# Tuesday, March 13, 2018

| Time at the lunch table | Caloric intake |
|---|---|
| 21.4 | 472 |
| 30.8 | 498 |
| 37.7 | 335 |
| 32.8 | 423 |
| 39.5 | 437 |
| 22.8 | 508 |
| 34.1 | 431 |
| 33.9 | 479 |
| 43.8 | 454 |
| 42.4 | 450 |
| 43.1 | 410 |
| 29.2 | 504 |
| 31.3 | 437 |
| 28.6 | 489 |
| 32.9 | 436 |

| | |
|---|---|
| 30.6 | 480 |
| 35.1 | 439 |
| 33.0 | 444 |
| 43.7 | 408 |

- Warm-up
  - Using the given data (*lunch19*)
    - Create a scatterplot
    - Find the regression line
    *in Calculator*
- Unit Overview
- Linear Regression inference
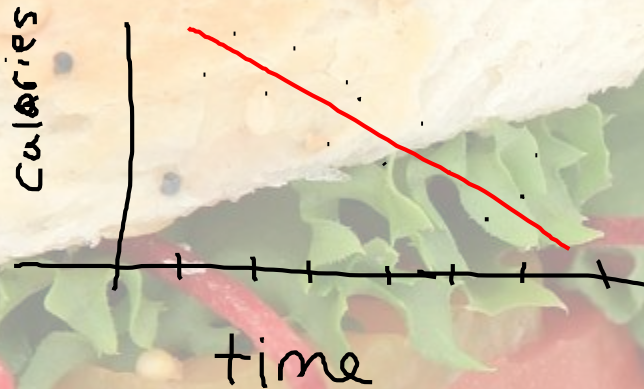
# Objectives

**Content Objective**: I will use the linear regression analysis to create confidence intervals and perform hypothesis testing.

**Social Objective**: I will participate in class activities.

**Language Objective**: I will clearly write down formulas, conditions & assumptions, and other notes.

# Warm-up

- Create a scatterplot
- Find the regression line

$$\frac{\Delta y}{\Delta x}$$

Calories = 574.98 − 3.69 time

When we spend 0 minutes at the lunch table, we predict 574.98 calories.

For every 1 minute increase at the lunch we predict calories will decrease by 3.69.

| Time at the lunch table | Caloric intake |
|---|---|
| 21.4 | 472 |
| 30.8 | 498 |
| 37.7 | 335 |
| 32.8 | 423 |
| 39.5 | 437 |
| 22.8 | 508 |
| 34.1 | 431 |
| 33.9 | 479 |
| 43.8 | 454 |
| 42.4 | 450 |
| 43.1 | 410 |
| 29.2 | 504 |
| 31.3 | 437 |
| 28.6 | 489 |
| 32.9 | 436 |
| 30.6 | 480 |
| 35.1 | 439 |
| 33.0 | 444 |
| 43.7 | 408 |

# How confident are you about the relationship?

- Confidence intervals about the slope…

- Hypothesis tests about the slope…

# But first... *Quantitative?*
## Assumptions and Conditions

- In Chapter 8 when we fit lines to data, we needed to check only the Straight Enough Condition.

- Now, when we want to make inferences about the coefficients of the line, we'll have to make more assumptions (and thus check more conditions).

- We need to be careful about the order in which we check conditions. If an initial assumption is not true, it makes no sense to check the later ones.
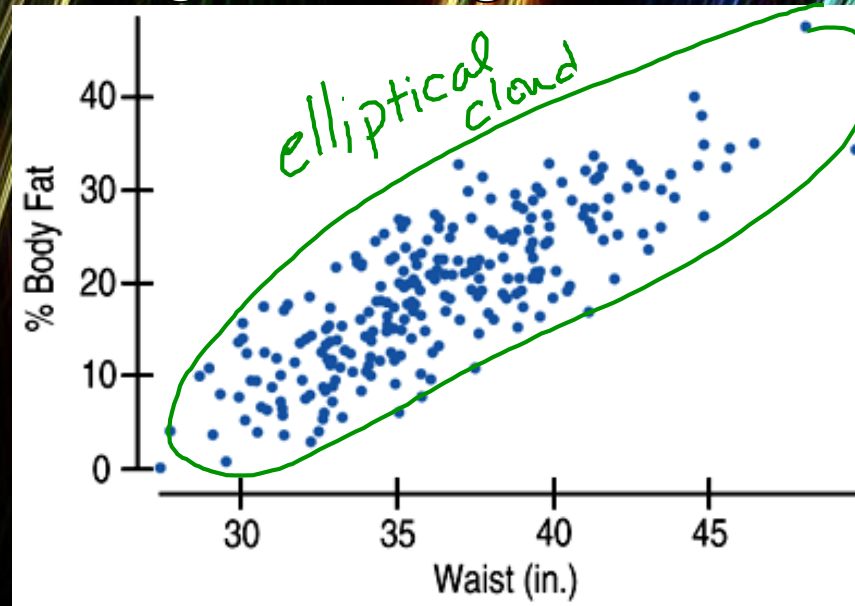
# Quantitative Data Condition

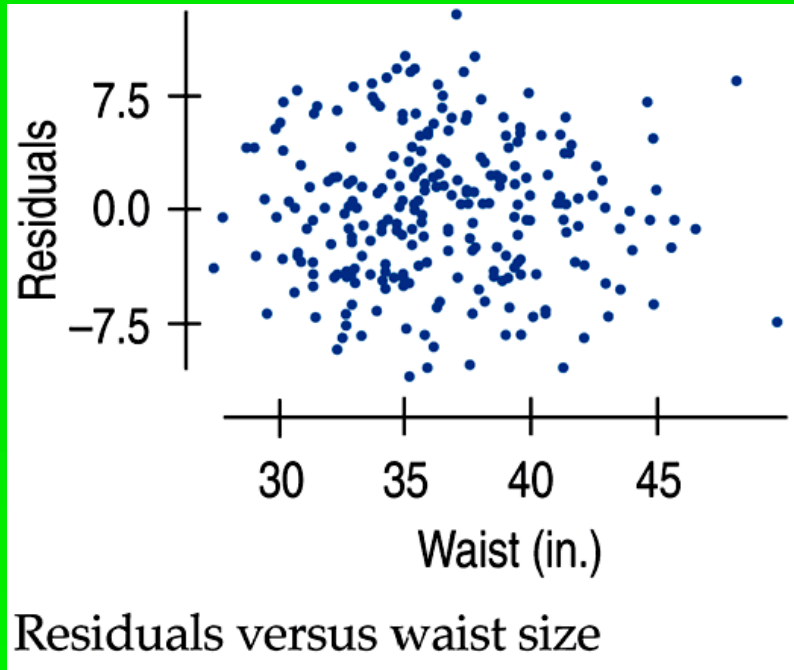The data must be quantitative for this to make sense.

# Straight Enough Condition

**Check the scatterplot—the shape must be linear or we can't use regression at all.**
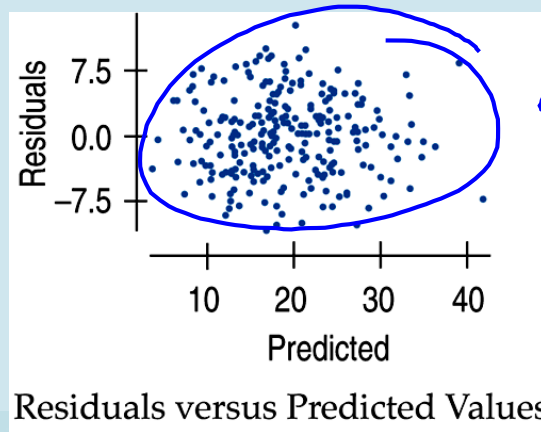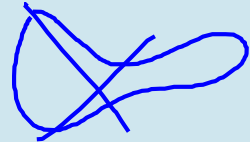
**Straight Enough Condition**

# Randomization Condition

Check the residual plot (part 1)—the residuals should appear to be randomly scattered.



Residuals versus waist size

# DOES THE PLOT THICKEN? CONDITION

Check the residual plot again - the spread of the residuals should be uniform.
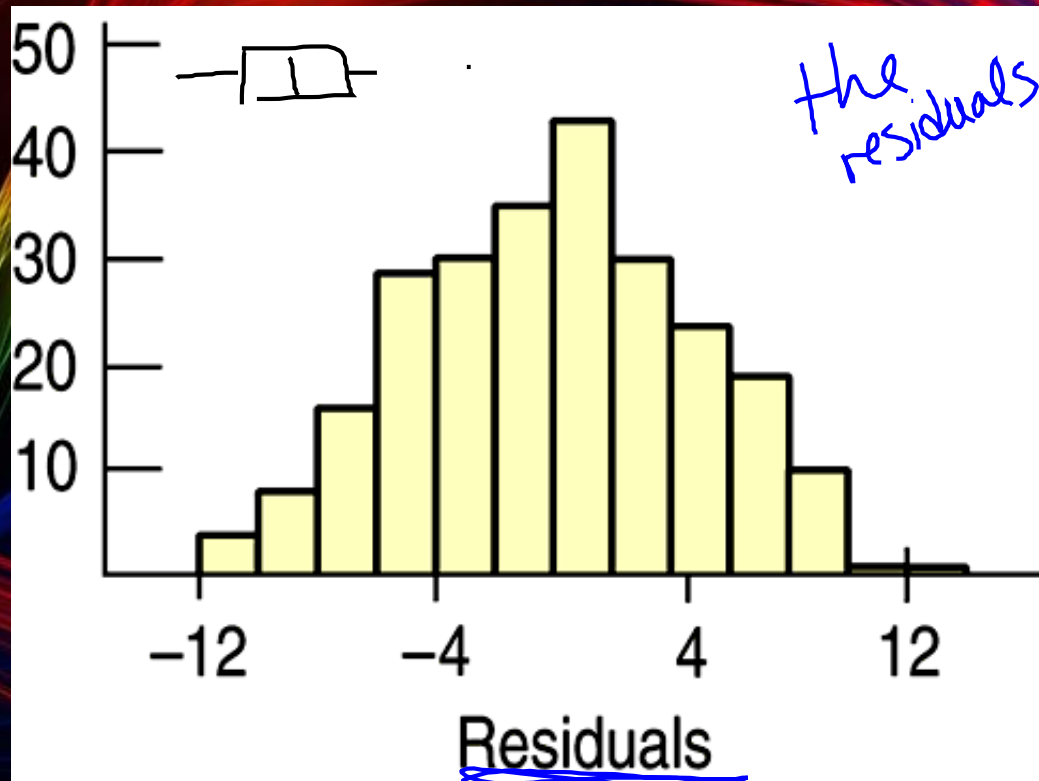


Residuals versus Predicted Values

# Nearly Normal Condition

**Check a histogram of the residuals. The distribution of the residuals should be unimodal and symmetric.**

# Outlier Condition:

Check for outliers.

*once with / one without outlier*

*stop or remove outlier*

# We need to finish our lunch

(table) Menu
L stat
L CI

$y = ax + b$
$b + ax$

- Confidence Interval

Lin Reg t-interval
x-list = time
y-list = calories
95% CL

(Slope)

$(-6.345, -1.049)$
Margin of Error = 2.648
Slope = -3.68
df = 17

# We need to finish our lunch

menu
↳ Stats
↳ tests

$\leftarrow .95 \rightarrow$
$0.025$

- **Hypothesis test**

slope of population relationship

$$H_a: \beta < 0$$

$$H_0: \beta = 0$$

Lin Reg  t-test

$$t_{17} = -2.946$$

$$p\text{-value} = 0.004$$

Due to a low p-value of 0.004, less than $\alpha$ of 0.025 we reject the null. There is statistical evidence that the slope relating time and calories is negative.

# MINITAB output

- The dataset "Healthy Breakfast" contains, among other variables, the *Consumer Reports* ratings of 77 cereals and the number of grams of sugar contained in each serving. (*Data source: Free publication available in many grocery stores. Dataset available through the* [Statlib Data and Story Library (DASL)](#).)

- Under the equation for the regression line, the output provides the least-squares estimate for the constant $b_0$ and the slope $b_1$. Since $b_1$ is the coefficient of the explanatory variable "Sugars," it is listed under that name. The calculated standard deviations for the intercept and slope are provided in the second column.

- We are comparing sugars and calories in each cereal.

# MINITAB output

| Predictor | Coef | StDev | T | P |
|-----------|--------|--------|------|-------|
| Constant | 80.81 | 56.04 | 1.44 | 0.187 |
| Calories | 2.4715 | 0.4072 | 6.07 | 0.000 |

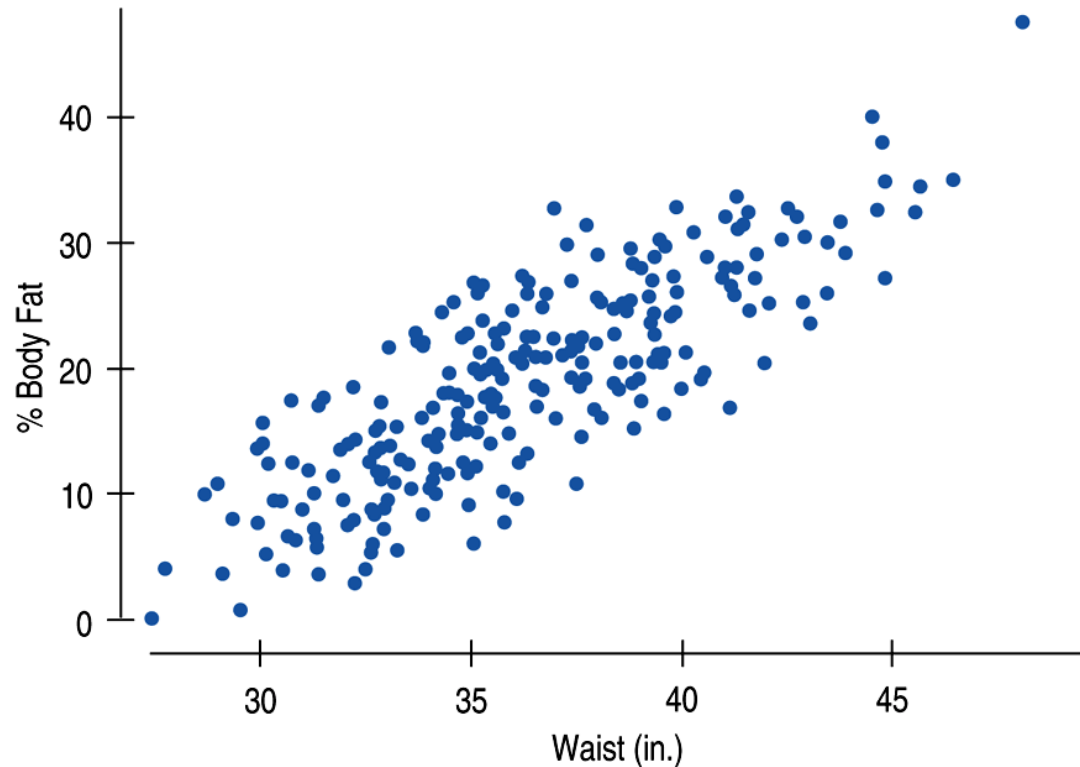S = 4.15116     R-Sq = 6.8%     R-Sq(adj) = 0.0%

# Homework

- Read chapter 27
  - take notes on formulas

# An Example: Body Fat and Waist Size

- Our chapter example revolves around the relationship between *% body fat* and *waist size* (in inches). Here is a scatterplot of our data set:

# The Population and the Sample

- When we found a confidence interval for a mean, we could imagine a single, true underlying value for the mean.

- When we tested whether two means or two proportions were equal, we imagined a true underlying difference.

- What does it mean to do inference for regression?

# The Population and the Sample (cont.)

- We know better than to think that even if we knew every population value, the data would line up perfectly on a straight line.

- In our sample, there's a whole distribution of *%body fat* for men with 38-inch waists:

# The Population and the Sample (cont.)

- This is true at each waist size.

- We could depict the distribution of *%body fat* at different *waist* sizes like this:

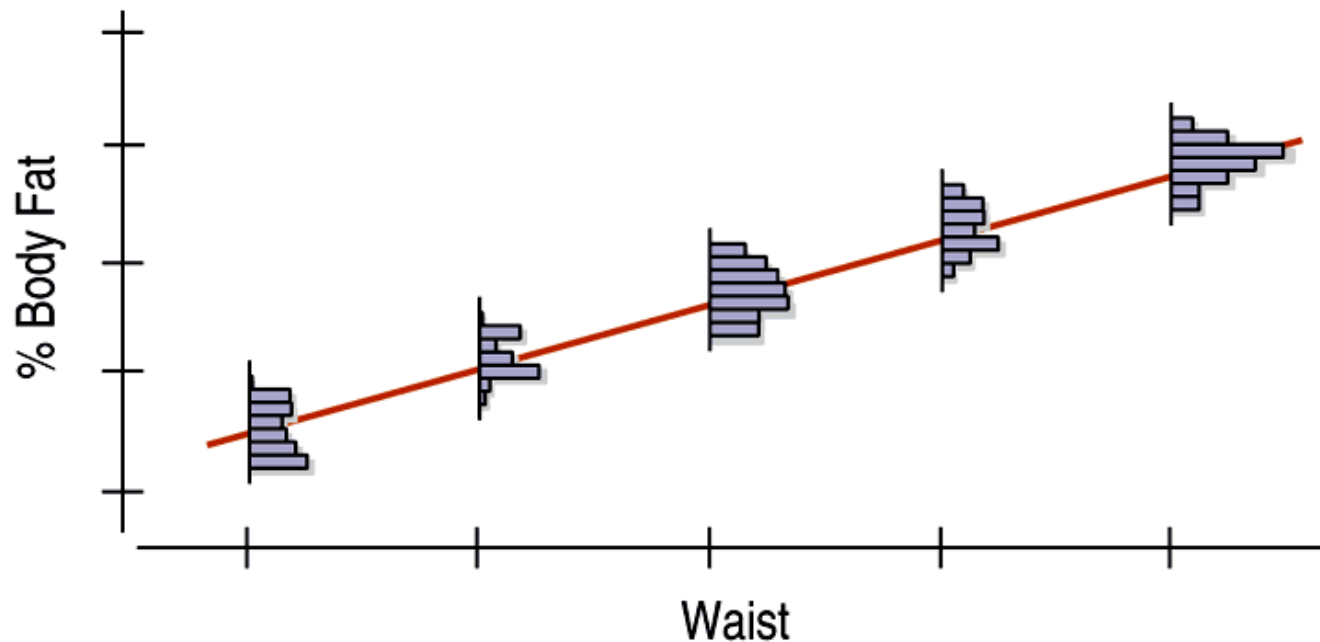# The Population and the Sample (cont.)

- The model assumes that the *means* of the distributions of *%body fat* for each *waist* size fall along the line even though the individuals are scattered around it.

- The model is not a perfect description of how the variables are associated, but it may be useful.

- If we had all the values in the population, we could find the slope and intercept of the *idealized regression line* explicitly by using least squares.

# The Population and the Sample (cont.)

- We write the idealized line with Greek letters and consider the coefficients to be *parameters*: $\beta_0$ is the intercept and $\beta_1$ is the slope.

- Corresponding to our fitted line of                    , we write

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

- Now, not all the individual *y*'s are at these means—some lie above the line and some below. Like all models, there are errors.

$$\mu_y=\beta_0+\beta_1 x$$

# The Population and the Sample (cont.)
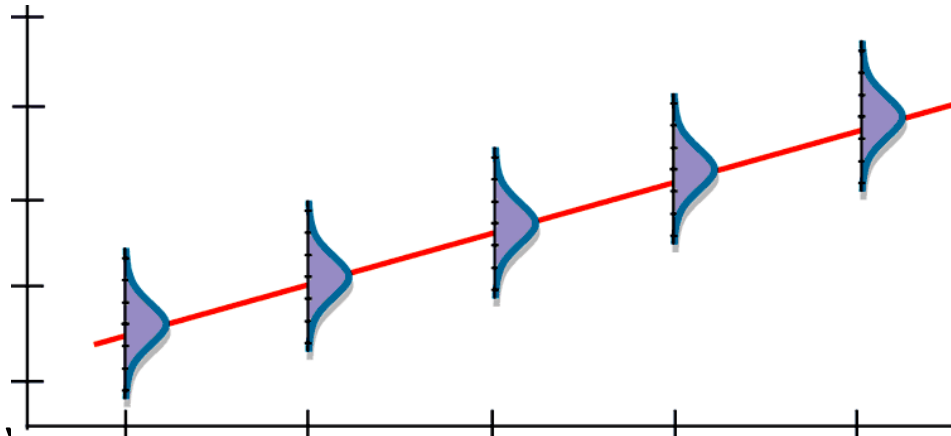
- Denote the errors by $\varepsilon$. These errors are random, of course, and can be positive or negative.

- When we add error to the model, we can talk about individual *y*'s instead of means:

This equation is now true for each data point (since there is an $\varepsilon$ to soak up the deviation) and gives a value of *y* for each *x*.

$$y = \beta_0 + \beta_1 x + \varepsilon$$

# Assumptions and Conditions (cont.)

- If all four assumptions are true, the idealized regression model would look like this:



- At each value of x there is a distribution of y-values that follows a Normal model, and each of these Normal models is centered on the line and has the same standard deviation.

# Which Come First:
# the Conditions or the Residuals?

- There's a catch in regression—the best way to check many of the conditions is with the residuals, but we get the residuals only *after* we compute the regression model.

- To compute the regression model, however, we should check the conditions.

- So we work in this order:
  1. Make a scatterplot of the data to check the Straight Enough Condition. (If the relationship isn't straight, try re-expressing the data. Or stop.)

# Which Come First:
# the Conditions or the Residuals? (cont.)

2. If the data are straight enough, fit a regression model and find the residuals, $e$, and predicted values, $\hat{y}$.

3. Make a scatterplot of the residuals against $x$ or the predicted values.
   - This plot should have no pattern. Check in particular for any bend, any thickening (or thinning), or any outliers.

4. If the data are measured over time, plot the residuals against time to check for evidence of patterns that might suggest they are not independent.

# Which Come First:
# the Conditions or the Residuals? (cont.)

5. If the scatterplots look OK, then make a histogram and Normal probability plot of the residuals to check the Nearly Normal Condition.

6. If all the conditions seem to be satisfied, go ahead with inference.

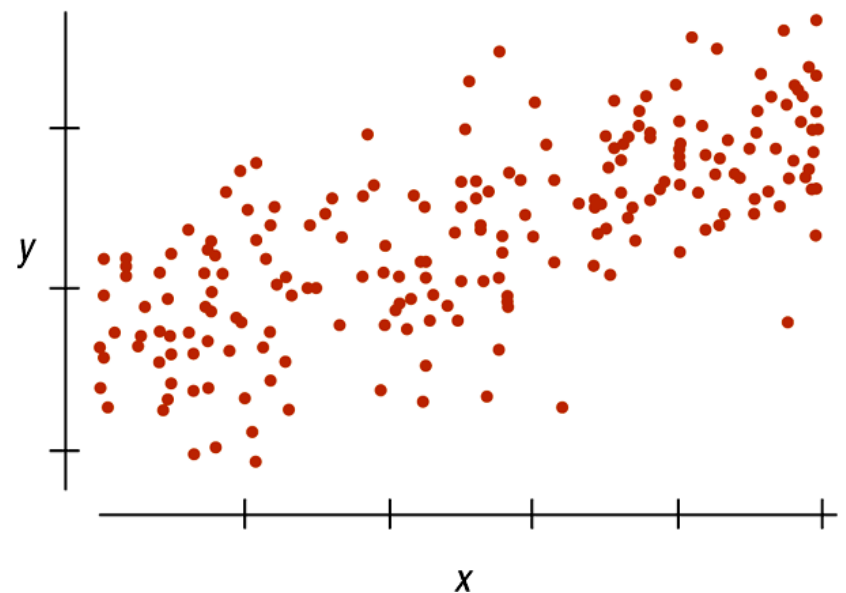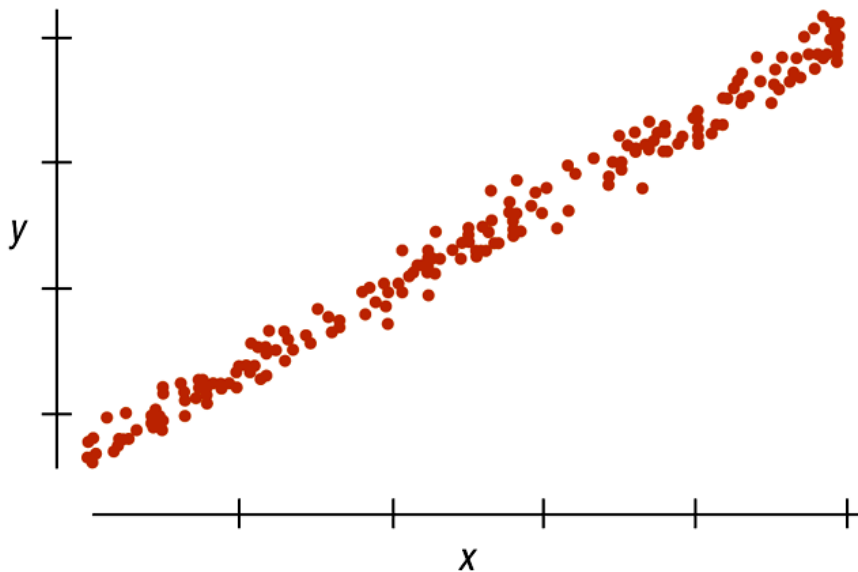# Intuition About Regression Inference

- We expect any sample to produce a $b_1$ whose expected value is the true slope, $\beta_1$.

- What about its standard deviation?

- What aspects of the data affect how much the slope and intercept vary from sample to sample?

# Intuition About Regression Inference (cont.)

- Spread around the line:
  - Less scatter around the line means the slope will be more consistent from sample to sample.
  - The spread around the line is measured with the residual standard deviation $s_e$.
  - You can always find $s_e$ in the regression output, often just labeled $s$.

# Intuition About Regression Inference (cont.)
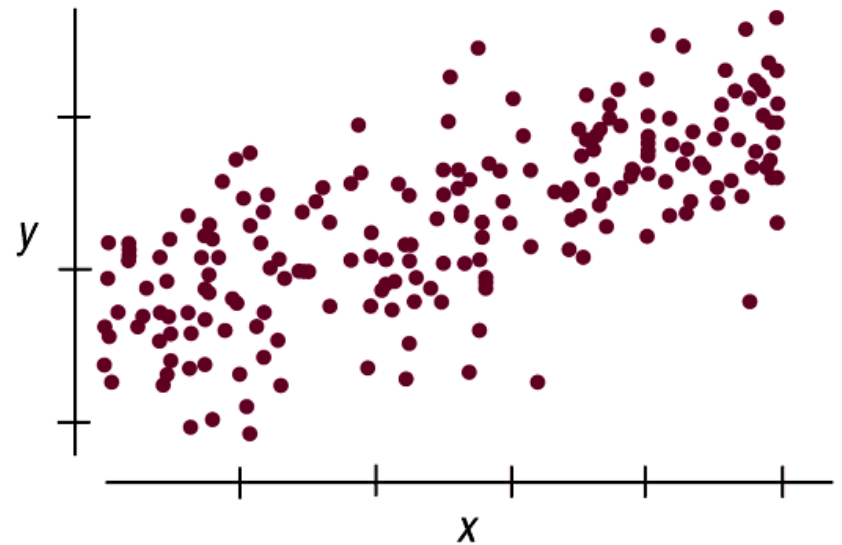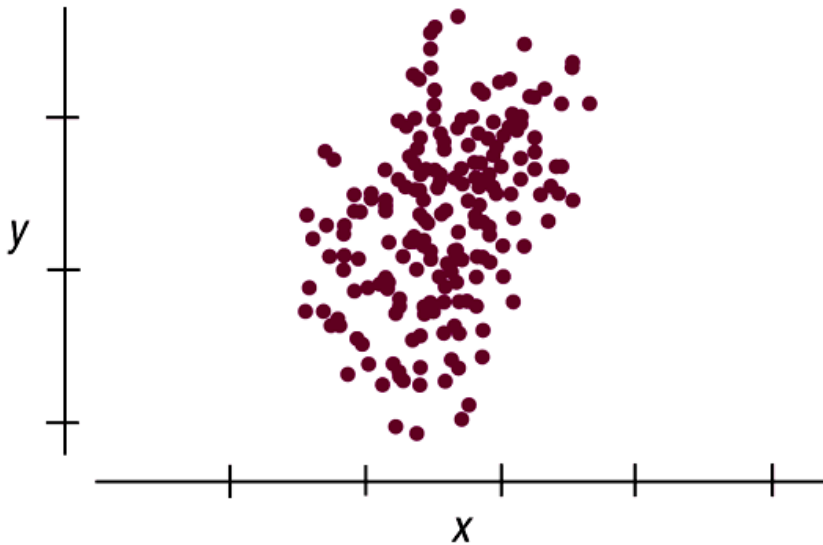
- Spread around the line:



Less scatter around the line means the slope will be more consistent from sample to sample.

# Intuition About Regression Inference (cont.)

- Spread of the *x*'s: A large standard deviation of *x* provides a more stable regression.

# Intuition About Regression Inference (cont.)

- Sample size: Having a larger sample size, *n*, gives more consistent estimates.

# Standard Error for the Slope

- Three aspects of the scatterplot affect the standard error of the regression slope:
  - spread around the line, $s_e$
  - spread of $x$ values, $s_x$
  - sample size, $n$.
- The formula for the standard error (which you will probably never have to calculate by hand) is:

$$SE(b_1) = \frac{s_e}{\sqrt{n-1}\, s_x}$$

# Sampling Distribution for Regression Slopes

- When the conditions are met, the standardized estimated regression slope

$$t = \frac{b_1 - \beta_1}{SE(b_1)}$$

follows a Student's $t$-model with $n - 2$ degrees of freedom.

# Sampling Distribution for Regression Slopes (cont.)

- We estimate the standard error with

$$SE(b_1) = \frac{s_e}{\sqrt{n-1}\,s_x}$$

where:

$$s_e = \sqrt{\frac{\sum(y-\hat{y})^2}{n-2}}$$

- $n$ is the number of data values
- $s_x$ is the ordinary standard deviation of the $x$-values.

# What About the Intercept?

- The same reasoning applies for the intercept.

- We can write $$\frac{b_0 - \beta_0}{SE(b_0)} : t_{n-2}$$
  but we rarely use this fact for anything.

- The intercept usually isn't interesting. Most hypothesis tests and confidence intervals for regression are about the slope.

# Regression Inference

- A null hypothesis of a zero slope questions the entire claim of a linear relationship between the two variables—often just what we want to know.

- To test $H_0: \beta_1 = 0$, we find

$$t_{n-2} = \frac{b_1 - 0}{SE(b_1)}$$

and continue as we would with any other *t*-test.

- The formula for a confidence interval for $\beta_1$ is

$$b_1 \pm t^*_{n-2} \times SE(b_1)$$

# *Standard Errors for Predicted Values

- Once we have a useful regression, how can we indulge our natural desire to predict, without being irresponsible?

- Now we have standard errors—we can use those to construct a confidence interval for the predictions, smudging the results in the right way to report our uncertainty honestly.

# *Standard Errors for Predicted Values (cont.)

- For our *%body fat* and *waist* size example, there are two questions we could ask:
  - Do we want to know the mean *%body fat* for *all* men with a *waist* size of, say, 38 inches?
  - Do we want to estimate the *%body fat* for a particular man with a 38-inch *waist*?

- The predicted *%body fat* is the same in both questions, but we can predict the *mean %body fat* for *all* men whose *waist* size is 38 inches with a lot more precision than we can predict the *%body fat* of a *particular individual* whose *waist* size happens to be 38 inches.

# *Standard Errors for Predicted Values (cont.)

- We start with the same prediction in both cases.
  - We are predicting for a new individual, one that was not in the original data set.
  - Call his *x*-value $x_v$ (38 inches).
  - The regression predicts *%body fat* as

$$\hat{y}_v = b_0 + b_1 x_v$$

# *Standard Errors for Predicted Values (cont.)

- Both intervals take the form

$$\hat{y}_v \pm t^*_{n-2} \times SE$$

- The *SE*'s will be different for the two questions we have posed.

# *Standard Errors for Predicted Values (cont.)

- The standard error of the *mean* predicted value is:

$$SE(\widehat{\mu}_v) = \sqrt{SE^2(b_1) \cdot (x_v - \bar{x})^2 + \frac{s_e^2}{n}}$$

- Individuals vary more than means, so the standard error for a single predicted value is larger than the standard error for the mean:
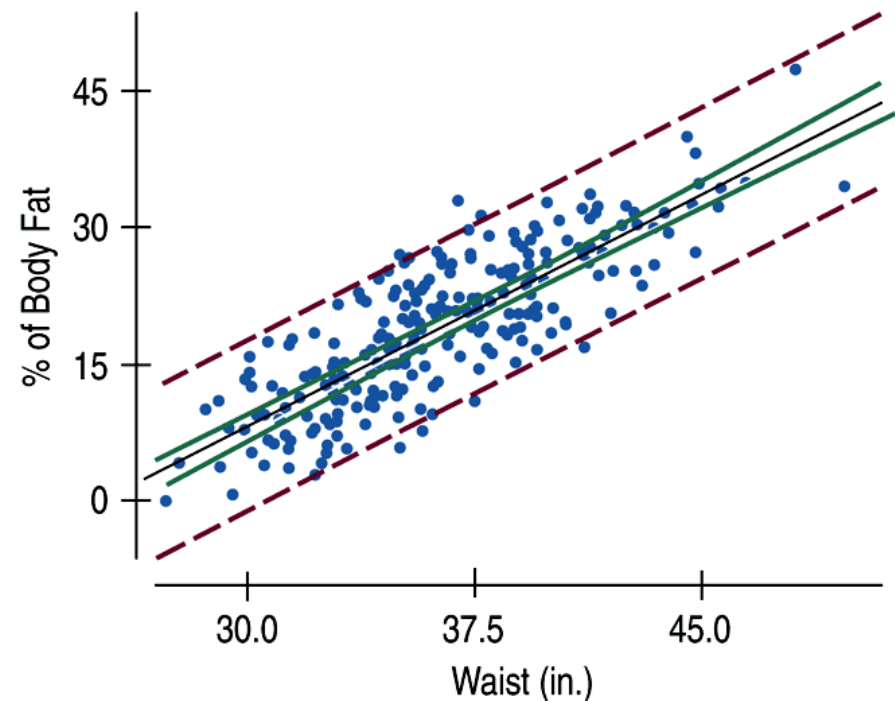
$$SE(\widehat{y}_v) = \sqrt{SE^2(b_1) \cdot (x_v - \bar{x})^2 + \frac{s_e^2}{n} + s_e^2}$$

# *Standard Errors for Predicted Values (cont.)

- Keep in mind the distinction between the two kinds of confidence intervals.
  - The narrower interval is a **confidence interval for the predicted mean value** at $x_v$
  - The wider interval is a **prediction interval for an individual** with that x-value.

# *Confidence Intervals for Predicted Values

- Here's a look at the difference between predicting for a mean and predicting for an individual.

- The solid green lines near the regression line show the 95% confidence interval for the mean predicted value, and the dashed red lines show the prediction intervals for individuals.

# What Can Go Wrong?

- Don't fit a linear regression to data that aren't straight.

- Watch out for the plot thickening.
  - If the spread in *y* changes with *x*, our predictions will be very good for some *x*-values and very bad for others.

- Make sure the errors are Normal.
  - Check the histogram and Normal probability plot of the residuals to see if this assumption looks reasonable.

# What Can Go Wrong? (cont.)

- Watch out for extrapolation.
    - It's always dangerous to predict for *x*-values that lie far from the center of the data.

- Watch out for high-influence points and outliers.

- Watch out for one-tailed tests.
    - Tests of hypotheses about regression coefficients are usually two-tailed, so software packages report two-tailed P-values.
    - If you are using software to conduct a one-tailed test about slope, you'll need to divide the reported P-value in half.

# What have we learned?

- We have now applied inference to regression models.

We've learned:

- Under certain assumptions, the sampling distribution for the slope of a regression line can be modeled by a Student's $t$-model with $n - 2$ degrees of freedom.

- To check four conditions, in order, to verify the assumptions. Most checks can be made by graphing the data and residuals.

# What have we learned?

- To use the appropriate $t$-model to test a hypothesis about the slope. If the slope of the regression line is significantly different from 0, we have strong evidence that there is an association between the two variables.

- To create and interpret a confidence interval or the true slope.

- We have been reminded yet again never to mistake the presence of an association for proof of causation.