

# Lecture 8: more Chapter 5, Section 3

## Relationships between Two Quantitative Variables; Regression

---

- Properties of Correlation
- Equation of Regression Line; Residuals
- Effect of Explanatory/Response Roles
- Unusual Observations
- Sample vs. Population
- Time Series; Additional Variables

# Looking Back: *Review*

---

## □ 4 Stages of Statistics

- Data Production (discussed in Lectures 1-4)
- Displaying and Summarizing
  - Single variables: 1 cat, 1 quan (discussed Lectures 5-8)
  - Relationships between 2 variables:
    - Categorical and quantitative (discussed in Lecture 9)
    - Two categorical (discussed in Lecture 10)
    - Two quantitative
- Probability
- Statistical Inference

# Review

---

- Relationship between 2 quantitative variables
  - Display with **scatterplot**
  - Summarize:
    - **Form**: linear or curved
    - **Direction**: positive or negative
    - **Strength**: strong, moderate, weak

If form is linear, **correlation**  $r$  tells direction and strength.

Also, equation of **least squares regression line** lets us predict a response  $\hat{y}$  for any explanatory value  $x$ .

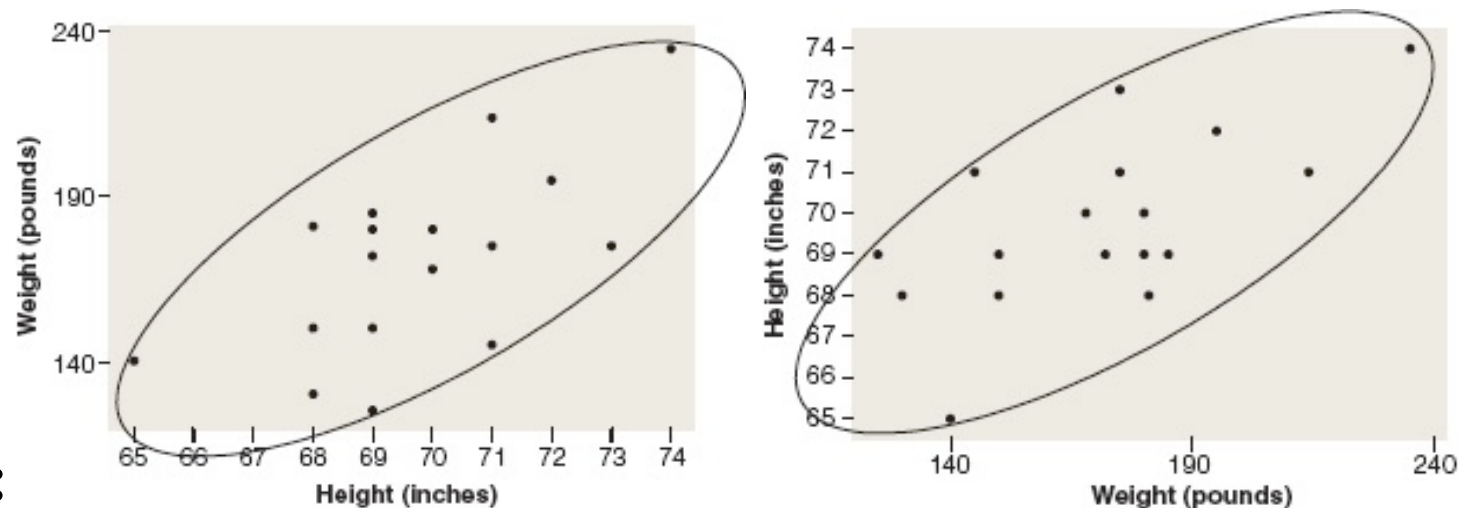
# More about Correlation $r$

---

- Tells direction and strength of linear relation between 2 quantitative variables
  - A strong curved relationship may have  $r$  close to 0
  - Correlation not appropriate for categorical data
- Unaffected by roles explanatory/response
- Unaffected by change of units
- Overstates strength if based on averages

# Example: *Correlation when Roles are Switched*

- **Background:** Male students' wt vs ht (left) or ht vs wt (right):



- **Questions:**

- How do directions and strengths compare, left vs. right?
- How do correlations  $r$  compare, left vs. right?

- **Responses:**

- \_\_\_\_\_
- \_\_\_\_\_

# More about Correlation $r$

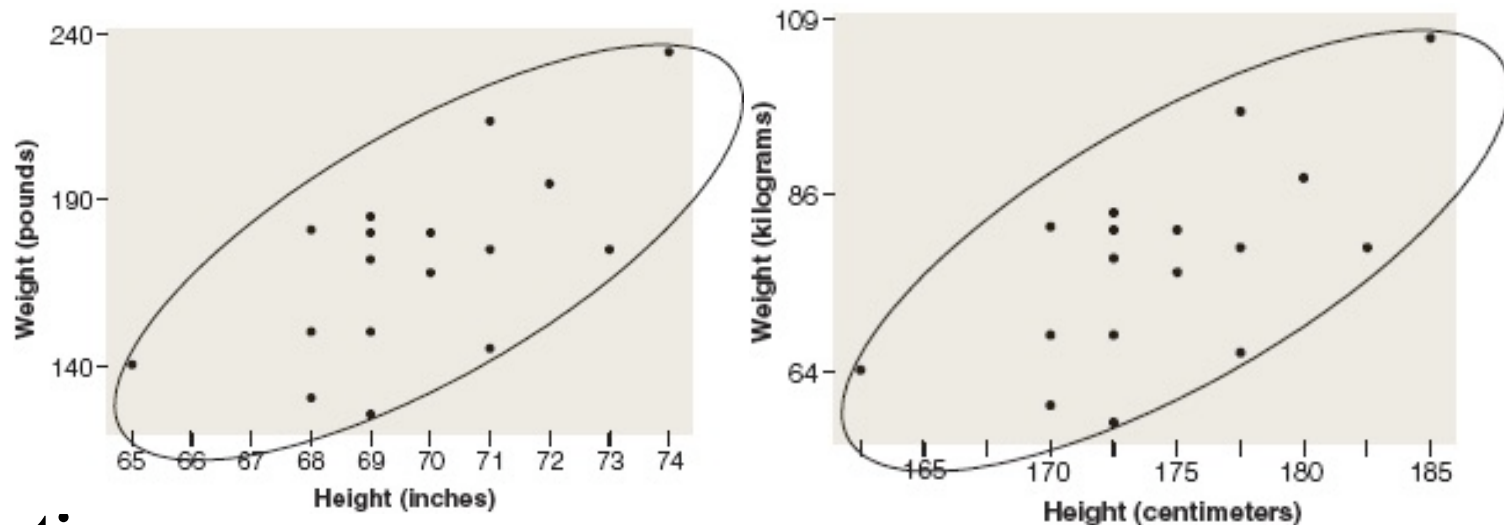
---

- Tells direction and strength of linear relation between 2 quantitative variables
  - A strong curved relationship may have  $r$  close to 0
  - Correlation not appropriate for categorical data
- Unaffected by roles explanatory/response
- Unaffected by change of units
- Overstates strength if based on averages

# Example: Correlation when Units are Changed

□ **Background:** For male students plot...

**Left:** wt (**lbs**) vs. ht (**in**) or **Right:** wt (**kg**) vs. ht (**cm**)



□ **Question:**

■ How do directions, strengths, and  $r$  compare, left vs. right?

□ **Response:**

## More about Correlation $r$

---

- Tells direction and strength of linear relation between 2 quantitative variables
  - A strong curved relationship may have  $r$  close to 0
  - Correlation not appropriate for categorical data
- Unaffected by roles explanatory/response
- Unaffected by change of units
- Overstates strength if based on averages



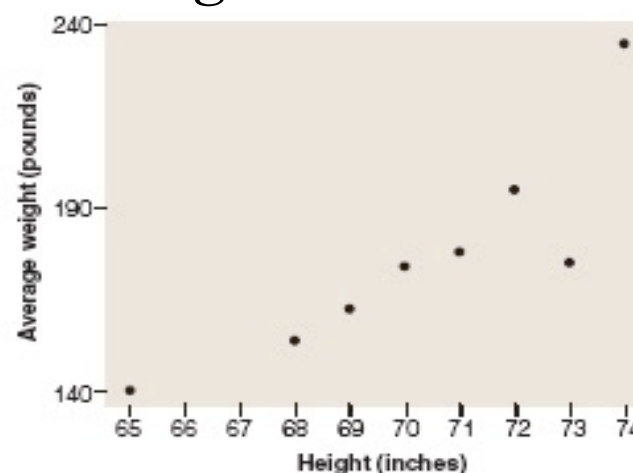
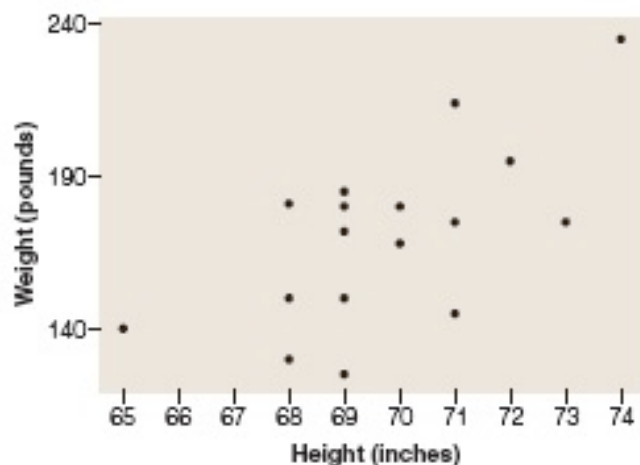
# Example: Correlation Based on Averages

- **Background:** For male students plot...

Ht	65	68	69	70	71	72	73	74
Wt	140	130 150 181	125 150 172 180 185	168 180	145 175 214	195	175	235
AvWt	140	153.7	162.4	174.0	178.0	195	175	235

**Left:** wt. vs. ht. or

**Right:** *average* wt. vs. ht.



- **Question:** Which one has  $r = +0.87$ ? (other  $r = +0.65$ )
- **Response:** Plot on \_\_\_\_\_ has  $r = +0.87$  (**stronger**).



## **Example:** *Correlation Based on Averages*

---

*In general, correlation based on averages tends to overstate strength because scatter due to individuals has been reduced.*

# Least Squares Regression Line

---

If form appears **linear**, then we picture points clustered around a straight line.

- **Questions** (Rhetorical):

1. Is there only one “best” line?
2. If so, how can we find it?
3. If found, how can we use it?

- **Responses:** (in reverse order)

3. If found, can use line to **make predictions**.

# Least Squares Regression Line

---

## ■ Response:

3. If found, can use line to **make predictions**.

Write equation of line  $\hat{y} = b_0 + b_1x$ :

- Explanatory value is  $x$
- Predicted response is  $\hat{y}$
- y-intercept is  $b_0$
- Slope is  $b_1$

and use the line to **predict** a response for any given explanatory value.

# Least Squares Regression Line

---

If form appears linear, then we picture points clustered around a straight line.

## ■ Questions:

1. Is there only one “best” line?
2. If so, how can we find it?
3. If found, how can we use it? *Predictions*

## ■ Response:

2. Find line that **makes best predictions.**

# Least Squares Regression Line

---

## ■ Response:

2. Find line that **makes best predictions:**

Minimize sum of squared *residuals* (prediction errors). Resulting line called **least squares line or regression line.**

*A Closer Look: The mathematician Sir Francis Galton called it the “regression” line because of the “regression to mediocrity” seen in any imperfect relationship: besides responding to  $x$ , we see  $y$  tending towards its average value.*

# Least Squares Regression Line

---

If form appears linear, then we picture points clustered around a straight line.

## ■ Questions:

1. Is there only one “best” line?
2. If so, how can we find it? *Minimize errors*
3. If found, how can we use it? *Predictions*

## ■ Response:

1. Methods of calculus → *unique “best” line*

# Least Squares Regression Line

---

If form appears linear, then we picture points clustered around a straight line.

## ■ Questions:

1. Is there only one “best” line?
2. If so, how can we find it?
3. If found, how can we use it?

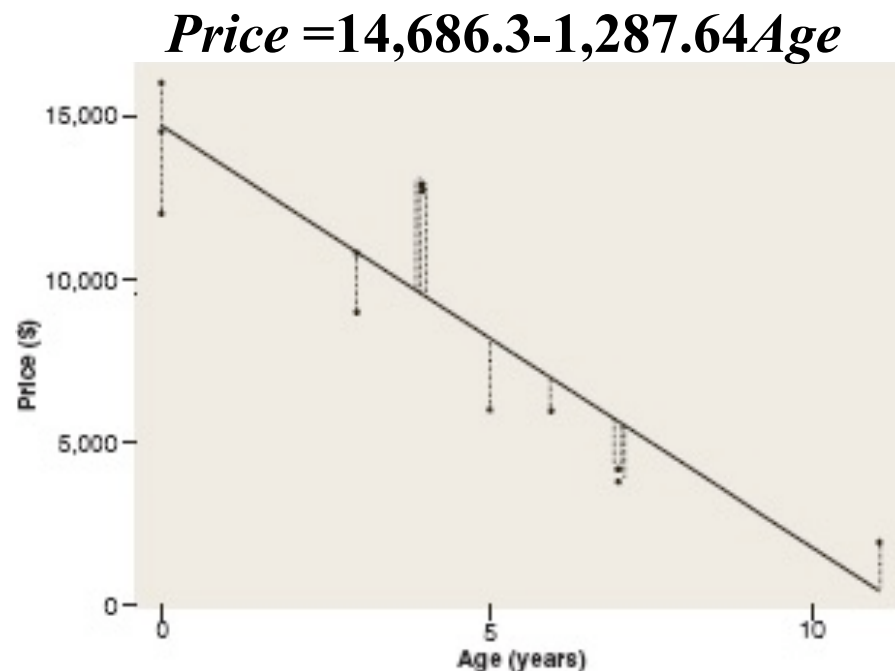
## ■ Response:

1. “Best” line has  $b_1 = r \frac{s_y}{s_x}$   $b_0 = \bar{y} - b_1 \bar{x}$



## Example: *Least Squares Regression Line*

- **Background:** Car-buyer wants to know if \$4,000 is a fair price for an 8-yr-old Grand Am; uses software to regress price on age for 14 used Grand Am's:



- **Question:** How can she use the line?
- **Response:** Predict for  $x=8$ ,  $\hat{y}$  \_\_\_\_\_.

# Least Squares Regression Line

---

Summarize linear relationship between explanatory ( $x$ ) and response ( $y$ ) values with line  $\hat{y} = b_0 + b_1x$  that minimizes sum of squared prediction errors (called *residuals*).

- **Slope:** predicted change in response  $y$  for every unit increase in explanatory value  $x$
- **Intercept:** where best-fitting line crosses  $y$ -axis (predicted response for  $x=0$ ?)

## Example: *Least Squares Regression Line*

---

- **Background:** Car-buyer used software to regress price on age for 14 used Grand Am's.

The regression equation is  
$$\text{Price} = 14690 - 1288 \text{ Age}$$

- **Question:** What do the slope (-1,288) and intercept (14,690) tell us?
- **Response:**
  - **Slope:** For each additional year in age, predict price \_\_\_\_\_
  - **Intercept:** Best-fitting line \_\_\_\_\_

## Example: *Extrapolation*

---

- **Background:** Car-buyer used software to regress price on age for 14 used Grand Am's.

The regression equation is  
$$\text{Price} = 14690 - 1288 \text{ Age}$$

- **Question:** Should we predict a new Grand Am to cost  $\$14,690 - \$1,288(0) = \$14,690$ ?
- **Response:**

# Definition

---

- **Extrapolation:** using the regression line to predict responses for explanatory values outside the range of those used to construct the line.

## Example: *More Extrapolation*

---

- **Background:** A regression of 17 male students' weights (lbs.) on heights (inches) yields the equation

$$\hat{y} = -438 + 8.7x$$

- **Question:** What weight does the line predict for a 20-inch-long infant?
- **Response:**

# Expressions for slope and intercept

---

Consider slope and intercept of the least squares regression line  $\hat{y} = b_0 + b_1x$

□ **Slope:**  $b_1 = r \frac{s_y}{s_x}$  so if  $x$  increases by a standard deviation, predict  $y$  to increase by  $r$  standard deviations

- **$|r|$  close to 1:**  $y$  responds closely to  $x$
- **$|r|$  close to 0:**  $y$  hardly responds to  $x$

# Expressions for slope and intercept

---

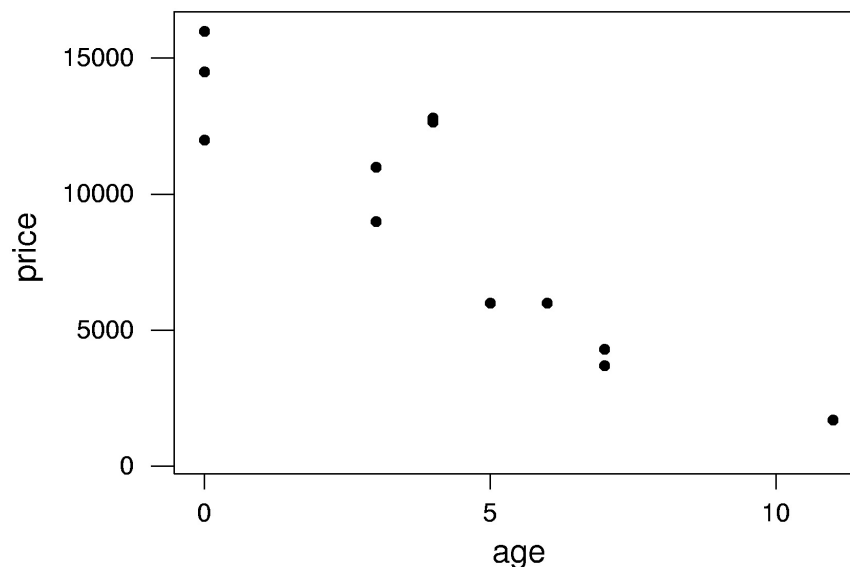
Consider slope and intercept of the least squares regression line  $\hat{y} = b_0 + b_1x$

- **Slope:**  $b_1 = r \frac{s_y}{s_x}$  so if  $x$  increases by a standard deviation, predict  $y$  to increase by  $r$  standard deviations
  - **Intercept:**  $b_0 = \bar{y} - b_1\bar{x}$  so when  $x = \bar{x}$  predict  $\hat{y} = b_0 + b_1\bar{x} = (\bar{y} - b_1\bar{x}) + b_1\bar{x} = \bar{y}$
- the line passes through the point of averages  $(\bar{x}, \bar{y})$
- See HW 5.53*



## Example: *Individual Summaries on Scatterplot*

- **Background:** Car-buyer plotted price vs. age for 14 used Grand Ams [(4, 13,000), (8, 4,000), etc.]



- **Question:** Guess the means and sds of age and price?
- **Response:** Age has approx. mean \_\_\_ yrs, sd \_\_\_ yrs; price has approx. mean \$\_\_\_\_\_, sd \$\_\_\_\_\_.

# Definitions

---

- **Residual:** error in using regression line  $\hat{y} = b_0 + b_1x$  to predict  $y$  given  $x$ . It equals the vertical distance *observed minus predicted* which can be written  $y_i - \hat{y}_i$

- **$s$ :** denotes typical residual size, calculated as

$$s = \sqrt{\frac{(y_1 - \hat{y}_1)^2 + \dots + (y_n - \hat{y}_n)^2}{n - 2}}$$

*Note:  $s$  just “averages” out the residuals  $y_i - \hat{y}_i$*

## Example: *Considering Residuals*

- **Background:** Car-buyer regressed price on age for 14 used Grand Ams [(4, 13,000), (8, 4,000), etc.].

The regression equation is

$$\text{price} = 14686 - 1290 \text{ age}$$

$$S = 2175$$

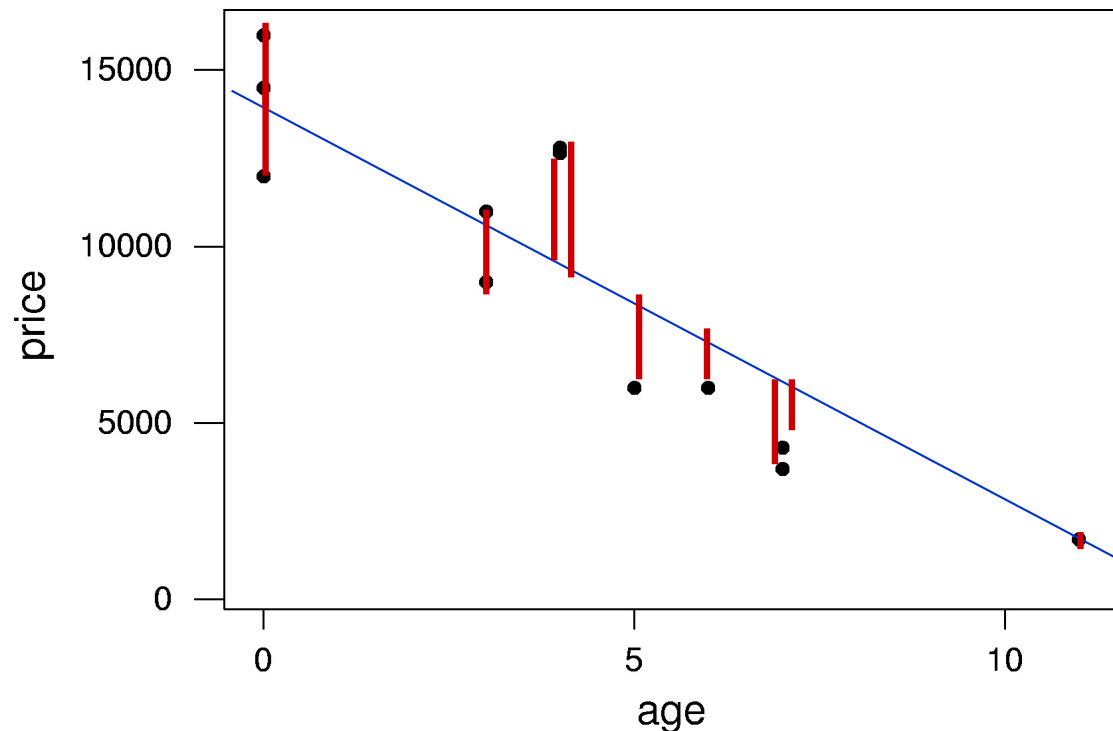
$$R\text{-Sq} = 78.5\%$$

$$R\text{-Sq}(\text{adj}) = 76.7\%$$

- **Question:** What does  $s = 2,175$  tell us?
- **Response:** Regression line predictions not perfect:
  - $x=4 \rightarrow \text{predict } \hat{y} =$   
actual  $y=13,000 \rightarrow \text{prediction error} =$
  - $x=8 \rightarrow \text{predict } \hat{y} =$   
actual  $y=4,000 \rightarrow \text{prediction error} =$
  - Typical size of 14 prediction errors is \_\_\_\_\_ (dollars)

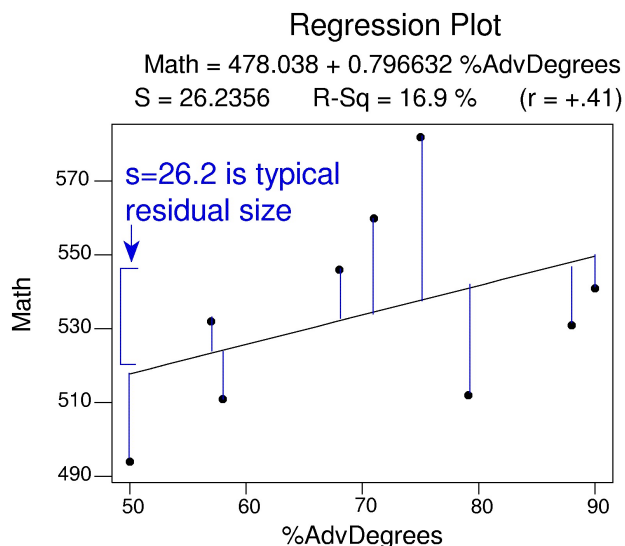
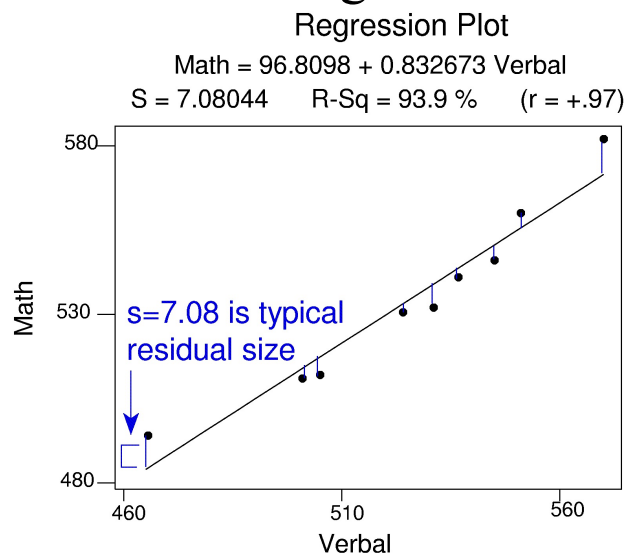
## Example: *Considering Residuals*

- Typical size of 14 prediction errors is  $s = 2,175$  (dollars): Some points' vertical distance from line more, some less;  $2,175$  is typical distance.



# Example: Residuals and their Typical Size $s$

- **Background:** For a sample of schools, regressed
  - average Math SAT on average Verbal SAT
  - average Math SAT on % of teachers w. advanced degrees

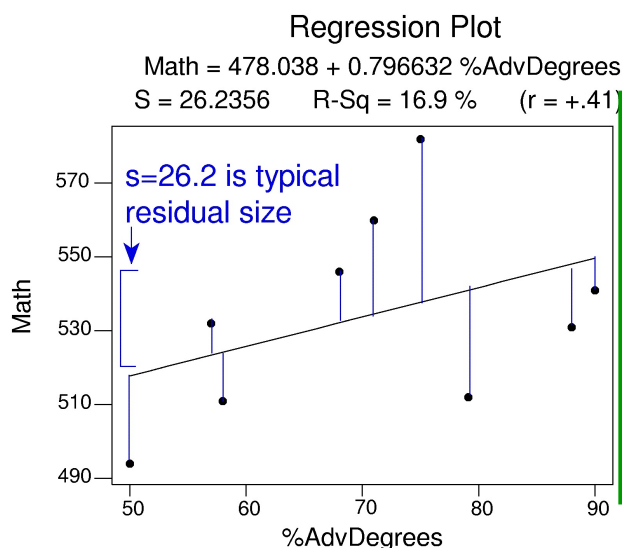
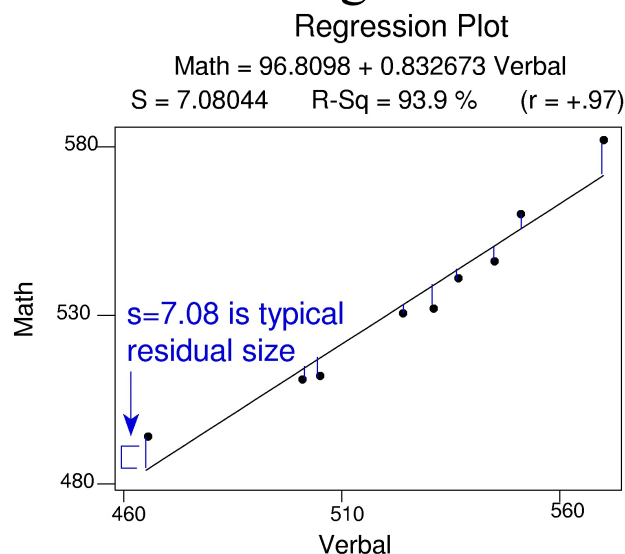


*A Closer Look: If output reports R-sq, take its square root (+ or - depending on slope) to find r.*

- **Question:** How are  $s = 7.08$  (left) and  $s = 26.2$  (right) consistent with the values of the correlation  $r$ ?
- **Response:** On left  $r = \sqrt{Rsq} = \sqrt{0.939} = 0.97$ ; relation is \_\_\_\_\_ and typical error size is \_\_\_\_\_ (only 7.08).

# Example: Residuals and their Typical Size $s$

- **Background:** For a sample of schools, regressed
  - average Math SAT on average Verbal SAT *Smaller  $s \rightarrow$  better predictions*
  - average Math SAT on % of teachers w. advanced degrees

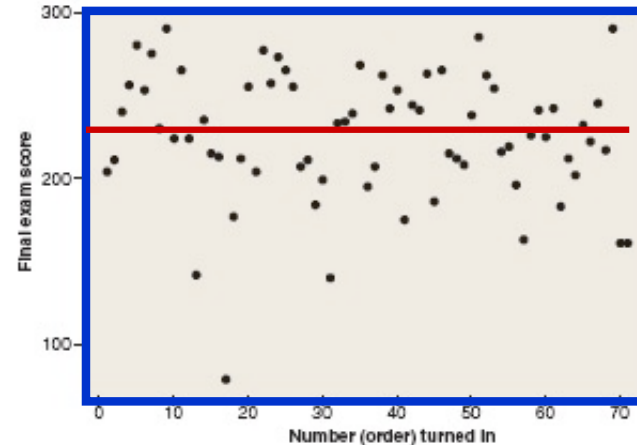
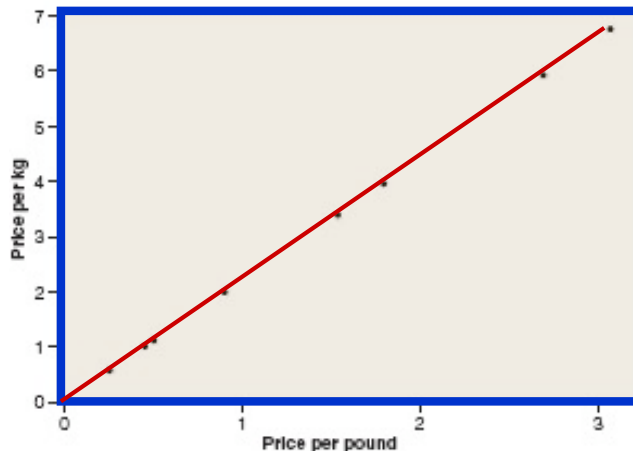


*Looking Back:  $r$  based on averages is overstated; strength of relationship for individual students would be less.*

- **Question:** How are  $s = 7.08$  (left) and  $s = 26.2$  (right) consistent with the values of the correlation  $r$ ?
- **Response:** On right  $r =$  \_\_\_\_\_ ; relation is \_\_\_\_\_ and typical error size is \_\_\_\_\_ ( $26.2$ ).

## Example: Typical Residual Size $s$ close to $s_y$ or 0

- **Background:** Scatterplots show relationships...
  - Price per kilogram vs. price per lb. for groceries
  - Students' final exam score vs. (number) order handed in



← *Regression line approx. same as line at average y-value.*

- **Questions:** Which has  $s = 0$ ? Which has  $s$  close to  $s_y$ ?
- **Responses:** Plot on left has  $s = \underline{\hspace{1cm}}$ : no prediction errors.  
Plot on right:  $s$  close to  $\underline{\hspace{1cm}}$ . (Regressing on  $x$  doesn't help; regression line is approximately horizontal.)



# Explanatory/Response Roles in Regression

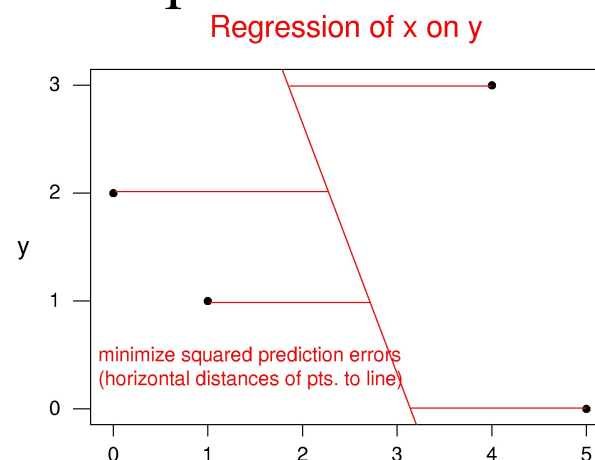
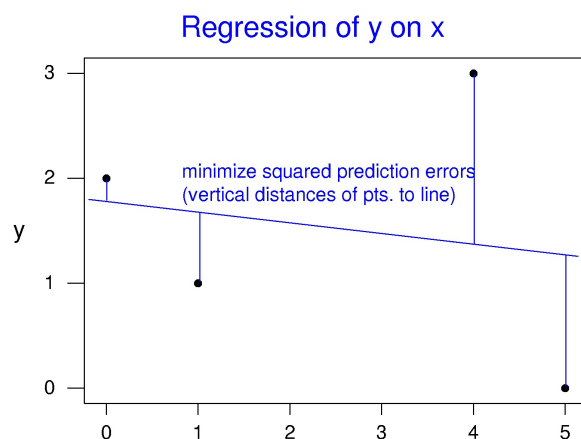
---

Our choice of roles, explanatory or response,  
does *not* affect the value of the correlation  $r$ ,  
but it *does* affect the regression line.



# Example: Regression Line when Roles are Switched

- **Background:** Compare regression of  $y$  on  $x$  (left) and regression of  $x$  on  $y$  (right) for same 4 points:



- **Question:** Do we get the same line regressing  $y$  on  $x$  as we do regressing  $x$  on  $y$ ?
- **Response:** The lines are very different.
  - Regressing  $y$  on  $x$ : \_\_\_\_\_ slope
  - Regressing  $x$  on  $y$ : \_\_\_\_\_ slope

*Context needed;  
consider variables  
and their roles  
before regressing.*

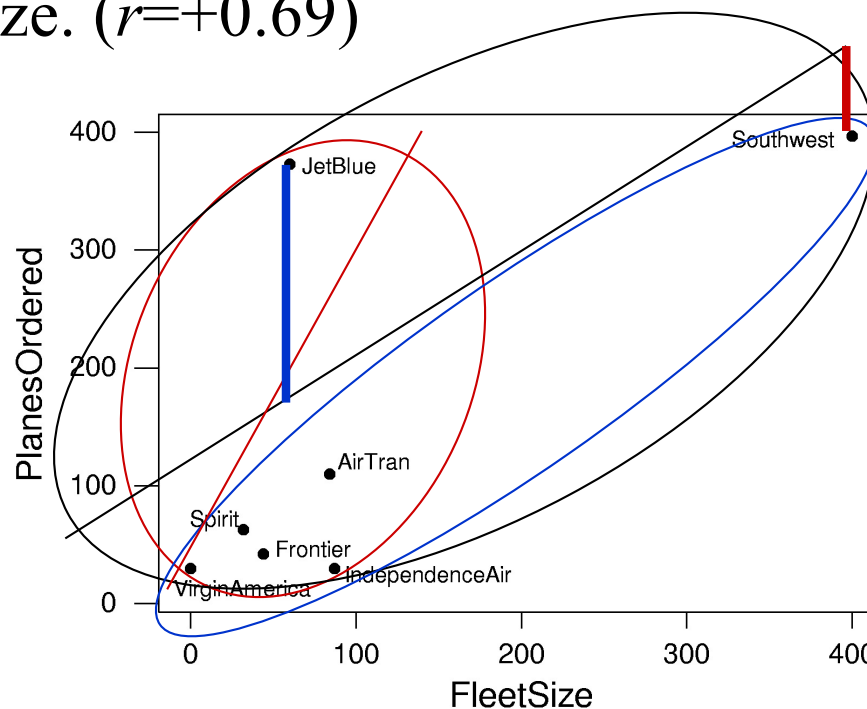
# Definitions

---

- **Outlier:** (in regression) point with unusually large residual
- **Influential observation:** point with high degree of influence on regression line.

## Example: *Outliers and Influential Observations*

- **Background:** Exploring relationship between orders for new planes and fleet size. ( $r=+0.69$ )



- **Question:** Are **Southwest** and **JetBlue** outliers or influential?
- **Response:**
  - **Southwest:** \_\_\_\_\_ (omit it  $\rightarrow$  slope changes a lot)
  - **JetBlue:** \_\_\_\_\_ (large residual; omit it  $\rightarrow r$  increases to  $+0.97$ )

## Example: *Outliers and Influential Observations*

- **Background:** Exploring relationship between orders for new planes and fleet size. ( $r = +0.69$ )

### Unusual Observations

Obs	FleetSiz	PlanesOr	Fit	SE Fit	Residual	St Resid
6	400	397.0	398.1	127.1	-1.1	-0.04 X
7	60	373.0	115.2	51.7	257.8	2.16R

R denotes an observation with a large standardized residual

X denotes an observation whose X value gives it large influence.

- **Question:** How does Minitab classify **Southwest** and **JetBlue**?

- **Response:**

- **Southwest:** \_\_\_\_\_ (marked \_\_\_\_ in Minitab)

- **JetBlue:** \_\_\_\_\_ (marked \_\_\_\_ in Minitab)

*Influential observations tend to be extreme in **horizontal** direction.*

# Definitions

---

- **Slope  $\beta_1$ :** how much response  $y$  changes in general (for entire **population**) for every unit increase in explanatory variable  $x$
- **Intercept  $\beta_0$ :** where the line that best fits all explanatory/response points (for entire **population**) crosses the  $y$ -axis

***Looking Back:** Greek letters often refer to population parameters.*

# Line for Sample vs. Population

---

- **Sample:** line best fitting **sampled** points: predicted response is

$$\hat{y} = b_0 + b_1x$$

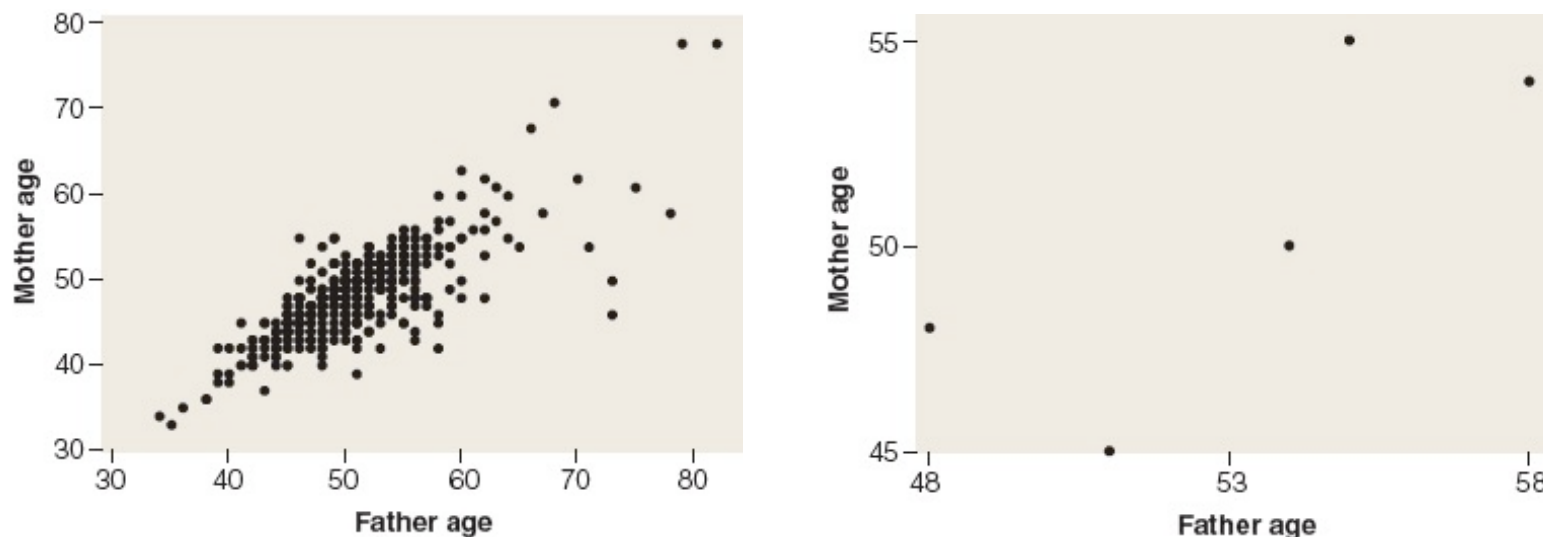
- **Population:** line best fitting **all** points in population from which given points were sampled: mean response is

$$\mu_y = \beta_0 + \beta_1x$$

A **larger sample** helps provide **more evidence** of a relationship between two quantitative variables in the general population.

## Example: *Role of Sample Size*

- **Background:** Relationship between ages of students' mothers and fathers; both scatterplots have  $r = +0.78$ , but sample size is over 400 (on left) or just 5 (on right):



- **Question:** Which plot provides more evidence of strong positive relationship in population?
- **Response:** Plot on \_\_\_\_\_

Can believe configuration on \_\_\_\_\_ occurred by chance.



# Time Series

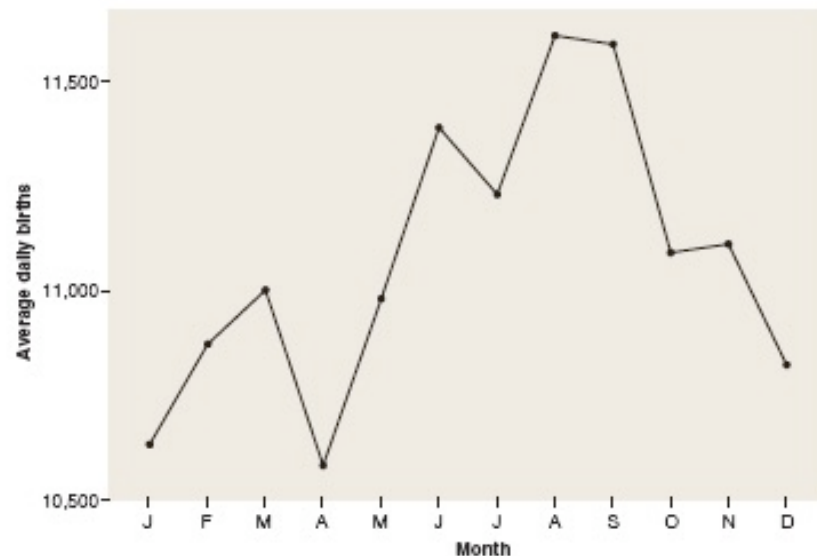
---

If explanatory variable is time, plot one response for each time value and “connect the dots” to look for general trend over time, also peaks and troughs.



## Example: *Time Series*

- **Background:** Time series plot shows average daily births each month in year 2000 in the U.S.:

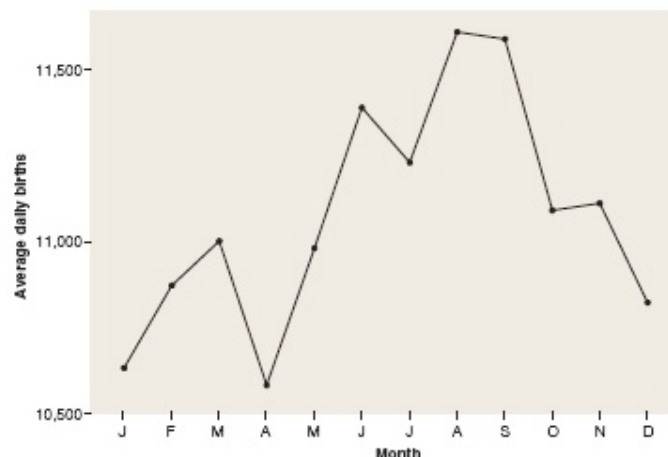


- **Question:** Where do you see a peak or a trough?

**Response:** Trough in \_\_\_\_\_, peak in \_\_\_\_\_

# Example: *Time Series*

- **Background:** Time series plot of average daily births in U.S.



- **Questions:** How can we explain why there are...
  - **Conceptions** in U.S.: fewer in July, more in December?
  - **Conceptions** in Europe: **more** in summer, **fewer** in winter?
- **Response:**

*A Closer Look: Statistical methods can't always explain "why", but at least they help understand "what" is going on.*

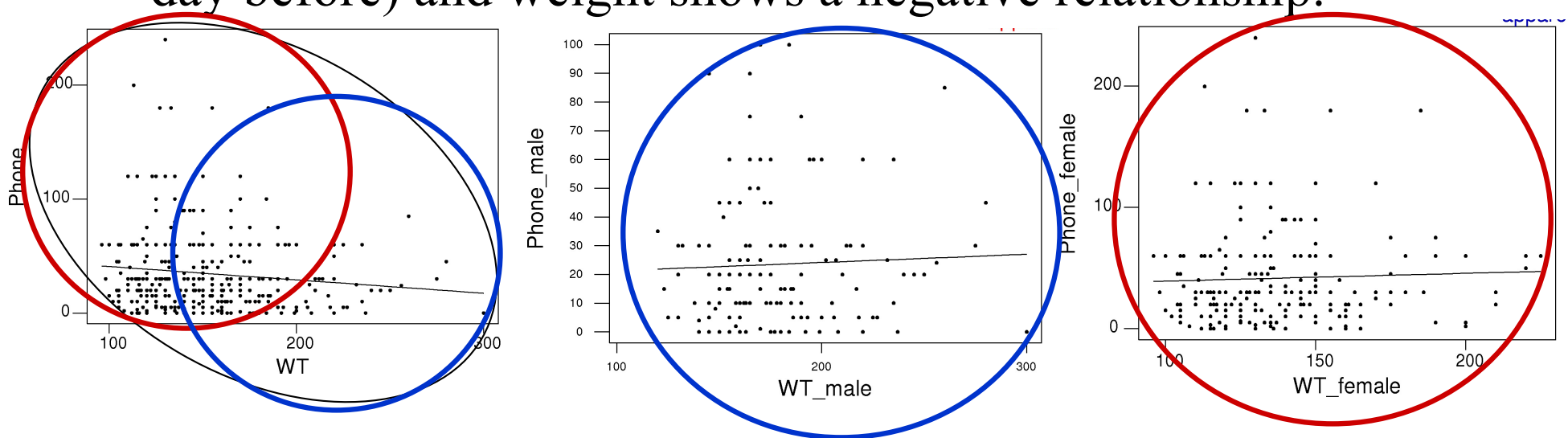
# Additional Variables in Regression

---

- **Confounding Variable:** Combining two groups that differ with respect to a variable that is related to both explanatory and response variables can affect the nature of their relationship.
- **Multiple Regression:** More advanced treatments consider impact of not just one but two or more quantitative explanatory variables on a quantitative response.

## Example: *Additional Variables*

- **Background:** A regression of phone time (in minutes the day before) and weight shows a negative relationship.



- **Questions:** Do heavy people talk on the phone less? Do light people talk more?
- **Response:** \_\_\_\_\_ is confounding variable → regress separately for \_\_\_\_\_ → no relationship

## Example: *Multiple Regression*

---

- **Background:** We used a car's age to predict its price.
- **Question:** What additional quantitative variable would help predict a car's price?
- **Response:**

# Lecture Summary

## *(Quantitative Relationships; Correlation)*

---

- Properties of  $r$ 
  - Unaffected by explanatory/response roles
  - Unaffected by change of units
  - Overstates strength if based on averages

# Lecture Summary (*Regression*)

---

- Equation of regression line
  - Interpreting slope and intercept
  - Extrapolation
- Residuals: typical size is  $s$
- Line affected by explanatory/response roles
- Outliers and influential observations
- Line for sample or population; role of sample size
- Time series
- Additional variables