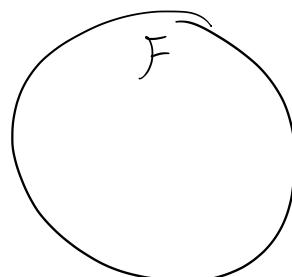


Ridge Regression

Generative Models

- Think about the model is generated and how data are generated from the model
 - prior
 - likelihood
- Bayes Rule : $P(\text{model} \mid \text{data}) \propto P(\text{data} \mid \text{model}) \cdot P(\text{model})$
 - posterior
 - likelihood
 - prior
- MAP : $\max_{\text{model}} \log p(\text{model} \mid \text{data})$
 - \uparrow
 - maximum a posteriori
 - $= \min_{\text{models}} -\log \text{likelihood} - \log \text{prior}$
 - supervised learning
 - loss function
 - regularization
- Maximum likelihood is a special case (where we don't use a prior, or use uniform prior so log prior is constant)
- Belief? Not necessary



$$\text{posterior} = \text{likelihood} \cdot \text{prior}$$

if gaussian \times gaussian then posterior is gaussian

if Bernoulli \times prior is Beta then posterior is Beta

Ridge Regression

$$\min_{\lambda} \frac{1}{n} \sum_i (y_i - f_{\lambda}(x_i))^2 + C \|\lambda\|_2^2 \quad \text{where } f_{\lambda}(x_i) = \sum_j \lambda_j x_{ij}$$

a Bayesian version of it:

$$\begin{aligned} \bar{\beta} &\sim N(\bar{0}, \sigma^2 I) \\ \bar{y} &\sim N(\bar{x}\bar{\beta}, \sigma^2 I) \end{aligned}$$

posterior: $p(\bar{\beta} | \bar{y}, \bar{x}) = p(\bar{y} | \bar{\beta}, \bar{x}) \cdot p(\bar{\beta}) \cdot \frac{1}{Z}$

$$= \frac{1}{Z} \frac{1}{(2\pi)^{n/2} \sigma^2} \exp\left(-\frac{1}{2\sigma^2} \|\bar{y} - \bar{x}\bar{\beta}\|^2\right) \cdot \frac{1}{(2\pi)^{p/2}} \exp\left(-\frac{1}{2\sigma^2} \|\bar{\beta}\|^2\right)$$

$$-\log p(\bar{\beta} | \bar{y}, \bar{x}) = -\log(\text{stuff}) + \frac{1}{2\sigma^2} \|\bar{y} - \bar{x}\bar{\beta}\|^2 + \frac{1}{2\sigma^2} \|\bar{\beta}\|^2 \leftarrow \text{Ridge regression}$$

Fact 1 : Ridge Regression has a generative + frequentist interpretation

Fact 2 : Least squares regression has a closed form solution.

$$\min_{\lambda} F(\lambda) \quad \text{where } F(\lambda) = \|\bar{y} - \bar{x}\bar{\lambda}\|^2 = \sum_i (y_i - \bar{x}_i \bar{\lambda})^2$$

$$\frac{dF(\lambda)}{d\lambda_j} = - \sum_i 2(y_i - \bar{x}_i \bar{\lambda}) x_{ij} = -2 \bar{x}_j^\top (\bar{y} - \bar{x}\bar{\lambda})$$

$$\nabla F(\bar{\lambda}) = -2 \bar{x}^\top (\bar{y} - \bar{x}\bar{\lambda}) = -2 [\bar{x}^\top \bar{y} - \bar{x}^\top \bar{x}\bar{\lambda}] = 0$$

$$\bar{x}^\top \bar{y} = \bar{x}^\top \bar{x} \bar{\lambda}^*$$

$$(\bar{x}^\top \bar{x})^{-1} \bar{x}^\top \bar{y} = \bar{\lambda}^*$$

$$\hat{y} = \bar{x} \bar{\lambda}^* = \underbrace{\bar{x} (\bar{x}^\top \bar{x})^{-1} \bar{x}^\top}_{\text{"hat matrix"}} \bar{y} = H \bar{y} = \text{"smoother" of } y$$

Fact 3 Ridge Regression has a closed form solution

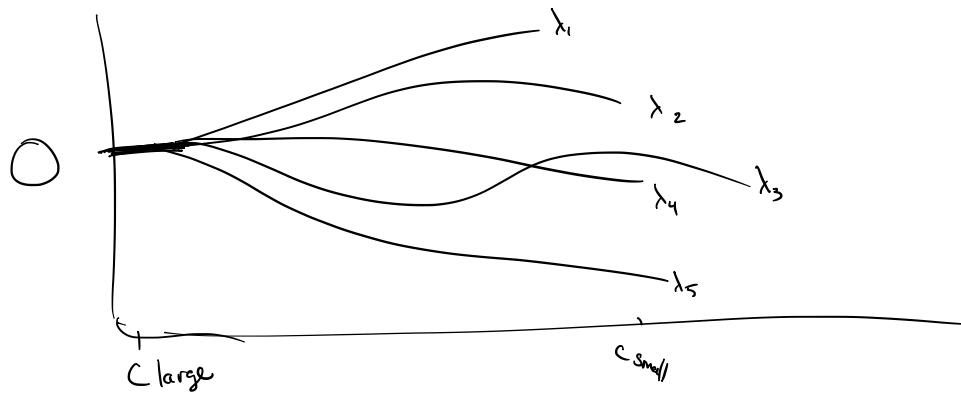
$$F(\lambda) = \|\mathbf{y} - \mathbf{X}\lambda\|_2^2 + C\|\lambda\|_2^2$$

$$\frac{\partial F(\lambda)}{\partial \lambda_j} = -2\mathbf{x}_j^T(\mathbf{y} - \mathbf{X}\lambda) + 2C\lambda_j$$

$$\nabla F(\lambda) = -2\mathbf{X}^T(\mathbf{y} - \mathbf{X}\lambda) + 2C\lambda = 2(-\mathbf{X}^T\mathbf{y} + \mathbf{X}^T\mathbf{X}\lambda + C\lambda) = 0$$

$$\mathbf{X}^T\mathbf{y} = (\mathbf{X}^T\mathbf{X} + C\mathbf{I})\lambda^*$$

$$(\mathbf{X}^T\mathbf{X} + C\mathbf{I})^{-1}\mathbf{X}^T\mathbf{y} = \lambda^*$$



Computation issues

$$\lambda^* = (\mathbf{X}^T\mathbf{X} + C\mathbf{I})^{-1}\mathbf{X}^T\mathbf{y}$$

$\mathcal{O}(p^3)$

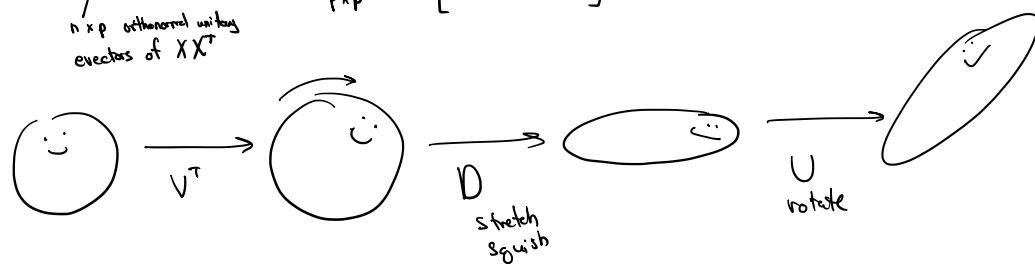
use SVD : $\mathbf{X} = \mathbf{U} \mathbf{D} \mathbf{V}^T$

$$\begin{bmatrix} \vdots & \vdots & & \vdots \\ u_1 & u_2 & \cdots & u_p \\ \vdots & \vdots & & \vdots \end{bmatrix}_{n \times p} \text{ orthonormal unitary vectors of } \mathbf{X}\mathbf{X}^T$$

$$\begin{bmatrix} d_1 & & & \\ & \ddots & & \\ & & d_p & \end{bmatrix}_{p \times p}$$

$$\begin{bmatrix} \vdots & \vdots & & \vdots \\ v_1^T & v_2^T & \cdots & v_p^T \\ \vdots & \vdots & & \vdots \end{bmatrix}_{p \times n} \text{ orthonormal vectors of } \mathbf{X}^T\mathbf{X}$$

$$\mathbf{A}^T\mathbf{A} = \mathbf{A}\mathbf{A}^T = \mathbf{I}$$



$$\text{Note } U^T U = I \quad V^T V = I$$

$$\begin{aligned} \text{Simplify } X^T X &= (UDV^T)^T (UDV^T) = V^T D^T U^T U D V^T \\ &= V D D V^T = V \begin{pmatrix} d_1^2 & & \\ & d_2^2 & \\ & & \ddots & d_n^2 \end{pmatrix} V^T = V D^2 V^T \end{aligned}$$

Sub into λ^*

$$\begin{aligned} \lambda^* &= (V D^2 V^T + CI)^{-1} X^T Y \\ &\quad C V \begin{pmatrix} 1 & & \\ & \ddots & \\ & & c \end{pmatrix} V^T = V \begin{pmatrix} c & & \\ & c & \\ & & c \end{pmatrix} V^T \\ &= (V D^2 V^T + V \begin{pmatrix} c & & \\ & c & \\ & & c \end{pmatrix})^{-1} X^T Y \\ &= \left(V \begin{pmatrix} d_1^2 + c & & \\ & d_2^2 + c & \\ & & \ddots & d_n^2 + c \end{pmatrix} V^T \right)^{-1} X^T Y \\ &= \left(V^T \begin{pmatrix} d_1^2 + c & & \\ & d_2^2 + c & \\ & & \ddots & d_n^2 + c \end{pmatrix} V \right)^{-1} X^T Y \\ &= V \begin{pmatrix} \frac{d_1^2 + c}{d_1^2 + c} & & \\ & \frac{d_2^2 + c}{d_2^2 + c} & \\ & & \ddots & \frac{d_n^2 + c}{d_n^2 + c} \end{pmatrix} V^T X^T Y \\ &\quad \underbrace{(UDV^T)^T}_{= V^T D^T U^T} \\ &= V \begin{pmatrix} \frac{d_1^2 + c}{d_1^2 + c} & & \\ & \frac{d_2^2 + c}{d_2^2 + c} & \\ & & \ddots & \frac{d_n^2 + c}{d_n^2 + c} \end{pmatrix} V^T U^T Y \\ &= V \begin{pmatrix} \frac{d_1^2 + c}{d_1^2 + c} & & \\ & \frac{d_2^2 + c}{d_2^2 + c} & \\ & & \ddots & \frac{d_n^2 + c}{d_n^2 + c} \end{pmatrix} U^T Y \\ &= V \text{ diag} \left(\frac{d_i^2 + c}{d_i^2 + c} \right) U^T Y = \lambda^* \end{aligned}$$

$$\begin{aligned} \hat{Y}_{\text{ridge}} &= X \lambda^* = U D V^T \text{ diag} \left(\frac{d_i^2}{d_i^2 + c} \right) U^T Y \\ &= U \text{ diag} \left(\frac{d_i^2}{d_i^2 + c} \right) U^T Y \\ &= \sum_j \bar{u}_j \frac{d_j^2}{d_j^2 + c} u_j^T \cdot Y \quad \begin{matrix} \uparrow & \uparrow & \uparrow \\ n \times 1 & & n \times 1 \end{matrix} \end{aligned}$$

Fact 4 This \uparrow is not too difficult to compute.

Fact 5: There's a relationship of ridge regression to principal components analysis (PCA)

PCA: finds orthogonal directions in which variance is maximized along first component, second component etc.

loading vector
 $\omega_{(1)} = \underset{\omega: \|\omega\|=1}{\operatorname{argmax}} (\chi_i \cdot \omega)^2 = \underset{\omega: \|\omega\|=1}{\operatorname{argmax}} \|\chi \omega\|^2$

$\langle \chi \omega, \chi \omega \rangle$
 $(\chi \omega)^T \chi \omega$
 $\omega^T \underbrace{\chi^T \chi}_{\text{occurs at largest eigenvalue of } \chi^T \chi \text{ when } \tilde{\omega} \text{ is corresponding eigenvector } v_i}$

pca_i = $\chi \cdot v_i$

matrix of all principal components = $\underline{\underline{\text{pca}}} = UDV^T$

pca_j = $\tilde{U}_j \cdot d_j$

Back to ridge

$$\hat{Y} = \chi \cdot \lambda^* = \chi \underbrace{V}_{T} \operatorname{diag}\left(\frac{d_j}{d_j^2 + C}\right) U^T Y$$

$$\chi V = UDV^T \xrightarrow{T} UD$$

call U "new data" χ

$$\begin{aligned} &= U D \operatorname{diag}\left(\frac{d_j}{d_j^2 + C}\right) U^T Y \\ &= \underbrace{U}_{\tilde{\chi}} \underbrace{\operatorname{diag}\left(\frac{d_j^2}{d_j^2 + C}\right)}_{\tilde{\lambda}^*} \underbrace{U^T}_{\tilde{Y}} Y \\ &= \underbrace{U}_{\tilde{\chi}} \underbrace{\operatorname{diag}\left(\frac{1}{\frac{d_j^2}{d_j^2 + C}}\right)}_{\tilde{\lambda}^*} \underbrace{DU^T}_{\tilde{Y}} Y \\ C = 0 &\Rightarrow \lambda_j^* = U_j^T Y \end{aligned}$$

$$C > 0 \Rightarrow \lambda_j^* = \frac{d_j^2}{d_j^2 + C} U_j^T Y$$

