

#### Additional Resources:

Section 8.4 in Diez, Barr, and Cetinkaya-Rundel, (2015), *OpenIntro Statistics* https://www.openintro.org/download.php?file=os3&redirect=/stat/textbook/os3.php

# **1.** Conducting Inference on $\beta_0, \beta_1, ..., \beta_p$ on logistic Regression model

Like with our linear regression models, we can formulate a logistic regression model for either the <u>sample</u> data or the <u>population</u> data.



We can also think of  $\hat{\beta}_i$  as a sample statistic that can be used as a point estimate for estimating  $\beta_i$ .

Like with multiple linear regression slopes and intercepts, we can do the following to conduct inference on  $\beta_0, \beta_1, ..., \beta_p$ .

- 1. Create a confidence interval for  $\beta_i$ .
- 2. Conduct hypothesis testing for  $\beta_i$ .
  - a. Using a p-value
  - b. Using a confidence interval.

## 1.1. CHECKING CONDITIONS FOR INFERENCE (AND MODEL FIT) ON $\beta_0, \beta_1, ..., \beta_p$ IN A LOGISTIC REGRESSION MODEL

#### First check logistic regression conditions for model fit and inference.

- 1. Independence of Observations Condition:
  - a. The sample is random.
  - b. n < 10% of the population size.
- 2. <u>Linearity Condition</u>: the explanatory variables should have a linear relationship with the **log-odds** of the response variable success level.

#### How to check in Python

For each numerical explanatory variable, plot this explanatory variable on the x-axis and the (0/1 response variable on the y-axis in a scatterplot. Then fit and plot a simple logistic regression curve for each of these scatterplots. If each of your curves are "S-shaped" for each of these plots, then we can say that this condition is met.

3. <u>No Multicollinearity Condition</u>: the explanatory variables should not have a strong linear relationship with each other.

# **1.2.** CREATING A $(1 - \alpha) \cdot 100\%$ Confidence Interval for $\beta_0, \beta_1, ..., or, \beta_p$ (individually) in a logistic Regression model

#### 1. First check logistic regression conditions for inference.

- a. Independence of Observations Condition
  - i. The sample is random.
  - ii. n < 10% of the population size.
- b. <u>Linearity Condition</u>: the explanatory variables should have a linear relationship with the **log-odds** of the response variable success level.
- c. <u>No Multicollinearity Condition</u>: the explanatory variables should not have a strong linear relationship with each other.

### **2.** The $(1 - \alpha) \cdot 100\%$ confidence interval for $\beta_i$ is calculated by:

 $(point \ estimate) \pm (critical \ value)(standard \ error)$ 

$$\widehat{\beta}_{\iota} \pm z^* SE_{\beta_i}$$

Notation:

3. The  $(1 - \alpha) \cdot 100\%$  confidence for  $e^{\beta_i}$  (the odds multiplier for  $x_i$ ) is calculated by:

$$(e^{\widehat{\beta}_l - z^*SE_{\beta_l}}, e^{\widehat{\beta}_l + z^*SE_{\beta_l}})$$

Remember from Unit 15, that  $e^{\beta_i}$  is easier to interpret than  $\beta_i$ . We call  $e^{\beta_i}$ , the odds multiplier for the explanatory variable  $x_i$  because if we were to increase this variable value by 1, we would expect the odds to increase by a multiple/factor of  $e^{\beta_i}$ .

After we have calculated a  $(1 - \alpha) \cdot 100\%$  confidence interval for  $\beta_i$ , we can just exponentiate the lower and upper bounds to get a confidence interval for  $e^{\beta_i}$ .

**Ex**: In the **Jupyter notebook**, create a 90% confidence interval for the slope that corresponds to the age variable in the population logistic regression model.

**1.3.** Conducting Hypothesis Testing on  $\beta_0$ ,  $\beta_1$ , ..., or,  $\beta_p$  (individually) in a logistic Regression model

#### 1. <u>Set up the hypotheses</u>

 $H_0: \beta_i = 0$  $H_A: \beta_i \neq 0$ 

Equivalent to ...

 $\begin{array}{l} H_0: e^{\beta_i} = 1 \\ H_A: e^{\beta_i} \neq 1 \end{array}$ 

#### 2. <u>Check the conditions for conducting inference on a population slope/intercept.</u>

- a. Independence of Observations Condition
  - i. The sample is random.
  - ii. n < 10% of the population size.
- b. <u>Linearity Condition</u>: the explanatory variables should have a linear relationship with the **log**-**odds** of the response variable success level.
- c. <u>No Multicollinearity Condition</u>: the explanatory variables should not have a strong linear relationship with each other.

 $\widehat{\beta}_{l}$ 

#### 3. Calculate the point estimate (observed sample statistic)

4. <u>Calculate the p-value</u>

$$p - value = 2P(Z \ge |\frac{\widehat{\beta_i} - 0}{SE_{\widehat{\beta_i}}}|)$$

### 5. Make a Decision

#### With a p-value

- a. If  $p value < \alpha$ , then we "reject the null hypothesis." And we say that "there IS sufficient evidence to suggest the alternative hypothesis."
- b. If  $p value \ge \alpha$ , then we "fail to reject the null hypothesis." And we say that "there IS NOT sufficient evidence to suggest the alternative hypothesis."

#### With a confidence interval for $\beta_i$

- a. If the **null value**  $\beta_i = 0$  is **NOT IN the confidence interval**, then we "reject the null hypothesis." And we say that "there IS sufficient evidence to suggest the alternative hypothesis."
- b. If **null value**  $\beta_i = 0$  is IN the confidence interval, then we "fail to reject the null hypothesis." And we say that "there IS NOT sufficient evidence to suggest the alternative hypothesis."

#### With a confidence interval for $e^{\beta_i}$

- a. If the **null value**  $e^{\beta_i} = 1$  is **NOT IN the confidence interval**, then we "reject the null hypothesis." And we say that "there IS sufficient evidence to suggest the alternative hypothesis."
- b. If **null value**  $e^{\beta_i} = 1$  is IN the confidence interval, then we "fail to reject the null hypothesis." And we say that "there IS NOT sufficient evidence to suggest the alternative hypothesis."

**Ex**: In the **Jupyter notebook**, create test whether there is sufficient evidence to suggest that the population slope for age is non-zero in the logistic regression model (which predicts the probability for approval for the president's foreign policy given age and sex.)

## **2. MODELING INTERACTION EFFECTS IN A LOGISTIC REGRESSION MODEL**

Like in linear regression models, we can model the interaction effect between two explanatory variables  $x_i$  and  $x_j$  by the product of these two variables  $x_i \cdot x_j$  in our logistic regression model and fitting a slope for this product as well as our other explanatory variables.

**Ex**: In the **Jupyter notebook**, create a logistic regression model that predicts the probability that a person in the sample approves of the president's foreign policy given:

- o sex,
- o age, and
- $\circ~$  the interaction between sex and age.

Then, write out the resulting logistic regression model for the sample.

Finally, use this model to predict the probability that 19 year old male supported the president's foreign policy (in 2017).