

Element of Statistical Learning Note

Zichong Wang

December 31, 2025

Contents

1 Chap 3: Linear Methods for Regression	2
1.1 Confidence region for $\hat{\beta}$	2
1.2 What is Linear	2
1.3 Gauss-Markov Theorem	2
1.4 QR decomposition	3
1.4.1 Application to Least Squared	3
1.4.2 About orthogonal matrix	4
1.4.3 SVD and orthogonal matrix	5
1.5 Multiple testing in forward selection	6
1.6 Ridge regression	6
1.6.1 How to use ridge regression	6
1.6.2 Equivalence of two forms	6
1.6.3 Bayesian view	8
1.6.4 Why avoid singular problem	8

1 Chap 3: Linear Methods for Regression

1.1 Confidence region for $\hat{\beta}$

We know that $\hat{\beta} \sim \mathcal{N}(\beta, \sigma^2(\mathbf{X}^T \mathbf{X})^{-1})$, thus $\hat{\beta} - \beta \sim \mathcal{N}(0, \sigma^2(\mathbf{X}^T \mathbf{X})^{-1})$.

Since for $\mathbf{z} \sim \mathcal{N}(0, \Sigma)$, we have $\mathbf{z}^T \Sigma^{-1} \mathbf{z} \sim \chi_k^2$, where $k = \text{rank}(\Sigma)$, we have

$$(\hat{\beta} - \beta)^T (\sigma^2(\mathbf{X}^T \mathbf{X})^{-1})^{-1} (\hat{\beta} - \beta) \sim \chi_{p+1}^2$$

Thus,

$$\frac{1}{\sigma^2} (\hat{\beta} - \beta)^T \mathbf{X}^T \mathbf{X} (\hat{\beta} - \beta) \sim \chi_{p+1}^2$$

And have approx confidence region

$$(\hat{\beta} - \beta)^T \mathbf{X}^T \mathbf{X} (\hat{\beta} - \beta) \leq \hat{\sigma}^2 \chi_{p+1, 1-\alpha}^2$$

Actually, $\hat{\sigma}^2 = \frac{1}{N-p-1} \sum_{i=1}^N (y_i - \hat{y}_i)^2 = \frac{1}{N-p-1} \text{RSS}$, and $\text{RSS}/\sigma^2 \sim \chi_{N-p-1}^2$,

$$\frac{(\hat{\beta} - \beta)^T \mathbf{X}^T \mathbf{X} (\hat{\beta} - \beta)/(p+1)}{\hat{\sigma}^2} \sim F_{p+1, N-p-1}$$

1.2 What is Linear

In the context of **linear model**, we are talking about linearity in parameters, meaning that the prediction \hat{y} is a linear combination of the parameters β_j . The \mathbf{X} itself can be non-linear transformations of the original features, e.g., polynomial terms, interaction terms, etc. $y = 1/(\beta_0 + \beta_1 x)$ and $y = \beta_0 e^{\beta_1 x}$ are not linear models, since they're not linear in parameters.

In the context of **Linear estimators**, we are talking about the estimator $\hat{\theta}$ (e.g. $\hat{\beta}$) can be written as a linear combination of the observed response values y_i , i.e. $\hat{\theta} = \mathbf{c}^T \mathbf{y}$. The weight \mathbf{c} depends only on \mathbf{X} , not on \mathbf{y} . A linear estimator **can be** a prediction at a new point, or the estimated coefficients $\hat{\beta}$ themselves.

1.3 Gauss-Markov Theorem

Why assume only know \mathbf{X} , but not \mathbf{y} ?

Note that though y_i as sample responses, are observable, the following statements and arguments including assumptions, proofs and the others assume under the only condition of knowing $\mathbf{X}_{i,j}$ but not y_i . — [1]

We have a *challenger* linear estimator $\tilde{\beta} = \mathbf{C} \mathbf{y}$, where $\mathbf{C} = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T + \mathbf{D}$, a modification

of OLS estimator. Ensure it's unbiased:

$$\begin{aligned}\mathbb{E}(\tilde{\beta}) &= \mathbb{E}(C\mathbf{y}) \\ &= C\mathbb{E}(\mathbf{y}) = ((\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T + D)\mathbf{X}\beta \\ &= \beta + D\mathbf{X}\beta \\ &= \beta\end{aligned}$$

Meaning $D\mathbf{X} = 0$.

Now, compute the variance:

$$\begin{aligned}\text{Var}(\tilde{\beta}) &= \text{Var}(C\mathbf{y}) \\ &= C\text{Var}(\mathbf{y})C^T \\ &= \sigma^2 CC^T \\ &= \sigma^2((\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T + D)((\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T + D)^T \\ &= \sigma^2((\mathbf{X}^T \mathbf{X})^{-1} + D\mathbf{D}^T) \\ &= \text{Var}(\hat{\beta}) + \sigma^2 D\mathbf{D}^T \\ &\geq \text{Var}(\hat{\beta})\end{aligned}$$

1.4 QR decomposition

Any real squared matrix \mathbf{A} can be decomposed as $\mathbf{A} = \mathbf{Q}\mathbf{R}$, where \mathbf{Q} is orthogonal matrix ($\mathbf{Q}^T \mathbf{Q} = \mathbf{I}$), and \mathbf{R} is upper triangular matrix.

Any rectangular matrix $\mathbf{A} \in \mathbb{R}^{m \times n}$ ($m \geq n$), we can decompose it as $\mathbf{A} = \mathbf{Q}\mathbf{R}$, where \mathbf{Q} is $m \times m$ orthogonal matrix, and \mathbf{R} is $m \times n$ upper triangular matrix (the last $m - n$ rows are all zero). It can be regarded as $\mathbf{A} = \mathbf{Q}\mathbf{R} = [\mathbf{Q}_1 \quad \mathbf{Q}_2] \begin{bmatrix} \mathbf{R}_1 \\ \mathbf{0} \end{bmatrix} = \mathbf{Q}_1 \mathbf{R}_1$, where $\mathbf{Q}_1 \in \mathbb{R}^{m \times n}$, $\mathbf{Q}_2 \in \mathbb{R}^{m \times (m-n)}$, $\mathbf{R}_1 \in \mathbb{R}^{n \times n}$ which is upper triangular.

If \mathbf{A} have k linearly independent columns, then first k columns of \mathbf{Q} form an orthonormal basis of the column space of \mathbf{A} . The fact that any column k of \mathbf{A} only depends on the first k columns of \mathbf{Q} corresponds to the triangular form of \mathbf{R} .

QR decomposition can be calculated using Gram-Schmidt process, or using Householder reflections. In practice, Householder reflections are more stable and efficient.

1.4.1 Application to Least Squared

In linear least squares problems, we aim to find a vector \mathbf{x} that minimizes the Euclidean norm of the residual for an overdetermined system $\mathbf{Ax} \approx \mathbf{b}$, where $\mathbf{A} \in \mathbb{R}^{m \times n}$ and $m \geq n$. The goal is to solve:

$$\min_{\mathbf{x}} \|\mathbf{Ax} - \mathbf{b}\|_2$$

Using the Full QR Decomposition, we substitute $\mathbf{A} = \mathbf{Q} \begin{bmatrix} \mathbf{R}_1 \\ \mathbf{0} \end{bmatrix}$:

$$\|\mathbf{Ax} - \mathbf{b}\|_2 = \left\| \mathbf{Q} \begin{bmatrix} \mathbf{R}_1 \\ \mathbf{0} \end{bmatrix} \mathbf{x} - \mathbf{b} \right\|_2$$

Since multiplying by an orthogonal matrix \mathbf{Q} preserves the Euclidean norm, we can left-multiply the entire expression by \mathbf{Q}^T :

$$\|\mathbf{A}\mathbf{x} - \mathbf{b}\|_2 = \left\| \begin{bmatrix} \mathbf{R}_1 \\ \mathbf{0} \end{bmatrix} \mathbf{x} - \mathbf{Q}^T \mathbf{b} \right\|_2$$

If we partition $\mathbf{Q}^T \mathbf{b}$ into two components— $\mathbf{c}_1 \in \mathbb{R}^n$ and $\mathbf{c}_2 \in \mathbb{R}^{m-n}$:

$$\left\| \begin{bmatrix} \mathbf{R}_1 \mathbf{x} - \mathbf{c}_1 \\ -\mathbf{c}_2 \end{bmatrix} \right\|_2^2 = \|\mathbf{R}_1 \mathbf{x} - \mathbf{c}_1\|_2^2 + \|\mathbf{c}_2\|_2^2$$

To minimize the total error, we must make the first term zero. The least squares solution is found by solving the square, upper-triangular system:

$$\mathbf{R}_1 \mathbf{x} = \mathbf{c}_1$$

Since \mathbf{R}_1 is upper triangular, we can efficiently solve this system using back substitution, **that's how we solve linear equations manually in algebra class!** The remaining term $\|\mathbf{c}_2\|_2$ represents the minimum residual norm (the "error" of the fit).

1.4.2 About orthogonal matrix

Let $\mathbf{Q} \in \mathbb{R}^{n \times n}$ be an orthogonal matrix. Write \mathbf{Q} in terms of its column vectors:

$$\mathbf{Q} = [\mathbf{q}_1 \ \mathbf{q}_2 \ \cdots \ \mathbf{q}_n], \quad \mathbf{q}_i^\top \mathbf{q}_j = \delta_{ij}.$$

Then

$$\mathbf{Q}^\top = \begin{bmatrix} \mathbf{q}_1^\top \\ \mathbf{q}_2^\top \\ \vdots \\ \mathbf{q}_n^\top \end{bmatrix}.$$

Computation of $\mathbf{Q}^\top \mathbf{Q}$.

$$\mathbf{Q}^\top \mathbf{Q} = \begin{bmatrix} \mathbf{q}_1^\top \\ \mathbf{q}_2^\top \\ \vdots \\ \mathbf{q}_n^\top \end{bmatrix} [\mathbf{q}_1 \ \mathbf{q}_2 \ \cdots \ \mathbf{q}_n] = \begin{bmatrix} \mathbf{q}_1^\top \mathbf{q}_1 & \mathbf{q}_1^\top \mathbf{q}_2 & \cdots & \mathbf{q}_1^\top \mathbf{q}_n \\ \mathbf{q}_2^\top \mathbf{q}_1 & \mathbf{q}_2^\top \mathbf{q}_2 & \cdots & \mathbf{q}_2^\top \mathbf{q}_n \\ \vdots & \vdots & \ddots & \vdots \\ \mathbf{q}_n^\top \mathbf{q}_1 & \mathbf{q}_n^\top \mathbf{q}_2 & \cdots & \mathbf{q}_n^\top \mathbf{q}_n \end{bmatrix}.$$

Using $\mathbf{q}_i^\top \mathbf{q}_j = \delta_{ij}$,

$$\mathbf{Q}^\top \mathbf{Q} = \begin{bmatrix} 1 & 0 & \cdots & 0 \\ 0 & 1 & \cdots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \cdots & 1 \end{bmatrix} = \mathbf{I}.$$

Computation of $\mathbf{Q}\mathbf{Q}^\top$.

$$\mathbf{Q}\mathbf{Q}^\top = [\mathbf{q}_1 \ \mathbf{q}_2 \ \cdots \ \mathbf{q}_n] \begin{bmatrix} \mathbf{q}_1^\top \\ \mathbf{q}_2^\top \\ \vdots \\ \mathbf{q}_n^\top \end{bmatrix} = \mathbf{q}_1 \mathbf{q}_1^\top + \mathbf{q}_2 \mathbf{q}_2^\top + \cdots + \mathbf{q}_n \mathbf{q}_n^\top = \sum_{k=1}^n \mathbf{q}_k \mathbf{q}_k^\top.$$

For any $\mathbf{x} \in \mathbb{R}^n$,

$$(\mathbf{Q}\mathbf{Q}^\top)\mathbf{x} = \sum_{k=1}^n \mathbf{q}_k(\mathbf{q}_k^\top \mathbf{x}) = \mathbf{x},$$

since $\{\mathbf{q}_1, \dots, \mathbf{q}_n\}$ is an orthonormal basis of \mathbb{R}^n . Hence

$$\mathbf{Q}\mathbf{Q}^\top = \mathbf{I}.$$

Therefore, for an orthogonal matrix \mathbf{Q} ,

$$\mathbf{Q}^\top \mathbf{Q} = \mathbf{Q}\mathbf{Q}^\top = \mathbf{I}.$$

1.4.3 SVD and orthogonal matrix

There are two kinds of SVD: full SVD and reduced SVD. Let $\mathbf{X} \in \mathbb{R}^{N \times p}$, and normally $N > p$.

- **Full SVD** in standard linear algebra.

$$\mathbf{X} = \mathbf{U}_{\text{full}} \mathbf{D}_{\text{full}} \mathbf{V}^T$$

- $\mathbf{U}_{\text{full}} \in \mathbb{R}^{N \times N}$ is orthogonal matrix, whose columns are eigenvectors of $\mathbf{X}\mathbf{X}^T$.
- $\mathbf{D}_{\text{full}} \in \mathbb{R}^{N \times p}$ is a rectangular diagonal matrix, and the last $N - p$ rows are all zero.
- $\mathbf{V} \in \mathbb{R}^{p \times p}$ is orthogonal matrix, whose columns are eigenvectors of $\mathbf{X}^T\mathbf{X}$, and are basis of row space of \mathbf{X} .

- **Reduced SVD** in linear regression. We simply ignore the zero parts of \mathbf{U} and \mathbf{D} .

\mathbf{D} and \mathbf{U} can be seen as

$$\mathbf{D}_{\text{full}} = \left[\begin{array}{ccc|c} d_1 & 0 & \dots & \\ 0 & \ddots & & \\ \dots & & d_p & \\ \hline 0 & 0 & 0 & \\ \vdots & \vdots & \vdots & \\ 0 & 0 & 0 & \end{array} \right] \left. \begin{array}{l} \\ \\ \\ \end{array} \right\} p \text{ (Non-zero part)} \left. \begin{array}{l} \\ \\ \\ \end{array} \right\} N - p \text{ (All zero part)}$$

$$\mathbf{U}_{\text{full}} = [\mathbf{U}_1 \quad | \quad \mathbf{U}_2]$$

Thus,

$$\begin{aligned} \mathbf{U}_{\text{full}} \cdot \mathbf{D}_{\text{full}} &= [\mathbf{U}_1 \quad \mathbf{U}_2] \cdot \begin{bmatrix} \mathbf{D}_p \\ \mathbf{0} \end{bmatrix} \\ &= \mathbf{U}_1 \cdot \mathbf{D}_p + \mathbf{U}_2 \cdot \mathbf{0} \\ &= \mathbf{U}_1 \cdot \mathbf{D}_p \end{aligned}$$

That's the reduced SVD:

$$\mathbf{X} = \mathbf{U} \mathbf{D} \mathbf{V}^T$$

where $\mathbf{U} \in \mathbb{R}^{N \times p}$ with orthonormal columns, $\mathbf{D} \in \mathbb{R}^{p \times p}$ diagonal with positive entries, and $\mathbf{V} \in \mathbb{R}^{p \times p}$ orthogonal.

In regression, we usually use reduced SVD, and $\mathbf{U}\mathbf{U}^T = \mathbf{H} \neq \mathbf{I}$, but $\mathbf{U}^T\mathbf{U} = \mathbf{I}$. For $\mathbf{y} \in \mathbb{R}^N$, $\mathbf{U}^T\mathbf{y}$ map \mathbf{y} from \mathbb{R}^N to \mathbb{R}^p , which is the coefficients in the basis of columns of \mathbf{U} . And $\mathbf{U}\mathbf{U}^T\mathbf{y}$ project \mathbf{y} onto the column space of \mathbf{X} .

1.5 Multiple testing in forward selection

ESL page 60:

Other more traditional packages base the selection on F -statistics, adding “significant” terms, and dropping “non-significant” terms. These are out of fashion, since they do not take proper account of the multiple testing issues.

Assume we have p candidate features, and already selected k features. When considering adding a new feature, we are actually performing $p - k$ hypothesis tests (each test for one feature). Even the rest $p - k$ features are all noise, with significance level α , we still have a probability of $1 - (1 - \alpha)^{p-k}$ to incorrectly add at least one noise feature.

1.6 Ridge regression

Answer questions: why two forms are equivalent? Why not equivariant under scaling of the inputs? What is a good practice for it? df of ridge? In the case of orthogonal inputs, why $\hat{\beta}^{\text{ridge}} = \hat{\beta}/(1 + \lambda)$?

Ridge regression shrinks the regression coefficients by imposing a penalty on their size:

$$\hat{\beta}^{\text{ridge}} = \arg \min_{\beta} \left\{ \sum_{i=1}^N (y_i - \beta_0 - \sum_{j=1}^p x_{ij}\beta_j)^2 + \lambda \sum_{j=1}^p \beta_j^2 \right\}$$

and is equivalent to

$$\begin{aligned} \hat{\beta}^{\text{ridge}} &= \arg \min_{\beta} \sum_{i=1}^N (y_i - \beta_0 - \sum_{j=1}^p x_{ij}\beta_j)^2 \\ &\text{subject to } \sum_{j=1}^p \beta_j^2 \leq t \end{aligned}$$

And there is a one-to-one correspondence between λ and t .

1.6.1 How to use ridge regression

Apparently, ridge regression is not equivariant under scaling of the inputs. For example, using OLS, measuring x in meters or in centimeters will not change the predictions, since the latter coefficients will just 100 times of the first. However, in ridge regression, the penalty term $\lambda \sum_{j=1}^p \beta_j^2$ will be affected by the scale of x_j . Which means, using centimeters instead of meters will make the penalty on β_j 10000 times larger, leading to different solutions. Thus, it's important to standardize the features (zero mean and unit variance) before applying ridge.

Usually, we calculate μ and σ from training set, and use them to standardize both training and test sets. Scaler can be regarded as part of the model, not data cleaning!

1.6.2 Equivalence of two forms

First, take a review of **KKT conditions**.

Consider a minimization problem with both equality and inequality constraints:

$$\begin{aligned} & \min_{\mathbf{x}} f(\mathbf{x}) \\ \text{subject to } & g_i(\mathbf{x}) \leq 0, \quad i = 1, \dots, m \\ & h_j(\mathbf{x}) = 0, \quad j = 1, \dots, l \end{aligned}$$

And the lagrangian function is:

$$L(\mathbf{x}, \boldsymbol{\lambda}, \boldsymbol{\mu}) = f(\mathbf{x}) + \sum_{i=1}^m \lambda_i g_i(\mathbf{x}) + \sum_{j=1}^l \mu_j h_j(\mathbf{x})$$

If \mathbf{x}^* is a local minimum, then there exist multipliers $\lambda_i^* \geq 0$ and μ_j^* such that the following conditions hold:

- **Stationary**

$$\nabla_{\mathbf{x}} L(\mathbf{x}^*, \boldsymbol{\lambda}^*, \boldsymbol{\mu}^*) = 0$$

- **Primal feasibility**

$$\begin{aligned} g_i(\mathbf{x}^*) &\leq 0, \quad i = 1, \dots, m \\ h_j(\mathbf{x}^*) &= 0, \quad j = 1, \dots, l \end{aligned}$$

The gradient can be regarded as the force to push a particle, primal stationary means the force of $\partial f(\mathbf{x}^*)$ is balanced by a linear sum of forces from constraints.

- **Dual feasibility**

$$\lambda_i^* \geq 0, \quad i = 1, \dots, m$$

All the $\partial g_i(\mathbf{x}^*)$ forces must be one-sided, pointing inwards into the feasible set for \mathbf{x} .

- **Complementary slackness**

$$\lambda_i^* g_i(\mathbf{x}^*) = 0, \quad i = 1, \dots, m$$

The force only activated when the particle is on the boundary of feasible set.

In ridge regression, we have:

$$\mathcal{L} = \|\mathbf{Y} - \mathbf{X}\beta\|_2^2 + \alpha(\|\beta\|_2^2 - t)$$

According to stationary condition:

$$\nabla_{\beta} \mathcal{L} = -2\mathbf{X}^T(\mathbf{Y} - \mathbf{X}\beta) + 2\alpha\beta = 0$$

On the other hand, solving the unconstrained form is

$$\nabla_{\beta} (\|\mathbf{Y} - \mathbf{X}\beta\|_2^2 + \lambda\|\beta\|_2^2) = -2\mathbf{X}^T(\mathbf{Y} - \mathbf{X}\beta) + 2\lambda\beta = 0$$

Thus, if we set $\lambda = \alpha$, the two forms are equivalent.

One step further, according to complementary slackness:

$$\alpha(\|\beta\|_2^2 - t) = 0$$

From unconstrained form, given λ , we can solve β , denote $\beta(\lambda)$, and define $t(\lambda) = \|\beta(\lambda)\|_2^2$. Then, $\beta(\lambda)$ is also the solution of constrained form with $t = t(\lambda)$ (apparently it's on the boundary, and the coefficient $\alpha = \lambda$).

From the constrained form, given t , we can also solve β , denote $\beta(t)$.

- If $\|\beta(t)\|_2^2 < t$, then according to complementary slackness, $\alpha = 0$, which means no penalty, $\lambda = 0$, back to OLS.
- If $\|\beta(t)\|_2^2 = t$, the boundary is effective, correspond to some $\alpha > 0$, and $\lambda = \alpha$.

1.6.3 Bayesian view

Assume prior $\beta \sim \mathcal{N}(0, \tau^2 \mathbf{I})$, and likelihood $\mathbf{Y} \sim \mathcal{N}(\mathbf{X}\beta, \sigma^2 \mathbf{I})$. Then the posterior is:

$$\begin{aligned} p(\beta | \mathbf{Y}) &\propto p(\mathbf{Y} | \beta) p(\beta) \\ &\propto \exp\left(-\frac{1}{2\sigma^2} \|\mathbf{Y} - \mathbf{X}\beta\|_2^2\right) \exp\left(-\frac{1}{2\tau^2} \|\beta\|_2^2\right) \\ &\propto \exp\left(-\frac{1}{2} \left(\frac{1}{\sigma^2} \|\mathbf{Y} - \mathbf{X}\beta\|_2^2 + \frac{1}{\tau^2} \|\beta\|_2^2\right)\right) \end{aligned}$$

Maximizing the posterior is equivalent to minimizing the negative log-posterior. The object is equivalent to ridge regression with $\lambda = \sigma^2/\tau^2$.

However, it seems that from bayesian view, we don't have a hard constrain on the size of β , just a prior that β is likely to be small. But the constrained form of ridge regression directly enforce $\|\beta\|_2^2 \leq t$. Is it contradictory?

No. Actually, the posterior distribution of β still have infinite support, meaning that β can still be large with small probability. Our MAP estimate of β is the mode of the posterior, it's just a **point estimate**. This point estimator can still satisfy the hard constraint $\|\beta\|_2^2 \leq t$ for some t .

1.6.4 SVD in ridge

The solution of ridge regression is:

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

When $\mathbf{X}^T \mathbf{X}$ is singular (not full rank), OLS estimator is not defined, and adding $\lambda \mathbf{I}$ ensures that $\mathbf{X}^T \mathbf{X} + \lambda \mathbf{I}$ is positive definite and invertible.

Take reduced SVD of $\mathbf{X} = \mathbf{U} \mathbf{D} \mathbf{V}^T$, $r = \text{Rank}(\mathbf{X})$, where $\mathbf{U} \in \mathbb{R}^{N \times r}$, $\mathbf{D} \in \mathbb{R}^{r \times r}$, $\mathbf{V} \in \mathbb{R}^{p \times r}$. Then,

$$\begin{aligned} \mathbf{X}^T \mathbf{X} + \lambda \mathbf{I} &= \mathbf{V} \mathbf{D}^T \mathbf{U}^T \mathbf{U} \mathbf{D} \mathbf{V}^T + \lambda \mathbf{I} \\ &= \mathbf{V} \mathbf{D}^T \mathbf{D} \mathbf{V}^T + \lambda \mathbf{V} \mathbf{V}^T \\ &= \mathbf{V} (\mathbf{D}^2 + \lambda \mathbf{I}) \mathbf{V}^T \end{aligned}$$

must be full rank and invertible.

It's also easy to see that in the case of orthogonal inputs (i.e., $\mathbf{X}^T \mathbf{X} = \mathbf{I}$), the ridge estimates are just a scaled version of the least squares estimates, that is

$$\hat{\beta}^{\text{ridge}} = \frac{1}{1 + \lambda} \hat{\beta}$$

Further more,

$$\begin{aligned} \mathbf{X} \beta^{\text{OLS}} &= \mathbf{X} (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{Y} \\ &= \mathbf{U} \mathbf{D} \mathbf{V}^T (\mathbf{V} \mathbf{D}^T \mathbf{U}^T \mathbf{U} \mathbf{D} \mathbf{V}^T)^{-1} \mathbf{V} \mathbf{D}^T \mathbf{U}^T \mathbf{Y} \\ &= \mathbf{U} \mathbf{D} \mathbf{D}^{-2} \mathbf{D}^T \mathbf{U}^T \mathbf{Y} \\ &= \mathbf{U} \mathbf{U}^T \mathbf{Y} \end{aligned}$$

and

$$\begin{aligned}
\mathbf{X}\beta^{\text{ridge}} &= \mathbf{X}(\mathbf{X}^T\mathbf{X} + \lambda\mathbf{I})^{-1}\mathbf{X}^T\mathbf{Y} \\
&= \mathbf{U}\mathbf{D}\mathbf{V}^T(\mathbf{V}(\mathbf{D}^2 + \lambda\mathbf{I})\mathbf{V}^T)^{-1}\mathbf{V}\mathbf{D}^T\mathbf{U}^T\mathbf{Y} \\
&= \mathbf{U}\mathbf{D}(\mathbf{D}^2 + \lambda\mathbf{I})^{-1}\mathbf{D}^T\mathbf{U}^T\mathbf{Y} \\
&= \mathbf{U}\text{diag}\left(\frac{d_1^2}{d_1^2 + \lambda}, \frac{d_2^2}{d_2^2 + \lambda}, \dots, \frac{d_r^2}{d_r^2 + \lambda}\right)\mathbf{U}^T\mathbf{Y}
\end{aligned}$$

In both way, $\mathbf{U}^T\mathbf{Y}$ gets the coordinates of \mathbf{Y} in the basis of columns of \mathbf{U} (the principal components of \mathbf{X}), and then send back to original space using \mathbf{U} . But ridge regression shrink the coordinates.

References

- [1] Wikipedia contributors, "Gauss–Markov theorem," *Wikipedia, The Free Encyclopedia*, https://en.wikipedia.org/wiki/Gauss%E2%80%93Markov_theorem (accessed Dec 29, 2025).