

Logistic Regression, Part III: Hypothesis Testing, Comparisons to OLS

Richard Williams, University of Notre Dame, <https://www3.nd.edu/~rwilliam/>

Last revised February 22, 2015

This handout steals heavily from Linear probability, logit, and probit models, by John Aldrich and Forrest Nelson, paper # 45 in the Sage series on Quantitative Applications in the Social Sciences; and Applied Logistic Regression Analysis Second Edition by Scott Menard, paper # 106 in that series. This handout primarily uses Stata; an older version of the handout that used SPSS may also be available.

WARNING: As Menard more or less points out, notation is wildly inconsistent across authors and programs when it comes to Logistic regression. I'm trying to more or less follow Menard, but you'll have to learn to adapt to whatever the author or statistical program happens to use.

Overview. In this handout, we'll examine hypothesis testing in logistic regression and make comparisons between logistic regression and OLS. A separate handout provides more detail about using Stata. The optional appendix to this handout provides more detail on how some of the key calculations are done.

There are a number of logical analogs between OLS and Logistic regression, i.e. the math is different but the functions served are similar. I will summarize these first, and then explain each of them in more detail:

OLS Regression	Logical Analog in Logistic Regression
Total Sums of Squares	$-2LL_0$, DEV_0 , D_0
Error/ Residual Sums of Squares	$-2LL_M$, DEV_M , D_M
Regression/Explained Sums of Squares	Model Chi Square, L^2 , G_M
Global F	Model Chi Square, L^2 , G_M
Incremental F Test	Chi-Square Contrast/ Incremental chi-square contrast
Incremental F Test and Wald test of the same hypotheses give identical results	Chi-square contrast between models and a Wald test of the same hypotheses generally do NOT give exactly identical results.

Using the same data as before, here is part of the output we get in Stata when we do a logistic regression of Grade on Gpa, Tuce and Psi.

```
. use https://www3.nd.edu/~rwilliam/statafiles/logist.dta, clear
. logit grade gpa tuce psi

Iteration 0:   log likelihood =  -20.59173
Iteration 1:   log likelihood = -13.496795
Iteration 2:   log likelihood = -12.929188
Iteration 3:   log likelihood = -12.889941
Iteration 4:   log likelihood = -12.889633
Iteration 5:   log likelihood = -12.889633

Logistic regression                               Number of obs   =           32
                                                  LR chi2(3)      =           15.40
                                                  Prob > chi2     =           0.0015
                                                  Pseudo R2      =           0.3740

Log likelihood = -12.889633
[Rest of output deleted]
```

Global tests of parameters. In OLS regression, if we wanted to test the hypothesis that all β 's = 0 versus the alternative that at least one did not, we used a global F test. In logistic regression, we use a *likelihood ratio chi-square test* instead. Stata calls this LR chi2. The value is 15.404. This is computed by contrasting a model which has no independent variables (i.e. has the constant only) with a model that does. Following is a general description of how it works; the appendix provides a detailed example.

The probability of the observed results given the parameter estimates is known as the *likelihood*. Since the likelihood is a small number less than 1, it is customary to use -2 times the log of the likelihood. -2LL is a measure of how well the estimated model fits the likelihood. A good model is one that results in a high likelihood of the observed results. This translates to a small number for -2LL (If a model fits perfectly, the likelihood is 1, and -2 times the log likelihood is 0).

-2LL is also called the Deviance, DEV, or simply D. Subscripts are often used to denote which model this particular deviance applies to. The smaller the deviance is, the better the model fits the data.

The “initial log likelihood function” is for a model in which only the constant is included. This is used as the baseline against which models with IVs are assessed. Stata reports LL_0 , -20.59173, which is the log likelihood for iteration 0. $-2LL_0 = -2 * -20.59173 = 41.18$.

$-2LL_0$, DEV_0 , or simply D_0 are alternative ways of referring to the deviance for a model which has only the intercept. This is analogous to the Total Sums of Squares, SST, in OLS Regression.

When GPA, PSI, and TUCE are in the model, $-2LL_M = -2 * -12.889633 = 25.78$. We can refer to this as DEV_M or simply D_M .

The $-2LL$ for a model, or DEV_M , indicates the extent to which the model fails to perfectly predict the values of the DV, i.e. it tells how much improvement is needed before the predictors provide the best possible prediction of the dependent variable. DEV_M is analogous to the Error Sums of Squares, SSE, in OLS regression.

The addition of these 3 parameters reduces $-2LL$ by 15.40, i.e.
 $DEV_0 - DEV_M = 41.183 - 25.779 = 15.40$. This is reflected in the *Model Chi-square*, which Stata labels as LR chi2.

The Model Chi-Square, also called Model L^2 or G_M , is analogous to the Regression (explained) Sums of Squares, SSR, in OLS regression. It is also the direct counterpart to the Global F Test in regression analysis. A significant value tells you that one or more betas differ from zero, but it doesn't tell you which ones.

$$G_M = L^2 = DEV_0 - DEV_M$$

The significance level for the model chi-square indicates that this is a very large drop in chi-square, ergo we reject the null hypothesis. The effect of at least one of the IVs likely differs from zero.

You can think of the Deviance as telling you how bad the model still is, while the Model L^2 , aka G_M tells you how good it is.

Incremental Tests / Likelihood Ratio Chi-Square Tests. There is also an analog to the incremental F test. Just like with OLS, we can compare constrained and unconstrained models. We use an incremental chi-square square statistic instead of an incremental F statistic. (More commonly, you see phrases like chi-square contrasts.) The difference between the deviances of constrained and unconstrained models has a chi-square distribution with degrees of freedom equal to the number of constraints.

Incremental chi-square test/ chi-square contrast (analog to incremental F test)

$$L^2 = DEV_{\text{Constrained}} - DEV_{\text{Unconstrained}}, \text{d.f.} = \text{number of constraints}$$

If the resulting chi-square value is significant, stick with the unconstrained model; if insignificant then the constraints can be justified. Alternatively, you'll get the same results using

$$L^2 = \text{Model } L^2_{\text{Unconstrained}} - \text{Model } L^2_{\text{Constrained}}, \text{d.f.} = \text{number of constraints}$$

The notation L^2 is used to signify that this is a Likelihood Ratio Chi Square test (as opposed to, say, a Pearson Chi-Square test, which has less desirable properties). Again, notation is wildly inconsistent across authors. G^2 is another notation sometime used.

WARNING: In OLS, an incremental F test and a Wald test give you the same results. In logistic regression, a chi-square contrast between models and a Wald test generally do NOT give

identical results. LR chi-square contrasts are considered better but in large samples it may not matter much.

Nested Models-Stata. In Stata, we can get incremental and global LR chi-square tests easily by using the `nestreg` command. We should include the `lr` option so we get likelihood ratio tests rather than Wald tests. The `quietly` option suppresses a lot of the intermediate information, but don't use it if you want to see those results.

```
. nestreg, lr quietly: logit grade gpa tuce psi
```

```
Block 1: gpa
Block 2: tuce
Block 3: psi
```

Block	LL	LR	df	Pr > LR	AIC	BIC
1	-16.2089	8.77	1	0.0031	36.4178	39.34928
2	-15.99148	0.43	1	0.5096	37.98296	42.38017
3	-12.88963	6.20	1	0.0127	33.77927	39.64221

With Stata, you can also use the `lrtest` command to do likelihood ratio contrasts between models, e.g.

```
. quietly logit grade gpa
. est store m1
. quietly logit grade gpa tuce
. est store m2
. quietly logit grade gpa tuce psi
. est store m3
. lrtest m1 m2
```

```
Likelihood-ratio test          LR chi2(1) =      0.43
(Assumption: m1 nested in m2)  Prob > chi2 =    0.5096
```

```
. lrtest m2 m3
```

```
Likelihood-ratio test          LR chi2(1) =      6.20
(Assumption: m2 nested in m3)  Prob > chi2 =    0.0127
```

Stepwise Logistic Regression-Stata. As with other Stata commands, you can use the `sw` prefix for stepwise regression. We can add the `lr` option so that likelihood-ratio, rather than Wald, tests are used when deciding the variables to enter next.

```
. sw, lr pe(.05) : logit grade gpa tuce psi
```

```
LR test                begin with empty model
p = 0.0031 < 0.0500   adding gpa
p = 0.0130 < 0.0500   adding psi
```

```
Logistic regression                Number of obs   =          32
                                   LR chi2(2)        =          14.93
                                   Prob > chi2       =          0.0006
Log likelihood = -13.126573         Pseudo R2     =          0.3625
```

grade	Coef.	Std. Err.	z	P> z	[95% Conf. Interval]	
gpa	3.063368	1.22285	2.51	0.012	.6666251	5.46011
psi	2.337776	1.040784	2.25	0.025	.2978755	4.377676
_cons	-11.60157	4.212904	-2.75	0.006	-19.85871	-3.344425

Tests of Individual Parameters. Testing whether any individual parameter equals zero proceeds pretty much the same way as in OLS regression. You can, if you want, do an incremental LR chi-square test. That, in fact, is the best way to do it, since the Wald test referred to next is biased under certain situations. For individual coefficients, Stata reports z values, which is b/s_b .

```
. logit grade gpa tuce psi, nolog
```

```
Logistic regression                Number of obs   =          32
                                   LR chi2(3)        =          15.40
                                   Prob > chi2       =          0.0015
Log likelihood = -12.889633         Pseudo R2     =          0.3740
```

grade	Coef.	Std. Err.	z	P> z	[95% Conf. Interval]	
gpa	2.826113	1.262941	2.24	0.025	.3507938	5.301432
tuce	.0951577	.1415542	0.67	0.501	-.1822835	.3725988
psi	2.378688	1.064564	2.23	0.025	.29218	4.465195
_cons	-13.02135	4.931325	-2.64	0.008	-22.68657	-3.35613

With Stata, you can also continue to use the `test` command. The `test` command does Wald tests, which aren't as good as LR tests but which may be adequate in large samples, e.g.

```
. * Test whether effects of gpa and tuce are both 0
. test gpa tuce
```

```
( 1) gpa = 0
( 2) tuce = 0
```

```
chi2( 2) = 6.35
Prob > chi2 = 0.0418
```

```
. * Test whether effects of gpa and psi are equal
. test gpa = psi
```

```
( 1) gpa - psi = 0
```

```
chi2( 1) = 0.11
Prob > chi2 = 0.7437
```

R² Analogs. As Menard points out in Applied Logistic Regression Analysis, Second Edition, several people have tried to come up with the equivalent of an R² measure for logistic regression. No one of these measures seems to have achieved widespread acceptance yet. One of the simplest and most popular formulas is

$$\text{Pseudo } R^2 = \text{Model } L^2 / \text{DEV}_0 = 1 - \text{DEV}_M / \text{DEV}_0 = 1 - \text{LL}_M / \text{LL}_0$$

where, as you'll recall, DEV₀ (or -2LL₀) pertains to the baseline model with intercept only. (Menard refers to this as R²_L; it is also called McFadden R²; Stata just calls it Pseudo R². Be careful when reading, since the term Pseudo R² gets applied to a lot of different statistics.) This statistic will equal zero if all coefficients are zero. It will come close to 1 if the model is very good. In the present case, for the model with gpa, psi and tuce included,

$$\text{Pseudo } R^2 = \text{Model } L^2 / \text{DEV}_0 = 15.404 / 41.183 = .374$$

Menard (p. 27) argues for the Pseudo R² statistic on the grounds that it is conceptually closest to OLS R² i.e. it reflects a proportionate reduction in the quantity actually being minimized, -2LL. However, as I explain in my categorical data class, you can make a logical case for most of the Pseudo R² measures.

Other ways of assessing "Goodness of Fit." There are other ways to assess whether or not the model fits the data. For example, there is the *classification table*. The command in Stata is `estat class` (you can also just use `lstat`)

```
. quietly logit grade gpa tuce psi
. estat class
```

Logistic model for grade

Classified	True		Total
	D	~D	
+	8	3	11
-	3	18	21
Total	11	21	32

Classified + if predicted Pr(D) >= .5
True D defined as grade != 0

Sensitivity	Pr(+ D)	72.73%
Specificity	Pr(- ~D)	85.71%
Positive predictive value	Pr(D +)	72.73%
Negative predictive value	Pr(~D -)	85.71%
False + rate for true ~D	Pr(+ ~D)	14.29%
False - rate for true D	Pr(- D)	27.27%
False + rate for classified +	Pr(~D +)	27.27%
False - rate for classified -	Pr(D -)	14.29%
Correctly classified		81.25%

In the classification table, cases with probabilities $\geq .50$ are predicted as having the event, other cases are predicted as not having the event. Ideally, you would like to see the two groups have very different estimated probabilities. In this case, of the 21 people who did not get A's, the model correctly predicted 18 would not but said that 3 would. Similarly, of the 11 who got A's, the model was right on 8 of them.

From the classification table, you can't tell how great the errors are. The 6 misclassified cases may have been within one or two percentage points of being classified correctly, or they may have been way off. For "rare" events, I am not sure how useful the table is. A 10% probability may be relatively high, but still not high enough to get the case classified as a 1 (e.g. there may be only 1 chance in a 1000 of the average 20 year old dying within the year; identifying those for whom the odds are 1 in 10 of dying may be quite useful.) Menard goes on at some length about other possible classification/prediction strategies.

Diagnostics. It can also be useful to run various diagnostics. These help to indicate areas or cases for which the model is not working well. Menard lists several statistics for looking at residuals. Menard also briefly discusses some graphical techniques that can be useful. Also see Hamilton's Statistics with Stata for some ideas.

In Stata, you can again use the `predict` command to compute various outliers. As was the case with OLS, Stata tends to use different names than SPSS and does some computations differently. Cases 2 and 27 seem to be the most problematic.

```
. * Generate standardized residuals
. predict p
(option pr assumed; Pr(grade))
. predict rstandard, rstandard
. extremes rstandard p grade gpa tuce psi
```

obs:	rstandard	p	grade	gpa	tuce	psi
27.	-2.541286	.8520909	0	3.51	26	1
18.	-1.270176	.5898724	0	3.12	23	1
16.	-1.128117	.5291171	0	3.1	21	1
28.	-.817158	.3609899	0	3.53	26	0
24.	-.7397601	.3222395	0	3.57	23	0
19.	.8948758	.6354207	1	3.39	17	1
30.	1.060433	.569893	1	4	21	0
15.	1.222325	.481133	1	2.83	27	1
23.	2.154218	.1932112	1	3.26	25	0
2.	3.033444	.1110308	1	2.39	19	1

Summary: Comparisons with OLS. There are many similarities between OLS and Logistic Regression, and some important differences. I'll try to highlight the most crucial points here.

OLS and its extensions	Logistic Regression
Estimated via least squares	Estimated via Maximum Likelihood.
Y is continuous, can take on any value	Y can only take on 2 values, typically 0 and 1
X's are continuous vars. Categorical variables are divided up into dummy variables	Same as OLS
X's are linearly related to Y; in the case of the LPM, X's are linearly related to P(Y=1)	X's are linearly related to log odds of event occurring. Log odds, in turn, are nonlinearly related to P(Y = 1).
Y's are statistically independent of each other, e.g., don't have serial correlation, don't include husbands and their wives as separate cases	Same as OLS
Robust standard errors can be used when error terms are not independent and identically distributed.	Same as OLS. Stata makes this easy (just add a <code>robust</code> parameter), SPSS does not.
There can be no perfect multicollinearity among the X's. High levels of multicollinearity can result in unstable sample estimates and large standard errors	Same as OLS. Techniques for detecting multicollinearity are also similar. In fact, as Menard points out, you could just run the corresponding OLS regression, and then look at the correlations of the IVs, the tolerances, variance inflation factors, etc. Or, use Stata's <code>collin</code> command.
Missing data can be dealt with via listwise deletion, pairwise deletion, mean substitution, multiple imputation	Pairwise deletion isn't an option. Can't do "mean substitution" on the DV. Otherwise, can use techniques similar to those that we've described for OLS.
Global F test is used to test whether any IV effects differ from 0. d.f. = K, N-K-1	Model chi-square statistic (also known as Model L^2 or G^2 or G_M) is used for same purpose. D.F. = number of IVs in the model = K.
Incremental F test is used to test hypotheses concerning whether subset of coefficients = 0. If you specify variables in blocks, the F change statistic will give you the info you need.	LR Chi-square statistic is used. $DEV_{Constrained} - DEV_{Unconstrained}$ $Model L^2_{Unconstrained} - Model L^2_{Constrained}$
T test or incremental F test is used to test whether an individual coefficient = 0	Can use a LR chi square test (preferable) or Wald statistic (probably usually ok, but not always).

Incremental F tests or T tests can be used to test equalities of coefficients within a model, equalities across populations, interaction effects.	Same basic procedures, substituting LR chi square tests for F tests.
Wald tests (as produced by the <code>test</code> command in stata) will produce the same results as incremental F tests. A nice thing about Wald tests is that they only require the estimation of the unconstrained model.	Wald tests can be performed, but they will generally NOT produce exactly the same results as LR tests. LR tests (which require the estimation of constrained and unconstrained models) are preferable, although in practice results will often be similar.
Can have interaction effects. Centering can sometimes make main effects easier to interpret. If you center the continuous vars, then the main effect of an IV like race is equal to the difference in the predicted values for an “average” black and white.	NOT quite the same as OLS. You can use interaction terms, but there are potential problems you should be aware of when interpreting results. See Allison (1999) or Williams (2009, 2010) for discussions. If you center, then the main effect of an IV like race is equal to the difference in the log odds for an “average” black and white.
Can do transformations of the IVs and DV to deal with nonlinear relationships, e.g. X^2 , $\ln(X)$, $\ln(Y)$.	Same as OLS for the IVs, but you of course can't do transformations of the dichotomous DV.
Can plot Y against X, examine residuals, plot X against residuals, to identify possible problems with the model	Similar to OLS. Can examine residuals.
Can do mindless, atheoretical stepwise regression	Similar to OLS
R^2 tells how much of total variance is “explained”.	Numerous Pseudo R^2 stats have been proposed. If you use one, make clear which one it is.
Can look at standardized betas.	There is actually a reasonable case for using standardized coefficients in logistic regression. Long & Freese's <code>spostado</code> routines include the <code>listcoef</code> command, which can do various types of standardization.
Can do path analysis. Can decompose association. Can estimate recursive and nonrecursive models. Programs like LISREL can deal with measurement error.	Most ideas of the “logic of causal order” still apply. But, many things, such as decomposition of effects, controlling for measurement error, estimating nonrecursive models, are much, much harder to do. There is work going on in this area, e.g. Lisrel, M-Plus, <code>gllamm</code> (an add-on routine to Stata).

Related Topics. Here is a super-quick look at other techniques for analyzing categorical data.

Probit. Probit models are an alternative to Logit models. They tend to produce almost identical results, and logit models are usually easier to work with. For some types of problems, there are more advanced probit techniques that can be useful.

Multinomial Logit. You can also have a dependent variable with more than two categories, e.g. the dependent variable might take the values Republican, Democrat, Other. The idea is that you talk about the probability of being in one group as opposed to another. In SPSS, use NOMREG, in Stata use `mlogit`.

Ordered Logit. Sometimes DVs are ordinal. Sometimes, it is ok to just treat them as interval-level and use OLS regression. But, other times an Ordered Logit routine is preferable. SPSS has PLUM. Stata has the built-in `ologit` and `oprobit`. Stata also has various user-written routines, including Williams's `oglm` and `gologit2`.

Appendix (Optional): Computing the log likelihood. This is adapted from J. Scott Long's Regression Models for Categorical and Limited Dependent Variables.

Define p_i as the probability of observing whatever value of y was actually observed for a given observation, i.e.

$$p_i = \begin{cases} \Pr(y_i = 1 | x_i) & \text{if } y_i = 1 \text{ is observed} \\ 1 - \Pr(y_i = 1 | x_i) & \text{if } y_i = 0 \text{ is observed} \end{cases}$$

If the observations are independent, the likelihood equation is

$$L(\boldsymbol{\beta} | \mathbf{y}, \mathbf{X}) = \prod_{i=1}^N p_i$$

The likelihood tends to be an incredibly small number, and it is generally easier to work with the log likelihood. Ergo, taking logs, we obtain the log likelihood equation:

$$\ln L(\boldsymbol{\beta} | \mathbf{y}, \mathbf{X}) = \sum_{i=1}^N \ln p_i$$

Before proceeding, let's see how this works in practice! Here is how you compute p_i and the log of p_i using Stata:

```
. quietly logit grade gpa tuce psi
. * Compute probability that y = 1
. predict pi
(option p assumed; Pr(grade))
. * If y = 0, replace pi with probability y = 0
. replace pi = 1 - pi if grade == 0
(21 real changes made)
. * compute log of pi
. gen lnpi = ln(pi)

. list grade pi lnpi, sep(8)
```

```
+-----+
| grade      pi      lnpi |
+-----+-----+-----+
| 1.         0  .9386242  -.0633401 |
| 2.         1  .1110308  -2.197947  |
| 3.         0  .9755296  -.0247748 |
|          --- Output deleted --- |
| 30.        1  .569893   -.5623066 |
| 31.        1  .9453403  -.0562103 |
| 32.        1  .6935114  -.3659876 |
+-----+-----+-----+
```

So, this tells us that the predicted probability of the first case being 0 was .9386. The probability of the second case being a 1 was .111. The probability of the 3rd case being a 0 was .9755; and so on. The likelihood is therefore

$$L(\boldsymbol{\beta} | \mathbf{y}, \mathbf{X}) = \prod_{i=1}^N p_i = .9386 * .1110 * .9755 * \dots * .6935 = .000002524$$

which is a really small number; indeed so small that your computer or calculator may have trouble calculating it correctly (and this is only 32 cases; imagine the difficulty if you have hundreds of thousands). Much easier to calculate is the log likelihood, which is

$$\ln L(\boldsymbol{\beta} | \mathbf{y}, \mathbf{X}) = \sum_{i=1}^N \ln p_i = -.0633 + -2.198 + \dots + -.366 = -12.88963$$

Stata's `total` command makes this calculation easy for us:

```
. total lnpi

Total estimation                Number of obs   =           32

-----+-----
            |          Total   Std. Err.   [95% Conf. Interval]
-----+-----
lnpi       |   -12.88963     3.127734    -19.26869    -6.510578
-----+-----
```

Note: The maximum likelihood estimates are those values of the parameters that make the observed data most likely. That is, the maximum likelihood estimates will be those values which produce the largest value for the likelihood equation (i.e. get it as close to 1 as possible; which is equivalent to getting the log likelihood equation as close to 0 as possible).