

Model (variable) selection ← Recall: ① Lack-of-fit (underfitting) ② overfitting

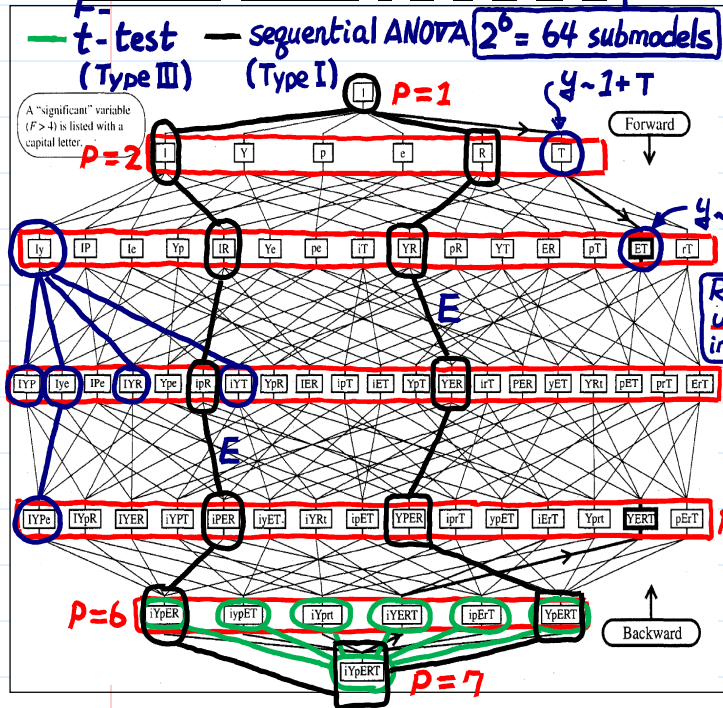
• **Q:** what is a "model selection" problem?

➢ consider the *full model*:

$$y = \beta_0 + \beta_1 g_1(x_1, \dots, x_m) + \beta_2 g_2(x_1, \dots, x_m) + \dots + \beta_{k-1} g_{k-1}(x_1, \dots, x_m) + \epsilon$$

For $1 \leq i \leq k-1$, should the term $\beta_i g_i$ be included in the *final fitted model*?

Example: 6 terms, I, Y, p, E, R, T (sub-)model: a subset set of all $k-1$ terms, e.g.,



always include intercept in these submodels
 $\{I, g_1, g_2\}$
 $\{I, g_2, g_4, g_5, g_{k-1}\}$
 $p = \#$ of parameters in a sub-model \leftarrow not include σ^2
 $\#$ of different sub-models = 2^{k-1}

hierarchical structure of all sub-models (see graph) **nested model**

➢ objective of model selection: select a "best" sub-model (or some good ones)

Q: what is a good sub-model? We usually hope a good model to have
 high R^2 $\rightarrow \hat{\sigma} \approx \sigma$ **lack of fit**
 not too many terms
 terms with significant t -tests

• **Q:** why bother to select a best subset of all terms? ← Why not just use the full model? p. 9-2

- simplicity: principle of Occam's Razor, removal of redundant terms results in a simpler model ← **overfitting: $Var(\hat{Y}) \uparrow$ (LN p. 6-7)** may be easier to interpret $\hat{\beta}$
- unnecessary terms will cost d.f. and add noise to the estimation of other quantities \Rightarrow less precise test/C.I. and tend to increase the standard error $\hat{\sigma}^2 = \frac{RSS}{n-p}$
- collinearity reduction: collinearity is caused by having too many terms trying to do same job \rightarrow significant t -test becomes insignificant ($(X^T X)^{-1}_{ii}$ may become large)
- save cost: if model is used for prediction, can save time and/or money by not measuring redundant terms

\therefore model selection is sensitive to outliers & influential points

- preliminary steps before performing variable selection
 - identify outliers and influential points --- may exclude them temporarily
 - add any terms, transformations, or (linear) combinations of the predictors or extra predictors that seem appropriate
- two types of variable selection procedures: *testing-based* and *criterion-based*

• testing-based procedure **impossible to perform if $k \geq n$ (supersaturated)**
 Recall: large \rightarrow insig, small \rightarrow sig.

① no need to be 0.05
 ② can be, say 0.15-0.2 **because of multiple testing**

\therefore collinearity cannot simultaneously remove terms with insignificant p -values **control the complexity of the chosen model**

backward elimination: (1) start with full model (all terms); (2) eliminate the term with the largest p -value greater than " α -to-remove" preset value; (3) refit the model and go to step (2); (4) stop when all p -values $< \alpha$ -to-remove

check graph in LN p. 9-1

forward selections: (1) start with no terms in the model ($y \sim I$);

check graph in LNp.9-1

(2) For terms not in the model, check their p -values if they are added to the model. Add the term with the smallest p -value less than " α -to-enter" preset value;

(3) refit the model and go to step (2); (4) stop when all the p -values $> \alpha$ -to-enter control the complexity of the chosen model \rightarrow Say, $0.15 \sim 0.2$, avoid stop too soon

stepwise regression: a combination of forward and backward and there are several variations on exactly how this is done. Roughly speaking, at any step, it can

check graph in LNp.9-1

(1) select a new term, according to " α -to-enter", or (2) remove a term from model, according to " α -to-remove", or (3) stop \leftarrow usually suggest α -to-enter $\leq \alpha$ -to-remove to avoid infinite loop.



drawbacks:

- may miss "optimal" model because of its "one-at-a-time" adding/dropping

α -values (α -to-enter and α -to-remove) should not be treated too literally: because of multiple testing occurring

no need to interpret as the significance level

- removal of less significant terms tends to increase the significance of the remaining terms \Rightarrow may lead to overstate the importance of the remaining terms

$\hat{\beta}$ might change & $\hat{\sigma}^2 = \frac{RSS}{n-p}$

e.g., In the example in LNp.9-1, if R & T have strong collinearity, after T is included, it might be very hard to explore submodels containing R . The procedure may only search part of the model space.

The procedure is not directly linked to final objectives of regression, such as prediction or interpretation. It's only based on statistical significance of testing in its selection.

physical significance

- for prediction purpose, testing-based procedure tends to pick smaller models than desired

• criterion-based procedure (k : # of all parameters, including intercept; p : # of parameters in a sub-model; m_p : a sub-model with p parameters): \leftarrow but not $\hat{\sigma}^2$

pick a criterion for judging the worth of a sub-model, consider all possible sub-models and pick those with best values of the criterion

of all possible sub-models $= 2^{k-1} \Rightarrow$ if k is large, computation may be too expensive, clever algorithm like "branch-and-bound" method can avoid it

criterion $\rightarrow m_p \rightarrow f(m_p) \in \text{model space}$
arg min max $f(m_p)$

adjusted R^2 (denoted by R_a^2) criterion: for a sub-model m_p , or stepwise search

RSS: from m_p
TSS: from null model $y \sim 1$

$R^2 = 1 - (RSS/TSS)$: not good, adding terms always increase R^2

$$R_a^2 = 1 - \{ [RSS/(n-p)] / [TSS/(n-1)] \} = 1 - [(n-1)/(n-p)](1-R^2) = 1 - (\hat{\sigma}_{m_p}^2 / \hat{\sigma}_{\text{null}}^2)$$

- will only increase when a term has some value \leftarrow same for all m_p 's
- larger R_a^2 is better [notice the connection between R_a^2 and $\hat{\sigma}_{m_p}^2$]

PRESS (Predicted Residual Sum of Square) criterion $\leftarrow R_a^2 \uparrow \Leftrightarrow \hat{\sigma}_{m_p}^2 \downarrow$

PRESS $= \sum_i \hat{\epsilon}_{(i)}^2$, where $\hat{\epsilon}_{(i)}$ are non-standardized jackknife residuals

entropy in information theory

- smaller value of PRESS is better $\leftarrow y_i \leftrightarrow \hat{y}_{(i)}$
- more expensive computation than R_a^2
- tends to pick bigger models (\Rightarrow may be desirable for prediction purpose) \leftarrow For prediction, tend to prefer more effects

Akaike Information Criterion (AIC) and Bayes Information Criterion (BIC)

substitute MLE

AIC $= -2(\text{maximized log-likelihood}) + 2p$

BIC $= -2(\text{maximized log-likelihood}) + \log(n)p$

usually $\log(n) > 2$

LNp.4-9

for linear model, $-2(\text{maximized log-likelihood}) = n \log(RSS_{m_p}/n) + \text{constant}$

- smaller value of AIC or BIC is better
- get a balance between model fit and model size: BIC penalizes larger models more heavily than AIC \Rightarrow BIC tends to prefer smaller models

➤ Mallow's C_p statistics: MSE of prediction, from m_p

$$E(\hat{\theta} - \theta)^2 = \text{Var}(\hat{\theta}) + \text{Bias}(\hat{\theta})^2 \quad (1/\sigma^2) \sum_i E[\hat{y}_i - E(y_i)]^2$$

would be a good criterion, which can be estimated by:

$$\frac{1}{\sigma^2} \sum_i \text{MSE}(\hat{y}_i) = E\left(\frac{RSS_{m_p}}{\sigma^2} + 2p - n\right) \rightarrow C_p = \frac{RSS_{m_p}}{\hat{\sigma}_{full}^2} + 2p - n,$$

where $\hat{\sigma}_{full}^2$ estimated from the model with all terms (full model)

and RSS_{m_p} is obtained from a sub-model m_p $\sum_i \text{MSE}(\hat{y}_i) = E\|H_i Y - X B_{full}\|^2 = E\|H_i(X B_{full} + \epsilon) - X B_{full}\|^2$

- cheap to compute

H_i : hat matrix under m_p

$$E(RSS_{m_p}) = E\|(I - H_i)Y\|^2 = E\|(I - H_i)X B_{full} + (I - H_i)\epsilon\|^2 = (n-p)\sigma^2 + \|(I - H_i)X B_{full}\|^2$$

closely related to R_a^2 and AIC, BIC

under full model: $RSS_{\{full\ model\}} = (n-k)\hat{\sigma}_{full}^2$, so $C_k = k$ for full model

- for sub-models that fit: $E(RSS_{m_p}) = (n-p)\sigma^2$, so $C_p \approx p$ if $\hat{\sigma}_{full}^2 \approx \hat{\sigma}_{m_p}^2 \approx \sigma^2$, then $C_p \approx \frac{(n-p)\hat{\sigma}_{full}^2}{\hat{\sigma}_{full}^2} + 2p - n = p$
- i.e., C_p close to p implies the sub-model fits

- for sub-models that do not fit: $E(RSS_{m_p}) \gg (n-p)\sigma^2$ and $C_p \gg p$

- it's usual to plot C_p against p . Models with small p and C_p around the $C_p = p$ line or less than p are desirable

different from other criteria in punishment term

Note: C_p , R_a^2 , AIC, BIC all trade-off fit in terms of RSS_{m_p} against complexity (p) \Rightarrow we prefer models with smaller RSS_{m_p} and smaller p ; however, $RSS_{m_p} \downarrow$ as $p \uparrow$

$$C_p = \frac{(n-k) + 2k - n}{\sigma^2}$$

- ambiguity about the best model is possible. When several candidate models exist:

check if models make similar predictions? if yes, can make decision on the basis of cost; if no, do not pick one model arbitrary. Report a range of models.

interpretations qualitatively similar? if not, avoid strong conclusion and report a range of models

examine which has the best diagnostics

- Notes:

e.g. \therefore collinearity

terms not in final model can still be correlated with the response

\Rightarrow not to say they are unrelated to the response;

\Rightarrow better to say they provide no additional explanatory effect beyond those terms included in final model \rightarrow an optimization

- It's important to keep in mind that model selection should not be divorced from the underlying purpose of investigation
- automatic variable selection are not guaranteed to be consistent with your goals. Use these methods as a guide only. can put weights on m_p 's
- these methods do not consider the natural hierarchy in some models: For example, in polynomial model, higher-order terms (such as $x_1^2, x_2^2, x_1 x_2$) should be considered only when corresponding lower-order terms (such as x_1 and x_2) have been included in the model \Rightarrow not all sub-models are candidate models

These criteria make sense from the statistical perspective. However, they may not make sense from the physical perspective.

(Note. These criteria do not use domain knowledge. Under such knowledge, some m_p 's might be preferred than the others.

can summarize consistent patterns in these models to get conclusion with stronger level of evidence (LNp.1-12)

cor(X_1, X_2) = 1 true model: $y \sim X_1$ model selection might choose $y \sim X_2$