

Linear Regression contd...

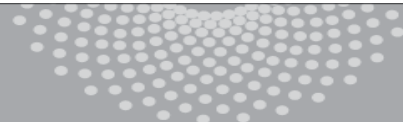
Aarti Singh

Machine Learning 10-315

Sept 28, 2020



MACHINE LEARNING DEPARTMENT



Carnegie Mellon.
School of Computer Science

Linear Regression

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

$$f(X_i) = X_i \beta$$



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

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

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

$$\mathbf{A} = \begin{bmatrix} X_1 \\ \vdots \\ X_n \end{bmatrix} = \begin{bmatrix} X_1^{(1)} & \dots & X_1^{(p)} \\ \vdots & \ddots & \vdots \\ X_n^{(1)} & \dots & X_n^{(p)} \end{bmatrix} \quad \mathbf{Y} = \begin{bmatrix} Y_1 \\ \vdots \\ Y_n \end{bmatrix}$$

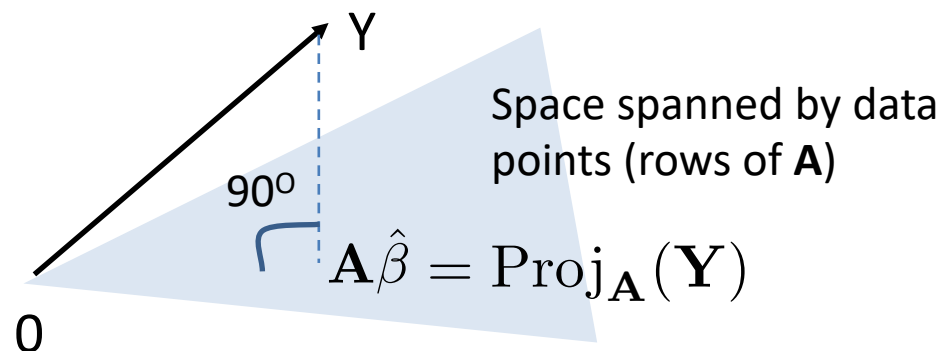
Linear regression solution satisfies Normal Equations

$$\underbrace{(\mathbf{A}^T \mathbf{A})}_{p \times p} \underbrace{\hat{\boldsymbol{\beta}}}_{p \times 1} = \underbrace{\mathbf{A}^T \mathbf{Y}}_{p \times 1}$$

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

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

Predicted labels for training points $\mathbf{A} \hat{\boldsymbol{\beta}} = \text{Proj}_{\mathbf{A}}(\mathbf{Y})$



Linear regression solution satisfies Normal Equations

$$\underbrace{(\mathbf{A}^T \mathbf{A})}_{p \times p} \underbrace{\hat{\boldsymbol{\beta}}}_{p \times 1} = \underbrace{\mathbf{A}^T \mathbf{Y}}_{p \times 1}$$

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

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

Later: When is $(\mathbf{A}^T \mathbf{A})$ invertible?

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

Now: What if $(\mathbf{A}^T \mathbf{A})$ is invertible but expensive (p very large)?

Gradient Descent

Even when $(\mathbf{A}^T \mathbf{A})$ is invertible, might be computationally expensive if \mathbf{A} is huge.

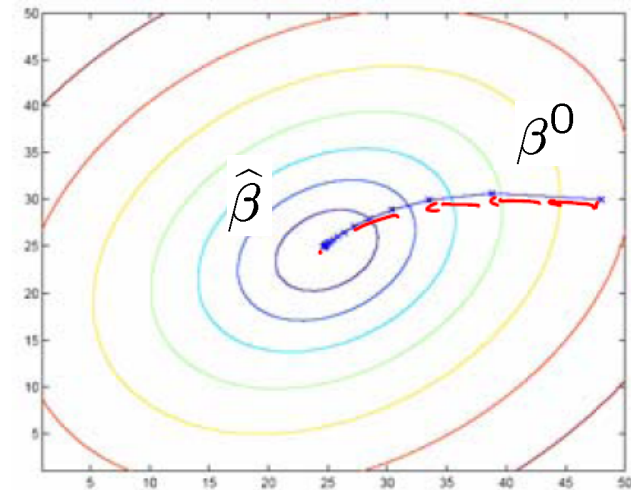
$$\hat{\beta} = \arg \min_{\beta} \frac{1}{n} (\mathbf{A}\beta - \mathbf{Y})^T (\mathbf{A}\beta - \mathbf{Y}) = \arg \min_{\beta} J(\beta)$$

Since $J(\beta)$ is convex, move along negative of gradient

Initialize: β^0

$$\begin{aligned} \text{Update: } \beta^{t+1} &= \beta^t - \frac{\alpha \partial J(\beta)}{2 \partial \beta} \Big|_t \\ &= \beta^t - \alpha \underbrace{\mathbf{A}^T (\mathbf{A}\beta^t - \mathbf{Y})}_{0 \text{ if } \hat{\beta} = \beta^t} \end{aligned}$$

step size



Stop: when some criterion met e.g. fixed # iterations, or $\frac{\partial J(\beta)}{\partial \beta} \Big|_{\beta^t} < \epsilon$.

Linear regression solution satisfies Normal Equations

$$\underbrace{(\mathbf{A}^T \mathbf{A})}_{p \times p} \underbrace{\hat{\boldsymbol{\beta}}}_{p \times 1} = \underbrace{\mathbf{A}^T \mathbf{Y}}_{p \times 1}$$

When is $(\mathbf{A}^T \mathbf{A})$ invertible ?

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

Null space argument

Linear regression solution satisfies Normal Equations

$$\underbrace{(\mathbf{A}^T \mathbf{A})}_{p \times p} \underbrace{\hat{\boldsymbol{\beta}}}_{p \times 1} = \underbrace{\mathbf{A}^T \mathbf{Y}}_{p \times 1}$$

When is $(\mathbf{A}^T \mathbf{A})$ invertible ?

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

$\text{Rank}(\mathbf{A}^T \mathbf{A}) =$ number of non-zero eigenvalues of $(\mathbf{A}^T \mathbf{A}) =$ number of non-zero singular values of $\mathbf{A} \leq \min(n, p)$ since \mathbf{A} is $n \times p$

So, $\text{rank}(\mathbf{A}^T \mathbf{A}), r \leq \min(n, p)$ not invertible if $r < p$ (e.g. $n < p$
i.e. high-dimensional setting)

Linear regression solution satisfies Normal Equations

$$\underbrace{(\mathbf{A}^T \mathbf{A})}_{p \times p} \underbrace{\hat{\boldsymbol{\beta}}}_{p \times 1} = \underbrace{\mathbf{A}^T \mathbf{Y}}_{p \times 1}$$

When is $(\mathbf{A}^T \mathbf{A})$ invertible?

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

If $\mathbf{A} = \mathbf{U} \mathbf{S} \mathbf{V}^T$, then normal equations $\underbrace{(\mathbf{S} \mathbf{V}^T)}_{r \times p} \hat{\boldsymbol{\beta}} = \underbrace{(\mathbf{U}^T \mathbf{Y})}_{r \times 1}$
 $\underbrace{\mathbf{S}}_{s-r \times r}$

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

Regularized Linear Regression

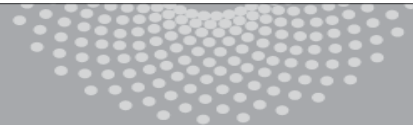
Aarti Singh

Machine Learning 10-315

Sept 28, 2020



MACHINE LEARNING DEPARTMENT



Carnegie Mellon.
School of Computer Science

Regularized Least Squares

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 β (small changes in input don’t translate to large changes in output)

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

Ridge Regression
(l2 penalty)

$$= \arg \min_{\beta} (\mathbf{A}\beta - \mathbf{Y})^T (\mathbf{A}\beta - \mathbf{Y}) + \lambda \|\beta\|_2^2 \quad \lambda \geq 0$$

$$\hat{\beta}_{\text{MAP}} = (\mathbf{A}^T \mathbf{A} + \lambda \mathbf{I})^{-1} \mathbf{A}^T \mathbf{Y}$$

Regularized Least Squares

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

Ridge Regression
(l2 penalty)

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

$$\lambda \geq 0$$

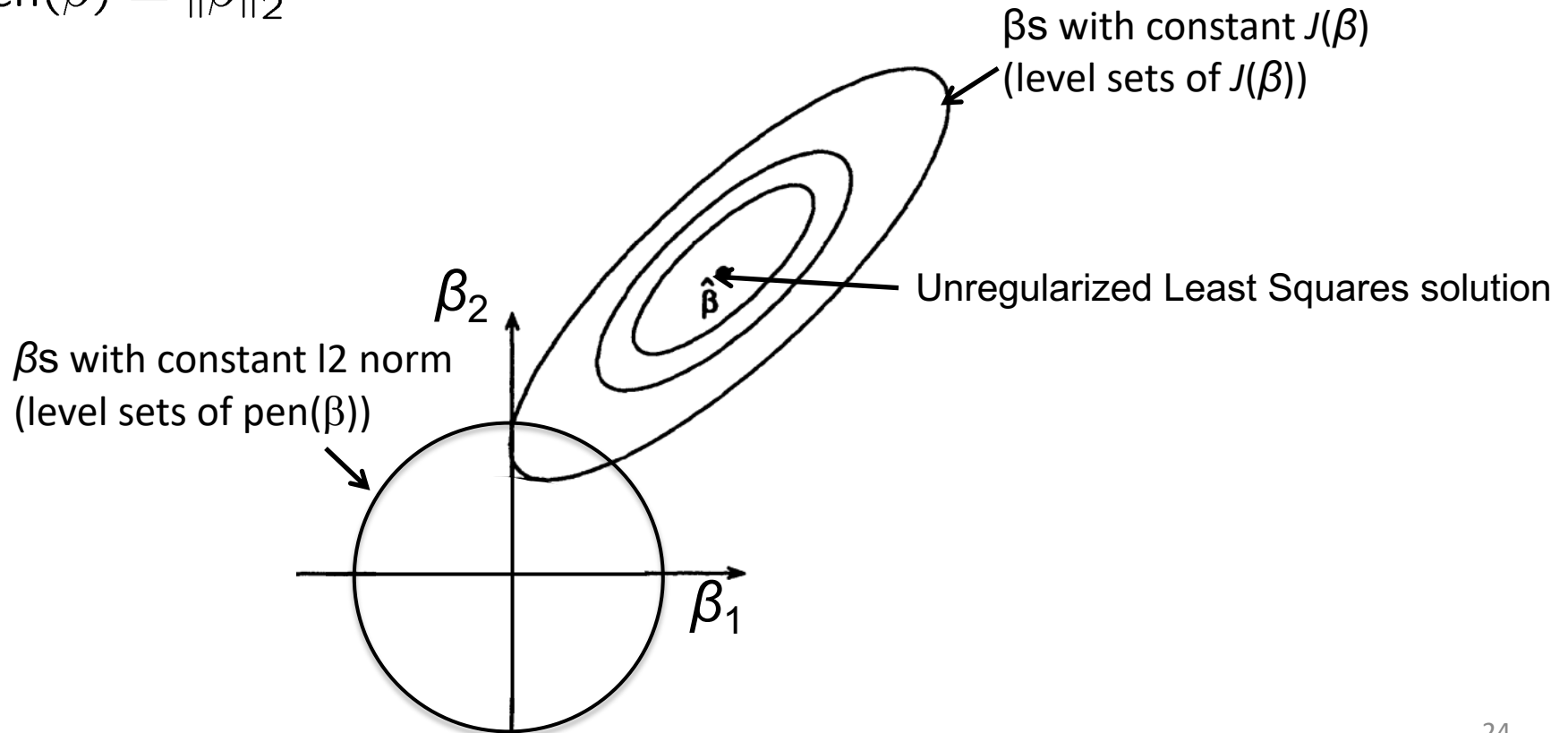
Is $(\mathbf{A}^T \mathbf{A} + \lambda \mathbf{I})$ invertible ?

Understanding regularized Least Squares

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

Ridge Regression:

$$\text{pen}(\beta) = \|\beta\|_2^2$$



Regularized Least Squares

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 β (small changes in input don’t translate to large changes in output)

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

Ridge Regression
(l2 penalty)

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

Lasso
(l1 penalty)

$$\lambda \geq 0$$

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

Regularized Least Squares

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 β (small changes in input don’t translate to large changes in output)

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

Ridge Regression
(l2 penalty)

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

Lasso
(l1 penalty)

$$\lambda \geq 0$$

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

Ridge Regression vs Lasso

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

Ridge Regression:

$$\text{pen}(\beta) = \|\beta\|_2^2$$

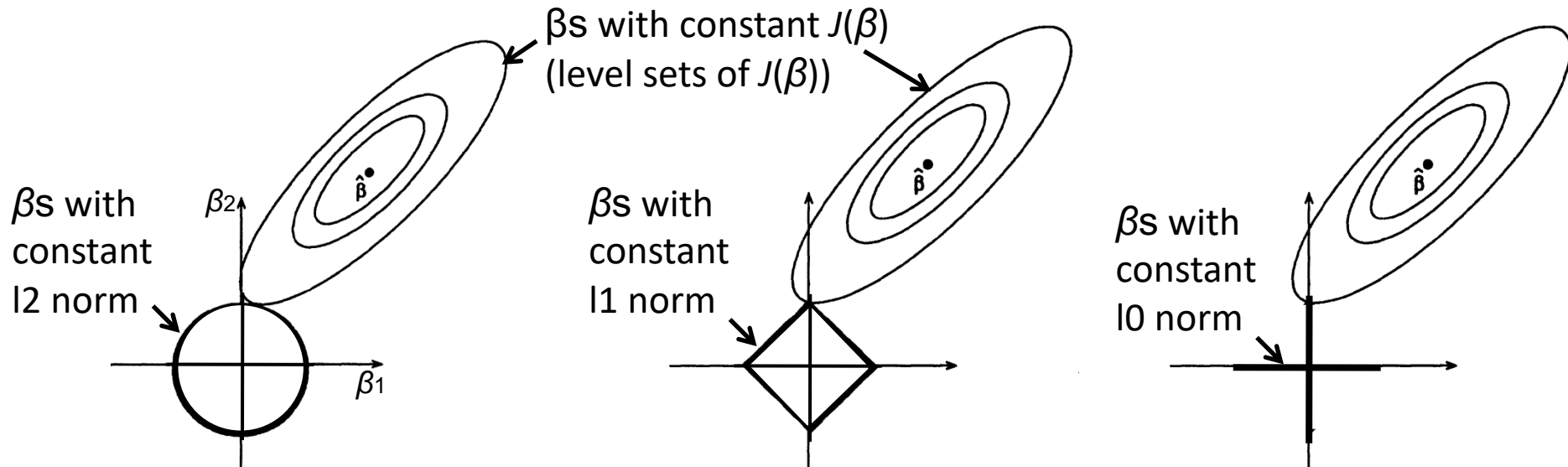
Lasso:

$$\text{pen}(\beta) = \|\beta\|_1$$

Ideally l0 penalty,

but optimization

becomes non-convex



Lasso (l1 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  
close all
```

```
n = 80; % datapoints  
p = 100; % features  
k = 10; % non-zero features
```

```
rng(20);  
X = randn(n,p);  
weights = zeros(p,1);  
weights(1:k) = randn(k,1)+10;  
noise = randn(n,1) * 0.5;  
Y = X*weights + noise;
```

```
Xtest = randn(n,p);  
noise = randn(n,1) * 0.5;  
Ytest = Xtest*weights + noise;
```

```
lassoWeights = lasso(X,Y,'Lambda',1,  
'Alpha', 1.0);  
Ylasso = Xtest*lassoWeights;  
norm(Ytest-Ylasso)
```

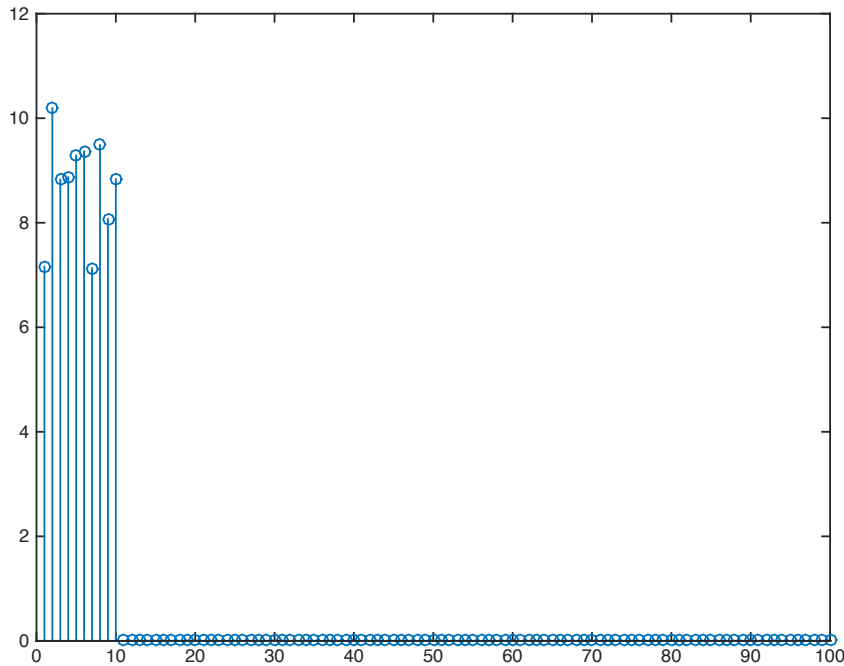
```
ridgeWeights = lasso(X,Y,'Lambda',1,  
'Alpha', 0.0001);  
Yridge = Xtest*ridgeWeights;  
norm(Ytest-Yridge)
```

```
stem(lassoWeights)  
pause  
stem(ridgeWeights)
```


Matlab example

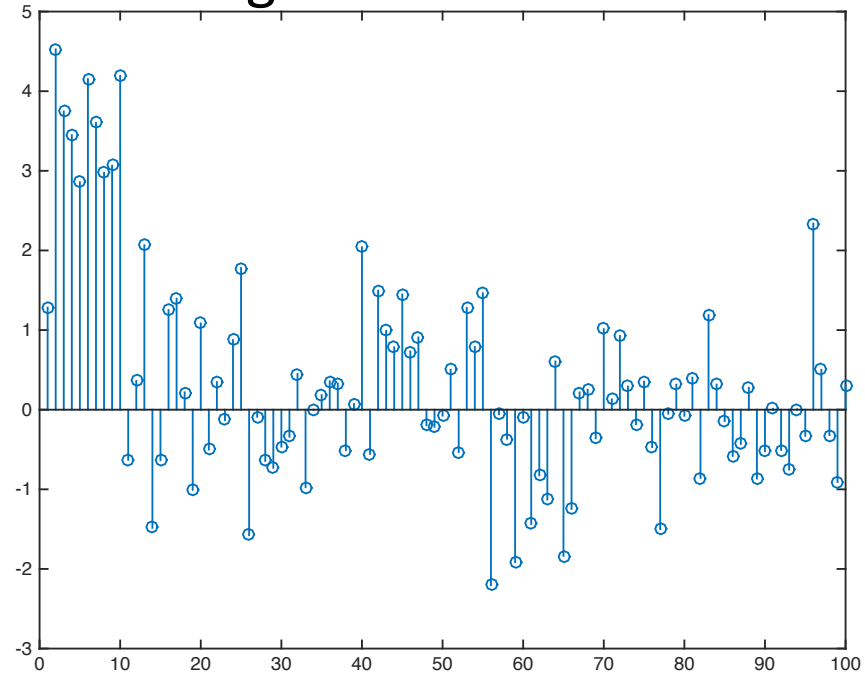
Test MSE = 33.7997

Lasso Coefficients



Test MSE = 185.9948

Ridge Coefficients



Least Squares and M(C)LE

Intuition: Signal plus (zero-mean) Noise model

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

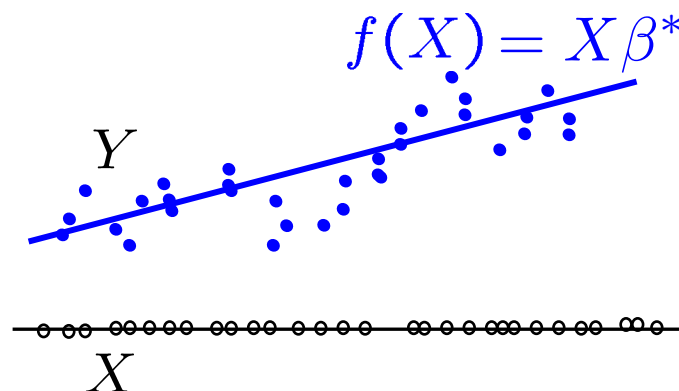
$$\epsilon \sim \mathcal{N}(0, \sigma^2 \mathbf{I}) \quad Y \sim \mathcal{N}(X\beta^*, \sigma^2 \mathbf{I})$$

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

Conditional log likelihood

➤ Breakout

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



Groups 1-10: [Jamboard 1 10](#)

Groups 11-20: [Jamboard 11 20](#)

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

Regularized Least Squares and M(C)AP

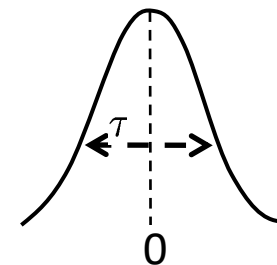
What if $(\mathbf{A}^T \mathbf{A})$ is not invertible ?

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

I) Gaussian Prior

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

$$p(\beta) \propto e^{-\beta^T \beta / 2\tau^2}$$



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

constant(σ^2, τ^2)

Ridge Regression

$$\hat{\beta}_{\text{MAP}} = (\mathbf{A}^T \mathbf{A} + \lambda \mathbf{I})^{-1} \mathbf{A}^T \mathbf{Y}$$

Regularized Least Squares and M(C)AP

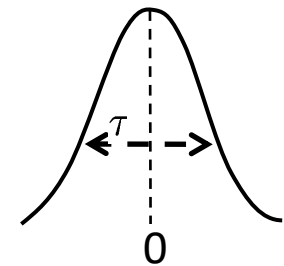
What if $(\mathbf{A}^T \mathbf{A})$ is not invertible ?

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

1) Gaussian Prior

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

$$p(\beta) \propto e^{-\beta^T \beta / 2\tau^2}$$



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

constant(σ^2, τ^2)

Ridge Regression

Prior belief that β is Gaussian with zero-mean biases solution to “small” β

Regularized Least Squares and M(C)AP

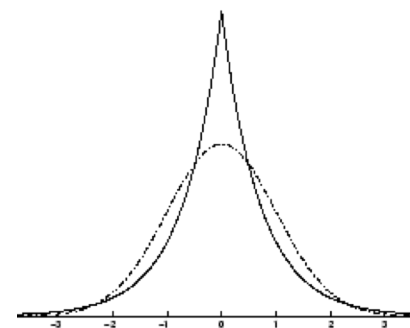
What if $(\mathbf{A}^T \mathbf{A})$ is not invertible ?

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

II) Laplace Prior

$$\beta_i \stackrel{iid}{\sim} \text{Laplace}(0, t)$$

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



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

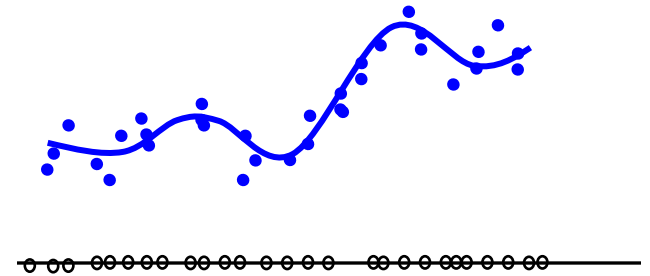
\downarrow
 constant(σ^2, t)

Lasso

Prior belief that β is Laplace with zero-mean biases solution to “sparse” β

Beyond Linear Regression

Polynomial regression
Regression with nonlinear features



Kernelized Ridge Regression (Later)

Local Kernel Regression (Later)

Polynomial Regression

degree m

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

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

$$\hat{\beta} = (\mathbf{A}^T \mathbf{A})^{-1} \mathbf{A}^T \mathbf{Y} \quad \text{or} \quad (\mathbf{A}^T \mathbf{A} + \lambda \mathbf{I})^{-1} \mathbf{A}^T \mathbf{Y} \quad \hat{f}_n(X) = \mathbf{X} \hat{\beta}$$

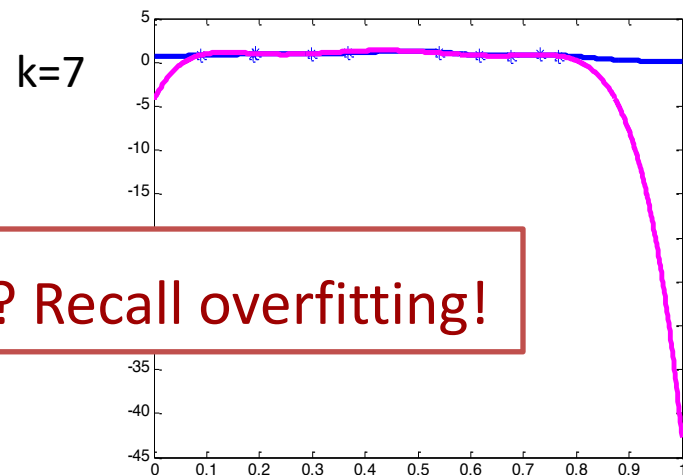
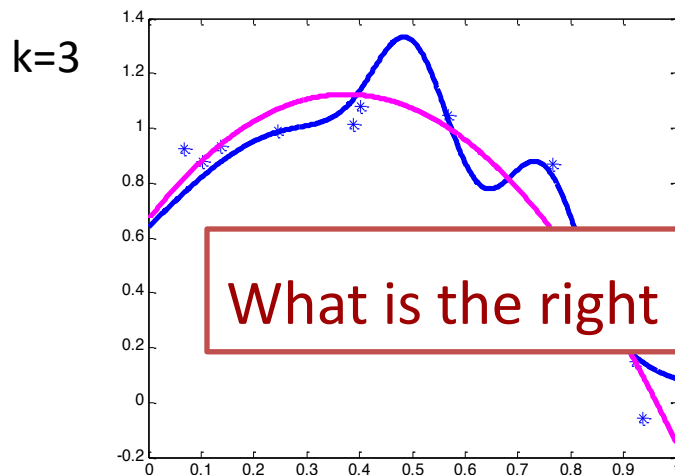
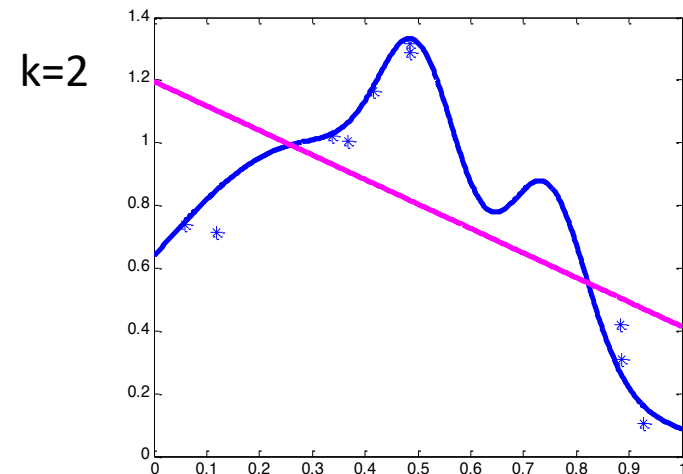
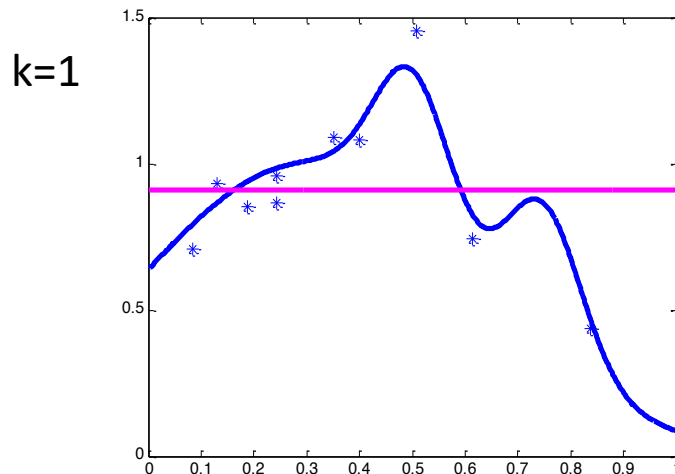
$$\text{where } \mathbf{A} = \begin{bmatrix} 1 & X_1 & X_1^2 & \dots & X_1^m \\ \vdots & & & \ddots & \vdots \\ 1 & X_n & X_n^2 & \dots & X_n^m \end{bmatrix}$$

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

+ ... terms up to degree m

Polynomial Regression

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



What is the right order? Recall overfitting!

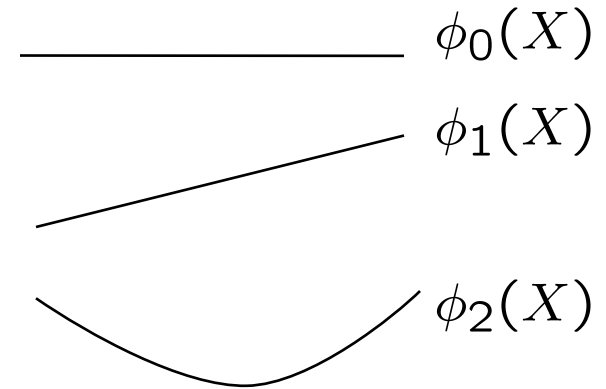
Regression with nonlinear features

$$f(X) = \sum_{j=0}^m \beta_j X^j = \sum_{j=0}^m \beta_j \phi_j(X)$$

Weight of
each feature



Nonlinear
features



In general, use any nonlinear features

e.g. e^X , $\log X$, $1/X$, $\sin(X)$, ...

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

or

$$(\mathbf{A}^T \mathbf{A} + \lambda \mathbf{I})^{-1} \mathbf{A}^T \mathbf{Y}$$

$$\mathbf{A} = \begin{bmatrix} \phi_0(X_1) & \phi_1(X_1) & \dots & \phi_m(X_1) \\ \vdots & & \ddots & \vdots \\ \phi_0(X_n) & \phi_1(X_n) & \dots & \phi_m(X_n) \end{bmatrix}$$

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

$$\mathbf{X} = [\phi_0(X) \ \phi_1(X) \ \dots \ \phi_m(X)]$$