

## Unit I Concept Summary

<p><b><u>The Who:</u></b> Cases or observations, “What you’re collecting data on”. Put in rows in StatCrunch.</p> <p><b><u>The What:</u></b> The variables, characteristics for each observation. Put in columns.</p>	<p style="text-align: center;"><b><u>Variables</u></b></p> <p><b><u>Quantitative:</u></b> Numerical and summary statistics like the mean and median make sense in context.</p> <p><b><u>Categorical:</u></b> Classify individuals. Usually words.</p> <p><b><u>Identifier:</u></b> A subgroup of categorical, these are labels that cannot really be analyzed statistically. May not be unique.</p>	
<p style="text-align: center;"><b><u>Graphs for Quantitative Variables</u></b></p> <ol style="list-style-type: none"> <li><b><u>Histogram</u></b>, to identify shape, visualize center, spread, unusual features.</li> <li><b><u>Boxplot</u></b>, shows minimum, quartiles, maximum and any official outliers.</li> <li><b><u>Dotplot</u></b>, good on StatCrunch to identify certain individuals, like largest, smallest, etc....</li> </ol>	<p style="text-align: center;"><b><u>Graphs for Categorical Variables</u></b></p> <ol style="list-style-type: none"> <li><b><u>Pie Charts</u></b>, stand alone or side-by-side to look for differences in groups.</li> <li><b><u>Bar Charts</u></b>, same idea, different display.</li> </ol>	<p style="text-align: center;"><b><u>Tables</u></b></p> <ol style="list-style-type: none"> <li><b><u>Frequency Tables</u></b>, to count up frequencies or % one variable at a time.</li> <li><b><u>Contingency Tables</u></b>, to count up frequencies or % two variables at a time.</li> </ol>
<p style="text-align: center;"><b><u>Describing the Distribution of a Quantitative Variable</u></b></p> <ol style="list-style-type: none"> <li><b><u>Shape</u></b>, symmetry, skewness, modes</li> <li><b><u>Center</u></b>, use the mean (roughly symmetric) or median (anything else), depending on shape.</li> <li><b><u>Spread</u></b>, use the standard deviation (roughly symmetric) or IQR (anything else), depending on shape.</li> <li><b><u>Official outliers, unusual shapes</u></b>, be specific, not generic.</li> </ol>	<p style="text-align: center;"><b><u>Analyzing Contingency Tables</u></b></p> <ol style="list-style-type: none"> <li>Row totals or column totals over the grand total give “<b><u>MARGINAL</u></b>” <b><u>proportions</u></b>.</li> <li>Intersections over the grand total give “<b><u>AND</u></b>” <b><u>proportions</u></b>.</li> <li>Intersections over column totals or row totals give “<b><u>CONDITIONAL</u></b>” <b><u>proportions</u></b>.</li> </ol>	<p style="text-align: center;"><b><u>Determining if Two Categorical Variables are Dependent on Each Other</u></b></p> <ol style="list-style-type: none"> <li>Using a contingency table, check for differences in marginal distributions and conditional distributions.</li> <li>Using side-by-side pie charts or stacked bar charts, look for substantial visual evidence of a difference in groups.</li> </ol>
<p style="text-align: center;"><b><u>Summary Statistics</u></b></p> <p><b>Sample Mean:</b> <math>\bar{y} = \frac{\sum y}{n}</math></p> <p><b>Sample Standard Deviation:</b></p> $s = \sqrt{\frac{\sum (y - \bar{y})^2}{n - 1}}$ <p><b>Sample Median:</b> The value in the middle position. If there isn’t one value, take the mean of the two values in the middle positions.</p> <p><b><math>Q_1</math>:</b> The first quartile, 25<sup>th</sup> percentile  <b><math>Q_3</math>:</b> The third quartile, 75<sup>th</sup> percentile  <b>IQR</b> = <math>Q_3 - Q_1</math></p> <p><b><math>P_{60}</math>:</b> The 60<sup>th</sup> percentile, etc...</p>	<p style="text-align: center;"><b><u>Summary Statistics</u></b></p> <p><b>Z-Score:</b> The number of standard deviations a data value lies from the mean.</p> $z = \frac{y - \bar{y}}{s}$ <p>Data values within 2 standard deviations of the mean are <b><u>not unusual</u></b>.</p> <p>Data values between 2 and 3 standard deviations from the mean <b><u>are unusual</u></b>.</p> <p>Data values more than 3 standard deviations away from the mean <b><u>are rare</u></b>.</p>	<p style="text-align: center;"><b><u>Summary Statistics</u></b></p> <p>Official outliers can be determined using the fences:</p> <p>Lower Fence = <math>Q_1 - 1.5(IQR)</math>  Upper Fence = <math>Q_3 + 1.5(IQR)</math></p> <p><b>Shifting:</b> Adding or subtracting the same number to every data value. Measures of position get shifted likewise. Measures of spread remain unchanged.</p> <p><b>Scaling:</b> Multiplying or dividing every data value by the same number. All summary statistics change likewise.</p>

## Unit I Concept Summary

<u>Definitions</u>	<u>Sampling Methods</u>
<p><b>Population:</b> The entire group of interest, usually obtaining data from one is impossible.</p> <p><b>Sample:</b> A subset of the population. In the real world, we typically deal with sample data.</p> <p><b>Observational Study:</b> A data collection method in which no manipulation or attempt is made to affect the outcome. In other words, obtaining data that already exists.</p> <p><b>Experiment:</b> A data collection method in which the experimenter manipulates treatments in order to see the effect on the response variable.</p>	<p><b>Simple Random Sample (SRS):</b> A sample collected in which every individual has an equal chance of being selected and every sample is equally likely to occur.</p> <p><b>Stratified Sample:</b> A sample collected in which the population is first broken down into groups (that are different). Then a SRS is taken from every group.</p> <p><b>Cluster Sample:</b> A sample collected in which the population is first broken down into groups (that are similar). Then clusters are randomly chosen and every member of the cluster is sampled.</p> <p><b>Systematic Sample:</b> A sample collected where every <math>k^{\text{th}}</math> individual is selected (like every 10<sup>th</sup>, e.g.).</p> <p><b>Multistage Sample:</b> A sample collected combining the above methods.</p> <p><b>Convenience Sample:</b> A sample collected with no randomization or usage of the above methods.</p>
<u>Designed Experiments</u>	
<p><b>Response Variable:</b> The variable of interest to the experimenter.</p> <p><b>Factors:</b> Variables controlled by the experimenter to see the effects on the response variable. Factor levels are randomly assigned to the experimental units or subjects.</p> <p><b>Blocking Factors:</b> Variables controlled by the experimenter to see the effects on the response variable. Blocking factors are not random, but pre-existing in the experimental units or subjects.</p> <p><b>Treatment:</b> The combination of factors and levels randomly assigned to the experimental units or subjects.</p> <p><b>Statistically Significant:</b> When the results of the experiment are too unusual to attribute to chance, the data / factor / experiment is said to be statistically significant.</p> <p><b>Control Group:</b> A treatment group in which no treatment, a baseline treatment, or a placebo treatment is applied.</p> <p><b>Replication:</b> Applying a treatment to multiple experimental units or subject.</p> <p><b>Lurking Variable:</b> A variable or factor that is the true reason for the results we see in an experiment (not attributed to a specific treatment).</p>	

## Unit II Correlation and Regression Summary

<p><b><u>Conditions for Linear Regression</u></b></p> <ol style="list-style-type: none"> <li>Both variables must be quantitative.</li> <li>The relationship must be linear when viewed on a scatterplot.</li> <li>The spread about the regression line must be equal for all <math>x</math>-values.</li> <li>All outliers and influential points must be investigated for accuracy and it must be determined if they are valid data points.</li> </ol> <p><b>Check a scatterplot and a Residuals vs. <math>X</math>-values plot on StatCrunch.</b></p> <ol style="list-style-type: none"> <li>Residual plots should have random scatter and equal spread left-to-right.</li> </ol>	<p><b><u>The Equation of the Linear Regression Line</u></b></p> $\hat{y} = b_1x + b_0$ <p>Slope Formula: <math>b_1 = r \left( \frac{s_y}{s_x} \right)</math></p> <p>y-Intercept Formula: <math>b_0 = \bar{y} - b_1\bar{x}</math></p> <p>Correlation Formula: <math>r = \frac{\sum(z_x z_y)}{n - 1}</math></p> <p><b><u>If you have data</u></b>, use StatCrunch features, not formulas.</p>	
<p><b><u>Interpreting Slope</u></b></p> <p>As the <math>x</math> variable increases by one unit, the slope tells us the change in the <math>y</math>-variable.</p> <p>This must be done in the context of the problem.</p>	<p><b><u>Interpreting <math>S_e</math></u></b></p> <p>“On average, our predicted ‘<math>y</math>-variable’s values are off by <math>\sim s_e</math> when using <math>x</math> = our ‘<math>x</math>-variable’”. We want this number to be small.</p> <p><b><u>Interpreting the y-Intercept</u></b></p> <p>The <math>y</math>-intercept is the value of the <math>y</math>-variable when <math>x = 0</math>.</p> <p>Be sure you have data at or very close to <math>x = 0</math>, and that the value of the <math>y</math>-intercept makes sense in the context of the problem. Otherwise, this value has no interpretable meaning in the context of the problem.</p>	
<p><b><u>Residuals</u></b></p> <ol style="list-style-type: none"> <li>We want residuals to be small.</li> <li><math>e</math> = Actual <math>y</math> – Predicted <math>y</math></li> <li><math>e = y - \hat{y}</math></li> <li>Residuals can be saved using the StatCrunch checkbox.</li> <li>Be able to interpret residuals in the context of the problem.</li> <li>Data points above the regression line have positive residuals and actual values are greater than the predicted values.</li> <li>Data points below the regression line have negative residuals and actual values are less than the predicted values.</li> <li>Residuals are in the same units as the <math>y</math>-variable.</li> </ol>	<p><b><u>Studentized Residuals</u></b></p> <ol style="list-style-type: none"> <li>Studentized residuals can be saved using the StatCrunch checkbox.</li> <li>Studentized residuals are in “standardized units”.</li> <li>Any Studentized residual exceeding <math>\pm 2</math> should be flagged. The original row of data should be checked for accuracy and appropriateness.</li> </ol>	<p><b><u>Cook’s Distances</u></b></p> <ol style="list-style-type: none"> <li>Cook’s distances can be saved using the StatCrunch checkbox.</li> <li>Any Cook’s distance exceeding <math>4 / n</math> should be flagged as an influential point.</li> <li>Any influential points should be checked for accuracy and appropriateness.</li> <li>Do not remove data points from the dataset without good reason. Outliers, data points with large Studentized residuals, or large Cook’s distances are investigated.</li> <li>One strategy for valid influential points is to run the regression twice, with and without the points, to see the effect.</li> </ol>
<p><b><u><math>R^2</math></u></b></p> <ol style="list-style-type: none"> <li><math>0\% \leq R^2 \leq 100\%</math>, with higher being better.</li> <li><math>R^2</math> measures the percentage of variation in the <math>y</math>-variable accounted for by the linear regression on the <math>x</math>-variable.</li> <li>When interpreting, you must put it into the context of the problem.</li> </ol>	<p><b><u><math>r</math></u></b></p> <ol style="list-style-type: none"> <li><math>-1 \leq r \leq 1</math></li> <li>Correlations from 0 up to 0.5 are weak.</li> <li>Correlations between 0.5 and 0.8 are moderate.</li> <li>Correlations between 0.8 and 1 are strong.</li> </ol>	<p><b><u>Prediction</u></b></p> <ol style="list-style-type: none"> <li>Plug in the desired <math>x</math>-value into the linear equation and solve for the predicted <math>y</math>-value.</li> <li>StatCrunch features should be used for the most accuracy.</li> <li>Only predict inside the <math>x</math>-range from the data. Predicting outside is extrapolation.</li> </ol>

## Unit II Probability Summary

<p style="text-align: center;"><b><u>Probability Rules</u></b></p> <ol style="list-style-type: none"> <li>1. <math>P(A^c) = 1 - P(A)</math></li> <li>2. <math>P(A \text{ or } B) = P(A) + P(B)</math> if A and B are disjoint.</li> <li>3. <math>P(A \text{ and } B) = P(A)P(B)</math> if A and B are independent.</li> <li>4. <math>P(A B) = P(A \text{ and } B) / P(B)</math></li> <li>5. <math>P(A \text{ or } B) = P(A) + P(B) - P(A \text{ and } B)</math> for all A, B.</li> <li>6. <math>P(A \text{ and } B) = P(A)P(B A)</math> for all A, B.</li> <li>7. A and B are independent if <math>P(A) = P(A B)</math>.</li> <li>8. A and B are disjoint if <math>P(A \text{ and } B) = 0</math>.</li> </ol>	<p style="text-align: center;"><b><u>Uniform Distributions (Continuous)</u></b></p> <ol style="list-style-type: none"> <li>1. Between lower and upper boundaries <math>a</math> and <math>b</math>, all probabilities of intervals of equal length are equally likely.</li> <li>2. Rectangular in shape. Height of rectangle is <math>\frac{1}{b-a}</math>.</li> <li>3. Probability = Area under rectangle.</li> <li>4. Percentiles can be found by putting the corresponding area to the left and finding the value on the <math>x</math>-axis.</li> <li>5. The mean is the midpoint: <math>\mu = \frac{a+b}{2}</math></li> <li>6. The probability function is <math>f(x) = \frac{1}{b-a}, a \leq x \leq b</math></li> </ol>
<p style="text-align: center;"><b><u>“At Least” Problems</u></b></p> <ol style="list-style-type: none"> <li>1. In words, determine the complement of what the problem is asking.</li> <li>2. Utilize the complement rule. Example given:   <math display="block">P(\text{at least one person is Canadian}) = 1 - P(\text{no one is Canadian})</math> </li> </ol>	<p style="text-align: center;"><b><u>Exponential Distributions (Continuous)</u></b></p> <ol style="list-style-type: none"> <li>1. Used to model probabilities for the time between events.</li> <li>2. Unimodal and skewed right, short times between events are more likely to occur.</li> <li>3. Probability = Area under the curve.</li> <li>4. Probabilities for intervals of values and percentiles can be found using Stat → Calculators → Exponential.</li> <li>5. Need to be told the mean, <math>\mu</math>, which is coincidentally the same as the standard deviation, <math>\sigma</math>.</li> </ol>
<p style="text-align: center;"><b><u>(Discrete) Distributions Presented in Table Form</u></b></p> <ol style="list-style-type: none"> <li>1. Will be given a list of quantitative outcomes and their respective probabilities.</li> <li>2. The mean can be calculated using:   <math display="block">\mu = \sum xP(x)</math> </li> </ol>	<p style="text-align: center;"><b><u>Normal Distributions (Continuous)</u></b></p> <ol style="list-style-type: none"> <li>1. Used to model probabilities with most values piling up near the mean and then tapering off in both directions.</li> <li>2. Unimodal and symmetric, defined by the mean, <math>\mu</math>, and standard deviation, <math>\sigma</math>, which will be given in the problem.</li> <li>3. Probability = Area under the curve.</li> <li>4. Probabilities for intervals of values and percentiles can be found using Stat → Calculators → Normal.</li> </ol>
<p style="text-align: center;"><b><u>Checking Data For Normality</u></b></p> <ol style="list-style-type: none"> <li>1. Make a histogram and look for the familiar bell-shaped distribution: unimodal and symmetric.</li> <li>2. Make a QQ plot and look for the data points to be generally straight on the diagonal.</li> </ol>	
<p style="text-align: center;"><b><u>Binomial Distributions (Discrete)</u></b></p> <ol style="list-style-type: none"> <li>1. Fixed number of trials, two outcomes per trial, <math>P(\text{Success}) = p</math> is constant, trials independent.</li> <li>2. Binomial random variables count up the number of successes, <math>x</math>, in <math>n</math> trials.</li> <li>3. <math>P(X = x) = \binom{n}{x} p^x (1-p)^{n-x}, 0 \leq x \leq n</math></li> <li>4. <math>\mu = np</math> and <math>\sigma = \sqrt{np(1-p)}</math></li> <li>5. Can use the StatCrunch calculator rather than the formula, Stat → Calculators → Binomial.</li> </ol>	

## Unit III Inference for Proportions Summary

<u>Conditions for Inference for Proportions</u>	<u>The Sampling Distribution of p-Hat</u>
<ol style="list-style-type: none"> <li>1. The variable or variables are categorical.</li> <li>2. The sample(s) are random or at least unbiased.</li> <li>3. The sample(s) size(s) are less than 10% of the population(s) size(s).</li> <li>4. We have or expect to have at least 10 successes and 10 failures in our sample(s).</li> </ol>	<p>As long as the conditions are met, the sample proportion <math>\hat{p}</math> will follow an approximate Normal distribution. The mean and standard deviation of this distribution are as follows:</p> $\mu_{\hat{p}} = p$ $\sigma_{\hat{p}} = \sqrt{\frac{p(1-p)}{n}}$ <p>This distribution is more for theoretical purposes, because in practice, we are investigating the value of <math>p</math>. We won't know it for real world problems, but this is the basis for our confidence intervals and hypothesis tests.</p>
<p style="text-align: center;"><b><u>One-Sample Confidence Interval for the Population Proportion</u></b></p> <p><math>\hat{p} \pm z \sqrt{\frac{\hat{p}(1-\hat{p})}{n}}</math> or use Stat → Proportion Stats → One Sample with data or summary → Confidence Interval.</p> <p>Margin of Error = <math>z \sqrt{\frac{\hat{p}(1-\hat{p})}{n}} = \frac{\text{Upper Bound} - \text{Lower Bound}}{2}</math></p> <p>Calculate Sample Size: <math>n = \frac{z^2 \hat{p}(1-\hat{p})}{(\text{ME})^2}</math> with educated guess for <math>\hat{p}</math> or default to 0.50 if necessary.</p>	
<p style="text-align: center;"><b><u>One-Sample Test for a Population Proportion</u></b></p> <p>Test the hypotheses <math>H_0 : p = p_0</math> vs <math>H_A : p &lt; \text{ or } &gt; \text{ or } \neq p_0</math> using the test statistic <math>z = \frac{\hat{p} - p_0}{\sqrt{\frac{p_0(1-p_0)}{n}}}</math>.</p> <p>Determine the P-value by shading under the Standard Normal model in the <math>H_A</math> direction.</p> <p>Or on StatCrunch, Stat → Proportion Stats → One Sample with data or summary → Hypothesis Test.</p> <p>Reject the null hypothesis in favor of the alternative hypothesis is the P-value is “low” or if the P-value is less than the stated significance level. See page 232 in the notes for the P-value diagram.</p>	
<p style="text-align: center;"><b><u>Two-Sample Intervals and Tests for a Difference in Proportions</u></b></p> <div style="display: flex; justify-content: space-between;"> <div style="width: 45%;"> <p><math>(\hat{p}_1 - \hat{p}_2) \pm z \sqrt{\frac{\hat{p}_1(1-\hat{p}_1)}{n_1} + \frac{\hat{p}_2(1-\hat{p}_2)}{n_2}}</math></p> <p><math>z = \frac{(\hat{p}_1 - \hat{p}_2)}{\sqrt{\frac{\hat{p}_{\text{pooled}}(1-\hat{p}_{\text{pooled}})}{n_1} + \frac{\hat{p}_{\text{pooled}}(1-\hat{p}_{\text{pooled}})}{n_2}}}</math></p> </div> <div style="width: 45%;"> <p><math>H_0 : p_1 = p_2</math> versus <math>H_A : p_1 &lt; \text{ or } &gt; \text{ or } \neq p_2</math></p> <p><math>\hat{p}_{\text{pooled}} = \frac{x_1 + x_2}{n_1 + n_2}</math></p> </div> </div> <p>StatCrunch, Stat → Proportion Stats → Two Sample with data or summary → Hypothesis Test or Confidence Interval.</p>	

Common Critical  $z$  Values:    90% →  $z = 1.645$     95% →  $z = 1.96$     98% →  $z = 2.326$     99% →  $z = 2.576$

## Unit III Inference for Means Summary

<u>Conditions for Inference for Means</u>	<u>The Sampling Distribution of <math>\bar{y}</math>-Bar</u>
<p>1. The variable or variables are quantitative and summarized using a mean.</p> <p>2. The sample(s) are random or at least unbiased.</p> <p>3. The sample(s) size(s) are less than 10% of the population(s) size(s).</p> <p>4. We have reason to believe the population is approximately Normal or our sample size(s) are at least size 30.</p>	<p>As long as the conditions are met, the sample mean <math>\bar{y}</math> will follow an approximate Normal distribution. The mean and standard deviation of this distribution are as follows:</p> $\mu_{\bar{y}} = \mu_y$ $\sigma_{\bar{y}} = \frac{\sigma_y}{\sqrt{n}}$ <p>This distribution is more for theoretical purposes, because in practice, we are investigating the true value of <math>\mu</math>. We won't know <math>\mu</math> or <math>\sigma</math> for real world problems, but this is the basis for our confidence intervals and hypothesis tests.</p>
<p style="text-align: center;"><b><u>Important</u></b></p> <p>When we run hypothesis tests or compute confidence intervals for means, the true value of <math>\sigma</math> will be unknown! We substitute the value of <math>s</math> from our collected data and as a result, we must use the Student's <math>t</math> model instead of the Standard Normal model.</p>	
<p style="text-align: center;"><b><u>One-Sample Confidence Interval for the Population Mean</u></b></p> <p><math>\bar{y} \pm t \left( \frac{s}{\sqrt{n}} \right)</math> with <math>df = n - 1</math> or use Stat <math>\rightarrow</math> T Stats <math>\rightarrow</math> One Sample with data or summary <math>\rightarrow</math> Confidence Interval.</p> <p>Margin of Error = <math>t \left( \frac{s}{\sqrt{n}} \right) = \frac{\text{Upper Bound} - \text{Lower Bound}}{2}</math></p> <p>Calculate Sample Size: <math>n = \left( \frac{z(\text{Best Guess for Standard Deviation})}{(\text{ME})} \right)^2</math>.</p>	
<p style="text-align: center;"><b><u>One-Sample Test for a Population Mean</u></b></p> <p>Test the hypotheses <math>H_0 : \mu = \mu_0</math> vs <math>H_A : \mu &lt; \text{ or } &gt; \text{ or } \neq \mu_0</math> using the test statistic <math>t = \frac{\bar{y} - \mu_0}{\left( \frac{s}{\sqrt{n}} \right)}</math>.</p> <p>Determine the P-value by shading under the Student's <math>t</math> model (<math>df = n - 1</math>) in the <math>H_A</math> direction.</p> <p>Or on StatCrunch, Stat <math>\rightarrow</math> T Stats <math>\rightarrow</math> One Sample with data or summary <math>\rightarrow</math> Hypothesis Test.</p>	
<p style="text-align: center;"><b><u>Two-Sample Intervals and Tests for a Difference in Means (Independent Samples, df from technology)</u></b></p> <p><math>H_0 : \mu_1 = \mu_2</math> versus <math>H_A : \mu_1 &lt; \text{ or } &gt; \text{ or } \neq \mu_2</math></p> $t = \frac{(\bar{y}_1 - \bar{y}_2)}{\sqrt{\frac{s_1^2}{n_1} + \frac{s_2^2}{n_2}}} \quad (\bar{y}_1 - \bar{y}_2) \pm t \sqrt{\frac{s_1^2}{n_1} + \frac{s_2^2}{n_2}}$ <p>StatCrunch, Stat <math>\rightarrow</math> T Stats <math>\rightarrow</math> Two Sample with data or summary <math>\rightarrow</math> Hypothesis Test or Confidence Interval.</p>	
<p style="text-align: center;"><b><u>Two-Sample Intervals and Tests for the Mean Difference (Dependent Sample)</u></b></p> <p>Take the differences for each paired data value and run a one-sample <math>t</math> test on the column of differences.</p>	



