

## Chapter 325

# Poisson Regression

---

### Introduction

Poisson regression is similar to regular multiple regression except that the dependent ( $Y$ ) variable is an observed count that follows the Poisson distribution. Thus, the possible values of  $Y$  are the nonnegative integers: 0, 1, 2, 3, and so on. It is assumed that large counts are rare. Hence, Poisson regression is similar to logistic regression, which also has a discrete response variable. However, the response is not limited to specific values as it is in logistic regression.

One example of an appropriate application of Poisson regression is a study of how the colony counts of bacteria are related to various environmental conditions and dilutions. Another example is the number of failures for a certain machine at various operating conditions. Still another example is vital statistics concerning infant mortality or cancer incidence among groups with different demographics.

Most books on regression analysis briefly discuss Poisson regression. We are aware of only one book that is completely dedicated to the discussion of the topic. This is the book by Cameron and Trivedi (1998). Most of the methods presented here were obtained from their book.

This program computes Poisson regression on both numeric and categorical variables. It reports on the regression equation as well as the goodness of fit, confidence limits, likelihood, and deviance. It performs a comprehensive residual analysis including diagnostic residual reports and plots. It can perform a subset selection search, looking for the best regression model with the fewest independent variables. It provides confidence intervals on predicted values.

---

### The Poisson Distribution

The Poisson distribution models the probability of  $y$  events (i.e. failure, death, or existence) with the formula

$$\Pr(Y = y | \mu) = \frac{e^{-\mu} \mu^y}{y!} \quad (y = 0, 1, 2, \dots)$$

Notice that the Poisson distribution is specified with a single parameter  $\mu$ . This is the mean incidence rate of a rare event per unit of *exposure*. Exposure may be time, space, distance, area, volume, or population size. Because exposure is often a period of time, we use the symbol  $t$  to represent the exposure. When no exposure value is given, it is assumed to be one.

The parameter  $\mu$  may be interpreted as the risk of a new occurrence of the event during a specified exposure period,  $t$ . The probability of  $y$  events is then given by

$$\Pr(Y = y | \mu, t) = \frac{e^{-\mu t} (\mu t)^y}{y!} \quad (y = 0, 1, 2, \dots)$$

The Poisson distribution has the property that its mean and variance are equal.

## The Poisson Regression Model

In Poisson regression, we suppose that the Poisson incidence rate  $\mu$  is determined by a set of  $k$  regressor variables (the  $X$ 's). The expression relating these quantities is

$$\mu = t \exp(\beta_1 X_1 + \beta_2 X_2 + \cdots + \beta_k X_k)$$

Note that often,  $X_1 \equiv 1$  and  $\beta_1$  is called the *intercept*. The regression coefficients  $\beta_1, \beta_2, \dots, \beta_k$  are unknown parameters that are estimated from a set of data. Their estimates are labeled  $b_1, b_2, \dots, b_k$ .

Using this notation, the fundamental Poisson regression model for an observation  $i$  is written as

$$\Pr(Y_i = y_i | \mu_i, t_i) = \frac{e^{-\mu_i t_i} (\mu_i t_i)^{y_i}}{y_i!}$$

where

$$\begin{aligned} \mu_i &= t_i \mu(\mathbf{x}_i' \boldsymbol{\beta}) \\ &= t_i \exp(\beta_1 X_{1i} + \beta_2 X_{2i} + \cdots + \beta_k X_{ki}) \end{aligned}$$

That is, for a given set of values of the regressor variables, the outcome follows the Poisson distribution.

## Solution by Maximum Likelihood Estimation

The regression coefficients are estimated using the method of maximum likelihood. The logarithm of the likelihood function is

$$\ln[L(\mathbf{y}, \boldsymbol{\beta})] = \sum_{i=1}^n y_i \ln[t_i \mu(\mathbf{x}_i' \boldsymbol{\beta})] - \sum_{i=1}^n t_i \mu(\mathbf{x}_i' \boldsymbol{\beta}) - \sum_{i=1}^n \ln(y_i!)$$

Note that some statistical packages ignore the last term since it does not involve the regression parameters. This will make their calculated log-likelihoods different from ours.

The likelihood equations may be formed by taking the derivatives with respect to each regression coefficient and setting the result equal to zero. Doing this leads to a set of nonlinear equations that admits no closed-form solution. Thus, an iterative algorithm must be used to find the set of regression coefficients that maximum the log-likelihood. Using the method of iteratively reweighted least squares, a solution may be found in five or six iterations. However, the algorithm requires a complete pass through the data at each iteration, so it is relatively slow for problems with a large number of rows. With today's computers, this is becoming less and less of an issue.

## Distribution of the MLE's

Applying the usual maximum likelihood theory, the asymptotic distribution of the maximum likelihood estimates (MLE's) is multivariate normal. That is,

$$\hat{\boldsymbol{\beta}} \sim N(\boldsymbol{\beta}, \boldsymbol{\beta} V_{\hat{\boldsymbol{\beta}}})$$

where

$$V_{\hat{\boldsymbol{\beta}}} = \left( \sum_{i=1}^n \mu_i \mathbf{x}_i \mathbf{x}_i' \right)^{-1}$$

Remember that in the Poisson model the mean and the variance are equal. In practice, the data almost always reject this restriction. Usually, the variance is greater than the mean—a situation called *overdispersion*. The

## Poisson Regression

increase in variance is represented in the model by a constant multiple of the variance-covariance matrix. That is, we use

$$V_{\hat{\beta}} = \phi \left( \sum_{i=1}^n \mu_i \mathbf{x}_i \mathbf{x}_i' \right)^{-1}$$

where  $\phi$  is estimated using

$$\hat{\phi} = \frac{1}{n-k} \sum_{i=1}^n \frac{(y_i - \hat{\mu}_i)^2}{\hat{\mu}_i}$$

NCSS provides the option of using  $\phi$  (phi) in the calculation of the variances of the regression coefficients.

## Goodness of Fit Tests

Overall performance of the model is measured by two chi-square tests. These are the Pearson statistic

$$P_P = \sum_{i=1}^n \frac{(y_i - \hat{\mu}_i)^2}{\hat{\mu}_i}$$

and the deviance, or  $G$ , statistic

$$D_P = \sum_{i=1}^n \left\{ y_i \ln \left( \frac{y_i}{\hat{\mu}_i} \right) - (y_i - \hat{\mu}_i) \right\}$$

Both of these statistics are approximately chi-square distributed with  $n - k$  degrees of freedom. When a test is rejected, there is a significant lack of fit. When a test is not rejected, there is no evidence of lack of fit.

The Pearson statistic is only chi-square distributed when you are analyzing grouped data, so if you are not using a frequency variable, you should not use the Pearson statistic as a goodness of fit test. The Pearson statistic is often used as a test of overdispersion.

## Deviance

The deviance is twice the difference between the maximum achievable log-likelihood and the log-likelihood of the fitted model. In multiple regression under normality, the deviance is the residual sum of squares. In the case of Poisson regression, the deviance is a generalization of the sum of squares. The formula for the deviance is

$$D(\mathbf{y}, \hat{\boldsymbol{\mu}}) = 2 \{ LL_{\mathbf{y}} - LL_{\hat{\boldsymbol{\mu}}} \}$$

## Pseudo R-Squared Measures

The  $R$ -squared statistic does not extend to Poisson regression models. Various pseudo  $R$ -squared tests have been proposed. These pseudo measures have the property that, when applied to the linear model, they match the interpretation of the linear model  $R$ -squared. In Poisson regression, the most popular pseudo  $R$ -squared measure is function of the log-likelihoods of three models

$$R^2 = \frac{LL_{fit} - LL_0}{LL_{max} - LL_0}$$

## Poisson Regression

where

$$LL_0 = \sum_{i=1}^n y_i \ln[t_i \hat{\mu}] - \hat{\mu} \sum_{i=1}^n t_i - \sum_{i=1}^n \ln(y_i!) \quad \text{where } \hat{\mu} = \frac{\sum_{i=1}^n y_i}{\sum_{i=1}^n t_i}$$

$$LL_{\max} = \sum_{i=1}^n y_i \ln[t_i y_i] - \sum_{i=1}^n t_i y_i - \sum_{i=1}^n \ln(y_i!)$$

$$LL_{fit} = \sum_{i=1}^n y_i \ln[t_i \hat{\mu}(\mathbf{x}_i' \boldsymbol{\beta})] - \sum_{i=1}^n t_i \hat{\mu}(\mathbf{x}_i' \boldsymbol{\beta}) - \sum_{i=1}^n \ln(y_i!)$$

Note that  $LL_0$  is the log-likelihood of the intercept-only model,  $LL_{fit}$  is the log-likelihood of the current model, and  $LL_{\max}$  is the maximum log-likelihood possible. The maximum log-likelihood occurs when the actual responses (the  $y_i$ 's) exactly equal the predicted responses (the  $\mu_i$ 's).

Notice that this value of  $R$ -squared varies between zero and one, with a perfect fit occurring at one. Also note that it assumes that there is an intercept in the model. This may be an actual explicit intercept or an implicit intercept (as when you use a complete set of indicator variables to represent a categorical variable).

---

## Residuals

As in any regression analysis, a complete residual analysis should be employed. This involves plotting the residuals against various other quantities such as the regressor variables (to check for outliers and curvature) and the response variable. Various residuals may be of interest. These will be presented next.

### Raw Residual

The raw residual is the difference between the actual response and the estimated value from the model. Because in the Poisson case, the variance is equal to the mean, we expect that the variances of the residuals are unequal. This can lead to difficulties in the interpretation of the raw residuals. However, they are still popular. The formula for the raw residual is

$$r_i = y_i - \hat{\mu}_i$$

### Pearson Residual

The Pearson residual corrects for the unequal variance in the residuals by dividing by the standard deviation. The formula for the Pearson residual is

$$p_i = \frac{y_i - \hat{\mu}_i}{\sqrt{\hat{\phi} \hat{\mu}_i}}$$

### Deviance Residual

The deviance residual is another popular residual. It is popular because the sum of squares of these residuals is the deviance statistic. The formula for the deviance residual is

$$d_i = \text{sign}(y_i - \hat{\mu}_i) \sqrt{2 \left\{ y_i \ln \left( \frac{y_i}{\hat{\mu}_i} \right) - (y_i - \hat{\mu}_i) \right\}}$$

## Poisson Regression

### Hat Values

The Hat matrix is used in residual diagnostics to measure the influence of each observation. The hat values,  $h_{ii}$ , are the diagonal entries of the Hat matrix which is calculated using

$$H = W^{1/2} X (X'WX)^{-1} X'W^{1/2}$$

where  $W$  is a diagonal matrix made up of  $\hat{\mu}_i$ .

The hat values should be studied themselves, to understand which observations have a large influence on the fitted regression coefficients. Large hat values are those that are larger than  $2k/n$ . They are also used to further standardize residuals as is shown next.

### Studentized Pearson Residual

The formula for the studentized Pearson residual is

$$sp_i = \frac{p_i}{\sqrt{1 - h_{ii}}}$$

### Studentized Deviance Residual

The formula for the studentized deviance residual is

$$sd_i = \frac{d_i}{\sqrt{1 - h_{ii}}}$$

## Subset Selection

Subset selection refers to the task of finding a small subset of the available regressor variables that does a good job of predicting the dependent variable. Because Poisson regression must be solved iteratively, the task of finding the best subset can be time consuming. Hence, techniques which look at all possible combinations of the regressor variables are not feasible. Instead, algorithms that add or remove a variable at each step must be used. Two such searching algorithms are available in this module: forward selection and forward selection with switching.

Before discussing the details of these two algorithms, it is important to comment on a couple of issues that can come up. The first issue is what to do about the binary variables that are generated for a categorical independent variable. If such a variable has six categories, five binary variables are generated. You can see that with two or three categorical variables, a large number of binary variables may result, which greatly increases the total number of variables that must be searched. To avoid this problem, the algorithms used here search on model terms rather than on the individual variables. Thus, the whole set of binary variables associated with a given term are considered together for inclusion in, or deletion from, the model. Its all or none. Because of the time consuming nature of the algorithm, this is the only feasible way to deal with categorical variables. If you want the subset algorithm to deal with them individually, you can generate the set of binary variables manually and designate them as Numeric Variables.

### Hierarchical Models

A second issue is what to do with interactions. Usually, an interaction is not entered in the model unless the individual terms that make up that interaction are also in the model. For example, the interaction term  $A*B*C$  is not included unless the terms  $A$ ,  $B$ ,  $C$ ,  $A*B$ ,  $A*C$ , and  $B*C$  are already in the model. Such models are said to be *hierarchical*. You have the option during the search to force the algorithm to only consider hierarchical models during its search. Thus, if  $C$  is not in the model, interactions involving  $C$  are not even considered. Even though the option for non-hierarchical models is available, we recommend that you only consider hierarchical models.

## Poisson Regression

### Forward Selection

The method of forward selection proceeds as follows.

1. Begin with no terms in the model.
2. Find the term that, when added to the model, achieves the largest value of  $R$ -squared. Enter this term into the model.
3. Continue adding terms until a preset limit on the maximum number of terms in the model is reached.

This method is comparatively fast, but it does not guarantee that the best model is found except for the first step when it finds the best single term. You might use it when you have a large number of observations so that other, more time consuming methods, are not feasible, or when you have far too many possible regressor variables and you want to reduce the number of terms in the selection pool.

### Forward Selection with Switching

This method is similar to the method of Forward Selection discussed above. However, at each step when a term is added, all terms in the model are switched one at a time with all candidate terms not in the model to determine if they increase the value of  $R$ -squared. If a switch can be found, it is made and the candidate terms are again searched to determine if another switch can be made.

When the search for possible switches does not yield a candidate, the subset size is increased by one and a new search is begun. The algorithm is terminated when a target subset size is reached or all terms are included in the model.

### Discussion

These algorithms usually require two runs. In the first run, you set the maximum subset size to a large value such as 10. By studying the Subset Selection reports from this run, you can quickly determine the optimum number of terms. You reset the maximum subset size to this number and make the second run. This two-step procedure works better than relying on some F-to-enter and F-to-remove tests whose properties are not well understood to begin with.

---

## Data Structure

At a minimum, datasets to be analyzed by Poisson regression must contain a dependent variable and one or more independent variables. For each categorical variable, the program generates a set of binary (0 and 1) variables that express the same information. For example, in the table below, the discrete variable AgeGroup will be replaced by the variables Ag2 through Ag6 (Ag1 is not needed).

Koch et. al. (1986) present the following data taken from the Third National Cancer Survey. This dataset contains the number of new melanoma cases in 1969-1971 among white males in two areas for various age groups. The size of the estimated population at risk is given in the variable Population.

## Poisson Regression

## Koch36 dataset

Melanoma	Area	AgeGroup	Population	AG1	AG2	AG3	AG4	AG5	AG6
61	0	<35	2880262	1	0	0	0	0	0
76	0	35-44	564535	0	1	0	0	0	0
98	0	45-54	592983	0	0	1	0	0	0
104	0	54-64	450740	0	0	0	1	0	0
63	0	65-74	270908	0	0	0	0	1	0
80	0	>74	161850	0	0	0	0	0	1
64	1	<35	1074246	1	0	0	0	0	0
75	1	35-44	220407	0	1	0	0	0	0
68	1	45-54	198119	0	0	1	0	0	0
63	1	54-64	134084	0	0	0	1	0	0
45	1	65-74	70708	0	0	0	0	1	0
27	1	>74	34233	0	0	0	0	0	1

---

## Missing Values

If missing values are found in any of the independent variables being used, the row is omitted. If only the value of the dependent variable is missing, that row will not be used during the estimation process, but its predicted value will be generated and reported on.

---

## Procedure Options

This section describes the options available in this procedure.

---

## Variables, Model Tab

This panel specifies the variables and model are used in the analysis.

---

### Variables

#### Dependent Y

Specify the dependent (response) variable. This is the variable to be predicted by the independent variables. The values in this variable should be non-negative integers (zero is okay).

#### Exposure T

Specify an optional variable containing exposure values. If this option is left blank, all exposures will be set to 1.0. This variable is specified when the exposures are different for each row.

The exposure is the amount of time, space, distance, volume, or population size from which the dependent variable is counted. For example, exposure may be the time in days, months, or years during which the values on that row were obtained. It may be the number of individuals at risk or the number of man-years from which the dependent variable is measured.

Each exposure must be a positive (non-zero) number or the row is ignored during the estimation phase.

#### Numeric X's

Specify the numeric (continuous) independent variables. By numeric, we mean that the values are numeric and at least ordinal. Nominal variables, even when coded with numbers, should be specified as Categorical Independent Variables. Although you may specify binary (0-1) variables here, they are better analyzed when you specify them as Categorical Independent Variables.

## Poisson Regression

If you want to create powers and cross-products of these variables, specify an appropriate model in the 'Custom Model' field under the Model tab.

If you want to create predicted values of  $Y$  for values of  $X$  not in your database, add the  $X$  values to the bottom of the database. They will not be used during estimation, but predicted values will be generated for them.

### Categorical X's

Specify categorical (nominal or group) independent variables in this box. By categorical we mean that the variable has only a few unique, numeric or text, values like 1, 2, 3 or Yes, No, Maybe. The values are used to identify categories.

Regression analysis is only defined for numeric variables. Since categorical variables are nominal, they cannot be used directly in regression. Instead, an internal set of numeric variables must be substituted for each categorical variable.

Suppose a categorical variable has  $G$  categories. NCSS automatically generates the  $G-1$  internal, numeric variables for the analysis. The way these internal variables are created is determined by the Recoding Scheme and, if needed, the Reference Value. These options can be entered separately with each categorical variable, or they can be specified using a default value (see Default Recoding Scheme and Default Reference Value below).

The syntax for specifying a categorical variable is  $VarName(CType; RefValue)$  where  $VarName$  is the name of the variable,  $CType$  is the recoding scheme, and  $RefValue$  is the reference value, if needed.

### CType

The recoding scheme is entered as a letter. Possible choices are B, P, R, N, S, L, F, A, 1, 2, 3, 4, 5, or E. The meaning of each of these letters is as follows.

- B for binary** (the group with the reference value is skipped).  
 Example: Categorical variable Z with 4 categories. Category D is the reference value.

Z	B1	B2	B3
A	1	0	0
B	0	1	0
C	0	0	1
D	0	0	0
- P for Polynomial** of up to 5th order (you cannot use this option with category variables with more than 6 categories).  
 Example: Categorical variable Z with 4 categories.

Z	P1	P2	P3
1	-3	1	-1
3	-1	-1	3
5	1	-1	-3
7	3	1	1
- R to compare each with the reference value** (the group with the reference value is skipped).  
 Example: Categorical variable Z with 4 categories. Category D is the reference value.

Z	C1	C2	C3
A	1	0	0
B	0	1	0
C	0	0	1
D	-1	-1	-1



## Poisson Regression

- **N** to compare each with the **next** category.  
Example: Categorical variable Z with 4 categories.  

Z	S1	S2	S3
1	1	0	0
3	-1	1	0
5	0	-1	1
7	0	0	-1
- **S** to compare each with the **average of all subsequent** values.  
Example: Categorical variable Z with 4 categories.  

Z	S1	S2	S3
1	-3	0	0
3	1	-2	0
5	1	1	-1
7	1	1	1
- **L** to compare each with the **prior** category.  
Example: Categorical variable Z with 4 categories.  

Z	S1	S2	S3
1	-1	0	0
3	1	-1	0
5	0	1	-1
7	0	0	1
- **F** to compare each with the **average of all prior** categories.  
Example: Categorical variable Z with 4 categories.  

Z	S1	S2	S3
1	1	1	1
3	1	1	-1
5	1	-2	0
7	-3	0	0
- **A** to compare each with the **average of all** categories (the Reference Value is skipped).  
Example: Categorical variable Z with 4 categories. Suppose the reference value is 3.  

Z	S1	S2	S3
1	-3	1	1
3	1	1	1
5	1	-3	1
7	1	1	-3
- **1** to compare each with the **first** category after sorting.  
Example: Categorical variable Z with 4 categories.  

Z	C1	C2	C3
A	-1	-1	-1
B	1	0	0
C	0	1	0
D	0	0	1

## Poisson Regression

- **2** to compare each with the **second** category after sorting.

Example: Categorical variable Z with 4 categories.

Z	C1	C2	C3
A	1	0	0
B	-1	-1	-1
C	0	1	0
D	0	0	1

- **3** to compare each with the **third** category after sorting.

Example: Categorical variable Z with 4 categories.

Z	C1	C2	C3
A	1	0	0
B	0	1	0
C	-1	-1	-1
D	0	0	1

- **4** to compare each with the **fourth** category after sorting.

Example: Categorical variable Z with 4 categories.

Z	C1	C2	C3
A	1	0	0
B	0	1	0
C	0	0	1
D	-1	-1	-1

- **5** to compare each with the **fifth** category after sorting.

Example: Categorical variable Z with 5 categories.

Z	C1	C2	C3	C4
A	1	0	0	0
B	0	1	0	0
C	0	0	1	0
D	0	0	0	1
E	-1	-1	-1	-1

- **E** to compare each with the **last** category after sorting.

Example: Categorical variable Z with 4 categories.

Z	C1	C2	C3
A	1	0	0
B	0	1	0
C	0	0	1
D	-1	-1	-1

### RefValue

A second, optional argument is the reference value. The reference value is one of the categories. The other categories are compared to it, so it is usually a baseline or control value. If neither a baseline or control value is evident, the reference value is the most frequent value.

For example, suppose you want to include a categorical independent variable, State, which has four values: Texas, California, Florida, and New York. Suppose the recoding scheme is specified as *Compare Each with Reference Value* with the reference value of *California*. You would enter

**State(R;California)**

## Poisson Regression

### Default Recoding Scheme

Select the default type of numeric variable that will be generated when processing categorical independent variables. The values in a categorical variable are not used directly in regression analysis. Instead, a set of numeric variables is automatically created and substituted for them. This option allows you to specify what type of numeric variable will be created. The options are outlined in the sections below.

The contrast type may also be designated within parentheses after the name of each categorical independent variable, in which case the default contrast type is ignored.

If your model includes interactions of categorical variables, this option should be set to 'Contrast with Reference' or Compare with All Subsequent' in order to match GLM results for factor effects.

- Binary** (the group with the reference value is skipped).  
 Example: Categorical variable Z with 4 categories. Category D is the reference value.
 

Z	B1	B2	B3
A	1	0	0
B	0	1	0
C	0	0	1
D	0	0	0
- Polynomial** of up to 5th order (you cannot use this option with category variables with more than 6 categories).  
 Example: Categorical variable Z with 4 categories.
 

Z	P1	P2	P3
1	-3	1	-1
3	-1	-1	3
5	1	-1	-3
7	3	1	1
- Compare Each with Reference Value** (the group with the reference value is skipped).  
 Example: Categorical variable Z with 4 categories. Category D is the reference value.
 

Z	C1	C2	C3
A	1	0	0
B	0	1	0
C	0	0	1
D	-1	-1	-1
- Compare Each with Next.**  
 Example: Categorical variable Z with 4 categories.
 

Z	S1	S2	S3
1	1	0	0
3	-1	1	0
5	0	-1	1
7	0	0	-1
- Compare Each with All Subsequent.**  
 Example: Categorical variable Z with 4 categories.
 

Z	S1	S2	S3
1	-3	0	0
3	1	-2	0
5	1	1	-1
7	1	1	1

## Poisson Regression

- **Compare Each with Prior**

Example: Categorical variable Z with 4 categories.

Z	S1	S2	S3
1	-1	0	0
3	1	-1	0
5	0	1	-1
7	0	0	1

- **Compare Each with All Prior**

Example: Categorical variable Z with 4 categories.

Z	S1	S2	S3
1	1	1	1
3	1	1	-1
5	1	-2	0
7	-3	0	0

- **Compare Each with Average**

Example: Categorical variable Z with 4 categories. Suppose the reference value is 3.

Z	S1	S2	S3
1	-3	1	1
3	1	1	1
5	1	-3	1
7	1	1	-3

- **Compare Each with First**

Example: Categorical variable Z with 4 categories.

Z	C1	C2	C3
A	-1	-1	-1
B	1	0	0
C	0	1	0
D	0	0	1

- **Compare Each with Second**

Example: Categorical variable Z with 4 categories.

Z	C1	C2	C3
A	1	0	0
B	-1	-1	-1
C	0	1	0
D	0	0	1

- **Compare Each with Third**

Example: Categorical variable Z with 4 categories.

Z	C1	C2	C3
A	1	0	0
B	0	1	0
C	-1	-1	-1
D	0	0	1

## Poisson Regression

- **Compare Each with Fourth**

Example: Categorical variable Z with 4 categories.

Z	C1	C2	C3
A	1	0	0
B	0	1	0
C	0	0	1
D	-1	-1	-1

- **Compare Each with Fifth**

Example: Categorical variable Z with 5 categories.

Z	C1	C2	C3	C4
A	1	0	0	0
B	0	1	0	0
C	0	0	1	0
D	0	0	0	1
E	-1	-1	-1	-1

- **Compare Each with Last**

Example: Categorical variable Z with 4 categories.

Z	C1	C2	C3
A	1	0	0
B	0	1	0
C	0	0	1
D	-1	-1	-1

### Default Reference Value

This option specifies the default reference value to be used when automatically generating indicator variables during the processing of selected categorical independent variables. The reference value is often the baseline, and the other values are compared to it. The choices are

- **First Value after Sorting – Fifth Value after Sorting**

Use the first (through fifth) value in alpha-numeric sorted order as the reference value.

- **Last Value after Sorting**

Use the last value in alpha-numeric sorted order as the reference value.

### Frequencies

This is an optional variable containing the frequency (observation count) for each row. Usually, you would leave this option blank and let each row receive the default frequency of one.

If your data have already been summarized, this option lets you specify how many actual rows each physical row represents.

---

## Regression Model

### Terms

This option specifies which terms (terms, powers, cross-products, and interactions) are included in the regression model. For a straight-forward regression model, select *1-Way*.

## Poisson Regression

The options are

- **Up to 1-Way**

This option generates a model in which each variable is represented by a single model term. No cross-products, interactions, or powers are added. Use this option when you want to use the variables you have specified, but you do not want to generate other terms.

This is the option to select when you want to analyze the independent variables specified without adding any other terms.

For example, if you have three independent variables A, B, and C, this would generate the model:

$$A + B + C$$

- **Up to 2-Way**

This option specifies that all individual variables, two-way interactions, and squares of numeric variables are included in the model. For example, if you have three numeric variables A, B, and C, this would generate the model:

$$A + B + C + A*B + A*C + B*C + A*A + B*B + C*C$$

On the other hand, if you have three categorical variables A, B, and C, this would generate the model:

$$A + B + C + A*B + A*C + B*C$$

- **Up to 3-Way**

All individual variables, two-way interactions, three-way interactions, squares of numeric variables, and cubes of numeric variables are included in the model. For example, if you have three numeric, independent variables A, B, and C, this would generate the model:

$$A + B + C + A*B + A*C + B*C + A*B*C + A*A + B*B + C*C + A*A*B + A*A*C + B*B*C + A*C*C + B*C*C$$

On the other hand, if you have three categorical variables A, B, and C, this would generate the model:

$$A + B + C + A*B + A*C + B*C + A*B*C$$

- **Up to 4-Way**

All individual variables, two-way interactions, three-way interactions, and four-way interactions are included in the model. Also included would be squares, cubes, and quartics of numeric variables and their cross-products.

For example, if you have four categorical variables A, B, C, and D, this would generate the model:

$$A + B + C + D + A*B + A*C + A*D + B*C + B*D + C*D + A*B*C + A*B*D + A*C*D + B*C*D + A*B*C*D$$

- **Interaction**

Mainly used for categorical variables. A saturated model (all terms and their interactions) is generated. This requires a dataset with no missing categorical-variable combinations (you can have unequal numbers of observations for each combination of the categorical variables). No squares, cubes, etc. are generated.

For example, if you have three independent variables A, B, and C, this would generate the model:

$$A + B + C + A*B + A*C + B*C + A*B*C$$

Note that the discussion of the Custom Model option discusses the interpretation of this model.

- **Custom Model**

The model specified in the *Custom Model* box is used.

## Poisson Regression

### Remove Intercept

Unchecked indicates that the intercept term,  $\beta_0$ , is to be included in the regression. Checked indicates that the intercept should be omitted from the regression model. Note that deleting the intercept distorts most of the diagnostic statistics ( $R^2$ , etc.). In most situations, you should include the intercept in the model.

### Replace Custom Model with Preview Model (button)

When this button is pressed, the Custom Model is cleared and a copy of the Preview model is stored in the Custom Model. You can then edit this Custom Model as desired.

### Maximum Order of Custom Terms

This option specifies that maximum number of variables that can occur in an interaction (or cross-product) term in a custom model. For example,  $A*B*C$  is a third order interaction term and if this option were set to 2, the  $A*B*C$  term would not be included in the model.

This option is particularly useful when used with the bar notation of a custom model to allow a simple way to remove unwanted high-order interactions.

### Custom Model

This options specifies a custom model. It is only used when the *Terms* option is set to *Custom*. A custom model specifies the terms (single variables and interactions) that are to be kept in the model.

### Interactions

An interaction expresses the combined relationship between two or more variables and the dependent variable by creating a new variable that is the product of the variables. The interaction between two numeric variables is generated by multiplying them. The interaction between two categorical variables is generated by multiplying each pair of indicator variables. The interaction between a numeric variable and a categorical variable is created by generating all products between the numeric variable and the indicator variables generated from the categorical variable.

### Syntax

A model is written by listing one or more terms. The terms are separated by a blank or plus sign. Terms include variables and interactions. Specify regular variables (main effects) by entering the variable names. Specify interactions by listing each variable in the interaction separated by an asterisk (\*), such as Fruit\*Nuts or  $A*B*C$ .

You can use the bar (|) symbol as a shorthand technique for specifying many interactions quickly. When several variables are separated by bars, all of their interactions are generated. For example,  $A|B|C$  is interpreted as  $A + B + C + A*B + A*C + B*C + A*B*C$ .

You can use parentheses. For example,  $A*(B+C)$  is interpreted as  $A*B + A*C$ .

Some examples will help to indicate how the model syntax works:

$$A|B = A + B + A*B$$

$$A|B A*A B*B = A + B + A*B + A*A + B*B$$

Note that you should only repeat numeric variable. That is,  $A*A$  is valid for a numeric variable, but not for a categorical variable.

$$A|A|B|B \text{ (Max Term Order=2)} = A + B + A*A + A*B + B*B$$

$$A|B|C = A + B + C + A*B + A*C + B*C + A*B*C$$

$$(A + B)*(C + D) = A*C + A*D + B*C + B*D$$

$$(A + B)|C = (A + B) + C + (A + B)*C = A + B + C + A*C + B*C$$

## Poisson Regression

### Subset Selection

#### Search Method

This option specifies the subset selection algorithm used to reduce the number of independent variables that used in the regression model. Note that since the solution algorithm is iterative, the selection process can be very time consuming. The Forward algorithm is much quicker than the Forward with Switching algorithm, but the Forward algorithm does not usually find as good of a model.

Also note that in the case of categorical independent variables, the algorithm searches among the original categorical variables, not among the generated individual binary variables. That is, either all binary variables associated with a particular categorical variable are included or not—they are not considered individually.

*Hierarchical models* are such that if an interaction is in the model, so are the terms that can be derived from it. For example, if  $A*B*C$  is in the model, so are  $A$ ,  $B$ ,  $C$ ,  $A*B$ ,  $A*C$ , and  $B*C$ . Statisticians usually adopt hierarchical models rather than non-hierarchical models. The subset selection procedure can be made to consider only hierarchical models during its search.

The subset selection options are:

- **None – No Search is Conducted**

No subset selection is attempted. All specified independent variables are used in the regression equation.

- **(Hierarchical) Forward**

With this algorithm, the term with the largest log likelihood is entered into the model. Next, the term that increases the log likelihood the most is added. This selection is continued until all the terms have been entered or until the maximum subset size has been reached.

If hierarchical models are selected, only those terms that will keep the model hierarchical are candidates for selection. For example, the interaction term  $A*B$  will not be considered unless both  $A$  and  $B$  are already in the model.

When using this algorithm, you must make one run that allows a large number of terms to find the appropriate number of terms. Next, a second run is made in which you decrease the maximum terms in the subset to the number after which the log likelihood does not change significantly.

- **(Hierarchical) Forward with Switching**

This algorithm is similar to the Forward algorithm described above. The term with the largest log likelihood is entered into the regression model. The term which increases the log likelihood the most when combined with the first term is entered next. Now, each term in the current model is removed and the rest of the terms are checked to determine if, when they are used instead, the likelihood function is increased. If a term can be found by this switching process, the switch is made and the whole switching operation is begun again. The algorithm continues until no term can be found that improves the likelihood. This model then becomes the best two-term model.

Next, the subset size is increased by one, the best third term is entered into the model, and the switching process is repeated. This process is repeated until the maximum subset size is reached. Hence, this model finds the optimum subset for each subset size. You must make one run to find an appropriate subset size by looking at the change in the log likelihood. You then reset the maximum subset size to this value and rerun the analysis.

If hierarchical models are selected, only those terms that will keep the model hierarchical are candidates for addition or deletion. For example, the interaction term  $A*B$  will not be considered unless both  $A$  and  $B$  are already in the model. Likewise, the term  $A$  cannot be removed from a model that contains  $A*B$ .



## Poisson Regression

### Stop search when number of terms reaches

Once this number of terms has been entered into the model, the subset selection algorithm is terminated. Often you will have to run the procedure twice to find an appropriate value. You would set this value high for the first run and then reset it appropriately for the second run, depending upon the values of the log likelihood.

Note that the intercept is counted in this number.

---

## Iterations Tab

---

### Estimation Options

The following options are used during the likelihood maximization process.

#### Maximum Iterations

Specifies the maximum number of iterations allowed during the iteration procedure. If this number is reached, the procedure is terminated prematurely. Typically, the maximum likelihood procedure converges in five or six iterations, so a value of twenty here should be ample.

#### Convergence Zero

This option specifies the convergence target for the maximum likelihood estimation procedure. When all of the maximum likelihood equations are less than this amount, the algorithm is assumed to have converged. In theory, all of the equations should be zero. However, about the best that can be achieved is 1E-13, so you should set this value to a number a little larger than this such as the default of 1E-9.

The actual value can be found by looking at the Maximum Convergence value on the Run Summary report.

---

## Reports Tab

The following options control which reports are displayed.

---

### Variance Adjustment

#### Use Dispersion Phi in SE's

Indicate whether to use the Phi multiplier in the calculation of the standard errors of the regression coefficients.

The Poisson model assumes that the mean and variance are identical. Usually, the variance is larger than the mean (called *overdispersion*). A correction can be applied to the standard errors by multiplying them by the Phi coefficient.

Note that this correction will not change the estimated regression coefficients.

---

## Alpha

### Alpha Level

Alpha is the significance level used in the hypothesis tests. One minus alpha is the confidence level of the confidence intervals. A value of 0.05 is most commonly used. This corresponds to a chance of error of 1 in 20. You should not be afraid to use other values since 0.05 became popular in pre-computer days when it was the only value available.

Typical values range from 0.001 to 0.20.

## Poisson Regression

---

### Select Reports – Summaries

#### Run Summary ... Means

Each of these options specifies whether the indicated report is calculated and displayed.

---

### Select Reports – Subset Selection

#### Subset Selection - Summary and Subset Selection - Detail

Indicate whether to display these subset selection reports.

---

### Select Reports – Estimation

#### Regression Coefficients ... Rate Coefficients

Indicate whether to display these estimation reports.

---

### Select Reports – Goodness-of-Fit

#### Lack-of-Fit Tests ... Log-Likelihood and R-Squared

Indicate whether to display these model goodness-of-fit reports.

---

### Select Reports – Row-by-Row Lists

#### Residuals ... Incidence

Indicate whether to display these list reports. Note that since these reports provide results for each row, they may be too long for normal use when requested on large databases.

#### Incidence Counts

Up to five incidence counts may be entered. The probabilities of these counts under the Poisson regression model will be displayed on the Incidence Report.

These values must be non-negative integers.

#### Exposure Value

Specify the exposure (time, space, distance, volume, etc.) value to be used as a multiplier on the Incidence Report. All items on that report are scaled to this amount. For example, if your data was scaled in terms of events per month but you want the Incidence report scaled to events per year, you would enter '12' here.

---

## Report Options Tab

These options control format of the reports.

---

### Variable Labels

#### Variable Names

This option lets you select whether to display only variable names, variable labels, or both.

## Poisson Regression

### Stagger label and output if label length is $\geq$

The names of the indicator variables can be too long to fit in the space provided. If the name contains more characters than the number specified here, only the name is shown on the first line of the report and the rest of the output is placed on the next line.

Enter *1* when you want the each variable's results printed on two lines.

Enter *100* when you want each variable's results printed on a single line.

---

## Decimal Places

### Precision

Specifies whether unformatted numbers (designated as decimal places = 'All') are displayed as single (7-digit) or double (13-digit) precision numbers in the output. All calculations are performed in double precision regardless of the Precision selected here.

### Single

Unformatted numbers are displayed with 7-digits. This is the default setting. All reports have been formatted for single precision.

### Double

Unformatted numbers are displayed with 13-digits. This option is most often used when the extremely accurate results are needed for further calculation. For example, double precision might be used when you are going to use the Multiple Regression model in a transformation.

### Double Precision Format Misalignment

Double precision numbers require more space than is available in the output columns, causing column alignment problems. The double precision option is for those instances when accuracy is more important than format alignment.

### Comments

1. This option does not affect formatted numbers such as probability levels.
2. This option only influences the format of the numbers as they presented in the output. All calculations are performed in double precision regardless of the Precision selected here.

### Y ... Chi-Square Decimals

Specify the number of digits after the decimal point to display on the output of values of this type. Note that this option in no way influences the accuracy with which the calculations are done.

Enter **All** to display all digits available. The number of digits displayed by this option is controlled by whether the **Precision** option is *Single* or *Double*.

---

## Plots Tab

These options control the attributes of the various plots.

---

### Select Plots

#### Incidence (Y/T) vs X Plot ... Resid vs X Plot

Indicate whether to display these plots. Click the plot format button to change the plot settings.

## Poisson Regression

### Edit During Run

This is the small check-box in the upper right-hand corner of the format button. If checked, the graphics format window for this plot will be displayed while the procedure is running so that you can format it with the actual data.

---

### Plot Options

#### Residual Plotted

This option specifies which of the five types of residuals are shown on the residual plots.

---

### Storage Tab

These options let you specify if, and where on the dataset, various statistics are stored.

*Warning: Any data already in these columns are replaced by the new data. Be careful not to specify columns that contain important data.*

---

### Data Storage Options

#### Storage Option

This option controls whether the values indicated below are stored on the dataset when the procedure is run.

- **Do not store data**  
No data are stored even if they are checked.
- **Store in empty columns only**  
The values are stored in empty columns only. Columns containing data are not used for data storage, so no data can be lost.
- **Store in designated columns**  
Beginning at the *Store First Item In* column, the values are stored in this column and those to the right. If a column contains data, the data are replaced by the storage values. Care must be used with this option because it cannot be undone.

#### Store First Item In

The first item is stored in this column. Each additional item that is checked is stored in the columns immediately to the right of this column.

Leave this value blank if you want the data storage to begin in the first blank column on the right-hand side of the data.

Warning: any existing data in these columns is automatically replaced, so be careful..

---

### Data Storage Options – Select Items to Store

#### Expanded X Values ... Covariance Matrix

Indicated whether to store these row-by-row values, beginning at the column indicated by the *Store First Item In* option. Note that several of these values include a different value for each group and so they require several columns when they are stored.

#### Expanded X Values

This option refers to the experimental design matrix. They include all binary and interaction variables generated.

---

## Example 1 – Poisson Regression using a Dataset with Indicator Variables

This section presents several examples. In the first example, the data shown earlier in the Data Structure section and found in the Koch36 dataset will be analyzed. Koch et. al. (1986) presented this dataset. It contains the number of new melanoma cases in 1969-1971 among white males in two areas for various age groups. The size of the estimated population at risk is given in the variable Population.

This dataset is instructive because it shows how easily categorical variables are dealt with. In this example, two categorical variables, Area and AgeGroup, will be included in the regression model. The dataset can also be used to validate the program since the results are given in Koch (1986).

You may follow along here by making the appropriate entries or load the completed template **Example 1** by clicking on Open Example Template from the File menu of the Poisson Regression window.

### 1 Open the Koch36 dataset.

- From the File menu of the NCSS Data window, select **Open Example Data**.
- Click on the file **Koch36.NCSS**.
- Click **Open**.

### 2 Open the Poisson Regression window.

- Using the Analysis menu or the Procedure Navigator, find and select the **Poisson Regression** procedure.
- On the menus, select **File**, then **New Template**. This will fill the procedure with the default template.

### 3 Specify the variables.

- On the Poisson Regression window, select the **Variables, Model** tab.
- Double-click in the **Dependent Y** box. This will bring up the variable selection window.
- Select **Melanoma** from the list of variables and click **Ok**. *Melanoma* will appear in the **Dependent Y** box.
- Double-click in the **T: Exposure Variable** box.
- Select **Population** from the list of variables and click **Ok**.
- Double-click in the **Categorical X's** box.
- Enter **Area(0) AgeGroup(<35)** in the **Categorical X's** box. The values in parentheses give the reference value for each variable.
- The rest of this panel can be left at the default values.

### 4 Specify the model.

- Set the **Terms** option to **1-Way**.
- Set the **Subset Selection** option to **None**.

### 5 Specify the reports.

- Select the **Reports** tab.
- Check all of the reports and plots. Normally, you would not want all of them, but we specify them now for documentation purposes.
- Set the **Incidence Counts** to **5 10 15 20 25**.
- Set the **Exposure Value** to **100000**.

### 6 Specify the decimals.

- Select the **Report Options** tab.
- Set the number of **decimal places for Probability** to **6**.

### 7 Run the procedure.

- From the Run menu, select **Run Procedure**. Alternatively, just click the green Run button.

## Poisson Regression

## Run Summary

Item	Value	Item	Value
Dependent Variable	Melanoma	Rows Used	12
Exposure Variable	Population	Sum of Frequencies	12
Frequency Variable	None	Iterations	5
Ind. Var's Available	2	Convergence Zero	1E-09
No. of X's in Model	6	Maximum Convergence	5.307754E-12
Pseudo R <sup>2</sup>	0.9931	Dispersion Phi	1.2230
Final Likelihood	-39.2199	Phi was not used to correct standard errors.	
Subset Method	None		

This report provides several details about the data and the MLE algorithm.

### Dependent, Exposure, and Frequency Variables

These variables are listed to provide a record of the variables that were analyzed.

### Ind. Var's Available

This is the number of independent variables that you have selected.

### No. of X's in Model

This is the number of actual X-variables generated from the terms in the model that was used in the analysis.

### Pseudo R<sup>2</sup>

This is the generalization of regular  $R^2$  in multiple regression. This value is discussed in detail in the Technical Details section of the chapter. Its formula is

$$R^2 = \frac{LL_{fit} - LL_0}{LL_{max} - LL_0}$$

### Final Likelihood

This is the value of the log likelihood that was achieved for this run.

### Subset Method

This is the type of subset selection that was run.

### Rows Used

This is the number of rows used by the estimation algorithm. Rows with missing values and filtered rows are not included. Always check this value to make sure that you are analyzing all of the data you intended to.

### Sum of Frequencies

This is the number of observations used by the estimation algorithm. If you specified a Frequency Variable, this will be greater than the number of rows. If not, they will be equal.

### Iterations

This is number of iterations used by the estimation algorithm. Usually, the algorithm will terminate in five or six iterations.

### Convergence Zero

The estimation algorithm continues until all of the likelihood equations are close to zero. This is *zero* to the algorithm. When the maximum convergence value is less than this amount, the algorithm has converged. Compare this value to the Maximum Convergence value.

### Maximum Convergence

The estimation algorithm continues until all of the likelihood equations are close to zero. This is largest value of all of these equations. It should be close to zero or the algorithm was terminated before it had converged.

## Poisson Regression

**Dispersion Phi**

This line gives the estimated value of the dispersion phi. It also indicates whether phi was used to adjust the standard errors of the regression coefficients and the predicted values.

**Model Summary**

Model	Model DF	Error DF	Log Likelihood	Deviance	AIC	Pseudo R <sup>2</sup>
Intercept	1	11	-484.0223	895.8197	897.8197	0.0000
Model	7	5	-39.2199	6.2149	20.2149	0.9931
Maximum	12	0	-36.1125	0.0000	24.0000	1.0000

This report is analogous to the analysis of variance table. It summarizes the goodness of fit of the model.

**Model**

This is the term(s) that are reported about on this row of the report. Note that the model line includes the intercept.

**Model DF**

This is the number of variables in the model.

**Error DF**

This is the number of observations minus the number of variables.

**Log Likelihood**

This is the value of the log-likelihood function for the intercept only model, the chosen model, and the saturated model that fits the data perfectly. By comparing these values, you obtain an understanding of how well you model fits the data.

**Deviance**

The deviance is the generalization of the sum of squares in regular multiple regression. It measures the discrepancy between the fitted values and the data.

**AIC**

This is Akaike's information criterion (AIC). It is equal to the deviance plus twice the number of parameters in the model. It combines a measure of the discrepancy between the fitted values and the data (the deviance) with a measure of the simplicity of the model (twice the number of parameters). It has been shown that using AIC to compare competing models with different numbers of parameters amounts to selecting the model with the minimum estimate of the mean squared error of prediction.

**Pseudo R<sup>2</sup>**

This is the generalization of regular  $R^2$  in multiple regression. This value is discussed in detail in the Technical Details section of the chapter. Its formula is

$$R^2 = \frac{LL_{fit} - LL_0}{LL_{max} - LL_0}$$

## Poisson Regression

## Means Report

Variable	Mean	Minimum	Maximum
Melanoma	68.667	27.000	104.000
Population	554422.917	34233.000	2880262.000

This report gives the mean, minimum, and maximum for each of the numeric variables in the analysis. Use it to check for obvious data errors.

## Regression Coefficients Section

Independent Variable	Regression Coefficient b(i)	Standard Error Sb(i)	Wald's Chi <sup>2</sup> H0: $\beta=0$	Prob Level	Lower 95.0% Confidence Limit	Upper 95.0% Confidence Limit
Intercept	-10.65831	0.09518	12538.43	0.000000	-10.84487	-10.47175
(Area=1)	0.81948	0.07103	133.11	0.000000	0.68027	0.95870
(AgeGroup="35-44")	1.79737	0.12093	220.92	0.000000	1.56036	2.03439
(AgeGroup="45-54")	1.91309	0.11844	260.90	0.000000	1.68095	2.14522
(AgeGroup="54-64")	2.24180	0.11834	358.89	0.000000	2.00987	2.47374
(AgeGroup="65-74")	2.36572	0.13152	323.56	0.000000	2.10795	2.62349
(AgeGroup=">74")	2.94468	0.13205	497.30	0.000000	2.68587	3.20349
Dispersion Phi		1.2230				

## Estimated Poisson Regression Model

Melanoma =

Exp( -10.6583092620666 + 0.819484586814042\*(Area=1) + 1.79737495802664\*(AgeGroup="35-44") + 1.91308772800918\*(AgeGroup="45-54") + 2.24180245796944\*(AgeGroup="54-64") + 2.36572417048965\*(AgeGroup="65-74") + 2.94467922306084\*(AgeGroup=">74") )

This report provides the estimated regression model and associated statistics. It provides the main results of the analysis.

## Validation

Koch (1986) gives the following estimates and standard errors.

Independent Variable	ML Estimate	Standard Error
Intercept	-10.66	0.01
Area	0.82	0.07
AG2	1.80	0.12
AG3	1.91	0.12
AG4	2.24	0.12
AG5	2.37	0.13
AG6	2.94	0.13

As you can see, these results match those provided by NCSS exactly—validating our algorithms. These results were also validated using SAS.

## Independent Variable

This item provides the name of the independent variable shown on this line of the report. The *Intercept* refers to the optional constant term. The *Dispersion Phi* is the estimated value of the phi coefficient.

Note that whether a line is skipped after the name of the independent variable is displayed is controlled by the *Stagger label and output if label length is  $\geq$*  option in the Format tab.



## Poisson Regression

### Regression Coefficient

These are the maximum-likelihood estimates of the regression coefficients,  $b_1, b_2, \dots, b_k$ . Their direct interpretation is difficult because the formula for the predicted value involves the exponential function.

### Standard Error

These are the asymptotic standard errors of the regression coefficients, the  $s_{b_i}$ . They estimate the precision of the regression coefficient. The standard errors are the square roots of the diagonal elements of this covariance matrix. The covariance matrix is obtained by inverting the observed information matrix evaluated at the maximum likelihood estimates.

If you Use Dispersion Phi option, the corrected standard error is shown. This is found by multiplying the simple standard error by the square root of phi. That is, the value displayed is  $s'_{b_i}$  where

$$s'_{b_i} = s_{b_i} \sqrt{\phi}$$

### Wald's Chi<sup>2</sup> H0: $\beta=0$

This is the one degree of freedom chi-square statistic for testing the null hypothesis that  $\beta_i = 0$  against the alternative that  $\beta_i \neq 0$ . The chi-square value is called a *Wald statistic*. This test has been found to follow the chi-square distribution only in large samples.

The test is calculated using

$$\chi_1^2 = \left( \frac{b_i}{s'_{b_i}} \right)^2$$

### Prob Level

The probability of obtaining a chi-square value greater than the above. This is the significance level of the test. If this value is less than some predefined alpha level, say 0.05, the variable is said to be statistically significant.

### Lower and Upper Confidence Limits

These provide a large-sample confidence interval for the values of the coefficients. The width of the confidence interval provides you with a sense of how precise the regression coefficients are. Also, if the confidence interval includes zero, the variable is not *statistically significant*. The formula for the calculation of the confidence interval is

$$b_i \pm z_{1-\alpha/2} s'_{b_i}$$

where  $1 - \alpha$  is the confidence coefficient of the confidence interval and  $z$  is the appropriate value from the standard normal distribution.

### Dispersion Phi

This is the estimate of the overdispersion correction multiplier, phi. Remember that in the Poisson model the mean and the variance are equal. In practice, the data almost always reject this restriction. Usually, the variance is greater than the mean—a situation called *overdispersion*. The increase in variance is represented in the model by a constant multiple of the variance-covariance matrix. That is, we use

$$V_{\hat{\beta}} = \phi \left( \sum_{i=1}^n \mu_i \mathbf{x}_i \mathbf{x}'_i \right)^{-1}$$

where  $\phi$  is estimated using

$$\hat{\phi} = \frac{1}{n-k} \sum_{i=1}^n \frac{(y_i - \hat{\mu}_i)^2}{\hat{\mu}_i}$$

## Poisson Regression

### Estimated Poisson Regression Model

This expression displays the estimated regression model in written format. It may be copied to the clipboard and used elsewhere. For example, you could copy it and paste it as a Variable Transformation.

Note that transformation must be less than 255 characters. Since this formula is often greater than 255 characters in length, you must use the FILE(filename) transformation. To do so, copy the formula to a text file using Notepad, Windows Write, or Word to receive the model text. Be sure to save the file as an unformatted text (ASCII) file. The transformation is FILE(filename) where *filename* is the name of the text file, including directory information. When the transformation is executed, it will load the file and use the transformation stored there.

### Rate Report

Independent Variable	Regression Coefficient b(i)	Rate Ratio Exp(b(i))	Lower 95.0% Confidence Limit	Upper 95.0% Confidence Limit
Intercept	-10.65831	0.00002	0.00002	0.00003
(Area=1)	0.81948	2.26933	1.97442	2.60830
(AgeGroup="35-44")	1.79737	6.03379	4.76055	7.64756
(AgeGroup="45-54")	1.91309	6.77397	5.37066	8.54396
(AgeGroup="54-64")	2.24180	9.41028	7.46233	11.86672
(AgeGroup="65-74")	2.36572	10.65175	8.23138	13.78381
(AgeGroup=">74")	2.94468	19.00457	14.67098	24.61823

This report provides the rate ratio for each independent variable.

#### Independent Variable

This item provides the name of the independent variable shown on this line of the report. The *Intercept* refers to the optional constant term.

#### Regression Coefficient

These are the maximum-likelihood estimates of the regression coefficients,  $b_1, b_2, \dots, b_k$ . Their direct interpretation is difficult because the formula for the predicted value involves the exponential function.

#### Rate Ratio

These are the exponentiated values of the regression coefficients. The formula used to calculate these is

$$RR_i = e^{b_i}$$

The rate ratio is mainly useful for interpretation of the regression coefficients of indicator variables. In this case, they estimate the incidence of the response variable (melanoma in this example) in the given category relative to the category whose indicator variable was omitted (usually called the *control* group).

#### Lower and Upper Confidence Limits

These provide a large-sample confidence interval for the rate ratios. The formula for the calculation of the confidence interval is

$$\exp(b_i \pm z_{1-\alpha/2} s'_{b_i})$$

where  $1 - \alpha$  is the confidence coefficient of the confidence interval and  $z$  is the appropriate value from the standard normal distribution.

## Lack-of-Fit Tests Section

Test	DF	Chi <sup>2</sup> Value	Prob Level
Pearson	5	6.12	0.295180
G Statistic	5	6.21	0.285867

These tests indicate whether there is a significant lack of fit to the data by the model.

This report provides the results of two goodness-of-fit tests. They indicate whether the current model adequately fits the data. The tests themselves are described in the Technical Details section of this chapter.

### Test

Indicates which of the two tests is shown on this line. Note that the *G* Statistic test is more accurate in small samples. The Pearson test is often used as a test for overdispersion.

### DF

Both of these tests are chi-square tests. This is the value of the degrees of freedom. It is equal to the number of observations minus the number of parameters in the regression model.

### Chi<sup>2</sup> Value

This is the value of the chi-square test statistic.

### Prob Level

This is the probability level of the test. The null hypothesis is that the model fits the data adequately. The alternative hypothesis is that the model is an inadequate representation of the data. If this probability level is less than some cutoff value such as 0.10 or 0.05, there is a significant lack of fit.

## Analysis of Deviance Section

Term Omitted	DF	Deviance	Increase From Model Deviance (Chi <sup>2</sup> )	Prob Level
All	1	968.0446		
Area	1	202.6602	124.22	0.000000
AgeGroup	5	875.1835	796.74	0.000000
None(Model)	7	78.4398		

This report is the Poisson regression analog of the analysis of variance table. It displays the results of a chi-square test used to test whether each of the individual terms in the regression are statistically significant after adjusting for all other terms in the model.

This report is not produced during a subset selection run.

Note that this report requires that a separate regression be run for each line. Thus, if the running time is too long, you might consider omitting this report.

### Term Omitted

This is the model term that is being tested. The test is formed by comparing the deviance statistic when the term is removed with the deviance of the complete model. Thus, the deviance when the term is left out of the model is shown.

The “All” line refers to the intercept-only model. The “None(Model)” refers to the complete model with no terms removed.

Note that it is usually not advisable to include an interaction term in a model when one of the associated main effects is missing—which is what happens here. However, in this case, we believe this to be a useful test.

## Poisson Regression

Note that the name may become very long, especially for interaction terms. These long names may misalign the report. You can force the rest of the items to be printed on the next line by using the *Stagger label and output* option in the Report Options tab. This should create a better looking report when the names are extra long.

### DF

This is the degrees of freedom of the  $\chi^2$  test displayed on this line.

### Deviance

The deviance is equal to minus two times the log likelihood achieved by the model being described on this line of the report. See the discussion given earlier in this chapter for a technical discussion of the deviance. A useful way to interpret the deviance is as the analog of the residual sum of squares in multiple regression. This value is used to create the difference in deviance that is used in the chi-square test.

### Increase From Model Deviance (Chi<sup>2</sup>)

This is the difference between the deviance for the model described on this line and the deviance of the complete model. This value follows the  $\chi^2$  distribution in medium to large samples. This value can be thought of as the analog of the residual sum of squares in multiple regression. Thus, you can think of this value as the increase in the residual sum of squares that occurs when this term is removed from the model.

Another way to interpret this test is as a redundancy test because it tests whether this term is redundant after considering all of the other terms in the model.

### Prob Level

This is the significance level of the chi-square test. This is the probability that a  $\chi^2$  value with degrees of freedom DF is equal to this value or greater. If this value is less than 0.05 (or other appropriate value), the term is said to be statistically significant.

---

## Log Likelihood & R<sup>2</sup> Report

Term(s) Omitted	DF	Log Likelihood	R <sup>2</sup> of Remaining Term(s)	Reduction From Model R <sup>2</sup>	Reduction From Saturated R <sup>2</sup>
All	1	-484.0223	0.0000		
Area	1	-101.3301	0.8544	0.1387	0.1456
AgeGroup	5	-437.5917	0.1037	0.8894	0.8963
None(Model)	7	-39.2199	0.9931	0.0000	0.0069
None(Saturated)	12	-36.1125	1.0000		0.0000

This report provides the log likelihoods and  $R^2$  values of various models. This report is not produced during a subset selection run.

Note that this report requires that a separate regression be run for each line. Thus, if the running time is too long, you might consider omitting this report.

### Term Omitted

This is the term that is omitted from the model. The “All” line refers to the intercept-only model. The “None(Model)” refers to the complete model with no terms removed. The “None(Saturated)” line gives the results for the saturated model.

Note that the name may become very long, especially for interaction terms. These long names may misalign the report. You can force the rest of the items to be printed on the next line by using *Stagger label and output if label length is ≥* option in the Report Options tab. This should create a better looking report when the names are extra long.

### DF

This is the degrees of freedom of the term displayed on this line.

## Poisson Regression

### Log Likelihood

This is the log likelihood of the model displayed on this line. Note that this is the log likelihood of the regression without the term listed.

### $R^2$ of Remaining Term(s)

This is the  $R^2$  of the model displayed on this line. Note that the model does not include the term listed at the beginning of the line.

Note that this is a pseudo  $R^2$  as discussed earlier in this chapter.

### Reduction From Model $R^2$

This is amount that  $R^2$  is reduced when the term is omitted from the regression model. This reduction is calculated from the  $R^2$  achieved by the full model.

This quantity is used to determine if removing a term causes a large reduction in  $R^2$ . If it does not, then the term can be safely removed from the model.

### Reduction From Saturated $R^2$

This is amount that  $R^2$  is reduced when the term is omitted from the regression model. This reduction is calculated from the  $R^2$  achieved by the saturated model. This item is included because it shows how removal of this term impacts the best  $R^2$  that is possible.

---

## Covariances of Regression Coefficients Section

The covariance matrix of the regression coefficients is not displayed as a report. However, it may be stored on the database for further investigation and use.

The covariance matrix is obtained by inverting the observed information matrix evaluated at the maximum likelihood estimates. If the Use Dispersion Phi option was checked, the original values are multiplied by phi.

---

## Residuals Report

Row	Melanoma (Y)	Predicted Value	Raw Residual	Pearson Residual	Deviance Residual	Population (T)
1	61	67.6998	-6.6998	-0.8143	-0.8283	2880262
2	76	80.0638	-4.0638	-0.4542	-0.4581	564535
3	98	94.4150	3.5850	0.3690	0.3667	592983
4	104	99.6974	4.3026	0.4309	0.4279	450740
5	63	67.8263	-4.8263	-0.5860	-0.5932	270908
6	80	72.2979	7.7021	0.9058	0.8904	161850
7	64	57.3002	6.6998	0.8851	0.8686	1074246
8	75	70.9362	4.0638	0.4825	0.4780	220407
9	68	71.5850	-3.5850	-0.4237	-0.4273	198119
10	63	67.3026	-4.3026	-0.5245	-0.5302	134084
11	45	40.1737	4.8263	0.7614	0.7469	70708
12	27	34.7021	-7.7021	-1.3075	-1.3609	34233

This report provides the predicted values and various types of residuals. Large residuals indicate data points that were not fit well by the regression model. You may consider removing rows with large residuals and refitting, but you must be certain that you have a good reason for doing so. You cannot remove them simply because they have large residuals.

### Row

The row number of the item. If you have excluded some rows by using a filter or if some of the rows had missing values, the row number identifies the original row on the database.

## Poisson Regression

### Y

This is the value of the dependent variable.

### Predicted Value

This is the predicted value of  $Y$ . It is the Poisson incidence rate,  $\hat{\mu}_i$ , estimated by

$$\begin{aligned} \hat{\mu}_i &= t_i \hat{\mu}(\mathbf{x}_i' \mathbf{b}) \\ &= t_i \exp(b_1 X_{1i} + b_2 X_{2i} + \dots + b_k X_{ki}) \end{aligned}$$

### Raw Residual

The raw residual is the different between the actual response and the estimated value from the model. The formula for the raw residual is

$$r_i = y_i - \hat{\mu}_i$$

### Pearson Residual

The Pearson residual corrects for the unequal variance in the residuals by dividing by the standard deviation. The formula for the Pearson residual is

$$p_i = \frac{y_i - \hat{\mu}_i}{\sqrt{\hat{\phi} \hat{\mu}_i}}$$

### Deviance Residual

The deviance residual is another popular residual. It is popular because the sum of squares of these residuals is the deviance statistic. The formula for the deviance residual is

$$d_i = \text{sign}(y_i - \hat{\mu}_i) \sqrt{2 \left\{ y_i \ln \left( \frac{y_i}{\hat{\mu}_i} \right) - (y_i - \hat{\mu}_i) \right\}}$$

where  $\text{sign}(x)$  is 1 if  $x$  is greater than or equal to 0 and -1 otherwise.

### T

The value of the exposure variable (if active) is provided for your reference.

---

## Predicted Values Report

Row	Melanoma (Y)	Predicted Value	Standard Error	Lower 95.0% Confidence Limit	Upper 95.0% Confidence Limit	Population (T)
1	61	67.6998	6.4440	55.0698	80.3297	2880262
2	76	80.0638	7.0419	66.2619	93.8657	564535
3	98	94.4150	7.8780	78.9743	109.8556	592983
4	104	99.6974	8.2257	83.5752	115.8195	450740
5	63	67.8263	6.7681	54.5610	81.0916	270908
6	80	72.2979	7.1850	58.2156	86.3802	161850
7	64	57.3002	5.5790	46.3656	68.2349	1074246
8	75	70.9362	6.3609	58.4691	83.4034	220407
9	68	71.5850	6.2636	59.3085	83.8615	198119
10	63	67.3026	5.9387	55.6630	78.9423	134084
11	45	40.1737	4.2609	31.8226	48.5249	70708
12	27	34.7021	3.7454	27.3612	42.0430	34233

This report provides the predicted values along with their standard errors and confidence limits.

If you want to generate predicted values and confidence limits for  $X$  values not on your database, you should add them to the bottom of the database, leaving  $Y$  blank (if you are using an exposure variable, set the value of  $T$  to a

## Poisson Regression

desired value). These rows will not be included in the estimation algorithm, but they will appear on this report with estimated  $Y$ 's.

### Row

The row number of the item. If you have excluded some rows by using a filter or if some of the rows had missing values, the row number identifies the original row on the database.

### Y

This is the value of the dependent variable.

### Predicted Value

This is the predicted value of  $Y$ . It is the predicted mean of the Poisson distribution,  $\hat{\mu}_i$ , estimated by

$$\begin{aligned}\hat{\mu}_i &= t_i \hat{\mu}(\mathbf{x}_i' \mathbf{b}) \\ &= t_i \exp(b_1 X_{1i} + b_2 X_{2i} + \cdots + b_k X_{ki})\end{aligned}$$

### Standard Error

The standard error of the predicted value is a measure of the precision of the estimated value. The formula for the standard error is

$$se_{\hat{\mu}_i} = \hat{\mu}_i \sqrt{\mathbf{x}_i' \mathbf{V}_{\hat{\beta}} \mathbf{x}_i}$$

where

$$\mathbf{V}_{\hat{\beta}} = \phi \left( \sum_{i=1}^n \hat{\mu}_i \mathbf{x}_i \mathbf{x}_i' \right)^{-1}$$

Note that if  $\phi$  is not used, it is set to one in the above formulas.

### Confidence Limits

These limits define a large-sample confidence interval for  $\mu_i$ . The formula is

$$\hat{\mu}_i \pm (z_{1-\alpha/2}) (se_{\hat{\mu}_i})$$

### T

The value of the exposure variable (if active) is provided for you reference.

## Residual Diagnostics Report

Row	Melanoma (Y)	Predicted Value	Raw Residual	Studentized Pearson Residual	Studentized Deviance Residual	Hat Diagonal
1	61	67.6998	-6.6998	-1.3095	-1.3321	0.6134
2	76	80.0638	-4.0638	-0.7361	-0.7425	0.6194
3	98	94.4150	3.5850	0.6303	0.6264	0.6573
4	104	99.6974	4.3026	0.7602	0.7548	0.6787
5	63	67.8263	-4.8263	-1.0285	-1.0411	0.6754
6	80	72.2979	7.7021	1.6939	1.6651	0.7140
7	64	57.3002	6.6998	1.3095	1.2852	0.5432
8	75	70.9362	4.0638	0.7361	0.7293	0.5704
9	68	71.5850	-3.5850	-0.6303	-0.6357	0.5481
10	63	67.3026	-4.3026	-0.7602	-0.7685	0.5240
11	45	40.1737	4.8263	1.0285	1.0089	0.4519
12	27	34.7021	-7.7021	-1.6939	-1.7632	0.4042
High Leverage Cutoff						1.166667

## Poisson Regression

This report provides the hat diagonals and studentized residuals. It allows you to study the leverage (influence) of each observation.

### Row

The row number of the item. If you have excluded some rows by using a filter or if some of the rows had missing values, the row number identifies the original row on the database.

### Y

This is the value of the dependent variable.

### Predicted Value

This is the predicted value of  $Y$ . It is the Poisson incidence rate,  $\hat{\mu}_i$ , estimated by

$$\begin{aligned}\hat{\mu}_i &= t_i \hat{\mu}(\mathbf{x}_i' \mathbf{b}) \\ &= t_i \exp(b_1 X_{1i} + b_2 X_{2i} + \cdots + b_k X_{ki})\end{aligned}$$

### Raw Residual

The raw residual is the difference between the actual response and the estimated value from the model. The formula for the raw residual is

$$r_i = y_i - \hat{\mu}_i$$

### Studentized Pearson Residual

The studentized Pearson residual is found by dividing the regular Pearson residual by the square root of one minus the hat diagonal. The formula is

$$(sp)_i = \frac{y_i - \hat{\mu}_i}{\sqrt{\hat{\phi} \hat{\mu}_i (1 - h_{ii})}}$$

### Studentized Deviance Residual

The studentized deviance residual is found by dividing the regular deviance residual by the square root of one minus the hat diagonal. The formula is

$$(sd)_i = \text{sign}(y_i - \hat{\mu}_i) \sqrt{\frac{2 \left\{ y_i \ln \left( \frac{y_i}{\hat{\mu}_i} \right) - (y_i - \hat{\mu}_i) \right\}}{1 - h_{ii}}}$$

### Hat Diagonal

This is the value of the influence measure,  $h_{ii}$ . The Hat matrix is used in residual diagnostics to measure the influence of each observation. The hat values,  $h_{ii}$ , are the diagonal entries of the Hat matrix which is calculated using

$$H = W^{1/2} X (X' W X)^{-1} X' W^{1/2}$$

where  $W$  is a diagonal matrix made up of  $\hat{\mu}_i$ .

The hat values should be studied to understand which observations have the greatest influence on the fitted regression coefficients. Large hat values are those that are larger than  $2k/n$ .



Poisson Regression

**Incidence Section when Exposure = 10000**

Row	Average Incidence Rate	Prob that Count is 5	Prob that Count is 10	Prob that Count is 15	Prob that Count is 20	Prob that Count is 25
1	2.3505	0.056990	0.000135	0.000000	0.000000	0.000000
2	14.1822	0.003313	0.062866	0.100093	0.030868	0.002778
3	15.9220	0.001037	0.035105	0.099684	0.054827	0.008800
4	22.1186	0.000011	0.001914	0.028111	0.079991	0.066422
5	25.0366	0.000001	0.000357	0.009747	0.051537	0.079521
6	44.6697	0.000000	0.000000	0.000000	0.000016	0.000457
7	5.3340	0.173603	0.024788	0.000297	0.000001	0.000000
8	32.1842	0.000000	0.000003	0.000332	0.006156	0.033343
9	36.1323	0.000000	0.000000	0.000036	0.001201	0.011606
10	50.1944	0.000000	0.000000	0.000000	0.000001	0.000034
11	56.8164	0.000000	0.000000	0.000000	0.000000	0.000001
12	101.3703	0.000000	0.000000	0.000000	0.000000	0.000000

This report gives the predicted incidence rate and Poisson probabilities for various counts.

**Row**

The row number of the item. If you have excluded some rows by using a filter or if some of the rows had missing values, the row number identifies the original row on the database.

**Average Incidence Rate**

This is the predicted incidence rate calculated using the formula

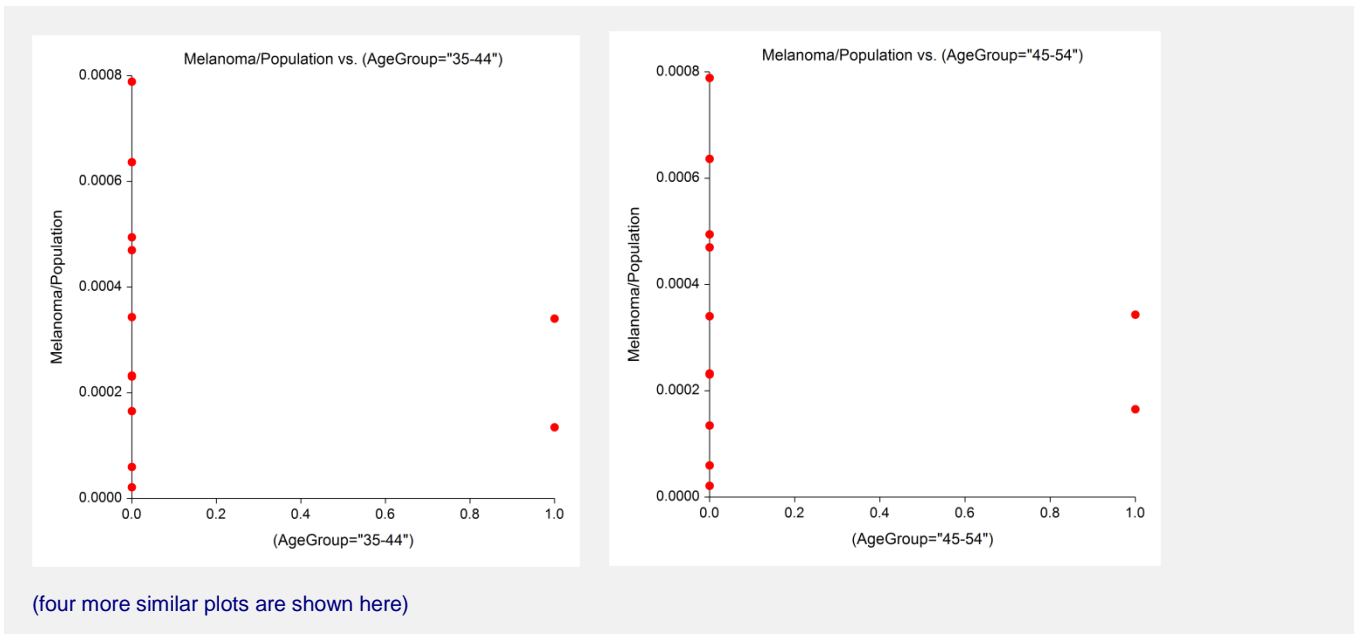
$$\hat{\mu}_i = T\hat{\mu}(\mathbf{x}'_i\mathbf{b})$$

Note that the calculation is made for a specific exposure value, not the value of *T* on the database. This allows you to make valid comparisons of the incidence rates.

**Prob that Count is Y**

Using the Poisson probability distribution, the probability of obtaining exactly *Y* events during the exposure amount given in the Exposure Value box is calculated for the values of *Y* specified in the Incidence Counts box.

**Incidence (Y/T) vs X's Plot(s)**

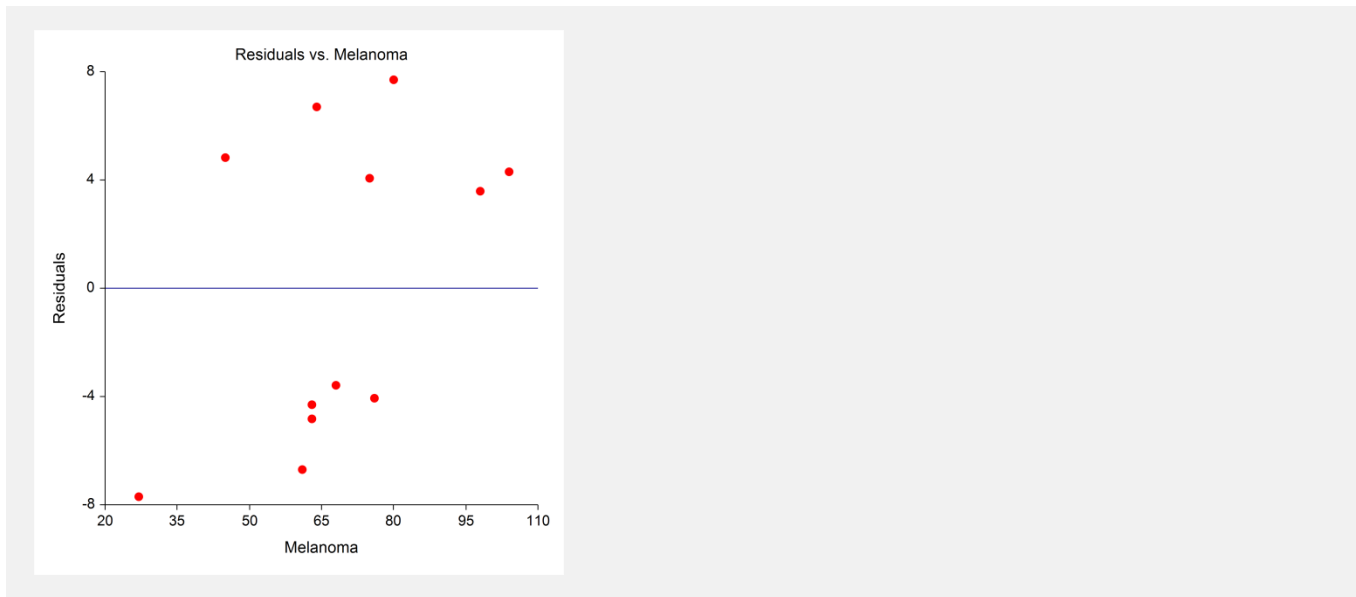


## Poisson Regression

These plots show each of the independent variables plotted against the incidence as measured by Y/T. They should be scanned for outliers and curvilinear patterns.

---

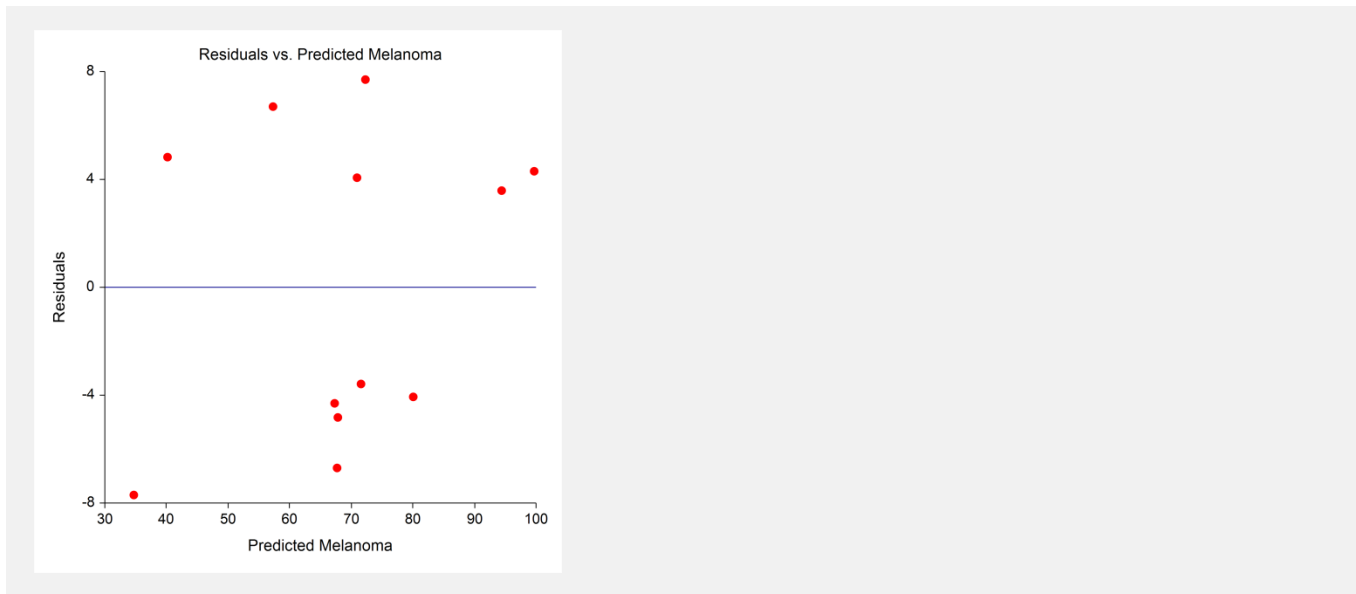
### Residuals vs Y Plot



This plot shows the residuals versus the dependent variable. It can be used to spot outliers.

---

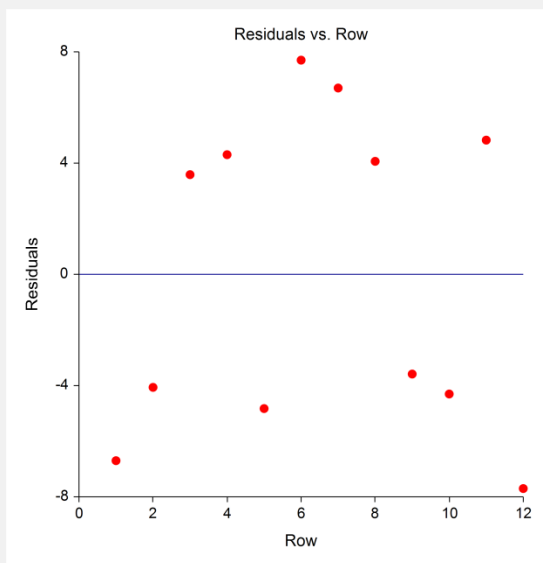
### Residuals vs Yhat Plot



This plot shows the residuals versus the predicted value (Yhat) of the dependent variable. It can show outliers.

---

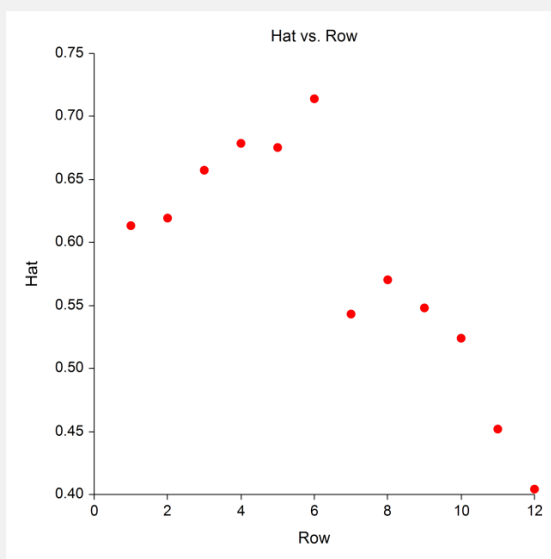
## Residuals vs Row Plot



This plot shows the residuals versus the row numbers. It is used to quickly spot rows that have large residuals.

---

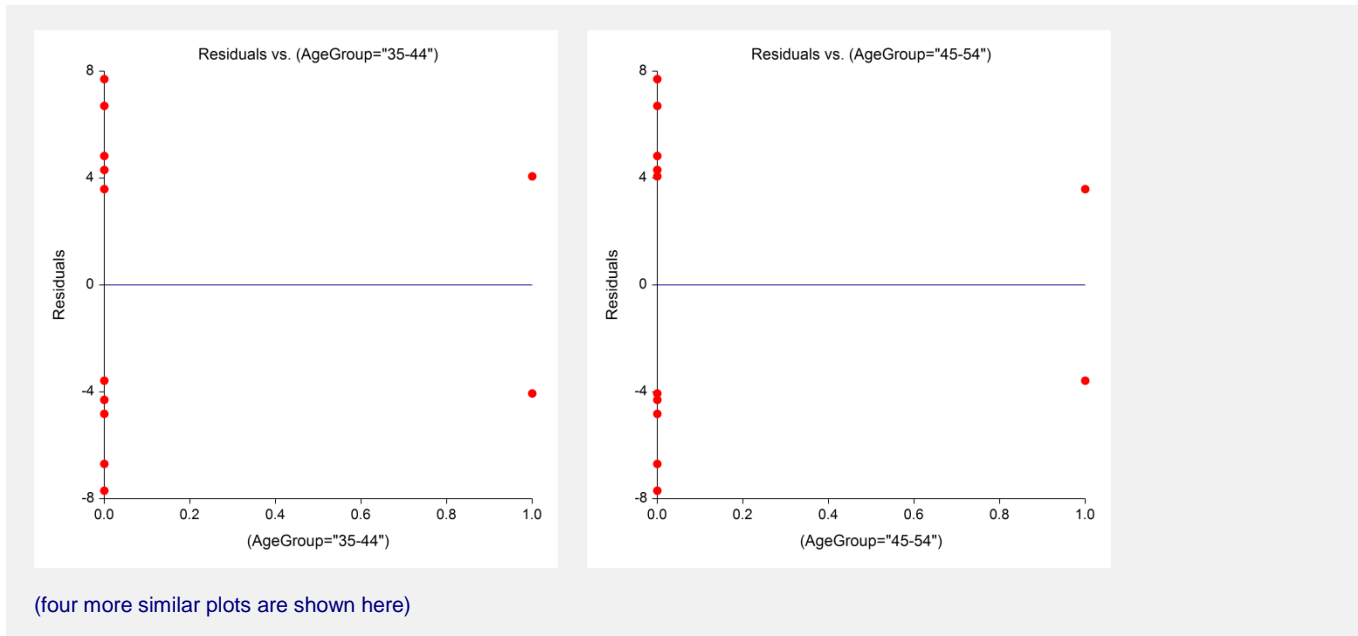
## Residuals vs Hat Plot



This plot shows the hat values versus the row numbers. It is used to quickly spot rows that have large hat values.

## Poisson Regression

## Residuals vs X's Plot(s)



These plots show the residuals plotted against the independent variables. They are used to spot outliers. They are also used to find curvilinear patterns that are not represented in the regression model.

## Example 2a – Subset Selection

This example will demonstrate how to select an appropriate subset of the independent variables that are available. The dataset to be analyzed consists of ten independent variables, a dependent variable, a frequency variable, and an exposure variable. The dependent variable was generated using independent variables X1, X2, and X3 using the formula

$$Count = Int[TimeExp(0.6 + 0.1X1 + 0.2X2 + 0.3X3)]$$

Variables X4, X5, and X6 were copies of X1 plus a small random component. Similarly, X7 and X8 were near copies of X2 and X9 and X10 were near copies of X3. These near copies of the original variables were added to cause confusion to the selection algorithm. The forty rows of data are stored in the PoisReg dataset.

Now we assume that we do not know how the data were generated. Our task is to find a subset of the ten independent variables that does a good job of fitting the data. We plan to make two runs. The goal of the first run will be to find an appropriate subset size. Then, in the second run, we will identify the variables in this subset and estimate the various regression statistics.

You may follow along here by making the appropriate entries or load the completed template **Example 2a** by clicking on Open Example Template from the File menu of the Poisson Regression window.

### 1 Open the PoisReg dataset.

- From the File menu of the NCSS Data window, select **Open Example Data**.
- Click on the file **PoisReg.NCSS**.
- Click **Open**.

### 2 Open the Poisson Regression window.

- Using the Analysis menu or the Procedure Navigator, find and select the **Poisson Regression** procedure.
- On the menus, select **File**, then **New Template**. This will fill the procedure with the default template.

### 3 Specify the variables.

- On the Poisson Regression window, select the **Variables, Model** tab.
- Set the **Dependent Y** to **Count**.
- Set the **T: Exposure Variable** to **Time**.
- Set the **X's: Numeric Independent Variables** to **X1-X10**.
- Set the **Frequencies** to **Cases**.

### 4 Specify the model.

- Set the **Terms** to **1-Way**.
- Set the **Search Method** to **Hierarchical Forward with Switching**.
- Set **Stop search when number of terms reaches** to **6**.
- The rest of this panel can be left at the default values.

### 5 Specify the reports.

- Select the **Reports** tab.
- Uncheck all of the reports and plots except **Run Summary**, **Subset Selection - Summary**, and **Subset Selection - Detail** (these should be checked).

### 6 Run the procedure.

- From the Run menu, select **Run Procedure**. Alternatively, just click the green Run button.

## Poisson Regression

## Run Summary

Item	Value	Item	Value
Dependent Variable	Count	Rows Used	40
Exposure Variable	Time	Sum of Frequencies	130
Frequency Variable	Cases	Iterations	20
Ind. Var's Available	10	Convergence Zero	1E-09
No. of X's in Model	5	Maximum Convergence	1.489973E-06
Pseudo R <sup>2</sup>	0.9980	Dispersion Phi	0.0138
Final Likelihood	-288.8153	Phi was not used to correct standard errors.	
Subset Method	Hierarchical Forward/Switching		

This report provides several details about the data and the MLE algorithm as it fit the best model found during the search. We note that, as expected, there were 40 rows used. The fact that 20 iterations were needed to solve the likelihood equations is a source of concern because this shows that the algorithm may not have converged. This may have been due to our fitting of a model that had too many terms.

## Subset Selection Summary Section

Number of Terms	Log Likelihood	R <sup>2</sup>	Deviance	AIC
1	-730.6939	0.0000	885.5007	887.5007
2	-434.0619	0.6700	292.2366	296.2366
3	-348.4077	0.8634	120.9282	126.9282
4	-288.8552	0.9979	1.8233	9.8233
5	-288.8343	0.9980	1.7815	11.7815
6	-288.8153	0.9980	1.7434	13.7434

This report will help us determine an appropriate subset size. By scanning each column, we can see that three variables are needed. All of these measures are functions of each other. However, they each offer insight into the appropriate subset size.

In this example, the four measures unanimously point to three as the appropriate subset size.

**Number of Variables**

This is the number of terms in the model including the intercept. Each line presents the results for the best model found for that subset size. The first line presents the results for the intercept-only model.

**Log Likelihood**

This is the value of the log likelihood function. Since the goal of maximum likelihood is to maximize this value, we want to select a subset size after which the log likelihood is not increased significantly.

In this example, after three terms are added (in addition to the intercept) the log likelihood does not change a great deal. The log likelihood points to a subset size of three terms plus the intercept for a total of four.

**R<sup>2</sup>**

This is the value of pseudo  $R^2$ —a measure of the adequacy of the model. Since our goal is to maximize this value, we want to select a subset size after which the this value is not increased significantly.

In this example, after four terms are included, the  $R^2$  is 0.9979 and it does not change a great deal. The  $R^2$  values point to a subset size of four.

**Deviance**

Deviance is a measure of the lack of fit. Hence, we want to select a subset size after which the deviance is not significantly decreased.

In this example, after four terms are included, the Deviance is 1.8233 and it does not change a great deal. The Deviance values point to a subset size of four.

## Poisson Regression

**AIC**

These are the Akaike information criterion values for each subset size. This criterion measures both the lack of fit and the size of the regression model. Our goal is to minimize this value.

In this example, the subset size of four gives the lowest value AIC and is thus the subset size implied by this statistic.

**Subset Selection Detail Section**

Step	Action	No. of Terms	No. of X's	Log Likelihood	R <sup>2</sup>	Term Entered	Term Removed
1	Add	1	1	-730.6939	0.0000	Intercept	
2	Add	2	2	-434.0619	0.6700	X3	
3	Add	3	3	-348.4423	0.8634	X2	
4	Switch	3	3	-348.4077	0.8634	X9	X3
5	Add	4	4	-289.2634	0.9970	X6	
6	Switch	4	4	-289.0943	0.9974	X8	X2
7	Switch	4	4	-288.8552	0.9979	X3	X9
8	Add	5	5	-288.8343	0.9980	X5	
9	Add	6	6	-288.8201	0.9980	X7	
10	Switch	6	6	-288.8153	0.9980	X2	X5

This report shows the progress of the subset selection algorithm through its various steps. It shows the original term added at each step and any switching that was done.

**Step**

This is the number of the step in the subset selection process.

**Action**

Two actions are possible at each step: Add or Switch. *Add* means that the subset size was increased and the term entered as added to the set of active regressor variables. *Switch* means that the subset size remained the same while one active regressor was removed and another was activated.

**No. of Terms**

This is the number of active terms (including the intercept) at the end of this step.

**No. of X's**

This is the number of active variables (excluding the intercept) at the end of this step. This reminds you of how many X variables were generated for each term involving a categorical variable.

**Log Likelihood**

This is the value of the log likelihood after this step was completed.

**R<sup>2</sup>**

This is the pseudo  $R^2$  value after this step was completed.

**Variable Entered**

This is the name of the regressor that was added to the list of active regressor variables.

**Variable Removed**

In switching steps, this is the name of the variable that was removed from the list of active regressor variables.

## Poisson Regression

## Example 2b – Subset Selection Continued

Example 2a completed the first step in the subset selection process by indicating that a subset of four terms is appropriate. Now, a second run must be made to find those terms.

The instructions provide here assume that you have just completed Example 2a. If you have not, you must complete it first since we will only tell you what needs to be changed.

You may follow along here by making the appropriate entries or load the completed template **Example 2b** by clicking on Open Example Template from the File menu of the Poisson Regression window.

### 1 Specify the model.

- On the Poisson Regression window, select the **Model tab**.
- Set the **Stop search when number of terms reaches to 4**.
- The rest of this panel can be left at the default values.

### 2 Specify the reports.

- Select the **Reports tab**.
- Uncheck all of the reports and plots except **Run Summary, Subset Selection - Summary, Subset Selection – Detail, Regression Coefficients, and Residuals** (these should be checked).

### 3 Run the procedure.

- From the Run menu, select **Run Procedure**. Alternatively, just click the Run button (the left-most button on the button bar at the top) or press the F9 function key.

## Run Summary Report

Item	Value	Item	Value
Dependent Variable	Count	Rows Used	40
Exposure Variable	Time	Sum of Frequencies	130
Frequency Variable	Cases	Iterations	10
Ind. Var's Available	10	Convergence Zero	1E-09
No. of X's in Model	3	Maximum Convergence	5.654499E-10
Pseudo R <sup>2</sup>	0.9979	Dispersion Phi	0.0142
Final Likelihood	-288.8552	Phi was not used to correct standard errors.	
Subset Method	Hierarchical Forward/Switching		

We note that the final model converged in only five iterations and the Maximum Convergence is less than Convergence Zero. This means that the algorithm terminated normally.

## Subset Selection Summary

Number of Terms	Log Likelihood	R <sup>2</sup>	Deviance	AIC
1	-730.6939	0.0000	885.5007	887.5007
2	-434.0619	0.6700	292.2366	296.2366
3	-348.4077	0.8634	120.9282	126.9282
4	-288.8552	0.9979	1.8233	9.8233

This report again shows us that a subset size of four is a reasonable choice.



## Poisson Regression

## Subset Selection Detail

Step	Action	No. of Terms	No. of X's	Log Likelihood	R <sup>2</sup>	Term Entered	Term Removed
1	Add	1	1	-730.6939	0.0000	Intercept	
2	Add	2	2	-434.0619	0.6700	X3	
3	Add	3	3	-348.4423	0.8634	X2	
4	Switch	3	3	-348.4077	0.8634	X9	X3
5	Add	4	4	-289.2634	0.9970	X6	
6	Switch	4	4	-289.0943	0.9974	X8	X2
7	Switch	4	4	-288.8552	0.9979	X3	X9

This report shows the algorithm's journey through the maze of possible models. During the process, three variables were switched in order to achieve a better model.

## Regression Coefficients Report

Independent Variable	Regression Coefficient b(i)	Standard Error Sb(i)	Wald's Chi <sup>2</sup> H0: $\beta=0$	Prob Level	Lower 95.0% Confidence Limit	Upper 95.0% Confidence Limit
Intercept	-0.12374	0.10638	1.35	0.2448	-0.33224	0.08476
X3	0.01047	0.00041	656.32	0.0000	0.00967	0.01127
X6	0.00345	0.00031	121.68	0.0000	0.00283	0.00406
X8	0.00677	0.00043	245.70	0.0000	0.00592	0.00761
Dispersion Phi		0.0142				

This report provides the details of the model that was selected. We note the X3, X6, and X8 were included in the model. We assume that X8 is taking the place of X2 and X6 is taking the place of X1. In fact, we ran a Poisson regression with X1, X2, and X3 in the model. The log likelihood for this model was -288.9466, which is slightly less than the -288.8552 achieved by our best model. This concludes our discussion of this example. Usually, we would go on to study the residual plots and complete the analysis by making a third run with only the variables X3, X6, and X8 specified.