

Warm up 2:00-2:10

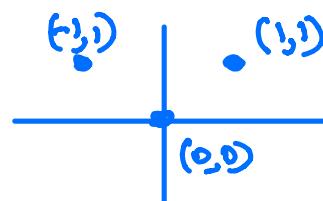
Exercise 11.6, 11

11. Let  $X$  have the uniform distribution on the three points  $-1, 0$ , and  $1$ . Let  $Y = X^2$ .

$$r = \frac{E(D_x D_y)}{\sigma_x \sigma_y}$$

a) Show that  $X$  and  $Y$  are uncorrelated.

Picture



we show that

$$r = 0$$

$$\text{i.e. } E(D_x D_y) = 0.$$

$$X \sim \text{Unif}(-1, 0, 1) \Rightarrow \mu_X = E(X) = \frac{-1 + 0 + 1}{3} = 0$$

$$\mu_Y = E(Y) = E(X^2) = \frac{1}{3}((-1)^2 + (0)^2 + (1)^2) = \frac{2}{3}$$

$$E(D_x D_y) = E((X - \mu_X)(Y - \mu_Y)) = E(X(Y - \frac{2}{3}))$$

$$= \frac{1}{3}(-1(1 - \frac{2}{3}) + 0(0 - \frac{2}{3}) + 1(1 - \frac{2}{3}))$$

$$= \frac{1}{3}(-\frac{1}{3} + 0 + \frac{1}{3}) = 0$$

$$\Rightarrow r = \frac{E(D_x D_y)}{\sigma_x \sigma_y} = \frac{0}{\sigma_x \sigma_y} = 0 \Rightarrow X, Y \text{ are uncorrelated}$$

$Y = X^2$  so  $Y$  is dependent on  $X$

So uncorrelated  $\not\Rightarrow$  indep.

But if you scatter plot is football shaped

then uncorrelated  $\Leftrightarrow$  independent.

## Last time

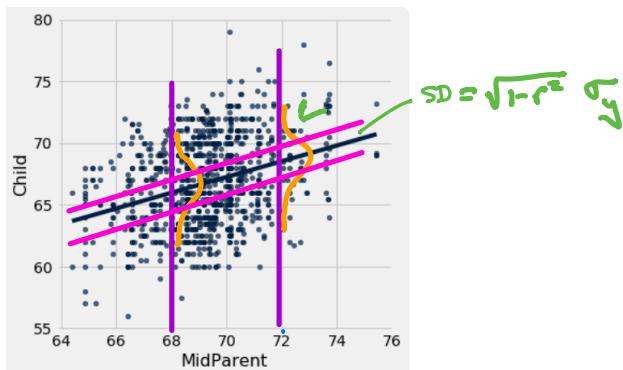
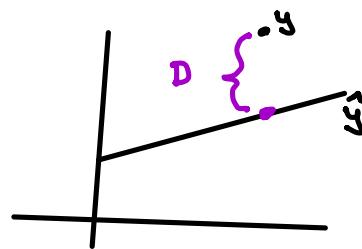
### Sec 11.5 The error in regression

let  $D = \hat{y} - y$  be the residual error

$$D \sim N(0, (1-r^2)\sigma_y^2)$$

$E(D)$        $\text{Var}(D)$

$\frac{2}{3}$  of your data are  
within two lines parallel  
to  $\hat{y}$   $\pm 1 SD = \sqrt{1-r^2} \sigma_y$  away.



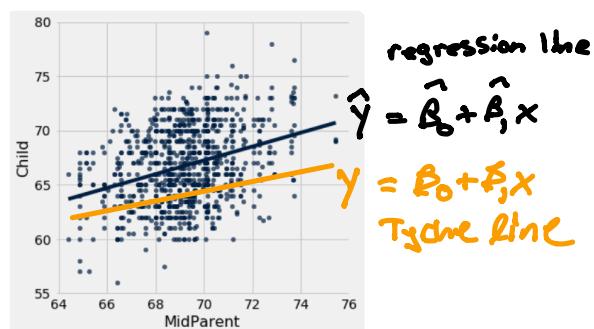
### Sec 12.1 the simple regression model

$x$  = Predictor variable  
 $y$  = response variable  
by  $y$  (random)

we assume for each of  $n$  observations

$$Y_i = \beta_0 + \beta_1 x_i + \varepsilon_i$$

- $x_1, x_2, \dots, x_n$  are fixed
- $\varepsilon_1, \varepsilon_2, \dots, \varepsilon_n \stackrel{\text{iid}}{\sim} N(0, \sigma^2)$
- only Tyche knows  $\beta_0, \beta_1$ , and  $\sigma^2$
- $y = \beta_0 + \beta_1 x$  is called the Tyche line,



Tyche line is the signal part of the response

The regression line estimates the Tyche line.

Today

- Sec 12.1 The simple linear regression model
- Sec 12.2 The distribution of the estimated slope.

## ① Sec 12.1 The simple linear regression model

Is the response variable dependent on the predictor variable(s)? In other words is the fitted line flat? This is the question we answer in chapter 12.

### Individual responses

Fix  $\beta_0, \beta_1, x_i$ :

$$y_i = \underbrace{\beta_0 + \beta_1 x_i}_{\substack{\text{response} \\ (\text{random})}} + \underbrace{\varepsilon_i}_{\substack{\text{signal} \\ (\text{constant})}} + \underbrace{\varepsilon_i}_{\substack{\text{error} \\ (\text{random}) N(0, \sigma^2)}}$$

$y_1, \dots, y_n$  are independent since  $\varepsilon_1, \dots, \varepsilon_n$  are independent.

### Average responses

Let  $\bar{y} = \frac{1}{n} \sum_{i=1}^n y_i$  be average response at index normals  
is normal

Show  $E(\bar{y}) = \beta_0 + \beta_1 \bar{x}$ ?  $y_i = \beta_0 + \beta_1 x_i + \varepsilon_i$

$$\begin{aligned} E(\bar{y}) &= E\left(\frac{1}{n} \sum_{i=1}^n \hat{y}_i\right) = \frac{1}{n} \sum_{i=1}^n E(\beta_0 + \beta_1 x_i + \varepsilon_i) \\ &= \frac{1}{n} \sum_{i=1}^n (\beta_0 + \beta_1 x_i + E(\varepsilon_i)) \\ &= \beta_0 + \beta_1 \bar{x} \quad \checkmark \end{aligned}$$

Furthermore,

$$Var(\bar{Y}) = \frac{1}{n^2} \sum_{i=1}^n Var(Y_i) = \frac{1}{n^2} \cdot n\sigma^2 = \frac{\sigma^2}{n}$$

so  $\bar{Y} \sim N(\beta_0 + \beta_1 \bar{x}, \frac{\sigma^2}{n})$

② Sec 12.2 the distribution of the estimated slope.

To make a hypothesis about whether  $\beta_1 = 0$  we need to know the distribution of  $\hat{\beta}_1$ .

$\hat{\beta}_1$  replaces  $\hat{\alpha}$  from chapter 11.

$$\text{recall } \hat{\alpha} = r \frac{\sigma_y}{\sigma_x} = \frac{E(D_x D_y)}{\sigma_x^2} = \frac{E((x - \mu_x)(y - \mu_y))}{\sigma_x^2}$$

### Estimated Slope

For our sample,

$$\hat{\beta}_1 = \frac{\frac{1}{n} \sum_{i=1}^n (x_i - \bar{x})(Y_i - \bar{Y})}{\frac{1}{n} \sum_{i=1}^n (x_i - \bar{x})^2} \quad \begin{matrix} \mu_x \\ \mu_y \\ \sigma_x^2 \end{matrix}$$

$\hat{\beta}_1$  is a linear combination of independent normals  $Y_1, \dots, Y_n$  and hence  $\hat{\beta}_1$  is normal

### Expectation of the estimated slope

$$\text{we have } E(Y_i) = \beta_0 + \beta_1 x_i + E(\varepsilon_i)$$

$$E(\bar{Y}) = \beta_0 + \beta_1 \bar{x}$$

$$\begin{aligned} \text{so } E(x_i - \bar{x}) &= \beta_0 + \beta_1 x_i - \beta_0 - \beta_1 \bar{x} \\ &= \beta_1 (x_i - \bar{x}) \end{aligned}$$

Show  $E(\hat{\beta}_1) = \beta_1$  ?

$$\begin{aligned} \hat{\beta}_1 &= \frac{\frac{1}{n} \sum_{i=1}^n (x_i - \bar{x})(Y_i - \bar{Y})}{\frac{1}{n} \sum_{i=1}^n (x_i - \bar{x})^2} \\ E(\hat{\beta}_1) &= \frac{\frac{1}{n} \sum_{i=1}^n (x_i - \bar{x}) E(Y_i - \bar{Y})}{\frac{1}{n} \sum_{i=1}^n (x_i - \bar{x})^2} = \beta_1 \frac{\frac{1}{n} \sum (x_i - \bar{x})^2}{\frac{1}{n} \sum (x_i - \bar{x})^2} \\ &= \beta_1 \quad \checkmark \end{aligned}$$

Hence  $\hat{\beta}_1$  is an unbiased estimate of  $\beta_1$ ,  
 we need  $\text{Var}(\hat{\beta}_1)$  to know the distribution  
 and can make a 95% CI for  $\beta_1$ .

Variance of the estimated slope :

FACT  $\text{Var}(\hat{\beta}_1) = \frac{\sigma^2}{n\text{var}(x)} = \frac{\sigma^2}{\sum_{j=1}^n (x_j - \bar{x})^2}$

Note as  $n \rightarrow \infty$ ,  $\sum_{j=1}^n (x_j - \bar{x})^2$  gets very large and  $\text{Var}(\hat{\beta}_1) \rightarrow 0$  so the difference between  $\hat{\beta}_1$  and  $\beta_1$  becomes very small with high probability.

Hence  $\hat{\beta}_1 \sim N(\beta_1, \frac{\sigma^2}{n\text{var}(x)})$

Exercise 12.4.)

- Recall that the intercept of the regression line is given by "the average of  $Y$  minus the slope times the average of  $x$ . That is,  $\hat{\beta}_0 = \bar{Y} - \hat{\beta}_1 \bar{x}$ . Is  $\hat{\beta}_0$  an unbiased estimator of  $\beta_0$ ?

$$E(\hat{\beta}_0) = E(\bar{Y}) - E(\hat{\beta}_1)\bar{x} = \beta_0 + \beta_1 \bar{x} - \beta_1 \bar{x} = \beta_0$$

$$E(\bar{Y}) = \beta_0 + \beta_1 \bar{x}$$

Standard error of the estimated slope :

We have  $SD(\hat{\beta}_1) = \frac{\sigma}{\sqrt{n\text{var}(x)}}$  where  $\sigma$  is SD of the errors,

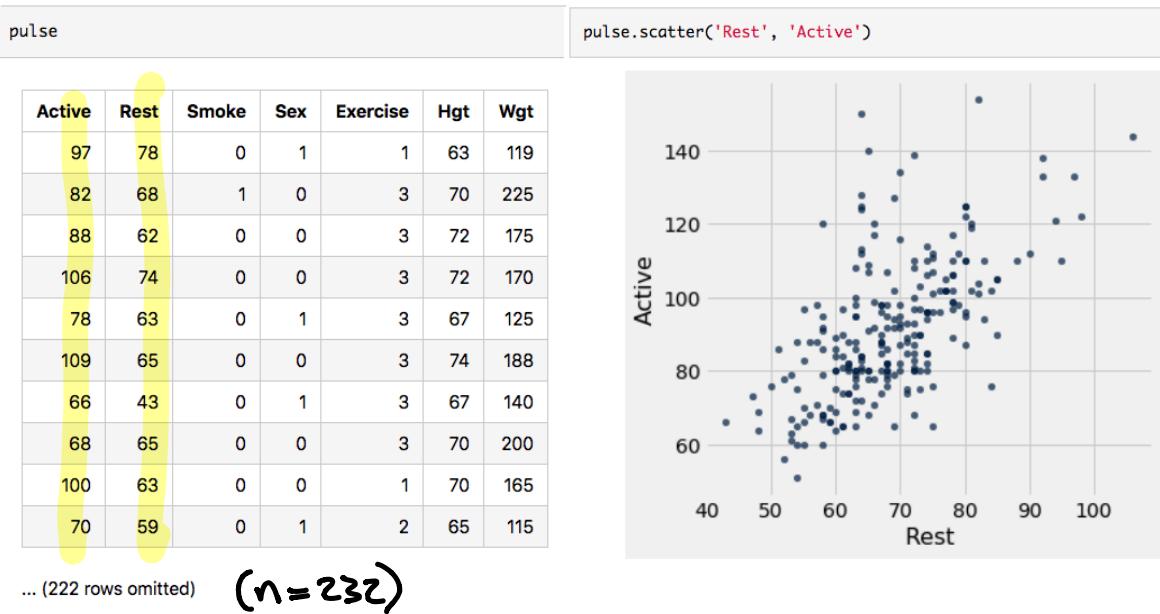
We estimate  $\sigma$  with the SD of the residuals. This is calculated in Python.

When the SD of an estimator is approximated by the data it is called the SE (standard error).

$$T = \frac{\hat{\beta}_1 - \beta_1}{SE(\hat{\beta}_1)} \text{ is approximately normal for large } n.$$

### Ex (Pulse Rates)

We wish to predict active pulse rates from resting pulse rates.



```
active = pulse.column(0)  
resting = pulse.column(1)
```

```
stats.linregress(x=resting, y=active)
```

Output:

$\hat{\beta}_1$	(1.142879681904831,
$\beta_0$	13.182572776013345,
$r$	0.6041870881060092,
P-val	1.7861044071652305e-24,
$SE(\hat{\beta}_1)$	0.09938884436389145)

$$n=232 \rightarrow \text{large} \Rightarrow T = \frac{\hat{\beta}_1 - \beta_1}{SE(\hat{\beta}_1)} \sim N(0, 1)$$

A 95% CI for  $\beta_1$  is

$$\hat{\beta}_1 \pm 2 SE(\hat{\beta}_1) = (0.944, 1.342)$$

this doesn't contain 0 so we reject the null  $H_0: \beta_1 = 0$

$$H_A: \beta_1 \neq 0$$

Under the null

$$T = \frac{\hat{\beta}_1 - 0}{SE(\hat{\beta}_1)} = 11.5$$

