

Lecture 24: Variable selection in linear models

Consider linear model $X = Z\beta + \varepsilon$, $\beta \in \mathbb{R}^p$ and $\text{Var}(\varepsilon) = \sigma^2 I_n$.

Like the LSE, the ridge regression estimator does not give 0 estimate to a component of β even if that component is 0.

Variable (or model) selection refers to eliminating covariates (columns of Z) corresponding to zero components of β .

Example 1. Linear regression models

- \mathcal{A} = a subset of $\{1, \dots, p\}$, indices of nonzero components of β
- The dimension of \mathcal{A} is $\dim(\mathcal{A}) = q \leq p$
- $\beta_{\mathcal{A}}$: sub-vector of β with indices in \mathcal{A}
- $Z_{\mathcal{A}}$: the corresponding sub-matrix of Z
- The number of models could be as large as 2^p
- Approximation to a response surface
 - The i th row of $Z_{\mathcal{A}} = (1, t_i, t_i^2, \dots, t_i^h)$, $t_i \in \mathcal{R}$
 - $\mathcal{A} = \{1, \dots, h\}$: a polynomial of order h
 - $h = 0, 1, \dots, p_n$

Example 2. 1-mean vs p -mean

- $n = pr$, $p = p_n$, $r = r_n$
- There are p groups, each has r identically distributed observations
- Select one model from two models
 - 1-mean model: all groups have the same mean μ_1
 - p -mean model: p groups have different means μ_1, \dots, μ_p
- $\mathcal{A} = \mathcal{A}_1$ or \mathcal{A}_p

$$Z = \begin{pmatrix} 1_r & 0 & 0 & \cdots & 0 \\ 1_r & 1_r & 0 & \cdots & 0 \\ \cdots & \cdots & \cdots & \cdots & \cdots \\ 1_r & 0 & 0 & \cdots & 1_r \end{pmatrix} \quad \beta = \begin{pmatrix} \mu_1 \\ \mu_2 - \mu_1 \\ \cdots \\ \mu_p - \mu_1 \end{pmatrix}$$

$$\begin{aligned} Z_{\mathcal{A}_p} &= Z & \beta_{\mathcal{A}_p} &= \beta \\ Z_{\mathcal{A}_1} &= \mathbf{1}_n & \beta_{\mathcal{A}_1} &= \mu_1 \end{aligned}$$

- In traditional studies, p is fixed and n is large, or p/n is small
- In modern applications, both p and n are large, and in some cases $p > n$, $p/n \rightarrow \infty$

Methods for variable selection

- Generalized Information Criterion (GIC)

Put a penalty on the dimension of the parameter: We minimize

$$\|X - Z_{\mathcal{A}}\beta_{\mathcal{A}}\|^2 + \lambda \hat{\sigma}^2 \dim(\beta_{\mathcal{A}}) \quad \text{over } \mathcal{A},$$

to obtain a suitable \mathcal{A} , and then estimate $\beta_{\mathcal{A}}$.

$\hat{\sigma}^2$ is a suitable estimator of the error variance σ^2

- The term $\|X - Z_{\mathcal{A}}\beta_{\mathcal{A}}\|^2$ measures goodness-of-fit of model \mathcal{A} , whereas the term $\lambda \hat{\sigma}^2 \dim(\beta_{\mathcal{A}})$ controls the “size” of \mathcal{A} .
 - If $\lambda_n = 2$, this is the C_p method, and close to the AIC
 - If $\lambda_n = \log n$, this is close to the BIC
- Regularization or penalized optimization simultaneously select variables and estimate θ by minimizing

$$\|X - Z\beta\|^2 + p_{\lambda}(\beta),$$

where $p_{\lambda}(\cdot)$ is a penalty function indexed by the penalty parameter $\lambda \geq 0$, which may depend on n and data.

Zero components of β are estimated as zeros and automatically eliminated.

● Examples of penalty functions

- Ridge regression: $p_\lambda(\beta) = \lambda \|\beta\|^2$;
- LASSO (least absolute shrinkage and selection operator):
 $p_\lambda(\beta) = \lambda \|\beta\|_1 = \lambda \sum_{j=1}^p |\beta_j|$, β_j is the j th component of β ;
- Adaptive LASSO: $p_\lambda(\beta) = \lambda \sum_{j=1}^p \tau_j |\beta_j|$, where τ_j 's are non-negative leverage factors chosen adaptively such that large penalties are used for unimportant β_j 's and small penalties for important ones;
- Elastic net: $p_\lambda(\beta) = \lambda_1 \|\beta\|_1 + \lambda_2 \|\beta\|^2$;
- Minimax concave penalty: $p_\lambda(\beta) = \sum_{j=1}^p (a\lambda - \beta_j)_+ / a$ for some $a > 0$;
- SCAD (smoothly clipped absolute deviation):
 $p_\lambda(\beta) = \sum_{j=1}^p \lambda \left\{ I(\beta_j \leq \lambda) + \frac{(a\lambda - \beta_j)_+}{(a-1)\lambda} I(\beta_j \geq \lambda) \right\}$ for some $a > 2$;
- There are also many modified versions of the previously listed methods.

● Resampling methods

Cross validation, bootstrap

● Thresholding

Compare $\hat{\beta}_j$ with a threshold (may depend on n and data) and eliminate estimates that are smaller than the threshold.

Assessment of variable/model selection procedures

\mathcal{A} = the set containing exactly indices of nonzero components of β

$\widehat{\mathcal{A}}$: a set of variables/model selected based on a selection procedure

The selection procedure is selection consistent if

$$\lim_{n \rightarrow \infty} P(\widehat{\mathcal{A}} = \mathcal{A}) = 1$$

Sometimes the following weaker version of consistency is desired.

Under model \mathcal{A} , $\mu = E(X|Z)$ is estimated by $\widehat{\mu}_{\mathcal{A}} = Z_{\mathcal{A}} \widehat{\beta}_{\mathcal{A}}$

We want to minimize the squared error loss

$$L_n(\mathcal{A}) = n^{-1} \|\mu - \widehat{\mu}_{\mathcal{A}}\|^2 \quad \text{over } \mathcal{A}$$

which is equivalent to minimizing the average prediction error

$$n^{-1} E \left[\|X^* - \widehat{\mu}_{\mathcal{A}}\|^2 \mid X, Z \right] \quad \text{over } \mathcal{A}$$

X^* : a future independent copy of X

The selection procedure is loss consistent if

$$L_n(\widehat{\mathcal{A}}) / L_n(\mathcal{A}) \rightarrow_p 1$$

Consistency of the GIC

Let \mathcal{M} denote a set of indices (model).

If $\mathcal{A} \subset \mathcal{M}$, then \mathcal{M} is a correct model; otherwise, \mathcal{M} is a wrong model.

The loss under model \mathcal{M} is equal to

$$L_n(\mathcal{M}) = \Delta_n(\mathcal{M}) + \varepsilon^\tau H_{\mathcal{M}} \varepsilon / n$$

$H_{\mathcal{M}} = Z_{\mathcal{M}}(Z_{\mathcal{M}}^\tau Z_{\mathcal{M}})^{-1} Z_{\mathcal{M}}^\tau$, $\Delta_n(\mathcal{M}) = \|\mu - H_{\mathcal{M}} \mu\|^2 / n$ (0 if \mathcal{M} is correct)

Let $\Gamma_{n,\lambda}(\mathcal{M}) = n^{-1} [\|X - Z_{\mathcal{M}} \beta_{\mathcal{M}}\|^2 + \lambda \hat{\sigma}^2 \dim(\beta_{\mathcal{M}})]$ (to be minimized)

$$\begin{aligned} \|X - Z_{\mathcal{M}} \beta_{\mathcal{M}}\|^2 &= \|X - H_{\mathcal{M}} X\|^2 = \|\mu - H_{\mathcal{M}} \mu + \varepsilon - H_{\mathcal{M}} \varepsilon\|^2 \\ &= n \Delta_n(\mathcal{M}) + \|\varepsilon\|^2 - \varepsilon^\tau H_{\mathcal{M}} \varepsilon + 2\varepsilon^\tau (I - H_{\mathcal{M}}) \mu \end{aligned}$$

When \mathcal{M} is a wrong model,

$$\begin{aligned} \Gamma_{n,\lambda}(\mathcal{M}) &= \frac{\|\varepsilon\|^2}{n} + \Delta_n(\mathcal{M}) - \frac{\varepsilon^\tau H_{\mathcal{M}} \varepsilon}{n} + \frac{\lambda \hat{\sigma}^2 \dim(\mathcal{M})}{n} + O_P\left(\frac{\Delta_n(\mathcal{M})}{n}\right) \\ &= \frac{\|\varepsilon\|^2}{n} + L_n(\mathcal{M}) + O_P\left(\frac{\lambda \dim(\mathcal{M})}{n}\right) + O_P\left(\frac{L_n(\mathcal{M})}{n}\right) \\ &= \frac{\|\varepsilon\|^2}{n} + L_n(\mathcal{M}) + o_P(L_n(\mathcal{M})) \end{aligned}$$

provided that

$$\liminf_{n \rightarrow \infty} \min_{\mathcal{M} \text{ is wrong}} \Delta_n(\mathcal{M}) > 0 \quad \text{and} \quad \frac{\lambda p}{n} \rightarrow 0$$

(The first condition implies that wrong is always worse than correct)

Among all wrong \mathcal{M} , minimizing $\Gamma_{n,\lambda}(\mathcal{M})$ is asymptotically the same as minimizing $L_n(\mathcal{M})$

Hence, the GIC is loss consistent when all models are wrong

The GIC selects the best wrong model, i.e., the best approximation to a correct model in terms of $\Delta_n(\mathcal{M})$, the leading term in the loss $L_n(\mathcal{M})$

For correct models, however, $\Delta_n(\alpha) = 0$ and $L_n(\mathcal{M}) = \varepsilon^\tau H_{\mathcal{M}} \varepsilon / n$

Correct models are nested, and \mathcal{A} has the smallest dimension and

$$\begin{aligned} \varepsilon^\tau H_{\mathcal{A}} \varepsilon &= \min_{\mathcal{M} \text{ is correct}} \varepsilon^\tau H_{\mathcal{M}} \varepsilon \\ \Gamma_{n,\lambda}(\mathcal{M}) &= \frac{\|\varepsilon\|^2}{n} - \frac{\varepsilon^\tau H_{\mathcal{M}} \varepsilon}{n} + \frac{\lambda \hat{\sigma}^2 \dim(\mathcal{M})}{n} \\ &= \frac{\|\varepsilon\|^2}{n} + L_n(\mathcal{M}) + \frac{\lambda \hat{\sigma}^2 \dim(\mathcal{M})}{n} - \frac{2\varepsilon^\tau H_{\mathcal{M}} \varepsilon}{n} \end{aligned}$$

If $\lambda \rightarrow \infty$, the dominating term in $\Gamma_{n,\lambda}(\mathcal{M})$ is $\lambda \hat{\sigma}^2 \dim(\mathcal{M})/n$.

Among correct models, the GIC selects a model by minimizing $\dim(\mathcal{M})$, i.e., it selects \mathcal{A} .

Combining the results, we showed that the GIC is selection consistent.

On the other hand, if $\lambda = 2$ (the C_p method, AIC), the term

$$\frac{2\hat{\sigma}^2 \dim(\mathcal{M})}{n} - \frac{2\varepsilon^\tau H_{\mathcal{M}} \varepsilon}{n}$$

is of the same order as $L_n(\mathcal{M}) = \varepsilon^\tau H_{\mathcal{M}} \varepsilon / n$ unless $\dim(\mathcal{M}) \rightarrow \infty$ for all but one correct model.

Under some conditions, the GIC with $\lambda = 2$ is loss consistent if and only if there does not exist two correct models with fixed dimensions.

Conclusion

- (1) The GIC with a bounded λ (C_p , AIC) is loss consistent when there is at most one fixed-dimension correct model; otherwise it is inconsistent.
- (2) The GIC with $\lambda \rightarrow \infty$ and $\lambda p/n \rightarrow 0$ (BIC) are selection consistent or loss consistent.

Example 2. 1-mean vs p -mean

\mathcal{A}_1 vs \mathcal{A}_p (always correct)

p_n groups, each with r_n observations

$$\Delta_n(\mathcal{A}_1) = \sum_{j=1}^p (\mu_j - \bar{\mu})^2 / p, \quad \bar{\mu} = \sum_{j=1}^p \mu_j / p$$

$n = p_n r_n \rightarrow \infty$ means that either $p_n \rightarrow \infty$ or $r_n \rightarrow \infty$

1. $p_n = p$ is fixed and $r_n \rightarrow \infty$

- The dimensions of correct models are fixed
- The GIC with $\lambda \rightarrow \infty$ and $\lambda/n \rightarrow 0$ is selection consistent
- The GIC with $\lambda = 2$ is inconsistent

2. $p_n \rightarrow \infty$ and $r_n = r$ is fixed

- Only one correct model has a fixed dimension
- The GIC with $\lambda_n = 2$ is loss consistent
- The GIC with $\lambda \rightarrow \infty$ is inconsistent, because $\lambda p_n / n = \lambda / r \rightarrow \infty$

3. $p_n \rightarrow \infty$ and $r_n \rightarrow \infty$

- Only one correct model has a fixed dimension
- The GIC is selection consistent, provided that $\lambda / r_n \rightarrow 0$

More on the case where $p_n \rightarrow \infty$ and $r_n = r$ is fixed

$$\hat{\sigma}^2 = S(\mathcal{A}_p)/n, \quad S(\mathcal{A}) = \|X - Z_{\mathcal{A}}\beta_{\mathcal{A}}\|^2.$$

It can be shown that

$$L_n(\mathcal{A}_1) = \Delta_n(\mathcal{A}_1) + \bar{e}^2 \rightarrow_p \Delta = \lim_{p \rightarrow \infty} \frac{1}{p} \sum_{j=1}^p \left(\mu_j - \frac{1}{p} \sum_{i=1}^p \mu_i \right)^2$$
$$L_n(\mathcal{A}_p) = \frac{1}{p} \sum_{i=1}^p \bar{e}_i^2 \rightarrow_p \frac{\sigma^2}{r}$$

where e_{ij} 's are iid, $E(e_{ij}) = 0$, $E(e_{ij}^2) = \sigma^2$, $\bar{e}_i = r^{-1} \sum_{j=1}^r e_{ij}$, and $\bar{e} = p^{-1} \sum_{i=1}^p \bar{e}_i$.

Then

$$\frac{L_n(\mathcal{A}_1)}{L_n(\mathcal{A}_p)} \rightarrow_p \frac{r\Delta}{\sigma^2}$$

The one-mean model is better if and only if $r\Delta < \sigma^2$.

The wrong model may be better!

The GIC with $\lambda_n \rightarrow \infty$ minimizes

$$\frac{S(\mathcal{A}_1)}{n} + \frac{\lambda_n}{n} \frac{S(\mathcal{A}_p)}{n-p} \quad \text{and} \quad \frac{S(\mathcal{A}_p)}{n} + \frac{\lambda_n}{r} \frac{S(\mathcal{A}_p)}{n-p}$$

Because

$$\frac{S(\mathcal{A}_1)}{n} = \Delta_n(\mathcal{A}_1) + \frac{1}{n} \sum_{i=1}^p \sum_{j=1}^r (e_{ij} - \bar{e}^2) \rightarrow_p \Delta + \sigma^2$$

$$\frac{S(\mathcal{A}_p)}{n} = \frac{1}{n} \sum_{i=1}^p \sum_{j=1}^r (e_{ij} - \bar{e}_i^2) \rightarrow_p \frac{(r-1)\sigma^2}{r}$$

and $\lambda_n/r \rightarrow \infty$, $P\{\text{GIC with } \lambda \rightarrow \infty \text{ selects } \mathcal{A}_1\} \rightarrow 1$

On the other hand, the C_p (GIC with $\lambda_n = 2$) is loss consistent, because the C_p minimizes

$$\frac{S(\mathcal{A}_1)}{n} + \frac{2}{n} \frac{S(\mathcal{A}_p)}{n-p} \quad \text{and} \quad \frac{S(\mathcal{A}_p)}{n} + \frac{2}{r} \frac{S(\mathcal{A}_p)}{n-p}$$

$$\frac{S(\mathcal{A}_1)}{n} + \frac{2}{n} \frac{S(\mathcal{A}_p)}{n-p} \rightarrow_p \Delta + \sigma^2,$$

$$\frac{S(\mathcal{A}_p)}{n} + \frac{2}{r} \frac{S(\mathcal{A}_p)}{n-p} \rightarrow_p \sigma^2 + \frac{\sigma^2}{r}$$

Asymptotically, the C_p selects \mathcal{A}_1 iff $\Delta < \sigma^2/r$, which is the same as the one-mean model is better.

Variable selection by thresholding

Can we do variable selection using p -values?

Or, can we simply select variables by using the values $\hat{\beta}_j, j = 1, \dots, p$?

Here $\hat{\beta}_j$ is the j th component of $\hat{\beta}$, the least squares estimator of β .

For simplicity, assume that $X|Z \sim N(Z\beta, \sigma^2 I)$.

Then

$$\hat{\beta}_j - \beta_j = \sum_{i=1}^n l_{ij} \varepsilon_i \Big| Z \sim N \left(0, \sigma^2 \sum_{i=1}^n l_{ij}^2 \right)$$

where ε_i and l_{ij} are the i th components of $\varepsilon = X - Z\beta$ and $(Z^\tau Z)^{-1} z_j$ is the j th row of Z

Because

$$1 - \Phi(t) \leq \frac{\sqrt{2\pi}}{t} e^{-t^2/2}, \quad t > 0$$

where Φ is the standard normal cdf,

$$P \left(|\hat{\beta}_j - \beta_j| > t \sqrt{\text{var}(\hat{\beta}_j | Z)} \Big| Z \right) \leq \frac{2\sqrt{2\pi}}{t} e^{-t^2/2}, \quad t > 0$$

Let J_j be the p -vector whose j th component is 1 and other components are 0:

$$l_{ij}^2 = [J_j^\tau (Z^\tau Z)^{-1} z_i]^2 \leq J_j^\tau (Z^\tau Z)^{-1} J_j z_i^\tau (Z^\tau Z)^{-1} z_i$$

$$\sum_{i=1}^n l_{ij}^2 \leq c_j \sum_{i=1}^n z_i^\tau (Z^\tau Z)^{-1} z_i = \rho c_j \leq \rho / \eta_n$$

where c_j is the j th diagonal element of $(Z^\tau Z)^{-1}$ and η_n is the smallest eigenvalue of $Z^\tau Z$.

Thus, for any j ,

$$P\left(|\hat{\beta}_j - \beta_j| > t\sigma\sqrt{\rho/\eta_n} \mid Z\right) \leq \frac{2\sqrt{2\pi}}{t} e^{-t^2/2}, \quad t > 0$$

and (letting $t = a_n/(\sigma\sqrt{\rho/\eta_n})$)

$$P\left(|\hat{\beta}_j - \beta_j| > a_n \mid Z\right) \leq C e^{-a_n^2 \eta_n / (2\sigma^2 \rho)}$$

for some constant $C > 0$,

$$P\left(\max_{j=1, \dots, p} |\hat{\beta}_j - \beta_j| > a_n \mid Z\right) \leq p C e^{-a_n^2 \eta_n / (2\sigma^2 \rho)}$$

Suppose that $p/n \rightarrow 0$ and $p/(\eta_n \log n) \rightarrow 0$ (typically, $\eta_n = O(n)$).

Then, we can choose a_n such that $a_n \rightarrow 0$ and $a_n^2(\eta_n \log n/p) \rightarrow \infty$ such that

$$P\left(\max_{j=1,\dots,p} |\hat{\beta}_j - \beta_j| > ca_n \mid Z\right) = O(n^{-s})$$

for any $c > 0$ and some $s \geq 1$; e.g.,

$$a_n = M \left(\frac{p}{\eta_n \log n} \right)^\alpha$$

for some constants $M > 0$ and $\alpha \in (0, \frac{1}{2})$.

What can we conclude from this?

Let $\mathcal{A} = \{j : \beta_j \neq 0\}$ and $\widehat{\mathcal{A}} = \{j : |\hat{\beta}_j| > a_n\}$

That is, $\widehat{\mathcal{A}}$ contains the indices of variables we select by thresholding $|\hat{\beta}_j|$ at a_n .

Selection consistency:

$$P\left(\widehat{\mathcal{A}} \neq \mathcal{A} \mid Z\right) \leq P\left(|\hat{\beta}_j| > a_n, j \notin \mathcal{A} \mid Z\right) + P\left(|\hat{\beta}_j| \leq a_n, j \in \mathcal{A} \mid Z\right)$$

The first term on the right hand side is bounded by

$$P\left(\max_{j=1,\dots,p} |\hat{\beta}_j - \beta_j| > a_n \mid Z\right) = O(n^{-s})$$

On the other hand, if we assume that

$$\min_{j \in \mathcal{A}} |\beta_j| \geq c_0 a_n$$

for some $c_0 > 1$, then

$$\begin{aligned} P\left(|\hat{\beta}_j| \leq a_n, j \in \mathcal{A} \mid \mathbf{Z}\right) &\leq P\left(|\beta_j| - |\hat{\beta}_j - \beta_j| \leq a_n, j \in \mathcal{A} \mid \mathbf{Z}\right) \\ &\leq P\left(c_0 a_n - |\hat{\beta}_j - \beta_j| \leq a_n, j \in \mathcal{A} \mid \mathbf{Z}\right) \\ &\leq P\left(\max_{j=1, \dots, p} |\hat{\beta}_j - \beta_j| \geq (c_0 - 1)a_n \mid \mathbf{Z}\right) \\ &= O(n^{-s}) \end{aligned}$$

Hence, we have consistency; in fact, the convergence rate is $O(n^{-s})$.

We can also obtain similar results by thresholding $|\hat{\beta}_j| / \sqrt{\sum_{i=1}^n l_{ij}^2}$.

This approach may not work if $p/n \not\rightarrow 0$.

If $p > n$, then $Z^T Z$ is not of full rank.

There exist several other approaches for the case where $p > n$; e.g., we replace $(Z^T Z)^{-1}$ by some matrix, or use ridge regression instead of LSE.