

$$(1) \text{ (mean, for all)}: \hat{Y} \pm t_{n-2} \times S_{Y|X} \sqrt{\frac{1}{n} + \frac{(x_0 - \bar{x})^2}{(n-1)S_x^2}}$$

X值离样本平均更远, interval更宽.
X值相同时, PI比CI更宽.

ANOVA TABLE

| | Source | DF | Sum of Squares | Mean Square | F statistic |
|---|------------|-------|----------------|---------------------------|-----------------------|
| Lec 7. 多元 | Regression | k | SSR | $MSR = \frac{SSR}{k}$ | |
| equation $\hat{Y} = \beta_0 + \beta_1 X_1 + \beta_2 X_2$ (有数值) | Error | n-k-1 | SSE | $MSE = \frac{SSE}{n-k-1}$ | $F = \frac{MSR}{MSE}$ |
| $S_{Y X} = \sqrt{MSE}$ <small>(n-1) def</small> $\hat{Y} = \beta_0 + \beta_1 X_1 + \dots + E$. | Total | n-1 | SSY | | |

解释 coefficient: $\beta_0 & \beta_1$ would have. For every additional A, the predicted of Y will increase by β_1 , holding β_0 constant.

Residual (实际 - 预测值) ? 选择最合适的模型

Reduced Model: $\hat{Y} = \beta_0 + E$ (R^2 对于人合模型有帮助) | $H_0: \beta_1 = \beta_2 = \beta_3 = 0$. vs $H_1: \text{at least } \beta_i \neq 0$.

解释: tend to deviate from predicted values by α in each direction.

Lec 8. $(V, F_{r, n-r, 0.05})$

$$SS(X_1, X_2, X_3) = SS(X_1) + SS(X_2|X_1) + SS(X_3|X_1, X_2).$$

对于 $SSR(A, B)$, AB 项不重要
对于 $SSR(A|B)$, AB 重要

$$SS(X_1, X_2) = SS(X_1) + SS(X_2|X_1) = SS(X_2) + SS(X_1|X_2)$$

$SSR(A)$ 没有单独的第 2 项

$$SS(X_3, X_4|X_1, X_2) = SS(X_1, X_2, X_3, X_4) - SS(X_1, X_2)$$

$$SS(X_1, X_2) = SS(X_1, X_2) - SS(X_2)$$

Model Reduced 哪个模型是正确的?

最大P值不会降低 SSR why

新加入元素 compete with former.

most variability has been explained by former.

有了 $\beta_1, \beta_2, \beta_3$ 后模型会更好.

enter order 不改变 SSR total | $q \leq n-p-q-1$, intercept | X_0 不影响

$$1^{\circ} \beta_1 = \beta_2 = \beta_3 = 0 \quad 4^{\circ} \beta_3 = 0, X_1, X_2 \text{ 在 model}$$

$$2^{\circ} \beta_1 = 0$$

$$3^{\circ} \beta_2 = \beta_3 = 0, X_1 \text{ 在 model}$$

$$5^{\circ} \beta_2 = 0, X_1 \text{ 在 model.}$$

$$\text{Overall } (X_1, \dots, X_k \text{ is better than empty}). \quad \begin{cases} Y = \beta_0 + \beta_1 X_1 + \dots + \beta_k X_k + E \\ Y = \beta_0 + E \end{cases} \quad F = \frac{SSR/k}{SSE/(n-k-1)} \quad \frac{k}{n-k-1} \text{ (df)}$$

Partial (X_i 是否显著有 $X_1 \dots X_p$ 无关).

$\beta_1 \neq 0$.

$$F = \frac{SS(X_1|X_1, \dots, X_p)}{MSE(X_1, \dots, X_p, X_1)} \quad \frac{1}{n-p-1-1} \text{ df}$$

有一个不为 0!

Multi Pa. ($X_i, X_j \dots$ 是不相关的已有 ...).

$$F = \frac{SS(X_i, X_j|X_1, \dots, X_p)}{MSE(X_1, \dots, X_p, X_i, X_j)} \quad \begin{cases} Y = \dots + \beta_p X_p + \beta_i X_i + \beta_j X_j + E \\ Y = \dots \end{cases}$$

$$\frac{ij \dots}{n-p-2-1} \text{ df}$$

Lec3 X-predictor Y-response

$$\text{Sum of Squares: } SSX = \sum_{i=1}^n (X_i - \bar{X})^2 \quad S_x = \sqrt{\frac{SSX}{n-1}} \quad S_{XY} = \frac{\sum (X_i - \bar{X})(Y_i - \bar{Y})}{n-1} = \frac{SSXY}{n-1}$$

$$\text{Correlation coefficient: } r = \frac{SSXY}{\sqrt{SSX \cdot SSY}} = \frac{S_{XY}}{S_x S_y} \in [-1, 1] \quad 0.4-0.7: \text{moderate}; 0.7-1: \text{strong}$$

Simple Linear Regression Line

$$\hat{Y} = \hat{\beta}_0 + \hat{\beta}_1 X$$

$$\text{Intercept} \rightarrow \text{slope} \cdot r = \frac{\sum (X_i - \bar{X})(Y_i - \bar{Y})}{\sum (X_i - \bar{X})^2}$$

Residual (误差): $\hat{e}_i = \text{实际值} - \text{预测值}$ (> 0 则 above line)

Lec4

$$R^2: \frac{SSY - SSE}{SSY} = r^2 = 1 - \frac{SSE}{SSY} = \frac{\text{Regression}}{\text{Total}}$$

$$R^2 = \frac{\sum (Y_i - \bar{Y})^2}{n-2}$$

Regression Error

$$SSY = SSR + SSE$$

Observed value measured to deviate from predicted value by $s_{Y|X}$

Coefficient of determination: % of the variability in the ... is explained by the ...

Certain? X. Linear model: $\hat{Y} = \hat{\beta}_0 + \hat{\beta}_1 X + E$ (不能有参数).

Independence: 一个变量的 X 是否影响另一个变量的 Y.

Linearity: Sample means 是否 approximately a straight line.

Normality: histogram of residuals 是否 approximately normal. { 是否 skewed.

Homoscedasticity: Residuals 是否 have same spread for each predicted value.

Impact of unusual obs:

Lec5: Slope

Significant predictor: $H_0: \beta_1 = 0$ vs. $H_a: \beta_1 \neq 0$

Test statistic: $t = \frac{\hat{\beta}_1 - \beta_1}{SE}$ CV: $t_{df, 0.95} = \pm 1.96$, 拒绝原假设 95% confidence that the predicted

whether $> \bar{x}$.

correlation = $\sqrt{R-\text{sqn}}$

Test: $t = \frac{\hat{\beta}_1 - \beta_1}{SE}$ CV: $t_{df, 0.95}$ X reject

\bar{x} is a plausible factor.

是否有 intercept: $H_0: \beta_0 = 0$ vs. $H_a: \beta_0 \neq 0$

$t = \frac{\hat{\beta}_0 - \beta_0}{SE}$ (1. But $t_{df, 0.95} \sqrt{\frac{1}{n} + \frac{1}{n-1}} s_{\beta_0}$)

// base \bar{x} , β_0 ? \bar{x} ,

同理

Lec6:

Test of correlation

$$t = \frac{r\sqrt{n-2}}{\sqrt{1-r^2}}$$

P小, 拒绝显著差异

outlier & regression(影响)

make the model fit worse

increase SSE \Rightarrow increase residual

standard deviation \Rightarrow decrease R square.

nonzero ($H_0: \rho = 0$ vs. $\rho \neq 0$)

$$Z = \frac{\frac{1}{2} \ln(\frac{1+r}{1-r}) - \frac{1}{2} \ln(\frac{1+\rho_0}{1-\rho_0})}{\sqrt{\frac{1}{n-3}}} // \frac{\frac{1}{2} \ln(\frac{1+r_1}{1-r_1}) - \frac{1}{2} \ln(\frac{1+r_2}{1-r_2})}{\sqrt{\frac{1}{n_1-3} + \frac{1}{n_2-3}}}$$

Fisher's Z.

P值，拒绝，差异显著。Test值 > CV // $H_0: \beta_1 = ?$ vs $H_A: \beta_1 \neq 0$

3个 dummy variable
 $\begin{pmatrix} 0 \\ 1 \\ 1 \\ 0 \end{pmatrix}$ -> 只在某中出现的

Equation. $\hat{Y} = \text{基础}(前俩) + \text{对应 DV 变量的 slope} + \text{intercept}$
Model: $\hat{Y} = \beta_0 + \beta_1 X + \beta_2 Z$

图: slope不一样, interaction就显著, DV 都改变, intercept 不变。

即, prediction of X, every additional 1, increases/decreases by XX.

confounding variable: 随解释不同时, 注意 extraneous variable

check: 让 β_1 在 [] , 若 $\beta_{1,2}$ 不在 [] , 则为 confounder

best set

扰乱关系 X 和 Y

control variable: 防止 confusion

判定是否 confound: FICCI

$P(C) \geq 0.975, 148$

2. 算 CI: $\beta_1 \pm (V \cdot SE)$

3. 看完模型后该元素 $\beta_1(E)$

是否在 CI, coefficient

4. 不在 CI 内 => 是 confound.

1. 先问模型是好? (CI 最窄)

是否 confound 和模型的 model 的 CI 相同。

2. Interaction.

Model $\hat{Y} = \beta_0 + \beta_3 X Z$

判别 relationship

$\hat{Y} = \beta_0 + \beta_1 X + \beta_2 Z + \beta_3 X Z E$

L12 所有直线 (must be parallel // significant intercept)

test dummy variables 的系数是否为 0.

$\begin{array}{c} z=1 \\ z=0 \\ z=-1 \end{array}$ 必须平行才能用 ANOVA.
 $\beta_1 + \beta_2$ 趋势: 是否 parallel
 β_3 方向: $\beta_1 + \beta_2$ 方向

2.1 Full: $\hat{Y} = \beta_0 + \beta_1 X + \beta_2 Z + E$ vs Reduced: $\hat{Y} = \beta_0 + \beta_1 X + E$

$F = \frac{\text{Full's SSR} - \text{Reduced's SSR}}{\text{Full's Res / df}}$ test Z 的线性

参数是否与 determined 线性参数不同

Adjusted Means.

DV 之间不能比, 只能比 up.

Full: $\hat{Y} = \beta_0 + \beta_1 X_1 + \beta_2 X_2 + \beta_3 Z_1 + \beta_4 Z_2 + E$ $\{H_0: \beta_3 = \beta_4 = 0\}$

Reduced: $\hat{Y} = \beta_0 + \beta_1 X_1 + \beta_2 X_2 + E$ $\{H_1: \text{at least one}\}$

$F = \frac{(\text{Sqr}(SS(\text{Full} - \text{Reduced}))^2 / \text{df}}{\text{Full's Res / df}}$

Full's Res / df 根据 Z 变量选择 DV.

求值: 根据 input 到一个 \bar{Y} , 找入 means 作图.

| Source | DF | SS | Mean Square | F-stat |
|------------|-------|-----|---------------------------|-----------------------|
| Regression | k | SSR | $MSR = \frac{SSR}{k}$ | $F = \frac{MSR}{MSE}$ |
| Error | n-k-1 | SSE | $MSE = \frac{SSE}{n-k-1}$ | |
| Total | n-1 | SSY | | |

$$R^2 = \frac{SSR}{SSY} \quad \text{Residual Standard Dev} = \sqrt{\frac{SSE}{n-k-1}} = \frac{\sqrt{SSE}}{\sqrt{n-k-1}}$$

① coincidence (是否 same slope & intercept) $H_0: \beta_1 = \beta_2 = 0$
P值, 拒绝 X coincident, significant. $Z \propto \frac{X_1 - X_2}{\sqrt{SSE}}$

② parallelism (是否 same slope but different intercept)

$H_0: \beta_2 = \beta_3 = 0$ // 小, 拒绝, not parallel // not reject L.

$F = \frac{(\beta_2, \beta_3 \text{ 对应的 } SSR)^2 / 2}{\text{Full model Res / DF}}$ β_0, β_1, E

③ $F = \frac{(\text{Full} + \beta_3 \text{ 对应的 } SSR - \text{Reduced}) / 2}{\text{Full} + \text{Res / DF}}$ $X \text{ and } Z \text{ on } X \text{ independent}$

Legal: 1 1 $\beta_0, \beta_1, \beta_2, \beta_3, E$
Illegal: 0 0 $X \text{ and } Y \text{ different } Z$

// Z=1 时, β_2 问题 intercept, β_3 问题 slope /

L13. 三因变量 (三者必须都满足), $\begin{cases} 1. \text{parallel} \\ 2. \text{intercept} \\ 3. \text{residuals} \end{cases}$

point 1. outlier: 有很大且极端的残差.

2. High leverage observation: X 距离很远, y 距离前值, Residuals 形如 R line ($> \frac{k+1}{n}$)

3. Influential point: 影响 slope / R line 的方向

Hi $\rightarrow y$ 和 Residuals 都很极端, Cook's distance $> \frac{4}{n}$

Leverage $(0, 1) \text{ cut-off: } \frac{2(k+1)}{n}$ // k 为有几个 predictor 元素, n 为 observations

High leverage \rightarrow 拥有极端值

Residuals 1. Standardized Residual: $z_i = \frac{y_i - \hat{y}_{\text{pred}}}{RSE}$

2. Studentized residual: $s_i = \frac{y_i - \hat{y}_{\text{pred}}}{\sqrt{MSE(1-h_i)}}$

3. $r_{(i)} = \frac{y_i - \hat{y}_i}{\sqrt{MSE(1-h_i)}}$ // If without i-th model,

outlier CV: $\text{agt}(1.025, df)$ \rightarrow $i: h_i$ $\propto S_{XIX}$

obs - pred = 1 $d_i = \frac{i}{k+1} \cdot \frac{h_i}{1-h_i} r_i^2$

Cook's distance = $\frac{4}{n_{\text{obs}}}$ // cut off: num of pred

look R² is influential, outcupt outlier, leverage & high L 引入多樣的 test for addition of β_i .

六有滿足那3個條件，才能safe to remove influential
normality straight-line
linearity 是否 obvious curve in resplot.
Homo: 這行不齊，是 random scatter

L14 為 straight, if homo & normal X
if non-X if linear X, but other \rightarrow linear fit

$\{\ln(Y)\}$: $\hat{Y} \rightarrow$, variance of res \neq 1 时, 应变量 > 0 .

\sqrt{Y} . 同 $\ln(Y)$: \hat{Y} ; \hat{Y} . 同 \hat{Y}

$\{\ln(Y)/\hat{Y}\}$: response is proportion or rate, $(0, 1) \sim (-\infty, \infty)$

Log transform.

通過 reg eq求 predicted response. 的, 再把 Y 正序

L15 多項式: $Y = \beta_0 + \beta_1 X + \beta_2 X^2 + \dots + \beta_k X^k + E$

修正linearity \Rightarrow lack of fit test ①

$\{\text{H}_0: \text{model合適} \rightarrow X \text{ lack of fit}$ | $F = \frac{\text{SS}_{\text{Lack of fit}} / (d-k-1)}{\text{pure SSPE} / (n-d)}$

1. Remove a subset. \rightarrow data-based.

2. create a linear combination of pred.

3. center predictors. \rightarrow structural

和 correlation 高的 pair, 但是与 response 有关 (cor)

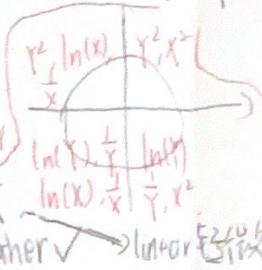
去除无关联的 pair, 2个高 pair 高向 cor.

P值 A/B/C.

为什么一个元素的 P 值那么高? (因为有向 with other elements).

另一个元素却大于 SSR, 很小的给这不多.

是两个元素都一样的东西不一样.



$$F = \frac{\text{extra ss for addition}}{\text{MSE Full model}}, df(1, n-k-1)$$

$$F = \frac{\text{ss}(X^2, \text{sq})}{\text{Res}/df}, df(1, n-k-1)$$

P1. 括弧, add 显著是升序型.

当 X lack of fit, 则为 best model.

如果 lack, 用 Julie, 带进模型.

如果 straight-line 和 quadratic 通过 lack of fit
看 addition of cubic term 是否 significant, F test.

L16 collinearity: predictors 很高高度相关 on + on

$$VIF = \frac{1}{1-R^2} \rightarrow r-square; \text{ remove.}$$

cluster of collinearity, R^2 越大越不好.

$$= \frac{\sum (\bar{Y} - \hat{Y})^2 / (d-k-1)}{\sum (Y - \hat{Y})^2 / (n-d)} = \frac{\text{lack of fit} / df}{\text{pure error} / df}, (V = q(123), df)$$

center: x. - only works when higher order terms are involved. Every term in this model is linear and centering \rightarrow x to \hat{x} to \hat{Y} to collinearity.

[Combination: $X/A \dots; B \dots$. They have different units of measurement and should not be combined using \dots bc the units would not make sense.

| ANOVA Table | | | | |
|-------------|-------|--|---------------------------------|---|
| Src | DF | Sum of Squares | Mean Square | F-stat |
| Regression | k | SSR | MSR = $\frac{SSR}{k}$ | $F = \frac{MSR}{MSE}$ |
| Error | N-k-1 | SSE | MSE = $\frac{SSE}{N-k-1}$ | |
| Total | N-1 | SSY | $SYIX = \sqrt{\frac{SSE}{n-2}}$ | |
| | | $SSR = \sum (\hat{Y}_i - \bar{Y})^2$, $SSE = \sum (Y_i - \hat{Y})^2$ | $n-2$ | $= \sqrt{MSE}$ |
| | | slope $\hat{\beta}_1 = \frac{\sum (X_i - \bar{X})(Y_i - \bar{Y})}{\sum (X_i - \bar{X})^2}$ | $R^2 = \frac{SSR}{SYIX}$ | residual standard deviation RSD. |
| | | $R^2 = \frac{\sum (Y_i - \bar{Y})^2}{\sum (Y_i - \hat{Y})^2}$ | $R^2\% \dots$ is explained by | |
| | | $\hat{\beta}_1 = \frac{SYIX}{\text{standard Deviation} / \sqrt{n-1}}$ | | |
| | | | | (I) (forall) $\hat{Y} = t_{n-2} \times SYIX \sqrt{\frac{1}{n} + \frac{(X_0 - \bar{X})^2}{(n-1)S_x^2}}$ |
| | | | | PI (single). $\hat{Y} = \sqrt{1 + \frac{1}{n} + \dots}$ |
| | | | | outlier 使模型更差, 提高 SSE, RSD, $\sqrt{R^2}$. 值离 mean 更远, interval 更宽 |
| | | | | 值相同, P1比U宽 |
| | | | | L2 ANACOVA 条件: 平行不能 YD1. |
| | | | | $H_0: \beta_3 = 0$. P值, test CV |
| | | | | $F = \frac{\text{Full}_{\text{res}} - R_{\text{res}}}{\text{Full}_{\text{res}} / df}$ 让 X 不重要. |
| | | | | adjust mean. Full 没有 XZ |
| | | | | $\begin{cases} \text{Full} = \beta_0 + \beta_1 X_1 + \beta_2 X_2 + \beta_3 Z_1 + \beta_4 Z_2 \\ R_{\text{res}} = \beta_0 + \beta_1 X_1 + \beta_2 X_2 + E \end{cases}$ |
| | | | | $H_0: \beta_3 = \beta_4 = 0$, vsat (east) |
| | | | | $F = \frac{(\text{Full}_{\text{res}} - R_{\text{res}}) / df}{\text{Full}_{\text{res}} / df}$ 来自 根据 output F |
| | | | | - T 代入 mean. Z 不影响 Z2 VG |
| | | | | $F = \frac{(\beta_3 / \beta_3 \text{ 对应的 } SSR) / 2}{\text{Full model Res/DF}}$ |
| | | | | $F = \frac{\text{Full}_{\text{res}} - R_{\text{res}}}{\text{Full Res/df}}$ |
| | | | | $F = \frac{(\text{Full}_{\text{res}} / 3) \text{ SSR} - R_{\text{res}} / 2 \text{ SSR}}{\text{Full Res/df}}$ |

Slope: $H_0: \beta_1 = \text{值}$ vs $H_1: \neq$. $t = \frac{\hat{\beta}_1 - \text{值}}{S_{\hat{\beta}_1}}$

C2. $\hat{\beta}_1 \pm t_{n-2} \left(\frac{S_{\hat{\beta}_1}}{S_{\hat{\beta}_1} \sqrt{n-1}} \right)$ // test < CV, 值 plausible

Two correlation. $p = p_{\text{just}}$

Intercept: $H_0: \beta_0 = \text{值}$, $H_1: \neq$. $t = \frac{\hat{\beta}_0 - \text{值}}{S_{\hat{\beta}_0}}$

C2: $\hat{\beta}_0 \pm t_{n-2} \left(S_{\hat{\beta}_0} \sqrt{\frac{1}{n} + \frac{\bar{x}^2}{n-1} S_x^2} \right)$

Zero correlation. $H_0: \rho = 0$ vs $\neq 0$.

/ Non-zero: $p = p_0$ vs $\neq p_0$

$t = \frac{\sqrt{n-2}}{\sqrt{1-p^2}}$

Z: $\frac{1}{2} \ln(\frac{1+p}{1-p}) - \frac{1}{2} \ln(\frac{1+p_0}{1-p_0})$

for持继add.

$Z = \frac{\text{Reg}(X_3) / 2}{\sqrt{\frac{1}{n-3} + \frac{1}{n_2-3}}}$

$SS(X_1, X_2, X_3) = SS(X_1) + SS(X_2 | X_1) + SS(X_3 | X_1, X_2)$. // 单SS, 第二行是那一个数.

$SS(X_3 | X_1, X_2) = SS(X_1, X_2, X_3) - SS(X_1, X_2)$

Reg Full, reduced $F = \frac{\text{Reg}(X_3) / 1}{\text{Residual}^2}$

test for addition of X_3

$SS(X_3 | X_1, X_2) = SS(X_1, X_2, X_3, X_4) - SS(X_1, X_2)$

$SS(A | B), B$ 重要

Reg Full, reduced $F = \frac{\text{Reg}(X_3) / 1}{\text{Residual}^2}$

test for addition of X_3

$SS(X_2 | X_1) \neq SS(X_2)$

$SS(X_2) = SS(X_1) + SS(X_2 | X_1) = SS(X_2) + SS(X_1 | X_2)$

Reg Full, reduced $F = \frac{\text{Reg}(X_3) / 1}{\text{Residual}^2}$

test for addition of X_3

X_2 在模型里, X_1 呢? $SS(X_1 | X_2) = SS(X_1, X_2) - SS(X_2)$.

Reg Full, reduced $F = \frac{\text{Reg}(X_3) / 1}{\text{Residual}^2}$

test for addition of X_3

Overall Reg Full: $\hat{Y} = \beta_0 + \beta_1 X_1 + \beta_2 X_2 + E$

Reg Full, reduced $F = \frac{\text{Reg}(X_3) / 1}{\text{Residual}^2}$

test for addition of X_3

$R^2 = Y = \hat{Y} + E$

Reg Full, reduced $F = \frac{\text{Reg}(X_3) / 1}{\text{Residual}^2}$

test for addition of X_3

$H_0: \beta_3 = \dots = \beta_k = 0$ vs at least one $\neq 0$.

Reg Full, reduced $F = \frac{\text{Reg}(X_3) / 1}{\text{Residual}^2}$

test for addition of X_3

Partial F $F = \frac{SS(X^* | X_1, \dots, X_p)}{MSE(X_1, \dots, X_p, X^*)}$

Reg Full, reduced $F = \frac{\text{Reg}(X_3) / 1}{\text{Residual}^2}$

test for addition of X_3

Mult. $t = \frac{\hat{\beta}_1 - \text{值}}{S_{\hat{\beta}_1}}$

Reg Full, reduced $F = \frac{\text{Reg}(X_3) / 1}{\text{Residual}^2}$

test for addition of X_3

正回归诊断! 所有条件都满足. standard. res $Z_i = \frac{Y - \hat{Y}}{S_{\text{res}}}$

Reg Full, reduced $F = \frac{\text{Reg}(X_3) / 1}{\text{Residual}^2}$

test for addition of X_3

1° outlier 很大且不在原直线上

Reg Full, reduced $F = \frac{\text{Reg}(X_3) / 1}{\text{Residual}^2}$

test for addition of X_3

2° High leverage obs. 距离远且不在原直线上

Reg Full, reduced $F = \frac{\text{Reg}(X_3) / 1}{\text{Residual}^2}$

test for addition of X_3

3° Influential point: 影响斜率. (X_4) 离原直线上很远

Reg Full, reduced $F = \frac{\text{Reg}(X_3) / 1}{\text{Residual}^2}$

test for addition of X_3

Outlier. (X_4) 离原直线上很远

Reg Full, reduced $F = \frac{\text{Reg}(X_3) / 1}{\text{Residual}^2}$

test for addition of X_3

Outlier. (X_4) 离原直线上很远

Reg Full, reduced $F = \frac{\text{Reg}(X_3) / 1}{\text{Residual}^2}$

test for addition of X_3

L16. collinearity $\lambda \approx 10^{-10}$

Reg Full, reduced $F = \frac{\text{Reg}(X_3) / 1}{\text{Residual}^2}$

test for addition of X_3

$\sqrt{F} = \frac{1}{1-\lambda^2} r^2$ r-square (DU对(原直))

Reg Full, reduced $F = \frac{\text{Reg}(X_3) / 1}{\text{Residual}^2}$

test for addition of X_3

lombing (原直) $\lambda = \frac{1}{\text{Residual}^2}$ (原直) $\lambda = \frac{1}{\text{Residual}^2}$

Reg Full, reduced $F = \frac{\text{Reg}(X_3) / 1}{\text{Residual}^2}$

test for addition of X_3

决定系数 $R^2 = \frac{\sum(Y - \hat{Y})^2}{\sum(Y - \bar{Y})^2}$

Reg Full, reduced $F = \frac{\text{Reg}(X_3) / 1}{\text{Residual}^2}$

test for addition of X_3

决定系数 $R^2 = \frac{\sum(Y - \hat{Y})^2}{\sum(Y - \bar{Y})^2}$

Reg Full, reduced $F = \frac{\text{Reg}(X_3) / 1}{\text{Residual}^2}$

test for addition of X_3

(2) 模型选择

Max Model: $\hat{Y} = \beta_0 + \beta_1 X_1 + \beta_2 X_2 + \dots$

model 最多多少 variable? $\beta_3 X_3 + \beta_4 X_4 + \beta_5 X_5 + E$ (全部)

$N > 5k \rightarrow$ predictor 数

$R^2 = 1 - \frac{\text{SSE}(n-k-1)}{\text{SSE}(n-1)}$ $\xrightarrow{\text{出现冗余}}$

variable 数

即对于一个 pred, 至少 5 个 obs.

$\frac{\text{SSE}}{n-1} \xrightarrow{\text{SSE}} \frac{\text{SSE}}{n-k-1}$ \downarrow $\frac{\text{SSE}}{\text{SSE}} \xrightarrow{\text{MSE}} \text{MSE}$

多个新 pred, $R^2 \uparrow$ (\uparrow 到 SSR) \downarrow SSR $\xrightarrow{\text{模型}}$

$C_p = \frac{\text{SSE}(\text{Reduced})}{\text{MSE}(\text{Full})} - [n - 2(p+1)]$

R^2 越大越好 / AIC = $n \ln(\text{SSE}) - n \ln(n) + 2p$

BIC = $n \ln(\text{SSE}) - n \ln(n) + p \ln(n)$

I.E. (18)

A. B 值越小越好 不是同模型 有几倍数的惩罚
 { 用 SSE 衡量模型. Forward: 添加 pred 直到 P 不是 pred, 不变 (opt)
 减去 $\ln(n)$ 是因为样本 对大模型有惩罚. $AIC(B)$ 会增加. $P = \underline{p_{\text{pred}+1}}$

Backward: 逐渐降 pred 让其减小. Stepwise: 两者结合, 在每一步做决定

F: 看单一模型, 找出里面最 Mio, 再看对元素有 X_i 的, 找出能降低 AIC, 再从已确定的选第三个. 找到最好的了. B: 从缺一的模型里选 BIC 最低的再缺二... 先做 B!

L21 时间序列

描述 linear trend: No. volatile (波动).

sly rocket (飙升). plummet (暴跌)

Positive auto correlation: 给定符号的残差跟前面相同符号

Negative autocorrelation: ... 跟前面相反, 自相关

高波动低相关: positive, independence condition 不成立

Consecutive positive residuals tend to be followed by positive residuals.

The same tends to happen once the residuals flip to being negative. (time series plot of residuals cross 0 次数统计).

Durbin-Watson Test 自相关是否存在 $H_0: \rho = 0$

值到 0 正自; 到 4 负自

$$\hat{Y} = E + \beta_1 L_1 + \beta_2 L_2 + \beta_3 L_3 \dots$$

(前向后)

Lagging: 将数据移位一定时间段以自适应 说明今天明天 highly correlated 但时间过去, 预测性降低

从每个 size 中挑出最大的, 算 diff, 对比.

LEC 19 交叉验证 { 一次删一个, 剩余继续拟合, 预测 Leave-one-out CV: 已删除的适合小样本 (<1000)

SSR = PRESS. $\text{RMSE} = \sqrt{\frac{\text{PRESS}}{n}}$ RMSE 越小越好, 说明模型精度.

k fold CV k=n 时与 LOOCV 相同

$$\text{RMSE} = \sqrt{\frac{\text{RMSE}_1^2 + \text{RMSE}_2^2 + \dots + \text{RMSE}_k^2}{k}}$$

LE (20) 遗嘱回归

$$\text{logit}(p) = \beta_0 + \beta_1 X_1 + \dots + \beta_k X_k \quad p = \frac{e^{\beta_0 + \beta_1 X_1 + \dots + \beta_k X_k}}{1 + e^{\beta_0 + \beta_1 X_1 + \dots + \beta_k X_k}}$$

logistic reg model \rightarrow A. 对应数值

$$\text{odds} = \frac{P}{1-P} \quad 95\% CI = \text{Coef} \pm 1.96 \text{ SE Coef}$$

| | Obs | Success | | Total | Accuracy = $\frac{T_P + T_N}{n}$ |
|---------|-----|---------|------|-------------------|---------------------------------------|
| | | Fail | Pass | | |
| Fail | | T_N | F_N | N _{pred} | Sensitivity = $\frac{T_P}{T_P + F_N}$ |
| Success | | F_P | T_P | P _{pred} | |
| Total | | N | P | n | Specificity = $\frac{T_N}{T_N + F_P}$ |

1° Accuracy % 确保预测模型 pass $\frac{T_P + T_N}{T_P + F_N}$

2° sensitivity. + % 被 misclassified. as rushes 1% 好坏, 3° SP: + % rush as pass.

$$OR_{A \text{ vs } B} = \frac{\text{odds}(A)}{\text{odds}(B)} \quad OR = e^{\beta_i}$$

[1] for slope: $\beta_1 \pm Z_{1-\alpha/2} S_{\beta_1}$ (不含 0, 显著)

[2] for odds ratio: $\exp(\beta_1 \pm Z_{1-\alpha/2} S_{\beta_1})$ 不显著

根据 D-W Test 判断是否相关: DW 接近 2 表示无相关.

P 值大, 无显著差异, 接受 H₀, 无相关.

→ 4

L22 missing data

MAR: 缺失数据与特征无关, 不同的变量.

MCAR: 值缺失的频率与所测变量值无关.

如何确定? (比其他变量缺失分布).

MCAR: 根据已知值的分布估算整体 mean/median.

自回归 model 例题: DW 2, P 值大, 无法拒绝 H_0 满足假设

? diagonal 有相同相关性

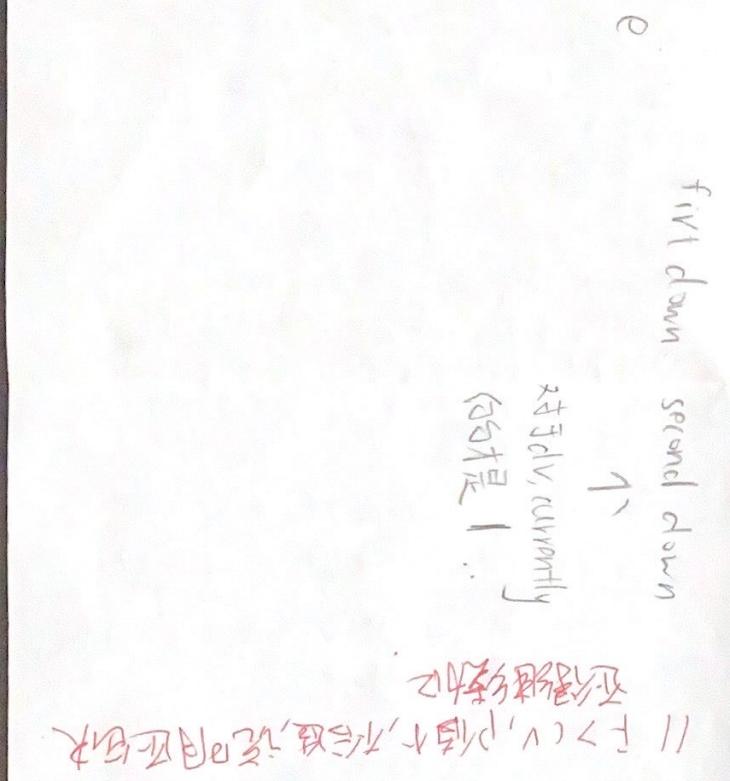
predictors 由 lagging 创建, 相同的值出现在不同的 row, 从而共享 correlation.

9. 12 11. 13.

MAR: 用 imputation 替代丢失值 with plausible alternative
 ↓
 no difference

MAR: 剔除 obs ; pred, 替代
 imputation: 完整 dataset, 更精确, 但有风险.
 ↳ 完全不同 lower levels.
 如果像 缺少对应类别的数据

如何 impute → symmetric: 用 mean 来 impute
 看 missing 与什么有关:
 表中, (No) (a) 差异很大, 将那两个格子给对立的, 其它差不多, 用的是 mean 值缺类别数据,
 把有 & mean 缺有 100% 人, 把没 ... 等 ...



只加 LOF, 根本是不 appropriate.
 只加 Add, 企图模型为 reduced add of X₁ & X₂. P₁ > (V, P₂ < V, P₃ < V) 说明 Add 重要

$$\begin{aligned}
 & \text{① Add of } X + X^2 = \beta_0 + \beta_1 X + \beta_2 X^2 \\
 & \text{② Local off-fit } Y = \beta_0 + \beta_1 X + E \quad | - , X - 2 \\
 & \text{③ Add } X^2 \quad Y = \beta_0 + \beta_1 X + E \quad | - X - 3 \\
 & \text{④ Local off-fit } Y = \beta_0 + \beta_1 X + E \quad | - , X - 2, X - 1 \\
 & \text{⑤ Add } X^3 \quad Y = \beta_0 + \beta_1 X + \beta_2 X^2 + E \quad | - X - 4 \\
 & \text{⑥ LOF} \quad Y = \beta_0 + \beta_1 X + \beta_2 X^2 + \beta_3 X^3 + E \quad | - X - 5
 \end{aligned}$$

为了计算 LOF T-test (CV, 没有)

$$\begin{aligned}
 & \text{⑦ Add } X^4 \quad Y = \beta_0 + \beta_1 X + \beta_2 X^2 + \beta_3 X^3 + \beta_4 X^4 + E \quad | - X - 6 \\
 & \text{⑧ Add } X^5 \quad Y = \beta_0 + \beta_1 X + \beta_2 X^2 + \beta_3 X^3 + \beta_4 X^4 + \beta_5 X^5 + E \quad | - X - 7 \\
 & \text{⑨ Add } X^6 \quad Y = \beta_0 + \beta_1 X + \beta_2 X^2 + \beta_3 X^3 + \beta_4 X^4 + \beta_5 X^5 + \beta_6 X^6 + E \quad | - X - 8 \\
 & \text{⑩ Add } X^7 \quad Y = \beta_0 + \beta_1 X + \beta_2 X^2 + \beta_3 X^3 + \beta_4 X^4 + \beta_5 X^5 + \beta_6 X^6 + \beta_7 X^7 + E \quad | - X - 9
 \end{aligned}$$

() more likely to
 comparing to ~ needing shelves
 b/c ... () < |, less likely to pass
 > more likely
 | equal

$$\text{Add F} = \frac{\text{Lack of fit} / df}{\text{Pure Error} / df}$$

$$\text{Add F} = \frac{\text{Lack of fit} / df}{\text{Error / df}}$$

$$\text{① LOF} \quad Y = \beta_0 + \beta_1 X + \beta_2 X^2 + \beta_3 X^3 + E \quad | - X - 1$$

$$\text{② Add } X^2 \quad Y = \beta_0 + \beta_1 X + E \quad | - X - 2, X - 1$$

$$\text{③ Add } X^3 \quad Y = \beta_0 + \beta_1 X + \beta_2 X^2 + E \quad | - X - 3$$