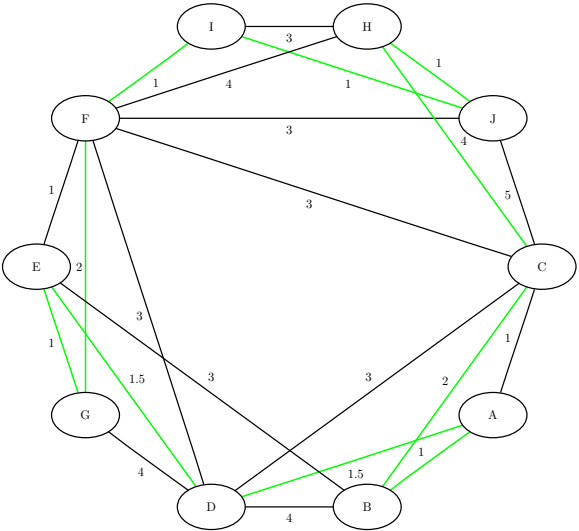
The data set TSPTour now contains the links in the optimal tour and is shown in [Figure 1.128](#).

**Figure 1.128** Links in the Optimal Traveling Salesman Tour

from	to	weight
A	B	1.0
B	C	2.0
C	H	4.0
H	J	1.0
I	J	1.0
F	I	1.0
F	G	2.0
E	G	1.0
D	E	1.5
A	D	1.5
		<b>16.0</b>

The minimum-cost links are shown in green in [Figure 1.129](#).

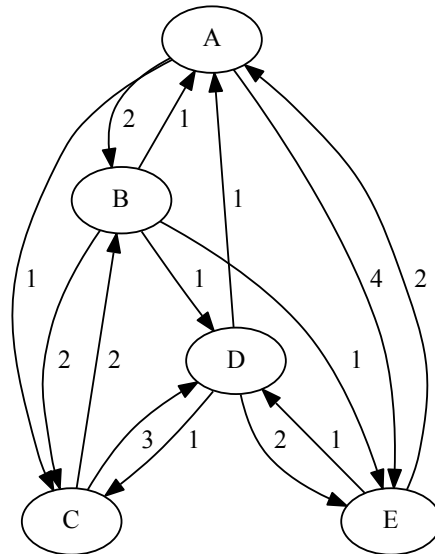
**Figure 1.129** Optimal Traveling Salesman Tour



## Traveling Salesman Problem Applied to a Simple Directed Graph

As another simple example, consider the weighted directed graph in Figure 1.130.

**Figure 1.130** A Simple Directed Graph



You can represent the links data set as follows:

```

data LinkSetIn;
  input from $ to $ weight @@;
  datalines;
A B 2.0 A C 1.0 A E 4.0 B A 1.0 B C 2.0
B D 1.0 B E 1.0 C B 2.0 C D 3.0 D A 1.0
D C 1.0 D E 2.0 E A 2.0 E D 1.0
;

```

The following statements calculate an optimal traveling salesman tour (on a directed graph) and output the results in the data sets TSPTour and NodeSetOut:

```

proc optgraph
  direction = directed
  loglevel = moderate
  data_links = LinkSetIn
  out_nodes = NodeSetOut;
  tsp
    out = TSPTour;
run;
%put &_OPTGRAPH_;
%put &_OPTGRAPH_TSP_;

```

The progress of the OPTGRAPH procedure is shown in Figure 1.131.

**Figure 1.131** PROC OPTGRAPH Log: Optimal Traveling Salesman Tour of a Simple Directed Graph

```

NOTE: -----
NOTE: -----
NOTE: Running OPTGRAPH version 14.1.
NOTE: -----
NOTE: -----
NOTE: The OPTGRAPH procedure is executing in single-machine mode.
NOTE: -----
NOTE: -----
NOTE: Reading the links data set.
NOTE: There were 14 observations read from the data set WORK.LINKSETIN.
NOTE: Data input used 0.00 (cpu: 0.00) seconds.
NOTE: Building the input graph storage used 0.00 (cpu: 0.00) seconds.
NOTE: The input graph storage is using 0.0 MBs (peak: 0.0 MBs) of memory.
NOTE: The number of nodes in the input graph is 5.
NOTE: The number of links in the input graph is 14.
NOTE: -----
NOTE: -----
NOTE: The TSP solver is starting using an augmented symmetric graph with 10 nodes and 19 links.
NOTE: -----
NOTE: -----
NOTE: Processing the traveling salesman problem.
NOTE: The initial TSP heuristics found a tour with cost 6 using 0.00 (cpu: 0.00) seconds.
NOTE: The MILP presolver value NONE is applied.
NOTE: The MILP solver is called.
NOTE: The Branch and Cut algorithm is used.

```

Node	Active	Sols	BestInteger	BestBound	Gap	Time
0	1	1	6.0000000	5.9001000	1.69%	0
0	0	1	6.0000000	6.0000000	0.00%	0

```

NOTE: Objective = 6.
NOTE: Processing the traveling salesman problem used 0.00 (cpu: 0.00) seconds.
NOTE: -----
NOTE: -----
NOTE: Creating nodes data set output.
NOTE: Creating traveling salesman data set output.
NOTE: Data output used 0.00 (cpu: 0.00) seconds.
NOTE: -----
NOTE: -----
NOTE: The data set WORK.NODESETOUT has 5 observations and 2 variables.
NOTE: The data set WORK.TSPTOUR has 5 observations and 3 variables.
STATUS=OK  TSP=OPTIMAL
STATUS=OPTIMAL  OBJECTIVE=6  RELATIVE_GAP=2.96E-16  ABSOLUTE_GAP=-1.77636E-15
PRIMAL_INFEASIBILITY=0  BOUND_INFEASIBILITY=0  INTEGER_INFEASIBILITY=0  BEST_BOUND=6  NODES=1
ITERATIONS=8  CPU_TIME=0.00  REAL_TIME=0.00

```

The data set NodeSetOut now contains a sequence of nodes in the optimal tour and is shown in [Figure 1.132](#).

**Figure 1.132** Nodes in the Optimal Traveling Salesman Tour

node	tsp_order
A	1
C	2
B	3
E	4
D	5

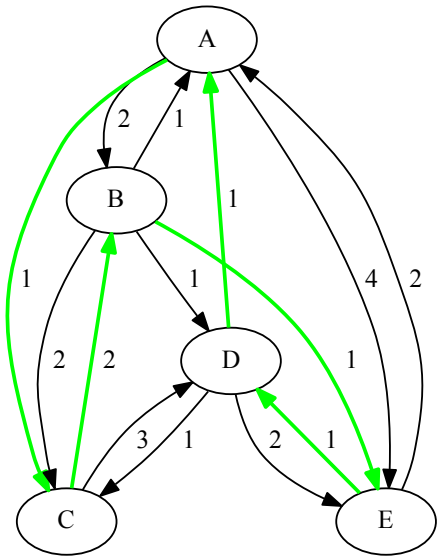
The data set TSPTour now contains the links in the optimal tour and is shown in [Figure 1.133](#).

**Figure 1.133** Links in the Optimal Traveling Salesman Tour

from	to	weight
A	C	1
C	B	2
B	E	1
E	D	1
D	A	1
		<b>6</b>

The minimum-cost links are shown in green in [Figure 1.134](#).

**Figure 1.134** Optimal Traveling Salesman Tour



## Macro Variables

The OPTGRAPH procedure defines one summary macro variable that reports the overall status and one detailed macro variable for each completed algorithm.

### Macro Variable `_OPTGRAPH_`

The OPTGRAPH procedure defines a macro variable named `_OPTGRAPH_`. This variable contains a character string that indicates the status of the OPTGRAPH procedure upon termination. The various terms of the variable are interpreted as follows:

#### STATUS

indicates the status of the procedure at termination. The STATUS term can take one of the following values:

OK	The procedure terminated normally.
OUT_OF_MEMORY	Insufficient memory was allocated to the procedure.
INTERRUPTED	The procedure was interrupted by the user during the input or setup phase.
ERROR	The procedure encountered an error.

#### BICONCOMP

indicates the status of the biconnected components algorithm at termination. This algorithm is described in the section “[Biconnected Components and Articulation Points](#)” on page 61. The BICONCOMP term can take one of the following values:

OK	The algorithm terminated normally.
ERROR	The algorithm encountered an error.

#### CENTRALITY

indicates the status of the centrality algorithms at termination. These algorithms are described in the section “[Centrality](#)” on page 65. The CENTRALITY term can take one of the following values:

OK	The algorithm terminated normally.
INTERRUPTED	The algorithm was interrupted by the user.
ERROR	The algorithm encountered an error.

#### CLIQUE

indicates the status of the clique-finding algorithms at termination. These algorithms are described in the section “[Clique](#)” on page 86. The CLIQUE term can take one of the following values:

OK	The algorithm terminated normally.
INTERRUPTED	The algorithm was interrupted by the user.
TIMELIMIT	The algorithm reached its execution time limit, which is specified in the <code>MAXTIME=</code> option in the <code>CLIQUE</code> statement.

SOLUTION_LIM	The algorithm reached its limit on the number of cliques found, which is specified in the <b>MAXCLIQUES=</b> option in the <b>CLIQUE</b> statement.
ERROR	The algorithm encountered an error.

**COMMUNITY**

indicates the status of the community algorithms at termination. These algorithms are described in the section “**Community Detection**” on page 89. The **COMMUNITY** term can take one of the following values:

OK	The algorithm terminated normally.
INTERRUPTED	The algorithm was interrupted by the user.
ERROR	The algorithm encountered an error.

**CONCOMP**

indicates the status of the connected components algorithm at termination. This algorithm is described in the section “**Connected Components**” on page 97. The **CONCOMP** term can take one of the following values:

OK	The algorithm terminated normally.
ERROR	The algorithm encountered an error.

**CORE**

indicates the status of the core decomposition algorithm at termination. This algorithm is described in the section “**Core Decomposition**” on page 102. The **CORE** term can take one of the following values:

OK	The algorithm terminated normally.
INTERRUPTED	The algorithm was interrupted by the user.
TIMELIMIT	The algorithm reached its execution time limit, which is specified in the <b>MAXTIME=</b> option in the <b>CORE</b> statement.
ERROR	The algorithm encountered an error.

**CYCLE**

indicates the status of the cycle detection algorithm at termination. This algorithm is described in the section “**Cycle**” on page 107. The **CYCLE** term can take one of the following values:

OK	The algorithm terminated normally.
TIMELIMIT	The algorithm reached its execution time limit, which is specified in the <b>MAXTIME=</b> option in the <b>CYCLE</b> statement.
SOLUTION_LIM	The algorithm reached its limit on the number of cycles found, which is specified in the <b>MAXCYCLES=</b> option in the <b>CYCLE</b> statement.
ERROR	The algorithm encountered an error.

**EIGEN**

indicates the status of the eigenvector solver at termination. This solver is described in the section “[Eigenvector Problem](#)” on page 112. The EIGEN term can take one of the following values:

OK	The solver terminated normally.
ITER_LIMIT_SOL	The solver reached the maximum number of iterations that is specified in the MAXITER= option and found a solution.
ITER_LIMIT_NOSOL	The solver reached the maximum number of iterations that is specified in the MAXITER= option and did not find a solution.
ERROR	The solver encountered an error.

**LAP**

indicates the status of the linear assignment solver at termination. This solver is described in the section “[Linear Assignment \(Matching\)](#)” on page 115. The LAP term can take one of the following values:

OPTIMAL	The solution is optimal.
INFEASIBLE	The problem is infeasible.
ERROR	The solver encountered an error.

**MCF**

indicates the status of the minimum-cost network flow solver at termination. This solver is described in the section “[Minimum-Cost Network Flow](#)” on page 116. The MCF term can take one of the following values:

OPTIMAL	The solution is optimal.
OPTIMAL_COND	The solution is optimal, but some infeasibilities (primal or bound) exceed tolerances because of scaling.
INFEASIBLE	The problem is infeasible.
UNBOUNDED	The problem is unbounded.
TIMELIMIT	The solver reached its execution time limit, which is specified in the <a href="#">MAX-TIME=</a> option in the <a href="#">MINCOSTFLOW</a> statement.
FAIL_NOSOL	The solver stopped because of errors and did not find a solution.
ERROR	The solver encountered an error.

**MINCUT**

indicates the status of the minimum-cut solver at termination. This solver is described in the section “[Minimum Cut](#)” on page 124. The MINCUT term can take one of the following values:

OPTIMAL	The solution is optimal.
INFEASIBLE	The problem is infeasible.
INTERRUPTED	The solver was interrupted by the user.
ERROR	The solver encountered an error.



**MST**

indicates the status of the minimum spanning tree solver at termination. This solver is described in the section “[Minimum Spanning Tree](#)” on page 128. The MST term can take one of the following values:

OPTIMAL	The solution is optimal.
INTERRUPTED	The algorithm was interrupted by the user.
ERROR	The solver encountered an error.

**REACH**

indicates the status of the reach algorithms at termination. These algorithms are described in the section “[Reach \(Ego\) Network](#)” on page 130. The REACH term can take one of the following values:

OK	The algorithm terminated normally.
INTERRUPTED	The algorithm was interrupted by the user.
ERROR	The algorithm encountered an error.

**SHORTPATH**

indicates the status of the shortest path algorithms at termination. These algorithms are described in the section “[Shortest Path](#)” on page 142. The SHORTPATH term can take one of the following values:

OK	The algorithm terminated normally.
INTERRUPTED	The algorithm was interrupted by the user.
ERROR	The algorithm encountered an error.

**SUMMARY**

indicates the status of the summary algorithms at termination. These algorithms are described in the section “[Summary](#)” on page 153. The SUMMARY term can take one of the following values:

OK	The algorithm terminated normally.
INTERRUPTED	The algorithm was interrupted by the user.
ERROR	The algorithm encountered an error.

**TRANSCL**

indicates the status of the transitive closure algorithm at termination. This algorithm is described in the section “[Transitive Closure](#)” on page 162. The TRANSCL term can take one of the following values:

OK	The algorithm terminated normally.
INTERRUPTED	The algorithm was interrupted by the user.
ERROR	The algorithm encountered an error.

**TSP**

indicates the status of the traveling salesman problem solver at termination. This algorithm is described in the section “[Traveling Salesman Problem](#)” on page 164. The TSP term can take one of the following values:

OPTIMAL	The solution is optimal.
OPTIMAL_AGAP	The solution is optimal within the absolute gap that is specified in the ABSOBJGAP= option.
OPTIMAL_RGAP	The solution is optimal within the relative gap that is specified in the RELOBJGAP= option.
OPTIMAL_COND	The solution is optimal, but some infeasibilities (primal, bound, or integer) exceed tolerances because of scaling.
TARGET	The solution is not worse than the target that is specified in the TARGET= option.
INFEASIBLE	The problem is infeasible.
UNBOUNDED	The problem is unbounded.
INFEASIBLE_OR_UNBOUNDED	The problem is infeasible or unbounded.
SOLUTION_LIM	The solver reached the maximum number of solutions specified in the MAXSOLS= option.
NODE_LIM_SOL	The solver reached the maximum number of nodes specified in the MAXNODES= option and found a solution.
NODE_LIM_NOSOL	The solver reached the maximum number of nodes specified in the MAXNODES= option and did not find a solution.
TIME_LIM_SOL	The solver reached the execution time limit specified in the MAXTIME= option and found a solution.
TIME_LIM_NOSOL	The solver reached the execution time limit specified in the MAXTIME= option and did not find a solution.
HEURISTIC_SOL	The solver used only heuristics and found a solution.
HEURISTIC_NOSOL	The solver used only heuristics and did not find a solution.
ABORT_SOL	The solver was stopped by the user but still found a solution.
ABORT_NOSOL	The solver was stopped by the user and did not find a solution.
OUTMEM_SOL	The solver ran out of memory but still found a solution.
OUTMEM_NOSOL	The solver ran out of memory and either did not find a solution or failed to output the solution due to insufficient memory.
FAIL_SOL	The solver stopped due to errors but still found a solution.
FAIL_NOSOL	The solver stopped due to errors and did not find a solution.

Each algorithm reports its own status information in an additional macro variable. The following sections provide more information about these macro variables.

### Macro Variable `_OPTGRAPH_BICONCOMP_`

The OPTGRAPH procedure defines a macro variable named `_OPTGRAPH_BICONCOMP_`. This variable contains a character string that indicates the status and some basic statistics about the results of the algorithm that PROC OPTGRAPH uses to calculate biconnected components. The various terms of the variable are interpreted as follows:

**STATUS**

indicates the status of the algorithm at termination. The STATUS term takes the same value as the term **BICONCOMP** in the `_OPTGRAPH_` macro as defined in the section “[Macro Variable \\_OPTGRAPH\\_](#)” on page 172.

**NUM\_COMPONENTS**

indicates the number of biconnected components found by the algorithm.

**NUM\_ARTICULATION\_POINTS**

indicates the number of articulation points found by the algorithm.

**CPU\_TIME**

indicates the CPU time (in seconds) taken by the algorithm.

**REAL\_TIME**

indicates the real time (in seconds) taken by the algorithm.

**Macro Variable `_OPTGRAPH_CENTRALITY_`**

The OPTGRAPH procedure defines a macro variable named `_OPTGRAPH_CENTRALITY_`. This variable contains a character string that indicates the status and some basic statistics about the results of the algorithm that PROC OPTGRAPH uses to calculate centrality. The various terms of the variable are interpreted as follows:

**STATUS**

indicates the status of the algorithm at termination. The STATUS term takes the same value as the term **CENTRALITY** in the `_OPTGRAPH_` macro as defined in the section “[Macro Variable \\_OPTGRAPH\\_](#)” on page 172.

**CPU\_TIME**

indicates the CPU time (in seconds) taken by the algorithm.

**REAL\_TIME**

indicates the real time (in seconds) taken by the algorithm.

**Macro Variable `_OPTGRAPH_CLIQUE_`**

The OPTGRAPH procedure defines a macro variable named `_OPTGRAPH_CLIQUE_`. This variable contains a character string that indicates the status and some basic statistics about the results of the algorithm that PROC OPTGRAPH uses to calculate cliques. The various terms of the variable are interpreted as follows:

**STATUS**

indicates the status of the algorithm at termination. The STATUS term takes the same value as the term **CLIQUE** in the `_OPTGRAPH_` macro as defined in the section “[Macro Variable \\_OPTGRAPH\\_](#)” on page 172.

**NUM\_CLIQUES**

indicates the number of cliques found by the algorithm.

**CPU\_TIME**

indicates the CPU time (in seconds) taken by the algorithm.

**REAL\_TIME**

indicates the real time (in seconds) taken by the algorithm.

**Macro Variable \_OPTGRAPH\_COMMUNITY\_**

The OPTGRAPH procedure defines a macro variable named `_OPTGRAPH_COMMUNITY_`. This variable contains a character string that indicates the status and some basic statistics about the results of the algorithm that PROC OPTGRAPH uses to calculate communities. The various terms of the variable are interpreted as follows:

**STATUS**

indicates the status of the algorithm at termination. The STATUS term takes the same value as the term `COMMUNITY` in the `_OPTGRAPH_` macro as defined in the section “Macro Variable `_OPTGRAPH_`” on page 172.

**RESOLUTION**

indicates the list of resolution levels specified in the `RESOLUTION_LIST=` option.

**NUM\_COMMUNITIES**

indicates the number of communities found by the algorithm at each resolution level.

**MODULARITY**

indicates the final modularity found by the algorithm at each resolution level.

**CPU\_TIME**

indicates the CPU time (in seconds) taken by the algorithm.

**REAL\_TIME**

indicates the real time (in seconds) taken by the algorithm.

**Macro Variable \_OPTGRAPH\_CONCOMP\_**

The OPTGRAPH procedure defines a macro variable named `_OPTGRAPH_CONCOMP_`. This variable contains a character string that indicates the status and some basic statistics about the results of the algorithm that PROC OPTGRAPH uses to calculate connected components. The various terms of the variable are interpreted as follows:

**STATUS**

indicates the status of the algorithm at termination. The STATUS term takes the same value as the term `CONCOMP` in the `_OPTGRAPH_` macro as defined in the section “Macro Variable `_OPTGRAPH_`” on page 172.

**NUM\_COMPONENTS**

indicates the number of connected components found by the algorithm.

**CPU\_TIME**

indicates the CPU time (in seconds) taken by the algorithm.

**REAL\_TIME**

indicates the real time (in seconds) taken by the algorithm.

**Macro Variable \_OPTGRAPH\_CORE\_**

The OPTGRAPH procedure defines a macro variable named `_OPTGRAPH_CORE_`. This variable contains a character string that indicates the status and some basic statistics about the results of the algorithm that PROC OPTGRAPH uses to calculate the core decomposition. The various terms of the variable are interpreted as follows:

**STATUS**

indicates the status of the algorithm at termination. The STATUS term takes the same value as the term **CORE** in the `_OPTGRAPH_` macro as defined in the section “[Macro Variable \\_OPTGRAPH\\_](#)” on page 172.

**CPU\_TIME**

indicates the CPU time (in seconds) taken by the algorithm.

**REAL\_TIME**

indicates the real time (in seconds) taken by the algorithm.

**Macro Variable \_OPTGRAPH\_CYCLE\_**

The OPTGRAPH procedure defines a macro variable named `_OPTGRAPH_CYCLE_`. This variable contains a character string that indicates the status and some basic statistics about the results of the algorithm that PROC OPTGRAPH uses to calculate cycles. The various terms of the variable are interpreted as follows:

**STATUS**

indicates the status of the algorithm at termination. The STATUS term takes the same value as the term **CYCLE** in the `_OPTGRAPH_` macro as defined in the section “[Macro Variable \\_OPTGRAPH\\_](#)” on page 172.

**NUM\_CYCLES**

indicates the number of cycles found by the algorithm.

**CPU\_TIME**

indicates the CPU time (in seconds) taken by the algorithm.

**REAL\_TIME**

indicates the real time (in seconds) taken by the algorithm.

**Macro Variable \_OPTGRAPH\_EIGEN\_**

The OPTGRAPH procedure defines a macro variable named `_OPTGRAPH_EIGEN_`. This variable contains a character string that indicates the status and some basic statistics about the results of the algorithm that PROC OPTGRAPH uses to calculate eigenvectors. The various terms of the variable are interpreted as follows:

**STATUS**

indicates the status of the algorithm at termination. The STATUS term takes the same value as the term **EIGEN** in the `_OPTGRAPH_` macro as defined in the section “[Macro Variable \\_OPTGRAPH\\_](#)” on page 172.

**CPU\_TIME**

indicates the CPU time (in seconds) taken by the algorithm.

**REAL\_TIME**

indicates the real time (in seconds) taken by the algorithm.

**Macro Variable \_OPTGRAPH\_LAP\_**

The OPTGRAPH procedure defines a macro variable named `_OPTGRAPH_LAP_`. This variable contains a character string that indicates the status and some basic statistics about the results of the algorithm that PROC OPTGRAPH uses to solve the linear assignment problem. The various terms of the variable are interpreted as follows:

**STATUS**

indicates the status of the solver at termination. The STATUS term takes the same value as the term **LAP** in the `_OPTGRAPH_` macro as defined in the section “[Macro Variable \\_OPTGRAPH\\_](#)” on page 172.

**OBJECTIVE**

indicates the total weight of the minimum linear assignment.

**CPU\_TIME**

indicates the CPU time (in seconds) taken by the solver.

**REAL\_TIME**

indicates the real time (in seconds) taken by the solver.

**Macro Variable \_OPTGRAPH\_MCF\_**

The OPTGRAPH procedure defines a macro variable named `_OPTGRAPH_MCF_`. This variable contains a character string that indicates the status and some basic statistics about the results of the algorithm that PROC OPTGRAPH uses to solve the minimum-cost network flow problem. The various terms of the variable are interpreted as follows:

**STATUS**

indicates the status of the solver at termination. The STATUS term takes the same value as the term **MCF** in the `_OPTGRAPH_` macro as defined in the section “[Macro Variable \\_OPTGRAPH\\_](#)” on page 172.

**OBJECTIVE**

indicates the total link weight of the minimum-cost network flow.

**CPU\_TIME**

indicates the CPU time (in seconds) taken by the solver.

**REAL\_TIME**

indicates the real time (in seconds) taken by the solver.

**Macro Variable \_OPTGRAPH\_MINCUT\_**

The OPTGRAPH procedure defines a macro variable named `_OPTGRAPH_MINCUT_`. This variable contains a character string that indicates the status and some basic statistics about the results of the algorithm that PROC OPTGRAPH uses to find the minimum cut. The various terms of the variable are interpreted as follows:

**STATUS**

indicates the status of the algorithm at termination. The STATUS term takes the same value as the term **MINCUT** in the `_OPTGRAPH_` macro as defined in the section “[Macro Variable \\_OPTGRAPH\\_](#)” on page 172.

**OBJECTIVE**

indicates the total link weight of the minimum cut.

**CPU\_TIME**

indicates the CPU time (in seconds) taken by the algorithm.

**REAL\_TIME**

indicates the real time (in seconds) taken by the algorithm.

**Macro Variable \_OPTGRAPH\_MST\_**

The OPTGRAPH procedure defines a macro variable named `_OPTGRAPH_MST_`. This variable contains a character string that indicates the status and some basic statistics about the results of the algorithm that PROC OPTGRAPH uses to solve the minimum spanning tree problem. The various terms of the variable are interpreted as follows:

**STATUS**

indicates the status of the solver at termination. The STATUS term takes the same value as the term **MST** in the `_OPTGRAPH_` macro as defined in the section “[Macro Variable \\_OPTGRAPH\\_](#)” on page 172.

**OBJECTIVE**

indicates the total link weight of the minimum spanning tree.

**CPU\_TIME**

indicates the CPU time (in seconds) taken by the solver.

**REAL\_TIME**

indicates the real time (in seconds) taken by the solver.

**Macro Variable \_OPTGRAPH\_REACH\_**

The OPTGRAPH procedure defines a macro variable named `_OPTGRAPH_REACH_`. This variable contains a character string that indicates the status and some basic statistics about the results of the algorithm that PROC OPTGRAPH uses to calculate reach networks. The various terms of the variable are interpreted as follows:

**STATUS**

indicates the status of the algorithm at termination. The STATUS term takes the same value as the term **REACH** in the `_OPTGRAPH_` macro as defined in the section “[Macro Variable \\_OPTGRAPH\\_](#)” on page 172.

**CPU\_TIME**

indicates the CPU time (in seconds) taken by the algorithm.

**REAL\_TIME**

indicates the real time (in seconds) taken by the algorithm.

**Macro Variable \_OPTGRAPH\_SHORTPATH\_**

The OPTGRAPH procedure defines a macro variable named `_OPTGRAPH_SHORTPATH_`. This variable contains a character string that indicates the status and some basic statistics about the results of the algorithm that PROC OPTGRAPH uses to calculate shortest paths. The various terms of the variable are interpreted as follows:

**STATUS**

indicates the status of the algorithm at termination. The STATUS term takes the same value as the term **SHORTPATH** in the `_OPTGRAPH_` macro as defined in the section “[Macro Variable \\_OPTGRAPH\\_](#)” on page 172.

**NUM\_PATHS**

indicates the number of shortest paths that the algorithm finds.

**CPU\_TIME**

indicates the CPU time (in seconds) taken by the algorithm.

**REAL\_TIME**

indicates the real time (in seconds) taken by the algorithm.

**Macro Variable \_OPTGRAPH\_SUMMARY\_**

The OPTGRAPH procedure defines a macro variable named `_OPTGRAPH_SUMMARY_`. This variable contains a character string that indicates the status and some basic statistics about the results of the algorithm that PROC OPTGRAPH uses to calculate summary statistics. The various terms of the variable are interpreted as follows:

**STATUS**

indicates the status of the algorithm at termination. The STATUS term takes the same value as the term **SUMMARY** in the `_OPTGRAPH_` macro as defined in the section “[Macro Variable \\_OPTGRAPH\\_](#)” on page 172.

**CPU\_TIME**

indicates the CPU time (in seconds) taken by the algorithm.

**REAL\_TIME**

indicates the real time (in seconds) taken by the algorithm.



## Macro Variable `_OPTGRAPH_TRANSCL_`

The OPTGRAPH procedure defines a macro variable named `_OPTGRAPH_TRANSCL_`. This variable contains a character string that indicates the status and some basic statistics about the results of the algorithm that PROC OPTGRAPH uses to calculate transitive closure. The various terms of the variable are interpreted as follows:

### STATUS

indicates the status of the algorithm at termination. The STATUS term takes the same value as the term `TRANSCL` in the `_OPTGRAPH_` macro as defined in the section “[Macro Variable `\_OPTGRAPH\_`](#)” on page 172.

### CPU\_TIME

indicates the CPU time (in seconds) taken by the algorithm.

### REAL\_TIME

indicates the real time (in seconds) taken by the algorithm.

## Macro Variable `_OPTGRAPH_TSP_`

The OPTGRAPH procedure defines a macro variable named `_OPTGRAPH_TSP_`. This variable contains a character string that indicates the status and some basic statistics about the results of the algorithm that PROC OPTGRAPH uses to solve the traveling salesman problem. The various terms of the variable are interpreted as follows:

### STATUS

indicates the status of the solver at termination. The STATUS term takes the same value as the term `TSP` in the `_OPTGRAPH_` macro as defined in the section “[Macro Variable `\_OPTGRAPH\_`](#)” on page 172.

### OBJECTIVE

indicates the objective value that the solver obtains at termination.

### RELATIVE\_GAP

specifies the relative gap between the best integer objective (`BestInteger`) and the objective of the best remaining node (`BestBound`) upon termination of the solver. The relative gap is equal to

$$| \text{BestInteger} - \text{BestBound} | / (1\text{E}-10 + | \text{BestBound} |)$$

### ABSOLUTE\_GAP

specifies the absolute gap between the best integer objective (`BestInteger`) and the objective of the best remaining node (`BestBound`) upon termination of the solver. The absolute gap is equal to

$$| \text{BestInteger} - \text{BestBound} |$$

### PRIMAL\_INFEASIBILITY

indicates the maximum (absolute) violation of the primal constraints by the solution.

**BOUND\_INFEASIBILITY**

indicates the maximum (absolute) violation by the solution of the lower or upper bounds (or both).

**INTEGER\_INFEASIBILITY**

indicates the maximum (absolute) violation of the integrality of integer variables that are returned by the solver.

**BEST\_BOUND**

specifies the best linear programming objective value of all unprocessed nodes in the branch-and-bound tree at the end of execution. A missing value indicates that the solver has processed either all or none of the nodes in the branch-and-bound tree.

**NODES**

specifies the number of nodes enumerated by the solver by using the branch-and-bound algorithm.

**ITERATIONS**

indicates the number of simplex iterations taken to solve the problem.

**CPU\_TIME**

indicates the CPU time (in seconds) taken by the algorithm.

**REAL\_TIME**

indicates the real time (in seconds) taken by the algorithm.

**NOTE:** The time reported in PRESOLVE\_TIME and SOLUTION\_TIME is either CPU time (default) or real time. The type is determined by the **TIMETYPE=** option.

---

## ODS Table Names

Each table that the OPTGRAPH procedure creates has a name associated with it, and you must use this name to refer to the table when you use ODS statements. These names are listed in [Table 1.59](#).

**Table 1.59** ODS Tables Produced by PROC OPTGRAPH

Table Name	Description	Required Statement or Option
PerformanceInfo	Information about the computing environment	Default output
ProblemSummary	Summary of the graph (or matrix) input	Default output
SolutionSummary	For each algorithm, summary of the solution status	Default output
Timing	Detailed real times for each phase of the procedure	<b>PERFORMANCE</b> with DETAILS option

The following code uses the example in the section “[Traveling Salesman Problem Applied to a Simple Undirected Graph](#)” on page 165 and calculates both an optimal traveling salesman tour and a minimum spanning tree. This code produces all four ODS output tables listed in [Table 1.59](#).

```
data LinkSetIn;
  input from $ to $ weight @@;
  datalines;
A B 1.0 A C 1.0 A D 1.5 B C 2.0 B D 4.0
B E 3.0 C D 3.0 C F 3.0 C H 4.0 D E 1.5
D F 3.0 D G 4.0 E F 1.0 E G 1.0 F G 2.0
F H 4.0 H I 3.0 I J 1.0 C J 5.0 F J 3.0
F I 1.0 H J 1.0
;

proc optgraph
  loglevel      = moderate
  data_links    = LinkSetIn
  out_nodes     = NodeSetOut;
  mst
    out         = MST;
  tsp
    out         = TSP;
  performance details;
run;
%put &_OPTGRAPH_;
%put &_OPTGRAPH_TSP_;
%put &_OPTGRAPH_MST_;
```

The performance information table in [Figure 1.135](#) provides information about the computing environment. The value for *Number of Threads* is the maximum number of threads that are used in processing.

**Figure 1.135** Performance Information Table

**The OPTGRAPH Procedure**

Performance Information	
Execution Mode	Single-Machine
Number of Threads	4

The procedure task timing table in [Figure 1.136](#) provides a breakdown of each task and the amount of real time spent in processing.

**Figure 1.136** Task Timing Table

Procedure Task Timing		
Task	Time (sec.)	Time
Input	0.00	6.78%
Setup	0.05	86.44%
Minimum Spanning Tree	0.00	1.70%
Traveling Salesman Problem	0.00	3.39%
Output	0.00	1.69%

The problem summary table in [Figure 1.137](#) provides a basic summary of the graph (or matrix) input.

**Figure 1.137** Problem Summary Table

Problem Summary	
Input Type	Graph
Number of Nodes	10
Number of Links	22
Graph Direction	Undirected

The solution summary tables in [Figure 1.138](#) and [Figure 1.139](#) provide a basic solution summary for each algorithm that is processed. The information in these tables matches the information that is provided in the macro variables for each algorithm, described in the section “[Macro Variables](#)” on page 172.

**Figure 1.138** Solution Summary Table for MST

Solution Summary	
Problem Type	Minimum Spanning Tree
Solution Status	Optimal
Objective Value	10
CPU Time	0.00
Real Time	0.00

**Figure 1.139** Solution Summary Table for TSP

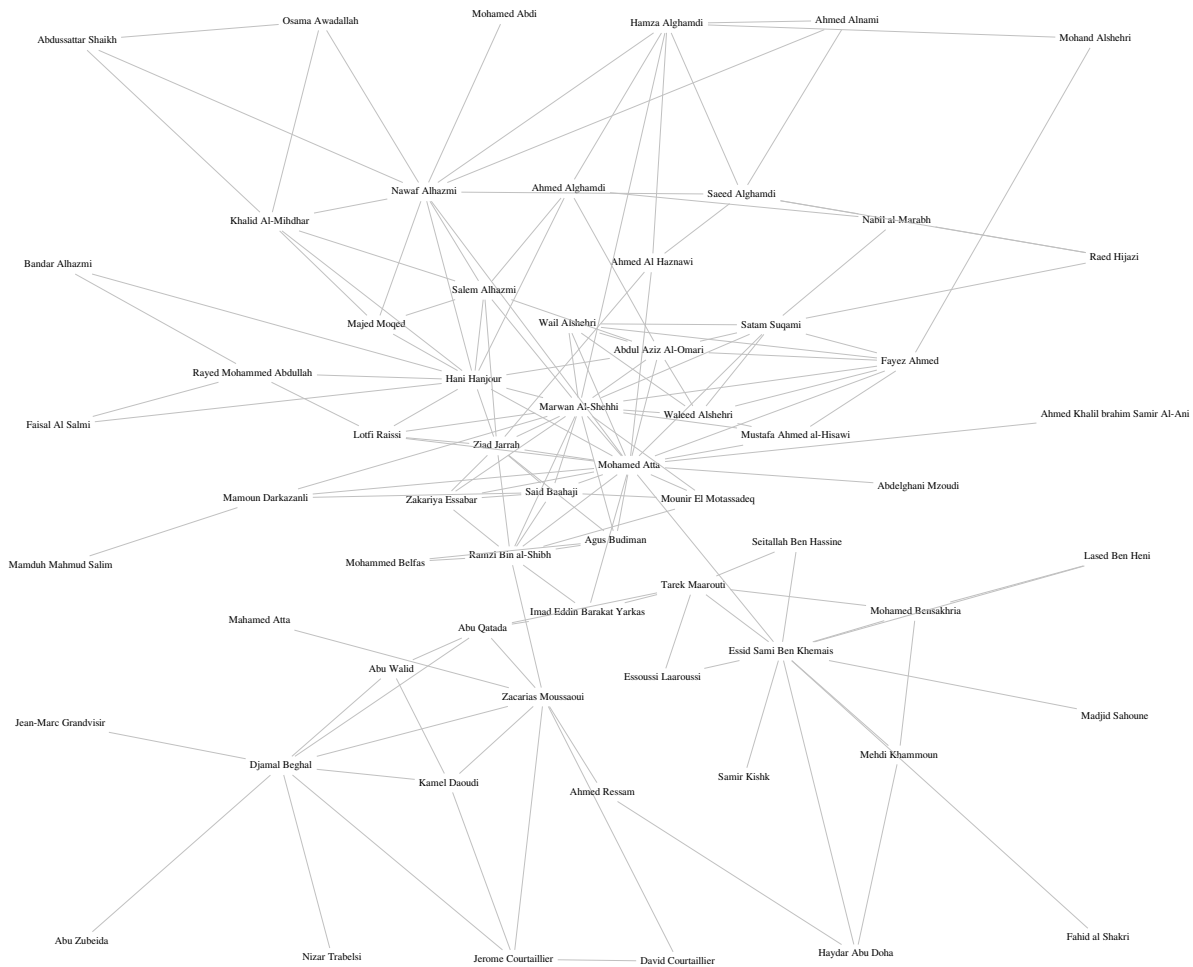
Solution Summary	
Problem Type	Traveling Salesman Problem
Solution Status	Optimal
Objective Value	16
Relative Gap	0
Absolute Gap	0
Primal Infeasibility	0
Bound Infeasibility	0
Integer Infeasibility	0
Best Bound	16
Nodes	1
Iterations	17
CPU Time	0.00
Real Time	0.00

## Examples: OPTGRAPH Procedure

### Example 1.1: Articulation Points in a Terrorist Network

This example considers the terrorist communications network from the attacks on the United States on September 11, 2001, described in Krebs 2002. Figure 1.140 shows this network, which was constructed after the attacks, based on collected intelligence information.

**Figure 1.140** Terrorist Communications Network from 9/11



The full network data include 153 links. The following statements show a small subset to illustrate the use of the BICONCOMP statement in this context:

```

data LinkSetInTerror911;
    input from & $32. to & $32.;
    datalines;
Abu Zubeida                Djamal Beghal
Jean-Marc Grandvisir       Djamal Beghal
Nizar Trabelsi             Djamal Beghal
Abu Walid                  Djamal Beghal
Abu Qatada                 Djamal Beghal
Zacarias Moussaoui         Djamal Beghal
Jerome Courtaillier        Djamal Beghal
Kamel Daoudi               Djamal Beghal
Abu Walid                  Kamel Daoudi
Abu Walid                  Abu Qatada
Kamel Daoudi               Zacarias Moussaoui
Kamel Daoudi               Jerome Courtaillier
Jerome Courtaillier        Zacarias Moussaoui
Jerome Courtaillier        David Courtaillier
Zacarias Moussaoui         David Courtaillier
Zacarias Moussaoui         Ahmed Ressam
Zacarias Moussaoui         Abu Qatada
Zacarias Moussaoui         Ramzi Bin al-Shibh
Zacarias Moussaoui         Mahamed Atta
Ahmed Ressam               Haydar Abu Doha
Mehdi Khammoun             Haydar Abu Doha
Essid Sami Ben Khemais     Haydar Abu Doha
Mehdi Khammoun             Essid Sami Ben Khemais
Mehdi Khammoun             Mohamed Bensakhria
...
;

```

Suppose that this communications network had been discovered before the attack on 9/11. If the investigators' goal was to disrupt the flow of communication between different groups within the organization, then they would want to focus on the people who are articulation points in the network.

To find the articulation points, use the following statements:

```

proc optgraph
    data_links = LinkSetInTerror911
    out_nodes = NodeSetOut;
    biconcomp;
run;

data ArtPoints;
    set NodeSetOut;
    where artpoint=1;
run;

```

The data set ArtPoints contains members of the network who are articulation points. Focusing investigations on cutting off these particular members could have caused a great deal of disruption in the terrorists' ability to communicate when formulating the attack.

**Output 1.1.1** Articulation Points of Terrorist Communications Network from 9/11

node	artpoint
Djamal Beghal	1
Zacarias Moussaoui	1
Essid Sami Ben Khemais	1
Mohamed Atta	1
Mamoun Darkazanli	1
Nawaf Alhazmi	1

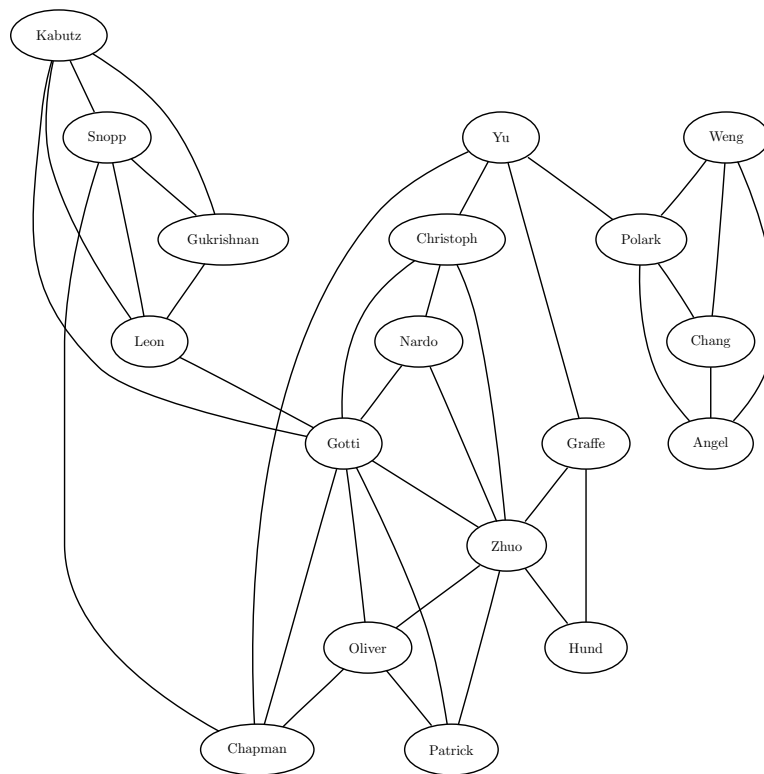
---

## Example 1.2: Influence Centrality for Project Groups in a Research Department

This example looks at an undirected graph that represents a few of the project groups in a hypothetical research department. A link between nodes A and B means that person A and B work together or that person A reports to person B. This graph represents six main projects.

- Department 1 (D1) consists of Snopp, Gukrishnan, Leon, and Kabutz. Snopp reports to Chapman.
- Department 2 (D2) consists of Oliver, Gotti, Patrick, and Zhuo. Oliver reports to Chapman.
- Department 3 (D3) consists of Gotti, Leon, and Kabutz. Gotti reports to Chapman.
- Department 4 (D4) consists of the following project teams who report to Yu. Yu reports to Chapman on this project.
  - Department 4a (D4a) consists of Polark, Chang, Weng, and Angel. Polark reports to Yu.
  - Department 4b (D4b) consists of Christoph, Nardo, Gotti, and Zhuo. Christoph reports to Yu.
  - Department 4c (D4c) consists of Graffe, Zhuo, and Hund. Graffe reports to Yu.

The links are shown in [Figure 1.141](#).

**Figure 1.141** Project Groups in a Research Department

The link weights measure the reporting magnitude. In general the higher the weight, the higher the contribution to the influence metric. Chapman is the director of the overall department, and Yu is the manager of a subgroup. The leads for the projects D1, D2, and D3 report to Chapman, and the leads for D4a, D4b, and D4c report to Yu. Reporting links to the director, Chapman, are given a link weight of 3, and reporting links to Yu are given a weight of 2. Links that represent people working together on a project all receive equal weight of 1. The node weights also represent some level of reporting: directors (4), managers (3), leads (2), and all others (1).

The project graph can be represented in the following link and node data sets:

```

data LinkSetInDept;
  input from $1-12 to $13-24 weight;
  datalines;
Yu      Chapman      3
Gotti   Chapman      3
Oliver  Chapman      3
Snopp   Chapman      3
Gukrishnan Leon      1
Snopp   Gukrishnan   1
Kabutz  Gukrishnan   1
Kabutz  Snopp         1
Snopp   Leon         1
Kabutz  Leon         1
Gotti   Oliver        1
Gotti   Patrick       1
Oliver  Patrick       1
  
```



```

Zhuo      Oliver      1
Zhuo      Gotti       1
Zhuo      Patrick     1
Kabutz    Gotti       1
Leon     Gotti       1
Polark    Yu          2
Polark    Chang       1
Chang     Angel       1
Polark    Angel       1
Weng      Polark      1
Weng      Chang       1
Weng      Angel       1
Christoph Yu          2
Christoph Nardo       1
Christoph Gotti       1
Christoph Zhuo        1
Nardo     Gotti       1
Nardo     Zhuo        1
Graffe    Yu          2
Graffe    Hund        1
Graffe    Zhuo        1
Zhuo      Hund        1
;

data NodeSetInDept;
    input node $1-12 weight;
    datalines;
Chapman    4
Yu         3
Gotti      2
Polark     2
Christoph  2
Oliver     2
Snopp     2
Zhuo       1
Nardo      1
Weng       1
Chang      1
Hund       1
Graffe     1
Leon      1
Gukrishnan 1
Kabutz     1
Patrick    1
Angel      1
;

```

The following statements calculate influence centrality (in addition to degree centrality):

```

proc optgraph
    loglevel      = moderate
    data_links    = LinkSetInDept
    data_nodes    = NodeSetInDept
    out_nodes     = NodeSetOut;

```

```

centrality
  degree    = out
  influence = weight;
run;
%put &_OPTGRAPH_;
%put &_OPTGRAPH_CENTRALITY_;

```

The progress of the procedure is shown in [Output 1.2.1](#).

### Output 1.2.1 PROC OPTGRAPH Log: Influence Centrality for Project Groups in a Research Department

---

```

NOTE: -----
NOTE: -----
NOTE: Running OPTGRAPH version 14.1.
NOTE: -----
NOTE: -----
NOTE: The OPTGRAPH procedure is executing in single-machine mode.
NOTE: -----
NOTE: -----
NOTE: Reading the nodes data set.
NOTE: There were 18 observations read from the data set WORK.NODESETINDEPT.
NOTE: Reading the links data set.
NOTE: There were 35 observations read from the data set WORK.LINKSETINDEPT.
NOTE: Data input used 0.00 (cpu: 0.02) seconds.
NOTE: Building the input graph storage used 0.00 (cpu: 0.00) seconds.
NOTE: The input graph storage is using 0.0 MBs (peak: 0.0 MBs) of memory.
NOTE: The number of nodes in the input graph is 18.
NOTE: The number of links in the input graph is 35.
NOTE: -----
NOTE: -----
NOTE: Processing centrality metrics.
NOTE: -----
NOTE: Processing degree centrality metrics.
NOTE: Processing centrality metrics used 0.0 MBs of memory.
NOTE: Processing degree centrality metrics used 0.00 (cpu: 0.00) seconds.
NOTE: -----
NOTE: Processing influence centrality metrics.
NOTE: Processing centrality metrics used 0.0 MBs of memory.
NOTE: Processing influence centrality metrics used 0.00 (cpu: 0.00) seconds.
NOTE: -----
NOTE: Processing centrality metrics used 0.00 (cpu: 0.00) seconds.
NOTE: -----
NOTE: -----
NOTE: Creating nodes data set output.
NOTE: Data output used 0.00 (cpu: 0.00) seconds.
NOTE: -----
NOTE: -----
NOTE: The data set WORK.NODESETOUT has 18 observations and 5 variables.
STATUS=OK  CENTRALITY=OK
STATUS=OK  CPU_TIME=0.00  REAL_TIME=0.00

```

---

The node data set NodeSetOut now contains the weighted influence centrality of the department's graph, including  $C_1$  (variable centr\_influence1\_wt) and  $C_2$  (variable centr\_influence2\_wt). This data set is shown in [Output 1.2.2](#).

**Output 1.2.2** Influence Centrality for Project Groups in a Research Department

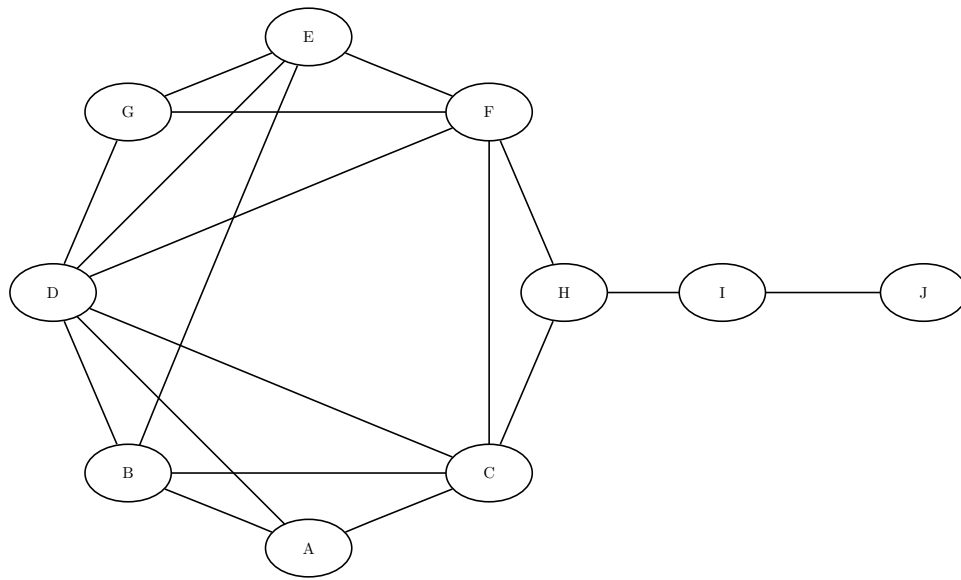
node	weight	centr_degree_out	centr_influence1_wt	centr_influence2_wt
Gotti	2	8	0.35714	1.57143
Zhuo	1	7	0.25000	1.17857
Oliver	2	4	0.21429	1.14286
Chapman	4	4	0.42857	1.10714
Christoph	2	4	0.17857	1.03571
Yu	3	4	0.32143	0.92857
Snopp	2	4	0.21429	0.82143
Leon	1	4	0.14286	0.82143
Patrick	1	3	0.10714	0.82143
Kabutz	1	4	0.14286	0.82143
Nardo	1	3	0.10714	0.78571
Polark	2	4	0.17857	0.64286
Graffe	1	3	0.14286	0.64286
Gukrishnan	1	3	0.10714	0.50000
Weng	1	3	0.10714	0.39286
Chang	1	3	0.10714	0.39286
Hund	1	2	0.07143	0.39286
Angel	1	3	0.10714	0.39286

As expected, the director Chapman has the highest first-order influence, since the weights of the reporting links to him are high. The highest second-order influence is Gotti, who reports to the director but is also involved in three different projects and therefore has a large sphere of influence. This example is revisited with other centrality metrics in other examples.

---

### Example 1.3: Betweenness and Closeness Centrality for Computer Network Topology

Consider a small network of 10 computers spread out across an office. Let a node represent a computer, and let a link represent a direct connection between the machines. For this example, consider the links as Ethernet connections that enable data to transfer between computers. If two computers are not connected directly, then the information must flow through other connected machines. Consider a topology as shown in [Figure 1.142](#). This is an example of the well-known *kite network*, which was popularized by David Krackhardt (1990) for better understanding of social networks in the workplace.

**Figure 1.142** Office Computer Network

Define the link data set as follows:

```
data LinkSetInCompNet;
  input from $ to $ @@;
  datalines;
A B  A C  A D  B C  B D
B E  C D  C F  C H  D E
D F  D G  E F  E G  F G
F H  H I  I J
;
```

To better understand the topology of the computer network, calculate the degree, closeness, and betweenness centrality. It is also interesting to look for articulation points in the computer network to identify places of vulnerability. All of these calculations can be done in one call to PROC OPTGRAPH as follows:

```
proc optgraph
  data_links = LinkSetInCompNet
  out_links  = LinkSetOut
  out_nodes  = NodeSetOut;
  centrality
    degree = out
    close  = unweight
    between = unweight;
  biconcomp;
run;
```

Output 1.3.1 shows the resulting node data set NodeSetOut sorted by closeness.

**Output 1.3.1** Node Closeness and Betweenness Centrality, Sorted by Closeness

node	centr_degree_out	centr_close_unwt	centr_between_unwt	artpoint
C	5	0.64286	0.23148	0
F	5	0.64286	0.23148	0
D	6	0.60000	0.10185	0
H	3	0.60000	0.38889	1
B	4	0.52941	0.02315	0
E	4	0.52941	0.02315	0
A	3	0.50000	0.00000	0
G	3	0.50000	0.00000	0
I	2	0.42857	0.22222	1
J	1	0.31034	0.00000	0

Output 1.3.2 shows the resulting node (NodeSetOut) and link data sets (LinkSetOut) sorted by betweenness.

**Output 1.3.2** Node Closeness and Betweenness Centrality, Sorted by Betweenness

Obs	node	centr_degree_out	centr_close_unwt	centr_between_unwt	artpoint
1	H	3	0.60000	0.38889	1
2	C	5	0.64286	0.23148	0
3	F	5	0.64286	0.23148	0
4	I	2	0.42857	0.22222	1
5	D	6	0.60000	0.10185	0
6	E	4	0.52941	0.02315	0
7	B	4	0.52941	0.02315	0
8	A	3	0.50000	0.00000	0
9	G	3	0.50000	0.00000	0
10	J	1	0.31034	0.00000	0

Obs	from	to	biconcomp	centr_between_unwt
1	H	I	2	0.44444
2	C	H	3	0.29167
3	F	H	3	0.29167
4	I	J	1	0.25000
5	E	F	3	0.12963
6	B	C	3	0.12963
7	A	C	3	0.12500
8	F	G	3	0.12500
9	C	D	3	0.09259
10	D	F	3	0.09259
11	A	D	3	0.08333
12	D	G	3	0.08333
13	C	F	3	0.07407
14	B	E	3	0.07407
15	B	D	3	0.05093
16	D	E	3	0.05093
17	A	B	3	0.04167
18	E	G	3	0.04167

The computers with the highest closeness centrality are *C* and *F*, because they have the shortest paths to all other nodes. These computers are key to the efficient distribution of information across the network. Assuming that the entire office has some centralized data that should be shared with all computers, machines *C* and *F* would be the best candidates for storing the data on their local hard drives. The computer with the highest betweenness centrality is *H*. Although machine *H* has only three connections, it is one of the most important machines in the office because it serves as the only way to reach computers *I* and *J* from the other machines in the office. Notice also that machine *H* is an articulation point because removing it would disconnect the office network. In this setting, computers with high betweenness should be carefully maintained and secured with UPS (uninterruptible power supply) systems to ensure they are always online.

---

### Example 1.4: Betweenness and Closeness Centrality for Project Groups in a Research Department

This example uses the same data as are used in the section “[Example 1.2: Influence Centrality for Project Groups in a Research Department](#)” on page 189, which illustrates influence centrality by considering the link weights that represent some measure of reporting magnitude. In “[Example 1.2: Influence Centrality for Project Groups in a Research Department](#)” on page 189, links between managers (or leads) and direct reports had higher link weights than links between non-managers. This interpretation makes sense in the context of influence centrality because weight and the metric are directly related. However, weight and the metric are inversely related for closeness and betweenness centrality.

This example considers the speed of the flow of information between people. In this sense, connections between managers and direct reports have *smaller values*, which cost less in the shortest path calculations. The following DATA step produces a new links data set, based on LinkSetInDept, which uses the inverse of the weight:

```
data LinkSetInDeptInv;
    set LinkSetInDept;
    weight = 1 / weight;
run;
```

The following statements calculate weighted (and unweighted) closeness and betweenness centrality. Notice that this example also uses the NTHREADS= option in the **PERFORMANCE** statement to specify two threads to allow the computation to be run in parallel. Since these data have 18 nodes, each thread can process 9 nodes simultaneously.

```
proc optgraph
    loglevel      = moderate
    data_links    = LinkSetInDeptInv
    out_links     = LinkSetOut
    out_nodes     = NodeSetOut;
    performance
        nthreads  = 2;
    centrality
        close     = both
        between   = both;
run;
%put &_OPTGRAPH_;
%put &_OPTGRAPH_CENTRALITY_;
```

The progress of the procedure is shown in [Output 1.4.1](#).

**Output 1.4.1** PROC OPTGRAPH Log: Closeness and Node Betweenness Centrality for Project Groups in a Research Department

```

NOTE: -----
NOTE: -----
NOTE: Running OPTGRAPH version 14.1.
NOTE: -----
NOTE: -----
NOTE: The OPTGRAPH procedure is executing in single-machine mode.
NOTE: -----
NOTE: -----
NOTE: Reading the links data set.
NOTE: There were 35 observations read from the data set WORK.LINKSETINDEPTINV.
NOTE: Data input used 0.00 (cpu: 0.00) seconds.
NOTE: Building the input graph storage used 0.00 (cpu: 0.00) seconds.
NOTE: The input graph storage is using 0.0 MBs (peak: 0.0 MBs) of memory.
NOTE: The number of nodes in the input graph is 18.
NOTE: The number of links in the input graph is 35.
NOTE: -----
NOTE: -----
NOTE: Processing centrality metrics.
NOTE: -----
NOTE: Processing weighted between/close centrality metrics using 2 threads.

      Algorithm      Nodes  Complete      Cpu      Real
      Algorithm      Nodes  Complete      Time      Time
      betwNL/close(wt)      18      100%      0.00      0.00
NOTE: Processing centrality metrics used 0.0 MBs of memory.
NOTE: Processing weighted between/close centrality metrics used 0.00 (cpu: 0.00) seconds.
NOTE: -----
NOTE: Processing unweighted between/close centrality metrics using 2 threads.

      Algorithm      Nodes  Complete      Cpu      Real
      Algorithm      Nodes  Complete      Time      Time
      betwNL/close(unwt)      18      100%      0.00      0.00
NOTE: Processing centrality metrics used 0.0 MBs of memory.
NOTE: Processing unweighted between/close centrality metrics used 0.00 (cpu: 0.00) seconds.
NOTE: -----
NOTE: Processing centrality metrics used 0.00 (cpu: 0.00) seconds.
NOTE: -----
NOTE: -----
NOTE: Creating nodes data set output.
NOTE: Creating links data set output.
NOTE: Data output used 0.00 (cpu: 0.00) seconds.
NOTE: -----
NOTE: -----
NOTE: The data set WORK.LINKSETOUT has 35 observations and 5 variables.
NOTE: The data set WORK.NODESETOUT has 18 observations and 5 variables.
STATUS=OK  CENTRALITY=OK
STATUS=OK  CPU_TIME=0.00  REAL_TIME=0.00

```

The node data set NodeSetOut shows the weighted and unweighted closeness and node betweenness centrality, as shown in [Output 1.4.2](#).

**Output 1.4.2** Closeness and Betweenness Centrality for Project Groups in a Research Department

node	centr_close_wt	centr_close_unwt	centr_between_wt	centr_between_unwt
Yu	0.87179	0.50000	0.50000	0.41262
Chapman	0.88696	0.50000	0.44118	0.23235
Gotti	0.81600	0.51515	0.20956	0.28444
Oliver	0.73913	0.44737	0.04044	0.02230
Snopp	0.75556	0.38636	0.16176	0.08088
Gukrishnan	0.46575	0.32692	0.00000	0.00000
Leon	0.50746	0.38636	0.00000	0.03885
Kabutz	0.50746	0.38636	0.00000	0.03885
Patrick	0.50000	0.37778	0.00000	0.00000
Zhuo	0.58286	0.47222	0.06618	0.15172
Polark	0.69388	0.38636	0.30882	0.30882
Chang	0.44156	0.29310	0.00000	0.00000
Angel	0.44156	0.29310	0.00000	0.00000
Weng	0.44156	0.29310	0.00000	0.00000
Christoph	0.68456	0.48571	0.05882	0.11275
Nardo	0.51777	0.42500	0.00000	0.00000
Graffe	0.67105	0.43590	0.08088	0.06642
Hund	0.45133	0.36957	0.00000	0.00000

The link data set LinkSetOut shows the weighted and unweighted link betweenness centrality, as shown in [Output 1.4.3](#).



**Output 1.4.3** Link Betweenness Centrality for Project Groups in a Research Department

from	to	weight	centr_between_wt	centr_between_unwt
Yu	Chapman	0.33333	0.39706	0.25576
Gotti	Chapman	0.33333	0.20221	0.09767
Oliver	Chapman	0.33333	0.14338	0.07623
Snopp	Chapman	0.33333	0.26471	0.16005
Gukrishnan	Leon	1.00000	0.00735	0.03431
Snopp	Gukrishnan	1.00000	0.11029	0.05637
Kabutz	Gukrishnan	1.00000	0.00735	0.03431
Kabutz	Snopp	1.00000	0.03676	0.03517
Snopp	Leon	1.00000	0.03676	0.03517
Kabutz	Leon	1.00000	0.00735	0.00735
Gotti	Oliver	1.00000	0.00000	0.03431
Gotti	Patrick	1.00000	0.05882	0.06066
Oliver	Patrick	1.00000	0.03676	0.02022
Zhuo	Oliver	1.00000	0.02574	0.03885
Zhuo	Gotti	1.00000	0.05515	0.10184
Zhuo	Patrick	1.00000	0.02941	0.04412
Kabutz	Gotti	1.00000	0.07353	0.12586
Leon	Gotti	1.00000	0.07353	0.12586
Polark	Yu	0.50000	0.41176	0.41176
Polark	Chang	1.00000	0.11029	0.11029
Chang	Angel	1.00000	0.00735	0.00735
Polark	Angel	1.00000	0.11029	0.11029
Weng	Polark	1.00000	0.11029	0.11029
Weng	Chang	1.00000	0.00735	0.00735
Weng	Angel	1.00000	0.00735	0.00735
Christoph	Yu	0.50000	0.13603	0.15870
Christoph	Nardo	1.00000	0.04779	0.04412
Christoph	Gotti	1.00000	0.02574	0.09620
Christoph	Zhuo	1.00000	0.03309	0.05147
Nardo	Gotti	1.00000	0.05515	0.05147
Nardo	Zhuo	1.00000	0.02206	0.02941
Graffe	Yu	0.50000	0.18015	0.12402
Graffe	Hund	1.00000	0.06985	0.04804
Graffe	Zhuo	1.00000	0.03676	0.08578
Zhuo	Hund	1.00000	0.05515	0.07696

Note that Chapman (director) and Yu (manager, reporting to Chapman) both have the highest weighted closeness centrality. However, Yu's weighted betweenness centrality is highest because he serves as more of a *gatekeeper* between his three groups (D4a, D4b, and D4c) and the rest of the department.

---

## Example 1.5: Eigenvector Centrality for Word Sense Disambiguation

In many languages, numerous words are polysemous (they carry more than one meaning). A common task in information retrieval is to assign the correct meaning to a polysemous word within a given context. Take the

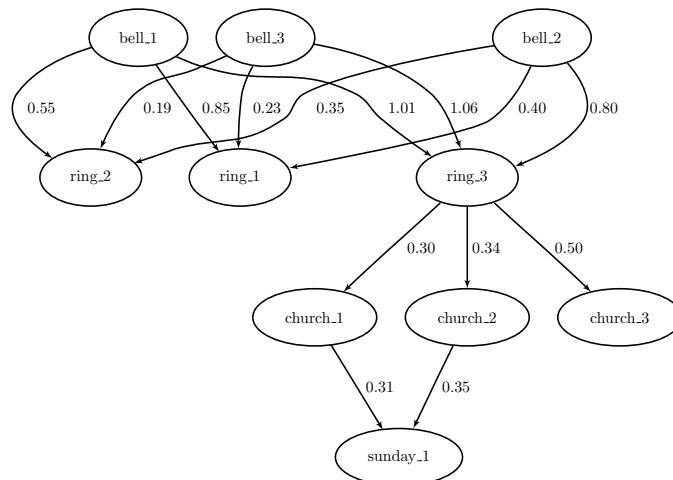
word “bass” as an example. It can mean either *a type of fish* (as in the sentence “I went fishing for some sea bass”) or *tones of low frequency* (as in the sentence “The bass part of the song is very moving”).

The following example from Mihalcea 2005 shows how eigenvector centrality can be used to disambiguate the word sense in the sentence “The church bells no longer ring on Sundays.” The following senses of words can be drawn from a dictionary:

- *church*
  1. one of the groups of Christians who have their own beliefs and forms of worship
  2. a place for public (especially Christian) worship
  3. a service conducted in a church
- *bell*
  1. a hollow device made of metal that makes a ringing sound when struck
  2. a push button at an outer door that gives a ringing or buzzing signal when pushed
  3. the sound of a bell
- *ring*
  1. make a ringing sound
  2. ring or echo with sound
  3. make (bells) ring, often for the purposes of musical edification
- *Sunday*
  1. first day of the week; observed as a day of rest and worship by most Christians

Using one of the similarity metrics defined in Sinha and Mihalcea 2007, you can generate a graph in which the nodes correspond to the word senses given above and the weights are determined by the similarity metric. The resulting graph is shown in Figure 1.143.

**Figure 1.143** Eigenvector Centrality for Word Sense Disambiguation



To identify the correct senses, you run eigenvector centrality on the graph and select the highest ranking sense for each word:

```
data LinkSetIn;
  input from $ to $ weight;
  datalines;
bell_1   ring_1   0.85
bell_1   ring_2   0.55
bell_1   ring_3   1.01
bell_2   ring_1   0.40
bell_2   ring_2   0.35
bell_2   ring_3   0.80
bell_3   ring_1   0.23
bell_3   ring_2   0.19
bell_3   ring_3   1.06
ring_3   church_1 0.30
ring_3   church_2 0.34
ring_3   church_3 0.50
church_1 sunday_1 0.31
church_2 sunday_1 0.35
;

proc optgraph
  data_links = LinkSetIn
  out_nodes  = NodeSetOut;
  centrality
    eigen    = weight;
run;

data NodeSetOut;
  length word $8 sense $1;
  set NodeSetOut;
  word  = scan(node,1,'_');
  sense = scan(node,2,'_');
run;

proc sort
  data = NodeSetOut
  out  = WordSenses;
  by word descending centr_eigen_wt;
run;

data WordSenses;
  set WordSenses(drop=centr_eigen_wt);
  by word;
  if first.word then output;
run;
```

The eigenvector scores and the implied word sense are shown in [Output 1.5.1](#).

**Output 1.5.1** Eigenvector Centrality for Word Sense Disambiguation

node	centr_eigen_wt
ring_3	1.00000
bell_1	0.77997
bell_3	0.59692
bell_2	0.53889
ring_1	0.48924
ring_2	0.35207
church_3	0.24081
church_2	0.17248
church_1	0.15222
sunday_1	0.05180

word	sense	node
bell	1	bell_1
church	3	church_3
ring	3	ring_3
sunday	1	sunday_1

**Example 1.6: Centrality Metrics for Project Groups in a Research Department**

The following statements use the WEIGHT2= option, and the project groups in a research department as depicted in Figure 1.141 on page 190. The data set contains the original weight and its inverse, which is used in the calculations of closeness and betweenness.

```
data LinkSetInDept2;
  input from $1-12 to $13-24 weight weightInv;
  datalines;
Yu          Chapman      3  0.33
Gotti       Chapman      3  0.33
Oliver      Chapman      3  0.33
Snopp       Chapman      3  0.33
Gukrishnan  Leon         1  1
Snopp       Gukrishnan   1  1
Kabutz      Gukrishnan   1  1
Kabutz      Snopp        1  1
Snopp       Leon         1  1
Kabutz      Leon         1  1
Gotti       Oliver       1  1
Gotti       Patrick      1  1
Oliver      Patrick      1  1
Zhuo        Oliver       1  1
Zhuo        Gotti        1  1
Zhuo        Patrick      1  1
Kabutz      Gotti        1  1
Leon       Gotti        1  1
Polark      Yu           2  0.50
Polark      Chang        1  1
Chang       Angel        1  1
Polark      Angel        1  1
```

```

Weng      Polark      1  1
Weng      Chang      1  1
Weng      Angel      1  1
Christoph Yu         2  0.50
Christoph Nardo      1  1
Christoph Gotti      1  1
Christoph Zhuo       1  1
Nardo     Gotti      1  1
Nardo     Zhuo       1  1
Graffe    Yu         2  0.50
Graffe    Hund       1  1
Graffe    Zhuo       1  1
Zhuo      Hund       1  1
;

proc optgraph
  data_nodes   = NodeSetInDept
  data_links   = LinkSetInDept2
  out_nodes    = NodeSetOut;
  performance
    nthreads   = 2;
  centrality
    clustering_coef
    degree     = out
    influence  = weight
    close      = weight
    between    = weight
    eigen      = weight
    weight2    = weightInv;
run;

```

The node data set NodeSetOut now shows the resulting centrality metrics given both weight interpretations.

**Output 1.6.1** Centrality for Project Groups in a Research Department

node	weight	centr_degree_out	centr_eigen_wt	centr_close_wt	centr_between_wt
Chapman	4	4	1.00000	0.88959	0.44118
Yu	3	4	0.62475	0.87404	0.50000
Gotti	2	8	0.70480	0.81849	0.20956
Polark	2	4	0.18777	0.69530	0.30882
Christoph	2	4	0.34168	0.68521	0.05882
Oliver	2	4	0.58858	0.74203	0.04044
Snopp	2	4	0.49133	0.75859	0.16176
Zhuo	1	7	0.32567	0.58319	0.06618
Nardo	1	3	0.18983	0.51813	0.00000
Weng	1	3	0.03591	0.44213	0.00000
Chang	1	3	0.03591	0.44213	0.00000
Hund	1	2	0.07667	0.45153	0.00000
Graffe	1	3	0.22852	0.67220	0.08088
Leon	1	4	0.21239	0.50822	0.00000
Gukrishnan	1	3	0.12674	0.46690	0.00000
Kabutz	1	4	0.21239	0.50822	0.00000
Patrick	1	3	0.22398	0.50074	0.00000
Angel	1	3	0.03591	0.44213	0.00000

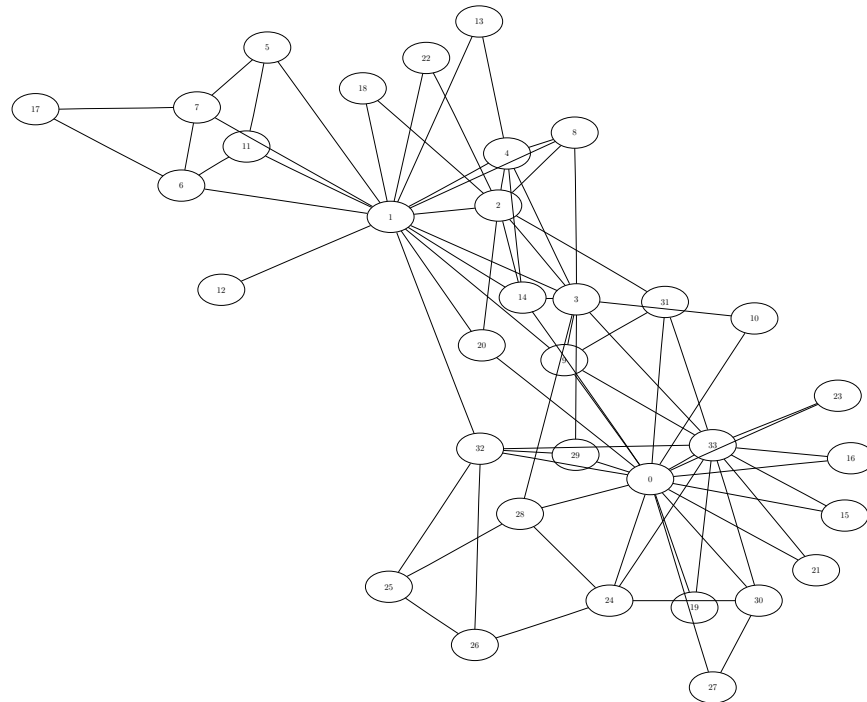
  

centr_influence1_wt	centr_influence2_wt	centr_cluster
0.42857	1.10714	0.16667
0.32143	0.92857	0.00000
0.35714	1.57143	0.28571
0.17857	0.64286	0.50000
0.17857	1.03571	0.50000
0.21429	1.14286	0.66667
0.21429	0.82143	0.50000
0.25000	1.17857	0.33333
0.10714	0.78571	1.00000
0.10714	0.39286	1.00000
0.10714	0.39286	1.00000
0.07143	0.39286	1.00000
0.14286	0.64286	0.33333
0.14286	0.82143	0.66667
0.10714	0.50000	1.00000
0.14286	0.82143	0.66667
0.10714	0.82143	1.00000
0.10714	0.39286	1.00000

## Example 1.7: Community Detection on Zachary's Karate Club Data

This example uses Zachary's Karate Club data (Zachary 1977), which describes social network friendships between 34 members of a karate club at a U.S. university in the 1970s. This is one of the standard publicly available data sets for testing community detection algorithms. It contains 34 nodes and 78 links. The graph is shown in Figure 1.144.

**Figure 1.144** Zachary's Karate Club Graph



The graph can be represented using the following links data set LinkSetIn:

```
data LinkSetIn;
  input from to weight @@;
  datalines;
0 9 1 0 10 1 0 14 1 0 15 1 0 16 1 0 19 1 0 20 1 0 21 1
0 23 1 0 24 1 0 27 1 0 28 1 0 29 1 0 30 1 0 31 1 0 32 1
0 33 1 2 1 1 3 1 1 3 2 1 4 1 1 4 2 1 4 3 1 5 1 1
6 1 1 7 1 1 7 5 1 7 6 1 8 1 1 8 2 1 8 3 1 8 4 1
9 1 1 9 3 1 10 3 1 11 1 1 11 5 1 11 6 1 12 1 1 13 1 1
13 4 1 14 1 1 14 2 1 14 3 1 14 4 1 17 6 1 17 7 1 18 1 1
18 2 1 20 1 1 20 2 1 22 1 1 22 2 1 26 24 1 26 25 1 28 3 1
28 24 1 28 25 1 29 3 1 30 24 1 30 27 1 31 2 1 31 9 1 32 1 1
32 25 1 32 26 1 32 29 1 33 3 1 33 9 1 33 15 1 33 16 1 33 19 1
33 21 1 33 23 1 33 24 1 33 30 1 33 31 1 33 32 1
;
```

The following statements use the RESOLUTION\_LIST= option to represent resolution levels (1, 0.5) in community detection on the Karate Club data. For more information about resolution levels, see the section “Resolution List” on page 92.

```
proc optgraph
  data_links      = LinkSetIn
  out_nodes       = NodeSetOut
  graph_internal_format = thin;
  community
    resolution_list = 1.0 0.5
    out_level       = CommLevelOut
    out_community   = CommOut
    out_overlap     = CommOverlapOut
    out_comm_links  = CommLinksOut;
run;
```

The data set NodeSetOut contains the community identifier of each node. It is shown in [Output 1.7.1](#).

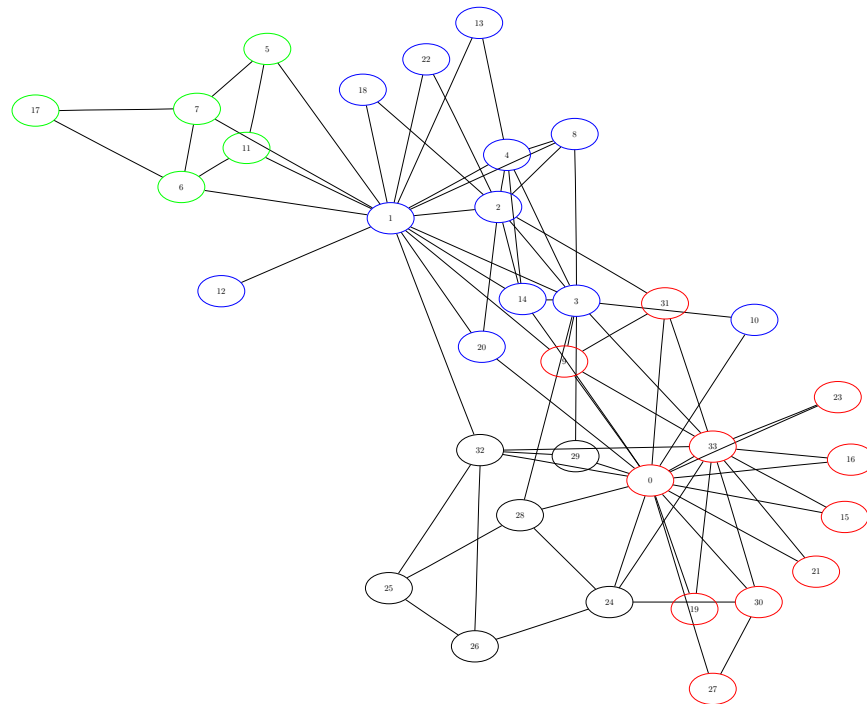
**Output 1.7.1** Community Nodes Output

node	community_1	community_2	node	community_1	community_2
0	0	0	33	0	0
9	0	0	2	1	1
10	1	1	1	1	1
14	1	1	3	1	1
15	0	0	4	1	1
16	0	0	5	3	1
19	0	0	6	3	1
20	1	1	7	3	1
21	0	0	8	1	1
23	0	0	11	3	1
24	2	0	12	1	1
27	0	0	13	1	1
28	2	0	17	3	1
29	2	0	18	1	1
30	0	0	22	1	1
31	0	0	26	2	0
32	2	0	25	2	0

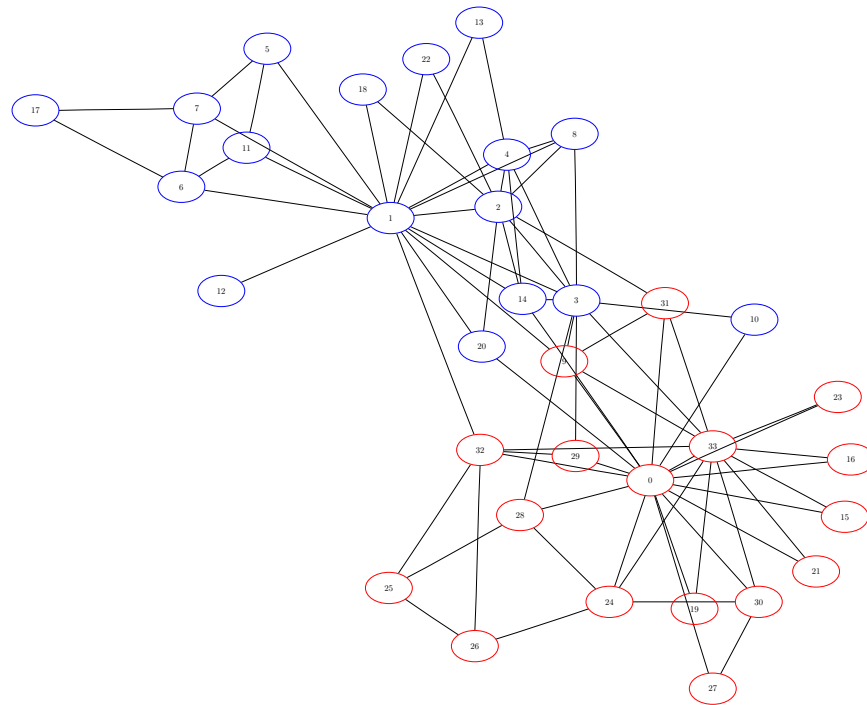
Column community\_1 contains the community identifier of each node when the resolution value is 1.0; column community\_2 contains the community identifier of each node when the resolution value is 0.5. Different node colors are used to represent different communities in [Figure 1.145](#) and [Figure 1.146](#). As you can see from the figures, four communities at resolution 1.0 are merged to two communities at resolution 0.5.



**Figure 1.145** Karate Club Communities (Resolution = 1.0)



**Figure 1.146** Karate Club Communities (Resolution = 0.5)



The data set `CommLevelOut` contains the number of communities and the corresponding modularity values found at each resolution level. It is shown in [Output 1.7.2](#).

**Output 1.7.2** Community Level Summary Output

level	resolution	communities	modularity
1	1.0	4	0.41880
2	0.5	2	0.37179

The data set CommOut contains the number of nodes contained in each community. It is shown in [Output 1.7.3](#).

**Output 1.7.3** Community Number of Nodes Output

level	resolution	community	nodes
1	1.0	0	11
1	1.0	1	12
1	1.0	2	6
1	1.0	3	5
2	0.5	0	17
2	0.5	1	17

The data set CommOverlapOut contains the intensity of each node that belongs to multiple communities. It is shown in [Output 1.7.4](#). Note that only the communities in the last resolution level (the smallest resolution value) are output in this data set. In this example, Node 0 belongs to two communities, with 82.3% of its links connecting to Community 0, and 17.6% of its links connecting to Community 1.

**Output 1.7.4** Community Overlap Output

node	community	intensity	node	community	intensity
0	0	0.82353	30	0	1.00000
0	1	0.17647	31	0	0.75000
9	0	0.60000	31	1	0.25000
9	1	0.40000	32	0	0.83333
10	0	0.50000	32	1	0.16667
10	1	0.50000	33	0	0.91667
14	0	0.20000	33	1	0.08333
14	1	0.80000	2	0	0.11111
15	0	1.00000	2	1	0.88889
16	0	1.00000	1	0	0.12500
19	0	1.00000	1	1	0.87500
20	0	0.33333	3	0	0.40000
20	1	0.66667	3	1	0.60000
21	0	1.00000	4	1	1.00000
23	0	1.00000	5	1	1.00000
24	0	1.00000	6	1	1.00000
27	0	1.00000	7	1	1.00000
28	0	0.75000	8	1	1.00000
28	1	0.25000	11	1	1.00000
29	0	0.66667	12	1	1.00000
29	1	0.33333	13	1	1.00000

**Output 1.7.4** *continued*

node	community	intensity
17	1	1.00000
18	1	1.00000
22	1	1.00000
26	0	1.00000
25	0	1.00000

The data set CommLinksOut shows how the communities are interconnected. It is shown in [Output 1.7.5](#). In this example, when the resolution value is 1, the link weight between Communities 0 and 1 is 7, and the link weight between Communities 1 and 2 is 4.

**Output 1.7.5** Community Links Output

level	resolution	from_community	to_community	link_weight
1	1.0	0	1	7
1	1.0	0	2	7
1	1.0	1	2	3
1	1.0	1	3	4
2	0.5	0	1	10

## Example 1.8: Recursive Community Detection on Zachary's Karate Club Data

This example illustrates the use of the RECURSIVE option in community detection on Zachary's Karate Club data. The data set appears in "[Example 1.7: Community Detection on Zachary's Karate Club Data](#)" on page 205. This example forces each community to contain no more than five nodes and the number of links between any pair of nodes within any community to be no greater than 2.

```
proc optgraph
  data_links          = LinkSetIn
  out_nodes           = NodeSetOut
  graph_internal_format = thin;
  community
    resolution_list    = 1.0
    recursive (max_comm_size = 5 max_diameter = 2 relation = AND)
    out_community      = CommOut;
run;
```

The data set NodeSetOut contains the community identifier of each node. It is shown in [Output 1.8.1](#).

**Output 1.8.1** Community Nodes Output

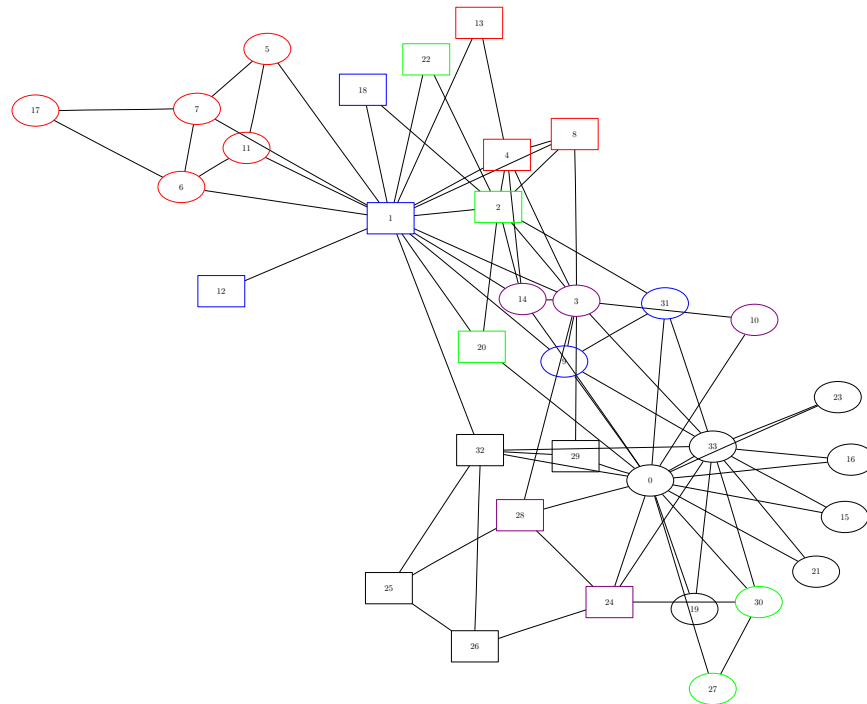
node	community_1	node	community_1
0	3	33	3
9	2	2	4
10	6	1	5
14	7	3	6
15	3	4	7
16	3	5	0
19	3	6	0
20	4	7	0
21	3	8	6
23	3	11	0
24	9	12	5
27	1	13	7
28	9	17	0
29	8	18	5
30	1	22	4
31	2	26	9
32	8	25	9

The data set CommOut contains the number of nodes contained in each community. It is shown in [Output 1.8.2](#).

**Output 1.8.2** Community Number of Nodes Output

level	resolution	community	nodes
1	1	0	5
1	1	1	2
1	1	2	2
1	1	3	7
1	1	4	3
1	1	5	3
1	1	6	3
1	1	7	3
1	1	8	2
1	1	9	4

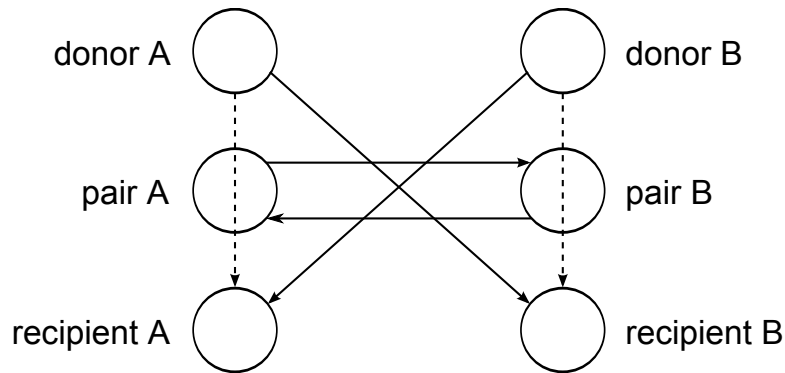
The community graph is shown in [Figure 1.147](#), with different node shapes and colors representing different communities.

**Figure 1.147** Karate Club Recursive Communities

As you can see from [Output 1.8.2](#), Community 3, whose nodes are drawn as black ellipses in [Figure 1.147](#), contains seven nodes even though the maximum number of nodes in any community is set to be 5. This is because Community 3 has a symmetric shape: Nodes 0 and 33 are in the center, and they symmetrically connect to Nodes 21, 15, 19, 16, and 23. Therefore, this community cannot be further split.

## Example 1.9: Cycle Detection for Kidney Donor Exchange

This example looks at an application of cycle detection to help create a kidney donor exchange. Suppose someone needs a kidney transplant and a family member is willing to donate one. If the donor and recipient are incompatible (because of blood types, tissue mismatch, and so on), the transplant cannot happen. Now suppose two donor-recipient pairs A and B are in this situation, but donor A is compatible with recipient B and donor B is compatible with recipient A. Then two transplants can take place in a two-way swap, shown graphically in [Figure 1.148](#). More generally, an  $n$ -way swap can be performed involving  $n$  donors and  $n$  recipients (Willingham 2009).

**Figure 1.148** Kidney Donor Exchange Two-Way Swap

To model this problem, define a directed graph as follows. Each node is an incompatible donor-recipient pair. Link  $(i, j)$  exists if the donor from node  $i$  is compatible with the recipient from node  $j$ . The link weight is a measure of the quality of the match. By introducing dummy links whose weight is 0, you can also include altruistic donors who have no recipients or recipients who have no donors. The idea is to find a maximum-weight node-disjoint union of directed cycles. You want the union to be node-disjoint so that no kidney is donated more than once, and you want cycles so that the donor from node  $i$  gives up a kidney if and only if the recipient from node  $i$  receives a kidney.

Without any other constraints, the problem could be solved as a linear assignment problem, as described in the section “[Linear Assignment \(Matching\)](#)” on page 115. But doing so would allow arbitrarily long cycles in the solution. Because of practical considerations (such as travel) and to mitigate risk, each cycle must have no more than  $L$  links. The kidney exchange problem is to find a maximum-weight node-disjoint union of short directed cycles.

One way to solve this problem is to explicitly generate all cycles whose length is at most  $L$  and then solve a set packing problem. You can use PROC OPTGRAPH to generate the cycles and then PROC OPTMODEL (see *SAS/OR User’s Guide: Mathematical Programming*) to read the PROC OPTGRAPH output, formulate the set packing problem, call the mixed integer linear programming solver, and output the optimal solution.

The following DATA step sets up the problem, first creating a random graph on  $n$  nodes with link probability  $p$  and Uniform(0,1) weight:

```
/* create random graph on n nodes with arc probability p
   and uniform(0,1) weight */
%let n = 100;
%let p = 0.02;
data LinkSetIn;
  do from = 0 to &n - 1;
    do to = 0 to &n - 1;
      if from eq to then continue;
      else if ranuni(1) < &p then do;
        weight = ranuni(2);
        output;
      end;
    end;
  end;
run;
```

The following statements use PROC OPTGRAPH to generate all cycles whose length is greater than or equal to 2 and less than or equal to 10:

```
/* generate all cycles with 2 <= length <= max_length */
%let max_length = 10;
proc optgraph
  loglevel          = moderate
  graph_direction   = directed
  data_links        = LinkSetIn;
  cycle
    minLength       = 2
    maxLength       = &max_length
    out              = Cycles
    mode             = all_cycles;
run;
%put &_OPTGRAPH_;
%put &_OPTGRAPH_CYCLE_;
```

PROC OPTGRAPH finds 224 cycles of the appropriate length, as shown in [Output 1.9.1](#).

#### Output 1.9.1 Cycles for Kidney Donor Exchange PROC OPTGRAPH Log

---

```
NOTE: -----
NOTE: -----
NOTE: Running OPTGRAPH version 14.1.
NOTE: -----
NOTE: -----
NOTE: The OPTGRAPH procedure is executing in single-machine mode.
NOTE: -----
NOTE: -----
NOTE: Reading the links data set.
NOTE: There were 194 observations read from the data set WORK.LINKSETIN.
NOTE: Data input used 0.00 (cpu: 0.00) seconds.
NOTE: Building the input graph storage used 0.00 (cpu: 0.00) seconds.
NOTE: The input graph storage is using 0.0 MBs (peak: 0.0 MBs) of memory.
NOTE: The number of nodes in the input graph is 97.
NOTE: The number of links in the input graph is 194.
NOTE: -----
NOTE: -----
NOTE: Processing cycle detection.
NOTE: The graph has 224 cycles.
NOTE: Processing cycle detection used 0.62 (cpu: 0.61) seconds.
NOTE: -----
NOTE: -----
NOTE: Creating cycle data set output.
NOTE: Data output used 0.00 (cpu: 0.00) seconds.
NOTE: -----
NOTE: -----
NOTE: The data set WORK.CYCLES has 2124 observations and 3 variables.
STATUS=OK  CYCLE=OK
STATUS=OK  NUM_CYCLES=224  CPU_TIME=0.61  REAL_TIME=0.62
```

---

From the resulting data set `Cycles`, use the following DATA step to convert the cycles into one observation per arc:

```
/* convert Cycles into one observation per arc */
data Cycles0(keep=c i j);
  set Cycles;
  retain last;
  c    = cycle;
  i    = last;
  j    = node;
  last = j;
  if order ne 1 then output;
run;
```

Given the set of cycles, you can now formulate a mixed integer linear program (MILP) to maximize the total cycle weight. Let  $C$  define the set of cycles of appropriate length,  $N_c$  define the set of nodes in cycle  $c$ ,  $A_c$  define the set of links in cycle  $c$ , and  $w_{ij}$  denote the link weight for link  $(i, j)$ . Define a binary decision variable  $x_c$ . Set  $x_c$  to 1 if cycle  $c$  is used in the solution; otherwise, set it to 0. Then, the following MILP defines the problem that you want to solve to maximize the quality of the kidney exchange:

$$\begin{aligned}
 &\text{maximize} && \sum_{c \in C} \left( \sum_{(i,j) \in A_c} w_{ij} \right) x_c \\
 &\text{subject to} && \sum_{c \in C: i \in N_c} x_c \leq 1 && i \in N && (\text{incomp\_pair}) \\
 &&& x_c \in \{0, 1\} && c \in C
 \end{aligned}$$

The constraint (incomp\_pair) ensures that each node (incompatible pair) in the graph is intersected at most once. That is, a donor can donate a kidney only once. You can use PROC OPTMODEL to solve this mixed integer linear programming problem as follows:

```
/* solve set packing problem to find maximum weight node-disjoint union
   of short directed cycles */
proc optmodel;
  /* declare index sets and parameters, and read data */
  set <num,num> ARCS;
  num weight {ARCS};
  read data LinkSetIn into ARCS=[from to] weight;
  set <num,num,num> TRIPLES;
  read data Cycles0 into TRIPLES=[c i j];
  set CYCLES = setof {<c,i,j> in TRIPLES} c;
  set ARCS_c {c in CYCLES} = setof {<(c),i,j> in TRIPLES} <i,j>;
  set NODES_c {c in CYCLES} = union {<i,j> in ARCS_c[c]} {i,j};
  set NODES = union {c in CYCLES} NODES_c[c];
  num cycle_weight {c in CYCLES} = sum {<i,j> in ARCS_c[c]} weight[i,j];

  /* UseCycle[c] = 1 if cycle c is used, 0 otherwise */
  var UseCycle {CYCLES} binary;

  /* declare objective */
  max TotalWeight
    = sum {c in CYCLES} cycle_weight[c] * UseCycle[c];
```



```

/* each node appears in at most one cycle */
con node_packing {i in NODES}:
    sum {c in CYCLES: i in NODES_c[c]} UseCycle[c] <= 1;

/* call solver */
solve with milp;

/* output optimal solution */
create data Solution from
    [c]={c in CYCLES: UseCycle[c].sol > 0.5} cycle_weight;
quit;
%put &_OROPTMODEL_;

```

PROC OPTMODEL solves the problem by using the mixed integer linear programming solver. As shown in [Output 1.9.2](#), it was able to find a total weight (quality level) of 26.02.

### Output 1.9.2 Cycles for Kidney Donor Exchange PROC OPTMODEL Log

---

```

NOTE: There were 194 observations read from the data set WORK.LINKSETIN.
NOTE: There were 1900 observations read from the data set WORK.CYCLES0.
NOTE: Problem generation will use 4 threads.
NOTE: The problem has 224 variables (0 free, 0 fixed).
NOTE: The problem has 224 binary and 0 integer variables.
NOTE: The problem has 63 linear constraints (63 LE, 0 EQ, 0 GE, 0 range).
NOTE: The problem has 1900 linear constraint coefficients.
NOTE: The problem has 0 nonlinear constraints (0 LE, 0 EQ, 0 GE, 0 range).
NOTE: The MILP presolver value AUTOMATIC is applied.
NOTE: The MILP presolver removed 46 variables and 35 constraints.
NOTE: The MILP presolver removed 901 constraint coefficients.
NOTE: The MILP presolver modified 0 constraint coefficients.
NOTE: The presolved problem has 178 variables, 28 constraints, and 999 constraint coefficients.
NOTE: The MILP solver is called.
NOTE: The parallel Branch and Cut algorithm is used.
NOTE: The Branch and Cut algorithm is using up to 4 threads.

```

Node	Active	Sols	BestInteger	BestBound	Gap	Time
0	1	3	22.7780692	868.9019355	97.38%	0
0	1	3	22.7780692	26.7803921	14.94%	0
0	1	4	23.2158345	26.1205372	11.12%	0
0	1	5	26.0202879	26.0202879	0.00%	0
0	0	5	26.0202879	26.0202879	0.00%	0

```

NOTE: The MILP solver added 31 cuts with 2815 cut coefficients at the root.
NOTE: Objective = 26.020287887.
NOTE: The data set WORK.SOLUTION has 6 observations and 2 variables.
STATUS=OK ALGORITHM=BAC SOLUTION_STATUS=OPTIMAL OBJECTIVE=26.020287887 RELATIVE_GAP=0
ABSOLUTE_GAP=0 PRIMAL_INFEASIBILITY=1.110223E-15 BOUND_INFEASIBILITY=0
INTEGER_INFEASIBILITY=4E-6 BEST_BOUND=26.020287887 NODES=1 ITERATIONS=117 PRESOLVE_TIME=0.01
SOLUTION_TIME=0.02

```

---

The data set Solution, shown in [Output 1.9.3](#), now contains the cycles that define the best exchange and their associated weight (quality).

**Output 1.9.3** Maximum Quality Solution for Kidney Donor Exchange

c cycle_weight	
12	5.84985
43	3.90015
71	5.44467
124	7.42574
222	2.28231
224	1.11757

**Example 1.10: Linear Assignment Problem for Minimizing Swim Times**

A swimming coach needs to assign male and female swimmers to each stroke of a medley relay team. The swimmers' best times for each stroke are stored in a SAS data set. The `LINEAR_ASSIGNMENT` statement evaluates the times and matches strokes and swimmers to minimize the total relay swim time.

The data are stored in matrix format, where the row identifier is the swimmer's name (variable name) and each swimming event is a column (variables: back, breast, fly, and free). In the following `DATA` step, the relay times are split into two categories, male and female:

```
data RelayTimes;
  input name $ sex $ back breast fly free;
  datalines;
Sue      F 35.1 36.7 28.3 36.1
Karen    F 34.6 32.6 26.9 26.2
Jan      F 31.3 33.9 27.1 31.2
Andrea   F 28.6 34.1 29.1 30.3
Carol    F 32.9 32.2 26.6 24.0
Ellen    F 27.8 32.5 27.8 27.0
Jim      M 26.3 27.6 23.5 22.4
Mike     M 29.0 24.0 27.9 25.4
Sam      M 27.2 33.8 25.2 24.1
Clayton  M 27.0 29.2 23.0 21.9
;

data RelayTimesF RelayTimesM;
  set RelayTimes;
  if      sex='F' then output RelayTimesF;
  else if sex='M' then output RelayTimesM;
run;
```

The following statements solve the linear assignment problem for both male and female relay teams:

```
proc optgraph
  data_matrix = RelayTimesF;
  linear_assignment
    out      = LinearAssignF
    id       = (name sex);
run;
%put &_OPTGRAPH_;
%put &_OPTGRAPH_LAP_;

proc optgraph
  data_matrix = RelayTimesM;
  linear_assignment
    out      = LinearAssignM
    id       = (name sex);
run;
%put &_OPTGRAPH_;
%put &_OPTGRAPH_LAP_;
```

The progress of the two PROC OPTGRAPH calls is shown in [Output 1.10.1](#) and [Output 1.10.2](#).

#### Output 1.10.1 PROC OPTGRAPH Log: Linear Assignment for Female Swim Times

---

```
NOTE: -----
NOTE: Running OPTGRAPH version 14.1.
NOTE: -----
NOTE: The OPTGRAPH procedure is executing in single-machine mode.
NOTE: -----
NOTE: The number of columns in the input matrix is 4.
NOTE: The number of rows in the input matrix is 6.
NOTE: The number of nonzeros in the input matrix is 24.
NOTE: Data input used 0.00 (cpu: 0.00) seconds.
NOTE: -----
NOTE: Processing the linear assignment problem.
NOTE: Objective = 111.5.
NOTE: Processing the linear assignment problem used 0.00 (cpu: 0.00) seconds.
NOTE: -----
NOTE: Data output used 0.00 (cpu: 0.00) seconds.
NOTE: -----
NOTE: The data set WORK.LINEARASSIGNF has 4 observations and 4 variables.
STATUS=OK  LAP=OPTIMAL
STATUS=OPTIMAL  OBJECTIVE=111.5  CPU_TIME=0.00  REAL_TIME=0.00
```

---

**Output 1.10.2** PROC OPTGRAPH Log: Linear Assignment for Male Swim Times

```

NOTE: -----
NOTE: Running OPTGRAPH version 14.1.
NOTE: -----
NOTE: The OPTGRAPH procedure is executing in single-machine mode.
NOTE: -----
NOTE: The number of columns in the input matrix is 4.
NOTE: The number of rows in the input matrix is 4.
NOTE: The number of nonzeros in the input matrix is 16.
NOTE: Data input used 0.00 (cpu: 0.00) seconds.
NOTE: -----
NOTE: Processing the linear assignment problem.
NOTE: Objective = 96.6.
NOTE: Processing the linear assignment problem used 0.00 (cpu: 0.00) seconds.
NOTE: -----
NOTE: Data output used 0.00 (cpu: 0.00) seconds.
NOTE: -----
NOTE: The data set WORK.LINEARASSIGNM has 4 observations and 4 variables.
STATUS=OK  LAP=OPTIMAL
STATUS=OPTIMAL  OBJECTIVE=96.6  CPU_TIME=0.00  REAL_TIME=0.00

```

The data sets LinearAssignF and LinearAssignM contain the optimal assignments. Note that in the case of the female data, there are more people (set  $S$ ) than there are strokes (set  $T$ ). Therefore, the solver allows for some members of  $S$  to remain unassigned.

**Output 1.10.3** Optimal Assignments for Best Female Swim Times

name	sex	assign	cost
Karen	F	breast	32.6
Jan	F	fly	27.1
Carol	F	free	24.0
Ellen	F	back	27.8
			<b>111.5</b>

**Output 1.10.4** Optimal Assignments for Best Male Swim Times

name	sex	assign	cost
Jim	M	free	22.4
Mike	M	breast	24.0
Sam	M	back	27.2
Clayton	M	fly	23.0
			<b>96.6</b>

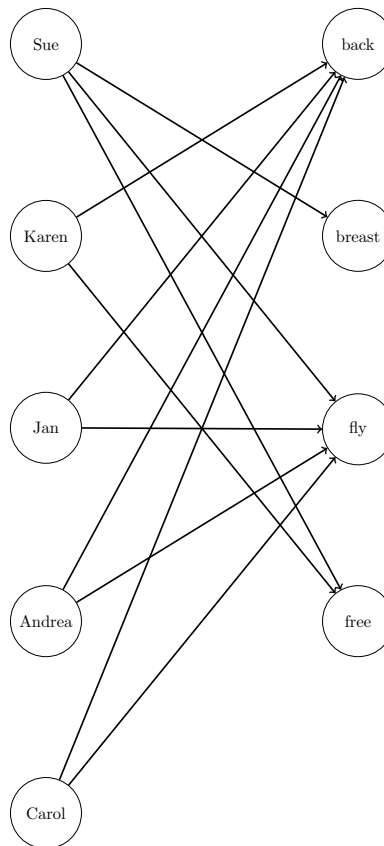
## Example 1.11: Linear Assignment Problem, Sparse Format versus Dense Format

This example looks at the problem of assigning swimmers to strokes based on their best times. However, in this case certain swimmers are not eligible to perform certain strokes. A missing (.) value in the data matrix identifies an ineligible assignment. For example:

```
data RelayTimesMatrix;
  input name $ sex $ back breast fly free;
  datalines;
Sue      F    .  36.7 28.3 36.1
Karen    F 34.6    .    . 26.2
Jan      F 31.3    . 27.1    .
Andrea   F 28.6    . 29.1    .
Carol    F 32.9    . 26.6    .
;
```

Recall that the linear assignment problem can also be interpreted as the minimum-weight matching in a bipartite graph. The eligible assignments define links between the rows (swimmers) and the columns (strokes), as in Figure 1.149.

**Figure 1.149** Bipartite Graph for Linear Assignment Problem



You can represent the same data in `RelayTimesMatrix` by using a links data set as follows:

```

data RelayTimesLinks;
  input name $ attr $ cost;
  datalines;
Sue      breast 36.7
Sue      fly    28.3
Sue      free   36.1
Karen    back   34.6
Karen    free   26.2
Jan      back   31.3
Jan      fly    27.1
Andrea   back   28.6
Andrea   fly    29.1
Carol    back   32.9
Carol    fly    26.6
;

```

This graph must be bipartite (such that  $S$  and  $T$  are disjoint). If it is not, PROC OPTGRAPH returns an error.

Now, you can use either input format to solve the same problem, as follows:

```

proc optgraph
  data_matrix = RelayTimesMatrix;
  linear_assignment
    out      = LinearAssignMatrix
    weight   = (back--free)
    id       = (name sex);
run;

proc optgraph
  graph_direction = directed
  data_links      = RelayTimesLinks;
  data_links_var
    from         = name
    to           = attr
    weight       = cost;
  linear_assignment
    out          = LinearAssignLinks;
run;

```

When you use the graph input format, the LINEAR\_ASSIGNMENT options WEIGHT= and ID= are not used directly.

The data sets LinearAssignMatrix and LinearAssignLinks now contain the optimal assignments, as shown in [Output 1.11.1](#) and [Output 1.11.2](#).

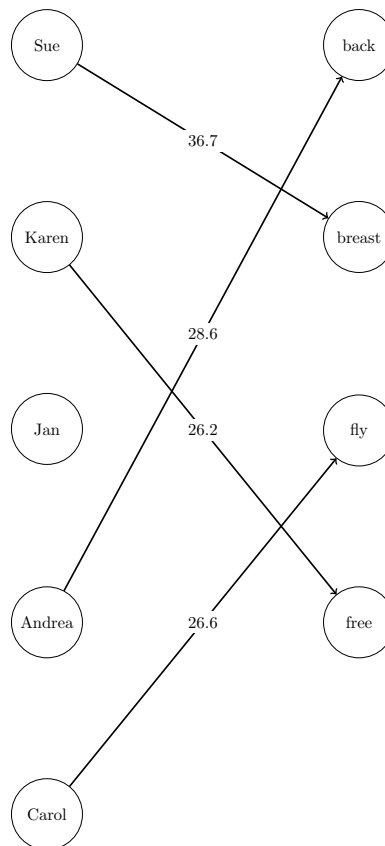
**Output 1.11.1** Optimal Assignments for Swim Times (Dense Input)

name	sex	assign	cost
Sue	F	breast	36.7
Karen	F	free	26.2
Andrea	F	back	28.6
Carol	F	fly	26.6
			<b>118.1</b>

**Output 1.11.2** Optimal Assignments for Swim Times (Sparse Input)

name	attr	cost
Sue	breast	36.7
Karen	free	26.2
Andrea	back	28.6
Carol	fly	26.6
		<b>118.1</b>

The optimal assignments are shown graphically in Figure 1.150.

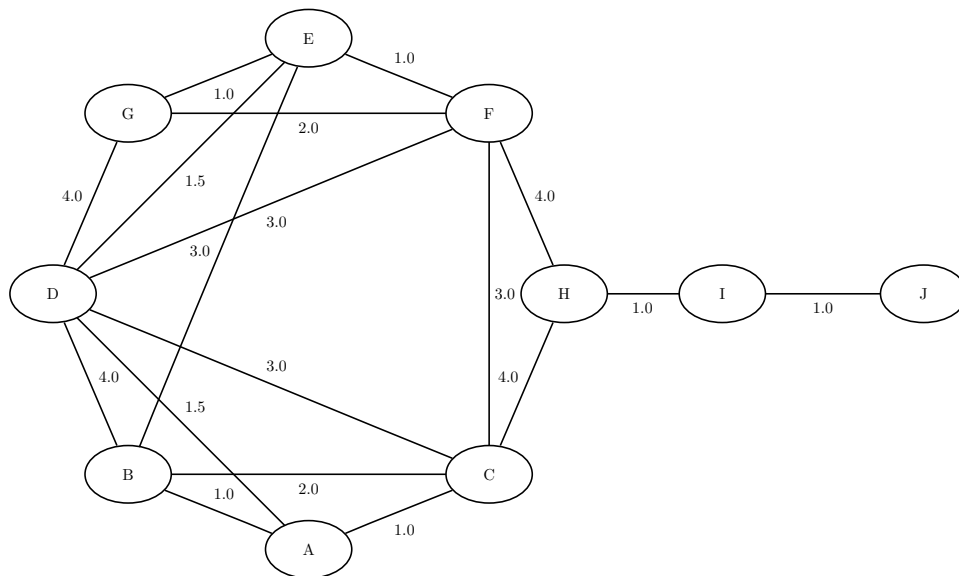
**Figure 1.150** Optimal Assignments for Swim Times

For large problems where a number of links are forbidden, the sparse format can be faster and can save a great deal of memory. Consider an example that uses the format of the DATA\_MATRIX= option with 15,000 columns ( $|S| = 15,000$ ) and 4,000 rows ( $|T| = 4,000$ ). To store the dense matrix in memory, PROC OPTGRAPH needs to allocate approximately  $|S| \cdot |T| \cdot 8/1024/1024 = 457$  MB. If the data have mostly ineligible links, then the sparse (graph) format that uses the DATA\_LINKS= option is much more efficient with respect to memory. For example, if the data have only 5% of the eligible links ( $15,000 \cdot 4,000 \cdot 0.05 = 3,000,000$ ), then the dense storage would still need 457 MB. The sparse storage for the same example needs approximately  $|S| \cdot |T| \cdot 0.05 \cdot 12/1024/1024 = 34$  MB. If the problem is fully dense (all links are eligible), then the dense format that uses the DATA\_MATRIX= option is the most efficient.

## Example 1.12: Minimum Spanning Tree for Computer Network Topology

Consider again the small network of computers described in the section “[Example 1.3: Betweenness and Closeness Centrality for Computer Network Topology](#)” on page 193. Suppose that this network has not yet been formed, but for structural reasons the connections between the machines shown in [Figure 1.142](#) are the only possible links. In designing the network, the goal is to make sure that each machine in the office can reach every other machine. To accomplish this goal, Ethernet lines must be constructed and run between the machines. The construction costs for each possible link are based approximately on distance and are shown in [Figure 1.151](#). Besides distance, the costs also reflect some restrictions due to physical boundaries. To connect all the machines in the office at minimal cost, you need to find a minimum spanning tree on the network of possible links.

**Figure 1.151** Potential Office Computer Network



Define the link data set as follows:

```
data LinkSetInCompNet;
  input from $ to $ weight @@;
  datalines;
A B 1.0  A C 1.0  A D 1.5  B C 2.0  B D 4.0
B E 3.0  C D 3.0  C F 3.0  C H 4.0  D E 1.5
D F 3.0  D G 4.0  E F 1.0  E G 1.0  F G 2.0
F H 4.0  H I 1.0  I J 1.0
;
```

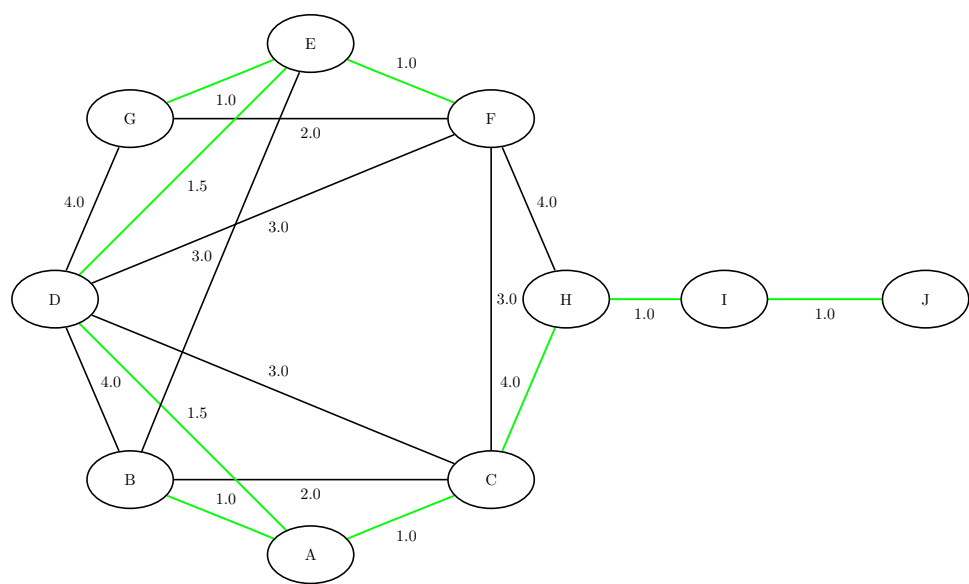
The following statements find a minimum spanning tree:

```
proc optgraph
  data_links = LinkSetInCompNet;
  minspantree
    out      = MinSpanTree;
run;
```



Output 1.12.1 shows the resulting data set MinSpanTree, which is displayed graphically in Figure 1.152 with the minimal cost links shown in green.

Figure 1.152 Minimum Spanning Tree for Office Computer Network



Output 1.12.1 Minimum Spanning Tree of a Computer Network Topology

from to weight		
I	J	1.0
A	C	1.0
E	F	1.0
E	G	1.0
H	I	1.0
A	B	1.0
A	D	1.5
D	E	1.5
C	H	4.0
13.0		

Example 1.13: Transitive Closure for Identification of Circular Dependencies in a Bug Tracking System

Most software bug tracking systems have some notion of *duplicate bugs* in which one bug is declared to be the same as another bug. If bug A is considered a duplicate (DUP) of bug B, then a fix for B would also fix A. You can represent the DUPs in a bug tracking system as a directed graph where you add a link  $A \rightarrow B$  if A is a DUP of B.

The bug tracking system needs to check for two situations when users declare a bug to be a DUP. The first situation is called a *circular dependence*. Consider bugs A, B, C, and D in the tracking system. The first user declares that A is a DUP of B and that C is a DUP of D. Then, a second user declares that B is a DUP of C, and a third user declares that D is a DUP of A. You now have a circular dependence, and no primary bug is

defined on which the development team should focus. You can easily see this circular dependence in the graph representation, because  $A \rightarrow B \rightarrow C \rightarrow D \rightarrow A$ . Finding such circular dependencies can be done using cycle detection, which is described in the section “Cycle” on page 107. However, the second situation that needs to be checked is more general. If a user declares that A is a DUP of B and another user declares that B is a DUP of C, this chain of duplicates is already an issue. The bug tracking system needs to provide one primary bug to which the rest of the bugs are duplicated. The existence of these chains can be identified by calculating the transitive closure of the directed graph that is defined by the DUP links.

Given the original directed graph  $G$  (defined by the DUP links) and its transitive closure  $G^T$ , any link in  $G^T$  that is not in  $G$  exists because of some chain that is present in  $G$ .

Consider the following data that define some duplicated bugs (called *defects*) in a small sample of the bug tracking system:

```
data DefectLinks;
  input defectId $ linkedDefect $ linkType $ when datetime16.;
  format when datetime16.;
  datalines;
D0096978 S0711218 DUPTO 20OCT10:00:00:00
S0152674 S0153280 DUPTO 30MAY02:00:00:00
S0153280 S0153307 DUPTO 30MAY02:00:00:00
S0153307 S0152674 DUPTO 30MAY02:00:00:00
S0162973 S0162978 DUPTO 29NOV10:16:13:16
S0162978 S0165405 DUPTO 29NOV10:16:13:16
S0325026 S0575748 DUPTO 01JUN10:00:00:00
S0347945 S0346582 DUPTO 03MAR06:00:00:00
S0350596 S0346582 DUPTO 21MAR06:00:00:00
S0539744 S0643230 DUPTO 10MAY10:00:00:00
S0575748 S0643230 DUPTO 15JUN10:00:00:00
S0629984 S0643230 DUPTO 01JUN10:00:00:00
;
```

The following statements calculate cycles in addition to the transitive closure of the graph  $G$  that is defined by the duplicated defects in DefectLinks. The output data set Cycles contains any circular dependencies, and the data set TransClosure contains the transitive closure  $G^T$ . To identify the chains, you can use PROC SQL to identify the links in  $G^T$  that are not in  $G$ .

```
proc optgraph
  loglevel          = moderate
  graph_direction   = directed
  data_links        = DefectLinks;
  data_links_var
    from            = defectId
    to              = linkedDefect;
  cycle
    out             = Cycles
    mode            = all_cycles;
  transitive_closure
    out             = TransClosure;
run;
%put &_OPTGRAPH_;
%put &_OPTGRAPH_CYCLE_;
%put &_OPTGRAPH_TRANSCL_;
```

```

proc sql;
  create table Chains as
  select defectId, linkedDefect from TransClosure
  except
  select defectId, linkedDefect from DefectLinks;
quit;

```

The progress of the procedure is shown in [Output 1.13.1](#).

**Output 1.13.1** PROC OPTGRAPH Log: Transitive Closure for Identification of Circular Dependencies in a Bug Tracking System

---

```

NOTE: -----
NOTE: -----
NOTE: Running OPTGRAPH version 14.1.
NOTE: -----
NOTE: -----
NOTE: The OPTGRAPH procedure is executing in single-machine mode.
NOTE: -----
NOTE: -----
NOTE: Reading the links data set.
NOTE: There were 12 observations read from the data set WORK.DEFECTLINKS.
NOTE: Data input used 0.00 (cpu: 0.00) seconds.
NOTE: Building the input graph storage used 0.00 (cpu: 0.00) seconds.
NOTE: The input graph storage is using 0.0 MBs (peak: 0.0 MBs) of memory.
NOTE: The number of nodes in the input graph is 16.
NOTE: The number of links in the input graph is 12.
NOTE: -----
NOTE: -----
NOTE: Processing cycle detection.
NOTE: The graph has 1 cycle.
NOTE: Processing cycle detection used 0.00 (cpu: 0.00) seconds.
NOTE: -----
NOTE: -----
NOTE: Processing the transitive closure.
NOTE: Processing the transitive closure used 0.00 (cpu: 0.00) seconds.
NOTE: -----
NOTE: -----
NOTE: Creating cycle data set output.
NOTE: Creating transitive closure data set output.
NOTE: Data output used 0.00 (cpu: 0.00) seconds.
NOTE: -----
NOTE: -----
NOTE: The data set WORK.CYCLES has 4 observations and 3 variables.
NOTE: The data set WORK.TRANSCLASURE has 20 observations and 2 variables.
STATUS=OK  CYCLE=OK  TRANSCL=OK
STATUS=OK  NUM_CYCLES=1  CPU_TIME=0.00  REAL_TIME=0.00
STATUS=OK  CPU_TIME=0.00  REAL_TIME=0.00
NOTE: Table WORK.CHAINS created, with 8 rows and 2 columns.

```

---

The data set Cycles contains one case of a circular dependence in which the DUPs start and end at S0152674.

**Output 1.13.2** Cycle in Bug Tracking System

cycle	order	node
1	1	S0152674
1	2	S0153280
1	3	S0153307
1	4	S0152674

The data set Chains contains the chains in the bug tracking system that come from the links in  $G^T$  that are not in  $G$ .

**Output 1.13.3** Chains in Bug Tracking System

defectId	linkedDefect
S0152674	S0152674
S0152674	S0153307
S0153280	S0152674
S0153280	S0153280
S0153307	S0153280
S0153307	S0153307
S0162973	S0165405
S0325026	S0643230

---

**Example 1.14: Reach Networks for Computation of Market Coverage of a Terrorist Network**

The problem of finding an efficient method for *covering* a market (a set of entities) is important in numerous industries. For example, consider that you are an advertising company with access to data that are collected from your customers' social networks. To keep costs at a minimum in some new promotion, you want to find a minimal set of customers to whom you need to advertise in order to reach the entire market. To solve this, you could first generate all the reach networks for each customer using PROC OPTGRAPH. These networks can then be used in a *set-covering problem*, which can be solved as an integer linear program using PROC OPTMODEL. Let  $N$  be the set of customers that you want to reach, and let the links  $A$  define the social network between those customers. If you use a one-hop reach network, you assume that if an advertisement is sent to customer  $i$ , then customer  $i$  will promote the advertisement to all his friends (those he is connected to in  $A$ ). If you use two-hop reach networks, you assume that customer  $i$ 's friends will also promote to their friends. So the question is: to which subset of customers should you advertise to reach all customers through the promotion mechanism?

This problem can be generalized as follows:

Given a graph  $G = (N, A)$ , choose a node set  $N^*$  of minimal size such that there is a path of length less than or equal to  $L$  to every node in  $N$  from a node in  $N^*$ .

To illustrate an application of this problem, consider again the terrorist communications network from the section “[Example 1.1: Articulation Points in a Terrorist Network](#)” on page 187. In this case, customers are alleged terrorists. Solving the covering problem here can tell you a subset of people to focus on in an investigation in order to cover all members of the network.

The following %GenerateReach macro runs PROC OPTGRAPH to generate the reach network for each person in the terrorist network for a variable hop limit:

```
%macro GenerateReach(limit=);
proc optgraph
  out_nodes      = NodeSetOut
  data_links     = LinkSetInTerror911;
  reach
    each_source
    out_nodes    = ReachNode
    maxreach     = &limit;
run;
%mend GenerateReach;
```

The following macro %SolverCover runs PROC OPTMODEL to solve the set-covering problem:

```
%macro SolverCover();
proc optmodel;
  string      tmpLabel;
  set<num>     NODE_ID;
  set<string>  NODE_LABEL init {};
  string      nodeIdToLabel{NODE_ID};
  num         nodeLabelToId{NODE_LABEL};

  set<num> REACH_SET{NODE_ID} init {};
  set<string,num> PAIRS;

  /* read data */
  read data NodeSetOut into NODE_ID=[_n_] nodeIdToLabel=node;
  read data ReachNode into PAIRS=[node reach];
  for{i in NODE_ID} do;
    tmpLabel = nodeIdToLabel[i];
    NODE_LABEL = NODE_LABEL union {tmpLabel};
    nodeLabelToId[tmpLabel] = i;
  end;
  for{<label,i> in PAIRS} do;
    REACH_SET[i] = REACH_SET[i] union {nodeLabelToId[label]};
  end;

  /* declare decision variables */
  var x {NODE_ID} binary;

  /* declare objective */
  minimize numNodes = sum{j in NODE_ID} x[j];

  /* cover constraint */
  con cover {i in NODE_ID}:
    sum{j in REACH_SET[i]} x[j] >= 1;

  /* solve */
  solve;

  create data Solution from [label]=
    (setof{j in NODE_ID : round(x[j].sol)=1}nodeIdToLabel[j]);
```

```
quit;
%mend SolverCover;
```

The following statements calculate the minimal cover for the one-hop limit:

```
%GenerateReach(limit=1);
%SolverCover();
```

In order to cover the network, assuming a one-hop limit, the investigators would need to investigate the people listed in the data set Solution, shown in [Output 1.14.1](#).

**Output 1.14.1** Minimal One-Hop Cover for Terrorist Communications Network

Obs	label
1	Djamal Beghal
2	Zacarias Moussaoui
3	Essid Sami Ben Khemais
4	Mohamed Atta
5	Agus Budiman
6	Mamduh Mahmud Salim
7	Fayez Ahmed
8	Nawaf Alhazmi
9	Hani Hanjour
10	Nabil al-Marabh

The following statements calculate the minimal cover for the two-hop limit:

```
%GenerateReach(limit=2);
%SolverCover();
```

If investigators assume a two-hop limit, they could focus their attention to the two people shown in [Output 1.14.2](#). Then by following their links (and their links' links) they could cover the entire network.

**Output 1.14.2** Minimal Two-Hop Cover for Terrorist Communications Network

Obs	label
1	Zacarias Moussaoui
2	Mohamed Atta

## Example 1.15: Traveling Salesman Tour through US Capital Cities

Consider a cross-country trip where you want to travel the fewest miles to visit all of the capital cities in all US states except Alaska and Hawaii. Finding the optimal route is an instance of the traveling salesman problem, which is described in section “[Traveling Salesman Problem](#)” on page 164.

The following PROC SQL statements use the built-in data set `maps.uscity` to generate a list of the capital cities and their latitude and longitude:

```
/* Get a list of the state capital cities (with lat and long) */
proc sql;
  create table Cities as
  select unique statecode as state, city, lat, long
  from maps.uscity
  where capital='Y' and statecode not in ('AK' 'PR' 'HI');
quit;
```

From this list, you can generate a links data set `CitiesDist` that contains the distances, in miles, between each pair of cities. The distances are calculated by using the SAS function `GEODIST`.

```
/* Create a list of all the possible pairs of cities */
proc sql;
  create table CitiesDist as
  select
    a.city as city1, a.lat as lat1, a.long as long1,
    b.city as city2, b.lat as lat2, b.long as long2,
    geodist(lat1, long1, lat2, long2, 'DM') as distance
  from Cities as a, Cities as b
  where a.city < b.city;
quit;
```

The following PROC OPTGRAPH statements find the optimal tour through each of the capital cities:

```
/* Find optimal tour using OPTGRAPH */
proc optgraph
  loglevel    = moderate
  data_links  = CitiesDist
  out_nodes   = TSPTourNodes;
  data_links_var
    from      = city1
    to        = city2
    weight    = distance;
  tsp
    out       = TSPTourLinks;
run;
%put &_OPTGRAPH_;
%put &_OPTGRAPH_TSP_;
```

The progress of the procedure is shown in [Output 1.15.1](#). The total mileage needed to optimally traverse the capital cities is 10,627.75 miles.

**Output 1.15.1** PROC OPTGRAPH Log: Traveling Salesman Tour through US Capital Cities

```

NOTE: -----
NOTE: -----
NOTE: Running OPTGRAPH version 14.1.
NOTE: -----
NOTE: -----
NOTE: The OPTGRAPH procedure is executing in single-machine mode.
NOTE: -----
NOTE: -----
NOTE: Reading the links data set.
NOTE: There were 1176 observations read from the data set WORK.CITIESDIST.
NOTE: Data input used 0.00 (cpu: 0.00) seconds.
NOTE: Building the input graph storage used 0.00 (cpu: 0.00) seconds.
NOTE: The input graph storage is using 0.1 MBs (peak: 0.2 MBs) of memory.
NOTE: The number of nodes in the input graph is 49.
NOTE: The number of links in the input graph is 1176.
NOTE: -----
NOTE: -----
NOTE: Processing the traveling salesman problem.
NOTE: The initial TSP heuristics found a tour with cost 10645.918753 using 0.08 (cpu: 0.08)
NOTE: The MILP presolver value NONE is applied.
NOTE: The MILP solver is called.
NOTE: The Branch and Cut algorithm is used.

```

Node	Active	Sols	BestInteger	BestBound	Gap	Time
0	1	1	10645.9187534	10040.5139714	6.03%	0
0	1	1	10645.9187534	10241.6970024	3.95%	0
0	1	1	10645.9187534	10262.9074205	3.73%	0
0	1	1	10645.9187534	10350.0790852	2.86%	0
0	1	1	10645.9187534	10381.4538838	2.55%	0
0	1	1	10645.9187534	10470.6122886	1.67%	0
0	1	1	10645.9187534	10492.6949171	1.46%	0
0	1	1	10645.9187534	10560.1837745	0.81%	0
0	1	1	10645.9187534	10576.0823291	0.66%	0
0	1	1	10645.9187534	10590.9748294	0.52%	0
0	1	1	10645.9187534	10607.8528157	0.36%	0
0	1	1	10645.9187534	10607.8528157	0.36%	0

```

NOTE: The MILP solver added 12 cuts with 3309 cut coefficients at the root.
NOTE: Objective of the best integer solution found = 10645.918753.
NOTE: Processing the traveling salesman problem used 0.10 (cpu: 0.09) seconds.
NOTE: -----
NOTE: -----
NOTE: The TSP solver is restarting using a reduced graph with 49 nodes and 97 links.
NOTE: -----
NOTE: -----
NOTE: Processing the traveling salesman problem.
NOTE: The initial TSP heuristics found a tour with cost 10645.918753 using 0.00 (cpu: 0.00)
NOTE: The MILP presolver value NONE is applied.
NOTE: The MILP solver is called.
NOTE: The Branch and Cut algorithm is used.

```

Node	Active	Sols	BestInteger	BestBound	Gap	Time
------	--------	------	-------------	-----------	-----	------



**Output 1.15.1** *continued*

0	1	1	10645.9187534	10040.5139714	6.03%	0
0	1	1	10645.9187534	10241.6970024	3.95%	0
0	1	1	10645.9187534	10262.9074205	3.73%	0
0	1	1	10645.9187534	10350.0790852	2.86%	0
0	1	1	10645.9187534	10381.4538838	2.55%	0
0	1	1	10645.9187534	10549.5506188	0.91%	0
0	1	1	10645.9187534	10576.0823291	0.66%	0
0	1	1	10645.9187534	10590.9748294	0.52%	0
0	1	1	10645.9187534	10607.8528157	0.36%	0
0	1	4	10645.9187534	10607.8528157	0.36%	0

NOTE: The MILP solver added 9 cuts with 51 cut coefficients at the root.

1	1	5	10627.7543183	10607.8528157	0.19%	0
2	0	5	10627.7543183	10627.7543183	0.00%	0

NOTE: Objective = 10627.754318.

NOTE: Processing the traveling salesman problem used 0.01 (cpu: 0.02) seconds.

NOTE: -----

NOTE: -----

NOTE: Creating nodes data set output.

NOTE: Creating traveling salesman data set output.

NOTE: Data output used 0.00 (cpu: 0.00) seconds.

NOTE: -----

NOTE: -----

NOTE: The data set WORK.TSPTOURNODES has 49 observations and 2 variables.

NOTE: The data set WORK.TSPTOURLINKS has 49 observations and 3 variables.

STATUS=OK TSP=OPTIMAL

STATUS=OPTIMAL OBJECTIVE=10627.754318 RELATIVE\_GAP=0 ABSOLUTE\_GAP=0 PRIMAL\_INFEASIBILITY=0

BOUND\_INFEASIBILITY=0 INTEGER\_INFEASIBILITY=0 BEST\_BOUND=10627.754318 NODES=3

ITERATIONS=366 CPU\_TIME=0.11 REAL\_TIME=0.11

The following PROC GPROJECT and PROC GMAP statements produce a graphical display of the solution:

```

/* Merge latitude and longitude */
proc sql;
  /* merge in the lat & long for city1 */
  create table TSPTourLinksAnno1 as
  select unique TSPTourLinks.*, cities.lat as lat1, cities.long as long1
  from TSPTourLinks left join cities
  on TSPTourLinks.city1=cities.city;
  /* merge in the lat & long for city2 */
  create table TSPTourLinksAnno2 as
  select unique TSPTourLinksAnno1.*, cities.lat as lat2, cities.long as long2
  from TSPTourLinksAnno1 left join cities
  on TSPTourLinksAnno1.city2=cities.city;
quit;

/* Create the annotated data set to draw the path on the map
   (convert lat & long degrees to radians, since the map is in radians) */
data anno_path;
  set TSPTourLinksAnno2;

```

```

length function color $8;
xsys='2'; ysys='2'; hsys='3'; when='a'; anno_flag=1;
function='move';
x=atan(1)/45 * long1;
y=atan(1)/45 * lat1;
output;
function='draw';
color="blue"; size=0.8;
x=atan(1)/45 * long2;
y=atan(1)/45 * lat2;
output;
run;

/* Get a map with only the contiguous 48 states */
data states;
    set maps.states (where=(fipstate(state) not in ('HI' 'AK' 'PR')));
run;

data combined;
    set states anno_path;
run;

/* Project the map and annotate the data */
proc gproject data=combined out=combined dupok;
    id state;
run;

data states anno_path;
    set combined;
    if anno_flag=1 then output anno_path;
    else                output states;
run;

/* Get a list of the endpoints locations */
proc sql;
    create table anno_dots as
    select unique x, y from anno_path;
quit;

/* Create the final annotate data set */
data anno_dots;
    set anno_dots;
    length function color $8;
    xsys='2'; ysys='2'; when='a'; hsys='3';
    function='pie';
    rotate=360; size=0.8; style='psolid'; color="red";
    output;
    style='pempty'; color="black";
    output;
run;

/* Generate the map with GMAP */
pattern1 v=s c=cxcccffcc repeat=100;
proc gmap data=states map=states anno=anno_path all;

```

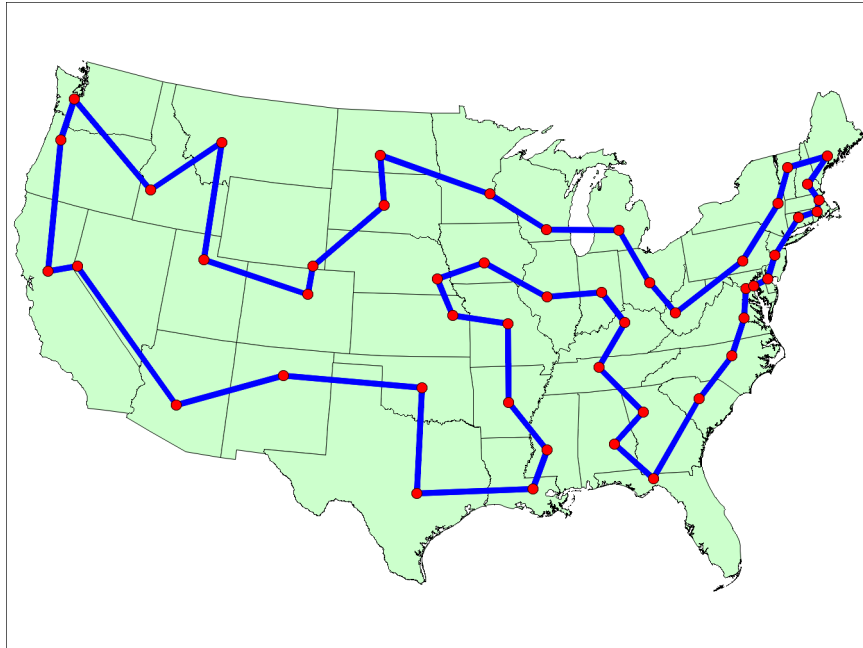
```

id state;
choro state / levels=1 nolegend coutline=black
anno=anno_dots des='' name="tsp";
run;

```

The minimal cost tour through the capital cities is shown on the US map in Figure 1.15.2.

**Output 1.15.2** Optimal Traveling Salesman Tour through US Capital Cities



The data set TSPTourLinks contains the links in the optimal tour. To display the links in the order they are to be visited, you can use the following DATA step:

```

/* Create the directed optimal tour */
data TSPTourLinksDirected(drop=next);
  set TSPTourLinks;
  retain next;
  if _N_ ne 1 and city1 ne next then do;
    city2 = city1;
    city1 = next;
  end;
  next = city2;
run;

```

The data set TSPTourLinksDirected is shown in Figure 1.153.

**Figure 1.153** Links in the Optimal Traveling Salesman Tour

City Name	City Name	distance	City Name	City Name	distance
Montgomery	Tallahassee	177.14	Denver	Salt Lake City	373.05
Tallahassee	Columbia	311.23	Salt Lake City	Helena	403.40
Columbia	Raleigh	182.99	Helena	Boise City	291.20
Raleigh	Richmond	135.58	Boise City	Olympia	401.31
Richmond	Washington	97.96	Olympia	Salem	146.00
Washington	Annapolis	27.89	Salem	Sacramento	447.40
Annapolis	Dover	54.01	Sacramento	Carson City	101.51
Dover	Trenton	83.88	Carson City	Phoenix	577.84
Trenton	Hartford	151.65	Phoenix	Santa Fe	378.27
Hartford	Providence	65.56	Santa Fe	Oklahoma City	474.92
Providence	Boston	38.41	Oklahoma City	Austin	357.38
Boston	Concord	66.30	Austin	Baton Rouge	394.78
Concord	Augusta	117.36	Baton Rouge	Jackson	139.75
Augusta	Montpelier	139.32	Jackson	Little Rock	206.87
Montpelier	Albany	126.19	Little Rock	Jefferson City	264.75
Albany	Harrisburg	230.24	Jefferson City	Topeka	191.67
Harrisburg	Charleston	287.34	Topeka	Lincoln	132.94
Charleston	Columbus	134.64	Lincoln	Des Moines	168.10
Columbus	Lansing	205.08	Des Moines	Springfield	243.02
Lansing	Madison	246.88	Springfield	Indianapolis	186.46
Madison	Saint Paul	226.25	Indianapolis	Frankfort	129.90
Saint Paul	Bismarck	391.25	Frankfort	Nashville-Davidson	175.58
Bismarck	Pierre	170.27	Nashville-Davidson	Atlanta	212.61
Pierre	Cheyenne	317.90	Atlanta	Montgomery	145.39
Cheyenne	Denver	98.33			<b>10,627.75</b>

## References

- Ahuja, R. K., Magnanti, T. L., and Orlin, J. B. (1993). *Network Flows: Theory, Algorithms, and Applications*. Englewood Cliffs, NJ: Prentice-Hall.
- Applegate, D. L., Bixby, R. E., Chvátal, V., and Cook, W. J. (2006). *The Traveling Salesman Problem: A Computational Study*. Princeton, NJ: Princeton University Press.
- Batagelj, V., and Zaversnik, M. (2003). “An  $O(m)$  Algorithm for Cores Decomposition of Networks.” *Computing Research Repository* cs.DS/0310049.
- Blondel, V. D., Guillaume, J. L., Lambiotte, R., and Lefebvre, E. (2008). “Fast Unfolding of Communities in Large Networks.” *Journal of Statistical Mechanics: Theory and Experiment* 10:10000–10014.
- Boitmanis, K., Freivalds, K., Ledins, P., and Opmanis, R. (2006). “Fast and Simple Approximation of the Diameter and Radius of a Graph.” In *Experimental Algorithms*, vol. 4007, edited by C. Alvarez, and M. Serna, 98–108. Berlin: Springer-Verlag. doi:10.1007/11764298\_9.
- Bron, C., and Kerbosch, J. (1973). “Algorithm 457: Finding All Cliques of an Undirected Graph.” *Communications of the ACM* 16:48–50.

- Cormen, T. H., Leiserson, C. E., and Rivest, R. L. (1990). *Introduction to Algorithms*. Cambridge, MA, and New York: MIT Press and McGraw-Hill.
- Fowler, J. H., and Joen, S. (2008). “The Authority of Supreme Court Precedent.” *Social Networks* 30:16–30. <http://jhfowler.ucsd.edu/judicial.htm>.
- Google (2011). “Google Maps.” Accessed March 16, 2011. <http://maps.google.com>.
- Harley, E. R. (2003). *Graph Algorithms for Assembling Integrated Genome Maps*. Ph.D. diss., University of Toronto.
- Johnson, D. B. (1975). “Finding All the Elementary Circuits of a Directed Graph.” *SIAM Journal on Computing* 4:77–84.
- Jonker, R., and Volgenant, A. (1987). “A Shortest Augmenting Path Algorithm for Dense and Sparse Linear Assignment Problems.” *Computing* 38:325–340.
- Kleinberg, J. (1998). “Authoritative Sources in a Hyperlinked Environment.” In *Proceedings of the Ninth Annual ACM-SIAM Symposium on Discrete Algorithms*, 668–677. Philadelphia: Society for Industrial and Applied Mathematics.
- Krackhardt, D. (1990). “Assessing the Political Landscape: Structure, Cognition, and Power in Organizations.” *Administrative Science Quarterly* 35:342–369.
- Krebs, V. (2002). “Uncloaking Terrorist Networks.” *First Monday* 7. [http://www.firstmonday.org/issues/issue7\\_4/krebs/](http://www.firstmonday.org/issues/issue7_4/krebs/).
- Kruskal, J. B. (1956). “On the Shortest Spanning Subtree of a Graph and the Traveling Salesman Problem.” *Proceedings of the American Mathematical Society* 7:48–50.
- Kumar, R., and Li, H. (1994). “On Asymmetric TSP: Transformation to Symmetric TSP and Performance Bound.” <http://home.eng.iastate.edu/~rkumar/PUBS/atsp.pdf>.
- Lancichinetti, A., and Fortunato, S. (2009). “Community Detection Algorithms: A Comparative Analysis.” *Physical Review E* 80:056117–056128.
- Landes, W. M., and Posner, R. A. (1976). “Legal Precedent: A Theoretical and Empirical Analysis.” *Journal of Law and Economics* 19:249–307.
- Mihalcea, R. (2005). “Unsupervised Large-Vocabulary Word Sense Disambiguation with Graph-Based Algorithms for Sequence Data Labeling.” In *Proceedings of the Conference on Human Language Technology and Empirical Methods in Natural Language Processing*, 411–418. Vancouver.
- Newman, M. E. J. (2010). *Networks: An Introduction*. Oxford: Oxford University Press.
- Raghavan, U. N., Albert, R., and Kumara, S. (2007). “Near Linear Time Algorithm to Detect Community Structures in Large-Scale Networks.” *Physical Review E* 76:36106–36117.
- Rochat, Y. (2009). “Closeness Centrality Extended to Unconnected Graphs: The Harmonic Centrality Index.” Paper presented at Sixth Applications of Social Network Analysis Conference, Zurich. [http://infoscience.epfl.ch/record/200525/files/\[EN\]ASNA09.pdf](http://infoscience.epfl.ch/record/200525/files/[EN]ASNA09.pdf).
- Ronhovde, P., and Nussinov, Z. (2010). “Local Resolution-Limit-Free Potts Model for Community Detection.” *Physical Review E* 81:46114–46129.

- Sinha, R., and Mihalcea, R. (2007). “Unsupervised Graph-Based Word Sense Disambiguation Using Measures of Word Semantic Similarity.” In *Proceedings of the IEEE International Conference on Semantic Computing*, 363–369. Los Alamitos, CA: IEEE Computer Society Press.
- Sleijpen, G. L. G., and van der Vorst, H. A. (2000). “A Jacobi-Davidson Iteration Method for Linear Eigenvalue Problems.” *SIAM Review* 42:267–293.
- Stoer, M., and Wagner, F. (1997). “A Simple Min-Cut Algorithm.” *Journal of the Association for Computing Machinery* 44:585–591.
- Tarjan, R. E. (1972). “Depth-First Search and Linear Graph Algorithms.” *SIAM Journal on Computing* 1:146–160.
- Traag, V. A., Van Dooren, P., and Nesterov, Y. (2011). “Narrow Scope for Resolution-Limit-Free Community Detection.” *Physical Review E* 84:016114 (1–9). doi:10.1103/PhysRevE.84.016114. <http://pre.aps.org/abstract/PRE/v84/i1/e016114>.
- Willingham, V. (2009). “Massive Transplant Effort Pairs 13 Kidneys to 13 Patients.” CNN Health. Accessed March 16, 2011. <http://www.cnn.com/2009/HEALTH/12/14/kidney.transplant/index.html>.
- Zachary, W. W. (1977). “An Information Flow Model for Conflict and Fission in Small Groups.” *Journal of Anthropological Research* 33:452–473.

# Index

- ABSOBJGAP= option
  - TSP statement, 45
- ALGORITHM= option
  - COMMUNITY statement, 26
  - CONCOMP statement, 28
- AUTH= option
  - CENTRALITY statement, 20
- BETWEEN= option
  - CENTRALITY statement, 21
- BETWEEN\_NORM= option
  - CENTRALITY statement, 21
- BICONCOMP option
  - SUMMARY statement, 42
- BICONCOMP statement
  - statement options, 20
- BY\_CLUSTER option
  - CENTRALITY statement, 21
  - REACH statement, 39
  - SUMMARY statement, 42
- CENTRALITY statement
  - statement options, 20
- CLIQUE statement
  - statement options, 25
- CLOSE= option
  - CENTRALITY statement, 21
- CLOSE\_NOPATH= option
  - CENTRALITY statement, 22
- CLUSTER= option
  - DATA\_NODES\_VAR statement, 33
- CLUSTERING\_COEF option
  - CENTRALITY statement, 22
- COMMUNITY statement
  - statement options, 26
- CONCOMP option
  - SUMMARY statement, 42
- CONCOMP statement
  - statement options, 28
- CONFLICTSEARCH= option
  - TSP statement, 45
- CORE statement
  - statement options, 29
- CUTOFF= option
  - TSP statement, 45
- CUTSTRATEGY= option
  - TSP statement, 45
- CYCLE statement
  - statement options, 30
- DATA\_LINKS= option
  - PROC OPTGRAPH statement, 17
- DATA\_LINKS\_VAR statement
  - statement options, 32
- DATA\_MATRIX= option
  - PROC OPTGRAPH statement, 17
- DATA\_NODES= option
  - PROC OPTGRAPH statement, 17
- DATA\_NODES\_SUB= option
  - PROC OPTGRAPH statement, 17
- DATA\_NODES\_VAR statement
  - statement options, 33
- DEGREE= option
  - CENTRALITY statement, 22
- DETAILS option
  - PERFORMANCE statement, 38
- DIAMETER\_APPROX= option
  - SUMMARY statement, 42
- DIGRAPH option
  - REACH statement, 39
- EACH\_SOURCE option
  - REACH statement, 39
- EIGEN= option
  - CENTRALITY statement, 22
- EIGEN\_ALGORITHM= option
  - CENTRALITY statement, 23
- EIGEN\_MAXITER= option
  - CENTRALITY statement, 23
- EIGENVALUES= option
  - EIGENVECTOR statement, 34
- EIGENVECTOR statement
  - statement options, 34
- EMPHASIS= option
  - TSP statement, 45
- FILTER\_SUBGRAPH= option
  - PROC OPTGRAPH statement, 18
- FROM= option
  - DATA\_LINKS\_VAR statement, 32
- GRAPH\_DIRECTION= option
  - PROC OPTGRAPH statement, 18
- GRAPH\_INTERNAL\_FORMAT= option
  - PROC OPTGRAPH statement, 18
- HEURISTICS= option
  - TSP statement, 46
- HUB= option

CENTRALITY statement, 23

ID= option  
 LINEAR\_ASSIGNMENT statement, 35

IGNORE\_SELF option  
 REACH statement, 39

INCLUDE\_SELFLINK option  
 PROC OPTGRAPH statement, 19

INFLUENCE= option  
 CENTRALITY statement, 23

LINEAR\_ASSIGNMENT statement  
 statement options, 35

LINK\_REMOVAL\_RATIO= option  
 COMMUNITY statement, 26

LINKS= option  
 CORE statement, 29

LOGFREQ= option  
 MINCOSTFLOW statement, 36  
 SHORTPATH statement, 40  
 TSP statement, 46

LOGFREQNODE= option  
 CENTRALITY statement, 24  
 SUMMARY statement, 42

LOGFREQTIME= option  
 CENTRALITY statement, 24  
 REACH statement, 39  
 SUMMARY statement, 43

LOGLEVEL= option  
 BICONCOMP statement, 20  
 CENTRALITY statement, 24  
 CLIQUE statement, 25  
 COMMUNITY statement, 26  
 CONCOMP statement, 29  
 CORE statement, 30  
 CYCLE statement, 30  
 EIGENVECTOR statement, 34  
 LINEAR\_ASSIGNMENT statement, 35  
 MINCOSTFLOW statement, 36  
 MINCUT statement, 37  
 MINSPANTREE statement, 37  
 PROC OPTGRAPH statement, 19  
 REACH statement, 39  
 SHORTPATH statement, 40  
 SUMMARY statement, 43  
 TRANSITIVE\_CLOSURE statement, 44  
 TSP statement, 46

LOWER= option  
 DATA\_LINKS\_VAR statement, 32

MAXCLIQUES= option  
 CLIQUE statement, 25

MAXCYCLES= option  
 CYCLE statement, 31

MAXITER= option

COMMUNITY statement, 27

EIGENVECTOR statement, 34

MAXLENGTH= option  
 CYCLE statement, 31

MAXLINKWEIGHT= option  
 CYCLE statement, 31

MAXNODES= option  
 TSP statement, 47

MAXNODEWEIGHT= option  
 CYCLE statement, 31

MAXNUMCUTS= option  
 MINCUT statement, 37

MAXREACH= option  
 REACH statement, 39

MAXSOLS= option  
 TSP statement, 47

MAXTIME= number  
 MINCOSTFLOW statement, 36

MAXTIME= option  
 CLIQUE statement, 25  
 CORE statement, 30  
 CYCLE statement, 31  
 TSP statement, 47

MAXWEIGHT= option  
 MINCUT statement, 37

MILP= option  
 TSP statement, 47

MINCOSTFLOW statement  
 statement options, 36

MINCUT statement  
 statement options, 36

MINLENGTH= option  
 CYCLE statement, 31

MINLINKWEIGHT= option  
 CYCLE statement, 31

MINNODEWEIGHT= option  
 CYCLE statement, 31

MINSPANTREE statement  
 statement options, 37

MODULARITY= option  
 COMMUNITY statement, 28

NEIGEN= option  
 EIGENVECTOR statement, 34

NODE= option  
 DATA\_NODES\_VAR statement, 33

NODESEL= option  
 TSP statement, 47

OPTGRAPH Procedure, 10

PERFORMANCE statement, 38

OUT= option  
 CLIQUE statement, 26  
 CYCLE statement, 32



- EIGENVECTOR statement, 34
- LINEAR\_ASSIGNMENT statement, 35
- MINCUT statement, 37
- MINSPAN TREE statement, 38
- SHORTPATH statement, 41
- SUMMARY statement, 43
- TRANSITIVE\_CLOSURE statement, 44
- TSP statement, 48
- OUT\_COMM\_LINKS= option
  - COMMUNITY statement, 27
- OUT\_COMMUNITY= option
  - COMMUNITY statement, 27
- OUT\_COUNTS1= option
  - REACH statement, 40
- OUT\_COUNTS2= option
  - REACH statement, 40
- OUT\_COUNTS= option
  - REACH statement, 40
- OUT\_LEVEL= option
  - COMMUNITY statement, 27
- OUT\_LINKS= option
  - PROC OPTGRAPH statement, 19
  - REACH statement, 40
- OUT\_NODES= option
  - PROC OPTGRAPH statement, 19
  - REACH statement, 40
- OUT\_OVERLAP= option
  - COMMUNITY statement, 27
- OUT\_PATHS= option
  - SHORTPATH statement, 41
- OUT\_WEIGHTS= option
  - SHORTPATH statement, 41
- PATHS= option
  - SHORTPATH statement, 41
- PERFORMANCE statement, 38
  - DETAILS option, 38
- PROBE= option
  - TSP statement, 48
- PROC OPTGRAPH statement
  - statement options, 17
- RANDOM\_FACTOR= option
  - COMMUNITY statement, 27
- RANDOM\_SEED= option
  - COMMUNITY statement, 27
- REACH statement
  - statement options, 39
- RECURSIVE (options)
  - COMMUNITY statement, 27
- RELOBJGAP= option
  - TSP statement, 48
- RESOLUTION\_LIST= option
  - COMMUNITY statement, 28
- SHORTPATH statement
  - statement options, 40
- SHORTPATH= option
  - SUMMARY statement, 43
- SINK= option
  - SHORTPATH statement, 41
- SOURCE= option
  - SHORTPATH statement, 41
- STANDARDIZED\_LABELS option
  - PROC OPTGRAPH statement, 19
- STRONGITER= option
  - TSP statement, 48
- STRONGLEN= option
  - TSP statement, 48
- SUBSIZESWITCH= option
  - CENTRALITY statement, 25
  - SUMMARY statement, 43
- SUMMARY statement
  - statement options, 42
- TARGET= option
  - TSP statement, 48
- TIMETYPE= option
  - PROC OPTGRAPH statement, 19
- TO= option
  - DATA\_LINKS\_VAR statement, 32
- TOLERANCE= option
  - COMMUNITY statement, 28
- TRANSITIVE\_CLOSURE statement
  - statement options, 44
- TSP statement
  - statement options, 44
- UPPER= option
  - DATA\_LINKS\_VAR statement, 32
- USEWEIGHT= option
  - SHORTPATH statement, 41
- VARSEL= option
  - TSP statement, 48
- WEIGHT2= option
  - CENTRALITY statement, 25
  - DATA\_NODES\_VAR statement, 33
  - SHORTPATH statement, 42
- WEIGHT= option
  - DATA\_LINKS\_VAR statement, 32
  - DATA\_NODES\_VAR statement, 33
  - LINEAR\_ASSIGNMENT statement, 35





# Gain Greater Insight into Your SAS<sup>®</sup> Software with SAS Books.

Discover all that you need on your journey to knowledge and empowerment.



[support.sas.com/bookstore](http://support.sas.com/bookstore)  
for additional books and resources.

