

Least Squares & Friends

Generative Models

- Think about how the model is generated and how the prior data are generated from the model.

likelihood

$$\text{Bayes Rule: } \underbrace{P(\text{model} | \text{data})}_{\text{posterior}} \propto \underbrace{P(\text{data} | \text{model})}_{\text{likelihood}} P(\text{model})_{\text{prior}}$$

$$\log \text{posterior} \leftarrow \underbrace{\log \text{likelihood}}_{\substack{\text{comes from} \\ \text{data}}} + \underbrace{\log \text{prior}}_{\substack{\text{comes} \\ \text{from prior}}}$$

$$\text{MAP: } \max_{\text{model}} \log p(\text{model} | \text{data}) \rightarrow \min_{\text{model}} \underbrace{-\log \text{likelihood}}_{\substack{\text{loss function}}} + \underbrace{-\log \text{prior}}_{\substack{\text{regularization}}}$$

- Maximum likelihood is a special case (where we don't use a prior, or use a uniform prior so log prior is constant.)
- Belief? Not necessary

Ridge Regression

$$\min_{\lambda} \frac{1}{n} \sum_i (y_i - f_x(x_i))^2 + C \|\lambda\|_2^2$$

$$\text{where } f_x(\bar{x}_i) = \sum_j \lambda_j x_{ij}$$

A Bayesian version of it:

$$\vec{\beta} \sim N(\bar{0}, \tau^2 \bar{I})$$

$$\bar{Y} \sim N(\bar{x}\bar{\beta}, \sigma^2 \bar{I})$$



$$\text{posterior} := p(\beta | \bar{Y}, \bar{X}) = p(\bar{Y} | \beta, \bar{X}) p(\beta) \cdot \frac{1}{Z}$$

"evidence"
doesn't depend on β → $\int p(\bar{Y} | \bar{\beta} \bar{X}) p(\bar{\beta}) d\bar{\beta}$

$$= \underbrace{\frac{1}{Z} \cdot \frac{1}{(2\pi)^{\frac{p}{2}} \sigma^p}}_{\text{constant}} \exp\left(-\frac{1}{2\sigma^2} \|\bar{Y} - \bar{X}\bar{\beta}\|_2^2\right) \cdot \underbrace{\frac{1}{(2\pi)^{\frac{p}{2}} \tau^p} \exp\left(-\frac{1}{2\tau^2} \|\beta\|^2\right)}_{\text{prior}}$$

$$-\log p(\bar{\beta} | \bar{Y}, \bar{X}) = -\log \text{stuff} + \underbrace{\frac{1}{2\sigma^2} \|Y - X\bar{\beta}\|_2^2}_{\text{residual sum of squares}} + \underbrace{\frac{1}{2\tau^2} \|\beta\|^2}_{\text{ridge term}}$$

Ridge Regression

Fact 1: Ridge regression has a generative & frequentist interpretation.

Fact 2 Least squares regression (no regularization) has a closed form solution.

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

$$\frac{\partial F(\bar{\lambda})}{\partial \lambda_j} = - \sum_i 2(y_i - \bar{x}_i \bar{\lambda}) x_{ij} = - 2 \bar{X}_{\cdot j}^T (\bar{Y} - \bar{X}\bar{\lambda})$$

in vector notation:

$$\nabla F(\bar{\lambda}) = - 2 \bar{X}^T (\bar{Y} - \bar{X}\bar{\lambda}) = - 2 [\bar{X}^T y - \bar{X}^T \bar{X}\bar{\lambda}] = 0$$

$$\bar{X}^T y = \bar{X}^T \bar{X} \bar{\lambda}^*$$

$$(\bar{X}^T \bar{X})^{-1} \bar{X}^T y = \bar{\lambda}^*$$

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

Fact 3 Ridge Regression has a closed form solution.

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

$$\frac{\partial F(\bar{\lambda})}{\partial \lambda_j} = - 2 \bar{X}_{\cdot j}^T (\bar{Y} - \bar{X}\bar{\lambda}) + 2C \lambda_j$$

$$\nabla F(\bar{\lambda}) = - 2 \bar{X}^T (\bar{Y} - \bar{X}\bar{\lambda}) + 2C \bar{\lambda}$$

$$\nabla F(\bar{\lambda}) = -2\bar{X}^T(\bar{Y} - \bar{X}\bar{\lambda}) + 2C\bar{\lambda}$$

$$= 2(-\bar{X}^T\bar{Y} + \bar{X}^T\bar{X}\bar{\lambda} + C\bar{\lambda}) = 0$$

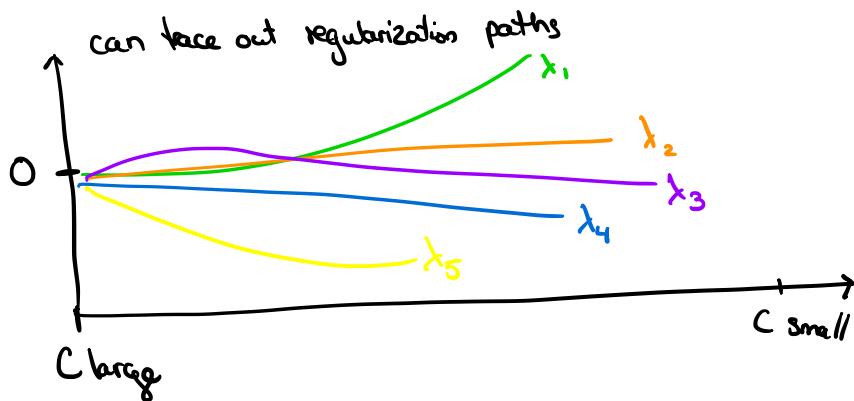
$$\bar{X}^T\bar{Y} = (\bar{X}^T\bar{X} + C\bar{I})\bar{\lambda}^*$$

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

minimize
↓

Hoerl & Kennard 1970

$$\text{Moore-Penrose inverse} \rightarrow \bar{X}^+ = \lim_{\delta \rightarrow 0} (\bar{X}^T\bar{X} + \delta I)^{-1}\bar{X}^T.$$



Computation issues

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

\uparrow
expensive $O(p^3)$

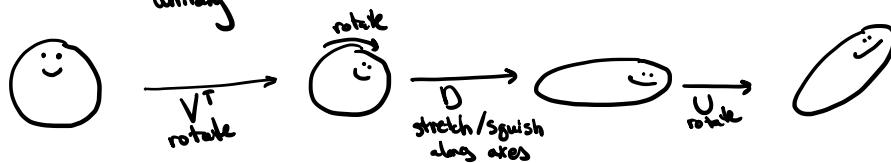
use SVD: $\bar{X} = \bar{U}\bar{D}\bar{V}^T \leftarrow O(np^2)$

$\begin{bmatrix} | & | & \dots & | \\ u_1 & u_2 & \dots & u_p \end{bmatrix}$
↑
eigenvectors of $\bar{X}\bar{X}^T$
 $n \times p$ orthonormal unitary

$$\begin{bmatrix} d_1 & & & \\ & \ddots & & \\ & & 0 & \\ & & & d_p \end{bmatrix}$$

$\begin{bmatrix} v_1^T \\ v_2^T \\ \vdots \\ v_p^T \end{bmatrix}$ } eigenvectors of $\bar{X}^T\bar{X}$
 $p \times p$ orthogonal
 $A^T A = AA^T = 1$

Note:
 $U^T U = 1$
 $V^T V = 1$



Simplify:

$$\bar{X}^T \bar{X} = (\bar{U} \bar{D} \bar{V}^T)^T (\bar{U} \bar{D} \bar{V}^T) \stackrel{\text{property of transpose}}{=} \bar{V} \underbrace{\bar{D} \bar{D}^T}_{\mathbf{1}} \underbrace{\bar{U}^T}_{\bar{U}} \bar{U} \bar{D} \bar{V}^T$$
$$= \bar{V} \bar{D} \bar{D} \bar{V}^T = \bar{V} \begin{pmatrix} d_1^2 & & \\ & \ddots & \\ & & d_p^2 \end{pmatrix} \bar{V}^T = \bar{V} \bar{D}^2 \bar{V}^T$$

Sub into λ^* :

$$\lambda^* = (\bar{V} \bar{D}^2 \bar{V}^T + \underline{C I})^{-1} \bar{X}^T \bar{Y}$$

$$C V \begin{pmatrix} \cdot & \cdot & \cdot \\ \cdot & \cdot & \cdot \\ \cdot & \cdot & \cdot \end{pmatrix} V^T = V \begin{pmatrix} c & & \\ & c & & \\ & & c & \\ & & & c \end{pmatrix} V^T$$

$$= (\bar{V} \bar{D}^2 \bar{V}^T + V \begin{pmatrix} c & & \\ & c & & \\ & & c & \\ & & & c \end{pmatrix} V^T)^{-1} \bar{X}^T \bar{Y}$$

$$= \left(V \begin{pmatrix} d_1^2 + c & & \\ & d_2^2 + c & & \\ & & \ddots & \\ & & & d_p^2 + c \end{pmatrix} V^T \right)^{-1} \bar{X}^T \bar{Y}$$

$$\stackrel{\text{property of inverse}}{=} (V^T)^{-1} \begin{pmatrix} d_1^2 + c & & \\ & d_2^2 + c & & \\ & & \ddots & \\ & & & d_p^2 + c \end{pmatrix} V^{-1} X^T Y$$

$$= V \begin{pmatrix} 1/(d_1^2 + c) & & \\ & 1/(d_2^2 + c) & & \\ & & \ddots & \\ & & & 1/(d_p^2 + c) \end{pmatrix} V^{-1} \underbrace{X^T Y}_{\cancel{V^{-1} \bar{V} \bar{D} \bar{U}^T Y}}$$

$$= V \begin{pmatrix} 1/(d_1^2 + c) & & \\ & \ddots & & \\ & & 1/(d_p^2 + c) \end{pmatrix} \cancel{V^{-1} \bar{V} \bar{D} \bar{U}^T Y} U^T Y$$

$$= V \begin{pmatrix} 1/(d_1^2 + c) & & \\ & \ddots & & \\ & & 1/(d_p^2 + c) \end{pmatrix} \begin{pmatrix} d_1 & & \\ & d_2 & & \\ & & \ddots & \\ & & & d_p \end{pmatrix} U^T Y$$

$$= V \begin{pmatrix} d_1 / (d_1^2 + c) & & \\ & \ddots & & \\ & & 1 / (d_p^2 + c) \end{pmatrix} U^T Y = V \text{ diag} \left(\frac{d_j}{d_j^2 + c} \right) U^T Y$$

Aside: $(\bar{U} \bar{D} \bar{V}^T)^T$
 $= \bar{V}^T \bar{D}^T \bar{U}^T$
 $= V D U^T$

Since we have x^* :

$$\begin{aligned}\hat{y}^{\text{ridge}} = Xx^* &= \underbrace{UDV^T}_{\text{SVD}} \text{diag} \left(\frac{d_j}{d_j^2 + c} \right) U^T y \\ &= U \text{diag} \left(\frac{d_j^2}{d_j^2 + c} \right) U^T y \\ &= \sum_j \underbrace{\bar{v}_j}_{\substack{n \times 1 \\ \text{factor}_j}} \frac{d_j^2}{d_j^2 + c} \underbrace{\bar{U}_j^T \cdot y}_{\substack{1 \times n \\ \sim n \times 1}} \\ &\quad \text{damping effect}\end{aligned}$$

Fact 4 This ↗ is not too difficult to compute.

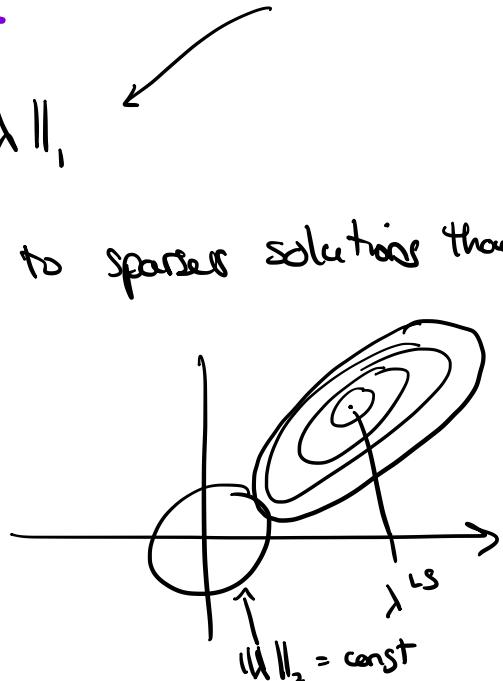
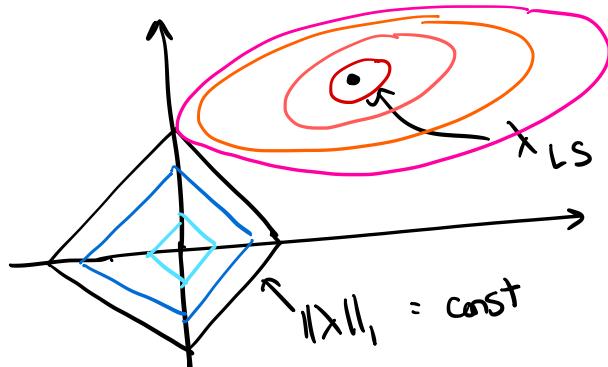
Computational issue averted.

(linear in n , linear in p after the svd)

Fact 5 None of this works for ℓ_1 penalties (lasso).

$$\sum_i (y_i - f(x_i))^2 + C\|\lambda\|_1$$

No closed form solution. Leads to sparser solutions though.



$$\sum_i (y_i - f(x_i))^2 \text{ s.t. } \|\lambda\|_1 \leq \alpha$$

$$\sum_i (y_i - f(x_i))^2 \text{ s.t. } \|\lambda\|_2 \leq \alpha$$

Kernel Least Squares

regular ridge

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

replace $\bar{\lambda} \rightarrow \bar{X}_r^T$ \leftarrow If I can get r , then I can get λ

$$F(r) = \|Y - \underbrace{\bar{X}\bar{X}^T}_{n \times n} r\|_2^2 + C\|X_r^T\|_2^2$$

$n \times p$ $p \times n$ $n \times 1$

$$F(r) = \|Y - K_r\|_2^2 + C \underbrace{\langle X_r^T, X_r^T \rangle}_{(X_r^T)^T (X_r^T)}$$

$r^T X^T X^T r$

$$\langle r, K_r \rangle$$

$$F(r) = \|Y - K_r\|_2^2 + C \langle r, K_r \rangle$$

$n \times n$ $n \times 1$

shorthand

$$\nabla F(r) \stackrel{\downarrow}{=} -2K(Y - K_r) + C2K_r = 0 \text{ at } r^*$$

$$-Y + K_r^* + Cr^* = 0$$

$$(K + CI)r^* = Y$$

$r^* = (K + CI)^{-1}Y$

and remember,

$X^* = X^T r^* = \sum_i X_i^T r_i^*$

Prediction at a new test point \tilde{x} :

$$\begin{aligned}
 f(\tilde{x}) &= \tilde{x} \lambda^* = \tilde{x} \sum_i x_i^T r_i^* = \sum_i \tilde{x} x_i^T r_i^* \\
 &= \sum_i \underbrace{k(\tilde{x}, x_i) r_i^*}_{K_{\tilde{x}, i}} = K_{\tilde{x}}^T r^* = \boxed{K_{\tilde{x}}^T (K + cI)^{-1} y}
 \end{aligned}$$

representer theorem
told us this would happen

Trick: $(P^{-1} + B^T R^{-1} B)^{-1} B^T R^{-1} = P B^T (B P B^T + R)^{-1}$

$$(P^{-1} + B^T R^{-1} B)^{-1} B^T R^{-1} (B P B^T + R) = P B^T$$

$$B^T R^{-1} (B P B^T + R) = (P^{-1} + B^T R^{-1} B) P B^T$$

$$B^T R^{-1} B P B^T + B^T = B^T + B^T R^{-1} B P B^T$$

$$\begin{aligned}
 & X^T r^* \\
 & X^T (X X^T + cI)^{-1} \quad \text{are they equal? } (X^T X + cI)^{-1} X^T \\
 & B^T (B B^T + R) = (P^{-1} + B^T R^{-1} B)^{-1} X^T R^{-1} \\
 & \quad \uparrow \quad \uparrow \\
 & P = I \quad cI \\
 & = c(I^{-1} + c X^T X)^{-1} X^T I^{-1} \\
 & = c \frac{1}{c} (I + \frac{1}{c} X^T X)^{-1} X^T I^{-1} \\
 & = (cI + X^T X)^{-1} X^T
 \end{aligned}$$

∵ equal
 Kernel ridge
 ||
 ridge
 in linear case