

χ^2 test:

H_0 : Observed distribution is compatible with a given distribution. [Multinomial]

O : Observed counts.

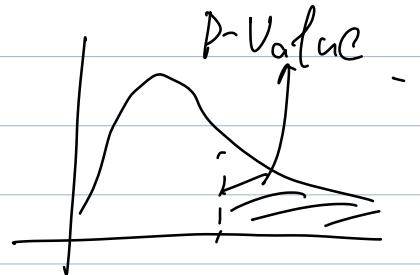
E : Expected counts.

K : Number of categories

n : Total number of outcomes.

$$\chi^2 = \sum_{i=1}^K (O_i - E_i)^2 / E_i$$

Under H_0 : $\chi^2 \sim \chi^2_{k-1}$



One Sample t-test

$H_0: \mu = \mu_0$, $H_1: \mu \neq \mu_0$

($\mu > \mu_0$)

($\mu < \mu_0$)

$$T = \frac{\bar{x} - \mu_0}{s/\sqrt{n}} \quad \text{with } s = \sqrt{s^2} \quad \begin{matrix} \text{Sample} \\ \text{Variance} \end{matrix}$$

\bar{x} : Sample mean

$T \sim \text{Student } t \text{ with } (n-1) \text{ df.}$

Test for 2 means:

i) Independent samples:

a) Unequal variances

Assumptions:

- Both are simple random samples
- $n_1 > 30$, $n_2 > 30$ or both samples are from Normal distribution.

$$H_0: \mu_1 - \mu_2 = \mu_0$$

$$H_1: \mu_1 - \mu_2 \neq \mu_0$$

Generally, $\mu_0 = 0$

$$T = \frac{(\bar{x}_1 - \bar{x}_2) - \mu_0}{\sqrt{\frac{s_1^2}{n_1} + \frac{s_2^2}{n_2}}}$$

$$df = \frac{(A+B)^2}{\frac{A^2}{n_1-1} + \frac{B^2}{n_2-1}}$$

$$A = \frac{s_1^2}{n_1}$$

$$B = \frac{s_2^2}{n_2}$$

b) $\sigma_1^2 = \sigma_2^2$

$$T = \frac{(\bar{x}_1 - \bar{x}_2) - \mu_0}{\sqrt{\frac{s_p^2}{n_1} + \frac{s_p^2}{n_2}}}$$

$$S_p = \frac{(n_1-1)s_1^2 + (n_2-1)s_2^2}{(n_1-1) + (n_2-1)} \quad T \sim t(n_1+n_2-2)$$

2) Matched Pairs

Assumptions:

- Sample of match pairs

- n : number of pairs. $n \geq 30$ or differences come from a Normal distribution.

$$H_0: \mu_d = 0 \quad H_1: \mu_d \neq 0$$

$$T = \frac{\bar{d}}{\frac{s_d}{\sqrt{n}}} \quad t(n-1) df$$

Density Estimation:

$$\hat{f}(x) = \frac{1}{nb} \sum_{j=1}^n K\left(\frac{x-x_j}{b}\right) ?$$

Default: $\hat{b} = 0.9 \min(\hat{\sigma}, R/1.34) n^{-1/5}$

Test for proportion:

$$H_0: \hat{p} = p$$

Page: Data explanation, challenges, references.

Oct. 14.

Simple linear Regression:

$$y_i = \beta_0 + \beta_1 x_i + \varepsilon_i \quad i=1, \dots, n$$

ε_i i.i.d. $N(0, \sigma^2)$

Least Square Estimation \rightarrow Don't need any assumption

Finding $\hat{\beta}_0$ and $\hat{\beta}_1$ such that the minimize:

$$\sum_{i=1}^n (y_i - \beta_0 - \beta_1 x_i)^2 = f(\beta_0, \beta_1)$$

$$\frac{\partial f(\beta_0, \beta_1)}{\partial \beta_0} = \sum_{i=1}^n -2(y_i - \beta_0 - \beta_1 x_i) = 0 \quad \textcircled{1}$$

$$\frac{\partial f(\beta_0, \beta_1)}{\partial \beta_1} = \sum_{i=1}^n -2x_i(y_i - \beta_0 - \beta_1 x_i) = 0 \quad \textcircled{D}$$

Using \textcircled{D}: $\sum_{i=1}^n y_i - n\beta_0 - \beta_1 \sum_{i=1}^n x_i = 0$

$$\Rightarrow \bar{y} - \beta_0 - \beta_1 \bar{x} = 0$$

$$\Rightarrow \hat{y} = \beta_0 + \beta_1 \bar{x} \quad \star$$

Using \textcircled{D}:

$$\sum_{i=1}^n x_i y_i - \sum_{i=1}^n x_i^2 \beta_1 - \sum_{i=1}^n x_i \beta_0 = 0 \quad \star \star$$

$$\Rightarrow \sum_{i=1}^n x_i y_i - \sum_{i=1}^n x_i^2 \cdot \beta_1 - \sum_{i=1}^n x_i \cdot \beta_0 = 0$$

By solving \star and $\star \star$:

$$\hat{\beta}_1 = \frac{\sum_{i=1}^n x_i y_i - n \bar{x} \bar{y}}{\sum_{i=1}^n x_i^2 - n \bar{x}^2} \quad \beta_0 = \bar{y} - \bar{x} \hat{\beta}_1$$

$$\hat{\beta}_1 = \frac{SS_{XY}}{SS_{XX}} \Rightarrow \begin{cases} SS_{XX} = \sum_{i=1}^n (x_i - \bar{x})^2 \\ SS_{XY} = \sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y}) \\ SS_{YY} = \sum_{i=1}^n (y_i - \bar{y})^2 \end{cases}$$

$$\hat{y}_i = \hat{\beta}_0 + \hat{\beta}_1 x_i \quad \varepsilon_i = y_i - \hat{\beta}_0 - \hat{\beta}_1 x_i$$

$$SSE = \sum_{i=1}^n (\hat{y}_i - y_i)^2$$

$$MSZ = SSE / (n-2) = \hat{\sigma}^2$$

$$E(\hat{\sigma}^2) = \sigma^2$$

Hypothesis testing.

$$H_0: \beta_1 = 0$$

$$H_1: \beta_1 \neq 0$$

$$T = \frac{\hat{\beta}_1 - 0}{\sqrt{\frac{s^2}{SS_{xx}}}}$$

$T \sim \text{student-t}$ with $n-2$ df.

Confidence interval for β_1 : at α level of significance:

The $(1-\alpha)\%$ CI (L, U) can be obtained with:

$$L = \hat{\beta}_1 - t_{\alpha/2, n-2} \frac{s^2}{\sqrt{SS_{xx}}}$$

$$U = \hat{\beta}_1 + t_{\alpha/2, n-2} \frac{s^2}{\sqrt{SS_{xx}}} \quad \hat{s}^2 = \frac{SS_{\epsilon}}{n-r}$$

$$R^2 = \frac{\sum_{i=1}^n (\hat{y}_i - \bar{y})^2}{\sum_{i=1}^n (y_i - \bar{y})^2} \rightarrow SSR$$

Confidence Interval for the mean response:

Interval for x^* :

$$\hat{y}|x^* = \hat{\beta}_0 + \hat{\beta}_1 x^*$$

$$M_{\hat{y}|x^*} = \hat{\beta}_0 + \hat{\beta}_1 x^*$$

Confidence interval: $(1-\alpha) 100\%$

$$\hat{y}|x^* \pm t_{\alpha/2, n-2} \cdot se \hat{y}|x^*$$

$$se \hat{y}|x^* = \hat{s} \sqrt{\frac{1}{n} + \frac{(x^* - \bar{x})^2}{\sum_{i=1}^n (x_i - \bar{x})^2}}$$

Prediction interval: prediction for a new observation at

$$y^{\text{new}} = \hat{\beta}_0 + \hat{\beta}_1 x^* + \varepsilon^{\text{new}}$$

Point estimation is the same,

$$\hat{y}_{\text{new}} = \hat{\beta}_0 + \hat{\beta}_1 x^*$$

But the interval is different:

$$\hat{y}_{\text{new}} \pm t_{\alpha/2, n-2} \cdot \text{Se } \hat{y}_{\text{new}}$$

$$\text{Se } \hat{y}_{\text{new}} = \sqrt{1 + \frac{1}{n} + \frac{(x^* - \bar{x})^2}{\sum_{i=1}^n (x_i - \bar{x})^2}}$$

H_0 for F test:

$$\beta_1 = \beta_2 = \dots = 0$$

H_1 for F test:

At least 1 of β_i is not 0.

Multiple Regression:

$$y_i = \beta_0 + \beta_1 x_{1i} + \beta_2 x_{2i} + \dots + \beta_p x_{pi} + \varepsilon_i$$

$$\varepsilon_i \stackrel{\text{iid}}{\sim} N(0, \sigma^2)$$

$$y = \begin{pmatrix} y_1 \\ \vdots \\ y_n \end{pmatrix} \quad \beta = \begin{pmatrix} \beta_0 \\ \vdots \\ \beta_p \end{pmatrix} \quad X = \begin{pmatrix} 1 & x_{11} & x_{21} & \dots & x_{p1} \\ \vdots & \vdots & \vdots & & \vdots \\ \vdots & \vdots & \vdots & & \vdots \\ 1 & x_{1n} & x_{2n} & \dots & x_{pn} \end{pmatrix}$$

$$\dim(X) = n \times (p+1) \quad \dim(\beta) = (p+1) \times 1$$

$$\dim(y) = n \times 1 \quad \dim(\varepsilon) = n \times 1 \quad \varepsilon \sim N(0, \sigma^2 I_n)$$

$$y = X\beta + \varepsilon \quad \varepsilon \sim N(0, \sigma^2 I_n)$$

1) LSZ : Find $\hat{\beta}$ that minimizes:

$$S(\beta) = \| y - X\beta \| = (y - X\beta)^T (y - X\beta)$$

2) MLZ:

$$y \sim N(x\beta, \sigma^2 I_n)$$

$$L(\beta, \sigma^2) = \prod_{i=1}^n \frac{1}{\sqrt{2\pi\sigma^2}} \exp\left(-\frac{(y_i - \beta_0 - \beta_1 x_{1,i} - \dots - \beta_p x_{p,i})^2}{2\sigma^2}\right)$$

or, equivalently:

$$L(\beta, \sigma^2) = \frac{1}{(2\pi\sigma^2)^{\frac{n}{2}}} \exp\left\{-\frac{(y - x\beta)^T(y - x\beta)}{2\sigma^2}\right\}$$

Maximizing $L(\beta, \sigma^2)$ is the same as maximizing:

$$\log L(\beta, \sigma^2) = -\frac{n}{2} \log(2\pi\sigma^2) - \frac{(y - x\beta)^T(y - x\beta)}{2\sigma^2}$$

For method 1:

LSZ: we find $\frac{\partial S(\beta)}{\partial \beta_0} = 0, \dots, \frac{\partial S(\beta)}{\partial \beta_p} = 0$

If X is a full rank matrix with $n > p+1$, and $\text{rank } X = p+1$. Then the solution is unique and given by:

$$(X^T X)^{-1} X^T Y.$$

For method 2:

MLZ: we find: $\frac{\partial l(\beta, \sigma^2)}{\partial \beta} = 0 \quad \boxed{\frac{\frac{\partial l(\beta, \sigma^2)}{\partial \sigma^2}}{\sigma^2} = 0}$ if σ^2 known

If σ^2 known: The MLZ with X a full rank matrix and $n > p+1$ is also $(X^T X)^{-1} X^T Y$.

But we also have σ^2 :

so. we also find $\frac{\partial l(\beta, \sigma^2)}{\partial \sigma^2} = 0$. not unbiased

This gives us:

$$\hat{\sigma}_{MLZ}^2 = \frac{(y - \hat{x}\hat{\beta})^T(y - \hat{x}\hat{\beta})}{n} = \frac{\hat{\epsilon}^T \hat{\epsilon}}{n}$$

So, we use:

$$\hat{\sigma}^2 = \frac{\hat{\epsilon}' \hat{\epsilon}}{n-(P+1)}$$

Properties of $\hat{\beta}$:

$$\begin{aligned} E(\hat{\beta}) &= E((X'X)^{-1} X' Y) \\ &= (X'X)^{-1} X' X \beta \\ &= \beta \end{aligned}$$

$$\begin{aligned} \text{Var}(\hat{\beta}) &= (X'X)^{-1} X' \text{Var}(Y) X (X'X)^{-1} \\ &= (X'X)^{-1} X' \sigma^2 I_n X (X'X)^{-1} \\ &= \sigma^2 (X'X)^{-1} \end{aligned}$$

$$\hat{\beta} \sim N(\beta, \sigma^2 (X'X)^{-1}) \quad \hat{\beta}_i \sim N(\beta_i, \sigma^2 (X'X)^{-1}_{ii})$$

F-test for Nested Model:

$$M_1: Y_i = \beta_0 + \beta_1 X_{1,i} + \dots + \beta_p X_{p,i} + \epsilon_i$$

$$M_2: Y_i = \beta_0 + \beta_1 X_{1,i} + \dots + \beta_p X_{p,i} + \beta_{p+1} X_{p+1,i} + \dots + \beta_{p+q} X_{p+q,i}$$

$$H_0: \beta_{p+1} = \beta_{p+2} = \dots = \beta_{p+q} = 0 \quad H_a: \beta_{p+1} = \dots = \beta_{p+q} \neq 0 + \epsilon_i$$

H_a: At least 1 of $\beta_{p+1}, \dots, \beta_{p+q}$ is not 0.

Fit M_1 and M_2 and get:

$$\underline{SSE_1} \quad \text{and} \quad \underline{SSE_2}$$

$$\begin{matrix} X_1 \\ \hat{\beta}^{(1)} \end{matrix} \quad \begin{matrix} X_2 \\ \hat{\beta}^{(2)} \end{matrix}$$

$$SSE_1 = (Y - X_1 \hat{\beta}^{(1)})^T (Y - X_1 \hat{\beta}^{(1)})$$

$$SSE_2 = (Y - X_1 \hat{\beta}^{(1)})^T (Y - X_2 \hat{\beta}^{(2)})$$

Note that: $SSE_1 > SSE_2$.

$$F = \frac{(SSE_1 - SSE_2)/q_0}{SSE_2/(n-(p+q_0+1))}$$

Under H_0 : $F \sim F$ -distribution with d.f.
of q_0 in the numerator, $n-(p+q_0+1)$ in the denominator.

What about comparing models that are not nested:

- AIC : Akaike's information

If you have a model M : \nearrow # of parameters in the model
 $AIC(M) = -2 \log L(M) + 2p(M)$
 \hookrightarrow likelihood of model M .
evaluated at the MLZ.

- BIC: Bayesian Information Criterion.

$BIC(M) = -2 \log(L(M)) + \log(n) \cdot p(M)$
 \hookrightarrow # of observations.

If you have models M_1, M_2, \dots, M_k , choose the one
that minimizes AIC or BIC

Adjusted R^2 :

$$R^2_{adj} = 1 - \left[\frac{SSE}{n-(p+1)} \right] / \left[\frac{\sum (y_i - \bar{y})^2}{n-1} \right]$$

$$SST = \sum (y_i - \bar{y})^2$$

Oct. 21: One-Way ANOVA.

$$Y_{ij} = \mu + \alpha_i + \epsilon_{ij}, \quad \epsilon_{ij} \stackrel{i.i.d.}{\sim} N(0, \sigma^2)$$

$i = 1 : a, \quad j = 1 : n_i$

↳ If n_j are same, the model is called balanced model.

$$n_i = n \text{ for all } i$$

1 factor with a levels. can we write this model as

$$y = x\beta + \epsilon \quad \epsilon \sim N(0, \sigma^2 I).$$

$$y = \begin{pmatrix} y_{11} \\ \vdots \\ y_{1n_1} \\ y_{21} \\ \vdots \\ y_{2n_2} \\ \vdots \\ y_{a1} \\ \vdots \\ y_{an_a} \end{pmatrix} \quad \beta = \begin{pmatrix} \mu \\ \alpha_1 \\ \vdots \\ \alpha_a \end{pmatrix} \quad \epsilon = \begin{pmatrix} \epsilon_{11} \\ \vdots \\ \epsilon_{1n_1} \\ \epsilon_{21} \\ \vdots \\ \epsilon_{2n_2} \\ \vdots \\ \epsilon_{a1} \\ \vdots \\ \epsilon_{an_a} \end{pmatrix}$$

$$X = \begin{pmatrix} 1 & 1 & 0 & \dots & 0 \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ 1 & 0 & 1 & 0 & \dots & 0 \\ \vdots & \vdots & \vdots & \vdots & \vdots & \vdots \\ 1 & 0 & 1 & 0 & \dots & 0 \\ \vdots & \vdots & \vdots & \vdots & \vdots & \vdots \\ 1 & 0 & \dots & \dots & 0 & 1 \\ \vdots & \vdots & \vdots & \vdots & \vdots & \vdots \\ 1 & 0 & \dots & \dots & 0 & 1 \end{pmatrix} \rightarrow \begin{array}{l} \text{Group 1} \\ \text{Group 2} \\ \vdots \\ \text{Group} \end{array}$$

X is not full rank.

↳ If X is not full rank, $M\beta$ is not unique
Add restrictions can solve this problem, or look at different parameterisations of the model.

(a) Consider the paramet:

$$y_{ij} = \mu_i + \varepsilon_{ij} \quad \begin{cases} i = 1:n \\ j = 1:n_i \\ \varepsilon_{ij} \text{ i.i.d } N(0, \sigma^2) \end{cases}$$

$$\beta = \begin{pmatrix} \mu_1 \\ \vdots \\ \mu_a \end{pmatrix} \quad \text{The resulting } X \text{ is a full rank.}$$

$X = \begin{pmatrix} 1 & 0 & \dots & 0 \\ \vdots & \ddots & \ddots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ 0 & \dots & 0 & 1 \end{pmatrix} \xrightarrow{\text{Group 1}}$

$\xrightarrow{\text{Group A.}}$

$$\hat{\beta} = (X^T X)^{-1} X^T Y$$

$$X^T X = \begin{pmatrix} n_1 & & & 0 \\ & n_2 & & \vdots \\ & & \ddots & 0 \\ 0 & \vdots & \ddots & n_a \end{pmatrix}$$

$$(X^T X)^{-1} = \begin{pmatrix} 1/n_1 & & & 0 \\ & 1/n_2 & & \vdots \\ & & \ddots & 0 \\ 0 & \vdots & \ddots & 1/n_a \end{pmatrix}$$

$$X^T Y = \begin{pmatrix} \sum_{j=1}^{n_1} y_{1,j} \\ \vdots \\ \sum_{j=1}^{n_a} y_{a,j} \end{pmatrix}$$

$$\hat{\beta} = \begin{pmatrix} \bar{y}_{1,.} \\ \vdots \\ \bar{y}_{a,.} \end{pmatrix}$$

$$\begin{aligned} \hat{\mu}_i &= \bar{y}_{i,.} \\ &= \sum_{j=1}^{n_i} y_{i,j} \end{aligned}$$

Hypothesis Testing:

$$H_0: \mu_1 = \mu_2 = \dots = \mu_a \quad \text{vs} \quad H_1: \text{at least one } \mu_i \text{ is different from the rest}$$

t -test If we fail to reject, then we are

$$Y_{ij} = \mu + \varepsilon_{ij} \quad \varepsilon_{ij} \sim i.i.d. N(0, \sigma^2)$$

$\hat{\mu} = \bar{y}_{..}$ t -test between $n\bar{M}$ and \bar{M} .

b) Add one restriction to this model.

$$Y_{ij} = \mu + \alpha_i + \varepsilon_{ij} \quad \varepsilon_{ij} \sim i.i.d. N(0, \sigma^2)$$

possible restrictions:

i) $\alpha_{i_0} = 0$ for one i_0

$\xrightarrow{n-1}$ $\alpha_1, \dots, \alpha_n$

ii) $\sum_{i=1}^a \alpha_i = 0$

$\xrightarrow{k-1}$ $\alpha_1, \dots, \alpha_k$

Test for 2 groups:

$$H_0: \mu_i = \mu_j \quad v.s. H_1: \mu_i \neq \mu_j$$

\downarrow $b-r$ var. est.

$$T = \frac{\bar{y}_i - \bar{y}_j}{\sqrt{MS\bar{z}(\bar{y}_i + \bar{y}_j)}} \quad df = N-a.$$

$\hookrightarrow MS\bar{z}$ for linear regression.

Threshold is:

$$T_{\alpha/2, N-a} \cdot \sqrt{MS\bar{z}(\bar{y}_i + \bar{y}_j)}.$$

$$\alpha = P(\text{Type I error})$$

$$= P(\text{Reject } H_0 \mid T > T_{\alpha/2, N-a})$$

$$= P(\text{Incorrect significant result})$$

$$Pr(\text{At least 1 significant is incorrect})$$

$$= 1 - P(\text{No significant result}) = 1 - (1-\alpha)^k$$

If $k = 20$, $\alpha = 0.05$

$$1 - (1-\alpha)^k \approx 0.64$$

Oct. 2f.

Randomized Block Models:

more than 2-dimensional paired t-test.

Multiple Block Factors:

Example: Latin square design.

| | Z_1 | Z_2 | Z_3 | Z_4 |
|-------|-------|-------|-------|-------|
| R_1 | A | B | C | D |
| R_2 | B | C | D | A |
| R_3 | C | D | A | B |
| R_4 | D | A | B | C |

$$Y_{ijk} = \mu + \alpha_i + \beta_j + \gamma_k + \varepsilon_{ijk}$$

↑ ↑ ↓
Block 1 Block 2 Factor

$i: 1: g$ $j: 1: g$ $k: 1: g$.

of parameters: $1 + (g-1) + (g-1) + (g-1) = 3g - 2$

of Observations: g^2

degree of freedom: $g^2 - 3g - 2$

Two-Way ANOVA Model:

Full model with interactions:

$$Y_{i,j,k} = \mu + \alpha_i + \beta_j + \gamma_{i,j} + \varepsilon_{i,j,k} \quad \varepsilon_{i,j,k} \stackrel{i.i.d.}{\sim} N(0, \sigma^2)$$

$i = 1:a \quad j = 1:b \quad k = 1:n_{i,j}$

Model with Interactions

$$Y_{i,j,k} = \mu + \alpha_i + \beta_j + \varepsilon_{i,j,k} \quad (\text{no interactions})$$

$$\frac{y}{N} = \frac{x\beta + \varepsilon}{N} \quad \hat{\beta} = (x^T x)^{-1} x^T \beta$$

$$\left(\begin{array}{c|c} y_{1,1,1} & \\ \vdots & \\ y_{1,1,n_{1,1}} & \\ \vdots & \\ y_{a,b,1} & \\ \vdots & \\ y_{a,b,n_{a,b}} & \end{array} \right) \quad N = \sum_{i=1}^a \sum_{j=1}^b n_{i,j}$$

These models are not full rank unless we make some restrictions.

usually: i. $\alpha_{i_0} = 0$ for some i_0 or

$$\sum_{i=1}^a \alpha_i = 0,$$

ii. $\beta_{j_0} = 0$ for some j_0 or

$$\sum_{j=1}^b \beta_j = 0$$

iii. $\gamma_{i_0 j_0} = 0$ for one i_0 and all j and $\gamma_{i_0 j_0} = 0$ for one j_0 and all i .

or

$$\sum_{i=1}^a \gamma_{i,j} = 0, \text{ for all } j \text{ and}$$

$$\sum_{j=1}^b \gamma_{i,j} = 0, \text{ for all } i$$

Estimates under restrictions are unique. For example:
in the balance case with $\eta_{ij} = \eta$ with 3 zero

sum constraints are:

$$M = \bar{y}_{...} \quad \hat{\alpha}_i = \bar{y}_{i...} - \bar{y}_{...} \quad \hat{\beta}_j = \bar{y}_{.j...} - \bar{y}_{...}$$

and $\hat{\gamma}_{ij} = \bar{y}_{ij...} - \bar{y}_{i...} - \bar{y}_{.j...} + \bar{y}_{...}$

Also consider: $y_{ijk} = M_{ij} + \varepsilon_{ijk}$
 $\hat{M}_{ij} = \bar{y}_{ij...}$

Oct. 28.

$$y_{ijk} = M + \alpha_i + \beta_j + \gamma_{ij} + \varepsilon_{ijk} \quad \left. \begin{array}{l} H_0: \gamma_{ij} = 0 \text{ for all } i, j \\ H_a: \text{At least } 1 \gamma_{ij} \neq 0 \end{array} \right\}$$

$$\left. \begin{array}{l} H_0: \alpha_i = 0 \text{ for all } i \\ H_a: \text{At least } 1 \alpha_i \neq 0 \end{array} \right\}$$

$$\left. \begin{array}{l} H_0: \beta_j = 0 \text{ for all } j \\ H_a: \text{At least } 1 \beta_j \neq 0 \end{array} \right\}$$

$$y_{ijk} = M_{ij} + \varepsilon_{ijk}.$$

$$H_0: \text{All } M_{ij} =$$

Restrictions: $\alpha_i = 0, \beta_j = 0$

$$\gamma_{ij} = 0, \gamma_{i..} = 0.$$

$\downarrow \quad \downarrow$
 for all $j \quad$ for all i .

$$\text{eq: } \hat{\mu}_{I,B} = \hat{\mu} + \hat{\alpha}_I + \hat{\beta}_B + \hat{\gamma}_{I,B}$$

ANOVA Tables:

- One way ANOVA $N: \text{Total # of obs.}$

| Source | DF | SS |
|----------|-------|--|
| Factor A | $a-1$ | $N \cdot \bar{y}_{..}^2 - (SS_A)$ |
| Error | $N-a$ | $\sum_{i=1}^a \sum_{j=1}^{n_i} (y_{ij} - \bar{y}_{ij})^2 (SS_E)$ |

- Two way ANOVA (Balanced Case)

$$y_{ijk} = \mu + \alpha_i + \beta_j + \gamma_{ij} + \epsilon_{ijk}$$

| Source | DF | SS | MS |
|-------------|--------------|--------|-------------------|
| Factor A | $a-1$ | SS_A | $SS_A/(a-1)$ |
| Factor B | $b-1$ | SS_B | $SS_B/(b-1)$ |
| Interaction | $(a-1)(b-1)$ | SS_I | $SS_I/(a-1)(b-1)$ |
| Error | $ab(n-1)$ | SS_E | $SS_E/ab(n-1)$ |

\hookrightarrow # of each group.

$$SS_A = b n \sum_i (\bar{y}_{i..} - \bar{y}_{...})^2$$

$$SS_B = a n \sum_j (\bar{y}_{.j} - \bar{y}_{...})^2$$

$$SS_I = n \sum_i \sum_j (\bar{y}_{ij} - \bar{y}_{i..} - \bar{y}_{.j} + \bar{y}_{...})^2$$

$$SS_E = \sum_i \sum_j \sum_k (y_{ijk} - \bar{y}_{...})^2$$

$$\textcircled{1} Y_{i,j} = \alpha_i + \beta_j X_{i,j} + \varepsilon_{i,j} \quad \varepsilon_{i,j} \stackrel{iid}{\sim} N(0, \sigma^2)$$

$\textcircled{1}$ $\textcircled{2} Y_{i,j} = \alpha + \delta_i + (\beta + \gamma_j) X_{i,j} + \varepsilon_{i,j} \quad \varepsilon_{i,j} \stackrel{iid}{\sim} N(0, \sigma^2)$

$\hat{\alpha}$ = Intercept

$\textcircled{2}$ $\hat{\alpha}_1 = \text{Intercept} + \text{Indicator.}$

Also:

$\hat{\alpha} = \text{Intercept} \quad \hat{\delta}_1 = 0 \quad \hat{\delta}_2 = \text{Indicator.}$

$\textcircled{1}: \hat{\beta}_1 = \text{Continuous}$

$\hat{\beta}_2 = \text{Continuous} + \text{Interaction.}$

Also:

$\hat{\beta} = \text{Continuous.} \quad \hat{\delta}_1 = 0 \quad \hat{\delta}_2 = \text{Interaction.}$

$$\begin{pmatrix} y_{1,1} \\ \vdots \\ y_{1,n_1} \\ y_{2,1} \\ \vdots \\ y_{2,n_2} \end{pmatrix} = \begin{pmatrix} 1 & 1 & 0 & x_{1,1} \\ \vdots & \vdots & \vdots & \vdots \\ 1 & 0 & x_{1,n_1} \\ 0 & 1 & x_{2,1} \\ \vdots & \vdots & \vdots \\ 0 & 1 & x_{2,n_2} \end{pmatrix} \begin{pmatrix} \alpha \\ \delta_1 \\ \delta_2 \\ \beta \end{pmatrix}$$

PS: If $\delta_1 = 0$, delete 2nd column.

↑
For $Y_{i,j} = (\alpha + \delta_i) + \beta X_{i,j} + \varepsilon_{i,j}$

Nov. 4th.

Logistic Regression.

y_1, \dots, y_n binary responses,

x_1, \dots, x_n continuous explanatory variable.

$$P(Y=y_i | X=x_i) = \theta(x_i)$$

$\hookrightarrow \theta(\cdot)$ is a function.

$$\text{logit}(\theta(x_i)) = \log\left(\frac{\theta(x_i)}{1-\theta(x_i)}\right) = \alpha + \beta x_i$$

$$\text{or: } \frac{\theta(x_i)}{1-\theta(x_i)} = e^{\alpha + \beta x_i}$$

\hookrightarrow odds: exponential function of x_i .

Interpretation:

$$\frac{\frac{\theta(x_i+1)}{1-\theta(x_i+1)}}{\frac{\theta(x_i)}{1-\theta(x_i)}} = \frac{e^{\alpha + \beta(x+1)}}{e^{\alpha + \beta x}} = e^\beta$$

1 unit increase in x , odds ratio will increase e^β .

We can have multiple explanatory variables:

$$\text{logit}(\theta_i) = \log\left(\frac{\theta_i}{1-\theta_i}\right) = \alpha + \sum_{j=1}^p \beta_j x_{ij}$$

What if we have a binary explanatory variable?

$$e^\beta = \frac{P(Y=1 | x=1)}{1-P(Y=1 | x=1)}$$

$$\frac{P(Y=1 | x=0)}{1-P(Y=1 | x=0)}$$

Once you fit the model, you can look at:

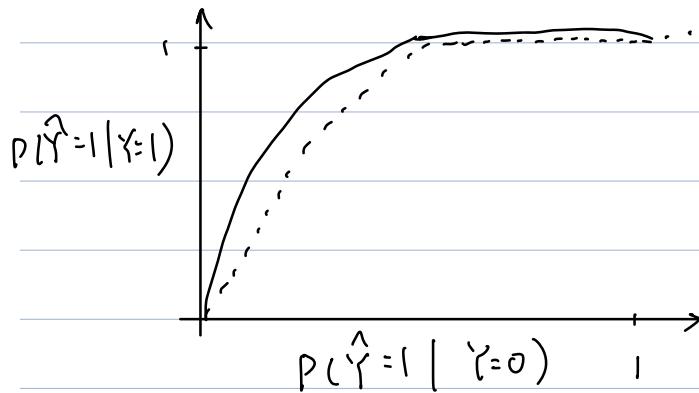
sensitivity:

$$P(Y=\hat{Y}=1 | Y=1)$$

specificity:

$$P(Y=\hat{Y}=0 | Y=0)$$

ROC: Receiving Operating characteristic Curve.



Estimation: Maximum Likelihood.

y_i ind. Bernoulli (θ_i)

Likelihood:

$$\mathcal{L}(\alpha, \beta) = \prod_{i=1}^n \theta_i^{y_i} (1-\theta_i)^{1-y_i}$$

$$\text{logit}(\theta_i) = \log\left(\frac{\theta_i}{1-\theta_i}\right) = \alpha + \beta x_i$$

$$\Rightarrow \frac{\theta_i}{1-\theta_i} = e^{\alpha + \beta x_i} \Rightarrow \theta_i = \frac{e^{\alpha + \beta x_i}}{1 + e^{\alpha + \beta x_i}}$$

$$= \ell(\alpha, \beta) = \sum_{i=1}^n y_i \log(\theta_i) + (1-y_i) \log(1-\theta_i)$$

$$= \sum_{i=1}^n y_i (\alpha + \beta x_i) + \sum_{i=1}^n \frac{1}{e^{\alpha + \beta x_i}}$$