# Regularized Linear Regression
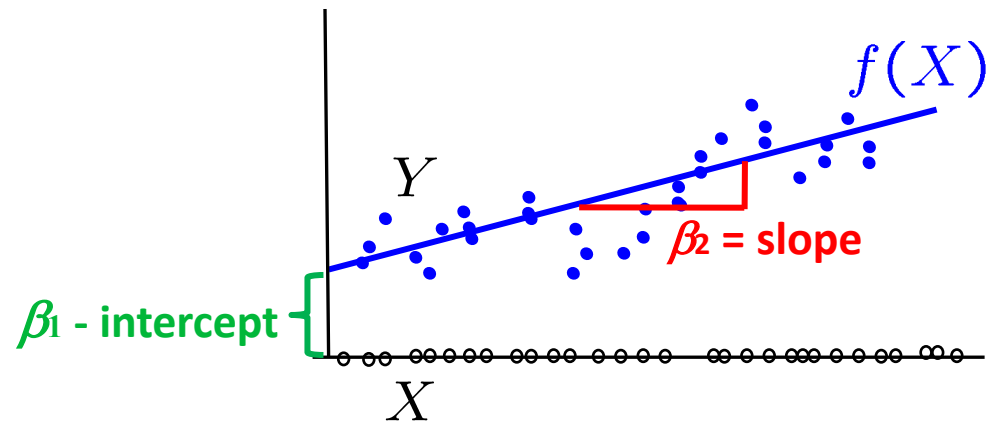
Aarti Singh

Machine Learning 10-315
Oct 28, 2019

# Linear Regression

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

**Least Squares Estimator**

$\mathcal{F}_L$ - Class of Linear functions

Uni-variate case:

$$f(X) = \beta_1 + \beta_2 X$$



$\beta_1$ - intercept

$\beta_2$ = slope

$f(X)$

$Y$

$X$

Multi-variate case:

$$f(X) = X\beta \qquad \text{where} \qquad X = [X^{(1)} \ldots X^{(p)}], \quad \beta = [\beta_1 \ldots \beta_p]^T$$

# Least Squares Estimator

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

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

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

$$\mathbf{A} = \begin{bmatrix} X_1 \\ \vdots \\ X_n \end{bmatrix} = \begin{bmatrix} X_1^{(1)} & \dots & X_1^{(p)} \\ \vdots & \ddots & \vdots \\ X_n^{(1)} & \dots & X_n^{(p)} \end{bmatrix} \qquad \mathbf{Y} = \begin{bmatrix} \mathbf{Y}_1 \\ \vdots \\ \mathbf{Y}_n \end{bmatrix}$$
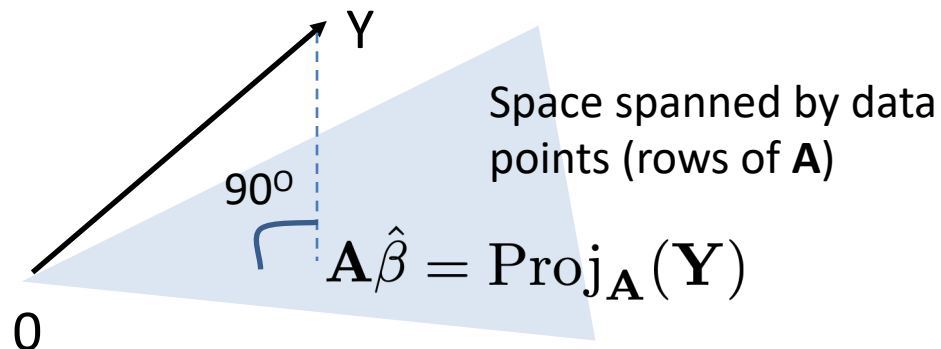
# Least Square solution satisfies Normal Equations

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

<span style="color:blue">p x p   p x1       p x1</span>

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

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

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

Y

90°

Space spanned by data points (rows of **A**)

$$\mathbf{A}\hat{\beta} = \mathrm{Proj}_{\mathbf{A}}(\mathbf{Y})$$

0

# Least Square solution satisfies Normal Equations

$$\frac{\partial J(\beta)}{\partial \beta}\bigg|_{\widehat{\beta}} = 0 \qquad \text{gives} \qquad (\mathbf{A}^T\mathbf{A})\widehat{\beta} = \mathbf{A}^T\mathbf{Y}$$

<span style="color:blue">p x p   p x1     p x1</span>

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

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

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

2) If dimension p is large, computing inverse is expensive O(p³)
   Gradient descent since objective is convex ($\mathbf{A}^\mathsf{T}\mathbf{A} \succeq 0$)

$$\boxed{\begin{aligned} \beta^{t+1} &= \beta^t - \frac{\alpha}{2}\frac{\partial J(\beta)}{\partial \beta}\bigg|_t \\ &= \beta^t - \alpha\,\mathbf{A}^T(\mathbf{A}\beta^t - Y) \end{aligned}}$$

27

# Least Square solution satisfies Normal Equations

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

p x p   p x1      p x1

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

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

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

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

# Least Square solution satisfies Normal Equations

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

p x p   p x1          p x1

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

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

If $\mathbf{A} = \mathbf{U}\mathbf{S}\mathbf{V}^\top$, then normal equations $(\mathbf{S}\mathbf{V}^\top)\hat{\beta} = (\mathbf{U}^\top\mathbf{Y})$

S - r x r                                          r x p   p x 1        r x 1

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

29

# 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)

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

<span style="color:red">Ridge Regression (l2 penalty)</span>

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

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

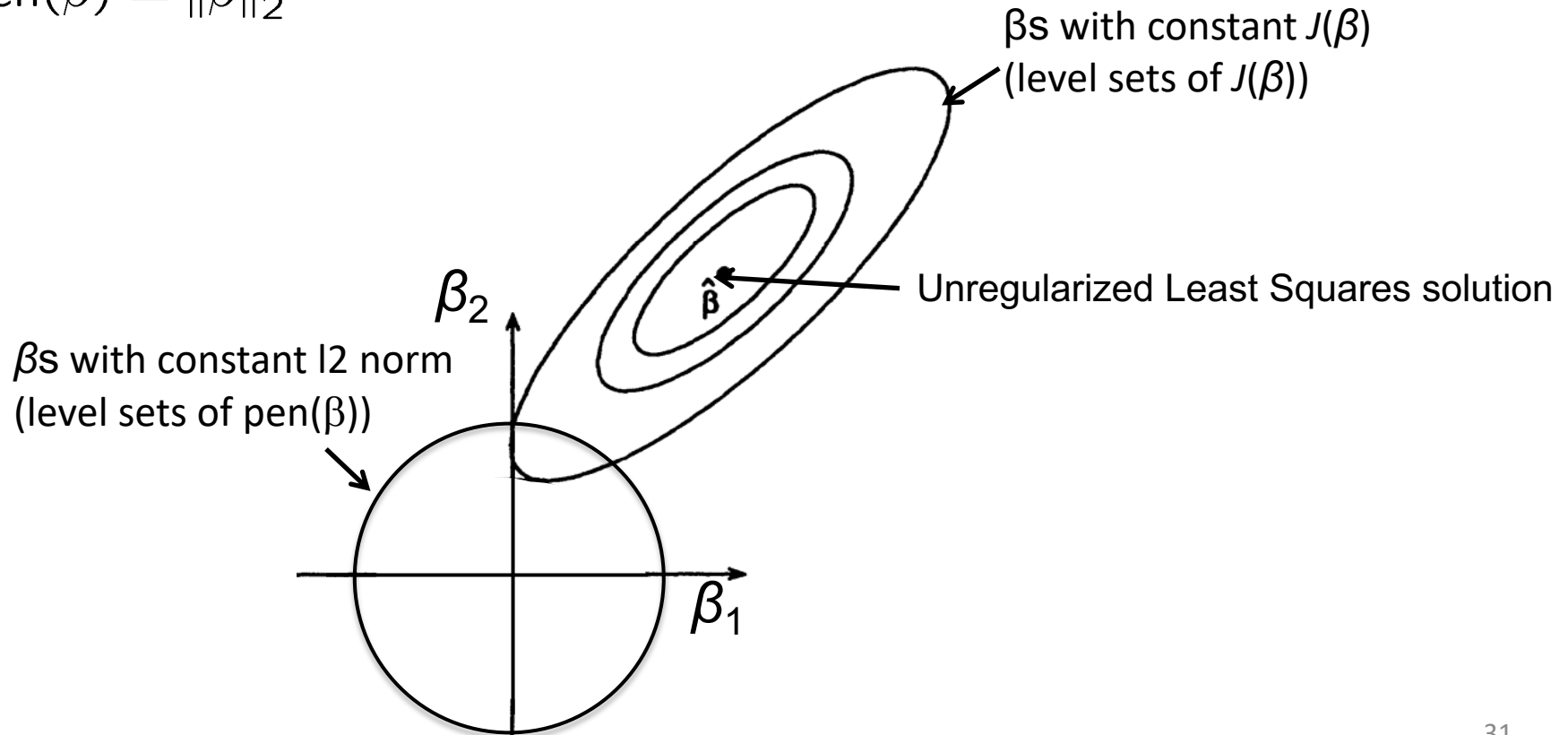Is $(\mathbf{A}^\top\mathbf{A} + \lambda\mathbf{I})$ invertible ?

# 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:
$$\mathrm{pen}(\beta) = \|\beta\|_2^2$$

βs with constant $J(\beta)$
(level sets of $J(\beta)$)

Unregularized Least Squares solution

$\hat{\beta}$

$\beta_2$

$\beta_1$

*β*s with constant l2 norm
(level sets of pen(β))

31

# 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)

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

<span style="color:red">Ridge Regression (l2 penalty)</span>

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

<span style="color:red">Lasso (l1 penalty)</span>

$$\lambda \geq 0$$

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

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

<span style="color:red">Ridge Regression (l2 penalty)</span>

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

<span style="color:red">Lasso (l1 penalty)</span>

$$\lambda \geq 0$$

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

# Ridge Regression vs Lasso

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

Ridge Regression:
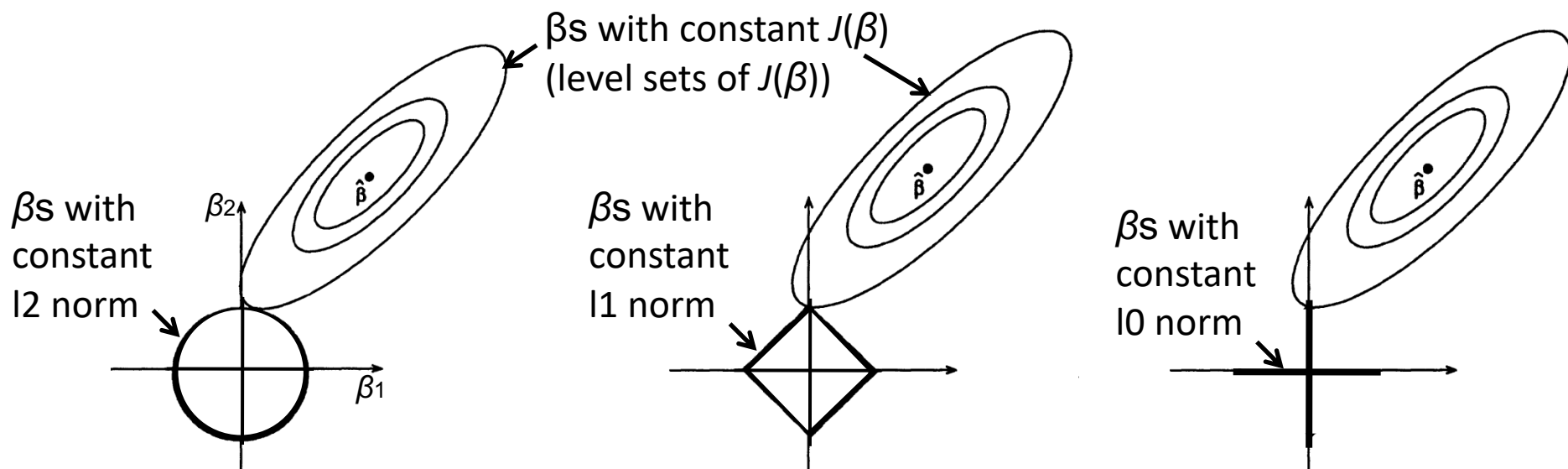$$\mathrm{pen}(\beta) = \|\beta\|_2^2$$

Lasso:
$$\mathrm{pen}(\beta) = \|\beta\|_1$$

Ideally l0 penalty, but optimization becomes non-convex

βs with constant $J(\beta)$
(level sets of $J(\beta)$)

βs with constant l2 norm

βs with constant l1 norm

βs with constant l0 norm

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!

34

# Matlab example

```matlab
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;
```

```matlab
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)
```
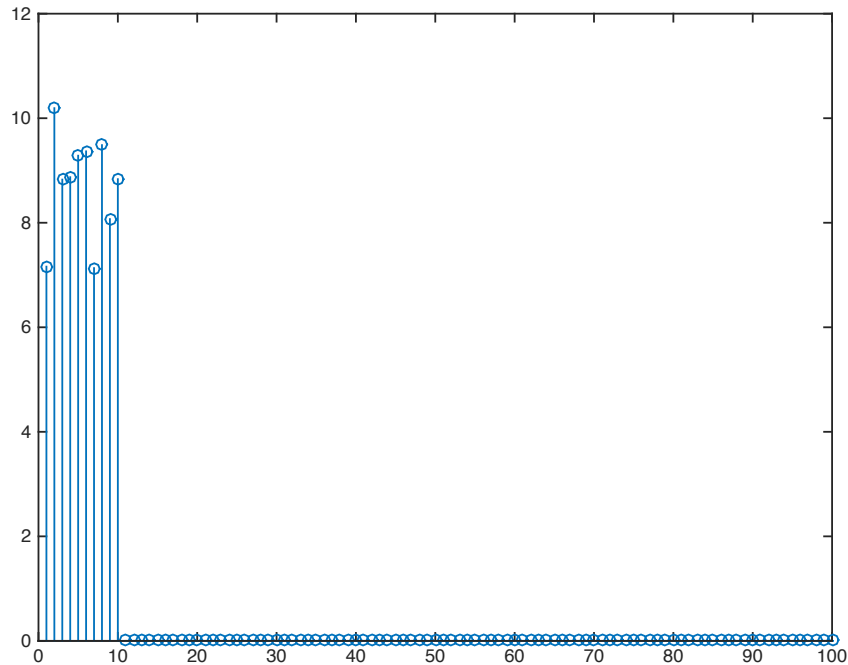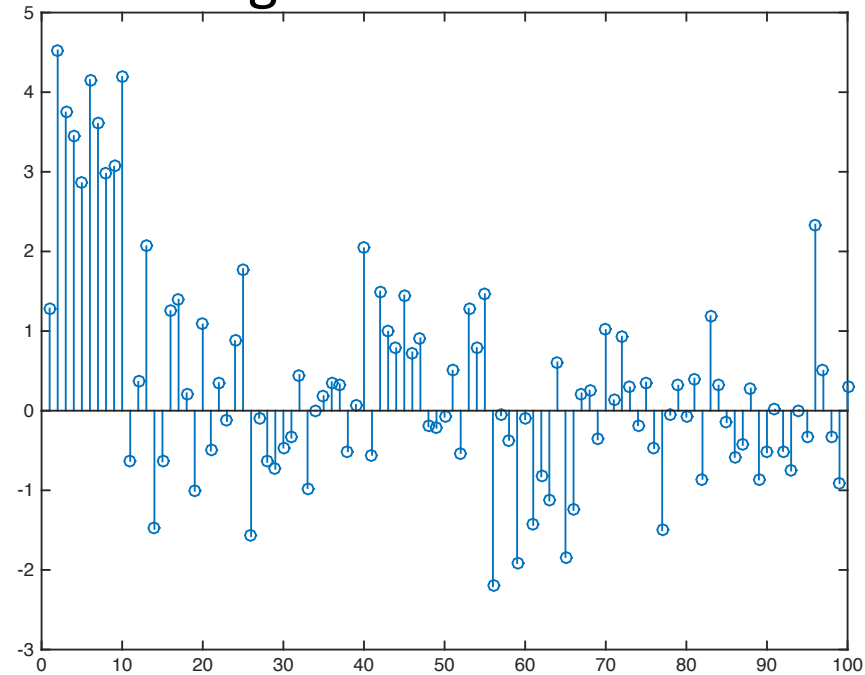
# Matlab example



Test MSE = 33.7997

Test MSE = 185.9948

Lasso Coefficients

Ridge Coefficients

# Regularized Least Squares – connection to MLE and MAP (Model-based approaches)

# Least Squares and M(C)LE
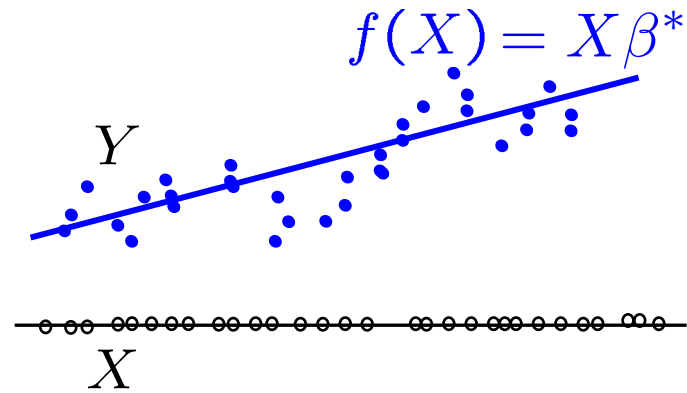
Intuition: Signal plus (zero-mean) Noise model

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

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

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

Conditional log likelihood

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

$f(X) = X\beta^*$

$Y$

$X$

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

# Regularized Least Squares and M(C)AP

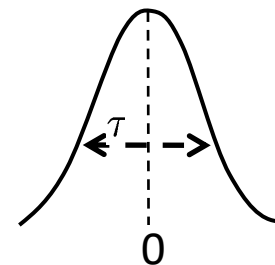What if $(\mathbf{A}^T \mathbf{A})$ is not invertible ?

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

I) Gaussian Prior

$$\beta \sim \mathcal{N}(0, \tau^2 \mathbf{I}) \qquad p(\beta) \propto e^{-\beta^T \beta / 2\tau^2}$$



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

$$\downarrow$$
$$\text{constant}(\sigma^2, \tau^2)$$

**Ridge Regression**

$$\widehat{\beta}_{\mathrm{MAP}} = (\boldsymbol{A}^\top \boldsymbol{A} + \lambda \boldsymbol{I})^{-1} \boldsymbol{A}^\top \boldsymbol{Y}$$
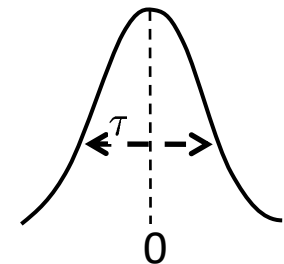
# Regularized Least Squares and M(C)AP

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

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

I) Gaussian Prior

$$\beta \sim \mathcal{N}(0,\tau^2\mathbf{I}) \qquad p(\beta) \propto e^{-\beta^T\beta/2\tau^2}$$



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

$$\downarrow$$

$$\text{constant}(\sigma^2,\tau^2)$$

**Ridge Regression**

Prior belief that β is Gaussian with zero-mean biases solution to "small" β
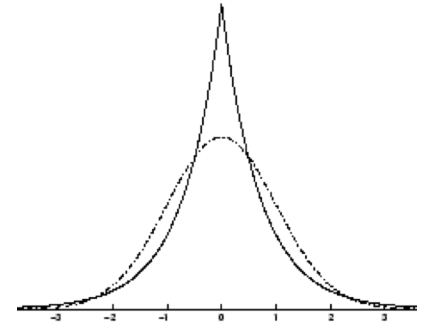
# Regularized Least Squares and M(C)AP

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

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

II) Laplace Prior

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



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

$$\downarrow$$

$$\mathsf{constant}(\sigma^2, t)$$

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