

The LR statistic

$$-2 [LR(\hat{\beta}, \hat{\sigma}^2) - L(\hat{\beta}, \hat{\sigma})] \left[-2 \log [LR \text{ statistic}] \right]$$

$$\sim \chi^2_k \quad (\text{asymptotically under } H_0)$$

Score Test:

- Has advantage that we don't have to calculate the unrestricted MLEs.

For a general wt function w ,

$$\frac{\partial L}{\partial \underline{\beta}} = -\frac{1}{2} \sum_{i=1}^n \left\{ \frac{1}{w_i} - \frac{\epsilon_i^2}{w_i^2} \right\} \frac{\partial w_i}{\partial \underline{\beta}}$$

Let, D : the matrix whose i th row is $\frac{\partial w_i}{\partial \underline{\beta}}$

$$\text{start } u = (u_1, \dots, u_n)' ; u_i = \frac{\epsilon_i^2}{\hat{\sigma}^2}$$

& $\hat{w}_i = \hat{\sigma}^2$ under the null

$$\frac{\partial L}{\partial \underline{\beta}} \Big|_{H_0} = -\frac{1}{2} \sum_{i=1}^n \frac{\hat{\sigma}^2}{\hat{\sigma}^4} \frac{\epsilon_i^2}{\hat{\sigma}^2} \cdot \frac{\partial \hat{w}_i}{\partial \underline{\beta}}$$

Differentiate and taking expectation

$$I_{H_0} = \begin{bmatrix} I_{\underline{\beta} \underline{\beta}} & 0 \\ 0 & I_{\underline{\alpha} \underline{\alpha}} \end{bmatrix}$$

$$\frac{\partial L}{\partial \underline{\beta}} = x' \sum^{-1} (y - x\underline{\beta}) ; \frac{\partial^2 L}{\partial \underline{\beta} \partial \underline{\beta}'} = -2 \sum^{-2} - (x' \sum^{-1} x)$$

$$\sum = \hat{\sigma}^2 I$$

$$\frac{\partial^2 L}{\partial \underline{\alpha} \partial \underline{\beta}'} = \frac{\partial}{\partial \underline{\beta}} \left(\frac{1}{2} \sum_{i=1}^n \frac{(y_i - x_i' \underline{\beta})^2}{w_i^2} \right) \frac{\partial w_i}{\partial \underline{\alpha}}$$

(first term is independent of β)

$$= \sum_i \frac{(y_i - x_i' \beta)}{w_i^2} \frac{\partial w_i}{\partial \beta} x_i'$$

\rightarrow zero expectation.

$$\frac{\partial^2 l}{\partial \beta \partial \beta'} = -\frac{1}{2} \sum_{i=1}^n \left(\frac{1}{w_i} - \frac{\epsilon^2}{w_i^2} \right) \frac{\partial w_i}{\partial \beta} \frac{\partial w_i}{\partial \beta'}$$

$$+ \frac{1}{2} \sum_{i=1}^n \left(\frac{1}{w_i^2} - \frac{2\epsilon^2}{w_i^3} \right) \frac{\partial w_i}{\partial \beta} \left(\frac{\partial w_i}{\partial \beta} \right)$$

$$E \left[\frac{1}{w_i} - \frac{\epsilon^2}{w_i^2} \right] = 0 \quad E \left[\frac{1}{w_i^2} - \frac{2\epsilon^2}{w_i^3} \right]$$

$$E(\epsilon^2) = \sigma^2 \quad \begin{cases} = -\sigma^{-2} & \text{under the null} \\ = \frac{1}{w_i^2} - \frac{2w_i}{w_i^3} = -\frac{1}{w_i^2} & \end{cases}$$

$$E \left(\frac{\partial^2 l}{\partial \beta \partial \beta'} \right) = -\frac{1}{2} D'D\sigma^{-2}$$

$$T_{H_0} = \begin{pmatrix} \sigma^{-2}(x'x) & 0 \\ 0 & \frac{1}{2}\sigma^{-2}(D'D) \end{pmatrix}$$

The score statistic

$$\left(\frac{\partial l}{\partial \beta} \right)' T_{H_0}^{-1} \left(\frac{\partial l}{\partial \beta} \Big|_{H_0} \right) \sim \chi_m^2$$

m : no of restrictions

imposed by the null

So the score statistic takes the form

$$\frac{1}{2} (\underline{u} - \underline{l}_m)' D (D'D)^{-1} D' (\underline{u} - \underline{l}_m)$$

If $\omega(z, \lambda) = \exp(-z/\lambda)$, $D = \hat{\sigma}^2 z^{100}$ under the null.

& the test statistic is

$$\left(\frac{1}{n} \sum_{i=1}^n (x_i' \beta)^2 \right) / \hat{\sigma}^2$$

- If variances are known of unknown parameter

$$\sigma_i^2 = \omega(z_i, \lambda, \beta) \quad \text{--- } \textcircled{*} \textcircled{*}$$

\hookrightarrow This is to allow for the possibility that the variance depends on the mean and possibly on other parameters.

$$\sigma_i^2 = \lambda, (x_i' \beta)^2 ; \lambda = (\lambda_1, \lambda_2)$$

Algorithm:

- ① obtain an estimate $\hat{\beta}$ (e.g., by LS method)

$$\text{② Compute } \hat{\Sigma} \text{ as } (\hat{X}' \hat{X})^{-1} \hat{X}' \hat{Y}$$

$$\hat{\Sigma} = \text{diag} \{ \omega(z_1, \lambda^*, \hat{\beta}^*) ; \dots ; \omega(z_n, \lambda^*, \hat{\beta}^*) \}$$

using a WLS program.

for λ^* , find λ such that $\epsilon_i^2 \approx \sigma_i^2$

- ③ with $\hat{\beta}$ fixed at $\hat{\beta}^*$, obtain an estimate

$$\lambda^* \text{ of } \lambda$$

Regressing Residuals:

$$e_i \approx \epsilon_i, \text{ then } E(e_i^2) \approx \text{Var}(\epsilon_i) = \omega(z_i, \lambda^*, \hat{\beta}^*)$$

- Obtain an estimate of λ by non-linear LS problem

$$\min_{\lambda} \sum_{i=1}^n [e_i^2 - \omega(z_i, \lambda, \hat{\beta}^*)]^2$$

Since $\text{Var}(e_i^2) \approx \text{Var}(e_i^2) = 2\omega(z_i, \beta^*)$,

an alternative method is to solve

$$\min_{\beta} \sum_{i=1}^n (\epsilon_i^2 - \omega(z_i, \beta^*))^2$$

for a non-negative solution $\hat{\beta}$.

Variance is a function of the mean:

If the variance is a smooth f^n of the mean, we can do better. Suppose that, $\sigma_i^2 = \omega(x_i, \beta)$:

ω -Known.

Algorithm: I: Obtain an estimate of β , say $\hat{\beta}_{OLS}$

II: Calculate $\hat{\Sigma} = \text{diag}\{\omega(x_1, \hat{\beta}), \dots, \omega(x_n, \hat{\beta})\}$

Note: III: Recompute $\hat{\beta} = (\hat{\Sigma}^{-1} \hat{x})^{-1} \hat{x}^T \hat{y}$

IV: Repeat II & III until convergence.

Carroll & Reppen showed that the estimate provided by this algorithm has the same asymptotic efficiency as $\hat{\beta}_{OLS}$.

This means, in this case, there is no cost, in not knowing the weights provided ω is known.

The same is true if ω is unknown but smooth.

Carroll showed if we plot the squared residuals from a LS fit v/s the OLS fitted values A smooth the plot, we can do just as well.

Let, \hat{w}_i be the smoothed value of e_i^2 .

Then if we use WLS with weights $1/\hat{w}_i$, we

obtain an estimate whose asymptotic efficiency relative
to $\hat{\beta}_{WLS}$ is 100%.

obtain an estimate whose asymptotic efficiency relative to $\hat{\beta}_{WLS}$ is 100%.

Transforming to equalize variances:

We want to find an increasing function f such that $f(y_i) = x_i' \beta + \epsilon_i$

where ϵ_i have equal variance.

Suppose $\text{Var}(y_i) = w(x_i)$; $x_i' \beta = \mu$.

Assuming w is known,

$$\text{Var}[f(y)] \approx \left(\frac{df}{d\mu} \right)^2 v(y) \quad E[f(y)] \approx f(\mu) + (y - \mu) f'(\mu)$$

The variance of the transformed responses $f(y_i)$ will be approximately constant if we choose f so that

$$\begin{aligned} [f'(\mu)]^2 w(\mu) &\text{ is constant} \\ f(\mu) &= \int \frac{du}{\{w(u)\}^{1/2}} \end{aligned}$$

Example: $y_i \sim \text{Poisson } (\mu)$; $w(\mu) = \mu$.

$f(\mu) = \mu^{1/2}$
we expect $y^{1/2}, \dots, y_n^{1/2}$ have approximately equal variance.

Example: $y_i \sim \text{Binomial } (m, p)$, $\mu = mp$; $w(\mu) = \mu(1 - \mu/m)$

$$\text{Then, } \int \frac{du}{\sqrt{\mu(1 - \mu/m)}} \propto \sin^{-1} \left(\left(\frac{\mu}{m} \right)^{1/2} \right)$$

If w is not known, we can experiment by the responses y_i using a power transformation.

We can transform with a sample power less than 1 and use the diagnostic plots to see if the variances have been made equal.

- we proceed by reducing the power until the wedge effect in the plot of squared residuals v/s fitted values disappear.

Box-Cox Transformation / Power family of Transformation

- A flexible family of transformation

When the values of the response variable are all positive.

$$g(y; \lambda) = \begin{cases} \frac{y^\lambda - 1}{\lambda} &; \lambda \neq 0 \\ \log y &; \lambda = 0 \end{cases}$$

Power transformation in y^λ when $\lambda = 0$; $y^1 = 1$; making all data pt equal.

Departure from Normality:

- Non-normal errors are detected by normal plot of residuals.

① first the LS residuals.

② Plot these argument quantiles of the standard normal distⁿ.

- Suppose the normality assumptions is satisfied.

$\underline{\sigma}$ has a singular distⁿ $N(\underline{0}, \sigma^2(I-H))$

$N(0, \sigma^2(I-H))$ is approximately $N(0, \sigma^2 I)$ provided h_{ij} 's are small.

Here residuals are approximately a random sample from a $N(0, \sigma^2)$ distⁿ.

Normal prob. plot is a special case of Q-Q plot.

So that, $E[e_{(i)}] \approx \sigma z(\alpha_i)$

$\alpha_i = \frac{i-0.5}{n} \Rightarrow z(\alpha) \text{ is the } \alpha \text{th quantile of the}$

standard normal distribution

$\int_{-\infty}^{\alpha} f(z) dz = \alpha$

- Normal error - roughly a straight line.

- Skewed error - will be a curve.

- Heavy-tailed distⁿ - S-shaped.

- Outliers - isolated pts.

Transforming the response

If the normal plot reveals non-normality, the standard remedy is to transform the response.

- Box-Cox.

- Box-Cox introduced this transformation to remedy several type of regression problems, so that all the regression assumptions would be satisfied.

(inv in exp variable, homogeneous variances, normal error).

Assume that there is a transformation parameter λ \Rightarrow

$$y_i^{(\lambda)} = g(y_i; \lambda) = x_i' \beta + e_i^{(\lambda)}$$

Under this transformation the likelihood f^{λ} for the original observations is

$$(2\pi\sigma^2)^{-n/2} \exp \left[-\frac{1}{2} \left(\underline{y}^{(\lambda)} - \underline{x}\beta \right) \left(\underline{y}^{(\lambda)} - \underline{x}\beta \right)' \right] |J|$$

for each $y_i > 0$

$$\text{Here, } |J| = \left| \prod_{i=1}^n \frac{dy_i^{(\lambda)}}{dy_i} \right| = \prod_{i=1}^n y_i^{\lambda-1}$$

For fixed λ , \circ is the likelihood corresponding to a standard LS problem except for the constant factor J .

For fixed λ , the max^m value of the likelihood f^{λ}

$$\text{is } (2\pi\hat{\sigma}^2)^{-n/2} \exp[-\hat{\sigma}^2/2] |J|$$

$$\hat{\sigma}^2 = \underline{y}^{(\lambda)}' (\mathbf{I} - \underline{x}(\underline{x}'\underline{x})^{-1}\underline{x}') \underline{y}^{(\lambda)}$$

$$= \text{RSS}(\lambda; y), \text{ say}$$

Apart from a constant, the max^m log likelihood is

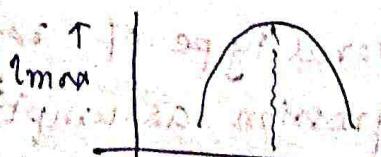
$$\begin{aligned} l_{\max}(\lambda) &= -\frac{1}{2} \ln \{ \text{RSS}(\lambda; y) \} + \frac{n}{2} \\ &\quad + (\lambda - 1) \sum_{i=1}^n \log y_i \end{aligned}$$

Box and Cox suggested

plotting $l_{\max}(\lambda)$ against λ

function for a trial series of

values & take that λ as λ which maximizes $l_{\max}(\lambda)$.



A more accurate value of $\hat{\lambda}$ can be obtained by

solving $\frac{d L_{\text{max}}(\lambda)}{d \lambda} = 0$

$$\Rightarrow \frac{y^\lambda - 1}{\lambda} = \frac{(y+c)^\lambda - 1}{\lambda}$$

Another family:

$$g(y; \lambda) = \begin{cases} \text{sign}(ay) \{1y+1\}/2, & \lambda \neq 0 \\ \text{sign}(y) \log \{1y+1\}, & \lambda = 0 \end{cases}$$

- defined ~~not~~ for all values of y .

- for each λ , g is a 1-1 monotone transformation

- general works well with symmetrically distributed observation with heavy tails.

Box-Cox is better ~~as~~ in case of skewed data

↳ can behave very badly in presence of outliers

$$y_i = x_i' \beta + \epsilon_i$$

Transforming both sides

If $E[y_i] = x_i' \beta$, but ~~as~~ the errors are non-normal and/or heteroscedastic.

- Then transforming the response will destroy the linear form of the mean.
- To avoid this; transform both sides,

$$g(y_i; \lambda) = g(x_i' \beta, \lambda) + \epsilon_i$$

where for some transformation $g(y; \lambda)$ and some value of λ , the errors ϵ_i are normally distributed with constant variance.

When $g(y; \lambda)$ is box-cox family, use ML method.

Then the transformed loglikelihood -

$$L(\underline{\beta}, \sigma^2, \lambda) = c - \frac{1}{2} \left\{ n \log \sigma^2 + \frac{1}{\sigma^2} \sum_{i=1}^n \left\{ g(y_i; \lambda) - g(x_i' \underline{\beta}; \lambda) \right\}^2 \right\}$$

$$+ (\lambda - 1) \sum_{i=1}^n \log y_i$$

$$\frac{\partial L}{\partial \sigma^2} = 0 \Rightarrow \hat{\sigma}^2(\underline{\beta}, \lambda) = \frac{1}{n} \sum_{i=1}^n (g(y_i; \lambda) - g(x_i' \underline{\beta}; \lambda))^2$$

$$\text{max } L(\underline{\beta}, \hat{\sigma}^2, \lambda) = c - \frac{1}{2} n \log \hat{\sigma}^2(\underline{\beta}, \lambda) + (\lambda - 1) \sum_{i=1}^n \log y_i$$

can be maximized by Fisher Scoring to

obtain the estimate of $\underline{\beta}$ & λ .

Alternatively, take $(\prod_{i=1}^n y_i)^{1/n} = y$

then the loglikelihood upto a constant then

$$- \frac{n}{2} \left[\log \hat{\sigma}^2(\underline{\beta}; \lambda) - 2 \lambda \log y \right]$$
$$= - \frac{n}{2} \log \left[\frac{\hat{\sigma}^2(\underline{\beta}; \lambda)}{(y)^2} \right]$$

is maximized by minimising

$$\sum_{i=1}^n \left[\frac{g(\underline{\beta}, \lambda) - g(x_i' \underline{\beta}; \lambda)}{(y)^2} \right]^2$$

- Standard error can be calculated from the information matrix.

This technique can be used to make error dist'

chosen to normal and also stabilize the variance.

- It is not necessary true that the same

transformation will achieve both the objections.

- If the transformation that ~~usually~~ makes the error distⁿ closer to normal does not also make the variance more homogeneous, it may be necessary to use weighting to achieve the constant variance (approx.) case.

Detecting Outliers and Dealing with Outliers

$$y_i = x_i' \beta + e_i \quad (\Delta + \gamma), \quad \hat{y}_i = x_i' \hat{\beta} + e_i$$

$\hat{\beta}$ estimates β .

$$\hat{y}_i = x_i' \hat{\beta} \quad ; \quad e_i = y_i - \hat{y}_i \Rightarrow y_i = \hat{y}_i + e_i$$

It is reasonable to consider e_i is an estimate of e_i .

- We can use the raw residuals e_i or the standardized residuals. $\frac{e_i}{\sigma_i}$ & $\frac{e_i}{s_i}$ are called raw residuals and $\frac{e_i}{\sigma_i}$ & $\frac{e_i}{s_i}$ are called standardized residuals.

\hookrightarrow identically distributed.

- Standard graphical plots such as box plot of residuals. \hookrightarrow we can identify an outliers, the points having large residuals.

- $|t_i| \sim t_{n-p-1}$; in absence of outliers, a reasonable definition of "large" is a pt for which $|t_i| > 2$.

\rightarrow The diagnostic approach works well which does not have high leverage.

- If it does we cannot expect the corresponding residual to reveal the presence of an outliers.

→ Suppose that the i th-response is

$$x_i y_i = \Delta_i \text{ rather than } y_i \text{ so, } x_i' \beta + \Delta_i + e_i$$

$$\text{Let, } \Delta = (0, 0, \dots, \Delta_i, 0, \dots, 0)'$$

$$\underline{\underline{(I-H)y}}; \text{ which has a slight problem}$$

$$\begin{aligned} \mathbb{E}[e] &= \mathbb{E}[(I-H)\Delta] = \mathbb{E}(x\beta + \Delta) \\ &= (I-H)\Delta \end{aligned}$$

$$\mathbb{E}[e_i] = (1-h_i) \Delta_i$$

If x_i is close to \bar{x} , h_i is small, we can expect the residual to reveal the outlier quite well.

→ Here $E[e_i]$ is close to Δ_i

→ If the data point is a high-leverage point, then hat matrix diagonal is close to 1,

so the residual will be much smaller than Δ_i

- Hat matrix diagonals are reasonable measures of leverage, can be interpreted in terms of M_{ii}

- the average of hat matrix diagonal \bar{h}/n , an arbitrary but reasonable definition of high-leverage pt is one which satisfy $h_i > 2\bar{h}/n$

Whether a point has high leverage depends on the ~~exp~~ variables included in the regression.

Q: Having identified that a pt is having (high) leverage, how can you detect whether the pt is an outlier.

- ① A ~~estimate~~ the effect that a data pt has on the regression by deleting and refitting the regression.
 - If the regression quantity change drastically, then the pt ~~is~~ in is a high influence pt or an outlier.

↳ leave-one-out-diagnostic $\left\{ \begin{array}{l} \text{regression coefficients,} \\ \text{fitted values, standard} \\ \text{errors, etc.} \end{array} \right.$

- ② Use a robust fitting method that is not affected by high-leverage pts., resulting in residuals, that better identify the outlier.

Leave-one-out Diagnostics

$$\text{Th: } \hat{\beta} - \hat{\beta}(i) = \frac{(x'x)^{-1} x_i e_i}{1-h_i}; \quad x_i: i\text{th row of } x.$$

Change in Estimated regression coefficients

i: $\hat{\beta}(i)$: LSE of $\hat{\beta}$ when i th obs is deleted.

$\hat{\beta} - \hat{\beta}(i) = \frac{(x'x)^{-1} x_i e_i}{1-h_i}$: the change is proportional to the size of the residual but is inflated if h_i is close to 1.

DFBETA

Let, $C = (x'x)^{-1} x'$ → catcher matrix

The (j,j) th element of C is significant if

$$e_{j,j} = [(x'x)^{-1} x_j]_j$$

So, the j th element of DFBETA is

$$\ell_{j,j}$$

$$\sum_i e_{j,i}^2 = \sum_i [(x'x)^{-1} x_i]_j^2$$

$$\begin{aligned} &= \sum_i \left[(x'x)^{-1} x_i x_i' (x'x)^{-1} \right]_j \\ &= \left[(x'x)^{-1} \sum_i x_i x_i' (x'x)^{-1} \right]_j \end{aligned}$$

$$\left[C C' = (x'x)^{-1} x' x (x'x)^{-1} = (x'x)^{-1} \right]_j$$

- Standardize the $(j+1)$ th element of DFBETA by an estimate of the standard error of $\hat{\beta}_j$, $[(x'x)^{-1}]_{j+1,j+1}^{1/2} s(i)$

$$= s(i) \left(\sum_i e_{j+1,i}^2 \right)^{1/2}$$

$$\text{DFBETA}_{j,j} = \frac{\hat{\beta}_j - \hat{\beta}(i)_j}{s(i) \left(\sum_i e_{j+1,i}^2 \right)^{1/2}}$$

In terms of externally studentized residuals h_i ,

$$\text{DFBETA } s_{i,j} = \frac{c_{j+1,i} e_i}{s(i) \left(\sum_i c_{j+1,i} \right)^{1/2} (1-h_i)}$$

$$= \frac{c_{j+1,i} \cdot t_i}{\left(\sum_i c_{j+1,i}^2 \right)^{1/2} (1-h_i)^{1/2}}$$

$$[\text{as } t_i = \frac{e_i}{s(i) \sqrt{1-h_i}}]$$

If the i th data pt is not an outlier and does not have high leverage we expect that

$$|t_i| < 2.$$

and y_i will not have a large effect on the value of $\hat{\beta}_j$. Since, $\hat{\beta}_j = \sum_{i=1}^n c_{j+1,i} y_i$

The $c_{j+1,i}$ is small in relation to $\left(\sum c_{j+1,i}^2 \right)^{-1/2}$

and approximately $n^{-1/2}$

A suitable cut off value for detecting large residuals using DFBETA $s_{i,j}$ is $2/\sqrt{n}$.

$$\beta_1, \dots, \beta_n, i=1, 2, \dots, n$$

$$n \rightarrow \hat{\beta}_n$$

$$n \rightarrow \hat{\beta}_2$$

DFFIT

Change in fitted values

The change in the i th obs fitted values is

$$\frac{x_i' \hat{\beta} - x_i' \hat{\beta}(i)}{1-h_i} = \frac{x_i' (x'x)^{-1} e_i}{1-h_i}$$

Lif DFFIT, $\frac{h_i e_i}{1-h_i}$

Standardize by the estimated standard error

$$s(i) h_i^{1/2}$$

$$\text{For DFFIT } s_{\text{DFFIT}} = \frac{h_i e_i}{s(i) h_i^{1/2} (1-h_i)}$$

$$= \frac{h_i^{1/2} e_i}{(1-h_i)}$$

$$\text{var}(x_i' \hat{\beta})$$

$$= h_i \sigma^2$$

$$= x_i' V(\hat{\beta}) x_i$$

$$= \sigma^2 x_i' (x'x)^{-1} x_i$$

$$= \sigma^2 h_i$$

$$= t_i \left(\frac{h_i (1-h_i)^{1/2}}{1-h_i} \right)$$

(for pts which are not outliers and do not have high leverage, $|t_i| < 2$ & h_i will not be far away from the average value k/n)

The cutoff pt DFFIT s_{DFFIT} is $2\sqrt{\frac{p}{n-p}}$ or

$$\text{even } 2\sqrt{\frac{p}{n}}$$

Cook's Distance:

(Cook's D):

[Distance between $\hat{\beta}(i)$ & $\hat{\beta}$]

$$z_1 = \frac{(\hat{\beta} - \beta)' x' x (\hat{\beta} - \beta)}{\sigma^2} \sim \chi_p^2$$

$$z_2 = \frac{RSS}{\sigma^2} = \frac{(n-p)s^2}{\sigma^2} \sim \chi_{n-p}^2$$

$$\frac{z_1/p}{z_2/(n-p)} = \frac{(\hat{\beta} - \beta)'(x'x)(\hat{\beta} - \beta)}{p s^2} \sim F_{p, n-p}$$

[$\hat{\beta}$ and s^2 are independent]

- So a $100(1-\alpha)\%$ confidence ellipsoid for β is

$$\{\beta : (\beta - \hat{\beta})' x'x (\beta - \hat{\beta}) \leq p s^2 F_{p, n-p}^{\alpha}\}$$

- Cook's distance measure D_i is defined as the distance of $\hat{\beta}_i$ and $\hat{\beta}(i)$

$$D_i = \frac{(\hat{\beta}(i) - \hat{\beta})' x'x (\hat{\beta}(i) - \hat{\beta})}{p s^2} \quad 1 \leq i \leq n$$

- based on confidence ellipsoid.

$$D_i = \frac{x_i'(x'x)^{-1}(x'x)^{-1}x_i e_i^2}{(1-h_i)^2 p s^2}$$

$$D_i = \frac{e_i^2 \cdot \frac{(1-h_i)^2 - 1}{p s^2}}{(1-h_i)^2} = \frac{(1-h_i)^2 - 1}{p (1-h_i)}$$

pts with large $|D_i|$ have considerable influence on the

LSE of β . Cook's suggestion $D_i > F_{p, n-p}^{0.10}$

- Similar to S_q of DFFITSS

(differing only in use of n_i instead of t_i and constant divisor p).

Another way of writing Cook's D

$$D_i = \frac{(\hat{y}_{(i)} - \hat{y})' (\hat{y}_{(i)} - \hat{y})}{p s^2}, \quad x \hat{\beta}(i) = \hat{y}_{(i)}$$

Covariance Ratio:

The estimated variance-covariance matrix of

$\hat{\beta}_0$ is $s^2(x'x)^{-1}$ & its leave-one-out version

$$\text{is } s(i)^2 [x(i)' x(i)]^{-1}$$

$$\text{COVRATIO}_i = \frac{\det [s(i)[x(i)' x(i)]^{-1}]}{\det [s^2(x'x)^{-1}]}$$

$$(n-p-1) s(i)^2 = (n-p) s(i)^2 - \frac{e_i^2}{1-h_i}; [s(i)^2 = \frac{e_i^2}{1-h_i}]$$

$$\Rightarrow \frac{s^2}{s(i)^2} = \frac{(n-p-1)}{n-p} + \frac{e_i^2}{s(i)^2 (1-h_i) (n-p)}$$

$$= \frac{n-p-1}{n-p} + \frac{e_i^2}{(n-p)}$$

$$\det [x(i)' x(i)] = \det [x'x - x_i x_i']$$

$$\begin{aligned} &= |x'x| \left[1 - \frac{x_i'(x'x)^{-1} x_i}{\det(A+uu')} \right] \left[= |A| (1 + u' A^{-1} u) \right] \\ &= |x'x| [1 - h_i] \end{aligned}$$

$$\text{So, COVRATIO}_i = \frac{s(i)^2 \{x(i)' x(i)\}^{-1}}{|x'x|}$$

$$= \frac{(s(i)^2)^p}{(s^2)^p} \left| \left(x(i)' x(i) \right)^{-1} \right|$$

$$= \frac{(n-p-1)}{(n-p)} \left(\frac{e_i^2}{(n-p)} \right)^{-p} \frac{1}{1-h_i}$$

$$(E - (i)^2) (E - (i)^2)$$

Extreme Cases:

① if $|t_{ii}| > 2$, but this case has the minimum leverage with $h_{ii} = \frac{1}{n}$ [$h_{ii} \geq \frac{1}{n}$ for regression models with intercept term]

$$\text{COVRATIO}_i = \left(1 + \frac{(t_{ii}^2 - 1)}{n-p} \right)^{-p} \frac{n}{(n-1)}$$

$$\approx \left[\frac{\left(\frac{n-1}{n} \right) \times (t_{ii}^2 - 1)^p}{n-p} \right] \cdot \frac{n}{n-1}$$

$$\approx \frac{1}{1 - \frac{(t_{ii}^2 - 1)^p}{n}}$$

When $n \gg p$. Therefore if $|t_{ii}| > 2$

$$\text{COVRATIO}_i \approx 1 - \frac{(t_{ii}^2 - 1)^p}{n} \ll 1 - \frac{3p}{n}$$

② The opposite situation in when t_{ii} is small but the case has high leverage.

Take $t_{ii} = 0$, $h_{ii} > 2p/n$

$$\begin{aligned} \text{COVRATIO}_i &= \left(1 - \frac{1}{n-p} \right)^{-p} (1-h_{ii})^{-1} \\ &\gg \left(1 - \frac{1}{n-p} \right)^{-p} \left(1 - \frac{2p}{n} \right)^{-1} \end{aligned}$$

$$\approx \left(1 + \frac{p}{n} \right) \left(1 + \frac{2p}{n} \right)$$

$$\approx \left(1 + \frac{p}{n} \right) \left(1 + \frac{2p}{n} \right) \text{ for large } n.$$

$\hat{h}_{ii} = \text{mult} \frac{3p}{n}$ [ignoring the higher order terms].

The cases having $|\text{COVRATIO}_i - 1| > \frac{3p}{n}$ are considered to have high influence.

[For high lever, COVRATIO will large]

Deleting Several Influential points

Consider there are $m (> 1)$ obs to be assessed for influence simultaneously.

Set, i denote the $m \times 1$ vector of indices specifying the pts to be assessed.

$$\text{Define, } D_i = \frac{(\hat{\beta}(i) - \hat{\beta}) \times (\hat{\beta}(i) - \hat{\beta})}{\text{ps}^2}$$

→ extension of single obs Cook's D.

Large value of D_i indicates that the set of m pts are influential.

- Selection of the set of pts to include in i , is

not obvious, as in some situations, the m pts are jointly influential but are not so.

- It is also difficult to choose all m -subsets of the n pts to investigate

Leave-many-out

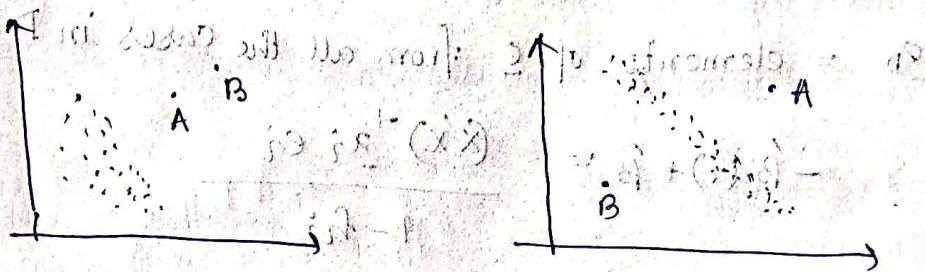
Single point deletion diagnostic will not be able to identify a pt as an outlier if it is masked by other.

→ leave-many-out diagnostic

Drawbacks

Practically applicable for small values of d .

It will be effective in identifying the situation only if all the points in a cluster are deleted.



If the cluster contains d pts, we must calculate leave- d -out diagnostics.

- This is a problem since d is not known in advance.

- The procedure is computationally ~~too~~ feasible for small values of d .

- Let D be the d subset of cases.

$X(D)$: $(n-d) \times p$ submatrix of X corresponding to the rows of the cases not in D .

X_D : $d \times p$ submatrix corresponding to the rows of the cases in D .

Then $[X(D)' X(D)]^{-1} = [x'x - X_D' X_D]^{-1}$

$$[X(D)' X(D)]^{-1} = [x'x - X_D' X_D]^{-1} = x'x - X_D' (I_d - H_D)^{-1} X_D (x'x)^{-1}$$

$A_{m \times m}, B_{n \times n} \rightarrow n s \quad U_{m \times n}, V_{n \times n}$

$$(A + UBV')^{-1} = A^{-1} - A^{-1}UB(B - BV'A^{-1}UB)^{-1}BVA^{-1}$$

$$H_D = X_D (x'x)^{-1} X_D'$$

$$\hat{\beta}_{(D)} = \hat{\beta} - (x'x)^{-1} x_D' (I_d - H_D)^{-1} e_D$$

e_D = elements of e from all the cases in D .

$$-\hat{\beta}(i) + \hat{\beta}_i = \frac{(x'x)^{-1} x_i e_i}{1-h_i}$$

The difference in fitted values for the cases in D .

$$\hat{y}_D - \hat{y}_D(D) = x_D(\hat{\beta}) - x_D(\hat{\beta}_{(D)})$$

$$\text{resp. for } i \text{ is } \hat{y}_i = H_D (S_D - H_D)^{-1} e_D$$

The analog of the COVRATIO can be shown

$$\left\{ \left[\frac{n-p-d}{n-p} + \frac{e_D' (S_d - H_D)^{-1} e_D}{n-p} \right] \det(I_d - H_D) \right\}$$

Test for outliers:

To test if a fixed set of K observations contains outliers, assuming the remaining $n-K$ observations are "clean"

Arrange the data so that the clean obs come first, followed by the possibly outlying obs.

Outlier Stuff Model:

$$\underline{y} = x \underline{\beta} + z \underline{\gamma} + e \quad [x \quad z \quad e]$$

$$z \equiv \begin{pmatrix} 0 & \overset{n-K}{\underset{K}{\rightarrow}} & x \\ I_K & \underset{K \times K}{\rightarrow} & x \end{pmatrix}$$

$\underline{\gamma}$: a K -vector containing the shifts for the outlying observation.

To test $\underline{\gamma} = \underline{0}$ i.e. $X = Y$ follows

Let H be the hat matrix for the regression

$$\underline{Y} = X\underline{\beta} + \underline{\epsilon}$$

Let, $H = \begin{pmatrix} H_{11} & H_{12} \\ H_{21} & H_{22} \end{pmatrix}$, $H_{11} \in (n-k) \times (n-k)$

Partition the residual $(I_n - H)\underline{Y} = \underline{\epsilon}$ as $\underline{\epsilon} = (\underline{\epsilon}_1', \underline{\epsilon}_2')'$

$$\underline{\epsilon}_1 = (n-k) \times 1$$

Then the LSE of $\underline{\gamma}$ is

$$\hat{\underline{\gamma}} = \underline{Z}' \underline{Z}^{-1} (I-H) \underline{Y}$$

$$= [(0, I_K) (I-H) \begin{pmatrix} 0 \\ I_K \end{pmatrix}]^{-1} (0, I_K) \underline{\gamma} \begin{pmatrix} \underline{\epsilon}_1 \\ \underline{\epsilon}_2 \end{pmatrix}$$

$$= 0 \cdot (I_K - H_{22})^{-1} \underline{\epsilon}_2 \text{ by inspection since } H_{22} \text{ is full rank}$$

$$F = \frac{(RSS_H - RSS)/2}{RSS/(n-p)} ; RSS_H = \underline{Y}'(I-H)\underline{Y}$$

$$RSS = \underline{Y}' R_G \underline{Y} = \underline{Y}' R_Y \underline{Y}$$

$$\hat{\underline{\gamma}}' R_Y \underline{Y} = \hat{\underline{\gamma}}' (0, I_K) (I-H) \underline{Y}$$

$$R = I-H$$

$$H = \underline{Y}' \underline{\epsilon}_2 \underline{\epsilon}_2' (I_K - H_{22})^{-1} \underline{\epsilon}_2$$

$$= \underline{\epsilon}_2' (I_K - H_{22})^{-1} \underline{\epsilon}_2$$

$$So, RSS_H - RSS = \underline{\epsilon}_2' (I_K - H_{22})^{-1} \underline{\epsilon}_2$$

So, the F-test for $\underline{\gamma} = \underline{0}$ is based on

$$F = \frac{[\underline{\epsilon}_2' (I_K - H_{22})^{-1} \underline{\epsilon}_2] / (K-n-p)}{[RSS - \underline{\epsilon}_2' (I_K - H_{22})^{-1} \underline{\epsilon}_2] / ((n-p-K))} \sim F_{K, n-p-K}$$

RSS \rightarrow RSS of model $y = x\beta + \epsilon$ $\hat{y} = \hat{x}\hat{\beta}$

$$\text{when } k=1, \text{ with } e_2' (I_{k+1} - H_{22})^{-1} e_2 = \frac{e_i^2}{1-h_i}$$

(corresponding \hat{y}_i to i th obs) $\hat{y}_i = \hat{x}_i \hat{\beta}$

$$F = \frac{(n-p-1) e_i^2 / (1-h_i)}{(n-p-1) \text{RSS} / e_i^2} \sim F_{1, n-p-1} \text{ (under } H_0)$$

$(\hat{y}_i - \hat{y})^2$ as $\hat{y}_i = \hat{y}(H_{ii})$ \hat{y} unbiased \hat{y} (when i th obs is not an outlier)

$$\text{Using } r_i^2 = \frac{e_i^2}{(1-h_i) s^2}$$

$$F = \frac{(n-p-1) r_i^2}{(n-p-r_i^2) / ((1-h_i) s^2)}$$

Masking and Swamping:

- Hat matrix diagonal and leave-one-out diagnostic are useful tools, but sometimes cannot identify outliers & high-leverage pt.

- Taken singly, points A and B are not influential, since either will have little effect on the fitted line, because remaining pts cont to affect the line.

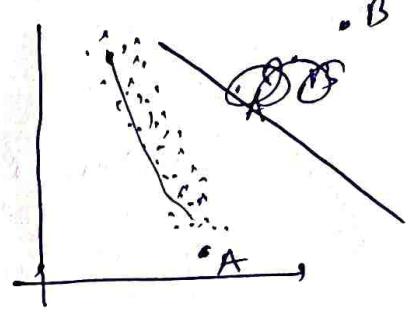
- Taken as a pair, they are influential.

- In this case, pts A & B mask each other.

- To use robust method in the least median squares to identify the cluster or use

leave-many-out (diagnostics - 22.1)

Pt A will have large residuals, since
 β is attracting the line away from A.



Swamping: whom pts that are not outlying can be mistaken as outliers.

- Sampling does not occur when robust methods are used.

$p=3$, two centered & scaled explanatory variables:

$$\underline{x_1^*}, \underline{x_2^*}$$

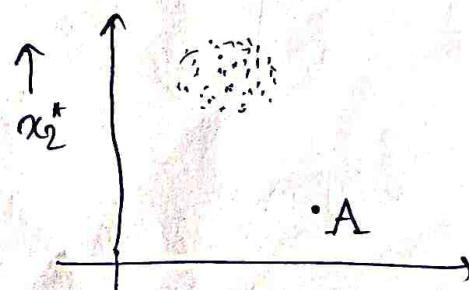
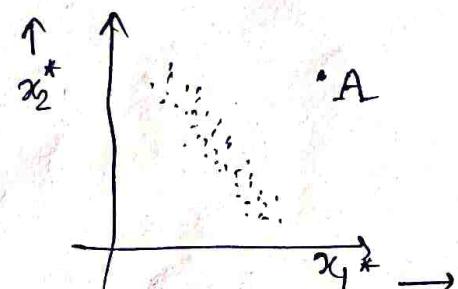
$$K(-i)$$

→ after removing i th obsⁿ,
then calculate & condⁿ no.

• Compare K and

$$K(i)$$

[K : CN on all obsⁿ].



Basic Time Series:

Defⁿ: (The Auto covariance function (ACVF)).

If $\{x_t : t \in \mathbb{Z}\}$ is a process $\Rightarrow \text{Var}(x_t) < \infty \forall t \in \mathbb{Z}$

then the ACVF of $\{x_t\}$ is

$$\gamma_x(r, s) = \text{Cov}(x_r, x_s), \text{ written as covariance}$$

$$= E[x_r - E(x_r)][x_s - E(x_s)],$$

$$r, s \in \mathbb{Z}$$

Stationary:

The TS $\{x_t : t \in \mathbb{Z}\}$ with index set $\mathbb{Z} = \{0, \pm 1, \pm 2, \dots\}$ is said to be (weak) stationary

if (i) $E|x_t|^2 < \infty \quad \forall t \in \mathbb{Z}$

(ii) $E[x_t] = m \quad \forall t \in \mathbb{Z}$

(iii) $\gamma_x(r, s) = \gamma_x(r+t, s+t) \quad \forall r, s, t \in \mathbb{Z}$

This is also known as Covariance stationary, stationary in wide sense or second order stationary.

If $\{x_t : t \in \mathbb{Z}\}$ is stationary $\gamma_x(r, s) = \gamma_x(r-s, 0)$

So we redefine the ACVF of a stationary TS as

$$\gamma_x(h) = \gamma_x(h, 0) = \text{Cov}(x_{t+h}, x_t), \quad h \in \mathbb{Z}$$

ACF (Auto Correlation Function) of $\{x_t\}$ is

$$p_x(h) = \frac{\gamma_x(h)}{\gamma_x(0)} = \text{Cov}(x_{t+h}, x_t) \quad \forall z, h \in \mathbb{Z}$$

Defⁿ: The TS $\{x_t : t \in \mathbb{Z}\}$ is strictly stationary if the joint distⁿ of $(x_{t_1}, x_{t_2}, \dots, x_{t_k})'$ and $(x_{t_1+h}, x_{t_2+h}, \dots, x_{t_k+h})'$ and the sense for all positive integer K and for all $t_1, t_2, \dots, t_k, h \in \mathbb{Z}$

Proposition:

If $\gamma(\cdot)$ is the ACVF of a stationary process, then

$$\gamma(0) > 0$$

$$|\gamma(h)| \leq \gamma(0) \quad \forall h \in \mathbb{Z}$$

$$\gamma(h) = \gamma(-h), \quad h \in \mathbb{Z}$$

MA(q): [Moving Average of order q]

$$X_t = Z_t + \theta_1 Z_{t-1} + \theta_2 Z_{t-2} + \dots + \theta_q Z_{t-q},$$

$$Z_t \sim WN(0, \sigma^2)$$

uncorrelated

AR(P): [AutoRegressive of Order P]

$$X_t - \phi_1 X_{t-1} - \phi_2 X_{t-2} - \dots - \phi_p X_{t-p} = Z_t;$$

$$Z_t \sim WN(0, \sigma^2)$$

Z_t is uncorrelated with X_s , $s < t$

ARMA(p,q)

$$X_t - \phi_1 X_{t-1} - \dots - \phi_p X_{t-p} = Z_t + \theta_1 Z_{t-1} + \dots + \theta_q Z_{t-q}$$

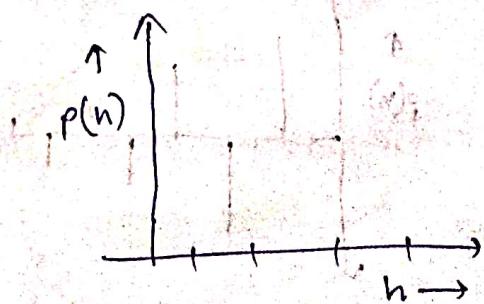
MA(1): $X_t = \theta Z_{t-1} + Z_t$

$$\gamma(X_t) = \gamma(0) = \sigma^2(1 + \theta^2)$$

$$\text{Cov}(X_t, X_{t+1}) = \gamma(1) = \theta \sigma^2$$

$$\text{Cov}(X_t, X_{t+h}) = \gamma(h) = 0, \quad \forall h \geq 2$$

$$\rho(1) = \theta / (1 + \theta^2) = \frac{\gamma(1)}{\gamma(0)}$$



Residuals are from an MA(1) process,

$$\text{Var}(\epsilon) = \sum = \sigma^2 \begin{pmatrix} 1+\theta^2 & \theta & 0 & \dots & 0 \\ \theta & 1+\theta^2 & \theta & \dots & 0 \\ 0 & \theta & 1+\theta^2 & \dots & 0 \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ 0 & \dots & 0 & \dots & 1+\theta^2 \end{pmatrix}$$

$$Y = X\beta + \epsilon ; \epsilon \sim (0, \Sigma)$$

$\Sigma_{(1)}$ from the residual by estimating θ .

$$\hat{\beta}_{(1)} = [X' \Sigma_{(1)}^{-1} X]^{-1} X' \Sigma_{(1)}^{-1} Y$$

Calculate the residual.

$$\text{Then } \hat{\theta} \Sigma_{(2)} \& \hat{\beta}_{(2)}$$

$$\hat{\beta}, \hat{\theta} \rightarrow \hat{z}_t , \epsilon_t = z_t + \theta z_{t-1}$$

$$\hat{\epsilon}_t \rightarrow \text{[of } \hat{z}_t \text{ is } \theta \text{ times previous residual]} : (1) AR$$

$$\underline{\text{AR}(1)}: X_t = \phi X_{t-1} + z_t , |\phi| < 1$$

$$\gamma(h) = \text{Cov}(X_t, X_{t+h})$$

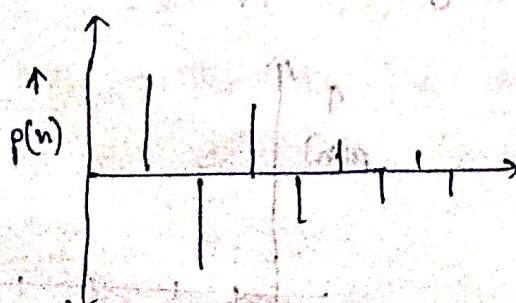
$$= \sigma^2 \sum_{j=0}^{\infty} \phi^j \phi^{j+h}$$

$$= \frac{\sigma^2 \phi^h}{1-\phi^2} ; h \geq 0$$

$$\rho(h) = \frac{\gamma(h)}{\gamma(0)}$$

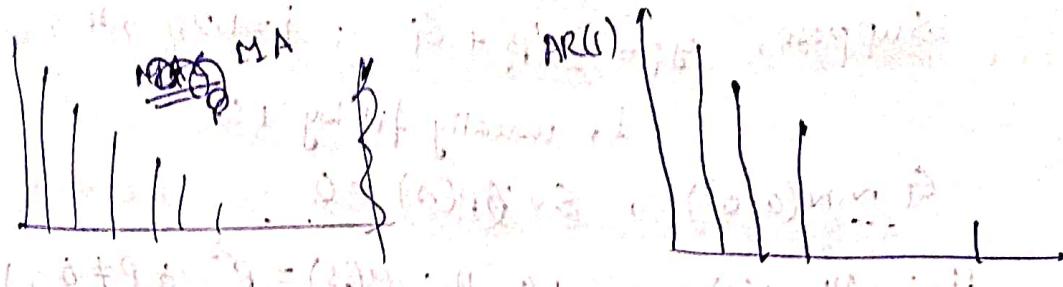
$$= \phi^h ; h \geq 0$$

$$\gamma(0) = \nu(x_t) = \frac{\sigma^2}{1-\phi^2}$$



PACF (Partial Auto Correlation ρ^h)

PACF (Partial Auto-correlation Function)



Serial Correlation, and Durbin Watson Test:

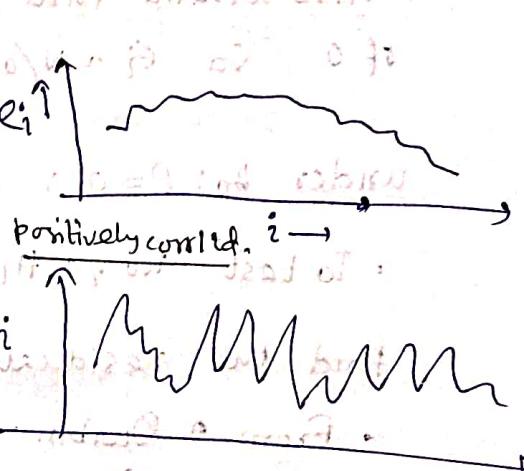
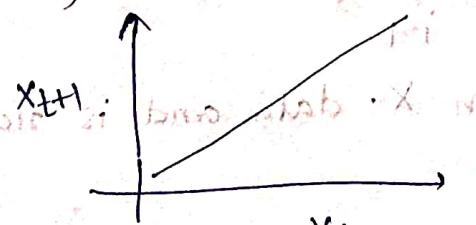
If the data are collected sequentially in time, then successive errors may be correlated.

If this is the case, a sequence plot of e_i against time order, ~~where~~ which is after e_i v/s i may show up the presence of any correlation.

For positively correlated error, a residual tends to have the same sign as its predecessor.

While for negatively correlated error, the signs

$$(x_t, x_{t+1})$$



Divide the ~~data~~ time ordered residuals into consecutive pairs and plotting are member of the pair against the other.

Serially correlated data show up as a linear trend in the ~~residual~~ plot.

Durbin-Watson Test:

Suppose, $y_i = x_i' \beta + \epsilon_i ; i=1,2,\dots,n$

↳ usually fit by LS.

$\epsilon \sim N(0, \sigma^2)$, so $E(\epsilon) = 0$

$H_0: \text{All } \rho(s) = 0 \text{ v/s } H_1: \rho(s) = \rho^s ; \rho \neq 0, |s| < 1$

An alternative that arises from the assumption that the error ϵ_{ij} are plotted & labelled are not all ρ .

$\epsilon_i = \rho \epsilon_{i-1} + z_i , z_i \sim N(0, \sigma^2)$ form shows

is indep of $\epsilon_{i-1}, \epsilon_{i-2}, \dots$

Also assume that mean & Var of ϵ_i are const, indep of i , so $\epsilon_i \sim N(0, \frac{\sigma^2}{1-\rho^2})$.

under $H_0: \rho = 0$; $\epsilon_i \sim N(0, \sigma^2) \rightarrow$ usual assumption

To test H_0 v/s H_1 , fit the model $y_i = x_i' \beta + \epsilon_i$ &

find the residuals e_i . Statistics plotted are e_i

From Durbin-Watson Statistic:

$$d = \sum_{i=2}^n (e_i - e_{i-1})^2 / \sum_{i=1}^n e_i^2$$

The distⁿ of d depends on X-data and is not independent of them.

The distⁿ of d lies between 0 & 4 & is symmetric about 2.

Percentage pts also depends on X-data & would have to be calculated each application to perform the test properly.

- Test is usually performed using tabulated bounds (d_L, d_U) instead of a single critical value, have to look into two critical values.
- d is used only for a lower tailed test $\rho > 0$.
- To test against $\rho < 0$, we need an upper tailed test. Can be handled as a lower tailed test using the statistic $4-d$.

1. One Sided Test: for $\rho > 0$

If $d < d_L$, d is significant, reject H_0 at level α

If $d > d_U$, d is not " , do not reject H_0 .

If $d_L \leq d \leq d_U$ → the test is inconclusive.

2. One sided Test: $\rho < 0$

Repeat (1) using $4-d$ instead of d .

3. Two sided equal tailed test $\rho \neq 0$

If $d < d_L$ or $4-d < d_L$, d is significant, reject H_0 at level α .

If $d > d_U$ or $4-d > d_U$, d is not significant, do not reject H_0 at level α .

Otherwise, the test is inconclusive.

A simplified version Test:

$H_0: \rho > 0$, If $d < d_L$, reject H_0 at level α , o.w.

do not reject.

$H_0: \rho < 0$, If $4-d < d_U$, reject H_0 at level α ,

o.w. do not reject.

$H_0: \rho = 0$, If $|d| < d_{\alpha/2}$ or $|d| > d_{1-\alpha/2}$, reject H_0 at level α .

of (D-W test). Follow writing assumption for basic

basic requirement of D-W test satisfies and
 $0 < \rho < 1$ if both series not prove linearly related

basic require a base few $(0, 1)$ following fact of ST -

series test basic several series are bivariate normal + $\text{E}(d) = 0$
and $\text{E}(d^2) = \sigma_d^2$

Diagonosing Collinearity:

The square of the condition no is an upper bound for the VIFs and

- If the condition no is small, there will be no regression coefficients with large variance.
- If the condⁿ no is large at least one eigen value must be small.
- The eigen value corresponding to the smallest eigen value should also be examined to determine the linear combinations that are causing the trouble.
- Can also calculate the variance proportions.

The variance proportion for the jth variable eigen value

λ_j .

The variance proportion for the jth var given value λ_j .

$$t_{j,r}^2 \lambda_r^{-1}$$

- Proportion near unity indicate that most of the var of a particular regression coeffs is due to a single small eigen value.

- High leverage pts can have a effect on collinearity.
- A more focused diagnostic is to compute the condition no $K(-i)$ of the regression with ith row removed.

- Comparison of $K(-i)$ and K (cn of the X matrix) will indicate any pts i.e. causing a drastic change.

Remedies of Collinearity:

- Collinearity arises when the data are deficient.

→ collect fresh data to repair the deficiencies

in the reg matrix.

- not always possible, when the exp variable are strongly correlated in the popl from which the data were collected.

- Another way is simply to discard the variables until the data are not collinear.

- The deleted variables may have a relationship with the response.

- Another way is to abandon the method of LS and use a biased estimation method, say ridge regression.

Ridge Regression:

- Introduced to deal with collinear data.

- LS method might result in a very poor estimates of the regression coefficients which the data are collinear.

- LSE is unbiased.

- In same situation, a biased estimator may have a smaller MSE than the var of an unbiased estimator.

→ ridge estimator

Consider the scaled and centered reg. model

$$y_i = \alpha + \gamma_1 x_{i1} + \gamma_2 x_{i2} + \dots + \gamma_{p-1} x_{ip} + \epsilon_i, \quad i=1, 2, \dots, n.$$

The ridge estimate of $\underline{\gamma}$ is defined as

$$(X^* X^* + K I)^{-1} X^* \underline{\gamma}$$

where K is a free parameter whose value must be specified. K is not necessarily an integer. When $K=0$, ridge estimate reduces to LSE.

Interpolation: It is basically the LSE of the data with augmented.

Imagine augmenting the observed data $(X^*, \underline{\gamma})$ with new data $(K I_{p-1}, 0)$.

The LSE calculated from the augmented data

$$\begin{pmatrix} X^* & \underline{\gamma} \\ K I_{p-1} & 0 \end{pmatrix} \text{ is } 2 \left[\begin{pmatrix} X^* & K I_{p-1} \\ K I_{p-1} & 0 \end{pmatrix} \right]^{-1} \begin{pmatrix} X^* & K I_{p-1} \\ K I_{p-1} & 0 \end{pmatrix} \underline{\gamma}$$

is just the ridge estimate.

$$\begin{aligned} \hat{\underline{\gamma}}(K) &= \left[X^* X^* + K I_{p-1} \right]^{-1} X^* \underline{\gamma} \\ &= [R_{xx} + K I_{p-1}]^{-1} R_{xx} R_{xx}^{-1} X^* \underline{\gamma} \\ &= [R_{xx} + K I_{p-1}]^{-1} R_{xx} \underline{\gamma} \\ &= \{R_{xx}(I_{p-1} + K R_{xx}^{-1})\}^{-1} R_{xx} \underline{\gamma} \\ &= (I_{p-1} + K R_{xx}^{-1})^{-1} \hat{\underline{\gamma}} \\ &= C_K \hat{\underline{\gamma}} \end{aligned}$$

So, the ridge estimate is clearly biased as $C_K \neq I_{p-1}$.

$$\begin{aligned} \text{MSE} &= E[\|\hat{\underline{\gamma}}(K) - \underline{\gamma}\|^2] = E[\|C_K \hat{\underline{\gamma}} - \underline{\gamma}\|^2] = E[\|C_K \hat{\underline{\gamma}} - C_K \underline{\gamma} + C_K \underline{\gamma} - \underline{\gamma}\|^2] \\ &= E[\|C_K(\hat{\underline{\gamma}} - \underline{\gamma})\|^2] + \end{aligned}$$

$$= E \left[\| c_k (\hat{\gamma} - \gamma) + (c_k - I_{p-1}) \hat{\gamma} \|^2 \right]$$

$$= E \left[\| c_k (\hat{\gamma} - \gamma) \|^2 \right] + \| (c_k - I_{p-1}) \hat{\gamma} \|^2$$

[cross product is 0 and $\hat{\gamma} = \gamma$]

$\hat{\gamma} - \gamma$ has mean zero and variance $\sigma^2 R_{xx}^{-1}$

$$I = E \left[\| c_k (\hat{\gamma} - \gamma) \|^2 \right] = E [(\hat{\gamma} - \gamma)' c_k' c_k (\hat{\gamma} - \gamma)]$$

$$= \sigma^2 \text{tr}(R_{xx}^{-1} c_k' c_k)$$

Spectral decomposition of R_{xx} , $R_{xx} = T \Lambda T'$

$$\Lambda = \text{diag}\{\lambda_1, \dots, \lambda_{p-1}\}$$

$$= \sigma^2 \text{tr}[T \Lambda^{-1} T']$$

$$= \sigma^2 \text{tr}[T \Lambda^{-1} T' T D T' T D T']$$

$$= \sigma^2 \text{tr}[T \Lambda^{-1} D^2 T']$$

$$= \sigma^2 \text{tr}[\Lambda^{-1} D^2 T' T]$$

$$= \sigma^2 \text{tr}[\Lambda^{-1} D^2]$$

$$= \sigma^2 \sum_{j=1}^{p-1} \frac{\lambda_j}{(k+\lambda_j)^2}$$

$$\text{Take } \alpha = T' \gamma$$

$$II = \gamma' (c_k - I_{p-1})' (c_k - I_{p-1}) \gamma$$

$$= \alpha' T^{-1} [T D T' - I_{p-1}]' [T D T' - I_{p-1}] \alpha \cdot (T')^{-1}$$

$$= \frac{1}{p-1} (\alpha' (D - I_{p-1})' (D - I_{p-1}) \alpha)$$

$$= \sum_{j=1}^{p-1} \alpha_j^2 \frac{k^2}{(k+\lambda_j)^2}$$

$$\text{So, } \text{MSE} = \sum_{j=1}^{p-1} \frac{\alpha_j^2 k^2 + \sigma^2 \lambda_j}{(k + \lambda_j)^2}$$

The derivative w.r.t. k

$$\sum_{j=1}^{p-1} \frac{2\lambda_j (\alpha_j^2 k - \sigma^2)}{(k + \lambda_j)^3} \rightarrow \text{negative for small values of } h,$$

values of h , so that for k sufficiently small, MSE decreases as k increases.

$$\text{Solve } \sum_{j=1}^{p-1} \frac{2\lambda_j (\alpha_j k - \sigma^2)}{(k + \lambda_j)^3} = 0$$

To find the value of k corresponding to the smallest MSE.

The minimising value depends on the unknowns α , σ^2 .

$$\hookrightarrow \text{use } \hat{\alpha} = T' \hat{\gamma} \text{ & } s^2$$

There is no guarantee that the result in \hat{k} will lead to the smaller MSE than LSE.

But generally, it works well.

If centered & scaled exp variables are orthogonal, then $R_{xx} = I_{p-1}$ and eigen values are all 1.

$$\text{Then } 2 \sum \frac{\hat{\alpha}_j^2 k - s^2}{(k+1)^3} = 0$$

$$\Rightarrow \hat{k} = \frac{(p-1)s^2}{\sum_{j=1}^{p-1} \hat{\alpha}_j^2} = \frac{(p-1)s^2}{\|\hat{\gamma}\|^2}$$

It gives good results in non-orthogonal regression.

- much better than LSE.

Dealing with Curvature.

conditional mean of the response y given the exp variable $\rightarrow E(y|x)$

We want to visualize the surface whether it can be adequately represented by a mean of x .

- Single x variable \rightarrow a plot y v/s x .
 \rightarrow smooth the plot. i.e. no kinks.

If the realisation is linear \rightarrow fit a linear model.

If not, try to transform y using a power transformation

- or try to use a polynomial model.

Two Exp Variable:

Plot y v/s $x = (x_1, x_2)$ using a 3d plot.

- using dynamic rotation, the relationship may be revealed.

More than two expl variables:

- not possible to visualize the surface directly.

A useful plot: residuals v/s fitted values.

\hookrightarrow are independent & under normality assumption.

- If the linear model is correct, then the plot should ~~not~~ display a horizontal pattern. band of pts.

- A curved regression surface (reveals) as a curved plot.

- Smoothing will enhance the interpretation.

Disadvantage:

The nature of the curvature is not revealed.

More seriously it is possible for the regression surface to be curved without the $e_i = y_i - \hat{y}_i$ plots revealing the curvature.

Partial residual Plot (PRP)

- Plots of suitably modified residuals v/s the exp variables

Suppose the true reg surface is

$$E(y|x) = \beta_0 + \beta_1 g(x_1) + \beta_2' x_2 \text{, i.e. } \textcircled{1} \quad x = (x_1, x_2)' \text{ and } g \text{ is unknown.}$$

$$E(y|x) = \beta_0 + \beta_1 x_1 + \beta_2' x_2$$

PRP are designed to reveal the form of g .

Suppose we fit $\beta_0 + \beta_1 x_1 + \beta_2' x_2$

when $\textcircled{*}$ is the true model.

$$Y = \beta_0 + \beta_1 x_1 + \beta_2' x_2 + [\beta_1 g(x_1) - \beta_1 x_1 + e]$$

so the errors are $e = \beta_1 [g(x_1) - x_1] + e'$ where e' is the error point.

Expect that the residual e from the linear fit would approximate $\beta_1 [g(x_1) - x_1] \approx \hat{\beta}_1 [g(x_1) - x_1]$

A plot of $e_i = e_i + \hat{\beta}_1 x_{1i} \approx \hat{\beta}_1 g(x_{1i})$ v/s x_{1i}

might reveal the shape of g .

e_i^* : partial residuals.

Features:

A LS fit through the origin of this will be $\hat{\beta}_3$.

The residuals from this fit are original residuals

$$e_i^* = y_i - \hat{\beta}_3 x_{ii} = y_i + \hat{\beta}_3 x_{ii} - \hat{\beta}_3 x_{ii} = e_i$$

$$(1.89) \quad \text{with zero bias}$$

matrix of the last two observations relative to $\hat{\beta}_3$ follows

so $\hat{\beta}_3$ is the LS estimate for y and x_3 configuration.

Drawbacks: ① The appearance of the plot may over emphasize the importance of the exp variable x_3 in the fit.

② Imp: If g is highly non-linear then the LSEs of β_2 are not very good estimates & this argument breaks down.

Example: $\beta_3 = 3$, scaled and centered x

$$\hat{\gamma}_1 = \frac{1}{1-n^2} \left[\underline{x}^{*(1)'} \underline{y} - r \underline{x}^{*(2)'} \underline{y} \right]$$

$$\underline{x}^{*(1)} = (x_{11}^*, x_{21}^*, \dots, x_{n1}^*)$$

$$r = \underline{x}^{*(1)'} \underline{x}^{*(2)}$$

$$H = P = \frac{1}{n} \ln \ln' + \underline{x}^{*(2)'} \underline{x}^{*(2)'} + \frac{1}{1-n^2} (\underline{x}^{*(1)} - r \underline{x}^{*(2)})$$

Now suppose that $y = \alpha \ln + \gamma_1 g + \gamma_2 x^{*(2)} + \epsilon$
 when $x_5 = (\ln, x^{*(1)}, x^{*(2)})$; $g = (g(x_{11}), \dots, g(x_{nn}))$

So, for $(I_n - H)x_5 = 0$

$$E[\underline{\epsilon}] = E[(I_n - H)\underline{y}] = (I_n - H)(\alpha \ln + \gamma_1 g + \gamma_2 x^{*(2)})$$

$$E(\hat{\gamma}_1) = \frac{1}{1-\gamma^2} [x^{*(1)} - \gamma x^{*(2)}] = E(y)$$

$$= \frac{1}{1-\gamma^2} [x^{*(1)} - \gamma x^{*(2)}] / (\alpha \ln + \gamma_1 g + \gamma_2 x^{*(2)})$$

$$= \frac{1}{1-\gamma^2} \gamma_1 (x^{*(1)} - \gamma x^{*(2)}) g' \quad [x^{*(j)}' \ln = 0; j=1,2]$$

$$\therefore E[e_i^*] = E[\underline{\epsilon}] + x^{*(1)} E[\hat{\gamma}_1] \quad [\text{since } x^{*(1)'} x^{*(2)} = 0]$$

$$= \gamma_1 (I_n - H) g + \frac{\gamma \gamma_1}{1-\gamma^2} x^{*(1)} \cdot [x^{*(1)} - \gamma x^{*(2)}] g$$

$$= (\bar{g} - \gamma_1 (\bar{g} - \bar{g} \ln)) - \frac{1}{1-\gamma^2} \gamma_1 x^{*(2)} (x^{*(2)} - \gamma x^{*(1)}) g$$

→ indicates that the plot does not reveal the shape of g , but the shape is contaminated by additional terms.

- The plot will suffer a large amount of contamination if the correlation between columns is larger.

Adding and dealing variables:

Added variable plot: Assume that a constant term has been included in the model.

Suppose that there is a reasonable fit line

relationship in the PRP for x_j^o .

The strength of the relationship is measured by the corrⁿ between x_{ij} and e_{ij}^* = $\rho_{ij} + \beta x_{ij}$

$$\tilde{\rho}_{ij} = \frac{\hat{\beta}_j \sum_{i=1}^n (x_{ij} - \bar{x}_j)^2}{\text{RSS} + \hat{\beta}_j^2 \sum_{i=1}^n (x_{ij} - \bar{x}_j)^2} \quad [\text{Show it}]$$

$$VIF_j = \sigma^2 \text{Var}(\hat{\beta}_j) \sum_{i=1}^n (x_{ij} - \bar{x}_j)^2$$

VIF_j will be large if x_j^o can be predicted accurately from the other explanatory variables.

If F_j is the f-statistic for testing $H_0: \beta_j = 0$,

$$\text{then } F_j = \frac{\hat{\beta}_j^2}{S^2 (x'x)_{j+1, j+1}} = \frac{\hat{\beta}_j^2 (x_{ij} - \bar{x}_j)^2}{S^2 VIF_j}$$

$$\text{Hence, } \tilde{\rho}_{ij}^2 = \frac{F_j \cdot VIF_j}{(n-p) + F_j \cdot VIF_j} \quad ; \quad \text{RSS} = (n-p) S^2$$

can get a better picture if the variable x_j to the regression should be added to the regression x_j' and the residuals from a fit excluding the variables x_j .

If $x = (\ln, \underline{x}^{(1)}, \dots, \underline{x}^{(p-1)})$ & $\underline{x}^{(j)}$ as x with $\underline{x}^{(j)}$ omitted.

Then consider all cases of adding either x_{ij} or x_j .

The test for $\underline{x}^{(j)} = (I_n - P_j) y$ or when $P_j = X^{(j)}(X^{(j)'}X^{(j)})^{-1}X^{(j)'}$

is $R_j = \frac{F_j}{F_0}$.

To add x_j we must first remove x_j , omitted.

- When considering whether to add x_j to the regression, we must assess how well we can predict the residuals $\underline{e}^{(j)}$ by using x_j .

Since x_j can be partly predicted by the other exp.

we are interested in how the residuals $\underline{e}^{(j)}$

can be predicted by the pt of x_j that is not predicted by the other exp variable.

The suggest at the relationship between

$\underline{e}^{(j)}$ and x_j residuals $(I_n - P_j) \underline{x}^{(j)}$

The plot of $\underline{e}^{(j)} = (I_n - P_j) y$ v/s $(I_n - P_j) \underline{x}^{(j)}$

→ A showing linear relationship $\underline{e}^{(j)} \rightarrow x^{(j)}$

indicating that the variable be included in the regression.

→ Added variable plot.

$P_j^2 \rightarrow$ Squared correlation between quantities in the added variable plot

$$P_j^2 = \frac{F_j}{n-p+f_j}; F_j \rightarrow \text{Statistic for testing } H_0: \beta_j = 0$$

Model Selection:

To decide which variables to include in the model & how to construct a predictor from the variables available

Assume that a large no. of potential exp var are available.

- Sometime, we assume that the true model is one of the form i.e. $y = \beta_0 + \beta_1 x_1 + \dots + \beta_K x_K + \epsilon$ — (1)
- The problem is to identify the variables (that are not) related to the response i.e. identifying the β 's that are zero.

Two approaches

(I) All possible Regression (APR) (Requires 2^k reg models)

- Define a criterion for goodness.

- Compute the criterion for each possible subset of variables.

Choose the subset which optimizes the criterion.

The criterion can be based on

- Standard goodness of fit criterion.

Or estimating the prediction error.

- On estimating some measure of distance between the could based on subset & the true model.

(II) To apply a seq of hypothesis tests to the problem and attempt to identify the non-zero's Betas in (1).

- forward selection, backward elimination, stepwise regression.
- Assumption: that ① is the true model
- Computationally less demanding, but there is no guarantee that model found will be optimal in terms of criteria in ①
- alternative to subject selection & LS is to use all exp variable & use a based estimation method that "shinks" the coefficients towards zero.

Prediction Error

- An initial data set $(y_i, x_i), i=1, 2, \dots, n$ containing of $n(p+1)$ dim multivariable observations.
- Assume that the first element of each x_i is 1, corresponding to the constant term.
- Training set: we use this set of observation to predict responses $y_{oi}, i=1, 2, \dots, m$, corresponding to m new vector $x_{oi}; i=1, 2, \dots, m$.
- $\underline{Y} = (y_1, \dots, y_n)'; \underline{y}_o = (y_{o1}, \dots, y_{om})'$
- Assume that $V(\underline{y}) = \sigma^2 I_n$; $V(\underline{y}_o) = \sigma^2 I_m$, $\underline{y}_o \perp \underline{x}$ are independent.
- $\underline{x} = \begin{pmatrix} x_1' \\ \vdots \\ x_n' \end{pmatrix} = \underline{x}_o, \underline{x}_o = \begin{pmatrix} x_{o1}' \\ \vdots \\ x_{om}' \end{pmatrix}, \underline{x} = \underline{x}_o$
- Calculate $\hat{\beta} = (\underline{x}'\underline{x})^{-1}\underline{x}'\underline{y}$ from the training set.
- The ISE of \underline{y}_o (also the estimate of $\mu_o = E(\underline{y}_o)$) is $\underline{x}_o \hat{\beta}$.

The sum of squared prediction error is

$$\sum_{i=1}^m (y_{oi} - \underline{x}_{oi}' \hat{\beta})^2 = \| \underline{y} - \underline{x}_o \hat{\beta} \|^2$$

$$= \| \underline{y}_o - \underline{\mu}_o + \underline{\mu}_o - \underline{x}_o \hat{\beta} \|^2$$

$$= \| \underline{y}_o - \underline{\mu}_o \|^2 + \| \underline{\mu}_o - \underline{x}_o \hat{\beta} \|^2 + 2 (\underline{y}_o - \underline{\mu}_o)' (\underline{\mu}_o - \underline{x}_o \hat{\beta})$$

Take expectation over the new data only and call it the prediction error (PE).

$$PE = E_{\underline{y}_o} \| \underline{y}_o - \underline{x}_o \hat{\beta} \|^2$$

$$= E_{\underline{y}_o} \| \underline{y}_o - \underline{\mu}_o \|^2 + \| \underline{\mu}_o - \underline{x}_o \hat{\beta} \|^2$$

$$E_{\underline{y}_o} [(\underline{y}_o - \underline{\mu}_o)' (\underline{\mu}_o - \underline{x}_o \hat{\beta})] = 0$$

$$\textcircled{*} \Rightarrow E_{\underline{y}_o} \left[\sum_{i=1}^m (y_{oi} - \underline{\mu}_{oi})^2 \right] + \| \underline{\mu}_o - \underline{x}_o \hat{\beta} \|^2$$

$$= m\sigma^2 + \| \underline{\mu}_o - \underline{x}_o \hat{\beta} \|^2$$

$$PE = m\sigma^2 + \| \underline{\mu}_o - \underline{x}_o \hat{\beta} \|^2$$

\downarrow measures how well the linear model reflects the underlying variability of the data.

$\underline{\mu}_o = \underline{x}_o' \hat{\beta}$ represents the mean response of the new data given the new predictors \underline{x}_o .

$\| \underline{\mu}_o - \underline{x}_o \hat{\beta} \|^2 \rightarrow$ model error (ME)

If $\underline{x}_o = \underline{x}$, then $m=n$ & $\underline{\mu}_o = E[\underline{y}] = \underline{\mu}$,

Then $\underline{\epsilon} = \underline{y} - \underline{\mu}$; $\underline{P} = \underline{x}(\underline{x}'\underline{x})^{-1}\underline{x}'$, say

$$ME = \| \underline{\mu} - \underline{x} \hat{\beta} \|^2 = \| \underline{\mu} - \underline{P} \underline{y} \|^2$$

$$= \| \underline{\mu} - \underline{I}(\underline{\mu} + \underline{\epsilon}) \|^2 = \| (\underline{I}_{n-p}) \underline{\mu} - \underline{P} \underline{\epsilon} \|^2$$

$$= \|(\mathbf{I}_n - \mathbf{P})\underline{\mu}\|^2 + \|\mathbf{P}(\mathbf{E}|_{\mathcal{X}})\| = \underline{\mu}'(\mathbf{I}_n - \mathbf{P})\underline{\mu} + \underline{\epsilon}'\mathbf{P}\underline{\epsilon}$$

[Cross product term consider as $\mathbf{P}^2 = \mathbf{P}$]

$$\begin{aligned} E(ME) &= \underline{\mu}'(\mathbf{I} - \mathbf{P})\underline{\mu} + \sigma^2 \text{tr}(\mathbf{P}) \\ &= \underline{\mu}'(\mathbf{I} - \mathbf{P})\underline{\mu} + \sigma^2 p. \end{aligned}$$

The corresponding

$$E(PE) = E[n\sigma^2 + ME] = (n+p)\sigma^2 + \underline{\mu}'(\mathbf{I}_n - \mathbf{P})\underline{\mu}$$

Define total bias & total variance of the prediction

$$X(\hat{\beta}) = (\mathbf{I} - \mathbf{P})\underline{\mu} + \underline{\epsilon}$$

$$\text{Total Bias} = \|\underline{\mu} - E(X(\hat{\beta}))\|$$

$$\text{Total Variance} = \text{tr}(V(X(\hat{\beta}))) = \sigma^2 p$$

$$\therefore E(ME) = \underline{\mu}'(\mathbf{I} - \mathbf{P})\underline{\mu} + \sigma^2 p$$

total variance

sq. of the total
bias which
variables of $\underline{\mu} \in \mathcal{C}(\mathbf{x})$

- As new variables are added to the model, the variability increases, but the bias decreases; unless the new variables are linearly dependent on the old variables, already included in the model.
- The best model, having the smallest expected ME (& hence the smallest expected PE) will be some compromise between these two conflicting requirements of small bias and small variability.

Suppose $\gamma = \beta_0 + \beta_1 x_1 + \dots + \beta_k x_k + \epsilon$

is exactly the true model, although some of the reg coeff may be equal to zero.

$$Y = X\beta + \epsilon, \quad X = (x_1, x_2), \quad \beta = (\beta_1, \beta_2)$$

$n \times p \quad k-p+1 \quad (k+1)x_1 \quad px_1$

① Including redundant variables does have a cost.

- the predictor using all variables will have a large variance than that based on x_1 .

② On the other hand if $\beta_2 \neq 0$ we use x_1 , the prediction is biased.

Suppose to predict $\underline{x}'\beta = \underline{x} = (x'_1, x'_2)'$ is $(k+1)$ vector

- Should we use $\underline{x}'\tilde{\beta}_1$ based on x_1 or unbiased predictor $\underline{x}'\hat{\beta}$ where $\hat{\beta}_1$ is the LSE of β based on X using all the exp. variable with $m=1$, the expected PE of the based predictor.

$$E(PE) = \sigma^2 + E[(\underline{x}'\tilde{\beta}_1 - \underline{x}'\beta)^2]$$

$$= \sigma^2 + \text{Var}(\underline{x}'\tilde{\beta}_1) + (\underline{x}' E(\tilde{\beta}_1) - \underline{x}'\beta)^2$$

$$\underline{x}' E(\tilde{\beta}_1) - \underline{x}'\beta = E(\tilde{\beta}_1) = \beta_1 + L\beta_2$$

$$= \underline{x}'\beta_1 + \underline{x}'L\beta_2 - \underline{x}'\beta_1 - \underline{x}'\beta_2$$

$$= (L'x_1 - x_2)' \beta_2$$

$$E(PE) = \sigma^2 + \text{Var}(\underline{x}'\tilde{\beta}_1) + E[(L'x_1 - x_2)' \beta_2]^2$$

→ for biased predictor.

For the unbiased predictors, when not 2nd diff matrix

$$\hat{E}(PE) = \sigma^2 + E(\underline{x}'\hat{\beta} - \underline{x}'\underline{\beta})^2$$

$$(\text{PE})^2 = \sigma^2 + \text{Var}(\underline{x}'\hat{\beta})$$

$$\text{PE}^2 = \sigma^2 + \text{Var}(\underline{x}'\hat{\beta}_1) + \sigma^2(L'\underline{x}_1 - \underline{x}_2)' M (L'\underline{x}_1 - \underline{x}_2)$$

$$\text{When } M = (\underline{x}_2' (I_n - P_1) \underline{x}_2)^{-1}; P_1 = \underline{x}_1 (\underline{x}_1' \underline{x}_1)^{-1} \underline{x}_1'$$

$$\text{write } L'\underline{x}_1 - \underline{x}_2 = h \text{ for } h = \frac{\underline{x}_1 - \underline{x}_2}{\|\underline{x}_1 - \underline{x}_2\|}$$

Comparing two PE's, the biased predictor $\underline{x}_1'\hat{\beta}_1$ has a smaller expected PE than the unbiased predictor if

$$(h'\hat{\beta}_2)^2 < \sigma^2 h' \mu h \quad \text{--- } \otimes \otimes$$

Since μ is p.d. we have \otimes

$$(h'\hat{\beta}_2)^2 < (h'\mu h) \hat{\beta}_2' \mu' \hat{\beta}_2 \quad \text{for any } b$$

$$\text{maximum of ratios of vectors } h \text{ w.r.t. } b \quad \max_{h, b \neq 0} \frac{(h'b)^2}{h'b} = b'L^{-1}b$$

$$\text{So if } \hat{\beta}_2' \mu' \hat{\beta}_2 = \hat{\beta}_2' X_2' (I_n - P_1) X_2 \hat{\beta}_2$$

then $\otimes \otimes$ is true for all h , \otimes

- So the biased predictor will have the smaller expected PE.

Choosing the Best subset

1. Goodness of fit - criterion

RSS is a measure of goodness of fit, but it is not a good absolute measure, since we have two sets of variables S_1 and S_2 , $S_1 \subset S_2$

Then the RSS for model based on S_2 will be smaller than for the model based on S_1 . (provided the variables are not known combination of those in S_2 & S_1)

- A way of correcting the RSS to accurate for the no of residuals is to use the estimated residual variance.

$$S^2 = \frac{RSS}{n-p}, \quad p \text{ is no of columns in } X.$$

- Another criterion is adjusted R^2

$$R^2_{adj} = 1 - \frac{n-1}{n-p} (1-R^2)$$

- Consider the model having the largest R^2

The use of R^2 is motivated by testing the accuracy of the model as the "full model" consider of all available exp variables.

- Suppose that there are K exp variables and would under consideration was $p-1 < K$ variable.

- R_p^2 & R_{K+1}^2 : coeff of determination

- F test for testing the ~~correct~~ accuracy of the small model as the larger model is

$$H_0: \beta_p = \beta_{p+1} = \dots = \beta_K = 0$$

The f-statistic.

$$F_p = \frac{\frac{R_{K+1}^2 - R_p^2}{1 - R_{K+1}^2}}{\frac{n-K-1}{K+1-p}} \cdot \frac{(n-K-1)}{(K+1-p)}$$

$$F = \frac{(RSS_{K+1} - RSS)/q}{RSS/(n-p)}$$

$$RSS = (1-R^2) \sum (y_i - \bar{y})^2$$

$$\text{So, } 1 - R_p^2 = (1 - R_{K+1}^2) \frac{(K+1-p) F_p + (n-K-1)}{n-K-1}$$

$$\text{and } \bar{R}_p^2 = (1 - (1 - R_p^2) \frac{n-1}{n-p}) \cdot S_0^2 (1-p) + (1-p)^2 =$$

$$= 1 - (1 - R_{K+1}^2) \frac{(K+1-p)F_p + (n-K-1)}{(n-K-1)} \cdot \frac{n-1}{n-p}$$

$$\text{If } F_p \gg 1 ; \frac{(K+1-p)F_p + (n-K-1)}{n-K-1} \gg \frac{K+1-p + n-K-1}{n-K-1}$$

$$\text{So, } \bar{R}_p^2 \leq 1 - (1 - R_{K+1}^2) \frac{n-1}{n-K-1}$$

which indicates \bar{R}_p^2 is minimized

large values of F_p are evidence in favour of K variables model. This motivates us to use adjusted R^2 as an model selection criterion.

Note:

$$\bar{R}_p^2 = 1 - (1 - R_p^2) \frac{n-1}{n-p} = 1 - (1 - R_p^2) \frac{(n-1)S_p^2}{(n-p)S_p^2} = \frac{(n-1)S_p^2}{(n-p)S_p^2}$$

$$; S_p^2 = \frac{\text{RSS}_p}{n-p} = \frac{\sum (y_i - \bar{y})^2}{n-p} \approx (p)^2$$

$$\bar{R}_p^2 = 1 - \frac{(n-1)S_p^2}{\sum (y_i - \bar{y})^2} = 1 - \frac{(n-1)S_p^2}{\text{RSS}_p}$$

So, the model is $\max^m \bar{R}^2$ is also the model with minimum S^2 .

Model Selection Criterion based on Prediction Error

To get an operational criterion

$$ME = \|\mu_0 - X_0 \hat{\beta}\|^2$$

To get an operational criterion we must estimate the expected ME.

Mallow's C_p

RSS_p : RSS from fitting a model with p -parameters by $n \times p$ regression matrix X p .

If $I_p = X_p (X_p' X_p)^{-1} X_p'$, then

$$E[\text{RSS}_p] = E[y' (I - I_p)' y] = \mu' (I - I_p) \mu + (n-p)\sigma^2$$

$$= E(ME) + (n-2p)\sigma^2$$

$$= \mu' (I-P)\mu + \sigma^2 p.$$

$$\frac{E(ME)}{\sigma^2} = \frac{E(RSS_p)}{\sigma^2} + 2p - n$$

If we had an estimate $\hat{\sigma}^2$, then we could use

it & $c_p = \frac{RSS_p}{\hat{\sigma}^2} + 2p - n \rightarrow$ as an estimate of $E(ME)/\sigma^2$, the scaled expected model error.

- If the model fits model in the sense that μ is well-approximated by vectors in $\mathcal{E}(X_p)$, then

$$\mu' (I_n - P_p) \mu = \| (I_n - P_p) \mu \|^2 \text{ will be small.}$$

$$\begin{aligned} E(c_p) &\approx \frac{E[RSS_p]}{\sigma^2} + 2p - n \\ &= \frac{\mu' (I_n - P_p) \mu}{\sigma^2} + \frac{(n-p)\sigma^2}{\sigma^2} + 2p - n \\ &\approx p \end{aligned}$$

So, the expected c_p will be close to p .

Mallow's suggested using C_p -plot, a plot of c_p v/s p for all possible models.

→ consider model for which $c_p \approx p$ (i.e. the model fits well)

- It is a common practise to use all K variables to estimate σ^2 .

$$\hat{\sigma}^2 = \frac{RSS_{K+1}}{n-K-1}$$

$$\Rightarrow C_{k+1} = (k+1)$$

$$\frac{RSS_p}{RSS_{k+1}} = \frac{n-k-1}{n-p} = \frac{1 - \bar{R}_p^2}{1 - \bar{R}_{k+1}^2}$$

If n is large compare to p , the smallest value of $C_p - p$ corresponds approximately to the largest \bar{R}_p^2 . (Show this)

Cross-Validation (CV)

- $(x_{oi}, y_{oi}) ; i=1, 2, \dots, m \rightarrow$ a new data following the same model as the training set data $(x_i, y_i), i=1, 2, \dots, n$

- An obvious measure of PE

$$\frac{1}{n} \sum_{i=1}^n [y_{oi} - \underline{x}_{oi}' \hat{\beta}]^2$$

$\hat{\beta}$: estimated using training set data.

- In practice, we don't have additional data
- Sometimes we can divide into two disjoint set of obs. → one for training set & the other for prediction.
- What we do →
 - Select a subset D of d observations.
 - estimate the prediction using $(n-d)$ obs.
 - Then calculate $[PE(D)]$

Repeat this process for selected d subsets D_1, D_2, \dots and average the resulting estimate $PE(D_i)$.

↳ This process is called cross-validation.

The most popular choice of d is 1. (This involves leaving out one obs in turn and calculating the estimate

from the remaining $(n-1)$ obs.).

The error in predicting the i th obs. as $y_i - \underline{x}_i' \hat{\beta}(i)$.

→ heads, the leave one-out or $cv(i)$ prediction error estimating.

$$cv(i) = \frac{1}{n} \sum_{i=1}^n (y_i - \underline{x}_i' \hat{\beta}(i))^2$$

$$y_i - \underline{x}_i' \hat{\beta}(i) = y_i - \underline{x}_i' [\hat{\beta} - (\underline{x}' \underline{x})^{-1} \underline{x}_i (y_i - \underline{x}_i' \hat{\beta})]$$

$$= y_i - \underline{x}_i' \hat{\beta} + \frac{h_i(y_i - \underline{x}_i' \hat{\beta})}{1-h_i}$$

$$= \frac{y_i - \underline{x}_i' \hat{\beta}}{1-h_i} [1 - h_i]$$

$$cv(i) = \frac{1}{n} \sum_{i=1}^n \frac{(y_i - \underline{x}_i' \hat{\beta})^2}{(1-h_i)^2}$$

Expected value of $cv(i)$:

$$\text{write } y_i = \{(I_n - P)\mu\}_i + \eta_i \quad \& \quad \epsilon = \underline{\epsilon} - \mu$$

$$y_i - \underline{x}_i' \hat{\beta} = \{y - x\hat{\beta}\}_i$$

$$= \{(I_n - P)y\}_i + \{(I_n - P)\epsilon\}_i$$

$$= \{(I_n - P)\mu\}_i + \{(I_n - P)\epsilon\}_i$$

$$= \eta_i + \{(I_n - P)\epsilon\}_i$$

$$E[(y_i - \underline{x}_i' \hat{\beta})^2] = E[(\eta_i + \{(I_n - P)\epsilon\}_i)^2]$$

$$= \eta_i^2 + E[\{(I_n - P)\epsilon\}_i^2]$$

Let, D_i : be a diagonal matrix whose i th diagonal element is 1 and rest are zero.

$$\begin{aligned} \left\{ (I_n - P) \in \right\}_i^2 &= \in' (I_n - P) D_i (I_n - P) \in \\ E \left[\left\{ (I_n - P) \in \right\}_i^2 \right] &= E \left[\in' (I_n - P) D_i (I_n - P) \in \right] \\ &= 0 + \sigma^2 \xrightarrow{\text{tr}} \text{tr} ((I_n - P) D_i (I_n - P)) \\ &= \sigma^2 \text{tr}(D_i (I_n - P)) \\ &= \sigma^2 (1 - h_i) \end{aligned}$$

~~$E \left[\left\{ (I_n - P) \in \right\}_i^2 \right]$~~

$$E[(y_i - x_i' \hat{\beta})^2] = \eta_i^2 + \sigma^2(1 - h_i)$$

$$E[nCV(i)] = \sum_{i=1}^n \left[\frac{\eta_i^2 + \sigma^2(1 - h_i)}{(1 - h_i)^2} \right]$$

The expected PE

$$\begin{aligned} E(PE) &= (n+p) \sigma^2 + \mu' (I - P) \mu \\ &= (n+p) \sigma^2 + \| \sigma (I_n - P) \mu \|^2 \\ &= (n+p) \sigma^2 + \sum_{i=1}^n \eta_i^2 \\ &= \sum_{i=1}^n [\sigma^2 (1 + h_i) + \eta_i^2] \end{aligned}$$

$$\begin{aligned} E[nCV(i) - PE] &= \sum_{i=1}^n \left[\frac{\eta_i^2 + (1 - h_i) \sigma^2}{(1 - h_i)^2} - \frac{\sigma^2 (1 + h_i) + \eta_i^2}{(1 - h_i)^2} \right] \\ &= \sum_{i=1}^n \eta_i^2 \frac{h_i(2 - h_i)}{(1 - h_i)^2} + \sigma^2 \sum_{i=1}^n \frac{h_i^2}{1 - h_i} \end{aligned}$$

Both the terms are positive $\Rightarrow nCV(i)$ tends to overestimate the PE.

$CV(d) = \sqrt{\frac{1}{d} \sum_{i=1}^d (x_i - \bar{x})^2}$

AIC \rightarrow A Information Criterion.

$$BIC \rightarrow \frac{3(1-\alpha)}{2} + \frac{3}{2}(1-\alpha) = \frac{3}{2}(1-\alpha)(\alpha+1)$$

$$CV \text{ for } d \geq 1 \left(\frac{d}{d+1} \right) \leq \left(\frac{d+1}{2} \right)^{\frac{1}{d}} = \left[\frac{d+1}{2} \right]^{\frac{1}{d+1}}$$

Suppose $D = \{i_1, \dots, i_g\}$ are the cases to be omitted,

(1-11) at 46th St. N.Y.C. T

$$\det \underline{Y}_D = \left\{ y_{ii}, \dots, \underline{\{y_{id}\}} \right\} \in \mathbb{R}^d \text{ and } \underline{Y}_D =$$

$\hat{\beta}_2$ = estimate of β_2 based on the observations not in D

$$\text{and } H_D = x_D (x'x)^{-1} x_D'$$

Since we are predicting the cases in D simultaneously and comparing them with \hat{Y}_D .

$$\underline{y}_D - x_D \hat{\beta}(D) = (I - H_D)^{-1} (\underline{y}_D - x_D \hat{\beta})$$

$$S_0, \quad CV(d) = \left(\frac{n}{d}\right)^{-1} \sum_D (y_D - x_D \hat{\beta})' (I - H_D)^2 (y_D - x_D \hat{\beta})$$

where Sum is taken over all d sub sets of observation.

- For even $d = 2$ or 3 , it is computationally very

Estimating Distributional Discrepancy: AIC

This criterion is based on a discrepancy between the true distⁿ of the data \underline{Y} and the distⁿ specified by the model under consideration.

$f(\underline{x})$: density of the true distn.

$\cdot f(y)$; " " specified model; named and noted

Discrepancy measure in the Kullback-Leibler discrepancy.

$$KL(f, g) = \int \log \frac{f(y)}{g(y)} f(y) dy = \int \log f(y) f(y) dy - \int \log g(y) f(y) dy$$

Note: This is not a true distance measure in the sense

$$KL(f, g) \neq KL(g, f)$$

$$KL(f, g) \geq KL(f, f) \Rightarrow 0$$

In practice, we are interested in a family of models $g(y; \theta)$ where θ ranges over some parameter space

→ use $KL(f, g; \theta)$ to measure fit in general

- The true model may or may not be of the form of $g(y; \theta)$ for some particular θ .
- We want to choose the model corresponding to the value of θ that minimizes $KL(f, g; \theta)$, or equivalently, that

$$-\int \log g(y; \theta) f(y) dy = -E[\log(g(y; \theta))]$$

This is equal to $KL(f, g, \theta)$ upto a constant

$$\int \log f(y) f(y) dy \rightarrow \text{inpt of } g.$$

- ① depends on two unknowns, the parameter θ & the true distⁿ. To estimate θ

Suppose we have a sample Y and we estimate θ using ML estimate

- Let, $\hat{\theta}(y)$ be the MLE that maximizes $g(y, \theta)$ as a function of θ .

→ leads to the modified criterion.

$$\Delta = - \int \log(s(x); \hat{\theta}(y)) f(x) dx$$

Then $E(\Delta) = - \iint \log(s(x); \hat{\theta}(y)) f(x) f(y) dx dy$

→ Still f is unknown.

To get operational criterion we need to estimate Δ .

- The standard estimate of Δ [basically 2Δ] is the Akaike Information Criterion (AIC)

$$AIC = -2 \log(s(\cdot; \hat{\theta}(\cdot))) + 2r$$

Where r is the dimension of the pos. vector $\hat{\theta}$.

In case of linear models,

- the true model $f(y)$ may be taken as multivariate normal with mean vector μ and variance matrix $\sigma^2 I_n$.

A typical considered model $s(y; \theta)$ is also multivariate normal with mean vector $x\beta$ and variance matrix $\sigma^2 I_n$.

By considerate model 3,

$$\log s(y; \theta) = \frac{n}{2} \ln(2\pi\sigma^2) + \frac{1}{2\sigma^2} (y - x\beta)' (y - x\beta)$$

& the MLEs are ; $\hat{\beta} = (x'x)^{-1} x'y$

$$\hat{\sigma}^2 = \frac{RSS}{n}$$

lets y_0 have the same dist as y and is

inpt of y

$$\text{Then } \Delta = E_{y_0} [-\ln s(y_0; \hat{\theta}(y))]$$

$$= n/2 \ln(2\pi\hat{\sigma}^2) + \frac{1}{2\hat{\sigma}^2} E_{y_0} [(y_0 - x\hat{\beta})' (y_0 - x\hat{\beta})]$$

$$= \frac{n}{2} \ln(2\pi\hat{\sigma}^2) + \frac{1}{2\hat{\sigma}^2} [n\sigma_0^2 + \| \underline{\mu} - \underline{x}\hat{\beta} \|^2]$$

Note that, MLE's $\hat{\beta}$ and $\hat{\sigma}^2$ are indep. and

$$\text{write } \lambda = \frac{\underline{\mu}'(\mathbf{I} - \mathbf{P})\underline{\mu}}{\sigma_0^2}$$

$$\begin{aligned} E_Y[2\Delta] &= E_Y[n \log(2\pi\hat{\sigma}^2) + \frac{1}{\hat{\sigma}^2}] \{ n\sigma_0^2 + \| \underline{\mu} - \underline{x}\hat{\beta} \|^2 \} \\ &= nE[\log 2\pi\hat{\sigma}^2] + \{ n\sigma_0^2 + E(\| \underline{\mu} - \underline{x}\hat{\beta} \|^2) \} E\left[\frac{1}{\hat{\sigma}^2}\right] \\ &= nE[\log 2\pi\hat{\sigma}^2] + \{ n\sigma_0^2 + (\sigma_0^2 + \sigma^2) \} E\left(\frac{1}{\hat{\sigma}^2}\right) \\ &= nE[\log 2\pi\hat{\sigma}^2] + (n+p+\lambda) E\left[\frac{\sigma^2}{\hat{\sigma}^2}\right] \end{aligned}$$

To estimate this we use AIC $= n + p + \lambda + O(n^{-1})$

$$\text{Take } r = p+1 \quad \frac{(n+p+\lambda)}{n} = \frac{n+p+\lambda}{n} + \frac{n+p+\lambda}{n^2} O(1)$$

$$AIC = n \log(2\pi\hat{\sigma}^2) + \frac{(\underline{y} - \underline{x}\hat{\beta})(\underline{y} - \underline{x}\hat{\beta})}{n} + 2(p+1)$$

$$(1) \hat{\sigma} = \sqrt{\frac{(n+p)\hat{\sigma}^2}{n-p-1}} = \hat{\sigma} \sqrt{1 + \frac{p}{n-p-1}}$$

$$= n \log(2\pi\hat{\sigma}^2) + n + 2(p+1)$$

$$= n \log(2\pi\hat{\sigma}^2) + n + 2(p+1) \quad (3)$$

Comparing two expression (2) & (3) & using the approximation $E\left[\frac{\sigma^2}{\hat{\sigma}^2}\right] = 1 + O(n^{-1})$

So upto $O(1)$, explicit value of AIC and $E[2\Delta]$ are the same. \hookrightarrow suggests that AIC is a reasonable of $E[2\Delta]$

$$[x_n = O\left(\frac{1}{n}\right) \Rightarrow nx_n = O(1)]$$