

**Binomial Test**: tests of hypotheses for medians.  $Z_B = \frac{B - 0.5n}{\sqrt{0.25n}}$ ,  $B = \# \text{ of obs greater than or equal to } \theta_{0.5}$

**Confidence Interval**:  $P(X_A < \theta_{0.5} < X_B) = 1 - \chi = \sum_{x=a}^{b-1} \binom{n}{x} (0.5)^n$  = the probability that at least "a" and at most "b-1" of the obs fall less than  $\theta_{0.5}$

**Large Sample Approximation**:  $\frac{A-np}{np(1-p)} = -Z_{(1-\alpha)/2}$ ,  $\frac{b-1-np}{np(1-p)} = Z_{(1-\beta)/2}$  | Power: In general, the binomial test will have higher power than the CLT test for heavy-tailed distributions.

**Power of CLT Test**:  $1 - \Phi(Z_{(1-\alpha)/2} - \frac{x-A}{\sigma/\sqrt{n}})$  | **Power of Binomial Test**:  $1 - \Phi(Z_{(1-\beta)/2} \sqrt{\frac{0.25}{p(1-p)}} - \frac{p-0.5}{\sqrt{p(1-p)/n}})$

**CH 1**  $\downarrow$  **CH 2**

**Permutation Test**: the test of the distribution of  $\binom{N}{m}$  differences of means.

**Process of Permutation Test**:

- Calculate  $\binom{N}{m}$ ,  $m = \# \text{ of observations designated to new method.}$
- Compute the difference, or any desirable computations, of all combinations.
- Compute  $\# \text{ of cases where the differences of a combination is greater than or equal to that of the targeted combination.}$

**Large Sample Approximation**:  $D_{\text{obs}} \pm Z_{(1-\alpha)/2} \sqrt{\frac{p(1-p)}{N}}$

**Wilcoxon Rank-Sum Test**: Let  $W_i$  be the sum of the ranks of the observations from treatment  $i$ . The test is a 2-sample permutation test based on  $W_i$ .

**Process of Wilcoxon Test**:

- Combine  $m+n$  obs and assign ranks to each.
- Compute  $W_1$  and  $W_2$  according to the preassigned ranks.
- Find all possible permutations of ranks.  $\binom{N}{m}$  combinations
- Compute  $W_i$ 's for each combination.
- $P_{\text{upper-tail}} = \frac{\# \text{ of rank sums} \geq \text{observed rank sum } W}{\binom{m+n}{m}}$

**Large Sample Approximation**:  $E(T_i) = m\bar{r}$ ,  $\text{Var}(T_i) = \frac{mn\bar{r}^2}{N-1} \Rightarrow \mu = \frac{\sum A_i}{N}$ ,  $\tau^2 = \frac{\sum (A_i - \mu)^2}{N} \Rightarrow Z = \frac{T_i - E(T_i)}{\sqrt{\text{Var}(T_i)}}$

**Mann-Whitney Statistics "U"**:  $U = \# \text{ of pairs } (X_i, X_j) \text{ for which } X_i < Y_j$ . \*  $i$  and  $j$  are from 2 different samples,  $H_0: F(x) = F(y)$

**Process of Mann-Whitney test**:

- make pairs of  $(X_i, Y_j)$  and count the # of pairs  $X_i < Y_j$
- Using the Mann-Whitney CI table, confirm whether  $U$  is within the CI.

**Hodges-Lehmann Estimate of  $\Delta$** : the median of all the pairwise differences of  $X_i - Y_j$  | Using "pooled" t-test  $\Rightarrow \bar{X}_i - \bar{Y}_j \pm t_{(1-\alpha)} S_e \sqrt{\frac{1}{m+n}}$ , with  $\text{df} = m+n-2$ ,  $S_e = \sqrt{\frac{(m-1)S_x^2 + (n-1)S_y^2}{m+n-2}}$

**Process**: form all  $m n$  pairwise differences of  $X_i - Y_j$

- arrange all pair-wise differences in an ascending order, and find  $k_a$  and  $k_b$  that satisfy,  $p_{\text{mid}}(k_a) < \delta \leq p_{\text{mid}}(k_b) \approx 1 - \alpha \%$ , referring to Mann-Whitney table

**Siegel-Tukey Test**: Test of equality of variations of two populations. EX)  $H_0: T_1 = T_2$

**Process of Siegel-Tukey Test**:

- arrange the observations in an ascending order.
- assign rank 1 to the smallest obs, rank 2 to the largest obs, rank 3 to the next largest obs, rank 4 to the next smallest, ... so on.
- apply Wilcoxon rank-sum test. The smaller rank sums are associated with the treatment that has the larger variability.

**Test on Deviances**: obtain  $\text{dev}_{ix}$  and  $\text{dev}_{jx}$  and compute RMD from the original data,  $\text{dev}_{ix} = X_i - M_1$ ,  $\text{dev}_{jx} = Y_j - M_2 \Rightarrow \text{RMD} = \frac{1}{2} \sum_{i=1}^m |\text{dev}_{ix}| / m$

alternatively,  $\max(\frac{1}{m} \sum_{i=1}^m |\text{dev}_{ix}| / m, \frac{1}{n} \sum_{j=1}^n |\text{dev}_{jx}| / n)$

- Permute the deviances among 2 treatments and obtain RMD for each permutation

**RMD**:  $\min(\frac{1}{m} \sum_{i=1}^m |\text{dev}_{ix}| / m, \frac{1}{n} \sum_{j=1}^n |\text{dev}_{jx}| / n)$

- Compute p-values,  $\text{RMD}_i \geq \text{RMD}_{\text{obs}}$  for upper-tail,  $\text{RMD}_i \leq \text{RMD}_{\text{obs}}$  for lower-tailed

**Kolmogorov-Smirnov Test**: can be applied even when you don't have any information about the population distributions.  $T_{KS} = \frac{m}{n} \max_i |\hat{F}_1(z) - \hat{F}_2(z)|$

**Process of K-S Test**:

- align the data in an ascending order and specify where each obs belongs to, and count one by one and check where the obs belong.
- calculate the absolute difference for each step and find the maximum absolute difference.

**Sampling Formulas**: let  $A_i$  be the ranks and  $M = \frac{\sum A_i}{N}$ ,  $E(T_i) = m\bar{r}$ ,  $\text{Var}(T_i) = \frac{mn\bar{r}^2}{N-1}$ ,  $\Gamma^2 = \frac{\sum A_i^2}{N} - M^2$ ,  $Z = \frac{T_i - E(T_i)}{\sqrt{\text{Var}(T_i)}}$

**Sampling Formulas for Wilcoxon Rank-Sum Test**:  $\sum_i i = \frac{N(N+1)}{2}$ ,  $\sum i^2 = \frac{N(N+1)(N+2)}{6}$ ,  $\Gamma^2 = \frac{(N-1)(N+1)}{12} \Rightarrow E(W) = \frac{m(N+1)}{2}$ ,  $\text{Var}(W) = \frac{mn(N+1)}{12}$

**CH-3**

**K-Sample Permutation Tests**: comparing the distributions of  $K$  different samples. EX)  $H_0: F_1(x) = F_2(x) = \dots = F_K(x)$  ~ Using F-statistic

**Sum of Squares for Treatments**:  $SST = \frac{1}{K} \sum_i n_i (\bar{X}_i - \bar{X})^2$ , where  $\bar{X} = (\frac{1}{K} \sum_{i=1}^K \sum_{j=1}^{n_i} X_{ij}) / N \Rightarrow \text{Mean Squares for Treatments } MST = \frac{SST}{K-1}$  & **Sum of Square for Error**  $SSE = \frac{1}{N} \sum_{i=1}^K (n_i - 1) S_i^2$

**Mean Squares for Error**:  $SSE = \frac{N-K}{N-k} \Rightarrow F = \frac{MST}{MSE}$ , with  $\text{df} = K-1$  for numerator and  $\text{df} = N-K$  for denominator

**Process of Permutation F-Test**:

- Obtain the F-statistic for the original data
- Obtain all possible permutations, and compute the F-statistic for each permutation.

**Alternative Forms of the Permutation F statistic**:  $SS_{\text{total}} = \sum_{i=1}^K \sum_{j=1}^{n_i} (X_{ij} - \bar{X})^2 = SST + SSE = C$ ,  $F = \frac{SST/K-1}{(C-SST)/(N-K)}$ , where  $SST = \frac{1}{K} \sum_i n_i \bar{X}_i^2 - N \bar{X}^2$ ,  $SSX = \frac{1}{K} \sum_i n_i \bar{X}_i^2$

**Kruskal-Wallis Statistic**:  $KW = \frac{12}{N(N+1)} \sum_{i=1}^K n_i (\bar{R}_i - \frac{N+1}{2})^2$ , most of KW critical values will be smaller than  $\chi^2$ 's critical values.

**Adjustment for Ties**:  $KW_{\text{ties}} = KW / \left\{ 1 - \frac{1}{N^3-N} \sum_{i=1}^K (t_i^3 - t_i) \right\}$

**Multiple Comparisons**: multiple comparisons using pairwise tests allow us to verify which treatment differs from the others, if any.

**EX**: Bonferroni Adjustment, Fisher's Protected Least Significant Difference, Tukey's Honest Significant Difference

**Family Wise Error Rate**:

- the probability of making one or more false rejections ( $FWER > \alpha$ )
- $FWER = 1 - P(\text{no false rejection for } K \text{ treatments})$
- $= 1 - P(\text{no rejection of } H_{0K})^{K(K-1)/2}$
- $= 1 - (1 - \alpha)^{K(K-1)/2}$

**Bonferroni Adjustment**: uses  $\alpha'$  rather than  $\alpha$ ,  $\alpha' = \frac{2\alpha}{K(K-1)}$

- If the data are normally distributed, we may use the t-test; otherwise, use the Wilcoxon rank-sum test or any other nonparametric approach.

**Fisher's Protected LSD**: run F-test for equality of means.

- if the F-statistic is significant at level  $\alpha$ , run all pairwise t-tests at level  $\alpha$ .

**Rank-based analogue of LSD for testing the equality of distributions**  $\Rightarrow |\bar{R}_i - \bar{R}_j| \geq Z_{(\alpha/2)} \sqrt{\frac{N(N+1)}{12} \left( \frac{1}{n_i} + \frac{1}{n_j} \right)}$

**Tukey's HSD Procedure**: measures the largest difference between sample means.  $Q = \sqrt{n} (\max_i (\bar{X}_i) - \min_i (\bar{X}_i)) / \sqrt{MSE} \Rightarrow |\bar{X}_i - \bar{X}_j| \geq Q(\alpha, K, K(n-1)) \sqrt{\frac{MSE}{n}} \Leftrightarrow |\bar{R}_i - \bar{R}_j| \geq Q(\alpha, K, 12n) \sqrt{\frac{MSE}{12n}}$

**Tukey-Kramer Procedure**: to declare treatments  $i$  and  $j$  to be statistically different.  $|\bar{X}_i - \bar{X}_j| \geq Q(\alpha, K, N-1) \sqrt{\frac{MSE}{2}} \left( \frac{1}{n_i} + \frac{1}{n_j} \right) \Leftrightarrow |\bar{R}_i - \bar{R}_j| \geq Q(\alpha, K, \infty) \sqrt{\frac{N(N+1)}{24} \left( \frac{1}{n_i} + \frac{1}{n_j} \right)}$

**Ordered Alternatives**: if treatments are not equal, it may be possible to anticipate the direction  $H_a: F_1(x) \geq F_2(x) \geq F_3(x) \geq \dots \geq F_K(x)$

**Jonckheere-Terpstra Test**: a test statistic of the form  $T$  in which  $T_{ij}$ 's are one-sided Mann-Whitney statistics,  $T_{ij} = \sum_{i>j} T_{ij}$

- To obtain the upper-tail p-value, compute JT<sub>obs</sub>, and obtain all possible samples without replacement of sizes  $n_i$  and  $n_j$ , or randomly select if  $N$  is large. fraction of the JT's greater than or equal to JT<sub>obs</sub>

$$\Rightarrow \text{Large Sample Approximation: } E(JT) = \frac{N^2 - \sum n_i^2}{4}, \quad \text{Var}(JT) = \{N^2(2N+3) - \sum n_i^2(2n_i+3)\}/72 \Rightarrow Z = \frac{JT - E(JT)}{\sqrt{\text{Var}(JT)}} \quad | \quad E(JT_{\text{Wilson}}) = \frac{\sum n_i(n_i+n_j+1)}{2}$$

CH 4 Paired-Comparison Permutation Test: Compute  $D_i$ 's and  $\bar{D}_i$ 's. Then compute  $D_i$  for all  $2^n$  permutations.  $\Rightarrow$  Upper-tail p-value = # of  $D_i$ 's  $\geq \bar{D}_{\text{obs}}$ ,  $H_0: F(x) = (-F(x))$

L Large-Sample Approximations:  $E(\bar{D}) = 0$ ,  $\text{Var}(\bar{D}) = \frac{1}{N} \sum_{i=1}^N |\bar{D}_i|^2$ ,  $Z = \bar{D}/\sqrt{\text{Var}(\bar{D})}$  |  $E(S+) = \frac{1}{2} \sum_{i=1}^N |\bar{D}_i|$ ,  $\text{Var}(S+) = \frac{1}{4} \sum_{i=1}^N |\bar{D}_i|^2$ ,  $Z = (S+ - E(S+))/\sqrt{\text{Var}(S+)}$  | Signed-Rank Test: Nonparametric test for paired comparisons based on ranks. We rank the absolute values of the differences, and then attach the signs to the ranks.  $\bar{S} = \frac{1}{n} \sum_{i=1}^n V_i D_i$ ,  $V_i = -1$  or  $1$

Nilcoxon Signed-Rank Test: Obtain  $SR_{\text{obs}}$  and all  $2^n$  possible  $SR_i$ .  $\Rightarrow P_{\text{upper}} = (\# \text{ of } SR_i \geq SR_{\text{obs}})/2^n$  |  $P_{\text{lower}} = (\# \text{ of } SR_i \leq SR_{\text{obs}})/2^n$  \* always one-tailed |  $SR_i = \sum_{j=1}^n V_{ij} D_{ij}$ ,  $V_{ij} = 0$  or  $1$

L Large-Sample Approximation:  $SR_i = \sum_{j=1}^n V_{ij} D_{ij} \rightarrow E(SR_i) = \frac{n(n+1)}{4}$  |  $\text{Var}(SR_i) = \frac{n(n+1)(2n+1)}{24} \rightarrow Z = (SR_i - E(SR_i))/\sqrt{\text{Var}(SR_i)}$  | Selecting Among Paired-Comparison Tests: Shifting by  $\Delta \Rightarrow$  signed-rank test

L Ranking with Zeros: the absolute values of the differences, including zeros, are ranked. - If the distributions have lighter tails  $\Rightarrow$  t-test

Features of Randomized Complete Block Design: - Experimental units are divided into blocks in a way that units or conditions within blocks are homogeneous.

- Blocks have the same # of units as there are treatments, and the treatments are randomly assigned to experimental units within blocks.

Permutation F-test for RCB: - Compute  $F_{\text{obs}}$  and F-statistics for all  $(K!)^b$  permutations.  $\Rightarrow P_{\text{upper}} = (\# \text{ of } F \geq F_{\text{obs}})/(K!)^b$  |  $SST^* = \sum (\bar{x}_i - \bar{x})^2$  or  $SSX^* = \sum (\bar{x}_i - \bar{x})^2$

Friedman's Test for RCB: - Assign ranks within blocks and obtain mean rank for each treatment.  $\Rightarrow FM = \frac{12b}{K(K+1)} \sum_{i=1}^K (\bar{R}_i - \frac{K+1}{2})^2$ , and refer to  $\chi^2$  table with  $K-1$  degrees of freedom.

Cochran's Q: used for experiments with binary outcomes | both use FM test | Ordered Alternatives for RCB: Page's Test: - measure of the association b/w the presumed order and the rank of the treatments.

↓ Kendall's W: level of agreement (Concordance) |  $E(SSR) = \sum E(\bar{R}_i - \frac{K+1}{2})^2 = \frac{b}{b-1} \sum \bar{V}_{bi}$  | Assume the obs are ranked,  $PG = \sum_i R_i$ ,  $\Rightarrow$  has the same process with the permutation test

Large-Sample Approximation:  $E(PG) = \frac{b(b-1)^2}{4}$ ,  $\text{Var}(PG) = \frac{b(b-1)K^2(K+1)^2}{144} \Rightarrow Z = \frac{PG - E(PG)}{\sqrt{\text{Var}(PG)}}$

CH 5 Pearson Correlation:  $r = \frac{1}{n} \sum (x_i - \bar{x})(y_i - \bar{y}) / \sqrt{\sum (x_i - \bar{x})^2 \sum (y_i - \bar{y})^2} \Rightarrow t_{\text{cor}} = \frac{n-2}{n-1} r \leftarrow t\text{-distribution with } n-2 \text{ df. } H_0: P = 0$  | For bivariate sampling,  $\beta_1 = P \frac{\sum V_i}{\sum X_i} \sim \hat{\beta}_1 = r \frac{s_y}{s_x}$

Slope of Least Squares Line:  $SSE = \sum (y_i - \hat{y}_i)^2 \leq \hat{\beta}_1 = \{\sum (x_i - \bar{x})(y_i - \bar{y})\} / \{\sum (x_i - \bar{x})^2\} \leq \hat{\beta}_0 = \bar{y} - \hat{\beta}_1 \bar{x}$  |  $t_{\text{slope}} = \hat{\beta}_1 \sqrt{\frac{\sum (x_i - \bar{x})^2}{MSE}} \leq MSE = \frac{SSE}{n-2}$  |  $\# \text{ of } \beta_i \text{'s} \geq \hat{\beta}_1^{\text{obs}}$

Process of Permutation test for Slope / Correlation: - Compute  $\hat{\beta}_1^{\text{obs}}$  or  $r_{\text{obs}}$ , and compute  $\hat{\beta}_i$ 's or  $r_i$ 's for all  $n!$  permutation (Permute Y's among the X's)  $\Rightarrow$

L Large Sample Approximation:  $Z_r = r \sqrt{n-1}$  | Spearman Rank Correlation: Same Process with Permutation test for Correlation except that it measures correlation with ranks.  $Z = \frac{r_n}{\sqrt{\text{Var}(r_n)}}$

Kendall's T: - Concordant means  $x_i < x_j$  implies  $y_i < y_j$  for  $i < j$  |  $T = 2P[(x_i - x_j)(y_i - y_j) > 0] - 1$  \*  $T = \frac{1}{2}$  implies no association | Large-Sample Approximation:  $E(T_r) = 0$ ,  $\text{Var}(T_r) = \frac{4n+10}{9(n-1)}$

Permutation Tests for Contingency Tables | Permutation  $\chi^2$  Test: - Compute all  $\binom{n}{r}$   $\chi^2$ -statistics, where  $n = \# \text{ of units}$ ,  $r = \# \text{ of rows}$ .  $\chi^2 = \sum_{i=1}^c \sum_{j=1}^r (a_{ij} - e_{ij})^2 / e_{ij}$ ,  $e_{ij} = \frac{n_{i,j}}{n} \Rightarrow$

=> and count the # of  $\chi^2$ -statistics greater than or equal to  $\chi^2_{\text{obs}}$  | Fisher's Exact Test: let  $X$  denote  $a_{11}$ , then  $PL(X \leq a_{11}) = \frac{n!}{a_{11}! n_{11}!} \frac{n_{11}!}{(n-a_{11})! (n-n_{11})!}$ , red-colored are fixed

McNemar's Test:  $P(T_{MN} \geq T_{MN}^{\text{obs}}) = \sum_{i=0}^n \binom{n}{i} (0.5)^n \Rightarrow Z_{MN} = \frac{T_{MN} - 0.5n}{\sqrt{0.25 \cdot n}}$  | Mantel-Haenszel Test:  $E(X_{ik}) = \frac{r_{ik} c_{ik}}{N_k}$ ,  $\text{Var}(X_{ik}) = \frac{r_{ik} c_{ik} c_{ik} r_{ik}}{N_k(N_k-1)}$   $\Rightarrow MH = (\sum (X_{ik} - E(X_{ik}))^2) / \sum \text{Var}(X_{ik})$

MH Estimate of the common odds ratio:  $\hat{\theta} = \frac{A}{B}$ ,  $A = \sum n_{11k} n_{21k} / N_k$ ,  $B = \sum n_{11k} n_{21k} / N_k \Rightarrow \text{Var}(\log \hat{\theta}) = \{\sum (n_{11k} + n_{21k})(n_{11k} n_{21k}) / N_k^2\} / 2A^2 + \{\sum (n_{11k} + n_{21k})(n_{12k} n_{22k}) / N_k^2\} / 2AB + \dots$

-->  $\{\sum (n_{12k} + n_{22k})(n_{12k} n_{22k}) / N_k^2\} / 2B^2$  | CH 8

Bootstrap Variance and Bias:  $B = E(\hat{\theta}) - \theta$ ,  $MSE = \text{Var} + B^2$ ,  $\hat{E} = \frac{1}{REP} \sum \hat{\theta}_{bi}$ ,  $\hat{B} = \hat{E} - \theta$ ,  $\text{Var} = \frac{1}{REP} \sum (\hat{\theta}_{bi} - \hat{E})^2$  | For symmetric dist, just use t-dist with df=n-1

Process of Obtaining a Bootstrap CI for  $\mu$ : - Compute  $\bar{X}_{\text{obs}}$  and  $S_{\text{obs}}$ , and compute all  $\bar{X}_b$ 's,  $S_b$ 's, and  $t_b$ 's  $\Rightarrow \bar{X} - t_{b, 0.025}(\frac{S}{\sqrt{n}}) < \mu < \bar{X} + t_{b, 0.025}(\frac{S}{\sqrt{n}})$  | The process is the same

CI for the Variance and Standard Deviation:  $\{(n-1)^2\} / \chi^2_{0.975} < S^2 < \{(n-1)^2\} / \chi^2_{0.025} \Leftrightarrow P(\chi^2_{0.025} < \frac{(n-1)^2}{S^2} < \chi^2_{0.975})$ , if normal,  $\# \text{ of } \chi^2 = \{(n-1)^2\} / S^2$  | Except for using  $\chi^2$ -statistics

Percentile: - draw bootstrap samples and compute b number of  $\hat{\theta}_b$ .  $\Rightarrow$  Find  $(1-\alpha/2)100^{\text{th}}$  and  $(\alpha/2)100^{\text{th}}$  percentiles of the bootstrap distribution.

Residual Method: - draw bootstrap samples and compute  $e_b = \hat{\theta}_b - \theta$ .  $\Rightarrow$  Find  $(1-\alpha/2)100^{\text{th}}$  and  $(\alpha/2)100^{\text{th}}$  percentiles of the bootstrap distribution.  $\Rightarrow Z_0 = \phi^{-1}\{\frac{1}{REP} \sum (\hat{\theta}_b < \theta)\}$

Bias-Corrected and Accelerated Method BCA: BCA adjusts for bias and skewness. If the bootstrap distribution is symmetric, then BCA, percentile, and residual method give the same conclusions.

- given  $a = \{\sum (\hat{\theta}_i - \theta)^3\} / \{\sum (\hat{\theta}_i - \theta)^2\}^{3/2}$ , where  $\hat{\theta}_i$  is the mean of  $\theta_i$ 's.  $\Rightarrow E[T(\theta)] = T(\theta) - \theta_0 [1 + QT(\theta)]$ ,  $\text{Var}[T(\theta)] = 1 - aQT(\theta) \Rightarrow \frac{T(\theta) + \theta_0 - \theta_p}{1 - a(\theta_0 - \theta_p)} < T(\theta) < \frac{T(\theta) + \theta_0 - \theta_p}{1 - a(\theta_0 + \theta_p)} \Rightarrow Z_U = \frac{\theta_0 - \theta_p}{1 - a(\theta_0 + \theta_p)} + Z_0$ ,  $Z_L = \frac{\theta_0 - \theta_p}{1 - a(\theta_0 - \theta_p)} + Z_0$

Bivariate Bootstrap Sampling: random sample with replacement of the pairs  $(x_i, y_i)$ ,  $i=1, 2, \dots, n$ . \* ns equally likely possible bootstrap samples

- draw a specified # of bivariate samples of size n and compute  $r_b$  for each bootstrap sample.  $\Rightarrow$  Obtain CI from  $r_b$ 's using BCA or percentile method.

Fixed-X Bootstrap Sampling: assume  $h(x)$  is a function of the independent variable, and  $Y$  may be expressed in terms of the regression model,  $Y = h(x) + \epsilon$

- Compute an estimate  $\hat{h}(x)$  and obtain  $e_i = Y_i - \hat{h}(x_i)$ .  $\Rightarrow$  Select n values of the  $e_i$ 's at random with replacement.  $\Rightarrow$  Compute  $\bar{Y}_i = \hat{h}(x_i) + e_{i,b}$

Bootstrap Inferences for the Slope of a Regression Line: - Compute the OLS estimates and the residuals  $e_i = Y_i - \hat{\theta}_0 - \hat{\theta}_1 x_i$ ,  $\Rightarrow SE(\hat{\theta}_1) = \sqrt{\frac{MSE}{\sum (x_i - \bar{x})^2}}$ ,  $t = \frac{\hat{\theta}_1 - \theta_1}{SE(\hat{\theta}_1)}$ ,  $t_{0.025} < \frac{\hat{\theta}_1 - \theta_1}{SE(\hat{\theta}_1)} < t_{0.975}$

=> - take samples of size n of the  $e_i$ 's with replacement, and form the bootstrap sample  $(x_i, e_{i,b})$ . Compute the least squares estimate of the slope  $\hat{\theta}_{1,b}$ , the estimated  $SE(\hat{\theta}_{1,b})$ , and  $t_b = \hat{\theta}_{1,b} / SE(\hat{\theta}_{1,b})$

CH 10 - Repeat above to obtain the bootstrap distribution of the  $t_b$ 's, and construct a CI or conduct a test of hypothesis for  $\theta_1$ .

Histogram: suggested width,  $d = 3.5(S/\sqrt{n})$ , where  $S = \text{sample SD}$  | Kernel Method:  $\hat{f}(x) = \frac{1}{n\Delta} \sum W\left(\frac{x-x_i}{\Delta}\right)$ , where  $W(z) = \text{kernel}$ ,  $\Delta = 1.06(S/\sqrt{n})$  |  $k$  smoothness

KNN Regression: For any  $x_0$ , we can always find K nearest values of  $x_i$  to  $x_0$ . We take average of those values and set it as the value of  $x_0$ .  $\Rightarrow \hat{f}(x_0) = \frac{1}{K} \sum Y_i$  |  $K$  smoothness

Loess Method:  $\sum_{x_i \in N_k(x_0)} [Y_i - \hat{f}(x_i)]^2 W(|x_0 - x_i|/\Delta_0)$ , where  $W(z) = \text{linear/quadratic/higher}$ ,  $N_k(x_0) = K$  nearest  $x_i$ 's,  $W(u) = (1-u)^3$ ,  $\Delta_0 = \max_{x_i \in N_k(x_0)} |x_i - x_0|$  | Span =  $K$ . Span  $\uparrow$  complexity  $\downarrow$ , Span  $\downarrow$  complexity  $\uparrow$

Kernel Method: uses points within the range of  $x_0 \pm h$ ,  $h$  = bandwidth  $\Rightarrow$   $h \uparrow$  complexity  $\downarrow$ ,  $h \downarrow$  complexity  $\uparrow$  |  $\hat{f}(x, y) = \frac{1}{n_h \Delta y} \sum_{i=1}^{n_h} W\left(\frac{x-x_i}{\Delta x}\right) W\left(\frac{y-y_i}{\Delta y}\right) \Rightarrow \hat{f}(x, y) = \int f(x, y) dy \approx \left\{ \frac{1}{n_h} \sum_{i=1}^{n_h} W\left(\frac{x-x_i}{\Delta x}\right) \right\} / \left\{ \sum_{i=1}^{n_h} W\left(\frac{y-y_i}{\Delta y}\right) \right\}$

$MSE = (Y_i - \hat{f}(x_i))^2 \Rightarrow CV_{(n)} = \frac{1}{n} \sum MSE$ ; \* using n is LOOCV and using n  $\rightarrow$  K is K-Fold CV. Knots: the points where the coefficients change. (usually set as quantiles) # of knots  $\uparrow$  complexity

Regression Splines (Piecewise Polynomials): fitting separate low-degree polynomials over different regions of  $X$ ,  $\Rightarrow Y_i = \theta_0 + \theta_1 x_i + \theta_2 x_i^2 + \dots + \theta_k x_i^k$ , where  $\theta_i$ 's differ in different parts of the range of  $X$ .

Cubic Splines, Natural Splines | Smoothing Splines: giving penalty to models with high variances.  $\sum (y_i - g(x_i))^2 + \lambda \int g''(x)^2 dx$

Generalized Additive Models:  $Y_i = \theta_0 + f_1(x_{i1}) + f_2(x_{i2}) + \dots + f_p(x_{ip}) = \theta_0 + \sum_{j=1}^p f_j(x_{ij}) + \epsilon_i$

Regression Trees: We choose to divide the regions (Predictor space) into high-dimensional rectangles or boxes. The goal is to find boxes  $R_j$ 's that minimizes RSS =  $\sum (Y_i - \hat{f}(x_i))^2$

Tree-Pruning: growing a very large tree and pruning it back in order to obtain a subtree. | Bagging: generating B different bootstrap datasets  $\Rightarrow \hat{f}(x) = \frac{1}{B} \sum \hat{f}^b(x) \Rightarrow$  Variable Importance

Random Forests: Unlike Bagging, each time a split in a tree is considered, a random sample of m predictors is chosen as split candidates from the full set of p predictors.

- In building a random forest, the algorithm is not even allowed to consider a majority of the available predictors.

- Suppose there is one very strong predictor  $x_1$  in the data. In the collection of trees, most or all of the trees will use  $x_1$  in the top split.

- We've seen that averaging highly correlated quantities doesn't lead to as large of reduction in variance as averaging many uncorrelated quantities

Boosting: each tree is grown using information from previously grown trees; each tree is fit on a modified version of the original data set. | d)  $\hat{f}(x) = \sum \lambda \hat{f}^b(x)$

- Set  $\hat{f}(x) = 0$  and  $r_i = y_i$  for all training data set.  $\Rightarrow$  Repeat B times: a) fitting  $\hat{f}^b$  with d+1 terminal nodes. b) update  $\hat{f}(x) \leftarrow \hat{f}(x) + \lambda \hat{f}^b(x)$  c) Update  $r_i \leftarrow r_i - \lambda \hat{f}^b(x_i)$

\* boosting can overfit if B is too large. - Very small  $\lambda$  can require using a very large B for good performance

- The number d of splits in each tree controls the complexity of the boosted model.