# Linear Regression (Matrix-vector form)

$$X = \begin{bmatrix} 1 & X_1 & X_2 & \dots & X_p \end{bmatrix}$$

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

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

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

$$J(\beta)$$

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

$n \times p$

$$\mathbf{Y} = \begin{bmatrix} \mathbf{Y}_1 \\ \vdots \\ \mathbf{Y}_n \end{bmatrix}$$

$X$

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

$p \times p$

If invertible, closed form expression or gradient descent

1

# 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

$A_{n \times p}$

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

n x p       S - r x r

$n \times r$ $\quad r \begin{bmatrix} \lambda_1 & 0 \\ 0 & \lambda_r \end{bmatrix} r \begin{bmatrix} & & \end{bmatrix}$   $p$

$S$

$U^T U = I = V^T V$

$\cancel{V^T}(V S U^T U S V^T)\widehat{\beta} = \cancel{V^T} V S U^T Y$

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

r x p    p x 1      r x 1

$r \leq \min(n, p)$

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

$n < p$

high-dim setting

2

# Regularized Linear Regression

Aarti Singh & Geoff Gordon

Machine Learning 10-701
Feb 22, 2021

# 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)

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

Ridge Regression
(l2 penalty)

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

$$\lambda \geq 0$$

$$\beta^T(A^TA+\lambda I)\beta \qquad \beta^TA^TA\beta \qquad \beta^T\beta$$

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

$$(A^TA+\lambda I)\,v = \lambda' v$$

$$A^TAv + \lambda v = \lambda' v$$

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

4

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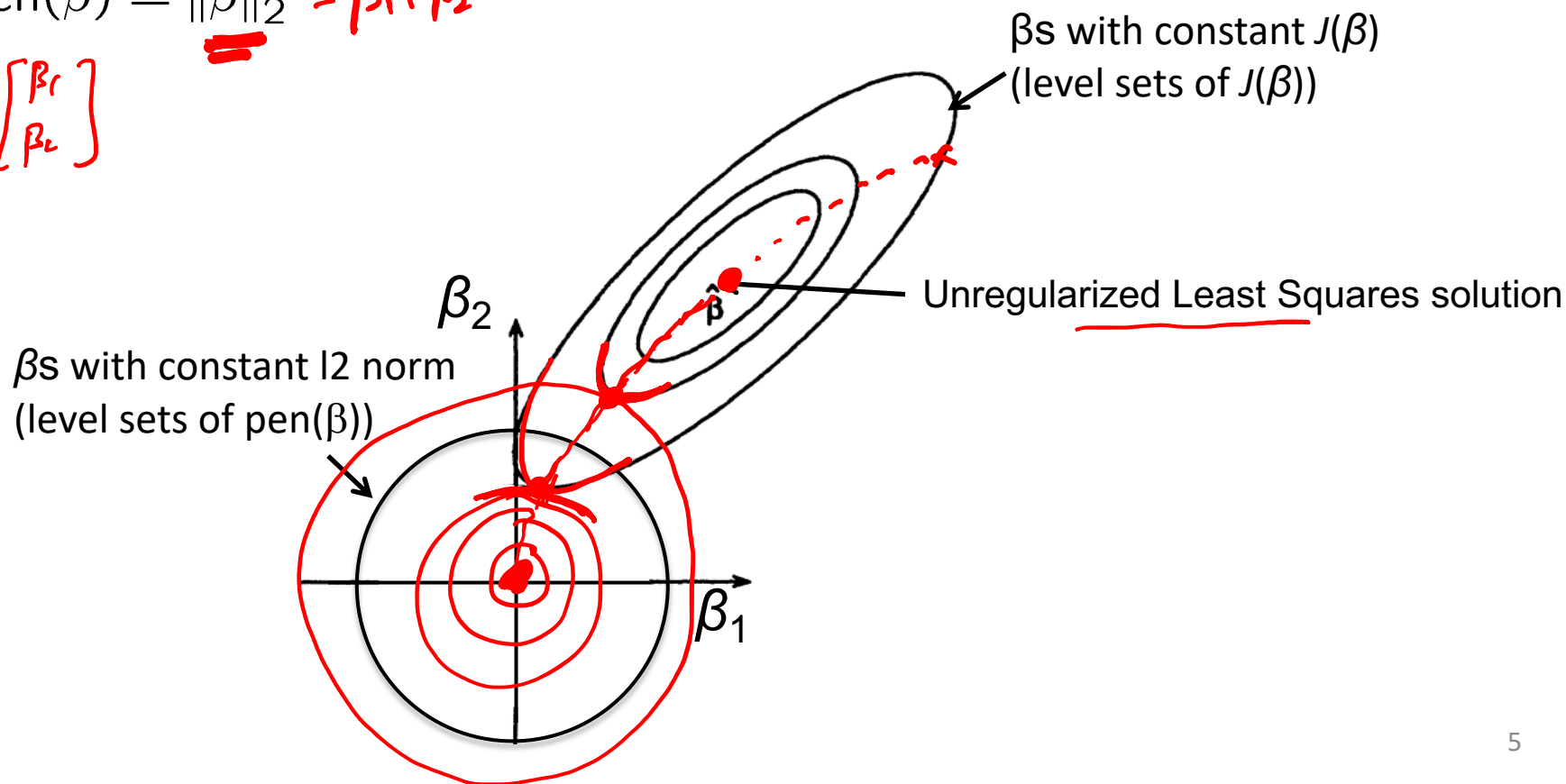
$$\min_{\beta}(\mathbf{A}\beta - \mathbf{Y})^T(\mathbf{A}\beta - \mathbf{Y}) + \lambda\text{pen}(\beta) = \min_{\beta} J(\beta) + \lambda\text{pen}(\beta)$$

$\lambda \geq 0$

$J^*$

Ridge Regression:

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

$\beta = \begin{bmatrix} \beta_1 \\ \beta_2 \end{bmatrix}$

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

Unregularized Least Squares solution

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

$\beta_2$

$\beta_1$

$\hat{\beta}$

# 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)

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

Ridge Regression
(l2 penalty)

$= \sum_i \beta_i^2$

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

Lasso
(l1 penalty)

$\lambda \geq 0$

$= \sum_i |\beta_i|$

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)

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

Ridge Regression
(l2 penalty)

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

Lasso
(l1 penalty)

$$\lambda \geq 0$$

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

# Ridge Regression vs Lasso

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

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

non-convex

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

Lasso:
$$\text{pen}(\beta) = \|\beta\|_1 = \sum_i |\beta_i| ✓$$

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

l0

$l_q < 1$

**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

```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)
```
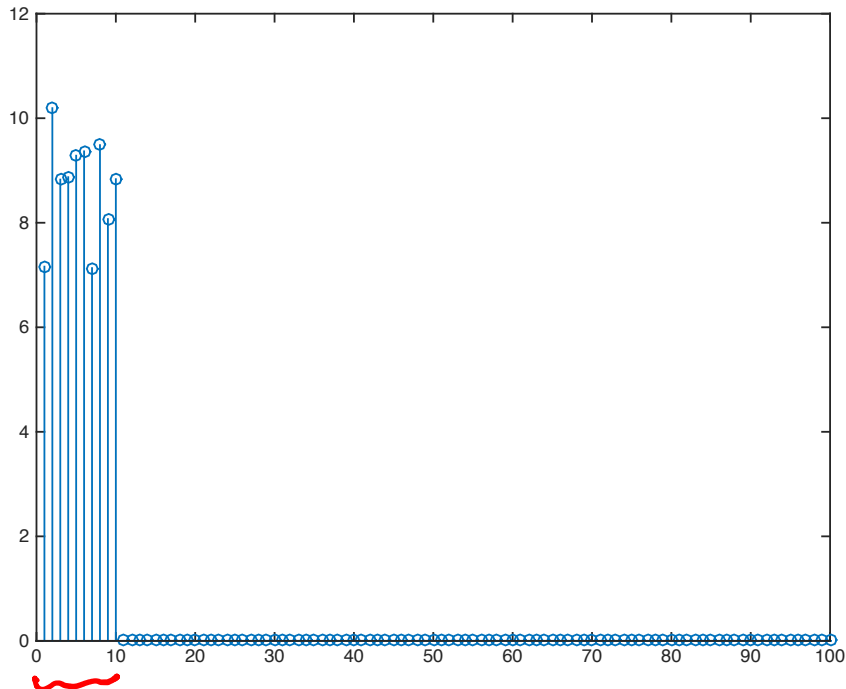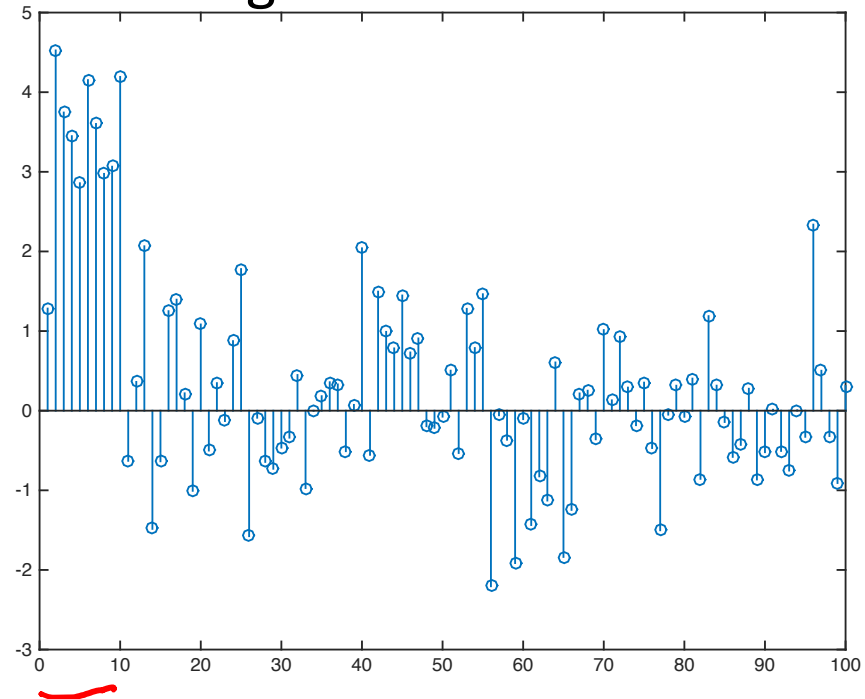
# Matlab example

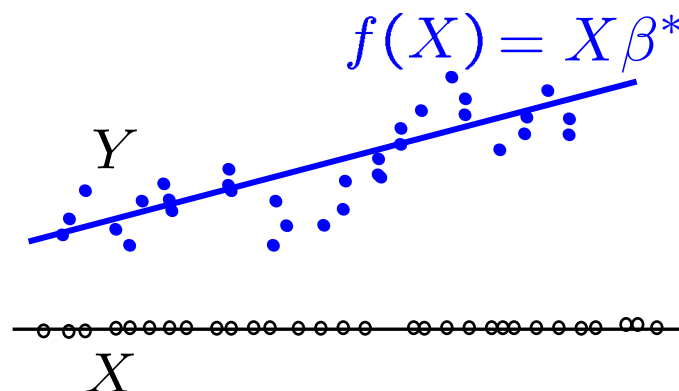Test MSE = 33.7997                    Test MSE = 185.9948


Lasso Coefficients


Ridge Coefficients

# Least Squares and M(C)LE

Intuition: Signal plus (zero-mean) Noise model

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

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

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

Conditional log likelihood

$$= \arg\max_{\beta} \log e^{-\frac{(Y_i - X_i\beta)^2}{2\sigma^2}}$$

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

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

$Y$

$X$

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

# Regularized Least Squares and M(C)AP

$$P(\theta|D) \propto P(\theta) P(D|\theta)$$

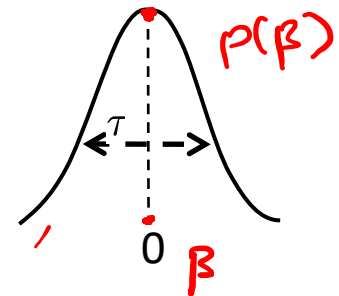What if $(\mathbf{A}^T\mathbf{A})$ is not invertible ?

$$\log p(D|\theta) \qquad + \log p(\theta)$$

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

I) Gaussian Prior

$$\Sigma$$

$$\beta^T\Sigma^{-1}\beta$$

$$p(\beta)$$

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

$$\|\beta\|^2$$

$$\log p(\beta) \propto -\|\beta\|^2$$

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

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

**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 ?

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

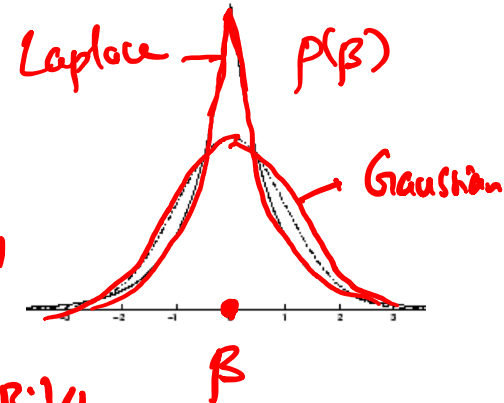$$-\sum_i |\beta_i| = -\|\beta\|_1$$

$$p(\beta) \propto e^{-|\beta|}$$

II) Laplace Prior

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

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

Laplace $\rightarrow$ $p(\beta)$

Gaussian

$$P(\beta) = \prod_i p(\beta_i) = \prod_i e^{-\frac{|\beta_i|}{t}}$$

$$= e^{-\sum_i |\beta_i|/t}$$

$\beta$

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

Lasso

constant$(\sigma^2,t)$

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

13

# Polynomial Regression

$$X = [1 \quad X \quad X^2 \quad X^3 \cdots X^m]$$

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

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

where $\mathbf{A} = \begin{bmatrix} 1 & X_1 & X_1^2 & \ldots & X_1^m \\ \vdots & & & \ddots & \vdots \\ 1 & X_n & X_n^2 & \ldots & X_n^m \end{bmatrix}$
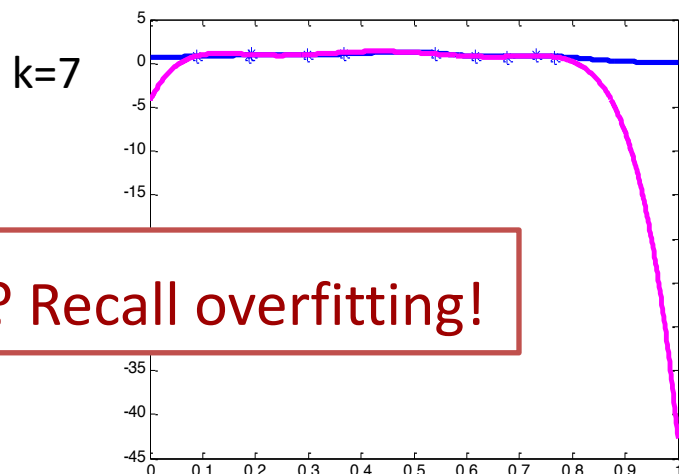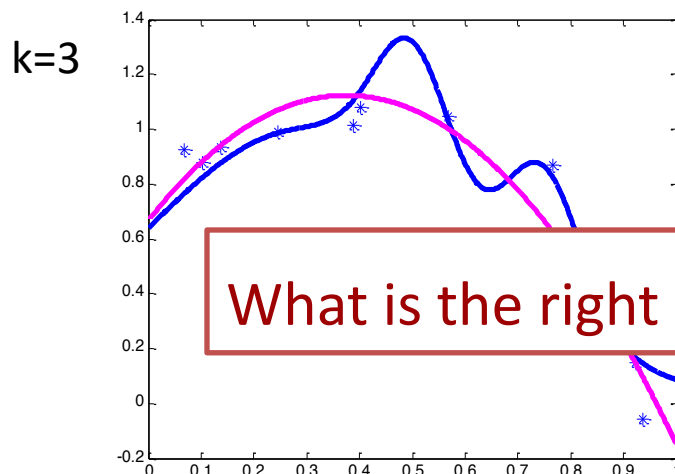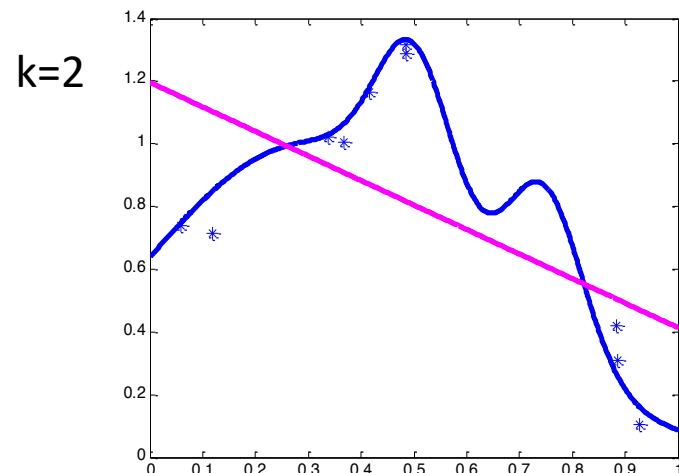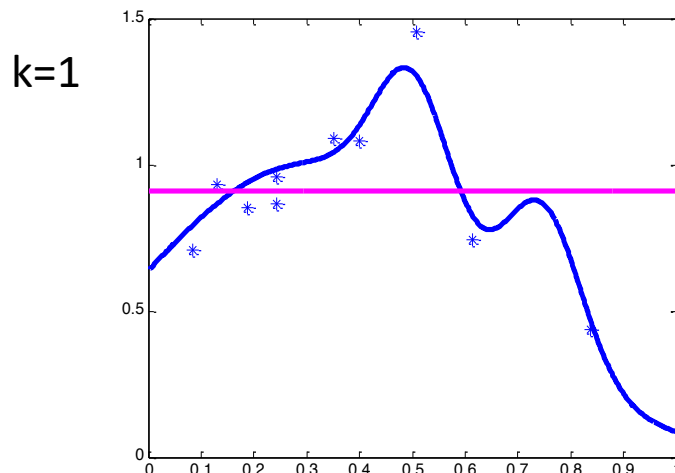
Multivariate (p-dim) case:
$$f(X) = \beta_0 + \beta_1 X^{(1)} + \beta_2 X^{(2)} + \cdots + \beta_p X^{(p)}$$
$$+ \sum_{i=1}^{p}\sum_{j=1}^{p} \beta_{ij} X^{(i)}X^{(j)} + \sum_{i=1}^{p}\sum_{j=1}^{p}\sum_{k=1}^{p} X^{(i)}X^{(j)}X^{(k)}$$
$$+ \ldots \text{terms up to degree m}$$

$p=2$

$X \rightarrow [1 \ X_1 \ X_2 \ X_1^2 \ X_2^2 \ X X_2]$

14

# Polynomial Regression

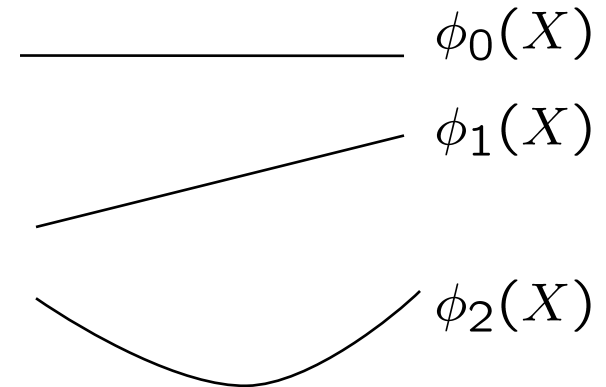Polynomial of order k, equivalently of degree up to k-1



What is the right order? Recall overfitting!

# Regression with nonlinear features

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

$\phi_0(X)$

$\phi_1(X)$

$\phi_2(X)$

Weight of each feature

Nonlinear features

In general, use any nonlinear features

e.g. $e^X$, log X, 1/X, sin(X), …

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

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

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

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