

Chatterjee

Advanced Topics in Stat. (Prof. Shirshendu Chatterjee)

Textbook: Applied Multivariate Analysis

Applied Linear Statistical model by kuther, etc.

Topics

Multiple Linear Regression and associated inference.

Logistic regression, ~~principle~~ principal component analysis,
clustering and classification, some design of experiment.

HW: Biweekly HW will be posted in Blackboard.

Solution of some of them will be posted.

Exams: Two Midterms (40%) ← Tentatively on the first
One project (20%) class of March and April.
Final Exam (40%) Due by the last class.

Makeup Exam (strict policy)

Formula sheets (3 pages for Midterms, 4 pages for Final
and calculators are allowed)

Ch. 7 from J & W.

Multiple linear regression

$y \leftrightarrow$ dependent variable (response)

$z_1 \dots z_k \leftrightarrow$ independent variable (predictors)

Assume Linear relation: $y = \beta_0 + \beta_1 z_1 + \dots + \beta_k z_k + \varepsilon$

Goal: Statistical inference about β_0, \dots, β_k

Prediction of response based on predictor variables.

We observe the response and predictors for n observations.

This is called the training data.

$$y_1 = \beta_0 + \beta_1 z_{11} + \beta_2 z_{12} + \dots + \beta_k z_{1k} + \varepsilon_1$$

$$y_2 = \beta_0 + \beta_1 z_{21} + \beta_2 z_{22} + \dots + \beta_k z_{2k} + \varepsilon_2$$

⋮

$$y_n = \beta_0 + \beta_1 z_{n1} + \beta_2 z_{n2} + \dots + \beta_k z_{nk} + \varepsilon_n$$

✓ check

$$\underline{y} = \begin{pmatrix} y_1 \\ y_2 \\ \vdots \\ y_n \end{pmatrix}, \underline{Z} = \begin{bmatrix} 1 & z_{11} & \dots & z_{1k} \\ 1 & z_{21} & \dots & z_{2k} \\ \vdots & \vdots & \ddots & \vdots \\ 1 & z_{n1} & \dots & z_{nk} \end{bmatrix}, \underline{\varepsilon} = \begin{pmatrix} \varepsilon_1 \\ \vdots \\ \varepsilon_n \end{pmatrix}, \underline{\beta} = \begin{pmatrix} \beta_0 \\ \vdots \\ \beta_k \end{pmatrix}$$

$$\underline{y} = \underline{Z} \underline{\beta} + \underline{\varepsilon}$$

$n \times 1$ $n \times (k+1)$ $(k+1) \times 1$

Assumption on errors

Errors have mean 0 and they are uncorrelated.

Remember that the random part in the model is the errors.

$$\mathbb{E}(\varepsilon_i) = 0 \text{ for all } i.$$

$$\text{cov}(\varepsilon_i, \varepsilon_j) = 0 \text{ for all } i \neq j.$$

$$\mathbb{E}(\varepsilon_i \varepsilon_j) = \mathbb{E}(\varepsilon_i) \mathbb{E}(\varepsilon_j)$$

$$\text{var}(\varepsilon_i) = \text{cov}(\varepsilon_i, \varepsilon_i) = \sigma^2. \quad (\text{homoskedasticity assumption})$$

In this case, we can use Least square method to estimate $\underline{\beta}$.

$$\text{Minimize } \sum_{i=1}^n (y_i - \beta_0 - \beta_1 z_{i1} - \dots - \beta_k z_{ik})^2 \quad (*)$$

$$\mathbb{E}(y_i) = \beta_0 + \beta_1 z_{i1} + \dots + \beta_k z_{ik}$$

The minimizer of (*) is called the LS estimator of $\underline{\beta}$.

Notation $\hat{\underline{\beta}}_{\text{LS}}$

the sum of squared error of differences.

$$= \sum_{i=1}^n (y_i - z_{i*} \hat{\underline{\beta}}_{\text{LS}})^2, \quad z_{i*} \text{ is the } i^{\text{th}} \text{ row.}$$

$$= (\underline{y} - \underline{Z} \hat{\underline{\beta}}_{\text{LS}})' (\underline{y} - \underline{Z} \hat{\underline{\beta}}_{\text{LS}}) = \underline{\varepsilon}' \underline{\varepsilon}$$

$$\frac{\partial}{\partial \beta_j} (\underline{y} - \underline{Z} \hat{\underline{\beta}}_{\text{LS}})' (\underline{y} - \underline{Z} \hat{\underline{\beta}}_{\text{LS}}) = 0.$$

$$2 \underline{Z}' (\underline{y} - \underline{Z} \hat{\underline{\beta}}_{\text{LS}}) = 0 \Rightarrow \underline{Z}' \underline{y} = \underline{Z}' \underline{Z} \hat{\underline{\beta}}_{\text{LS}}$$

$$\begin{aligned} \underline{Z}' \underline{Z} \hat{\underline{\beta}}_{\text{LS}} &= \underline{y}' \\ \underline{Z}' \underline{Z} &= \underline{Z}' \underline{Z} \\ \underline{Z}' \underline{Z} &= (\underline{Z}' \underline{Z})^{-1} \underline{Z}' \underline{Z} \underline{Z} \\ \underline{Z}' \underline{Z} &= \underline{Z}' \underline{Z} \end{aligned}$$

$$\hat{\underline{\beta}}_{\text{LS}} = (\underline{Z}' \underline{Z})^{-1} \underline{Z}' \underline{y}$$

If $\text{rank}(Z) = k+1$, then $\hat{\beta}_{LS} = (Z'Z)^{-1}Z'y$.

\hat{y} = predicted value of y

$$= Z\hat{\beta} = Z(Z'Z)^{-1}Z'y = P_Z y, \text{ where}$$

P_Z is the projection matrix on the column space of Z .
 $Z(Z'Z)^{-1}Z'$

Recall that projection matrices are idempotent,

i.e., $P_Z^2 = P_Z$ Also, $I - P_Z$ is idempotent.

Also, $P_Z(I - P_Z) = 0$: $P_Z \perp (I - P_Z)$

Thus, column space of P_Z and $I - P_Z$ are orthogonal.

$$E(\hat{\beta}_{LS}) = E((Z'Z)^{-1}Z'y)$$

$$= (Z'Z)^{-1}Z'E(y)$$

$$= (Z'Z)^{-1}Z'\hat{\beta} = \hat{\beta}$$

$$\text{cov}(\hat{\beta}) = \text{cov}((Z'Z)^{-1}Z'y) = (Z'Z)^{-1}Z' \frac{\text{cov}(y)}{\text{cov}(y)} Z(Z'Z)^{-1}$$

$$\begin{aligned} &= (Z'Z)^{-1}Z' (\sigma^2 I) Z(Z'Z)^{-1} \\ &= \sigma^2 (Z'Z)^{-1} \end{aligned}$$

Unbiased estimator for σ^2

$$\text{The estimated error } \hat{\varepsilon} = y - \hat{y}$$

$$\begin{aligned} &= y - P_Z y = (I - P_Z)y = \hat{\varepsilon} \\ \text{So, } \sum_{i=1}^n \hat{\varepsilon}_i^2 &= (\hat{\varepsilon}' \hat{\varepsilon}) = y'(I - P_Z)^2 y \quad (I - P_Z)^2 = (I - P_Z) \\ &= y'(I - P_Z)y \end{aligned}$$

$$E(\hat{\varepsilon}' \hat{\varepsilon}) = E[y'(I - P_Z)y]$$

$$= E[\text{tr}(y'(I - P_Z)y)]$$

$$= E[\text{tr}((I - P_Z)y y')]$$

$$= \text{tr}(E((I - P_Z)y y'))$$

$$= \text{tr}[(I - P_Z)E(y y')]$$

$$= \text{tr}[(I - P_Z)(\sigma^2 I + Z\hat{\beta}\hat{\beta}'Z')] \quad \text{By } (I - P_Z)Z\hat{\beta}\hat{\beta}'Z' = 0$$

$$= \sigma^2 \text{tr}[(I - P_Z) + 0]$$

$$= \sigma^2 \text{tr}(I - P_Z) = \sigma^2 \text{rank}(I - P_Z) = \sigma^2(n - k - 1)$$

Since, $(I - P_Z)P_Z = 0$,

$$\text{rank}(P_Z) + \text{rank}(I - P_Z) = n$$

Also, $\text{rank}(\hat{\varepsilon}) = \text{rank}(P_Z) = \# \text{columns of } Z \text{ assuming that}$

$$= k+1$$

Z has full column rank.

$$\text{So, } \text{rank}(I - P_Z) = n - (k+1).$$

So, an unbiased estimator of σ^2 is

$$\frac{\hat{\Sigma}' \hat{\Sigma}}{n-k-1} = \frac{\mathbf{Y}' (\mathbf{I} - \mathbf{P}_Z) \mathbf{Y}}{n-k-1} = \hat{\sigma}^2$$

Under the "mean 0 uncorrelated" assumption on the errors, we found unbiased estimator of β and σ^2 .

Math B 11800

2/6/18, Tue

$$S(\beta) = (\mathbf{Y} - \mathbf{Z}\beta)' (\mathbf{Y} - \mathbf{Z}\beta) = \sum_{i=1}^n (y_i - z_i\beta)^2$$

$$\beta_{\text{LS}} = (\beta_0, \beta_1, \dots, \beta_r) = \sum_{i=1}^n (y_i - \beta_0 - \beta_1 z_{i1} - \dots - \beta_r z_{ir})^2$$

$$\begin{bmatrix} \frac{\partial S}{\partial \beta_0} \\ \vdots \\ \frac{\partial S}{\partial \beta_r} \end{bmatrix} = \frac{\partial S}{\partial \beta}$$

Last time we looked at the least square estimate of β

$$\hat{\beta}_{\text{LS}} = (\mathbf{Z}' \mathbf{Z})^{-1} \mathbf{Z}' \mathbf{Y}$$

• Decomposition of sum of squares.

$$\text{Total Sum of Square} = \sum_{i=1}^n y_i^2$$

$$\begin{array}{l} \text{Model} \\ \mathbf{Y} = \mathbf{Z}\beta + \varepsilon \end{array}$$

$\uparrow \quad \uparrow$
 $n \times 1 \quad n \times (r+1) \times 1$

$$\text{Total SS} = \mathbf{Y}' \mathbf{Y} - \mathbf{Y}' (\mathbf{I} - \mathbf{P}_Z) \mathbf{Y} = \mathbf{Y}' \mathbf{Y} - \mathbf{Y}' \mathbf{P}_Z \mathbf{Y}$$

$$\sum_{i=1}^n y_i^2 = \sum_{i=1}^n (y_i - \bar{y})^2 + n\bar{y}^2$$

$$\sum_{i=1}^n (y_i - \bar{y})^2 = \sum_{i=1}^n y_i^2 + \sum_{i=1}^n \bar{y}^2 - 2 \sum_{i=1}^n y_i \bar{y}$$

$$= \sum_{i=1}^n y_i^2 - n\bar{y}^2$$

(total SS (corrected) or total SS about the mean.)

$$\bar{y} = \frac{1}{n} \mathbf{1}' \mathbf{y}, \quad \mathbf{1}_n = \begin{bmatrix} 1 \\ \vdots \\ 1 \end{bmatrix}$$

in terms of Matrix?

$$\begin{pmatrix} y_1 - \bar{y} \\ \vdots \\ y_n - \bar{y} \end{pmatrix} = \begin{pmatrix} y_1 \\ \vdots \\ y_n \end{pmatrix} - \begin{pmatrix} \bar{y} \\ \vdots \\ \bar{y} \end{pmatrix} = \underline{y} - \underline{\frac{1}{n} \mathbf{1} \mathbf{1}' \underline{y}} = \underline{y} - \frac{1}{n} \mathbf{1} \mathbf{1}' \underline{y}$$

$$P_{\mathbf{z} \mathbf{z}} = \frac{1}{n} (\mathbf{1} \mathbf{1}')^{-1} \mathbf{1} \mathbf{1}' = \frac{1}{n} \mathbf{1} \mathbf{1}'$$

$$\text{So, } (\mathbf{I} - P_{\mathbf{z}})^2 = (\mathbf{I} - P_{\mathbf{z}})$$

$$\sum_{i=1}^n (y_i - \bar{y})^2 = \begin{pmatrix} y_1 - \bar{y} \\ \vdots \\ y_n - \bar{y} \end{pmatrix}' \begin{pmatrix} y_1 - \bar{y} \\ \vdots \\ y_n - \bar{y} \end{pmatrix}$$

Total SS about the mean

$$= \underline{y}' (\mathbf{I} - P_{\mathbf{z}})^2 \underline{y} = \underline{y}' (\mathbf{I} - P_{\mathbf{z}}) \underline{y}$$

$$\sum_{i=1}^n (y_i - \bar{y})^2 = \underline{y}' (\mathbf{I} - P_{\mathbf{z}}) \underline{y} + \underline{y}' (P_{\mathbf{z}} - P_{\mathbf{z}}) \underline{y}$$

$\hat{\beta}^T \hat{\beta}$

(if bigger \Rightarrow prediction is failed.) (Model is good or not.)

$$\text{Total SS about mean} = \underline{y}' (\mathbf{I} - P_{\mathbf{z}}) \underline{y}$$

$$\underline{y}' (\mathbf{I} - P_{\mathbf{z}}) \underline{y} = \underline{y}' (\mathbf{I} - P_{\mathbf{z}})^2 \underline{y}$$

$$\text{"residual SS"} = \underline{y}' (\mathbf{I} - P_{\mathbf{z}}) (\mathbf{I} - P_{\mathbf{z}}) \underline{y}$$

$$(\mathbf{I} - P_{\mathbf{z}}) \underline{y} = \underline{y} - P_{\mathbf{z}} \underline{y} = \underline{y} - \hat{\underline{y}} \quad (\text{because } P_{\mathbf{z}} \underline{y} = \mathbf{z} \hat{\beta} = \hat{\underline{y}})$$

$$= \hat{\underline{y}}$$

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

residual SS

$$= \sum_{i=1}^n \hat{z}_i^2 = \sum_{i=1}^n (y_i - \hat{\beta}_0 - \hat{\beta}_1 z_{i1} - \dots - \hat{\beta}_r z_{ir})^2$$

"Regression SS"

$$\star \underline{y}' (P_{\mathbf{z}} - P_{\mathbf{z}}) \underline{y} = ?$$

$$P_{\mathbf{z}} \mathbf{1} \mathbf{1}' = \mathbf{1} \mathbf{1}', \text{ then}$$

$$P_{\mathbf{z}} P_{\mathbf{z}} = P_{\mathbf{z}} \frac{1}{n} \mathbf{1} \mathbf{1}' \mathbf{1} \mathbf{1}' = \frac{1}{n} P_{\mathbf{z}} \mathbf{1} \mathbf{1}' \mathbf{1} \mathbf{1}'$$

$$P_{\mathbf{z}} = \frac{1}{n} \mathbf{1} \mathbf{1}' \mathbf{1} \mathbf{1}' = P_{\mathbf{z}}$$

In general,
they are not equal.
But here they are equal

$$\text{so, } (P_{\mathbf{z}} - P_{\mathbf{z}})^2 = P_{\mathbf{z}}^2 + P_{\mathbf{z}}^2 - P_{\mathbf{z}} P_{\mathbf{z}} - P_{\mathbf{z}} P_{\mathbf{z}}$$

$$(P_{\mathbf{z}} - P_{\mathbf{z}}) = P_{\mathbf{z}} + P_{\mathbf{z}} - P_{\mathbf{z}} - P_{\mathbf{z}}$$

$P_{\mathbf{z}} P_{\mathbf{z}} \neq P_{\mathbf{z}} P_{\mathbf{z}}$
"Not commutative"

$$\text{"Regression SS"} = P_{\mathbf{z}} - P_{\mathbf{z}} \quad \text{"idempotent"}$$

$$\underline{y}' (P_{\mathbf{z}} - P_{\mathbf{z}}) \underline{y} = \underline{y}' (P_{\mathbf{z}} - P_{\mathbf{z}})(P_{\mathbf{z}} - P_{\mathbf{z}}) \underline{y} = \sum_{i=1}^n (\hat{y}_i - \bar{y})^2$$

$$(P_{\mathbf{z}} - P_{\mathbf{z}}) \underline{y} = \hat{\underline{y}} - \bar{y} = \begin{pmatrix} \hat{y}_1 - \bar{y} \\ \vdots \\ \hat{y}_n - \bar{y} \end{pmatrix}$$

$$\begin{matrix} P_{\mathbf{z}} \underline{y} - P_{\mathbf{z}} \underline{y} \\ \hat{\underline{y}} - \bar{y} \end{matrix}$$

Also observe $(I - P_Z)(P_Z - P_{\bar{Z}}) = 0$.

Also, $(I - P_{\bar{Z}})P_{\bar{Z}} = 0$

\therefore , the decomposition of \underline{Y} into $(I - P_{\bar{Z}})\underline{Y}$ and $P_Z\underline{Y}$ is an orthogonal decomposition.

Also, $(I - P_{\bar{Z}})\underline{Y} = (I - P_Z)\underline{Y} + (P_Z - P_{\bar{Z}})\underline{Y}$

This is also orthogonal decomposition.

A measure of goodness of the regression is

$$\begin{aligned} R^2 &= \frac{\text{Reg. SS}}{\text{Tot. SS about Mean}} = \frac{\underline{Y}'(P_Z - P_{\bar{Z}})\underline{Y}}{\underline{Y}'(I - P_{\bar{Z}})\underline{Y}} \\ &= \frac{\sum_{i=1}^n (\hat{Y}_i - \bar{Y})^2}{\sum_{i=1}^n (Y_i - \bar{Y})^2} = 1 - \frac{\sum_{i=1}^n (Y_i - \hat{Y}_i)^2}{\sum_{i=1}^n (Y_i - \bar{Y})^2} \\ &= 1 - \frac{\text{Res. SS}}{\text{Total SS about mean}} \end{aligned}$$

$\therefore R^2$ lies between 0 and 1.

not good fit

If R^2 is close to 1, then the regression model is a good fit.

(predictor has no information for \underline{Y})

If R^2 is close to 0, then the model is that the predictors cannot predict the response well.

$\therefore R^2 = 1$ only when $\hat{Y}_i = 0$ for all i .

$R = \sqrt{R^2}$ is called

Multiple correlation coefficient.

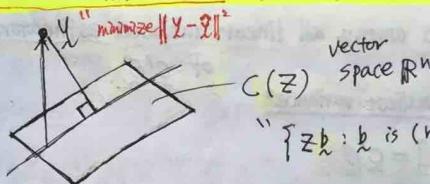
This means $Y_i = \hat{Y}_i = \hat{\beta}_0 + \hat{\beta}_1 Z_{i1} + \dots + \hat{\beta}_r Z_{ir}$.

$R^2 = 0$, when $\hat{\beta}_0 = \bar{Y}$, $\hat{\beta}_1 = \dots = \hat{\beta}_r = 0$.

Geometric Point of view for the LS method.

The LS problem:

minimize $\|\underline{Y} - \underline{Z}\underline{\beta}\|^2 = (\underline{Y} - \underline{Z}\underline{\beta})'(\underline{Y} - \underline{Z}\underline{\beta})$ with respect to $\underline{\beta}$



" $\{\underline{Z}\underline{b} : \underline{b} \text{ is } (r+1) \times 1 \text{ vector}\}$ "

So, the minimization problem is equivalent to orthogonal projection of \underline{Y} onto $C(Z)$.

So, if $\underline{Z}\underline{b}$ is the closest to \underline{Y} among all points of $C(Z)$, then $\underline{Y} - \underline{Z}\underline{b}$ is orthogonal to $C(Z)$.

So, $\underline{Z}'(\underline{Y} - \underline{Z}\underline{b}) = 0$ for any $\underline{z} \in C(Z)$.

So, $\underline{Z}'(\underline{Y} - \underline{Z}\underline{b}) = 0 \Rightarrow \underline{b} = (\underline{Z}'\underline{Z})^{-1}\underline{Z}'\underline{Y}$

Hence, $\underline{Z}\underline{b} = P_Z\underline{Y}$ by $\underline{Z}\underline{b} = \underline{Z}(\underline{Z}'\underline{Z})^{-1}\underline{Z}'\underline{Y}$

• Best Linear Unbiased Estimator (BLUE)

column vector

For any vector $\underline{\beta}$, suppose we want to estimate $\underline{\beta}'\underline{\beta}$

An estimator will be called linear if it has the form $\underline{a}'\underline{y}$ for some \underline{a}

Thm (Gauss)

$\underline{\beta}'\hat{\underline{\beta}}$ is the best among all linear unbiased estimators of $\underline{\beta}'\underline{\beta}$.

Best in the sense "smallest variance"

$$\mathbb{E}[\underline{\beta}'\hat{\underline{\beta}}] = \underline{\beta}'\mathbb{E}[\hat{\underline{\beta}}] = \underline{\beta}'\underline{\beta}$$

so, $\underline{\beta}'\hat{\underline{\beta}}$ is unbiased for $\underline{\beta}'\underline{\beta}$. $\underline{a}'\underline{Z}\hat{\underline{\beta}}$ so $\underline{a}'\underline{Z} = \underline{\beta}'$

$$\underline{\beta}'\hat{\underline{\beta}} = \underline{\beta}'(\underline{Z}'\underline{Z})^{-1}\underline{Z}'\underline{Y} = \underline{a}'\underline{Y}, \text{ where } \underline{a} = \underline{Z}(\underline{Z}'\underline{Z})^{-1}\underline{\beta}$$

linear so, $\underline{\beta}'\hat{\underline{\beta}}$ is also linear estimator.

If $\underline{a}'\underline{Y}$ is unbiased for $\underline{\beta}'\underline{\beta}$, then $\mathbb{E}[\underline{a}'\underline{Y}] = \mathbb{E}[\underline{a}'\underline{Z}\underline{\beta} + \underline{a}'\underline{\varepsilon}]$

$$\Rightarrow \mathbb{E}[\underline{a}'\underline{Y}] = \underline{\beta}'\underline{\beta} \text{ for all } \underline{\beta} \quad \underline{a}'\underline{Z} = \underline{\beta}'$$

$$\text{so, } \underline{a}'\underline{Z}\underline{\beta} = \underline{\beta}'\underline{\beta} \text{ for all } \underline{\beta} \quad \text{so, } \underline{a}'\underline{Z} = \underline{\beta}'$$

$$\text{by } (*) \quad \mathbb{E}[\underline{Y}] \text{ or } (\underline{C}\underline{a}')\underline{\beta} = 0, \forall \underline{\beta}$$

$$\text{var}(\underline{a}'\underline{Y}) = \underline{a}'\text{var}(\underline{Y})\underline{a} = \sigma^2 \underline{a}'\underline{a}$$

$$\text{var}(\underline{\beta}'\hat{\underline{\beta}}) = \text{var}(\underline{a}'\underline{Y}) = \sigma^2 \underline{a}'\underline{a}$$

$$\begin{aligned} \underline{a}'\underline{a} &= (\underline{a} - \underline{a}^*)'(\underline{a} - \underline{a}^*) \\ &= (\underline{a} - \underline{a}^*)'(\underline{a} - \underline{a}^*) + \underline{a}'\underline{a}^* + 2(\underline{a} - \underline{a}^*)'\underline{a}^* \geq \underline{a}'\underline{a}^*. \end{aligned}$$

Math B7800

2/6/18, Thur

$$\underline{Y} = (\underline{I} - \underline{P}_{\underline{Z}})\underline{Y} + \underline{P}_{\underline{Z}}\underline{Y} \text{ "orthogonal decomposition"}$$

BLUE (continued)

$$\underline{a}^* = \underline{Z}(\underline{Z}'\underline{Z})^{-1}\underline{\beta}$$

$$\underline{a} \text{ satisfies } \underline{a}'\underline{Z} = \underline{\beta}'$$

Need to check $(\underline{a} - \underline{a}^*)'\underline{a}^* = 0$.

$$(\underline{a} - \underline{a}^*)'\underline{Z}(\underline{Z}'\underline{Z})^{-1}\underline{\beta} = 0, \text{ since } (\underline{a} - \underline{a}^*)'\underline{Z} = \underline{\beta}' - \underline{\beta}' = 0$$

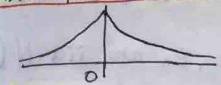
Inference about Regression

so far we didn't assume any distribution for the error ε_i 's.
we only assumed $\mathbb{E}(\varepsilon_i) = 0$ and $\text{cov}(\varepsilon_i, \varepsilon_j) = \begin{cases} 0 & \text{if } i \neq j \\ \sigma^2 & \text{if } i = j \end{cases}$

Now, we need to assume distribution of the error.

1. Suppose the errors are iid double exponential dist.

$$\text{with pdf } f(x) = \frac{1}{2\sigma} e^{-|x|/\sigma}, x \in \mathbb{R}$$



Suppose $\varepsilon_1, \dots, \varepsilon_n$ are iid with pdf $f(x)$.

Q: Find MLE of $\underline{\beta}, \sigma^2$.

$$\text{Model } \underline{Y} = \underline{Z}\underline{\beta} + \underline{\varepsilon}$$

Suppose y_1, y_2, \dots, y_n are independent.

$$\mathbb{E}(y_i) = z_{i*}\underline{\beta} = \beta_0 + \beta_1 z_{i1} + \dots + \beta_r z_{ir}$$

dist. of y_i ?

PDF of y_i ?

PDF of $\varepsilon_i = y_i - z_{i*}\beta$ is $f(x)$ i.e. Double exponential dist.

$$f(y_i) = \frac{1}{2} e^{-|y_i - z_{i*}\beta|}$$

So, joint pdf of y_1, \dots, y_n .

$$f(y_1, \dots, y_n) = \prod_{i=1}^n f(y_i) = \frac{1}{2^n} e^{-\sum_{i=1}^n |y_i - z_{i*}\beta|} = L(\beta)$$

"Not continuous!" "likelihood function"

Maximizing $L(\beta)$ is equivalent to

$$\text{Minimizing } \sum_{i=1}^n |y_i - \beta_0 - \beta_1 z_{i1} - \dots - \beta_r z_{ir}| \text{ w.r.t } \beta.$$

The estimator that minimizes the last term is called LAD (Least absolute deviation) estimator of β .

Note Suppose the errors are normally distributed.

$$\text{So, } \varepsilon_i \sim N_n(0, \sigma^2 I_n)$$

$$\begin{aligned} \text{Normal Dist} \\ f(x) = \frac{1}{\sqrt{2\pi}\sigma} e^{-\frac{1}{2}(\frac{x-\mu}{\sigma})^2} \end{aligned}$$

So, $\varepsilon_1, \dots, \varepsilon_n$ are iid $N(0, \sigma^2)$

* Q: Find MLE of β, σ^2

Distribution of y_i is $N(z_{i*}\beta, \sigma^2)$

$$\text{Pdf of } y_i \quad f(y_i) = \frac{1}{\sqrt{2\pi}\sigma} e^{-\frac{1}{2\sigma^2}(y_i - z_{i*}\beta)^2}$$

Joint pdf of y_1, \dots, y_n

$$f(y_1, \dots, y_n) = \left(\frac{1}{\sqrt{2\pi}\sigma}\right)^n e^{-\sum_{i=1}^n \frac{1}{2\sigma^2}(y_i - z_{i*}\beta)^2} = L(\beta, \sigma^2)$$

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

$$- \frac{1}{2\sigma^2} \sum_{i=1}^n (y_i - z_{i*}\beta)^2$$

$$\frac{\partial}{\partial \beta} l(\beta, \sigma^2) = -\frac{1}{2\sigma^2} \underbrace{\frac{\partial}{\partial \beta} \sum_{i=1}^n (y_i - z_{i*}\beta)^2}_{=0} = 0$$

$$\xrightarrow{\text{minimizer}} \Rightarrow 0$$

$$\text{So, } \hat{\beta} = \hat{\beta}_{LS} = (\mathbf{Z}'\mathbf{Z})^{-1} \mathbf{Z}' \mathbf{y}$$

MLE of β

$$\frac{\partial}{\partial \sigma^2} l(\beta, \sigma^2) = -\frac{n}{2\sigma^2} + \frac{1}{2\sigma^4} \sum_{i=1}^n (y_i - z_{i*}\beta)^2 = 0$$

$$\text{So, } \hat{\sigma}^2 = \frac{1}{n} \sum_{i=1}^n (y_i - z_{i*}\hat{\beta})^2 = \frac{1}{n} \sum_{i=1}^n \hat{\varepsilon}_i^2 = \frac{1}{n} \hat{\varepsilon}' \hat{\varepsilon}$$

Hence, $\hat{\beta}, \hat{\sigma}^2$ are the MLE.

$\hat{\sigma}^2$ is MLE, but not unbiased for σ^2 .

$$\text{The unbiased estimator for } \sigma^2 \text{ is } \frac{\hat{\varepsilon}' \hat{\varepsilon}}{n-r-1} = s^2$$

$\hat{\beta}$ is unbiased for β .

$$\begin{aligned} \hat{\beta} &= (\mathbf{Z}'\mathbf{Z})^{-1} \mathbf{Z}' \mathbf{y} \sim N_{r+1} \left(\underbrace{(\mathbf{Z}'\mathbf{Z})^{-1} \mathbf{Z}' \beta}_{I}, \underbrace{(\mathbf{Z}'\mathbf{Z})^{-1} \mathbf{Z}' (\sigma^2 \mathbf{I})}_{\mathbf{Z}' (\mathbf{Z}'\mathbf{Z})^{-1}} \right) \\ \mathbf{y} &\sim N_n(\mathbf{Z}\beta, \sigma^2 \mathbf{I}) \\ &= N_{r+1}(\beta, \sigma^2 (\mathbf{Z}'\mathbf{Z})^{-1}) \end{aligned}$$

Fact $\hat{\beta}$ and $\hat{\sigma}^2$ (or s^2) are independent. ?

We will show $\hat{\beta}$ and $\hat{\sigma}^2$ are independent.

$$\hat{\Sigma} = (I - P_Z)Y, \quad \hat{\beta} = (Z'Z)^{-1}Z'Y.$$

$$\text{cov}(\hat{\Sigma}, \hat{\beta}) = (I - P_Z)(\sigma^2 I)Z(Z'Z)^{-1} = 0. \leftarrow$$

since $(I - P_Z)Z = 0$

So, $\hat{\Sigma}, \hat{\beta}$ are independent.

Since $\hat{\sigma}^2$ and s^2 are both function of $\hat{\Sigma}$ they are independent with $\hat{\beta}$.

Q: Distribution of $\hat{\sigma}^2$ and s^2 ?

$$\text{Fact: } \hat{\Sigma}'\hat{\Sigma} \sim \chi^2_{n-r-1}.$$

$$R^2 = \frac{Y'(I-P_Z)Y}{Y'(I-P_Z)Y}$$

$$1-R^2 = \frac{Y'(Z-P_Z)Y}{Y'(Z-P_Z)Y}$$

Linear Algebra Digression

Suppose A is symmetric real matrix $(n \times n)$.

If $\lambda_1, \lambda_2, \dots, \lambda_n$ are the eigenvalues of A , then eigenvalues of A^2 are $\lambda_1^2, \dots, \lambda_n^2$.

If $A = U\Delta U'$ is the spectral decomposition

i.e., U is orthogonal and $\Delta = \begin{pmatrix} \lambda_1 & & \\ & \ddots & \\ & & \lambda_n \end{pmatrix}$,

then $A^2 = U\Delta U'U\Delta U' = U\Delta^2 U'$, $\Delta^2 = \begin{pmatrix} \lambda_1^2 & & \\ & \ddots & \\ & & \lambda_n^2 \end{pmatrix}$. This shows that $\lambda_1^2, \dots, \lambda_n^2$ are the eigenvalues of A .

If A is idempotent, then each eigenvalue is either 0 or 1.

because if $A = U\Delta U'$, then $A = A^2$

$$\Rightarrow U\Delta U' = U\Delta^2 U'$$

$$\Rightarrow \Delta = \Delta^2$$

$$\Rightarrow \lambda_i = \lambda_i^2 \quad \forall i \in \mathbb{N}$$

$$\text{So, } \lambda_i = \text{either 0 or 1}.$$

$$\text{So, } A = U \begin{bmatrix} I & 0 \\ 0 & 0 \end{bmatrix} U', \quad U = \begin{bmatrix} U_1 & U_2 \\ n \times k & n \times (n-k) \end{bmatrix}$$

$$= U_1 U_1' \quad \text{rank}(A) = k.$$

$$\hat{\Sigma}'\hat{\Sigma} = Y'(I - P_Z)Y. \quad \text{since } (I - P_Z) \text{ is idempotent and rank}(I - P_Z) = n - r - 1$$

$$\text{So, } I - P_Z = UU', \quad \text{where } U \text{ is } n \times (n-r-1).$$

$$\text{and } U'U = I.$$

$$\text{So, } Y'(I - P_Z)Y = Y'UU'Y \quad \text{So, } \frac{Y'Y}{\sigma^2} \sim \chi^2_{n-r-1}$$

$$\text{Now, see that } U'Y \sim N(U'Z\beta, \sigma^2 I) \quad \frac{Y'(I - P_Z)Y}{\sigma^2} = \frac{\hat{\Sigma}'\hat{\Sigma}}{\sigma^2}$$

$$\text{Also, } UU'Z = (I - P_Z)Z = 0. \quad \left\{ \begin{array}{l} U'Y \sim N_{n-r-1}(0, \sigma^2 I) \\ U'Z = 0 \end{array} \right.$$

$$\text{So, } U'UU'Z = 0$$

$$\text{So, } U'Z = 0.$$

2/13/18, Tue

Math B7800

Model: $\mathbf{y} = \mathbf{z}\beta + \varepsilon$, where $\varepsilon \sim N_n(0, \sigma^2 I_n)$.

Last time $\frac{\mathbf{y}'(\mathbf{I} - \mathbf{P}_z)\mathbf{y}}{\sigma^2} \sim \chi^2_{n-r-1}$ rank of $(\mathbf{I} - \mathbf{P}_z)$.

Also, $\hat{\beta}$ and $\mathbf{y}'(\mathbf{I} - \mathbf{P}_z)\mathbf{y}$ are independent.

Goal: Obtain Confidence Region for β . Then find Scheffe CI and Bonferroni CI for β_i .

Recall: $\hat{\beta} = (\mathbf{z}'\mathbf{z})^{-1}\mathbf{z}'\mathbf{y} \sim N_{r+1}(\beta, \sigma^2(\mathbf{z}'\mathbf{z})^{-1})$

Recall: For symmetric positive definite matrix A , we can define

\mathcal{Q} : Confidence Region for β .

So, $(\hat{\beta} - \beta) \sim N_{r+1}(0, \sigma^2(\mathbf{z}'\mathbf{z})^{-1})$.

$\frac{1}{\sigma}(\mathbf{z}'\mathbf{z})^{1/2}(\hat{\beta} - \beta) \sim N_{r+1}(0, I_{r+1})$

So, $\left[\frac{1}{\sigma}(\mathbf{z}'\mathbf{z})^{1/2}(\hat{\beta} - \beta) \right]' \frac{1}{\sigma}(\mathbf{z}'\mathbf{z})^{1/2}(\hat{\beta} - \beta) = \mathbf{u}'\mathbf{u} \sim \chi^2_{r+1}$.

$\hookrightarrow \frac{1}{\sigma^2}(\hat{\beta} - \beta)'(\mathbf{z}'\mathbf{z})(\hat{\beta} - \beta) \sim \chi^2_{r+1}$

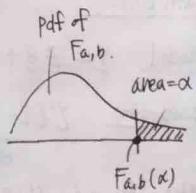
So, using independent of $\hat{\beta}$ and $\mathbf{y}'(\mathbf{I} - \mathbf{P}_z)\mathbf{y}$.

$F = \frac{\frac{1}{r+1} \frac{1}{\sigma^2} (\hat{\beta} - \beta)' (\mathbf{z}'\mathbf{z}) (\hat{\beta} - \beta)}{\frac{1}{n-r-1} \frac{1}{\sigma^2} \mathbf{y}' (\mathbf{I} - \mathbf{P}_z) \mathbf{y}} \sim F_{r+1, n-r-1}$

$$= \frac{1}{r+1} \frac{1}{\sigma^2} (\hat{\beta} - \beta)' (\mathbf{z}'\mathbf{z}) (\hat{\beta} - \beta)$$

Notation: $F_{r+1, n-r-1}(x)$.

$$F_{a,b}(x)$$



so, $P(F \leq F_{r+1, n-r-1}(x)) = 1 - \alpha$.

Thus, $\{\beta : (\hat{\beta} - \beta)'(Z'Z)(\hat{\beta} - \beta) \leq (r+1)s^2 F_{r+1, n-r-1}(\alpha)\}$

is a $100(1-\alpha)\%$ Confidence Region for β .

Next topic is CI for β_i "only i , one interval."

Distribution of $\hat{\beta}_i$?

$$\hat{\beta}_i \sim N(\beta_i, s^2(Z'Z)^{-1}_{i+1, i+1}) \text{ for } i=0, 1, \dots, r.$$

Note that $(Z'Z)^{-1}_{i+1, i+1}$ refers to the $(i+1), (i+1)$ entry of $(Z'Z)^{-1}$.

Now, $\frac{\hat{\beta}_i - \beta_i}{\sqrt{s^2(Z'Z)^{-1}_{i+1, i+1}}} \sim N(0, 1)$.

Since s^2 is unknown, we cannot use it for a CI of β_i . We replace s^2 by s^2 .

$$\begin{aligned} \frac{\hat{\beta}_i - \beta_i}{\sqrt{s^2(Z'Z)^{-1}_{i+1, i+1}}} &\sim t_{n-r-1}. \\ &= \frac{(\hat{\beta}_i - \beta_i)\sqrt{s^2(Z'Z)^{-1}_{i+1, i+1}}}{\sqrt{s^2/\hat{s}^2}} \sim N(0, 1) \end{aligned}$$

So, a $100(1-\alpha)\%$ CI for β_i is

$$\hat{\beta}_i \pm \sqrt{s^2(Z'Z)^{-1}_{i+1, i+1}} t_{n-r-1} \left(\frac{\alpha}{2}\right)$$

$$\hat{\beta}_i \pm \sqrt{\text{var}(\hat{\beta}_i)} t_{n-r-1} \left(\frac{\alpha}{2}\right)$$

$$\text{var}(\hat{\beta}_i) = s^2(Z'Z)^{-1}_{i+1, i+1}$$

$$\text{var}(\hat{\beta}_i) = s^2(Z'Z)^{-1}_{i+1, i+1}$$

estimator

• Simultaneous CI for β_i , $i=0, 1, \dots, k$.

Two standard methods

- (1) Bonferroni, (2) Scheffe's method

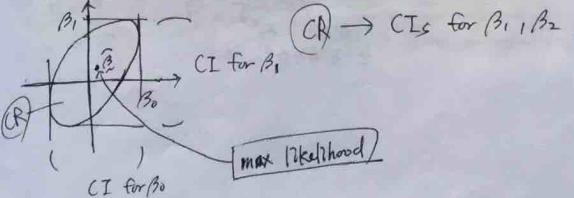
(1) Bonferroni simultaneous CI for $\beta_0, \beta_1, \dots, \beta_k$.

with confidence $1-\alpha$ is

$$\hat{\beta}_i \pm \sqrt{\text{var}(\hat{\beta}_i)} t_{n-r-1} \left(\frac{\alpha}{2(k+1)}\right)$$

(2) Scheffe's method, we need to project the $(1-\alpha)100\%$ CR onto different axes.

Suppose β has 2 components.



Diagression (Linear Algebra)

Suppose A is a symmetric matrix.

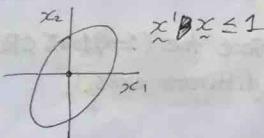
$$\max_{\underline{x}} \underline{x}' A \underline{x} = \lambda_1(A) \quad (\text{largest eigenvalue of } A) \\ : \|\underline{x}\|_2 = 1$$

$$= \max_{\substack{\underline{x}: \|\underline{x}\|_2 = 1}} \underline{x}' A \underline{x} = \lambda_1(A). \quad \text{positive definite.}$$

Suppose A is symmetric, B is symmetric and pd.

$$\max_{\underline{x}: \underline{x}' B \underline{x}} (\underline{x}' A \underline{x}) = \max_{\substack{y: \|y\|_2 = 1}} y' B^{-1/2} A B^{-1/2} y \\ \text{Define } y = B^{1/2} \underline{x} \quad = \lambda_1(B^{-1/2} A B^{-1/2})$$

Next, for two matrices C and D for which CD , DC are both defined, then all nonzero eigenvalues of CD and DC are same.



Finding the projection onto x_1 -axes.

$$\Leftrightarrow \max_{\substack{\underline{x}: \underline{x}' B \underline{x} \leq 1}} \underline{x}_1^2$$

$$\max_{\substack{\underline{x}: \underline{x}' B \underline{x} \leq 1}} \underline{x}_1^2$$

Finding the projection onto x_1 -axis.

$$\Leftrightarrow \max_{\substack{\underline{x}: \underline{x}' B \underline{x} \leq 1}} \underline{x}_1^2$$

$$= \max_{\substack{\underline{x}: \underline{x}' B \underline{x} \leq 1}} \underline{x}' A \underline{x}, \text{ where } A = \begin{pmatrix} 1 & 0 \\ 0 & 0 \end{pmatrix} = \begin{pmatrix} 1 & 0 \\ 0 & 0 \end{pmatrix}$$

$$= \underline{e}_1 \underline{e}_1'$$

$$= \lambda_1(B^{-1/2} A B^{-1/2})$$

$$= \lambda_1(B^{-1/2} A B^{-1/2})$$

$$= \lambda_1(B^{-1/2} \underline{e}_1 \underline{e}_1' B^{-1/2})$$

$$= \lambda_1(\underline{e}_1' B^{-1} \underline{e}_1) = \underline{e}_1' B^{-1} \underline{e}_1 = (B^{-1})_{1,1}$$

scalar.

So, we see that

$$\max_{\substack{\underline{x}: \underline{x}' B \underline{x} \leq 1}} \underline{x}_1^2 = (B^{-1})_{1,1}.$$

Recall 100(1- α)% CR for β was

$$(\beta - \hat{\beta})^* (\underline{z}' \underline{z}) (\beta - \hat{\beta}) \leq (r+1) s^2 F_{r+1, n-r-1}(\alpha).$$

Write in the term $\underline{x}' B \underline{x} \leq 1$, where

$$\underline{x} = \beta - \hat{\beta}, \quad B = (\underline{z}' \underline{z}) \frac{1}{(r+1) s^2 F_{r+1, n-r-1}(\alpha)}.$$

$$\max_{\beta: \beta \text{ is MCR.}} |\beta_i - \hat{\beta}_i|^2 = (\mathbf{B}^{-1})_{i+1, i+1}, i=0, 1, \dots, r \\ = (r+1) s^2 F_{r+1, n-r-1}(\alpha) (\mathbf{Z}' \mathbf{Z})_{i+1, i+1}^{-1}$$

so, the Scheffe simultaneous CI for β_0, \dots, β_r is

$$\hat{\beta}_i \pm \sqrt{(r+1) s^2 F_{r+1, n-r-1}(\alpha) (\mathbf{Z}' \mathbf{Z})_{i+1, i+1}^{-1}}$$

$$= \hat{\beta}_i \pm \sqrt{(r+1) F_{r+1, n-r-1}(\alpha)} \widehat{\text{var}}(\hat{\beta}_i)$$

Scheffe simultaneous CI for all linear combinations of the components of β , namely $\underline{\alpha}' \underline{\beta}$ for all vectors $\underline{\alpha}$ is given by

$$\underline{\alpha}' \hat{\beta} \pm \sqrt{(r+1) F_{r+1, n-r-1}(\alpha)} \widehat{\text{var}}(\underline{\alpha}' \hat{\beta}), \text{ where}$$

$$\widehat{\text{var}}(\underline{\alpha}' \hat{\beta}) = \underline{\alpha}' \widehat{\text{var}}(\hat{\beta}) \underline{\alpha} \\ = s^2 \underline{\alpha}' (\mathbf{Z}' \mathbf{Z})^{-1} \underline{\alpha}$$

$$\underline{\beta}_2 = \mathbf{Z} (\mathbf{Z}' \mathbf{Z})^{-1} \mathbf{Z}' \\ \underline{\beta}_{21} = \mathbf{Z}_1 (\mathbf{Z}' \mathbf{Z})^{-1} \mathbf{Z}_1'$$

$$\mathbf{Z} = \begin{bmatrix} 1 & 2 \\ 3 & 4 \end{bmatrix}$$

$$\mathbf{Z}_1 = \begin{bmatrix} 1 & 0 \\ 3 & 0 \end{bmatrix}$$

Math B7800

2/15/18, Thur

March 06, Tue - Exam 01. (3 pages formula sheet)

Likelihood Ratio test

$$\underline{\beta} = \begin{bmatrix} \beta_1 \\ \vdots \\ \beta_{r+1} \\ \beta_2 \end{bmatrix}_{(r+1) \times 1} \quad \mathbf{Z} = \begin{bmatrix} \mathbf{Z}_1 & | & \mathbf{Z}_2 \end{bmatrix}_{n \times (r+1)} \quad \begin{bmatrix} \mathbf{Z}_1 \\ \mathbf{Z}_2 \end{bmatrix}_{n \times (r+1)} \quad \begin{bmatrix} \mathbf{Z}_1 & | & \mathbf{Z}_2 \end{bmatrix}_{n \times (r+1)}$$

$$\text{Then } \underline{y} = \mathbf{Z} \underline{\beta} + \underline{\xi} = \begin{bmatrix} \mathbf{Z}_1 & | & \mathbf{Z}_2 \end{bmatrix} \begin{bmatrix} \beta_1 \\ \vdots \\ \beta_{r+1} \\ \beta_2 \end{bmatrix} + \underline{\xi} \\ = \mathbf{Z}_1 \beta_1 + \mathbf{Z}_2 \beta_2 + \underline{\xi}$$

$$\text{"Test"} \quad H_0: \beta_2 = 0 \quad \text{vs} \quad H_1: \beta_2 \neq 0.$$

Recall the likelihood ratio test (LRT).

$$\Delta = \frac{\sup_{H_0} L(\beta, \sigma^2)}{\sup_{H_1} L(\beta, \sigma^2)} \text{ is the likelihood Ratio.}$$

We reject H_0 if Δ is small.

We express Δ as a monotone function of a statistic T (of the data) whose null distribution is known (under H_0).

Then $\{\Delta < c\} \iff \{T < c' \text{ or } > c'\}$ for any constant c .

Then, for any constant c , there is another constant c' s.t.
 $\{\Delta < c\} \Leftrightarrow \{T < c'\}$ if Δ is increasing
 $\Leftrightarrow \{T > c'\}$ if Δ is decreasing.

Then we will reject H_0 if $\begin{cases} T < c' & \text{in case 1} \\ T > c' & \text{in case 2.} \end{cases}$

Since the null distribution of T will be known, we can find c'
so that the "level condition" is met.

In our case,

$$L(\beta, \sigma^2) = (2\pi)^{-n/2} (\sigma^2)^{-n/2} \exp \left[-\frac{1}{2\sigma^2} (\mathbf{y} - \mathbf{z}\beta)' (\mathbf{y} - \mathbf{z}\beta) \right]$$

The maximize of $L(\beta, \sigma^2)$ are

$$\hat{\beta} = (\mathbf{z}'\mathbf{z})^{-1} \mathbf{z}'\mathbf{y}, \quad \hat{\sigma}^2 = \frac{(\mathbf{y} - \mathbf{z}\hat{\beta})' (\mathbf{y} - \mathbf{z}\hat{\beta})}{n}$$

$$\sup L(\beta, \sigma^2) = L(\hat{\beta}, \hat{\sigma}^2) = (2\pi)^{-n/2} (\hat{\sigma}^2)^{-n/2} e^{-\frac{n}{2}}$$

If H_0 is true, then $\mathbf{y} = \mathbf{z}_1\beta_1 + \xi$.

$$L\left(\begin{pmatrix} \beta_1 \\ 0 \end{pmatrix}, \sigma^2\right) = (2\pi)^{-n/2} (\sigma^2)^{-n/2} \exp \left[-\frac{1}{2\sigma^2} (\mathbf{y} - \mathbf{z}_1\beta_1)' (\mathbf{y} - \mathbf{z}_1\beta_1) \right]$$

The maximizer of $L(\beta_1, \sigma^2)$ are

$$\hat{\beta}_1 = (\mathbf{z}_1'\mathbf{z}_1)^{-1} \mathbf{z}_1' \mathbf{y}, \quad \hat{\sigma}^2 = \frac{(\mathbf{y} - \mathbf{z}_1\hat{\beta}_1)' (\mathbf{y} - \mathbf{z}_1\hat{\beta}_1)}{n}$$

$$\text{So, } \sup_{H_0} L(\beta, \sigma^2) = L\left(\begin{pmatrix} \hat{\beta} \\ 0 \end{pmatrix}, \hat{\sigma}^2\right)$$

$$= (2\pi)^{-n/2} (\hat{\sigma}^2)^{-n/2} e^{-\frac{n}{2}}$$

$$\text{So, } \Delta = \left(\frac{\hat{\sigma}^2}{\sigma^2} \right)^{-n/2} = \left(\frac{\mathbf{y}'(I - P_{\mathbf{z}})\mathbf{y}}{\mathbf{y}'(I - P_{\mathbf{z}})\mathbf{y}} \right)^{-n/2}$$

"decreasing"

So, we will reject H_0 if $\frac{\mathbf{y}'(I - P_{\mathbf{z}})\mathbf{y}}{\mathbf{y}'(I - P_{\mathbf{z}})\mathbf{y}} > c$, for some c .

Recall that $(I - P_{\mathbf{z}}), (I - P_{\mathbf{z}_1})$ are idempotent.

$$\text{So, } \frac{\mathbf{y}'(I - P_{\mathbf{z}})\mathbf{y}}{\sigma^2} \sim \chi_{n-r+1}^2. \quad \text{"not independent"}$$

$$\text{under } H_0, \quad \frac{\mathbf{y}'(I - P_{\mathbf{z}_1})\mathbf{y}}{\sigma^2} \sim \chi_{n-q+1}^2 \quad \text{rank}(\mathbf{z}_1) = q+1$$

These two χ^2 will be independent if $(I - P_{\mathbf{z}})(I - P_{\mathbf{z}_1}) = 0$.

$$= I - P_{\mathbf{z}} - P_{\mathbf{z}_1} + P_{\mathbf{z}_1} \neq 0$$

But, "Not independent" (as $P_{\mathbf{z}} \cdot \mathbf{z}_1 = \mathbf{z}_1$)

$$\frac{\mathbf{y}'(I - P_{\mathbf{z}_1})\mathbf{y}}{\mathbf{y}'(I - P_{\mathbf{z}})\mathbf{y}} - 1 = \frac{\mathbf{y}'(P_{\mathbf{z}} - P_{\mathbf{z}_1})\mathbf{y}}{\mathbf{y}'(I - P_{\mathbf{z}})\mathbf{y}} \quad \text{independent}$$

$$\text{check } (P_{\mathbf{z}} - P_{\mathbf{z}_1})(I - P_{\mathbf{z}}) = P_{\mathbf{z}} - P_{\mathbf{z}}^2 - P_{\mathbf{z}_1} + P_{\mathbf{z}_1} = 0.$$

Under H_0 ,

$$\frac{\frac{1}{r-q} \chi' (P_z - P_{z_1}) \chi / \sigma^2}{\frac{1}{n-r-1} \chi' (I - P_z) \chi / \sigma^2} \sim F_{r-q, n-r-1}.$$

$$\text{So, } \frac{1}{r-q} \frac{1}{\sigma^2} \chi' (P_z - P_{z_1}) \chi \sim F_{r-q, n-r-1}.$$

So, we reject H_0 at level of α if

$$\chi' (P_z - P_{z_1}) \chi > (r-q) \sigma^2 F_{r-q, n-r-1}(\alpha)$$

test Δ , decreasing, $\Rightarrow \text{fInv}(1-\alpha, r-q, n-r-1)$

$$\beta = \begin{bmatrix} \beta_1 \\ \beta_2 \end{bmatrix} \quad * \text{ Confidence Region for } \beta_2$$

Recall $\hat{\beta} \sim N_{p+1}(\beta, \sigma^2(Z'Z)^{-1})$. Then $\hat{\beta}_2 \sim ?$

$$\hat{\beta}_2 \sim N_{r-q}(\beta_2, \sigma^2(Z'_2(I-P_{z_1})Z_2)^{-1})$$

$$Z = [Z_1 | Z_2]$$

$$Z'Z = \begin{bmatrix} Z_1' \\ Z_2' \end{bmatrix} \begin{bmatrix} Z_1 & Z_2 \end{bmatrix} = \begin{bmatrix} Z_1'Z_1 & Z_1'Z_2 \\ Z_2'Z_1 & Z_2'Z_2 \end{bmatrix}$$

$$\text{Inverse of block matrix. } \begin{bmatrix} A & B \\ C & D \end{bmatrix}^{-1} = \begin{bmatrix} * & * \\ * & (D-CA^{-1}B)^{-1} \end{bmatrix}$$

$$(Z'Z)^{-1} = \begin{bmatrix} * & * \\ * & (Z_2'Z_2 - Z_2'Z_1(Z_1'Z_1)^{-1}Z_1'Z_2)^{-1} \end{bmatrix} (Z_2'(I-P_{z_1})Z_2)^{-1}$$

$$\frac{(\hat{\beta}_2 - \beta_2)' Z_2' (I - P_{z_1}) Z_2 (\hat{\beta}_2 - \beta_2)}{\sigma^2} \sim \chi^2_{r-q}.$$

$$\text{Then } \frac{1}{r-q} \frac{1}{\sigma^2} (\hat{\beta}_2 - \beta_2)' Z_2' (I - P_{z_1}) Z_2 (\hat{\beta}_2 - \beta_2) \sim F_{r-q, n-r-1}$$

So, a $100(1-\alpha)\%$ CR for β_2 is

$$\{ \beta_2 : (\hat{\beta}_2 - \beta_2)' Z_2' (I - P_{z_1}) Z_2 (\hat{\beta}_2 - \beta_2) \leq (r-q) \sigma^2 F_{r-q, n-r-1}(\alpha) \}$$

H_0 belongs to this ellipse.
is true

Prediction using the fitted regression model.

We have fitted a linear regression model on

$$Y = \begin{pmatrix} y_1 \\ \vdots \\ y_n \end{pmatrix}, \text{ and } Z = \begin{bmatrix} 1 & z_{11} & \dots & z_{1r} \\ 1 & z_{21} & \dots & z_{2r} \\ \vdots & \vdots & \ddots & \vdots \\ 1 & z_{n1} & \dots & z_{nr} \end{bmatrix}$$

For a future observations with known predictor variables.

(say $Z_0 = [1 z_{01} \dots z_{0r}]'$) how to obtain point and interval estimates of the corresponding (1) response and (2) mean response.

Predictor.

$$\tilde{Z}_0 = \begin{bmatrix} 1 \\ z_{01} \\ \vdots \\ z_{0r} \end{bmatrix}$$

If y_0 is the corresponding response, then

$$y_0 = \beta_0 + \beta_1 z_{01} + \dots + \beta_r z_{0r} + \varepsilon_0. \quad \tilde{z}_0' \tilde{\beta} + \varepsilon_0 = \\ = \tilde{z}_0' \tilde{\beta} + \varepsilon_0. \quad (1 \ z_{01} \ \dots \ z_{0r}) \begin{pmatrix} \beta_0 \\ \vdots \\ \beta_r \end{pmatrix} + \varepsilon_0$$

$\mathbb{E}[y_0] = \tilde{z}_0' \tilde{\beta}$. This is a parameter. We estimate this by

$\tilde{z}_0' \hat{\beta}$ (This is the BLUE). CI for $\tilde{z}_0' \tilde{\beta}$ is

$$\tilde{z}_0' \hat{\beta} \pm t_{n-r-1} \left(\frac{\alpha}{2}\right) \sqrt{s^2 \tilde{z}_0' (\tilde{z}' \tilde{z})^{-1} \tilde{z}_0}$$

this is for an average prediction.

To predict y_0 , the point estimate $\tilde{z}_0' \hat{\beta}$ is unbiased.

Interval estimate here is called Prediction Interval.

$$\text{var}(\hat{y}_0) = \text{var}(\tilde{z}_0' \hat{\beta}) + \sigma^2$$

Prediction interval : $\tilde{z}_0' \hat{\beta} \pm t_{n-r-1} \left(\frac{\alpha}{2}\right) \sqrt{s^2 + s^2 \tilde{z}_0' (\tilde{z}' \tilde{z})^{-1} \tilde{z}_0}$

Math B7800

2/22/18, Thur

Prediction Interval

Set up We have our regression model $y = \tilde{z}' \tilde{\beta} + \varepsilon$

number of parameters
↓
response predictors

We observe values of y and \tilde{z} for n observations.

We arrange them to have

$$\tilde{y} = \begin{bmatrix} y_1 \\ \vdots \\ y_n \end{bmatrix}, \quad \tilde{z} = \begin{bmatrix} \tilde{z}_1' \\ \vdots \\ \tilde{z}_n' \end{bmatrix} = \begin{bmatrix} 1 & z_{11} & \dots & z_{1r} \\ \vdots & \vdots & \ddots & \vdots \\ 1 & z_{n1} & \dots & z_{nr} \end{bmatrix}$$

Using these we estimate $\tilde{\beta}$: $\hat{\beta} = (\tilde{z}' \tilde{z})^{-1} \tilde{z}' \tilde{y}$

Now, suppose we have a future observation for which the predictor variables are known.

But, we do not know the value of the response.

$$y_0 = \tilde{z}_0' \tilde{\beta} + \varepsilon_0 \quad \text{error for future obs.}$$

↑
predictor
for future obs.
↓
response
for future obs.

$\mathbb{E}(y_0) = \tilde{z}_0' \tilde{\beta}$ This is a parameter, so we can find estimate and CI for it.
↑
mean responses.

Prediction Interval

Goal: to find an interval that will contain y_0 (response of a future observation) with a given confidence.

Consider $y_0 - \hat{z}_0' \hat{\beta}$ and find its distribution.

Recall, under normality assumption, $\hat{\beta} \sim N_{r+1}(\beta, \sigma^2(z'z)^{-1})$.

$$\text{Then } \underbrace{z_0' \hat{\beta}}_{\text{scalar}} \sim N(z_0' \beta, \sigma^2 z_0' (z'z)^{-1} z_0)$$

variance from ε_0 (future error)

$$y_0 \sim N(z_0' \beta, \sigma^2)$$

$$\text{so, } y_0 - z_0' \hat{\beta} \sim N(0, \sigma^2(1 + z_0' (z'z)^{-1} z_0))$$

$$\frac{y_0 - z_0' \hat{\beta}}{\sqrt{\sigma^2(1 + z_0' (z'z)^{-1} z_0)}} \sim N(0, 1)$$

$$\frac{y_0 - z_0' \hat{\beta}}{\sqrt{\sigma^2(1 + z_0' (z'z)^{-1} z_0)}} \sim t_{n-r-1}$$

$$\text{so, } P\left(\left|\frac{y_0 - z_0' \hat{\beta}}{\sqrt{\sigma^2(1 + z_0' (z'z)^{-1} z_0)}}\right| \leq t_{n-r-1}(\alpha/2)\right) = 1-\alpha.$$

so, 100(1-\alpha)% PI for y_0 .

$$\hat{z}_0' \hat{\beta} \pm \sqrt{\sigma^2(1 + z_0' (z'z)^{-1} z_0)} t_{n-r-1}(\alpha/2)$$

Model Verification*

Recall, $\hat{\varepsilon}_j = y_j - \hat{y}_j$

$$\hat{\varepsilon} = \begin{bmatrix} \hat{\varepsilon}_1 \\ \vdots \\ \hat{\varepsilon}_n \end{bmatrix} = y - \hat{y} = (I - P_z)y$$

$$\text{Var}(\hat{\varepsilon}) = \sigma^2(I - P_z)$$

Normalized error estimate

$$\hat{\varepsilon}_j^* = \frac{\varepsilon_j}{\sqrt{\text{Var}(\hat{\varepsilon}_j)}} = \frac{\varepsilon_j}{\sqrt{\sigma^2(I - P_z)_{jj}}} \quad \text{normalized variance.}$$

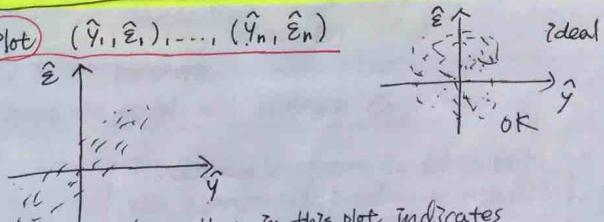
Residual Plot

If we plot the histogram of $\hat{\varepsilon}_1^*, \dots, \hat{\varepsilon}_n^*$, then we expect the histogram to be close to that of $N(0, 1)$.



Violations indicate violations of normality assumption on

Plot $(\hat{y}_1, \hat{\varepsilon}_1), \dots, (\hat{y}_n, \hat{\varepsilon}_n)$

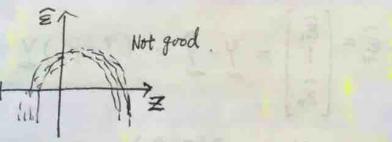


Any pattern in this plot indicates that the Lin. model assumption is not true.

Plot $\hat{\varepsilon}$ against $Z_{*2}, \dots, Z_{*(n+1)}$

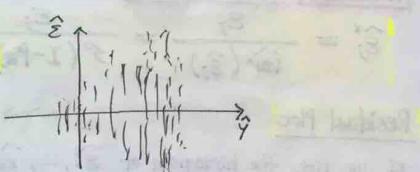
Expect no pattern if assumption are right.

But there is pattern, that indicates violation of the linear model namely corresponding Z variable needs to be transformed before using.



Violatility plot

In the plot $\hat{\varepsilon}$ vs \hat{X} , such a pattern violates the "equal variance" assumption.



Q-Q plot

Q-Q plot for the normalized residuals

To check "normality assumption".

Model selections

(Predictors) When the number of predictors is large, we need to select "suitable small" subset of predictors.

First, naive approach is to consider all subsets and choose the one which maximizes R^2 (RSS)

This doesn't work, as R^2 is an increasing function of r (# of predictors).

Adjusted R^2 , $\bar{R}^2 = 1 - (1-R^2) \frac{n-1}{n-r-1}$

One method is to use \bar{R}^2 to select model.

A "better" measure is Mallow's

$$C_p = \frac{\text{residual SS using } p \text{ predictors}}{\text{residual variances } (s^2)} - (n-2p)$$

$\left. \begin{array}{l} p = \# \text{ variables including } \\ \text{intercept} \\ n = \# \text{ observations} \end{array} \right\}$

C_p is not monotone.

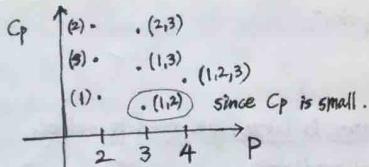
Choose the model with minimum C_p .

$$C_p = \left(\frac{\text{residual sum of squares for subset model with } p \text{ parameters, including an intercept}}{\text{residual variance for full model (residual SS)}} \right) - (n-2p)$$

$p = r+1$

A plot of the pairs (p, C_p) , one for each subset of predictors.

Ex Suppose $y = \beta_0 + \beta_1 z_1 + \beta_2 z_2 + \beta_3 z_3 + \epsilon$



These two methods work when # predictors is moderate, as we need to consider all possible subsets.

Next, we will see step-wise regression.

Math B7800

2/27/18, Tue

Predictor selection

Stepwise selection method

Recall

$$\text{If } \beta_2 = \begin{pmatrix} \beta_{(1)} \\ \beta_{(2)} \end{pmatrix} \xrightarrow{(q+1) \times 1} \xrightarrow{(r-q) \times 1}$$

Test $H_0: \beta_{(2)} = 0 \text{ vs } H_1: \beta_{(2)} \neq 0$

$$\Sigma = \left[\begin{array}{c|c} Z_{(1)} & Z_{(2)} \\ \hline n \times (q+1) & n \times (r-q) \end{array} \right] \quad \begin{array}{l} \text{if } H_0, \text{ exclude } Z_{(2)} \\ \text{if } H_1, \text{ we need to include } Z_{(2)} \end{array}$$

We use the LRT which boils down to

$$F = \frac{Y'(P_Z - P_{Z'}) Y \frac{1}{r-q}}{Y'(I - P_Z) Y \frac{1}{n-r-1}} \quad \begin{array}{l} \text{"We reject } H_0 \text{ at level } \alpha \text{ "} \\ \text{when } F > F_{r-q, n-r-1}(\alpha) \end{array}$$

We will refer to the model: $y = \beta_0 + \beta_1 z_1 + \dots + \beta_r z_r + \epsilon$ as "regression function"

Stepwise method

Step 1: Scan through all the variables (one at time)

and which contributes the most to the Regression SS.

$y = \beta_0 + \beta_1 z_1$ Find which one has R^2 value the most.

$$y = \beta_0 + \beta_2 z_2 \quad \vdots \quad Z = \begin{bmatrix} 1 & z_{11} \\ 1 & z_{12} \\ \vdots & \vdots \\ 1 & z_{nr} \end{bmatrix} \rightarrow \begin{array}{l} Y'(P_Z - P_1) Y \\ Y'(I - P_1) Y \end{array} \quad \begin{array}{l} \text{compute} \\ \text{* use AIC, R}^2, Cp \end{array}$$

$$y = \beta_0 + \beta_3 z_3 \quad Z = \begin{bmatrix} 1 & z_{11} \\ 1 & z_{13} \\ \vdots & \vdots \\ 1 & z_{nr} \end{bmatrix}$$

Suppose k is the one which maximizes

Step 2 Test $\beta = \begin{bmatrix} \beta_0 \\ \beta_k \end{bmatrix}$

Test $H_0: \beta_k = 0$ vs $H_1: \beta_k \neq 0$

$$Z = \begin{bmatrix} 1 & z_{1k} \\ \vdots & \vdots \\ 1 & z_{nk} \end{bmatrix}, Z_1 = \begin{bmatrix} 1 \\ \vdots \\ 1 \end{bmatrix}$$

$$F = \frac{Y'(P_Z - P_{Z_1})Y \frac{1}{1-\alpha}}{Y'(I-P_Z)Y \frac{1}{n-1}} > F_{1, n-2}(\alpha) \quad (*)$$

reject $H_0 \Rightarrow \beta_k \neq 0$

then include k in the model.

If $(*)$ holds, then "include" Z_k in the regression function.

Current model is $Y = \beta_0 + \beta_k Z_k + \epsilon$.

Step 3 We scan through the remaining variables (one at time) and find out which one contributes to the Reg. SS the most. (keeping the existing regression function)

$$\text{Current model: } Y = \beta_0 + \beta_k Z_k + \beta_l Z_l \rightarrow Z = \begin{bmatrix} 1 & z_{1k} & z_{1l} \\ \vdots & \vdots & \vdots \\ 1 & z_{nk} & z_{nl} \end{bmatrix}, Z_1 = \begin{bmatrix} 1 & z_{1k} \\ \vdots & \vdots \\ 1 & z_{nk} \end{bmatrix}$$

$$Y = \beta_0 + \beta_k Z_k + \beta_{k1} Z_{k1}$$

$$Y = \beta_0 + \beta_k Z_k + \dots + \beta_{kh} Z_{kh} \rightarrow Z = \begin{bmatrix} 1 & z_{1k} & z_{1h} \\ \vdots & \vdots & \vdots \\ 1 & z_{nk} & z_{nh} \end{bmatrix}$$

$$Y = \beta_0 + \beta_k Z_k + \beta_r Z_r$$

Find out which model has "max value" of $Y'(P_Z - P_{Z_1})Y$.

Suppose L is the one that maximizes: Reg. SS

Then $\beta = \begin{bmatrix} \beta_0 \\ \beta_k \\ \beta_L \end{bmatrix}$ $H_0: \beta_L = 0$ vs $H_1: \beta_L \neq 0$.

$$F = \frac{Y'(P_Z - P_{Z_1})Y \frac{1}{2-1}}{Y'(I-P_Z)Y \frac{1}{n-2-1}}$$

$$Z = \begin{bmatrix} 1 & z_{1k} & z_{1L} \\ \vdots & \vdots & \vdots \\ 1 & z_{nk} & z_{nL} \end{bmatrix}, Z_1 = \begin{bmatrix} 1 & z_{1k} \\ \vdots & \vdots \\ 1 & z_{nk} \end{bmatrix}$$

Include Z_L if $F > F_{1, n-3}(\alpha)$

If included, then current model is

$$Y = \beta_0 + \beta_k Z_k + \beta_L Z_L + \epsilon$$

Step 4 When a new variable is included, we need to check whether one of the existing ones fail the F test or not.

Suppose at a stage, we include Z_{im} to $(Z_{i1}, Z_{i2}, \dots, Z_{im-1})$

then test:

$$\beta_0 = \begin{bmatrix} \beta_0 \\ \beta_{i1} \\ \vdots \\ \beta_{im} \end{bmatrix}$$

$$H_0: \beta_{ii} = 0 \text{ vs } H_1: \beta_{ii} \neq 0 \quad Z_k \quad Z_L \dots$$

$$\vdots \quad \beta_k = 0 \quad \beta_L = 0$$

$$H_0: \beta_{im} = 0 \text{ vs } H_1: \beta_{im} \neq 0$$

"Repeat Steps 3,4"

If in one of the tests, the H_0 is accepted, then the corresponding variable goes out. $\beta_k = 0$

At the end, we will have a subset $\{i_1, \dots, i_p\}$ of $\{1, \dots, r\}$ such that no additional variable can be included, (i.e., F-test is insignificant) and no existing variable leaves (i.e., the corresponding F test is significant).
accept $H_0 \rightarrow$ exclude the corresponding variable.

Then we stop.

"don't drop."
 to reject H_0
 ⇒ include variable.
 the corresponding.

$$AIC = n \ln \left(\frac{\text{residual SS for subset model with } p \text{ parameters, including an intercept}}{n} \right) + 2p$$

overall, we want to select models from those having the smaller values of AIC.

$$\begin{pmatrix} 2 & 3 \\ 1 & 1 \\ 2 & 1 \\ 1 & 3 \end{pmatrix}$$

Information Criterion Method

The famous information criterion assigned to subsets of variables is the Akaike's Information Criterion (AIC).

For a subset of size p ,

$$AIC = n \log \left(\frac{\text{Residual SS for that subset}}{n} \right) + 2p$$

↓ decreasing ↑ increasing

one needs to minimize AIC to select variables.

↳ Multicollinearity

So far we have assumed that Z has full rank.

then $Z'Z$ was nonsingular.

If this condition does not hold, then we say that the data has Multicollinearity.

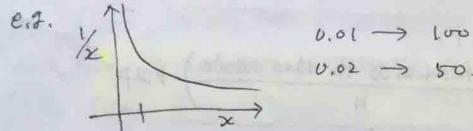
In some cases, $Z\alpha = \zeta$ for some $\alpha \neq 0$.

In this case, $Z'Z$ is singular.

In many cases, $Z'Z$ is "nearly singular", meaning it has some "very small" eigenvalues.

In the first case, we choose a linearly independent subset of columns of Z .

In the second case, $(Z'Z)^{-1}$ becomes numerically unstable.



Then remove one of the highly correlated columns of Z .

3/11/18, Thur

② Multivariate "Multiple" Regression.

For each individual, we have a vector of response ($1 \times m$)

The model

$$Y^t = [y_1, \dots, y_m] = \underbrace{[1, z_1, \dots, z_r]}_{1 \times (r+1)} \begin{bmatrix} \beta_{01} & \dots & \beta_{0m} \\ \vdots & \ddots & \vdots \\ \beta_{r1} & \dots & \beta_{rm} \end{bmatrix} + \varepsilon^t$$

$1 \times m$

$$\mathbb{E}[\varepsilon] = \underset{m \times 1}{0}, \quad \mathbb{E}[\varepsilon \varepsilon^t] = \Sigma_{m \times m}$$

We observe the responses and predictors from n individuals.

values of

We organize them:

$$Y^t = \underbrace{\tilde{z}^t \beta + \varepsilon^t}_{\vdots} \quad Y = \begin{bmatrix} Y^1 \\ \vdots \\ Y^n \end{bmatrix}_{n \times m}, \quad \tilde{Z} = \begin{bmatrix} \tilde{z}^1 \\ \vdots \\ \tilde{z}^n \end{bmatrix} = \begin{bmatrix} 1 & z_{11} & \dots & z_{1r} \\ \vdots & \vdots & \ddots & \vdots \\ 1 & z_{n1} & \dots & z_{nr} \end{bmatrix}_{n \times (r+1)}$$

$$\beta = \begin{bmatrix} \beta_{01} & \dots & \beta_{0m} \\ \vdots & \ddots & \vdots \\ \beta_{r1} & \dots & \beta_{rm} \end{bmatrix}_{(r+1) \times m}, \quad \varepsilon = \begin{bmatrix} \varepsilon^1 \\ \vdots \\ \varepsilon^n \end{bmatrix}_{n \times m}$$

Multivariate multiple regression model: $\underline{Y = Z\beta + \varepsilon}$

$$= \begin{bmatrix} \varepsilon_{11} & \dots & \varepsilon_{1m} \\ \vdots & \ddots & \vdots \\ \varepsilon_{n1} & \dots & \varepsilon_{nm} \end{bmatrix}$$

Each observation has m -response variables which are correlated.

$$\underset{n \times m}{Y'} = \underset{1 \times (n+1)}{\underbrace{Z' \beta + \varepsilon'}} \quad (\text{the model}) \quad E[\varepsilon'] = \Omega'$$

$$\text{COV}(\varepsilon') = \sum_{m \times m}$$

$$= \begin{bmatrix} \sigma_{11} & \dots & \sigma_{1m} \\ \vdots & \ddots & \vdots \\ \sigma_{m1} & \dots & \sigma_{mm} \end{bmatrix}$$

" n = # responses for each individual observation"

We observe the values of the predictors and responses from n individuals.

We arrange them into $\checkmark Y = \begin{bmatrix} Y'_1 \\ \vdots \\ Y'_n \end{bmatrix} = [y_{11}, \dots, y_{1m}]$

$$Z = \begin{bmatrix} Z'_1 \\ \vdots \\ Z'_n \end{bmatrix} = \begin{bmatrix} 1 & z_{11} & \dots & z_{1r} \\ \vdots & \vdots & \ddots & \vdots \\ 1 & z_{n1} & \dots & z_{nr} \end{bmatrix}, \quad \beta = \begin{bmatrix} \beta_{01} & \dots & \beta_{0m} \\ \vdots & \ddots & \vdots \\ \beta_{r1} & \dots & \beta_{rm} \end{bmatrix}$$

$$\varepsilon = \begin{bmatrix} \varepsilon'_1 \\ \vdots \\ \varepsilon'_n \end{bmatrix} = [\varepsilon_{11}, \dots, \varepsilon_{1m}]$$

$\varepsilon'_1, \dots, \varepsilon'_n$ are independent.

$$Y = Z\beta + \varepsilon$$

$$Y = [Y_{(1)} | \dots | Y_{(m)}], \quad \beta = [\beta_{(1)} | \dots | \beta_{(m)}]$$

$$\begin{bmatrix} Y_{11} \\ \vdots \\ Y_{1r} \end{bmatrix} \quad \begin{bmatrix} Y_{1m} \\ \vdots \\ Y_{nm} \end{bmatrix}$$

$$\varepsilon = [\varepsilon_{(1)} | \dots | \varepsilon_{(m)}]$$

$$\begin{bmatrix} \varepsilon_{11} \\ \vdots \\ \varepsilon_{1r} \end{bmatrix} \quad \begin{bmatrix} \varepsilon_{1m} \\ \vdots \\ \varepsilon_{nm} \end{bmatrix}$$

$$\text{So, } [Y_{(1)} | \dots | Y_{(m)}] = Z[\beta_{(1)} | \dots | \beta_{(m)}] + [\varepsilon_{(1)} | \dots | \varepsilon_{(m)}]$$

$$\text{So, } Y_{(i)} = Z\beta_{(i)} + \varepsilon_{(i)}, \text{ for } i = 1, 2, \dots, m.$$

*goal: Estimate β , i.e., $\hat{\beta}$.

one method is to estimate $\beta_{(i)}$ by $\hat{\beta}_{(i)}$ (the least square estimate from the univariate case).

$$\hat{\beta}_{(i)} = (Z'Z)^{-1}Z'Y_{(i)}$$

$$\text{Using this, } \hat{\beta} = [\hat{\beta}_{(1)} | \dots | \hat{\beta}_{(m)}]$$

$$= (Z'Z)^{-1}Z' [Y_{(1)} | \dots | Y_{(m)}]$$

$$= (Z'Z)^{-1}Z' Y_{n \times m}$$

$$\text{For this } \hat{\beta}, \quad \hat{Y} = Z\hat{\beta} = Z(Z'Z)^{-1}Z'Y = P_Z Y.$$

$$\text{and } \hat{\varepsilon} = Y - \hat{Y} = (I - P_Z)Y$$

$$P_Z Y \quad Z(Z'Z)^{-1}Z'$$

least square The error for an estimator B of $\hat{\beta}$, defined to least square

$$f(B) = \text{tr} \left[\underbrace{(Y - ZB)'(Y - ZB)}_{m \times n} \right]_{n \times m}$$

$$A = [A_{(1)} | \dots | A_{(m)}]$$

$$(A'A)_{i,j}$$

Fact $\hat{\beta}$ minimizes $f(\beta)$, i.e.,

$$f(\hat{\beta}) \leq f(B) \text{ for any } B.$$

$$f(B) = \sum_{i=1}^m (Y_{(i)} - Z\beta_{(i)})'(Y_{(i)} - Z\beta_{(i)})$$

The i th summand is minimized at $\beta_{(i)} = \hat{\beta}_{(i)}$ (from our knowledge about the univariate case).

$$\text{so, } f(B) \geq \sum_{i=1}^m (Y_{(i)} - Z\hat{\beta}_{(i)})'(Y_{(i)} - Z\hat{\beta}_{(i)}) = f(\hat{\beta})$$

Now, suppose we consider a different least square error

$$g(B) = |(Y - ZB)'(Y - ZB)|$$

Fact (from Lin. Alg.)

If C and D are symmetric non-negative definite (nnnd) matrices of same size, then

determinate = product of eigenvalues.

case 1:
determinate
not absolute value.

proof If C and D are both positive semi-definite (nnnd but not pd) then $|C| = |D| = 0$ and $|C+D| \geq 0$ (some eigenvalue is zero)

case 2: one of the matrices is pd. say D is pd.

$$\text{Then } C+D = D^{1/2}(D^{-1/2}CD^{-1/2})D^{1/2} + D. \quad \text{For any nnnd matrix } A, \\ |I+A| \geq 1+|A|$$

$$= D^{1/2}(D^{-1/2}CD^{-1/2} + I)D^{1/2}$$

$$\text{so, } |C+D| = |D^{1/2}| \underbrace{|I+D^{-1/2}CD^{-1/2}|}_{\geq 1+|D^{-1/2}CD^{-1/2}|} \cdot |D|^{1/2} \geq |D^{1/2}|(1+|D^{-1/2}CD^{-1/2}|)|D|^{1/2} \\ = |D| + |C|$$

$$\text{write } Y - ZB = (Y - Z\hat{\beta}) + Z(\hat{\beta} - B).$$

$$\text{so, } (Y - ZB)'(Y - ZB)$$

$$= (Y - Z\hat{\beta})'(Y - Z\hat{\beta}) + (Z(\hat{\beta} - B))'Z(\hat{\beta} - B)$$

$$+ (Y - Z\hat{\beta})'Z(\hat{\beta} - B) + [Z(\hat{\beta} - B)]'(Y - Z\hat{\beta})$$

$$= 0$$

$$\nwarrow \text{ by } \hat{\beta} = (Z'Z)^{-1}Z'Y \Leftrightarrow (Z'Z)\hat{\beta} = Z'Y$$

$$(\hat{\beta} - B)' \underbrace{Z' (Y - Z\hat{\beta})}_{= 0}$$

$$\text{so, } g(B) = |(Y - ZB)'(Y - ZB)|$$

$$= |(Y - Z\hat{\beta})'(Y - Z\hat{\beta}) + (Z\hat{\beta} - ZB)'(Z\hat{\beta} - ZB)|$$

$$\geq |(Y - Z\hat{\beta})'(Y - Z\hat{\beta})| + |(Z\hat{\beta} - ZB)'(Z\hat{\beta} - ZB)|$$

$$\text{so, } g(B) \geq g(\hat{\beta}) + |(\hat{\beta} - B)'Z'Z(\hat{\beta} - B)|$$

$$\geq g(\hat{\beta}) \quad (\text{as } A'A \text{ is nnnd for any } A).$$

$$\Rightarrow |A'A| \text{ is non-negative. unknown}$$

Under the additional normality assumption, i.e., $Z \sim N_m(0, \Sigma)$. We want to find MLE for β, Σ .

Def of MLE, Likelihood func --

The likelihood function

$$L(\beta, \Sigma) = \text{joint density of } Y's$$

$$= \prod_{i=1}^n (\text{pdf of } Y_i)$$

Recall $Y' = Z'\beta + \varepsilon'$

If $\varepsilon' \sim N_m(0, \Sigma)$, then $Y' \sim N_m(Z'\beta, \Sigma)$.

pdf of Y : $f_Y(Y) = (2\pi)^{-\frac{n}{2}} |\Sigma|^{-\frac{1}{2}} \exp\left(-\frac{1}{2} (Y - Z'\beta)' \Sigma^{-1} (Y - Z'\beta)\right)$

$$\Rightarrow \prod_{i=1}^n f_{Y_i}(y_i) = (2\pi)^{-\frac{nm}{2}} |\Sigma|^{-\frac{n}{2}} \exp\left(-\frac{1}{2} \sum_{i=1}^n (y_i - Z'\beta)' \Sigma^{-1} (y_i - Z'\beta)\right)$$

$$\text{So, } L(\beta, \Sigma) = (2\pi)^{-\frac{nm}{2}} |\Sigma|^{-\frac{n}{2}} \exp\left(-\frac{1}{2} \text{tr}\left[\sum_{i=1}^n (y_i - Z'\beta)' \Sigma^{-1} (y_i - Z'\beta)\right]\right)$$

Fact \downarrow by $\text{tr}(CD) = \text{tr}(DC)$

$$= (2\pi)^{-\frac{nm}{2}} |\Sigma|^{-\frac{n}{2}} \exp\left(-\frac{1}{2} \text{tr}[\Sigma^{-1} (Y - Z\beta)' (Y - Z\beta)]\right).$$

We will maximize in 2 steps.

First, we wrt Σ , then wrt β .

Step 1 β is fixed. call $(Y - Z\beta)' (Y - Z\beta) = D$.

define $h(\Sigma) = |\Sigma|^{-\frac{n}{2}} \exp(-\frac{1}{2} \text{tr}(\Sigma^{-1} D))$

$$\text{tr}(\Sigma^{-1} D) = \text{tr}(\Sigma^{-\frac{1}{2}} D \Sigma^{-\frac{1}{2}}) = \sum_{i=1}^m \lambda_i, \text{ where } \lambda_1, \dots, \lambda_m \text{ are}$$

the eigenvalues of $\Sigma^{-\frac{1}{2}} D \Sigma^{-\frac{1}{2}}$.

Then, $\prod_{i=1}^m \lambda_i = |\Sigma^{-\frac{1}{2}} D \Sigma^{-\frac{1}{2}}| = \frac{|D|}{|\Sigma|}$

$$\text{then, } h(\Sigma) = \frac{1}{|D|^{\frac{n}{2}}} \left(\frac{|D|}{|\Sigma|}\right)^{\frac{n}{2}} \exp\left(-\frac{1}{2} \sum_{i=1}^m \lambda_i\right)$$

$$= \frac{1}{|D|^{\frac{n}{2}}} \prod_{i=1}^m \lambda_i e^{-\frac{1}{2} \lambda_i}$$

A is nnd if $\underline{x}' A \underline{x} \geq 0$ for all \underline{x}

If A is symmetric, then nnd \Leftrightarrow all eigen values ≥ 0

A is p.d. if $\underline{x}' A \underline{x} > 0$ unless $\underline{x} = 0$.

If A is symmetric, then p.d. \Leftrightarrow all eig values > 0 .

positive semidefinite = $\{\text{nnd}\} - \{\text{p.d.}\}$.

Recall multivariate multiple linear regression.

model $\underline{Y} = \underline{Z} \beta + \varepsilon$ — $n \times m$

$$\begin{matrix} n \times m \\ 1 \\ n \times (r+1) \\ (r+1) \times m \end{matrix}$$

$$\hat{\beta} = (\underline{Z}' \underline{Z})^{-1} \underline{Z}' \underline{Y}$$

$$\hat{\Sigma} = \frac{1}{n} \hat{\varepsilon}' \hat{\varepsilon}$$

Result If the rows of ε are iid $N_m(0, \Sigma)$, then the MLE of β and

Σ will be $\hat{\beta} = (\underline{Z}' \underline{Z})^{-1} \underline{Z}' \underline{Y}$, $\hat{\Sigma} = \frac{1}{n} \hat{\varepsilon}' \hat{\varepsilon}$, where

$$\hat{\varepsilon} = \underline{Y} - \hat{\underline{Y}} = \underline{Y} - \underline{Z} \hat{\beta} = \underline{Y} - \underline{Z} (\underline{Z}' \underline{Z})^{-1} \underline{Z}' \underline{Y} = (\underline{I} - \underline{P}_Z) \underline{Y}$$

Proof Recall the likelihood function

$$L(\beta, \Sigma) = (2\pi)^{-\frac{nm}{2}} |\Sigma|^{-\frac{n}{2}} \exp\left(-\frac{1}{2} \text{tr}(\Sigma' D)\right), \text{ where}$$

$$D = (\underline{Y} - \underline{Z} \beta)' (\underline{Y} - \underline{Z} \beta)$$

For fixed β (so, D is also fixed), we will maximize

$$f(\Sigma) = |\Sigma|^{-\frac{n}{2}} \exp\left(-\frac{1}{2} \text{tr}(\Sigma' D)\right).$$

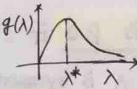
$$= \frac{1}{|D|^{\frac{n}{2}}} \left(\frac{|D|}{|\Sigma|} \right)^{\frac{n}{2}} \exp\left[-\frac{1}{2} \text{tr}(\Sigma'^{-\frac{1}{2}} D \Sigma'^{-\frac{1}{2}})\right]$$

symmetric, then eigenvalues are real value.

If the eigenvalues of $\Sigma^{-1/2} D \Sigma^{-1/2}$ are $\lambda_1, \dots, \lambda_m$,
then $\text{tr}(\Sigma^{-1/2} D \Sigma^{-1/2}) = \sum_{i=1}^m \lambda_i$

$$\prod_{i=1}^m \lambda_i = \left| \sum^{-1/2} D \Sigma^{-1/2} \right| = \left| \Sigma^{-1/2} \right| |D| \left| \Sigma^{-1/2} \right| = \frac{|D|}{|\Sigma|}$$

Then $f(\Sigma) = \frac{1}{|D|^{n/2}} \left(\prod_{i=1}^m \lambda_i \right)^{n/2} \exp\left(-\frac{1}{2} \sum_{i=1}^m \lambda_i\right)$



$$= \frac{1}{|D|^{n/2}} \prod_{i=1}^m g(\lambda_i), \text{ where } g(\lambda) = \lambda^{n/2} e^{-\lambda/2}$$

We will maximize g w.r.t. λ . If λ^* maximizes g , then $f(\Sigma)$ will be maximized by taking $\lambda_i = \lambda^*$ for all $i = 1, 2, \dots, m$.

$$\frac{d}{d\lambda} \log(g(\lambda)) = \frac{n}{2\lambda} - \frac{1}{2}.$$

$$\text{so, } \frac{d}{d\lambda} (\log g(\lambda)) = 0 \iff \lambda = n, \text{ so } \lambda^* = n.$$

To find maximizer of $f(\Sigma)$, we need to find Σ so that all eigenvalues of $\Sigma^{-1/2} D \Sigma^{-1/2}$ are equal to n .

Thus, we must have $\Sigma^{-1/2} D \Sigma^{-1/2} = n I \Rightarrow n \Sigma = D$

$$\Sigma = \frac{1}{n} D$$

$$\text{so, } f(\Sigma) \leq f\left(\frac{1}{n} D\right).$$

$$\text{so, } L(\beta, \Sigma) \leq L(\beta, \frac{1}{n} D) = L\left(\beta, \frac{1}{n} (\gamma - z\beta)' (\gamma - z\beta)\right)$$

$$= (2\pi)^{-\frac{nm}{2}} |\frac{1}{n} D|^{-\frac{nm}{2}} \exp\left[-\frac{1}{2} \text{tr}\left((\frac{1}{n} D)^{-1} D\right)\right]$$

$$= (2\pi)^{-\frac{nm}{2}} |D|^{-\frac{nm}{2}} n^{n/2} e^{-\frac{mn}{2}}$$

because we have already seen that $\beta = \hat{\beta}$ minimizes $h(\beta) = |(\gamma - z\beta)' (\gamma - z\beta)|$

$$\max_{\beta, \Sigma} L(\beta, \Sigma) = C \left| \hat{\Sigma} \right|^{-\frac{nm}{2}}, \quad C \text{ is constant}$$

$$(2\pi)^{-\frac{nm}{2}} n^{n/2} e^{-\frac{mn}{2}}$$

properties of $\hat{\beta}$.

$$\mathbb{E}[\hat{\beta}] = \mathbb{E}\left[(z'z)^{-1} z' \gamma\right] = (z'z)^{-1} z' \mathbb{E}(\gamma) = (z'z)^{-1} z' z \beta = \beta.$$

$$\text{Recall } \gamma = [Y_{(1)} | \dots | Y_{(m)}] = z\beta + \varepsilon.$$

$$\mathbb{E}[\hat{\beta}_{(i)}] = \beta_{(i)}, \text{ for } i=1, \dots, m = z_{(i)} \beta_{(i)} + [\varepsilon_{(1)} | \dots | \varepsilon_{(m)}]$$

$$\text{cov}(\hat{\beta}_{(j)}, \hat{\beta}_{(k)}) =$$

$$\text{cov}((z'z)^{-1} z' Y_{(j)}, (z'z)^{-1} z' Y_{(k)})$$

$$= (z'z)^{-1} z' \text{cov}(Y_{(j)}, Y_{(k)}) z (z'z)^{-1}$$

$$= \text{cov}(\varepsilon_{(j)}, \varepsilon_{(k)})$$

$$= \sigma_{jk} I_{n \times n}$$

$$\Sigma = \begin{bmatrix} \varepsilon_{11} & \dots & \varepsilon_{1m} \\ \vdots & & \vdots \\ \varepsilon_{m1} & \dots & \varepsilon_{mm} \end{bmatrix}$$

$m \times m$
of responses.

$$\varepsilon_{(j)} = \begin{bmatrix} \varepsilon_{j1} \\ \vdots \\ \varepsilon_{jn} \end{bmatrix}, \quad \varepsilon_{(k)} = \begin{bmatrix} \varepsilon_{k1} \\ \vdots \\ \varepsilon_{kn} \end{bmatrix}$$

$$= (z'z)^{-1} z' (\sigma_{jk} I_n) z (z'z)^{-1}$$

$$= \sigma_{jk} (z'z)^{-1}$$

$$\text{Also, } \hat{\gamma}' \hat{\Sigma} = (\mathbf{P}_Z \mathbf{Y})' (\mathbf{I} - \mathbf{P}_Z) \mathbf{Y} = \mathbf{Y}' \underbrace{(\mathbf{I} - \mathbf{P}_Z)}_{=0} \mathbf{Y}$$

$$\star \text{cov}(\hat{\gamma}_{(j)}, \hat{\gamma}_{(k)}) = 0 \text{ for } j, k.$$

$$= \text{cov}(\mathbf{P}_Z \mathbf{Y}_{(j)}, (\mathbf{I} - \mathbf{P}_Z) \mathbf{Y}_{(k)})$$

$$= \mathbf{P}_Z \text{cov}(\mathbf{Y}_{(j)}, \mathbf{Y}_{(k)}) (\mathbf{I} - \mathbf{P}_Z) = \mathbf{P}_Z (\sigma_{jk} \mathbf{I}) (\mathbf{I} - \mathbf{P}_Z) = 0.$$

$$\text{Recall } \hat{\Sigma} = \frac{1}{n} \hat{\epsilon}' \hat{\epsilon} = \frac{1}{n} (\mathbf{Y} - \mathbf{Z}\hat{\beta})' (\mathbf{Y} - \mathbf{Z}\hat{\beta}).$$

$$\mathbb{E}[\hat{\Sigma}] =$$

to check
unbiased

$$\text{Find } \mathbb{E}[\hat{\epsilon}_{(j)}' \hat{\epsilon}_{(k)}]$$

$$= \mathbb{E}[\mathbf{Y}_{(j)}' (\mathbf{I} - \mathbf{P}_Z)^2 \mathbf{Y}_{(k)}]$$

$$= \mathbb{E}[\mathbf{Y}_{(j)}' (\mathbf{I} - \mathbf{P}_Z) \mathbf{Y}_{(k)}] = \mathbb{E}[\text{tr}(\mathbf{Y}_{(j)}' (\mathbf{I} - \mathbf{P}_Z) \mathbf{Y}_{(k)})]$$

$$= \mathbb{E}[\text{tr}((\mathbf{I} - \mathbf{P}_Z) \mathbf{Y}_{(k)} \mathbf{Y}_{(j)}')] = \text{tr} \mathbb{E}[(\mathbf{I} - \mathbf{P}_Z) \mathbf{Y}_{(k)} \mathbf{Y}_{(j)}']$$

$$= \text{tr}((\mathbf{I} - \mathbf{P}_Z) \{ \text{cov}(\mathbf{Y}_{(k)}, \mathbf{Y}_{(j)}) + \mathbf{Z}\beta_{(k)} \beta_{(j)}' \mathbf{Z}' \})$$

$$= \text{tr}[(\mathbf{I} - \mathbf{P}_Z) \sigma_{jk} \mathbf{I} + 0] \quad \text{, } \text{tr}(\mathbf{I} - \mathbf{P}_Z)$$

$$= \sigma_{jk} \text{tr}(\mathbf{I} - \mathbf{P}_Z) = \sigma_{jk} \underset{\substack{\text{idempotent} \\ \text{matrix}}}{\text{rank}}(\mathbf{I} - \mathbf{P}_Z) = \sigma_{jk}(n-r-1)$$

$$\text{So, } \mathbb{E}[\hat{\epsilon}' \hat{\epsilon}] = (n-r-1) \Sigma$$

So, unbiased estimator of Σ is

$$\star S = \frac{1}{n-r-1} \hat{\epsilon}' \hat{\epsilon} = \frac{1}{n-r-1} (\mathbf{Y} - \mathbf{Z}\hat{\beta})' (\mathbf{Y} - \mathbf{Z}\hat{\beta})$$

$$= \frac{n}{n-r-1} \hat{\Sigma} = \frac{1}{n-r-1} \mathbf{Y}' (\mathbf{I} - \mathbf{P}_Z) \mathbf{Y}.$$

Also, $n \hat{\Sigma} = \hat{\epsilon}' \hat{\epsilon} \sim W_{m, n-r-1}(\Sigma)$ "Wishart Distribution"

Also, $\hat{\Sigma}$ and $\hat{\beta}$ are independent.

under normality assumption,

$$\hat{\beta}_{(j)} \sim N_{r+1}(\beta_{(j)}, \sigma_{jj}(\mathbf{Z}' \mathbf{Z})^{-1})$$

3/29 ~ No class

3/13/18, Tue

Math B7800

Likelihood Ratio Test

$$\beta_{(r+1) \times m} = \begin{bmatrix} \beta_1 \\ \beta_2 \end{bmatrix}_{(r+1) \times m}^{(q+1) \times m}$$

* Test: $H_0: \beta_2 = 0$ vs $H_1: \beta_2 \neq 0$

Partition Z matrix accordingly: $Z_{n \times (r+1)} = [Z_1 | Z_2]$

$$E[Y] = Z\beta = [Z_1 | Z_2] \begin{bmatrix} \beta_1 \\ \beta_2 \end{bmatrix} = Z_1\beta_1 + Z_2\beta_2$$

"Under $H_0: \beta_2 = 0$ "

$E[Y] = Z_1\beta_1$ ($H_0 \Rightarrow$ the predictors in Z_2 are redundant.)

Recall Likelihood Ratio Test.

$$\Delta = \frac{\max_{H_0} L(\beta, \Sigma)}{\max L(\beta, \Sigma)} \text{ "the likelihood ratio"}$$

so, $L(\beta, \Sigma)$ is maximized at "MLE of β, Σ " $\hat{\beta}, \hat{\Sigma}$

$$\hat{\beta} = \hat{\beta} = (Z'Z)^{-1}Z'Y \quad \text{and} \quad \hat{\Sigma} = \frac{1}{n}(Y - Z\hat{\beta})(Y - Z\hat{\beta})'$$

$$\text{and } L(\hat{\beta}, \hat{\Sigma}) = (2\pi)^{-\frac{n}{2}} |\hat{\Sigma}|^{-\frac{1}{2}} e^{-\frac{1}{2} \hat{\Sigma}^{-1} (Y - Z\hat{\beta})^T (Y - Z\hat{\beta})} = \max L(\beta, \Sigma)$$

under H_0 , the model becomes $Y = Z_1\beta_1 + \varepsilon$

so, under H_0 , $L(\beta, \Sigma)$ is maximized at

$$\hat{\beta} = \left[\frac{(Z_1'Z_1)^{-1}Z_1'Y}{0} \right], \quad \hat{\Sigma}_1 = \frac{1}{n}(Y - Z_1\hat{\beta})'(Y - Z_1\hat{\beta})$$

matlab

$$|\hat{\Sigma}| = \det(\hat{\Sigma})$$

$$e^{-\frac{nm}{2}} = \exp\left(-\frac{nm}{2}\right)$$

$$-\begin{bmatrix} 1 & 2 \\ 3 & 4 \\ 0 & 0 \\ 0 & 0 \end{bmatrix}$$

$$L([\hat{\beta}_0], \hat{\Sigma}_1) = (2\pi)^{-\frac{m}{2}} |\hat{\Sigma}_1|^{-\frac{m}{2}} e^{-\frac{\text{tr}(\hat{\Sigma}_1)}{2}}$$

$$\text{so, } \Delta_1 = \frac{L([\hat{\beta}_0], \hat{\Sigma}_1)}{L(\hat{\beta}_0, \hat{\Sigma})} = \left(\frac{|\hat{\Sigma}|}{|\hat{\Sigma}_1|} \right)^{\frac{m}{2}}$$

so, $\Delta_1 = \frac{|\hat{\Sigma}|}{|\hat{\Sigma}_1|}$, We will reject H_0 , when Δ_1 is small or equivalently $\frac{|\hat{\Sigma}|}{|\hat{\Sigma}_1|}$ is small.

$$U = -[n-r-1-\frac{1}{2}(m-r+q+1)] \ln \left(\frac{|\hat{\Sigma}|}{|\hat{\Sigma}_1|} \right) \sim \chi^2_{m(r-q)}$$

Reject H_0 at 5% level if $U > \chi^2_{m(r-q)}$.

$$\text{Note: } \frac{|\hat{\Sigma}|}{|\hat{\Sigma}_1|} = \frac{|\hat{\Sigma}|}{|\hat{\Sigma} + (\hat{\Sigma}_1 - \hat{\Sigma})|}$$

since $(\hat{\Sigma}_1 - \hat{\Sigma})$ and $\hat{\Sigma}_1$ are independent (we proved this only in the univariate case), we need to compare $\hat{\Sigma}_1$ and $\hat{\Sigma}_1 - \hat{\Sigma}$, in particular, consider the ratio matrix: $(\hat{\Sigma}_1 - \hat{\Sigma})(\hat{\Sigma}_1)^{-1}$.

Let the positive eigenvalues of $(\hat{\Sigma}_1 - \hat{\Sigma})(\hat{\Sigma}_1)^{-1}$ are $\eta_1 \geq \eta_2 \geq \dots \geq \eta_s > 0$.

$$\text{• Wilks' lambda: } \prod_{i=1}^s \frac{1}{1+n_i} = \frac{|\hat{\Sigma}|}{|\hat{\Sigma} + (\hat{\Sigma}_1 - \hat{\Sigma})|} = \frac{|E|}{|E + H|} \quad \begin{cases} E = \hat{\Sigma} \\ H = \hat{\Sigma}_1 - \hat{\Sigma} \end{cases} \quad (1)$$

$$\text{• Pillai's trace: } \sum_{i=1}^s \frac{\eta_i}{1+n_i} = \text{tr}((\hat{\Sigma}_1 - \hat{\Sigma})((\hat{\Sigma}_1 - \hat{\Sigma}) + \hat{\Sigma})^{-1}) \quad (2)$$

$$= \text{tr}((\hat{\Sigma}_1 - \hat{\Sigma})(\hat{\Sigma}_1)^{-1}) \quad (L_2)$$

$$\text{• Hotelling's trace: } \sum_{i=1}^s \eta_i = \text{tr}((\hat{\Sigma}_1 - \hat{\Sigma})(\hat{\Sigma}_1)^{-1}) \quad (3)$$

$$\text{• Roy's greatest root: } \frac{\eta_1}{1+n_1} \quad (L_3) \quad (4)$$

Fact (Linear Algebra)

For two matrices E and H , (both invertible, symmetric, and same size)

$$\frac{|E|}{|E + H|} = \frac{1}{|E| |E + H|} = \frac{1}{|I + E^{-1}H|}$$

$$= \prod_{i=1}^s \frac{1}{1+n_i} \quad \begin{matrix} \text{if the eigenvalues of } E^{-1}H \text{ are} \\ n_1, n_2, \dots, n_s, 0 \end{matrix}$$

Apply this with $E = \hat{\Sigma}$, $H = \hat{\Sigma}_1 - \hat{\Sigma}$ to obtain (1)

$$\text{For (2), } \sum_{i=1}^s \frac{\eta_i}{1+n_i} = \sum_{i=1}^s \frac{1}{1+\frac{1}{n_i}} \quad \text{"m-dimension"}$$

$$\begin{aligned} \text{trace}(H(E+E^{-1})) &= \text{tr}(H[(E^{-1}+I)H^{-1}]) \\ &= \text{tr}(HH^{-1}(E^{-1}+I)^{-1}) = \text{tr}((I+E^{-1})^{-1}) \end{aligned}$$

If the eigenvalues of EH^{-1} are n_1, \dots, n_s , then

$$\text{the eigenvalues of } EH^{-1} \text{ are } \frac{1}{n_1}, \dots, \frac{1}{n_s}, \text{ then}$$

$$\text{the eigenvalues of } (I+EH^{-1}) \text{ are } 1 + \frac{1}{n_1}, \dots, 1 + \frac{1}{n_s},$$

$$\text{then the eigenvalues of } (I+EH^{-1})^{-1} \text{ are } \frac{1}{1+\frac{1}{n_1}}, \dots, \frac{1}{1+\frac{1}{n_s}}.$$

$$\text{then } \text{tr}((I+EH^{-1})^{-1}) = \sum_{i=1}^s \frac{1}{1+\frac{1}{n_i}} \quad \text{this satisfies (2)}$$

We will reject H_0 if L_1 is small, or L_2 is large, or
 L_3 is large, or L_4 is large.

$$\begin{bmatrix} 1 & 2 & 3 \\ 4 & 5 & 6 \\ 7 & 8 & 9 \end{bmatrix}$$

$\lambda = 2.6 + 3c$

Fact (Linear Algebra)

Suppose A and B are two matrices (same size and squared).
 A is n.n.d and B is p.d. (so, B is non-singular) $\xrightarrow{\text{invertible}} \det(B) \neq 0$.
 Then the eigenvalues of AB^{-1} are non-negative.

Fact For two matrices C and D , the non-zero eigenvalues of CD and DC are same.

using this, all non-zero eigenvalues of AB^{-1} and $\underbrace{B^{-1/2}AB^{-1/2}}_{\substack{\text{"symmetric"} \\ \text{n.n.d}}} \xrightarrow{\text{n.n.d}}$ are same.

$$\beta = \begin{bmatrix} \beta_{01} & \cdots & \beta_{0m} \\ \beta_{11} & \cdots & \beta_{1m} \\ \vdots & \ddots & \vdots \\ \beta_{r1} & \cdots & \beta_{rm} \end{bmatrix}_{(r+1) \times m}$$

$$\hat{\beta} = (Z'Z)^{-1}Z'Y$$

$$\hat{\beta}' = Y'Z(Z'Z)^{-1}$$

Prediction

For a future observation with (unknown) response vector y_0 and (known) predictor values \underline{z}_0 ,

$$\mathbb{E}[y_0] = \underline{z}_0 \beta, \quad \mathbb{E}[y_0] = \beta' \underline{z}_0$$

$$\underline{z}_0 = \begin{bmatrix} 1 \\ z_{01} \\ \vdots \\ z_{0r} \end{bmatrix}_{(r+1)+1}$$

$\xrightarrow{\text{confidence region}} \text{Prediction region}$.

1. Point estimator of both y_0 and $\mathbb{E}(y_0)$

$$\hat{\beta}' \underline{z}_0 = Y'Z(Z'Z)^{-1} \underline{z}_0$$

2. Confidence region for $\mathbb{E}(y_0)$ $\xrightarrow{\substack{y_0 \text{ is a vector} \\ \text{not a single value}}} \text{So, not confidence interval.}$

3. Prediction region for y_0

$$y_0 \sim N_m(\beta' \underline{z}_0, \Sigma), \quad \hat{\beta}' \underline{z}_0 \sim N_m(\beta' \underline{z}_0, (\underline{z}_0' (Z'Z)^{-1} \underline{z}_0) \Sigma)$$

Recall

$$\beta = [\beta_{(1)} | \cdots | \beta_{(m)}]$$

$$\mathbb{E}[y_0] = \begin{bmatrix} \underline{z}_0' \beta_{(1)} \\ \vdots \\ \underline{z}_0' \beta_{(m)} \end{bmatrix}$$

$$\text{cov}(\beta_{(j)}, \beta_{(k)}) = \sigma_{jk} (Z'Z)^{-1}$$

$$\text{so, cov}(\underline{z}_0' \beta_{(j)}, \underline{z}_0' \beta_{(k)}) = \underline{z}_0' \sigma_{jk} (Z'Z)^{-1} \underline{z}_0 = (\underline{z}_0' (Z'Z)^{-1} \underline{z}_0) \sigma_{jk}$$

$$\text{Hotelling } T^2 = \frac{(\hat{\beta}' \underline{z}_0 - \beta' \underline{z}_0)' S^{-1} (\hat{\beta}' \underline{z}_0 - \beta' \underline{z}_0)}{\sqrt{\underline{z}_0' (Z'Z)^{-1} \underline{z}_0} \sqrt{\underline{z}_0' (Z'Z)^{-1} \underline{z}_0}}$$

$$\text{, where } S = \frac{1}{n-r-1} (Y - Z\hat{\beta})'(Y - Z\hat{\beta})$$

$$100(1-\alpha)\% \text{ confidence region} = \left\{ \beta' \underline{z}_0 : T^2 \leq \frac{m(n-r-1)}{n-r-m} F_{m, n-r-m}(\alpha) \right\}$$

The 100 $(1-\alpha)\%$ simultaneous confidence intervals for $E(Y_i) = \mathbf{z}_i' \hat{\beta}$ are

$$\mathbf{z}_i' \hat{\beta} \pm \sqrt{\left(\frac{m(n-r-1)}{n-r-m}\right) F_{m,n-r-m}(\alpha)} \sqrt{\mathbf{z}_i' (\mathbf{Z}'\mathbf{Z})^{-1} \mathbf{z}_i} S_{ii}, \quad i=1, \dots, m$$

where $\hat{\beta}_{(i)}$ is the i th column of $\hat{\beta}$ and S_{ii} is the i th diagonal element of S . Beta-hat $(:, i)$

$$3. Y_0 - \hat{\beta}' \mathbf{z}_0 = (\beta - \hat{\beta})' \mathbf{z}_0 + \varepsilon_0 \sim N_m(0, (1 + \mathbf{z}_0' (\mathbf{Z}'\mathbf{Z})^{-1} \mathbf{z}_0) \Sigma)$$

The 100 $(1-\alpha)\%$ prediction ellipsoid for Y_0 becomes

$$(\mathbf{y}_0 - \hat{\beta}' \mathbf{z}_0)' S^{-1} (\mathbf{y}_0 - \hat{\beta}' \mathbf{z}_0) \leq (1 + \mathbf{z}_0' (\mathbf{Z}'\mathbf{Z})^{-1} \mathbf{z}_0) \left[\left(\frac{m(n-r-1)}{n-r-m} \right) F_{m,n-r-m}(\alpha) \right]$$

The 100 $(1-\alpha)\%$ simultaneous prediction intervals for the individual responses Y_{0i} :

$$\mathbf{z}_0' \hat{\beta}_{(i)} \pm \sqrt{\left(\frac{m(n-r-1)}{n-r-m}\right) F_{m,n-r-m}(\alpha)} \sqrt{(1 + \mathbf{z}_0' (\mathbf{Z}'\mathbf{Z})^{-1} \mathbf{z}_0) S_{ii}}, \quad i=1, \dots, m$$

$\mathbf{z}_0' \text{Beta-hat}(:, i)$

$$\mathbf{z}_0 = [1 \ \mathbf{z}(3, 2:\text{end})]'$$

$$\therefore Y_3 = \mathbf{z}_0' \text{Beta-hat}(:, 3).$$

Math B7800

3/15/18, Thur

Recall suppose A is pd matrix ($n \times n$) and symmetric.
positive-definite

$$\max \{ (\mathbf{b}' \mathbf{x})^2 : \mathbf{x}' A \mathbf{x} \leq 1 \} = \mathbf{b}' A^{-1} \mathbf{b}$$

Proof $(\mathbf{b}' \mathbf{x})^2 = \underbrace{\mathbf{b}' A^{-1/2}}_{\text{vector}} \underbrace{A^{1/2} \mathbf{x}}_{\text{vector}} = (A^{-1/2} \mathbf{b})' (A^{1/2} \mathbf{x})$

Then

$$\begin{aligned} (\mathbf{b}' \mathbf{x})^2 &\leq \left[(A^{-1/2} \mathbf{b})' (A^{-1/2} \mathbf{b}) \right] \left[(A^{1/2} \mathbf{x})' (A^{1/2} \mathbf{x}) \right] \\ &= (\mathbf{b}' A^{-1} \mathbf{b}) \underbrace{(\mathbf{x}' A \mathbf{x})}_{\leq 1} \leq \mathbf{b}' A^{-1} \mathbf{b} \end{aligned}$$

by Cauchy-Schwartz Inequality

For two $n \times 1$ vector \mathbf{x} and \mathbf{y} , $(\mathbf{x}' \mathbf{y}) \leq (\mathbf{x}' \mathbf{x})(\mathbf{y}' \mathbf{y})$
Equality holds if $\mathbf{y} = c \mathbf{x}$ for some constant c .

Equality holds if $A^{1/2} \mathbf{x} = c A^{-1/2} \mathbf{b}$ for some constant c ,
and c should be such that $\mathbf{x}' A \mathbf{x} = 1$.

$$(*) \mathbf{x} = c A^{-1} \mathbf{b} \quad \text{so, } \mathbf{x}' A \mathbf{x} = c^2 \mathbf{b}' A^{-1} \mathbf{b} = 1$$

$$(c \mathbf{b}' A^{-1}) A (c A^{-1} \mathbf{b}) \uparrow \quad \text{so, } c = \pm \frac{1}{\sqrt{\mathbf{b}' A^{-1} \mathbf{b}}}$$

The value of \mathbf{x} that attains the max is $\pm \frac{1}{\sqrt{\mathbf{b}' A^{-1} \mathbf{b}}} A^{-1} \mathbf{b}$.

Recall the model

$$Y = Z\beta + \varepsilon, \quad Y_{n \times m} = \begin{bmatrix} y_1 \\ y_2 \\ \vdots \\ y_n \end{bmatrix}, \quad Z_{n \times (r+1)} = \begin{bmatrix} z_1' \\ \vdots \\ z_n' \end{bmatrix}$$

$$\beta = [\beta_{(1)} \mid \cdots \mid \beta_{(m)}] \quad \begin{bmatrix} \beta_{10} & \beta_{20} \\ \beta_{11} & \beta_{21} \\ \beta_{12} & \beta_{22} \end{bmatrix}$$

$$\varepsilon_{n \times m} = \begin{bmatrix} \varepsilon_1' \\ \vdots \\ \varepsilon_m' \end{bmatrix} \text{ such that } \varepsilon_1, \dots, \varepsilon_m \text{ are iid } N_m(\mathbf{0}, \Sigma)$$

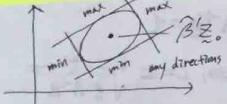
(assumption) $\rightarrow \hat{\beta} = (Z'Z)^{-1}Z'Y$

Now, we have a future observation with unknown response y_0 ,
but known predictors z_0 .

$$y_0 \sim N(\beta'z_0, \Sigma), \quad E[y_0] = \beta'z_0.$$

We saw that $100(1-\alpha)\%$ CR for $\beta'z_0$ is

$$\{z: (z - \hat{\beta}'z_0)' S^{-1} (z - \hat{\beta}'z_0) \leq \frac{m(n-r-1)}{n-r-m} F_{m, n-r-m}(\alpha)\}$$



Q: Obtain simultaneous Scheffe CI for linear combinations of the components of $\beta'z_0$ ($m \times 1$).

so, we want CI for $b'\beta'z_0$ (for all $m \times 1$ vectors b).

$\hat{\beta}'z_0 = \{z: (\hat{z} - \hat{\beta}'z_0)' A (\hat{z} - \hat{\beta}'z_0) \leq 1\}$, where

$$A = \frac{S^{-1}}{\frac{m(n-r-1)}{n-r-m} F_{m, n-r-m}(\alpha)}$$

If $z = \hat{z} - \hat{\beta}'z_0$, then the constraint is $z'Az \leq 1$. Under this constraint, $(b'z)^2 \leq b'A^{-1}b$ by Cauchy-Schw. Inequality.
So, $(b'(\hat{z} - \hat{\beta}'z_0))^2 \leq b' \frac{m(n-r-1)}{n-r-m} S F_{m, n-r-m}(\alpha) b$

so, \hat{z} represents $\beta'z_0$.

so, the simultaneous CI for $b'\beta'z_0$ is

$$b'\beta'z_0 \pm \sqrt{\frac{m(n-r-1)}{n-r-m} F_{m, n-r-m}(\alpha) b'b}$$

$P(\text{the interval } I_b \text{ contains } b'\beta'z_0 \text{ for all } b) = 1 - \alpha$.
simultaneous.

• General regression function "joint distribution of y and \mathbf{z} "

Suppose (y, z_1, \dots, z_r) have a joint distribution with mean

$$\mu = \begin{bmatrix} \mu_y \\ \mu_{\mathbf{z}} \\ \vdots \\ \mu_{z_r} \end{bmatrix} \quad \text{and covariance } \Sigma = \begin{bmatrix} \sigma_{yy} & \sigma'_{yz} \\ \sigma_{yz} & \Sigma_{zz} \end{bmatrix}_{r \times r}$$

Goal: predict y based on z_1, \dots, z_r ← random variables.

We predict y based on z_1, \dots, z_r via some function $g(z_1, \dots, z_r)$.
prediction error minimized.

prediction error = $y - g(z_1, \dots, z_r)$, $g: \mathbb{R}^r \rightarrow \mathbb{R}$.

mean prediction squared error = $E[(y - g(z_1, \dots, z_r))^2]$

What function minimizes $g(z_1, \dots, z_r)$?

Step 1 $E(y - a)^2 = h(a)$, $h(a)$ is minimized at $a = E(y)$.
 $E[y] \nearrow$ minimize

Step 2 $E[(y - a)^2 | \mathbf{z}] = h(a)$, $h(a)$ is minimized at $E(y|\mathbf{z}) = \frac{E(y, \mathbf{z})}{E(\mathbf{z})}$

Step 3 $E[(y - a)^2 | z_1, \dots, z_r]$ is minimized at $E(y|z_1, \dots, z_r) = g^*(z_1, \dots, z_r)$

$E[(y - g(z_1, \dots, z_r))^2] = E\left\{ E[(y - g(z_1, \dots, z_r))^2 | z_1, \dots, z_r] \right\} \quad (*)$

④ $E[(y - g(z_1, \dots, z_r))^2 | z_1, \dots, z_r] \geq E[(y - g^*(z_1, \dots, z_r))^2 | z_1, \dots, z_r]$

Taking expectation both sides,

$E[(y - g(z_1, \dots, z_r))^2] \geq E[(y - g^*(z_1, \dots, z_r))^2] \quad \text{by } (*)$

So, the minimizer $E(y|z_1, \dots, z_r)$ is called the regression function of y on z_1, \dots, z_r . "Jointly normal \Rightarrow linear regression"

Q: When $\begin{pmatrix} y \\ z_1 \\ \vdots \\ z_r \end{pmatrix} \sim N_{r+1}(\mu, \Sigma)$, then $E(y|z_1, \dots, z_r) = ?$

Then the conditional distribution of y given z_1, \dots, z_r is

$$N(\mu_y + \sigma'_{yz} \Sigma_{zz}^{-1} (\mathbf{z} - \mu_z), \sigma_{yy} - \sigma'_{yz} \Sigma_{zz}^{-1} \sigma_{yz})$$

So, $E[y|z_1, \dots, z_r] = \mu_y + \sigma'_{yz} \Sigma_{zz}^{-1} (\mathbf{z} - \mu_z) \leftarrow \text{constant} \Rightarrow \text{linear}$
 $= \beta_0 + \beta_1' \mathbf{z}$, where

$$\beta_0 = (\mu_y - \sigma'_{yz} \Sigma_{zz}^{-1} \mu_z), \beta_1 = \Sigma_{zz}^{-1} \sigma_{yz}$$

* result If the joint distribution is not normal, then $E(y|z_1, \dots, z_r)$ need not to be linear.

Ex Suppose (y, z) has joint pdf, $f_{y,z}(y, z) = z e^{-yz-1}$, $0 \leq y, z < \infty$.

Find $E(y|z) = ?$

$$f_{y|z}(y|z) = \frac{f_{y,z}(y, z)}{f_z(z)} = \frac{z e^{-yz-1}}{\int_0^\infty z e^{-yz-1} dy} = z e^{-yz}$$

"marginal" →

$$E[y|z] = \int_0^\infty y \cdot f_{y|z}(y|z) = \int_0^\infty y z e^{-yz} dy = \frac{1}{z}$$

So, $E[y|z]$ is not linear in z .

In general, (without normality) among all linear predictors of y based on z_1, \dots, z_r , the one that minimizes the mean squared error is $\beta_0 + \beta' z$, where $\beta_0 = (M_y - \Sigma_{yz}^{-1} M_z)$, $\beta = \Sigma_{zz}^{-1} \Sigma_{yz}$.

Def A linear predictor of y based on z_1, \dots, z_r has the form $b_0 + b' z$ for some constants $b_0, b = \begin{pmatrix} b_1 \\ \vdots \\ b_r \end{pmatrix}$

If $f(b_0, b) = E[(y - b_0 - b' z)^2] \geq f(\beta_0, \beta)$,
for any b_0, b .

Math B7800

3/20/18, Tue

• Best Linear Predictor

Recall (y, z_1, \dots, z_r) has a joint distribution with

$$\text{mean } \mu = \begin{bmatrix} M_y \\ M_z \end{bmatrix}_{(r+1) \times 1}, \text{ covariance } \Sigma = \begin{bmatrix} \Sigma_{yy} & \Sigma_{yz} \\ \Sigma_{yz} & \Sigma_{zz} \end{bmatrix}_{(r+1) \times (r+1)}$$

Goal: Find a linear predictor of y based on z_1, \dots, z_r , which has smallest possible mean squared error.

A linear predictor has the form $b_0 + b' z$, $z = \begin{bmatrix} z_1 \\ \vdots \\ z_r \end{bmatrix}_{r \times 1}$.

Its mean squared error: $g(b_0, b) = E[(y - b_0 - b' z)^2]$

Goal: Minimize $g(b_0, b)$ w.r.t. b_0, b .

Result: $g(b_0, b) \geq g(\beta_0, \beta)$, where $\beta_0 = M_y - \Sigma_{yz}^{-1} M_z$, $\beta = \Sigma_{zz}^{-1} \Sigma_{yz}$.

$$\begin{aligned} \text{proof } g(b_0, b) &= E[(y - b_0 - b' z)^2] \stackrel{(a+b+c)^2 = a^2 + b^2 + c^2 + 2ab + 2bc + 2ac}{=} E[(y - M_y - b'(z - M_z) - b_0 + M_y - b'M_z)^2] \\ &= E[(y - M_y)^2] + E[(b'(z - M_z))^2] + (b_0 - M_y + b'M_z)^2 \\ &\quad \xrightarrow{\text{by}} + 0 + 0 - 2E[(y - M_y)b'(z - M_z)] \quad b' \Sigma_{zz} b \\ &\quad \xrightarrow{2E[(y - M_y)(b'(z - M_z))]} 2E[(b'(z - M_z))(b_0 - M_y + b'M_z)] \\ &\quad \xrightarrow{2E[(b'(z - M_z))(b_0 - M_y + b'M_z)]} \\ &\quad = \Sigma_{yy} + b' \sum_{zz} b - 2b' \Sigma_{yz} - \textcircled{*} \end{aligned}$$

$$\frac{\partial g(\mu_y - k' M_z, k)}{\partial k} = 0 + 2 \sum_{zz} k - 2 \Sigma_{yz}$$

Set $\frac{\partial g}{\partial k} = 0$, we get $k = \underbrace{\sum_{zz}^{-1} \Sigma_{yz}}_{\beta} \leftarrow \text{"minimizer of } k\text{"}$

so, $g(b_0, k) \geq g(\mu_y - k' M_z, k)$
 $\geq g(\mu_y - \beta' M_z, \beta) = g(\beta_0, \beta)$
 $\beta_0 = \mu_y - \Sigma_{yz}' \sum_{zz}^{-1} M_z, \beta = \sum_{zz}^{-1} \Sigma_{yz}$

• Smallest mean squared error among linear predictors:

$$\begin{aligned} E[(y - \beta_0 - \beta' z)^2] &= E[(y - (\mu_y - \Sigma_{yz}' \sum_{zz}^{-1} M_z) - \Sigma_{yz}' \sum_{zz}^{-1} z)^2] \\ &= E[(y - \mu_y + \Sigma_{yz}' \sum_{zz}^{-1} M_z - \Sigma_{yz}' \sum_{zz}^{-1} z)^2] \\ &= E[((y - \mu_y) - \Sigma_{yz}' \sum_{zz}^{-1} (z - \mu_z))^2] \\ &= g(\beta_0, \beta) = \sigma_{yy} + (\beta' \sum_{zz} \beta - 2 \beta' \Sigma_{yz}) \quad \begin{matrix} \beta' \\ \Sigma_{yz} \sum_{zz}^{-1} \Sigma_{yz} \\ \sum_{zz} \end{matrix} \\ \text{by plugin } \textcircled{*} \rightarrow &= \sigma_{yy} + \Sigma_{yz}' \sum_{zz}^{-1} \Sigma_{yz} - 2 \Sigma_{yz}' \sum_{zz}^{-1} \Sigma_{yz} \\ &= \sigma_{yy} - \Sigma_{yz}' \sum_{zz}^{-1} \Sigma_{yz}. \\ &= \sigma_{yy} \left[1 - \frac{\Sigma_{yz}' \sum_{zz}^{-1} \Sigma_{yz}}{\sigma_{yy}} \right] \\ &= \sigma_{yy} [1 - \rho_{y(z)}^2], \text{ where } \rho_{y(z)} = \pm \sqrt{\frac{\Sigma_{yz}' \sum_{zz}^{-1} \Sigma_{yz}}{\sigma_{yy}}} \end{aligned}$$

called the "population multiple correlation coefficient between y and z ".

Result \star

$$\max_{b_0, k} \text{corr}(y, b_0 + k' z) = \text{corr}(y, \beta_0 + \beta' z) = \rho_{y(z)}$$

$$\text{proof } [\text{corr}(y, b_0 + k' z)]^2 = \frac{[\text{cov}(y, b_0 + k' z)]^2}{\text{var}(y) \text{var}(b_0 + k' z)} = \frac{(k' \Sigma_{yz})^2}{\sigma_{yy} k' \sum_{zz} k}$$

Recall using Cauchy-Schwarz Inequality.

$$\begin{aligned} (k' \Sigma_{yz})^2 &= (k' \sum_{zz}^{-1} \sum_{zz} \Sigma_{yz})^2 \\ &\leq (k' \sum_{zz} k) (\Sigma_{yz}' \sum_{zz}^{-1} \Sigma_{yz}) \end{aligned}$$

$$\text{so, } \frac{(k' \Sigma_{yz})^2}{k' \sum_{zz} k} \leq \Sigma_{yz}' \sum_{zz}^{-1} \Sigma_{yz}$$

$$\text{so, } [\text{corr}(y, b_0 + k' z)]^2 \leq \frac{\Sigma_{yz}' \sum_{zz}^{-1} \Sigma_{yz}}{\sigma_{yy}}$$

If y, z has a joint dist., we can predict y based on z , and vice-versa.

$$= (\text{corr}(y, \beta_0 + \beta' z))^2 = \rho_{y(z)}^2$$

Now, suppose we have $\begin{bmatrix} y \\ z \end{bmatrix}$ is jointly $N_{n+1}(\mu, \Sigma)$, we observe

a sample of size n , $(y_1, z_1), \dots, (y_n, z_n)$.

Based on these samples, we find MLE of Best linear predictor, smallest mean squared error, $\rho_{y(z)}$.

We compute the sample mean and sample covariance,

$$\hat{\mu} = \begin{pmatrix} \bar{y} \\ \bar{z} \end{pmatrix}, S = \begin{bmatrix} S_{yy} & S'_{yz} \\ S'_{yz} & S_{zz} \end{bmatrix}, \text{ where } S_{yy} = \frac{1}{n-1} \sum_{i=1}^n (y_i - \bar{y})^2.$$

\downarrow $(y_i - \bar{y})(y_i - \bar{y})' / n-1$

$$(\hat{S}_{yz})_{k,l} = \frac{1}{n-1} \sum_{i=1}^n (y_i - \bar{y})(z_{ik} - \bar{z}_k), k=1, \dots, r$$

\downarrow $(y_i - \bar{y})(z_{il} - \bar{z}_l)' / n-1$

$$(\hat{S}_{zz})_{k,l} = \frac{1}{n-1} \sum_{i=1}^n (z_{ik} - \bar{z}_k)(z_{il} - \bar{z}_l), k, l = 1, 2, \dots, r$$

\downarrow $(z_{ik} - \bar{z}_k)(z_{il} - \bar{z}_l)' / n-1$

\star MLE of μ, S are $\hat{\mu}, \frac{n-1}{n} S$.

* Recall (Invariance property)

If $\hat{\theta}$ is the MLE of θ , then for any function g ,

$g(\hat{\theta})$ is an MLE of $g(\theta)$.

Recall $\beta_0 = My - \hat{S}'_{yz} \hat{S}_{zz}^{-1} \hat{M}_z, \beta_1 = \hat{S}_{zz}^{-1} \hat{S}_{yz}$.

so, $\hat{\beta}_0 = \bar{y} - \hat{S}'_{yz} \hat{S}_{zz}^{-1} \bar{z}$ and $\hat{\beta}_1 = \hat{S}_{zz}^{-1} \hat{S}_{yz}$

* MLE for Best Linear Predictor

$$= \hat{\beta}_0 + \hat{\beta}_1' \bar{z} = \bar{y} + \hat{S}'_{yz} \hat{S}_{zz}^{-1} (\bar{z} - \bar{z})$$

$\leftarrow \bar{z} \text{ is a variable}$
 "observable"

$$\hat{\rho}_{y(z)}^2 = \frac{\hat{S}'_{yz} \hat{S}_{zz}^{-1} \hat{S}_{yz}}{S_{yy}}$$

\star MLE for smallest mean squared error
 $S_{yy} - \hat{S}'_{yz} \hat{S}_{zz}^{-1} \hat{S}_{yz}$ is

$$\frac{n-1}{n} (S_{yy} - \hat{S}'_{yz} \hat{S}_{zz}^{-1} \hat{S}_{yz})$$

$$\hat{\rho}_{y(z)}^2 = \frac{\hat{S}'_{yz} \hat{S}_{zz}^{-1} \hat{S}_{yz}}{S_{yy}}$$

Multiple response case

$\begin{bmatrix} Y \\ Z \end{bmatrix} \sim \begin{bmatrix} m \times 1 \\ r \times 1 \end{bmatrix}$ have a joint distribution with mean $\mu = \begin{pmatrix} M_y \\ M_z \end{pmatrix}$ and

$$\text{cov} \begin{pmatrix} Y \\ Z \end{pmatrix} = \sum = \begin{bmatrix} \sum_{yy} & \sum_{yz} \\ \sum_{yz} & \sum_{zz} \end{bmatrix} \begin{matrix} m \times m \\ m \times r \\ r \times m \\ r \times r \end{matrix}$$

A linear predictor of Y based on Z has the form

$$\underbrace{b_0}_{m \times 1} + \underbrace{b' Z}_{m \times r} \quad \text{where } b_0 \text{ is a constant vector and } b \text{ is } m \times r \text{ matrix.}$$

We minimize mean square and cross products of errors

$$\mathbb{E}[(Y - b_0 - b' Z)(Y - b_0 - b' Z)'] \quad \leftarrow \mathbb{E} Z' \quad A = \text{matrix}$$

* The best linear predictor of Y based on Z is

$$My - \hat{S}'_{yz} \hat{S}_{zz}^{-1} \hat{M}_z + \hat{S}'_{yz} \hat{S}_{zz}^{-1} Z = \beta_0 + \beta' Z$$

β

$$= My + \hat{S}'_{yz} \hat{S}_{zz}^{-1} (\bar{z} - \bar{M}_z)$$

"smallest" (in the matrix sense: $A \geq B$ for matrices A, B
 $\text{if } A - B \text{ is n.n.d.}")$

* mean square and products of errors

$$= \sum_{yy} - \hat{S}'_{yz} \hat{S}_{zz}^{-1} \hat{S}_{yz} = \sum_{yy} \cdot z$$

$$= \mathbb{E}[(Y - \beta_0 - \beta' Z)(Y - \beta_0 - \beta' Z)']$$

3/22/18, Thur

"Partial correlation coefficient" between y_k, y_l eliminating
the effect of \underline{z} is the (k, l) corr. coeff. corresponding to the covariance
matrix $\Sigma_{yy \cdot z}$.

e.g. If $\text{cov}(y_k, y_l)$ is small, y_k, y_l are dependent on \underline{z} .
w/ eliminating the effect of \underline{z} .
& $\text{corr}(y_k, y_l)$ is large
w/ the effect of \underline{z} .

If $\text{corr}(y_k, y_l)$ is large, then y_k, y_l are not dependent on \underline{z} .
w/ eliminating the effect of \underline{z} .
& $\text{cov}(y_k, y_l)$ is large
w/ the effect of \underline{z} they are ~~moderately~~ strongly correlated each other.

$$\Sigma_{yy \cdot z} = \begin{bmatrix} \sigma_{y_1 y_1 \cdot z} & \sigma_{y_1 y_2 \cdot z} & \dots & \sigma_{y_1 y_m \cdot z} \\ \vdots & \ddots & & \vdots \\ \sigma_{y_m y_1 \cdot z} & \dots & \dots & \sigma_{y_m y_m \cdot z} \end{bmatrix}$$

so, Partial corr. coeff. between y_k, y_l eliminating the effect of \underline{z}

$$= \frac{\sigma_{y_k y_l \cdot z}}{\sqrt{\sigma_{y_k y_k \cdot z} \sigma_{y_l y_l \cdot z}}}$$

Math B1800

4/10, Tue "Exam 2"

or 4/12, Thur

Chapter 14 of the other text book.

Non-linear Regression response

In linear regression model, $Y = \beta_0 + \beta_1 x + \epsilon$ predictor error

$E[Y] = \beta_0 + \beta_1 x$ (a linear function both in parameters and predictors)

In nonlinear regression, the relation between $E[Y]$ and the parameters β_0, β_1 , and the predictor x is no longer linear.

"generalized linear model"

An important and useful subclass is the generalized linear model (GLM).

Here, the relation between $E[Y]$ and β_0, β_1 , and x is via a link function η . s.t. $\eta(E[Y]) = \beta_0 + \beta_1 x$.

e.g.

Logistic regression, Probit response function,

log-log mean response function.

In many examples, the response variable takes two possible values.

Ex 1 While predicting the possibility of heart disease based on diet and lifestyle measures.

Ex 2 While predicting whether a married women will be in the work force based on family background, education, household income, etc.

Ex3 while predicting whether a home owner will buy liability insurance.

Ex4 while predicting whether a company will have a legal division or not.

Ex5 while deciding whether to give loan/credit card.

Now, we assume y takes two values: 0 and 1.

Suppose $\pi = P(Y=1)$, so $P(Y=0) = 1-\pi$.

$$\mathbb{E}[Y] = \pi, \quad \text{var}(Y) = \pi(1-\pi)$$

If we try to use linear model, $y_i = \beta_0 + \beta_1 x_i + \varepsilon_i$

Issues:
 1. $\text{Var}(\varepsilon_i) = \text{Var}(y_i) = \pi_i(1-\pi_i) = \mathbb{E}(y_i)(1-\mathbb{E}(y_i))$

(variance depends on mean)

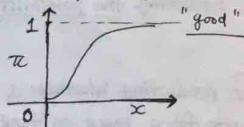
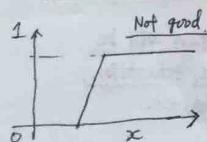
2. Possible values of ε_i : $1 - \mathbb{E}(y_i), 0 - \mathbb{E}(y_i)$
 $\Rightarrow 1 - \pi_i, -\pi_i$.

(Normal approximation for ε_i is not good).

*3. $\pi_i = \mathbb{E}(y_i) = \beta_0 + \beta_1 x_i \leftarrow \mathbb{E}(y_i) = 1\pi_i + 0(1-\pi_i) = \pi_i = P(Y_i=1)$

Since π_i is a probability, $|\beta_0 + \beta_1 x_i| \leq 1$.

(π_i is dependent on x_i)



Ex y^c = Duration of pregnancy

X = Index for alcohol consumption.

Here, we can use linear model $y^c = \beta_0^c + \beta_1^c x + \varepsilon^c$, where $\varepsilon^c \sim N(0, \sigma^2)$

It is considered to be a premature delivery if $y^c \leq 38$ weeks,
 otherwise full term delivery.

Define

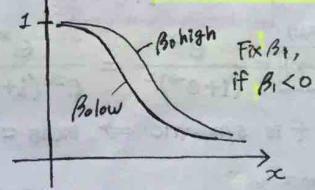
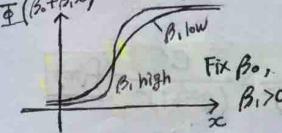
$$y = \begin{cases} 1 & \text{if } y^c \leq 38 \text{ (preterm)} \\ 0 & \text{if } y^c > 38 \text{ (full term)} \end{cases}$$

$$\begin{aligned} \pi &= P[y=1] = P[y^c \leq 38] = P(\beta_0^c + \beta_1^c x + \varepsilon^c \leq 38) \\ &= P[\varepsilon^c \leq 38 - \beta_0^c - \beta_1^c x] = P\left[\frac{\varepsilon^c}{\sigma} \leq \frac{38 - \beta_0^c - \beta_1^c x}{\sigma}\right] \\ &= P[Z \leq \beta_0 - \beta_1 x] = \Phi(\beta_0 - \beta_1 x) \end{aligned}$$

$$\text{So, } \Phi^{-1}(\pi) = \beta_0 + \beta_1 x$$

$$\Phi^{-1}(\mathbb{E}(y)) = \beta_0 + \beta_1 x$$

This is called "probit mean response function" (probit transformation
 $\Phi(\beta_0 + \beta_1 x)$ is $x \mapsto \Phi^{-1}(x)$).

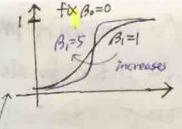


Properties of the curve

1. Always lie between 0 and 1.

2. Monotone function.

3. For fixed β_0 , if β_1 increases, then curve becomes more sigmoid. (S자 모양의) (이중 S자 모양)



4. If β_1 is fixed and β_0 is changed, then the curve shifts. (이중 S자 모양)

5. If $P(Y=1) = \pi = \Phi(\beta_0 + \beta_1 x)$, then

$$1 - \pi = P(Y=0) = 1 - \Phi(\beta_0 + \beta_1 x) \quad (\text{Using the property}) \\ = \Phi(-\beta_0 - \beta_1 x) \quad (\Phi(z) + \Phi(-z) = 1)$$

"(read) of meaning"

* This model is hard to interpret value of β_1 (parameter) (probit model)

Logistic Regression

Logistic distribution is a continuous distribution with pdf

pdf $f(x) = \frac{e^x}{(1+e^x)^2}, x \in \mathbb{R}$

cdf $F(x) = \frac{e^x}{1+e^x}, x \in \mathbb{R}$

$$f(-x) = \frac{e^{-x}}{(1+e^{-x})^2} = \frac{e^x}{e^{2x}(1+e^{-x})^2} = \frac{e^x}{(e^x+1)^2} = f(x)$$

So, f is symmetric \Rightarrow mean = 0.

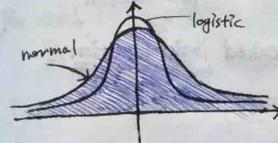
$$e^{2x} (e^{-x} + 2e^{-x} + 1) \\ = (1 + 2e^{-x} + e^{2x}) = (e^x + 1)^2$$

Variance = $\int_{-\infty}^{\infty} x^2 \frac{e^x}{(1+e^x)^2} dx$ pdf
 $= 2 \int_0^{\infty} x^2 \frac{e^x}{(1+e^x)^2} dx$ "even function"

"by part"
 $\int u dv = uv - \int v du$
 $U = x^2, dU = \frac{2x}{1+e^x} dx$
 $V = -\frac{1}{(1+e^x)}$
 $\Rightarrow 0$

$$= 2 \left[x^2 \left(-\frac{1}{1+e^x} \right) \Big|_0^\infty + \int_0^\infty \frac{2x}{1+e^x} dx \right]$$

$$= 4 \int_0^\infty \frac{x}{1+e^x} dx \quad \text{check} \quad = 4 \int_0^\infty \log(1+e^{-x}) dx \quad \text{check}$$



Now, suppose $y^c = \beta_0^c + \beta_1^c x + \varepsilon^c$, where

ε^c has logistic distribution with mean 0, variance σ^2 .

$$y = \begin{cases} 1 & \text{if } y^c \leq 38 \\ 0 & \text{if } y^c > 38 \end{cases}$$

$$\pi = P(Y=1) = P(y^c \leq 38) = P(\varepsilon^c \leq 38 - \beta_0^c - \beta_1^c x)$$

$$= P\left(\frac{\varepsilon^c}{\sigma^2} \frac{\pi}{\sqrt{3}} \leq \frac{\pi}{\sqrt{3}} 38 - \underbrace{\frac{\pi}{\sqrt{3}} \beta_0^c}_{\text{var}=1} - \underbrace{\frac{\pi}{\sqrt{3}} \beta_1^c x}_{\beta_0} \beta_1\right)$$

$$= \frac{e^{\beta_0 + \beta_1 x}}{1 + e^{\beta_0 + \beta_1 x}}$$

← cdf

$$y = \frac{e^x}{1+e^x}, \text{ then } x = \log\left(\frac{y}{1-y}\right)$$

$$\begin{aligned} & \text{Left: } y + e^x y = e^x \\ & \text{Right: } \log y + \log(1-y) = x \end{aligned}$$

$$\frac{y}{1-y} = \frac{e^x}{(1+e^x)-e^x} = e^x$$

$$\text{So, } \log\left(\frac{\pi}{1-\pi}\right) = \beta_0 + \beta_1 x$$

The function $y \rightarrow \log\left(\frac{y}{1-y}\right)$ is called logit function.

$$\frac{\pi}{1-\pi} = \frac{P(y=1)}{P(y=0)} \leftarrow \text{odd ratio.}$$

Math B7800

4/10/18, Tue

"4/26 - Exam II"

(Tue) Thur.

Generalized Linear Model

Logistic Regression / GLM

We have a binary response y and predictor x .

y takes two values: 0 or 1.

" $y = E(y) + \text{error}$ "

For the i th observation, $\pi_i = E(y_i) = P(y_i=1)$.

$y_i = \text{response}$, $x_i = \text{predictor}$

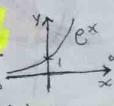
For GLM, $\pi_i = F(\beta_0 + \beta_1 x_i)$, where F is a function (usually F is the cdf of a distribution).

Different choices of F give different GLM. (pdf)

(For a random variable X , the cdf $F_X(x) = P(X \leq x)$)

* If $F(x) = \Phi(x)$, where Φ is the cdf of $N(0,1)$. Then we have the probit regression model.

$$\pi_i = \Phi(\beta_0 + \beta_1 x_i) \Leftrightarrow \Phi^{-1}(\pi_i) = \beta_0 + \beta_1 x_i$$

* When $F(x) = \frac{e^x}{1+e^x} = 1 - \frac{1}{1+e^x}$ "increasing" ↗ 

as $x \rightarrow -\infty$, $F(x) = 0$

as $x \rightarrow \infty$, $F(x) = 1$

of logistic distribution

Using this F , we get logistic regression model.

$$\star \pi_i = \frac{e^{\beta_0 + \beta_1 x_i}}{1 + e^{\beta_0 + \beta_1 x_i}} \Leftrightarrow \frac{1}{\pi_i} = 1 + \frac{1}{e^{\beta_0 + \beta_1 x_i}} \Leftrightarrow \frac{\pi_i}{1 - \pi_i} = e^{\beta_0 + \beta_1 x_i}$$

$F(\beta_0 + \beta_1 x_i)''$

$P(y_i=1)$

"odds" ↗ $\Leftrightarrow \log \frac{\pi_i}{1 - \pi_i} = \beta_0 + \beta_1 x_i$

$$\Leftrightarrow \log \frac{\pi_i}{1-\pi_i} = \beta_0 + \beta_1 x_i$$

is called odds = $\frac{P(Y_i=1)}{P(Y_i=0)} = \frac{\pi_i}{1-\pi_i}$

"odd ratio".

If $F(x) = 1 - e^{-e^x}$ → increasing
is the cdf of Gumbel dist.

(If $y = 1 - e^{-e^x}$, then $x = \log(-\log(1-y))$)

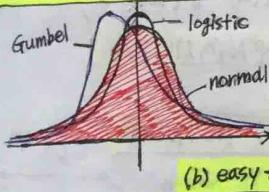
then $F^{-1}(\pi) = \log(-\log(1-\pi))$.

④ log-log regression function.

$$\text{Here } \pi_i = 1 - e^{\beta_0 + \beta_1 x_i} = F(\beta_0 + \beta_1 x_i)$$

$$\Leftrightarrow \log(-\log(1-\pi_i)) = \beta_0 + \beta_1 x_i$$

The pdf of the "three distribution".



Logistic regression is the most popular among the three GLM.

The reason is (a) availability of software for computation;
(b) easy to interpret the parameters.

For logistic regression, $\log \frac{\pi}{1-\pi} = \beta_0 + \beta_1 x$. odd = $\frac{\pi}{1-\pi} = \frac{P(Y_i=1)}{P(Y_i=0)}$

β_1 captures the change in odds for each unit change in x .
Suppose π and $\tilde{\pi}$ correspond to the predictor value x and $x+1$.

$$\text{Then } \log \frac{\pi}{1-\pi} = \beta_0 + \beta_1 x$$

$$\log \frac{\tilde{\pi}}{1-\tilde{\pi}} = \beta_0 + \beta_1 (x+1)$$

$$\text{Subtract } \log \frac{\tilde{\pi}}{1-\tilde{\pi}} - \log \frac{\pi}{1-\pi} = \beta_1$$

$$= \log \frac{\pi/1-\pi}{\tilde{\pi}/1-\tilde{\pi}}$$

$$\star \text{Odd Ratio (OR)} = \frac{\tilde{\pi}/(1-\tilde{\pi})}{\pi/(1-\pi)} = e^{\beta_1}$$

If $\hat{\beta}_1$ is the MLE of β_1 , then $\hat{OR} = e^{\hat{\beta}_1}$.

★ How to find MLE of β_0, β_1 in logistic regression.

Need to find the "likelihood function" "joint dist. of y ".

$L(\beta_0, \beta_1) = \text{Joint distribution of the responses.}$

We have n -individuals with known response and predictor values

		What is the pmf of y_1 ?
y	x	
y_1	x_1	$P(Y_1=1) = \pi_1 \Leftrightarrow P(Y_1=y) = \pi_1^y (1-\pi_1)^{1-y}$
\vdots	\vdots	
y_n	x_n	$P(Y_1=0) = 1-\pi_1 \quad g_1(y) \uparrow$
		Independent

Then the joint pmf of y_1, \dots, y_n

$$g(y_1, \dots, y_n) = \prod_{i=1}^n g_i(y_i) = \prod_{i=1}^n \pi_i^{y_i} (1-\pi_i)^{1-y_i}$$

$$\log[g(y_1, \dots, y_n)] = \sum_{i=1}^n [y_i \log \pi_i + (1-y_i) \log(1-\pi_i)]$$

$$L(\beta_0, \beta_1) = \sum_{i=1}^n \left[y_i \log \frac{\pi_i}{1-\pi_i} + \log(1-\pi_i) \right]$$

$$= \sum_{i=1}^n \left[y_i (\beta_0 + \beta_1 x_i) - \log(1 + e^{\beta_0 + \beta_1 x_i}) \right]$$

$$\pi_i = \frac{e^{\beta_0 + \beta_1 x_i}}{1 + e^{\beta_0 + \beta_1 x_i}}$$

To obtain MLE $\hat{\beta}_0$ and $\hat{\beta}_1$ of

β_0 and β_1 , we need to maximize

$$\log L(\beta_0, \beta_1) = \sum_{i=1}^n y_i (\beta_0 + \beta_1 x_i) - \log(1 + e^{\beta_0 + \beta_1 x_i})$$

"No closed form expression"

Ex (Effect of programming experience in successfully completing projects in time)

$$y_i = \begin{cases} 1 & \text{if } i\text{th candidate completed project in time} \\ 0 & \text{otherwise} \end{cases}$$

x_i = amount of programming experience (in months), i is index of person
The dataset is for 25 candidates (available in book).
($i=1, \dots, 25$)

Logistic regression was used $\hat{\beta}_1 = 0.615$, $\hat{\beta}_0 = -3.05$.

so the fitted logistic regression model:

$$\log \frac{\pi}{1-\pi} = -3.05 + 0.615x = \beta_0 + \beta_1 x$$

$$\text{In this case, } \widehat{OR} = e^{\hat{\beta}_1} = e^{0.615} = 1.175.$$

For two persons, P_1, P_2 , where P_2 has one month more experience.

$$\frac{\text{Odds}_2}{\text{Odds}_1} = 1.175 \Leftrightarrow \frac{\text{Odds}_2 - \text{Odds}_1}{\text{Odds}_1} = 0.175 = 17.5\%$$

so for each additional month of experience, the odds for success increase by 17.5%.
↑ ratio $P(y_i=1)$ per $P(y_i=0)$.

I have two future candidates having 10 months and 15 months of experiences.

$$\pi\left(\frac{15}{10}\right)$$

Relative difference of odds

$$\widehat{OR} = e^{5\hat{\beta}_1}, 5 = 15 - 10 \quad // \quad \frac{\text{Odds } 15}{\text{Odds } 10} \\ = e^{0.8075} = 2.24$$

$$\frac{\text{Odds}_2 - \text{Odds}_1}{\text{Odds}_1} = 1.24$$

$$e^{0.615 \cdot 5} = \frac{e^{0.615 \cdot 5}}{e^{0.615 \cdot 5}} = \boxed{2.24}$$

$$\frac{e^{0.615 \cdot 15}}{e^{0.615 \cdot 10}} = \boxed{2.24}$$

• Logistic Regression

last time Example showing the interpretation of the parameter β_1 , connection with odd Ratio.

so far, we consider the case where one y value is present for each x value. But there can be multiple y values corresponding to one x value.

Ex (coupon distribution)

A company gives out 5 different coupons (\$5, \$10, \$15, \$20, \$30) to 200 customers for each category.

count how many customers have redeemed those coupons.

(#) coupon amount	# customer given	# customers who redeemed
5	200	55
10	200	74
15	200	95
20	200	122
30	200	144

For $i=1, 2, \dots, 5$, $Y_i \sim \text{Bin}(200, \pi_i)$, where $\pi_i = \frac{e^{\beta_0 + \beta_1 x_i}}{1 + e^{\beta_0 + \beta_1 x_i}}$

$$\text{So, } \log \frac{\pi_i}{1 - \pi_i} = \beta_0 + \beta_1 x_i = \pi_i^{n=200}$$

pmf of Y_i is $g_i(y_i) = P(Y_i = y_i) = \binom{n}{y_i} \pi_i^{y_i} (1 - \pi_i)^{n-y_i}$

$$P(Y_i \leq y_i) = \sum_{k=0}^{y_i} \binom{n}{k} \pi_i^k (1 - \pi_i)^{n-k} = \binom{200}{y_i} \pi_i^{y_i} (1 - \pi_i)^{200-y_i}$$

cdf

• Joint pmf of y_1, \dots, y_5

$$f(y_1, \dots, y_5) = \prod_{i=1}^5 \left(\frac{200}{y_i} \right)^{y_i} \pi_i^{y_i} (1-\pi_i)^{200-y_i}$$

$$\log g(y_1, \dots, y_5) = \sum_{i=1}^5 \left[\log \left(\frac{200}{y_i} \right) + y_i \log \pi_i + (200-y_i) \log (1-\pi_i) \right]$$

$$= \sum_{i=1}^5 \left[\log \left(\frac{200}{y_i} \right) + y_i \log \frac{\pi_i}{1-\pi_i} + 200 \log (1-\pi_i) \right]$$

so, the log-likelihood function "maximize this function"

$$\ell(\beta_0, \beta_1) = \log L(\beta_0, \beta_1)$$

$$= C + \sum_{i=1}^5 \left[y_i (\beta_0 + \beta_1 x_i) - 200 \log (1 + e^{\beta_0 + \beta_1 x_i}) \right]$$

Next, we maximize the above function with respect to β_0, β_1 using numerical method.

The maximizer will be the MLE. find $\hat{\beta}_0, \hat{\beta}_1$

• The Interpretation of β_1

$$\log \frac{\pi_i}{1-\pi_i} = \beta_0 + \beta_1 x_i$$

so, if $\tilde{x}_i = x_i + 1$ and $\tilde{\pi}_i, \pi_i$ are the corresponding probabilities.

$$\log \frac{\tilde{\pi}_i}{1-\tilde{\pi}_i} - \log \frac{\pi_i}{1-\pi_i} = \beta_1 (\tilde{x}_i - x_i) = \boxed{\beta_1}$$

$$\beta_0 + \beta_1 (x+1) \quad \beta_0 + \beta_1 x$$

$$= \log \frac{\tilde{\pi}_i / (1-\tilde{\pi}_i)}{\pi_i / (1-\pi_i)}$$

$$\text{so, } OR = e^{\beta_1}, \widehat{OR} = e^{\hat{\beta}_1} = \boxed{\log(OR)}$$

EX Suppose $\hat{\beta}_1 = 0.01$. What is the % increase in odds if \$12 is offered instead of \$4.

$$\begin{aligned} \text{The odds ratio } \widehat{OR} &= e^{\hat{\beta}_1} = e^{0.01} = 1.083. \\ &\text{estimated} \quad \text{Odds}_{12} = 1.083 \quad \text{so, } \frac{\text{Odds}_{12} - \text{Odds}_4}{\text{Odds}_4} = \frac{1.083 - 1}{1.083 - 1} = 0.083 \end{aligned}$$

So, the odds will increase by 8.3% for any x_i .

π_i depends on x_i
chance

(odds does not depend on x_i)

• Multiple logistic Regression

Here, we have $\underbrace{p-1 \text{ predictors}}$ and one response.

Again, y takes two values: 0 or 1

x_i 's are either categorical or continuous variables.

EX (Disease occurrence)

$$y = \begin{cases} 1 & \text{if individual infected} \\ 0 & \text{if not} \end{cases}$$

Four predictors

Two continuous (Age, cholesterol level)
 Two categorical (Socio economic, city status, sector)

		Socio economic status		City sector	
		Middle-X3	X4-lower	Sector 1	X5-sector
Upper	0	0		0	
Middle	1	0		1	
Lower	0	1			

so, the predictor $\tilde{x} = (1, x_1, \dots, x_5)'$

so, we have information about n individuals

we know $(y_i, \tilde{x}_i), i=1, \dots, n$.

The logistic regression model is as follows.

$$y_i \sim \text{Ber}(\pi_i), \text{ where } \pi_i = \frac{e^{\tilde{x}_i \beta}}{1 + e^{\tilde{x}_i \beta}}, \beta = \begin{bmatrix} \beta_0 \\ \beta_1 \\ \vdots \\ \beta_5 \end{bmatrix}$$

so, the pmf of y_i is $g_i(y_i) = \pi_i^{y_i} (1 - \pi_i)^{1-y_i}$

so, the joint pmf $g(y_1, \dots, y_n) = \prod_{i=1}^n \pi_i^{y_i} (1 - \pi_i)^{1-y_i}$

$$\begin{aligned} \text{so, } \log g(y_1, \dots, y_n) &= \sum_{i=1}^n \left[y_i \log \frac{\pi_i}{1 - \pi_i} + \log(1 - \pi_i) \right] \\ &= \sum_{i=1}^n \left[y_i (\tilde{x}_i \beta) - \log(1 + e^{\tilde{x}_i \beta}) \right] \end{aligned}$$

Interpretation of the parameters

e^{β_i} is the odd ratio between two cases where everything except the i th predictor is kept fixed and the i th predictor value is increased by 1.

Suppose $\hat{\beta}_1$ is such that $e^{\hat{\beta}_1} = 1.0309$

this means that the odds for having the disease will go up by 3.09% for each year of age.

Suppose

$$e^{\hat{\beta}_5} = 4.82 \quad (\text{if change "sector", odds of disease 5 times. e.g. 5} \rightarrow 25)$$

then the odds of having the disease becomes approx. five time in sector 2 than the odds in sector 1.

Inference about the parameters

The inference in this case is for large sample size

The MLE $\hat{\beta}$ approximately follows normal distribution with mean β and covariance $-H(\hat{\beta})$, where

$$H_{ij} = \frac{\partial^2 \log L(\beta)}{\partial \beta_i \partial \beta_j} \Big|_{\hat{\beta}} \quad (\text{If } \hat{\beta} \text{ given, can you find covariance matrix?})$$

This will allow us to obtain confidence intervals and simultaneous CI for different components of β .

Math B7800

4/17/18, Tue

We are discussing logistic regression.

Ex In a simple logistic regression, ($\log \frac{\pi}{1-\pi} = \beta_0 + \beta_1 x$) if the estimate of β_0 and β_1 are -3 and 1.2, i.e., $\hat{\beta}_0 = -3$, $\hat{\beta}_1 = 1.2$ then What is the probability that among 100 future individuals having $x=5$, there will be at least 60 individuals with positive response.

$$P(\text{an individual having } x=5 \text{ gives positive response}) = \frac{e^{\beta_0 + \beta_1 x}}{1 + e^{\beta_0 + \beta_1 x}} = \pi(5)$$

so, # individuals among the future 100 individuals who will give positive response. - *

$$\sim \text{Bin}(100, \pi(5))$$

$$P(\# * \geq 60) = \sum_{k=60}^{100} \binom{100}{k} \pi(5)^k (1-\pi(5))^{100-k}$$

$$\widehat{P}(\# * \geq 60) \quad \leftarrow \text{plug in } \hat{\beta}_0 = -3, \hat{\beta}_1 = 1.2$$

Recall: If $\hat{\theta}$ is an MLE of θ , then for any function g , $g(\hat{\theta})$ is MLE of $g(\theta)$.

so, MLE of the probability is

$$\sum_{k=60}^{100} \binom{100}{k} (\hat{\pi}(5))^k (1-\hat{\pi}(5))^{100-k}, \text{ where}$$

$$\hat{\pi}(5) = \frac{e^{-3+1.2 \times 5}}{1 + e^{-3+1.2 \times 5}} = \frac{e^3}{1 + e^3}$$

M376

Q2: use normal approximation to approximate the estimate.

Recall: $\text{Bin}(n, p) \approx N(np, npq)$, $q = 1-p$.

so the normal approximation is $P(N(100\hat{\pi}(5), 100\cdot\hat{\pi}(5)(1-\hat{\pi}(5))) \geq 60)$

$$= P(N(95.26, 4.5) \geq 60)$$

$$= 1 - \Phi\left(\frac{60 - 95.26}{\sqrt{4.5}}\right)$$

standard deviation

Inference for the parameters in logistic regression.

The inference is correct for large sample size $\hat{\beta}$ is approximately normal with mean β and covariance matrix $[-H(\hat{\beta})]^{-1}$

Where $H(\beta)$ is the matrix of second derivatives

$$H_{ij}(\beta) = \frac{\partial^2 \log L(\beta)}{\partial \beta_i \partial \beta_j}, i, j = 0, 1, \dots, p-1$$

this is correct.
check the last note.

Using this approximation,

We obtain CI for β_k , simultaneous CI for the parameters $\beta_0, \dots, \beta_{k-1}$.

Also, we can test $H_0: \beta_k = 0 \text{ vs } H_1: \beta_k \neq 0$ - (a)

(b) $H_0: \beta_k = 0 \text{ vs } H_1: \beta_k > 0$

(c) $H_0: \beta_k = 1 \text{ vs } H_1: \beta_k \neq 1$

Q: CI for β_k ?

Standard normal $N(0, 1)$.

$$100(1-\alpha)\% \text{ CI for } \beta_k = \hat{\beta}_k \pm Z_{\alpha/2} \sqrt{(-H(\hat{\beta}))_{kk}^{-1}} \quad - (a)$$

• How to test $H_0: \beta_k = 0 \text{ vs } H_1: \beta_k > 0$.

Reject H_0 at level α if $\hat{\beta}_k > Z_\alpha \sqrt{(-H(\hat{\beta}))_{kk}^{-1}}$ - (b)
one side

• How to test $H_0: \beta_k = 1 \text{ vs } H_1: \beta_k \neq 1$.

Reject H_0 at level α if $|\hat{\beta}_k - 1| > Z_{\alpha/2} \sqrt{(-H(\hat{\beta}))_{kk}^{-1}}$ - (c)

• $H_0: \beta_k = 1 \text{ vs } H_1: \beta_k < 1$?

Reject H_0 at level α if $\hat{\beta}_k < 1 - Z_\alpha \sqrt{(-H(\hat{\beta}))_{kk}^{-1}}$

• Likelihood Ratio test

$q < p$,

$$H_0: \beta_q = \beta_{q+1} = \dots = \beta_{p-1} = 0 \text{ vs } H_1: H_0 \text{ is false}$$

$$\Delta = \text{likelihood ratio} = \frac{\sup_{H_0} L(\beta)}{\sup_{H_1} L(\beta)}$$

Note that under H_0 , we also have a logistic regression model

$$\log \frac{\pi}{1-\pi} = \beta_0 + \beta_1 x_1 + \dots + \beta_{q-1} x_{q-1}$$

using similar approach, we can find MLE of $\beta_0, \dots, \beta_{q-1}$ under H_0 .

call the MLE $\hat{\beta}_{0,q}$.

$$\text{so, } \Delta = \frac{L([\hat{\beta}_{0,q}])}{L(\hat{\beta})} \quad \text{small } \rightarrow \text{reject } H_0.$$

$$-2\log \Delta \approx \chi^2_{p-q}.$$

reject H_0 if $-2\log \Delta > \chi^2_{p-q, \alpha}$
at level α

• Model selection

Define $AIC = -2\log L(\hat{\beta}) + 2(\# \text{ variables})$ ← minimize

We compare all possible subsets of predictors, when the total # of predictors is ≤ 30 or 40.

If the number of predictors is more than 40,

We should use stepwise methods for finding the optional model.

• Principal Component Analysis (PCA)

population principal component

Suppose $\underline{x} = \begin{bmatrix} x_1 \\ \vdots \\ x_p \end{bmatrix}$ is a vector consisting of p features of an individual.

We want to choose linear combinations (any combination of \underline{x})

$$y_1 = a_1' \underline{x} \quad (a_1, \dots, a_p \text{ are } p \times 1 \text{ vector of constants})$$

so that y_1, \dots, y_p are uncorrelated and

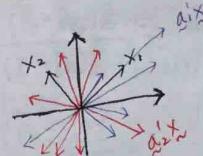
"population" y_1, \dots, y_p have max possible variance.

Suppose Σ is covariance matrix of $\underline{x} = [x_1, \dots, x_p]'$

Then the constant vectors will depend on Σ .

These y_i 's are called (population) principal component of \underline{x}

y_1 is the first principal component, y_2 is the second and so on.



Math B7800

4/19/18, Thur

Question 3 from Exam 1

Let $n=50$, $R^2=0.6$.
 $H_0: \beta_1 = \dots = \beta_{10} = 0$ vs $H_1: H_0$ is false.

$$Y = X\beta + \varepsilon$$

$$= \beta_0 + \beta_1 X_1 + \dots + \beta_{10} X_{10} + \varepsilon$$

$$\begin{aligned} Z &= \begin{bmatrix} 1 & Z_{1,1} & \dots & Z_{1,10} \\ \vdots & \vdots & & \vdots \\ 1 & Z_{50,1} & \dots & Z_{50,10} \end{bmatrix} = \begin{bmatrix} Z_1 \\ Z_2 \end{bmatrix} & \beta &= \begin{bmatrix} \beta_0 \\ \beta_1 \\ \vdots \\ \beta_{10} \end{bmatrix} = \begin{bmatrix} \beta_0 \\ \beta_{(1)} \\ \vdots \\ \beta_{(2)} \end{bmatrix} \end{aligned}$$

$$\text{So, } Z_1 = \begin{bmatrix} 1 \\ \vdots \\ 1 \end{bmatrix}, \quad Z_2 = \begin{bmatrix} Z_{1,1} & \dots & Z_{1,10} \\ \vdots & \ddots & \vdots \\ Z_{50,1} & \dots & Z_{50,10} \end{bmatrix}$$

So, $H_0: \beta_{(2)} = 0$ vs $H_1: \beta_{(2)} \neq 0$.

LRT: $\frac{Y'(P_Z - P_{Z_1})Y}{Y'(I - P_Z)Y}$ Reject H_0 if LRT is large
 $(P_{Z_1} = P_1)$ We would reject H_0 at 5% level

$$= \frac{Y'(P_Z - P_1)Y}{Y'(I - P_Z)Y} = R^2 \quad \text{If } R^2 > \frac{q}{40} F_{q,40}(0.05).$$

• Principal Component Analysis (PCA)

Suppose \underline{x} is a random vector with mean $\underline{\mu}$ and covariance matrix Σ .

\underline{x} represents p features of the individuals in a population.
 p is usually large.

Goal: choose fewer (less than p) linear combinations of features without losing much information. $\rightarrow \underline{y}_1, \dots, \underline{y}_n$.

Steps Find linear combinations $\underline{a}_1, \dots, \underline{a}_p$ s.t.

$$\underline{y} = \begin{bmatrix} \underline{a}_1' \\ \vdots \\ \underline{a}_p' \end{bmatrix} \underline{x} \quad \text{normalized} \quad \underline{y}_i = \underline{a}_i' \underline{x}$$

\underline{y}_i must have uncorrelated components and max possible variance.

Step 1 • First principal component maximize $\underline{a}' \Sigma \underline{a} = \text{var}(\underline{a}' \underline{x})$ w.r.t. \underline{a} s.t. $\underline{a}' \underline{a} = 1$.

We saw $\max_{\substack{\underline{a}: \underline{a}' \underline{a} = 1}} \underline{a}' \Sigma \underline{a} = \lambda_1 \leftarrow \text{largest eigenvalue of } \Sigma$

$$\max_{\substack{\underline{a}: \underline{a}' \underline{a} = 1}} \underline{a}' \Sigma \underline{a} \quad \text{by spectral decomposition.}$$

\underline{a}_1 = eigenvector of Σ corresponding to λ_1 achieves the max.

• Recall: Spectral decomposition of real symmetric matrix.

$A = U \Delta U'$ for any symmetric real matrix A , where U is orthogonal and Δ is diagonal.

The diagonal entries of Δ are the eigenvalues of A and $U' = U^{-1}$. $UU' = I$.

$$\text{If } U = \begin{bmatrix} \underline{u}_1 & \cdots & \underline{u}_p \end{bmatrix}_{p \times p}$$

$$U \Delta U' = \begin{bmatrix} \underline{u}_1 & \cdots & \underline{u}_p \end{bmatrix} \begin{bmatrix} \lambda_1 & & 0 \\ 0 & \ddots & \\ & & \lambda_p \end{bmatrix} \begin{bmatrix} \underline{u}_1 \\ \vdots \\ \underline{u}_p \end{bmatrix}$$

$$= \sum_{i=1}^p \lambda_i \underline{u}_i \underline{u}_i' = A$$

Let $\Sigma = U \Delta U'$ be the spectral decomposition where

$$\Delta = \begin{bmatrix} \lambda_1 & & 0 \\ 0 & \ddots & \\ & & \lambda_p \end{bmatrix} \quad \text{s.t. } \lambda_1 \geq \cdots \geq \lambda_p \geq 0 \leftarrow \text{because } \Sigma \text{ is always hnd.}$$

Note: $\underline{u}_1, \dots, \underline{u}_p$ form a basis of \mathbb{R}^p .

so, any $\underline{a} \in \mathbb{R}^p$ can be expressed as $\underline{a} = \underline{x}' \underline{u}$

$$\underline{a} = x_1 \underline{u}_1 + x_2 \underline{u}_2 + \cdots + x_p \underline{u}_p \quad \text{for some constants } x_1, \dots, x_p.$$

$$\underline{a} \cdot \underline{u}_i = x_1 \cdot 1 + x_2 \cdot 0 + \cdots + x_p \cdot 0 = x_i$$

Similarly, $x_i = (\underline{a} \cdot \underline{u}_i)$ for $i = 1, \dots, p$.

Then any $\underline{a} \in \mathbb{R}^p$ can be written as $\underline{a} = \sum_{i=1}^p (\underline{a} \cdot \underline{u}_i) \underline{u}_i$

$$\text{Then } \underline{\alpha}' \underline{\alpha} = \left(\sum_{i=1}^P (\underline{\alpha} \cdot \underline{u}_i) \underline{u}_i' \right) \left(\sum_{j=1}^P (\underline{\alpha} \cdot \underline{u}_j) \underline{u}_j' \right)$$

$$= \sum_{i=1}^P (\underline{\alpha} \cdot \underline{u}_i)^2$$

$$\underline{\alpha}' \Sigma \underline{\alpha} = \left(\sum_{i=1}^P (\underline{\alpha} \cdot \underline{u}_i) \underline{u}_i' \right) \left(\sum_{j=1}^P \lambda_j \underline{u}_j \underline{u}_j' \right) \left(\sum_{k=1}^P (\underline{\alpha} \cdot \underline{u}_k) \underline{u}_k' \right)$$

$$(\underline{\alpha} \cdot \underline{u}_i) \underline{u}_i' \lambda_j \underline{u}_j \underline{u}_j' (\underline{\alpha} \cdot \underline{u}_k) \underline{u}_k'$$

$$= \dots \quad \underline{u}_i' \underline{u}_j \underline{u}_j' \underline{u}_k = 0 \text{ if } i \neq j \text{ or } j \neq k$$

$$\Rightarrow = \sum_{i=1}^P (\underline{\alpha} \cdot \underline{u}_i) \lambda_i (\underline{\alpha} \cdot \underline{u}_i) = \sum_{i=1}^P \lambda_i (\underline{\alpha} \cdot \underline{u}_i)^2 \quad \text{---} \circledast$$

So, what does the problem reduce to?

$$\text{call } x_i = (\underline{\alpha} \cdot \underline{u}_i) \quad \text{convex combination}$$

$$\max_{x_1, \dots, x_p} \frac{\lambda_1 x_1^2 + \dots + \lambda_p x_p^2}{x_1^2 + \dots + x_p^2} = \lambda_1 \frac{x_1^2}{x_1^2 + \dots + x_p^2} + \dots + \lambda_p \frac{x_p^2}{x_1^2 + \dots + x_p^2}$$

$$= x_1 = \lambda_1 \alpha_1 + \dots + \lambda_p \alpha_p, \text{ where } \alpha_1 + \dots + \alpha_p = 1.$$

This shows that $y_1 = \underline{u}_1' \underline{x}$ and $\text{var}(y_1) = \lambda_1$.

$$\text{Second principal component: } y_2 = \underline{\alpha}_2' \underline{x} \text{ s.t. } \text{cov}(y_2, y_1) = \underline{\alpha}_2' \sum \underline{u}_1 \underbrace{\underline{u}_1'}_{\text{eigen vector of } \Sigma}$$

$$= \underline{\alpha}_2' (\lambda_1 \underline{u}_1)$$

$$= \lambda_1 \underline{\alpha}_2' \underline{u}_1 = 0$$

$$\text{so, } \underline{\alpha}_2' \underline{u}_1 = 0.$$

Now, we need to maximize $\frac{\underline{\alpha}_2' \sum \underline{u}_2}{\underline{\alpha}_2' \underline{u}_2}$ subject to the constraint $\underline{\alpha}_2' \underline{u}_1 = 0$.

$$\underline{\alpha}_2 = \sum_{i=1}^P \frac{x_i}{\underline{\alpha}_2' \underline{u}_2} (\underline{\alpha}_2 \cdot \underline{u}_i) \underline{u}_i \quad \text{(normalize)}$$

$$\text{If } \underline{\alpha}_2' \underline{u}_1 = 0, \text{ then } \underline{\alpha}_2 = \sum_{i=2}^P (\underline{\alpha}_2 \cdot \underline{u}_i) \underline{u}_i$$

$$\text{Also, } \underline{\alpha}_2' \sum \underline{u}_2 = \sum_{i=1}^P \lambda_i (\underline{\alpha}_2 \cdot \underline{u}_i)^2 = \sum_{i=2}^P \lambda_i (\underline{\alpha}_2 \cdot \underline{u}_i)^2$$

$$\frac{\underline{\alpha}_2' \sum \underline{u}_2}{\underline{\alpha}_2' \underline{u}_2} = \frac{\lambda_2 (\underline{\alpha}_2 \cdot \underline{u}_2)^2 + \dots + \lambda_p (\underline{\alpha}_2 \cdot \underline{u}_p)^2}{(\underline{\alpha}_2 \cdot \underline{u}_2)^2 + \dots + (\underline{\alpha}_2 \cdot \underline{u}_p)^2}$$

So, using similar argument,

$$\max_{\substack{\underline{\alpha}_2 \neq 0 \\ \underline{\alpha}_2' \underline{u}_1 = 0}} \frac{\underline{\alpha}_2' \sum \underline{u}_2}{\underline{\alpha}_2' \underline{u}_2} = \lambda_2$$

$\underline{\alpha}_2 = \underline{u}_2$ satisfies

$$\underline{\alpha}_2' \underline{u}_1 = 0 \text{ and } \frac{\underline{\alpha}_2' \sum \underline{u}_2}{\underline{\alpha}_2' \underline{u}_2} = \lambda_2$$

Thus, $y_2 = \underline{u}_2' \underline{x}$, and $\text{var}(y_2) = \lambda_2$.

\uparrow
eigenvector
of Σ

\uparrow
eigenvalue
of Σ

~~Fact~~

Following similar argument, $Y_i = U_i' \bar{X}$, $i=1, \dots, p$

$$\text{Var}(Y_i) = \lambda_i \quad \text{trace}(\Delta)$$

$$\sum_{i=1}^p \text{Var}(Y_i) = \sum_{i=1}^p \lambda_i = \text{trace}(\Sigma) = \sum_{i=1}^p \sigma_{ii}$$

$$\begin{aligned} \text{tr}(\Sigma) &= \text{tr}(P\Delta P') \\ &= \text{tr}(\Delta P') = \text{tr}(\Delta). \end{aligned}$$

$\sum_{i=1}^p \frac{\text{Var}(X_i)}{\sum_{ii}} \left(\begin{array}{l} \text{total variance} \\ \text{doesn't change} \end{array} \right)$

• proportion of total explained by i th principal component (PC)

$$= \frac{\lambda_i}{\lambda_1 + \dots + \lambda_p} \quad (\text{this is called portion of variance explained by } i\text{th PC})$$

We will choose first k PC if the corresponding explained variance $\frac{\lambda_1 + \dots + \lambda_k}{\lambda_1 + \dots + \lambda_p}$ is close to 1.

If \bar{X} has mean \bar{M} , covariance Σ , then the (population) PCs are

$$Y_1 = U_1' (\bar{X} - \bar{M}), \dots, Y_p = U_p' (\bar{X} - \bar{M}),$$

where U_1, \dots, U_p are the eigenvectors of Σ .

• PC for standardized random vector.

Suppose $\bar{X}_{p \times 1}$ has mean $\bar{0}$, covariance Σ

Define $Z_i = \frac{X_i - \bar{M}_i}{\sqrt{\sigma_{ii}}}$, $\bar{Z} = \begin{bmatrix} Z_1 \\ \vdots \\ Z_p \end{bmatrix}$

$$\begin{aligned} \text{cov}(Z_i, Z_j) &= \text{cov}\left(\frac{X_i - \bar{M}_i}{\sqrt{\sigma_{ii}}}, \frac{X_j - \bar{M}_j}{\sqrt{\sigma_{jj}}}\right) \\ &= \frac{\sigma_{ij}}{\sqrt{\sigma_{ii}} \sqrt{\sigma_{jj}}} = P_{ij} \end{aligned}$$

$$\mathbb{E}(\bar{Z}) = \bar{0}$$

$$\text{cov}(\bar{Z}) = P = \begin{bmatrix} 1 & P_{12} & P_{13} & \cdots & P_{1p} \\ P_{21} & 1 & & & | \\ \vdots & & & & | \\ P_{p1} & P_{p2} & \cdots & \cdots & 1 \end{bmatrix}$$

Let $V = \text{Diag}(\Sigma)$

$$= \begin{bmatrix} \sigma_{11} & 0 \\ 0 & \ddots & 0 \\ & \ddots & \sigma_{pp} \end{bmatrix}$$

$$= (V^{1/2})^{-1} \Sigma (V^{1/2})^{-1}$$

The PCs of \bar{Z} are $\tilde{Y}_1, \dots, \tilde{Y}_p$, where

$$\tilde{Y}_i = \tilde{U}_i \bar{Z}, \text{ where } \tilde{U}_1, \dots, \tilde{U}_p \text{ are eigenvectors of } P$$

$$= \tilde{U}_i (V^{1/2})^{-1} (\bar{X} - \bar{M})$$

Math B7800

09224

4/24/18, Tue

Principal Components

Last time population PC for standardized covariance.

Special covariance matrices

case1: when Σ is diagonal, $\Sigma = \begin{bmatrix} \sigma_{11} & 0 \\ 0 & \dots \\ 0 & \sigma_{pp} \end{bmatrix}$

Here, the eigenvalues of Σ are $\sigma_{11}, \dots, \sigma_{pp}$, we order them in decreasing order $\sigma_{(1)} \geq \dots \geq \sigma_{(p)}$.

$\mathbf{x} = \begin{pmatrix} x_1 \\ \vdots \\ x_p \end{pmatrix}$ • First principal component = $(X_{(1)} - \mu_{(1)})$

• The i th PC is $(X_{(i)} - \mu_{(i)})$

Ex If $\Sigma = \begin{bmatrix} 2 & 0 \\ 0 & 3 \end{bmatrix}$ and \mathbf{x} has mean μ and cov Σ .

then $PC_1 = x_2 - \mu_2$

$PC_2 = x_1 - \mu_1$

$PC_3 = x_2 - \mu_2$

corresponding correlation matrix = I . \leftarrow all eigenvalues are equal important.

Case 2: All pairwise correlations of the components are equal to ρ .

Suppose $\text{var}(X_i) = \sigma_{ii}$.

then $\text{cov}(X_i, X_j) = \rho \sqrt{\sigma_{ii}} \sqrt{\sigma_{jj}}$ if $i \neq j$

so, $\Sigma = \begin{bmatrix} \sigma_{11} & \rho \sqrt{\sigma_{11}} \sqrt{\sigma_{22}} & \dots & \rho \sqrt{\sigma_{11}} \sqrt{\sigma_{pp}} \\ \vdots & \ddots & \ddots & \vdots \\ \rho \sqrt{\sigma_{pp}} \sqrt{\sigma_{11}} & \dots & \sigma_{pp} \end{bmatrix}$

corresponding correlation matrix = $\begin{bmatrix} 1 & \rho & \rho & \dots & \rho \\ \rho & 1 & \rho & \dots & \rho \\ \rho & \rho & 1 & \dots & \rho \end{bmatrix} = \rho$

Step 1: $\begin{bmatrix} 1 & \dots & 1 \\ \vdots & \ddots & \vdots \\ 1 & \dots & 1 \end{bmatrix}_{n \times n} = \frac{1}{\sqrt{n}} \frac{1}{\sqrt{n}}$, $\frac{1}{\sqrt{n}}$ is $n \times 1$ vector of all 1

$\frac{1}{\sqrt{n}} \frac{1}{\sqrt{n}} = n$
the eigenvalues of $\frac{1}{\sqrt{n}} \frac{1}{\sqrt{n}}$ are $(n, 0, 0, \dots, 0)$

$(\frac{1}{\sqrt{n}} \frac{1}{\sqrt{n}}) \cdot \frac{1}{\sqrt{n}} = n \cdot \frac{1}{\sqrt{n}}$. So, $\frac{1}{\sqrt{n}}$ is an eigenvector of $\frac{1}{\sqrt{n}} \frac{1}{\sqrt{n}}$ corresponding to eigenvalue n .

If \underline{u} is an eigenvector of $\frac{1}{\sqrt{n}} \frac{1}{\sqrt{n}}$ corresponding to 0, then

$$(\frac{1}{\sqrt{n}} \frac{1}{\sqrt{n}}) \underline{u} = 0 = \frac{1}{\sqrt{n}} (\frac{1}{\sqrt{n}} \underline{u}) = (\frac{1}{\sqrt{n}} \underline{u}) \frac{1}{\sqrt{n}}. \text{ So, } \frac{1}{\sqrt{n}} \underline{u} = 0$$

Def the eigenspace of a matrix A corresponding to eigenvalue λ is $\{\underline{x} : A\underline{x} = \lambda \underline{x}\}$.

In our case, the eigenspace of $\frac{1}{\sqrt{n}} \frac{1}{\sqrt{n}}$ corresponding to eigenvalue λ is $\text{Null}(\frac{1}{\sqrt{n}} \frac{1}{\sqrt{n}})$.
 $\dim(N(A)) = \dim \text{col}(A) - \text{rank}(A)$.

In this case, $\dim(\text{eigenspace}) = n - 1$.

Any basis of $N(\frac{1}{\sqrt{n}} \frac{1}{\sqrt{n}})$, or in particular, any set of mutually orthogonal vectors $\underline{u}_1, \dots, \underline{u}_{n-1}$ satisfying $\frac{1}{\sqrt{n}} \underline{u}_i = 0$, $i=1, \dots, n-1$, and be used as $P\underline{c}_2, \dots, P\underline{c}_n$ for $\frac{1}{\sqrt{n}} \frac{1}{\sqrt{n}}$.

$$\left[\frac{1}{\sqrt{2}}, \frac{-1}{\sqrt{2}}, 0, \dots, 0 \right] \text{ "normalized"} \quad \left[\frac{1}{\sqrt{2 \times 3}}, \frac{1}{\sqrt{2 \times 3}}, \frac{-2}{\sqrt{2 \times 3}}, 0, \dots, 0 \right] \text{ "mutually orthogonal"} \quad \left[\frac{1}{\sqrt{(i-1)i}}, \frac{1}{\sqrt{(i-1)i}}, \dots, \frac{1}{\sqrt{(i-1)i}}, \frac{-(i-1)}{\sqrt{(i-1)i}}, 0, 0, \dots, 0 \right] = \underline{e}_i \text{ orthogonal.}$$

for $i = 2, 3, \dots, n$.

then $\underline{e}_1 = \frac{1}{\sqrt{n}} \frac{1}{\sqrt{n}}$, $\underline{e}_2, \dots, \underline{e}_n$ are independent eigenvectors of $\frac{1}{\sqrt{n}} \frac{1}{\sqrt{n}}$.

In the example, the correlation matrix was

$$\begin{bmatrix} 1 & p & \dots & p \\ p & 1 & \dots & \vdots \\ \vdots & \ddots & \ddots & \vdots \\ p & \dots & \dots & 1 \end{bmatrix}_{p \times p} = p \mathbf{1} \mathbf{1}' / p + (1-p) \mathbf{I}_p$$

If $\lambda_1, \dots, \lambda_p$ are the eigenvalues of $A_{p \times p}$ with corresponding eigenvectors u_1, \dots, u_p , then what are the eigenvalues and eigenvectors of $a\mathbf{I}_p + bA$?

$a+b\lambda_1, a+b\lambda_2, \dots, a+b\lambda_p$ are the eigenvalues with eigenvectors u_1, \dots, u_p .

$$Au_i = \lambda_i u_i$$

$$bAu_i = b\lambda_i u_i$$

$$au_i + bAu_i = au_i + b\lambda_i u_i = (a+b\lambda_i)u_i$$

$$(aI + bA)u_i$$

so, the eigen values of $p \mathbf{1} \mathbf{1}' / p + (1-p) \mathbf{I}_p$ are

$$p \cdot p, 1-p, 1-p, \dots, 1-p.$$

$$\frac{1}{1+(p-1)p}$$

$$\text{so, } PC_1 = \left(\frac{1}{\sqrt{p}}, \frac{1}{\sqrt{p}}, \dots, \frac{1}{\sqrt{p}} \right)' (\bar{x} - \mu)$$

$$\vdots \\ PC_i = \left(\frac{1}{\sqrt{(i-1)i}}, \frac{1}{\sqrt{(i-1)i}}, \dots, \frac{1}{\sqrt{(i-1)i}}, \frac{-(i-1)}{\sqrt{(i-1)i}}, 0, \dots, 0 \right)' (\bar{x} - \mu)$$

for $i=2, \dots, p$.

proportion of variance explained by the first component.

total variance : sum of eigenvalues.

total variance = p .

$$\text{so, the proportion} = \frac{1+(p-1)p}{p} = p + \frac{1-p}{p}$$

If p is large and/or p is large, then this proportion is quite large.

If $p=0.8, p=5$, then the proportion is 0.84.

so, it is enough to consider the first component.

Sample PC (Principal Components)

Suppose x_1, \dots, x_n are n iid sample from a population with mean μ_{px1} and covariance Σ_{pxp} .

Obtain sample mean \bar{x} and sample cov $S = \frac{1}{n-1} \sum_{i=1}^n (x_i - \bar{x})(x_i - \bar{x})'$. Then obtain eigenvalues $\hat{\lambda}_1 \geq \dots \geq \hat{\lambda}_p$ of S with corresponding eigenvectors $\hat{e}_1, \dots, \hat{e}_p$.

then the i th sample PC is $\hat{e}_i' (x - \bar{x})$

If S is the sample covariance matrix, then for any constant vector a , what is the sample variance of $a' x_1, \dots, a' x_n$?

$$\text{Ans } a' S a$$

For two constant vectors a, b , sample cov between $a' x, b' x$ is $a' S b$.

keeping these in mind, the sample PC maximizes sample variance subject to uncorrelated components.

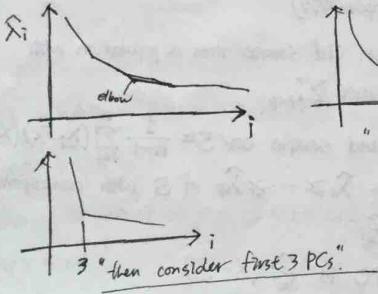
so, we want to find linear combinations $\hat{c}_1, \dots, \hat{c}_p$ s.t.

if $y_i = \hat{c}_i(x - \bar{x})$, then the components $y = \begin{pmatrix} y_1 \\ \vdots \\ y_p \end{pmatrix}$

have zero sample covariance and max sample variance.

Number of PC

look at the scree plot (λ_i) and look for 'elbow'.



Math B7800

5/1/18, Tue

Final exam is similar with Exam 1. with 4 page sheets.
(5/22) calculator.

Sample PC

x_1, x_2, \dots, x_n are iid $p \times 1$ random vectors from a population with mean μ and covariance matrix Σ .

Obtain sample mean \bar{x} and sample covariance matrix S .

Obtain eigenvalues, eigenvectors $(\hat{\lambda}_i, \hat{e}_i)$ $i=1, \dots, p$ for the matrix S .
 $\hat{\lambda}_1 \geq \dots \geq \hat{\lambda}_p$

then the sample PC are given by

$$\hat{e}_i'(x - \bar{x}), i=1, \dots, p.$$

$$\frac{1}{n-1} \sum_{i=1}^n (x_i - \bar{x})(x_i - \bar{x})'$$

Interpretation

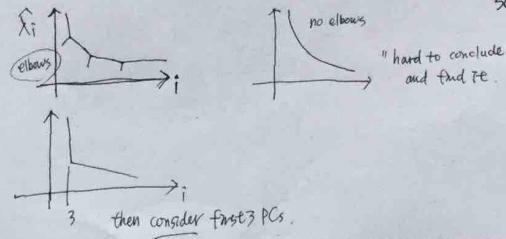
Total sample variance = trace(S)

$$= \sum_{i=1}^p \hat{\lambda}_i$$

$$= \sum_{i=1}^p \underbrace{\text{sample variance of } \hat{e}_i'(\bar{x} - \mu)}_{= \hat{\lambda}_i}$$

The PC corresponding to small $\hat{\lambda}_i$ are dropped.

Either use Scree plot or a thumb rule applied to standardized sample PC.



Standardized sample PC

Initially we have $\bar{x}_1, \dots, \bar{x}_n$.

Sample mean $\bar{\bar{x}}$, sample cov = S , $D = \text{Diag}(S)$.

$$\tilde{x}_i = D^{-1/2}(\bar{x}_i - \bar{\bar{x}}) \quad \text{standardized } \bar{x}_i$$

$\tilde{x}_1, \dots, \tilde{x}_n$ are standardized version of \bar{x} 's.

$$\bar{\bar{x}} = \bar{x} \quad \text{and} \quad \frac{1}{n-1} \sum_{i=1}^n \tilde{x}_i \tilde{x}_i' = \text{sample cov matrix of } \tilde{x}_1, \dots, \tilde{x}_n.$$

$$= R \quad (\text{sample correlation matrix} \\ p \times p \quad \text{for } \bar{x}_1, \dots, \bar{x}_n)$$

$\text{trace}(R) = P$, If $\hat{\lambda}_i, i=1, \dots, p$ are the eigenvalues of R .

$$\text{then } \sum_{i=1}^p \hat{\lambda}_i = P$$

The PC for which $\hat{\lambda} \leq 1$ are dropped.

Standardized sample PC are

$\hat{e}'_i z = \hat{e}'_i D^{-1/2}(\bar{x} - \bar{\bar{x}})$, where \hat{e}_i is an eigenvector
corresponding to $\hat{\lambda}_i$ and $\hat{e}_1, \dots, \hat{e}_p$ are orthogonal.

(3) The distribution of $\hat{\lambda}_i$ and $\hat{e}_i, i=1, \dots, p$ are asymptotically independent.

X Suppose the spectral decomposition of $S_{3 \times 3}$

$$\begin{aligned} \hat{\lambda}_1 &= 10 & \hat{e}_1 &= \begin{bmatrix} 1/\sqrt{3} \\ 1/\sqrt{3} \\ 1/\sqrt{3} \end{bmatrix}, & \hat{e}_2 &= \begin{bmatrix} 1/\sqrt{2} \\ -1/\sqrt{2} \\ 0 \end{bmatrix}, & \hat{e}_3 &= \begin{bmatrix} 1/\sqrt{6} \\ 1/\sqrt{6} \\ -2/\sqrt{6} \end{bmatrix} \\ \hat{\lambda}_2 &= 5 & & & & & \\ \hat{\lambda}_3 &= 1 & & & & & \end{aligned}$$

Q: Find 95% CI for $\lambda_1, \lambda_2, \lambda_3$. (n is large), $n=80$.

$$\text{For } \lambda_1, \quad \frac{-1.96}{\sqrt{2 \cdot 10^2}} \leq \frac{\sqrt{80}(10-\lambda_1)}{\sqrt{2 \cdot 10^2}} \leq \frac{1.96}{\sqrt{2 \cdot 10^2}} \quad \alpha = 0.05$$

$$? \leq \lambda_1 \leq ?$$

Sample PC estimate population PC.

Recall $y_i = e_i^T (X - \mu)$, $i=1, \dots, p$ where the population PC.

Recall $\text{cov}(X) = \Sigma$, Σ has eigenvalues $\lambda_1 \geq \dots \geq \lambda_p$, with corresponding eigenvectors e_1, \dots, e_p .

$$\text{cov}(Y) = \begin{bmatrix} \lambda_1 & 0 \\ 0 & \lambda_p \end{bmatrix} = \Lambda, \quad \text{cov}(y_i, y_j) = e_i^T \Sigma e_j = \begin{cases} 0 & \text{if } i \neq j \\ \lambda_i & \text{if } i = j \end{cases}$$

large sample inference

$$(1) \sqrt{n} \left[\begin{pmatrix} \hat{\lambda}_1 \\ \vdots \\ \hat{\lambda}_p \end{pmatrix} - \begin{pmatrix} \lambda_1 \\ \vdots \\ \lambda_p \end{pmatrix} \right] \approx N_p(0, 2\Lambda^2)$$

$$(2) \sqrt{n} (\hat{e}_k - e_k) \approx N_p(0, E_k), \text{ where } \underbrace{\hat{e}_k}_{\text{rank 1}} \leftarrow \text{rank is 1.} \quad \underbrace{xx'}_{\text{rank is 1.}}$$

$$E_k = \sum_{\substack{i=1 \\ i \neq k}}^p \frac{\lambda_i \lambda_k}{(\lambda_i - \lambda_k)^2} e_i e_i^T \leftarrow \text{rank } p-1 \Rightarrow \text{singular matrix.}$$

Note that rank(E_k) is $p-1$.

This is fine because \hat{e}_k satisfies $\sum_{i=1}^p (\hat{e}_{ki})^2 = 1$, and so the vector \hat{e}_k lies on $p-1$ dimensional sphere.

Q: Find 99% CI for e_{21} .

$$-Z_{\alpha/2} \leq \frac{\sqrt{80} (1/\sqrt{2} - e_{21})}{\sqrt{(\hat{E}_2)_{1,1}}} \leq Z_{\alpha/2}$$

$$\hat{E}_2 = \sum_{\substack{i=1 \\ i \neq 2}}^3 \frac{\lambda_i \lambda_k}{(\lambda_i - \lambda_k)^2} e_i e_i^T$$

$$\xrightarrow{\alpha=0.01} = \frac{\lambda_1 \lambda_2}{(\lambda_1 - \lambda_2)^2} e_1 e_1^T + \frac{\lambda_2 \lambda_3}{(\lambda_2 - \lambda_3)^2} e_2 e_2^T$$

$$(\hat{E}_2)_{1,1} = \frac{10 \cdot 5}{5^2} \frac{1}{3} + \frac{5}{16} \frac{1}{6}$$

Q: Find 95% CR for (e_{21}, e_{22}) .

$$\sqrt{80} \begin{pmatrix} \hat{e}_{21} - e_{21} \\ \hat{e}_{22} - e_{22} \end{pmatrix} \sim N_2(0, A)$$

$$A = (\hat{E}_2)_{1,1} \quad \hat{E}_2 = \frac{10 \times 5}{5^2} \hat{e}_1 \hat{e}_1^T + \frac{5 \times 1}{16} \hat{e}_2 \hat{e}_2^T$$

$$\sqrt{80} A^{-1/2} \begin{pmatrix} \hat{e}_{21} - e_{21} \\ \hat{e}_{22} - e_{22} \end{pmatrix} \sim N_2(0, I_2)$$

$$80 \begin{pmatrix} \hat{e}_{21} - e_{21} \\ \hat{e}_{22} - e_{22} \end{pmatrix}^T A^{-1} \begin{pmatrix} \hat{e}_{21} - e_{21} \\ \hat{e}_{22} - e_{22} \end{pmatrix} \approx \chi^2_2$$

$$\text{CR} := \left\{ \begin{pmatrix} e_{21} \\ e_{22} \end{pmatrix} : T \leq \chi^2_{2, 0.05} \right\}$$

Test x_1, \dots, x_n iid with mean μ and cov Σ and correlation matrix P .

H_0 : All pairwise correlations are same $P = P_0 = \begin{bmatrix} 1 & p & \dots & p \\ p & 1 & \dots & p \\ \vdots & \vdots & \ddots & \vdots \\ p & p & \dots & 1 \end{bmatrix}$

H_1 : H_0 is false.

Likelihood Ratiotest at level α :

First obtain sample correlation matrix R , $R_{ij} = \frac{s_{ij}}{\sqrt{s_{ii}s_{jj}}}$

\bar{r}_k = average of the off-diagonal entries in the k^{th} row of R
 $k = 1, \dots, p$.

$$\bar{r}_k = \frac{1}{p-1} \sum_{i=1, i \neq k}^p R_{ik}$$

$$\bar{r} = \text{average of all off diagonal} = \frac{1}{p(p-1)} \sum \sum_{i \neq k} R_{ik}$$

$$\Gamma = \frac{(p-1)(1-(1-\bar{r})^2)}{p - (p-2)(1-\bar{r})^2},$$

$$T = \frac{p-1}{(1-\bar{r})^2} \left[\sum \sum_{i \neq k} (\bar{r}_{ik} - \bar{r})^2 - \bar{r} \sum_k (\bar{r}_k - \bar{r})^2 \right]$$

$$\stackrel{H_0}{\sim} \chi^2_{(p+1)(p-2)/2}$$

Math B7800

5/3/18 / Thur

Ch.11 Classification

In this set up, we suppose that the data (observation vector, feature vector) comes from one of the classes (or groups).

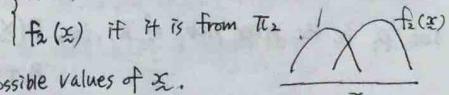
The class labels of the data vectors are unknown.

The set of possible class is known.

We would like to predict the class labels.

Two classes Suppose T_1 and T_2 denote the two possible classes that the data vectors \underline{x} can come from.

The pdf of \underline{x} is $f_1(\underline{x})$ if it is from T_1 , $f_1(\underline{x})$

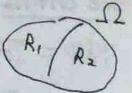


Ω = set of all possible values of \underline{x} .

We want to divide Ω into two disjoint regions

R_1 and R_2

so that $\underline{x} \in R_i$ will be classified to T_i .



Misclassification Probability

$$\begin{aligned} P(2|1) &= P(\text{observation from } \pi_1 \text{ is classified to } \pi_2) \\ &= P(\text{observation classified to } \pi_2 \mid \text{it came from } \pi_1) \\ &= \int_{R_2} f_1(x) dx \end{aligned}$$

Similarly,

$$P(1|2) = \int_{R_1} f_2(x) dx$$

In many cases, we have prior knowledge about the probability that x comes from π_1 or π_2 .

Let P_i be the prior prob. of π_i , i.e., $P(x \in \pi_i) = P_i$

$$P_1 + P_2 = 1$$

One criteria to minimize

TMP = Total misclassification prob.

$$= P(x \text{ is misclassified})$$

$$= P(2|1)P_1 + P(1|2)P_2$$

$$= P_1 \int_{R_2} f_1(x) dx + P_2 \int_{R_1} f_2(x) dx$$

~~Cost~~

Cost involved in Misclassification

		π_1	π_2	"Actual"
		○	C(1 2)	Average of Expected cost of Misclassification
<u>Your classification</u>	π_1	○	C(2 1)	$\underline{ECM} = \mathbb{E}(\text{cost})$
	π_2	○		
				$= P_1 P(2 1) C(2 1) + P_2 P(1 2) C(1 2)$
				$= P_1 \int_{R_2} f_1(x) dx C(2 1) + P_2 \int_{R_1} f_2(x) dx C(1 2)$

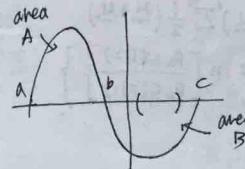
Assume that we know $P_1, P_2, C(1|2), C(2|1), f_1(x), f_2(x)$.

Goal: find R_1, R_2 so that ECM is minimized.

$$\int_{R_1} f_1(x) dx + \int_{R_2} f_2(x) dx = 1$$

$$\text{so, } \underline{ECM} = P_1 \left(1 - \int_{R_1} f_1(x) dx \right) C(2|1) + P_2 \int_{R_1} f_2(x) dx C(1|2)$$

$$= \int_{R_1} \left[P_2 C(1|2) f_2(x) - P_1 C(2|1) f_1(x) \right] dx + P_1 C(2|1)$$



$\int_R h(x) dx$ is minimum or max when $R = [b, c]$ when $R = [a, b]$.

$$-B \leq \int_R f(x) dx \leq A$$

So, using similar argument, ECM will be minimized

$$\text{If } R_1 = \left\{ \tilde{x} : \frac{f_1(\tilde{x})}{f_2(\tilde{x})} > \frac{C(1|2)p_2}{C(2|1)p_1} \right\}$$

↓ ratio ↓ ratio ↓ ratio
 ↑ pdf of cost of prior.

Suppose π_1 is $N_p(\mu_1, \Sigma)$ where $\mu_1 \neq \mu_2$
and π_2 is $N_p(\mu_2, \Sigma)$ but same covariance.

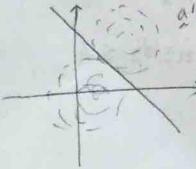
Assume μ_1, μ_2, Σ are known cov.

$$\text{Then } f_1(\tilde{x}) = \frac{1}{(2\pi)^{1/2} |\Sigma|^{1/2}} \exp\left[-\frac{1}{2} (\tilde{x} - \mu_1)' \Sigma^{-1} (\tilde{x} - \mu_1)\right]$$

$$f_2(\tilde{x}) = \frac{1}{(2\pi)^{1/2} |\Sigma|^{1/2}} \exp\left[-\frac{1}{2} (\tilde{x} - \mu_2)' \Sigma^{-1} (\tilde{x} - \mu_2)\right]$$

$$\frac{f_1(\tilde{x})}{f_2(\tilde{x})} = \exp\left[-\frac{1}{2} \{ (\tilde{x} - \mu_1)' \Sigma^{-1} (\tilde{x} - \mu_1) - (\tilde{x} - \mu_2)' \Sigma^{-1} (\tilde{x} - \mu_2) \}\right]$$

$$= \exp\left[-(\mu_1 - \mu_2)' \Sigma^{-1} \left(\tilde{x} - \frac{\mu_1 + \mu_2}{2} \right) \right]$$

$$\text{so, } R_1 = \left\{ \tilde{x} : (\underbrace{\mu_1 - \mu_2)' \Sigma^{-1} \tilde{x}}_{a' \tilde{x} - b \geq 0} - (\mu_1 - \mu_2)' \Sigma^{-1} \frac{1}{2} (\mu_1 + \mu_2) \geq \ln \left[\frac{p_2 C(1|2)}{p_1 C(2|1)} \right] \right\}$$


If μ_1, μ_2, Σ are not known, but we have a "training sample"

(n_1 observations $\tilde{x}_{1i}, i=1, \dots, n_1$ from π_1 and
 n_2 observations $\tilde{x}_{2j}, j=1, \dots, n_2$ from π_2)

Then we estimate μ_1 by \bar{x}_1
 μ_2 by \bar{x}_2 .

Σ by Spotted, where

$$\text{Spotted} = \frac{n_1 - 1}{n_1 - 1 + n_2 - 1} S_1 + \frac{n_2 - 1}{n_1 - 1 + n_2 - 1} S_2 \quad \text{and}$$

$$S_j = \frac{1}{n_j - 1} \sum_{i=1}^{n_j} (\tilde{x}_{ji} - \bar{x}_j)(\tilde{x}_{ji} - \bar{x}_j)'$$

We replace μ_1, μ_2, Σ by $\bar{x}_1, \bar{x}_2, \text{Spotted}$ in R_1 .

Note: $n_1 - 1 + n_2 - 1 \geq p$.

If suppose π_1, π_2 are $N_p(\mu_1, \Sigma_1), N_p(\mu_2, \Sigma_2)$, then

$$\frac{f_1(x)}{f_2(x)} = \frac{|\Sigma_2|^{1/2}}{|\Sigma_1|^{1/2}} \exp\left[-\frac{1}{2} (\tilde{x} - \mu_1)' \Sigma_1^{-1} (\tilde{x} - \mu_1) + \frac{1}{2} (\tilde{x} - \mu_2)' \Sigma_2^{-1} (\tilde{x} - \mu_2)\right] = \frac{f_1(x)}{f_2(x)}$$

so, $\boxed{R_1}$

$$R_1 = \left\{ \tilde{x} : \frac{1}{2} \tilde{x}' (\Sigma_2^{-1} - \Sigma_1^{-1}) \tilde{x} + (\mu_1' \Sigma_1^{-1} - \mu_2' \Sigma_2^{-1}) \tilde{x} - \frac{1}{2} (\mu_1' \Sigma_1^{-1} \mu_1 - \mu_2' \Sigma_2^{-1} \mu_2) + \frac{1}{2} \ln \frac{|\Sigma_2|}{|\Sigma_1|} \geq \ln \left[\frac{p_2 C(1|2)}{p_1 C(2|1)} \right] \right\}$$

This is a quadratic separator.

When $\mu_1, \mu_2, \Sigma_1, \Sigma_2$ are not known, but we have ~~a training sample~~
a training sample, then estimate μ_i by \bar{x}_i and
 Σ_i by S_i .

Math BN800

5/8/18, Tue

Classification

Recall that in this set up the observations are coming from one of g possible groups: $\pi_1, \pi_2, \dots, \pi_g$.

case 1: The pdf for different groups are known.

case 2: The pdfs' are unknown.

case 2.1: The pdfs' are parametric with unknown parameters.

case 2.2: The pdfs' are nonparametric.

In case 1, we minimize either TMP or ECM

<u>(Total Misclassification)</u>	<u>Expected Cost of misclassification</u>
probability	(Expected Cost of misclassification)

For $g=2$

$$\text{TMP} = P_1 \int_{R_2} f_1(x) dx + P_2 \int_{R_1} f_2(x) dx$$

P_i = prior probability of π_i .

$$\text{ECM} = P_1 C(2|1) \int_{R_2} f_1(x) dx + P_2 C(1|2) \int_{R_1} f_2(x) dx$$

		predict	
		π_1	π_2
Actual	π_1	0	$C(2 1)$
	π_2	$C(1 2)$	0

The choice of R_1 and R_2 that minimizes ECM

$$R_1 = \{ \underline{x} : \frac{f_1(\underline{x})}{f_2(\underline{x})} > \frac{C(1|2) P_2}{C(2|1) P_1} \}$$

Note: if $C(2|1) = C(1|2)$, then

minimizing ECM and TMP are equivalent.

$$R_2 = \{ \underline{x} : \frac{f_1(\underline{x})}{f_2(\underline{x})} < \frac{C(1|2) P_2}{C(2|1) P_1} \}$$

optimum error rate (OER) = smallest possible error of a classification rule.

~~Ex~~ Suppose π_1 and π_2 denote $N_p(\mu_1, \Sigma), N_p(\mu_2, \Sigma)$ with known parameters. Suppose $C(2|1) = C(1|2), P_1 = P_2 = \frac{1}{2}$.

Find OER?

The classification rule that minimizes ECM or TMP is given by

$$R_1 = \{ \underline{x} : \frac{f_1(\underline{x})}{f_2(\underline{x})} \geq 1 \}, R_2 = \{ \underline{x} : \frac{f_1(\underline{x})}{f_2(\underline{x})} < 1 \}$$

$$\frac{f_1(\underline{x})}{f_2(\underline{x})} = \exp \left[-\frac{1}{2} (\underline{x} - \mu_1)' \Sigma^{-1} (\underline{x} - \mu_1) + (\underline{x} - \mu_2)' \Sigma^{-1} (\underline{x} - \mu_2) \right]$$

$$\frac{f_1(\underline{x})}{f_2(\underline{x})} \geq 1 \Leftrightarrow \log \frac{f_1(\underline{x})}{f_2(\underline{x})} \geq 0 \Leftrightarrow (\mu_1 - \mu_2)' \Sigma^{-1} \underline{x} - (\mu_1 - \mu_2)' \Sigma^{-1} \frac{1}{2} (\mu_1 + \mu_2) \geq 0$$

$$so, R_1 = \{ \underline{x} : (\mu_1 - \mu_2)' \Sigma^{-1} \underline{x} \geq (\mu_1 - \mu_2)' \Sigma^{-1} \frac{1}{2} (\mu_1 + \mu_2) \}$$

$$R_2 = \{ \underline{x} : (\mu_1 - \mu_2)' \Sigma^{-1} \underline{x} < (\mu_1 - \mu_2)' \Sigma^{-1} \frac{1}{2} (\mu_1 + \mu_2) \}$$

$$TMP = P_1 P(2|1) + P_2 P(1|2)$$

$$P(2|1) = P(\underline{x} \in R_2 | \underline{x} \sim N(\mu_1, \Sigma))$$

$$P(\underline{x} \in R_2), \text{ where } \underline{x} \sim N_p(\mu_1, \Sigma)$$

$$so, (\mu_1 - \mu_2)' \Sigma^{-1} \underline{x} \sim N_1((\mu_1 - \mu_2)' \Sigma^{-1} \mu_1, (\mu_1 - \mu_2)' \Sigma^{-1} (\mu_1 - \mu_2))$$

$$P(\underline{x} \in R_2) = P((\mu_1 - \mu_2)' \Sigma^{-1} \underline{x} < (\mu_1 - \mu_2)' \Sigma^{-1} \frac{1}{2} (\mu_1 + \mu_2))$$

$$= P \left[\frac{(\mu_1 - \mu_2)' \Sigma^{-1} (\underline{x} - \mu_1)}{\sqrt{(\mu_1 - \mu_2)' \Sigma^{-1} (\mu_1 - \mu_2)}} < \frac{(\mu_1 - \mu_2)' \Sigma^{-1} \left[\frac{1}{2} (\mu_1 + \mu_2) - \mu_1 \right]}{\sqrt{(\mu_1 - \mu_2)' \Sigma^{-1} (\mu_1 - \mu_2)}} \right]$$

$$= P \left[\underline{z} < -\frac{1}{2} \sqrt{(\mu_1 - \mu_2)' \Sigma^{-1} (\mu_1 - \mu_2)} \right]$$

$$= \Phi(-\frac{1}{2} \Delta)$$

$$P(1|2) = P(\underline{x} \in R_1) \text{ if } \underline{x} \sim N_p(\mu_2, \Sigma).$$

$$= P \left[(\mu_1 - \mu_2)' \Sigma^{-1} \underline{x} \geq (\mu_1 - \mu_2)' \Sigma^{-1} \frac{1}{2} (\mu_1 + \mu_2) \right], \text{ where}$$

$$(\mu_1 - \mu_2)' \Sigma^{-1} \underline{x} \sim N((\mu_1 - \mu_2)' \Sigma^{-1} \mu_2, (\mu_1 - \mu_2)' \Sigma^{-1})$$

$$= P \left[\frac{(\mu_1 - \mu_2)' \Sigma^{-1} (\underline{x} - \mu_2)}{\sqrt{(\mu_1 - \mu_2)' \Sigma^{-1} (\mu_1 - \mu_2)}} \geq \frac{(\mu_1 - \mu_2)' \Sigma^{-1} \left[\frac{1}{2} (\mu_1 + \mu_2) - \mu_2 \right]}{\sqrt{(\mu_1 - \mu_2)' \Sigma^{-1} (\mu_1 - \mu_2)}} \right] (\mu_1 - \mu_2)$$

$$= P[\underline{z} \geq \frac{1}{2} \Delta] = 1 - \Phi(\frac{1}{2} \Delta) = \Phi(-\frac{1}{2} \Delta). \leftarrow "OER"$$

$$= \mathbb{P}[Z \geq \frac{1}{2}\Delta] = 1 - \mathbb{P}(\frac{1}{2}\Delta) = \mathbb{P}(-\frac{1}{2}\Delta)$$

$$\text{So, OER} = \mathbb{P}(-\frac{\Delta}{2})$$

$\Delta = \sqrt{(\mu_1 - \mu_2)' \Sigma^{-1} (\mu_1 - \mu_2)}$ is "distance" between μ_1, μ_2 .

This is called Mahalanobis (PC Mahalanobis)
distance between two means.

Case 2.1 pdfs are parametric with unknown parameters.

We have a training sample for which the group label is known.

Then we estimate the unknown parameters using the
training sample.

Then use the classification rule used in case 1 with
the parameters replaced with their estimates.

Here, we use \hat{R}_1, \hat{R}_2 instead of R_1 and R_2 .

The corresponding rates are called Actual Error Rate (AER).

$$\underline{\text{AER}} = P_1 \int_{\hat{R}_2} f_1(x) dx + P_2 \int_{\hat{R}_1} f_2(x) dx$$

AER is not known, as it involves the unknown pdfs
 $f_1(x), f_2(x)$.

We estimate AER from the training data.

AER (APER, apparent error rate)

$$= \frac{N_{1m} + N_{2m}}{N_1 + N_2} = \text{proportion of misclassification}$$

within training sample.

		predicted	
		π_1	π_2
Actual	π_1	N_{1c}	N_{1m}
	π_2	N_{2c}	N_{2m}

Logistic Regression

Starting with training data for which group labels are known.
We can treat the group labels as the response variable.

For the i th training sample,

$$y_i = \begin{cases} 1 & \text{if it is from } \pi_1 \\ 2 & \text{if it is from } \pi_2 \end{cases}$$

The data: x_1, \dots, x_n ($P \times 1$)

logistic model

$$\text{logit}(P(Y=1)) = \beta_0 + \beta_1 x_1 + \dots + \beta_p x_p$$

$$= \log \left(\frac{P(Y=1)}{P(Y=0)} \right)$$

First estimate: $\beta_0, \beta_1, \dots, \beta_p$.

Then allocate x to π_1 if

$$\beta_0 + \hat{\beta}_1 x_1 + \dots + \hat{\beta}_p x_p > 0.$$

Then allocate x to T_{l_1} if

$$\hat{\beta}_0 + \hat{\beta}_1 x_1 + \dots + \hat{\beta}_p x_p > 0.$$

When there are $g \geq 2$ populations T_{l_1}, \dots, T_{l_g} with

$$\text{pdf } f_1(x), \dots, f_g(x) \rightarrow$$

$C(k|i)$ = cost of misclassification of an observation coming from T_{l_i} to T_{l_k} .

p_1, p_2, \dots, p_g are prior probabilities.

$$ECM(1) = C(2|1)f_2(x) + C(3|1)f_3(x) + \dots + C(g|1)f_g(x)$$

$$ECM(i) = \sum_{k=1}^g C(k|i)f_k(x)$$

$$ECM = p_1 ECM(1) + p_2 ECM(2) + \dots + p_g ECM(g).$$

Allocate x to T_{l_k} if

$$\sum_{\substack{i=1 \\ i \neq k}}^g p_i C(k|i) f_i(x) \text{ is the smallest.}$$

Math B7800

5/10/18, Thur

No date and time for variables on the project.

Last time

Classification with g groups/classes.

↓
supervised learning

$T_{l_1}, T_{l_2}, \dots, T_{l_g}$ are g groups with prior probabilities p_1, p_2, \dots, p_g .

The pdfs are $f_1(x), \dots, f_g(x)$

case 1 : pdfs are known.

$$ECM = p_1 ECM(1) + \dots + p_g ECM(g).$$

cost of misclassification

$C(k|i)$ = cost for misclassifying an individual from T_{l_i} to T_{l_k} .

$$ECM(1) = C(2|1) \cdot p(2|1) + C(3|1) \cdot p(3|1) + \dots + C(g|1) \cdot p(g|1)$$

$$ECM(i) = \sum_{k=1}^g C(k|i) p(k|i), i = 1, 2, \dots, g.$$

$p(k|i)$ = prob. of
classifying
an observation
from T_{l_i} to T_{l_k}

$$ECM = \sum_{i=1}^g p_i \sum_{k=1}^g C(k|i) p(k|i)$$

$$= p_1 ECM(1) + \dots + p_g ECM(g)$$

The classification rule that minimizes ECM will allocate x to T_{l_k} if $\sum_{i=1}^g f_i(x) p_i C(k|i)$ is minimum.

\leq

When $g=2$

we allocate x to π_i if $f_i(x)P_iC(2|1)$
 $\geq f_2(x)P_2C(1|2)$ minimum.
 $\Leftrightarrow \frac{f_1(x)}{f_2(x)} \geq \frac{C(1|2)P_2}{C(2|1)P_1}$

When the costs $C(k|i)$ are all equal, then it is enough
to compare $\beta_k = \sum_{i=1}^g P_i f_i(x)$.

Choose k such that $\beta_k = \min_i \beta_i$

$$\begin{aligned}\beta_k &= \sum_{i=1}^g P_i f_i(x) - P_k f_k(x) \\ &= C - P_k f_k(x).\end{aligned}$$

so, minimizing $\beta_k \Leftrightarrow$ maximizing $P_k f_k(x)$

Posterior prob. of a group

posterior prob. for group i

$$= P(\pi_i | x) = P(X \text{ comes from } \pi_i | x=x)$$

$$= \frac{P_i f_i(x)}{\sum_{k=1}^g P_k f_k(x)}$$

so, the allocation rule that maximizes the posterior prob.

will allocate x to π_i if $P_i f_i(x) = \max_k P_k f_k(x)$

When the costs are equal, then minimizing ECM

\Leftrightarrow maximizing posterior.

Ex If π_i is $N_p(\mu_i, \Sigma_i)$, $i=1, \dots, g$.

Suppose equal costs, then Allocate \bar{x} to π_i If $P_i f_i(\bar{x})$ is the max.

Equivalently, $\log(P_i f_i(\bar{x}))$ is the max.

$$\Leftrightarrow \log(P_i) - \frac{1}{2} \log(|\Sigma_i|) - \frac{1}{2} (\bar{x} - \mu_i)' \Sigma_i^{-1} (\bar{x} - \mu_i) \text{ "quadratic"}$$

Ex If π_i is $N_p(\mu_i, \Sigma)$, $i=1, \dots, g$,
then allocate \bar{x} to π_i if

$$\log(P_i) + \mu_i' \bar{x} - \frac{1}{2} \mu_i' \Sigma^{-1} \mu_i \text{ "linear"}$$

is the max.

Chapter 12 Clustering (unsupervised learning)

No training data. → No fitted model. Example M12

Setup: we have a bunch of "objects", which we need to classify into different groups / clusters

Ex1 Community detection.

Ex2 Segregating genes.

Ex3 Segregating patients / healthy

First, assume that our observations are points in \mathbb{R}^p .

$x_1, x_2, \dots, x_n \in \mathbb{R}^p$ e.g. $p=2$

Two kinds of algorithms

1. Hierarchical

2. Non-Hierarchical

1. Let $D_{n \times n}$ be the distance matrix. $D_{ij} = \|x_i - x_j\|$

Step 0 we have singleton clusters $\{x_1\}, \{x_2\}, \dots, \{x_n\}$

D = Distance matrix.

Step 1 find clusters U and V s.t. $\text{dist}(U, V)$ is the minimum among all distances of pairs of clusters.

Step 2 Merge clusters U and V , so delete U, V and add $U \cup V$.

[e.g. after step 0, if $\|x_1 - x_2\| = \min_{i \neq j} \|x_i - x_j\|$

then new cluster structure is $\{x_1, x_2\}, \{x_3\}, \dots, \{x_n\}$

Step 3 Adjust the distance matrix.

Delete the rows and columns corresponding to clusters U and V .

Add a new row and column for $U \cup V$ in $D_{n-1 \times n-1}$

Step 4 Repeat step 2 and 3 until you get one cluster

$\{x_1, \dots, x_n\}$

Distance between clusters

(a) single link $\text{dist}(U, x) = \min \{d(x, y) : y \in U\}$

(b) average link $\text{dist}(U, x) = \text{dist}(x, \text{centroid of } U)$

2. Non-

• Nonheirarchical Clustering

(k means clustering algorithm)

Assume that there are K clusters. y_1, \dots, y_K

Goal: Find out the location of K centers and allocation of the points x_1, \dots, x_n to

$A(x_1), \dots, A(x_n) \in \{y_1, \dots, y_K\}$. so that

$\sum_{i=1}^n \|x_i - A(x_i)\|$ is minimum among all possible centers
and all possible allocation rules. "Np hard"
so, "approximation"

Step 1 Start with any K groups.

Step 2 Obtain the K centroids these groups.
(means)

Step 3 Allocate each vector to its closest centroid.
this defines new K groups.

Repeat Step 2 and 3 until the sum of the distances
between the vectors and centroids stabilize.