

# Linear Regression contd...

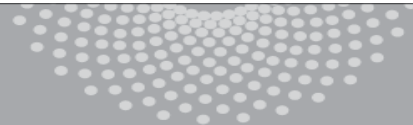
Aarti Singh

Machine Learning 10-315

Sept 28, 2020



**MACHINE LEARNING** DEPARTMENT



**Carnegie Mellon.**  
School of Computer Science

# Linear Regression

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

*Exy*

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

*parameters*



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

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

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

*J(β)*

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

# Linear regression solution satisfies Normal Equations

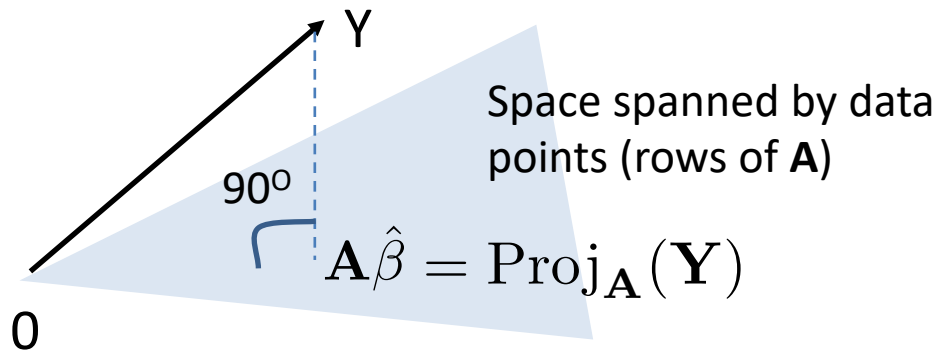
$$(\mathbf{A}^T \mathbf{A}) \hat{\beta} = \mathbf{A}^T \mathbf{Y}$$

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

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

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

Predicted labels for training points  $\hat{\mathbf{A}} \hat{\beta} = \text{Proj}_{\mathbf{A}}(\mathbf{Y})$       $\hat{f}_n^L(\mathbf{A})$   
" "  
 $\hat{\beta}$   
 $\text{Proj}_{\mathbf{A}} = \mathbf{A}(\mathbf{A}^T \mathbf{A})^{-1} \mathbf{A}^T$   
← projects onto



# Linear regression solution satisfies Normal Equations

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

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

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

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

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

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.

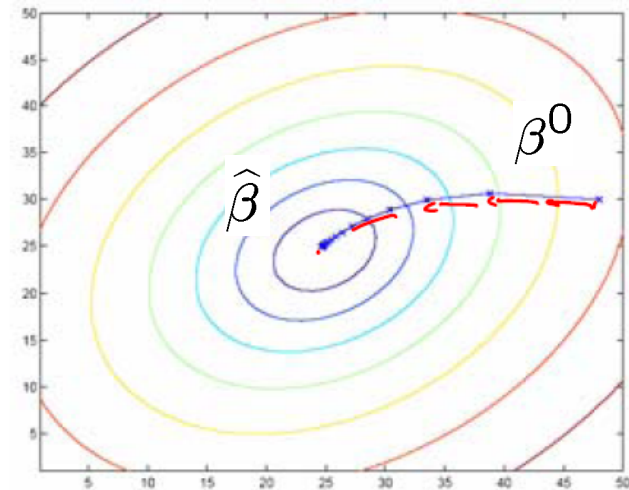
$$\hat{\beta} = \arg \min_{\beta} \frac{1}{n} (\mathbf{A}\beta - \mathbf{Y})^T (\mathbf{A}\beta - \mathbf{Y}) = \arg \min_{\beta} J(\beta)$$

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

Initialize:  $\beta^0$

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

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

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

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

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

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

$\rightarrow$   $\text{null}(M) = \{v : Mv = 0\}$

Null space argument

$\text{null space}(\mathbf{A}^T \mathbf{A}) = \text{null space}(\mathbf{A})$  [Claim]

1) Consider  $v$  s.t.  $\mathbf{A}^T \mathbf{A} v = 0$

$\Rightarrow v^T \mathbf{A}^T \mathbf{A} v = 0$

$\Rightarrow z^T z = 0$

$z^T z = \sum_i z_i^2 = 0$   
 $\Rightarrow z = 0$   
 $\equiv \mathbf{A} v = 0$

2) Consider  $v$  s.t.  $\mathbf{A} v = 0$

$\Rightarrow \mathbf{A}^T \mathbf{A} v = 0$

$\text{rank} = p - \text{dim}(\text{null})$

# Linear regression solution satisfies Normal Equations

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

\* collinear rows or columns  
or •  $n < p$   
high-dim setting  
not enough data

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  $\text{null}(\mathbf{A}^T \mathbf{A}) = \text{null}(\mathbf{A})$

$$\Rightarrow \text{rank}(\mathbf{A}^T \mathbf{A}) = \text{rank}(\mathbf{A})$$

Same number of linearly independent rows & columns

How many?  $\mathbf{A}_{n \times p} = \begin{bmatrix} x_1^{(1)} & \dots & x_1^{(p)} \\ \vdots & \dots & \vdots \\ x_n^{(1)} & \dots & x_n^{(p)} \end{bmatrix}$

$$\text{rank}(\mathbf{A}) \leq \min(n, p) \stackrel{?}{\neq} p \text{ (full rank)}$$

\*

# 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

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

$\rightarrow$   $r$  independent equations  
 $p$  unknowns ( $\hat{\beta}$ )

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  $\underbrace{(\mathbf{S} \mathbf{V}^T)}_{r \times p} \hat{\beta} = \underbrace{(\mathbf{U}^T \mathbf{Y})}_{r \times 1}$

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

# Regularized Linear Regression

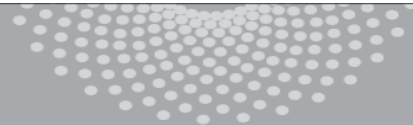
Aarti Singh

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  $\beta$  (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)

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

$$\frac{\partial J(\beta)}{\partial \beta} = 2\mathbf{A}^T\mathbf{A}\beta - 2\mathbf{A}^T\mathbf{Y} + 2\lambda\beta \quad \underbrace{\quad}_{=\beta^T\beta}$$

$$= 2(\mathbf{A}^T\mathbf{A} + \lambda\mathbf{I})\beta - 2\mathbf{A}^T\mathbf{Y} \quad | \quad \hat{\beta}_{\text{MAP}} = 0$$

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

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

all rows & col<sup>m</sup> are independent  
eval  $(\mathbf{A}^T\mathbf{A} + \lambda\mathbf{I}) = \text{eval}(\mathbf{A}^T\mathbf{A}) + \lambda > 0$

Yes!  $\lambda > 0$

$$(\mathbf{A}^T\mathbf{A} + \lambda\mathbf{I})\mathbf{v} = \mathbf{A}^T\mathbf{A}\mathbf{v} + \lambda\mathbf{v} = (\eta + \lambda)\mathbf{v}$$

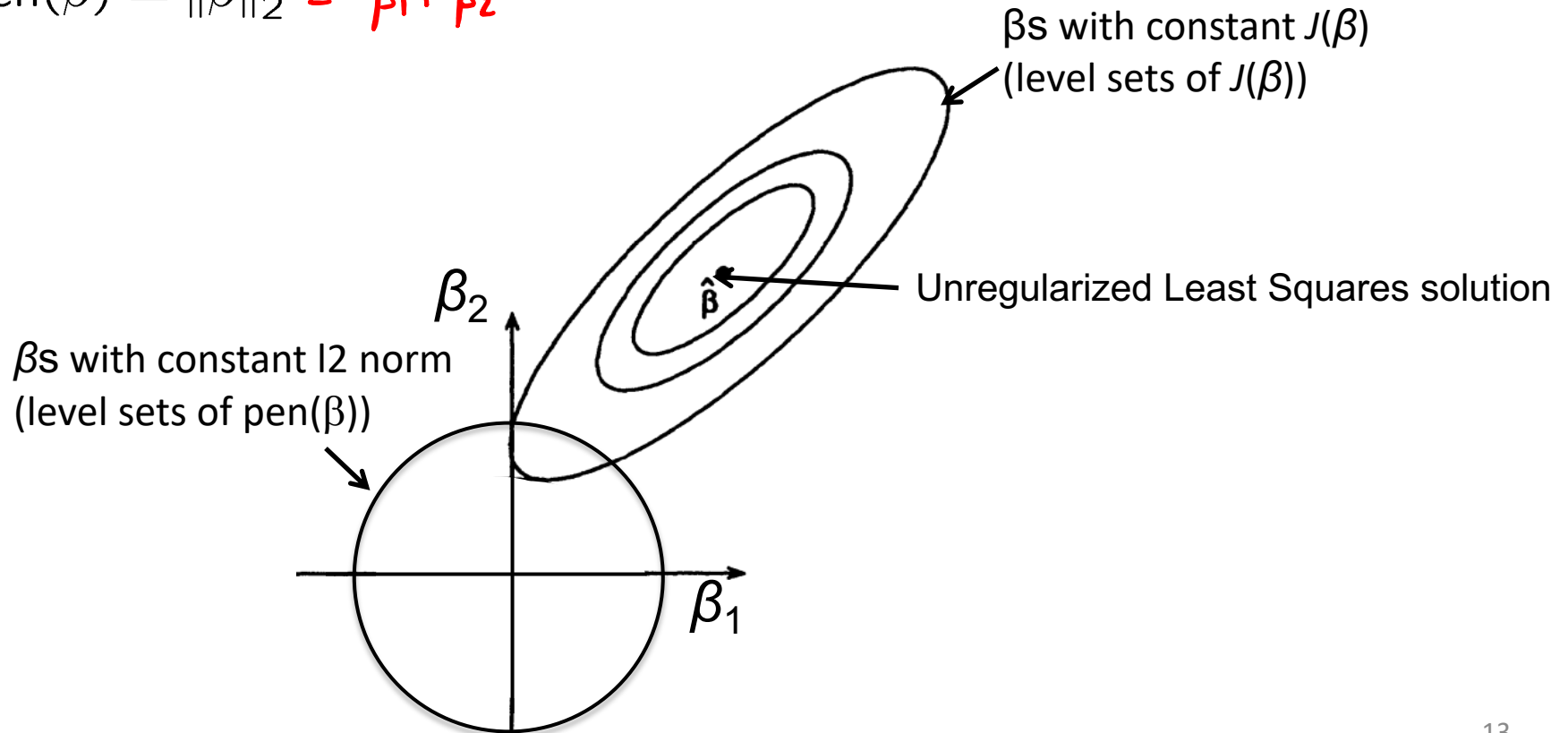
Is  $(\mathbf{A}^T\mathbf{A} + \lambda\mathbf{I})$  invertible?

# 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:

$$\text{pen}(\beta) = \|\beta\|_2^2 = \beta_1^2 + \beta_2^2$$



# 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  $\beta$  (small changes in input don’t translate to large changes in output)

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

$\sum_i \beta_i^2$   
“ $\sum_i \beta_i^2$ ”  
Ridge Regression  
(l2 penalty)

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

$\sum_i |\beta_i|$   
“ $\sum_i |\beta_i|$ ”  
Lasso  
(l1 penalty)

$\lambda \geq 0$

Many  $\beta$  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  $\beta$  (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:

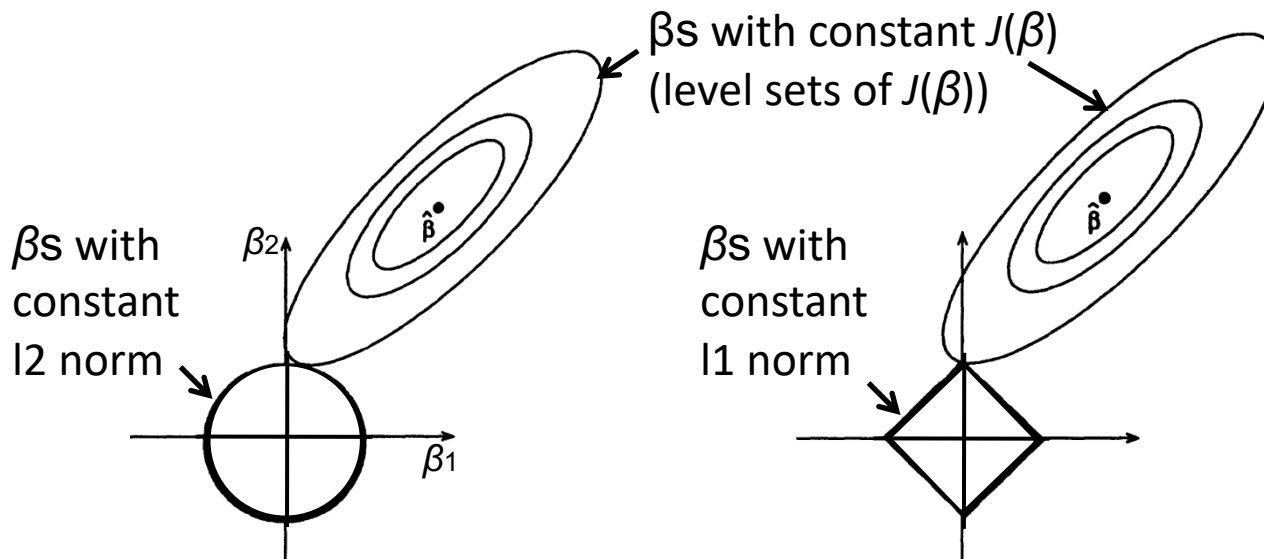
$$\text{pen}(\beta) = \|\beta\|_2^2$$

Lasso:

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

Ideally l0 penalty,  
but optimization  
becomes non-convex

$$\|\beta\|_0 = \sum_i \mathbb{1}_{\{\beta_i \neq 0\}}$$



**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)
```

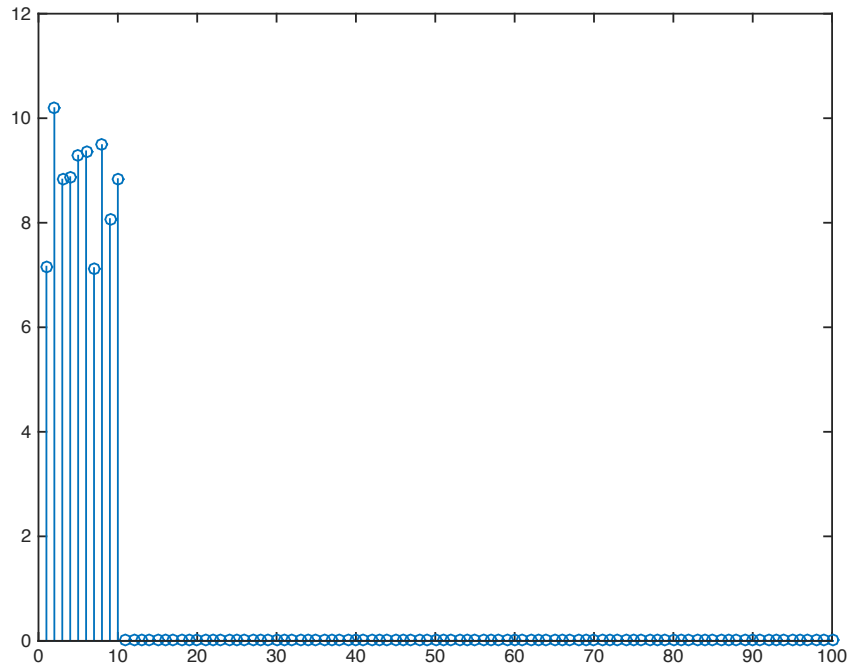
```
ridgeWeights = lasso(X,Y,'Lambda',1,  
'Alpha', 0.0001);  
Yridge = Xtest*ridgeWeights;  
norm(Ytest-Yridge)
```

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

# Matlab example

Test MSE = 33.7997

## Lasso Coefficients



Test MSE = 185.9948

## Ridge Coefficients

