



## ST1232 cheat sheet

Statistics for Life Sciences (National University of Singapore)

<b>Hypothesis Testing</b> 1. Assumptions: - Data comes from randomization - Shape of population distribution is normal - Sample size large enough 2. Hypothesis: - $H_0$ : null hypothesis - $H_1$ : parameter falls in alternative range of values 3. Test statistic: - How far point estimate falls from the parameter value given in null hypothesis - Measured by number of standard errors between point estimate and parameter value in $H_0$ 4. Interpreting test statistic: - Presume $H_0$ is true - But if sample test statistic falls out in the tail of sampling distribution, value is unlikely - P-value is the probability of "what we saw or something more extreme" given $H_0$ is true	<b>Significance level:</b> - a number such that we reject $H_0$ if p-value is less than or equal to that number - if we reject $H_0$ , the results are statistically different <b>Hypothesis Testing for proportions</b> 1. Assumptions: - Variable collected is <b>categorical</b> - Data comes from randomization - Sample size $n$ is sufficiently <b>large</b> that the sampling distribution of sample proportion $\hat{p}$ is approximately normal if $np_0(1-p_0) \geq 5$ 2. Hypothesis - $H_0: p = p_0$ - $H_1: p \neq p_0$ 3. Test statistic: - $Z = \frac{\hat{p} - p_0}{\sqrt{\frac{p_0(1-p_0)}{n}}}$ $N(p_0, \frac{p_0(1-p_0)}{n})$ when $H_0$ is true $Z \sim (0,1)$	<b>Hypothesis testing for mean (One sample t-test)</b> 1. Assumptions: - Variable collected is <b>quantitative</b> - Data comes from randomization - Population distribution is approximately normal - $n$ is <b>small</b> $\rightarrow$ one sided t-test 2. Hypothesis: - $H_0: \mu = \mu_0$ 3. Test statistic: - Distance between sample mean $\bar{X}$ and $H_0$ value of population mean - Measured in terms of standard errors $T = \frac{\bar{X} - \mu_0}{s/\sqrt{n}}$ with $n-1$ degrees of freedom 4. Conclusion: - If $H_0$ is true - Mean is not significantly different from $\mu_0$	<b>Normality assumption</b> 1. Histogram with a Normal pdf overlaid - Symmetric, unimodal 2. Q-Q plot - Standardized sample quantiles against theoretical quantiles of $N(0,1)$ distribution - If fall on a straight line, the data is normal - Right tail below straight line: fatter - Left tail below straight line: thinner <b>Errors in conclusions</b> 1. Type 1 error: occurs if we reject $H_0$ when in fact it is true - when $H_0$ is assumed to be true 2. Type 2 error: occurs if we do not reject $H_0$ when it is false - When $H_1$ is assumed to be true 3. $\uparrow$ sample size $\rightarrow$ type 2 $\downarrow$ 4. $\downarrow$ significance level, type 1 $\downarrow$ but type 2 $\uparrow$	<b>2 independent sample t-test, Equal variance</b> 1. Assumptions: - <b>Independent</b> samples with randomized experiment - Sample standard deviations not twice... - Population distribution is approximately normal and $n$ is <b>small</b> 2. Hypothesis: - $H_0: \mu_1 = \mu_2$ 3. Test statistic: - Point estimate of difference between the population means $\bar{X} - \bar{Y}$ - Estimate of common variance: $s_p^2 = \frac{(n_1-1)s_1^2 + (n_2-1)s_2^2}{n_1+n_2-2}$ $T = \frac{(\bar{X} - \bar{Y}) - 0}{se}$ with $n_1+n_2-2$ degrees of freedom - $se = s_p \sqrt{\frac{1}{n_1} + \frac{1}{n_2}}$	<b>2 independent sample t-test, Unequal variance</b> 1. Assumptions: - Quantitative variable for 2 groups, independent, random - Sample standard deviation twice... - Approximately normal and $n$ is <b>small</b> 2. Hypothesis: - $H_0: \mu_1 = \mu_2$ 3. Test statistic: - Point estimate of difference between the population means $\bar{X} - \bar{Y}$ $T = \frac{(\bar{X} - \bar{Y}) - 0}{se}$ with complicated number of degrees of freedom (df) - $se = \sqrt{\frac{s_1^2}{n_1} + \frac{s_2^2}{n_2}}$ 4. Conclusion: - If do not reject $H_0$ , no statistically difference between means			
<b>Dependent sample test</b> -Testing if $\mu_1 = \mu_2$ 1. Construct new variable $-D_i = X_i - Y_i$ 2. Hypothesis: - $H_0: \mu_D = 0$ 3. Test statistic: $T = \frac{\bar{D} - \mu_0}{s/\sqrt{n}}$ with $n-1$ degrees of freedom 4. Pairing up reduced standard error in $\bar{D}$	<b>ANOVA: Compare more than 2 Group means</b> 1. Assumptions: - Population distribution for the response for k groups are normal - Largest group sd less than twice of smallest sd - Data randomized - N observations in total: n observations in k group, $N=nk$ 2. Hypothesis: - $H_0: \mu_1 = \mu_2 = \dots = \mu_k$ - $H_1$ : At least 2 means are not equal 3. Test statistic: - Involves within grp &	<b>Multiple comparison test</b> -To control Experiment error rate $\alpha$ : the probability of making at least 1 Type 1 error when all m hypothesis are true <b>1. Bonferroni correction:</b> - Not doing ANOVA - perform each test at $\alpha/m$ significance level - for each CI at $(1-\alpha/m)100\%$ level - p-value $\hat{a}/m$ <b>2. Tukey:</b> - replacing use of t-distribution for each pairwise distribution with a multiplier from another	<b>Wilcoxon signed rank For matched pairs</b> 1. Assumptions: - 2 groups are independent - Data random - Difference of observations <b>can be ranked</b> - Population distribution of difference is <b>symmetric</b> 2. Hypothesis: - $H_0: q_{0.5} = 0$ 3. Test statistic: - Sum of the ranks of the difference that were	<b>2-sample Wilcoxon rank sum (Mann-Whitney U)</b> 1. Assumptions: - Independent observations $X_1, X_2, \dots, X_{n1}$ and $Y_1, Y_2, \dots, Y_{n2}$ from 2 groups - Random data Hypothesis: - $H_0$ : 2 samples are from same distribution 2. Test statistic: - Minimum of $R_x$ and $R_y$ , $R_x$	<b>Kruskal-Wallis</b> -If <b>more than 2 groups of independent samples</b> -Not normal and unequal variance 1. Assume: -Independent observations from <b>k groups where k <math>\geq 2</math></b> -Random data 2. Hypothesis: - $H_0$ : The distributions of all groups are the same	<b>Simple linear Regression</b> -F-test equals t-test 1. Assumptions: - Random data - R/S between X and Y is linear $\epsilon \sim N(0, \sigma^2)$ normally distributed subpopulation for each explanatory variable and same variance $Y \sim N(\beta_0 + \beta_1 X)$ 2. t-test Hypothesis:	<b>Multiple linear Regression</b> -more than one explanatory variable -t-test for individual coefficient (quantitative) - <b>F-test for overall significance (categorical)</b> -Use <b>adjusted <math>R^2</math></b> to compare models <b>Indicator variables</b> $\hat{Y} = \beta_0 + \beta_1 X_1 + \beta_2 X_2$ when reference=female (0) & male (1) 1. Hypothesis: - $H_0$ : all coefficients=0 - $H_1$ : at least 1 not zero	<b>Point estimates</b> - $\hat{\beta}_0$ (constant) and $\hat{\beta}_1$ are point estimates - $\hat{Y}$ is a point estimate of the <b>mean of response</b> at a particular X value -when X increases by 1 unit, mean response $\hat{Y}$ increases by $\hat{\beta}_1$ unit

<p><b>Compare within 2 groups</b></p> <p>-H<sub>0</sub>: <math>\mu_i = \mu_j</math></p> <p>if F-test is significant</p> <p>-</p> $T = \frac{\bar{Y}_i - \bar{Y}_j}{se}$ <p>where se=</p> $s_p \sqrt{\frac{1}{n_i} + \frac{1}{n_j}}$ <p>with N-k degrees of freedom</p>	<p>between grp variability</p> <p>- F=</p> <p><b>Btwn grp vari</b></p> <p><b>within grp vari</b></p> <p>with numerator k-1 and denominator N-k degrees of freedom</p> <p>4. p-value: always the areas to the right of it; larger F → against H<sub>0</sub></p>	<p>distribution (a studentized range distribution)</p> <p>- input: no. of k grps; output: all significant pairwise comparisons</p> <p><b>3. Dunnett:</b></p> <p>-k groups with one control, conduct m=k-1 comparisons</p> <p>-same as Tukey but will provide shorter CI</p>	<p>positive (W+)</p> <p>- If <math>W_+ \approx W_-</math>, do not reject H<sub>0</sub></p> <p><b>For one sample</b></p> <p>1. Hypothesis:</p> <p>- H<sub>0</sub>:</p> $q_{0.5} = m_0$ <p>2. Test statistic:</p> <p>-sum of ranks</p> <p><math>\hat{m}_0</math> known as (W+)</p>	<p>sample while R<sub>v</sub> is the sum of ranks of Y-sample</p> <p>- H<sub>0</sub> is correct if <math>R_x \approx R_v</math></p>	<p>3. Test stats:</p> <p>-Uses the variability of sample-mean ranks</p> <p>-Under H<sub>0</sub>, follows a <math>\chi^2_{k-1}</math> distribution (chi-square)</p>	<p>- H<sub>0</sub>: <math>\beta_1 = 0</math> and t-distribution follows n-2 degrees of freedom</p> <p>3. F-test Hypothesis:</p> <p>- H<sub>0</sub>: <math>\beta_1 = 0</math> and F distribution with 1 and n-2 degrees of freedom</p>	<p>2. <b>Adjusted R<sup>2</sup></b></p> <p>=0.068 means only 6.8% of variation in the marks change is explained by regression model</p> <p><b>With interaction term</b></p> $\hat{Y} = \beta_1 + \beta_1 X_1 + \beta_2$	<p><b>Estimating <math>\sigma^2</math></b></p> <p>-how variable the response measurements are around fitted line</p> <p>- <math>\hat{\sigma}^2</math> is computed using residuals</p> $e_i = Y_i - \hat{Y}$ <p>-Variance is <b>minimized at</b></p> $X_0 = \bar{X}$ <p><b>, smaller the better</b></p>
<p><b>Scatterplot</b></p> <p>1. Linearity violated:</p> <ul style="list-style-type: none"><li>- Add higher order terms in X eg. X<sup>2</sup></li></ul> <p>2. Variance not constant:</p> <ul style="list-style-type: none"><li>- Transform the response by taking ln(Y), <math>\sqrt{Y}</math> or (1/Y)</li><li>- 1 unit increase in X, mean response <math>\hat{Y}</math> increases by <math>e^{\hat{\beta}_1}</math> times</li></ul> <p><b>Residuals plot</b></p> <ul style="list-style-type: none"><li>- Detect non-normality</li><li>- Check for non-constant variance and the need to transform Y</li><li>- Check for the need to add higher order terms in X</li></ul> <p><b>Explanatory data analysis</b></p> <p><b>1. Frequency table:</b></p> <ul style="list-style-type: none"><li>-For categorical or <b>quantitative variable</b></li><li>-Modal category: highest frequency</li></ul>	<p><b>Plots to make</b></p> <p>1. Plot r<sub>i</sub>'s against <math>\hat{Y}_i</math> on</p> <p>2. Plot r<sub>i</sub>'s against X<sub>i</sub></p> <p>3. Create histogram of r<sub>i</sub></p> <p>4. Create QQ plot using r<sub>i</sub></p> <ul style="list-style-type: none"><li>• Transform Y when there is funnel shape</li><li>• Add higher order X term when there is a curve band</li><li>• Non-normal when points are outside 3 and -3</li></ul> <p>*Note that r<sub>i</sub> are not independent</p> <p><b>Cook's distance plot</b></p> <ul style="list-style-type: none"><li>- To check for influential point</li></ul> <p><b>R<sup>2</sup></b>: The proportion of total variation of the response that is explained by the fitted regression.</p> <p><b>4. Boxplot</b></p> <p>-Minimum, lower quartile, median, upper quartile and maximum</p> <p>-For <b>comparing between quantitative vs categorical variable</b></p> <p>-Skewness of distributions (mention if obvious)</p>	<p><b>Chi-square test for 2 × 2 with continuity correction</b></p> <p>1. Assumptions:</p> <ul style="list-style-type: none"><li>- 2 categorical variables</li><li>- all expected cell count &gt;5</li><li>- Data random</li></ul> <p>2. Hypothesis:</p> <ul style="list-style-type: none"><li>- H<sub>0</sub>: 2 Variables are independent</li><li>- H<sub>1</sub>: 2 Variables are dependent</li></ul> <p>3. Test statistic:</p> $\chi^2 = \sum \frac{( observed count - expected )^2}{expected}$ <p>4. P-value:</p> <ul style="list-style-type: none"><li>- Right tail probability of <math>\chi^2</math> distribution with 1 degree of freedom</li></ul> <p><b>Pearson chi-square</b></p> <ul style="list-style-type: none"><li>- r rows and c columns that define 2 categorical random variables</li><li>- use <math>\chi^2_{(r-1)(c-1)}</math> distribution</li></ul> <p><b>1. Mean</b></p>	<p><b>Fisher's exact</b></p> <p>1. Assume:</p> <ul style="list-style-type: none"><li>- <math>\geq 20\%</math> expected cell count &lt;5</li><li>- 2 binary categorical variables</li><li>- Data random</li></ul> <p>2. Hypothesis:</p> <ul style="list-style-type: none"><li>- H<sub>0</sub>: 2 variables are independent</li></ul> <p>3. Test stats:</p> <ul style="list-style-type: none"><li>- The first cell count</li></ul> <p><b>Linear by linear</b></p> <ul style="list-style-type: none"><li>-There is at least 1 ordinal variable</li></ul> <p><b>6.pth-quantile</b></p>	<p><b>Strength of association</b></p> <p>-How different is P(disease exposure) is from P(disease no exposure)</p> <p>- <math>\hat{p}_1 = \frac{a}{a+b}</math> is point estimate of P(disease exposure)</p> <p>- <math>\hat{p}_2 = \frac{c}{c+d}</math> is a point estimate of P(disease no exposure)</p> <p><b>Prospective study:</b></p> <ul style="list-style-type: none"><li>-randomly assign the exposure variable to subjects or record their exposure variable status</li></ul> <p><b>Retrospective study:</b></p> <ul style="list-style-type: none"><li>-case control study</li><li>-cannot obtain <math>\hat{p}_1</math> and <math>\hat{p}_2</math> from 2 × 2 table</li></ul> <p><b>Probability</b></p>	<p><b>Difference of proportion:</b></p> <p>- <math>\hat{p}_1 - \hat{p}_2</math></p> <p><b>Relative risk:</b></p> $\frac{\hat{p}_1}{\hat{p}_2}$ <p><b>Odds of success:</b></p> $-odds = \frac{p}{p-1} \text{ where } p \text{ is the probability of success and } p-1 \text{ is failure}$ <p>-odds ratio= <math>\frac{a/b}{c/d} = \frac{ad}{bc}</math></p> <p>-The odds of having disease in expose is <math>\frac{a}{b}</math> times the odds of having disease in unexposed</p> <p><b>Computing CI for OR</b></p> <p>-estimated se for ln(OR) is <math>se(\ln(OR)) = \sqrt{\frac{1}{a} + \frac{1}{b} + \frac{1}{c} + \frac{1}{d}}</math></p> <p>-Formula for a (1- <math>\alpha</math>) 100% CI:</p> $e^{\ln(\frac{ad}{bc}) \pm q_{1-\alpha/2} \sqrt{\frac{1}{a} + \frac{1}{b} + \frac{1}{c} + \frac{1}{d}}}$ <p><b>Random variable:</b></p>	<p><b>Logistic regression</b></p> <p>-Response variable Y is a binary random variable (categorical)</p> <p>1. Assumptions:</p> <ul style="list-style-type: none"><li>-Data is randomized</li><li>-Observations are independent</li><li>-Denote p<sub>i</sub>=P(Y<sub>i</sub>=1) (success probability) given a value of X<sub>i</sub></li></ul> <p>2. Odds of success:</p> $\frac{P(Y_i=1)}{P(Y_i=0)} = \frac{p_i}{1-p_i}$ <p>- <math>\ln\left(\frac{p_i}{1-p_i}\right) = \beta_0 + \beta_1 X_i</math></p> <p>- <math>\ln\left(\frac{p_i}{1-p_i}\right) = \beta_0 + \beta_1 X_{i1} + \beta_2 X_{i2} + \dots</math></p> <p>3. Distributional assumption:</p> <ul style="list-style-type: none"><li>- <math>Y_i \sim \text{Bin}(1, p_i)</math></li><li>- <math>p_i = \frac{e^{\beta_0 + \beta_1 X_i}}{1 + e^{\beta_0 + \beta_1 X_i}}</math></li></ul> <p><b>Maximum likelihood estimation:</b></p> <ul style="list-style-type: none"><li>-Test fitness of model using R<sup>2</sup></li></ul> <p><b>Binomial Distribution</b></p>	<p><b>Model interpretation</b></p> <p>-For every increase in X<sub>i</sub>, the ln odds increase by <math>\beta_1</math> times</p> <p>-For every increase in X<sub>i</sub>, the odds of success will increase by <math>e^{\beta_1}</math> times</p> <p><b>Wald Hypothesis Test</b></p> <p>-Test if an individual term is significantly different from zero</p> <p>1. Hypothesis:</p> <ul style="list-style-type: none"><li>-H<sub>0</sub>: <math>\beta_1 = 0</math></li></ul> <p>2. Test statistic:</p> $-s.e.\left(\frac{\hat{\beta}_1}{\hat{\sigma}_1}\right)^2 \text{ with sampling distribution } \chi^2_1$ <p><b>Confidence Interval of Wald Test</b></p> <ul style="list-style-type: none"><li>- <math>\hat{\beta}_1 \pm q_{1-\frac{\alpha}{2}} \times s.e.\left(\hat{\beta}_1\right)</math></li></ul> <p>-odds ratio corresponding to <math>\beta</math> is <math>(e^L, e^U)</math></p> <p>-95% CI for odds ratio corresponding to 1 unit increase in X<sub>i</sub> is <math>(e^L, e^U)</math></p> <p><b>Confidence interval</b></p>	

<p>-Look at proportion of modal category, compare it via the difference</p> <p><b>2.Bar Plot:</b></p> <p>-To display <b>single categorical variable</b></p> <p>-Mention group of categories with high/low proportions</p> <p>-Mention any apparent trend in proportions if there is ordering</p> <p><b>3.Histogram</b></p> <p>-use bars to portray the frequency or relative frequency of the possible outcomes for <b>quantitative variable</b></p> <p>-Overall pattern: Gap, outlier?</p> <p>-Skewed or symmetric?</p> <p>-Unimodal? Bimodal?</p>	<p>-Spread of data</p> <p>-Median and IQR</p> <p>-No. of outliers</p> <p>-Outlier: 1.5 <math>\times</math> IQR below and above upper/lower quartile</p> <p><b>5.Contingency table</b></p> <p>-To display <b>2 categorical variables</b></p> <p><b>6.Scatterplot</b></p> <p>-For <b>association</b> between <b>quantitative vs quantitative</b> variable</p> <p>-correlation:</p> $r = \frac{1}{n-1} \sum_{i=1}^n \frac{(X_i - \bar{X})(Y_i - \bar{Y})}{s_x s_y}$ <p>-r is between 1 and -1</p>	<p>-Sum of observations divided by number of observations</p> $\bar{X} = \frac{1}{n} \sum_{i=1}^n X_i$ <p><b>2.Median</b></p> <p>-Middle value of the observations when observations ordered from smallest to largest</p> <p>- <math>\left(\frac{n+1}{2}\right)</math> th largest if odd, average of <math>\left(\frac{n}{2}\right)</math> th and <math>X_{0.5}</math> th if even <math>\rightarrow</math></p> <p><b>*If data is highly skewed, report median and if symmetric, report mean</b></p> <p><b>3.Range</b></p> <p>-Difference between largest and smallest observations</p> <p>-Measure of spread but sensitive to extreme observations</p> <p><b>4.Variance</b></p> <p>-Average of the squared deviation from the mean</p> $s^2 = \frac{1}{n-1} \sum_{i=1}^n (X_i - \bar{X})^2$ <p><b>5.Standard deviation</b></p> <p>-Represents average distance from mean</p> <p>-68% observations fall between 1 sd of mean, between <math>\bar{X} - s</math> and <math>\bar{X} + s</math></p> <p>-95% fall between 2 sd of mean</p> <p>-All fall between 3 sd of mean</p>	<p>-value that p percent of the observations fall below or at that value (<math>\hat{q}_{0.5} = X_{0.5}</math>)</p> <p><b>7.IQR</b></p> <p>-Distance between upper and lower quartile</p> $(\hat{q}_{0.75} - \hat{q}_{0.25})$ <p>-50% of sample points fall within IQR</p> <p>-how spread the 'middle' data is</p> <p><b>Lurking variable</b></p> <p>-Variable that is usually unobserved that influence the association between variables of primary interest</p> <p><b>Confounding</b></p> <p>-2 explanatory variables are associated with a response variable but also associated with each other</p> <p><b>Poor sources of bias</b></p> <p>-Sampling bias</p> <p>-Non-response bias</p> <p>-Response bias</p>	<p><b>1.Mutually exclusive (disjoint)</b></p> <p>-A and B has no common outcome</p> $A \cap B = \emptyset$ <p>-If <math>P(A) \geq 0</math>, <math>P(S)=1</math>,</p> $P(A \cup B) = P(A) + P(B)$ <p><b>2.Not mutually exclusive</b></p> $P(A \cup B) = P(A) + P(B) - P(A \cap B)$ $P(A) = P(A \cap B) + P(A \cap B^c)$ <p><b>3.Independent event</b></p> <p>-2 events do not influence one another</p> $P(A \cap B) = P(A) \cdot P(B)$ $P(A B) = P(A)$ <p><b>4.Conditional probability</b></p> <p>-Probability of A given B</p> $P(A B) = \frac{P(A \cap B)}{P(B)}$ <p><b>5.Bayes Theorem</b></p> <p><b>6.Sensitivity</b></p> <p>-Probability that the test is positive given that the person has disease <math>P(A B)</math></p> <p><b>7.Specificity</b></p> <p>-Probability that the test is negative given no disease <math>P(A^c B^c)</math></p> <p><b>8.Prevalance</b></p> <p>-Having Disease <math>P(B)</math></p>	<p>-Numerical measurement of outcome of experiment via random sampling</p> $p_z = \frac{e^{-20} 20^z}{z!}$ <p>for barplot</p> <p><b>1.Mean:</b></p> <p>-sum of probabilities multiplied by possibilities</p> $\mu = \sum_x x p_x$ <p><b>2.Expected value of X, E(X):</b> The mean of the probability distribution of a random variable X</p> $E(\bar{X}) = \frac{1}{n} \sum_{i=1}^n X_i = \mu$ <p><b>*mean of variables</b></p> <p>=mean of each random variable</p> <p><b>3.Standard deviation</b></p> <p>-measures variability of a random variable from the mean</p> $\sigma^2 = \sum_x (X - \mu)^2 p$ <p><b>4.Variance</b></p> <p>-large standard deviation has larger variability</p> $Var(\bar{X}) = \frac{1}{n^2} \sum_{i=1}^n \sigma^2$ <p><b>Probability density function pdf</b></p> <p>-curve that determines probabilities of intervals</p> <p><b>Mean of continuous variable X</b></p> $\mu = \int x f(x) dx$ <p><b>Variance of continuous variable</b></p> $\sigma^2 = \int (x - \mu)^2 f(x) dx$ <p><b>Pth quantile of cont variable</b></p> $P(X \leq q_p) = p$ <p><b>Combination:</b></p> $C_k^n = \frac{n!}{k!(n-k)!}$	<p>-n trials have 2 possible outcomes, independent</p> <p>-same probability of success</p> $X \sim \text{Bin}(n, p)$ $P(X=x) = C_x^n p^x (1-p)^{n-x}$ $E(X) = np, Var(X) = np(1-p)$ <p><b>Normal distribution</b></p> <p>-Symmetric, bell-shaped</p> $X \sim N(\mu, \sigma^2)$ $\bar{X} = \frac{X_1 + \dots + X_n}{n} \sim N\left(\mu, \frac{\sigma^2}{n}\right)$ $Z = \frac{\bar{X} - \mu}{\frac{\sigma}{\sqrt{n}}} \sim N(0,1)$ <p><b>Probability distribution</b></p> <p><math>P(X=1)=p, P(X=0)=1-p</math></p> <p><b>Sampling distribution:</b></p> <p>-Probability distribution that specifies the probabilities for possible values that statistic can take</p> <p><b>Sampling proportion <math>\hat{p}</math> for random sample size n:</b></p> <p>Mean=<math>\mu</math>, sd=<math>\sqrt{\frac{p(1-p)}{n}}</math></p> <p><math>p, \left(\sqrt{\frac{p(1-p)}{n}}\right)^2</math> if <math>\hat{p} \sim N(\mu, \sigma^2)</math></p> <p>n is sufficiently large that <math>n\hat{p}(1-\hat{p}) \geq 5</math>, the sampling distribution is <b>approximately normal</b></p> <p><b>CLT:</b> <math>\bar{X} \sim N\left(\mu, \frac{\sigma^2}{n}\right)</math></p> <p>-independent observations</p> <p>-n <math>\geq 30</math></p> <p>-X is approximately normal</p>	<p><b>*a long run interpretation</b></p> <p><b>Margin of error:</b></p> <p>-measures how accurate the point estimate is likely to be estimating a parameter <math>\rightarrow</math> sd</p> <p><b>Standard error:</b></p> <p>-an estimated deviation of a sampling distribution</p> <p><b>Population proportion:</b> <math>p</math></p> <p><b>Sample proportion:</b> <math>\hat{p}</math></p> <p><math>\rightarrow</math> for cat</p> <p>When <math>n\hat{p}(1-\hat{p}) \geq 5</math></p> $\hat{p} \pm q_{1-\frac{\alpha}{2}} \times \sqrt{\frac{\hat{p}(1-\hat{p})}{n}}$ <p>99%=2.58, 80%=1.28, 95%=1.96</p> <p><b>Width:</b> higher the CI, smaller <math>\alpha</math> and wider the width</p> $2 \times q_{1-\frac{\alpha}{2}} \times \sqrt{\frac{\hat{p}(1-\hat{p})}{n}}$ $n = \left( \frac{2 \times q_{1-\frac{\alpha}{2}}}{D} \right)^2 p(1-p)$ <p><b>*use p=0.5 (highest variability to obtain highest n)</b></p> <p><b>t-distribution:</b></p> <p>-probability on the degree of freedom(df)</p> <p>-thicker tail, more variability than <math>N(0,1)</math></p> <p>-larger the df, closer to <math>N(0,1)</math></p> $\bar{X} \pm t_{n-1, 1-\frac{\alpha}{2}} \times \frac{s}{\sqrt{n}}$ $2 \times t_{n-1, 1-\frac{\alpha}{2}} \times \frac{s}{\sqrt{n}} \leq D$ $n \geq \left( \frac{t_{n-1, 1-\frac{\alpha}{2}} \times s}{D} \right)^2$
---	--	--	---	--	---	--	--