

# How to use the SAS® Software to Evaluate Screening Tests Using Predictive Values in Conjunction with ROC Curves

Richard Severino, The Queen’s Medical Center, Honolulu, HI

## ABSTRACT

Development of new screening tests is ongoing in the medical and biotechnology fields, be it for a medical condition or in an attempt to predict outcome. The screen usually involves a measurement that is then compared to a cutoff point in order to classify the test subject into one of two categories. In developing screening tools, researchers often concentrate on the sensitivity and specificity of the test, and use ROC curves to evaluate the discriminating power of the screening test often ignoring the predictive values of the test. The predictive values of the screen are equally important as they measure the accuracy of the prediction made on a subject whose true condition is unknown. By examining the predictive values, a cutoff point may be chosen such that it minimizes the false positive and false negative rates. A method for obtaining and examining these values using base SAS and SAS/STAT is presented.

## INTRODUCTION

A screening or diagnostic test is designed to predict the presence or absence of a condition, or even predict one of two possible outcomes. If a person is tested for a given condition and the test is positive, what is the probability that the person really has the condition? On the other hand, if the test is negative, what is the probability that the person is really free of the condition?

Plotting an ROC curve has become a very popular tool for evaluating the accuracy of such tests and prediction models (Stead and MacDonald, 1997). The ROC curve gives an indication of how well the test performs when classifying a person or test sample whose condition is known. However, an ROC curve does not give any indication of how accurate the test is when classifying a person or test sample whose condition is not known.

The question of practical importance is: “Given the test result is positive, what is the probability that the subject really has the condition?” The answer to this question is provided by the predictive values of the test.

## SOME PROPERTIES OF SCREENING TESTS

The screening or diagnostic test is one which results in classifying the subject of the test into one of two categories. For example, a home pregnancy test will result in concluding that the woman who took the test is either pregnant or is not. If the test indicates that the condition being screened for exists then the test is said to be positive, otherwise it is said to be negative.

Table 1 shows a cross classification of the results of a screening test and the known condition of the subjects tested. (If you are not certain of what kind of data is required to construct a table such as this, Stead and MacDonald (1997) provide a good example and explanation.) In Table 1 There are  $a + c$  subjects

who actually have the condition being screened for. Of these  $a+c$  subjects,  $a$  of them tested positive for the condition.

**Table 1. Cross Tabulation of Test Results with Actual Condition**

		Condition		Totals (Test)
		Present	Absent	
Test	Positive	$a$	$b$	$a+b$
	Negative	$c$	$d$	$c+d$
Totals (Condition)		$a+c$	$b+d$	$N$

A screening test has four important properties listed in Table 2 along with their definitions. The third column in Table 2 shows how to compute the value of each of the properties directly from Table 1. Although they are often reported as percentages, these properties are probabilities and therefore their values range from 0 to 1.

**Table 2. Properties of a Screening Test**

Property	Definition	Computation from Table 1.
Sensitivity	Probability that a subject known to have the condition will test positive	$\frac{a}{a+c}$
Specificity	Probability that a subject who is known to be free of the condition will test negative	$\frac{d}{b+d}$
Predictive Value Positive	Probability that a subject who tests positive does have the condition	$\frac{a}{a+b}$
Predictive Value Negative	Probability that a subject who tests negative is free of the condition	$\frac{d}{c+d}$

A screening test is usually performed by comparing an observed value or characteristic to a reference or cutoff point. If the observed value is above the cutoff point the test is positive ( or negative, depending on the test). A table such as Table 1 can be formed for different cutoff points. The probabilities defined in Table 2 can then be computed for each cutoff point considered . The cutoff point which yields the best combination of values for the properties in Table 2 is then chosen as the reference point to use when conducting the test.

## THE ROC CURVE

Stead and MacDonald (1997) refer to the sensitivity as the true positive rate (TP). The false positive rate (FP) which they define is none other than the quantity (1- specificity). They use percentages rather than proportions. They describe how to obtain an ROC curve which is a plot of the sensitivity against (1-specificity). The ROC curve gives a graphical representation of how well the test performs with respect to sensitivity and specificity. They also present a method to facilitate choosing the cutoff point which corresponds to the desired point on the ROC curve. Fletcher, Fletcher and Wagner (1988) provide a more detailed discussion of ROC curves.

A test with high sensitivity and specificity is desirable. Unfortunately, in practice, you usually have to sacrifice one for the other (Fletcher, Fletcher and Wagner, 1988). Generally, the faster the ROC curve rises from the origin and then curves to right the more accurate the test is considered. The area under the curve is one measure of the overall accuracy or predictive power of the test (SAS Institute Inc., 1995).

**PREDICTIVE VALUES**

While high sensitivity and specificity are good starting points for a test, they tell you nothing of the accuracy of the test when it is applied to a subject whose condition is unknown. This is where an ROC curve is no longer useful.

When you apply the test to a subject whose condition is unknown and the test is positive, then you must ask yourself “what is the probability that this subject really has the condition?” The answer to this question is given by the predictive value positive. Similarly, if the test is negative then the predictive value negative will provide you with the probability that the subject is truly free of the condition given the test was negative.

All four properties in Table 2 can be written as conditional probabilities (Fleiss, 1981). The predictive value positive PVP and the predictive value negative (PVN) can be written as follows (Rosner, 1986):

$$PVP = \frac{PREV * SENSITIVITY}{PREV * SENSITIVITY + (1 - PREV) * (1 - SPECIFICITY)} \quad (1)$$

$$PVN = \frac{(1 - PREV) * SPECIFICITY}{(1 - PREV) * SPECIFICITY + PREV * (1 - SENSITIVITY)} \quad (2)$$

where *PREV* is the prevalence of the condition in question. In Table 1, the prevalence is  $(a+c)/N$ .

You could think of the quantity PVP as being the true positive rate and the quantity (1-PVN) as being the false negative rate.

So, not only do the predictive values differ by definition from the sensitivity and specificity of a test, but they also take into account the prevalence of the condition which is being tested while the sensitivity and specificity do not. It is worth noting at this point that the cutoff point which yields the best combination of sensitivity and specificity for a test will not necessarily yield the best combination of predictive values.

Because the PVP and PVN take into account the prevalence of the condition, you can evaluate the predictive values of the test for different values of prevalence. This is especially important if the prevalence of the condition in the sample that is used to evaluate a test is much higher than the prevalence of the condition in the general population. In the next section, you will see that depending on the value of the prevalence, there can be quite a difference in the predicted values while the sensitivity and specificity remain unchanged.

**USING SAS TO OBTAIN PREDICTIVE VALUES**

**METHOD WITH EXAMPLE**

Obtaining predictive values is not difficult, but rather tedious. Stead and MacDonald (1997) presented an approach to obtaining an ROC curve using only DATA steps, PROC SORT, PROC FREQ, and PROC GPLOT of the SAS system. With some modifications and the substitution of PROC PLOT for PROC GPLOT, the same approach is used to obtain and plot the predictive values of the test in addition to the sensitivity and specificity. The SAS Macro in appendix A (Macro A1) will accomplish this task creating a SAS data set such as the one in Table 3 and the plots in Figures 1 and 2.

**Table 3. Data Set Needed to Plot an ROC Curve and the Predictive Values (not all variables in the data set are shown).**

CUTOFF	SENSITIVITY <sup>†</sup>	1-SPECIFICITY <sup>†</sup>	PVP <sup>††</sup>	PVN <sup>††</sup>
0.50	100	100	0.50000	0.00000
0.55	100	92	0.52083	1.00000
0.60	94	86	0.52222	0.70000
0.65	94	78	0.54651	0.78571
0.70	92	72	0.56098	0.77778
0.75	84	68	0.55263	0.66667
0.80	76	62	0.55072	0.61290
0.85	66	60	0.52381	0.54054
0.90	58	60	0.49153	0.48780
0.95	52	52	0.50000	0.50000
1.00	48	50	0.48980	0.49020
1.05	42	42	0.50000	0.50000
1.10	42	38	0.52500	0.51667
1.15	36	34	0.51429	0.50769
1.20	24	34	0.41379	0.46479
1.25	16	28	0.36364	0.46154
1.30	14	22	0.38889	0.47561
1.35	12	20	0.37500	0.47619
1.40	8	16	0.33333	0.47727
1.45	4	8	0.33333	0.48936
1.50	4	0	1.00000	0.51020

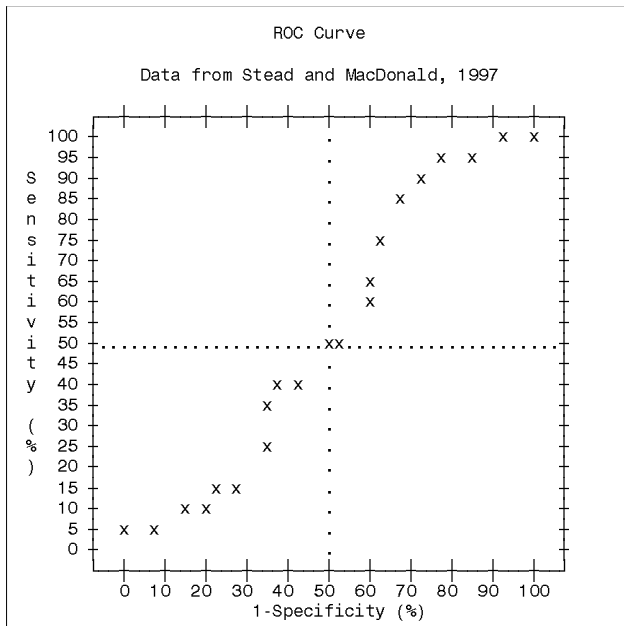
<sup>†</sup> represented as a percentage

<sup>††</sup> represented as a probability (proportion)

Once the sensitivity, specificity and initial predictive values have been obtained, the procedure to obtain predictive values for a range of prevalence values is not very different. Macro A2 (the second SAS Macro in Appendix A) will take the SAS data set in Table 3 as input along with the upper and lower values of the prevalence range and the increment size. This macro also takes as input the upper limit, lower limit and increment size for the cutoff values which will be used in PROC PLOT. The macro also allows you to specify the name of an ASCII output file to which it will write all the data generated so that it can be used by other

software if necessary.

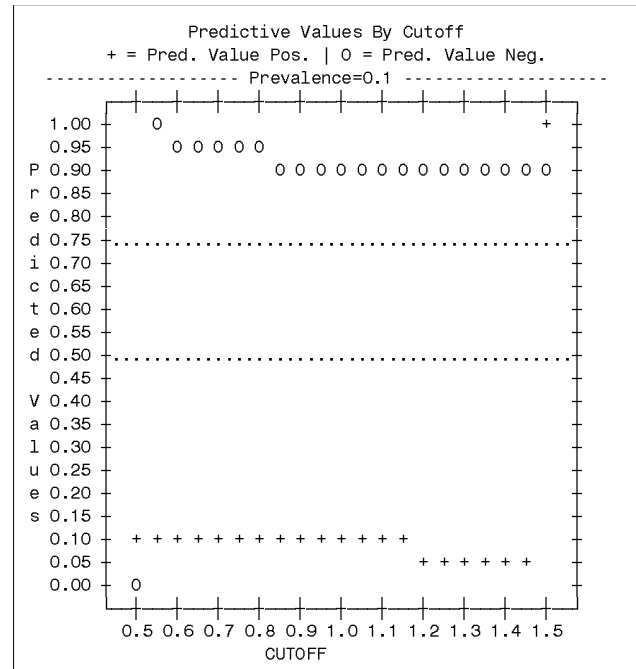
**Figure 1. ROC Curve**



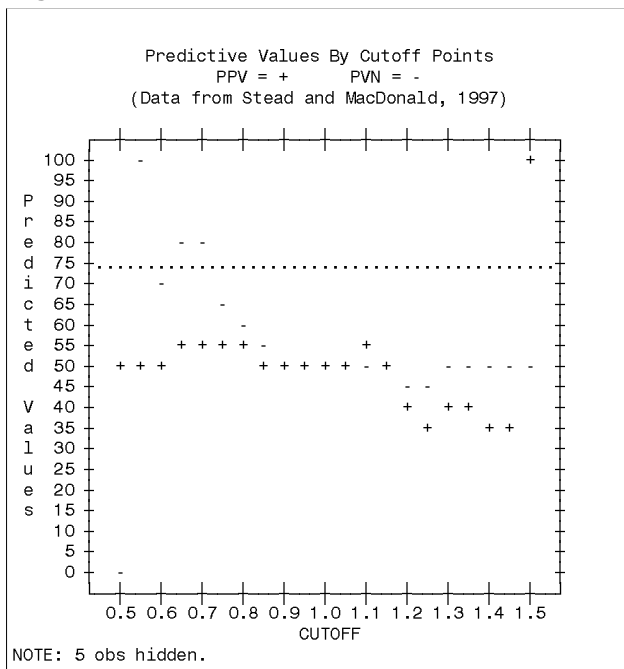
specified as well as the range.

For example, the sample data set given by Stead and Macdonald (1997), consists of 50 cases that are in fact positive, and 50 cases that are negative. So the prevalence in the sample is 0.5 (50%). If the prevalence of the condition being screened in the population is 0.1 (or 10%) then, as can be seen by comparing Figures 2 and 3, the predicted values are quite different.

**Figure 3. Predictive Values when Prevalence=0.1**



**Figure 2. Predictive Values**



In order to evaluate the predictive values for different levels of prevalence, a DO LOOP is used to cycle through a range of values for prevalence. For each distinct prevalence, the predicted values are computed for the cutoff points previously determined. Here the predicted values are computed using equations (1) and (2). The range of values for prevalence is specified by setting the lowest and highest values and the size of the incremental change desired. If you only want to specify one prevalence value, then specify the same value for the highest and lowest values. The number of distinct prevalences used depends on the increment

If the prevalence in the population is close to 0.10, and even if only a 50% predicted value positive is required, then it is clear from Figure 3 that none of the cutoffs considered will result in an acceptable test. On the other hand, if the predicted value negative is critical and we can disregard the predicted value positive, then the cutoff points from .6 to .9 may result in a satisfactory test with PVN=95%. Note that while the cutoff point of 0.5 corresponds to a PVN=100%, this is because it happens to be the lowest observed value of test. The predicted values at the lowest and highest cutoff points must be examined carefully if these cutoff points happen to coincide with the minimum and maximum observed values of the test.

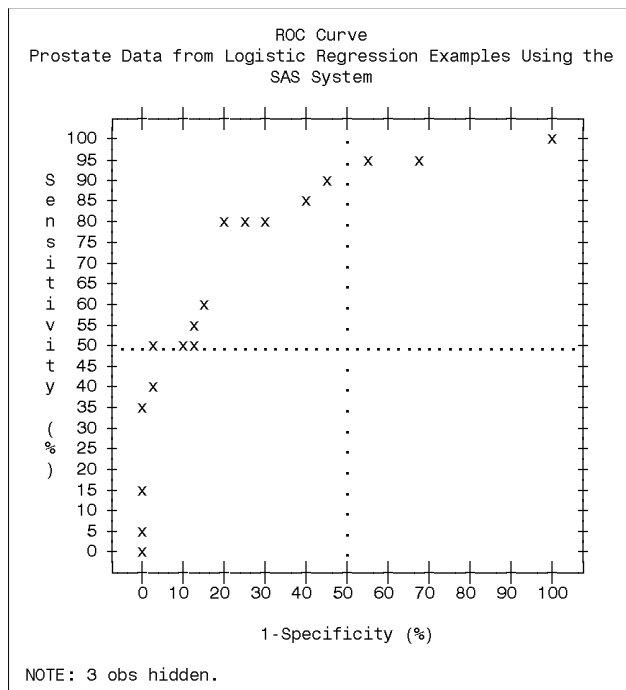
**ANOTHER EXAMPLE**

The data for this example is the Prostate data set from the first edition of Logistic Regression Examples Using the SAS System (1995). In this example, a logistic regression model was used to predict the probability of nodal involvement based on three predictor variables. Using PROC LOGISTIC, the predicted probabilities are saved with the original data to a new data set. The screening test then consists of predicting nodal involvement based on the predicted probability obtained from the logistic regression. For example, if the cutoff point is 0.45 then we predict that all subjects with predicted probability greater than 0.45 have nodal involvement, i.e. the test is positive.

It is worth mentioning here that if you are using a the predicted probabilities from a logistic regression model as your test measurement, then there are at least two options in PROC LOGISTIC which would automatically generate the sensitivity, specificity and predictive values, allowing you to generate an ROC curve as well as the type of plots shown here. In fact there are some advantages to using those options instead of the procedures presented here, but then there are some disadvantages. For a good discussion of those options, I suggest you read Logistic Regression Examples Using the SAS System (1995).

The ROC curve in Figure 4 is obtained using macro A1. This ROC curve indicates that the test is quite accurate with respect to sensitivity and specificity. If plotted using PROC Gplot, it would be possible to label each point with a symbol corresponding to the cutoff point. You can however plot the sensitivity and specificity against the cutoff points as was done with the predictive values in Figure 2. That will make it easy to see which cutoff point yields the best combination of sensitivity and specificity for the test under consideration.

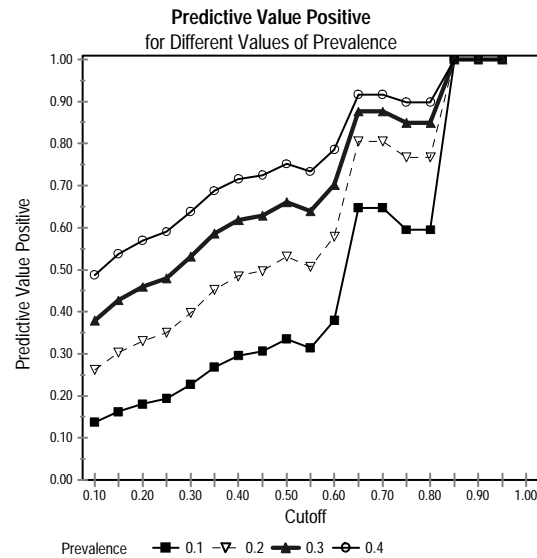
**Figure 4. ROC curve for the Prostate Data.**



The predictive values are obtained for a range of prevalence values using macro A2. Figure 5 shows the predictive values positive plotted against the cutoff points for various values of prevalence. This plot was created with Corel Quattro Pro version 7 using the ascii data written out by macro A2. Of course those of you lucky enough to have access to SAS/GRAPH software won't need to use anything else.

Notice how much the predictive value positive changes as the prevalence changes. If the prevalence in the population is 0.1, then the predictive value positive is not that great until the cutoff point is so high that the test won't be positive for any subjects other than extreme cases. Of course you can not make any conclusions until you also examine the predictive value negative information. A similar plot can be obtained for the predictive values negative.

**Figure 5. Predictive Values for the Prostate Data.**



You must decide which of the predictive values is more critical. This is usually done by examining the effects of misclassification. What would happen if you predict that a patient has nodal involvement? If the test you are using has a predictive value positive of 0.70, then there 30% of the time when you predict nodal involvement, you will be wrong and the patient will undergo some procedure that is not necessary. If the procedure has no risk, then it doesn't really matter (other than economically), but if the procedure is risky you don't want to put someone through it unnecessarily.

**CONCLUSION**

While ROC curves provide an initial evaluation of the discriminatory power of a screening test, the predictive values are more reflective of how accurate the test will be when applied to subjects whose condition is unknown. The predictive values depend on the prevalence of the condition while the sensitivities and specificities do not. By examining the predictive values with different values of prevalence, you can more easily find the cutoff point which will minimize the false negative and false positive rates of the test.

SAS® and SAS/GRAPH are registered trademarks of the SAS Institute Inc. in the USA and other countries. ® indicates USA registration.

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

**REFERENCES**

Fleiss, Joseph L., *Statistical Methods for Rates and Proportions, Second Edition*, John Wiley and Sons, 1981.  
 Fletcher, Robert H., Fletcher, Suzanne W. and Wagner, Edward

H., *Clinical Epidemiology, the Essentials, Second Edition*, Williams and Wilkins, 1988.

Rosner, Bernard, *Fundamentals of Biostatistics*, PWS Publishers, 1986.

SAS Institute Inc., *Logistic Regression Examples Using the SAS System, Version 6, First Edition*, Cary, NC: SAS Institute Inc., 1995.

Stead, Andrew G., and MacDonald, Karen G., *Constructing ROC Curves with the SAS System*, Proceedings of the 22<sup>nd</sup> Annual SAS Users' Group Conference, March, 1997.

APPENDIX A

Macro A1.

```

/* *****
The following macro
    ROC_PV(DATAIN,LOWLIM,UPLIM,INC)
will read a SAS data set which contains at
least the following:
    TESTVAL a quantitative measure thought to
            be predictive of the existence of
            a given condition
    CNDITION a variable indicating the presence
            or absence of the condition
            associated with TESTVAL
The macro will then compute sensitivities,
specificities and Predictive Values (Positive
and Negative).
The macro will then create a plot of the ROC
curve as well plots of the predictive values by
cutoff points.
The macro will also print the data.
DATAIN   Input SAS data set
LOWLIM   Lowest cutoff point
UPLIM   Highest cutoff point
INC      Increment for range of cutoff values
NINC     the number of cutoff points (based
on LOWLIM, UPLIM, and INC)
I        Loop Index
CUTOFF   Cutoff for TEST
CNDITION Actual Observed Condition
TEST     Predicted condition based on test
cutoff value
-----
This Macro is based on the Macro written by
A. Stead and appearing in "Constructing ROC
Curves with the SAS System" By Andrew G. Stead
and Karen G. MacDonald, published in the 1997
SUGI 22 Proceedings
Modified and Extended By
                                Richard Severino 06/97
Revised: 08/97
***** */
%MACRO ROC_PV(DATAIN,LOWLIM,UPLIM,INC);
OPTIONS MTRACE MPRINT CENTER PAGENO=1 LINESIZE=120;
DATA ROC;
SET &DATAIN;

```

```

LOWLIM=&LOWLIM; UPLIM=&UPLIM; INC=&INC;
NINC=(UPLIM-LOWLIM)/INC ;
DO I=1 TO NINC+1 ;
    CUTOFF=LOWLIM + (I-1)*(INC);
    IF TESTVAL>CUTOFF THEN TEST=1; ELSE TEST=0;
    OUTPUT;
END;
DROP I;
RUN;

PROC SORT; BY CUTOFF;
RUN;

PROC FREQ; BY CUTOFF;
TABLE TEST*CNDITION / OUT=PCTS1 OUTPCT NOPRINT;
RUN;
/* ----- */

DATA SENSIT; SET PCTS1; /* GET SENSITIVITY */
IF CNDITION=1 AND TEST=1 ;
SENS_PCT = PCT_COL;
SENS = PCT_COL/100 ;
DROP PCT_COL PCT_ROW;

LABEL SENS_PCT="Sensitivity (%)"
SENS="Sensitivity";

RUN;
/* ----- */

DATA SPECIF; SET PCTS1; /* GET SPECIFICITY */
IF CNDITION=0 AND TEST=0 ;
SPEC_PCT = PCT_COL;
SPEC = PCT_COL/100;
DROP PCT_COL PCT_ROW;

LABEL SPEC_PCT = "Specificity (%)"
SPEC = "Specificity";

RUN;
/* ----- */

DATA OMSPEC; SET PCTS1; /* GET 1-SPECIFICITY */
IF CNDITION=0 AND TEST=1 ;
_1MSPEC = PCT_COL;
DROP PCT_COL PCT_ROW;

LABEL _1MSPEC="1-Specificity (%)" ;

RUN;
/* ----- */

DATA PVPOS; SET PCTS1; /* GET PRED. VAL. POSITIVE */
IF TEST=1 AND CNDITION=1 ;
PVP = PCT_ROW;
DROP PCT_COL PCT_ROW;

LABEL PVP="Pred. Value Positive (%)" ;

RUN;
/* ----- */

DATA PVNEG; SET PCTS1; /* GET PRED. VAL. NEGATIVE */
IF TEST=0 AND CNDITION=0;
PVN = PCT_ROW;
DROP PCT_COL PCT_ROW;

LABEL PVN="Pred. Value Negative (%)" ;

RUN;

/* ----- *
| MERGE:
| SENSITIVITY, SPECIFICITY AND 1-SPECIFICITY
| INTO ONE DATA SET
* ----- */

DATA ROC2;
MERGE SENSIT SPECIF OMSPEC; BY CUTOFF;

IF SENS = . THEN SENS = 0.0;
IF SENS_PCT = . THEN SENS_PCT = 0.0;
IF SPEC = . THEN SPEC = 0.0;
IF _1MSPEC = . THEN _1MSPEC = 0.0;

RUN;

PROC PRINT DATA=ROC2 ;
TITLE 'ROC_PV (ROC2 DATA SET) DATA';
RUN;
/* ----- */
/* ----- *
| MERGE:
| PREDICTIVE VALUES POS. AND NEG.
| INTO ONE DATA SET
* ----- */

DATA PV1; MERGE PVPOS PVNEG; BY CUTOFF;
IF PVP = . THEN PVP = 0.0;
IF PVN = . THEN PVN = 0.0;

RUN;

PROC PRINT DATA=PV1;

```

```
TITLE 'PREDICTIVE VALUES (PV DATA)';
RUN;
/* ----- */
/* ***** */
|                CREATE PLOTS                |
* ***** */

PROC PLOT DATA=ROC2 NOLEGEND;
  TITLE 'ROC CURVE';
  PLOT SENS_PCT*1MSPEC = 'X'
      / VAXIS=0 TO 100 BY 5 VREF=50 VREFCHAR='.'
        HAXIS=0 TO 100 BY 5 HREF=50 HREFCHAR='.'
        BOX HEXPAND; ;
RUN;

PROC PLOT DATA=PV1 VPERCENT=50 NOLEGEND;
  TITLE 'PREDICTIVE VALUES BY CUTOFF POINTS';
  PLOT PVP*CUTOFF='+' PVN*CUTOFF='0'
      / VAXIS= 0 TO 100 BY 5 VREF=75 VREFCHAR='.'
        HAXIS= &LOWLIM TO &UPLIM BY &INC
        BOX HEXPAND;
RUN;
/* ***** */
%MEND;
```

Macro A2.

```
/* ***** */
  The following macro
  PV( DATAIN, PREVLO, PREVHI, INC,
      CUTLO, CUTHI, CUTINC, PV_OUT )
  will read a SAS data set which contains
  sensitivities and specificities the test. The
  macro will then compute Predictive Values (Pos-
  itive and Negative) for a range of Prevalence
  values.

  The macro will then create plots of predictive
  values and will write an ASCII data file so that
  the data can be plotted using other software if
  necessary.

  DATAIN   Input SAS data set
  PREVLO    Lowest Prevalence Value
  PREVHI    Highest Prevalence Value
  INC       Increment for range of Prevalence
            values
  CUTLO     Lowest cutoff point to be plotted
  CUTHI     Highest cutoff point to be plotted
  CUTINC    Increment for plotting the range of
            cutoff points
  PV_OUT    path and filename for an ASCII file
            which will be written containing Pre-
            valence, Cutoff, Sensitivity,
            Specificity and Predictive values
  -----
  Written By Richard Severino                06/97
  Revised: 08/97
  -----
/* ***** */

%MACRO PV( DATAIN, PREVLO, PREVHI, INC,
          CUTLO, CUTHI, CUTINC, PV_OUT );

OPTIONS MTRACE MPRINT CENTER PAGENO=1 LINESIZE=120;

FILENAME PV2_OUT &PV_OUT ;

DATA PV2;
  SET &DATAIN;
  PREVLO=&PREVLO;
  PREVHI=&PREVHI;
  INC2=&INC;

  NINC2 = ( PREVHI - PREVLO )/INC2 ;

DO I=1 TO NINC2+1 ;
  PREV = PREVLO + (I-1)*INC2 ;

IF SENS > 0.0 THEN
  PVP = (PREV*SENS)/( (PREV*SENS)+(1-PREV)*(1-SPEC) ) ;
```

```
ELSE
  PVP = 0.0 ;

IF SPEC > 0.0 THEN
  PVN=((1-PREV)*SPEC)/(((1-PREV)*SPEC)+(PREV*(1-SENS)));
ELSE
  PVN = 0.0 ;

  OUTPUT;
END;
DROP I;

IF PVP = . THEN PVP = 0.0 ;

LABEL PREV='Prevalence' ;
RUN;

PROC SORT;
  BY PREV CUTOFF;
RUN;

/* ----- */
: WRITE SENSITIVITY, SPECIFICITY :
: PREDICTIVE VALUES, CUTOFF AND :
: PREVALENCE DATA TO AN ASCII FILE :
* ----- */

DATA TEMP; SET PV2;
FILE PV2_OUT ;
PUT @2 PREV 4.2 @8 CUTOFF 6.4
    @15 (SENS SPEC PVP PVN) (6.4 7.4 7.4 7.4) ;
RUN;
/* ----- */

PROC PRINT DATA=PV2;
RUN;

/* ===== */
| PLOT THE PREDICTIVE VALUES |
* ===== */

PROC PLOT DATA=PV2 NOLEGEND; BY PREV;
  TITLE1 'PREDICTIVE VALUES BY CUTOFF';
  TITLE2 ' ';
  TITLE3 '+ = PRED. VALUE POS. | 0 = PRED. VALUE NEG.';

  PLOT PVP*CUTOFF = '+' PVN*CUTOFF = '0'
      / VAXIS = 0 TO 1.00 BY .05 VREF=.50 .75 VREFCHAR='.'
        HAXIS = &CUTLO TO &CUTHI BY &CUTINC
        BOX OVERLAY HEXPAND;
RUN;
/* ===== */

%MEND ;
```

**AUTHOR CONTACT INFORMATION**

Richard Severino  
 The Queen's Medical Center  
 Research Planning and Development  
 1301 Punchbowl Street  
 Honolulu, HI 96813  
  
 Telephone: (808)-547-4427  
 e-mail: severino@hawaii.edu, rseverino@queens.org