

# Unit 11: Logistic Regression – Part 1

## Case Studies:

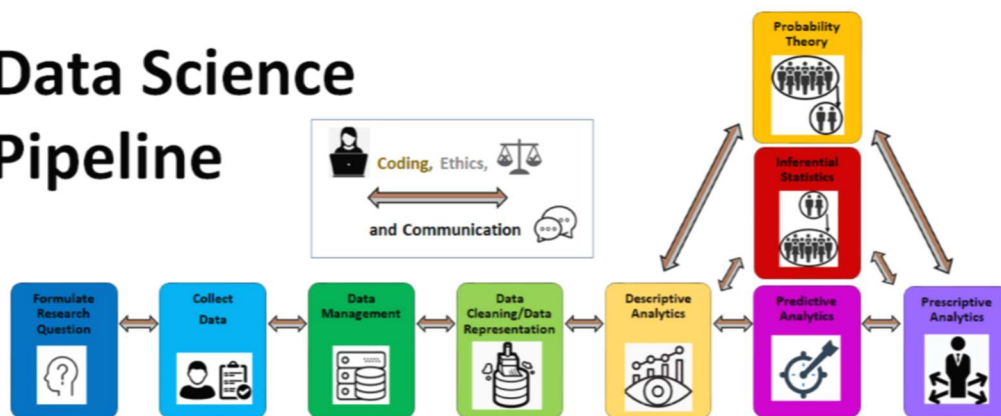
- To introduce the concept of simple logistic regression, we will examine the association between sex and approval for the direction the country is going in (in 2017)?
  - Response: Approve/disapprove of the direction
  - Explanatory: sex
- To introduce the concept of multiple logistic regression, we will examine the association between sex and age and approval for the direction the country is going in (in 2017)?
  - Response: Approve/disapprove of the direction
  - Explanatory(s):
    - Sex
    - Age



disapprove

approve

## Data Science Pipeline



## Summary of Concepts:

- Analyses for Associations
- Association Analyses Summary:** (Numerical(s) and Categorical(s) Explanatory Variable(s)-> Categorical Response Variable (with 2 levels)
- Why should we not use linear regression to model a categorical response variable?
- What curve can we fit instead when we have a categorical response variable with 2 levels?
- Odds vs. Probability
- Fitting a Simple Logistic Regression Model
- Fitting a Multiple Logistic Regression Model
- Making predictions with a logistic regression equation.
- Interpreting  $\hat{\beta}_0, \hat{\beta}_1, \hat{\beta}_2, \dots, \hat{\beta}_p$  in a Logistic Regression Model

## Additional Resources:

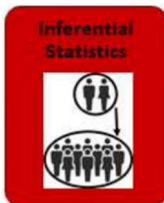
Section 8.4 in Diez, Barr, and Cetinkaya-Rundel, (2015), *OpenIntro Statistics* <https://www.openintro.org/download.php?file=os3&redirect=/stat/textbook/os3.php>

# 1. ANALYSES FOR ASSOCIATIONS

## Questions to consider, when selecting an analysis to test an association.



1. Which variable is the **response variable** in this association?
  - a. Is it a **categorical** or **numerical** variable?
  - b. If it's a categorical variable, **how many levels** does it have?
2. Which variable(s) is the **explanatory variable** in this association?
  - a. Is it a **categorical** or **numerical** variable?
  - b. If it's a categorical variable, **how many levels** does it have?
3. How would you **quantify this association**?
  - a. Difference between two summary statistics? What two summary statistics?
  - b. With a model? What kind of model?



4. Are you interested in an association in a **sample** or a **population**?
5. When is it **appropriate to use this test** for association?

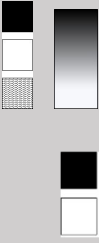



6. Can you use this model/test to **make predictions**?
  - a. How would you quantify the performance of your predictions?

## 2. ASSOCIATION ANALYSIS SUMMARY:

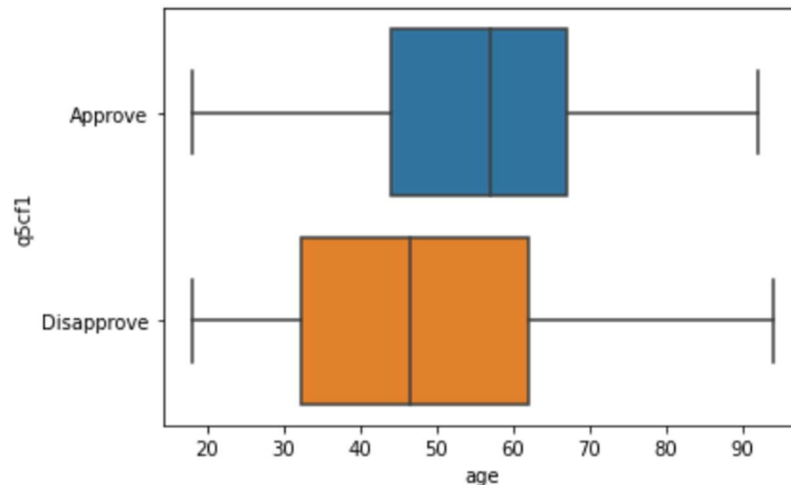
RESPONSE: CATEGORICAL (TWO LEVELS)

EXPLANATORY: NUMERICAL(S), CATEGORICAL(S) (ANY NUMBER OF LEVELS)

Research Questions about Associations	Type of Variables Involved in the Association Test	 <p><b><u>Explanatory Variables:</u></b> Numerical Variables and/or Categorical Variables</p> <p><b><u>Response Variable:</u></b> Categorical Variable (two levels)</p>
	Example	<p>Is there an association between <b>sex</b> and <b>age</b> and <b>opinion</b> on the direction that the country is going in (satisfied/dissatisfied)?</p> 
	Type of Association (Way to Quantify Association)	<p><b>Logistic Regression Model</b> (<i>linear relationship between explanatory variables and response variable (y)</i>)</p>
Descriptive Analytics	How to <u>Describe</u> an Association in a <u>Sample</u> ?	<p><u>Model:</u></p> <ul style="list-style-type: none"> <li><math>\log\left(\frac{\hat{p}}{1-\hat{p}}\right) = \hat{\beta}_0 + \hat{\beta}_1 x_1 + \hat{\beta}_2 x_2 + \dots \hat{\beta}_p x_p</math></li> </ul>
	When is this analysis ( <i>for the sample</i> ) appropriate to use?	See upcoming unit.
Inferential Statistics	How to <u>Infer</u> an Association for a <u>Population</u> ?	See upcoming unit.
	When is this analysis ( <i>for the population</i> ) appropriate to use?	See upcoming unit.
Predictive Analytics	Making Predictions	See upcoming unit.
	How to quantify the performance of your prediction(s)?	See upcoming unit.

### 3. WHY SHOULD WE NOT USE LINEAR REGRESSION TO MODEL A CATEGORICAL RESPONSE VARIABLE?

Suppose we wanted to fit a simple linear regression model, this time using **approval (or disapproval) of the president's foreign policy** as the response variable and **age** as the explanatory variable.

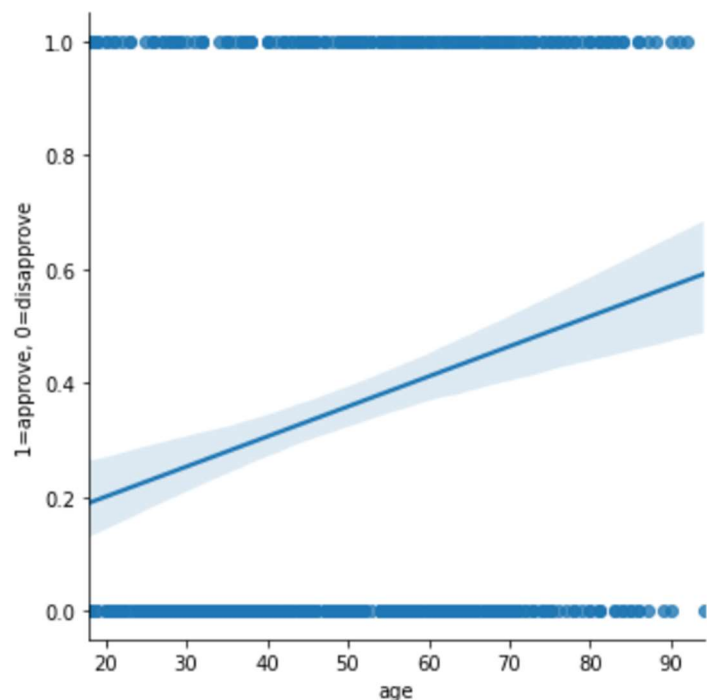


What we could *try* to do is create a new numerical variable **y** in which:

- $y = \underline{\hspace{2cm}}$ , when a survey respondent approves
- $y = \underline{\hspace{2cm}}$ , when a survey respondent disapproves.

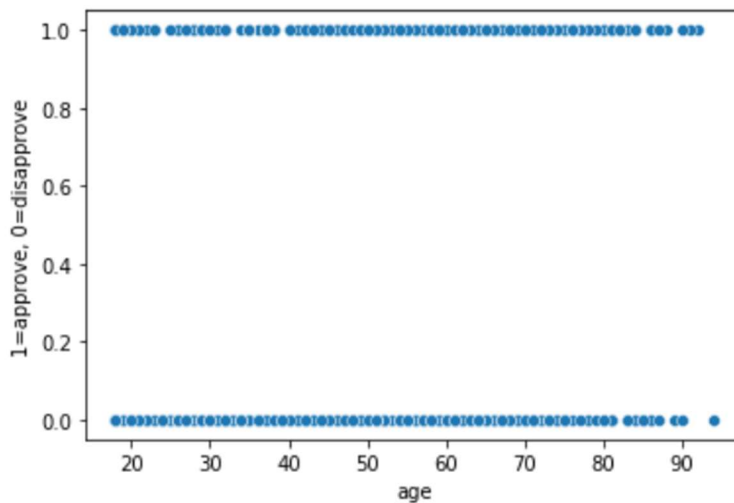
Would we expect a best fit line to provide a good fit of the data? What are some issues we might have when using this best fit line when making predictions?

	age	sex	q5cf1	y
1	70.0	Female	Disapprove	0
2	69.0	Female	Disapprove	0
4	70.0	Female	Disapprove	0
6	89.0	Female	Disapprove	0
7	92.0	Female	Approve	1



#### 4. WHAT CURVE CAN WE FIT INSTEAD WHEN WE HAVE A CATEGORICAL RESPONSE VARIABLE (WITH 2 LEVELS)?

What kind of curve would be a better fit of the data?



The curve that we just drew is called a \_\_\_\_\_. This function has the following properties:

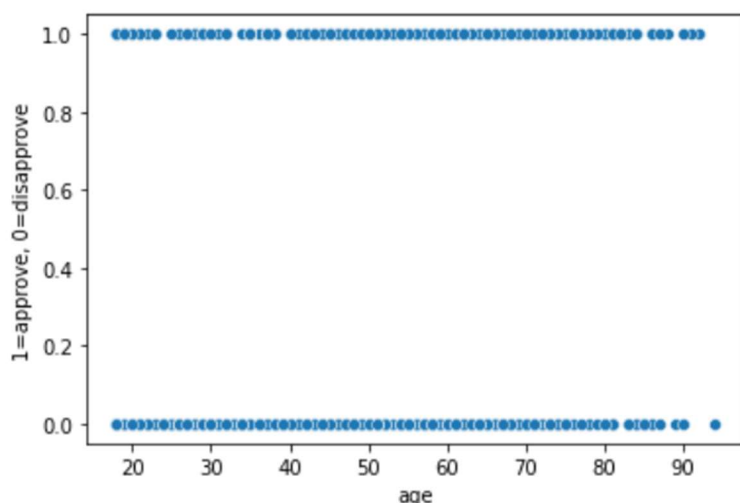
- It's S-shaped and has two horizontal asymptotes at  $y = \underline{\hspace{2cm}}$  and  $y = \underline{\hspace{2cm}}$ .

- Curves of *this type* are defined by the formula  $S(x) = \frac{1}{1+e^{-x}} = \frac{e^x}{e^x+1}$

By multiplying a  $\widehat{\beta}_1$  to x and then adding a  $\widehat{\beta}_0$ , we can horizontally stretch out the curve and horizontally shift the curve in an attempt to better fit our data.

$$s(\widehat{\beta}_0 + \widehat{\beta}_1 x) = \frac{e^{\widehat{\beta}_0 + \widehat{\beta}_1 x}}{e^{\widehat{\beta}_0 + \widehat{\beta}_1 x} + 1}$$

Even with this better fitting logistic curve, in most datasets we will never have all data points (which have either a y=1 or a y=0) fit perfectly on this curve.



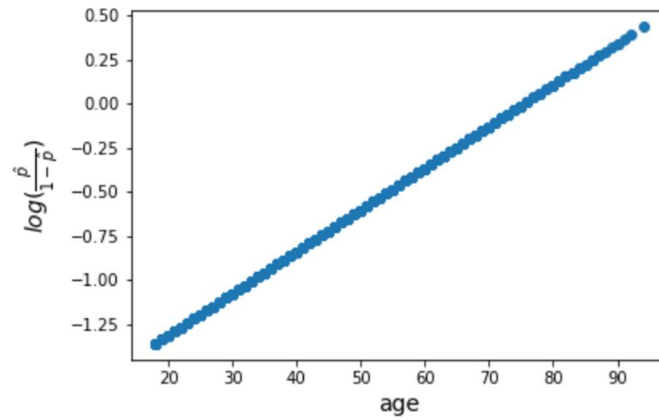
So how should we interpret predictions made with this curve then if our predictions are in between (0,1)?

$$P(\widehat{Y} = 1) = \widehat{p} = \frac{e^{\widehat{\beta}_0 + \widehat{\beta}_1 x}}{e^{\widehat{\beta}_0 + \widehat{\beta}_1 x} + 1}$$

In this case, we call  $\widehat{p}$  the **predictive probability**, which is the predicted probability that the corresponding response variable will be the \_\_\_\_\_ (ie. y=1), for the given explanatory variable value x.

With some algebraic manipulation, we can convert this equation into what we call our **logistic regression model** for the sample data.

$$\log\left(\frac{\hat{p}}{1-\hat{p}}\right) = \hat{\beta}_0 + \hat{\beta}_1 x$$



## 5. ODDS VS. PROBABILITY

In general, if  $p$  is the probability of some “success” level happening, then we call:

$$\frac{p}{1-p} = \text{the odds of a success.}$$

In general, if  $p$  is the probability of some “success” level happening, then we call:

$$\log \left( \frac{p}{1-p} \right) = \text{the log-odds of a success (odds against failure)}$$



**Ex:** For instance, if the probability of a horse winning a race is  $p = 1/10$ , then

- the probability of a horse losing a race is  $1 - p = \underline{\hspace{2cm}}$
- thus, the odds of the horse winning the race is:

What does it mean to say that the odds of a given horse winning are 1 to 9? Think of it as "chances" in a box of tickets. For every 1 "chance" of winning there are 9 "chances" of losing.

In general, we define **odds for success** are “number of successes” to “number of failures.”

We can equivalently say that the odds against a given horse winning are 9 to 1. We can still think of it as "chances" in a box of tickets. For every 1 "chance" of winning there are 9 "chances" of losing, but we flip the order in which we state these odds...

In general, we define **odds against success** are “number of failures” to “number of successes.”

Ex: What are the odds against the horse winning?

In general, if  $p$  is the probability of some “success” level happening, then we call:

$$\frac{1-p}{p} = \text{the odds against a success (odds of failure)}$$

### Converting Odds (in “number of successes” to “number of failures” form) to Probability of Success

$$\text{probability of success} = p = \frac{\text{Number of Chances of Success}}{\text{Total Number of Chances}}$$

$$\text{probability of failure} = 1 - p = \frac{\text{Number of Chances of Failure}}{\text{Total Number of Chances}}$$

### Converting Odds (in numerical form) to Probabilities

$$\text{odds of success} = \frac{p}{1-p} \quad \text{and} \quad p = \frac{\text{odds of success}}{1 + \text{odds of success}}$$

$$\text{odds of failure} = \frac{1-p}{p} \quad \text{and} \quad p = \frac{\text{odds of success}}{1 + \text{odds of success}}$$

Ex: If the odds of a horse winning a race are 1 to 24, what is the probability that the horse wins the race?

Ex: If there is a 40% chance of rain, what are the odds it won't rain?

## 6. FITTING A SIMPLE LOGISTIC REGRESSION CURVE

So how do we find these optimal values  $\hat{\beta}_0, \hat{\beta}_1$  that best fit our data?

Efficient coefficient estimates are obtained for this model by the **method of maximum likelihood**.

We can use Python to estimate these optimal values  $\hat{\beta}_0, \hat{\beta}_1$  that best fit our data.

Logit Regression Results

Dep. Variable:	y	No. Observations:	1503			
Model:	Logit	Df Residuals:	1501			
Method:	MLE	Df Model:	1			
Date:	Sun, 25 Oct 2020	Pseudo R-squ.:	0.01091			
Time:	20:52:45	Log-Likelihood:	-955.55			
converged:	True	LL-Null:	-966.09			
Covariance Type:	nonrobust	LLR p-value:	4.412e-06			
	coef	std err	z	P> z	[0.025	0.975]
Intercept	-0.9211	0.082	-11.208	0.000	-1.082	-0.760
sex[T.Male]	0.5027	0.110	4.563	0.000	0.287	0.719

**Fit the logistic regression model.**

### Notation:

- Make sure to put a hat over the predictive probability  $p$  to indicate that this is a prediction.
- Put the variables into words.

## 7. FITTING A MULTIPLE LOGISTIC REGRESSION CURVE

Similarly, if we have multiple explanatory variables (or a categorical variable with >2 levels, which translates to more than one indicator variables), then we can fit multiple slopes  $\hat{\beta}_1, \dots, \hat{\beta}_p$  like we did with multiple linear regression.

$$\log\left(\frac{\hat{p}}{1-\hat{p}}\right) = \hat{\beta}_0 + \hat{\beta}_1 x_1 + \dots + \hat{\beta}_p x_p$$

We also use the **method of maximum likelihood** to find these optimal values  $\hat{\beta}_0, \hat{\beta}_1, \dots, \hat{\beta}_p$  that best fit our data. Python will also find these optimal values  $\hat{\beta}_0, \hat{\beta}_1, \dots, \hat{\beta}_p$  for us.

```
Optimization terminated successfully.  
Current function value: 0.612754  
Iterations 5
```

Logit Regression Results

Dep. Variable:	y	No. Observations:	691
Model:	Logit	Df Residuals:	688
Method:	MLE	Df Model:	2
Date:	Tue, 06 Apr 2021	Pseudo R-squ.:	0.06252
Time:	10:40:34	Log-Likelihood:	-423.41
converged:	True	LL-Null:	-451.65
Covariance Type:	nonrobust	LLR p-value:	5.457e-13

	coef	std err	z	P> z	[0.025	0.975]
Intercept	-2.3609	0.288	-8.190	0.000	-2.926	-1.796
sex[T.Male]	0.8856	0.167	5.296	0.000	0.558	1.213
age	0.0260	0.005	5.466	0.000	0.017	0.035

Categorical explanatory variables are translated into indicator variables in the same way that they are for multiple linear regression models.

**Ex:** Use the Python summary output to formulate the logistic regression model.

$$\log\left(\frac{\hat{p}}{1 - \hat{p}}\right) = -2.3609 + .8858sex[T.Male] + 0.026age$$

## 8. MAKING PREDICTIONS WITH A LOGISTIC REGRESSION MODEL

Use this logistic regression model to predict the following:

- a. The log-odds that a 20-year-old female supports the president's foreign policy *in the sample*.
- b. The odds that a 20-year-old female supports the president's foreign policy *in the sample*.
- c. The probability that a 20-year-old female supports the president's foreign policy *in the sample*.

## 9. INTERPRETING THE VALUES OF $\widehat{\beta}_0, \widehat{\beta}_1, \widehat{\beta}_2 \dots, \widehat{\beta}_p$

### Easier Interpretation of the Logistic Regression Model

$$\log\left(\frac{\widehat{p}}{1-\widehat{p}}\right) = \widehat{\beta}_0 + \widehat{\beta}_1 x_1 + \dots + \widehat{\beta}_p x_p$$

We can exponentiate both sides of the logistic regression model above to help us interpret what the values of  $\widehat{\beta}_0, \widehat{\beta}_1, \widehat{\beta}_2 \dots, \widehat{\beta}_p$  mean.

$$\left(\frac{\widehat{p}}{1-\widehat{p}}\right) = e^{\widehat{\beta}_0 + \widehat{\beta}_1 x_1 + \dots + \widehat{\beta}_p x_p}$$

Remember that  $\left(\frac{\widehat{p}}{1-\widehat{p}}\right)$  represents the predicted odds of the “success level” happening for a given set of values for  $x_1, x_2, \dots, x_p$ .

### Interpreting $\widehat{\beta}_0$

By plugging in 0 for all  $x_i$  (for  $i=1, \dots, p$ ), we get:

$$\left(\frac{\widehat{p}}{1-\widehat{p}}\right) = e^{\widehat{\beta}_0},$$

Thus, we call  $e^{\widehat{\beta}_0}$  the **baseline odds** as it is the \_\_\_\_\_ when  $x_1 = 0, x_2 = 0, \dots, x_p = 0$ .

Ex:  $\widehat{\beta}_0 = -2.3609$

So, we say that the baseline odds of someone in the sample supporting the president’s foreign policy are  $e^{\widehat{\beta}_0} = 0.0943$  (or in other words, 943 to 10,000). This would represent a person who is (female and 0 years old, which is nonsensical).

### Interpreting $\hat{\beta}_i$ ( $i = 1, \dots, p$ ), when $x_i$ is numerical:

We can use properties of exponents to represent the logistic regression equation as

$$odds = \left( \frac{\hat{p}}{1 - \hat{p}} \right) = e^{\hat{\beta}_0} \cdot e^{\hat{\beta}_1 x_1} \cdot e^{\hat{\beta}_2 x_2} \cdot \dots \cdot e^{\hat{\beta}_p x_p}$$

If we were to increase  $x_i$  by one and hold everything else constant, we would get:

$$\frac{odds_{new}}{odds_{old}} = \frac{\left( \frac{\hat{p}}{1 - \hat{p}} \right)_{new}}{\left( \frac{\hat{p}}{1 - \hat{p}} \right)_{old}} = \frac{e^{\hat{\beta}_0} \cdot e^{\hat{\beta}_1 x_1} \cdot \dots \cdot e^{\hat{\beta}_i (x_i + 1)} \cdot \dots \cdot e^{\hat{\beta}_p x_p}}{e^{\hat{\beta}_0} \cdot e^{\hat{\beta}_1 x_1} \cdot \dots \cdot e^{\hat{\beta}_i x_i} \cdot \dots \cdot e^{\hat{\beta}_p x_p}} = e^{\hat{\beta}_i}$$

Thus we call,  $e^{\hat{\beta}_i}$ , the **odds multiplier** as this is the *multiple* we would expect the odds to increase (or decrease if  $e^{\hat{\beta}_i} < 0$ ) by if we increased  $x_i$  by 1 (holding everything else constant).

Ex:  $\hat{\beta}_1 = 0.0260$

So, we say that we would expect the odds of someone in the sample supporting the president's foreign policy to increase by a factor of  $e^{\hat{\beta}_1} = 1.026$ , all else held equal.

### Interpreting $\hat{\beta}_i$ ( $i = 1, \dots, p$ ), when $x_i$ is an indicator variable:

Ex: Using our fitted logistic regression model

$$\log\left(\frac{\hat{p}}{1-\hat{p}}\right) = -2.3609 + .08858sex[T.Male] + 0.026age,$$

If we were to find the difference in the log-odds of a male supporting the president's foreign policy vs. a female supporting the president's foreign policy (everything else held equal), we would get:

$$\begin{aligned}\log\left(\frac{\hat{p}_{male}}{1-\hat{p}_{male}}\right) - \log\left(\frac{\hat{p}_{female}}{1-\hat{p}_{female}}\right) \\ = (-2.3609 + .8858(1) + 0.026age) - (-2.3609 + .8858(0) + 0.026age) = .8858\end{aligned}$$

By using properties of logarithms, we also get

$$\log\left(\frac{\frac{\hat{p}_{male}}{1-\hat{p}_{male}}}{\frac{\hat{p}_{female}}{1-\hat{p}_{female}}}\right) = .8858$$

And by exponentiating both sides, we get:

$$\frac{odds_{male}}{odds_{female}} = \frac{\frac{\hat{p}_{male}}{1-\hat{p}_{male}}}{\frac{\hat{p}_{female}}{1-\hat{p}_{female}}} = e^{.8858} = 2.424$$

So this shows that  $e^{.8858} = 2.424$  represents an **odds ratio**, which in this case represents that we expect that the odds that a male supports the president's foreign policy is 2.424 higher than the odds that a female supports the president's foreign policy.

### Interpretation in General:

If  $x_i$  is an indicator variable, we expect that the odds that the level of the categorical explanatory variable in which  $x_i = 1$  is a success will be  $e^{\hat{\beta}_i}$  higher (or lower if  $e^{\hat{\beta}_i} < 0$ ) than the reference level for that categorical variable level.