

STA 104: Applied Statistical Methods: Nonparametric Statistics

Course Material Summary

University of California at Davis

Last Edit Date: 02/21/2022

Disclaimer and Term of Use:

1. We do not guarantee the accuracy and completeness of the summary content. Some of the course material may not be included, and some of the content in the summary may not be correct. You should use this file properly and legally. We are not responsible for any results from using this file.
2. Although most of the content in this summary is originally written by the creator, there may be still some of the content that is adapted (derived) from the slides and codes from *Professor Maxime Pouokam*. We use those as references and quotes in this file. Please [contact us](#) to delete this file if you think your rights have been violated.
3. This work is licensed under a [Creative Commons Attribution 4.0 International License](#).

STA 104 Summary

Topic	Content
Parametric vs non-parametric	<p>Non-parametric statistics uses techniques that do not require typical assumptions of traditional techniques.</p> <p>In a traditional test, we assume:</p> <ol style="list-style-type: none"> 1) Random sample was taken, or equivalently the X_i values are independent. 2) The sample size $n \geq 30$ (CLT), and the population is normal <p>When we do not have these assumptions above, the distributions based on CLT cannot be used, which means we need to assume a named distribution. This is what non-parametric do. It is often called “distribution free”.</p> <p>When assumptions are NOT violated, the parametric tests have more power. When assumptions are violated, the non-parametric tests have more power.</p> <p>* $Power = 1 - P(\text{Type II error}) = P(\text{Reject } H_0 \mid H_0 \text{ False})$</p>
Test for single median	
Binomial test	<p>We use median because we do not have these assumptions in parametric tests. Median is actually a proportion, which is 50%. θ_m denotes all possible medians. θ_m^* is the hypothesized median. n is the sample size.</p> <p>Assumption: A random sample was taken from the population (equivalently, observations are independent).</p> <p>Step 1: State the null and alternative hypothesis $H_0: \theta_m = \theta_m^* \text{ v.s. } H_1: \theta_m \neq \theta_m^* (H_0: p = 0.5 \ H_1: p \neq 0.5) \text{ (two - sided)}$ $H_0: \theta_m \leq \theta_m^* \text{ v.s. } H_1: \theta_m > \theta_m^* (H_0: p \leq 0.5 \ H_1: p > 0.5) \text{ (right tail, one - sided)}$ $H_0: \theta_m \geq \theta_m^* \text{ v.s. } H_1: \theta_m < \theta_m^* (H_0: p \geq 0.5 \ H_1: p < 0.5) \text{ (left tail, one - sided)}$</p> <p>Step 2: Calculate the test statistic $B^+ = \# \text{ of } X_i > \theta_m^*, B^+ \sim \text{Binomial}(n, \frac{1}{2})$</p> <p>Step 3: Calculate p-value $H_1: \theta_m \neq \theta_m^* (H_0: p = 0.5 \ H_1: p \neq 0.5) \Rightarrow p - \text{value} = 2(\min \{P(X \geq B^+), P(X \leq B^+)\})$ $H_1: \theta_m > \theta_m^* (H_0: p \leq 0.5 \ H_1: p > 0.5) \Rightarrow p - \text{value} = P(X \geq B^+)$ $H_1: \theta_m < \theta_m^* (H_0: p \geq 0.5 \ H_1: p < 0.5) \Rightarrow p - \text{value} = P(X \leq B^+)$ Interpretation: If the true median equals to θ_m^*, we can observe our data or more extreme with probability p-value.</p> <p>Step 4: Reject H_0 if p-value $< \alpha$</p>
Normal Approximation to Binomial test	<p>When we have a reasonable sample size, we may assume $B^+ = N(np, \sqrt{np(1-p)})$ based on CLT, where $p = 0.5$ under H_0. We now can use a z distribution.</p>

	<p>Assumption: A random sample was taken from the population (equivalently, observations are independent), and there are at least 5 observations above and below the hypothesized median.</p> <p>Step 1: State the null and alternative hypothesis $H_0: \theta_m = \theta_m^*$ v.s. $H_1: \theta_m \neq \theta_m^*$ ($H_0: p = 0.5$ $H_1: p \neq 0.5$) (two – sided) $H_0: \theta_m \leq \theta_m^*$ v.s. $H_1: \theta_m > \theta_m^*$ ($H_0: p \leq 0.5$ $H_1: p > 0.5$) (right tail, one – sided) $H_0: \theta_m \geq \theta_m^*$ v.s. $H_1: \theta_m < \theta_m^*$ ($H_0: p \geq 0.5$ $H_1: p < 0.5$) (left tail, one – sided)</p> <p>Step 2: Calculate the test statistic $Z_S = \frac{S - n(0.5)}{\sqrt{n(0.25)}}$, where $S = B^+ = \# \text{ of } X_i > \theta_m^*$, $B^+ \sim \text{Binomial}(n, \frac{1}{2})$</p> <p>Step 3: Calculate p-value $H_1: \theta_m \neq \theta_m^*$ ($H_0: p = 0.5$ $H_1: p \neq 0.5$) \Rightarrow $p\text{-value} = 2P(Z > Z_S)$ $H_1: \theta_m > \theta_m^*$ ($H_0: p \leq 0.5$ $H_1: p > 0.5$) \Rightarrow $p\text{-value} = P(Z > Z_S)$ $H_1: \theta_m < \theta_m^*$ ($H_0: p \geq 0.5$ $H_1: p < 0.5$) \Rightarrow $p\text{-value} = P(Z < Z_S)$ Interpretation: If the true median equals to θ_m^*, we can observe our data or more extreme with probability p-value.</p> <p>Step 4: Reject H_0 if p-value $< \alpha$</p>
Confidence Interval for median	<p>Find a $(1 - \alpha)100\%$ confidence interval for the median, using the normal approximation to binomial.</p> <p>Step 1: Get the location Lower bound location $= -z_{1-\frac{\alpha}{2}} * (\sqrt{0.25n}) + 0.5n$ Upper bound location $= +z_{1-\frac{\alpha}{2}} * (\sqrt{0.25n}) + 0.5n + 1$</p> <p>Step 2: Find the number in the rounded location to get the confidence interval $(X_{\text{lower bound location}}, X_{\text{upper bound location}})$</p>
Estimation for Percentile and CDF	<p>Find a $(1 - \alpha)100\%$ confidence interval for the CDF at x.</p> <p>Step 1: Get the proportion $\hat{F}(x) = \hat{p} = \frac{\# \text{ of } X_i \leq x}{n} \sim N(p, \sqrt{(p(1-p))/n})$</p> <p>Step 2: Get the lower and upper bound Lower bound: $\hat{p} - z_{1-\frac{\alpha}{2}} * \sqrt{(p(1-p))/n}$ Upper bound: $\hat{p} + z_{1-\frac{\alpha}{2}} * \sqrt{(p(1-p))/n}$</p>

	<p>Step 3: Convert the lower and upper bound to percentile (lower bound percentile, upper bound percentile) $\times 100\%$</p>
Confidence Intervals for percentiles	<p>Find a $(1 - \alpha)100\%$ confidence interval for the $(p^*)100^{th}$ percentile.</p> <p>Step 1: Get the location Lower bound location $= n(p^*) - z_{1-\frac{\alpha}{2}} * \sqrt{p^*(1-p^*)n}$ Upper bound location $= n(p^*) + 1 + z_{1-\frac{\alpha}{2}} * \sqrt{p^*(1-p^*)n}$</p> <p>Step 2: Find the number in the rounded location and get the confidence interval $(X_{lower\ bound\ location}, X_{upper\ bound\ location})$</p> <p># When we get a location equals to 0 or n+1, we should use 1 or n as our location.</p>
Tests for two groups	
Comparing two means	<p>The goal is to determine whether two means are statistically different. Assumptions for parametric test are:</p> <ol style="list-style-type: none"> 1) Random sample from both groups 2) Groups are independent 3) \bar{X}_1 and \bar{X}_2 are normal
Permutation test for two groups	<p>Let $F_1(x)$ = CDF for group 1, $F_2(x)$ = CDF for group 2. If the distributions for the groups are equal, $F_1(x) = F_2(x)$. Both groups are from the same population.</p> <p>Assumption: A random sample was taken from each group, groups independent.</p> <p>Step 1: state the null and alternative hypothesis $H_0: F_1(x) = F_2(x)$ v.s $H_1: F_1(x) \geq F_2(x)$ or $F_1(x) \leq F_2(x)$ (two – sided) # distributions are different $H_1: F_1(x) \leq F_2(x)$ (right tail, one – sided) # group 1 tends to be larger than group 2 $H_1: F_1(x) \geq F_2(x)$ (left tail, one – sided) # group 2 tends to be larger than group 1</p> <p>Step 2: Calculate the observed statistic and all permutations $D^{OBS} = \bar{X}_1 - \bar{X}_2$ or $D^{OBS} = total_1 - total_2$ or $D^{OBS} = median_1 - median_2$ Permutations $= \binom{m+n}{m} = \binom{m+n}{n} = \frac{(m+n)!}{m!n!}$</p> <p>Step 3: Calculate the permutation p-value $H_1: F_1(x) \geq F_2(x)$ or $F_1(x) \leq F_2(x) \Rightarrow \frac{\# of D_i \geq D^{OBS} }{permutations}$ $H_1: F_1(x) \leq F_2(x) \Rightarrow \frac{\# of D_i \geq D^{OBS}}{permutations}$</p>

	$H_1: F_1(x) \geq F_2(x) \Rightarrow \frac{\# \text{ of } D_i \leq D^{OBS}}{\text{permutations}}$ <p>Step 4: If p-value < α, reject H_0.</p> <p># When we have asymmetric distributions, we use the median to compare outliers. Otherwise, we use total or mean.</p> <p>## If the sample sizes of each group are the same, the test results from total and mean are the same. Otherwise, the results will be different.</p>
Approximate Permutation Test	<p>Steps for an approximate permutation test (for coding):</p> <ol style="list-style-type: none"> 1) Record D^{OBS} 2) Create one vector of all observations 3) Randomly shuffle the (m + n) observations, and assign first m to group 1, last n to group 2 4) Compute D_i = observed difference (in means/medians/totals) 5) Repeat step 3 and 4, $R > 2000$ times 6) Based on these R random values of D_i, we have an approximate p-values are: $H_1: F_1(x) \geq F_2(x) \text{ or } F_1(x) \leq F_2(x) \Rightarrow (\# \text{ of } D_i \geq D^{OBS}) / R$ $H_1: F_1(x) \leq F_2(x) \Rightarrow (\# \text{ of } D_i \geq D^{OBS}) / R$ $H_1: F_1(x) \geq F_2(x) \Rightarrow (\# \text{ of } D_i \leq D^{OBS}) / R$ 7) If p-value < α, reject H_0
Confidence Interval for p-value	<p>A $(1 - \alpha)100\%$ CI for a p-value p^* is:</p> $p^* \pm z_{1-\frac{\alpha}{2}} \sqrt{p^*(1-p^*)/R}$
Normal Approximation to permutation	<p>We use the overall mean \bar{x}^* and overall standard deviation S^* in our test statistics. For this test, we need $n + m \geq 30$.</p> <p>Assumption: A random sample was taken from each group, groups independent.</p> <p>Step 1: state the null and alternative hypothesis $H_0: F_1(x) = F_2(x)$ v.s $H_1: F_1(x) \geq F_2(x) \text{ or } F_1(x) \leq F_2(x)$ (two – sided) # distributions are different $H_1: F_1(x) \leq F_2(x)$ (right tail, one – sided) # group 1 tends to be larger than group 2 $H_1: F_1(x) \geq F_2(x)$ (left tail, one – sided) # group 2 tends to be larger than group 1</p> <p>Step 2: Calculate the test statistic $Z_S = \frac{\bar{x}_1 - \bar{x}^*}{S^*/\sqrt{m}} \text{ or } Z_S = \frac{\bar{x}_2 - \bar{x}^*}{S^*/\sqrt{n}}$ <p>where $\bar{x}^* = \frac{1}{m+n} \sum x_i$, $S^* = \sqrt{\frac{1}{m+n-1} \sum (x_i - \bar{x}^*)^2}$</p></p> <p>Step 3: Calculate the permutation p-value</p>

	$H_1: F_1(x) \geq F_2(x) \text{ or } F_1(x) \leq F_2(x) \Rightarrow p\text{-value} = 2P(Z > Z_S)$ $H_1: F_1(x) \leq F_2(x) \Rightarrow p\text{-value} = P(Z > Z_S)$ $H_1: F_1(x) \geq F_2(x) \Rightarrow p\text{-value} = P(Z < Z_S)$ <p>Step 4: If p-value < α, reject H_0.</p>
Wilcoxon Rank Sum (WRS) test	<p>Assumption: A random sample was taken from each group, groups independent.</p> <p>Step 1: State the null and alternative hypothesis $H_0: F_1(x) = F_2(x)$ v.s $H_1: F_1(x) \geq F_2(x) \text{ or } F_1(x) \leq F_2(x)$ (two – sided) # distributions are different $H_1: F_1(x) \leq F_2(x)$ (right tail, one – sided) # group 1 tends to be larger than group 2 $H_1: F_1(x) \geq F_2(x)$ (left tail, one – sided) # group 2 tends to be larger than group 1</p> <p>Step 2: Calculate test statistic 1) Combine the m + n values into one group 2) Calculate the rank for each data point: $R(x_i) = \# \text{ of data } \leq x_i, i = 1, \dots, m + n$ Note: If there are ties, average the ranks of the tied observations, and assign the tied values as their ranks 3) Calculate the total rank in group 1 (arbitrary choice of groups). This is our test statistic, $W_{OBS} = \sum_{group\ 1} R(x_i)$.</p> <p>Step 3: Calculate the exact p-value Permutations = $\binom{m+n}{n}$, W_i = sum of rank in group 1 $H_1: F_1(x) \geq F_2(x) \text{ or } F_1(x) \leq F_2(x) \Rightarrow 2 * \min\left(\frac{\# \text{ of } W_i \geq W_{OBS}}{\text{permutations}}, \frac{\# \text{ of } W_i \leq W_{OBS}}{\text{permutations}}\right)$ $H_1: F_1(x) \leq F_2(x) \Rightarrow \frac{\# \text{ of } W_i \geq W_{OBS}}{\text{permutations}}$ $H_1: F_1(x) \geq F_2(x) \Rightarrow \frac{\# \text{ of } W_i \leq W_{OBS}}{\text{permutations}}$</p> <p>Step 4: If p-value < α, reject H_0.</p> <p># WRS tends to have higher power when the distribution is skewed, outliers are present, since assigning ranks essentially removes all influence of both issues.</p> <p>## Permutation tests tend to have higher power when the distribution is thought to be symmetric, and when using the mean.</p>
Large Sample Approximation to WRS	Let $N = m + n$, and $R(x_1), \dots, R(x_N)$ be the corresponding combined ranks of the two groups. Let S_1 = sum of ranks in group 1.

	<p>Under the assumption that the distributions are equal, every $R(x_i)$ should have been equally likely to come from both groups.</p> <p>Assumption: A random sample was taken from each group, independent groups, combined sample size at least 30.</p> <p>Step 1: State the null and alternative hypothesis $H_0: F_1(x) = F_2(x)$ v.s $H_1: F_1(x) \geq F_2(x)$ or $F_1(x) \leq F_2(x)$ (two – sided) # distributions are different $H_1: F_1(x) \leq F_2(x)$ (right tail, one – sided) # group 1 tends to be larger than group 2 $H_1: F_1(x) \geq F_2(x)$ (left tail, one – sided) # group 2 tends to be larger than group 1</p> <p>Step 2: Our test-statistic is $Z_S = \frac{W_{OBS} - E[S_1]}{\sqrt{\sigma_S^2}}$ where $W_{OBS} = \sum_{group\ 1} R(x_i)$, $E[S_1] = m\mu_R$, $\sigma_S^2 = \frac{mn\sigma_R^2}{N-1}$ where $\mu_R = \frac{1}{N} \sum R(x_i)$, $\sigma_R^2 = \frac{1}{N} \sum (R(x_i) - \bar{x}_R)^2$ If $N \geq 30$, we have $S_1 \sim N(m\mu_R, \frac{mn\sigma_R^2}{N-1})$ # If there are no ties, then $E[S_1] = \frac{m(N+1)}{2}$, $\sigma_R^2 = \frac{mn(N+1)}{12}$</p> <p>Step 3: Get the p-value $H_1: F_1(x) \geq F_2(x)$ or $F_1(x) \leq F_2(x) \Rightarrow P(Z > Z_S)$ $H_1: F_1(x) \leq F_2(x) \Rightarrow P(Z > Z_S)$ $H_1: F_1(x) \geq F_2(x) \Rightarrow P(Z < Z_S)$</p> <p>Step 4: If p-value < α, reject H_0.</p>
Mann-Whitney Test (alternative to WRS)	<p>Let X_1, \dots, X_m be our sample from group 1. Let Y_1, \dots, Y_n be our sample from group 2.</p> <p>Assumption: A random sample was taken from each group, groups independent.</p> <p>Step 1: State the null and alternative hypothesis $H_0: F_1(x) = F_2(x)$ v.s $H_1: F_1(x) \geq F_2(x)$ or $F_1(x) \leq F_2(x)$ (two – sided) # distributions are different $H_1: F_1(x) \leq F_2(x)$ (right tail, one – sided) # group 1 tends to be larger than group 2 $H_1: F_1(x) \geq F_2(x)$ (left tail, one – sided) # group 2 tends to be larger than group 1</p> <p>Step 2: Calculate test statistic</p>

	<p> $U_{MW} = (\# \text{ of pairs } (X_i < Y_i)) + \frac{1}{2} (\# \text{ of pairs } (X_i = Y_i))$ # If group 1 is lower than group 2, U_{MW} will be closed to the maximum # of pairs. ## If group 1 is larger than group 2, U_{MW} will be closed to 1. ### Number of possible pairs = $m \cdot n$ #### The test statistic is in U or Mann-Whitney Distribution. </p> <p>Step 3: Calculate p-value</p> <p>Let $U_{1-\frac{\alpha}{2}} = \left(1 - \frac{\alpha}{2}\right) 100^{th}$ percentile of U (upper)</p> <p>Let $U_{\frac{\alpha}{2}} = \left(\frac{\alpha}{2}\right) 100^{th}$ percentile of U (lower)</p> <p>$H_1: F_1(x) \geq F_2(x) \text{ or } F_1(x) \leq F_2(x) \Rightarrow \text{If } U_{MW} > U_{1-\frac{\alpha}{2}} \text{ or } U_{MW} < U_{\frac{\alpha}{2}} \Rightarrow < \alpha$</p> <p>$H_1: F_1(x) \leq F_2(x) \Rightarrow \text{If } U_{MW} < U_{\frac{\alpha}{2}} \Rightarrow < \frac{\alpha}{2}$</p> <p>$H_1: F_1(x) \geq F_2(x) \Rightarrow \text{If } U_{MW} > U_{1-\frac{\alpha}{2}} \Rightarrow < \frac{\alpha}{2}$</p> <p># Need to look at the table of Mann-Whitney Distribution.</p> <p>Step 4: If p-value < α, reject H_0.</p>
Kolmogorov Smirnov (KS) Test	<p>Assumption: A random sample was taken from each group, groups independent. Distributions should be continuous.</p> <p>Step 1: State the null and alternative hypothesis</p> <p>$H_0: F_1(x) = F_2(x) \text{ v.s. } H_1: F_1(x) \geq F_2(x) \text{ or } F_1(x) \leq F_2(x) \text{ (two-sided)}$ # distributions are different</p> <p>Step 2: Calculate test statistic</p> <p>Let $\hat{F}_1(x)$ = empirical CDF of group 1</p> <p>Let $\hat{F}_2(x)$ = empirical CDF of group 2</p> <p>Then,</p> <ol style="list-style-type: none"> 1) Combine the data from both groups 2) Calculate $\hat{F}_1(x)$ for both groups observations Calculate $\hat{F}_2(x)$ for both groups observations 3) Calculate the difference between $\hat{F}_1(x) - \hat{F}_2(x)$ for all observations 4) Our test-statistic is then $K_S = \max \hat{F}_1(x) - \hat{F}_2(x)$ <p>Step 3: Calculate p-value</p> <p>The p-value is a permutation p-value:</p> <p>$(\# \text{ of } \hat{F}_1(x) - \hat{F}_2(x) \geq K_S) / \binom{m+n}{n}$</p> <p>or divided by R if it's a random permutation test</p>

	Step 4: If p-value < α, reject H_0.																
Confidence Interval for shift parameter	<p>Step 1: Find all n*m pairwise differences, $X_i - Y_i$</p> <p>Step 2: Order the pairwise differences, call them pwd(1), pwd(2), ..., pwd(n*m)</p> <p>Step 3: We want the locations, call them ka and kb $P(\text{pwd}(ka) \leq \Delta \leq \text{pwd}(kb)) = 1 - \alpha \Rightarrow ka \leq U \leq kb - 1$ # kb - 1 because of discrete data</p> <p>Step 4: Confidence interval is ($ka = U_{\frac{\alpha}{2}} + 1$, $kb = U_{1-\frac{\alpha}{2}}$)</p> <p># If CI of Δ has both bounds > 0, then group 1 has larger distribution/measurement than group 2. ## If CI of Δ has both bounds < 0, then group 1 has smaller distribution/measurement than group 2. ### If CI of Δ contains 0, then there is no significant difference between group 1 and group 2.</p>																
Choose an appropriate test	<table border="1"> <thead> <tr> <th>Distribution</th><th>Statistic</th><th>Winner</th></tr> </thead> <tbody> <tr> <td>Symmetric</td><td>Mean</td><td>Permutation</td></tr> <tr> <td>Symmetric</td><td>Median</td><td>Wilcoxon Rank Sum</td></tr> <tr> <td>Asymmetric</td><td>Mean</td><td>Wilcoxon Rank Sum</td></tr> <tr> <td>Asymmetric</td><td>Median</td><td>Permutation</td></tr> </tbody> </table>		Distribution	Statistic	Winner	Symmetric	Mean	Permutation	Symmetric	Median	Wilcoxon Rank Sum	Asymmetric	Mean	Wilcoxon Rank Sum	Asymmetric	Median	Permutation
Distribution	Statistic	Winner															
Symmetric	Mean	Permutation															
Symmetric	Median	Wilcoxon Rank Sum															
Asymmetric	Mean	Wilcoxon Rank Sum															
Asymmetric	Median	Permutation															
Tests for three or more groups																	
ANOVA (non-parametric, permutation based)	<p>Notation: Assume we have K groups</p> <p>Let $X_{ij} = j^{th}$ observation from i^{th} group</p> <p>Let n_i = sample size of i^{th} group</p> <p>Let \bar{X}_i = sample mean of i^{th} group</p> <p>Let S_i^2 = sample variance of i^{th} group</p> <p>Let N = overall sample size = $\sum_{i=1}^k n_i$</p> <p>Let \bar{X} = overall sample mean = $\frac{\sum_{i=1}^k n_i \bar{X}_i}{N}$</p> <p>The following measure the difference between groups:</p> <p>SST = Sum of squared treatment = $\sum_{i=1}^k n_i (\bar{X}_i - \bar{X})^2$</p> <p>$MST = \frac{SST}{k-1}$</p> <p>The following measure the variances within each group:</p> <p># The idea is the same as parametric ANOVA. We compare the difference in means to the overall mean to the spread of each group.</p>																

	<p> $SSE = \sum_{i=1}^k (n_i - 1) S_i^2 = \text{Sum of square errors}$ $MSE = \frac{SSE}{N-k}$ </p> <p>Test statistic:</p> <p> $F_s = \frac{MST}{MSE}$ # When F_s is large => variance between groups is larger than within groups => means are significantly different ## When F_s is small => variance between groups is smaller than within groups => means are not significantly different </p> <p>Assumptions (traditional):</p> <ol style="list-style-type: none"> 1) Random samples are taken from all k groups 2) All k groups are independent 3) $\sigma_1 = \sigma_2 = \dots = \sigma_k$ equal variance (Levene's Test) 4) $\epsilon_{ij} \sim N(0, \sigma_\epsilon^2)$ independent and identically distributed (QQ plot and Shapiro-Wilks Test) <p># When the assumptions do not hold, we do not know what the distribution of F_s. But, we can find the permutation distribution.</p> <p>Assumptions (non-parametric):</p> <p>A random sample was taken from each group, groups independent.</p> <p>Step 1: State the null and alternative hypothesis</p> <p> $H_0: F_1(x) = F_2(x) = \dots = F_k(x)$ v.s. $H_1: F_i(x) \leq F_j(x)$ or $F_i(x) \geq F_j(x)$ for some $i \neq j$ </p> <p>Step 2: Calculate the observed test statistic</p> <p> $F_{OBS} = \frac{MST}{MSE}$ </p> <p>Step 3: Find the permutation p-value:</p> <p> Possible permutations = $\frac{N!}{n_1!n_2!\dots n_k!}$ We can also use random permutations: <ol style="list-style-type: none"> 1) Randomly assign the N observations into the k groups, $R > 4000$ times 2) Calculate the R values of F_s, denote F_i 3) Our p-value is $(\# \text{ of } F_i \geq F_{OBS})/R$ </p> <p>Step 4: If p-value < α, reject H_0.</p>
Kroskall-Wallis (KW) Test (permutation based)	<p>Kroskall-Wallis test uses ranks rather than the actual X_{ij} values. Has confidence interval.</p> <p>Assumptions:</p>

	<p>A random sample was taken from each group, groups independent.</p> <p>Step 1: State the null and alternative hypothesis $H_0: F_1(x) = F_2(x) = \dots = F_k(x) \text{ v.s. } H_1: F_i(x) \leq F_j(x) \text{ or } F_i(x) \geq F_j(x) \text{ for some } i \neq j$</p> <p>Step 2: Calculate the test statistic $KW_{OBS} = \frac{1}{S_R^2} \sum_{i=1}^k n_i \left(\bar{R}_i - \frac{N+1}{2} \right)^2$ Where S_R^2 = variance of ranks, regardless of groups, \bar{R}_i = mean rank of each group # This form of KW test works whenever ties are present or not</p> <p>Step 3: Calculate the approximate permutation p-value p-value = (# of $KW_i \geq KW_{OBS}$)/R</p> <p>Step 4: If p-value < α, reject H_0.</p>
<p>Large Sample Approximation to Kroskall-Wallis Test</p>	<p>If the n_i's are large, but an assumption of ANOVA is violated, we may use a large sample approximation.</p> <p>Motivation: In traditional ANOVA, we know that SST/σ_e^2 is distributed X^2 with df = k - 1.</p> <p>Now, replace X_{ij} with R_{ij}, which is the corresponding ranks, we can see:</p> $SST_R = \sum_{i=1}^k n_i \left(\bar{R}_i - \frac{N+1}{2} \right)^2$ <p>But the normalizing constant for the X^2 distribution has changed (since we are using R_{ij}) $E[c(SST_R)] = k - 1$ (since we know $E[X_{k-1}^2] = k - 1$), which gives $c = 1/S_R^2$ This gives our test statistic as:</p> $KW = \frac{1}{S_R^2} \sum_{i=1}^k n_i \left(\bar{R}_i - \frac{N+1}{2} \right)^2 \sim X_{k-1}^2$ <p>Assumptions: A random sample was taken from each group, groups independent, combined sample size at least 30.</p> <p>Step 1: State the null and alternative hypothesis $H_0: F_1(x) = F_2(x) = \dots = F_k(x) \text{ v.s. } H_1: F_i(x) \leq F_j(x) \text{ or } F_i(x) \geq F_j(x) \text{ for some } i \neq j$</p> <p>Step 2: Calculate the test statistic $KW = \frac{1}{S_R^2} \sum_{i=1}^k n_i \left(\bar{R}_i - \frac{N+1}{2} \right)^2$ Where S_R^2 = variance of ranks, regardless of groups, \bar{R}_i = mean rank of each group</p>

	<p>Step 3: Calculate the p-value $p\text{-value} = P(X_{k-1}^2 > KW)$</p> <p>Step 4: If $p\text{-value} < \alpha$, reject H_0.</p>
Asymptotic Bonferroni and Tukey cutoffs (Corrections for multiple comparisons)	<p>Assumptions: A random sample was taken from each group, groups independent, combined sample size at least 30.</p> <p>Bonferroni (BON) cutoff: $BON = Z_{1-\frac{\alpha}{2g}} \sqrt{S_R^2 \left(\frac{1}{n_i} + \frac{1}{n_j} \right)}$</p> <p>Tukey (HSD) cutoff: $HSD = q_\alpha(k, df = N - k) \sqrt{S_R^2 \left(\frac{1}{n_i} + \frac{1}{n_j} \right)}$</p> <p>Parametric version of Tukey: We reject H_0 if $\bar{X}_i - \bar{X}_j \geq q_\alpha(k, df = N - k) \sqrt{MSE \left(\frac{1}{n_i} + \frac{1}{n_j} \right)}$</p> <p>If $\bar{R}_i - \bar{R}_j > BON \text{ or } HSD$, we may conclude that the specific group have significant different average ranks.</p>
Permutation cutoff for Bonferroni and Tukey	<p>Assumptions: A random sample was taken from each group, groups independent.</p> <p>There are $\binom{k}{2}$ possible permutations. Compare the p-values to $\frac{\alpha}{g}$.</p> <p>Step 1: Randomly shuffle each observation into a group, $R > 4000$.</p> <p>Step 2: Pick a comparison measure, T_{ij}. Common values are $\bar{X}_i - \bar{X}_j$, $\bar{R}_i - \bar{R}_j$, $median_i - median_j$, $\frac{\bar{X}_i - \bar{X}_j}{\sqrt{MSE \left(\frac{1}{n_i} + \frac{1}{n_j} \right)}}$</p> <p>Step 3: For each R permutation, calculate $Q_{ij} = \max T_{ij}$</p> <p>Step 4: Let $q^*(\alpha)$ be the $(1-\alpha)100^{th}$ percentile of Q_{ij}. Then groups i and j are significant different if $T_{ij}^{OBS} > q^*(\alpha)$. We can also find the p-value = $(\# \text{ of } Q_{ij} \geq T_{ij}^{OBS}) / R$. If $p\text{-value} < \alpha$, groups i and j are significant different.</p>

Kroskall-Wallis v.s. Permutation	<p>The KW test will have higher power than a permutation test when:</p> <ol style="list-style-type: none"> 1) Outliers are present 2) The distribution of one or more groups is skewed 3) The distribution of one or more groups has “heavy tails”
Test for linear relationship	
Parametric test for correlation	<p>Assumptions:</p> <ol style="list-style-type: none"> 1) Pairs are independent (random selection of pairs) 2) (x_i, y_i) are distributed bivariate normal, where $r = \frac{1}{n-1} \sum_{i=1}^k \left(\frac{x_i - \bar{x}}{s_x} \right) \left(\frac{y_i - \bar{y}}{s_y} \right)$ <p>Let ρ denote the population correlation between numeric variables X and Y. We measure n pairs of data, (x_i, y_i).</p> <p>Step 1: State the null and alternative hypothesis</p> <p>$H_0: \rho = 0 \quad v.s. \quad H_1: \rho \neq 0$ $H_0: \rho \geq 0 \quad v.s. \quad H_1: \rho < 0$ $H_0: \rho \leq 0 \quad v.s. \quad H_1: \rho > 0$</p> <p>Step 2: Calculate the test statistic</p> $t_s = r \sqrt{\frac{n-2}{1-r^2}}$ <p>Step 3: Calculate the p-value</p> <p>$H_1: \rho \neq 0 \Rightarrow p\text{-value} = 2P(t > t_s)$ $H_1: \rho < 0 \Rightarrow p\text{-value} = P(t < t_s)$ $H_1: \rho > 0 \Rightarrow p\text{-value} = P(t > t_s)$</p> <p>Step 4: If p-value < α, reject H_0.</p> <p>We can also create linear regression line and a test for the slope:</p> <p>True model: $Y_i = \beta_0 + \beta_1 x_i + \epsilon_i$ Least square line: $\hat{y} = \hat{\beta}_0 + \hat{\beta}_1 \bar{x}$, $\hat{\beta}_1 = r \left(\frac{s_y}{s_x} \right)$, $\beta_1 = \bar{y} - \hat{\beta}_1 \bar{x}$</p> <p>Assumptions:</p> <ol style="list-style-type: none"> 1) Pairs are randomly sampled/independent 2) $\epsilon_i \sim N(0, \sigma_\epsilon^2)$ independent and identically distributed

	<p>Step 1: State the null and alternative hypothesis</p> $H_0: \beta_1 = 0 \text{ v.s. } H_1: \beta_1 \neq 0$ $H_0: \beta_1 \geq 0 \text{ v.s. } H_1: \beta_1 < 0$ $H_0: \beta_1 \leq 0 \text{ v.s. } H_1: \beta_1 > 0$ <p>Step 2: Calculate the test statistic</p> $t_s = \hat{\beta}_1 \sqrt{\frac{\sum(x_i - \bar{x})^2}{MSE}} \sim t(df = n - 2), \text{ where } MSE = \frac{\sum(y_i - \hat{y}_i)^2}{n - 2}$ <p>Step 3: Calculate the p-value</p> $H_1: \beta_1 \neq 0 \Rightarrow p\text{-value} = 2P(t > t_s)$ $H_1: \beta_1 < 0 \Rightarrow p\text{-value} = P(t < t_s)$ $H_1: \beta_1 > 0 \Rightarrow p\text{-value} = P(t > t_s)$ <p>Step 4: If p-value < α, reject H_0.</p> <p># If β_1 or $\rho = 0 \Rightarrow$ no linear relationship between X and Y ## If β_1 or $\rho < 0 \Rightarrow$ negative linear relationship between X and Y ### If β_1 or $\rho > 0 \Rightarrow$ positive linear relationship between X and Y</p>
Permutation test for the slope	<p>Common reasons we use a non-parametric test:</p> <ol style="list-style-type: none"> 1) Outliers present (violates normality) 2) Non constant variance (violates normality) 3) Small sample size (may not be able to conclude normal) <p>Assumptions: A random sample of pairs of data was taken.</p> <p>Step 1: State the null and alternative hypothesis</p> $H_0: \beta_1 = 0 \text{ v.s. } H_1: \beta_1 \neq 0$ $H_0: \beta_1 \geq 0 \text{ v.s. } H_1: \beta_1 < 0$ $H_0: \beta_1 \leq 0 \text{ v.s. } H_1: \beta_1 > 0$ <p>Step 2: Calculate the observed test hypothesis</p> $\hat{\beta}_1^{OBS} = \text{estimated least-squares slope} = r \frac{s_y}{s_x}$ <p>Step 3: Calculate the permutation p-value</p>

	<p>There are n ways to pair the first y_i with an x_i, then $n-1$ ways to pair the second y_i with an x_i, etc. There are $n!$ possible permutations.</p> <p>1) Permute the data and calculate $\hat{\beta}_1^i$</p> <p>2) Repeat for either</p> <ul style="list-style-type: none"> i) All $n!$ permutations ii) $R > 3000$ random permutations <p>3) The actual or estimated permutation p-values are</p> $H_1: \beta_1 \neq 0 \Rightarrow \frac{\# \text{ of } \hat{\beta}_1^i \geq \hat{\beta}_1^{OBS} }{n!} \text{ (actual) or } \frac{\# \text{ of } \hat{\beta}_1^i \geq \hat{\beta}_1^{OBS} }{R} \text{ (estimated)}$ $H_1: \beta_1 < 0 \Rightarrow \frac{\# \text{ of } \hat{\beta}_1^i \leq \hat{\beta}_1^{OBS}}{n!} \text{ (actual) or } \frac{\# \text{ of } \hat{\beta}_1^i \leq \hat{\beta}_1^{OBS}}{R} \text{ (estimated)}$ $H_1: \beta_1 > 0 \Rightarrow \frac{\# \text{ of } \hat{\beta}_1^i \geq \hat{\beta}_1^{OBS}}{n!} \text{ (actual) or } \frac{\# \text{ of } \hat{\beta}_1^i \geq \hat{\beta}_1^{OBS}}{R} \text{ (estimated)}$ <p>Step 4: If p-value $< \alpha$, reject H_0.</p>
Large Sample approximation to permutation test for the slope	<p>Assumptions: A random sample of pairs of data was taken, combined sample size at least 30.</p> <p>Step 1: State the null and alternative hypothesis</p> $H_0: \rho = 0 \text{ v.s. } H_1: \rho \neq 0$ $H_0: \rho \geq 0 \text{ v.s. } H_1: \rho < 0$ $H_0: \rho \leq 0 \text{ v.s. } H_1: \rho > 0$ <p>Step 2: Calculate the test statistic</p> $Z_s = \frac{r-0}{1/\sqrt{n-1}} = r\sqrt{n-1} \sim N(0,1/\sqrt{n-1})$ <p>Step 3: Calculate the p-value</p> $H_1: \rho \neq 0 \Rightarrow 2P(Z > Z_s)$ $H_1: \rho < 0 \Rightarrow P(Z < Z_s)$ $H_1: \rho > 0 \Rightarrow P(Z > Z_s)$ <p>Step 4: If p-value $< \alpha$, reject H_0.</p>
Spearman's Rank Correlation	<p>Let $R(X_i)$ = rank for x_i, $i = 1, \dots, n$; $R(Y_i)$ = rank for y_i, $i = 1, \dots, n$. $\bar{R}(x)$ = average rank of x_i, $S_{R(Y)}$ = standard deviation of rank of Y_i</p> <p>Step 1: State the null and alternative hypothesis</p> $H_0: \rho_s = 0 \text{ v.s. } H_1: \rho_s \neq 0$ $H_0: \rho_s \leq 0 \text{ v.s. } H_1: \rho_s > 0$

	<p>$H_0: \rho_s \geq 0$ v.s. $H_1: \rho_s < 0$</p> <p>Step 2: Calculate the test statistic</p> $r_s = \frac{1}{n-1} \sum_{i=1}^k \left(\frac{R(x_i) - \bar{R}(x)}{S_{R(x)}} \right) \left(\frac{R(y_i) - \bar{R}(y)}{S_{R(y)}} \right)$ <p>Step 3: Calculate the p-value</p> <p>$H_1: \rho_s \neq 0 \Rightarrow 2P(r_s^* > r_s)$</p> <p>$H_1: \rho_s > 0 \Rightarrow P(r_s^* > r_s)$</p> <p>$H_1: \rho_s < 0 \Rightarrow P(r_s^* < r_s)$</p> <p># If $H_1: \rho_s < 0$, $P(r_s^* < -c) = P(r_s^* > c)$</p> <p>Step 4: If p-value < α, reject H_0.</p>
Kendall's Tau	<p>Kendall's Tau does not use ranks directly, but also does not use the original data.</p> <p>Suppose we look at a pair of (x_i, y_i), say (x_1, y_1) and (x_2, y_2).</p> <p>1) If as X increases, Y tends to increase, then we should see $x_1 > x_2 \Rightarrow y_1 > y_2$.</p> <p>2) If as X increases, Y tends to decrease, then we should see $x_1 > x_2 \Rightarrow y_1 < y_2$.</p> <p>We use this to describe “discordant” and “concordant” pairs.</p> <p>Concordant pairs: If $X_i < X_j \Rightarrow Y_i < Y_j$, or equivalently $(X_i - X_j)(Y_i - Y_j) > 0$ (or $X_i > X_j \Rightarrow Y_i > Y_j$)</p> <p>Discordant pairs: If $X_i < X_j \Rightarrow Y_i > Y_j$, or equivalently $(X_i - X_j)(Y_i - Y_j) < 0$ (or $X_i < X_j \Rightarrow Y_i > Y_j$)</p> <p>If most pairs are concordant \Rightarrow positive linear relationship</p> <p>If most pairs are discordant \Rightarrow negative linear relationship</p> <p>The “population” value of Kendall's Tau is</p> <p>$\tau = 2P[(X_i - X_j)(Y_i - Y_j) > 0] - 1$, which is a rescaled probability of concordant pairs.</p> <p>If all pairs are concordant, $\tau = 1$. If all pairs are discordant, $\tau = -1$. If exactly half are concordant, half are discordant, $\tau = 0$.</p> <p>There are $\binom{n}{2}$ total pairs $(X_i, X_j), (Y_i, Y_j)$ then</p> <p>$U_{ij} = 1$ if $(X_i - X_j)(Y_i - Y_j) > 0$ (concordant)</p> <p>$U_{ij} = \frac{1}{2}$ if $(X_i - X_j)(Y_i - Y_j) = 0$ (tied)</p> <p>$U_{ij} = 0$ if $(X_i - X_j)(Y_i - Y_j) < 0$ (discordant)</p>

	<p>Let $V_i = \sum_{j=i+1}^n U_{ij} = \#$ of concordant pairs for i^{th} value (x_i, y_i). $\# j = i + 1$ ensures that we are never comparing the same pair.</p> $r_\tau = \frac{2[\sum_{i=1}^{n-1} V_i]}{\binom{n}{2}} - 1$
Exact Hypothesis Test for τ	<p>Step 1: State the null and alternative hypothesis $H_0: \tau = 0 \quad v.s. \quad H_1: \tau \neq 0$ $H_0: \tau \leq 0 \quad v.s. \quad H_1: \tau > 0$ $H_0: \tau \geq 0 \quad v.s. \quad H_1: \tau < 0$</p> <p>Step 2: Calculate test statistic $r_\tau = \frac{2[\sum_{i=1}^{n-1} V_i]}{\binom{n}{2}} - 1$</p> <p>Step 3: Calculate the p-value $H_1: \tau \neq 0 \Rightarrow 2P(r_\tau^* > r_\tau)$ $H_1: \tau > 0 \Rightarrow P(r_\tau^* > r_\tau)$ $H_1: \tau < 0 \Rightarrow P(r_\tau^* < r_\tau)$</p> <p>Step 4: If p-value $< \alpha$, reject H_0.</p>
Permutation test for τ	<p>Same step 1 and 2 as exact hypothesis test for τ</p> <p>Step 3: Calculate the p-value $H_1: \tau \neq 0 \Rightarrow (\# r_\tau^* \geq r_{\tau OBS})/R$ $H_1: \tau > 0 \Rightarrow (\# r_\tau^* \geq r_{\tau OBS})/R$ $H_1: \tau < 0 \Rightarrow (\# r_\tau^* \leq r_{\tau OBS})/R$</p> <p>Step 4: If p-value $< \alpha$, reject H_0.</p>
Asymptotic Approximation for τ	<p>The following formula can be used with or without ties. Let $s_i = \#$ of ties for the i^{th} tied value of X Let $t_i = \#$ of ties for the i^{th} tied value of Y</p>
When to use which correlation	<p>1) When there are no outliers, and the distribution is approximately symmetric (but with low sample size) use a permutation test for the slope.</p> <p>2) When outliers are present in the data, use Spearman's or Kendall's.</p>

	3) Kendall and Spearman tend to have similar results, but Spearman tends to have higher power at low sample sizes, and Kendall has higher power in large sample sizes.
--	--