

VE472 Lecture 5

Jing Liu

UM-SJTU Joint Institute

Summer

- It can be shown that the least squares estimator is unbiased and consistent, furthermore, it is the best in the sense described by the following theorem.

Theorem 0.1 (Gauss-Markov Theorem)

Given $\mathbb{E}[\varepsilon] = 0$ and $\text{Var}[\varepsilon] = \sigma^2 \mathbf{I}$, then the least squares estimator, if it exists,

$$\hat{\beta}_{\text{lse}} = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{y}$$

is the best linear unbiased estimator of β in the sense of minimum variance for

$$\mathbf{y} = \mathbf{X}\beta + \varepsilon$$

also unbiased

that is, for any estimator of the form $\tilde{\beta} = \mathbf{A}_{(k+1) \times n} \mathbf{y}$ such that $\mathbb{E}[\tilde{\beta}] = \beta$, then

$$\text{Var}[\alpha^T \tilde{\beta}] \geq \text{Var}[\alpha^T \hat{\beta}_{\text{lse}}] \quad \text{where } \alpha \in \mathbb{R}^{k+1}$$

smallest variance

- This is one of the main reasons why $\hat{\beta}_{\text{lse}}$ is widely used for small data.

Q: What happens if \mathbf{X} is not full rank?

Q: Will the last theorem still hold when \mathbf{X} is not full rank? no

Q: Why do we expect the least squares estimate is no longer unique?

$$\arg \min_{\mathbf{b} \in \mathbb{R}^{k+1}} \|\mathbf{y} - \mathbf{X}\mathbf{b}\|^2$$

一种方法: $\min \|\mathbf{b}\|$ (L4)
minimum norm solution:
eliminate unnecessary variables

- In this situation, we say β is non-identifiable, and additional constraint(s) need to be introduced to reach a unique estimate of β .
- The matrix \mathbf{X} being less than full rank is not specific to big data but it does become more prominent when \mathbf{X} becomes more complex in modern era.
- The simplest type of constraint is to set some of the coefficients β_j to zero, that is, we select or exclude some variables from the k independent variables.
- Variable selection is traditionally done according to some criterion, e.g.

AIC/BIC

- When building a **predictive model**, variable selection is done according to

$$\text{MSE}(\hat{Y}_i) = \mathbb{E} \left[\left(\hat{Y}_i - Y_i \right)^2 \right]$$

which is different from the MSE of an estimator $\hat{\theta}$ of $\theta \in \mathbb{R}$

$$\text{MSE}(\hat{\theta}) = \mathbb{E} \left[\left(\hat{\theta} - \theta \right)^2 \right] = \text{Var}[\hat{\theta}] + \underbrace{\left(\mathbb{E}[\hat{\theta}] - \theta \right)^2}_{\text{Bias}}$$

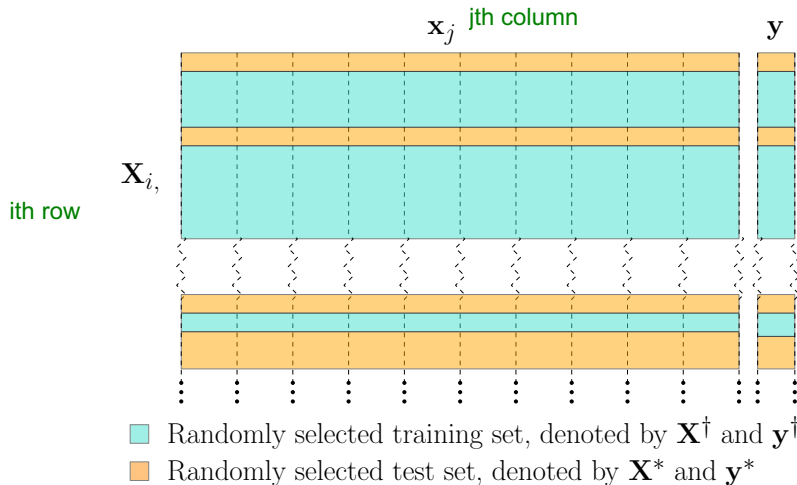
since $\text{MSE}(\hat{Y}_i)$, a.k.a. **prediction error**, involves two random variables.

- In the context of regression analysis, it can be shown that

$$\text{MSE}(\hat{Y}_i) = \mathbb{E} \left[\left(\hat{Y}_i - Y_i \right)^2 \right] = \text{Var}[\hat{Y}_i] + \underbrace{\left(\mathbb{E}[\hat{Y}_i] - \mathbb{E}[Y_i] \right)^2}_{\text{Bias}} + \sigma^2$$

- In practice, we have to estimate $\text{MSE}(\hat{Y}_i)$, e.g. via simple training-test split.

- Simple training-test split is widely used when n is sufficiently large.



- Let $n - m$ and m denote the number of cases in the training and test set.

- Given a simple training-test split, we "train" various linear models

$$\mathbf{y}^\dagger = \mathbf{X}_\ell^\dagger \mathbf{b}_\ell + \boldsymbol{\varepsilon}_\ell$$

where only some of the k independent variables, say $p \leq k$, are included, e.g.

$$y_i = \beta_0 + \beta_1 x_{i1} + 0x_{i2} + \cdots + \beta_{k-1} x_{i,k-1} + 0x_{ik} + e_i$$

- The MSE of the ℓ th model is estimated by computing the followings

$$\mathbf{b}_\ell = \arg \min_{\mathbf{b} \in \mathbb{R}^{p+1}} \left\| \mathbf{y}^\dagger - \mathbf{X}_\ell^\dagger \mathbf{b} \right\|^2$$

$$\hat{\text{MSE}} = \frac{1}{m} \left\| \mathbf{y}^* - \mathbf{X}_\ell^* \mathbf{b}_\ell \right\|^2 \quad \text{.calculate using test set}$$

beta_k = 0 / 1

- Note there are 2^k models in total, and each \mathbf{b}_ℓ defines just one of them, the model with the smallest $\hat{\text{MSE}}$ is considered to be the best predictive model.

problem: might have too many models

Q: Is there any way to reduce MSE without conducting variable selection?

- Note the following two aspects of this approach when comes to big datasets:

1. When the dataset is complex, i.e. a **big k value**, the number of models

$$2^k$$

grows really quickly, and significantly outstrip computing power we have.

2. When the dataset is large, i.e., **as n grows**, **$\hat{\beta}_{\text{lse}}$ might not be the best**

$$\text{MSE}(\hat{Y}_i) = \mathbb{E}\left[\left(\hat{Y}_i - Y_i\right)^2\right] = \text{Var}[\hat{Y}_i] + \left(\mathbb{E}[\hat{Y}_i] - \mathbb{E}[Y_i]\right)^2 + \sigma^2$$

tolerate some bias, reduce variance

- The Gauss-Markov theorem **only** guarantees this approach would lead us to the minimum MSE estimators of β amongst all **linear unbiased** estimators.

$$\hat{Y}_i = \hat{\beta}_0 + \hat{\beta}_1 x_{i1} + \hat{\beta}_2 x_{i2} + \cdots + \hat{\beta}_k x_{ik} = \mathbf{x}_i^T \hat{\beta}$$

where \mathbf{x}_i^T denotes the row vector $\begin{bmatrix} 1 & x_{i1} & \cdots & x_{ik} \end{bmatrix}$ that we predict Y_i with

- However, unbiasedness is not relevant if the dataset is large enough, being consistent will guarantee the quality of the prediction for sufficiently large n .
- In other words, for a large dataset, we can tolerate some bias as long as it is consistent, so we should find a predictive model that minimises the estimated

$$\text{MSE}(\hat{Y}_i) = \mathbb{E} \left[\left(\hat{Y}_i - Y_i \right)^2 \right] = \text{Var}[\hat{Y}_i] + \left(\mathbb{E}[\hat{Y}_i] - \mathbb{E}[Y_i] \right)^2 + \sigma^2$$

without restricting ourselves to unbiased estimators to form a predictor.

- Just like you have seen earlier, depending on the size of dataset, n , there might be a significant improvement in $\hat{\text{MSE}}$ if we give up unbiasedness.
- **Shrinkage methods**, which assign the importance of each X_j in predicting y

$$\mathbb{E}[Y_i] = \mathbf{x}_{i,\cdot}^T \boldsymbol{\beta} = \beta_0 + \beta_1 x_{i1} + \beta_2 x_{i2} + \cdots + \beta_k x_{ik}$$

as well as estimating the coefficient β_j , offer ways to avoid variable selection as well as estimators which lead to smaller $\hat{\text{MSE}}$ s by having some small bias.

Q: Where can we get some bias?

- Recall the following property of variance for fixed scalars w_1 , α_1 , w_2 and α_2

$$\begin{aligned}\text{Var} [w_1 Z_1 \alpha_1 + w_2 Z_2 \alpha_2] &= w_1^2 \text{Var} [Z_1] \alpha_1^2 + w_2^2 \text{Var} [Z_2] \alpha_2^2 \\ &\quad + 2w_1 w_2 \text{Cov} [Z_1, Z_2] \alpha_1 \alpha_2\end{aligned}$$

Z_1 Z_2 如果 not correlated, $\text{Cov} = 0$

where Z_1 and Z_2 are two arbitrary random variables with finite variance.

- Notice the variance shrinks to 0 as the "size" of w_1 and w_2 shrinks.
- Given the linear model, $\mathbf{y} = \mathbf{X}\boldsymbol{\beta} + \boldsymbol{\varepsilon}$, the variance of following predictor,

\mathbf{w} : weight vector

$$\hat{Y}_i = w_0 \hat{\gamma}_0 1 + w_1 \hat{\gamma}_1 x_{i1} + w_2 \hat{\gamma}_2 x_{i2} + \cdots + w_k \hat{\gamma}_k x_{ik}$$

which will be used as a prediction for Y_i given x_{ij} values, shrinks to zero in a similar fashion to the case above if the ℓ_2 norm of \mathbf{w} shrinks to zero.

- Conceptually, \mathbf{w} can be thought as a vector of some kind of weighting factor for each X_j in predicting Y_i , and $\hat{\gamma}$ is some estimator of $\boldsymbol{\beta}$ with respect to \mathbf{w} .

- In practice, we do not separate w_j from $\hat{\gamma}_j$, for $j = 0, \dots, k$,

$$\hat{Y}_i = w_0 \hat{\gamma}_0 1 + w_1 \hat{\gamma}_1 x_{i1} + w_2 \hat{\gamma}_2 x_{i2} + \dots + w_k \hat{\gamma}_k x_{ik}$$

- Instead of explicitly introducing \mathbf{w} , we form an estimator $\hat{\beta}$ that takes $\|\hat{\beta}\|$ into account for shrinking it shrinks the variance, and use the linear predictor

$$\hat{Y}_i = \mathbf{x}_i^T \hat{\beta}$$

- One natural approach is still to minimise $\|\mathbf{y} - \mathbf{X}\mathbf{b}\|^2$ with respect to \mathbf{b} , but with the constraint $\|\mathbf{b}\|^2 \leq c_\lambda$, where c_λ defines the level of shrinkage.
- This approach is known as **ridge regression**,

$$\begin{aligned} & \arg \min_{\mathbf{b} \in \mathbb{R}^{k+1}} \|\mathbf{y} - \mathbf{X}\mathbf{b}\|^2 \\ & \text{s. t.} \quad \|\mathbf{b}\|^2 \leq c_\lambda \end{aligned}$$

and its solution is known as the **ridge estimate**, often denoted by **$\mathbf{b}_{\text{ridge}}$** .

Theorem 0.2

The constrained minimisation problem in ridge regression for some $c_\lambda > 0$

$$\begin{aligned} \arg \min_{\mathbf{b} \in \mathbb{R}^{k+1}} & \|\mathbf{y} - \mathbf{X}\mathbf{b}\|^2 \\ \text{s. t.} & \quad \|\mathbf{b}\|^2 \leq c_\lambda \end{aligned}$$

can be converted into the following penalised/regularised least squares problem

$$\arg \min_{\mathbf{b} \in \mathbb{R}^{k+1}} \left\{ \|\mathbf{y} - \mathbf{X}\mathbf{b}\|^2 + \lambda \|\mathbf{b}\|^2 \right\} \quad \text{where } 0 \leq \lambda < \infty$$

$\|\mathbf{b}\|$ 越小, variance 越小, 同时不能让 bias 太大

where $0 \leq \lambda < \infty$ depends on c_λ , and the solution of this problem is given by

$$\mathbf{b}_{\text{ridge}} = (\mathbf{X}^T \mathbf{X} + \lambda \mathbf{I})^{-1} \mathbf{X}^T \mathbf{y}$$

$\mathbf{X}^T \mathbf{X} + \lambda \mathbf{I}$ is always invertible

- Recall our motivation of using $\mathbf{b}_{\text{ridge}}$ is to avoid variable selection, and attain a smaller $\text{MSE}(\hat{Y})$ by giving up unbiasedness and relying on consistency.

Theorem 0.3

Given $\mathbb{E}[\boldsymbol{\varepsilon}] = 0$ and $\text{Var}[\boldsymbol{\varepsilon}] = \sigma^2 \mathbf{I}$, then the ridge estimator,

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

for the unknown parameter $\boldsymbol{\beta}$ in the linear model,

$$\mathbf{y} = \mathbf{X}\boldsymbol{\beta} + \boldsymbol{\varepsilon}$$

is biased and the bias is given by

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

but $\hat{\boldsymbol{\beta}}_{\text{ridge}}$ is consistent in the following sense,

$$\lim_{n \rightarrow \infty} \Pr \left[\left\| \hat{\boldsymbol{\beta}}_{\text{ridge}} - \boldsymbol{\beta} \right\| \geq \delta \right] = 0 \quad \text{for all } \delta > 0$$

where n denotes the number of cases in the dataset.

- Note all it does in comparison to $\mathbf{b}_{\text{lse}} = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{y}$ is that it creates a "ridge" down the diagonal of the matrix that needs to be "inverted"

$$\mathbf{b}_{\text{ridge}} = (\mathbf{X}^T \mathbf{X} + \lambda \mathbf{I})^{-1} \mathbf{X}^T \mathbf{y}$$

- Of course, this distorts the original least squares estimate depending on λ

$$\arg \min_{\mathbf{b} \in \mathbb{R}^{k+1}} \left\{ \|\mathbf{y} - \mathbf{X}\mathbf{b}\|^2 + \lambda \|\mathbf{b}\|^2 \right\} \quad \text{where } 0 < \lambda < \infty$$

- It is clear that we have the following limiting cases of the ridge estimate
lse is special case of ridge

$$\begin{array}{lll} \mathbf{b}_{\text{ridge}} \rightarrow \mathbf{b}_{\text{lse}} & \text{as} & \lambda \rightarrow 0 \\ \mathbf{b}_{\text{ridge}} \rightarrow \mathbf{0} & \text{as} & \lambda \rightarrow \infty \end{array}$$

lambda太大 bias太大

- The two limiting cases represent the zero bias and the zero variance solution.

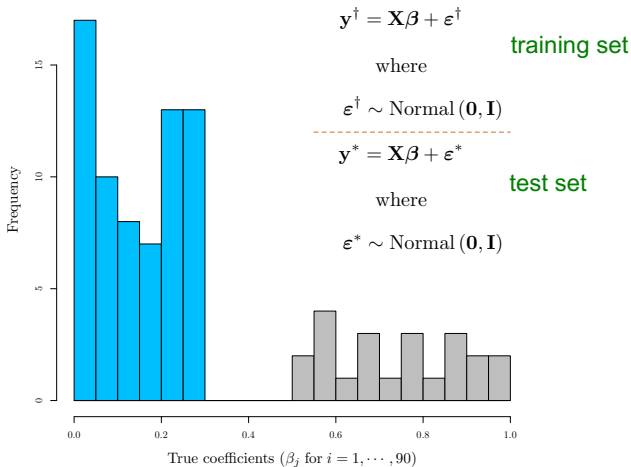
$$\text{MSE}(\hat{Y}_i) = \mathbb{E} \left[\left(\hat{Y}_i - Y_i \right)^2 \right] = \text{Var}[\hat{Y}_i] + \left(\mathbb{E}[\hat{Y}_i] - \mathbb{E}[Y_i] \right)^2 + \sigma^2$$

- It is reasonable to expect there is some optimal λ that achieves the balance.

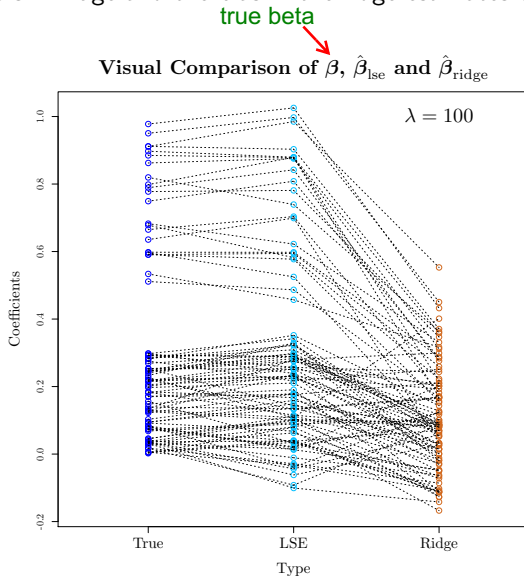
- Training and test sets are simulated with the following true coefficients.

Training and test sets based on simulated β

True coefficients are randomly generated, 22 large ones and 68 small ones



- Notice the shrinkage and the bias in the ridge estimates of the coefficients.

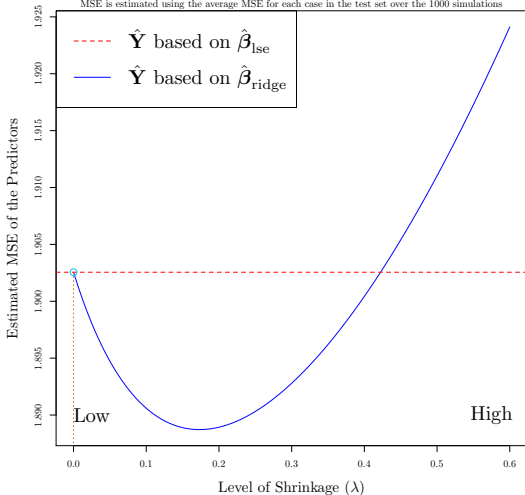


big coefficients
are pushed down
(penalty on norm
of b)

- As we expected, there is an optimal λ according to the 1000 simulations.

Simulation Study of Predictors \hat{Y}

MSE is estimated using the average MSE for each case in the test set over the 1000 simulations

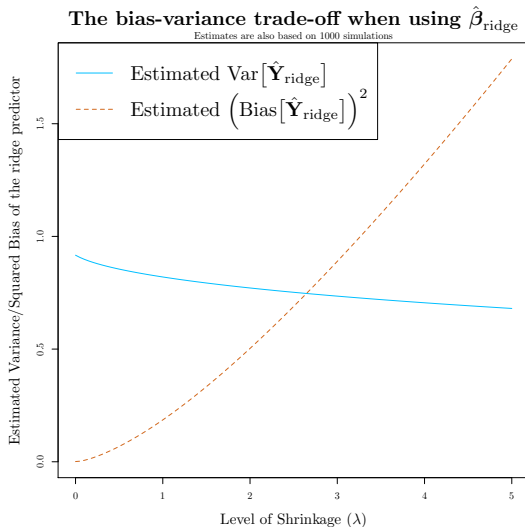


找best lambda
的方法只有试

use b_{lse}

use ridge

- As λ increases, the variance shrinks but the squared bias increases sharply.



bias变大很快，
variance慢慢变小

- The simulation result is a direct consequence of the following theorem.

Theorem 0.4 reason of the previous slide

Let $\hat{\theta}_\ell$ for $\ell = 1, 2$ denote two distinct estimators of vector θ with second order moments and generalised mean square errors:

$$\mathbf{M}_\ell = \mathbb{E} \left[\left(\hat{\theta}_\ell - \theta \right) \left(\hat{\theta}_\ell - \theta \right)^T \right] \quad \text{and} \quad \text{GMSE} \left[\hat{\theta}_\ell \right] = \mathbb{E} \left[\left(\hat{\theta}_\ell - \theta \right)^T \mathbf{A} \left(\hat{\theta}_\ell - \theta \right) \right]$$

generalized MSE

where $\mathbf{A} \succeq 0$, that is, \mathbf{A} is a positive semi-definite matrix, then $\mathbf{M}_1 - \mathbf{M}_2 \succeq 0$ if and only if

$$\text{GMSE} \left[\hat{\theta}_1 \right] - \text{GMSE} \left[\hat{\theta}_2 \right] \geq 0 \quad \text{for all} \quad \mathbf{A} \succeq 0 \quad \text{positive definite matrix}$$

If β is fixed but unknown and $\mathbf{y} = \mathbf{X}\beta + \epsilon$, where $\mathbb{E}[\epsilon] = 0$ and $\text{Var}[\epsilon] = \sigma^2 \mathbf{I}$, then there is an optimal $\lambda_{\text{opt}} > 0$ so that

$$\mathbf{M}_{\text{lse}} - \mathbf{M}_{\text{ridge}} \succ 0 \quad \text{and} \quad \boxed{\text{GMSE} \left[\hat{\beta}_{\text{ridge}}(\lambda_{\text{opt}}) \right] < \text{GMSE} \left[\hat{\beta}_{\text{lse}} \right]} \quad \text{for all} \quad \mathbf{A} \succeq 0$$

where $\hat{\beta}_{\text{lse}} = \left(\mathbf{X}^T \mathbf{X} \right)^{-1} \mathbf{X}^T \mathbf{y}$ and $\hat{\beta}_{\text{ridge}}(\lambda) = \left(\mathbf{X}^T \mathbf{X} + \lambda \mathbf{I} \right)^{-1} \mathbf{X}^T \mathbf{y}$.

- Notice the last Theorem only guarantees the existence of the optimal λ_{opt} .
- Suppose the data is big in the sense that k is large, but n is relatively small.

Q: How to determine the shrinkage parameter λ_{opt} ?

- Given a dataset $\{\mathbf{X}_{n \times (k+1)}, \mathbf{y}_{n \times 1}\}$, and a simple training-test split

$$\{\mathbf{X}_{(n-m) \times (k+1)}^\dagger, \mathbf{X}_{m \times (k+1)}^*, \mathbf{y}_{(n-m) \times (k+1)}^\dagger, \mathbf{y}_{m \times (k+1)}^*\}$$

for any value of $0 \leq \lambda < \infty$, we can estimate $\text{MSE}(\hat{Y}_i)$ for arbitrary i by

试出来

$$\mathbf{b}_{\text{ridge}} = \arg \min_{\mathbf{b} \in \mathbb{R}^{k+1}} \left\{ \|\mathbf{y}^\dagger - \mathbf{X}^\dagger \mathbf{b}\|^2 + \lambda \|\mathbf{b}\|^2 \right\}$$

$$\hat{\text{MSE}} = \frac{1}{m} \|\mathbf{y}^* - \mathbf{X}^* \mathbf{b}_{\text{ridge}}\|^2$$

- Of course, we can not only do this for a value of λ , but a set of values of λ , thus choose the λ that numerically minimises $\hat{\text{MSE}}$ to be λ_{opt} .