

# Stat 515: Introduction to Statistics

## Chapter 6

# Recall Definitions from Ch 2

- **Statistic:** numerical summary of a sample
  - Mean( $\bar{x}$ ), proportion( $\hat{p}$ ), median, mode, standard deviation( $s$ ), variance( $s^2$ ), Q1, Q3, IQR, etc.
  - We use US alphabet letters to denote these
- **Parameter:** numerical summary of a population
  - Mean( $\mu_x$ ), proportion( $\rho$ ), median, mode, standard deviation( $\sigma$ ), variance( $\sigma^2$ ), Q1, Q3, IQR, etc.
  - We usually don't know these values
  - We use Greek letters to denote these

# Sampling Distributions

- A **sampling distribution** is the **probability distribution** that specifies probabilities for the possible values of the sample mean or proportion.
- A **sampling distribution** is a special case of a probability distribution where the outcome of an experiment that we are interested in is a sample statistic such as a **sample proportion**( $\hat{p}$ ) or **sample mean** ( $\bar{x}$ )
  - It's the same as what we were doing before, but now instead of singular observations we're looking at groups

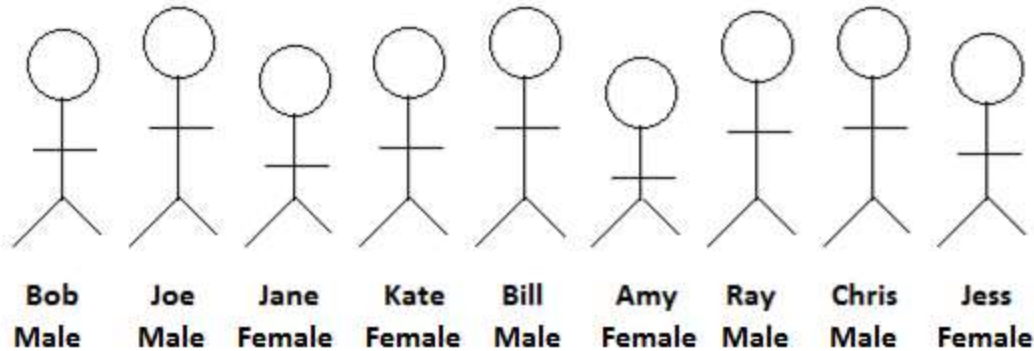
# Sampling Distributions

- This is confusing.
  - Remember, before we talked about events and random variables in  $n$  trials
  - Now, we're talking about  $m$  groups of  $n$  trials which yield  $m$  sample means or  $m$  sample proportions
    - $\bar{x}_i = \frac{\sum x}{n}$  for  $i = 1, 2, \dots, m$
    - $\hat{p}_i = \frac{x}{n}$  for  $i = 1, 2, \dots, m$

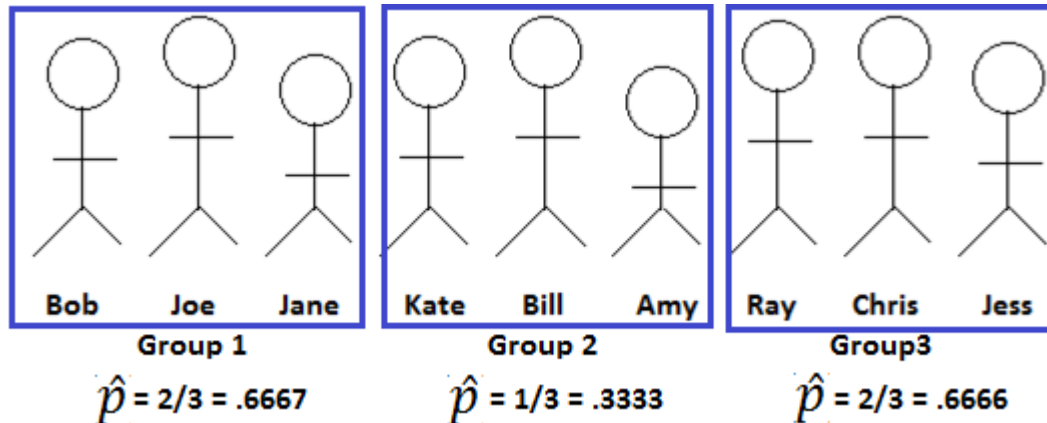
# Sampling Distributions

- Variable: Gender of Students

- Before, we measured individuals:



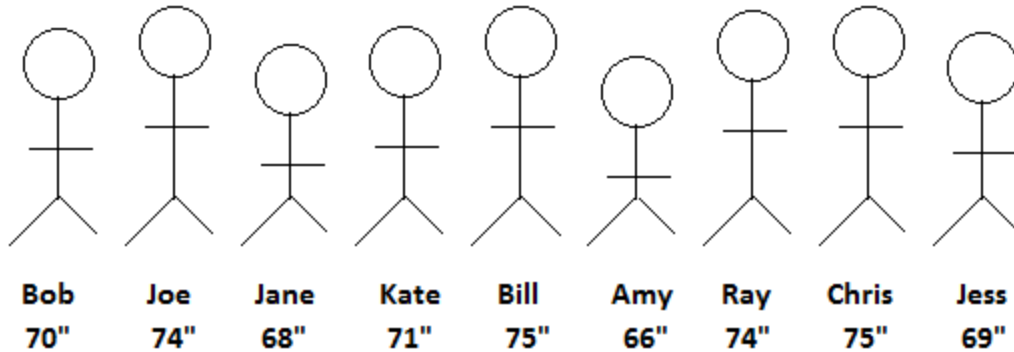
- Now, we have one measurement across groups:



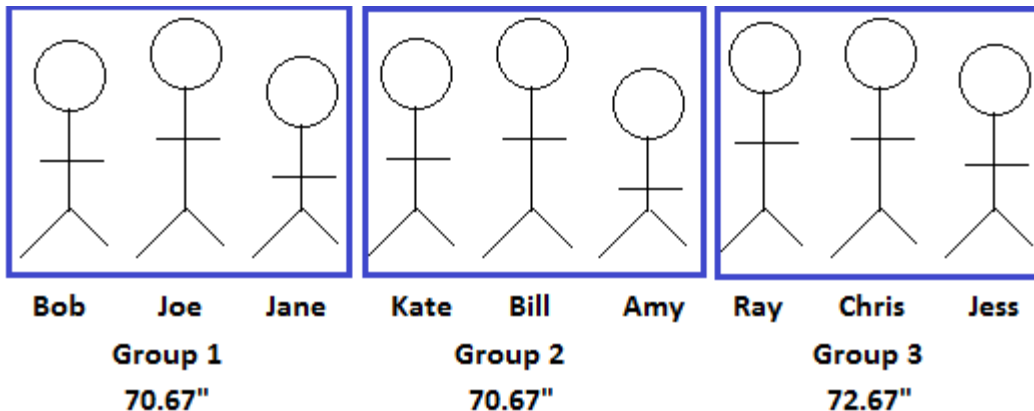
# Sampling Distributions

- Variable: Heights of Americans

– Before, we measured individuals:

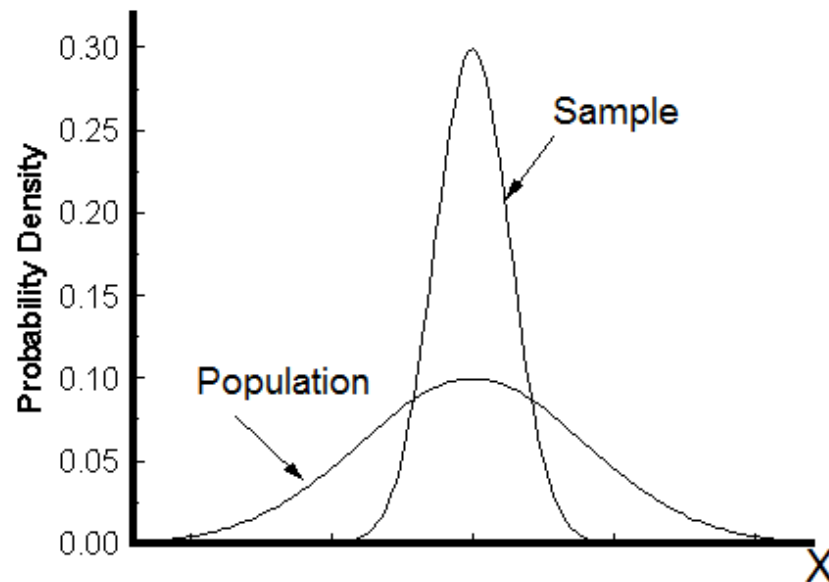


– Now, we have one measurement across groups:



# Sampling Distribution - Graphs

- Sample vs. Population: the sampling distribution is narrower than the population because grouping the data reduces the variation; pay attention to the standard error equations



# Sampling Distributions: Proportions

- This first sampling distribution we'll talk about is the **sampling distribution for the sample proportion  $\hat{p}$** .
- The idea is that there is some **true population proportion out there,  $\rho$** , but in most cases it isn't feasible to know it
  - We may not have enough time or money to poll the population
  - It may be infeasible to get a population measure



# Sampling Distributions: Proportions

- We look at **sample proportions**,  $\hat{p}$ , the proportion of observations in our sample that have a certain characteristic among our sample
  - Think “x out of n” then  $\hat{p} = \frac{x}{n}$
- We’ve looked at this before in the **descriptive statistics** but now we’re going to talk about **all possible sample proportions from repeated random samples from the population** and their distribution (mean and standard deviation)

# Sampling Distributions: Proportions

- **Before we had categorical observations:**  $x_1, x_2, x_3, \dots, x_n$ 
  - We would summarize all  $x$ 's with one **sample proportion, one  $\hat{p}$**
  - $\hat{p} = \frac{\text{number of } x \text{ with desired trait}}{\text{total sample size}}$   
= the proportion of our sample with the desired trait

# Sampling Distributions: Proportions

- **Now we have m groups of n subjects with categorical observations:**  
 $\{x_{1,1}, x_{1,2}, x_{1,3}, \dots, x_{1,n}\}, \{x_{2,1}, x_{2,2}, x_{2,3}, \dots, x_{2,n}\}, \dots, \{x_{m,1}, x_{m,2}, x_{m,3}, \dots, x_{m,n}\}$
- **Now, we find summary statistics for each group**  
 $\widehat{p}_1, \widehat{p}_2, \widehat{p}_3, \widehat{p}_4, \dots, \widehat{p}_m$

– We have m sample proportions , one  $\widehat{p}$  for each group

–  $\widehat{p}_1 = \frac{\text{number of } x \text{ with desired trait in group 1}}{\text{total sample size of group 1}}$

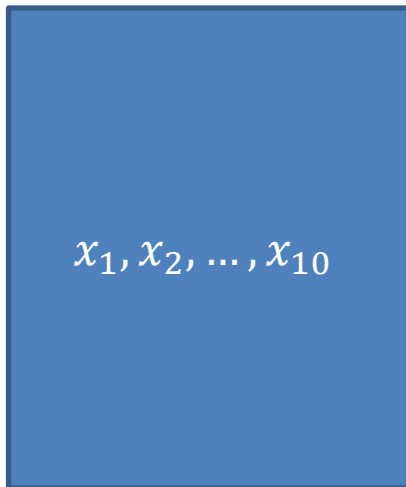
–  $\widehat{p}_2 = \frac{\text{number of } x \text{ with desired trait in group 2}}{\text{total sample size of group 2}} \dots$

–  $\widehat{p}_m = \frac{\text{number of } x \text{ with desired trait in group } m}{\text{total sample size of group } m}$

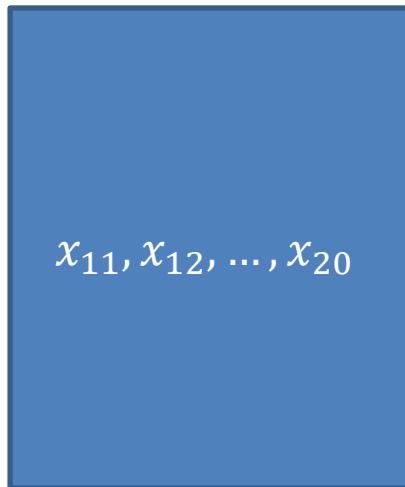
# Sampling Distributions: Proportions

- You could think of each group as a barrel and we're only interested in the proportion of each barrel; we are no longer interested in the individual responses like we might have been before
- The example below shows how we could summarize 40 observations by splitting them into four representative sample proportions

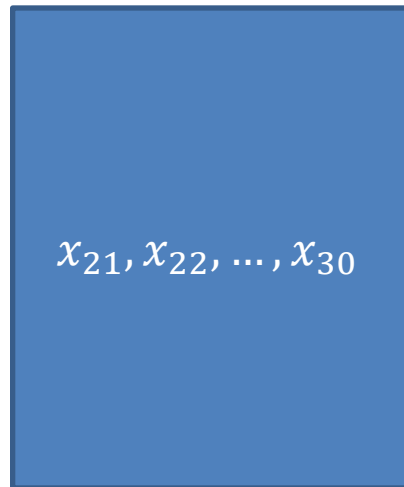
$\widehat{p}_1$



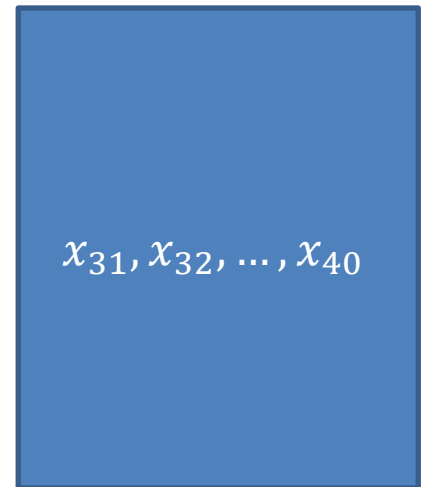
$\widehat{p}_2$



$\widehat{p}_3$



$\widehat{p}_4$



# Sampling Distribution – Mean and SD

- The **mean of the sampling distribution** for a sample proportion will always equal the population proportion:  $\mu_{\hat{p}} = \rho$ 
  - Even though we know the mean is the population proportion, we note that some  $\hat{p}$  will be lower and some will be higher

# Sampling Distribution – Mean and SD

- **Think about it this way:**
  - **Q:** If the population proportion of females in the United States is 51% what would you expect the number of females to be in a random sample of 100 Americans?
  - **A:** 51%, or 51 of 100, is our best guess; think of the binomial expectation.
- Later, we'll do this the other way around and we will call  $\hat{p}$  the **point estimate for  $\rho$**  since it's our best guess for the population proportion if we don't know it

# Sampling Distribution – Mean and SD

- The **standard error**, the standard deviation of all possible sample proportions, is:

$$\begin{aligned}\sigma_{\hat{p}} &= \sqrt{\frac{\rho(1 - \rho)}{n}} \\ &= \mathbf{St. Dev}(\hat{p}_1, \hat{p}_2, \hat{p}_3, \hat{p}_4, \dots, \hat{p}_m)\end{aligned}$$

# Sampling Distribution – Mean and SD

- **Think about it this way:**

- **Q:** If our best guess for  $\rho$  is  $\hat{p}$  we need a **measure of reliability** for our estimate
- **A:** We'll talk more about this later, but our standard error calculator is a big part of this

- Recall:  $\sigma_{\hat{p}} = \sqrt{\frac{\rho(1-\rho)}{n}}$

- Later, in the case we don't know  $\rho$  we're estimating it with our **point estimate**  $\hat{p}$ 
  - Consider:

$$\sqrt{\frac{\hat{p}(1 - \hat{p})}{n}}$$



# Sampling Distribution – Mean and SD

- $\mu_{\hat{p}} = \rho$ 
  - Even though we know the mean is the population proportion, we note that some  $\hat{p}$  will be lower and some will be higher
- $\sigma_{\hat{p}} = \sqrt{\frac{\rho(1-\rho)}{n}}$ 
  - Aside:
    - What if we increase  $n$ ?
      - The standard deviation shrinks
    - What if we decrease  $n$ ?
      - The standard deviation grows

# Sampling Distribution:

- Now that we know the mean and standard deviation of the sample proportions we can calculate z-scores to find some probabilities associated with sample proportions just like we did before.

$$\mu_{\hat{p}} = \rho$$
$$\sigma_{\hat{p}} = \sqrt{\frac{\rho(1-\rho)}{n}}$$
$$z = \frac{\text{observation} - \text{mean}}{\text{st.dev}} = \frac{\hat{p} - \mu_{\hat{p}}}{\sigma_{\hat{p}}} = \frac{\hat{p} - \rho}{\sqrt{\frac{\rho(1-\rho)}{n}}}$$

# Sampling Distribution:

$$P(\hat{p} > c) = 1 - P\left(z < \frac{c - \mu_{\hat{p}}}{\sigma_{\hat{p}}}\right) = 1 - P\left(z < \frac{c - \rho}{\sqrt{\frac{\rho(1 - \rho)}{n}}}\right)$$

$$P(\hat{p} < c) = P\left(z < \frac{c - \mu_{\hat{p}}}{\sigma_{\hat{p}}}\right) = P\left(z < \frac{c - \rho}{\sqrt{\frac{\rho(1 - \rho)}{n}}}\right)$$

$$\begin{aligned} P(c_1 < \hat{p} < c_2) &= P\left(z < \frac{c_2 - \mu_{\hat{p}}}{\sigma_{\hat{p}}}\right) - P\left(z < \frac{c_1 - \mu_{\hat{p}}}{\sigma_{\hat{p}}}\right) \\ &= P\left(z < \frac{c_2 - \rho}{\sqrt{\frac{\rho(1 - \rho)}{n}}}\right) - P\left(z < \frac{c_1 - \rho}{\sqrt{\frac{\rho(1 - \rho)}{n}}}\right) \end{aligned}$$

# Sampling Distributions - Example

- Say, we know that **16% of Americans approve of Congress (Gallup)**.
- **What is the sampling distribution of the sample proportion** of Americans that approve of Congress for  $n=100$ ?
  - Note, we aren't interested in the yes or no's individually but the proportion among the ten
  - Here,  $X$ =the proportion of the one hundred Americans in each group

# Sampling Distributions - Example

- Say, we know that **16% of Americans approve of Congress (Gallup)**.
- **What is the sampling distribution of the sample proportion** of Americans that approve of Congress for  $n=100$ ?
  - $n$  = sample size = **sample size of one hundred**= 100
  - $p$  = population proportion = **16%** =.16

# Sampling Distributions – Example

- Let's find the sampling distribution mean:
- **The mean of all sample proportions of n=100**  
 $= \mu_{\hat{p}} = \rho = 16\% = .16$ 
  - Some  $\hat{p}$  will be lower and some will be higher but **the mean of all sample proportions of n=100 will be .6**

# Sampling Distributions – Example

- Let's find the sampling distribution st. error:
- **The st. deviation of all sample proportions of n=100**

$$\begin{aligned} \text{= Standard Error} &= \sigma_{\hat{p}} = \sqrt{\frac{\rho(1-\rho)}{n}} \\ &= \sqrt{\frac{.16(1 - .16)}{100}} = .0367 \end{aligned}$$

# Sampling Distributions – Example

- Let's find the sampling distribution :

$$\mu_{\hat{p}} = \rho = 16\% = .16$$

$$\sigma_{\hat{p}} = \sqrt{\frac{\rho(1 - \rho)}{n}} = \sqrt{\frac{.16(1 - .16)}{100}} = .0367$$



# Sampling Distributions – Example

- The probability that **most**, of our sample of  $n=100$ , approve of Congress:

R

$$\begin{aligned} P(\hat{p} > .5) &= 1 - P(\hat{p} < .5) \\ &= 1 - pnorm(.5, .16, .0367) \\ &\approx 0 \end{aligned}$$

Z-table:

$$\begin{aligned} P(\hat{p} > .5) &= P\left(z > \frac{.5 - .16}{.0367}\right) = P(Z > 9.26) \\ &= 1 - P(Z \leq 9.26) \approx 1 - 1 \\ &= 0 \end{aligned}$$

**Note:** the z-table is a less accurate approximation than R

# Sampling Distributions – Example

- The probability that **less than 10%**, of our sample of  $n=100$ , approve of Congress:

R

$$\begin{aligned} P(\hat{p} < .1) &= \text{pnorm}(.1, .16, .0367) \\ &= .0510 \end{aligned}$$

Z-table:

$$\begin{aligned} P(\hat{p} < .1) &= P\left(z < \frac{.1 - .16}{.0367}\right) = P(Z < -1.63) \\ &= .0516 \end{aligned}$$

**Note:** the z-table is a less accurate approximation than R

# Sampling Distributions – Example

- The probability that **between 5 and 19 percent**, of our sample of  $n=100$ , approve of Congress:

R

$$\begin{aligned}P(.05 < \hat{p} < .19) &= P(\hat{p} < .19) - P(\hat{p} < .05) \\ &= \text{pnorm}(.19, .16, .0367) - \text{pnorm}(.05, .16, .0367) \\ &= .7917991\end{aligned}$$

Z-table:

$$\begin{aligned}P(.05 < \hat{p} < .19) &= P(\hat{p} < .19) - P(\hat{p} < .05) \\ &= P\left(z < \frac{.19 - .16}{.0367}\right) - P\left(z < \frac{.05 - .16}{.0367}\right) \\ &= P(Z < .82) - P(Z < -3.00) \\ &= .7939 - .0013 \\ &= .7926\end{aligned}$$

# Sampling Distributions – Example

- **Note:** we had to assume normality of  $\hat{p}$  to use the Z-score transformation and R code I provided to solve the previous probabilities
- Later, we will see how we are able to make that assumption – unlocking all of the nice methodologies of the Normal distribution

# Sampling Distributions: Means

- This second sampling distribution we'll talk about is the **sampling distribution for the sample mean**.
- The idea is that there is some **true population mean out there,  $\mu$** , but it might not be feasible to know it
  - We may not have enough time or money to poll the population
  - It may be infeasible to get a population measure

# Sampling Distributions: Means

- Instead, we look at **sample mean,  $\bar{x}$** , the mean of quantitative observations
- We've looked at this before in the **descriptive statistics** but now we're going to talk about **all possible sample means from repeated random samples from our population**

# Sampling Distributions: Means

- **Before we had quantitative observations:**  $x_1, x_2, x_3, \dots, x_n$ 
  - We would summarize all  $x$ 's with one **sample mean, one  $\bar{x}$**
  - $\bar{x} = \frac{\text{the sum of } x\text{'s}}{\text{the total sample size}} = \frac{\sum x}{n}$   
  
= the mean of the observations in our sample

# Sampling Distributions: Means

- **Now we have m groups of n subjects with categorical observations:**  
 $\{x_{1,1}, x_{1,2}, x_{1,3}, \dots, x_{1,n}\}, \{x_{2,1}, x_{2,2}, x_{2,3}, \dots, x_{2,n}\}, \dots, \{x_{m,1}, x_{m,2}, x_{m,3}, \dots, x_{m,n}\}$
- Now, we find summary statistics for each group

$$\bar{x}_1, \bar{x}_2, \bar{x}_3, \bar{x}_4, \dots, \bar{x}_m$$

– We have m sample means, one  $\bar{x}$  for each group

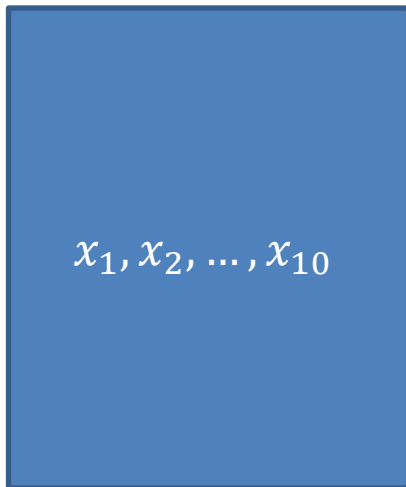
- $\bar{x}_1 = \frac{\text{the sum of } x' \text{'s from group 1}}{\text{the total sample size of group 1}} = \frac{\sum x}{n}$
- $\bar{x}_2 = \frac{\text{the sum of } x' \text{'s from group 2}}{\text{the total sample size of group 2}} = \frac{\sum x}{n} \dots$
- $\bar{x}_m = \frac{\text{the sum of } x' \text{'s from group } m}{\text{the total sample size of group } m} = \frac{\sum x}{n}$



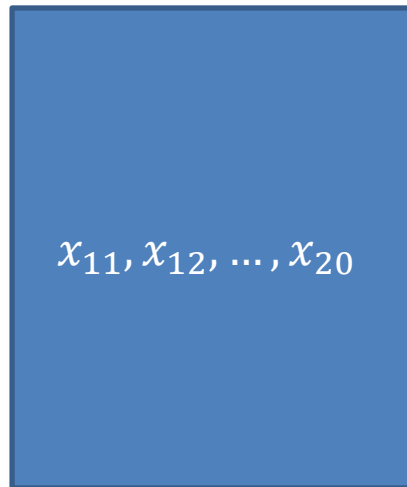
# Sampling Distributions: Means

- You could think of each group as a barrel and we're only interested in the mean of each barrel; we are no longer interested in the individual responses
- The example below shows how we could summarize 40 observations, into four representative sample means

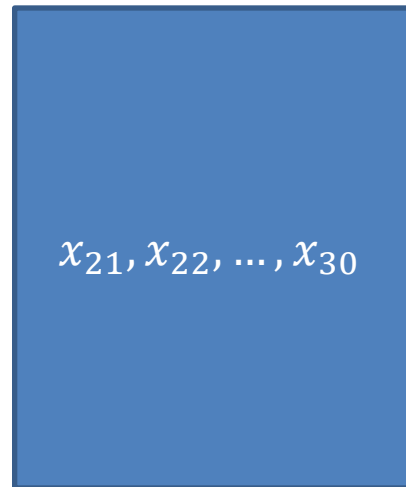
$\bar{x}_1$



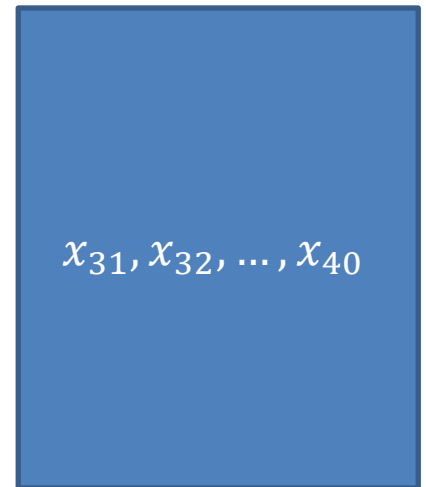
$\bar{x}_2$



$\bar{x}_3$



$\bar{x}_4$



# Sampling Distribution – Mean and SD

- The mean of the sampling distribution for a sample mean will always equal the population mean:  $\mu_{\bar{x}} = \mu_x$ 
  - This is the mean of all possible sample means, but we note that some  $\bar{x}$  will be lower and some will be higher

# Sampling Distribution – Mean and SD

- **Think about it this way:**
  - **Q:** If the population mean of time Americans spend on social media is 100 minutes with a standard deviation of 25 minutes what would you expect the average time a sample of 35 Americans spent on social media?
  - **A:** 100 minutes is our best guess.
- Later, we'll do this the other way around and we will call  $\bar{x}$  the **point estimate for  $\mu_x$**  since it's our best guess for the population mean if we don't know it

# Sampling Distribution – Mean and SD

- The standard error, the standard deviation of all possible sample means, is:

$$\begin{aligned}\sigma_{\bar{x}} &= \frac{\sigma_x}{\sqrt{n}} \\ &= \mathbf{St. Dev}(\bar{x}_1, \bar{x}_2, \bar{x}_3, \bar{x}_4, \dots, \bar{x}_m)\end{aligned}$$

# Sampling Distribution – Mean and SD

- **Think about it this way:**
  - **Q:** If our best guess for  $\mu$  is  $\bar{x}$  we need a **measure of reliability** for our estimate
  - **A:** We'll talk more about this later, but our standard error calculator is a big part of this
- Later, in the case we don't know  $\mu_x$  or  $\sigma_x$  we're estimating it with our **point estimate  $\bar{x}$** 
  - Recall:  $\sigma_{\bar{x}} = \frac{\sigma_x}{\sqrt{n}}$
  - Consider:  $\frac{s_x}{\sqrt{n}}$  [**Note: we estimate  $\sigma_x = s_x$** ]

# Sampling Distribution – Mean and SD

- $\mu_{\bar{x}} = \text{mean of all sample means} = \mu_x$ 
  - Even though we know the mean is the population mean, we note that some  $\bar{x}$  will be lower and some will be higher
- $\sigma_{\bar{x}} = \frac{\sigma_x}{\sqrt{n}}$
- Aside:
  - What if we increase n?
    - The standard deviation shrinks
  - What if we decrease n?
    - The standard deviation grows

# Sampling Distribution:

- Now that we know the mean and standard error of the sample means we can calculate z-scores to find some probabilities associated with sample means just like we did before.

$$\mu_{\bar{x}} = \mu_x$$
$$\sigma_{\bar{x}} = \frac{\sigma_x}{\sqrt{n}}$$

$$z = \frac{\textit{observation} - \textit{mean}}{\textit{st. dev}} = \frac{\bar{x} - \mu_{\bar{x}}}{\sigma_{\bar{x}}} = \frac{\bar{x} - \mu_x}{\frac{\sigma_x}{\sqrt{n}}}$$

# Sampling Distribution:

$$P(\bar{x} > c) = 1 - P\left(z < \frac{c - \mu_{\bar{x}}}{\sigma_{\bar{x}}}\right) = 1 - P\left(z < \frac{c - \mu_x}{\frac{\sigma_x}{\sqrt{n}}}\right)$$

$$P(\bar{x} < c) = P\left(z < \frac{c - \mu_{\bar{x}}}{\sigma_{\bar{x}}}\right) = P\left(z < \frac{c - \mu_x}{\frac{\sigma_x}{\sqrt{n}}}\right)$$

$$\begin{aligned} P(c_1 < \bar{x} < c_2) &= P\left(z < \frac{c_2 - \mu_{\bar{x}}}{\sigma_{\bar{x}}}\right) - P\left(z < \frac{c_1 - \mu_{\bar{x}}}{\sigma_{\bar{x}}}\right) \\ &= P\left(z < \frac{c_2 - \mu_x}{\frac{\sigma_x}{\sqrt{n}}}\right) - P\left(z < \frac{c_1 - \mu_x}{\frac{\sigma_x}{\sqrt{n}}}\right) \end{aligned}$$



# Sampling Distributions - Example

- Say, we know that **the average American spends 100 minutes on social media per day with a standard deviation of 25 minutes.**
- **What is the sampling distribution of the sample mean** of time Americans spend on social media for  $n=35$ ?
  - Note, we aren't interested in the individuals but the group of thirty five
  - Here,  $X$ =the proportion of the ten Americans in each group

# Sampling Distributions - Example

- Say, we know that **the average American spends 100 minutes on social media per day with a standard deviation of 25 minutes.**
- **What is the sampling distribution of the sample mean of time Americans spend on social media for  $n=35$ ?**
  - $n$  = sample size = **sample size of thirty five** = 35
  - $\mu_x$  = population mean = 100
  - $\sigma_x$  = population standard deviation = 25

# Sampling Distributions – Example

- Let's find the sampling distribution mean:
- **The mean of all sample means of n=35**  
 $= \mu_{\bar{x}} = \mu_x = 100$ 
  - Some  $\bar{x}$  will be lower and some will be higher but **the mean of all sample means of n=35 will be 100**

# Sampling Distributions – Example

- Let's find the sampling distribution st. error:
- **The st. deviation of all sample means of n=35**  
= Standard Error

$$= \sigma_{\bar{x}} = \frac{\sigma_x}{\sqrt{n}} = \frac{25}{\sqrt{35}} = 4.2258$$

# Sampling Distributions – Example

- Let's find the sampling distribution :

$$\mu_{\bar{x}} = \mu_x = 100$$

$$\sigma_{\bar{x}} = \frac{\sigma_x}{\sqrt{n}} = \frac{25}{\sqrt{35}} = 4.2258$$

# Sampling Distributions – Example

- The probability that a sample of  $n=35$  spend **more than two hours** on social media on average:

R

$$\begin{aligned}P(\bar{x} > 120) &= 1 - P(\bar{x} < 120) \\ &= 1 - pnorm(120, 100, 4.2258) \\ &= .000001107046\end{aligned}$$

Z-table:

$$\begin{aligned}P(\bar{x} > 120) &= P\left(z > \frac{120 - 100}{4.2258}\right) = P(Z > 4.73) \\ &= 1 - P(Z < 4.73) \approx 1 - 1 \\ &= 0\end{aligned}$$

**Note:** the z-table is a less accurate approximation than R

# Sampling Distributions – Example

- The probability that a sample of  $n=35$  spend **less than one hour** on social media on average:

R

$$\begin{aligned}P(\bar{x} < 60) &= P(\bar{x} < 60) \\ &= \text{pnorm}(60, 100, 4.2258) \\ &= 1.458519 * 10^{-21}\end{aligned}$$

Z-table:

$$\begin{aligned}P(\bar{x} > 60) &= P\left(z < \frac{60 - 100}{4.2258}\right) = P(Z < -9.47) \\ &= P(Z < -9.47) \\ &\approx 0\end{aligned}$$

**Note:** the z-table is a less accurate approximation than R

# Sampling Distributions – Example

- The probability that a sample of  $n=35$  spend **between 1 and 1.5 hours** on social media on average:

R

$$\begin{aligned}P(60 < \bar{x} < 90) &= P(\bar{x} < 90) - P(\bar{x} < 60) \\ &= \text{pnorm}(90, 100, 4.2258) - \text{pnorm}(60, 100, 4.2258) \\ &= .008980629\end{aligned}$$

Z-table:

$$\begin{aligned}P(60 < \bar{x} < 90) &= P\left(\frac{90 - 100}{4.2258} < z < \frac{60 - 100}{4.2258}\right) \\ &= P(Z < -2.37) - P(Z < -9.47) \\ &\approx .0089 - 0 \\ &= 0\end{aligned}$$

**Note:** the z-table is a less accurate approximation than R



# Sampling Distributions – Example

- **Note:** we had to assume normality of  $\bar{x}$  to use the Z-score transformation and R code I provided to solve the previous probabilities
- Later, we will see how we are able to make that assumption – unlocking all of the nice methodologies of the Normal distribution

# Recall: Law of Large Numbers

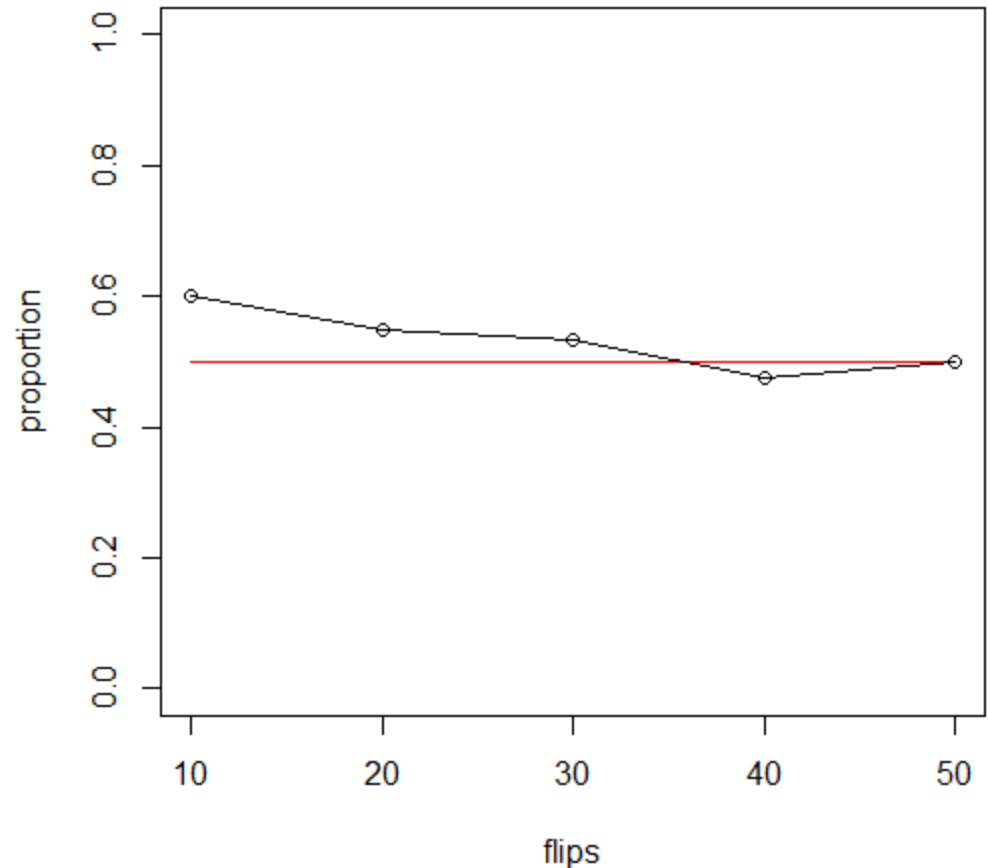
- **(LLN 1)** – As the sample size increases the sample estimates ( $\bar{x}$  or  $\hat{p}$ ) approach the population values ( $\mu$  or  $\sigma$ )
- **(LLN 2)** – As the number of trials increase the **proportion** of occurrences of any given outcome approaches the probability in the long run.

# Recall: Law of Large Numbers

- 10 flips: 6 heads were flipped
  - *Total proportion*  $= \frac{x}{n} = \frac{6}{10} = .60 = 60\% \text{ heads}$
- 10 more flips: 5 heads were flipped
  - *Total proportion*  $= \frac{x}{n} = \frac{5+6}{10+10} = \frac{11}{20} = .55 = 55\% \text{ heads}$
- 10 more flips: 5 heads were flipped
  - *Total proportion*  $= \frac{x}{n} = \frac{11+5}{20+10} = \frac{16}{30} = .5333 = 53.33\% \text{ heads}$
- 10 more flips: 3 heads were flipped
  - *Total proportion*  $= \frac{x}{n} = \frac{16+3}{30+10} = \frac{19}{40} = .475 = 47.5\% \text{ heads}$
- 10 more flips: 6 heads were flipped
  - *Total proportion*  $= \frac{x}{n} = \frac{19+6}{40+10} = \frac{25}{50} = .5 = 50\% \text{ heads}$

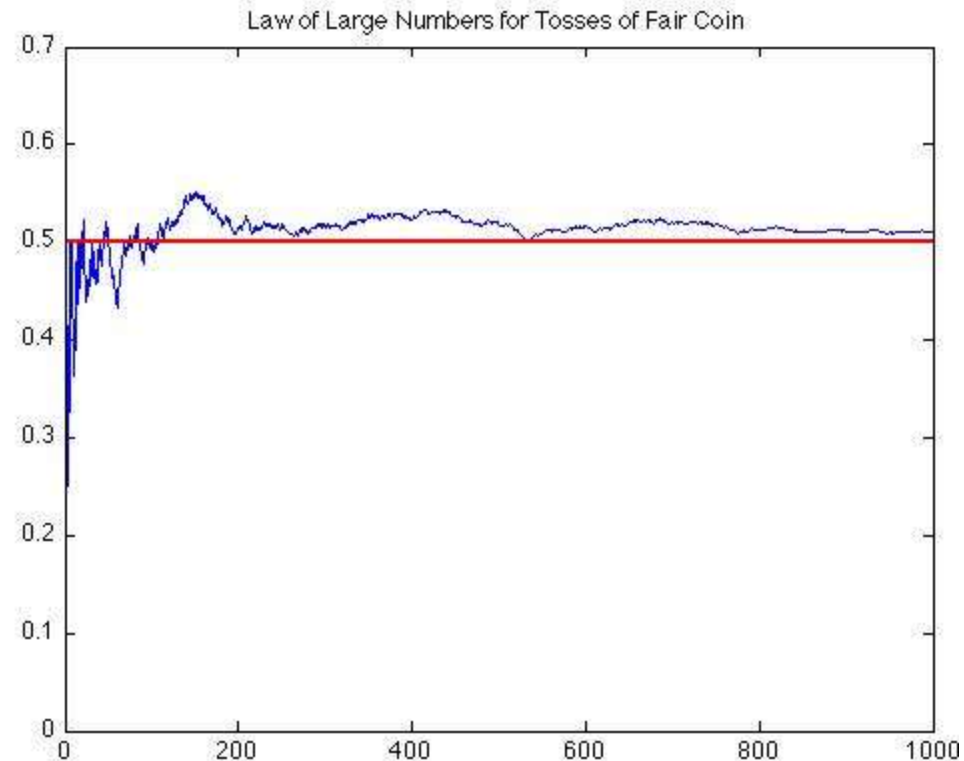
# Recall: Law of Large Numbers

- **(LLN)** – As the number of flips increase the **proportion** of heads approaches the **probability** of seeing a heads,  $P(\text{heads}) = .5$ , which is the red line.



# Recall: Law of Large Numbers

- At first the proportion is all over the place – you can see the large spikes in the graph
- Importantly, we see that the proportion of coins that landed on heads levels off and gets closer and closer to 50%, the probability, which is where we expect it to go ‘in the long run!’



# Central Limit Theorem: Proportions

- For random sampling with a **large sample size  $n$ , the sampling distribution of the sample proportion** is approximately a normal distribution
  - $n * p \geq 15$  and  $n * (1 - p) \geq 15$
- Introduction:
  - [https://www.youtube.com/watch?v=Pujol1yC1\\_A](https://www.youtube.com/watch?v=Pujol1yC1_A)

# Central Limit Theorem: Means

- For random sampling with a **large sample size  $n$ , the sampling distribution of the sample mean** is approximately a normal distribution
  - For us, 30 is close enough to infinity
- Introduction:
  - [https://www.youtube.com/watch?v=Pujol1yC1\\_A](https://www.youtube.com/watch?v=Pujol1yC1_A)

# Central Limit Theorem: Means

- 1) For any population the sampling distribution of  $\bar{x}$  is bell shaped when the sample size  $n$  is large, when  $n$  is thirty or more
- 2) The sampling distribution of  $\bar{x}$  is bell-shaped when the population distribution is distribution is bell-shaped, regardless of sample size
- 3) We do not know the shape of the sampling distribution of  $\bar{x}$  if the sample size is small and the population distribution isn't bell-shaped



# Central Limit Theorem

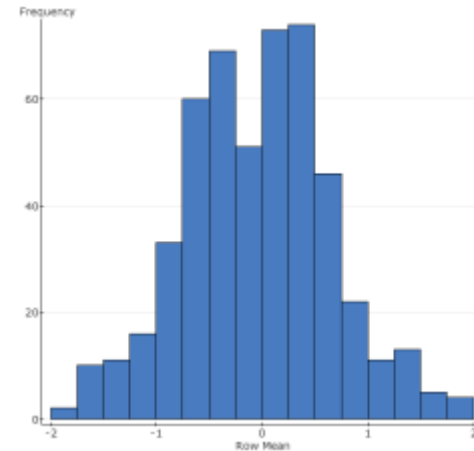
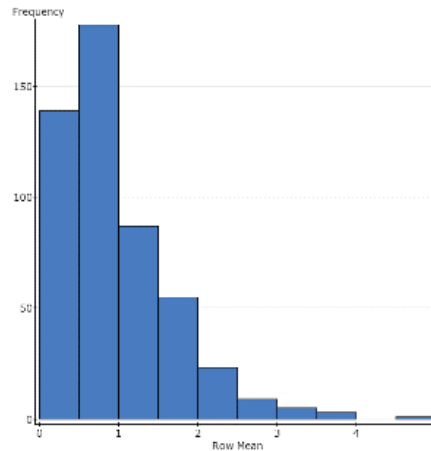
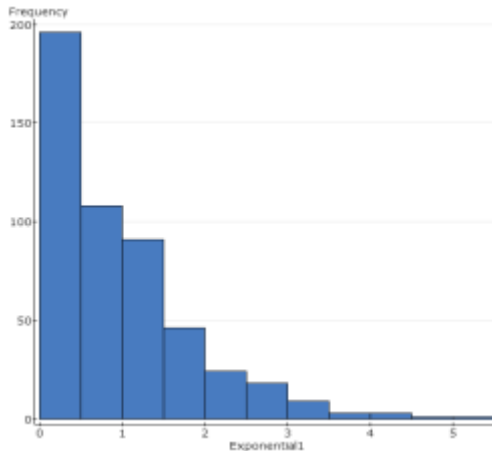
For any population the sampling distribution of  $\bar{x}$  is bell shaped when the sample size  $n$  is large, when  $n$  is thirty or more

**Note:** for small sample size we can't say this.

Population

$\bar{x}$  when  $n=2$

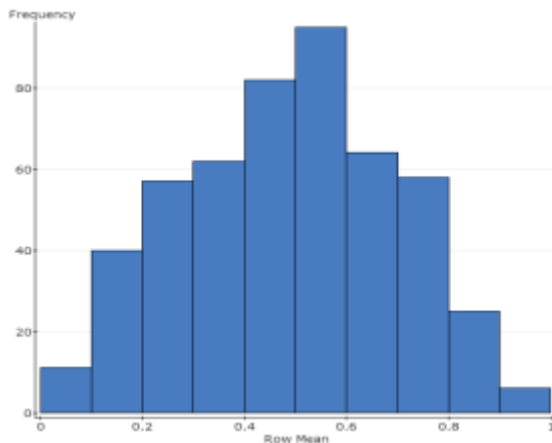
$\bar{x}$  when  $n=30$



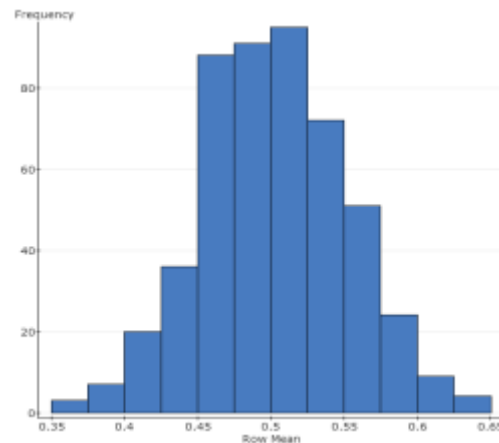
# Central Limit Theorem

The sampling distribution of  $\bar{x}$  is bell-shaped when the population distribution is distribution is bell-shaped, regardless of sample size

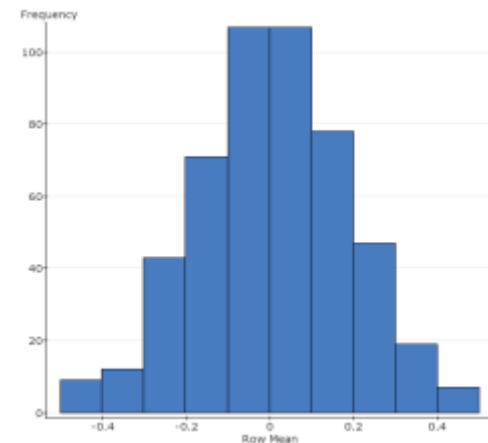
Population



$\bar{x}$  when n=2



$\bar{x}$  when n=30



# Chapter 6 Supplement (not in text book)

# Central Limit Theorem

- **Recall:** The sampling distribution of  $\bar{x}$  is bell-shaped even for  $n=1$  if  $X$  follows the normal distribution
- If a simple random sample is drawn from the population then  $z = \frac{\bar{x} - \mu_x}{\sigma_{\bar{x}}}$  follows the standard normal distribution
- The difficulty in this is that we rarely know the value of the parameters:  $\mu_x$  or  $\sigma_x$  in  $\sigma_{\bar{x}} = \frac{\sigma_x}{\sqrt{n}}$

# Central Limit Theorem

- To combat this difficulty we will cover a few more named distributions and sampling distributions
  1.  $s^2$  and the  $\chi^2$  (chi-squared) distribution
  2.  $\bar{x}, s^2$  and the Student's t-distribution
  3. Two variances and the F distribution

$s^2$  and the  $\chi^2$  (chi-squared) distribution

- The chi-squared density can be defined as follows:

$$f_{X(x)} = \frac{1}{2^{\frac{n}{2}} \Gamma\left(\frac{n}{2}\right)} x^{\frac{n}{2}-1} e^{-\frac{x}{2}} I(X > 0)$$

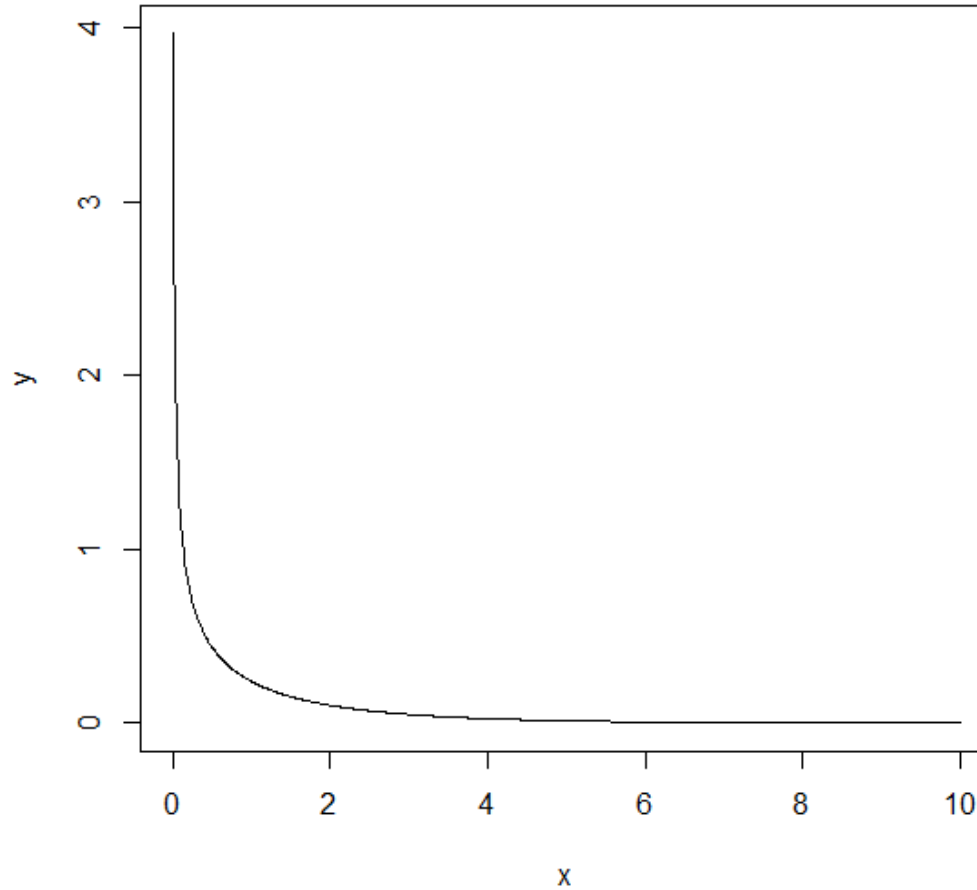
- Mean =  $n$
- Variance =  $2n$

## $\chi^2$ (chi-squared) distribution in R

- $P(X = x) = 0$  as the probability of any one value is always zero
- $P(X \leq x) = \text{pchisq}(x, n)$
- $P(X \geq x) = 1 - \text{pchisq}(x, n)$
- $P(x_1 < X < x_2) = \text{pchisq}(x_2, n) - \text{pchisq}(x_1, n)$

# $s^2$ and the $\chi^2$ (chi-squared) distribution

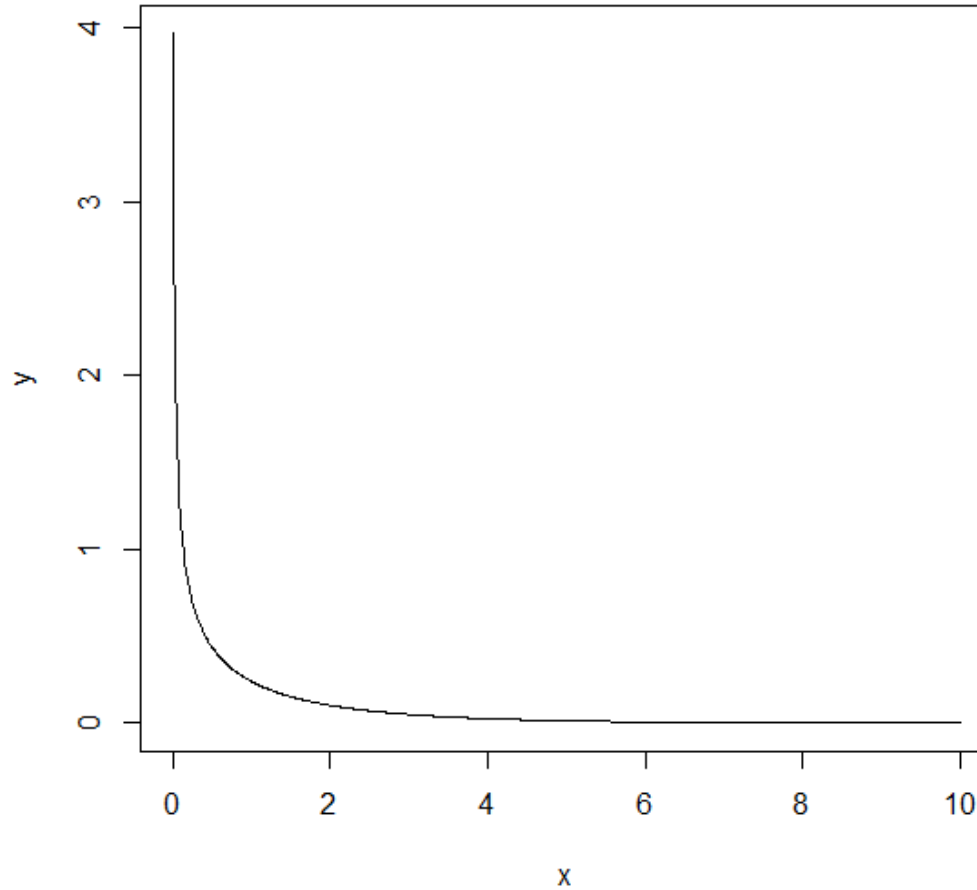
n=1





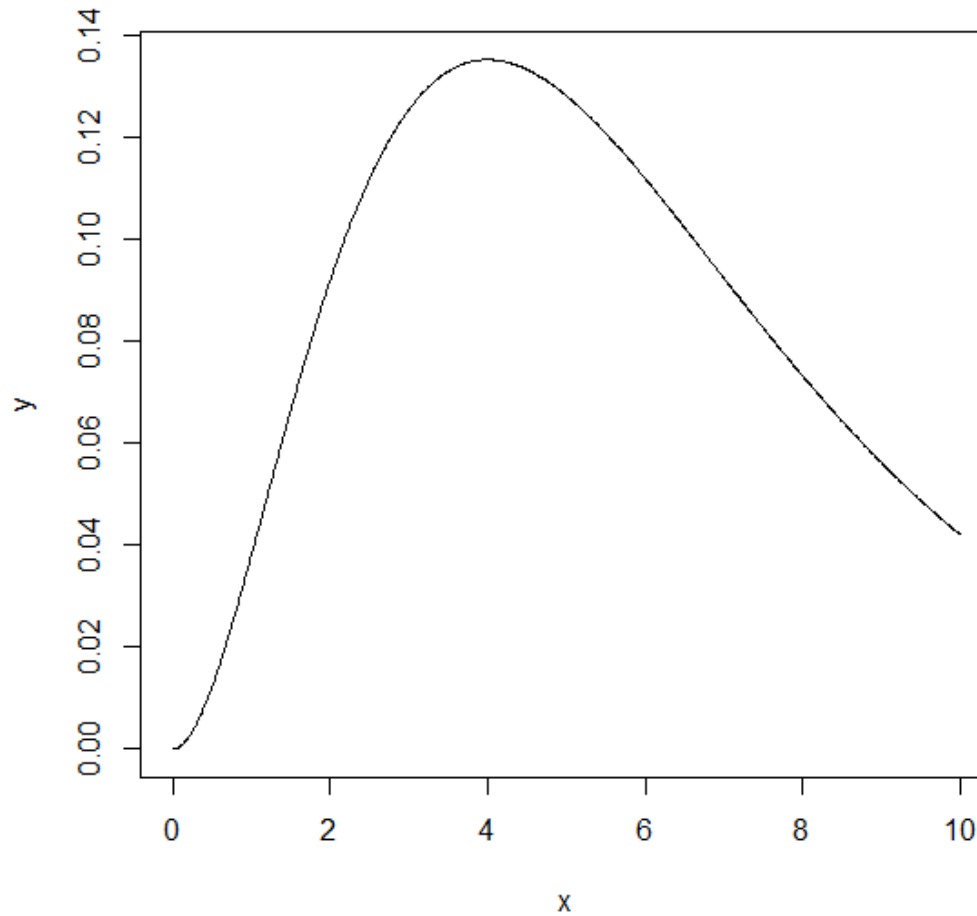
# $s^2$ and the $\chi^2$ (chi-squared) distribution

$n=1$



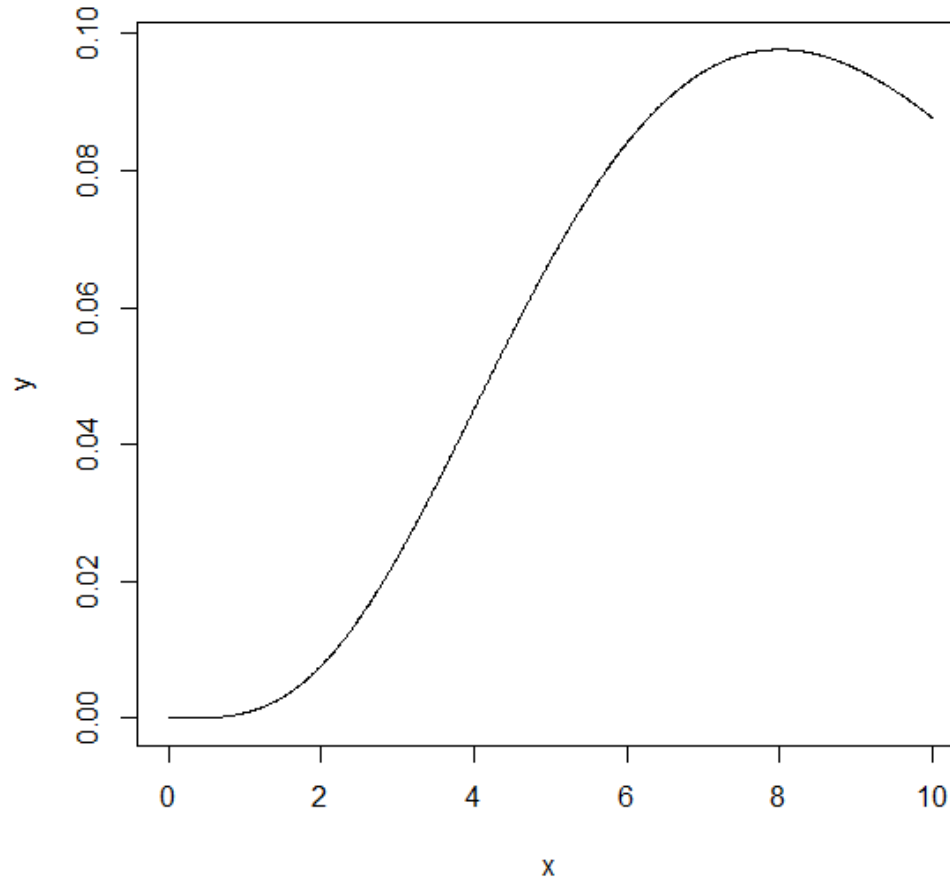
# $s^2$ and the $\chi^2$ (chi-squared) distribution

n=6



# $s^2$ and the $\chi^2$ (chi-squared) distribution

n=10



$s^2$  and the  $\chi^2$  (chi-squared) distribution

- The chi-squared distribution can be defined as follows:

$$X^2 = z_1^2 + z_2^2 + \cdots + z_n^2$$

- **$X^2$  follows the chi-squared distribution with  $n$  degrees of freedom** where  $z_i$  are independent random variables that follow the standard normal distribution

## $s^2$ and the $\chi^2$ (chi-squared) distribution

- Consider  $x_1, x_2, \dots, x_n$  independent from a normal distribution
- If we look at the formula for the sample variance from chapter two we see that we are summing  $(x_i - \bar{x})^2$  (n normal random variables)

## $s^2$ and the $\chi^2$ (chi-squared) distribution

- We're close to showing  $s^2$  follows the  $\chi^2$  distribution but  $(x_i - \bar{x})^2$  are normal – not the standard normal.
- Now, if we can write it as a sum of squared standard normal random variables we can say that  $s^2$  follows the  $\chi^2$  distribution
  - This requires a “cute” trick

$s^2$  and the  $\chi^2$  (chi-squared) distribution

$$s^2 = \frac{\sum (x_i - \bar{x})^2}{n - 1}$$
$$\left( \frac{n - 1}{\sigma_x^2} \right) s^2 = \frac{\sum (x_i - \bar{x})^2}{n - 1} \left( \frac{n - 1}{\sigma_x^2} \right)$$
$$\left( \frac{(n - 1)s^2}{\sigma_x^2} \right) = \sum \left( \frac{x_i - \bar{x}}{\sigma_x} \right)^2$$

$s^2$  and the  $\chi^2$  (chi-squared) distribution

$$X_{n-1}^2 = \left( \frac{(n-1)s^2}{\sigma_x^2} \right) = \sum \left( \frac{x_i - \bar{x}}{\sigma_x} \right)^2$$

- Again, consider  $\left( \frac{x_i - \bar{x}}{\sigma_x} \right)$  on the right hand side
- If we change  $\bar{x}$  to  $\mu$  we have the usual  $z = \frac{x_i - \mu_x}{\sigma_x}$
- **Note:** We lose one degree of freedom, going from  $n$  to  $n-1$ , because we are using  $\bar{x}$  instead of  $\mu$



## $\bar{x}$ , $s^2$ and the Student's t-distribution

- Still dealing with the difficulty of needing to know the population standard deviation for the Central limit theorem we talk about the t-distribution

## $\bar{x}$ , $s^2$ and the Student's t-distribution

- The t density can be defined as follows:

$$f_{X(x)} = \frac{\Gamma\left(\frac{n+1}{2}\right)}{\sqrt{\pi n} \Gamma\left(\frac{n}{2}\right)} \left(1 + \frac{x^2}{n}\right)^{-\frac{n+1}{2}} I(X \in \mathbb{R})$$

- Mean = 0
- Variance =  $\frac{n}{n-2}$

# Properties of the t-distribution

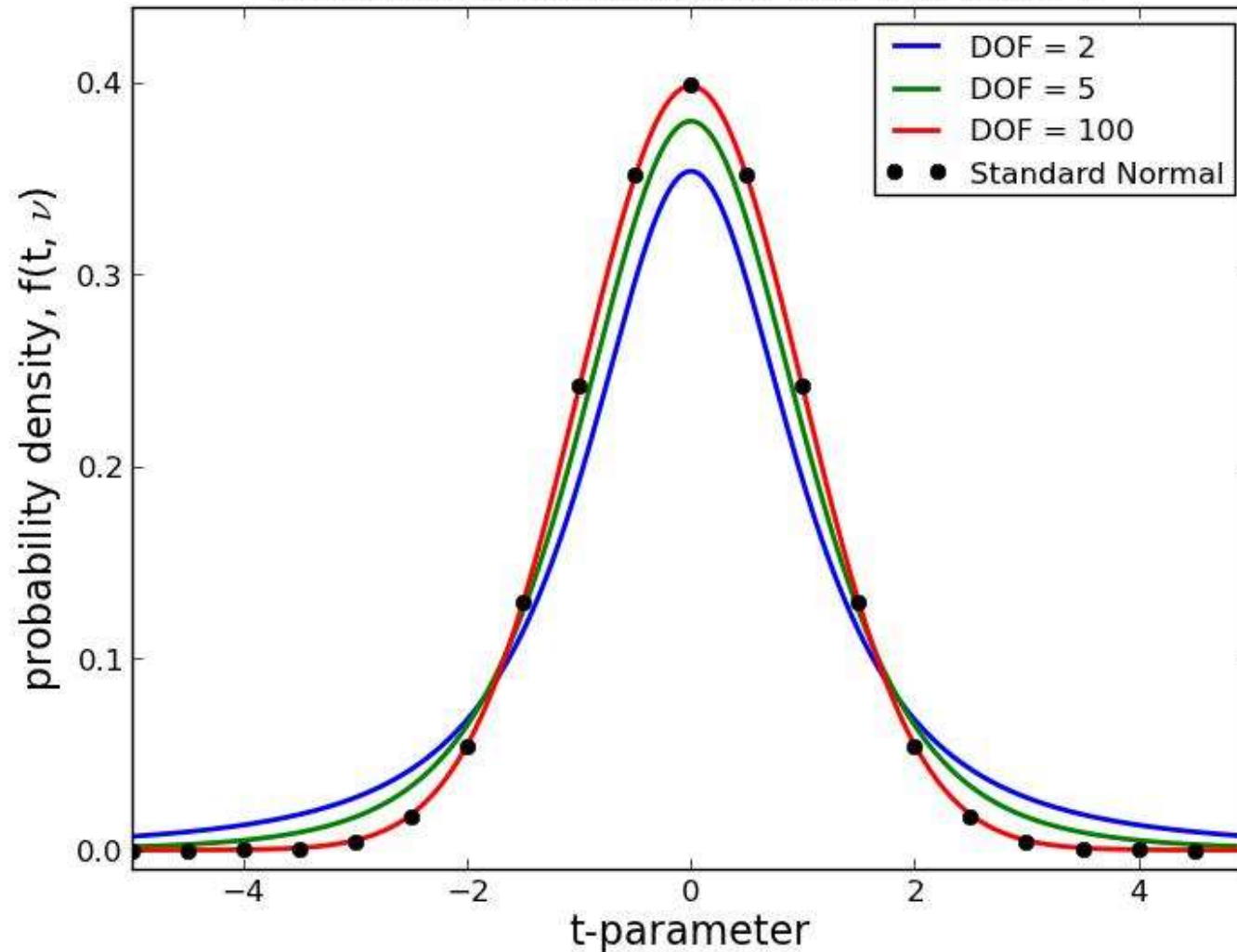
1. The t-distribution is different for different degrees of freedom
2. The t-distribution is centered and symmetric at 0
3. The area under the curve is 1 and  $\frac{1}{2}$  on either side of 0
4. The probability approaches 0 as we move away from 0
5. The t-distribution has fatter tails than the standard normal
6. As the sample size increases t gets close to z

## $\bar{x}, s^2$ and the Student's t-distribution in R

- $P(X = x) = 0$  as the probability of any one value is always zero
- $P(X \leq x) = \text{pt}(x, n)$
- $P(X \geq x) = 1 - \text{pt}(x, n)$
- $P(x_1 < X < x_2) = \text{pt}(x_2, n) - \text{pt}(x_1, n)$

# The t-distribution

Student distribution for various  $\nu$



## $\bar{x}$ , $s_x^2$ and the Student's t-distribution

- The t distribution can be defined as follows:

$$t_n = \frac{Z}{\sqrt{\frac{\chi_n^2}{n}}}$$

- $t_n$  follows the t distribution with n degrees of freedom where Z follows the standard normal distribution and  $\chi^2$  is chi-squared and divided by its degrees of freedom

# $\bar{x}$ , $s_x^2$ and the Student's t-distribution

- Relating it back to  $s_x^2$  and  $\chi_n^2$

$$t_n = \frac{Z}{\sqrt{\frac{\chi_n^2}{n}}} = \frac{\left( \frac{\bar{x} - \mu_x}{\frac{\sigma_x}{\sqrt{n}}} \right)}{\sqrt{\frac{\left( \frac{(n-1)s_x^2}{\sigma_x^2} \right)}{n-1}}}$$

# $\bar{x}$ , $s_x^2$ and the Student's t-distribution

- Relating it back to  $s_x^2$  and  $\chi_n^2$

$$\begin{aligned} t_{n-1} &= \frac{\left( \frac{\bar{x} - \mu_x}{\frac{\sigma_x}{\sqrt{n}}} \right)}{\sqrt{\frac{s_x^2}{\sigma_x^2}}} = \frac{\left( \frac{\bar{x} - \mu_x}{\frac{\sigma_x}{\sqrt{n}}} \right)}{\frac{s_x}{\sigma_x}} = \frac{\left( \frac{\bar{x} - \mu_x}{\frac{\sigma_x}{\sqrt{n}}} \right) \frac{\sigma_x}{s_x}}{\frac{\sigma_x}{s_x}} \\ &= \frac{\bar{x} - \mu_x}{(s_x / \sqrt{n})} \end{aligned}$$



## $\bar{x}$ , $s_x^2$ and the Student's t-distribution

- Note the similarity to the z-score: the only difference here is that we estimate  $\sigma_x$  with  $s_x$
- Later we will use the t-distribution when we make inference on the mean and don't know the population standard deviation

$$t_{n-1} = \frac{\bar{x} - \mu_x}{(s_x/\sqrt{n})}$$

# Two variances and the F distribution

- This result is something we will use in Chapter 11 and deals with the relationship between the F and the t distribution.

# Two variances and the F distribution

- The F density can be defined as follows:

$$f_{x(x)} = \frac{\left( \sqrt{\frac{(d_1 x)^2 d_2^2}{(d_1 x + d_2)^{(d_1 + d_2)}}} \right)}{xB \left( \frac{d_1}{2}, \frac{d_2}{2} \right)} I(x \geq 0)$$

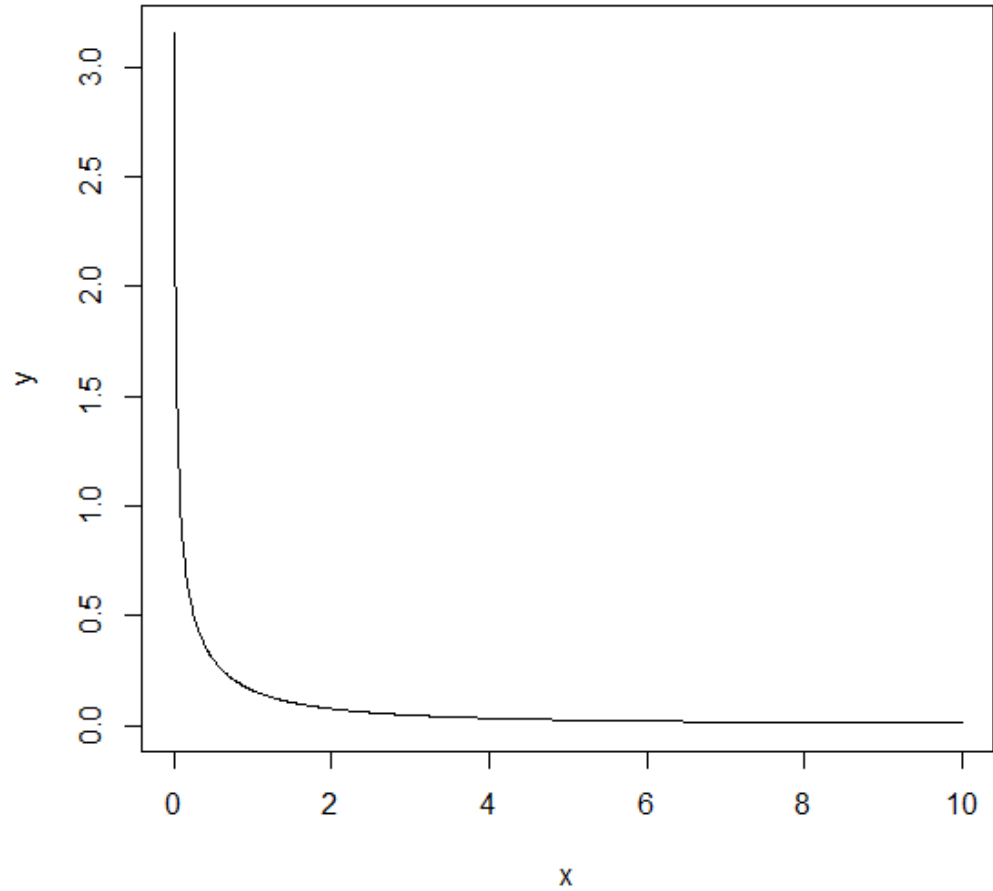
- Mean =  $\frac{d_2}{d_2 - 2}$  for  $d_2 > 2$
- Variance =  $\frac{2 * d_2^2 (d_1 + d_2 - 2)}{d_1 (d_2 - 2)^2 (d_2 - 4)}$  for  $d_2 > 4$

# Two variances and the F distribution in R

- $P(X = x) = 0$  as the probability of any one value is always zero
- $P(X \leq x) = \text{pf}(x, n_x, n_y)$
- $P(X \geq x) = 1 - \text{pf}(x, n_x, n_y)$
- $P(x_1 < X < x_2) = \text{pf}(x_2, n_x, n_y) - \text{pf}(x_1, n_x, n_y)$

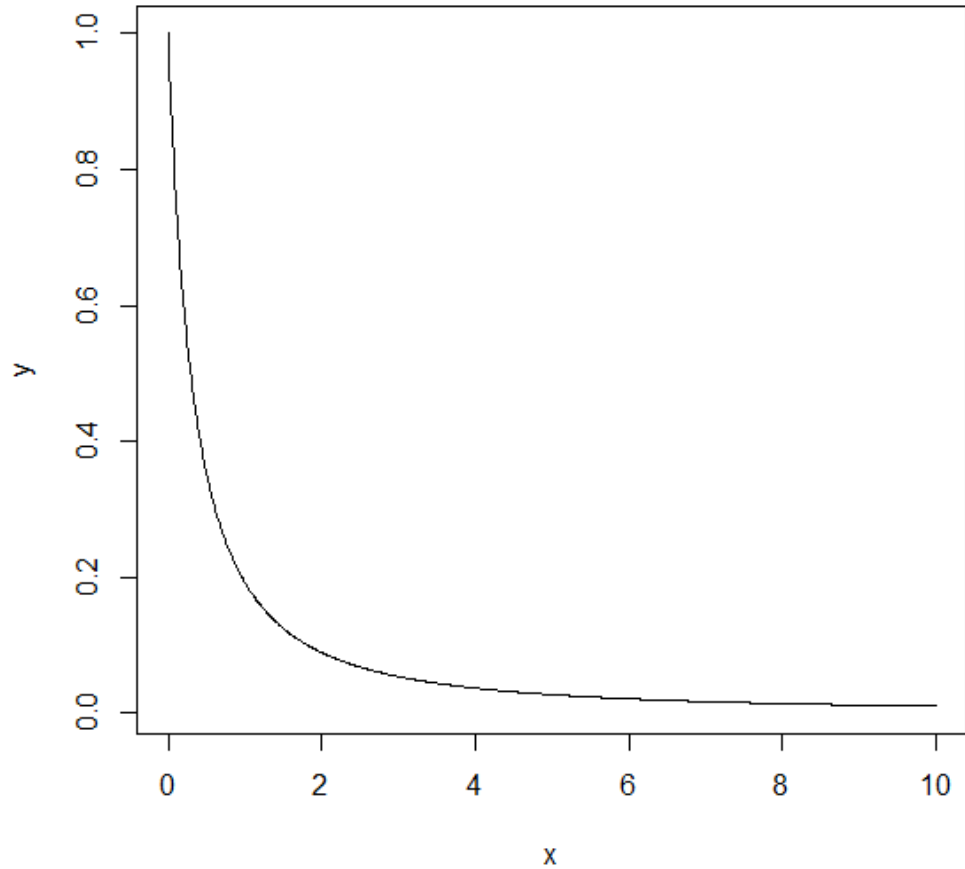
# Two variances and the F distribution

$$n_x = 1, n_y = 1$$



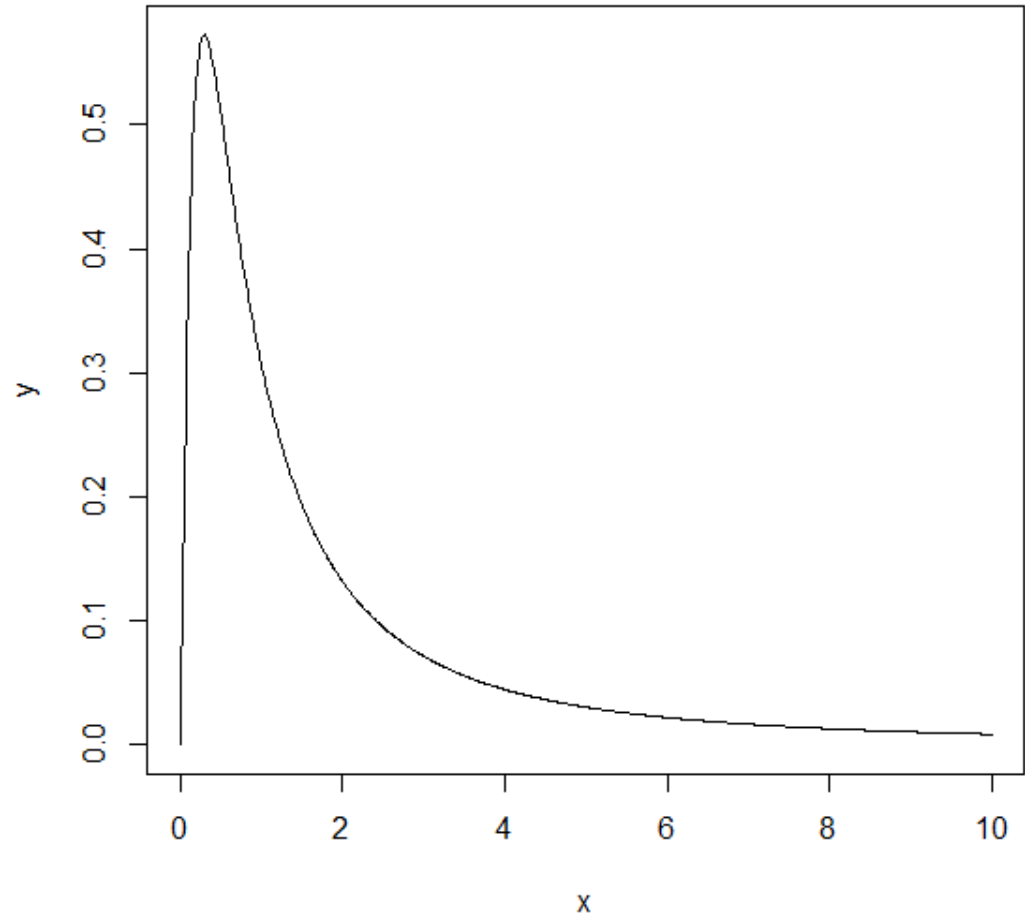
# Two variances and the F distribution

$$n_x = 2, n_y = 1$$



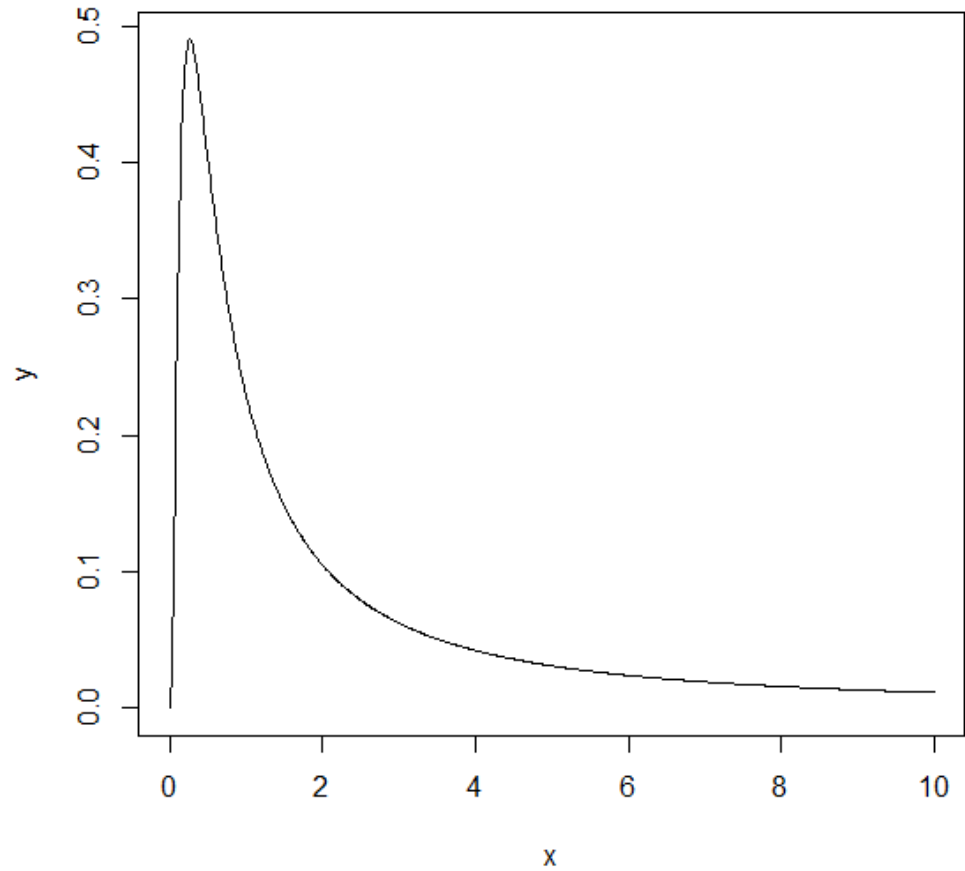
# Two variances and the F distribution

$$n_x = 5, n_y = 2$$



# Two variances and the F distribution

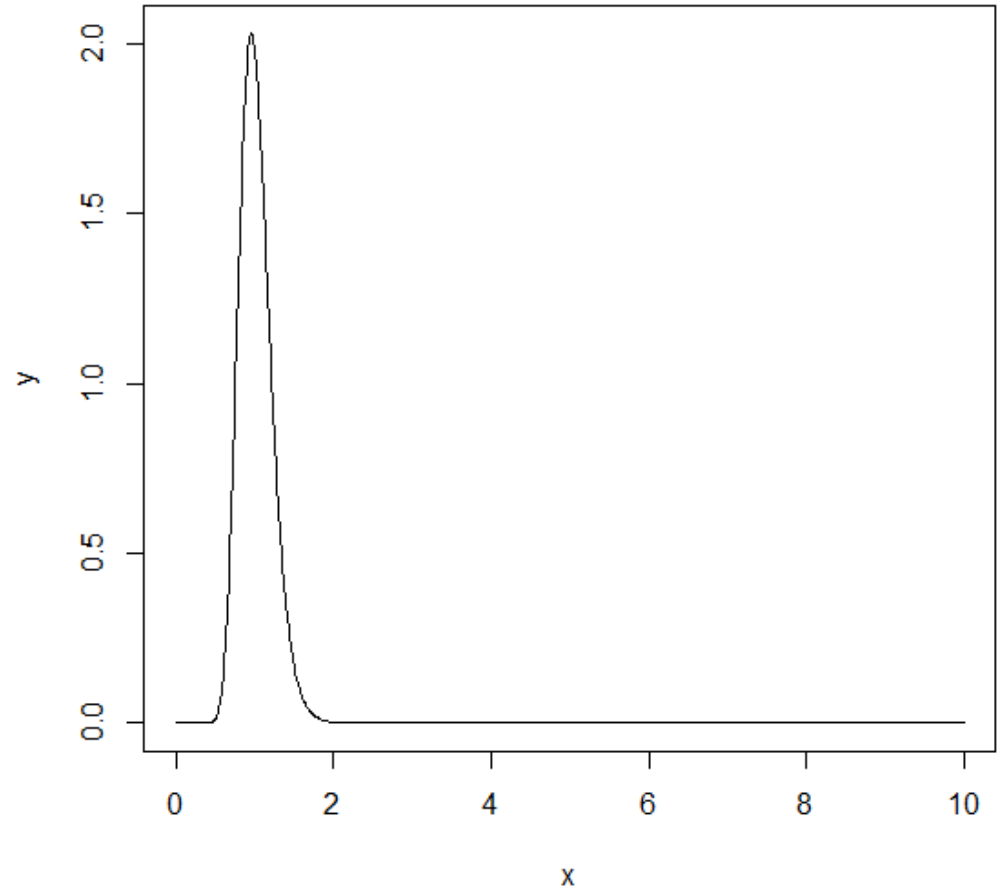
$$n_x = 10, n_y = 1$$





# Two variances and the F distribution

$$n_x = 100, n_y = 100$$



# Two variances and the F distribution

- The F distribution can be defined as follows:

- $$F_{n_x, n_y} = \frac{\left(\frac{X_x^2}{n_x}\right)}{\left(\frac{X_y^2}{n_y}\right)}$$

- Where  $X_x^2$  with  $n_x$  degrees of freedom and  $X_y^2$  with  $n_y$  degrees of freedom are independent  $\chi^2$  random variables

# Two variances and the F distribution

- Relating it back to  $\bar{x}$ ,  $s_x^2$ , and the t-distribution

$$F_{n_x-1, n_y-1} = \frac{\left( \frac{(n_x - 1)s_x^2}{\sigma_x^2} \right)}{\left( \frac{(n_y - 1)s_y^2}{\sigma_y^2} \right)} = \frac{\frac{s_x^2}{\sigma_x^2}}{\frac{s_y^2}{\sigma_y^2}} = \frac{\frac{s_x^2}{s_y^2}}{\frac{\sigma_x^2}{\sigma_y^2}}$$

# Two variances and the F distribution

- Thus allowing us to compare the variances of two populations

- We say  $\frac{\left(\frac{s_x^2}{s_y^2}\right)}{\left(\frac{\sigma_x^2}{\sigma_y^2}\right)}$  follows  $F_{n_x-1, n_y-1}$ , the F

distribution with  $n_x - 1$  and  $n_y - 1$  degrees of freedom

# Summaries!

- Bunnies, Rabbits and the NY times
  - <https://www.youtube.com/watch?v=jvoxEYmQHNM>

# Sampling Distribution for the Sample Proportion Summary

Shape of sample	Center of sample	Spread of sample
<p>The shape of the distribution is bell shaped if</p> <p><math>n * \rho \geq 15</math> <b>and</b> <math>n * (1 - \rho) \geq 15</math></p>	$\mu_{\hat{p}} = \rho$	$\sigma_{\hat{p}} = \sqrt{\frac{\rho(1 - \rho)}{n}}$

# Sampling Distribution for the Sample Mean Summary

Shape, Center and Spread of Population	Shape of sample	Center of sample	Spread of sample
Population is normal with mean $\mu$ and standard deviation $\sigma$ .	Regardless of the sample size $n$ , the shape of the distribution of the sample mean is normal	$\mu_{\bar{x}} = \mu$	$\sigma_{\bar{x}} = \frac{\sigma_x}{\sqrt{n}}$
Population is not normal with mean $\mu$ and standard deviation $\sigma$ .	As the sample size $n$ increases, the distribution of the sample mean becomes approximately normal	$\mu_{\bar{x}} = \mu$	$\sigma_{\bar{x}} = \frac{\sigma_x}{\sqrt{n}}$

$s_x^2$  and the  $\chi^2$  (chi-squared) distribution

- The chi-squared density can be defined as follows:

$$f_{X(x)} = \frac{1}{2^{\frac{n}{2}} \Gamma\left(\frac{n}{2}\right)} x^{\frac{n}{2}-1} e^{-\frac{x}{2}} I(X > 0)$$

- Mean =  $n$
- Variance =  $2n$



## $\chi^2$ (chi-squared) distribution in R

- $P(X = x) = 0$  as the probability of any one value is always zero
- $P(X \leq x) = \text{pchisq}(x, n)$
- $P(X \geq x) = 1 - \text{pchisq}(x, n)$
- $P(x_1 < X < x_2) = \text{pchisq}(x_2, n) - \text{pchisq}(x_1, n)$

$s_x^2$  and the  $\chi^2$  (chi-squared) distribution

- The chi-squared distribution can be defined as follows:

$$X^2 = z_1^2 + z_2^2 + \cdots + z_n^2$$

- **$X^2$  follows the chi-squared distribution with n degrees of freedom** where  $z_i$  are independent random variables that follow the standard normal distribution

# The Chi-Squared Distribution

$$X_{n-1}^2 = \left( \frac{(n-1)s^2}{\sigma_x^2} \right) = \sum \left( \frac{x_i - \bar{x}}{\sigma_x} \right)^2$$

- $\left( \frac{(n-1)s^2}{\sigma_x^2} \right)$  follows a  $\chi^2$  distribution with  $n-1$  degrees of freedom

## $\bar{x}$ , $s_x^2$ and the Student's t-distribution

- The t density can be defined as follows:

$$f_{X(x)} = \frac{\Gamma\left(\frac{n+1}{2}\right)}{\sqrt{\pi n} \Gamma\left(\frac{n}{2}\right)} \left(1 + \frac{x^2}{n}\right)^{-\frac{n+1}{2}} I(X \in \mathbb{R})$$

- Mean = 0
- Variance =  $\frac{n}{n-2}$

## $\bar{x}, s_x^2$ and the Student's t-distribution in R

- $P(X = x) = 0$  as the probability of any one value is always zero
- $P(X \leq x) = \text{pt}(x, n)$
- $P(X \geq x) = 1 - \text{pt}(x, n)$
- $P(x_1 < X < x_2) = \text{pt}(x_2, n) - \text{pt}(x_1, n)$

## $\bar{x}$ , $s_x^2$ and the Student's t-distribution

- The t distribution can be defined as follows:

$$t_n = \frac{Z}{\sqrt{\frac{\chi_n^2}{n}}}$$

- $t_n$  follows the t distribution with n degrees of freedom where Z follows the standard normal distribution and  $\chi^2$  is chi-squared and divided by its degrees of freedom

$\bar{x}$ ,  $s_x^2$  and the Student's t-distribution

$$t_{n-1} = \frac{Z}{\sqrt{\frac{\chi_n^2}{n}}} = \frac{\bar{x} - \mu_x}{(s_x/\sqrt{n})}$$

- $\frac{\bar{x} - \mu_x}{(s_x/\sqrt{n})}$  follows a t distribution with n-1 degrees of freedom
- Note the similarity to the z-score: the only difference here is that we estimate  $\sigma_x$  with  $s_x$

# Two variances and the F distribution

- The F density can be defined as follows:

$$f_{x(x)} = \frac{\left( \sqrt{\frac{(d_1 x)^2 d_2^2}{(d_1 x + d_2)^{(d_1 + d_2)}}} \right)}{xB \left( \frac{d_1}{2}, \frac{d_2}{2} \right)} I(x \geq 0)$$

- Mean =  $\frac{d_2}{d_2 - 2}$  for  $d_2 > 2$
- Variance =  $\frac{2 * d_2^2 (d_1 + d_2 - 2)}{d_1 (d_2 - 2)^2 (d_2 - 4)}$  for  $d_2 > 4$



# Two variances and the F distribution in R

- $P(X = x) = 0$  as the probability of any one value is always zero
- $P(X \leq x) = \text{pf}(x, n_x, n_y)$
- $P(X \geq x) = 1 - \text{pf}(x, n_x, n_y)$
- $P(x_1 < X < x_2) = \text{pf}(x_2, n_x, n_y) - \text{pf}(x_1, n_x, n_y)$

# Two variances and the F distribution

- The F distribution can be defined as follows:

- $$F_{n_x, n_y} = \frac{\left(\frac{X_x^2}{n_x}\right)}{\left(\frac{X_y^2}{n_y}\right)}$$

- Where  $X_x^2$  with  $n_x$  degrees of freedom and  $X_y^2$  with  $n_y$  degrees of freedom are independent  $\chi^2$  random variables

# Two variances and the F distribution

$$F_{n_x-1, n_y-1} = \frac{\left(\frac{S_x^2}{S_y^2}\right)}{\left(\frac{\sigma_x^2}{\sigma_y^2}\right)}$$

$\frac{\left(\frac{S_x^2}{S_y^2}\right)}{\left(\frac{\sigma_x^2}{\sigma_y^2}\right)}$  follows, the F distribution with  $n_x - 1$  and  $n_y - 1$  degrees of freedom