

Regression Analysis

Chapter 10: Model Selection

Shaobo Jin

Department of Mathematics

Motivation

It is common that we are unsure about which regressors should be used in the regression model and how they should be used.

- ① Should I fit $\beta_0 + \beta_1 x_1$ or $\beta_0 + \beta_1 x_1 + \beta_2 x_2$? **Nested models**
- ② Should I fit $\beta_0 + \beta_1 x$ or $\beta_0 + \beta_1 x + \beta_2 x^2$, or even a data-driven $f(x)$?
- ③ Should I fit $\beta_0 + \beta_1 x_1$ or $\beta_0 + \beta_2 x_2$? **Non-nested models**

Information Criterion

Suppose that $\mathbf{y} \mid \mathbf{X}$ follows some distribution with density $f(\mathbf{y} \mid \mathbf{X})$.

- We want our model to be complex enough to fit the data well.
 - One way to achieve so is to maximize $f(\mathbf{y} \mid \mathbf{X})$ or $\log f(\mathbf{y} \mid \mathbf{X})$.
- On the other hand, we don't want the model to be too complex to avoid [overfitting](#).
- We cannot simply choose a complex model in order to have a better fit to the current data.
- An information criterion is often of the form

$$-c \log f(\mathbf{y} \mid \mathbf{X}) + \text{penalty of model complexity}$$

for some constant c .

Kullback-Leibler Divergence

Suppose that the true data generating process is $g(y)$, but we assume $f(y | \mathbf{x})$, where f and g may not be the same.

- The **Kullback-Leibler (KL)** divergence is

$$\text{KL}(g, f) = \text{E} \left[\log \left(\frac{g(y)}{f(y | \mathbf{x})} \right) \right] = \int \log \left(\frac{g(y)}{f(y | \mathbf{x})} \right) g(y) dy.$$

- We can show that $\text{KL}(g, f) \geq 0$.
- A large $\text{KL}(g, f)$ means that f is far away from g .

We can rewrite $\text{KL}(g, f)$ as

$$\text{KL}(g, f) = \text{E} [\log g(y)] - \int \log [f(y | \mathbf{x})] g(y) dy.$$

AIC: Minimizing Distance of the Fit from the Truth

- If the unknown parameters are estimated by **maximum likelihood**, a nearly “unbiased” estimator of

$$\text{KL}(g, f) - \mathbb{E}[\log g(y)] = - \int \log[f(y | \mathbf{x})] g(y) dy$$

is

$$\log[f(\mathbf{y} | \mathbf{X})] - m,$$

where m is the number of parameters in our model.

- The **Akaike information criterion (AIC)** is defined as

$$\text{AIC} = -2 \log[f(\mathbf{y} | \mathbf{X})] + 2m.$$

- We prefer the model with the smallest AIC or a parsimonious model that has AIC near the minimum.
- In practice, AIC tends to be conservative, in the sense that it tends to select more explanatory variables.

OBS! AIC in Book

The book defines AIC to be

$$n \log \left(\frac{\text{RSS}}{n} \right) + 2p = n \log \left(\frac{\sum_{i=1}^n (y_i - \hat{y}_i)^2}{n} \right) + 2p.$$

Keep in mind that this expression requires the normality assumption!

Under the normality assumption, the AIC becomes

$$\text{AIC} = n \log \left(\frac{\sum_{i=1}^n (y_i - \hat{y}_i)^2}{n} \right) + 2(p + 1) + n \log(2\pi) + n.$$

BIC: Consistent Model Selection

- Bayesian information criterion penalizes a complex model much more than AIC.

$$\text{BIC} = -2 \log [f(\mathbf{y} \mid \mathbf{X})] + \log(n) m.$$

- We prefer the model with the smallest BIC or a parsimonious model that has BIC near the minimum.
- BIC is consistent in model selection in the sense that

$$P(\text{Choose the true model if it is a candidate}) \rightarrow 1, \text{ as } n \rightarrow \infty.$$

- In contrast, AIC is not consistent.

Stepwise Regression Using AIC

If we have $p - 1$ covariates, then we can fit in total 2^{p-1} models, which is computationally intensive for a large p .

- **Forward selection** starts with an intercept model and sequentially adds the term that gives the lowest AIC.
- **Backward elimination** starts with a model with all explanatory variables and sequentially removes the term that has the greatest decrease in AIC.
- **Stepwise regression** is a combination of both by checking whether adding one term or deleting one term will yield the smallest AIC.

Cross Validation (CV)

CV focuses on the prediction property.

Algorithm 1: One version of cross validation

```
1 Randomly split the data set into  $K$  nonoverlapping groups (K-fold CV) or  
   split the data set into  $n$  groups (leave-one-out CV, aka jackknife);  
2 for  $k = 1$  in  $1 : K$  do  
3   Take the  $k$ th group as test set and the remaining groups as training  
   set ;  
4   while for each model do  
5     Fit it on the training set and evaluate it on the test set ;  
6     Retain the model performance (e.g., MSE, misclassification error,  
       log-likelihood) ;  
7   end  
8 end  
9 Summarize the model performance (e.g., average across  $K$  groups) ;  
0 Choose the model that performs the best ;  
1 Refit the chosen model using the entire data set ;
```

More on Cross Validation

Suppose that the true model is one of our candidate models.

- It has been proved that, for **leave-one-out** CV,

$$\begin{aligned}\lim_{n \rightarrow \infty} P(\text{Choose a model where a nonzero } \beta \text{ is missing}) &= 0, \\ \lim_{n \rightarrow \infty} P(\text{Choose the most parsimonious true model}) &\neq 1.\end{aligned}$$

Hence, the leave-one-out CV is conservative that select a model of excessive size.

- It has been proved that you need **leave- n_v -out** CV if you want

$$\lim_{n \rightarrow \infty} P(\text{Choose the most parsimonious true model}) = 1,$$

where $n_v/n \rightarrow 1$ as $n \rightarrow \infty$.

More on Leave-One-Out Cross Validation

Suppose now that all candidate models are wrong. Consider

$$L_n = n^{-1} (\mathbf{E}(\mathbf{y} \mid \mathbf{X}) - \hat{\mathbf{y}})^T (\mathbf{E}(\mathbf{y} \mid \mathbf{X}) - \hat{\mathbf{y}}).$$

It has been proved that, for any $\epsilon > 0$, we have

$$\lim_{n \rightarrow \infty} P \left(\left| \frac{L_n \text{ of the model selected by CV}}{\text{Minimum } L_n \text{ among all candidate models}} - 1 \right| < \epsilon \right) = 1.$$

Cross Validation and AIC

In fact, we can show that leave-one-out CV is asymptotically equivalent to AIC if the model performance is measured by

$$\sum_{i=1}^n \log f(y_i; \hat{\theta}_{-i}),$$

where $\hat{\theta}_{-i}$ is the estimator of θ after removing the i th observation.

Mallows Criterion

The [Mallows criterion](#) for linear regression is

$$C_p = \sum_{i=1}^n (y_i - \hat{y}_i)^2 + 2\sigma^2 p,$$

where σ^2 is assumed to be known. If it is unknown, it is replaced by an unbiased estimator.

- 1 If σ^2 is known, model selections by C_p and AIC are equivalent under the normality assumption.
- 2 C_p is an unbiased estimator of the out-of-sample prediction mean squared error, if the model is correct.

Some Recommendations

- 1 The **marginality principle** means that when you include higher-order terms into the model, always include all lower-order (and same-order terms).

$$g(\mu_i) = \beta_0 + (\beta_1 x_{i1} + \beta_2 x_{i2}) + (\beta_3 x_{i1}^2 + \beta_4 x_{i1} x_{i2} + \beta_5 x_{i2}^2),$$

$$g(\mu_i) = \beta_0 + (\beta_1 x_{i1} + \beta_2 x_{i2}) + \beta_4 x_{i1} x_{i2},$$

$$g(\mu_i) = \beta_0 + \beta_3 x_{i1}^2 + \beta_4 x_{i1} x_{i2} + \beta_5 x_{i2}^2.$$

- 2 For a qualitative predictor with more than two categories, we select the whole variable rather than individual dummy variables created from the qualitative predictor. Otherwise, the result may depend on the choice of reference category.

Ridge Regression

Suppose that we have demeaned both \mathbf{y} and \mathbf{x} . The **ridge regression** minimizes the penalized sum-of-squares

$$(\mathbf{y} - \mathbf{X}\boldsymbol{\beta})^T (\mathbf{y} - \mathbf{X}\boldsymbol{\beta}) + \lambda \boldsymbol{\beta}^T \boldsymbol{\beta},$$

where $\lambda > 0$ is a tuning parameter.

- ① The penalized sum-of-squares is equivalent to

$$\min_{\boldsymbol{\beta}} (\mathbf{y} - \mathbf{X}\boldsymbol{\beta})^T (\mathbf{y} - \mathbf{X}\boldsymbol{\beta}), \quad \text{s.t. } \boldsymbol{\beta}^T \boldsymbol{\beta} \leq t,$$

for some t .

- ② The ridge estimator has a closed form expression

$$\hat{\boldsymbol{\beta}} = (\mathbf{X}^T \mathbf{X} + \lambda \mathbf{I})^{-1} \mathbf{X}^T \mathbf{y}.$$

- ③ Unfortunately, the ridge regression does not achieve variable selection.

Lasso

Suppose that we have demeaned both y and x . The **least absolute shrinkage and selection operator** (**lasso**) minimizes

$$(\mathbf{y} - \mathbf{X}\boldsymbol{\beta})^T (\mathbf{y} - \mathbf{X}\boldsymbol{\beta}) + \lambda \sum_j |\beta_j|.$$

- ① The penalized sum-of-squares is equivalent to

$$\min_{\boldsymbol{\beta}} (\mathbf{y} - \mathbf{X}\boldsymbol{\beta})^T (\mathbf{y} - \mathbf{X}\boldsymbol{\beta}), \quad \text{s.t.} \quad \sum_j |\beta_j| \leq t,$$

for some t .

- ② The lasso estimator will achieve **sparsity**: the estimates of some parameters can be exactly zero.

Elastic Net

Suppose that we have demeaned both y and x . The **elastic net** is a hybrid of both ridge and lasso:

$$\frac{1}{n} (\mathbf{y} - \mathbf{X}\boldsymbol{\beta})^T (\mathbf{y} - \mathbf{X}\boldsymbol{\beta}) + \lambda \left[\alpha \sum_j |\beta_j| + \frac{1-\alpha}{2} \boldsymbol{\beta}^T \boldsymbol{\beta} \right],$$

where $\lambda > 0$ is the tuning parameter and $0 \leq \alpha \leq 1$ is also a tuning parameter.

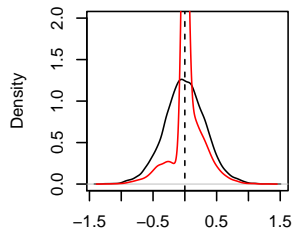
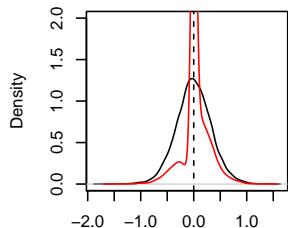
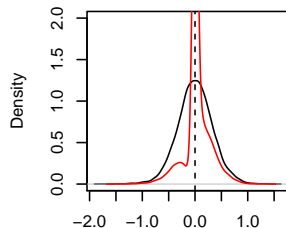
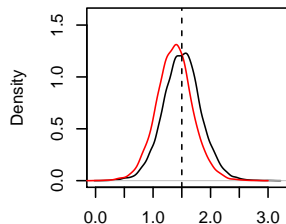
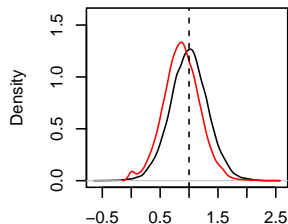
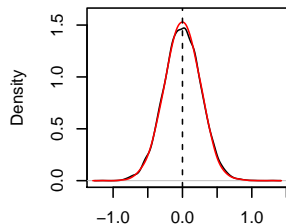
- ① $\alpha = 1$: lasso.
- ② $\alpha = 0$: ridge.

How to Choose λ

- Estimators under penalization are not invariant under rescaling, i.e., a model with meters can be radically different from a model with centimeters. For this reason, covariates are commonly standardized before entering the model (or it is a **must** to some extent).
- A grid search is commonly used: the optimal λ is chosen from a large number of pre-specified λ values.
- Commonly used criteria include AIC, BIC, CV, etc.

Consequence of Penalization

- **bias-variance trade-off** : If $\lambda = 0$, the usual estimator is obtained. If $\lambda > 0$, the **bias-variance trade-off** occurs.
- **Shrinkage**:
 - $\hat{\beta}$ obtained for a positive λ is shrunk towards zero comparing with the ML estimator with $\lambda = 0$. $\hat{\beta}$ is actually a function of λ .
 - The ridge never produces exact zero estimates, but the lasso and the elastic net can produce exact zero estimates.
- **Variable selection**: Lasso and elastic net conduct simultaneous parameter estimation and variable selection.
- **Highly correlated regressors**:
 - Lasso randomly selects one regressor from the set of highly correlated regressors.
 - Ridge and elastic net tend to perform better than lasso.
- **High dimensional**: They work when $n < p$.

Simulation with $n = 20$ and 5 Covariates

Generalization

- The lasso, ridge, and elastic net estimators are biased. In contrast, SCAD and MCP are supposed to be unbiased if the true values of regression coefficients are large.
- The **group lasso** can group the variables first and penalize the group of variables.
 - When dummy variables are created for qualitative variables, the dummy variables from the same qualitative variable are penalized together.
 - Variable selection is conducted to the qualitative variable, not only any individual dummy variables.
 - In contrast, the lasso may only include some dummy variables from the same variable, which depends on the choice of the reference group.
 - Lasso does not care about the marginality principle either.