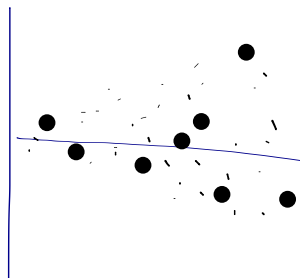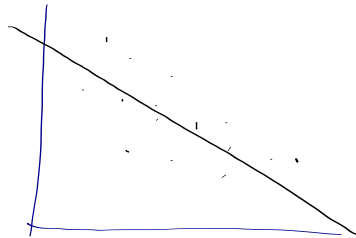Find your new seat

TODAY
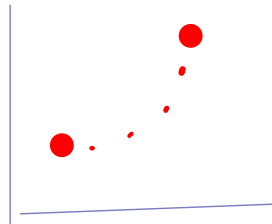
# The Role of s and r² in Regression
## (pages 188-192)
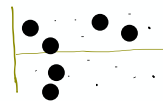
# LAST Class

We can't use r
to justify linearity.

$r = 0.89$

Why
?

strong r
but......
Form is non-linear

So....

We use Residual Plots
to determine if a LSRL
is appropriate.

We use Residual Plots
to determine if a LSRL
is appropriate.

$\hat{y} =$

once we establish that
a linear model is
appropriate. ••••

there are two tools
available to tell us how
good the predictions can be.

---

We use Residual Plots
to determine if a LSRL
is appropriate.

$S$ and $r^2$
determine how good
the predictions will be.
( How well does the line work)

$S_x$

We use Residual Plots
to determine if a LSRL
is appropriate.

$S$ and $r^2$
determine how good
the predictions will be.

( How well does the line work)

$S$    Standard Deviations of
       the residuals

$r^2$   Coefficient of
       Determination

---

$S$ and $r^2$
determine how good
the predictions will be.

( How well does the line work)

Aim Today

Interpet

$S$    Standard Deviations of
       the residuals

$r^2$   Coefficient of
       Determination

## Normally

Experience First - Formalize later

Today - Formalize right away

---

*handout*

The **standard deviation of the residuals** $s$, measures the size of a <u>typical</u> residual. That is, $s$ measures the typical distance between the actual $y$ values and the predicted $y$ values.

$y$

$\hat{y}$

$rms$

$s$ is sometimes called the typical prediction error.

Small $s$ is good

different than $s_x$

What's the typical residual size?

Typical prediction error

Can't just average them 🙂

---

**Note:** The sum of residuals will up to 0.

| L2 | L3 | L4 | 4 |
|---|---|---|---|
| .5 | -1.398 | 1.898 | |
| 1 | .207 | .793 | |
| 1 | 1.812 | -.812 | |
| 1.7 | 3.417 | -1.717 | |
| 4 | 5.022 | -1.022 | |
| 5 | 6.627 | -1.627 | |
| 9 | 8.232 | .768 | |

L4(1)=1.898

Residuals

1-Var Stats L4

1-Var Stats
x̄=2.2222222ᴇ-4
Σx=.002
Σx²=14.574056
Sx=1.349724766
σx=1.272532713
↓n=9

Sum of residuals

> The **standard deviation of the residuals s** measures the size of a typical residual. That is, **s** measures the typical distance between the actual *y* values and the predicted *y* values.
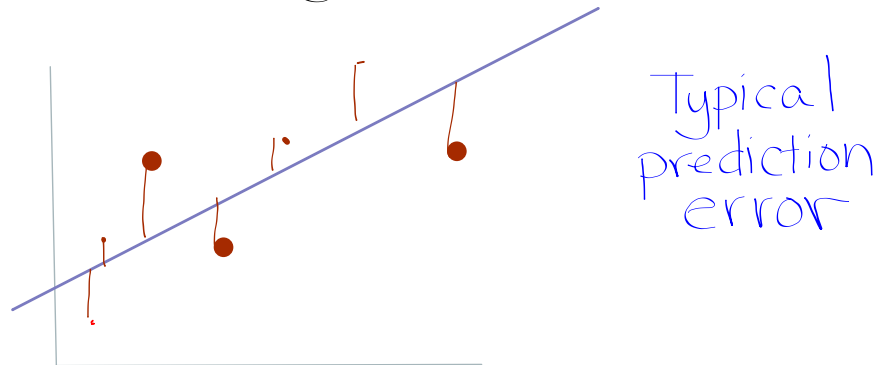
$$s = \sqrt{\frac{\sum residuals^2}{n-2}} = \sqrt{\frac{\sum (y_i - \hat{y})^2}{n-2}}$$

Most likely you will be given this value.

$$S_x = \sqrt{\frac{\sum (x_i - \bar{x})^2}{n-1}}$$

We divide by *n-2* rather than *n-1*. We used *n-1* for *s* when we estimated the mean (used $\bar{x}$ for μ). Now we are estimating both slope and the y-intercept, so we use *n-2*. We subtract one more for each parameter we estimate.

# Coefficient of Determination

$r^2$ measures the fraction of the variability in the *y* variable that is accounted for by the LSRL using *x*.

---

$$r^2 = 1 - \frac{\sum \text{residuals}^2}{\sum (y_i - \bar{y})^2}$$

The **coefficient of determination $r^2$** measures the percent reduction in the sum of squared residuals when using the least-squares regression line to make predictions, rather than the mean value of *y*.

In other words, $r^2$ measures the percent of the variability in the response variable that is accounted for by the least-squares regression line.

$r^2$ tells us how much better the LSRL does at predicting values of *y* than simply guessing the mean *y* for each value in the dataset.

# Backpacking

do #1

We'll do #2 and #3 as a class

---

**Backpacking -** Ninth-grade students at the Webb Schools go on a backpacking trip each fall. Students are divided into hiking groups of size 8 by selecting names from a hat. Before leaving, students and their backpacks are weighed. The data here are from one hiking group.

| Body weight (lb) | 120 | 187 | 109 | 103 | 131 | 165 | 158 | 116 |
|---|---|---|---|---|---|---|---|---|
| Backpack weight (lb) | 26 | 30 | 26 | 24 | 29 | 35 | 31 | 28 |

Analyze the data using stapplet.com.

1. Using www.stapplet.com find the LSRL of the data. Write it below (in you know what form!)

**Backpacking** - Ninth-grade students at the Webb Schools go on a backpacking trip each fall. Students are divided into hiking groups of size 8 by selecting names from a hat. Before leaving, students and their backpacks are weighed. The data here are from one hiking group.

| Body weight (lb) | 120 | 187 | 109 | 103 | 131 | 165 | 158 | 116 |
|---|---|---|---|---|---|---|---|---|
| Backpack weight (lb) | 26 | 30 | 26 | 24 | 29 | 35 | 31 | 28 |

Analyze the data using stapplet.com.

1. Using www.stapplet.com find the LSRL of the data. Write it below (in you know what form!)

$$\hat{y} = 16.265 + 0.091\,x$$

$$\widehat{Backpack\ weight} = 16.265 + 0.091\,(Body\ Weight)$$

2. Find and interpret **S**, *the standard deviation of the residuals.*

3. Find and interpret the value of $r^2$, the coefficient of determination.

2. Find and interpret **S**, *the standard deviation of the residuals.*

$S = 2.27$   The actual backpack weight is typically about 2.27 lb. away from the weight predicted by the LSRL with $x$ = the body weight.

3. Find and interpret the value of $r^2$, the coefficient of determination.

---

2. Find and interpret **S**, *the standard deviation of the residuals.*

$S = 2.27$   The actual backpack weight is typically about 2.27 lb. away from the weight predicted by the LSRL with $x$ = the body weight.

3. Find and interpret the value of $r^2$, the coefficient of determination.

2. Find and interpret **S**, *the standard deviation of the residuals.*

$S = 2.27$  The actual backpack weight is typically about 2.27 lb. away from the weight predicted by the LSRL with x = the body weight.

3. Find and interpret the value of $r^2$, the coefficient of determination.

$r^2 = 0.632$
About 63.2% of the variability in backpack weight is accounted for by the LSRL with x = body weight.

---

## s and $r^2$

**Big Ideas:**

Standard Deviation of Residuals (s)

Interpretation ●

Coefficient of Determination $r^2$

Interpretation ●

---

Root mean Squares and $r^2$

**Big Ideas:**

Standard Deviation of Residuals (s)

Interpretation ●

"The actual y-context is typically about S away from the #predicted by the LSRL". with x= context

Coefficient of Determination $r^2$

Interpretation ●

"About $r^2$ % of the variability in y-context is accounted for by the LSRL when x= x-context

## Mickey's last bungee jump

Partners • A ↔ B

↓ does a     ↘ does b

After the response is written, switch papers and check each other. •

---

(a) Interpret the value of *s*.

(b) Interpret the value of $r^2$.

**(a) Interpret the value of s.**

The distance travelled by Mickey is typically about 4.11 cm away from the distance predicted by the LSRL with $x$ = #rubber bands.

**(b) Interpret the value of $r^2$.**

About 98% of the variability in dist. travelled by Mickey is accounted for by the LSRL with $x$ = #rubber bands.

---

## Interpreting Computer Regression Output
(pages 192-194)

You are not expected to able to use the software but you are expected to interpret the output.

From the output, be sure you can find the:

slope     a                    $\hat{y} = a + bx$

y-intercept  b

          S

          $r^2$

---

**Minitab**

                    Slope              y intercept

| Predictor | Coef | SE Coef | T | P |
|---|---|---|---|---|
| Constant | 38257 | 2446 | 15.64 | 0.000 |
| Miles Driven | -0.16292 | 0.03096 | -5.26 | 0.000 |

                                              $r^2$

S = 5740.13          R-Sq = 66.4%    R-Sq(adj) = 64.0%

          Standard deviation of the residuals

**JMP**

Summary of Fit

| | | |
|---|---|---|
| RSquare | 0.664248 | $r^2$ |
| RSquare Adj | 0.640266 | Standard deviation |
| Root Mean Square Error | 5740.131 | of the residuals |
| Mean of Response | 27833.69 | |
| Observations (or Sum Wgts) | 16 | |

Parameter Estimates

| Term | Estimate | Std Error | t Ratio | Prob>|t| |
|---|---|---|---|---|
| Intercept | 38257.135 | 2445.813 | 15.64 | <.0001 |
| Miles Driven | -0.162919 | 0.030956 | -5.26 | 0.0001 |

*y* intercept — Slope

---

# Can we predict a school's average SAT math score?
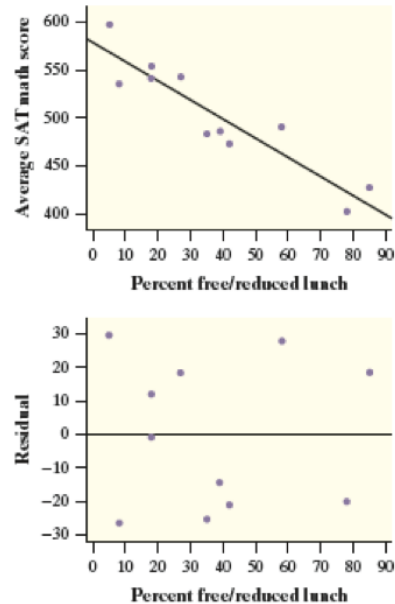*Interpreting regression output*

check in w/ table mates
as you progress

**Can we predict a school's average SAT math score?** *Interpreting regression output*

A random sample of 11 high schools was selected from all the high schools in Michigan. The percent of students who are eligible for free/reduced lunch and the average SAT math score of each high school in the sample were recorded.

Students with household income below a certain threshold are eligible for free/reduced lunch.

Here are a scatterplot with the least-squares regression line added, a residual plot, and some computer output:



| Predictor | Coef | SE Coef | T | P |
|---|---|---|---|---|
| Constant | 577.9 | 12.5 | 46.16 | 0.000 |
| Foot length | −1.993 | 0.276 | −7.22 | 0.000 |

S = 23.3168   R-Sq = 85.29%   R-Sq(adj) = 83.66%

$r^2 = 85.29$
$r =$

**(a) Is a line an appropriate model to use for these data? Explain how you know the answer.**

Because the scatter plot shows a linear association and the residual plot shows no leftover pattern, the line is approp.

↑ shows a random scatter

**(b) Find the correlation.**

$r = \pm\sqrt{0.8529} = \pm 0.924$

because the relationship is negative

$r = -0.924$

r  values tell us
two things

$r = 0.56$
0.67                Strength

but also direction

$r = -0.67$

---

| Predictor | Coef | SE Coef | T | P |
|---|---|---|---|---|
| Constant | 577.9 | 12.5 | 46.16 | 0.000 |
| Foot length | −1.993 | 0.276 | −7.22 | 0.000 |

S = 23.3168   R-Sq = 85.29%   R-Sq(adj) = 83.66%

$r^2 = 85.29$

(a) Is a line an appropriate model to use for these data? Explain how you know the answer.

Because the scatter plot shows a linear association and the residual plot shows no leftover pattern, the line is appropr.

shows a random scatter

(b) Find the correlation.

$r = \pm\sqrt{0.8529} = \pm 0.924$

because the relationship is negative

$r = -0.924$

(c) What is the equation of the least-squares regression line that describes the relationship between percent free/reduced lunch and average SAT math score? Define any variables that you use.

$$\hat{y} = 577.9 - 1.993x \quad \text{where}$$

$\hat{y}$ is the predicted average SAT math Score.
and $x$ is percent free/reduced lunch.

(d) By about how much do the actual average SAT math scores typically vary from the values predicted by the least-squares regression line with $x$ = percent free/reduced lunch?

$S = 23.3168$ so the actual average SAT math scores typically vary by about 23.3168 from the values predicted by the regression line using $x$ = percent/free lunch.

See your
LCQ

Assignment:

**3.2** 55, 57, 59, 67

pp. 188-194