## **Regularized Linear Regression**

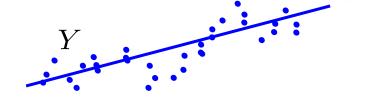
Aarti Singh

Machine Learning 10-315 Sept 22, 2021



### Linear Regression (Matrix-vector form)

$$\widehat{f}_n^L = \arg\min_{f \in \mathcal{F}_L} \frac{1}{n} \sum_{i=1}^n (f(X_i) - Y_i)^2 \qquad \qquad Y$$





$$\widehat{\beta} = \arg\min_{\beta} \frac{1}{n} \sum_{i=1}^{n} (X_i \beta - Y_i)^2$$

$$= \arg\min_{\beta} \frac{1}{n} (\mathbf{A}\beta - \mathbf{Y})^T (\mathbf{A}\beta - \mathbf{Y})$$

$$\widehat{f}_n^L(X) = X\widehat{\beta}$$

$$\mathbf{A} = \begin{bmatrix} X_1 \\ \vdots \\ X_n \end{bmatrix} = \begin{bmatrix} X_1^{(1)} & \dots & X_1^{(p)} \\ \vdots & \ddots & \vdots \\ X_n^{(1)} & \dots & X_n^{(p)} \end{bmatrix} \quad \mathbf{Y} = \begin{bmatrix} \mathbf{Y}_1 \\ \vdots \\ \mathbf{Y}_n \end{bmatrix}$$

## Least Square solution satisfies Normal Equations

$$\left. \frac{\partial J(\beta)}{\partial \beta} \right|_{\widehat{\beta}} = 0$$
 gives  $(\mathbf{A}^T \mathbf{A}) \widehat{\beta} = \mathbf{A}^T \mathbf{Y}$ 

If  $(\mathbf{A}^T \mathbf{A})$  is invertible,

1) If dimension p not too large, analytical solution:

$$\widehat{\beta} = (\mathbf{A}^T \mathbf{A})^{-1} \mathbf{A}^T \mathbf{Y} \qquad \widehat{f}_n^L(X) = X \widehat{\beta}$$

2) If dimension p is large, computing inverse is expensive  $O(p^3)$  Gradient descent since objective is convex ( $A^TA \succeq 0$ )

$$\beta^{t+1} = \beta^t - \frac{\alpha}{2} \frac{\partial J(\beta)}{\partial \beta} \Big|_t$$
$$= \beta^t - \alpha \mathbf{A}^T (\mathbf{A} \beta^t - Y)$$

# Linear regression solution satisfies Normal Equations

$$(\mathbf{A}^T\mathbf{A})\widehat{\boldsymbol{\beta}} = \mathbf{A}^T\mathbf{Y}$$

$$\mathbf{p} \times \mathbf{p} \quad \mathbf{p} \times \mathbf{1}$$

$$\mathbf{p} \times \mathbf{1}$$

When is  $(\mathbf{A}^T\mathbf{A})$  invertible ? Recall: Full rank matrices are invertible. What is rank of  $(\mathbf{A}^T\mathbf{A})$  ?

If 
$$\mathbf{A} = \mathbf{U}\mathbf{S}\mathbf{V}^{ op}$$
, then

normal equations 
$$(\mathbf{S}\mathbf{V}^{\top})\hat{\beta} = (\mathbf{U}^{\top}\mathbf{Y})$$
 requations in p unknowns.

$$(\mathbf{V} \mathbf{S} \mathbf{V}^{\mathsf{T}} \mathbf{V} \mathbf{S} \mathbf{V}^{\mathsf{T}}) \hat{\boldsymbol{\beta}} = \mathbf{V} \mathbf{S} \mathbf{U}^{\mathsf{T}} \mathbf{Y}$$

$$\mathbf{V}^{\mathsf{T}} \mathbf{S}^{\mathsf{T}} \mathbf{V}^{\mathsf{T}} \hat{\boldsymbol{\beta}} = \mathbf{V}^{\mathsf{T}} \mathbf{V}^{\mathsf{T}} \mathbf{Y}$$

$$= (\mathbf{U}^{\mathsf{T}} \mathbf{Y})$$

$$\mathbf{S}^{\mathsf{T}} \mathbf{S}^{\mathsf{T}} \mathbf{V}^{\mathsf{T}} \hat{\boldsymbol{\beta}} = \mathbf{S}^{\mathsf{T}} \mathbf{S} \mathbf{U}^{\mathsf{T}} \mathbf{Y}$$

$$\mathbf{S}^{\mathsf{T}} \mathbf{S}^{\mathsf{T}} = \mathbf{U}^{\mathsf{T}} \mathbf{Y}$$

Under-determined if r < p, hence no unique solution.

What if 
$$(\mathbf{A}^T\mathbf{A})$$
 is not invertible ?  $\checkmark$ 

r equations , p unknowns – underdetermined system of linear equations many feasible solutions

Need to constrain solution further

e.g. bias solution to "small" values of  $\beta$  (small changes in input don't translate to large changes in output)

$$\begin{split} \widehat{\beta}_{\text{MAP}} &= \arg\min_{\beta} \ \sum_{i=1}^{n} (Y_i - X_i \beta)^2 + \lambda \|\beta\|_2^2 \quad \text{Ridge Regression} \\ &= \arg\min_{\beta} \ (\mathbf{A}\beta - \mathbf{Y})^T (\mathbf{A}\beta - \mathbf{Y}) \ + \lambda \|\beta\|_2^2 \quad \lambda \geq 0 \\ \widehat{\beta}_{\text{MAP}} &= (\mathbf{A}^\top \mathbf{A} + \lambda \mathbf{I})^{-1} \mathbf{A}^\top \mathbf{Y} \end{split}$$

$$\widehat{\beta}_{\mathsf{MAP}} = \arg\min_{\beta} \sum_{i=1}^{n} (Y_i - X_i \beta)^2 + \lambda \|\beta\|_2^2$$

= 
$$\arg\min_{\beta} \frac{1}{n} (\mathbf{A}\beta - \mathbf{Y})^T (\mathbf{A}\beta - \mathbf{Y}) + \lambda \|\beta\|_2^2$$

$$\frac{\partial J(\beta)}{\partial \beta} = 2\overrightarrow{A}\overrightarrow{A}\beta - 2\overrightarrow{A}\overrightarrow{Y} + \lambda \frac{\partial \|\beta\|^2}{\partial \beta} = 2\overrightarrow{A}\beta$$

$$= 2(\overrightarrow{A}A + \lambda I)\beta^{-2}\overrightarrow{A}\overrightarrow{Y}$$

$$= 2\lambda\beta$$

= O =) 
$$\hat{\beta} = (AA + \lambda L)'A'Y$$

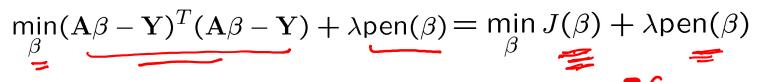
Is 
$$(\mathbf{A}^{ op}\mathbf{A} + \lambda \mathbf{I})$$
 invertible ?

Ridge Regression (12 penalty)

$$\lambda > 0$$

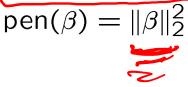
$$eval(A^TA+\lambda F) = eval(A^TA) + \lambda$$

#### **Understanding regularized Least Squares**



720





 $\beta$ s with constant  $J(\beta)$   $\sim$  (level sets of  $J(\beta)$ )

Unregularized Least Squares solution

 $\beta$ s with constant I2 norm (level sets of pen( $\beta$ ))

 $\beta_2$ 

What if  $(\mathbf{A}^T \mathbf{A})$  is not invertible ?

r equations , p unknowns – underdetermined system of linear equations many feasible solutions

Need to constrain solution further

e.g. bias solution to "small" values of  $\beta$  (small changes in input don't translate to large changes in output)

$$\widehat{\beta}_{\mathsf{MAP}} = \arg\min_{\beta} \sum_{i=1}^{n} (Y_i - X_i \beta)^2 + \lambda \|\beta\|_2^2$$

Ridge Regression (12 penalty)

$$\widehat{\beta}_{\mathsf{MAP}} = \arg\min_{\beta} \sum_{i=1}^{n} (Y_i - X_i \beta)^2 + \lambda \|\beta\|_1$$



Many  $\beta$  can be zero – many inputs are irrelevant to prediction in high-dimensional settings (typically intercept term not penalized)

What if  $(\mathbf{A}^T \mathbf{A})$  is not invertible ?

r equations , p unknowns – underdetermined system of linear equations many feasible solutions

Need to constrain solution further

e.g. bias solution to "small" values of  $\beta$  (small changes in input don't translate to large changes in output)

$$\widehat{\beta}_{\mathsf{MAP}} = \arg\min_{\beta} \sum_{i=1}^{n} (Y_i - X_i \beta)^2 + \lambda \|\beta\|_2^2$$

Ridge Regression (12 penalty)

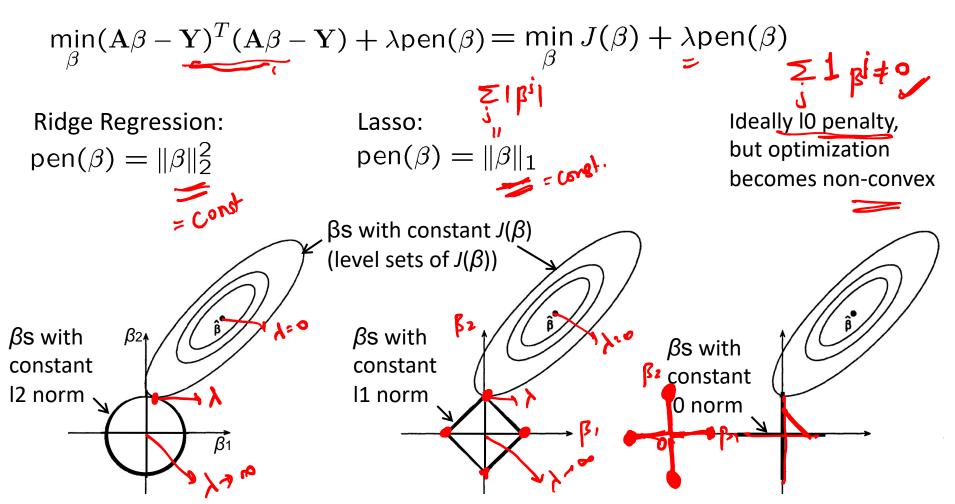


$$\widehat{\beta}_{\mathsf{MAP}} = \arg\min_{\beta} \sum_{i=1}^{n} (Y_i - X_i \beta)^2 + \lambda \|\beta\|_1$$

 $\begin{array}{c} \lambda \geq 0 \\ \text{Lasso} \\ \text{(I1 penalty)} \end{array}$ 

No closed form solution, but can optimize using sub-gradient descent (packages available)

## Ridge Regression vs Lasso



Lasso (11 penalty) results in sparse solutions – vector with more zero coordinates Good for high-dimensional problems – don't have to store all coordinates, interpretable solution!



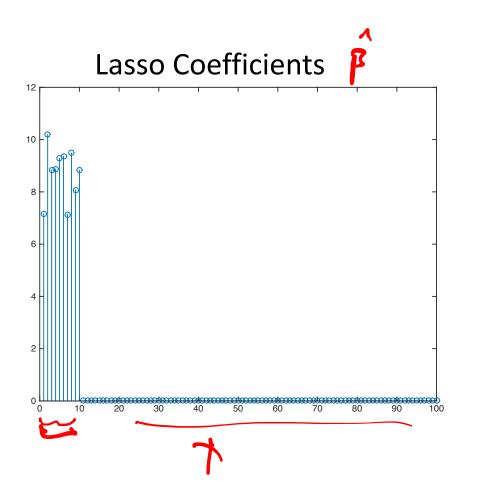
## Matlab example

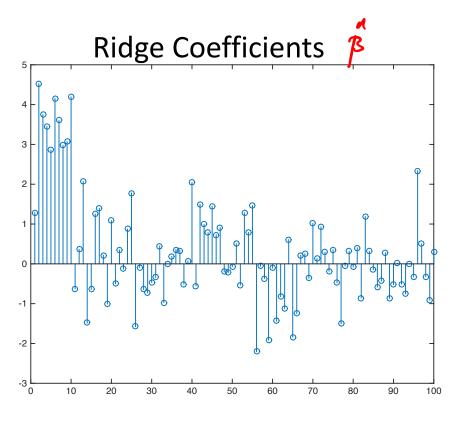
```
clear all
                                       lassoWeights = lasso(X,Y,'Lambda',1,
close all
                                       'Alpha', 1.0);
                                       Ylasso = Xtest*lassoWeights;
n = 80;
        % datapoints
                                       norm(Ytest-Ylasso)
p = 100; % features
k = 10; % non-zero features
                                       ridgeWeights = lasso(X,Y,'Lambda',1,
                                       'Alpha', 0.0001);
rng(20);
                                       Yridge = Xtest*ridgeWeights;
X = randn(n,p);
                                       norm(Ytest-Yridge)
weights = zeros(p,1);
weights(1:k) = randn(k,1)+10;
                                       stem(lassoWeights)
noise = randn(n,1) * 0.5;
                                       pause
Y = X*weights + noise;
                                       stem(ridgeWeights)
Xtest = randn(n,p);
noise = randn(n,1) * 0.5;
Ytest = Xtest*weights + noise;
```

## Matlab example

Test MSE = 33.7997

Test MSE = 185.9948





# 7 (41- Xip) Least Squares and M(C)LE

Intuition: Signal plus (zero-mean) Noise model

$$Y = f^*(X) + \epsilon = X\beta^* + \epsilon$$

$$\epsilon \sim \mathcal{N}(0, \sigma^2 \mathbf{I}) \quad Y \sim \mathcal{N}(X\beta^*, \sigma^2 \mathbf{I})$$

$$\widehat{\beta}_{\text{total}} = - \arg\max_{\mathbf{I}} \log n(\{Y_i\}^n) + \beta \sigma^2 \{X_i\}^n$$

 $\widehat{\beta}_{\text{MLE}} = \arg\max_{\beta} \log p(\{Y_i\}_{i=1}^n | \beta, \sigma^2, \{X_i\}_{i=1}^n) \prod_{i=1}^N p(Y_i | X_i; \beta, \sigma^i)$ 

Conditional log likelihood

= ary max 
$$\frac{1}{2}$$
 by  $\frac{1}{2}$   $\frac$ 

Least Square Estimate is same as Maximum Conditional Likelihood Estimate under a Gaussian noise model!