

# Tutorial 9 HT

## Research Methods for Political Science - PO3110

---

Andrea Salvi

2 April 2019

Trinity College Dublin,

<https://andrsalvi.github.io/research-methods/>

1. A few more things on Logistic Regression

## A few more things on Logistic Regression

---

- Reported in the "Variables in the equation box".

# Wald Statistic

- Reported in the "Variables in the equation box".
- SPSS Reports  $z^2$ , that is Wald Squared.

# Wald Statistic

- Reported in the "Variables in the equation box".
- SPSS Reports  $z^2$ , that is Wald Squared.
- As we did for the  $t$ , we calculate the new value as  $z^2 = (\frac{\beta}{SE})^2$
- SPSS conveniently compares that to a the relevant critical value of the  $\chi^2$  distribution so to obtain a p-value.

# Wald Statistic

- Reported in the "Variables in the equation box".
- SPSS Reports  $z^2$ , that is Wald Squared.
- As we did for the  $t$ , we calculate the new value as  $z^2 = (\frac{\beta}{SE})^2$
- SPSS conveniently compares that to a the relevant critical value of the  $\chi^2$  distribution so to obtain a p-value.
- As we've seen for linear regression,  $H_0 : \beta = 0$

- Reported in the "Variables in the equation box".
- SPSS Reports  $z^2$ , that is Wald Squared.
- As we did for the  $t$ , we calculate the new value as  $z^2 = (\frac{\beta}{SE})^2$
- SPSS conveniently compares that to a the relevant critical value of the  $\chi^2$  distribution so to obtain a p-value.
- As we've seen for linear regression,  $H_0 : \beta = 0$
- **Interpretation:** If  $p \leq 0.05$  we can conclude that there is a statistically significant relationship between the corresponding IV and the probability of DV = 1 (event happening).



# Wald Statistic

- Reported in the "Variables in the equation box".
- SPSS Reports  $z^2$ , that is Wald Squared.
- As we did for the  $t$ , we calculate the new value as  $z^2 = (\frac{\beta}{SE})^2$
- SPSS conveniently compares that to a the relevant critical value of the  $\chi^2$  distribution so to obtain a p-value.
- As we've seen for linear regression,  $H_0 : \beta = 0$
- **Interpretation:** If  $p \leq 0.05$  we can conclude that there is a statistically significant relationship between the corresponding IV and the probability of DV = 1 (event happening).
- Might be worth double-checking with bootstrapping when you get very large  $\beta$ s. That leads to inflate SE, which in turns might lead to a misleading Walt statistic (smaller than it should be).

- Reported in the "Variables in the equation box".
- SPSS Reports  $z^2$ , that is Wald Squared.
- As we did for the  $t$ , we calculate the new value as  $z^2 = \left(\frac{\beta}{SE}\right)^2$
- SPSS conveniently compares that to a the relevant critical value of the  $\chi^2$  distribution so to obtain a p-value.
- As we've seen for linear regression,  $H_0 : \beta = 0$
- **Interpretation:** If  $p \leq 0.05$  we can conclude that there is a statistically significant relationship between the corresponding IV and the probability of DV = 1 (event happening).
- Might be worth double-checking with bootstrapping when you get very large  $\beta$ s. That leads to inflate SE, which in turns might lead to a misleading Walt statistic (smaller than it should be).
- This leads to an ARTIFICIAL increase in the probability of rejecting the predictor as being statistically significantly different than zero. (Type II error – false negative).

# Pseudo R Squared

- Remember we used  $R^2$  as a measure of the fit of the model and F as a test of overall fit.

# Pseudo R Squared

- Remember we used  $R^2$  as a measure of the fit of the model and F as a test of overall fit.
- We cannot use the same measures/tests here, because the coefficients and standard errors in the model are calculated here using maximum-likelihood estimation.

# Pseudo R Squared

- Remember we used  $R^2$  as a measure of the fit of the model and F as a test of overall fit.
- We cannot use the same measures/tests here, because the coefficients and standard errors in the model are calculated here using maximum-likelihood estimation.
- **CANNOT** be interpreted in terms of 'percentage variance explained'

# Pseudo R Squared

- Remember we used  $R^2$  as a measure of the fit of the model and F as a test of overall fit.
- We cannot use the same measures/tests here, because the coefficients and standard errors in the model are calculated here using maximum-likelihood estimation.
- **CANNOT** be interpreted in terms of 'percentage variance explained
- It is calculated by dividing the model chi-square (based on the log-likelihood) by the baseline -2LL (the log-likelihood of the model before any predictors were entered).

# Pseudo R Squared

- Remember we used  $R^2$  as a measure of the fit of the model and F as a test of overall fit.
- We cannot use the same measures/tests here, because the coefficients and standard errors in the model are calculated here using maximum-likelihood estimation.
- **CANNOT** be interpreted in terms of 'percentage variance explained
- It is calculated by dividing the model chi-square (based on the log-likelihood) by the baseline  $-2LL$  (the log-likelihood of the model before any predictors were entered).
- Nagelkerke R Squared is the proportional reduction in the absolute value of the log-likelihood measure and as such it is a measure of how much the "badness" of fit improves as a result of the inclusion of the predictor variables.

# Pseudo R Squared

- Remember we used  $R^2$  as a measure of the fit of the model and F as a test of overall fit.
- We cannot use the same measures/tests here, because the coefficients and standard errors in the model are calculated here using maximum-likelihood estimation.
- **CANNOT** be interpreted in terms of 'percentage variance explained
- It is calculated by dividing the model chi-square (based on the log-likelihood) by the baseline  $-2LL$  (the log-likelihood of the model before any predictors were entered).
- Nagelkerke R Squared is the proportional reduction in the absolute value of the log-likelihood measure and as such it is a measure of how much the "badness" of fit improves as a result of the inclusion of the predictor variables.



# Pseudo R Squared

- Cox Snell as well as Nagelkerke R square can be interpreted on a scale from 0 (poor fit, indicating that the predictors are useless at predicting the outcome variable) to 1 (best fit, indicating that the model predicts the outcome variable perfectly).

# Pseudo R Squared

- Cox Snell as well as Nagelkerke R square can be interpreted on a scale from 0 (poor fit, indicating that the predictors are useless at predicting the outcome variable) to 1 (best fit, indicating that the model predicts the outcome variable perfectly).
- Cox Snell never reaches its theoretical maximum of 1 (equation in field if you are curious).The Nagelkerke measure performs an ad hoc adjustment to CS so that it can reach 1.
- Independently, in the logistic regression context, these measures tell us very little.
  1. A pseudo R-squared only has meaning when compared to another pseudo R-squared of the same type, on the same data, predicting the same outcome.
  2. In this situation, the higher pseudo R-squared indicates which model better predicts the outcome.

- A maximum likelihood estimator is an estimator that makes the observed the data most likely to have occurred.

# Model Fit

- A maximum likelihood estimator is an estimator that makes the observed data most likely to have occurred. We maximize the likelihood function:  $P(data|model)$  Which is the probability getting the observed data given the model.
- Model Summary box

# Model Fit

- A maximum likelihood estimator is an estimator that makes the observed data most likely to have occurred. We maximize the likelihood function:  $P(data|model)$  Which is the probability getting the observed data given the model.
- Model Summary box
- SPSS does not report the Maximum Likelihood but the **-2LogLikelihood (deviance statistic)**: Scarcely informative in absolute terms. Very useful for comparisons between use it to compare different models.

# Model Fit

- A maximum likelihood estimator is an estimator that makes the observed data most likely to have occurred. We maximize the likelihood function:  $P(\text{data}|\text{model})$  Which is the probability getting the observed data given the model.
- Model Summary box
- SPSS does not report the Maximum Likelihood but the **-2LogLikelihood (deviance statistic)**: Scarcely informative in absolute terms. Very useful for comparisons between use it to compare different models.
- Larger values of deviance indicate a poorer fit. You can either look at them across the different models or:

# Model Fit

- A maximum likelihood estimator is an estimator that makes the observed the data most likely to have occurred. We maximize the likelihood function:  $P(data|model)$  Which is the probability getting the observed data given the model.
- Model Summary box
- SPSS does not report the Maximum Likelihood but the **-2LogLikelihood (deviance statistic)**: Scarcely informative in absolute terms. Very useful for comparisons between use it to compare different models.
- Larger values of deviance indicate a poorer fit. You can either look at them across the different models or:
- SPSS always compares that with the baseline model (Omnibus test for model coefficients box). How does that work?

- Omnibus Test box test the  $H_0$  = no difference in terms of fit between the baseline model and our model.



- Omnibus Test box test the  $H_0$  = no difference in terms of fit between the baseline model and our model.
- The likelihood ratio statistic is simply the deviance of the baseline model MINUS the deviance of the new model.

- Omnibus Test box test the  $H_0$  = no difference in terms of fit between the baseline model and our model.
- The likelihood ratio statistic is simply the deviance of the baseline model MINUS the deviance of the new model.
- Since this ratio converges asymptotically to a chi squared distribution, SPSS calculates the p-values from there.

- Omnibus Test box test the  $H_0$  = no difference in terms of fit between the baseline model and our model.
- The likelihood ratio statistic is simply the deviance of the baseline model MINUS the deviance of the new model.
- Since this ratio converges asymptotically to a chi squared distribution, SPSS calculates the p-values from there.
- Looking at the p-values we can now determine whether our new model is a significant improvement upon the baseline model.