

# LECTURE 13: DUMMY VARIABLES

ECON 480 - ECONOMETRICS - FALL 2018

---

Ryan Safner

November 7, 2018

Dummy Variables

Recoding Dummies

## DUMMY VARIABLES

---

- You now have the “minimal toolkit” for data analysis with multivariate OLS regression

## OVERVIEW OF THE REMAINDER OF THE COURSE

- You now have the “minimal toolkit” for data analysis with multivariate OLS regression
- The remainder of the course is about two types of extensions to the toolkit:

## OVERVIEW OF THE REMAINDER OF THE COURSE

- You now have the “minimal toolkit” for data analysis with multivariate OLS regression
- The remainder of the course is about two types of extensions to the toolkit:
  1. “Data Wrangling”

## OVERVIEW OF THE REMAINDER OF THE COURSE

- You now have the “minimal toolkit” for data analysis with multivariate OLS regression
- The remainder of the course is about two types of extensions to the toolkit:
  1. “Data Wrangling”
    - Altering variables or data for useful analysis

## OVERVIEW OF THE REMAINDER OF THE COURSE

- You now have the “minimal toolkit” for data analysis with multivariate OLS regression
- The remainder of the course is about two types of extensions to the toolkit:
  1. “Data Wrangling”
    - Altering variables or data for useful analysis
    - Dummy variables for categorical data (**R** calls them **factors**)



## OVERVIEW OF THE REMAINDER OF THE COURSE

- You now have the “minimal toolkit” for data analysis with multivariate OLS regression
- The remainder of the course is about two types of extensions to the toolkit:
  1. “Data Wrangling”
    - Altering variables or data for useful analysis
    - Dummy variables for categorical data (**R** calls them **factors**)
    - Transforming variable scales or models to fit nonlinear data

- You now have the “minimal toolkit” for data analysis with multivariate OLS regression
- The remainder of the course is about two types of extensions to the toolkit:
  1. “Data Wrangling”
    - Altering variables or data for useful analysis
    - Dummy variables for categorical data (**R** calls them **factors**)
    - Transforming variable scales or models to fit nonlinear data
  2. Advanced identification strategies and unique problems

- You now have the “minimal toolkit” for data analysis with multivariate OLS regression
- The remainder of the course is about two types of extensions to the toolkit:
  1. “Data Wrangling”
    - Altering variables or data for useful analysis
    - Dummy variables for categorical data (**R** calls them **factors**)
    - Transforming variable scales or models to fit nonlinear data
  2. Advanced identification strategies and unique problems
    - Time Series data

- You now have the “minimal toolkit” for data analysis with multivariate OLS regression
- The remainder of the course is about two types of extensions to the toolkit:
  1. “Data Wrangling”
    - Altering variables or data for useful analysis
    - Dummy variables for categorical data (**R** calls them **factors**)
    - Transforming variable scales or models to fit nonlinear data
  2. Advanced identification strategies and unique problems
    - Time Series data
    - Panel data

- You now have the “minimal toolkit” for data analysis with multivariate OLS regression
- The remainder of the course is about two types of extensions to the toolkit:
  1. “Data Wrangling”
    - Altering variables or data for useful analysis
    - Dummy variables for categorical data (**R** calls them **factors**)
    - Transforming variable scales or models to fit nonlinear data
  2. Advanced identification strategies and unique problems
    - Time Series data
    - Panel data
    - Fixed effects and random effects models

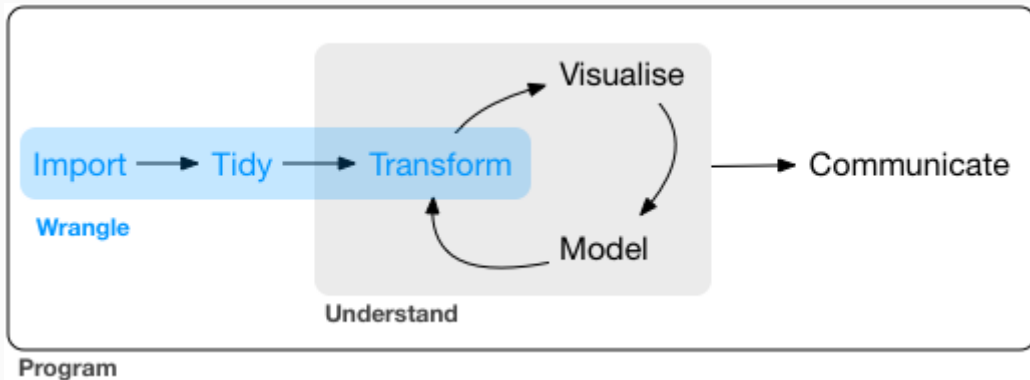
- You now have the “minimal toolkit” for data analysis with multivariate OLS regression
- The remainder of the course is about two types of extensions to the toolkit:
  1. “Data Wrangling”
    - Altering variables or data for useful analysis
    - Dummy variables for categorical data (**R** calls them **factors**)
    - Transforming variable scales or models to fit nonlinear data
  2. Advanced identification strategies and unique problems
    - Time Series data
    - Panel data
    - Fixed effects and random effects models
    - Difference-in-difference models

- You now have the “minimal toolkit” for data analysis with multivariate OLS regression
- The remainder of the course is about two types of extensions to the toolkit:
  1. “Data Wrangling”
    - Altering variables or data for useful analysis
    - Dummy variables for categorical data (**R** calls them **factors**)
    - Transforming variable scales or models to fit nonlinear data
  2. Advanced identification strategies and unique problems
    - Time Series data
    - Panel data
    - Fixed effects and random effects models
    - Difference-in-difference models
    - Instrumental variables models

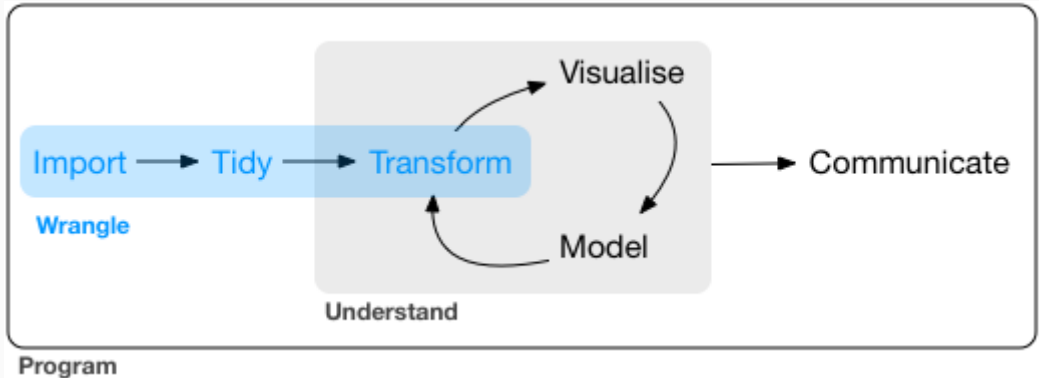
- You now have the “minimal toolkit” for data analysis with multivariate OLS regression
- The remainder of the course is about two types of extensions to the toolkit:
  1. “Data Wrangling”
    - Altering variables or data for useful analysis
    - Dummy variables for categorical data (**R** calls them **factors**)
    - Transforming variable scales or models to fit nonlinear data
  2. Advanced identification strategies and unique problems
    - Time Series data
    - Panel data
    - Fixed effects and random effects models
    - Difference-in-difference models
    - Instrumental variables models
    - Linear probability, logit, and probit models



- "Data wrangling" is a term for altering and cleaning data from raw form (often unusable) to a form that is useful for analysis (e.g. plotting and regressions)



- "Data wrangling" is a term for altering and cleaning data from raw form (often unusable) to a form that is useful for analysis (e.g. plotting and regressions)
- A **significant portion** of data analysis is initial data wrangling



- Recall **categorical variables** place an individual into one of several possible categories

- Recall **categorical variables** place an individual into one of several possible categories
  - e.g. sex, season, political party

- Recall **categorical variables** place an individual into one of several possible categories
  - e.g. sex, season, political party
  - may be responses to questions

- Recall **categorical variables** place an individual into one of several possible categories
  - e.g. sex, season, political party
  - may be responses to questions
  - can be quantitative (e.g. age, zip code)

Cut	Fair	Good	Very Good	Premium	Ideal
Count	1610	4906	12082	13791	21551
Proportion	0.030	0.091	0.224	0.256	0.400

Cut characteristics of 53,940 diamonds

- Recall **categorical variables** place an individual into one of several possible categories
  - e.g. sex, season, political party
  - may be responses to questions
  - can be quantitative (e.g. age, zip code)

Cut	Fair	Good	Very Good	Premium	Ideal
Count	1610	4906	12082	13791	21551
Proportion	0.030	0.091	0.224	0.256	0.400

Cut characteristics of 53,940 diamonds

- Also recall R calls this type of data a **factor**

### Example

Do men earn higher wages on average than women?

- Using basic statistics, we can test for a statistically significant difference in group means with a ***t*-test**<sup>1</sup>

---

<sup>1</sup>See the **Handout** on Blackboard for this example.



## Example

Do men earn higher wages on average than women?

- Using basic statistics, we can test for a statistically significant difference in group means with a ***t*-test**<sup>1</sup>
- Let:

---

<sup>1</sup>See the **Handout** on Blackboard for this example.

### Example

Do men earn higher wages on average than women?

- Using basic statistics, we can test for a statistically significant difference in group means with a **t-test**<sup>1</sup>
- Let:
  - $\bar{Y}_M$  the average earnings of a sample of  $n_M$  men

---

<sup>1</sup>See the **Handout** on Blackboard for this example.

### Example

Do men earn higher wages on average than women?

- Using basic statistics, we can test for a statistically significant difference in group means with a **t-test**<sup>1</sup>
- Let:
  - $\bar{Y}_M$  the average earnings of a sample of  $n_M$  men
  - $\bar{Y}_W$  the average earnings of a sample of  $n_W$  women

---

<sup>1</sup>See the **Handout** on Blackboard for this example.

## Example

Do men earn higher wages on average than women?

- Using basic statistics, we can test for a statistically significant difference in group means with a **t-test**<sup>1</sup>
- Let:
  - $\bar{Y}_M$  the average earnings of a sample of  $n_M$  men
  - $\bar{Y}_W$  the average earnings of a sample of  $n_W$  women
  - Difference in group averages  $d = \bar{Y}_M - \bar{Y}_W$

---

<sup>1</sup>See the **Handout** on Blackboard for this example.

## Example

Do men earn higher wages on average than women?

- Using basic statistics, we can test for a statistically significant difference in group means with a **t-test**<sup>1</sup>
- Let:
  - $\bar{Y}_M$  the average earnings of a sample of  $n_M$  men
  - $\bar{Y}_W$  the average earnings of a sample of  $n_W$  women
  - Difference in group averages  $d = \bar{Y}_M - \bar{Y}_W$
- The hypothesis test is:

---

<sup>1</sup>See the **Handout** on Blackboard for this example.

## Example

Do men earn higher wages on average than women?

- Using basic statistics, we can test for a statistically significant difference in group means with a **t-test**<sup>1</sup>
- Let:
  - $\bar{Y}_M$  the average earnings of a sample of  $n_M$  men
  - $\bar{Y}_W$  the average earnings of a sample of  $n_W$  women
  - Difference in group averages  $d = \bar{Y}_M - \bar{Y}_W$
- The hypothesis test is:
- $H_0 : d = 0$

---

<sup>1</sup>See the **Handout** on Blackboard for this example.

## Example

Do men earn higher wages on average than women?

- Using basic statistics, we can test for a statistically significant difference in group means with a **t-test**<sup>1</sup>
- Let:
  - $\bar{Y}_M$  the average earnings of a sample of  $n_M$  men
  - $\bar{Y}_W$  the average earnings of a sample of  $n_W$  women
  - Difference in group averages  $d = \bar{Y}_M - \bar{Y}_W$
- The hypothesis test is:
  - $H_0 : d = 0$
  - $H_1 : d \neq 0$

---

<sup>1</sup>See the **Handout** on Blackboard for this example.

- In a regression, we can easily compare across groups via a **dummy variable**<sup>2</sup>

---

<sup>2</sup>Also called a **binary variable** or **dichotomous variable**



- In a regression, we can easily compare across groups via a **dummy variable**<sup>2</sup>
  - Dummy variable *only* = 0 or = 1, depending on if a condition is met

---

<sup>2</sup>Also called a **binary variable** or **dichotomous variable**

- In a regression, we can easily compare across groups via a **dummy variable**<sup>2</sup>
  - Dummy variable *only* = 0 or = 1, depending on if a condition is met
  - Signifies whether an observation belongs to a category or not

---

<sup>2</sup>Also called a **binary variable** or **dichotomous variable**

- In a regression, we can easily compare across groups via a **dummy variable**<sup>2</sup>
  - Dummy variable *only* = 0 or = 1, depending on if a condition is met
  - Signifies whether an observation belongs to a category or not

### Example

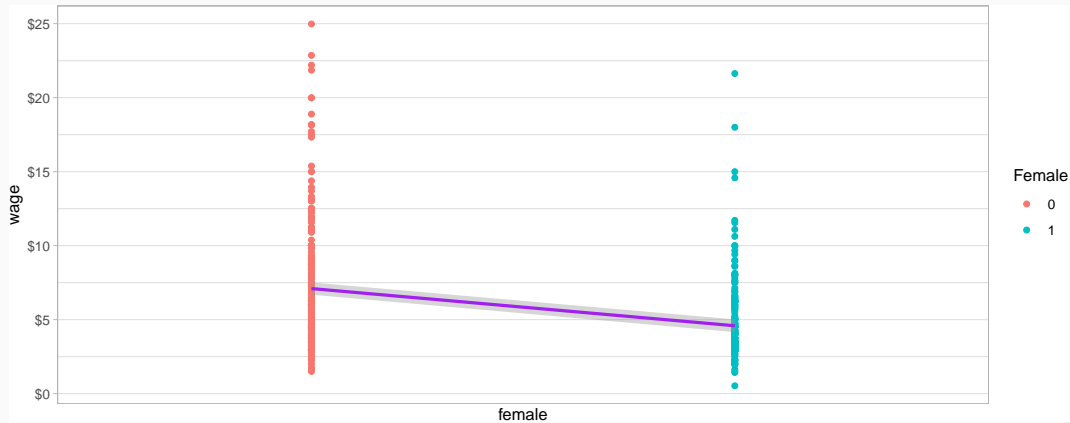
$$\widehat{Wage}_i = \hat{\beta}_0 + \hat{\beta}_1 Female_i \quad \text{where } Female_i = \begin{cases} 1 & \text{if } i \text{ is Female} \\ 0 & \text{If } i \text{ is Male} \end{cases}$$

- Again,  $\hat{\beta}_1$  makes less sense as the “slope” of a line in this context

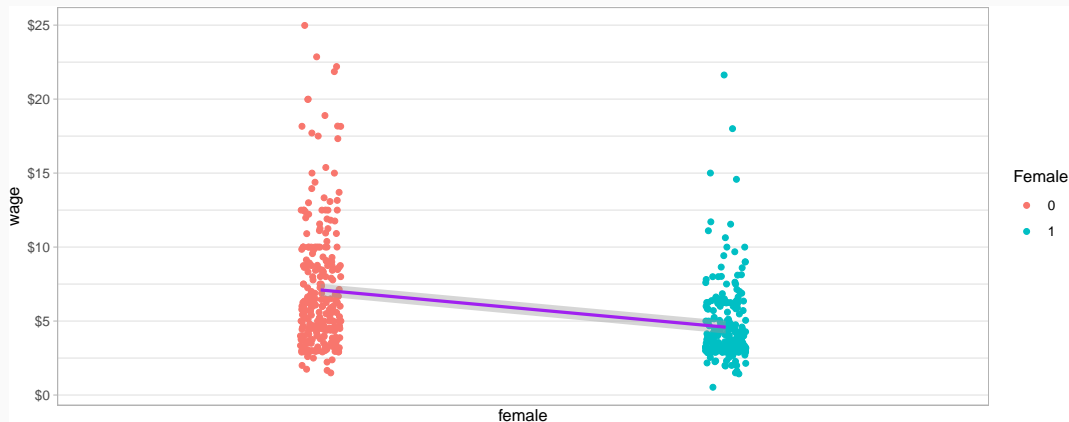
---

<sup>2</sup>Also called a **binary variable** or **dichotomous variable**

## COMPARING GROUPS IN REGRESSION: SCATTERPLOT

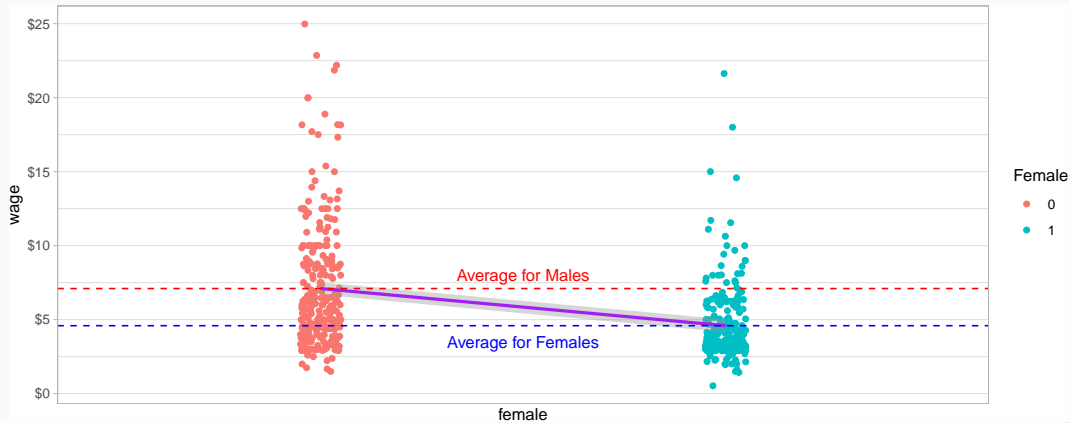


## COMPARING GROUPS IN REGRESSION: SCATTERPLOT WITH JITTERING



- use `geom_jitter()` instead of `geom_point()` to “jitter” the data to avoid overplotting

## COMPARING GROUPS IN REGRESSION: SCATTERPLOT WITH JITTERING II



$$\hat{Y}_i = \hat{\beta}_0 + \hat{\beta}_1 D_i \quad \text{where } D_i = \{0, 1\}$$

- When  $D_i = 0$  (Control group):

$$\hat{Y}_i = \hat{\beta}_0 + \hat{\beta}_1 D_i \quad \text{where } D_i = \{0, 1\}$$

- When  $D_i = 0$  (Control group):
  - $\hat{Y}_i = \hat{\beta}_0$



$$\hat{Y}_i = \hat{\beta}_0 + \hat{\beta}_1 D_i \quad \text{where } D_i = \{0, 1\}$$

- When  $D_i = 0$  (Control group):
  - $\hat{Y}_i = \hat{\beta}_0$
  - $E[Y|D_i = 0] = \hat{\beta}_0 \iff$  the mean of  $Y$  when  $D_i = 0$

$$\hat{Y}_i = \hat{\beta}_0 + \hat{\beta}_1 D_i \quad \text{where } D_i = \{0, 1\}$$

- When  $D_i = 0$  (Control group):
  - $\hat{Y}_i = \hat{\beta}_0$
  - $E[Y|D_i = 0] = \hat{\beta}_0 \iff$  the mean of  $Y$  when  $D_i = 0$
- When  $D_i = 1$  (Treatment group):

$$\hat{Y}_i = \hat{\beta}_0 + \hat{\beta}_1 D_i \quad \text{where } D_i = \{0, 1\}$$

- When  $D_i = 0$  (Control group):
  - $\hat{Y}_i = \hat{\beta}_0$
  - $E[Y|D_i = 0] = \hat{\beta}_0 \iff$  the mean of  $Y$  when  $D_i = 0$
- When  $D_i = 1$  (Treatment group):
  - $\hat{Y}_i = \hat{\beta}_0 + \hat{\beta}_1 D_i$

$$\hat{Y}_i = \hat{\beta}_0 + \hat{\beta}_1 D_i \quad \text{where } D_i = \{0, 1\}$$

- When  $D_i = 0$  (Control group):
  - $\hat{Y}_i = \hat{\beta}_0$
  - $E[Y|D_i = 0] = \hat{\beta}_0 \iff$  the mean of  $Y$  when  $D_i = 0$
- When  $D_i = 1$  (Treatment group):
  - $\hat{Y}_i = \hat{\beta}_0 + \hat{\beta}_1 D_i$
  - $E[Y|D_i = 1] = \hat{\beta}_0 + \hat{\beta}_1 \iff$  the mean of  $Y$  when  $D_i = 1$

$$\hat{Y}_i = \hat{\beta}_0 + \hat{\beta}_1 D_i \quad \text{where } D_i = \{0, 1\}$$

- When  $D_i = 0$  (Control group):
  - $\hat{Y}_i = \hat{\beta}_0$
  - $E[Y|D_i = 0] = \hat{\beta}_0 \iff$  the mean of  $Y$  when  $D_i = 0$
- When  $D_i = 1$  (Treatment group):
  - $\hat{Y}_i = \hat{\beta}_0 + \hat{\beta}_1 D_i$
  - $E[Y|D_i = 1] = \hat{\beta}_0 + \hat{\beta}_1 \iff$  the mean of  $Y$  when  $D_i = 1$
- So the *difference* in group means:

$$\hat{Y}_i = \hat{\beta}_0 + \hat{\beta}_1 D_i \quad \text{where } D_i = \{0, 1\}$$

- When  $D_i = 0$  (Control group):

- $\hat{Y}_i = \hat{\beta}_0$
- $E[Y|D_i = 0] = \hat{\beta}_0 \iff$  the mean of  $Y$  when  $D_i = 0$

- When  $D_i = 1$  (Treatment group):

- $\hat{Y}_i = \hat{\beta}_0 + \hat{\beta}_1 D_i$
- $E[Y|D_i = 1] = \hat{\beta}_0 + \hat{\beta}_1 \iff$  the mean of  $Y$  when  $D_i = 1$

- So the *difference* in group means:

$$\begin{aligned} &= E[Y_i|D_i = 1] - E[Y_i|D_i = 0] \\ &= (\hat{\beta}_0 + \hat{\beta}_1) - (\hat{\beta}_0) \\ &= \hat{\beta}_1 \end{aligned}$$

### Example

$$\widehat{Wage}_i = \hat{\beta}_0 + \hat{\beta}_1 Female_i$$
$$Female_i = \begin{cases} 1 & \text{if } i \text{ is Female} \\ 0 & \text{If } i \text{ is Male} \end{cases}$$

### Example

$$\widehat{Wage}_i = \hat{\beta}_0 + \hat{\beta}_1 Female_i$$
$$Female_i = \begin{cases} 1 & \text{if } i \text{ is Female} \\ 0 & \text{If } i \text{ is Male} \end{cases}$$

- Mean wage for males:



### Example

$$\widehat{Wage}_i = \hat{\beta}_0 + \hat{\beta}_1 Female_i$$
$$Female_i = \begin{cases} 1 & \text{if } i \text{ is Female} \\ 0 & \text{if } i \text{ is Male} \end{cases}$$

- Mean wage for males:  $E[Wage | Female = 0] = \hat{\beta}_0$

### Example

$$\widehat{Wage}_i = \hat{\beta}_0 + \hat{\beta}_1 Female_i$$
$$Female_i = \begin{cases} 1 & \text{if } i \text{ is Female} \\ 0 & \text{If } i \text{ is Male} \end{cases}$$

- Mean wage for males:  $E[Wage | Female = 0] = \hat{\beta}_0$
- Mean wage for females:

### Example

$$\widehat{Wage}_i = \hat{\beta}_0 + \hat{\beta}_1 Female_i$$
$$Female_i = \begin{cases} 1 & \text{if } i \text{ is Female} \\ 0 & \text{if } i \text{ is Male} \end{cases}$$

- Mean wage for males:  $E[Wage|Female = 0] = \hat{\beta}_0$
- Mean wage for females:  $E[Wage|Female = 1] = \hat{\beta}_0 + \hat{\beta}_1$

### Example

$$\widehat{Wage}_i = \hat{\beta}_0 + \hat{\beta}_1 Female_i$$
$$Female_i = \begin{cases} 1 & \text{if } i \text{ is Female} \\ 0 & \text{if } i \text{ is Male} \end{cases}$$

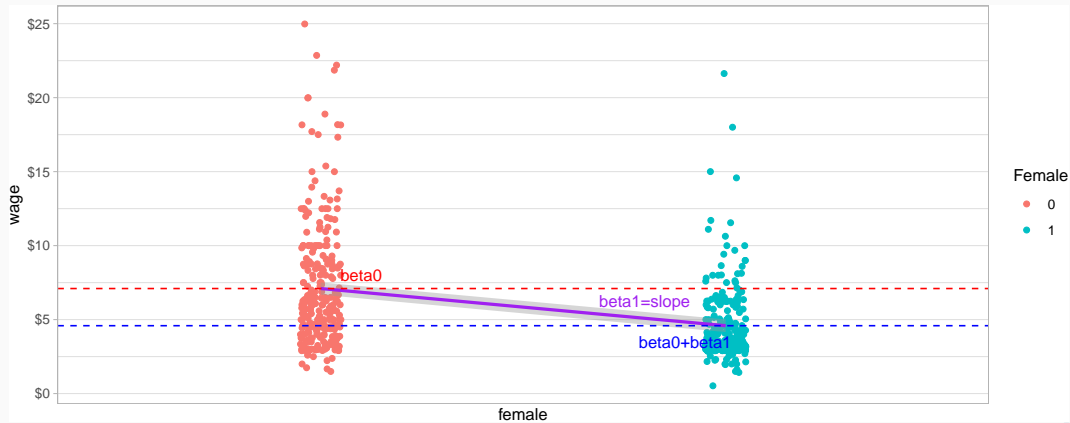
- Mean wage for males:  $E[Wage|Female = 0] = \hat{\beta}_0$
- Mean wage for females:  $E[Wage|Female = 1] = \hat{\beta}_0 + \hat{\beta}_1$
- Difference in wage between males & females:

### Example

$$\widehat{Wage}_i = \hat{\beta}_0 + \hat{\beta}_1 Female_i$$
$$Female_i = \begin{cases} 1 & \text{if } i \text{ is Female} \\ 0 & \text{if } i \text{ is Male} \end{cases}$$

- Mean wage for males:  $E[Wage|Female = 0] = \hat{\beta}_0$
- Mean wage for females:  $E[Wage|Female = 1] = \hat{\beta}_0 + \hat{\beta}_1$
- Difference in wage between males & females:  $\hat{\beta}_1$

## DUMMY VARIABLES AS GROUP MEANS II





- OLS Regression:  $\widehat{Wage}_i = 7.10 - 2.51 \text{ Female}_i$   
(0.21) (0.30)



## DUMMY REGRESSION VS. GROUP MEANS

- OLS Regression:  $\widehat{\text{Wage}}_i = 7.10 - 2.51 \text{ Female}_i$   
(0.21) (0.30)
- Simple tabulation of group means:

Sex	Avg. Wage ( $\bar{Y}$ )	SE(avg) ( $s_Y$ )	n
Female	4.59	0.16	252
Male	7.10	0.21	274
Difference	-2.51	(0.30)	-

- OLS Regression:  $\widehat{\text{Wage}}_i = 7.10 - 2.51 \text{ Female}_i$   
(0.21) (0.30)
- Simple tabulation of group means:

Sex	Avg. Wage ( $\bar{Y}$ )	SE(avg) ( $s_Y$ )	n
Female	4.59	0.16	252
Male	7.10	0.21	274
Difference	-2.51	(0.30)	-

- Differences in means:  $\bar{Y}_F - \bar{Y}_M = 4.59 - 7.10 = -2.51$

## DUMMY REGRESSION VS. GROUP MEANS

- OLS Regression:  $\widehat{\text{Wage}}_i = 7.10 - 2.51 \text{ Female}_i$   
(0.21) (0.30)
- Simple tabulation of group means:

Sex	Avg. Wage ( $\bar{Y}$ )	SE(avg) ( $s_Y$ )	n
Female	4.59	0.16	252
Male	7.10	0.21	274
Difference	-2.51	(0.30)	-

- Differences in means:  $\bar{Y}_F - \bar{Y}_M = 4.59 - 7.10 = -2.51$
- $SE(\bar{Y}_F - \bar{Y}_M) = \sqrt{\frac{s_M^2}{n_M} + \frac{s_F^2}{n_F}} = \sqrt{\frac{0.21^2}{274} + \frac{0.16^2}{252}} \approx 0.30$

```
# Our data comes from WAGE1.dta which you can find in Blackboard under data

# Load WAGE1 as wages
library("foreign") # to load .dta Stata files
wages<-read.dta("../Data/WAGE1.dta")

# there's a lot of variables in wages, let's only look at wage and female for now
wages<-subset(wages, select=c("wage","female"))
```

```
# just get a sense of the data
```

```
head(wages)
```

```
##   wage female
## 1 3.10      1
## 2 3.24      1
## 3 3.00      0
## 4 6.00      0
## 5 5.30      0
## 6 8.75      0
```

- We want to look at the data under certain **conditions**

---

<sup>3</sup>Later, I will show you how to do this in **dplyr**, a popular package that makes data wrangling easier

- We want to look at the data under certain **conditions**
- Can do this in base R by **subsetting** data using square brackets [ ]<sup>3</sup>

---

<sup>3</sup>Later, I will show you how to do this in **dplyr**, a popular package that makes data wrangling easier

- We want to look at the data under certain **conditions**
- Can do this in base R by **subsetting** data using square brackets `[]`<sup>3</sup>
- Syntax: `data[df$variable condition]` where **condition** is likely:

---

<sup>3</sup>Later, I will show you how to do this in **dplyr**, a popular package that makes data wrangling easier



- We want to look at the data under certain **conditions**
- Can do this in base R by **subsetting** data using square brackets `[]`<sup>3</sup>
- Syntax: `data[df$variable condition]` where **condition** is likely:
  - A logical test, i.e. `>`, `<`, `!=`, `<=`, `>=`, `==` some value

---

<sup>3</sup>Later, I will show you how to do this in **dplyr**, a popular package that makes data wrangling easier

```
# look at average wage for men  
summary(wages$wage[wages$female==0])
```

```
##      Min. 1st Qu.  Median    Mean 3rd Qu.    Max.  
##    1.500   4.143   6.000   7.099   8.765  24.980
```

```
sd(wages$wage[wages$female==0]) # get sd
```

```
## [1] 4.160858
```

```
# look at average wage for women  
summary(wages$wage[wages$female==1])
```

```
##      Min. 1st Qu.  Median    Mean 3rd Qu.    Max.  
##    0.530   3.000   3.750   4.588   5.510   21.630
```

```
sd(wages$wage[wages$female==1]) # get sd
```

```
## [1] 2.529363
```

```
dummyreg<-lm(wage~female, data=wages)
summary(dummyreg)
```

```
##
## Call:
## lm(formula = wage ~ female, data = wages)
##
## Residuals:
##      Min       1Q   Median       3Q      Max
## -5.5995 -1.8495 -0.9877  1.4260 17.8805
##
## Coefficients:
##              Estimate Std. Error t value Pr(>|t|)
## (Intercept)   7.0995     0.2100  33.806 < 2e-16 ***
## female       -2.5118     0.3034  -8.279 1.04e-15 ***
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 3.476 on 524 degrees of freedom
```

$$\widehat{Wage}_i = 7.10 - 2.51 \text{ Female}_i \\ (0.21) \quad (0.30)$$

$$\widehat{\text{Wage}}_i = 7.10 - 2.51 \text{ Female}_i$$

(0.21)   (0.30)

- Does this mean we've accurately measured the gender-wage gap as \$2.51/hr?

$$\widehat{Wage}_i = 7.10 - 2.51 \text{ Female}_i$$

(0.21)   (0.30)

- Does this mean we've accurately measured the gender-wage gap as \$2.51/hr?
- Are there variables for which the following is true?

$$\text{corr}(\text{wage}, Z) \neq 0$$

$$\text{corr}(\text{female}, Z) \neq 0$$

$$\widehat{Wage}_i = 7.10 - 2.51 \text{ Female}_i$$

(0.21)   (0.30)

- Does this mean we've accurately measured the gender-wage gap as \$2.51/hr?
- Are there variables for which the following is true?

$$\begin{aligned} \text{corr}(\text{wage}, Z) &\neq 0 \\ \text{corr}(\text{female}, Z) &\neq 0 \end{aligned}$$

- `female` is probably endogenous, must include other control variables

## RECODING DUMMIES

---



- What if instead of *female* we had used:

### Example

$$\widehat{Wage}_i = \hat{\beta}_0 + \hat{\beta}_1 Male_i$$

$$Male_i = \begin{cases} 1 & \text{if } i \text{ is Male} \\ 0 & \text{If } i \text{ is Female} \end{cases}$$

- What if instead of *female* we had used:

### Example

$$\widehat{Wage}_i = \hat{\beta}_0 + \hat{\beta}_1 Male_i$$
$$Male_i = \begin{cases} 1 & \text{if } i \text{ is Male} \\ 0 & \text{If } i \text{ is Female} \end{cases}$$

- *female* is a variable already in the data, we need to generate the *male* variable

- Again, a very useful R function is

```
ifelse(conditions, do.this.if.true, do.this.if.false)
```

- Again, a very useful R function is

```
ifelse(conditions, do.this.if.true, do.this.if.false)
```

- So let's create a `male` variable in our `wages` dataframe that we define as `1` if `female==0` and `0` otherwise (i.e. if `female==1`)

- Again, a very useful R function is

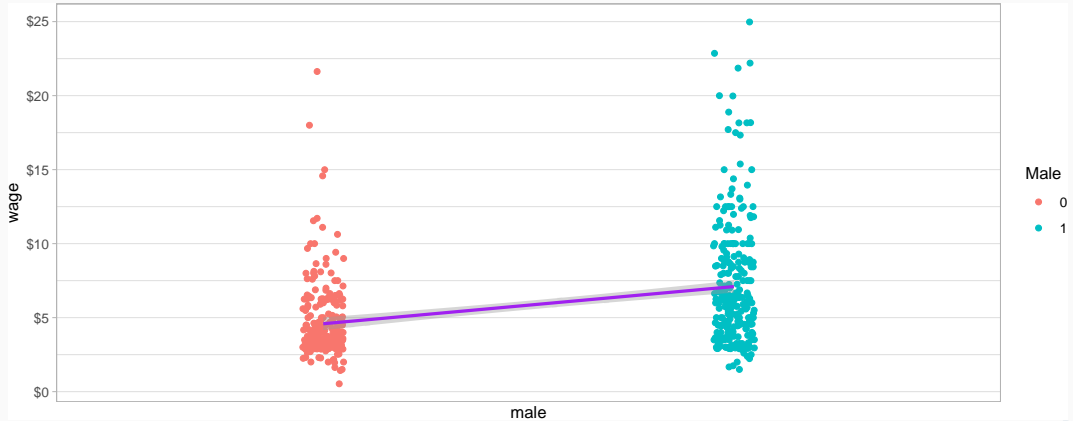
```
ifelse(conditions, do.this.if.true, do.this.if.false)
```

- So let's create a `male` variable in our `wages` dataframe that we define as 1 if `female==0` and 0 otherwise (i.e. if `female==1`)

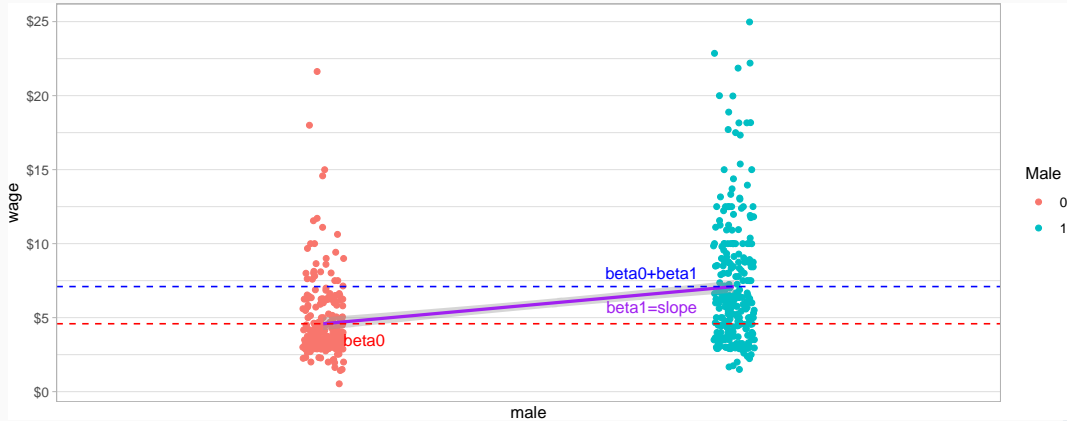
```
wages$male<-ifelse(wages$female==0,1,0)  
head(wages) # verify that it worked
```

```
##   wage female male  
## 1 3.10      1    0  
## 2 3.24      1    0  
## 3 3.00      0    1  
## 4 6.00      0    1  
## 5 5.30      0    1  
## 6 8.75      0    1
```

## SCATTERPLOT WITH MALE



## SCATTERPLOT WITH MALE II



## THE DUMMY REGRESSION WITH MALE

```
mreg<-lm(wage~male, data=wages)
summary(mreg)
```

```
##
## Call:
## lm(formula = wage ~ male, data = wages)
##
## Residuals:
##      Min       1Q   Median       3Q      Max
## -5.5995 -1.8495 -0.9877  1.4260 17.8805
##
## Coefficients:
##              Estimate Std. Error t value Pr(>|t|)
## (Intercept)   4.5877     0.2190  20.950 < 2e-16 ***
## male          2.5118     0.3034   8.279 1.04e-15 ***
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 3.476 on 524 degrees of freedom
```

$$\widehat{Wage}_i = 4.59 + 2.51 Male_i$$

(0.21) (0.30)



## THE DUMMY REGRESSION: MALE OR FEMALE

```
library("stargazer")
stargazer(dummyreg, mreg, type="latex",
          header=FALSE, float=FALSE)
```

	Dependent variable:	
	wage	
	(1)	(2)
female	-2.512*** (0.303)	
male		2.512*** (0.303)
Constant	7.099*** (0.210)	4.588*** (0.219)
Observations	526	526
R <sup>2</sup>	0.116	0.116
Adjusted R <sup>2</sup>	0.114	0.114
Residual Std. Error (df = 524)	3.476	3.476
F Statistic (df = 1; 524)	68.537***	68.537***

Note:

\* p<0.1; \*\* p<0.05; \*\*\* p<0.01

- Note it doesn't matter if we use male or female  
males always earn \$2.51 more than females

## THE DUMMY REGRESSION: MALE OR FEMALE

```
library("stargazer")
stargazer(dummyreg, mreg, type="latex",
          header=FALSE, float=FALSE)
```

	Dependent variable:	
	wage	
	(1)	(2)
female	-2.512*** (0.303)	
male		2.512*** (0.303)
Constant	7.099*** (0.210)	4.588*** (0.219)
Observations	526	526
R <sup>2</sup>	0.116	0.116
Adjusted R <sup>2</sup>	0.114	0.114
Residual Std. Error (df = 524)	3.476	3.476
F Statistic (df = 1; 524)	68.537***	68.537***

Note:

\* p<0.1; \*\* p<0.05; \*\*\* p<0.01

- Note it doesn't matter if we use **male** or **female** males always earn \$2.51 more than females
- Compare the constant (mean for the **D=0** group)