

## 8.4 – Interaction terms

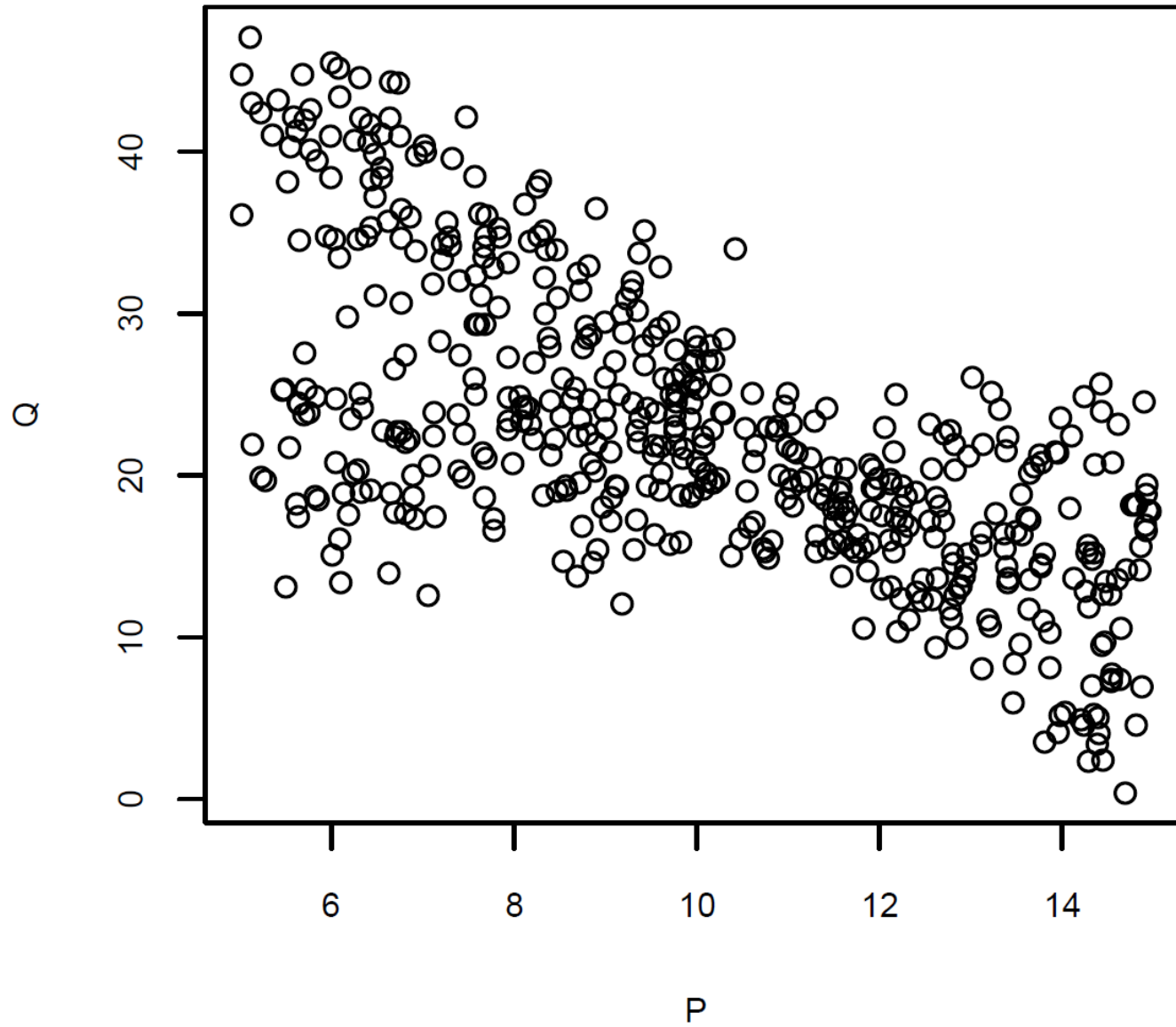
- A type of non-linear effect
- Allows for different effects for different groups (when using a dummy)

## 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))
```

```

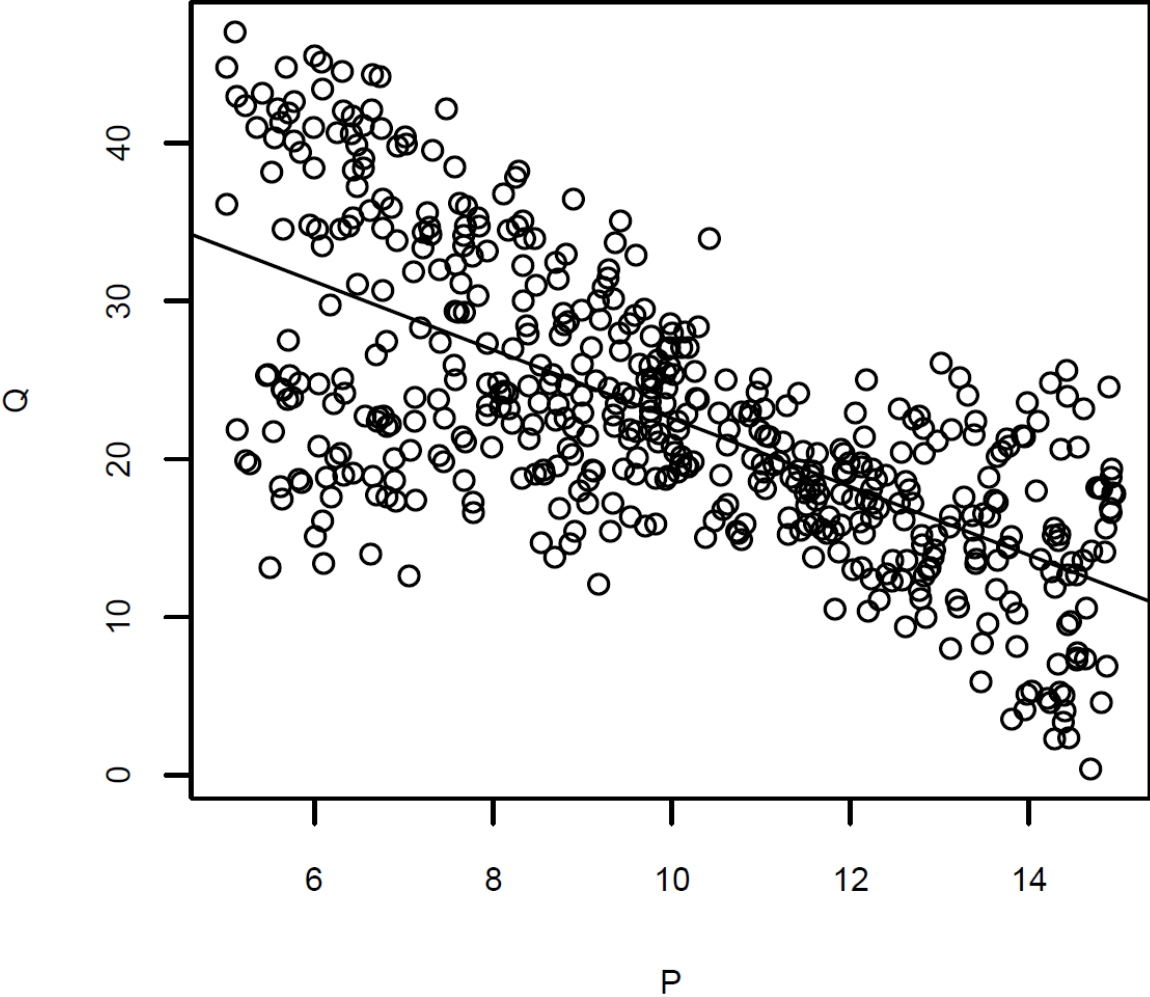
Coefficients:
              Estimate Std. Error t value Pr(>|t|)
(Intercept)  44.2152     1.0776   41.03  <2e-16 ***
P            -2.1634     0.1041  -20.78  <2e-16 ***

```

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

Add the line:

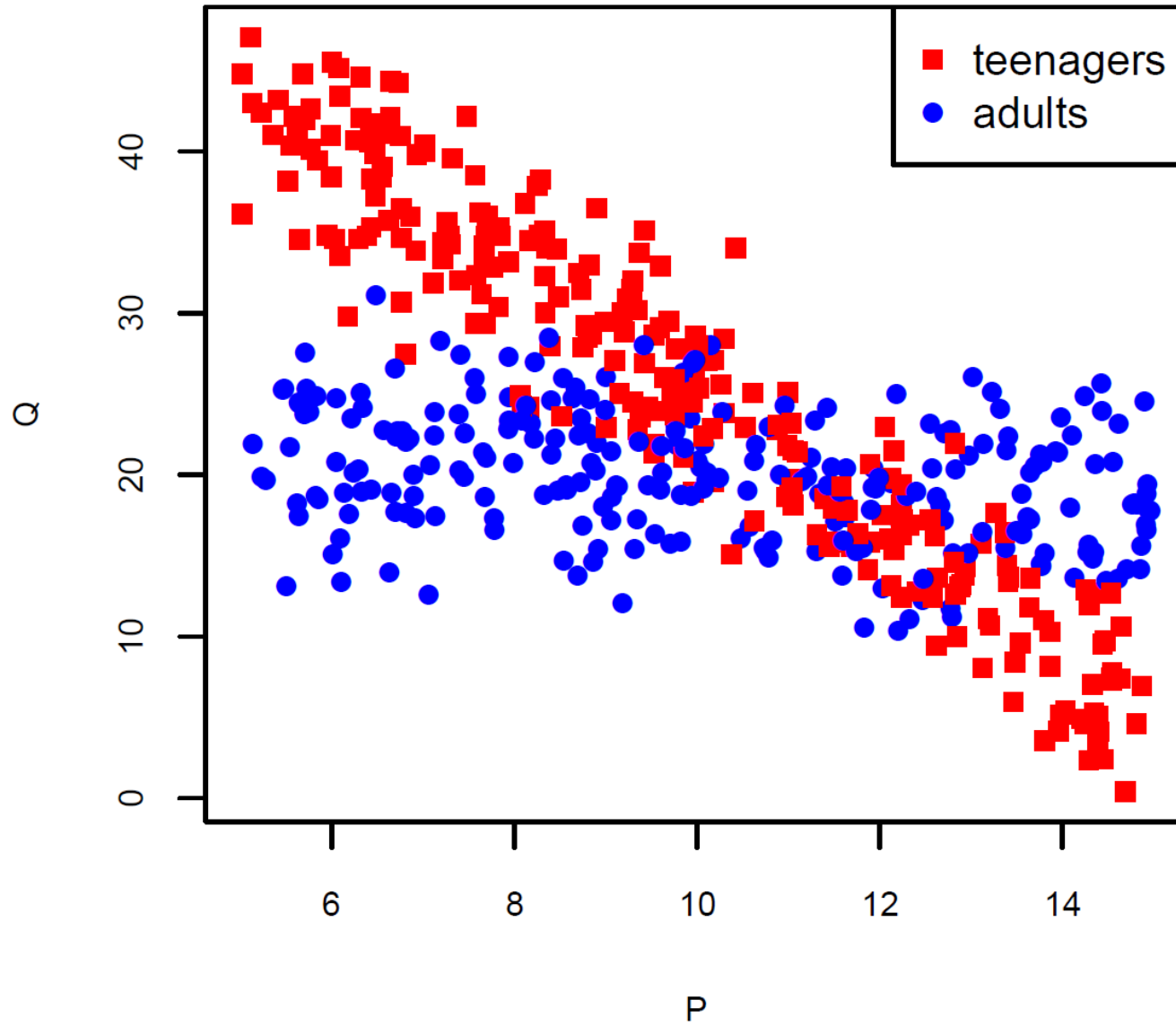
Figure 8.2: Marijuana data, with estimated regression line from  $Q = \beta_0 + \beta_1 P + \epsilon$  added to the plot.



- 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 ***
```

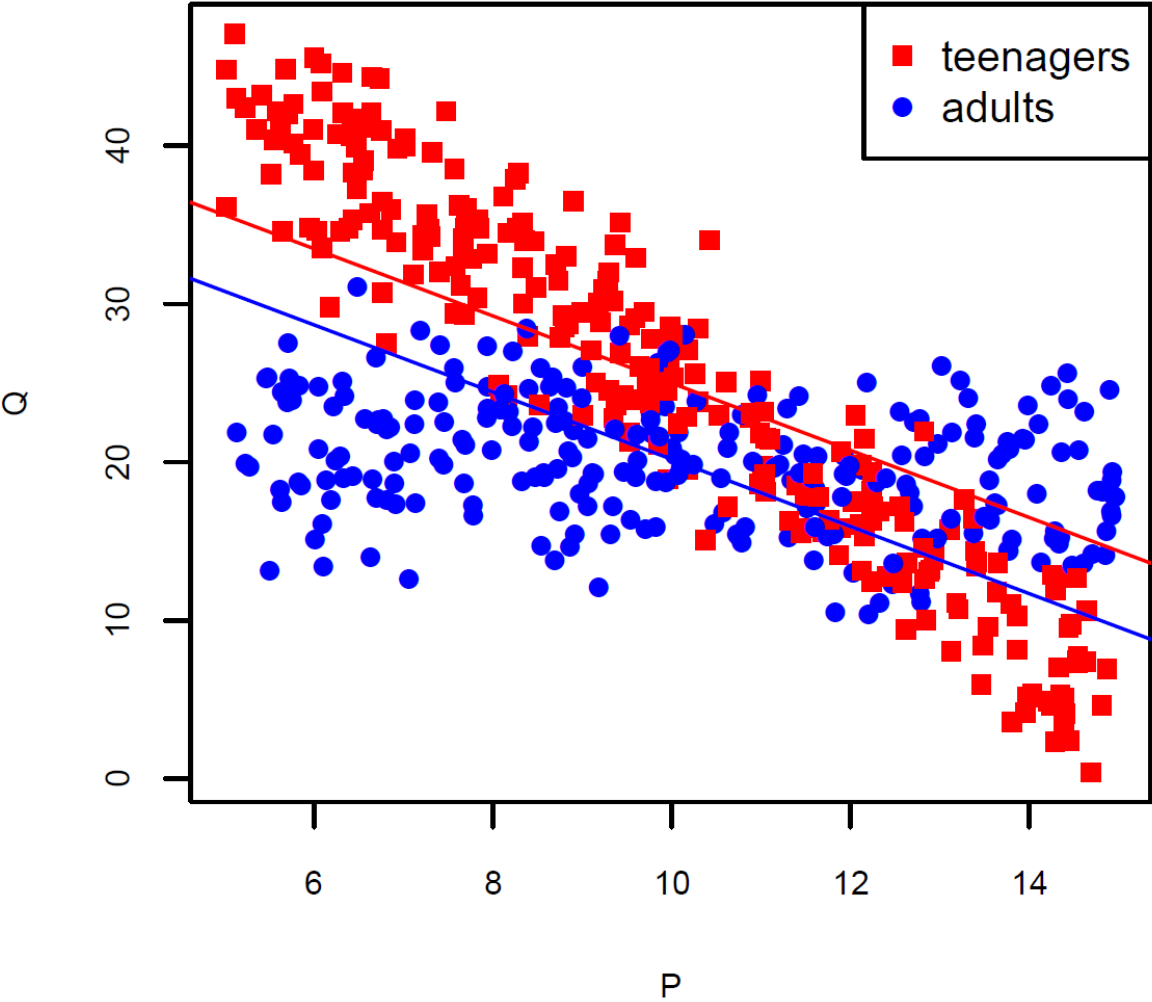
Interpretation?

- 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  $\beta_0$  and the slope is  $\beta_1$ .

## 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  $\beta_3$  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:

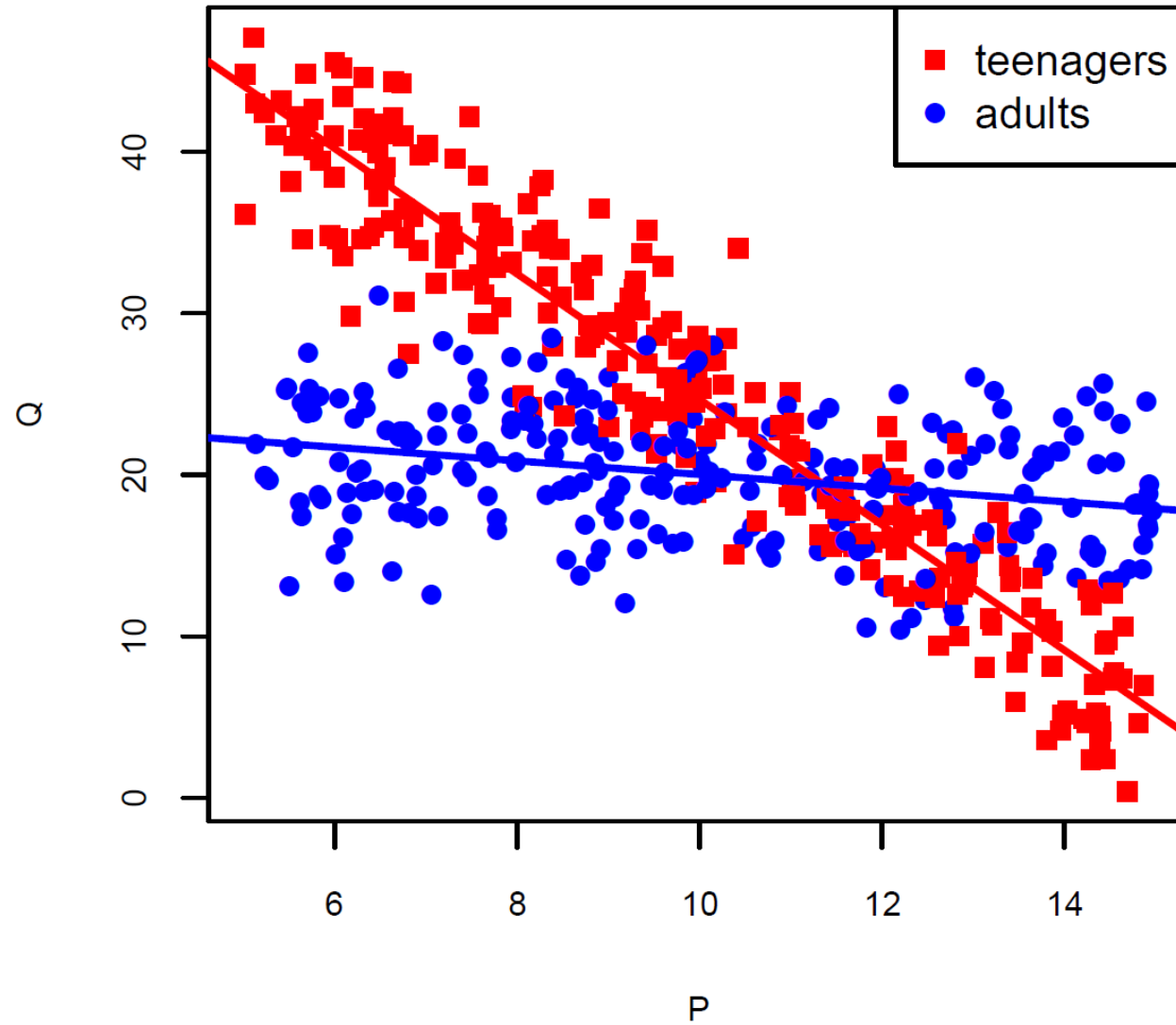
```
summary(lm(Q ~ P + adult + adult_P))
```

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.



## Dummy-dummy interaction: differences-in-differences

- CPS data again
- *bach* = 1 if individual has a university degree, = 0 otherwise
- start with basic model:

```
summary(lm(log(wage) ~ female + bach))
```

Coefficients:

	Estimate	Std. Error	t value	Pr(> t )	
(Intercept)	2.07175	0.03108	66.657	< 2e-16	***
female	-0.22886	0.04240	-5.397	1.02e-07	***
bach	0.39177	0.04976	7.873	1.97e-14	***

Interpretation of results?

- There is a different wage for men and women
- There is a different wage for *bach* and no *bach*
- There is no different effect of *bach* for women vs. men

We might want to allow for the effect of a degree on *wage* to be different for men and women (a difference-in-difference).

We could estimate the model:

$$\log(wage) = \beta_0 + \beta_1 female + \beta_2 bach + \beta_3 (female \times bach) + \epsilon$$

where  $\beta_3$  is the additional percentage increase in wages for women with an education, versus men with an education. In R, we create the dummy-dummy interaction term by:

```
fem_bach <- female*bach
```

and include it in our regression:

```
summary(lm(log(wage) ~ female + bach + fem_bach))
```



Coefficients :					
	Estimate	Std. Error	t value	Pr(> t )	
(Intercept)	2.08291	0.03292	63.280	< 2e-16	***
female	-0.25309	0.04849	-5.219	2.58e-07	***
bach	0.34500	0.06736	5.122	4.25e-07	***
fem_bach	0.10292	0.09994	1.030	0.304	

## Interpretation?

- $b_1 \rightarrow$  Women without a degree make 25% less than men without a degree
- $b_2 \rightarrow$  Men with a degree make 35% more than men without a degree
- $b_2 + b_3 \rightarrow$  Women with a degree make 45% more than women without a degree
- $b_3 \rightarrow$  The “difference-in-difference”. The effect of a degree on *wage* is 10% more for women than for men

#### 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:

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

Similarly, testing  $\beta_3 = 0$  in the model:

$$\log(wage) = \beta_0 + \beta_1 female + \beta_2 bach + \beta_3 (female \times bach) + \epsilon$$

is a test of whether there is a different effect of education for women than for men. From the regression output in the previous section, we see that the  $p$ -value for the estimated coefficient on `fem_bach` is 0.304. We fail to reject the null that there is no difference in the effect of education between men and women.

## Extra stuff

- Continuous-continuous interactions
- Add other variables
- Polynomial model – multiple interactions – need  $F$ -test