

Regression Diagnosis

We have postulate several conditions:

i) $E(Y)$ is L.F of $(X_i)^p$ (linearity)

ii) ε_i are indep. (independence)

iii) $E(\varepsilon_i) = 0$, $\sigma^2(\varepsilon_i) = \sigma^2$. (homogeneity)

iv) ε_i is normal dist. $1 \leq i \leq n$. (normality)

\Rightarrow Whether these hypothesis are reasonable or not. If not, how can we correct it?

(1) Residual and Plot:

$$\textcircled{1} \hat{e} = Y - X\hat{\beta} = (I - H)Y, \quad H = X(X^T X)^{-1}X^T.$$

So \hat{e} is connected with $H = (h_{ij})_{n \times n}$

prop. i) $h_{ii} \in [0, 1]$. Besides, $h_{ii} = 1 \Rightarrow h_{ji} = 0, j \neq i$

ii) $\sum_i h_{ii} = p+1$ (p is number of variables)

p.f. i) By $M^2 = M \Rightarrow h_{ii} = \sum_j h_{ij}^2 \geq h_{ii}^2$

ii) $\sum h_{ii} = \text{tr}(H) = \text{tr}(X^T X)^{-1}(X^T X)$

Def: $r_j = \hat{e}_j / \sqrt{1 - h_{jj}} S$ is normalized residual, where

$$S = \sqrt{SSE / (n - p - 1)}, \quad \hat{e} = (\hat{e}_1, \dots, \hat{e}_n)$$

Prop: i) $r_i / (n - p - 1) \sim B(\frac{1}{2}, \frac{n - p - 2}{2})$

ii) $E(r_i) = 0$, $\text{Cov}(r_i, r_j) = - \frac{h_{ij}}{\sqrt{1 - h_{ii}} \sqrt{1 - h_{jj}}}$

iii) Generally, $\tilde{\epsilon}_i$ isn't indept with σ^2 .

But if the hypothesis hold, then $\{\epsilon_i\}$:

iiia $N(0,1)$, approximately.

① Plot:

We choose r_i as y -axis of plot

And x_i as x -axis of plot (or $\hat{\eta}$)

Rmk: If the model is reasonable. Then the positions of data points is random.

Since $r_i \sim N(0,1)$, i.i.d.

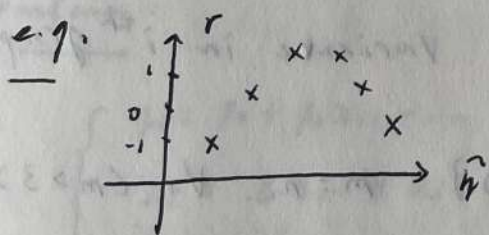
If there's some rule of dist. of points

Then we can suspect the model fails.

(2) Model Diagnosis:

① Linearity:

i) Criteria: η is L.F. of $\{x_i\}$.



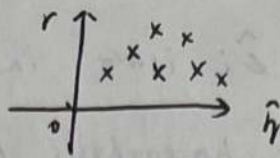
the dist. of points have a rule of quadratic function.

\Rightarrow So we will suspect it's wrong.

ii) Recorrection:

Consider: $\eta_i = \beta_0 + \beta_1 x_i + \beta_2 x_i^2 + \epsilon_i$

\Rightarrow We have plot:



② Homogeneity:

$$\text{If } \begin{cases} \eta_{ij} = \beta_0 + \beta_1 x_{ij1} + \dots + \beta_p x_{ijp} + \epsilon_{ij}, & 1 \leq i \leq k, 1 \leq j \leq n_i \\ \epsilon_{ij} \sim N(0, \sigma_i^2), & \text{indep.} \end{cases}$$

$$\text{Test: } H_0: \sigma_1^2 = \sigma_2^2 = \dots = \sigma_k^2.$$

i) Diagnose by Plot:

If it's not homogeneous. Then plot $r - \hat{\eta}$ have a rule (some trend)

ii) Diagnose by trivial:

If we can repeat the experiments. Then we can use the following three tests.

a) Martley Test:

$$F = \frac{\max_{1 \leq i \leq k} S_i^2}{\min_{1 \leq i \leq k} S_i^2}, \quad \text{where } S_i^2 = \frac{1}{n_i - 1} \sum_{j=1}^{n_i} (\eta_{ij} - \hat{\eta}_{ij})^2$$

$$\hat{\eta}_{ij} = x_{ij}^T \hat{\beta}, \quad (S_i^2 \text{ is variance in } i^{\text{th}} \text{ group})$$

$$R = \{ F > F_{1-\alpha}(k, m-1) \}, \quad m = n_i, \forall i, (m > 3).$$

b) Cochran Test:

$$G = \frac{\max_{1 \leq i \leq k} S_i^2}{\sum_{i=1}^k S_i^2}, \quad R = \{ G > G_{1-\alpha}(k, m-1) \},$$

$$m = n_i, \forall i, m > 3.$$

a) Bartlett Test:

$$\chi^2 = \frac{1}{c} \left[f_2 \ln s^2 - \sum_{i=1}^k (n_i - 1) s_i^2 \right] \xrightarrow{H_0} \chi^2_{(k-1)}.$$

Where $f_2 = \sum (n_i - 1)$, $s^2 = \sum (n_i - 1) s_i^2 / f_2$, $c = \frac{\sum \frac{1}{n_i - 1} - \frac{1}{f_2}}{3(n-1)} + 1$

$$R = \{ \chi^2 > \chi^2_{1-\alpha}(k-1) \}.$$

iii) Correction:

a) By generalised LSE:

Directly consider $\varepsilon \sim N(0, \sigma^2 \text{diag}\{\sigma_1^2/\sigma^2, \dots, \sigma_k^2/\sigma^2\})$

b) By Transformation:

If $E(Y) = m = m(x_1, \dots, x_p)$, $\text{Var}(Y) = g(m(x_1, \dots, x_p))$

Set $Z = f(Y)$. Find f st. $\text{Var}(Z) = \text{const.}$

By Taylor Expansion: $Z \approx f(m) + f'(m)(Y - m)$

$$\Rightarrow \text{Var}(Z) = f'(m)^2 g(m) = c.$$

Solve $f'(y) = \int \sqrt{c/g(m)} dm$. Consider Z .

③ Independence:

e.g.
$$\begin{cases} \eta_i = \beta_0 + \beta_1 x_{i1} + \dots + \beta_k x_{ik} + \varepsilon_i, & 1 \leq i \leq n. \\ \varepsilon_i = \rho \varepsilon_{i-1} + u_i \\ u_i \stackrel{i.i.d.}{\sim} N(0,1) \end{cases} \quad \text{(one-order case)}$$

$$\Rightarrow \text{Test: } H_0: \rho = 0 \text{ v.s. } H_1: \rho \neq 0.$$

(By Time-Series analysis)

④ Normality:

i) Diagnosis:

a) since if $Y \sim N(X\beta, \sigma^2 I_n)$. Then $r_i \sim N(0, 1)$.

Calculate the ratio of data's fallen into $(-1, 1)$, $(-2, 2)$, $(-3, 3)$.

b) Pearson - χ^2 test.

ii) Correction:

We want to find a transformation on Y .

$$\text{st. } \begin{pmatrix} y^{(1)} \\ \vdots \\ y^{(n)} \end{pmatrix} \sim N_n(X\beta, \sigma^2 I_n)$$

$$\text{Def: Box-Cox Transform: } y^{(\lambda)} = \begin{cases} (y_i^{\lambda} - 1)/\lambda, & \lambda \neq 0 \\ \ln y_i, & \lambda = 0. \end{cases}$$

\Rightarrow Find λ . st. It's most likely normal dist.

\Rightarrow Consider MLE $L(\beta, \sigma^2, \lambda) = |J| \cdot$

$$(2\pi\sigma^2)^{-\frac{n}{2}} \exp\left(-\frac{1}{2\sigma^2} (Y^{(\lambda)} - X\beta)^T (Y^{(\lambda)} - X\beta)\right)$$

$$\text{where } |J| = \prod |y_i|^{\lambda-1}.$$

$$\Rightarrow \begin{cases} \hat{\beta}_\lambda = (X^T X)^{-1} X^T Y^{(\lambda)} \\ \hat{\sigma}_\lambda^2 = \frac{1}{n} Y^{(\lambda)T} (I - M) Y^{(\lambda)} \end{cases}$$

$$\Rightarrow \text{Solve } \max_{\lambda} L(\lambda, \hat{\beta}_\lambda, \hat{\sigma}_\lambda^2) = \max_{\lambda} |J| (2\pi\sigma^2)^{-\frac{n}{2}}$$

By numerical method, find $\hat{\lambda}$.

(3) Data Diagnosis:

① Outlier:

It means that a data point departing the model a lot

$$\text{Consider } \begin{cases} \eta_i = x_i^T \beta + \varepsilon_i, & i \neq j \\ \eta_j = x_j^T \beta + \eta + \varepsilon_j. \end{cases} \quad \varepsilon_i \stackrel{i.i.d.}{\sim} N(0, \sigma^2)$$

$$\Rightarrow \text{i.e. } Y = (X \quad e_j(n)) \begin{pmatrix} \beta \\ \eta \end{pmatrix} + \varepsilon.$$

$$H_0: \eta = 0 \text{ v.s. } H_1: \eta \neq 0.$$

$$F = \frac{Y^T (M - M_0) Y / 1}{Y^T (I - M) Y / (n-p-1)} = \frac{(n-p-1) r_j^2}{n-p-1} \sim F(1, n-p-1)$$

$$R = \{ F > F_{1-\alpha}(1, n-p-1) \}.$$

② Influential Observation:

It means the data point influences the estimate statistics a lot.

Rmk: i) It can provide more information than other points.

ii) It may be an outlier or may not be.

③ Analysis:

i) Cook Distance:

$$\text{Define: } IF_i = \hat{\beta}_{(-i)} - \hat{\beta}, \quad \hat{\beta}_{(-i)} = (X_{(-i)}^T X_{(-i)})^{-1} X_{(-i)}^T Y.$$

Def: $D_i(m, c)$ is Cook distance between $\hat{\beta}_{(-i)}$ and

$$\hat{\beta} : D_i(m, c) = (\hat{\beta}_{(-i)} - \hat{\beta})^T M (\hat{\beta}_{(-i)} - \hat{\beta}) / c$$

Rmk: If $D_i(m, c) \uparrow$. Then it means i^{th} group of data is influential.

prop. $D_i(m, c) = r_i^2 \cdot \frac{S^2}{c} \cdot p_i(m)$, $p_i(m) = \frac{x_i^T (X^T X)^{-1} M (X^T X)^{-1} x_i}{1 - h_{ii}}$

Rmk: i) r_i^2 is the evaluation of fitting degree.

ii) $p_i(m)$ describes the position of data point \vec{x}_i .

iii) When $D_i(m, c)$ is large $\Rightarrow x_i$ is influential observation (Departure a lot)

When $h_{ii} \approx 1$, x_i is high leverage point.

When r_i is large, x_i is outlier.

ii) AP - Statistics:

Denote: $Z = (X \ y)$, $Z(-I)$ is matrix depleted the rows with index in I .

Def: AP-statistics: $R_I = |Z(-I)^T Z(-I)| / |Z^T Z|$.

Rmk: R_I is smaller \Rightarrow It's more likely the data depleted is influential.