



Analyzing count data and correcting for over-dispersion

Contents

1. Example: Analyzing counts from a complex sampling design	1
1.1 Setting up the analysis through the GUI.....	1
1.2 Discussion of results – Poisson-log model	7
1.3 Estimated outcomes for different groups.....	8
2. Ignoring stratification and clustering in the sample.....	9
2.1 Setting up the analysis.....	9
2.2 Discussion of results	10
3. Example: Correcting for over-dispersion in an analysis of counts	11
3.1 Setting up the analysis.....	11
3.2 Discussion of results – negative Binomial model.....	13

1. Example: Analyzing counts from a complex sampling design

A question that a researcher may want to address is whether ethnicity and gender effects are associated with the number of substance abuse diagnoses. An appropriate statistical model for this type of count variable is a GLIM with a Poisson distribution and a log link function.

1.1 Setting up the analysis through the GUI

The first step is to open the LSF shown above in the LISREL LSF window. Use the **Open** option on the **File** menu of the root window of LISREL to load the **Open** dialog box. Select the **Lisrel Data (*.lsf)** option from the **Files of type** drop-down list box and browse for the file **cntdiag.psf** in the **Complex Survey GLIM examples** folder. Click on the **Open** button to open the file **cntdiag.lsf**.

	cntdiag	sex	race_d	CENREG	FACTYPE	A2TWA0
1	1.00	0.00	0.00	4.00	4.00	190.70
2	0.00	0.00	0.00	4.00	4.00	44.30
3	0.00	0.00	0.00	4.00	4.00	44.30
4	0.00	1.00	0.00	4.00	4.00	44.30
5	0.00	0.00	0.00	4.00	4.00	44.30
6	0.00	0.00	0.00	4.00	4.00	44.30
7	0.00	0.00	0.00	4.00	4.00	44.30
8	0.00	0.00	0.00	4.00	4.00	44.30
9	0.00	0.00	0.00	4.00	4.00	44.30
10	0.00	0.00	1.00	4.00	2.00	371.90
11	0.00	0.00	1.00	4.00	2.00	371.90
12	0.00	0.00	1.00	4.00	2.00	371.90
13	0.00	0.00	1.00	4.00	2.00	371.90
14	0.00	0.00	1.00	4.00	2.00	371.90
15	0.00	0.00	0.00	4.00	2.00	371.90

The next step is to enter the model specifications into the sequence of the four SurveyGLIM GUI dialog boxes. The **Title and Options** dialog box is the first dialog box and is accessed by selecting the **Title and Options** option on the **SurveyGLIM** menu above. In order to identify the analysis, enter the string **Poisson-Log Model for ADSS Data** into the **Title** string field to produce the following **Title and Options** dialog box.

Title and Options ×

Title:

Maximum Number of Iterations:

Convergence Criterion:

Missing Data Value:

Suppress Iterative Details Variance Adjustment

Response Variable Ordering _____

Ascending Descending

Reference Category Code _____

0 -1

Optimization Method _____

Fisher-Scoring Newton-Raphson

Additional Output _____

Residual file Data file

To build syntax, proceed to the Survey Design screen and click the Finish button

The default options will be used for this example. Click the **Next** button to access the **Distributions and Links** dialog box. Since we intend to fit a Poisson-log model, select the **Poisson** option from the **Distribution type** drop-down list box. For this example, we will estimate the scale parameter of the model by using the Pearson χ^2 estimate. Select the **Pearson** option from the **Estimate scale?** drop-down list box to produce the following **Distributions and Links** dialog box.

Distributions and Links X

Distribution type: Poisson

Link function: Log

Include intercept? Yes No

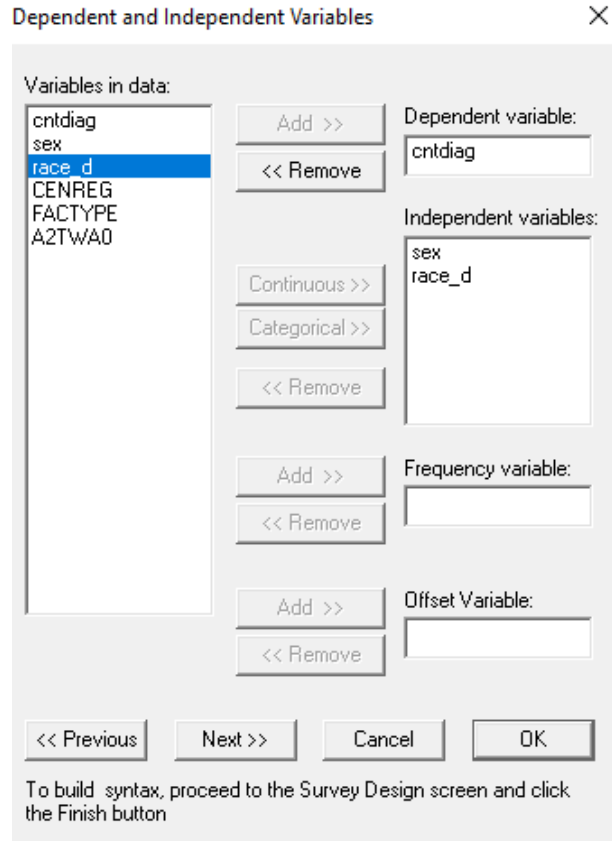
Estimate dispersion? Yes Fixed value:

Estimate scale? Pearson

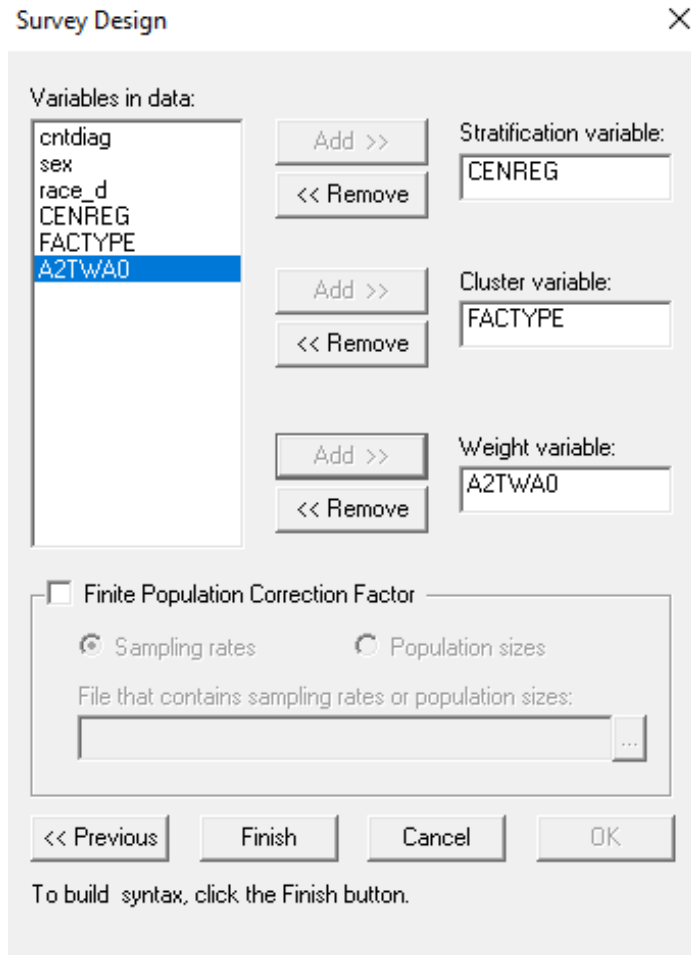
<< Previous Next >> Cancel OK

To build syntax, proceed to the Survey Design screen and click the Finish button

Proceed to the **Dependent and Independent Variables** dialog box by clicking on the **Next** button. Specify the response variable cntdiag by selecting it from the **Variables in data** list box and clicking on the **Add** button of the **Dependent variable** section. In a similar fashion, add the covariates sex and race_d to the **Independent variables** list box to produce the following **Dependent and Independent Variables** dialog box.



Since the data are not frequency table data and no offset variable is used for this example, go to the **Survey Design** dialog box by clicking on the **Next** button. The strata are the census regions (CENREG) and are specified by selecting the variable CENREG from the **Variables in data** list box and clicking on the **Add** button of the **Stratification variable** section. Similarly, add the PSU variable FACTYPE and the design weight variable A2TWA0 to the **Cluster variable** and **Weight variable** boxes respectively to produce the following **Survey Design** dialog box.



As no finite population information is available, we proceed to click on the **Finish** button to open the following text editor window for **cntdiag.prl**.

```

LISREL for Windows - [cntdiag.PRL]
File Edit Options Window Help
GlimOptions Converge=0.0001 MaxIter=100 MissingCode=-999999 Response=Ascending
RefCatCode=-1 IterDetails=No
Method=Fisher;
Title=Poisson-Log model for ADSS data;
SY='cntdiag.lsf';
Distribution=POI;
Link=LOG;
Intercept=Yes;
Scale=Pearson;
DepVar=cntdiag;
CoVars=sex race_d;
Stratum=CENREG;
Cluster=FACTYPE;
Weight=A2TWA0;
Ready

```

We are now ready to submit the GLIM analysis. This is achieved by clicking on the **Run Prelis** toolbar icon to produce the text editor window for **cntdiag.out**.

1.2 Discussion of results – Poisson-log model

A portion of the results of the Poisson-log GLIM analysis is shown in the following text editor window.

The screenshot shows a text editor window titled "cntdiag.OUT" with the following content:

Statistic	Value	Den. DF	Num. DF	P Value
Adjusted Wald F	2.8314	2	7	0.125599
Wald Chi-square	6.4718	2		0.125599

Note: The Wald F Test and Chi-square Statistics are statistics to test the null hypothesis that all the regression weights are equal to zero.

Estimated Regression Weights

Parameter	Estimate	Standard Error	z Value	P Value
intcept	0.3302	0.0557	5.9248	0.0000
sex	0.0619	0.0709	0.8726	0.3829
race_d	0.1167	0.0620	1.8818	0.0599
SCALE	0.7479			

Note: The scale parameter estimate is based on the Pearson Chi-square value
 $\phi = \text{Square Root of } (\text{The Pearson Chi-square value/degrees of freedom})$

SurveyGLIM reports the Adjusted Wald F and χ^2 test statistic values for testing the null hypothesis that all the regression weights are equal to zero which may be expressed as (*cf.* American Institutes for Research & Cohen, 2003)

$$F_w = \frac{\left(\sum_{h=1}^H n_h - H - r + 1 \right)}{\left(\sum_{h=1}^H n_h - H \right) * r} \hat{\beta}' \hat{\Upsilon}^{-1} \hat{\beta}$$

and

$$X_w^2 = \hat{\beta}' \hat{\Upsilon}^{-1} \hat{\beta}$$

respectively where H denotes the number of strata, $\sum_{h=1}^H n_h$ denotes the number of PSUs, r denotes the number of covariates of the model, $\hat{\boldsymbol{\beta}}$ denotes the estimate of the parameter vector, $\boldsymbol{\beta}$, of regression weights and $\hat{\mathbf{Y}}$ denotes the estimated asymptotic covariance matrix of the estimators of the elements of $\boldsymbol{\beta}$. If the null hypothesis is correct, F_w and X_w^2 approximately follow an F distribution with r and $\sum_{h=1}^H n_h - H - r + 1$ degrees of freedom and a χ^2 distribution with r degrees of freedom respectively.

Both the values of the Wald F and χ^2 test statistics are not statistically significant if a significance level of 5% is used. Hence, there is insufficient evidence to conclude that both gender and race influence the number of diagnoses of a client. This finding is supported by the non-significant z test statistic values for the significance of the individual parameters.

The scale parameter estimate is less than unity which indicates under-dispersion for the response variable. In other words, the sample variance of the variable cntdiag is less than its mean.

1.3 Estimated outcomes for different groups

The fitted model follows from the output file above as

$$\hat{E}[\text{cntdiag}_k] = \exp(0.33 + 0.06 * \text{sex}_k + 0.12 * \text{race}_k)$$

Although gender and race did not significantly affect the number of diagnoses, the following examples illustrate how the fitted model can be used to calculate the mean of number of diagnoses for various subgroups when there are statistically significant differences among them. This fitted model implies that the mean number of diagnoses for a white female client ($\text{sex}_k = 1$ and $\text{race}_k = 1$) is given by

$$\exp(0.33 + 0.06 + 0.12) = \exp(0.51) = 1.67$$

Similarly, the mean number of diagnoses for a nonwhite female client ($\text{sex}_k = 1$ and $\text{race}_k = 0$) is 1.48.

It also follows from the output above that $\exp(\hat{\beta}_1) = \exp(0.06) = 1.06$ is the multiplicative effect of gender on the fitted number of diagnoses for a client. This implies that, on the average, female clients have a 6% higher estimated mean number of diagnoses than male clients. Similarly, it follows that $\exp(\hat{\beta}_2) = \exp(0.12) = 1.13$ which implies that, on the average, the fitted number of diagnoses is 13% higher for white clients than for nonwhite clients.

2. Ignoring stratification and clustering in the sample

2.1 Setting up the analysis

The stratification and clustering can be ignored by not specifying the stratification and cluster variables on the **Survey Design** dialog box. However, it is recommended to change the title of the analysis to distinguish it from the previous analysis. This is done by selecting the **Title and Options** option on the **SurveyGLIM** menu to go to the **Title and Options** dialog box and then by entering the string **Fitting a Poisson-Log model with design weights only** in the **Title** string field. Since our model remains the same, click on the **Next** buttons of the **Title and Options**, the **Distributions/Links** and the **Dependent and Independent Variables** dialog boxes respectively to go to the **Survey Design** dialog box. Remove the stratification and cluster variables by clicking on the **Remove** buttons of the **Stratification variable** and **Cluster variable** sections to produce the following **Survey Design** dialog box.

Survey Design

Variables in data:

- cntdiag
- sex
- race_d
- CENREG**
- FACTYPE
- A2TWA0

Add >> << Remove

Stratification variable:

Add >> << Remove

Cluster variable:

Add >> << Remove

Weight variable:

A2TWA0

Finite Population Correction Factor

Sampling rates Population sizes

File that contains sampling rates or population sizes:

<< Previous Finish Cancel OK

To build syntax, click the Finish button.

As this completes our modifications, click on the **Finish** button to open the following text editor window for **cntdiag.prl**.

```

GlimOptions Converge=0.0001 MaxIter=100 MissingCode=-999999 Response=Ascending
Method=Fisher;
Title=Poisson-Log model for ADSS data;
SY='cntdiag.lsf';
Distribution=POI;
Link=LOG;
Intercept=Yes;
Scale=Pearson;
DepVar=cntdiag;
CoVars=sex race_d;
Weight=A2TWA0;|

```

As before, submit the analysis by clicking on the **Run Prelis** toolbar icon to produce the text editor window for **cntdiag.out**.

2.2 Discussion of results

A portion of the text editor window for **cntdiag.out** is shown below.

```

Estimated Regression Weights

Parameter      Estimate      Standard
-----      -
intcept        0.3302        0.0009
sex            0.0619        0.0016
race_d         0.1167        0.0015
SCALE          0.7479
z Value        P Value
-----      -
intcept        367.2546      0.0000
sex            39.3322      0.0000
race_d         75.4532      0.0000

Note: The scale parameter estimate is based on the Pearson Chi-square value
      phi = Square Root of (The Pearson Chi-square value/degrees of freedom)

```

The results above indicate that although the parameter estimates are identical to those obtained when the design of the complex survey was taken into account, the standard error estimates are significantly smaller (*cf.* Brogan, 1998). As a consequence, both gender and race appear to have a statistically significant effect on the number of substance abuse diagnoses at a $p < 0.00001$ level of confidence. This is a reversal of the results obtained when the complex sampling design was taken into account. As this example indicates, inferences based on an analysis that does not correct for the reduced precision of a complex sampling design can be very misleading.

3. Example: Correcting for over-dispersion in an analysis of counts

The results for the Poisson-log model indicated the presence of under-dispersion. Although the negative Binomial distribution is intended for dealing with over-dispersion, we will use it here for illustrative purposes.

3.1 Setting up the analysis

In order to fit the Negative Binomial-log model interactively to the data in **cntdiag.psf**, we only need to re-specify the sampling distribution. As in the previous analysis, start by modifying the title to **Fitting a Negative Binomial-Log model** by accessing the **Title and Options** dialog box and clicking the **Next** button to go to the **Distributions and Links** dialog box. Select the **Negative Binomial** option from the **Distribution** drop-down list box to produce the following **Distributions and Links** dialog box.

Distributions and Links

Distribution type: Negative Binomial

Link function: Log

Include intercept? Yes No

Estimate dispersion? Yes Fixed value:

Estimate scale? Pearson

<< Previous Next >> Cancel OK

To build syntax, proceed to the Survey Design screen and click the Finish button

Since the rest of the model remains the same, click on the **Next** buttons of the **Distributions and Links** and the **Dependent and Independent Variables** dialog boxes respectively to go to the **Survey Design** dialog box. Specify the complex survey design again by selecting the variables **CENREG** and **FACTYPE** from the **Variables in data** list box and clicking on the **Add** buttons of the **Stratification variable** and **Cluster variable** sections respectively to produce the following **Survey Design** dialog box.

X

Survey Design

Variables in data:

cntdiag	Add >>	Stratification variable:
sex	<< Remove	CENREG
race_d		
CENREG		
FACTYPE	Add >>	Cluster variable:
A2TWA0	<< Remove	FACTYPE
	Add >>	Weight variable:
	<< Remove	A2TWA0

Finite Population Correction Factor

Sampling rates
 Population sizes

File that contains sampling rates or population sizes:

...

To build syntax, click the Finish button.

Click on the **Finish** button to open the following text editor window for **cntdiag.prl**.

```

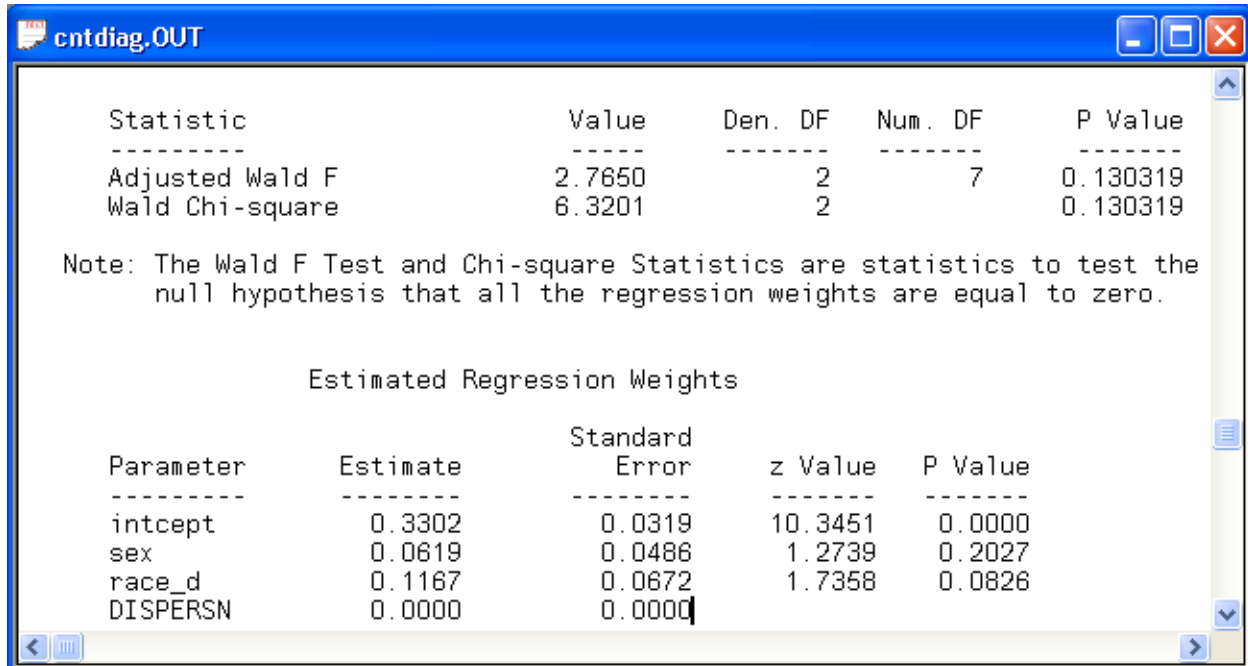
File Edit Options Window Help
GlimOptions Converge=0.0001 MaxIter=100 MissingCode=-999999 Response=Ascending
Method=Fisher;
Title=Poisson-Log model for ADSS data;
SY='cntdiag.lsf';
Distribution=NBIN;
Link=LOG;
Intercept=Yes;
Dispersion=Yes;
DepVar=cntdiag;
CoVars=sex race_d;
Stratum=CENREG;
Cluster=FACTYPE;
Weight=A2TWA0;
Ready CAP SCRI

```

Submit the analysis by clicking on the **Run Prelis** toolbar icon to open the text editor window for the corresponding output file **cntdiag.out**.

3.2 Discussion of results – negative Binomial model

A portion of the text editor window for **cntdiag.out** is shown below.



```
Statistic          Value      Den. DF  Num. DF  P Value
-----          -
Adjusted Wald F    2.7650      2        7      0.130319
Wald Chi-square    6.3201      2        7      0.130319

Note: The Wald F Test and Chi-square Statistics are statistics to test the
      null hypothesis that all the regression weights are equal to zero.

      Estimated Regression Weights

Parameter      Estimate      Standard      z Value      P Value
-----      -
intcept        0.3302        0.0319        10.3451      0.0000
sex            0.0619        0.0486         1.2739      0.2027
race_d         0.1167        0.0672         1.7358      0.0826
DISPERSN       0.0000        0.0000
```

A comparison of these results with those obtained for the Poisson-log model shows that the estimates are the same, but that the standard error estimates are different. However, the conclusions are the same as those made based on the results for the Poisson-log model.

The zero estimate of the dispersion parameter of the Negative Binomial distribution indicates that over-dispersion seen with the Poisson distribution does not apply to this particular analysis. This finding is in agreement with the Poisson scale estimate less than unity, which indicated the presence of under-dispersion rather than over-dispersion.