# **Regularized Linear Regression**

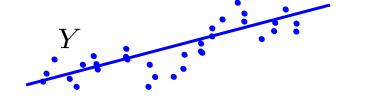
Aarti Singh

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

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

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

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 
$$(\mathbf{S}\mathbf{V}^{\top})\hat{\beta} = (\mathbf{U}^{\top}\mathbf{Y})$$

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

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)

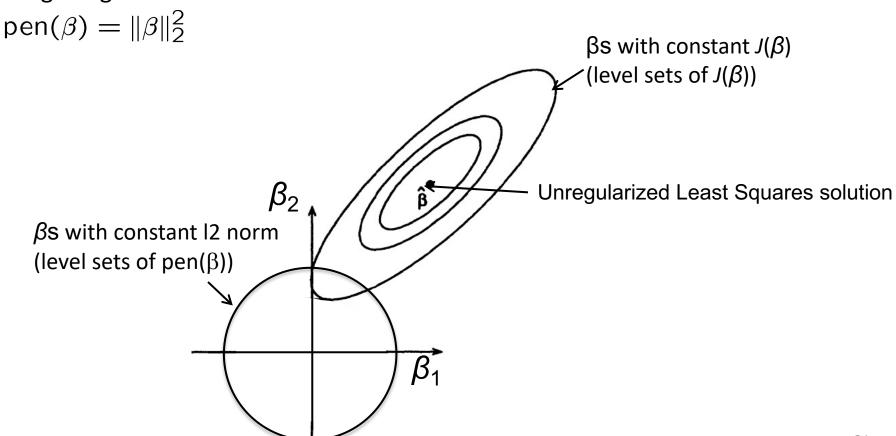
$$\begin{split} \widehat{\beta}_{\mathsf{MAP}} &= \arg\min_{\beta} \sum_{i=1}^{n} (Y_i - X_i \beta)^2 + \lambda \|\beta\|_2^2 & \mathsf{Ridge Regression} \\ &= \arg\min_{\beta} \ \ (\mathbf{A}\beta - \mathbf{Y})^T (\mathbf{A}\beta - \mathbf{Y}) \ + \lambda \|\beta\|_2^2 & \lambda \geq 0 \\ \widehat{\beta}_{\mathsf{MAP}}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

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

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

$$\min_{\beta} (\mathbf{A}\beta - \mathbf{Y})^T (\mathbf{A}\beta - \mathbf{Y}) + \lambda \mathrm{pen}(\beta) = \min_{\beta} J(\beta) + \lambda \mathrm{pen}(\beta)$$

#### Ridge Regression:



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)

$$\widehat{\beta}_{\mathsf{MAP}} = \arg\min_{\beta} \sum_{i=1}^{n} (Y_i - X_i \beta)^2 + \lambda \|\beta\|_2^2 \qquad \begin{array}{l} \mathsf{Ridge \ Regression} \\ \mathsf{(I2 \ penalty)} \end{array}$$

$$\widehat{\beta}_{\mathsf{MAP}} = \arg\min_{\beta} \sum_{i=1}^n (Y_i - X_i \beta)^2 + \lambda \|\beta\|_1 \qquad \text{Lasso} \tag{I1 penalty}$$

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)

$$\widehat{\beta}_{\mathsf{MAP}} = \arg\min_{\beta} \sum_{i=1}^{n} (Y_i - X_i \beta)^2 + \lambda \|\beta\|_2^2 \qquad \begin{array}{l} \mathsf{Ridge \ Regression} \\ \mathsf{(12 \ penalty)} \end{array}$$

$$\widehat{\beta}_{\mathsf{MAP}} = \arg\min_{\beta} \sum_{i=1}^{n} (Y_i - X_i \beta)^2 + \lambda \|\beta\|_1 \qquad \text{Lasso} \\ \text{(I1 penalty)}$$

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

#### Ridge Regression vs Lasso

$$\min_{\beta} (\mathbf{A}\beta - \mathbf{Y})^T (\mathbf{A}\beta - \mathbf{Y}) + \lambda \mathrm{pen}(\beta) = \min_{\beta} J(\beta) + \lambda \mathrm{pen}(\beta)$$

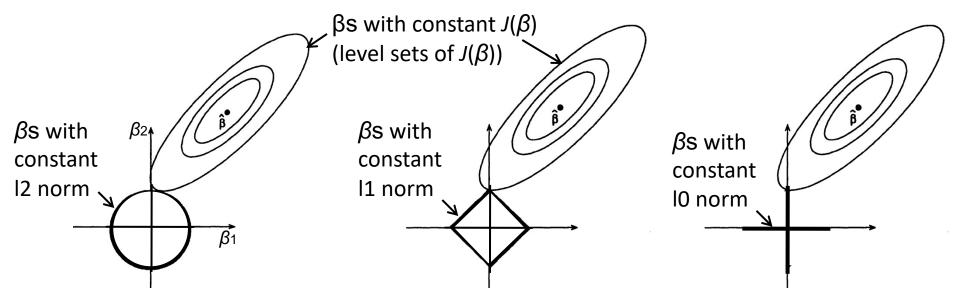
Ridge Regression:

$$pen(\beta) = \|\beta\|_2^2$$

Lasso:

$$pen(\beta) = \|\beta\|_1$$

Ideally IO penalty, but optimization becomes non-convex



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

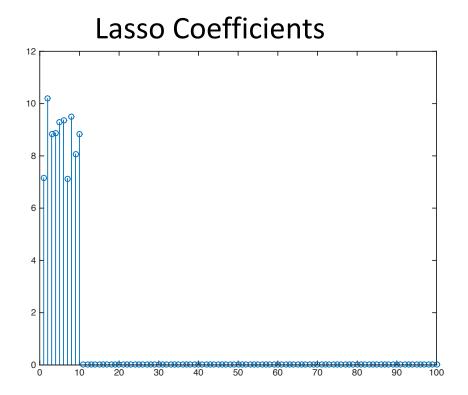
```
clear all
                                      lassoWeights = lasso(X,Y,'Lambda',1,
                                      'Alpha', 1.0);
close all
                                      Ylasso = Xtest*lassoWeights;
n = 80; % datapoints
                                      norm(Ytest-Ylasso)
p =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

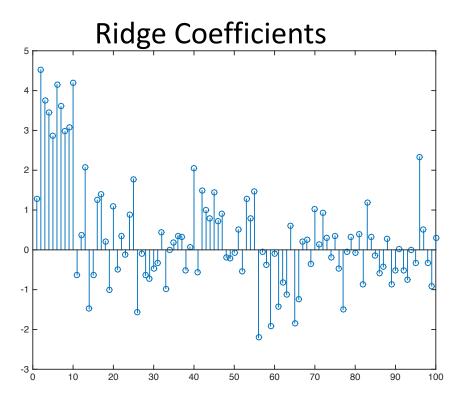
Ytest = Xtest\*weights + noise;

## Matlab example

Test MSE = 33.7997

Test MSE = 185.9948





## Least Squares and M(C)LE

Intuition: Signal plus (zero-mean) Noise model

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

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

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

Conditional log likelihood

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

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

#### Regularized Least Squares and M(C)AP

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

$$\widehat{\beta}_{\text{MAP}} = \arg\max_{\beta} \log p(\{Y_i\}_{i=1}^n | \beta, \sigma^2, \{X_i\}_{i=1}^n + \log p(\beta)$$
 Conditional log likelihood log prior

I) Gaussian Prior

$$\beta \sim \mathcal{N}(0, \tau^2 \mathbf{I})$$

$$eta \sim \mathcal{N}(0, au^2\mathbf{I})$$
  $p(eta) \propto e^{-eta^Teta/2 au^2}$ 

$$\widehat{\beta}_{\mathsf{MAP}} = \arg\min_{\beta} \sum_{i=1}^n (Y_i - X_i \beta)^2 + \lambda \|\beta\|_2^2 \qquad \underset{\mathsf{constant}(\sigma^2, \tau^2)}{\mathsf{Ridge Regression}}$$

#### Regularized Least Squares and M(C)AP

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

$$\widehat{\beta}_{\text{MAP}} = \arg\max_{\beta} \log p(\{Y_i\}_{i=1}^n | \beta, \sigma^2, \{X_i\}_{i=1}^n + \log p(\beta)$$
 Conditional log likelihood log prior

II) Laplace Prior

$$eta_i \stackrel{iid}{\sim} \mathsf{Laplace}(\mathsf{0},t) \qquad \qquad p(eta_i) \propto e^{-|eta_i|/t}$$

$$p(\beta_i) \propto e^{-|\beta_i|/t}$$

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

## **Polynomial Regression**

degree m

Univariate (1-dim)  $f(X) = \beta_0 + \beta_1 X + \beta_2 X^2 + \cdots + \beta_m X^m = \mathbf{X}\beta$  case:

where 
$$\mathbf{X} = [1 \ X \ X^2 \dots X^m]$$
 ,  $\beta = [\beta_1 \dots \beta_m]^T$ 

$$\widehat{\beta} = (\mathbf{A}^T \mathbf{A})^{-1} \mathbf{A}^T \mathbf{Y}$$

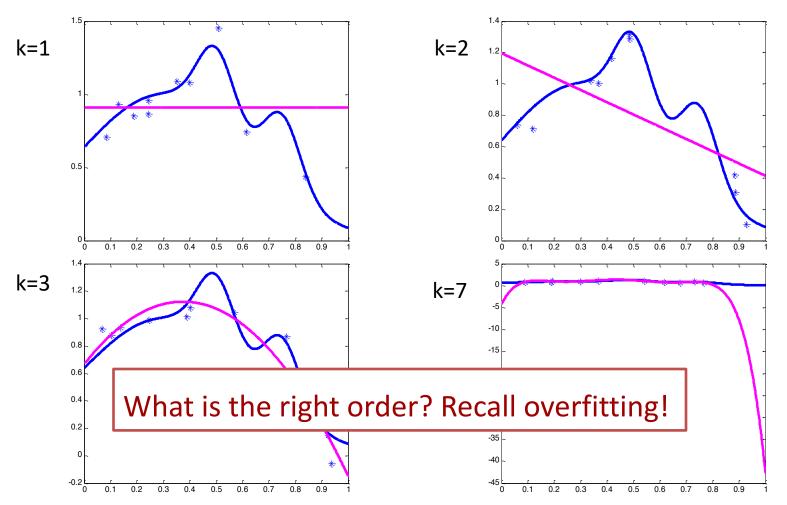
$$\widehat{f}_n(X) = \mathbf{X}\widehat{\beta}$$

where 
$$\mathbf{A}=\left[\begin{array}{ccccc} 1 & X_1 & X_1^2 & \dots & X_1^m \\ \vdots & & \ddots & \vdots \\ 1 & X_n & X_n^2 & \dots & X_n^m \end{array}\right]$$

Multivariate (p-dim) 
$$f(X) = \beta_0 + \beta_1 X^{(1)} + \beta_2 X^{(2)} + \dots + \beta_p X^{(p)}$$
 case: 
$$+ \sum_{i=1}^p \sum_{j=1}^p \beta_{ij} X^{(i)} X^{(j)} + \sum_{i=1}^p \sum_{j=1}^p \sum_{k=1}^p X^{(i)} X^{(j)} X^{(k)} + \dots \text{ terms up to degree m}$$

# **Polynomial Regression**

Polynomial of order k, equivalently of degree up to k-1



#### Regression with nonlinear features

$$f(X) = \sum_{j=0}^{m} \beta_j X^j = \sum_{j=0}^{m} \beta_j \phi_j(X)$$
Weight of each feature features 
$$\phi_0(X)$$

$$\phi_1(X)$$

In general, use any nonlinear features

e.g. e<sup>X</sup>, log X, 1/X, sin(X), ...

$$\widehat{\beta} = (\mathbf{A}^T \mathbf{A})^{-1} \mathbf{A}^T \mathbf{Y} \qquad \mathbf{A} = \begin{bmatrix} \phi_0(X_1) \ \phi_1(X_1) \ \dots \ \phi_m(X_1) \\ \vdots \ \phi_0(X_n) \ \phi_1(X_n) \ \dots \ \phi_m(X_n) \end{bmatrix}$$

$$\widehat{f}_n(X) = \mathbf{X}\widehat{\beta}$$
  $\mathbf{X} = [\phi_0(X) \ \phi_1(X) \ \dots \ \phi_m(X)]$