

Statistical Analysis Cheatsheet

Compiled by Gejun Zhu (zhug3@miamioh.edu) in preparation for the analysis comprehensive exam by using William Chen's formula sheet template.

Last Updated August 11, 2015

Regression Analysis

Simple Linear Regression

Model : $Y_i = \beta_0 + \beta_1 X_i + \epsilon_i, \epsilon_i \stackrel{iid}{\sim} N(0, \sigma^2), Y_i \sim N(\beta_0 + \beta_1 X_i, \sigma^2)$.

$$E_i = Y_i - \hat{Y}_i, \sum E_i = 0, \sum Y_i = \sum \hat{Y}_i, \sum X_i E_i = 0, \sum \hat{Y}_i E_i = 0.$$

$$B_1 = \frac{\sum_{i=1}^n (X_i - \bar{X})(Y_i - \bar{Y})}{\sum_{i=1}^n (X_i - \bar{X})^2} = \frac{S_{XY}}{S_{XX}}, S_{XY} = \sum XY - \frac{\sum X \sum Y}{n}, B_0 = \bar{Y} - B_1 \bar{X}.$$

$$B_1 = \sum_{i=1}^n Y_i \frac{(X_i - \bar{X})}{\sum_{i=1}^n (X_i - \bar{X})^2} = \sum_{i=1}^n k_i Y_i; \sum k_i = 0, \sum k_i X_i = 1, \sum k_i^2 = \frac{1}{S_{XX}}.$$

$$B_1 \sim N(\beta_1, \frac{\sigma^2}{S_{XX}}); \frac{B_1 - \beta_1}{\sqrt{V(B_1)}} \sim N(0, 1), \frac{S_{B_1}^2}{V(B_1)} = \frac{MSE/S_{XX}}{\sigma^2/S_{XX}} = \frac{SSE}{(n-2)\sigma^2},$$

$$\frac{SSE}{\sigma^2} \sim \chi_{n-2}^2 \Rightarrow \frac{B_1 - \beta_1}{S_{B_1}} = \frac{B_1 - \beta_1}{\sqrt{MSE/S_{XX}}} \sim T_{n-2}, \text{CI: } b_1 \pm t_{1-\alpha/2, n-2} \cdot s_{b_1}, \frac{SSR}{\sigma^2} \sim \chi_{p-1}^2.$$

$$\text{Inference on } E(Y_h): \hat{Y}_h = B_0 + B_1 X_h, \hat{Y}_h \sim N(\beta_0 + \beta_1 X_h, \sigma^2 [\frac{1}{n} + \frac{(X_h - \bar{X})^2}{S_{XX}}])$$

$$\text{Prediction on a new observation: } \hat{y} \pm t_{1-\alpha/2, n-2} \sqrt{mse[1 + \frac{1}{n} + \frac{(X_h - \bar{X})^2}{S_{XX}}]}.$$

$$SST = \sum (Y_i - \bar{Y})^2 = \sum (Y_i - \hat{Y}_i)^2 + \sum (\hat{Y}_i - \bar{Y})^2 = SSE + SSR$$

$$\text{If } var(Y_i) = \sigma^2, \text{ and } Y_i\text{'s are uncorrelated, then } Cov(\sum a_i Y_i, \sum b_i Y_i) = \sigma^2 \sum a_i b_i.$$

$$B_1 \text{ and } \bar{Y} \text{ are uncorrelated, } Cov(B_1, \bar{Y}) = 0 \text{ because}$$

$$Cov(B_1, \bar{Y}) = Cov(\sum k_i Y_i, \sum \frac{1}{n} Y_i) = \frac{\sigma^2}{n} \sum k_i = 0.$$

ANOVA Table - Analysis of variance for simple linear regression

Source	SS	DF	MS	Expected
Regression	SSR	1	MSR=SSR/1	$\sigma^2 + \beta_1^2 S_{XX}$
Error	SSE	n-2	MSE=SSE/(n-2)	σ^2
Total	SST	n-1		

$$\text{Under } H_0: \beta_1 = 0, F^* = \frac{MSR}{MSE} \sim F_{1, n-2}, R^2 = \frac{SSR}{SST} = 1 - \frac{SSE}{SST}, B_1 = r \sqrt{\frac{S_{YY}}{S_{XX}}}.$$

$$E(MSR) = \sigma^2 + \beta_1^2 S_{XX}, SSR = B_1^2 S_{XX}, E(MSE) = E(\frac{SSE}{n-2}) = \frac{\sigma^2}{n-1} E(\frac{SSE}{\sigma^2}) = \sigma^2$$

$$\text{Studentized residuals: } E_i^* = \frac{E_i}{\sqrt{V(\hat{E}_i)}}$$

Assumptions - LINE

- Linearity: No curvature in the residual plot; (high-order, log/square)
- Independence: $\epsilon_i \stackrel{iid}{\sim} N(0, \sigma^2)$;
- Normality of error: QQ plot; (GLM, poisson regression...)
- Equal Variance: standardized residual inside [-3, 3]. (weighted obs)

Matrix Approach - Matrix form

$$\mathbf{Y}_{n \times 1} = \mathbf{X}_{n \times p} \boldsymbol{\beta}_{p \times 1} + \boldsymbol{\epsilon}_{n \times 1}, \boldsymbol{\epsilon} \sim MN(\mathbf{0}, \sigma^2 \mathbf{I}). \boldsymbol{\beta} = (\mathbf{X}'\mathbf{X})^{-1} \mathbf{X}'\mathbf{Y}.$$

$$\mathbf{Y} \sim MN(\boldsymbol{\mu}, \boldsymbol{\Sigma}), \text{ then } \mathbf{AY} + \mathbf{b} \sim MN(\mathbf{A}\boldsymbol{\mu} + \mathbf{b}, \mathbf{A}\boldsymbol{\Sigma}\mathbf{A}').$$

$$\mathbf{B} \sim MN(\boldsymbol{\beta}, \sigma^2 (\mathbf{X}'\mathbf{X})^{-1}), \hat{\mathbf{Y}} = \mathbf{XB} = \mathbf{HY} \sim MN(\mathbf{X}\boldsymbol{\beta}, \sigma^2 \mathbf{H}), \mathbf{H} = \mathbf{X}(\mathbf{X}'\mathbf{X})^{-1} \mathbf{X}'.$$

$$S^2(\mathbf{B}) = MSE(\mathbf{X}'\mathbf{X})^{-1}, \mathbf{E} = \mathbf{Y} - \hat{\mathbf{Y}} = (\mathbf{I} - \mathbf{H})\mathbf{Y} \sim MN(\mathbf{0}, \sigma^2 (\mathbf{I} - \mathbf{H})).$$

$$\sum Y_i^2 = \mathbf{Y}'\mathbf{Y}, SSTO = \mathbf{Y}'(\mathbf{I} - \frac{1}{n})\mathbf{Y}, SSE = \mathbf{Y}'(\mathbf{I} - \mathbf{H})\mathbf{Y}, SSR = \mathbf{Y}'(\mathbf{H} - \frac{1}{n})\mathbf{Y}.$$

$$\hat{\mathbf{Y}} = \mathbf{HY}, \hat{\mathbf{Y}}' = \mathbf{Y}'\mathbf{H}, 0 = \sum \hat{Y}_i E_i = \sum \hat{Y}_i Y_i - \sum \hat{Y}_i^2.$$

$$\text{Distribution of } \hat{Y}_h: \hat{Y}_h = \mathbf{X}_h' \mathbf{B} \sim MN(\mathbf{X}_h' \boldsymbol{\beta}, \sigma^2 \mathbf{X}_h' (\mathbf{X}'\mathbf{X})^{-1} \mathbf{X}_h)$$

$$S^2(\hat{Y}_h) = MSE(\mathbf{X}_h' (\mathbf{X}'\mathbf{X})^{-1} \mathbf{X}_h)$$

$$pred = Y_{h(new)} - \hat{Y}_h, pred \sim N(0, \sigma^2 (1 + \mathbf{X}_{h(new)}' (\mathbf{X}'\mathbf{X})^{-1} \mathbf{X}_{h(new)}))$$

Multiple Regression - Multiple regression

$$\mathbf{Y} = \mathbf{X}\boldsymbol{\beta} + \boldsymbol{\epsilon}, \text{ where } \boldsymbol{\epsilon} \sim MN(\mathbf{0}, \sigma^2 \mathbf{I}).$$

FWER - Bonferroni and Holm. Bonferroni: compare p-value with α/g ; Holm: sort p-values and multiple g, g-1, ..., 1 in order and compare with α finally.

$$SSTO = SSR(X_1) + SSE(X_1) = SSR(X_1, X_2) + SSE(X_1, X_2), \text{ where}$$

$$SSR(X_1, X_2) = SSR(X_1) + SSR(X_2|X_1), SSE(X_1, X_2) = SSE(X_2) - SSR(X_1|X_2).$$

$$\text{In general, } SSR(X_q, X_{q+1}, \dots, X_{p-1}|X_1, X_2, \dots, X_{q-1}) =$$

$$SSE(X_1, X_2, \dots, X_{q-1}) - SSE(X_1, X_2, \dots, X_{p-1}) = SSE_R - SSE_F.$$

$$\text{Partial F-test: } H_0: \beta_q = \beta_{q+1} = \dots = \beta_{p-1} = 0, H_a: \text{At least one}$$

$$\beta_k \neq 0, k = q, q+1, \dots, p-1. \text{ Test statistic}$$

$$F^* = \frac{\frac{SSE_R - SSE_F}{df_R - df_F}}{MSE_F} = \frac{\frac{SSR(X_q, X_{q+1}, \dots, X_{p-1}|X_1, X_2, \dots, X_{q-1})}{p-q}}{\frac{SSE(X_1, X_2, \dots, X_{p-1})}{n-p}}. \text{ If } H_0 \text{ is true, then}$$

$$F^* \sim F_{p-q, n-p}.$$

$$\text{coefficient of partial determination is the proportion of the variation in Y}$$

$$\text{"explained" by an indep. variable when other indep. variables are in the model.}$$

$$R_{Y1|2}^2 = \frac{SSE(X_2) - SSE(X_1, X_2)}{SSE(X_2)} = \frac{SSR(X_1|X_2)}{SSE(X_2)} = 1 - \frac{SSE(X_1, X_2)}{SSE(X_2)}$$

$$R_{Y1|2,3,4} = 1 - \frac{SSE(X_1, X_2, X_3, X_4)}{SSE(X_2, X_3, X_4)}$$

$$\text{multicollinearity diagnostic: Variance Inflation Factor (VIF) = } (1 - R_k^2)^{-1}, \text{ where}$$

$$R_k^2 = \text{coefficient of determination when } X_k \text{ is regressed upon other predictors.}$$

$$\text{If } VIF > 1, \text{ variance of } B_k \text{ is inflated due to correlations b/w } X_k \text{ and other predictors. If } X_k \text{ is uncorrelated with other predictors, then } R_k^2 = 0 \text{ and } VIF_k = 1.$$

Model Diagnostics - More about model diagnostics

Added-variable Plots (1) regress Y on predictors except X_k and obtain the

residuals; (2) regress X_k on other predictors and obtain residuals; (3) plot (1) vs (2);

Leverage: A measure of how unusual an X is. (diagonal values of Hat matrix,

$$\sum h_{ii} = tr(H) = tr[X(X'X)^{-1}X'] = tr[X'X(X'X)^{-1}] = tr[I_{p \times p}] = p)$$

Influence: An influence point is its exclusive causes substantial changes to the

fitting data. Just because a point has high leverage doesn't mean it has high influence.

Measures of influence include:

$$DFFITs_i = \frac{\hat{Y}_i - \hat{Y}_i(i)}{\sqrt{MSE(i)h_{ii}}}$$

$$\text{Cook's Distance} = \frac{\hat{\epsilon}_i^2}{pMSE} \left[\frac{h_{ii}}{(1-h_{ii})^2} \right]$$

$$DFBETAS_{k(i)} = \frac{b_k - b_k(i)}{\sqrt{MSE(i) c_{kk}}} \text{ where } c_{kk} \text{ is } (k+1)^{th} \text{ diagonal in } (X'X)^{-1}.$$

Design of experiments

CRD with one factor

Model one factor with $a \geq 2$ levels. $H_0: \mu_1 = \mu_2, \dots, \mu_a$ or $\hat{\tau}_i = 0$.

- $y_{ij} = \mu + \tau_i + \epsilon_{ij}, i = 1, 2, \dots, a, j = 1, 2, \dots, n_j, \epsilon_{ij} \sim N(0, \sigma^2)$;
- LSE estimator $\hat{\mu} + \hat{\tau}_i = \bar{y}_{i.}$, if $\sum n_i \hat{\tau}_i = 0$ or $\hat{\mu} \hat{=} 0$ or $\hat{\tau}_a \hat{=} 0$;
- $MS_{trt} = \frac{SS_{trt}}{a-1} = \frac{\sum_{i=1}^a n_i (\bar{y}_{i.} - \bar{y}_{..})^2}{a-1}, MSE = \frac{SSE}{N-a} = \frac{\sum_{i=1}^a \sum_{j=1}^{n_i} (y_{ij} - \bar{y}_{i.})^2}{N-a}$;
- $E(MSE) = \sigma^2, E(MS_{trt}) = \sigma^2 + \frac{\sum_{i=1}^a n_i \tau_i^2}{a-1}, S_p^2 = \frac{(n_1-1)S_1^2 + (n_2-1)S_2^2}{n_1+n_2-2}$;
- $SS_T = \sum_{i=1}^a \sum_{j=1}^{n_i} y_{ij}^2 - \frac{(\sum y_{..})^2}{N}, SS_{trt} = \sum_{i=1}^a \frac{y_{i.}^2}{n_i} - \frac{(\sum y_{..})^2}{N}$;
- Fact: Under $H_0, SSE/\sigma^2 \sim \chi_{N-a}^2, SS_{trt}/\sigma^2 \sim \chi_{a-1}^2$, independent;
- $\frac{SS_{trt}/(a-1)\sigma^2}{SSE/(N-a)\sigma^2} \sim F_{a-1, N-a}$; rej $H_0 > F(\alpha, a-1, N-a), p = P(F_{a-1, N-a} > F_0)$;
- $E(\bar{y}_{i.}) = \mu_i, V(\bar{y}_{i.}) = \sigma^2/n_i, \frac{\bar{y}_{i.} - \mu_i}{\sqrt{MSE/n_i}} \sim T_{N-a}$;
- CI: $\bar{y}_{i.} \pm t_{\alpha/2, N-a} \sqrt{MSE/n_i}, \bar{y}_{s.} - \bar{y}_{t.} \pm t_{\alpha/2, N-a} \sqrt{MSE/n_s + MSE/n_t}$.
- Linear contrasts: $\Gamma = \sum c_i \mu_i, C = \sum c_i \hat{y}_{i.}$ with $\sum c_i = 0. E(C) = \Gamma,$

$$V(C) = \sigma^2 \sum \frac{c_i^2}{n_i}. \text{ CI: } \sum c_i \hat{y}_{i.} \pm t_{\alpha/2, N-a} \sqrt{MSE \sum \frac{c_i^2}{n_i}}, t = \frac{\sum c_i \hat{y}_{i.} - c}{\sqrt{MSE \sum \frac{c_i^2}{n_i}}}.$$

ANOVA Table Analysis of variance for three factor fixed effects model.

Source	DF	Expected Mean Square
A	$a - 1$	$\sigma^2 + \frac{bcn \sum \tau_i^2}{a-1}$
AB	$(a-1)(b-1)$	$\sigma^2 + \frac{cn \sum \sum (\tau\beta)_{ij}^2}{(a-1)(b-1)}$
ABC	$(a-1)(b-1)(c-1)$	$\sigma^2 + \frac{n \sum \sum \sum (\tau\beta\gamma)_{ijk}^2}{(a-1)(b-1)(c-1)}$
Error	$abc(n-1)$	σ^2

Basic Blocking Designs

Model two factors - the treatment factor τ_i and the block factor β_j .

$$Y_{ij} = \mu + \tau_i + \beta_j + \epsilon_{ij}, i = 1, 2, \dots, a, j = 1, 2, \dots, b. (\sum \tau_i = 0 \text{ and } \sum \beta_j = 0)$$

$$H_0: \tau_0 = \tau_1 = \dots = \tau_a = 0 \text{ or } \mu_1 = \mu_2 = \dots = \mu_a, H_a: \text{Not } H_0 \text{ (at least two means differs)} . E(\bar{Y}_{i.}) = \mu + \tau_i$$

A balanced incomplete block design (BIBD) includes a treatment factor with a levels, a blocking factor with b levels, each block includes k experimental units, which implies a total of bk runs. This means that each treatment occurs $r = bk/a$ times. Each treatment occurs either 0 or 1 times, and each pair of treatments occurs together in a block exactly λ times. $N = bk$. (1) $ar = bk$; (2) $r(k-1) = \lambda(a-1)$; (3) $b \geq a$.

Source	DF	Sum of Squares
Treatments	$a - 1$	$\sum_i \frac{y_{i.}^2}{b} - \frac{y_{..}^2}{N}$
Blocks	$b - 1$	$\sum_j \frac{y_{.j}^2}{a} - \frac{y_{..}^2}{N}$
Error	$N - a - b + 1$	$SS_{total} - SS_{trts} - SS_{blocks}$
Total	$N - 1$	$\sum \sum y_{ij}^2 - \frac{y_{..}^2}{N}$

$$E(MS_{trt}) = \sigma^2 + \frac{b \sum \tau_i^2}{a-1}, E(MS_{blk}) = \sigma^2 + \frac{a \sum \beta_j^2}{b-1}, E(MSE) = \sigma^2$$

$$F_0 = MSE/MS_{trt}, p\text{-value} = P(F_{a-1, (b-1)(a-1)} > F_0).$$

$$Q_i = y_{i.} - \frac{1}{k} \sum_j n_{ij} y_{.j}, \hat{\tau}_i = \frac{k Q_i}{\lambda a}, \hat{\mu} = \frac{y_{..}}{N} = \frac{y_{..}}{bk}, LS\text{Mean}(\mu_i) = \hat{\mu} + \hat{\tau}_i$$

2^k Factorial Designs

Model $Y_{ijk} = \mu + \tau_i + \beta_j + (\tau\beta)_{ij} + \epsilon_{ijk}, i = 1, 2, \dots, a, j = 1, 2, \dots, b, k = 1, 2, \dots, n$

$$\sum \tau = 0, \sum \beta = 0, \sum_i (\tau\beta)_{ij} = 0, \sum_j (\tau\beta)_{ij} = 0$$

$$\hat{\mu} = \bar{y}_{...}, \hat{\tau}_i = \bar{y}_{i..} - \bar{y}_{...}, \hat{\beta}_j = \bar{y}_{.j.} - \bar{y}_{...}, \hat{\tau}\hat{\beta}_{ij} = \bar{y}_{ij.} - \bar{y}_{i..} - \bar{y}_{.j.} + \bar{y}_{...}$$

$$\text{Overall test: } \mu_{11} = \mu_{12} = \dots = \mu_{ab}, \text{ test statistic } F_0 = \frac{MS_{trt}}{MSE};$$

$$\text{Interaction test: } (\alpha\beta)_{ij} = 0 \text{ for all } ij, \text{ test statistic } F_0 = \frac{MS_{AB}}{MSE}.$$

Two-level Fractional Factorial Designs

Design resolution - A fractional factorial design's resolution is the length of the shortest word and its defining relation.

2^{k-p} terms, 2^p alias.

Random Effects and Mixed Models

Model $Y_{ij} = \mu + \tau_i + \epsilon_{ij}, i = 1, 2, \dots, a; j = 1, 2, \dots, n_i$, where τ_i are assumed to be independent $N(0, \sigma_\tau^2)$ random variables.

$$H_0: \sigma_\tau^2 = 0 \text{ vs. } H_a: \sigma_\tau^2 > 0, \text{ test stat: } F_0 = \frac{MS_{trt}}{MSE}, F_0 \sim F_{a-1, N-a} \text{ under } H_0.$$

Some facts: $Y_{ij} \sim N(\mu, \sigma_\tau^2 + \sigma^2)$ (1) if $i \neq k$ - different treatment levels,

$$Cov(Y_{ij}, Y_{kj}) = 0 \text{ since } \tau_i \text{ and } \tau_k \text{ are independent and } E(\tau_i \tau_k) = E(\tau_i)E(\tau_k) = 0;$$

$$(2) \text{ if } k \neq l - \text{same treatment different obs, } Cov(Y_{ij}, Y_{kl}) = \sigma_\tau^2.$$

$$E(MSE) = \sigma^2, E(MS_{trt}) = \sigma^2 + n_0 \sigma_\tau^2 \text{ where } n_0 = n \text{ if all } n_i = n \text{ and}$$

$$n_0 = \frac{1}{a-1} [N - \frac{\sum n_i^2}{N}].$$

$$\text{Estimates: } \hat{\sigma}^2 = MSE \text{ and } \hat{\sigma}_\tau^2 = \frac{MS_{trt} - MSE}{n_0}.$$

$$\text{Confidence interval for } \frac{\sigma_\tau^2}{\sigma_\tau^2 + \sigma^2}: \frac{MS_{trts}/(n\sigma_\tau^2 + \sigma^2)}{MSE/\sigma^2} \sim F_{a-1, N-a},$$

$$(F_{1-\alpha/2, a-1, N-a} \leq \frac{MS_{trts}}{MSE} \frac{\sigma_\tau^2}{n\sigma_\tau^2 + \sigma^2} \leq F_{\alpha/2, a-1, N-a}) = 1 - \alpha, P(L \leq \frac{\sigma_\tau^2}{\sigma_\tau^2 + \sigma^2} \leq U) = 1 - \alpha,$$

$$L = \frac{1}{n} \left(\frac{MS_{trts}}{MSE} \frac{1}{F_{\alpha/2, a-1, N-a}} - 1 \right), U = \frac{1}{n} \left(\frac{MS_{trts}}{MSE} \frac{1}{F_{1-\alpha/2, a-1, N-a}} - 1 \right),$$

$$\frac{L}{L+1} \leq \frac{\sigma_\tau^2}{\sigma_\tau^2 + \sigma^2} \leq \frac{U}{U+1}$$

Two-factor factorial with random factors: $Y_{ijk} = \mu + \tau_i + \beta_j + (\tau\beta)_{ij} + \epsilon_{ijk}, i = 1, 2, \dots, a, j = 1, 2, \dots, b, k = 1, 2, \dots, n, V(\tau_i) = \sigma_\tau^2, V(\beta_j) = \sigma_\beta^2, V[(\tau\beta)_{ij}] = \sigma_{\tau\beta}^2, \text{ and } V(\epsilon) = \sigma^2.$

Expected mean squares: $E(MS_A) = \sigma^2 + n\sigma_{\tau\beta}^2 + bn\sigma_{\tau}^2$;
 $E(MS_B) = \sigma^2 + n\sigma_{\tau\beta}^2 + an\sigma_{\beta}^2$; $E(MS_{AB}) = \sigma^2 + n\sigma_{\tau\beta}^2$; $E(MSE) = \sigma^2$.
Two-factor mixed model: Factor A is fixed; factor B is random.
 $Y_{ijk} = \mu + \tau_i + \beta_j + (\tau\beta)_{ij} + \epsilon_{ijk}$, where

- $i = 1, 2, ..., a, j = 1, 2, ..., b, k = 1, 2, ..., n$;
- τ_i is a fixed effect with $\sum \tau_i = 0$;
- $\beta_j \sim N(0, \sigma_{\beta}^2)$, $(\tau\beta)_{ij} \sim N(0, \sigma_{\tau\beta}^2)$, and $\epsilon_{ijk} \sim N(0, \sigma^2)$.

$Y_{ijk} \sim N(\mu + \tau_i, \sigma^2 + \sigma_{\beta}^2 + \sigma_{\tau\beta}^2)$.
Expected mean squares: $E(MSE) = \sigma^2$, $E(MS_A) = \sigma^2 + n\sigma_{\tau\beta}^2 + bn\frac{\sum \tau_i^2}{a-1}$,
 $E(MS_B) = \sigma^2 + n\sigma_{\tau\beta}^2 + an\sigma_{\beta}^2$, $E(MS_{AB}) = \sigma^2 + n\sigma_{\tau\beta}^2$
Variance components estimates: $\hat{\sigma}^2 = MSE$, $\hat{\sigma}_{\tau\beta}^2 = \frac{MS_{AB}-MSE}{n}$, $\hat{\sigma}_{\beta}^2 = \frac{MS_B-MS_{AB}}{an}$
Hypothesis Tests: (1) $H_0 : \sigma_{\tau\beta}^2 = 0$ vs. $H_a : \sigma_{\tau\beta}^2 > 0$ using $F = \frac{MS_{AB}}{MSE}$; (2)
 $H_0 : \sigma_{\beta}^2 = 0$ vs. $H_a : \sigma_{\beta}^2 > 0$ using $F = \frac{MS_B}{MS_{AB}}$; (3) $H_0 : \tau_i = 0$ vs. $H_a : not\ H_0$
using $F = \frac{MS_A}{MS_{AB}}$.

Approximate F-test: degree of freedom $\nu = \frac{(\sum c_i MS_i)^2}{\sum \frac{c_i^2 MS_i^2}{\nu_i}}$

Nested Designs

Model $Y_{ijk} = \mu + \tau_i + \beta_{j(i)} + \epsilon_{ijk}$, $i = 1, 2, ..., a; j = 1, 2, ..., b, k = 1, 2, ..., n$.
A random; B(A) random: $Cov(Y_{ijk}, Y_{mno}) = \sigma_{\beta}^2 + \sigma_{\tau}^2$ if $i = m, j = n, k \neq o$;
 $Cov(Y_{ijk}, Y_{mno}) = \sigma_{\tau}^2$ if $i = m, j \neq n$; $Cov(Y_{ijk}, Y_{mno}) = 0$ if $i \neq m$;
A fixed; B(A) random: $Cov(Y_{ijk}, Y_{mno}) = \sigma_{\beta}^2$ if $i = m, j = n$; 0 otherwise.

Generalized Linear Models

Textbook 1

- An othognoal matrix $C_{k \times k}$ has the property $C'C = CC' = I$, i.e. $C' = C^{-1}$. The eigenvalues of $A_{k \times k}$ are the same as $C'AC$.
- P and Q are nonsingular, then $rank(AQ) = rank(PA) = rank(A)$.
- $A_{n \times n}$, symmetric, then $x_i'x_j = 0$ for $i \neq j$. $P_{n \times n}$ nonsingular, then $Tr(P^{-1}AP) = Tr(A)$.
- $A_{n \times n}$, symmetric, then A can be factorized as $A = P\Lambda P^{-1}$, where $\Lambda_{ii} = \lambda_i$, P is an orthogonal matrix, i.e. $PP' = I$.
- $A_{n \times n}$, symmetric, idempotent, then $r(A) = tr(A) = r(P'AP) = tr(P'AP)$.

- $z = a'Y$, $\frac{\partial z}{\partial Y} = a$; $z = Y'Y$, $\frac{\partial z}{\partial Y} = 2Y$; $z = Y'AY$, $\frac{\partial z}{\partial Y} = AY + A'Y$.
- $E(Y) = \mu$, $E(a'Y) = a'E(Y) = a'\mu$; $V(Y) = V$, $V(a'Y) = a'V(Y)a$, $V(AY) = AV(Y)A'$.
- $E(Y'AY) = tr(AV) + \mu'A\mu$.
- If $Y_{k \times 1} \sim N(\mu, I)$, then $Y'Y \sim \chi_{k,\lambda}^2 = \frac{1}{2}(\mu'\mu)$.
- $Y_{n \times 1} \sim N(\mu, I)$, $A = A'$, then $Y'AY \sim \chi_{k,\lambda}^2$ with $k = r(A)$, $\lambda = \frac{1}{2}(\mu'A\mu)$ iff $A = A^2$.
- $Y_{n \times 1} \sim N(\mu, \sigma^2 I)$, $A = A'$, then $Y'AY \sim \chi_{k,\lambda}^2$ with $k = r(A)$, $\lambda = \frac{1}{2\sigma^2}(\mu'A\mu)$ iff $A = A^2$.
- $Y_{n \times 1} \sim N(\mu, V)$, $A = A'$, then $Y'AY \sim \chi_{k,\lambda}^2$ with $k = r(AV) = r(A)$, $\lambda = \frac{1}{2}(\mu'A\mu)$ iff $AV = (AV)^2$.
- $Y_{n \times 1} \sim N(\mu, V)$, then $Y'V^{-1}Y \sim \chi_{k,\lambda'}^2$ with $k = n$, $\lambda = \frac{1}{2}(\mu'V^{-1}\mu)$.
- $Y_{n \times 1} \sim N(\mu, V)$, then AY and BY are independent iff $AVB' = 0$.
- $Y_{n \times 1} \sim N(\mu, V)$, $A_{n \times n} = A'$, $B_{m \times n}$, then $Y'AY$ and BY are independent iff $BVA = 0$.
- $Y_{n \times 1} \sim N(\mu, V)$, $A_{n \times n} = A'$, $B_{n \times n} = B'$, then $Y'AY$ and $Y'BY$ are independent iff $AVB = 0$.
- $B = (X'X)^{-1}XY$, $\hat{Y} = XB = X(X'X)^{-1}XY = HY$, $E(B) = \beta$, $var(B) = \sigma^2(X'X)^{-1}$, $E(\hat{Y}) = X\beta$, $var(\hat{Y}) = \sigma^2H$.
- $SSE = Y'(I - H)Y$ with $df = n - p$, $SSR = Y'(H - \frac{1}{n}J)Y$ with $df = p - 1$, $SST = Y'(I - \frac{1}{n}J)Y$ with $df = n - 1$.
- If $Y = X\beta + \epsilon$, $\epsilon \sim N(0, \sigma^2 I)$, then $B = (X'X)^{-1}XY \sim N(\beta, \sigma^2(X'X)^{-1})$.
- $\frac{(n-p)s^2}{\sigma^2} = \frac{(n-p)MSE}{\sigma^2} = \frac{SSE}{\sigma^2} = \frac{1}{\sigma^2}Y'(I - H)Y \sim \chi_{n-p}^2$.
- B and $\frac{SSE}{\sigma^2}$ are independent.
- $\frac{b_j - \beta_j}{\sqrt{var(b_j)}} = \frac{b_j - \beta_j}{\sigma\sqrt{c_{jj}}} \sim N(0, 1)$, c_{jj} is j th diag entry of $(X'X)^{-1}$.
- $\frac{(b_j - \beta_j) / (\sigma\sqrt{c_{jj}})}{\sqrt{\frac{SSE}{\sigma^2} / (n - p)}} \sim t_{n-p} \Rightarrow b_j \pm t_{n-p}\sqrt{MSEc_{xx}}$
- $LB \sim N(L\beta, \sigma^2L(X'X)^{-1}L')$.

- Let $M = (LB)'\left(\sigma^2\left(L(X'X)^{-1}L'\right)^{-1}\right)^{-1}(LB) \sim \chi_{r,\lambda}^2$, where $\lambda = \frac{1}{2\sigma^2}(LB)'\left(L(X'X)^{-1}L'\right)^{-1}(LB)$.

- $E(M) = r\sigma^2 + (L\beta)'\left(L(X'X)^{-1}L'\right)^{-1}(L\beta)$

- $F^* = \frac{(Lb)'\left(L(X'X)^{-1}L'\right)^{-1}(Lb)/r}{\frac{SSE}{\sigma^2}/(n-p)} = \frac{MSQ}{MSE} \sim F_{r,n-p}$ under $H_0 : L\beta = 0$.

- $\frac{SSR}{\sigma^2} \sim \chi_{p,\lambda'}^2$ where $\lambda = \frac{1}{2\sigma^2}\beta'(X'X)\beta$.

- $MSQ(L\beta) = (LB)'\left(\sigma^2\left(L(X'X)^{-1}L'\right)^{-1}\right)^{-1}(LB) = \frac{SSR}{\sigma^2}$.

- $A = X(X'X)^{-1}X' - X_2(X_2'X_2)^{-1}X_2'$ is idempotent; $r(A) = r$.

Textbook 2

- standardized residual $r_i = \frac{y_i - \hat{\mu}_i}{\hat{\sigma}}$.
- $f(y; \theta) = exp\{a(y)b(\theta) + c(\theta) + d(y)\}$
- glm = exp family + link func (mono + diff)
- $E[a(y)] = -c'(\theta)/b(\theta)$
- $var[a(y)] = \frac{b''(\theta)c'(\theta) - c''(\theta)b'(\theta)}{[b'(\theta)]^3}$
- score info: $U = \frac{\partial l(\theta; y)}{\partial \theta} = a(y)b'(\theta) + c'(\theta)$
- $E(U) = 0$; $J = var(U) = -E(U') = \frac{b''(\theta)c'(\theta)}{b'(\theta)} - c''(\theta)$
- $\frac{U - 0}{\sqrt{J}} \sim N(0, 1)$, $U'J^{-1}U = \frac{U^2}{J} \sim \chi^2(1)$
- wald statistic $(b - \beta)'J(b)(b - \beta) \sim \chi^2(p)$
- $\lambda = \frac{L(b_{max}; y)}{L(b; y)}$, b_{max}/b , - MLE for saturated/reduced model
- $D = 2[l(b_{max}; y) - l(b; y)]$
- $AIC = -2l(\hat{\tau}; y) + 2p$; $BIC = -2l(\hat{\tau}; y) + 2p \times (\text{\#of obs})$