Paper 232-2012

Getting to the Good Part of Data Analysis: Data Access, Manipulation, and Customization Using JMP®

Audrey Ventura, SAS Institute Inc., Cary, NC

ABSTRACT

Effective data analysis requires easy access to your data no matter what format it comes in. JMP can handle a wide variety of formats. Once the data is in JMP, you can choose from a variety of options to reshape the data with just a few clicks. Finally, customize your data with labels, colors, and data roles so that graphs and charts automatically look the way you want them to. This paper walks through two or three story lines that demonstrate how JMP can easily import, reshape, and customize data (even large datasets) in ways that allow your data to be displayed in vibrant visualizations that will wow your audience.

INTRODUCTION

Tasks such as data cleaning or reshaping data might be up to 80% of the time spent with a particular data set. Data cleaning is the task of managing missing values, adjusting formats, fixing dates and times, adding new columns or formulas, and removing duplicate or erroneous rows. Reshaping data is the task of molding your data into the format that is appropriate for the type of analysis you want to do and might consist of joining multiple data sources, transposing rows and columns, or splitting or stacking columns. Once all this is done, you have finally gotten to the "good part", where you can ask questions of your data.

JMP is an interactive statistical visualization and discovery tool that was developed by SAS. JMP has been available since 1989. The current version of JMP, released in March 2012, is JMP 10. One of the central pieces of JMP is its data table, which is a rich data manipulation tool. JMP makes it easy to get to the "good part" of data analysis.

This paper will summarize the data access capabilities of JMP, and then it will step through three examples starting with raw data to show various ways JMP makes data manipulation easy.

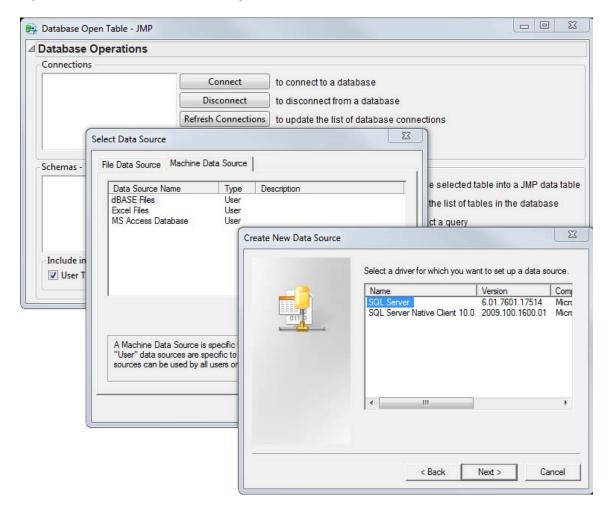
- Data Access
- Example 1 N.C. State Fair Attendance Internet Open, Value Ordering
- Example 2 U.S. Unemployment Column Name Add-In, Standardize Column Attributes, Stack
- Example 3 Baseball Attendance Concatenate, Join, Standardize Column Attributes, Recode

DATA ACCESS

It is worth noting before discussing the examples that JMP can access a wide variety of data types. The examples will showcase only a few ways of accessing data to demonstrate its flexibility. A common way

is copying and pasting data into JMP from another application such as Excel or a text file. JMP can also import many common file types such as Excel, SAS, or text files. JMP 9 introduced mapping capabilities and the ability to read ESRI shapefiles. If you are using a JMP competitor product, such as Minitab or SPSS, it is easy for you to experiment with JMP as you can import those data files directly into JMP. Finally, JMP can connect to most types of databases, since more and more data is stored on database servers.

Figure 1. Database Connections Using JMP



EXAMPLE 1 N.C. STATE FAIR ATTENDANCE DATA

Suppose you are interested in attending the N.C. State Fair, which spans 11 days. Your schedule is flexible and you are interested in attending on a day with low attendance so that the lines will be shorter. A search finds the N.C. State Fair¹ web site, which had attendance data for the last 25 years. You can use **File->Internet Open** in JMP to import the data.

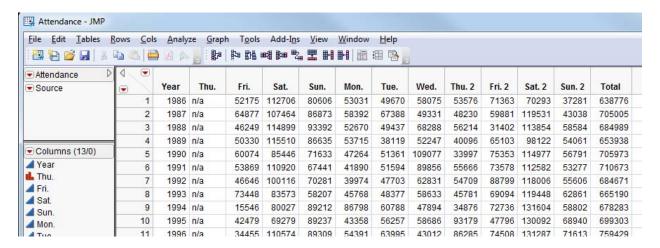
¹ Data from the N.C. State Fair web site is available at http://www.ncstatefair.org/2011/About/Attendance.htm.

Figure 2. Internet Open Dialog Box



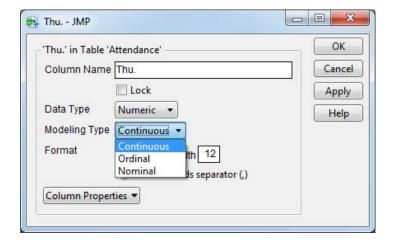
The data is imported nicely, with one exception in the first Thursday column, which has been imported as character data. JMP interpreted the "n/a" value as character, but this is easily changed. Also note, the columns named Thu 2, Fri 2, Sat 2 and Sun 2 refer to the second Thursday, second Friday, and so on, since the fair spans two weekends.

Figure 3. Attendance Data from the N.C. State Fair



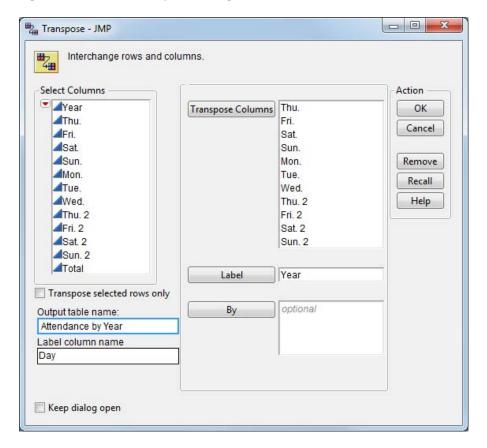
Two actions are needed to analyze attendance per day. First, change Thu to a continuous column through the **Column Info** dialog box.

Figure 4. Change Data Type and Modeling Type in the Column Info Dialog Box



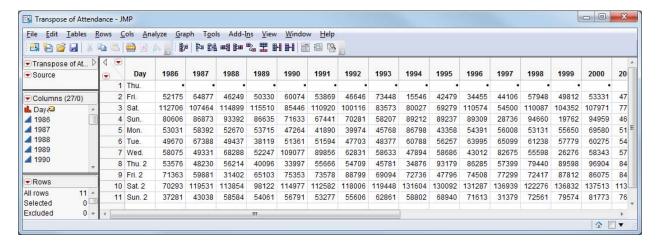
Second, transpose the data so that years are represented in column names and attendance by day is represented on each row. You can access JMP Transpose functionality through **Tables->Transpose**.

Figure 5. The JMP Transpose Dialog Box



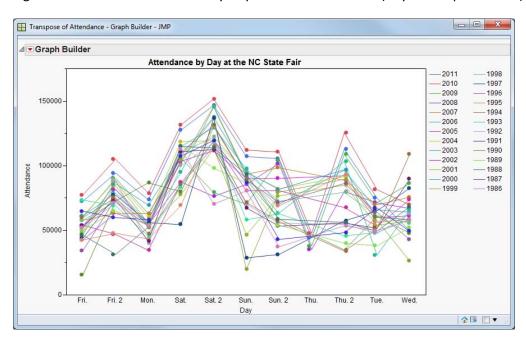
Once you press OK, you will have a new data table with the rows and columns transposed.

Figure 6. Transposed Attendance Data



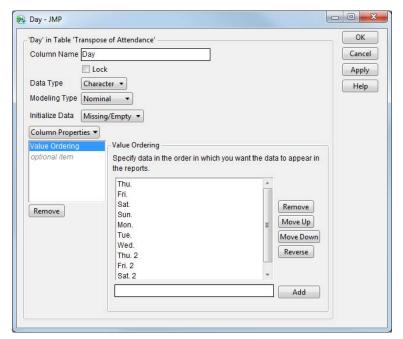
Now you have reached the "good part" and the data is in a format that can be analyzed. Figure 7 shows a line plot of attendance for each day. However, there is one small issue shown in Figure 7: the default in JMP is to alphabetize the days of the week along the x-axis. The graph is not displaying incorrect data, but the days of the week are not in the order you would like to see.

Figure 7. Line Plot of Attendance by Day at the N.C. State Fair (Days Are Alphabetized)



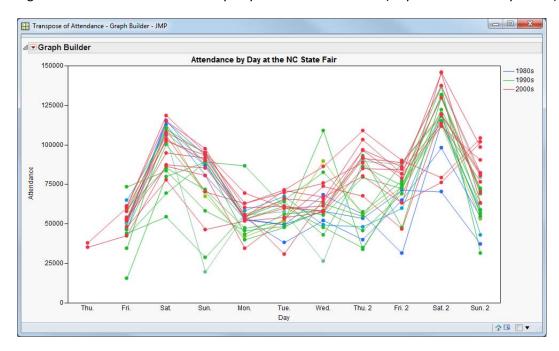
Revisit the data table to add a **Value Ordering** column property to the Day column by way of the **Column Info** dialog box. Value Ordering tells the column what order you would like to display the values.

Figure 8. Example of Adding the Value Ordering Column Property



Once value ordering has been added, now you can recreate the Graph Builder report using the proper order of the days. Additional customization will show you trends of the last three decades.

Figure 9. Line Plot of Attendance by Day at the N.C. State Fair (Days Are in Weekday Order)



You can see from this graph that Saturdays are, not surprisingly, the days with the heaviest attendance. Attendance on the first Thursday is low, but there are also very few data points. Tuesday might be the

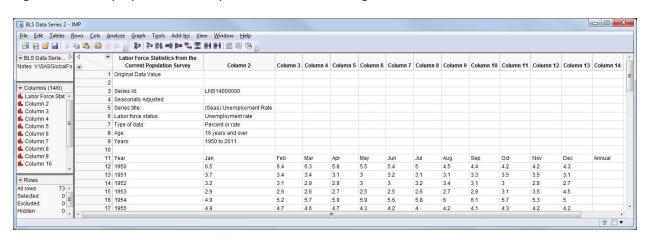
best day to attend the fair based on lowest average attendance in recent years. You could further analyze this data by looking at attendance trends for the last 10 years only.

EXAMPLE 2 U.S. UNEMPLOYMENT DATA (Seasonally Adjusted)

For the next example, suppose you are interested in comparing the recent rates of high unemployment to trends in unemployment rates over the last 60 years. The Bureau of Labor Statistics² makes employment data publicly available on their web site. This example uses information from their Current Population Surveys³ for seasonally adjusted unemployment rates.

The tables available on the Bureau of Labor Statistics web site are not able to be imported using Internet Open in JMP (which you used in Example 1), due to the way queries are processed and returned. However, the tables are easily downloaded into a Microsoft Excel format that can be imported into JMP. The first attempt does not work well, because of 10 rows of notes at the top of the spreadsheet. (See Figure 10.) JMP imports data as columns, and was not expecting the superfluous rows. You have two options. One option would be to close the JMP table and simply copy-and-paste from Excel to JMP. The second option is to proceed and clean up the extra rows in JMP. In this case, you choose the second option.

Figure 10. Unemployment Data Imported into JMP Showing 10 Rows of Notes



JMP can manage the small difficulty posed by the extra rows, as it is an easy thing to delete the 10 rows of notes. Next, you can use the Column Names Add-In⁴ to move the header row from row 1 of the data table into the column names area. The Column Names add-in is available on the JMP File Exchange⁵.

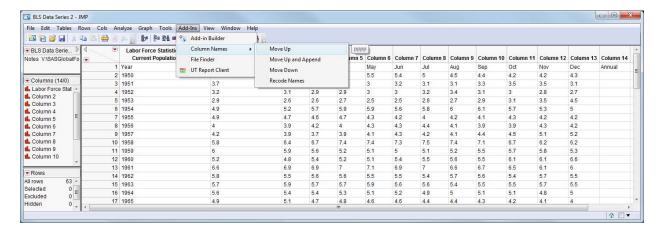
² Bureau of Labor Statistics Historical "A" tables are available at http://www.bls.gov/webapps/legacy/cpsatab1.htm.

³ Current Population Surveys from BLS are available at http://www.bls.gov/cps/#data.

⁴ Column Names Add-In is accessible on the JMP File Exchange at http://support.sas.com/demosdownloads/downarea_t4.jsp?productID=110491&jmpflag=Y.

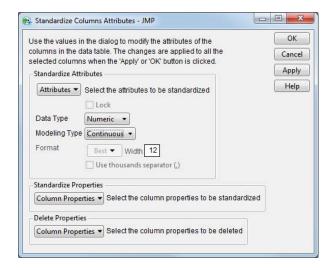
JMP File Exchange can be accessed from http://www.jmp.com/community/.

Figure 11. Using the Column Names Add-In



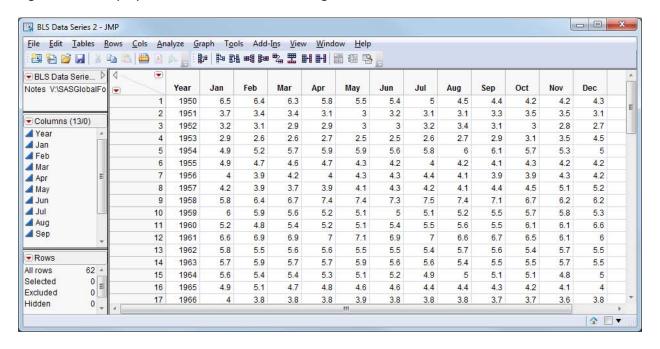
The Annual column is empty for seasonally adjusted data, so it can be deleted. There is also still a small issue to correct in that all the data is considered character data. Use Standardize Attributes functionality in JMP to change the properties of many columns at once. In this case change all the columns to continuous using **Cols->Standardize Attributes**. In this dialog box, change the data type to numeric and change the modeling type to continuous.

Figure 12. Standardize Column Attributes Dialog Box



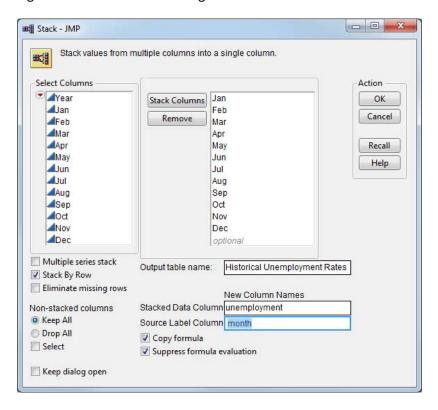
Now the data table is ready to be analyzed.

Figure 13. Unemployment Data, After Data Cleaning



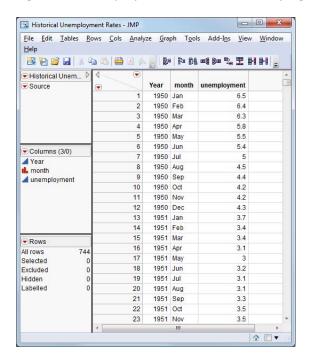
You would like to plot the data over time, but the data is not in the correct format. What you need to do next is *stack* the data, so that each year has 12 rows of data, one for each month. JMP makes this easy with its Stack functionality, which you can invoke through **Tables->Stack**.

Figure 14. The JMP Stack Dialog Box



The resulting table is what you expect. Instead of 62 rows for 62 years of data, we have 744 (12 x 62) rows of data.

Figure 15. Unemployment Data, After Reshaping Using the Stack Functionality



From this newly shaped data, you can now create a visualization of unemployment rate over time. In the JMP Graph Builder platform, drag and drop columns into the X and Y roles to see the rise and fall in unemployment rate over time. Note that the Month and Year values are nested in the x-axis so you did not have to create a merged column that contained both month and year data.

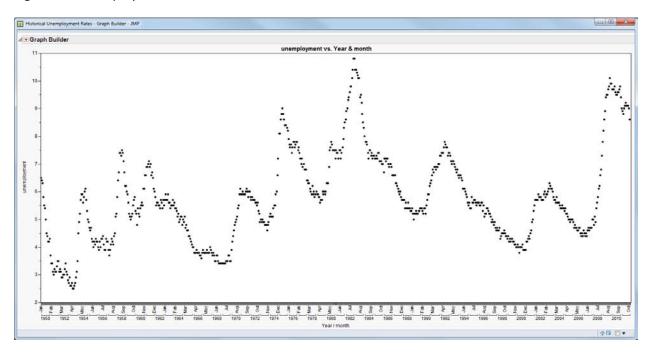


Figure 16. Unemployment Rates over Time

The graph shows the rise and fall of unemployment rates over a 60-year span of time. You can clearly see the rapid rise in unemployment rates at the end of 2008. It is interesting to note that the period of highest unemployment rates occurred in the early 1980s. It is also heartening to see that unemployment appears to be decreasing in recent months.

One more note about this data. If you are an experienced SAS user, you might have already known how to use the DATA step to reshape the data. In this example you used Stack functionality that is available in JMP on the data shown in Figure 13. An alternative would be to use the JMP ability to integrate directly with SAS to create a SAS data set instead. Save the data table from Figure 13 as a SAS data set (for example, C:\SASGlobalForum2011\unemployment.sas7bdat). If you have access to SAS from your machine, you can open the SAS Program Editor in JMP and submit the code to SAS. This code shown below would create a data set named new.sas7bdat that looks identical to the data in Figure 15.

Figure 17. SAS Program Editor within JMP

```
StackviaSAS - JMP

File Edit Tables Analyze Graph Tools Add-Ins View Window Help

1 libname sgf 'C:\SASGlobalForum2011';
2 proc transpose data=sgf.unemployment out=sgf.new(RENAME=(Col1=UnempRate));
3 var jan feb mar apr may jun jul aug sep oct nov dec;
4 by year;
5 run;
6
```

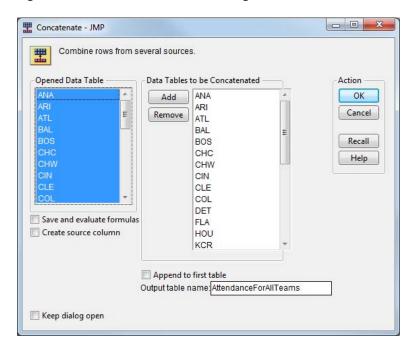
EXAMPLE 3 BASEBALL ATTENDANCE

Many major cities have baseball teams, and some cities have seen teams come and go, and come back again. Suppose you want to analyze attendance per city, to see if you can determine if fans in certain cities are more loyal than others. If you want to analyze data on a per-city basis, you will need information about the different ballparks within each city as well as attendance for each team. Also, you might want to consider the capacity of each ballpark, since sizes vary widely; for example, the Dodgers' stadium seats 56,000 people while Oakland's stadium only seats 37,000.

Attendance data is available online at Baseball-Reference.com⁶ for each team. Suppose you saved attendance data for each team into a JMP data table. To more easily work with the data, you should use the Concatenate functionality in JMP to combine all the attendance information into one big table. Use **Tables->Concatenate** to invoke the dialog box.

⁶ Data is accessible from http://www.baseball-reference.com.

Figure 18. The JMP Concatenate Dialog Box



The capacity of each ballpark is available online at BallParksofBaseball.com⁷ and can be imported into JMP using Internet Open or copy-and-paste. Shown below are the two tables you just created – AttendanceForAllTeams, which contains average attendance information for each team and each year, including which ballpark the team played in each year; and BallParkCapacity, which contains a single row for past and present ballparks.

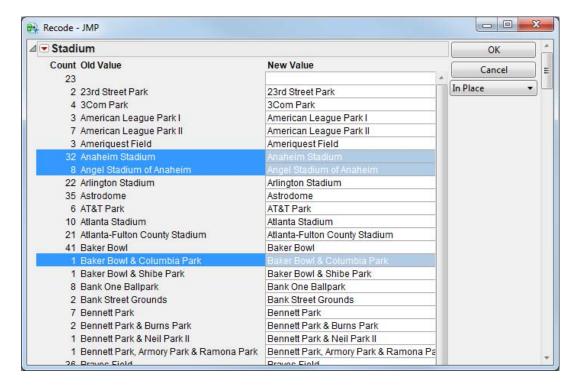
⁷ Data can be accessed from http://www.ballparksofbaseball.com.

AttendanceForAllTeams - JMP File Edit Tables Rows Cols AttendanceForA... Playoffs Lq Source 2011 Los Angeles Angels of Anaheim AL West 76 Angel Stadium of Anaheim 86 2010 Los Angeles Angels of Anaheim AL West 80 82 Angel Stadium of Anaheim **Columns** (15/0) 1 Lost ALCS (4-2) Angel Stadium of Anaheim 2009 Los Angeles Angels of Anaheim AL West 97 65 ▲ Rk Year
Tm
Lg 1 Lost LDS (3-1) 2008 Los Angeles Angels of Anaheim AL West 100 62 Angel Stadium of Anaheim 2007 Los Angeles Angels of Anaheim AL West 94 68 1 Lost LDS (3-0) Angel Stadium of Anaheim 2006 Los Angeles Angels of Anaheim AL West 89 73 Angel Stadium of Anaheim 95 1 Lost ALCS (4-1) Angel Stadium of Anaheim 2005 Los Angeles Angels of Anaheim AL West 67 ▲ L ▲ Finish 1 Lost LDS (3-0) Angel Stadium of Anaheim 2004 Anaheim Angels AL West 92 70 2003 Anaheim Angels AL West 77 ♣ Playoffs 10 2002 Anaheim ♣ Stadium 2001 Anaheim BallParkCapacity - JMP _ D X 11 11 12 2000 Anaheim Rows Cols Analyze Graph Tools Add-Ins View 12 File Edit Tables Rows 13 13 1999 Anaheim All rows 2.446 14 14 1998 Anaheim 15 15 1997 Anaheim Excluded 0 City Teams Ballpark Capacity Hidden 1 Arlington Texas Rangers Arlington Stadium 43,521 2 Arlington Texas Rangers Rangers Ballpark in Arlington 49,166 3 Atlanta Atlanta Braves Atlanta-Fulton County Stadium 52.013 4 Atlanta Turner Field 49,381 Atlanta Braves 5 Baltimore Oriole Park at Camden Yards 45,971 Baltimore Orioles Columns (4/0) 6 Baltimore Baltimore Orioles Memorial Stadium 54,000 ♣ Teams 7 Baltimore Baltimore Orioles Oriole Park 8.000 Ballpark 8 Boston Boston Braves Braves Field 42.000 Capacity 9 Boston Boston Braves South End Grounds 5.000 10 Boston Boston Red Sox Fenway Park 39.928 11 Boston Boston Red Sox Huntington Ave. Grounds 11.500 32,000 12 Brooklyn Brooklyn Dodgers Ebbets Field Rows 13 Chicago Chicago Cubs West Side Park 16,000 All rows 14 Chicago Chicago Cubs Wrigley Field 41,118 Selected 15 Chicago Chicago White Sox Comiskey Park I 52,000 Excluded Chicago White Sox South Side Park 15,000 16 Chicago **↑** □ ▼

Figure 19. Tables Showing Concatenated Attendance for All Teams and the Capacity for Each Ballpark

Your next task is to join the two tables so that capacity information for each stadium can be used. The matching columns you will use are the Stadium/Ballpark columns. It is important to examine the values for these two matching columns. Use **Cols->Recode** to look at the values and make judgments about the values of the names. For example, the *Anaheim Stadium* and the *Angel Stadium of Anaheim* are actually the same, but stadiums go through name changes frequently. Recode allows you to modify all values as appropriate. Another example, *Baker Bowl & Columbia Park*, shows that in one year, a team played in two different stadiums. For these cases you would have to decide how to handle the row. The Recode step is necessary as part of data cleaning in any software application and can take time. JMP helps minimize the amount of time spent here.

Figure 20. The JMP Recode Dialog Box



Once you have determined the correct values for stadium names, your next task is to join the two tables. Invoke the Join functionality that is available in JMP using **Tables->Join**. In the dialog box, match Stadiums=Ballparks. You may also select which columns you would like to keep in the resulting table.

Figure 21. The JMP Join Dialog Box



The resulting table contains 2,527 rows, but a few hundred rows contain some missing data because of mismatching in the matching columns we selected. This is primarily due to very old ballparks listed in the attendance data sets (some from the 1800s!) that did not have a corresponding entry in the capacity table.

2945228 36361

2927399 36141

(- 0 - K Attendance-Capacity-Joined - JMP File Edit Tables Bows Cols Analyze Graph Tools Add-Ins View Window Help PERMETER BER 四日日日 - B ▼ Attendance-Capa Attendance Attend/G 665 Arlington 666 Arlington 667 Arlington 43,521 43,521 Arlington Stadium Arlington Stadium Texas Rangers 1973 Texas Rangers Arlington Stadium Arlington Stadium Arlington Stadium Arlington Stadium 84 Arlington Stadium Texas Rangers 43,521 1974 Texas Rangers 1193902 14924 668 Arlington 669 Arlington Texas Rangers Texas Rangers 1975 Texas Rangers 1976 Texas Rangers 79 Arlington Stadium 76 Arlington Stadium 1127924 1164982 43 521 1976 Texas Rangers 1977 Texas Rangers 670 Arlingtor Texas Rangers 94 Arlington Stadium 15441 Columns (11/0) Artington Stadium 43,521 1250722 671 Arlington Texas Rangers Arlington Stadium 43,521 1978 Texas Rangers 87 Arlington Stadium 1447963 17658 43,521 43,521 43,521 83 Arlington Stadium 76 Arlington Stadium 57 Arlington Stadium Arlington Stadium Arlington Stadium 1519671 1980 Texas Rangers 1981 Texas Rangers 14977 15180 Texas Rangers Arlington Stadium 675 Arlington 676 Arlington 677 Arlington 678 Arlington Texas Rangers Texas Rangers Texas Rangers Arlington Stadium Arlington Stadium Arlington Stadium 1982 Texas Rangers 1983 Texas Rangers 1984 Texas Rangers 1985 Texas Rangers 64 Arlington Stadium 77 Arlington Stadium 69 Arlington Stadium 43.521 1154432 1363469 1102471 1112497 43 521 16833 43,521 43,521 43,521 62 Arlington Stadium Texas Rangers Arlington Stadium 13906 679 Arlington Texas Rangers Arlington Stadium 43.521 1986 Texas Rangers 87 Arlington Stadium 1692002 20889 680 Arlington 681 Arlington 682 Arlington Texas Rangers Texas Rangers Arlington Stadium Arlington Stadium 1987 Texas Rangers 1988 Texas Rangers 75 Arlington Stadium 70 Arlington Stadium 21766 19530 1763053 1581901 Texas Rangers Arlington Stadium 43,521 1989 Texas Rangers 83 Arlington Stadium 25234 Arlington Stadium Arlington Stadium Arlington Stadium 83 Arlington Stadium 85 Arlington Stadium 77 Arlington Stadium 683 Arlington Texas Rangers 43,521 1990 Texas Rangers 2057911 684 Arlington 685 Arlington 686 Arlington Texas Rangers Texas Rangers 43,521 28367 43,521 27139 27711 Rows
All rows
Selected
Excluded
Hidden
Labelled 2244616 Texas Rangers Arlington Stadium 1993 Texas Rangers 86 Arlington Stadium 2,527 687 Arlington Texas Rangers The Ballpark in Arlington 49 166 1994 Texas Rangers 52 The Ballpark in Arlington 2503198 39733 688 Arlington 689 Arlington 1995 Texas Rangers 1996 Texas Rangers 1985910 2889020

Figure 22. Resulting JMP Data Table after Attendance and Capacity Tables Have Been Joined

You might want to measure attendance by how full the stadium is. Create a new column named Percent of Capacity with a formula of Attend/G divided by Capacity.

49,166

49,166

1997 Texas Rangers

1998 Texas Rangers

77 The Ballpark in Arlington

88 The Ballpark in Arlington

95 The Ballpark in Arlington 71 The Ballpark in Arlington

Figure 23. Example of the Column Info Dialog Box with a Formula

The Ballpark in Arlington

The Ballpark in Arlington

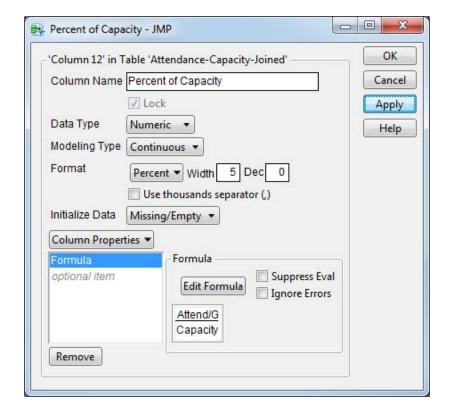
690 Arlington

691 Arlington

692 Arlington 693 Arlington

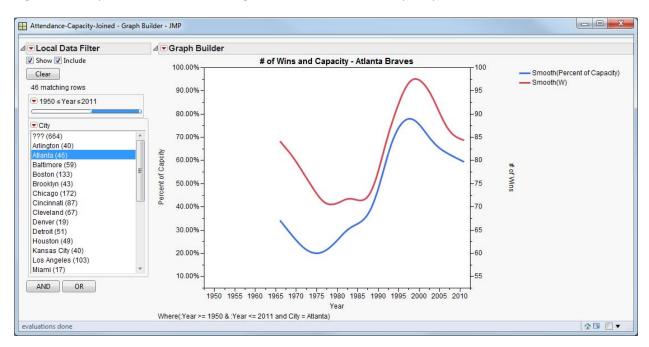
Texas Rangers

Texas Rangers



Using the JMP **Graph Builder** platform, create a graph showing # of Wins and Percent of Capacity over the years. You can use the **Local Data Filter** to limit the years to 1970 or later, since data from those years contain at least 24 teams per year and is thus more interesting. Also use **Local Data Filter** to view each city by itself to see how the Percent of Capacity is affected by the # of Wins. Below are two examples. The Atlanta Braves franchise started in 1966. Its curves follow each other closely, indicating that average attendance in Atlanta is strongly correlated with how well the team is doing. Conversely, Boston Red Sox fans continue to grow more loyal over time, even when the Red Sox hit a rough streak in the 1980s and early 1990s. This might imply that Red Sox fans are more loyal. You could investigate this idea further by using some of the more analytical platforms in JMP.

Figure 24. Graph Builder Chart Showing # of Wins and Percent Capacity over Time for Atlanta



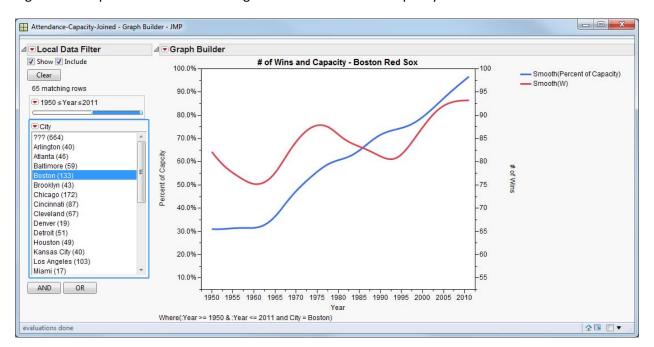


Figure 25. Graph Builder Chart Showing # of Wins and Percent Capacity over Time for Boston

CONCLUSION

JMP 10 continues the tradition of offering varied and easy-to-use ways of accessing and reshaping data. Knowing the ways in which JMP can help analyze data will allow you to quickly take raw data from any source, reshape, and customize it in order to get to the "good part" of data analysis. JMP will help you uncover the answers within your data.

CONTACT INFORMATION

Your comments and questions are valued and encouraged. Contact the author:

Audrey Ventura SAS Campus Drive SAS Institute Inc.

E-mail: Audrey. Ventura@jmp.com

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.