

①

Error Analysis.

→ After filling the model if we want to cross check that the model fitting is justifiable or not - then error analysis will help us.

→ Outlier analysis.

Model

$$\mathbf{Y} \sim N(\mathbf{X}\boldsymbol{\beta}, \sigma^2 \mathbf{I}_n)$$

Estimated error.

$$\hat{\mathbf{e}} = \mathbf{y} - \hat{\mathbf{y}} = (\mathbf{I} - \mathbf{P}_x) \mathbf{y} \sim N(0, \sigma^2 (\mathbf{I} - \mathbf{P}_x))$$

$$\hat{\sigma}^2 = \frac{\hat{\mathbf{e}}^T \hat{\mathbf{e}}}{n-k-1} \text{ and } \frac{\hat{\mathbf{e}}^T \hat{\mathbf{e}}}{\sigma^2} \sim \chi_{n-k-1}^2. \quad \mathbf{P}_x = \mathbf{H} = \begin{pmatrix} \mathbf{h}_{ij} \\ (k+1) \times (k+1) \end{pmatrix}$$

$$\text{cov}(e_i, e_j) = \begin{cases} \sigma^2 (1 - h_{ii}) & \text{if } i=j \\ \sigma^2 (-h_{ij}) & \text{if } i \neq j \end{cases}$$

$$\frac{e_i}{\sqrt{\sigma^2 (1 - h_{ii})}} \sim N(0, 1) \quad \text{Marginally} \Rightarrow t_i = \frac{e_i}{\sqrt{\hat{\sigma}^2 (1 - h_{ii})}} \sim t_{n-k-1}.$$

~~is large~~

$$\frac{e_i}{\sqrt{\hat{\sigma}^2(1-h_{ii})}} \stackrel{a}{\sim} N(0,1) \quad \begin{matrix} n \gg k \\ \text{because} \end{matrix} \quad \left. \right\} \xrightarrow{d} \text{converges in distribution.}$$

$$T = \frac{Z}{\sqrt{Y/n}} \sim t_n \quad \left| \begin{matrix} Z \sim N(0,1) \\ Y \sim \chi^2_k \end{matrix} \right. \text{ind.}$$

$$t_n \equiv C(0,1)$$

$\gamma_n \xrightarrow{p} c$ WLLN as $n \uparrow \infty$. w.p.1.

$$\gamma_n \sim \chi^2_n.$$

$$Y = \sum_{i=1}^n z_i^2 \quad z_i \stackrel{iid}{\sim} N(0,1).$$

$$\frac{\gamma}{n} = \frac{1}{n} \sum_{i=1}^n z_i^2$$

$$E(z_i) = 1.$$

$$\left\{ \begin{array}{l} X_n \xrightarrow{d} X \\ Y_n \xrightarrow{p} c \\ Y_n X_n \xrightarrow{d} cX. \end{array} \right.$$

- Slutsky's theorem.
???

②

Standardized residual.

$$d_i = \frac{e_i}{\sqrt{MSE_{\text{Error}}}}$$

$$MSE_{\text{Error}} = \frac{SSE_{\text{Error}}}{n-k-1}$$

$$e_i = y_i - \hat{y}_i$$

(3).

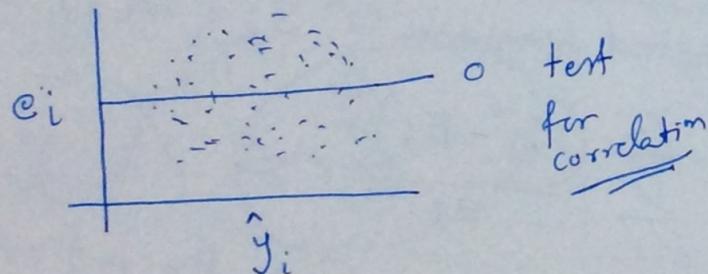
Studentized residual.

$$r_i = \frac{e_i}{\sqrt{MSE_{\text{Error}}(1-h_{ii})}}$$

For simple linear regression.

$$r_i = \frac{e_i}{\sqrt{MSE_{\text{Error}} \left[1 - \left(\frac{1}{n} + \frac{(x_i - \bar{x})^2}{\sum} \right) \right]}}$$

- ① residuals are marginally normally distributed ($k \ll n$)
- ② residuals are uncorrelated to the predicted values.



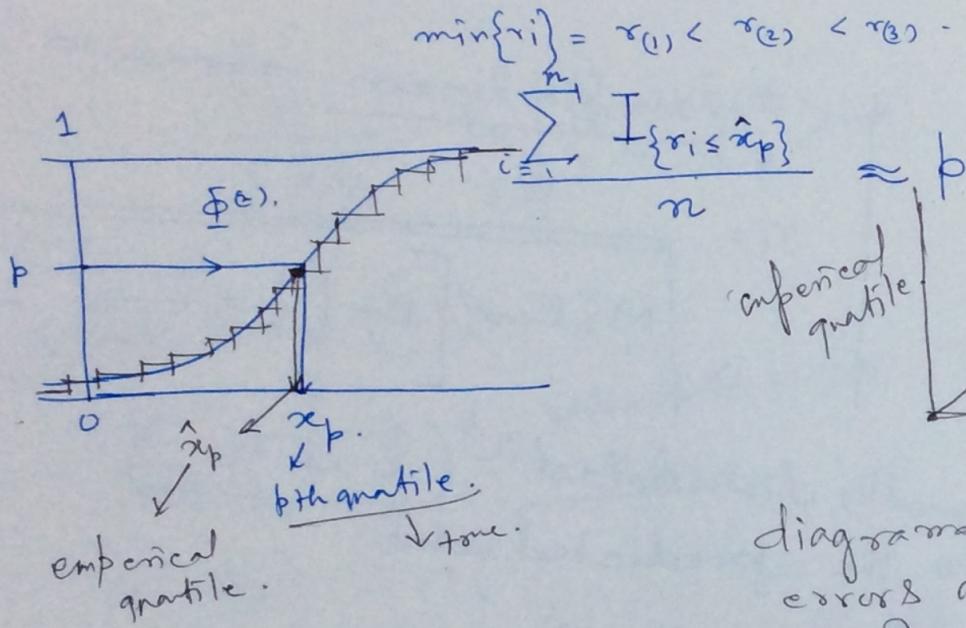
test
for
correlation

Correlation zero
should be reflected.

$r_i \stackrel{a}{\sim} N(0, 1)$
 → diagram based. ✓ (q-q-plot)
 → pdf based.
 → cdf based. } testing

$$\textcircled{1} \quad p\text{-th quantile } x_p \text{ if } p = \Phi(x_p) = P(X \leq x_p) \quad \underline{x \sim N(0,1)}$$

Let r_i 's be sorted in increasing order.



$$r_{(n)} = \max\{r_i\}$$

$$I_A = \begin{cases} 1 & \text{if } A \text{ is satisfied} \\ 0 & \text{otherwise.} \end{cases}$$

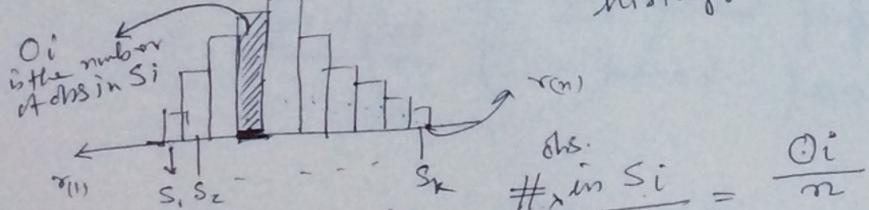
q-q plot:

true quantile.

diagrammatic approach to check whether errors are ~~for~~ following normal distribution or not. (qqplot).

\textcircled{2} bdf based:

construct the histogram.



$$\rightarrow P(\text{data} \in S_i) = \frac{\# \text{obs in } S_i}{n} = \frac{O_i}{n}$$

$$\rightarrow P(X \in S_i) = \text{calculate from true cdf} = f_i \text{ say.}$$

expected frequency is $S_i = n \cdot f_i = E_i$

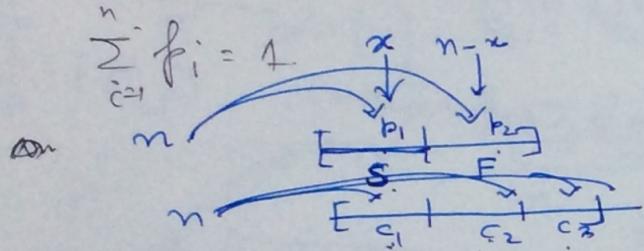
$$\sum_{i=1}^K \frac{(O_i - E_i)^2}{E_i} \sim \chi^2_{K-1}$$

When H_0 is true.

(5)

$$(s_1, s_2, \dots, s_k)$$

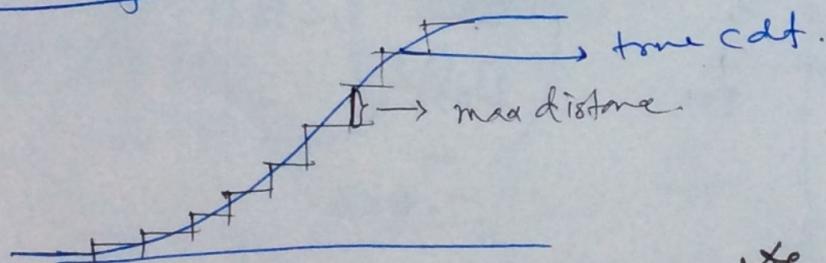
n obs.

multinomial (n, p_1, p_2, \dots, p_k)

When the p_k or parameters of the null distribution are estimated from data, then χ^2 statistic follows.

$$\sum_{i=1}^k \frac{(o_i - \hat{E}_i)^2}{\hat{E}_i} = \chi^2_{k-1 - r}$$

② Cdf based test. [only for continuous r.v.s]



$$\begin{aligned}\hat{E}_i &= n \hat{p}_i \\ &= n \hat{p}_i (\hat{\mu}, \hat{\sigma}^2)\end{aligned}$$

$$\begin{aligned}F_n(x) &= \frac{\sum_{i=1}^n I_{\{x_i \leq x\}}}{n} \quad \text{empirical CDF.} \\ F(x) &\text{ is true CDF.}\end{aligned}$$

we study $\sqrt{n}(F_n - F)$

and $n \uparrow \infty$.

Multinomial \rightarrow Multivariate normal with k components, but dim. will be

$\chi^2_{(k-1)}$ hence χ^2_{k-1} df.

When expected frequencies are completely known as the true CDF is completely specified

We study the distribution of

$$\max_x \sqrt{n} |F_n(x) - F(x)| \xrightarrow{d} N(0, t(t)) \quad \text{as } n \rightarrow \infty$$

\xrightarrow{d}

Supremum:

$$U \sim U[0,1] \quad \frac{\sum_{i=1}^n I\{u_i < t\}}{n}$$

$$N(0, t(t)) \sim \sqrt{n} \left(\frac{\sum_{i=1}^n I\{u_i \leq t\}}{n} - at \right)$$

$$G_n(y) = \frac{\sum_{i=1}^n I\{y_i \leq y\}}{n}$$

If we know $X \sim F(\cdot)$ which is continuous.

$$F(X) \sim U[0,1]$$

transfer the data. $F(x_i) = y_i$

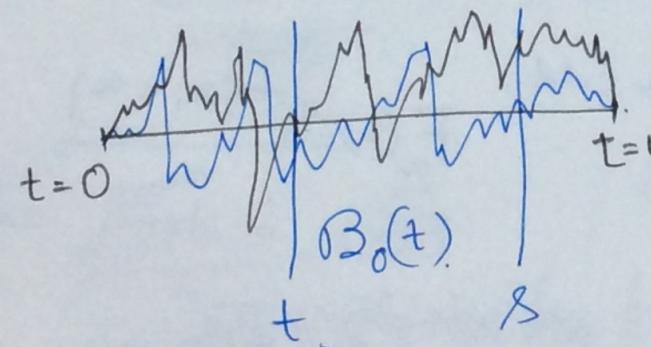
$$Y = F(X)$$

Kolmogorov-Smirnov test
KS-test.

$$\sup_{0 \leq y \leq 1} \sqrt{n} |G_n(y) - y| \xrightarrow{d} \sup_{0 \leq t \leq 1} |B_0(t)|$$

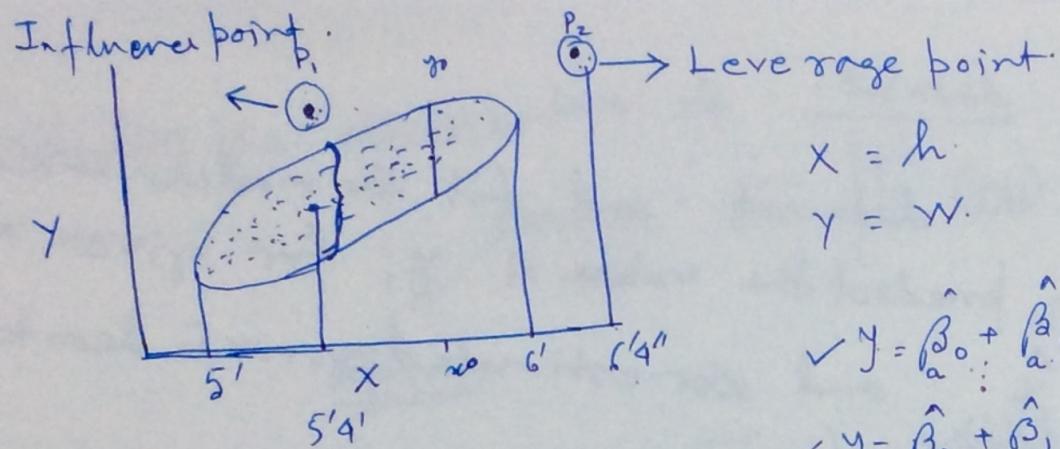
Kolmogorov distribution.
Brownian bridge on $[0,1]$.

$B_0(t)$ is known as the standard Brownian bridge on $[0,1]$.



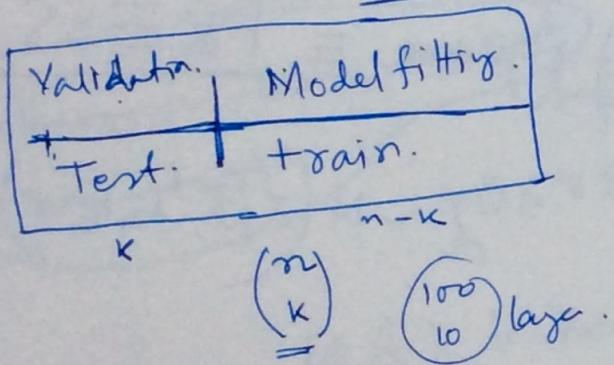
$$\left(\frac{B_0(t)}{B_0(s)} \right) \sim N \left[\begin{pmatrix} 0 \\ 0 \end{pmatrix} \begin{pmatrix} \frac{t(1-t)}{1-t} & \frac{t(1-s)}{1-t} \\ \frac{t(1-s)}{1-t} & \frac{s(1-t)}{1-t} \end{pmatrix} \right]$$

$t < s$



$$P_i = (y_i, x_i)$$

case 1 $|x_0 - x_i| < \delta \rightarrow$ prefer to include P_i }
case 2 $|x_0 - x_i| \gg \delta \rightarrow$ prefer not to include P_i } ②.



Leave-one-out / Jack-knife.
 x_1, x_2, \dots, x_n .
 ~~x_1, x_2, \dots, x_n~~ .
 This method is repeated for each data point.

- ① if the number of data point is less then we may not discard the outlier.
 ② If the data is logically incorrect then we discard it.

$$(m-2) \hat{\sigma}^2 = \frac{s_{yy} - \frac{s_{xx}}{s_{xx}} \hat{s}_{yy}^2}{\text{without filling the model}}$$

↓ using $(n-1)$ data point for model fitting.
 validation with one data point.

$$D = \{(x_i, y_i) \mid i=1, 2, \dots, n\}.$$

(8)

If we remove (x_i, y_i) from the data set and fit the model based on $(n-1)$ data points and then predict the value of y_i for given x_i then we denote it as $\hat{y}_{(i)}$ and estimated error is denoted as

$$e_{(i)} = y_i - \hat{y}_{(i)}$$

where $\hat{y}_{(i)}$ is the predicted value with model based on $(n-1)$ data points.
not w.r.t. (x_i, y_i)

$$\sum_{i=1}^n e_{(i)}^2 = \text{Predicted residual sum of square} = \underline{\text{PRESS}}.$$

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

$$Y \sim N(\underline{x}\underline{\beta}, \underline{I_n} \sigma^2)$$

~~$$e_i = (y_i - \hat{y}_i) \sim N(0, \sigma^2 (1-h_{ii})).$$~~

$$\text{It can be shown } e_{(i)} = \frac{e_i}{1-h_{ii}} \sim N\left(0, \frac{\sigma^2}{1-h_{ii}}\right).$$

Montgomery C.7 appendix.

$$\frac{e_{(i)}}{\sqrt{v(e_{(i)})}} = \frac{e_i}{\sqrt{\sigma^2(1-h_{ii})}} = \frac{e_i}{\sqrt{v(e_i)}} \quad \textcircled{2}.$$

$$\hat{\sigma}^2 = \sum_{i=1}^n \frac{e_i^2}{n-k-1}$$

σ^2 unknown.

$\hat{\sigma}^2 \stackrel{?}{=} \text{MSE}_{\text{error}}$. not preferable to be for jack-knife method.

$$\left\{ \begin{array}{l} \sigma^2 \triangleq S_{(i)}^2 = \frac{(n-k-1) \text{MSE}_{\text{error}} - \frac{\hat{e}_i^2}{1-h_{ii}}}{n-k-2} \\ \text{when } (x_i, y_i) \text{ not used.} \end{array} \right.$$

$$T = \frac{e_i}{\sqrt{S_{(i)}^2 (1-h_{ii})}} \sim t_{n-k-2}$$

C.8

if $n \gg k$.
then approximate it by Normal dist.

$$(A + \underline{u} \underline{v}^T)^{-1} = A^{-1} - \frac{A^{-1} \underline{u} \underline{v}^T A^{-1}}{1 + \underline{v}^T A^{-1} \underline{u}}$$

$$\boxed{\begin{aligned} & X_{(i)}^T X_{(i)} \\ & = (X^T X - \underline{x}_i \underline{x}_i^T) \end{aligned}}$$

To test.

H₀: ith obs is not outlier.

H₁: ith obs is an outlier.

Leverage point in a data set is observed because of \tilde{x} values.

(10)

it is identified with the quantity - $h_{ii} = \tilde{x}_i^T (X^T X)^{-1} \tilde{x}_i$.

$\text{tr}(P_x) = \text{tr}(H) = \sum_{i=1}^n h_{ii} > \frac{(k+1)}{\frac{m}{2}}$ we may suspect
that there can be at least one leverage point(s) in the dataset:

Large magnitude h_{ii} or/and $|e_i|$ can be considered as
an influential point which has a prominent impact in estimation

of $\hat{\beta}$ or $\hat{\epsilon}_i$.

$$VIF = \text{Variance inflation factor} = \frac{1}{1-h_{ii}}$$

$\hat{\beta}_{j(i)}$ is estimated value of $\hat{\beta}_j$ when (\tilde{x}_i, y_i) is not used
in estimation. i.e. $(m-1)$ observations are used

Impact of influence point on \hat{y}_i

Cook's distance: [1977, 79]

(11)
NOT an F statistic

$$D_i = \frac{(\hat{\beta}_{(i)} - \hat{\beta})^T X^T X (\hat{\beta}_{(i)} - \hat{\beta}) / (k+1)}{\sum_{i=1}^n (\hat{y}_i - \tilde{y})^2 / (n-k-1)} = \frac{(\hat{y}_{(i)} - \tilde{y})^T (\hat{y}_{(i)} - \tilde{y}) / (k+1)}{X^T (I - P_X) X / (n-k-1)}$$

\hat{y} is the predicted value of y for given x when all data are used to estimate $\hat{\beta}$.
 \tilde{y} is the predicted value of \tilde{y} for given x when i th obs is not used to estimate $\hat{\beta}$.

$$\hat{y}_{(i)} = X \hat{\beta}_{(i)}$$

$$\boxed{\hat{\beta}_{(i)} - \hat{\beta} = - (X^T X)^{-1} \tilde{x}_i e_i / (1 - h_{ii})}$$

Ch 6
Matrices.

$$D_i = \frac{e_i^2 h_{ii} / (k+1)}{(1 - h_{ii})^2 \text{MSE}_{\text{Error}}} = \left(\frac{e_i}{\sqrt{(1 - h_{ii}) \text{MSE}_{\text{Error}}}} \right)^2 \frac{h_{ii}}{(1 - h_{ii})(k+1)}$$

In prediction how much impact is present of the data. (x_i, y_i) is measured by Cook's distance. It can be expressed as a fraction elements which are available from complete data estimation.
For $n \gg k$ D_i is compared with 1. If $D_i > 1$ we may consider it has an impact in prediction.

DF BETA

Belsley, Kuh, Welsch (1980).

(12)

$$\hat{\beta}_j - \hat{\beta}_{j(i)} \sim ?$$

$$\hat{\beta} \sim N(\underline{\beta}, C\sigma^2)$$

$$\hat{\beta}_j - \hat{\beta}_{j(i)} = \frac{\gamma_{j,i} e_i}{1 - h_{ii}}$$

Denote:

$$R = (X^T X)^{-1} X^T$$

$$R^T R = (X^T X)^{-1} = (C_{ij})$$

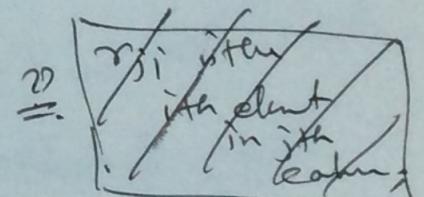
 r_j^T is the j th row of R .

everything is known from old.
estimation process with all
data points.

$$\frac{\hat{\beta}_j - \hat{\beta}_{j(i)}}{\sqrt{\gamma_{j,i}^T r_j s_{(i)}^2}} = \frac{r_{ji}}{\sqrt{\gamma_{j,i}^T r_j}} \cdot \left(\frac{e_i}{\sqrt{s_{(i)}^2 (1 - h_{ii})}} \right) \frac{!}{\sqrt{1 - h_{ii}}}$$

$$t_{(i)} = \frac{e_i}{\sqrt{s_{(i)}^2 (1 - h_{ii})}} \sim t_{n-k-2}$$

t test to check H₀: E $(\hat{\beta}_j - \hat{\beta}_{j(i)}) = 0$.
H₁: E $(\hat{\beta}_j - \hat{\beta}_{j(i)}) \neq 0$.



(13)

Cov Ratio.

$$\text{Cov Ratio} = \frac{|\hat{\Sigma}(\hat{\beta}_{(i)})|}{|\hat{\Sigma}(\hat{\beta})|} = \frac{|\hat{\sigma}_{(i)}^2 (x_{(i)}^T x_{(i)})^{-1}|}{|\hat{\sigma}^2 (x^T x)^{-1}|}$$

$$= \frac{|\sigma_{(i)}^2 (x_{(i)}^T x_{(i)})^{-1}|}{|\text{MSE}_{\text{error}} (x^T x)^{-1}|} = \left(\frac{\sigma_{(i)}^2}{\text{MSE}_{\text{error}}} \right)^{k+1} \frac{|(x_{(i)}^T x_{(i)})^{-1}|}{|T(x^T x)^{-1}|}$$

$$= \left(\frac{\sigma_{(i)}^2}{\text{MSE}_{\text{error}}} \right)^{k+1} \frac{1}{1 - h_{ii}}$$

$$\hat{\beta}_{(i)} \sim N(\beta, \sigma^2 (x_{(i)}^T x_{(i)})^{-1})$$

$$\begin{pmatrix} y_1 \\ y_2 \\ \vdots \\ y_n \end{pmatrix} = \begin{pmatrix} X \\ \times \end{pmatrix}_{n \times k} \begin{pmatrix} \beta_0 \\ \beta_1 \\ \vdots \\ \beta_k \end{pmatrix} + \tilde{\epsilon}.$$

✓ Cov Ratio. $\gtrsim 1 \pm 3\left(\frac{k+1}{n}\right)$. \rightarrow ith obs has impact. (14)

✓ $|DF\text{Beta}_{(i)}| > 2\sqrt{n}$. \Rightarrow ith obs has impact.

✓ DFFIT_i = $\frac{\hat{y}_i - \hat{y}_{(i)}}{\sqrt{s_{(i)}^2 h_{ii}}} = \left(\frac{h_{ii}}{1-h_{ii}}\right)^{1/2} t_{(i)}$.

✓ $|DFFIT_i| > \frac{2(k+1)}{n}$ then ith obs has impact.

$x_i^T (X^T X)^{-1} x_i$ is large enough.

or. $T_x(H) > \frac{2(k+1)}{n}$. then there may exists leverage points.

Variable selection & Model building.

Assume the true model $\tilde{Y} = X\tilde{\beta} + \epsilon$ $\tilde{\beta} \in \mathbb{R}^{k+1}$, $\tilde{Y} \in \mathbb{R}^n$, $\epsilon \in \mathbb{R}^n$

But when we fit the model we use $p < k+1$ variables instead $(k+1)$.

$$\tilde{Y} = [X_p : X_r] \begin{pmatrix} \tilde{\beta}_p \\ \tilde{\beta}_r \end{pmatrix} + \epsilon$$

$$\Rightarrow \tilde{Y} = X_p \tilde{\beta}_p + X_r \tilde{\beta}_r + \epsilon \text{ where } p+r = k+1.$$

$$\tilde{\beta}_p = (\beta_0 \ \beta_1 \ \dots \ \beta_{k-r})^T \quad \tilde{\beta}_r = (\beta_{k-r+1} \ \dots \ \beta_k)^T$$

Fitted model:

$$\tilde{Y} = X_p \tilde{\beta}_p + \epsilon.$$

$$\hat{\tilde{\beta}}_p = (X_p^T X_p)^{-1} X_p^T \tilde{Y} \rightarrow \text{Reduced model. } x$$

$$\hat{\beta} \begin{pmatrix} \tilde{\beta}_p \\ \tilde{\beta}_r \end{pmatrix} \hat{\beta} = \hat{\beta} = (X^T X)^{-1} X^T \tilde{Y} \rightarrow \text{complete model.}$$

Q1

whether $\hat{\beta}_{\tilde{p}}$ is an unbiased estimator of β_p under reduced model?

(15)

$$\hat{\beta}_{\tilde{p}} = \underbrace{(x_p^T x_p)^{-1} x_p^T y}_{\text{fitted.}} \text{ under reduced model.}$$

$$\begin{aligned} E(\hat{\beta}_{\tilde{p}}) &= (x_p^T x_p)^{-1} x_p^T (x \beta) \\ &= (x_p^T x_p)^{-1} x_p^T [x_p \beta_p + x_r \beta_r] \\ &= (x_p^T x_p)^{-1} (x_p^T x_p) \beta_p + (x_p^T x_p)^{-1} (x_p^T x_r) \beta_r \end{aligned}$$

True model:
 $y \sim N(x\beta, \sigma^2 I_n)$.

Assume $(x_p^T x_p)$ is invertable.

$$= \beta_p + (x_p^T x_p)^{-1} (x_p^T x_r) \beta_r$$

\Rightarrow In general $E(\hat{\beta}_{\tilde{p}}) \neq \beta_p$ but if $x_p^T x_p = 0$ ie. the columns of x_p are orthogonal to the columns of x_r then only

$$E(\hat{\beta}_{\tilde{p}}) = \beta_p$$

→ orthogonal polynomial.

→ Principal component regression.

Q2

Is $\hat{\sigma}_p^2$, under reduced fitted model, an unbiased estimator of σ^2 ? (17)

$$\hat{\sigma}_p^2 = \frac{\mathbf{y}^T (\mathbf{I}_n - \mathbf{x}_p (\mathbf{x}_p^T \mathbf{x}_p)^{-1} \mathbf{x}_p^T) \mathbf{y}}{n-p}.$$

True:
 $\mathbf{y} \sim N(\tilde{\mathbf{x}}\beta, \sigma^2 \mathbf{I}_n)$

$$\begin{aligned} E(\hat{\sigma}_p^2) &= \frac{\sigma^2(n-p) + \tilde{\beta}^T \mathbf{x}^T (\mathbf{I}_n - \mathbf{x}_p (\mathbf{x}_p^T \mathbf{x}_p)^{-1} \mathbf{x}_p) \mathbf{x} \beta}{n-p} \\ &= \sigma^2 + \frac{\tilde{\beta}^T \mathbf{x}^T (\mathbf{I}_n - \mathbf{x}_p (\mathbf{x}_p^T \mathbf{x}_p)^{-1} \mathbf{x}_p) \mathbf{x} \beta}{n-p} \\ &= \sigma^2 + \frac{\tilde{\beta}_{nr}^T \mathbf{x}_{nr}^T (\mathbf{I}_n - \mathbf{x}_p (\mathbf{x}_p^T \mathbf{x}_p)^{-1} \mathbf{x}_p) \mathbf{x}_{nr} \tilde{\beta}_{nr}}{n-p}. \end{aligned}$$

$\geq \sigma^2$ as $(\mathbf{I}_n - \mathbf{x}_p (\mathbf{x}_p^T \mathbf{x}_p)^{-1} \mathbf{x}_p^T)$
 is a p.s.d matrix.

~~$\mathbf{z} \sim N(\mu, \sigma^2 \mathbf{I}_n)$~~
 $E(\mathbf{z}^T \mathbf{A} \mathbf{z})$
 $= \#_r(\mathbf{A} \mathbf{I}_n \sigma^2) + \#_{nr}(\mathbf{A} \mathbf{x}_{nr})$

A is square, symmetric
 & idempotent.

$$\mathbf{x} \tilde{\beta} = \mathbf{x}_p \tilde{\beta}_p + \mathbf{x}_{nr} \tilde{\beta}_{nr}$$

If $\mathbf{e}(\mathbf{x}_p) \perp \mathbf{e}(\mathbf{x}_{nr})$

then $(\mathbf{I}_n - \mathbf{x}_p (\mathbf{x}_p^T \mathbf{x}_p)^{-1} \mathbf{x}_p^T)$ is the projection matrix of $\mathbf{e}(\mathbf{x}_{nr})$.

Q3 Wem: $\hat{\beta}_p \xrightarrow{\text{reduced mod}} \hat{\beta}$ has less variance compared the estimator of β_p under complete model? (18)

$$D(\hat{\beta}_p) = \sigma^2 (X_p^T X_p)^{-1}$$

$$\begin{aligned} D(\hat{\beta}) &= \sigma^2 (X^T X)^{-1} = \sigma^2 \left([\bar{X}_p : \bar{X}_r]^T [\bar{X}_p : \bar{X}_r] \right)^{-1} \\ &= \sigma^2 \left(\begin{matrix} X_p^T X_p & X_p^T X_r \\ -X_r^T X_p & X_r^T X_r \end{matrix} \right)^{-1} \boxed{(A + BCD)^{-1}} \\ &= A^{-1} - A^{-1} B (C^{-1} + D A^{-1} B)^{-1} D A^{-1} \end{aligned}$$

$$A^{-1} = \begin{bmatrix} A_{11} & A_{12} \\ A_{21} & A_{22} \end{bmatrix}^{-1}$$

$$= \begin{bmatrix} A_{11} - A_{12} A_{22}^{-1} A_{21} & (-[A_{11} - A_{12} A_{22}^{-1} A_{21}]^{-1} A_{12} A_{22}^{-1}) \\ -A_{21} A_{11}^{-1} [A_{22} - A_{22} A_{11}^{-1} A_{12}] & A_{22} - A_{21} A_{11}^{-1} A_{12} \end{bmatrix}$$

