

Why is ridge regression a shrinkage method? Suppose the design matrix \mathbf{X} has ON^a columns, $\mathbf{X}^T \mathbf{X} = \mathbf{I}_p$. Then the ridge estimate/prediction is a shrinkage version of the LS estimate/prediction.

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

$$\begin{aligned} \hat{\boldsymbol{\beta}}^{\text{ridge}} &= (\mathbf{X}^T \mathbf{X} + \lambda \mathbf{I})^{-1} \mathbf{X}^T \mathbf{y} \\ &= \frac{1}{1 + \lambda} \mathbf{X}^T \mathbf{y} = \frac{1}{1 + \lambda} \hat{\boldsymbol{\beta}}^{\text{LS}} \end{aligned}$$

$$\hat{\mathbf{y}}_{\text{LS}} = \mathbf{X} \hat{\boldsymbol{\beta}}^{\text{LS}}$$

$$\hat{\mathbf{y}}_{\text{ridge}} = \mathbf{X} \hat{\boldsymbol{\beta}}^{\text{ridge}} = \frac{1}{1 + \lambda} \mathbf{y}_{\text{LS}}$$

^aOrthonormal (ON): orthogonal with norm one.

In case the columns of \mathbf{X} are not orthogonal, we can reformulate the regression on an orthogonal version of \mathbf{X} , known as the principal components analysis or SVD. Similarly we can see that the ridge estimate/prediction is a shrinkage version of the LS estimate/prediction. (You can skip the next 6 slides.)

Consider a singular value decomposition (SVD) of \mathbf{X} :

$$\mathbf{X}_{n \times p} = \mathbf{U}_{n \times p} \mathbf{D}_{p \times p} \mathbf{V}_{p \times p}^T,$$

where

- $\mathbf{U}_{n \times p}$: columns \mathbf{u}_j 's form an ON basis for $C(\mathbf{X})$, $\mathbf{U}^T \mathbf{U} = \mathbf{I}_p$.
- $\mathbf{V}_{p \times p}$: columns \mathbf{v}_j 's form an ON basis for \mathbb{R}^p with $\mathbf{V}^T \mathbf{V} = \mathbf{I}_p$.
- $\mathbf{D}_{p \times p}$: diagonal matrix with diagonal entries $d_1 \geq d_2 \geq \cdots \geq d_p \geq 0$ being the singular values of \mathbf{X} .

For ease of exposition we assume $n > p$ and $\text{rank}(\mathbf{X}) = p$. Therefore

$d_p > 0$.

The Geometric interpretation of SVD:

$$\mathbf{X}_{n \times p} = \mathbf{U}_{n \times p} \mathbf{D}_{p \times p} \mathbf{V}_{p \times p}^T$$

Map a unit circle in \mathbf{R}^p to an ellipse in \mathbf{R}^n

$$\mathbf{X}_{n \times p} \mathbf{v}_{j_{p \times 1}} = \mathbf{U}_{n \times p} \mathbf{D}_{p \times p} \mathbf{V}_{p \times p}^T \mathbf{v}_{j_{p \times 1}} = d_j \mathbf{u}_j.$$

Consider a singular value decomposition (SVD) of \mathbf{X} :

$$\mathbf{X}_{n \times p} = \mathbf{U}_{n \times p} \mathbf{D}_{p \times p} \mathbf{V}_{p \times p}^T,$$

where

- $\mathbf{U}_{n \times p}$: columns \mathbf{u}_j 's form an ON basis for $C(\mathbf{X})$, $\mathbf{U}^T \mathbf{U} = \mathbf{I}_p$.
- $\mathbf{V}_{p \times p}$: columns \mathbf{v}_j 's form an ON basis for \mathbb{R}^p with $\mathbf{V}^T \mathbf{V} = \mathbf{I}_p$.
- $\mathbf{D}_{p \times p}$: diagonal matrix with diagonal entries $d_1 \geq d_2 \geq \dots \geq d_p \geq 0$ being the singular values of \mathbf{X} .

For ease of exposition we assume $n > p$ and $\text{rank}(\mathbf{X}) = p$. Therefore $d_p > 0$.

- PCA: write $\mathbf{X} = \mathbf{FV}^T$ where each columns of $\mathbf{F}_{n \times p} = \mathbf{UD}$ is the so-called **principal components** and each column of \mathbf{V} is the **principal component directions** of \mathbf{X} ;

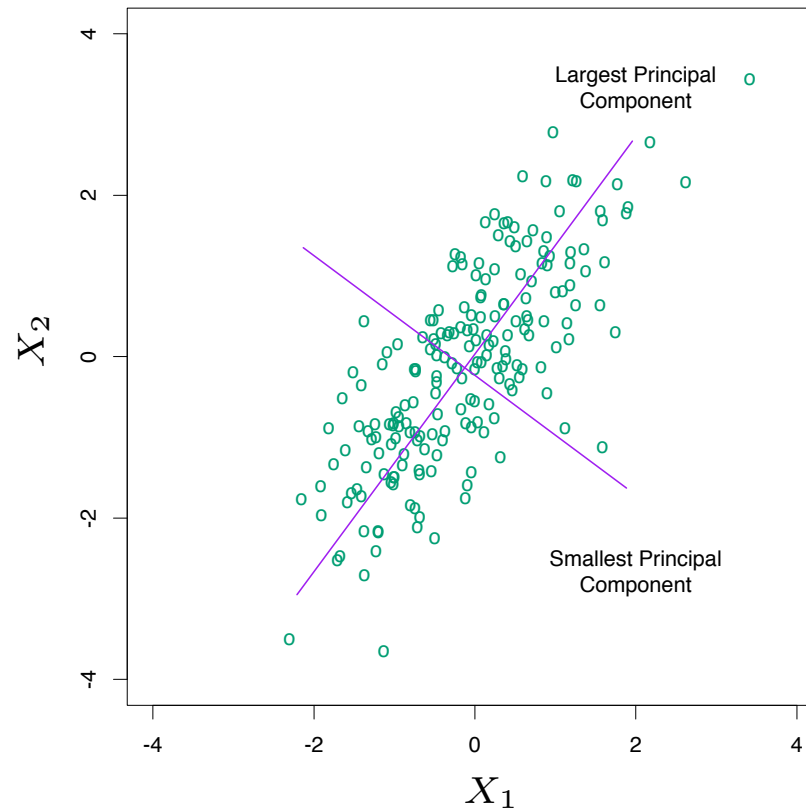


FIGURE 3.9. *Principal components of some input data points. The largest principal component is the direction that maximizes the variance of the projected data, and the smallest principal component minimizes that variance. Ridge regression projects \mathbf{y} onto these components, and then shrinks the coefficients of the low-variance components more than the high-variance components.*

Write

$$\mathbf{y} - \mathbf{X}\boldsymbol{\beta} = \mathbf{y} - \mathbf{UDV}\boldsymbol{\beta} = \mathbf{y} - \mathbf{F}\boldsymbol{\alpha}.$$

there is a one-to-one correspondence between $\boldsymbol{\beta}_{p \times 1}$ and $\boldsymbol{\alpha}_{p \times 1}$ and

$\|\boldsymbol{\beta}\|^2 = \|\boldsymbol{\alpha}\|^2$. So

$$\min_{\boldsymbol{\beta} \in \mathbb{R}^p} \|\mathbf{y} - \mathbf{X}\boldsymbol{\beta}\|^2 + \lambda \|\boldsymbol{\beta}\|^2 \iff \min_{\boldsymbol{\alpha} \in \mathbb{R}^p} \|\mathbf{y} - \mathbf{F}\boldsymbol{\alpha}\|^2 + \lambda \|\boldsymbol{\alpha}\|^2.$$

$$\hat{\boldsymbol{\alpha}}^{\text{LS}} = \mathbf{D}^{-1} \mathbf{U}^T \mathbf{y}, \quad \hat{\alpha}_j^{\text{LS}} = \frac{1}{d_j} \mathbf{u}_j^T \mathbf{y}$$

$$\hat{\boldsymbol{\alpha}}^{\text{ridge}} = \text{diag}\left(\frac{d_j}{d_j^2 + \lambda}\right) \mathbf{U}^T \mathbf{y}, \quad \hat{\alpha}_j^{\text{ridge}} = \frac{d_j^2}{d_j^2 + \lambda} \hat{\alpha}_j^{\text{LS}}$$

So the ridge estimate $\hat{\boldsymbol{\alpha}}^{\text{ridge}}$ shrinks the LS estimate $\hat{\boldsymbol{\alpha}}^{\text{LS}}$ by the factor $d_j^2/(d_j^2 + \lambda)$: directions with smaller eigen values get more shrinkage.

- The LS prediction

$$\mathbf{F}\hat{\boldsymbol{\alpha}}^{\text{LS}} = (\mathbf{U}\mathbf{D})\mathbf{D}^{-1}\mathbf{U}^T\mathbf{y} = \mathbf{U}\mathbf{U}^T\mathbf{y} = \sum_{j=1}^p (\mathbf{u}_j^T \mathbf{y}) \mathbf{u}_j.$$

- The ridge prediction

$$\mathbf{F}\hat{\boldsymbol{\alpha}}^{\text{ridge}} = \mathbf{U}\text{diag}\left(\frac{d_j}{d_j^2 + \lambda}\right)\mathbf{U}^T\mathbf{y} = \sum_{j=1}^p \frac{d_j^2}{d_j^2 + \lambda} (\mathbf{u}_j^T \mathbf{y}) \mathbf{u}_j$$

- So the ridge prediction $\hat{\mathbf{y}}_{\text{ridge}}$ shrinks the LS prediction $\hat{\mathbf{y}}_{\text{LS}}$ by factor $d_j^2/(d_j^2 + \lambda)$: directions with smaller eigen values get more shrinkage.