# Regularized Linear Regression
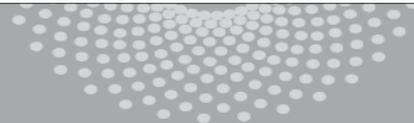
Aarti Singh

Machine Learning 10-701

Mar 22, 2023

# Mean square error regression

Optimal predictor:
$$f^* = \arg\min_f \mathbb{E}[(f(X) - Y)^2]$$

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

**Class of predictors**

$\mathcal{F}$ - Class of Linear functions
- Class of Polynomial functions
- Class of nonlinear functions

# Least Square solution satisfies Normal Equations

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

<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:

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

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

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

# Linear regression 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 $(\mathbf{A}^T\mathbf{A})$ ?

# Linear regression 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

S - r x r

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

r x p   p x 1      r x 1

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

# 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)

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

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

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

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

# Ridge Regression

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

Ridge Regression
(l2 penalty)

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

$$\lambda \geq 0$$

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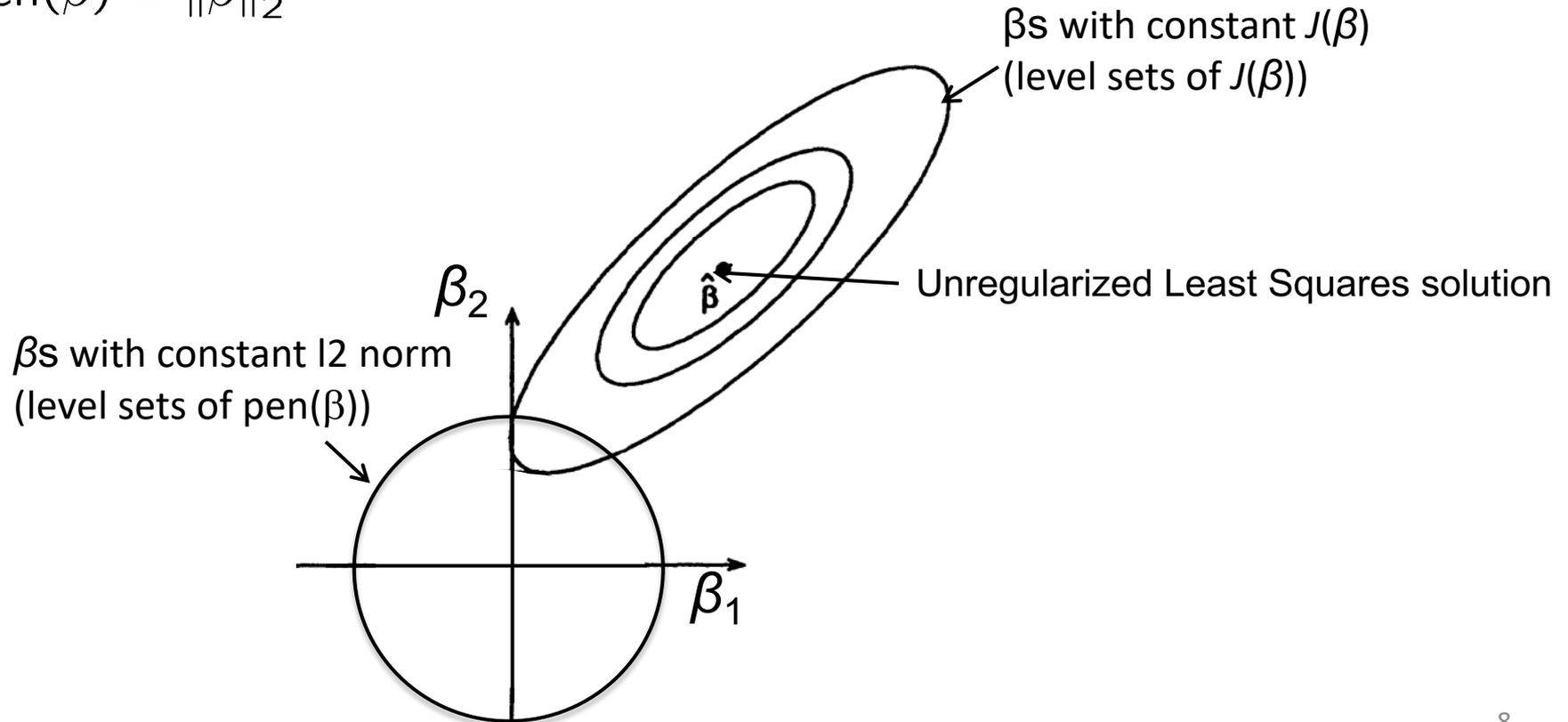
Is $(\mathbf{A}^\top \mathbf{A} + \lambda \mathbf{I})$ invertible ?

# 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

$\beta_2$

$\hat{\boldsymbol{\beta}}$

βs with constant l2 norm
(level sets of pen(β))

$\beta_1$

# 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 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>

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

$$\lambda \geq 0$$

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

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

# Ridge Regression vs Lasso

$$\min_{\beta}(\mathbf{A}\beta - \mathbf{Y})^T(\mathbf{A}\beta - \mathbf{Y}) + \lambda\mathrm{pen}(\beta) = \min_{\beta} J(\beta) + \lambda\mathrm{pen}(\beta)$$
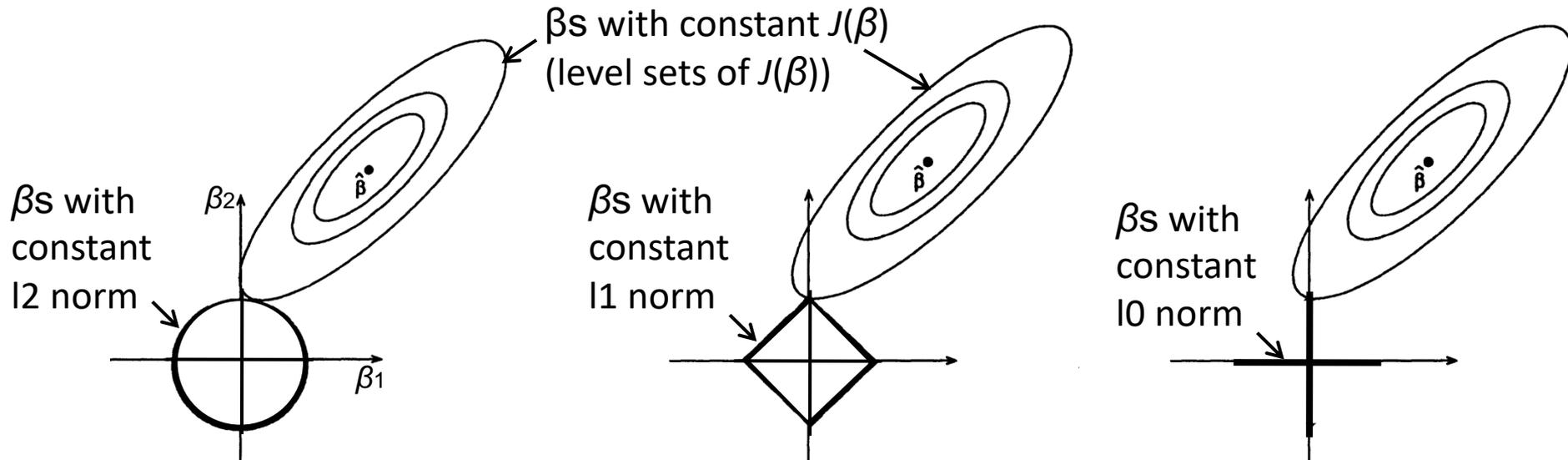
Ridge Regression:
$$\mathrm{pen}(\beta) = \|\beta\|_2^2$$

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!**

11

# 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;

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)
```
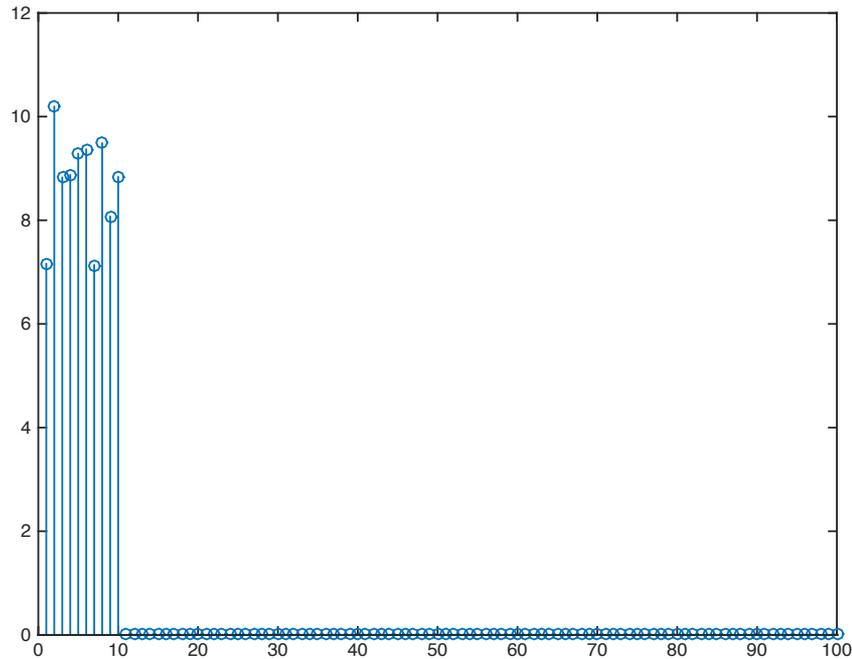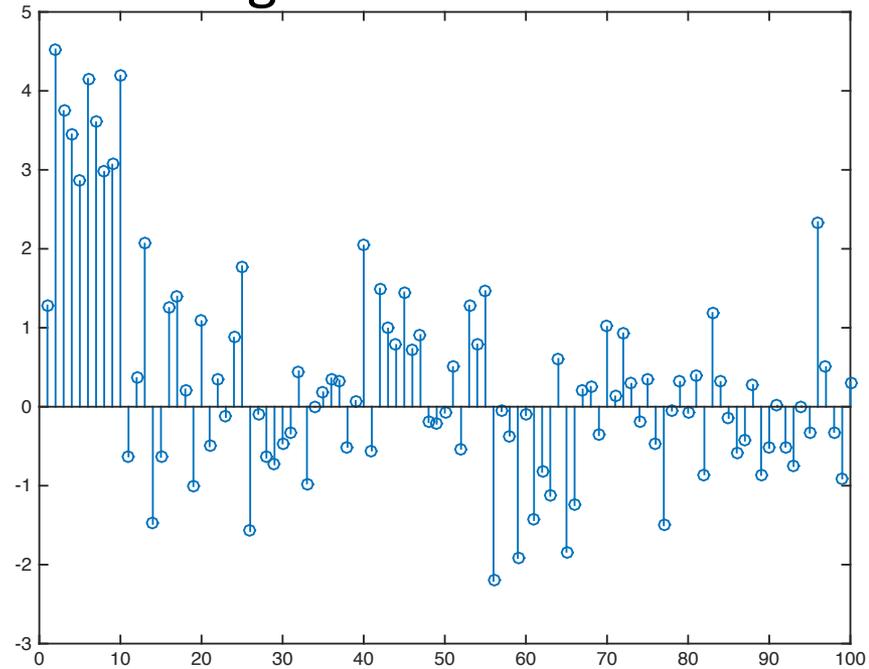
# Matlab example

Test MSE = 33.7997

Test MSE = 185.9948



Lasso Coefficients

Ridge Coefficients

# Least Squares and M(C)LE

Intuition: Signal plus (zero-mean) Noise model

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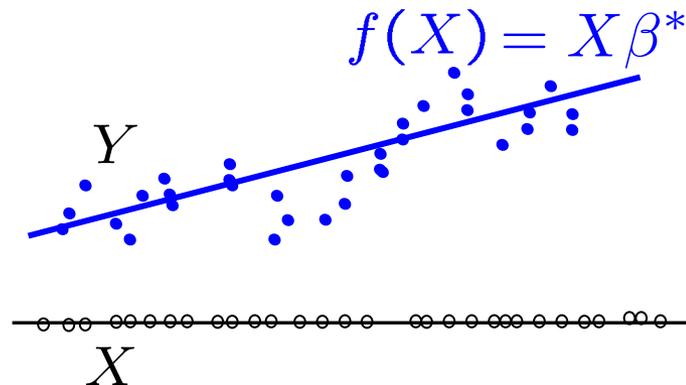
$f(X) = X\beta^*$

$Y$

$X$

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

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

Conditional log likelihood

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

**Least Square Estimate is same as Maximum Conditional Likelihood Estimate under a Gaussian model !**
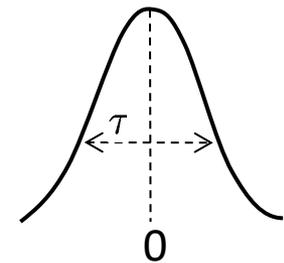
# Regularized Least Squares and M(C)AP

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**

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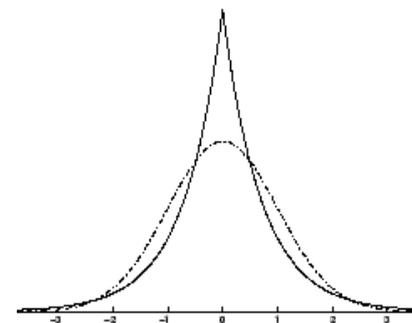
# Regularized Least Squares and M(C)AP

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

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

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

$$\downarrow$$

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

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

# Polynomial Regression

degree m

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

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

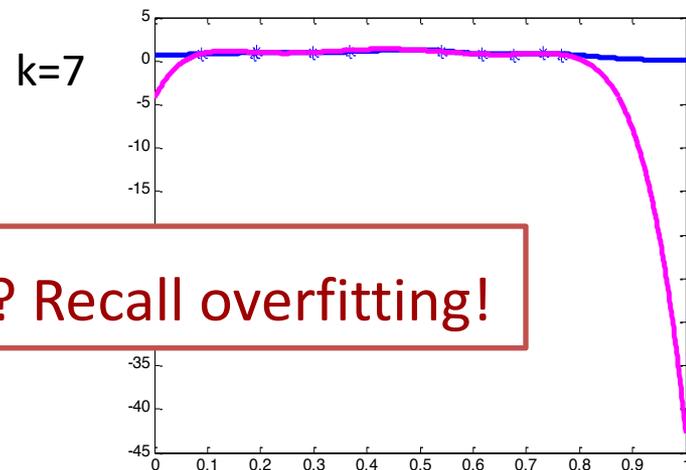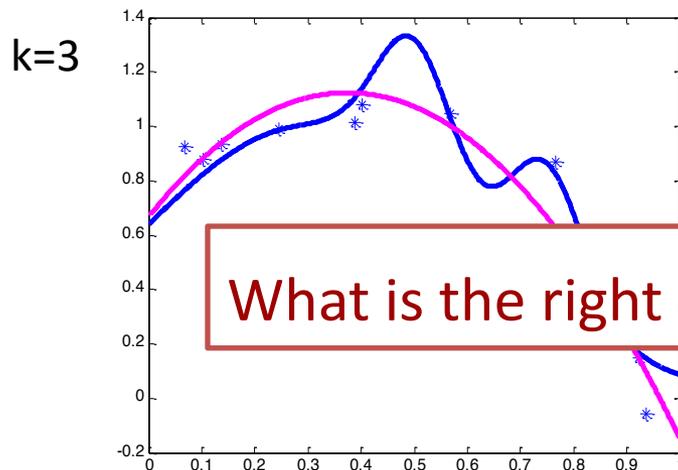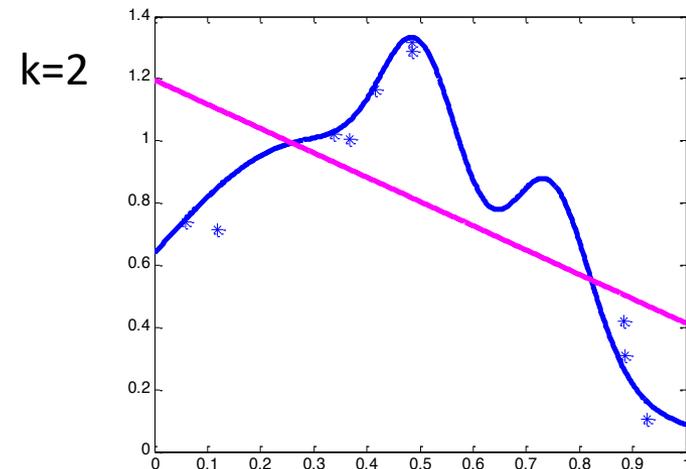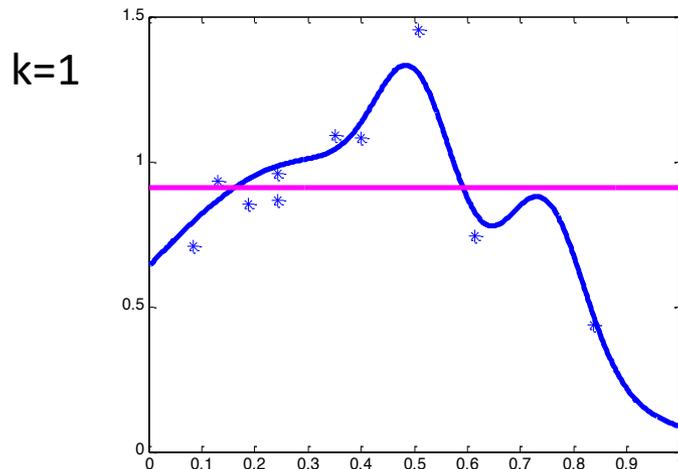$$\widehat{\beta} = (\mathbf{A}^T\mathbf{A})^{-1}\mathbf{A}^T\mathbf{Y} \qquad \widehat{f}_n(X) = \mathbf{X}\widehat{\beta}$$

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

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

# Polynomial Regression

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



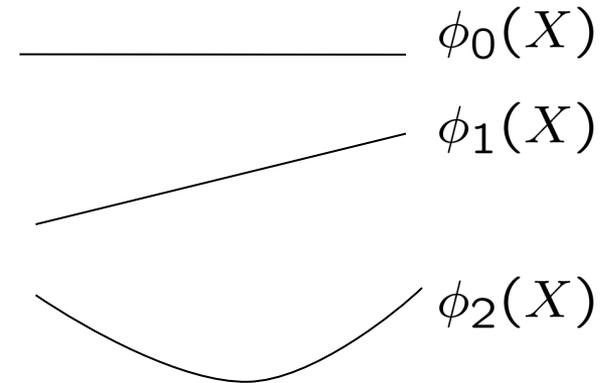What is the right order? Recall overfitting!

# 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

$\phi_0(X)$

$\phi_1(X)$

$\phi_2(X)$

In general, use any nonlinear features

e.g. $e^X$, log X, 1/X, sin(X), …

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

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

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

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

# Poll

- The maximum likelihood estimate of model parameter α for the random variable y ~N(α $x_1 x_2^3$ , $\sigma^2$), where $x_1$ and $x_2$ are random variables, can be learned using linear regression on n iid samples of ($x_1$, $x_2$, y)

  – True
  – False

# Can we kernelize linear regression?

# Linear (Ridge) regression

$$\min_{\beta} \sum_{i=1}^{n} (Y_i - X_i\beta)^2 + \lambda\|\beta\|_2^2 \qquad \widehat{\beta} = (\mathbf{A}^T\mathbf{A} + \lambda\mathbf{I})^{-1}\mathbf{A}^T\mathbf{Y}$$

Recall

$$\mathbf{A} = \begin{bmatrix} X_1 \\ \vdots \\ X_n \end{bmatrix} = \begin{bmatrix} X_1^{(1)} & \ldots & X_1^{(p)} \\ \vdots & \ddots & \vdots \\ X_n^{(1)} & \ldots & X_n^{(p)} \end{bmatrix}$$

Hence $\mathbf{A}^T\mathbf{A}$ is a p x p matrix whose entries denote the (sample) correlation between the features

NOT inner products between the data points – the inner product matrix would be $\mathbf{A}\mathbf{A}^T$ which is n x n (also known as Gram matrix)

Using dual formulation, we can write the solution in terms of $\mathbf{A}\mathbf{A}^T$

# Ridge regression

$$\min_{\beta} \sum_{i=1}^{n} (Y_i - X_i\beta)^2 + \lambda\|\beta\|_2^2 \qquad \widehat{\beta} = (\mathbf{A}^T\mathbf{A} + \lambda\mathbf{I})^{-1}\mathbf{A}^T\mathbf{Y}$$

<u>Similarity with SVMs</u>

Primal problem:

$$\min_{\beta, z_i} \sum_{i=1}^{n} z_i^2 + \lambda\|\beta\|_2^2$$
$$\text{s.t. } z_i = Y_i - X_i\beta$$

SVM Primal problem:

$$\min_{w, \xi_i} C \sum_{i=1}^{n} \xi_i + \frac{1}{2}\|w\|_2^2$$
$$\text{s.t. } \xi_i = \max(1 - Y_i X_i\, w, 0)$$

Lagrangian:

$$\sum_{i=1}^{n} z_i^2 + \lambda\|\beta\|^2 + \sum_{i=1}^{n} \alpha_i(z_i - Y_i + X_i\beta)$$

$\alpha_i$ – Lagrange parameter, one per training point

# Kernelized ridge regression

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

Using dual, can re-write solution as:

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

How does this help?
- Only need to invert n x n matrix (instead of p x p or m x m)
- More importantly, kernel trick!

$\mathbf{A}\mathbf{A}^T$ involves only inner products between the training points
BUT still have an extra $\mathbf{A}^T$

Recall the predicted label is $\widehat{f}_n(X) = \mathbf{X}\widehat{\beta}$

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

$\mathbf{X}\mathbf{A}^T$ contains inner products between test point $\mathbf{X}$ and training points!

# Kernelized ridge regression

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

Using dual, can re-write solution as:

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

How does this help?
- Only need to invert n x n matrix (instead of p x p or m x m)
- More importantly, kernel trick!

$$\widehat{f}_n(X) = \mathbf{K}_X (\mathbf{K} + \lambda \mathbf{I})^{-1} \mathbf{Y} \quad \text{where} \quad \begin{aligned} \mathbf{K}_X(i) &= \boldsymbol{\phi}(X) \cdot \boldsymbol{\phi}(X_i) \\ \mathbf{K}(i,j) &= \boldsymbol{\phi}(X_i) \cdot \boldsymbol{\phi}(X_j) \end{aligned}$$

Work with kernels, never need to write out the high-dim vectors

Ridge Regression with (implicit) nonlinear features $\boldsymbol{\phi}(X)$! $\quad f(X) = \phi(X)\beta$