

8.4 – Interaction terms

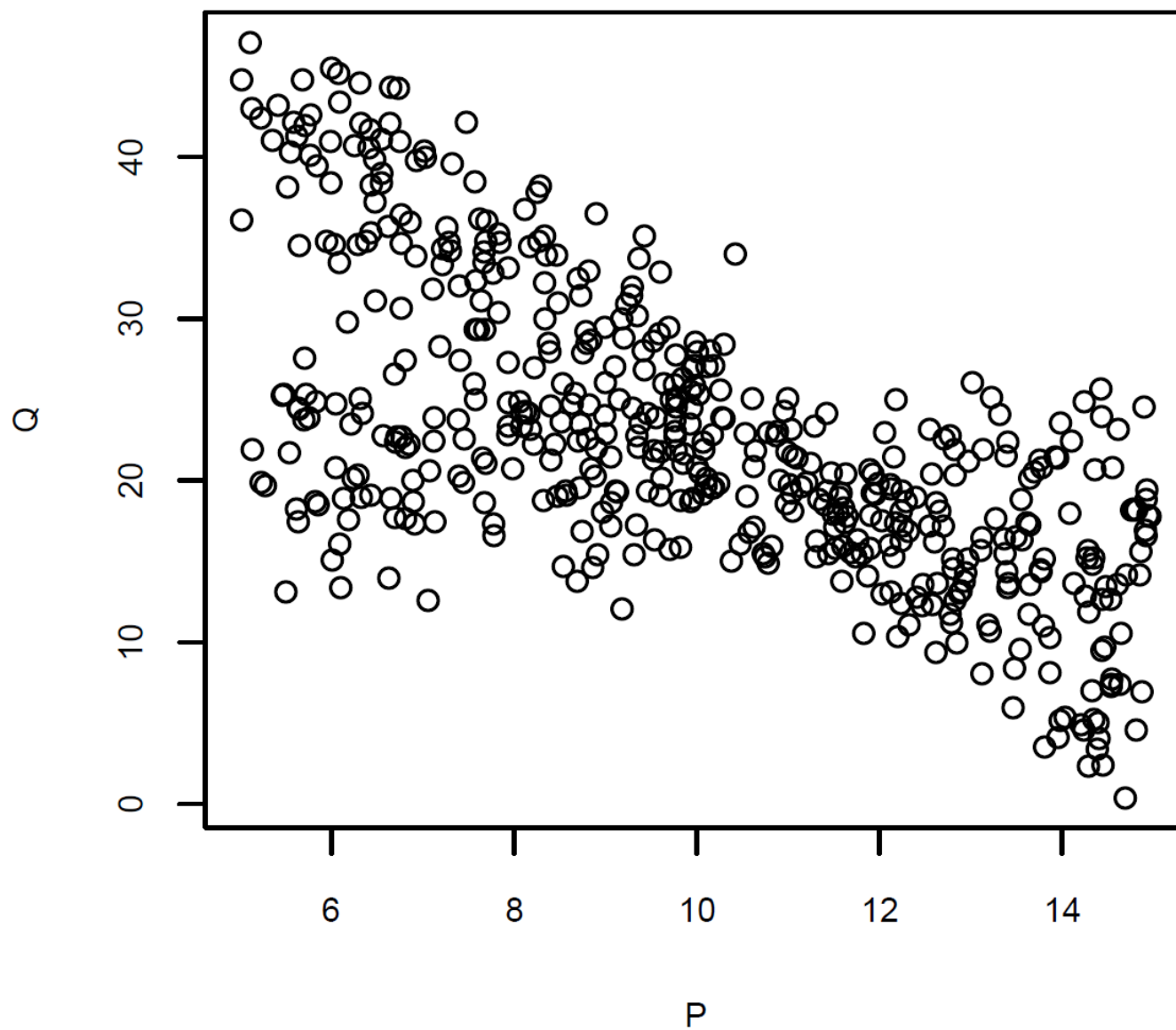
- A type of non-linear effect
- Allows for different effects for different groups (when using a dummy)

A hypothetical data set – demand for marijuana

Suppose that 500 marijuana users are surveyed in different locations, and the variables in the data are:

- Q - the quantity of marijuana consumed, in grams, per month
- P - the average price per gram in the individual's location
- $adult = 1$ if the individual is an adult, $= 0$ if the individual is a teenager

Figure 8.1: Plot of the hypothetical demand for marijuana data.



- Notice anything?
- Ignore the *adult* dummy variable, estimate a regression

```
summary(lm(Q ~ P))
```

```

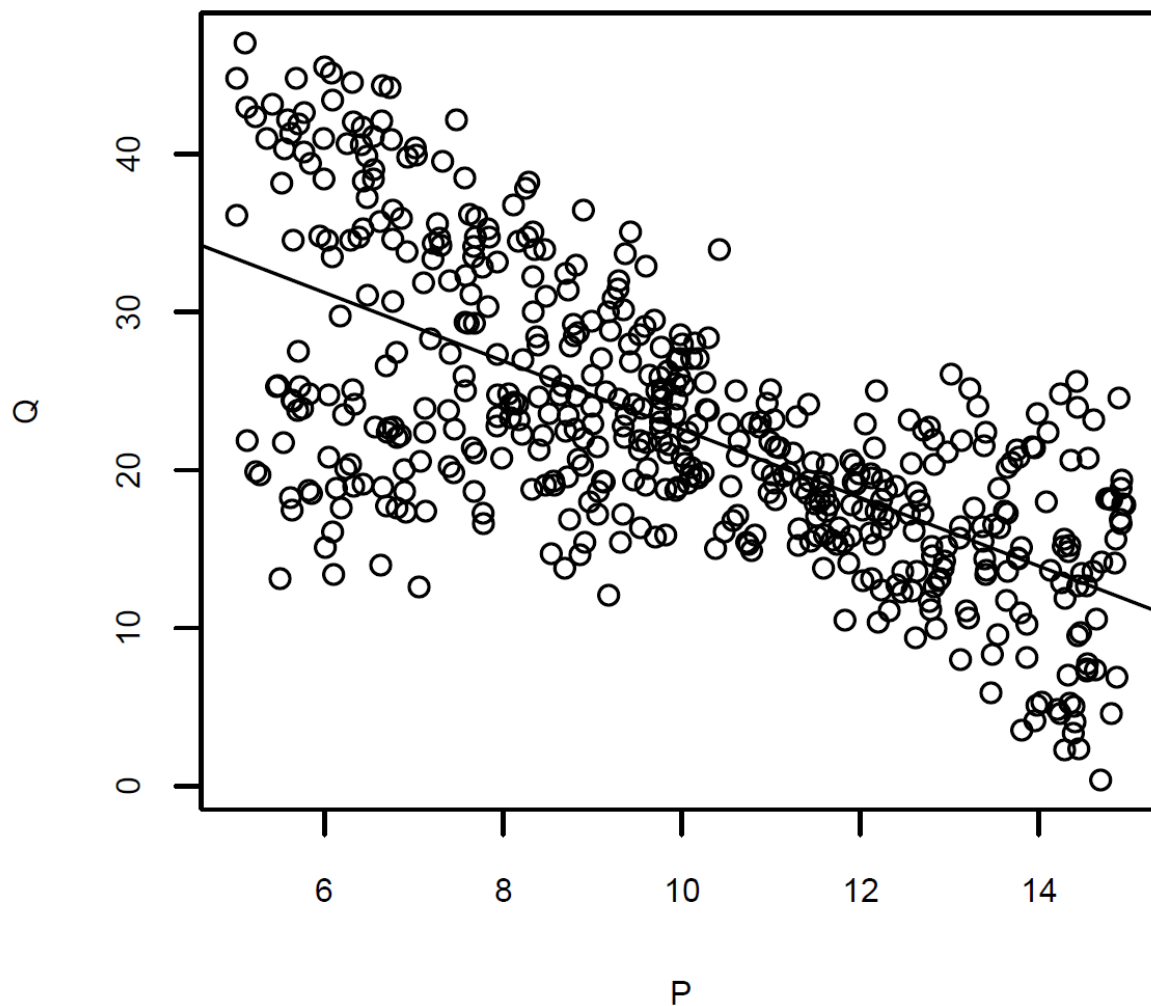
Coefficients:
              Estimate Std. Error t value Pr(>|t|)
(Intercept)   44.2152     1.0776   41.03  <2e-16 ***
P             -2.1634     0.1041  -20.78  <2e-16 ***

```

Increase in price of \$1 leads to decrease in consumption of 2.16 grams/month.

Add the line:

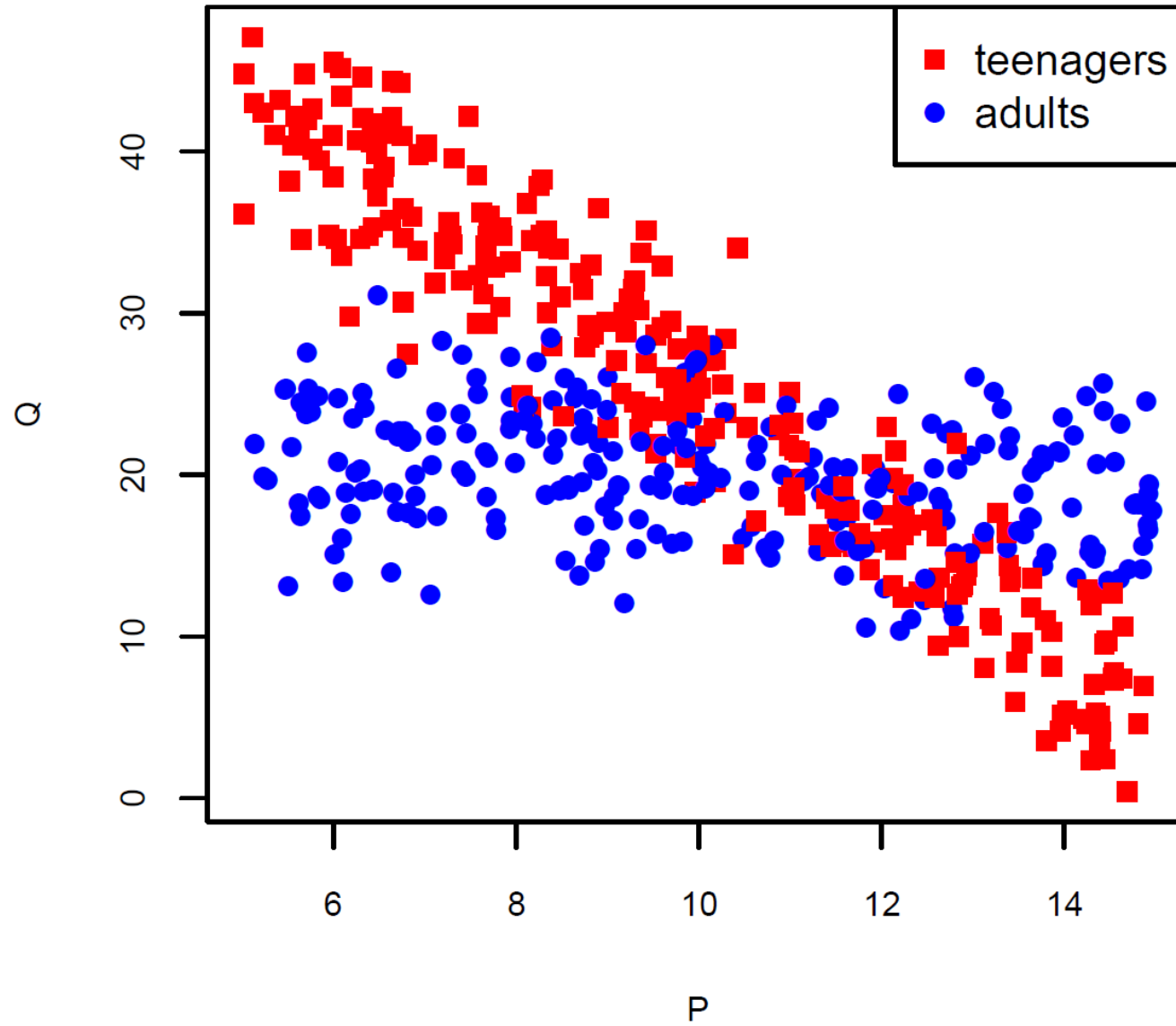
Figure 8.2: Marijuana data, with estimated regression line from $Q = \beta_0 + \beta_1 P + \epsilon$ added to the plot.



- We're getting an “average” regression line for the two groups
- Ideally, we would like a separate regression slope for each
- Why might the slope (marginal effect) be different between groups

Plot the data by group (teenagers and adults):

Figure 8.3: Marijuana data plotted by age group.



Let's add the dummy variable to the regression:

```
summary(lm(Q ~ P + adult))
```

Coefficients:

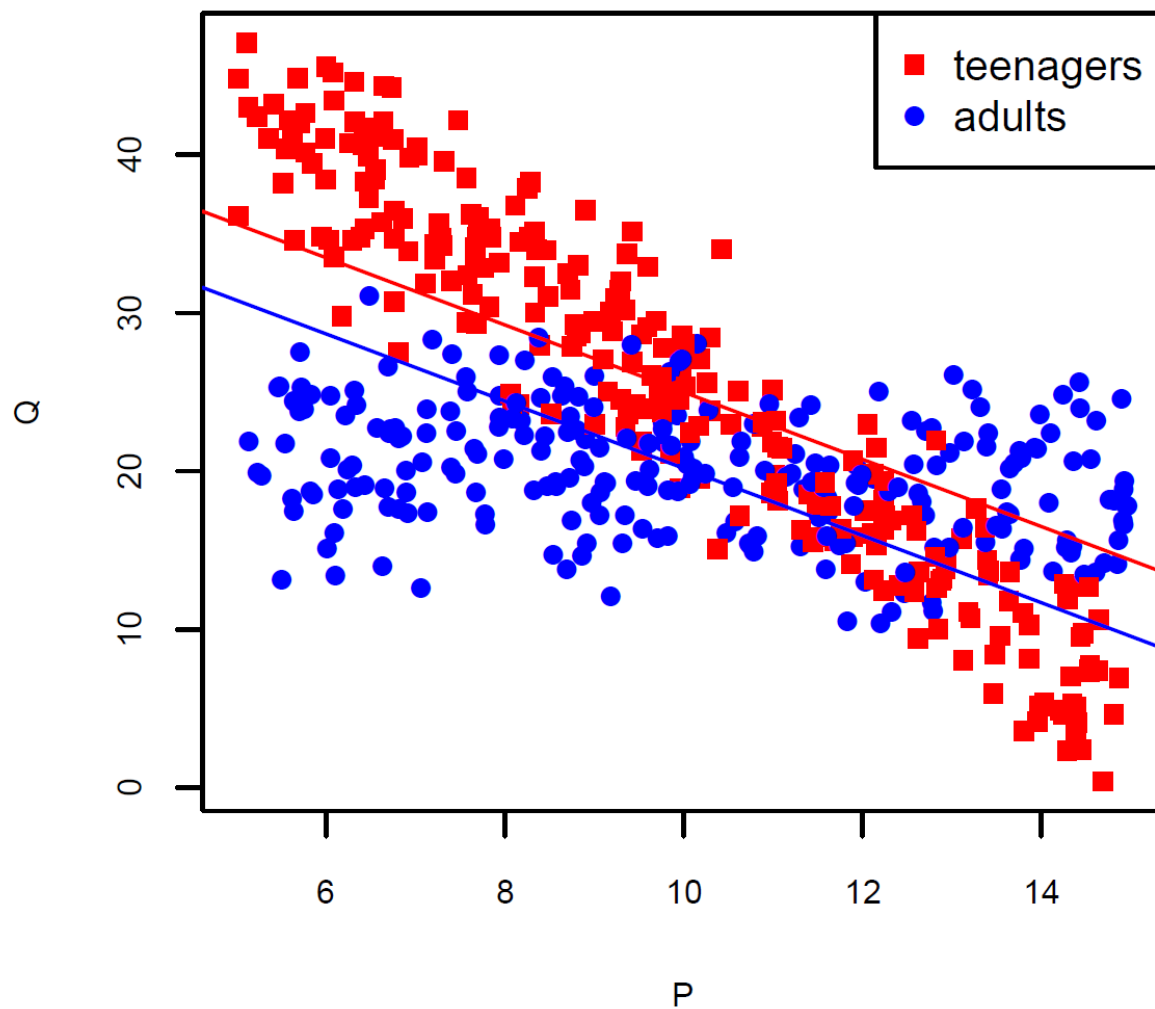
	Estimate	Std. Error	t value	Pr(> t)	
(Intercept)	46.21319	1.02971	44.880	<2e-16	***
P	-2.12242	0.09712	-21.854	<2e-16	***
adult	-4.81124	0.54975	-8.752	<2e-16	***

Interpretation?

- Adults consume 4.81 g less
- Slope?

Does the dummy variable do the trick? See the regression lines plotted:

Figure 8.4: With the addition of the dummy variable, each group has a different intercept, but the same slope.



Two separate regression lines, but only the intercepts differ (slope the same). In order to get what we want, we need an *interaction term*. In this case, it will be a *dummy-continuous* interaction.

Ideally, we want to allow the effect of P on Q to be different for adults and teenagers. How to do this?

Estimate the population model:

$$Q = \beta_0 + \beta_1 P + \beta_2 adult + \beta_3 (adult \times P) + \epsilon \quad (8.2)$$

where $adult \times P$ is the interaction term, and is a new variable that is created by multiplying the other two variables together. To see how model 8.2 allows for two separate lines, consider what the population model is for teenagers ($adult = 0$), and for adults ($adult = 1$).

Population model for teenagers

Let's substitute in the value $adult = 0$ into equation 8.2 and get the population model for teenagers:

$$\begin{aligned} Q &= \beta_0 + \beta_1 P + \beta_2(0) + \beta_3(0 \times P) + \epsilon \\ &= \beta_0 + \beta_1 P + \epsilon \end{aligned} \tag{8.3}$$

From equation 8.3, we can see that the intercept is β_0 and the slope is β_1 .

Population model for adults

Substituting in the value $adult = 1$ into equation 8.2, we get the population model for adults:

$$\begin{aligned} Q &= \beta_0 + \beta_1 P + \beta_2(1) + \beta_3(1 \times P) + \epsilon \\ &= (\beta_0 + \beta_2) + (\beta_1 + \beta_3)P + \epsilon \end{aligned} \tag{8.4}$$

For adults, the intercept is $\beta_0 + \beta_2$ and the slope is $\beta_1 + \beta_3$. The marginal effect of price on consumption differs by β_3 between the two groups.

Estimation with an interaction term

To include a dummy-continuous interaction term in our regression, we simply create a new variable by multiplying the dummy variable (*adult*) and the continuous variable *P* together:

```
adult_P <- adult*P
```

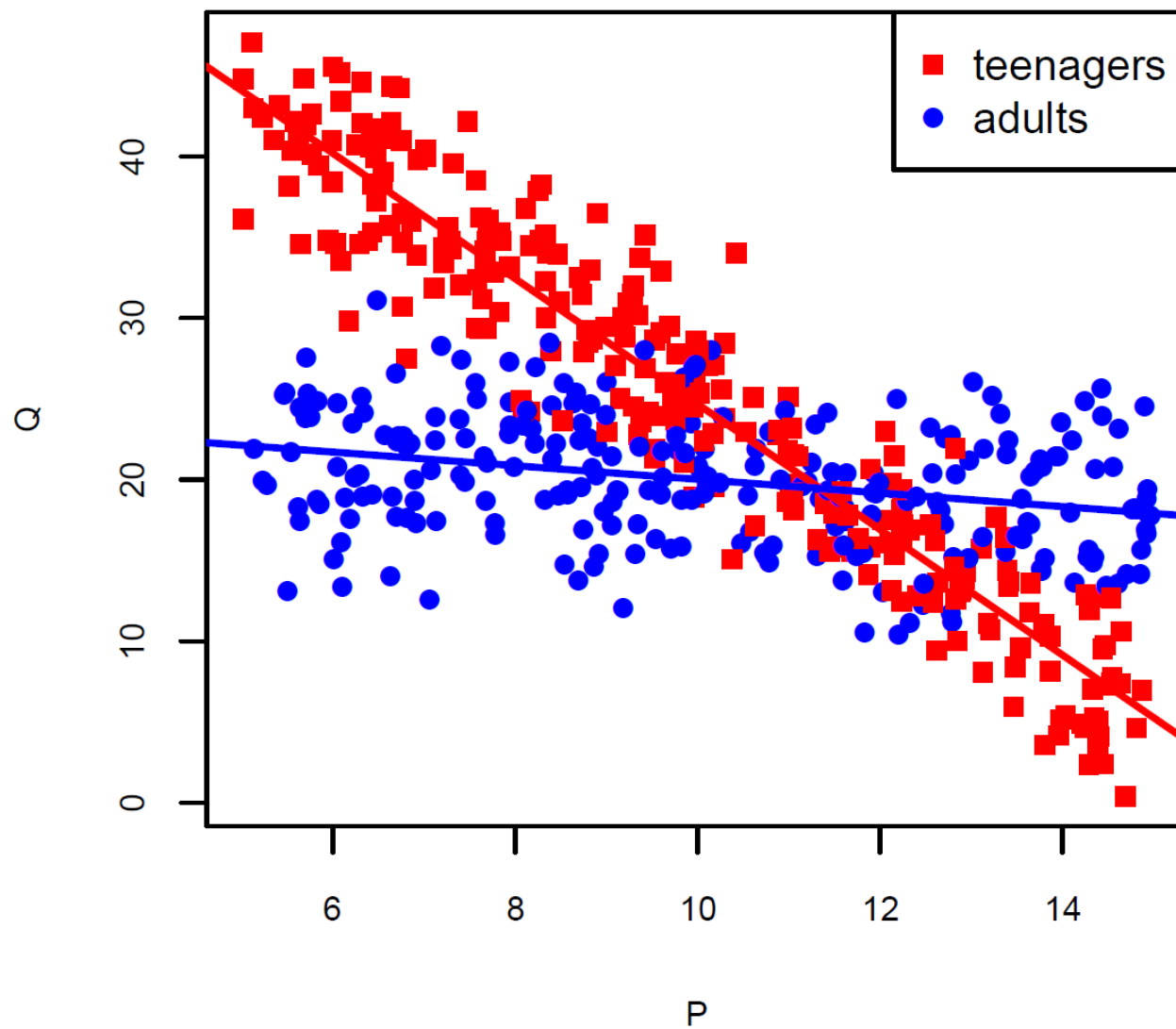
and include the new variable in the regression:

```
summary(lm(Q ~ P + adult + adult_P))
```

```
Coefficients:
              Estimate Std. Error t value Pr(>|t|)
(Intercept)   63.48944    0.85166   74.55  <2e-16 ***
P             -3.88168    0.08339  -46.55  <2e-16 ***
adult         -39.25222    1.21030  -32.43  <2e-16 ***
adult_P        3.45993    0.11695   29.58  <2e-16 ***
```

The estimated value of 3.46 (on the `adult_P` dummy-continuous interaction term) means that the decrease in consumption due to an increase in price of \$1 is 3.46 grams/month less for adults than it is for teenagers. That is, the effect of price on quantity is -3.88 for teenagers, and $(-3.88 + 3.46 = -0.42)$ for adults. The demand curve is much steeper for teenagers.

Figure 8.5: Two separate regression lines for the two different groups.



8.4.4 Hypothesis tests involving dummy interactions

An important use of dummy interaction terms is to test whether there is a different effect between two groups. In the marijuana example, the interaction term measures the difference in the slope of the demand curve between the two groups. To test the hypothesis that the sensitivity of marijuana consumption to changes in price is the same for teenagers as it is for adults, we could test the hypothesis:

$$H_0 : \beta_3 = 0$$

$$H_A : \beta_3 \neq 0$$

in the model:

$$Q = \beta_0 + \beta_1 P + \beta_2 adult + \beta_3 (adult \times P) + \epsilon$$

Differences-in-differences (DiD) [Not yet in textbook!]

Dummy-dummy interactions can be used for something called “Differences-in-differences” (DiD) estimation.

Example: increasing the minimum wage (image by Stable Diffusion)



- In 1992, New Jersey's minimum wage rose from \$4.25 to \$5.05 per hour.
- Card and Krueger (1994) surveyed 410 fast-food restaurants before and after the increase, and asked about things like the number of employees.

Download Card and Krueger data:

```
did <- read.csv("https://rtgodwin.com/data/card.csv")
```

Some variables to look at for now:

EMP – number of full-time employees

TIME – a dummy equal to 0 for before the wage increase, 1 for after the increase

STATE – a dummy equal to 0 for Pennsylvania, equal to 1 for New Jersey

Difference in the number of employees before and after the wage increase:

```
mean(did$EMP[did$STATE == 1 & did$TIME == 1]) -  
  mean(did$EMP[did$STATE == 1 & did$TIME == 0])  
[1] 0.4666667
```

The difference is not significant:

```
dids <- subset(did, STATE==1)
summary(lm(EMP ~ TIME, data=dids))
```

Coefficients:

	Estimate	Std. Error	t value	Pr(> t)
(Intercept)	20.4306	0.5289	38.627	<2e-16 ***
TIME	0.4667	0.7480	0.624	0.533

Signif. codes: 0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1

Residual standard error: 9.298 on 616 degrees of freedom

Multiple R-squared: 0.0006315, Adjusted R-squared: -0.0009909

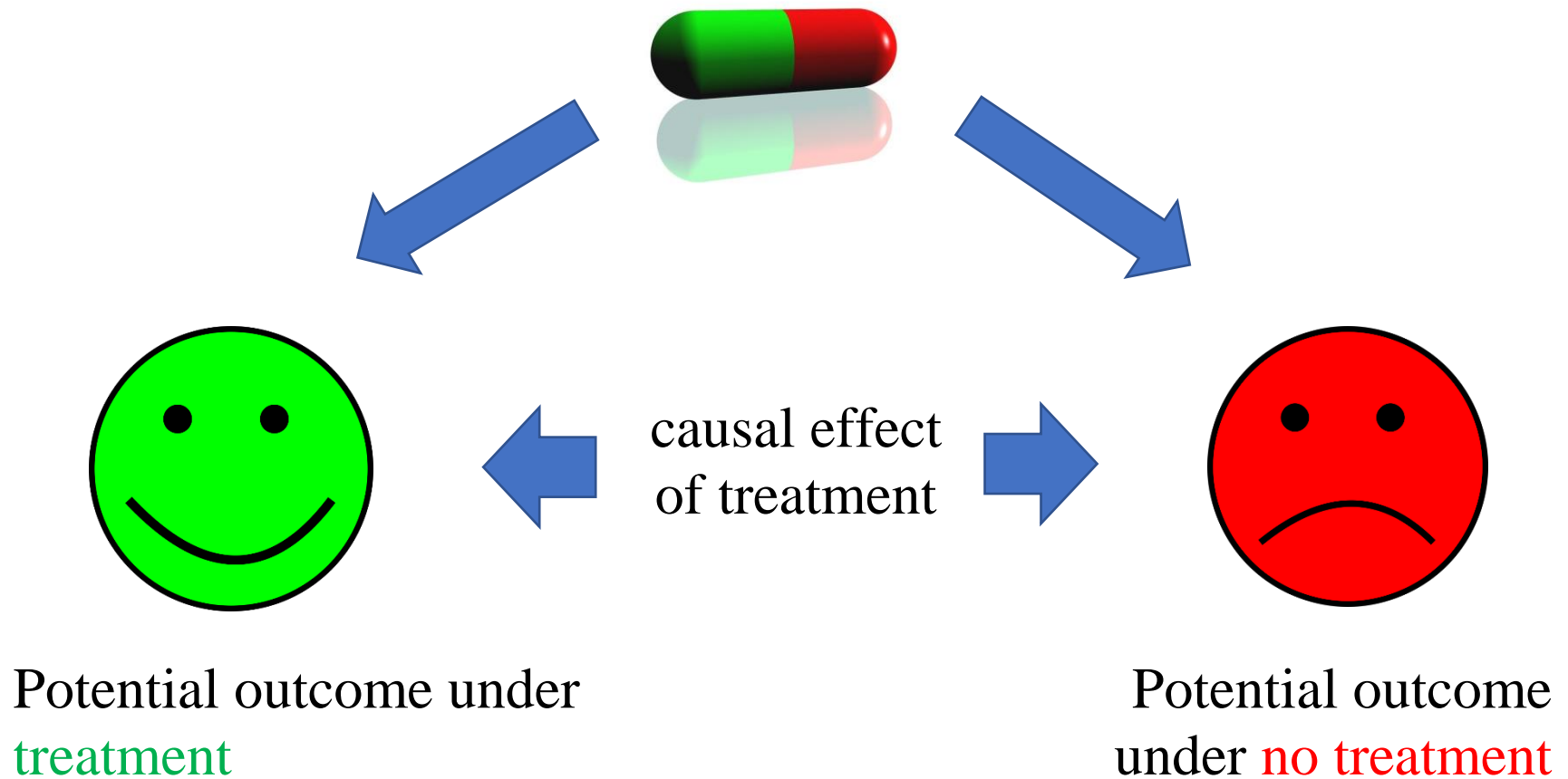
F-statistic: 0.3892 on 1 and 616 DF, p-value: 0.5329

So, the causal effect of the increase in minimum wage on employment is estimated to be an increase of 0.47 workers on average, but this increase is not statistically significant.

What is the problem with calling this a “causal effect”?

Next: “The Fundamental Problem of Causal Inference”

Fundamental problem of causal inference



Suppose we want to know the *difference* that a cause (treatment) makes.

That is, we want to know:

$$E[y_1 - y_0]$$

- y_1 – outcome with treatment
- y_0 – outcome without treatment

Treatment is broadly defined:

- Treatment with a drug - (y_1 and y_0 blood pressure with/without the drug)
- Addictions treatment (methadone) – (y_1 and y_0 probability of success)
- Health insurance - (y_1 and y_0 the number of visits to the doctor with or without insurance)
- Education (y_1 and y_0 the wage with/without an education)
- Job training
- Monetary policy
- Student debt
- Information
- **Increase in minimum wage** (y_1 and y_0 the employment rate)

Fundamental problem of causal inference

Because an “individual” can’t be in both states (treated and untreated), we can’t observe both y_1 and y_0 .

We can never observe a causal effect!

- One of the two outcomes will occur, and is factual.
- The other outcome(s) is imagined, or counterfactual.
- We only ever observe either y_1 **or** y_0 .

Maybe we could observe a causal effect?

Wooldridge calls it a problem of “missing data”.

How could we observe the missing data?

- Time travel
- Parallel universe

Barring the above, we have to think in *counterfactuals* and try to find ways to estimate what the unobserved outcome (y_1 or y_0) would have looked like so that we can calculate $y_1 - y_0$.

Estimation of a causal effect

Unit	Treated:	Outcome under treatment y_1	Outcome under no treatment y_0
1	yes	✓	?
2	yes	✓	?
3	no	?	✓
4	no	?	✓



causal effect estimate

Back to minimum wage example

EMP (y)	number of full-time employees
TIME	0 for before the wage increase
	1 for after the increase
STATE	0 for Pennsylvania (no wage increase – “control”)
	1 for New Jersey (wage increase – “treatment”)

The naïve approach is to take the difference between New Jersey’s employment before and after the wage increase:

$$\bar{y}_{at\ TIME=1} - \bar{y}_{at\ TIME=0} = 0.4667$$

But for this to be the causal effect, need to assume that the level of employment would have stayed constant over the 6 months!

```
dids <- subset(did, STATE==1)
summary(lm(EMP ~ TIME, data=dids))
```

Coefficients:

	Estimate	Std. Error	t value	Pr(> t)
(Intercept)	20.4306	0.5289	38.627	<2e-16 ***
TIME	0.4667	0.7480	0.624	0.533

Signif. codes: 0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1

Residual standard error: 9.298 on 616 degrees of freedom

Multiple R-squared: 0.0006315, Adjusted R-squared: -0.0009909

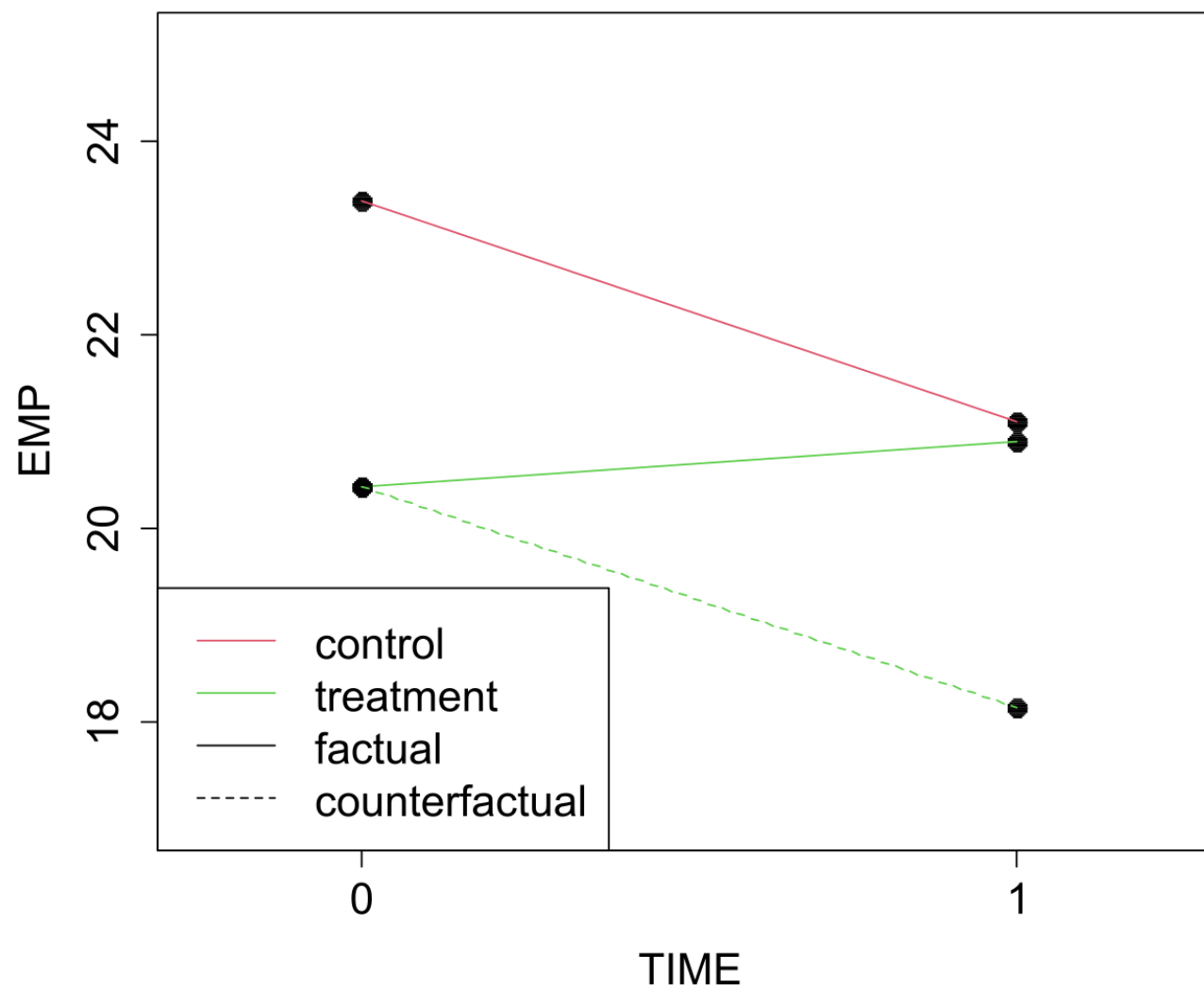
F-statistic: 0.3892 on 1 and 616 DF, p-value: 0.5329

Table 1: Average employment by STATE and TIME

	TIME = 0	TIME = 1	Difference
New Jersey STATE = 1 (treatment)	20.431	20.897	0.466
Pennsylvania STATE = 0 (control)	23.380	21.096	-2.283
Difference	-2.949	-0.199	2.750

- Parallel trends assumption: the *difference* in employment that occurred for the control group would have also occurred for the treatment group (if they hadn't have been treated): -2.283
- The *difference* in employment that actually did occur under treatment was 0.466
- The *difference-in-difference* is $0.466 - (-2.283) = 2.750$

Average number of employees before and after wage increase, by state



We can get the DiD estimator by differencing the sample means between groups. But often, we want to include other “X” variables in the model in order to avoid OVB. If we estimate the model:

$$EMP = \beta_0 + \beta_1 TIME + \beta_2 STATE + \beta_3 (TIME \times STATE) + \epsilon$$

Then b_3 is the DiD estimator!

- Other “X” variables can be added to the model
- $TIME \times STATE$ is an **interaction term**
- β_1 is the effect of $TIME$ for the control group
- β_2 is the difference in EMP at $TIME = 0$
- β_3 is the difference in the effect of $TIME$ between the two groups

$$EMP = \beta_0 + \beta_1 TIME + \beta_2 STATE + \beta_3 (TIME \times STATE) + \epsilon$$

Plug in values for the dummies to get the interpretation of the β :

<i>TIME</i>	<i>STATE</i>	<i>EMP</i>	difference
0	0	β_0	β_1 (for control)
1	0	$\beta_0 + \beta_1$	
0	1	$\beta_0 + \beta_2$	$\beta_1 + \beta_3$ (for treatment)
1	1	$\beta_0 + \beta_1 + \beta_2 + \beta_3$	

Difference over time for **control**: β_1

Difference over time for **treatment**: $\beta_1 + \beta_3$

Difference-in-difference: $(\beta_1 + \beta_3) - \beta_1 = \beta_3$


```
summary(lm(EMP ~ TIME + STATE + I(TIME * STATE), data = did))
```

Coefficients:

	Estimate	Std. Error	t value	Pr(> t)
(Intercept)	23.380	1.098	21.288	<2e-16 ***
TIME	-2.283	1.553	-1.470	0.1419
STATE	-2.949	1.224	-2.409	0.0162 *
I(TIME * STATE)	2.750	1.731	1.588	0.1126

Signif. codes: 0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1

Residual standard error: 9.511 on 764 degrees of freedom

Multiple R-squared: 0.007587, Adjusted R-squared: 0.00369

F-statistic: 1.947 on 3 and 764 DF, p-value: 0.1206

Average number of employees before and after wage increase, by state

