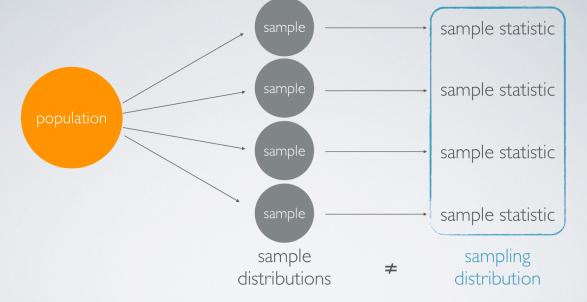


CLT and Sampling

https://gallery.shinyapps.io/CLT_mean/

<p>Sample vs Sampling Distribution</p>	<p>Sample:</p> <ul style="list-style-type: none"> - Distribution of data from the samples <p>Sampling</p> <ul style="list-style-type: none"> - Distribution of a statistic from several samples (ex. Average of all possible samples) 
<p>Central Limit Theorem;</p> <p>Because this theorem is “central” to statistic logics</p>	<p>CLT:</p> <ul style="list-style-type: none"> - Distribution of sample statistics is nearly normal, centered at the population mean, with standard error equal to population standard deviation divided by square root of sample size - N means the shape of the distribution - Mean = center - SE = spread of the graph $\bar{x} \sim N \left(\text{mean} = \mu, SE = \frac{\sigma}{\sqrt{n}} \right)$ <p style="text-align: center;"><i>shape center spread</i></p> <ul style="list-style-type: none"> - Note: if we don't have the population standard deviation (sigma), then we use the sample standard deviation (s) - Note: we are re-calibrating a population distribution by taking a sample and calibrating around the sample statistic (mean).
<p>Central Limit Theorem Conditions:</p>	<ol style="list-style-type: none"> 1. Independent: <ol style="list-style-type: none"> a. Random sample/assignment b. If sampling without replacement, $n < 10\%$ of population* 2. Sample size/skew:

	<p>a. Either the population distribution is normal, or if the population distribution is skewed, the sample size if large $n>30$</p> <p>*More on sampling without replacement:</p> <ul style="list-style-type: none"> - Sampling without replacement - don't take more than 10% of the population to avoid dependence (ex. Population has you and your extended family, if you take 60% of the population, chances are you and your extended family are in there and then the sample would not be independent)
Central Limit In Practice	<ol style="list-style-type: none"> 1. Distribution of a population 2. Distribution of a sample from the population (distribution most resembles 1) 3. Distribution of sample mean from size of 7 of sample (more normal than 2 given CLT, but more skewed than 4) 4. Distribution of sample mean from size of 100 (distribution is most normal given CLT since size if bigger) <p>*sample statistics is more normally distributed than the sample itself</p> <div style="border: 1px solid black; padding: 10px;"> <p>Four plots: Determine which plot (A, B, or C) is which.</p> <p>(1) The distribution for a population ($\mu = 10, \sigma = 7$), (2) a single random sample of 100 observations from this population, (3) a distribution of 100 sample means from random samples with size 7, and (4) a distribution of 100 sample means from random samples with size 49.</p> </div>
Central Limit Probability Example:	<p>Suppose my iPod has 3,000 songs. The histogram below shows the distribution of the lengths of these songs. We also know that, for this iPod, the mean length is 3.45 minutes and the standard deviation is 1.63 minutes. Calculate the probability that a randomly selected song lasts more than 5 minutes.</p> <div style="border: 1px solid black; padding: 10px; text-align: center;"> $X = \text{length of one song}$ $P(X > 5) = \frac{350 + 100 + 25 + 20 + 5}{3000} = \frac{500}{3000} \approx 0.17$ </div>

I'm about to take a trip to visit my parents and the drive is 6 hours. I make a random playlist of 100 songs. What is the probability that my playlist lasts the entire drive?

$$6 \text{ hours} = 360 \text{ minutes}$$

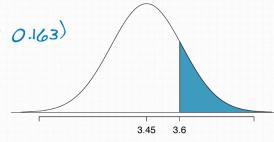
$$P(X_1 + X_2 + \dots + X_{100} > 360 \text{ min}) = ?$$

$$P(\bar{X} > 3.6) = ?$$

$$\bar{X} \sim N(\text{mean} = \mu = 3.45, SE = \frac{\sigma}{\sqrt{n}} = \frac{1.63}{\sqrt{100}} = 0.163)$$

$$Z = \frac{3.6 - 3.45}{0.163} = 0.92$$

$$P(Z > 0.92) = 0.179$$



Step 1: Population distribution is shown above as a bar chart

Step 2: To find the probability of a sample of 100 songs with total playtime of more than 6 hours, we create a distribution of sample statistics of the mean using CLT then find the z-score to calculate the probability that the average is greater than 3.6min

Given:

- n=3000
- mean=3.45
- sd=1.63

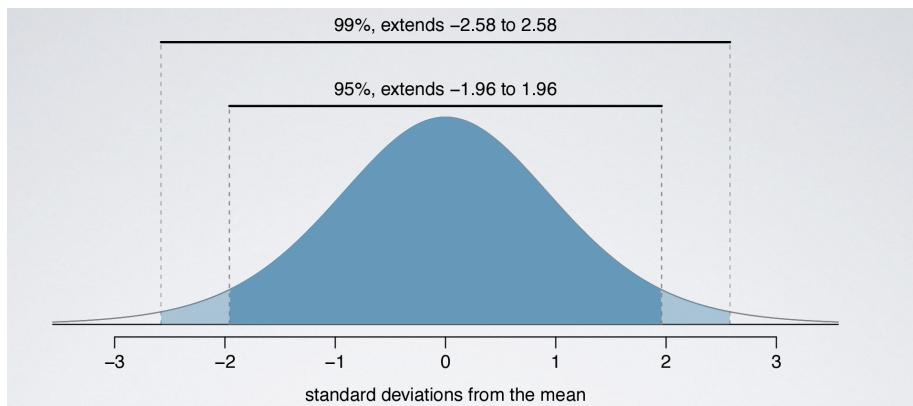
Find:

- $P(\bar{x} > 3.6)$
- Z-score = (observed value - mean)/standard error
- Why standard error rather than deviation? Because we are calculating for the sample statistic (mean) and the standard deviation of 1.63 applies to the length of songs.

Confidence Intervals (for a mean)

Confidence Interval	Plausible range of values for the population parameter Remember: x-bar is sample mean u is population mean
Confidence Interval and Central Limit Theorem:	The standard error lets us know how close the sample statistic's distribution is to the population mean. “The mean of the sample will be within x-standard errors of the population mean” Approximate confidence interval of 64% = mean +- SE Approximate confidence interval of 95% = mean +- 2SE

	<p>Approximate confidence interval of 99.7% = mean \pm 3SE</p> <p>Confidence interval for a population mean: Computed as the sample mean plus/minus a margin of error (critical value corresponding to the middle XX% of the normal distribution times the standard error of the sampling distribution).</p> $\bar{x} \pm z^* \frac{s}{\sqrt{n}}$																																																																																																																		
Confidence Interval Conditions: (sample as CLT)	<p>Conditions for this confidence interval:</p> <ol style="list-style-type: none"> 1. Independence: Sampled observations must be independent. <ul style="list-style-type: none"> ► random sample/assignment ► if sampling without replacement, $n < 10\%$ of population 2. Sample size/skew: $n \geq 30$, larger if the population distribution is very skewed. 																																																																																																																		
Finding the Confidence with Critical Values	<p>Critical Value = the specific standard deviation amount (z^*)</p> <ol style="list-style-type: none"> 1. Find the confidence of 95% (margin of error 5%) (Margin of Error)/2 $(1-0.95)/2$ 2. Find corresponding critical value for z <ol style="list-style-type: none"> a. Using z-table b. Using R (<code>qnorm(0.025)</code>) 3. Resulting z is the specific standard deviation of 1.96 <div style="display: flex; align-items: center;"> <div style="flex: 1;"> </div> <div style="flex: 1; text-align: right;"> <table border="1" style="width: 100%; border-collapse: collapse;"> <thead> <tr> <th></th> <th colspan="4">Second decimal place</th> <th></th> </tr> <tr> <th>0.07</th> <th>0.06</th> <th>0.05</th> <th>0.04</th> <th>0.00</th> <th>Z</th> </tr> </thead> <tbody> <tr><td>0.0003</td><td>0.0003</td><td>0.0003</td><td>0.0003</td><td>0.0003</td><td>-3.4</td></tr> <tr><td>0.0004</td><td>0.0004</td><td>0.0004</td><td>0.0004</td><td>0.0005</td><td>-3.3</td></tr> <tr><td>0.0005</td><td>0.0006</td><td>0.0006</td><td>0.0006</td><td>0.0007</td><td>-3.2</td></tr> <tr><td>0.0008</td><td>0.0008</td><td>0.0008</td><td>0.0008</td><td>0.0010</td><td>-3.1</td></tr> <tr><td>0.0011</td><td>0.0011</td><td>0.0011</td><td>0.0012</td><td>0.0013</td><td>-3.0</td></tr> <tr><td>0.0015</td><td>0.0015</td><td>0.0016</td><td>0.0016</td><td>0.0019</td><td>-2.9</td></tr> <tr><td>0.0021</td><td>0.0021</td><td>0.0022</td><td>0.0023</td><td>0.0026</td><td>-2.8</td></tr> <tr><td>0.0028</td><td>0.0029</td><td>0.0030</td><td>0.0031</td><td>0.0035</td><td>-2.7</td></tr> <tr><td>0.0038</td><td>0.0039</td><td>0.0040</td><td>0.0041</td><td>0.0047</td><td>-2.6</td></tr> <tr><td>0.0051</td><td>0.0052</td><td>0.0054</td><td>0.0055</td><td>0.0062</td><td>-2.5</td></tr> <tr><td>0.0068</td><td>0.0069</td><td>0.0071</td><td>0.0073</td><td>0.0082</td><td>-2.4</td></tr> <tr><td>0.0089</td><td>0.0091</td><td>0.0094</td><td>0.0096</td><td>0.0107</td><td>-2.3</td></tr> <tr><td>0.0116</td><td>0.0119</td><td>0.0122</td><td>0.0125</td><td>0.0139</td><td>-2.2</td></tr> <tr><td>0.0150</td><td>0.0154</td><td>0.0158</td><td>0.0162</td><td>0.0179</td><td>-2.1</td></tr> <tr><td>0.0192</td><td>0.0197</td><td>0.0202</td><td>0.0207</td><td>0.0228</td><td>-2.0</td></tr> <tr><td>0.0244</td><td>0.0250</td><td>0.0256</td><td>0.0262</td><td>0.0287</td><td>-1.9</td></tr> <tr><td>0.0307</td><td>0.0314</td><td>0.0322</td><td>0.0329</td><td>0.0359</td><td>-1.8</td></tr> </tbody> </table> </div> </div>		Second decimal place					0.07	0.06	0.05	0.04	0.00	Z	0.0003	0.0003	0.0003	0.0003	0.0003	-3.4	0.0004	0.0004	0.0004	0.0004	0.0005	-3.3	0.0005	0.0006	0.0006	0.0006	0.0007	-3.2	0.0008	0.0008	0.0008	0.0008	0.0010	-3.1	0.0011	0.0011	0.0011	0.0012	0.0013	-3.0	0.0015	0.0015	0.0016	0.0016	0.0019	-2.9	0.0021	0.0021	0.0022	0.0023	0.0026	-2.8	0.0028	0.0029	0.0030	0.0031	0.0035	-2.7	0.0038	0.0039	0.0040	0.0041	0.0047	-2.6	0.0051	0.0052	0.0054	0.0055	0.0062	-2.5	0.0068	0.0069	0.0071	0.0073	0.0082	-2.4	0.0089	0.0091	0.0094	0.0096	0.0107	-2.3	0.0116	0.0119	0.0122	0.0125	0.0139	-2.2	0.0150	0.0154	0.0158	0.0162	0.0179	-2.1	0.0192	0.0197	0.0202	0.0207	0.0228	-2.0	0.0244	0.0250	0.0256	0.0262	0.0287	-1.9	0.0307	0.0314	0.0322	0.0329	0.0359	-1.8
	Second decimal place																																																																																																																		
0.07	0.06	0.05	0.04	0.00	Z																																																																																																														
0.0003	0.0003	0.0003	0.0003	0.0003	-3.4																																																																																																														
0.0004	0.0004	0.0004	0.0004	0.0005	-3.3																																																																																																														
0.0005	0.0006	0.0006	0.0006	0.0007	-3.2																																																																																																														
0.0008	0.0008	0.0008	0.0008	0.0010	-3.1																																																																																																														
0.0011	0.0011	0.0011	0.0012	0.0013	-3.0																																																																																																														
0.0015	0.0015	0.0016	0.0016	0.0019	-2.9																																																																																																														
0.0021	0.0021	0.0022	0.0023	0.0026	-2.8																																																																																																														
0.0028	0.0029	0.0030	0.0031	0.0035	-2.7																																																																																																														
0.0038	0.0039	0.0040	0.0041	0.0047	-2.6																																																																																																														
0.0051	0.0052	0.0054	0.0055	0.0062	-2.5																																																																																																														
0.0068	0.0069	0.0071	0.0073	0.0082	-2.4																																																																																																														
0.0089	0.0091	0.0094	0.0096	0.0107	-2.3																																																																																																														
0.0116	0.0119	0.0122	0.0125	0.0139	-2.2																																																																																																														
0.0150	0.0154	0.0158	0.0162	0.0179	-2.1																																																																																																														
0.0192	0.0197	0.0202	0.0207	0.0228	-2.0																																																																																																														
0.0244	0.0250	0.0256	0.0262	0.0287	-1.9																																																																																																														
0.0307	0.0314	0.0322	0.0329	0.0359	-1.8																																																																																																														
Confidence Level Accuracy vs Precision	<p>The higher the confidence level, the wider the spread to “capture more range of data” so we can have more confidence that we captured the right number.</p>																																																																																																																		



However, this means decreased precision, aka the usefulness of the data. For example, a larger weather range may be more accurate, but the extremes in the range means we don't know what temperature to dress for.

CL ↑ width ↑ accuracy ↑

precision ↓

For higher precision and higher accuracy, increase sample size

Calculate the required sample size for target margin of error

$$ME = z^* \frac{s}{\sqrt{n}} \rightarrow n = \left(\frac{z^* s}{ME} \right)^2$$

Calculating sample size for target margin of error Example:

A group of researchers want to test the possible effect of an epilepsy medication taken by pregnant mothers on the cognitive development of their children. As evidence, they want to estimate the IQ scores of three-year-old children born to mothers who were on this medication during pregnancy.

Previous studies suggest that the SD of IQ scores of three-year-old children is 18 points.

How many such children should the researchers sample in order to obtain a 90% confidence interval with a margin of error less than or equal to 4 points?

$$ME \leq 4 \text{ pts} \quad 4 = 1.65 \frac{18}{\sqrt{n}} \rightarrow n = \left(\frac{1.65 \times 18}{4} \right)^2 = 55.13$$

$z^* = 1.65$ We need at least 56 such children in the sample to obtain a maximum margin of error of 4 points.

$$\sigma = 18$$

	z^* is calculated using the chart or by r where $qnorm(0.05)$
Interpreting Confidence Interval Examples:	<p>The General Social Survey asks: "For how many days during the past 30 days was your mental health, which includes stress, depression, and problems with emotions, not good?" Based on responses from 1,151 US residents, the survey reported a 95% confidence interval of 3.40 to 4.24 days in 2010. Interpret this interval in context of the data.</p> <p><i>We are 95% confident that Americans on average have 3.40 to 4.24 bad mental health days per month</i></p> <p>The General Social Survey asks: "For how many days during the past 30 days was your mental health, which includes stress, depression, and problems with emotions, not good?" Based on responses from 1,151 US residents, the survey reported a 95% confidence interval of 3.40 to 4.24 days in 2010.</p> <p>In this context, what does a 95% confidence level mean?</p> <p><i>95% of random samples of 1,151 Americans will yield CIs that capture the true population mean of number of bad mental health days per month.</i></p>
Confidence Interval Calculation	<p>A sample of 50 college students were asked how many exclusive relationships they've been in so far. The students in the sample had an average of 3.2 exclusive relationships, with a standard deviation of 1.74. In addition, the sample distribution was only slightly skewed to the right. Estimate the true average number of exclusive relationships based on this sample using a 95% confidence interval.</p> <p>Step 1: Does this meet the Conditions of Confidence Intervals?</p> <p>1. random sample & $50 < 10\%$ of all college students <i>We can assume that the number of exclusive relationships one student in the sample has been in is independent of another.</i></p> <p>2. $n > 30$ & not so skewed sample <i>We can assume that the sampling distribution of average number of exclusive relationships from samples of size 50 will be nearly normal.</i></p> <p>Step 2: Find CI $CI = \text{mean} \pm \text{margin of error}$</p> <p>Step 3: Find margin of error then put it together to find CI</p> $SE = \frac{s}{\sqrt{n}} = \frac{1.74}{\sqrt{50}} \approx 0.246$ $\begin{aligned} \bar{x} \pm z^* SE &= 3.2 \pm 1.96 (0.246) \\ &= 3.2 \pm 0.48 \\ &= (2.72, 3.68) \end{aligned}$ <p>Step 5: Interpret the Answer</p>

We are 95% confident that college students on average have been in 2.72 to 3.68 exclusive relationships.