**Section 8.4**
**Inference for Linear Regression**
**Based on content in OpenIntro Stats, 4th Ed**

## Gear up for Inference

▶ Inference in this class has been about this: Make a decision about a parameter based on a test statistic generated from good data.

▶ Inference for linear regression is about this too.

▶ We assume two variables $x$ and $y$ have a linear association plus some noise:
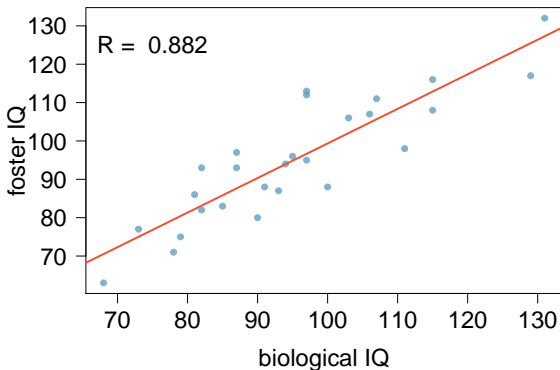
$$y = \beta_0 + \beta_1 x + \epsilon.$$

▶ In this theoretical description, $\beta_0$ and $\beta_1$ are parameters, a sort of theoretical $y$-intercept ($\beta_0$) and theoretical slope ($\beta_1$) describing the association.

▶ We make a decision about $\beta_1$ by gathering data, generating a test statistic, and analyzing it (finding a p-value).

▶ Our test statistic will be calculated based on the equation of the least-squares regression line calculated from the data:

$$\hat{y} = b_0 + b_1 x$$

Chapter 8: Regression
└─ Inference for linear regression
   └─ Understanding regression output from software

# Nature or nurture?

In 1966 Cyril Burt published a paper called "The genetic determination of differences in intelligence: A study of monozygotic twins reared together and apart". The data consist of IQ scores for [an assumed random sample of] 27 identical twins, one raised by foster parents, the other by the biological parents.

Chapter 8: Regression
└─ Inference for linear regression
    └─ Understanding regression output from software

Which of the following is false?

```
Coefficients:
                Estimate Std. Error t value Pr(>|t|)
(Intercept)      9.20760    9.29990   0.990    0.332
bioIQ            0.90144    0.09633   9.358  1.2e-09

Residual standard error: 7.729 on 25 degrees of freedom
Multiple R-squared: 0.7779, Adjusted R-squared: 0.769
F-statistic: 87.56 on 1 and 25 DF,  p-value: 1.204e-09
```

(a) Additional 10 points in the biological twin's IQ is associated with additional 9 points in the foster twin's IQ, on average.

(b) Roughly 78% of the foster twins' IQs can be accurately predicted by the model.

(c) The linear model is $\widehat{fosterIQ} = 9.2 + 0.9 \times bioIQ$.

(d) Foster twins with IQs higher than average IQs tend to have biological twins with higher than average IQs as well.

Chapter 8: Regression
└─ Inference for linear regression
  └─ Understanding regression output from software

Which of the following is <u>false</u>?

```
Coefficients:
              Estimate Std. Error t value Pr(>|t|)
(Intercept)    9.20760    9.29990   0.990    0.332
bioIQ          0.90144    0.09633   9.358  1.2e-09

Residual standard error: 7.729 on 25 degrees of freedom
Multiple R-squared: 0.7779, Adjusted R-squared: 0.769
F-statistic: 87.56 on 1 and 25 DF,  p-value: 1.204e-09
```

(a) Additional 10 points in the biological twin's IQ is associated with additional 9 points in the foster twin's IQ, on average.

(b) *Roughly 78% of the foster twins' IQs can be accurately predicted by the model.*

(c) The linear model is $\widehat{fosterIQ} = 9.2 + 0.9 \times bioIQ$.

(d) Foster twins with IQs higher than average IQs tend to have biological twins with higher than average IQs as well.

Chapter 8: Regression
└─ Inference for linear regression
  └─ Understanding regression output from software

# Testing for the slope

Assuming that these 27 twins comprise a representative sample of all twins separated at birth, we would like to test if these data provide convincing evidence that the IQ of the biological twin is a significant predictor of IQ of the foster twin. What are the appropriate hypotheses?

(a) $H_0 : b_0 = 0$; $H_A : b_0 \neq 0$

(b) $H_0 : \beta_0 = 0$; $H_A : \beta_0 \neq 0$

(c) $H_0 : b_1 = 0$; $H_A : b_1 \neq 0$

(d) $H_0 : \beta_1 = 0$; $H_A : \beta_1 \neq 0$

Chapter 8: Regression
└─Inference for linear regression
  └─Understanding regression output from software

# Testing for the slope

Assuming that these 27 twins comprise a representative sample of all twins separated at birth, we would like to test if these data provide convincing evidence that the IQ of the biological twin is a significant predictor of IQ of the foster twin. What are the appropriate hypotheses?

(a) $H_0 : b_0 = 0$; $H_A : b_0 \neq 0$

(b) $H_0 : \beta_0 = 0$; $H_A : \beta_0 \neq 0$

(c) $H_0 : b_1 = 0$; $H_A : b_1 \neq 0$

(d) $H_0 : \beta_1 = 0$; $H_A : \beta_1 \neq 0$

Chapter 8: Regression
└─ Inference for linear regression
   └─ Understanding regression output from software

# Testing for the slope (cont.)

|              | Estimate | Std. Error | t value | Pr(>|t|) |
|-------------:|---------:|-----------:|--------:|---------:|
| (Intercept)  | 9.2076   | 9.2999     | 0.99    | 0.3316   |
| bioIQ        | 0.9014   | 0.0963     | 9.36    | 0.0000   |

Chapter 8: Regression
└─ Inference for linear regression
    └─ Understanding regression output from software

# Testing for the slope (cont.)

|             | Estimate | Std. Error | t value | Pr(>|t|) |
|-------------|----------|------------|---------|----------|
| (Intercept) | 9.2076   | 9.2999     | 0.99    | 0.3316   |
| bioIQ       | 0.9014   | 0.0963     | 9.36    | 0.0000   |

▶ We always use a $t$-test in inference for regression.

Chapter 8: Regression
└─Inference for linear regression
  └─Understanding regression output from software

# Testing for the slope (cont.)

|             | Estimate | Std. Error | t value | Pr(>|t|) |
|-------------|----------|------------|---------|----------|
| (Intercept) | 9.2076   | 9.2999     | 0.99    | 0.3316   |
| bioIQ       | 0.9014   | 0.0963     | 9.36    | 0.0000   |

▶ We always use a $t$-test in inference for regression.

*Remember:* Test statistic, $T = \frac{point\ estimate - null\ value}{SE}$

Chapter 8: Regression
└─Inference for linear regression
   └─Understanding regression output from software

# Testing for the slope (cont.)

|             | Estimate | Std. Error | t value | Pr(>\|t\|) |
|-------------|----------|------------|---------|-----------|
| (Intercept) | 9.2076   | 9.2999     | 0.99    | 0.3316    |
| bioIQ       | 0.9014   | 0.0963     | 9.36    | 0.0000    |

- ▶ We always use a $t$-test in inference for regression.

  *Remember:* Test statistic, $T = \frac{point\ estimate - null\ value}{SE}$

- ▶ Point estimate $= b_1$, the observed slope.

Chapter 8: Regression
└─ Inference for linear regression
  └─ Understanding regression output from software

# Testing for the slope (cont.)

|  | Estimate | Std. Error | t value | Pr(>|t|) |
|---|---|---|---|---|
| (Intercept) | 9.2076 | 9.2999 | 0.99 | 0.3316 |
| bioIQ | 0.9014 | 0.0963 | 9.36 | 0.0000 |

▶ We always use a $t$-test in inference for regression.

*Remember:* Test statistic, $T = \frac{point\ estimate - null\ value}{SE}$

▶ Point estimate $= b_1$, the observed slope.

▶ $SE_{b_1}$ is the standard error associated with the slope (given in the table!)

Chapter 8: Regression
└─ Inference for linear regression
   └─ Understanding regression output from software

# Testing for the slope (cont.)

|  | Estimate | Std. Error | t value | Pr(>|t|) |
|---|---|---|---|---|
| (Intercept) | 9.2076 | 9.2999 | 0.99 | 0.3316 |
| bioIQ | 0.9014 | 0.0963 | 9.36 | 0.0000 |

▶ We always use a $t$-test in inference for regression.

  *Remember: Test statistic,* $T = \frac{point\ estimate - null\ value}{SE}$

▶ Point estimate $= b_1$, the observed slope.

▶ $SE_{b_1}$ is the standard error associated with the slope (given in the table!)

▶ Degrees of freedom associated with the slope is $df = n - 2$, where $n$ is the sample size.

  (We lose 1 degree of freedom for each parameter we estimate, and in simple linear regression we estimate 2 parameters, $\beta_0$ and $\beta_1$.)

Chapter 8: Regression
└─ Inference for linear regression
   └─ Understanding regression output from software

# Testing for the slope (cont.)

|             | Estimate | Std. Error | t value | Pr(>|t|) |
|-------------|----------|------------|---------|----------|
| (Intercept) | 9.2076   | 9.2999     | 0.99    | 0.3316   |
| biolQ       | 0.9014   | 0.0963     | 9.36    | 0.0000   |

$$T \;=\; \frac{0.9014 - 0}{0.0963} = 9.36$$

Chapter 8: Regression
└─ Inference for linear regression
  └─ Understanding regression output from software

# Testing for the slope (cont.)

|              | Estimate | Std. Error | t value | Pr(>|t|) |
|-------------:|---------:|-----------:|--------:|---------:|
| (Intercept)  | 9.2076   | 9.2999     | 0.99    | 0.3316   |
| biolQ        | 0.9014   | 0.0963     | 9.36    | 0.0000   |

$$T = \frac{0.9014 - 0}{0.0963} = 9.36$$
$$df = 27 - 2 = 25$$

Chapter 8: Regression
  └─Inference for linear regression
       └─Understanding regression output from software

# Testing for the slope (cont.)

|  | Estimate | Std. Error | t value | Pr(>|t|) |
|---|---|---|---|---|
| (Intercept) | 9.2076 | 9.2999 | 0.99 | 0.3316 |
| bioIQ | 0.9014 | 0.0963 | 9.36 | 0.0000 |

$$
\begin{aligned}
T &= \frac{0.9014 - 0}{0.0963} = 9.36 \\
df &= 27 - 2 = 25 \\
p - value &= P(|T| > 9.36) < 0.01
\end{aligned}
$$

In fact, p-value is:

```
> 2*(1-pt(9.36,25))
[1] 1.197331e-09
```

Chapter 8: Regression
└─ Inference for linear regression
   └─ CI for the slope

# Confidence interval for the slope

Remember that a confidence interval is calculated as *point estimate* $\pm$ *ME* and the degrees of freedom associated with the slope in a simple linear regression is $n - 2$. Which of the below is the correct 95% confidence interval for the slope parameter? Note that the model is based on observations from 27 twins.

|             | Estimate | Std. Error | t value | Pr($>$|t|) |
|-------------|----------|------------|---------|------------|
| (Intercept) | 9.2076   | 9.2999     | 0.99    | 0.3316     |
| biolQ       | 0.9014   | 0.0963     | 9.36    | 0.0000     |

(a) $9.2076 \pm 1.65 \times 9.2999$

(b) $0.9014 \pm 2.06 \times 0.0963$

(c) $0.9014 \pm 1.96 \times 0.0963$

(d) $9.2076 \pm 1.96 \times 0.0963$

Chapter 8: Regression
└─ Inference for linear regression
 └─ CI for the slope

# Confidence interval for the slope

Remember that a confidence interval is calculated as *point estimate* $\pm$ *ME* and the degrees of freedom associated with the slope in a simple linear regression is $n - 2$. Which of the below is the correct 95% confidence interval for the slope parameter? Note that the model is based on observations from 27 twins.

|             | Estimate | Std. Error | t value | Pr($>$\|t\|) |
|-------------|----------|------------|---------|----------|
| (Intercept) | 9.2076   | 9.2999     | 0.99    | 0.3316   |
| biolQ       | 0.9014   | 0.0963     | 9.36    | 0.0000   |

$$n = 27 \qquad df = 27 - 2 = 25$$

(a) $9.2076 \pm 1.65 \times 9.2999$

(b) $0.9014 \pm 2.06 \times 0.0963$

(c) $0.9014 \pm 1.96 \times 0.0963$

(d) $9.2076 \pm 1.96 \times 0.0963$

Chapter 8: Regression
└─Inference for linear regression
  └─CI for the slope

# Confidence interval for the slope

Remember that a confidence interval is calculated as *point estimate* $\pm$ *ME* and the degrees of freedom associated with the slope in a simple linear regression is $n - 2$. Which of the below is the correct 95% confidence interval for the slope parameter? Note that the model is based on observations from 27 twins.

|             | Estimate | Std. Error | t value | Pr($>$|t|) |
|-------------|----------|------------|---------|------------|
| (Intercept) | 9.2076   | 9.2999     | 0.99    | 0.3316     |
| biolQ       | 0.9014   | 0.0963     | 9.36    | 0.0000     |

(a) $9.2076 \pm 1.65 \times 9.2999$

(b) $0.9014 \pm 2.06 \times 0.0963$

(c) $0.9014 \pm 1.96 \times 0.0963$

(d) $9.2076 \pm 1.96 \times 0.0963$

$n = 27 \qquad df = 27 - 2 = 25$

$95\% : t_{25}^{\star} = 2.06$

Chapter 8: Regression
└─ Inference for linear regression
   └─ CI for the slope

# Confidence interval for the slope

Remember that a confidence interval is calculated as *point estimate* $\pm$ *ME* and the degrees of freedom associated with the slope in a simple linear regression is $n - 2$. Which of the below is the correct 95% confidence interval for the slope parameter? Note that the model is based on observations from 27 twins.

|  | Estimate | Std. Error | t value | Pr(>|t|) |
|---|---|---|---|---|
| (Intercept) | 9.2076 | 9.2999 | 0.99 | 0.3316 |
| biolQ | 0.9014 | 0.0963 | 9.36 | 0.0000 |

(a) $9.2076 \pm 1.65 \times 9.2999$

(b) $0.9014 \pm 2.06 \times 0.0963$

(c) $0.9014 \pm 1.96 \times 0.0963$

(d) $9.2076 \pm 1.96 \times 0.0963$

$$n = 27 \qquad df = 27 - 2 = 25$$

$$95\% : t_{25}^{\star} = 2.06$$

$$0.9014 \pm 2.06 \times 0.0963$$

Chapter 8: Regression
└─ Inference for linear regression
  └─ CI for the slope

# Confidence interval for the slope

Remember that a confidence interval is calculated as *point estimate* $\pm$ *ME* and the degrees of freedom associated with the slope in a simple linear regression is $n - 2$. Which of the below is the correct 95% confidence interval for the slope parameter? Note that the model is based on observations from 27 twins.

|             | Estimate | Std. Error | t value | Pr(>|t|) |
|-------------|----------|------------|---------|----------|
| (Intercept) | 9.2076   | 9.2999     | 0.99    | 0.3316   |
| bioIQ       | 0.9014   | 0.0963     | 9.36    | 0.0000   |

(a) $9.2076 \pm 1.65 \times 9.2999$

(b) $0.9014 \pm 2.06 \times 0.0963$

(c) $0.9014 \pm 1.96 \times 0.0963$

(d) $9.2076 \pm 1.96 \times 0.0963$

$$n \;=\; 27 \qquad df = 27 - 2 = 25$$
$$95\% : \; t_{25}^{\star} \;=\; 2.06$$
$$0.9014 \;\pm\; 2.06 \times 0.0963$$
$$(0.7 \quad , \quad 1.1)$$

Chapter 8: Regression
└─ Inference for linear regression
   └─ CI for the slope

# Recap

- Inference for the slope for a single-predictor linear regression model:
    - Hypothesis test:

    $$T = \frac{b_1 - null\ value}{SE_{b_1}} \qquad df = n - 2$$

    - Confidence interval:

    $$b_1 \pm t^{\star}_{df=n-2} SE_{b_1}$$

- The null value is often 0 since we are usually checking for *any* relationship between the explanatory and the response variable.

- The regression output gives $b_1$, $SE_{b_1}$, and *two-tailed* p-value for the $t$-test for the slope where the null value is 0.

- We rarely do inference on the intercept, so we'll be focusing on the estimates and inference for the slope.

Chapter 8: Regression
└─Inference for linear regression
  └─CI for the slope

# Caution

- Always be aware of the type of data you're working with: random sample, non-random sample, or population.

Chapter 8: Regression
└─ Inference for linear regression
  └─ CI for the slope

# Caution

- Always be aware of the type of data you're working with: random sample, non-random sample, or population.
- Statistical inference, and the resulting p-values, are meaningless when you already have population data.

Chapter 8: Regression
└─ Inference for linear regression
  └─ CI for the slope

# Caution

- ▶ Always be aware of the type of data you're working with: random sample, non-random sample, or population.
- ▶ Statistical inference, and the resulting p-values, are meaningless when you already have population data.
- ▶ If you have a sample that is non-random (biased), inference on the results will be unreliable.

Chapter 8: Regression
└─ Inference for linear regression
   └─ CI for the slope

# Caution

▶ Always be aware of the type of data you're working with: random sample, non-random sample, or population.

▶ Statistical inference, and the resulting p-values, are meaningless when you already have population data.

▶ If you have a sample that is non-random (biased), inference on the results will be unreliable.

▶ The ultimate goal is to have independent observations.