

L10. Model selection - Criteria, and Search Procedures

Criteria for Model selection

- Adj- R^2
- Mallon's Cp
- t test p values
- AIC and BIC

Automatic Search Procedures

- Best subset
- Stepwise
- Forward, Backward

or both

↑
commonly seen in practice.

So far by talking about t test, anova, R^2 and adj- R^2 ,

we have a taste that choosing what predictors to be left in the model is not always obvious.

- When we use different criterias, it might lead to different candidates.

Another difficulty might be when we have a large number of predictors, there are too many candidate models. For ex, 10
~~parameters~~
~~predictors~~ $\rightarrow 2^{10} = 1024$ possible regression models.

so some automatic search proc procedures are needed to help us.

- Best subset would fit all possible models but could be a huge computational burden
- Stepwise is simpler but might "miss" the best model. since it's just adding / dropping predictors in one-way, and stops when it meets the criteria.

2. Selection Criteria

$$(i) \text{adj-}R^2 = 1 - \left(\frac{n-1}{n-p}\right) \cdot \frac{SSE}{SST}$$

takes a balance of ~~gain in SSE~~ decrease in SSE and decrease in df = n-p when adding more predictors

(ii) Mallows Cp.

dist. version, not sample version.



The Mean Square Error for \hat{Y}_i is defined as

$$E(\hat{Y}_i - \mu_i)^2$$

where $\mu_i = E(Y_i)$

The bias of $\hat{Y}_i - \mu_i$ can be broken into

$$\hat{Y}_i - \mu_i = (\underbrace{E(\hat{Y}_i) - \mu_i}_{\text{when the model}}) + (\underbrace{\hat{Y}_i - E(\hat{Y}_i)}_{\text{based on a sample data}})$$

is incorrect. The
is a bias

\hat{Y}_i based on a
sample data has
a bias from $E(\hat{Y}_i)$

Then it can be shown that

$$E(\hat{Y}_i - \mu_i)^2 = [E(\hat{Y}_i) - \mu_i]^2 + \text{Var}(\hat{Y}_i)$$

And the total mean square error for all observations is

$$\sum_{i=1}^n [E(\hat{Y}_i) - \mu_i]^2 + \sum_{i=1}^n \text{Var}(\hat{Y}_i)$$

Defined the "standardized" version as

$$P_p = \frac{1}{\sigma^2} \left[\sum_{i=1}^n \left[E(Y_i) - \hat{Y}_i \right]^2 + \sum_{i=1}^n \text{Var}(\hat{Y}_i) \right]$$

[OPTIONAL]

$$\text{SSE} = \sum (Y_i - \hat{Y}_i)^2$$

$$= (\vec{Y} - \vec{X}\vec{B})^\top (\vec{Y} - \vec{X}\vec{B})$$

$$= \vec{Y}^\top (I_n - \vec{X}(\vec{X}^\top \vec{X})^{-1} \vec{X}^\top) \vec{Y}$$

$\underbrace{\quad}_{\vec{A}}$

By using a result for quadratic forms:

$$E(\vec{Y}^\top \vec{A} \vec{Y}) = E(\vec{Y}^\top) \vec{A} E(\vec{Y}) + \text{tr}(\Sigma A)$$

where $\Sigma = \text{Var}(\vec{Y})$

$$E(\text{SSE}) = E(\vec{Y}^\top) (I_n - \vec{X}(\vec{X}^\top \vec{X})^{-1} \vec{X}^\top) E(\vec{Y})$$

$$+ \text{tr}[I_n - \vec{X}(\vec{X}^\top \vec{X})^{-1} \vec{X}^\top] \sigma^2$$

$$= \sum_{i=1}^n (E(Y_i) - E(\hat{Y}_i))^2 + \sigma^2 (n - \text{tr}[(\vec{X}^T \vec{X})^{-1}])$$

$$= \sum_{i=1}^n (E(\hat{Y}_i) - \mu_i)^2 + \sigma^2(n-p)$$

$= \text{tr}(2P)$
 $= P$

$$\text{Var}(\hat{Y}_i) = \text{cov}(Y_i, \hat{Y}_i) = \text{cov}\left(\sum_{j=1}^n h_{ij} Y_j, \sum_{k=1}^n h_{ik} Y_k\right)$$

$$= \sum_{j=1}^n \sum_{k=1}^n h_{ij} h_{ik} \text{cov}(Y_j, Y_k)$$

$$= \sum_{j=1}^n h_{ij}^2 \text{Var}(Y_j)$$

$$= \sigma^2 \sum_{j=1}^n h_{ij}^2 = \sigma^2 h_{ii}$$

$$\sum_{j=1}^n h_{ij}^2 = h_{ii}$$

because $H^T H = H$

$$\sum_{i=1}^n \text{Var}(\hat{Y}_i) = \sigma^2 \sum h_{ii}$$

$$= \sigma^2 \text{tr}(H)$$

$$\text{tr}(H) =$$

$$\text{tr}(X(X^T X)^{-1} X^T)$$

$$= \sigma^2 P$$

$$= \text{tr}(X^T X (X^T X)^{-1})$$

$$= \text{tr}(I_p)$$

$$= P$$

Therefore,

$$\begin{aligned}P_p &= \frac{1}{\sigma^2} [E(\text{SSE}) - (n-p)\sigma^2 + p\sigma^2] \\&= \frac{E(\text{SSE})}{\sigma^2} - (n-2p)\end{aligned}$$

when the fitted model is accurate, $E(\text{SSE}) = (n-p)\sigma^2$

then $P_p = P$ \square

- In practice, we have a total of $(p-1)$ predictors (p parameters) and Mallon's C_p is calculated for each model with $k-1 \leq p-1$ predictors. (k parameters):

$$C_p = \frac{\text{SSE - fitted model with } (k-1) \text{ predictors}}{\text{MSE - full model with } (p-1) \text{ predictors}} - (n-2k)$$

- If the model with the chosen $k-1$ predictors is "good" model, AKA. MSE_{k-1} is close to MSE_{p-1} C_p should be around k .

- In practice, since C_p is a sample estimate, it could be less than $p \cdot |k+1|$
 When $C_p < 1$, it's doubtful, might suggest overfitting.
- The commonly used $C_p = p \cdot \hat{M}$ when including all possible predictors, so it can not be used to evaluate the full model, only the model with less than $p-1$ predictors.
- It's recommended to use Mallows' C_p to choose some initial candidate models.

- (iii) use t-test p-values: less common.
- (iv) AIC and BIC - Note: Textbook gave ~~a~~ version with modified constants, what we introduced here is more common.
 - Akaike's Information Criterion
 - Bayesian Information Criterion (or SBC)

They ~~are~~ both measure the likelihood fit of the model

amount of information LOST by a given model:

the less information a model loses,
 the higher the quality of that model

$$(a) \text{AIC} = 2K - 2 \ln(\underset{\uparrow}{L(\beta)})$$

↑
Number of parameters ↑
maximum likelihood function

- Recall: For MLR, the likelihood function

$$L(\beta) = \prod_{i=1}^n \frac{1}{\sqrt{2\pi}} \exp \left[-\frac{(y_i - \mu_i(\beta))^2}{2\sigma^2} \right]$$

To maximize the likelihood function, it's equivalent to
maximize

$$\ell(\beta, \sigma) = \ln(L) = -n \log \sigma - \frac{n}{2} \log 2\pi - \frac{1}{2\sigma^2} \sum_{i=1}^n (y_i - (\hat{x}\beta))^2$$

and $\hat{\beta}_{MLE} = \hat{\beta}_{OLS}$, $\hat{\sigma}^2_{MLE} = \frac{SSE}{n}$, $\ln(\hat{L})$

Note! $\ln(\hat{L})$ is negative

the maximized value is the minimized

- $\ln(\hat{L})$, which is positive.

In fact, $\rightarrow \ln(\hat{L}) \sim \chi^2_{df=n-p}$ non-centered chi-square.

- In summary, a smaller value of $-\ln(\hat{L})$

indicates a better goodness of fit:

$$\rightarrow \ln(\hat{L}) = n \log \hat{\sigma}^2 + n \ln \frac{1}{\hat{\sigma}^2} + n \ln 2\pi + \frac{SSE}{\hat{\sigma}^2}$$

where $\hat{\sigma}^2 = \frac{SSE}{n}$

- AIC consider the penalty of adding predictors:

$$AIC = -2\ln(\hat{L}) + 2K$$

\uparrow
penalty coefficient on number of parameters

- With the same number of predictors,

smallest AIC has the biggest likelihood.

- With the same likelihood,

~~AIC~~ AIC would choose (smaller) the model with less predictors.

- (b) Through complicated calculation, Gideon E. Schwarz gave a Bayesian argument to give a heavier penalty on the number of parameters, under the consideration of sample size

$$BIC = -2\ln(\hat{L}) + \underbrace{\ln n \cdot K}_{\text{penalty coefficient on number of parameters}}$$

- AIC and BIC theoretically have different set-up:
AIC is based on likelihood function, BIC is estimating a function of a posterior probability of a model being true. Both criteria are based on various assumptions and asymptotic approximations.
- In practice, the only difference is in the size of the penalty. BIC penalizes model complexity more heavily.
- AIC is better when we hate false negative finding more
BIC is better when we hate false positive finding more

II. Automatic Search Procedures.

No matter which criteria we choose, there are two most common automatic search procedures:

(a) Best Subsets:

It fits all 2^P models and print the best fitting model(s)

with one predictor, two predictors, three predictors and so on

(the one with only Intercept is normally ignored)

— Good to give a list of candidates

$$\text{Ex. } Y \sim X_1 + X_2 + X_3 + X_4$$

Number of Predictors	adj-R ²	Mallows Cp	X ₁	X ₂	X ₃	X ₄
1	64.5	138.7				X
1	64.5	142.5				X
2	97.4	2.7	X	X		
2	97.7	5.5	X		X	
3	97.6	3.0	X	X		X
3	97.6	3.0	X	X	X	
4	97.4	3.0	X	X	X	X

adj-R²
 Mallows Cp
 choiil

(b) Stepwise - less seen in practice.

Basic idea: adding or removing variables based on t test or anova test ~~result~~ p-value

Before starting: α_E : to-enter significant level

α_R : to-remove significant level

Procedure: Using $Y \sim X_1 + X_2 + X_3$ as example.

(i) Fit SLR: $Y \sim X_1$

$Y \sim X_2$

$Y \sim X_3$

If none of

the pvalue $\leq \alpha_E$

Pick $Y \sim \text{Null}$

STOP

For those with p-value $\leq \alpha_E$

choose the one with smallest p-value.

i.e. X_1

Pick $Y \sim X_1$

(i) Fit: $Y \sim X_1 + X_2$

$Y \sim X_1 + X_3$

If none of

the 2nd pvalue $\leq \alpha_E$

STOP

For the 2nd predictor with p-value $\leq \alpha_E$

pick the 2nd predictor with the smallest
 p-value to enter the model. $p\text{-value of } X_1 > \alpha_R$
 $\xrightarrow{\text{Remove } X_1} Y \sim X_2$
 then check p-value of X_1

\downarrow
 p-value of $X_1 \leq \alpha_R$.

$$Y \sim X_1 + X_2$$

(c) StepAIC, stepBIC

- R comes with stepAIC with options "both", "forward", "backward"
- Python doesn't have build-in automatic search. (yet)

$$Y \sim X_1 + X_2 + X_3$$

- Backward Starting with $Y \sim X_1 + X_2 + X_3$

(i) Drop one:

$$X_3 : Y \sim X_1 + X_2 \quad 63.8 \quad (\text{lowest})$$

$$X_2 : Y \sim X_1 + X_3 \quad 65.8$$

$$X_1 : Y \sim X_2 + X_3 \quad 75.8$$

(ii) Start with $Y \sim X_1 + X_2$

Drop one:

$$X_1 : Y \sim X_2 \quad 73.1$$

$$X_2 : Y \sim X_1 \quad 78.4$$

Stop. Pick $Y \sim X_1 + X_2$

• Forward starting with $y \sim \text{Null}$ and add one:

$$\begin{array}{ll} \text{(i)} \quad y \sim x_1 & 73.1 \\ y \sim x_2 & 78.4 \\ \hline y \sim x_3 & 70.1 \end{array}$$

$$\text{(ii)} \quad y \sim x_2 + x_3 \quad 75.8$$

$$\boxed{y \sim x_1 + x_3 \quad 65.8}$$

stop: pick $y \sim x_1 + x_3$

• Direction = both.: each step would add back the one deleted from previous step to see if it increase/ decreasing AIC/BIC.

(i) Starting with $y \sim x_1 + x_2 + x_3 + x_4$ $AIC = 56.28$

$$-x_1: y \sim x_2 + x_3 + x_4 \quad 54.41$$

$$-x_4: y \sim x_2 + x_3 \quad 49.10$$

$$-x_2: y \sim x_1 + x_3 + x_4 \quad 54.88$$

(ii) Starting: $y \sim x_2 + x_3$

$$+x_1: y \sim x_1 + x_2 + x_3 \quad 56.28$$

$$-x_2: y \sim x_3$$

$$-x_3: y \sim x_1 + x_2 + x_4 \quad 58.22$$

$$-x_4: y \sim x_2$$

$$+x_4: y \sim x_2 + x_3 + x_4$$

(ii) starting with $y \sim x_2 + x_3 + x_4$

$$+x_1: y \sim x_1 + x_2 + x_3 + x_4 \quad 56.28 \quad \text{add it back}$$

$$-x_2: y \sim x_3 + x_4 \quad 52.11$$

$$-x_3: y \sim x_2 + x_4 \quad 53.49$$

... until find
the
lowest AIC.