Agenda

Section 12.1
Day 1

Tomorrow

MORE AP Review

MONDAY
Continue with
Section 12.1



Chapter 3: Describing Relationships
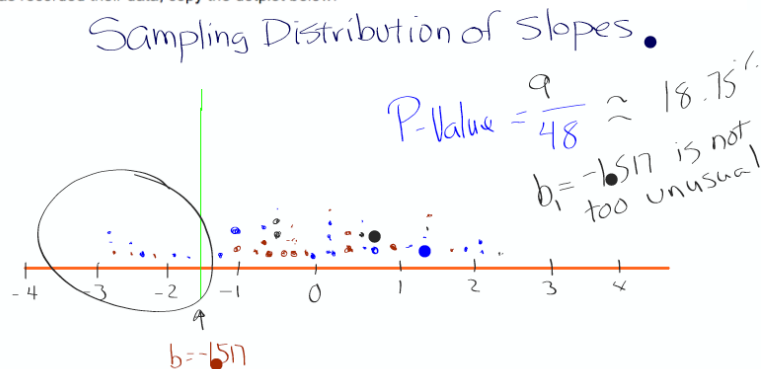
We need to find out how likely it is that this slope occurs purely by chance if row has no affect on scores.

You have now calculated three different possible values for the slope based on random assignment. Take these 3 values to the dotplot on the whiteboard in the front of the room. When everyone in class has recorded their data, copy the dotplot below:
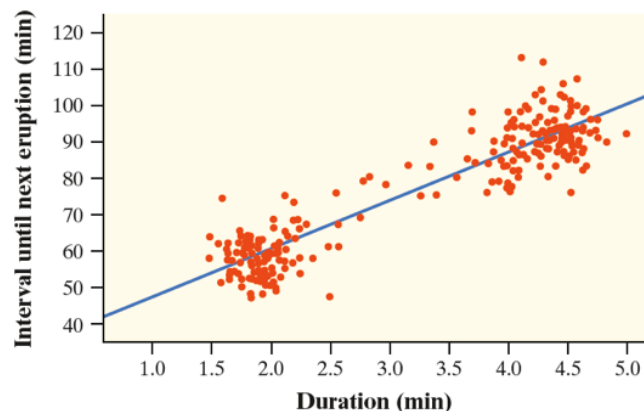
Sampling Distribution of Slopes.

$$P\text{-Value} = \frac{9}{48} \approx 18.75\%$$

$b_1 = -1.517$ is not too unusual

b = -1.517

Different samples of the same size from the same population will give different estimates for the slope

## Learning Targets

✓ CHECK the conditions for performing inference about the slope $\beta_1$ of the population (true) regression line.

✓ INTERPRET the values of $\beta_0$, $\beta_1$, $\sigma$, and $SE_{b_1}$ in context, and DETERMINE these values from computer output.
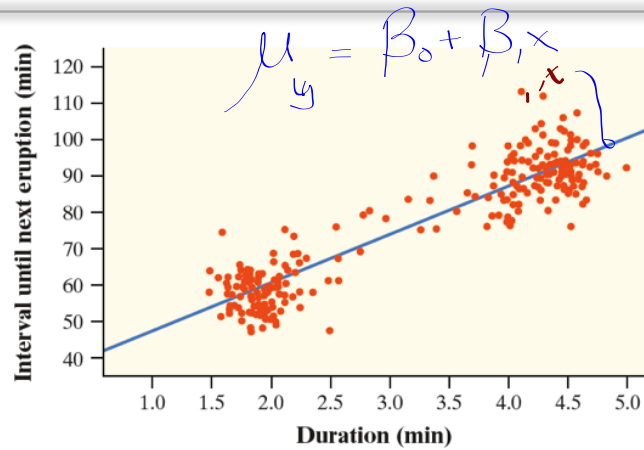
## Inference for Linear Regression



Scatterplot of the duration and interval between eruptions of Old Faithful for all 263 eruptions in a single month. The **population least-squares regression line** is shown in blue.

# Inference for Linear Regression

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

Sample

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



Scatterplot of the duration and interval between eruptions of Old Faithful for all 263 eruptions in a single month. The **population least-squares regression line** is shown in blue.
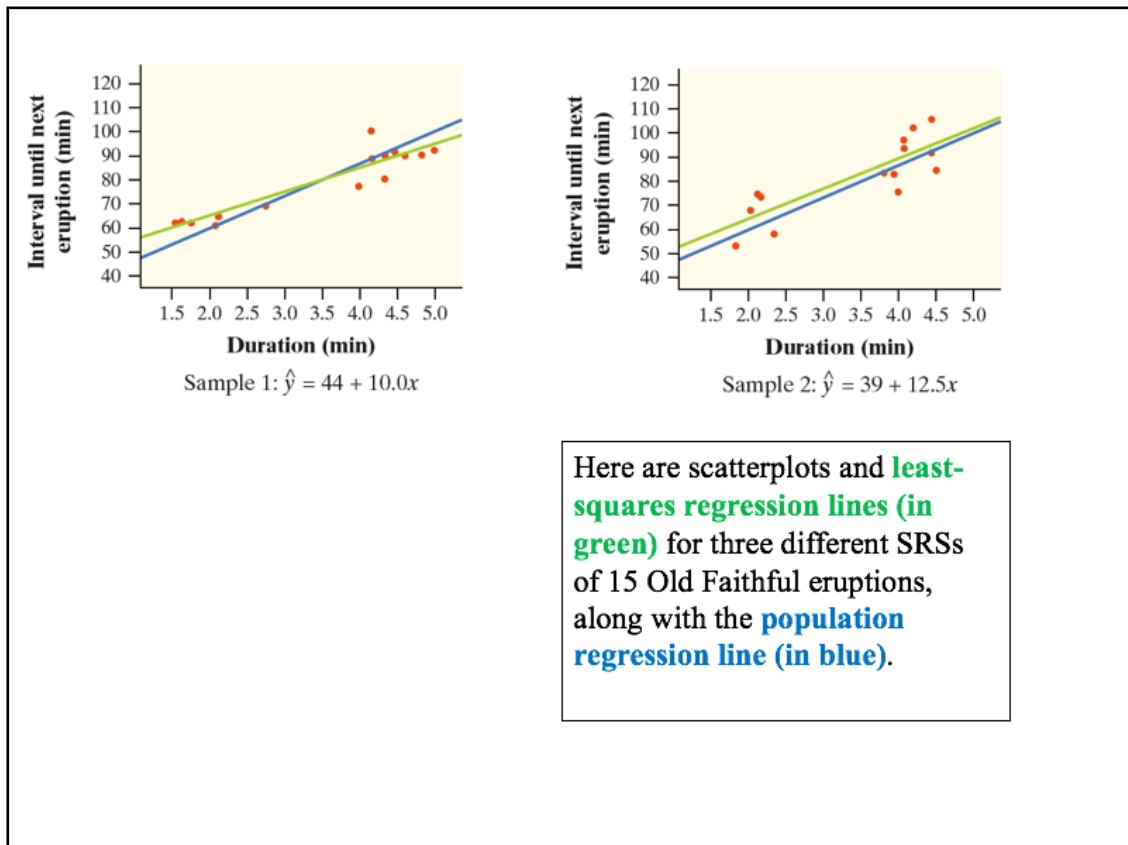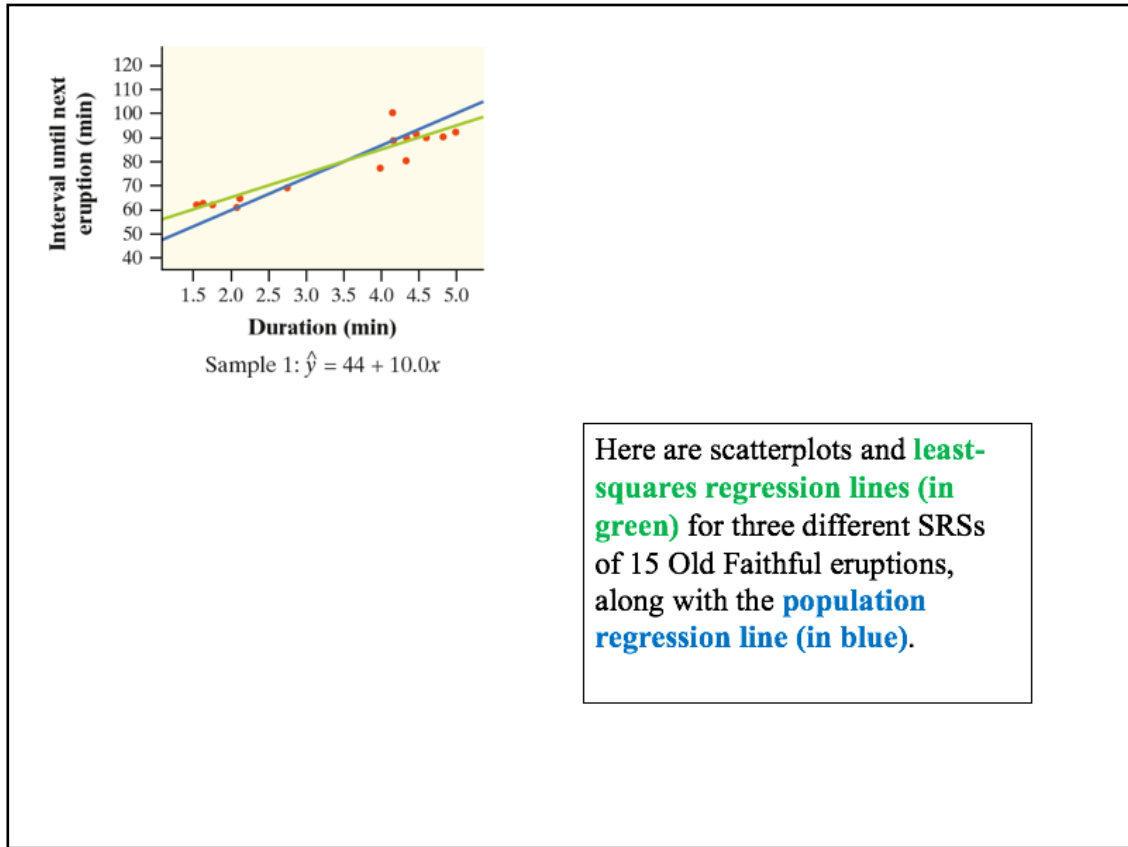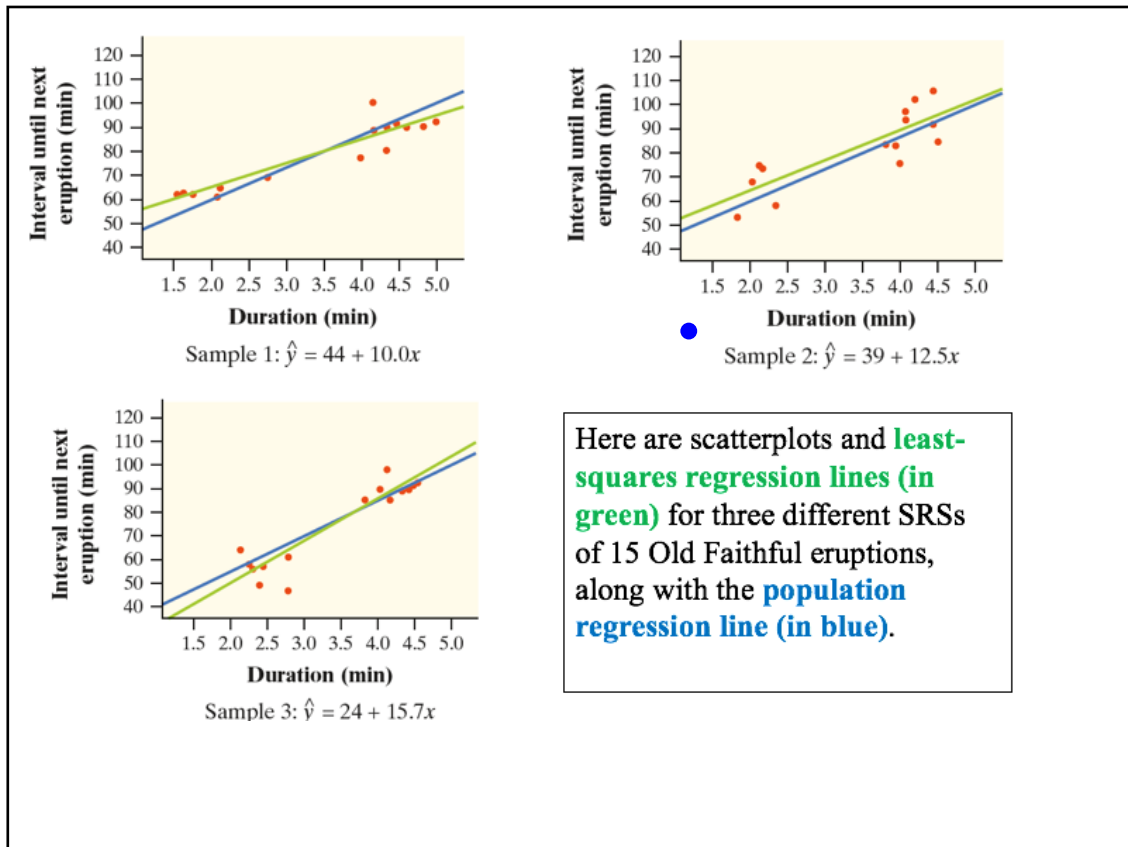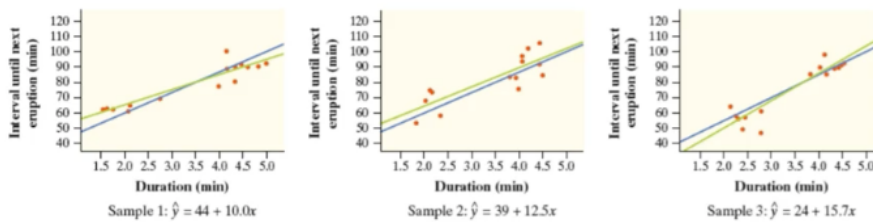
Here are scatterplots and **least-squares regression lines (in green)** for three different SRSs of 15 Old Faithful eruptions, along with the **population regression line (in blue).**

Sample 1: $\hat{y} = 44 + 10.0x$

Here are scatterplots and **least-squares regression lines (in green)** for three different SRSs of 15 Old Faithful eruptions, along with the **population regression line (in blue).**



Sample 1: $\hat{y} = 44 + 10.0x$

Sample 2: $\hat{y} = 39 + 12.5x$

Sample 1: $\hat{y} = 44 + 10.0x$

Sample 2: $\hat{y} = 39 + 12.5x$

Sample 3: $\hat{y} = 24 + 15.7x$

Here are scatterplots and **least-squares regression lines (in green)** for three different SRSs of 15 Old Faithful eruptions, along with the **population regression line (in blue).**
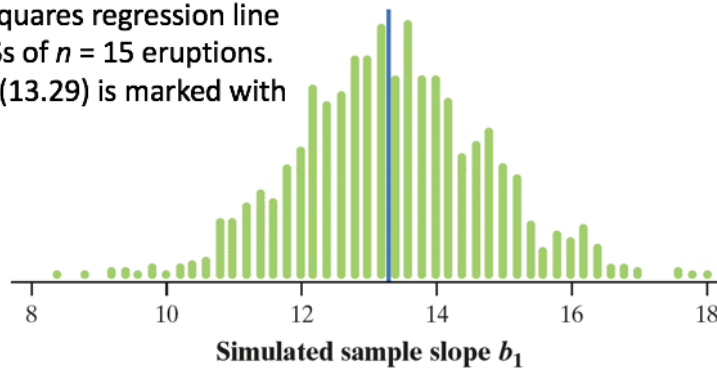
## Different samples of the same size from the same population will give different estimates for the slope.



Sample 1: $\hat{y} = 44 + 10.0x$

Sample 2: $\hat{y} = 39 + 12.5x$
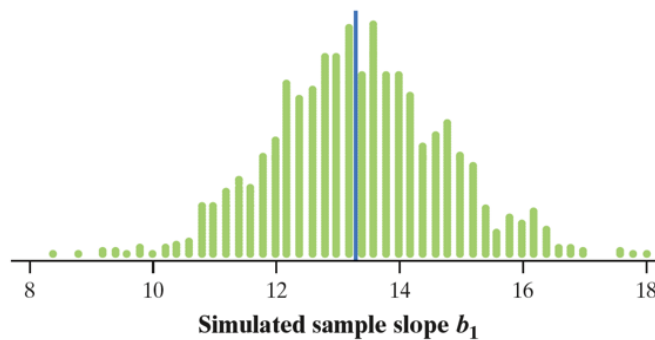
Sample 3: $\hat{y} = 24 + 15.7x$

The pattern is described by its sampling distribution

The figure shows a dotplot of the sample slope $b_1$ of the least-squares regression line in 1000 simulated SRSs of $n = 15$ eruptions. The population slope (13.29) is marked with a **blue vertical line**.
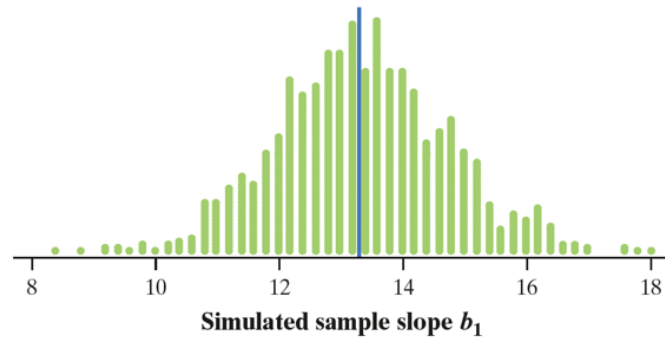


Simulated sample slope $b_1$

## Sampling Distribution of $b_1$

which estimates the population slope $\beta_1$



Simulated sample slope $b_1$
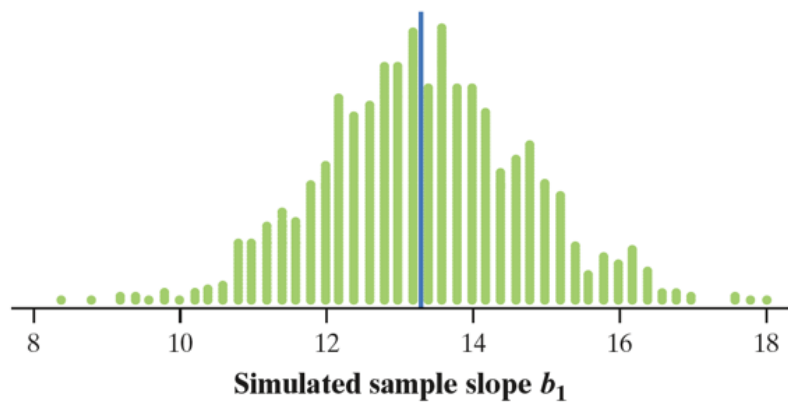
**Shape:** Approximately Normal



**Shape:** Approximately Normal

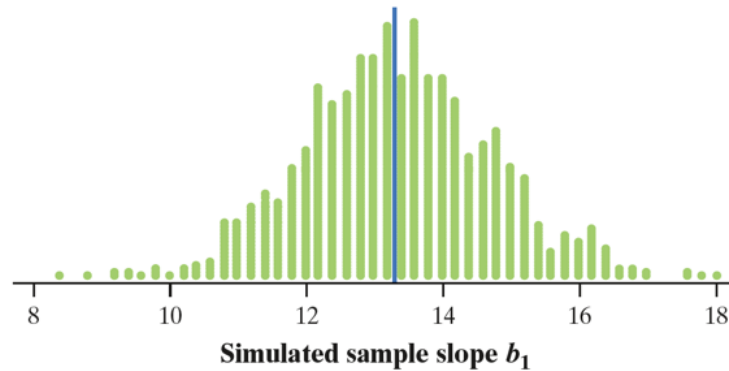**Center:** $\mu_{b_1} = \beta_1 = 13.29$ ($b_1$ is an unbiased estimator of $\beta_1$.)

**Shape:** Approximately Normal

**Center:** $\mu_{b_1} = \beta_1 = 13.29$ ($b_1$ is an unbiased estimator of $\beta_1$.)

**Variability:** $\sigma_{b_1} = \dfrac{\sigma}{\sigma_x\sqrt{n}} = \dfrac{6.47}{1.18\sqrt{15}} = 1.42$, where $\sigma_x$ is the standard deviation of duration for the 263 eruptions.



Simulated sample slope $b_1$

---

When certain conditions are met, we can anticipate the shape, center, and variability of the sampling distribution of the sample slope.



(b)                Simulated sample slope $b_1$

|  | We used | to estimate |
|---|---|---|
| Ch. 8 | $\hat{P}$ | $P$ |
|  | $\overline{X}$ | $\mu$ |
| Ch. 10 | $(\hat{P_1} - \hat{P_1})$ | $(P_1 - P_2)$ |
|  | $(\overline{X_1} - \overline{X_2})$ | $(\mu_1 - \mu_2)$ |
| Ch. 12 | $\hat{y} = b_0 + b_1 x$ | $\mu_y = \beta_0 + \beta_1 x$ |

Sample regression line       population regression line

we'll focus on estimating slope $\beta_1$

---

Big Idea: If the data come from a random sample or randomized experiment, the least-squares regression line is just an *estimate* of the population (true) least-squares regression line.

- Population line: $\mu_y = \beta_0 + \beta_1 x$
- Sample line: $\hat{y} = b_0 + b_1 x$

handout

---

A regression line calculated from every value in the population is called a **population regression line** (true regression line). The equation of a population regression line is $\mu_y = \beta_0 + \beta_1 x$ where

- $\mu_y$ is the mean *y*-value for a given value of *x*.
- $\beta_0$ is the population *y* intercept.
- $\beta_1$ is the population slope.

Population

A regression line calculated from a sample is called a **sample regression line** (estimated regression line). The equation of a sample regression line is $\hat{y} = b_0 + b_1 x$ where

- $\hat{y}$ is the estimated mean *y*-value for a given value of *x*.
- $b_0$ is the sample *y* intercept.
- $b_1$ is the sample slope.

Sample

### Sampling Distribution of a Slope

Choose an SRS of $n$ observations $(x, y)$ from a population of size $N$ with least-squares regression line $\mu_y = \beta_0 + \beta_1 x$. Let $b_1$ be the slope of the sample regression line. Then:

- The **mean** of the sampling distribution of $b_1$ is $\mu_{b_1} = \beta_1$.
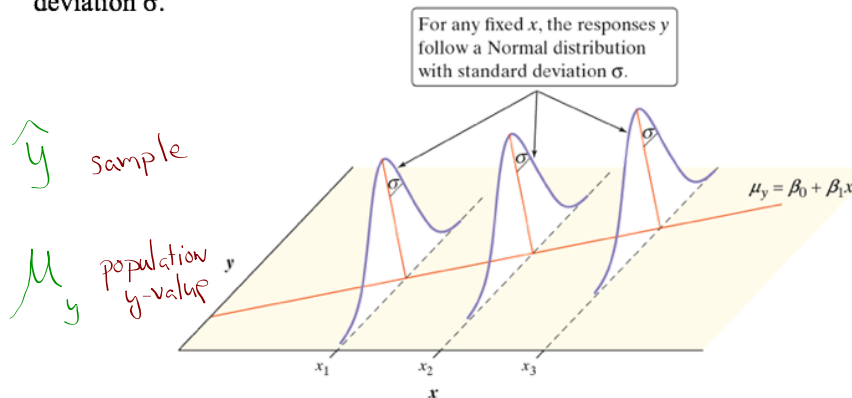- The **standard deviation** of the sampling distribution of $b1$ is

$$\sigma_{b_1} = \frac{\sigma}{\sigma_x \sqrt{n}}$$

  as long as the *10% condition* is satisfied: $n < 0.10\ N$.
- The sampling distribution of $b_1$ will be **approximately Normal** if the values of the response variable $y$ follow a Normal distribution for each value of the explanatory variable $x$ (the *Normal condition*).
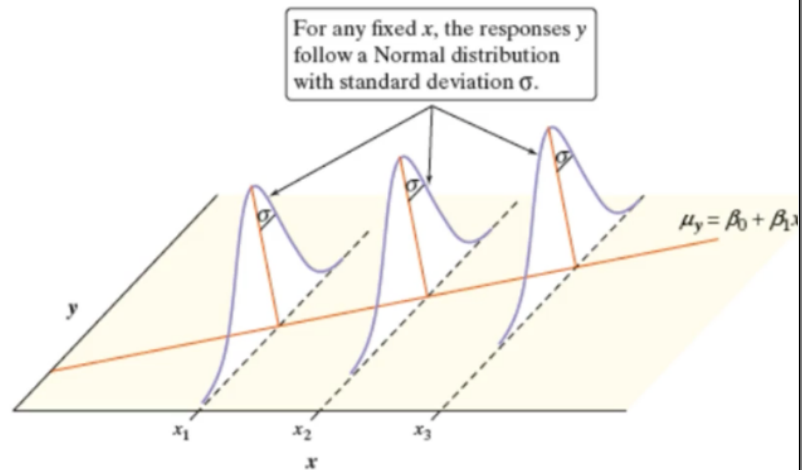
# Conditions for Regression Inference

When the conditions for inference are met, the regression model looks like the one shown here. The line is the population (true) regression line, which shows how the mean response $\mu_y$ changes as the explanatory variable $x$ changes. For any fixed value of $x$, the observed response $y$ varies according to a Normal distribution having mean $\mu_y$ and standard deviation $\sigma$.



$\hat{y}$  sample

$\mu_y$  population y-value

- The conditions are based on the model for linear regression:

> For any fixed $x$, the responses $y$ follow a Normal distribution with standard deviation $\sigma$.

- **L**inear
- **I**ndependent
- **N**ormal
- **E**qual SD
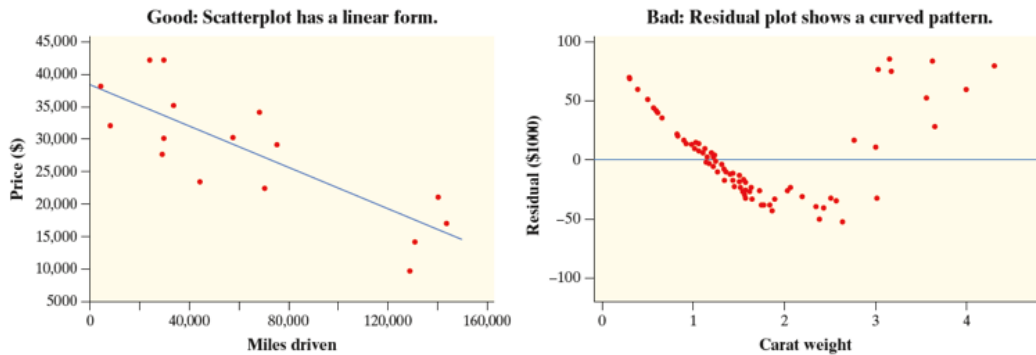- **R**andom

$\mu_y = \beta_0 + \beta_1 x$

---

| Conditions for Regression Inference |
|---|

Suppose we have $n$ observations on a quantitative explanatory variable $x$ and a quantitative response variable $y$. Our goal is to study or predict the behavior of $y$ for given values of $x$.

- **L**inear: The actual relationship between $x$ and $y$ is linear. For any fixed value of $x$, the mean response m$y$ falls on the population (true) regression line $\mu_y = \beta_0 + \beta_1 x$.
- **I**ndependent: Individual observations are independent of each other. When sampling without replacement, check the *10% condition*.
- **N**ormal: For any fixed value of $x$, the response $y$ varies according to a Normal distribution.
- **E**qual SD: The standard deviation of $y$ (call it $\sigma$) is the same for all values of $x$.
- **R**andom: The data come from a random sample from the population of interest or a randomized experiment.

## Here's a summary of how to check the conditions one by one.

**Linear:** Examine the scatterplot to see if the overall pattern is roughly linear. Make sure there are no leftover curved patterns in the residual plot.
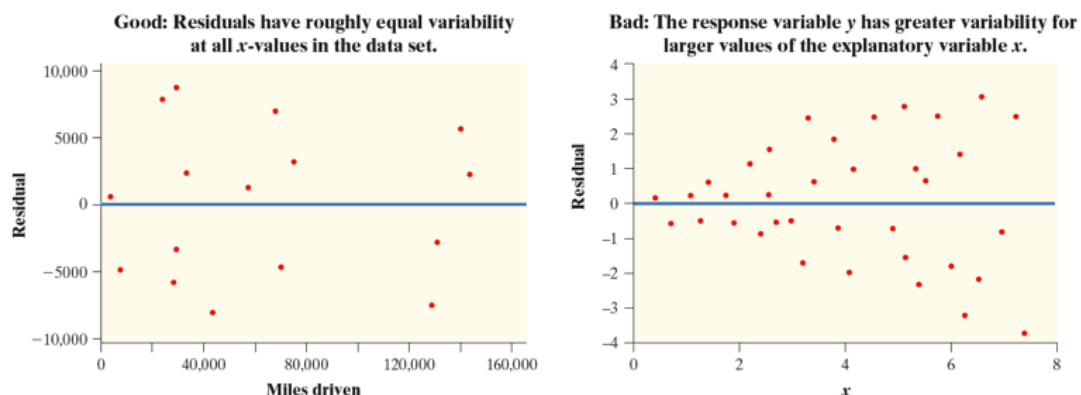
**Good: Scatterplot has a linear form.**

**Bad: Residual plot shows a curved pattern.**

Here's a summary of how to check the conditions one by one.

**Independent:** Knowing the value of the response variable for one individual shouldn't help predict the value of the response variable for other individuals. If sampling is done without replacement, remember to check that the sample size is less than 10% of the population size (*10% condition*).

**Normal:** Make a histogram, dotplot, stemplot, boxplot, or Normal probability plot of the residuals and check for strong skewness or outliers. Ideally, we would check the Normality of the residuals at each possible value of $x$. Because we rarely have enough observations at each $x$-value, however, we make one graph of all the residuals to check for Normality.

**Equal SD:** Look at the scatter of the residuals above and below the "residual = 0" line in the residual plot. The variability of the residuals in the vertical direction should be roughly the same from the smallest to the largest $x$-value.

Good: Residuals have roughly equal variability at all $x$-values in the data set.

Bad: The response variable $y$ has greater variability for larger values of the explanatory variable $x$.

Random: See if the data came from a random sample from the population of interest or a randomized experiment. If not, we can't make inferences about a larger population or about cause and effect.

## Check Your Understanding

Mrs. Barrett's class did a fun experiment using paper helicopters. After making 70 helicopters using the same template, students randomly assigned 14 helicopters to each of five drop heights: 152 cm, 203 cm, 254 cm, 307 cm, and 442 cm.

Teams of students released the 70 helicopters in a random order and measured the flight times in seconds. The class used computer software to carry out a least-squares regression analysis for these data. Some output from this regression analysis is shown here. We checked conditions for performing inference earlier.
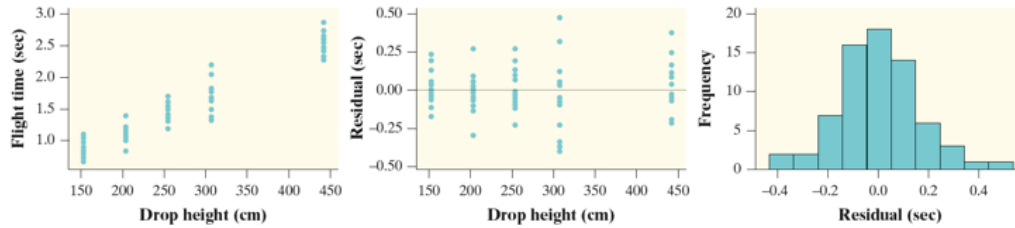
```
Regression Analysis: Flight time versus Drop height
Predictor                Coef      SE Coef      T       P
Constant               −0.03761    0.05838    −0.64   0.522
Drop height (cm)       0.0057244  0.0002018   28.37   0.000
  S = 0.168181      R-Sq = 92.2%       R-Sq(adj) = 92.1%
```
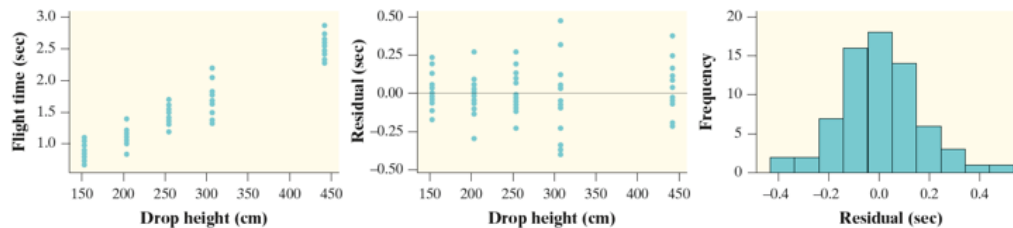
(a) With your group, discuss whether all of the conditions are met for regressions inference

**Problem:**



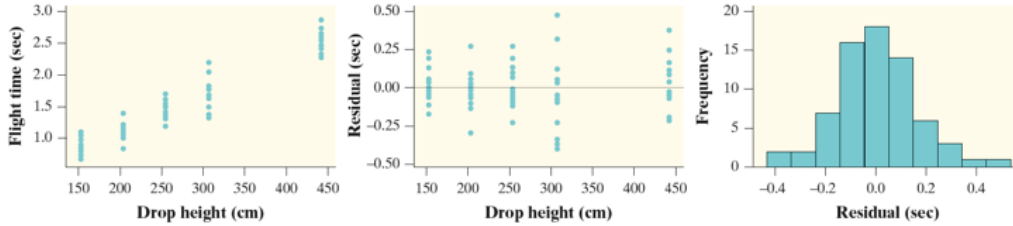Check whether the conditions for performing inference about the regression model are met.

- **Linear:** The scatterplot shows a clear linear form and there is no leftover curved pattern in the residual plot. ✓

- **Independent:** Because the helicopters were released in a random order and no helicopter was used twice, knowing the result of one observation should not help us predict the value of another observation. ✓

---

**Problem:**



Check whether the conditions for performing inference about the regression model are met.

- **Normal:** There is no strong skewness or outliers in the histogram of the residuals. ✓

- **Equal SD:** The residual plot shows a similar amount of scatter about the residual = 0 line for each drop height. However, flight times seem to vary a little more for the helicopters that were dropped from a height of 307 cm. ✓

## Problem:



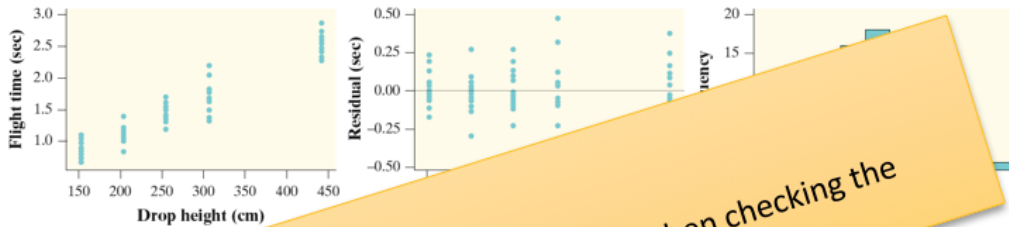Check whether the conditions for performing inference about the regression model are met.

- Random: The helicopters were randomly assigned to the five possible drop heights. ✓

(Remember the LINER acronym.)

---

## Problem:



Check whether th~~~
regres~~~

**CAUTION**: Don't overreact to minor issues in the graphs when checking the Normal and Equal SD conditions

~~~the

~~~andomly assigned to the five possible

(Remember the LINER acronym.)

**(b) What is the estimate for $\beta_0$? Interpret this value.**

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

Time to land

Drop height

y-intercept ($b_0$)

**Regression Analysis: Flight time versus Drop height**

| Predictor | Coef | SE Coef | T | P |
|---|---|---|---|---|
| Constant | -0.03761 | 0.05838 | -0.64 | 0.522 |
| Drop height (cm) | 0.0057244 | 0.0002018 | 28.37 | 0.000 |

S = 0.168181      R-Sq = 92.2%      R-Sq(adj) = 92.1%

slope ($b_1$)

$SE_{b_1}$

---

**(b) What is the estimate for $\beta_0$? Interpret this value.**

$b_0 = -0.03761$

If the helicopter is dropped from 0 cm, it is predicted to take -0.03761 seconds to land. This has no meaning.

**(b) What is the estimate for $\beta_1$? Interpret this value.**
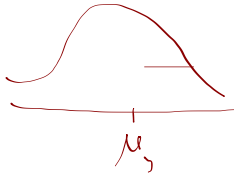
$b_1 = 0.0057244$

.0057

for every additional cm of drop height, the landing time rises by 0.0057 seconds

(c) What is the estimate for σ? Interpret this value.

$S = 0.16818$

the actual flight times typically vary by 0.168 seconds from predicted by LSRL.

(d) Give the standard error of the slope $SE_{b1}$. Interpret this value.



$SE_{b_1} =$

---

(d) Give the standard error of the slope $SE_{b1}$. Interpret this value.

$SE_{b_1} = 0.0002018$

If we repeated the random assignment many times, the slope of the sample LSRL typically varies by 0.0002 from the true slope.

See your test

12.1..... 1, 3, 5

and study pp. 769-776

I will be unavailable before school tomorrow.