

# Linear Regression

Aarti Singh

Co-instructor: Barnabas Póczos

Machine Learning 10-401

Mar 24, 2016



**MACHINE LEARNING** DEPARTMENT



# Discrete to Continuous Labels

## Classification

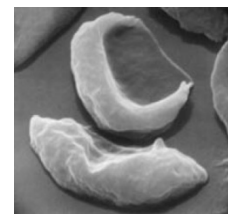


**X = Document**



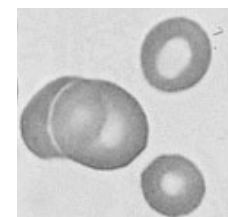
Sports  
Science  
News

**Y = Topic**



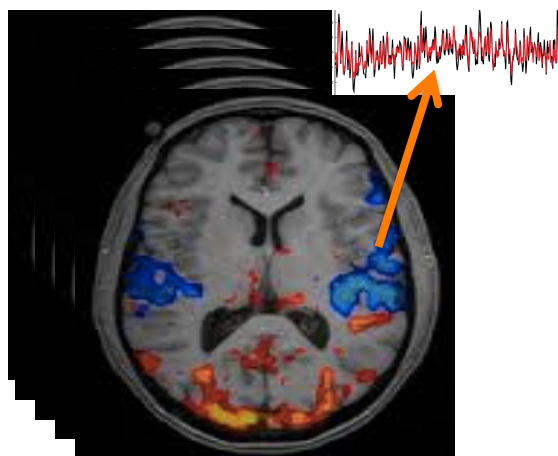
Anemic cell  
Healthy cell

**Y = Diagnosis**



**X = Cell Image**

## Regression



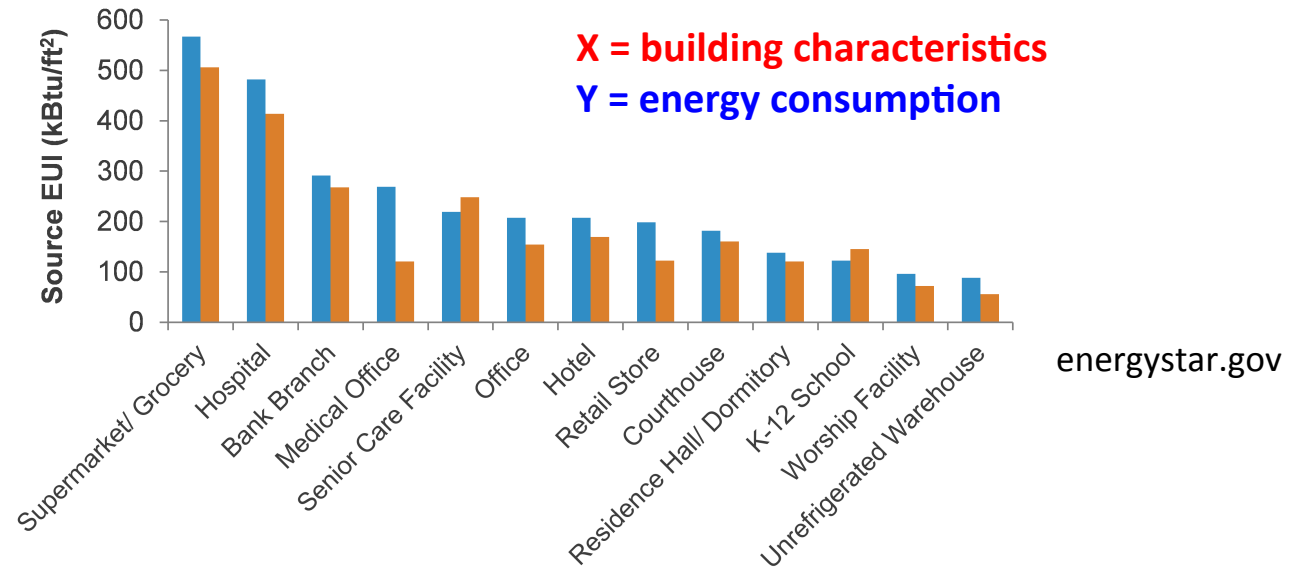
**X = Brain Scan**



**Y = Age of a subject**

# Regression Tasks

## Estimating Energy Usage



## Estimating Contamination



# Supervised Learning

**Goal:** Construct a **predictor**  $f : X \rightarrow Y$  to minimize loss function (performance measure)

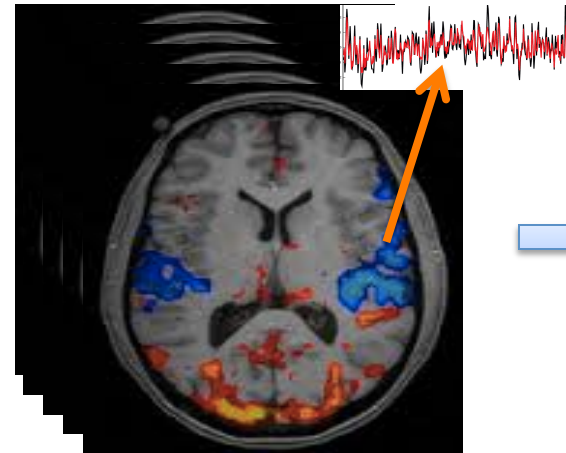


Sports  
Science  
News

**Classification:**

$$P(f(X) \neq Y)$$

**Probability of Error**



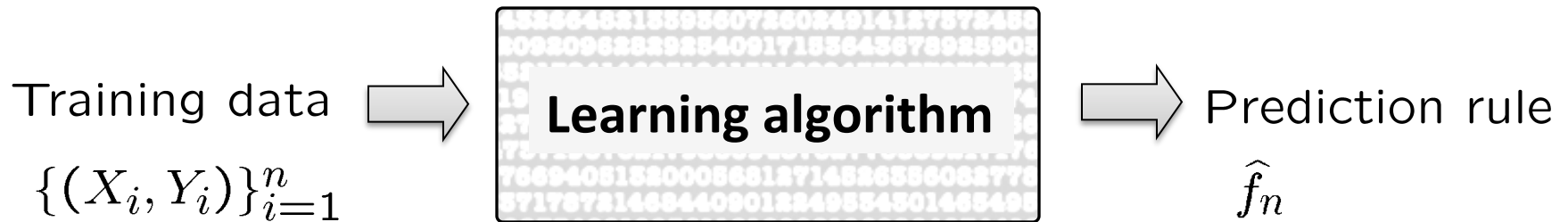
**Y = Age of  
a subject**

**Regression:**

$$\mathbb{E}[(f(X) - Y)^2]$$

**Mean Squared Error**

# Regression algorithms



Linear Regression

Regularized Linear Regression – Ridge regression, Lasso

Polynomial Regression

Kernelized Ridge Regression

Gaussian Process Regression

Kernel regression, Regression Trees, Splines, Wavelet estimators, ...

# Replace Expectation with Empirical Mean

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

Empirical Minimizer:  $\hat{f}_n = \arg \min_{f \in \mathcal{F}} \underbrace{\left( \frac{1}{n} \sum_{i=1}^n (f(X_i) - Y_i)^2 \right)}_{\text{Empirical mean}}$

Law of Large Numbers:

$$\frac{1}{n} \sum_{i=1}^n [\text{loss}(Y_i, f(X_i))] \xrightarrow{n \rightarrow \infty} \mathbb{E}_{XY} [\text{loss}(Y, f(X))]$$

# Restrict class of predictors

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

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

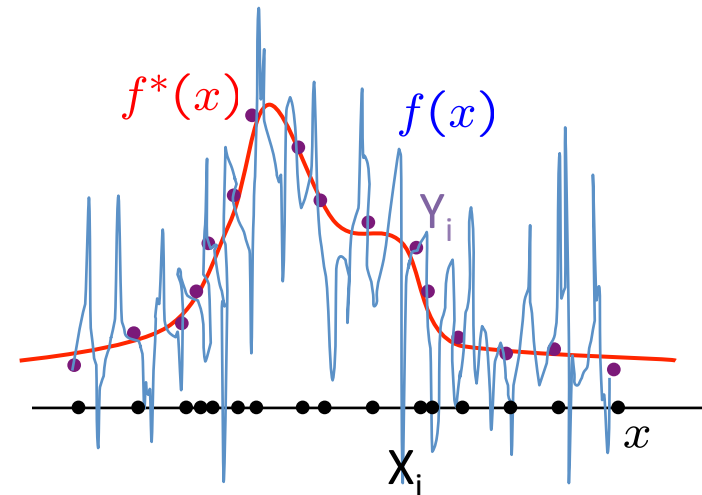
Class of predictors

Why?

Overfitting!

Empirical loss minimized by any function of the form

$$f(x) = \begin{cases} Y_i, & x = X_i \text{ for } i = 1, \dots, n \\ \text{any value,} & \text{otherwise} \end{cases}$$



# Restrict class of predictors

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

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

Class of predictors

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



# Linear Regression

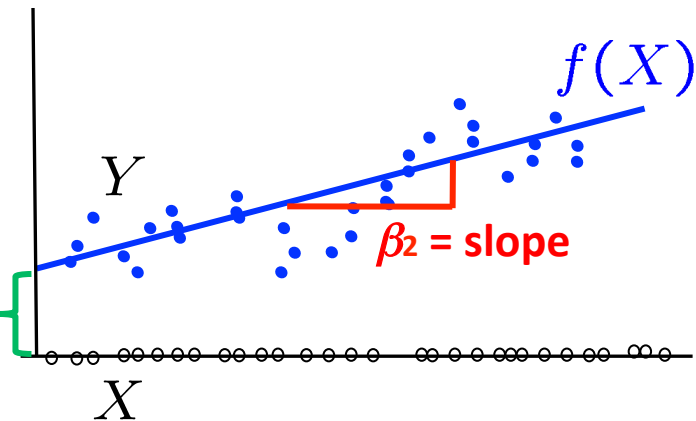
$$\hat{f}_n^L = \arg \min_{f \in \mathcal{F}_L} \frac{1}{n} \sum_{i=1}^n (f(X_i) - Y_i)^2 \quad \text{Least Squares Estimator}$$

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

Uni-variate case:

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

$\beta_1$  - intercept



Multi-variate case:

$$f(X) = f(X^{(1)}, \dots, X^{(p)}) = \beta_1 X^{(1)} + \beta_2 X^{(2)} + \dots + \beta_p X^{(p)}$$

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

# Least Squares Estimator

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



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

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

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

# Least Squares Estimator

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

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

$$\left. \frac{\partial J(\beta)}{\partial \beta} \right|_{\hat{\beta}} = 0$$

# Least Square 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,

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

# Linear Regression

Aarti Singh

Co-instructor: Barnabas Poczos

Machine Learning 10-401

Mar 29, 2016



**MACHINE LEARNING** DEPARTMENT



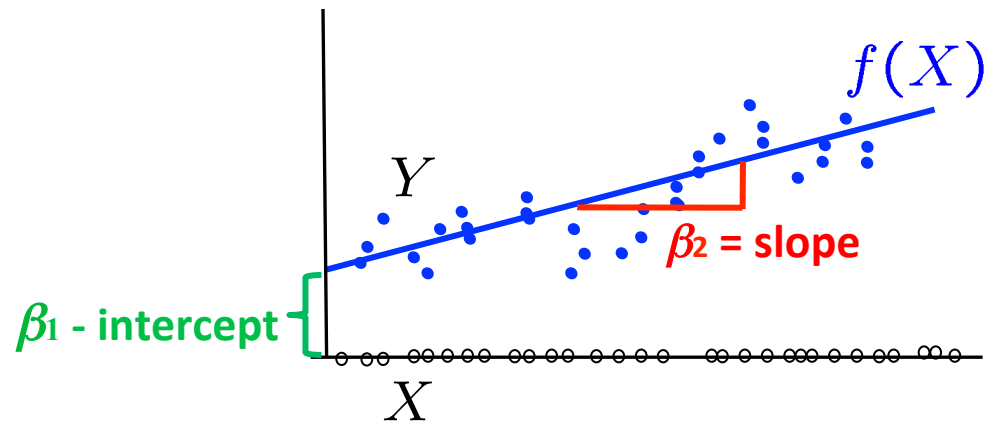
# Linear Regression

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

Least Squares Estimator

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

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



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

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

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

# Least Square 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,

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

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

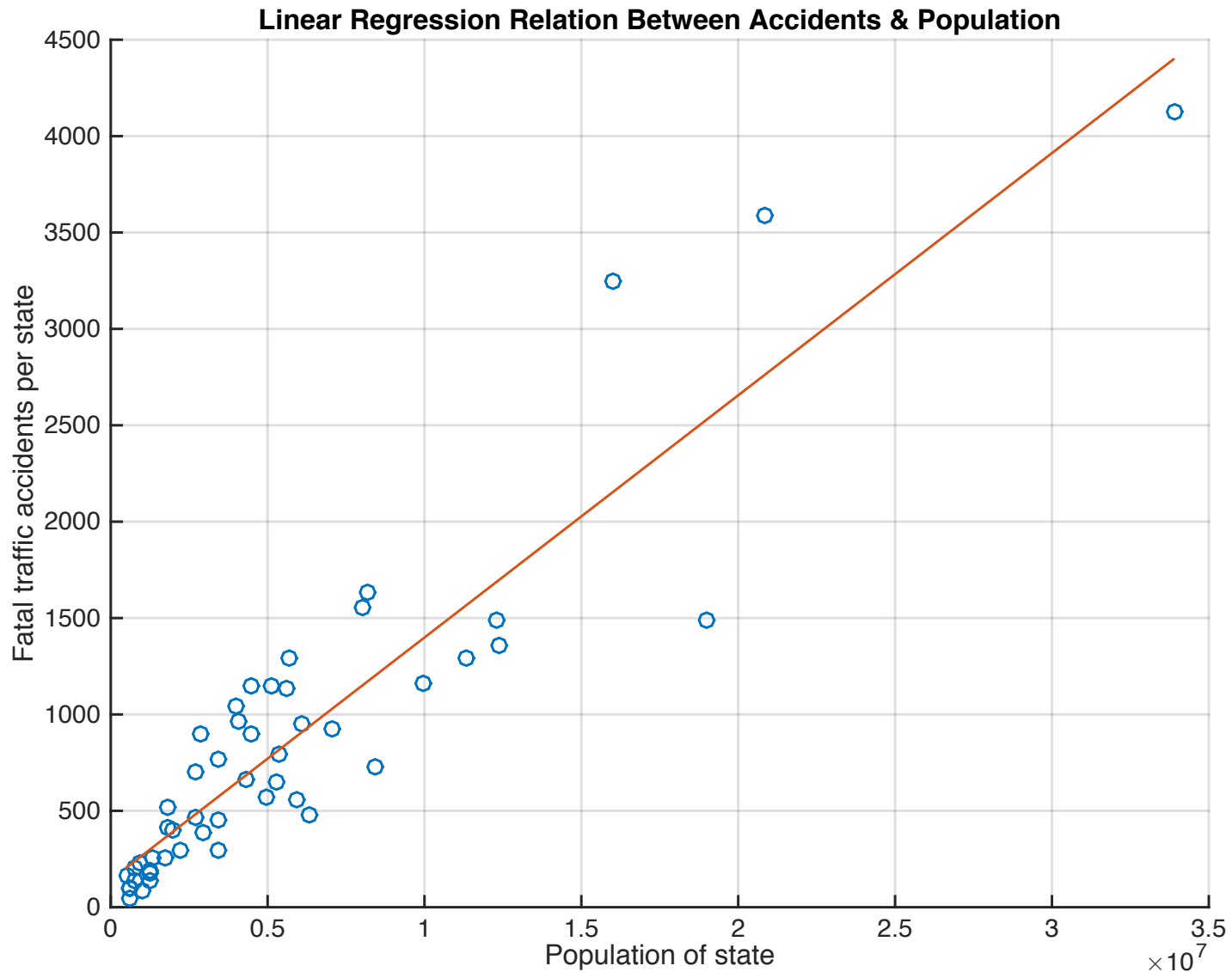
# Matlab example – linear regression

```
load accidents
x = hwydata(:,14);           %Population of states
y = hwydata(:,4);           %Accidents per state
scatter(x,y)
hold on
X = [ones(length(x),1) x];

b = X\y;
yhat = X*b;
plot(x,yhat)
xlabel('Population of state')
ylabel('Fatal traffic accidents per state')
title('Linear Regression Relation Between Accidents &
Population')
```



# Matlab example – linear regression



# Least Square 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,

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

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)

# Least Square 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,

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

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

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

# Least Square 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,

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

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

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

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

Constrain solution i.e. Regularization (later)

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

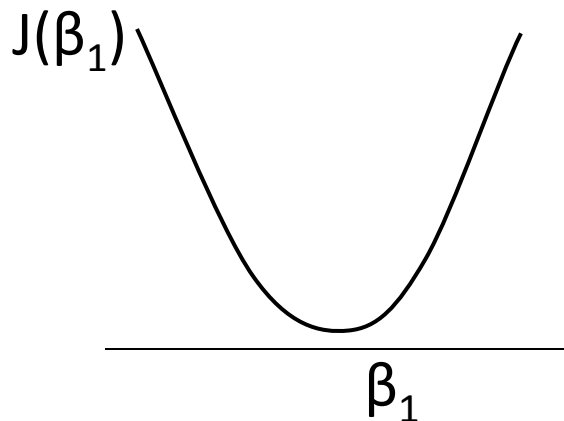
# Gradient Descent

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

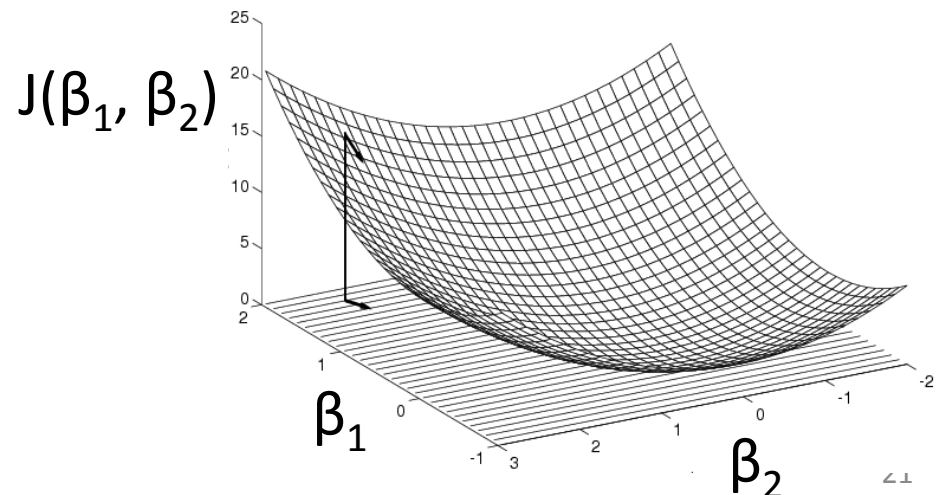
$$\hat{\beta} = \arg \min_{\beta} \frac{1}{n} (\mathbf{A}\beta - \mathbf{Y})^T (\mathbf{A}\beta - \mathbf{Y}) = \arg \min_{\beta} J(\beta)$$

Treat as optimization problem

Observation:  $J(\beta)$  is convex in  $\beta$ .



**How to find the minimizer?**



# Gradient Descent

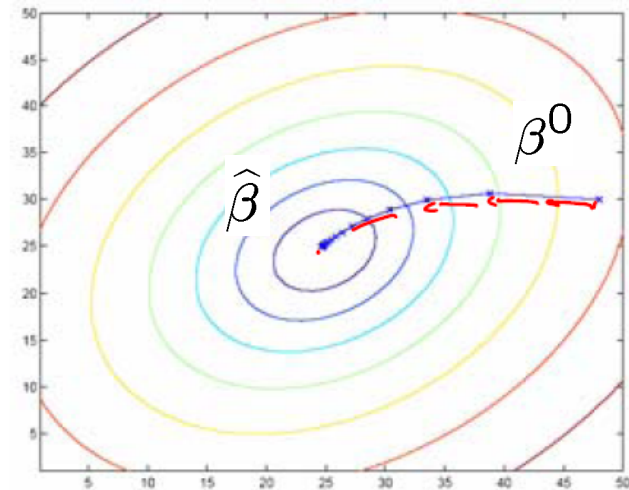
Even when  $(\mathbf{A}^T \mathbf{A})$  is invertible, might be computationally expensive if  $\mathbf{A}$  is huge.

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

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

Initialize:  $\beta^0$

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



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

# 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)

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

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

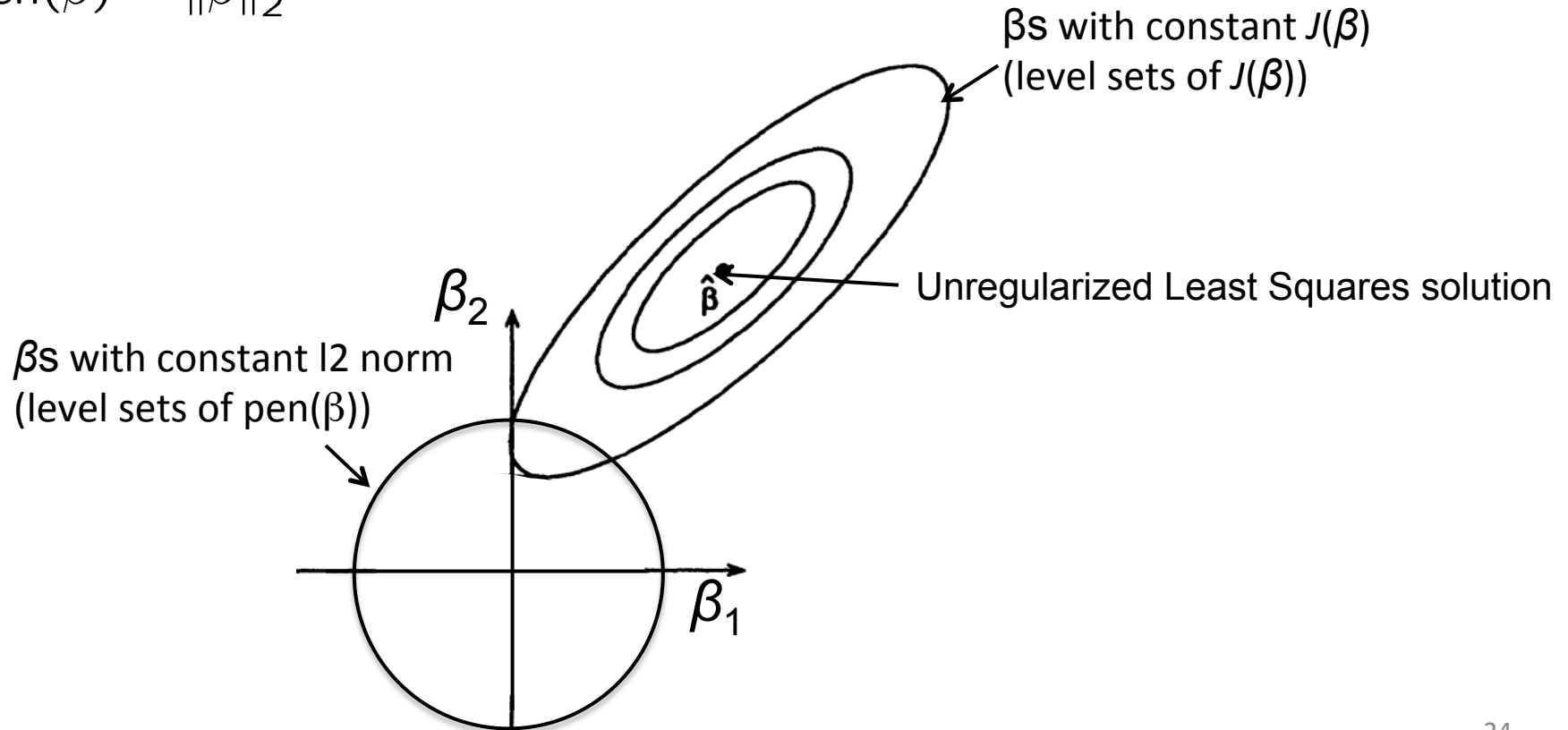
Is  $(\mathbf{A}^T \mathbf{A} + \lambda \mathbf{I})$  invertible ?

# 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  $\beta$  (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)

$$\lambda \geq 0$$

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:

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

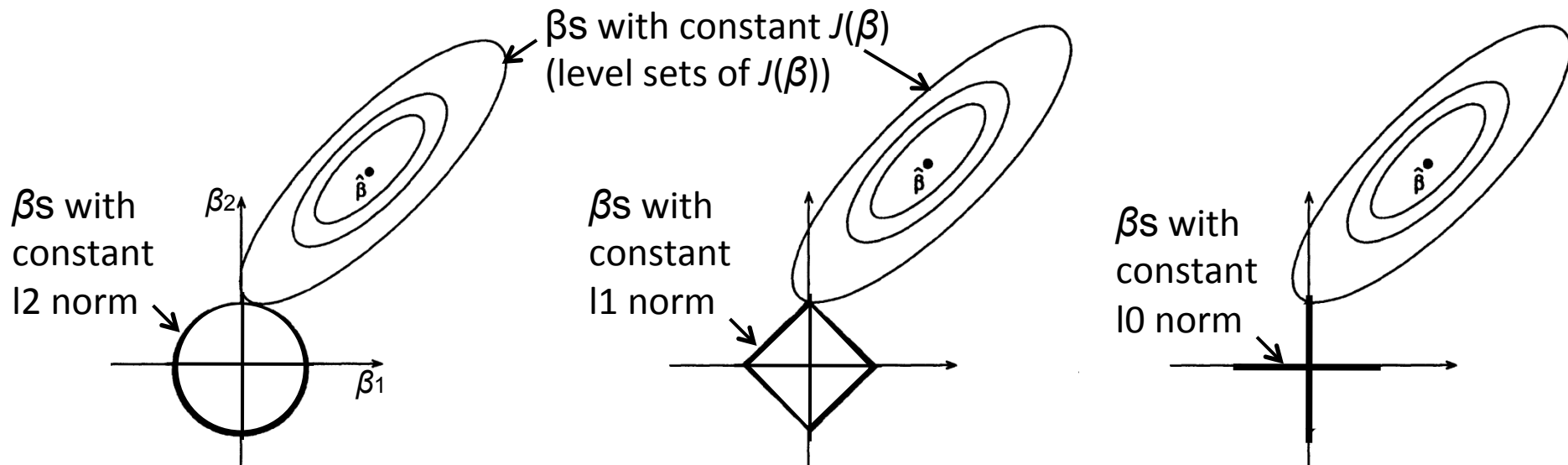
Lasso:

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

Ideally l0 penalty,

but optimization

becomes non-convex



**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)
```

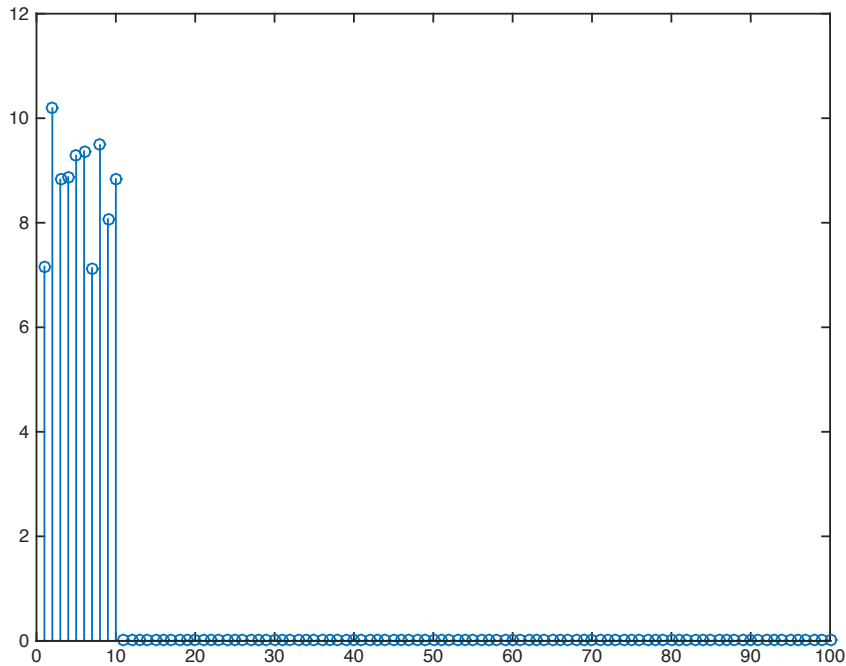
```
ridgeWeights = lasso(X,Y,'Lambda',1,  
    'Alpha', 0.0001);  
Yridge = Xtest*ridgeWeights;  
norm(Ytest-Yridge)
```

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

# Matlab example

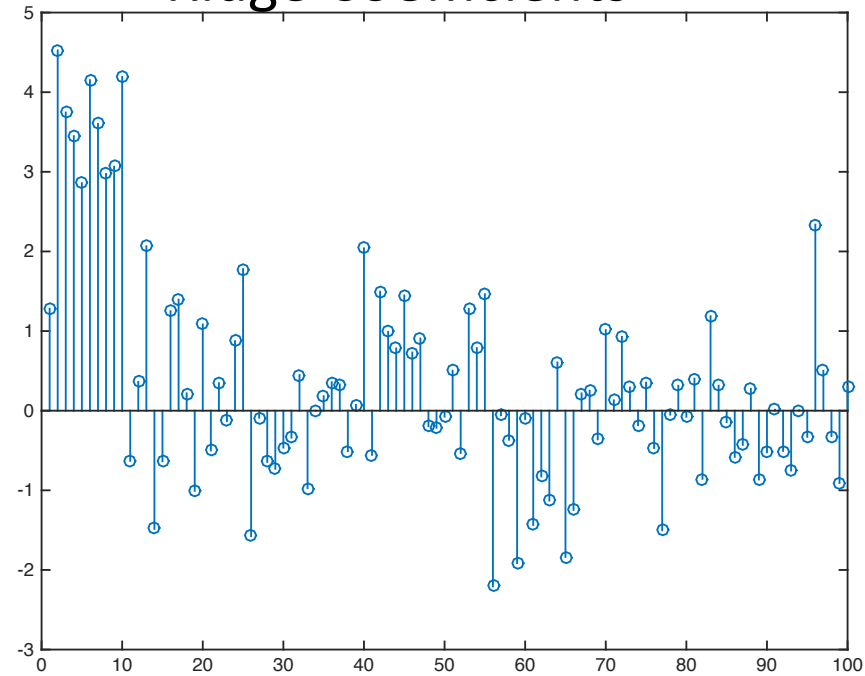
Test MSE = 33.7997

Lasso Coefficients



Test MSE = 185.9948

Ridge Coefficients



# Regularized Least Squares Regression

Aarti Singh

Co-instructor: Barnabas Poczos

Machine Learning 10-401

Mar 31, 2016



**MACHINE LEARNING** DEPARTMENT



# 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)

$$\lambda \geq 0$$

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

# **Regularized Least Squares – connection to MLE and MAP**



# Least Squares and M(C)LE

Intuition: Signal plus (zero-mean) Noise model

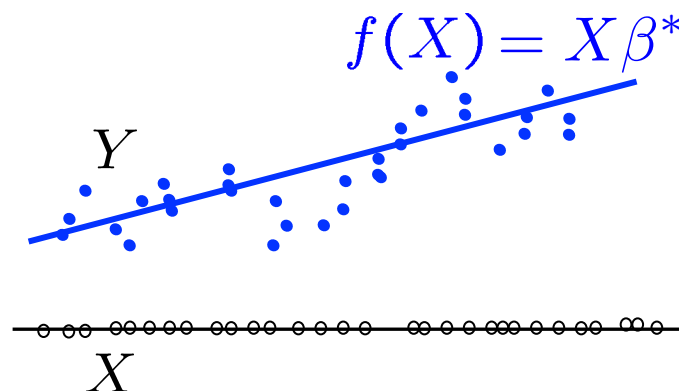
$$Y = f^*(X) + \epsilon = X\beta^* + \epsilon$$

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

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

Conditional log likelihood

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



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

# Regularized Least Squares and M(C)AP

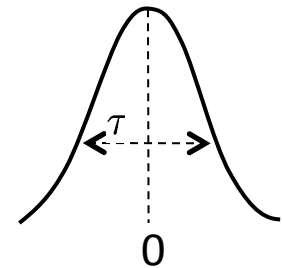
What if  $(\mathbf{A}^T \mathbf{A})$  is not invertible ?

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

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( $\sigma^2, \tau^2$ )

**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

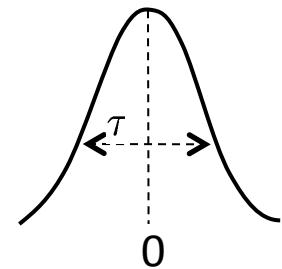
What if  $(\mathbf{A}^T \mathbf{A})$  is not invertible ?

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

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( $\sigma^2, \tau^2$ )

**Ridge Regression**

Prior belief that  $\beta$  is Gaussian with zero-mean biases solution to “small”  $\beta$

# Regularized Least Squares and M(C)AP

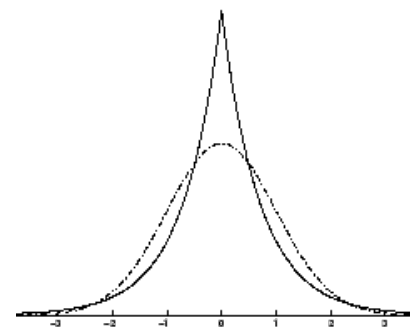
What if  $(\mathbf{A}^T \mathbf{A})$  is not invertible ?

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

II) Laplace Prior

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

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



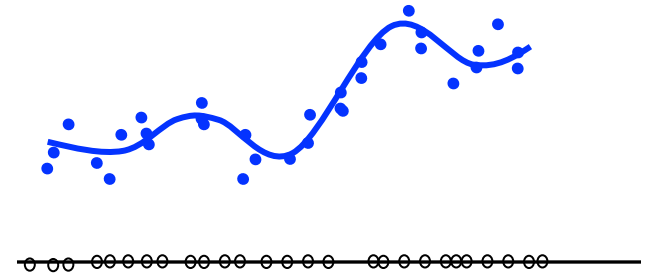
$$\hat{\beta}_{\text{MAP}} = \arg \min_{\beta} \sum_{i=1}^n (Y_i - X_i \beta)^2 + \underbrace{\lambda \|\beta\|_1}_{\text{constant}(\sigma^2, t)}$$

**Lasso**

Prior belief that  $\beta$  is Laplace with zero-mean biases solution to “sparse”  $\beta$

# Beyond Linear Regression

Polynomial regression  
Regression with nonlinear features



Kernelized Ridge Regression

Local Kernel Regression

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

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

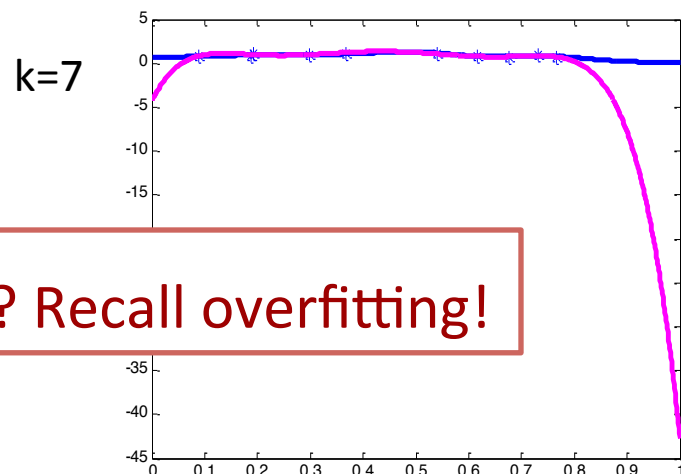
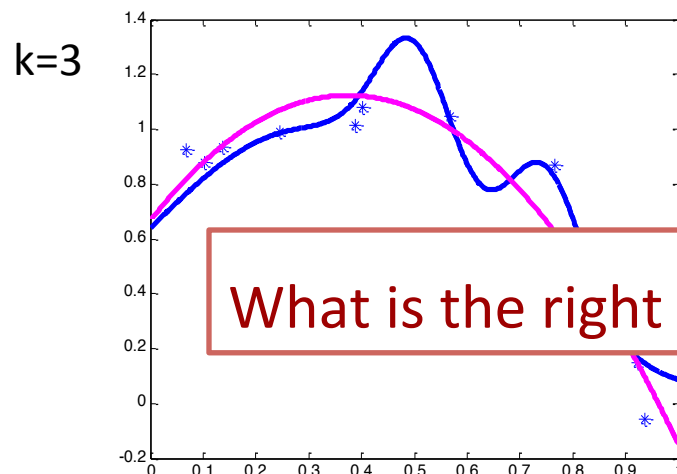
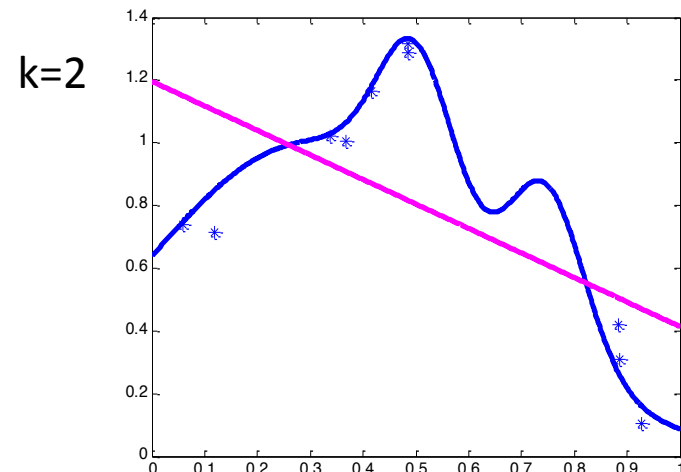
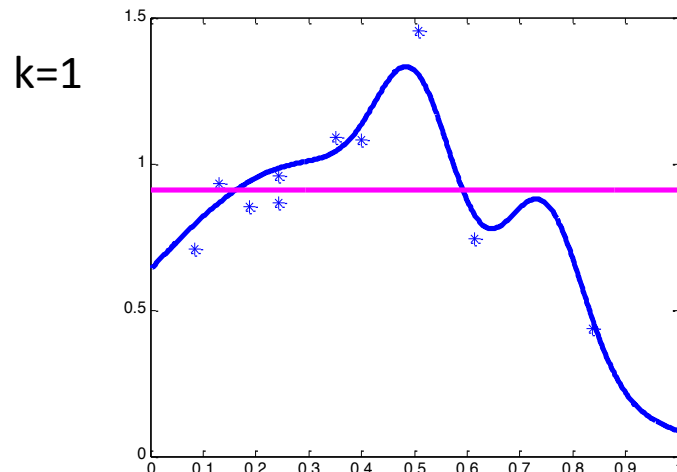
$$\text{where } \mathbf{A} = \begin{bmatrix} 1 & X_1 & X_1^2 & \dots & X_1^m \\ \vdots & & & \ddots & \vdots \\ 1 & X_n & X_n^2 & \dots & X_n^m \end{bmatrix}$$

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

$$+ \sum_{i=1}^p \sum_{j=1}^p \beta_{ij} X^{(i)} X^{(j)} + \sum_{i=1}^p \sum_{j=1}^p \sum_{k=1}^p X^{(i)} X^{(j)} X^{(k)} \\ + \dots \text{terms up to degree m}$$

# Polynomial Regression

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

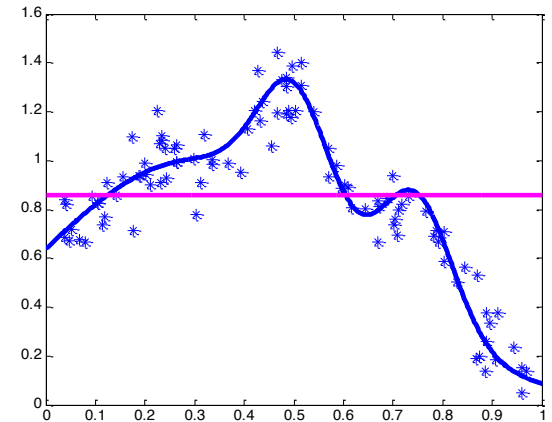
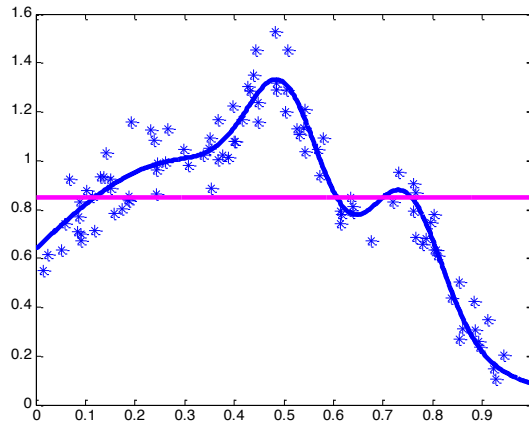
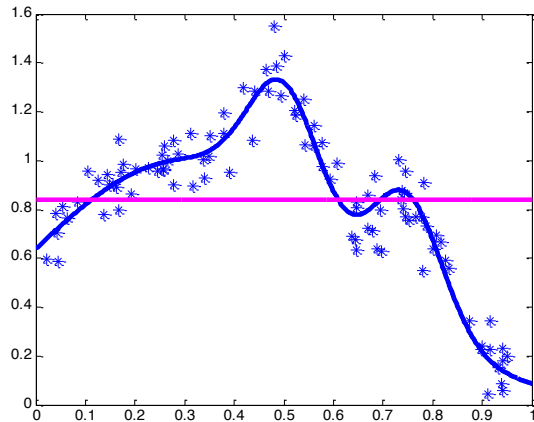


What is the right order? Recall overfitting!

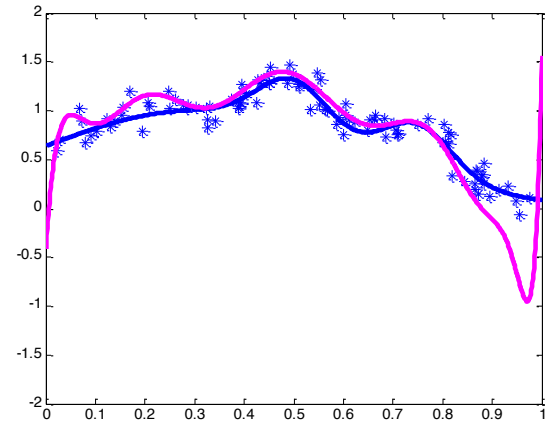
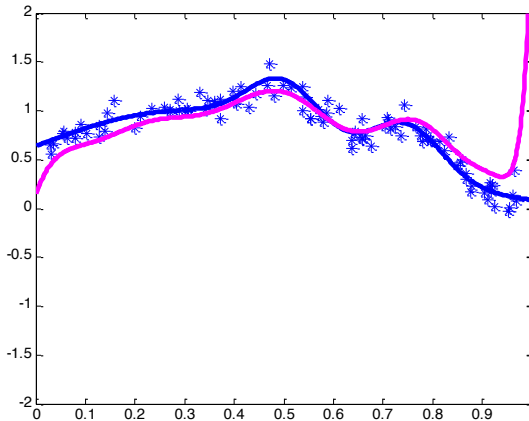
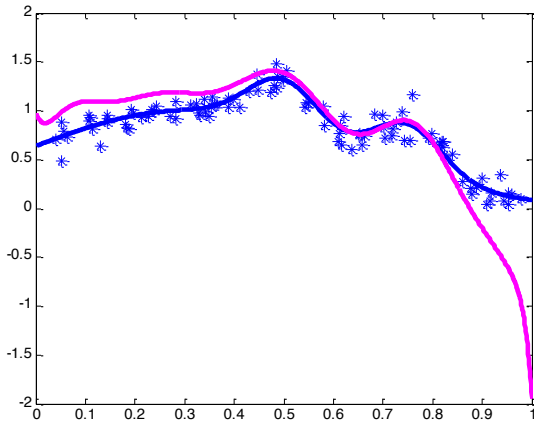
# Bias – Variance Tradeoff

3 Independent training datasets

Large bias, Small variance – poor approximation but robust/stable

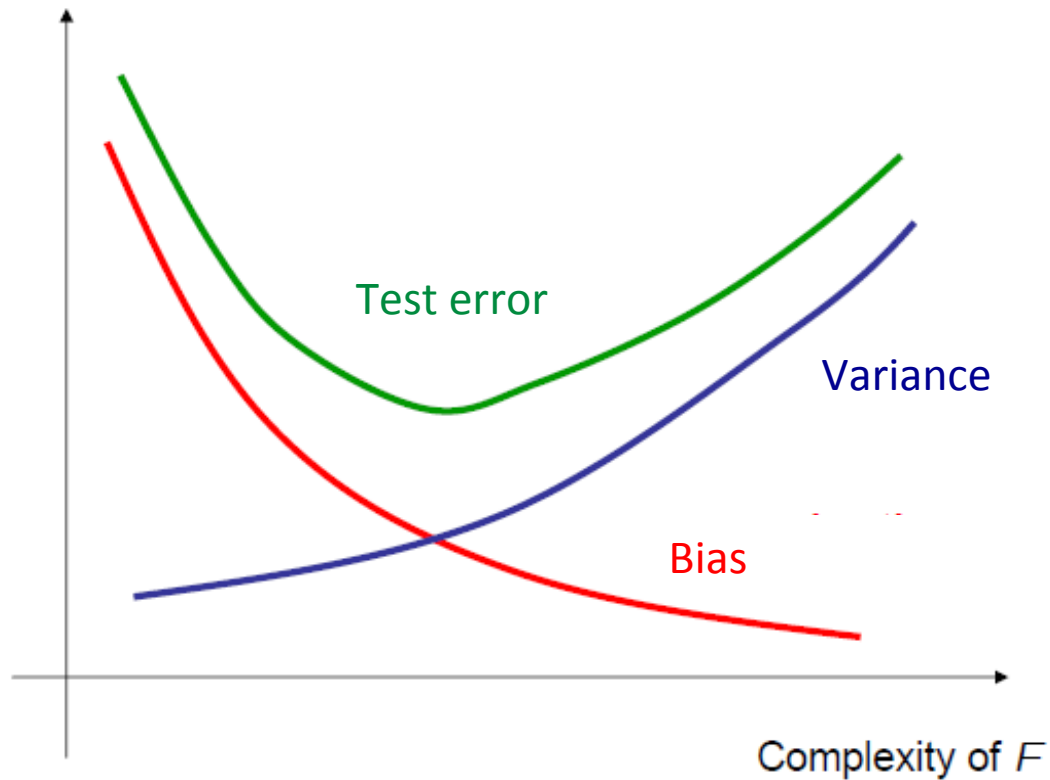


Small bias, Large variance – good approximation but unstable

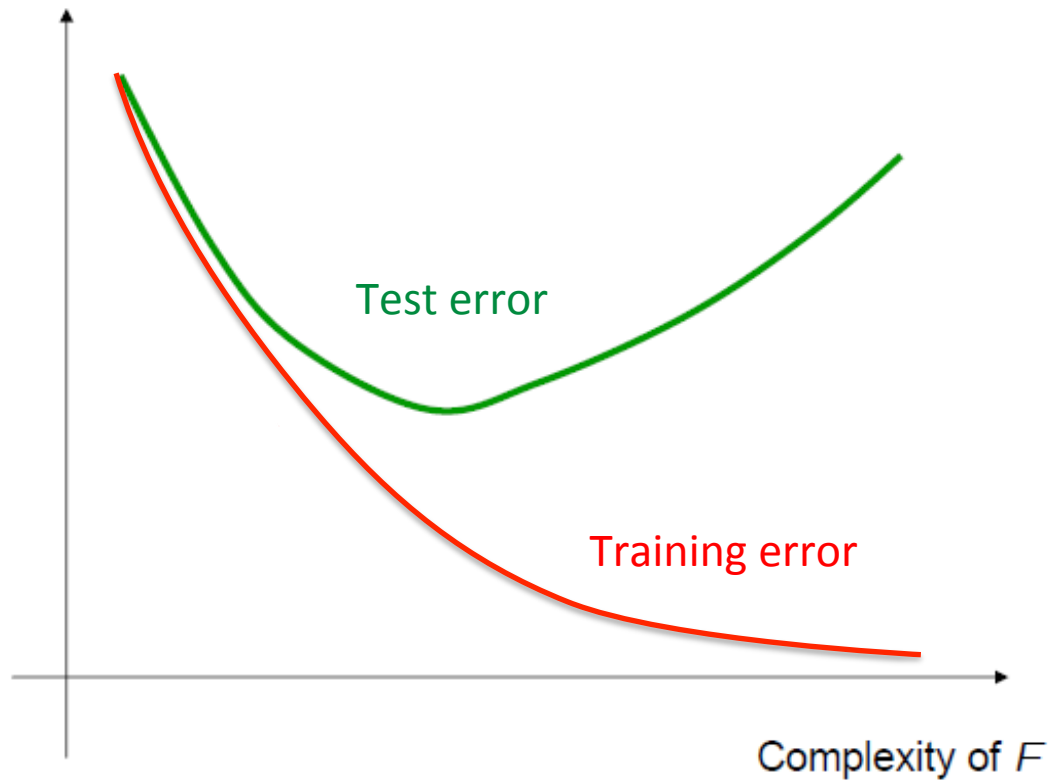




# Effect of Model Complexity



# Effect of Model Complexity



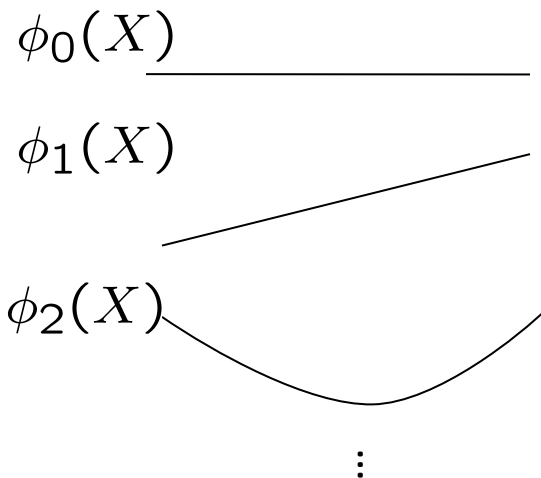
# Regression with basis functions

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

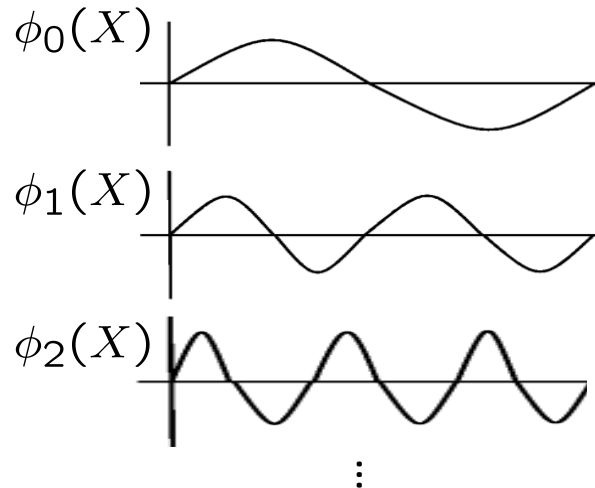
Basis coefficients

Basis functions (Linear combinations yield meaningful spaces)

Polynomial Basis

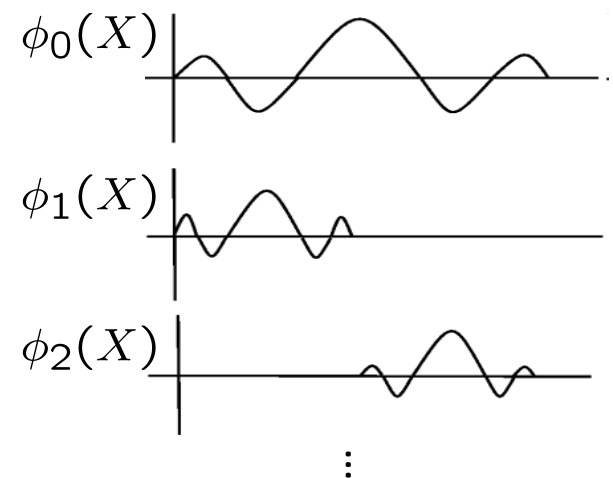


Fourier Basis



Good representation for  
periodic functions

Wavelet Basis



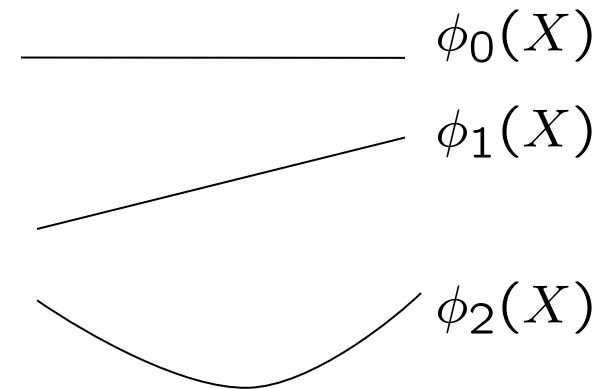
Good representation for  
local functions

# 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

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

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

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

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

**Can we use kernels?**

# Ridge regression (dual)

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

## Similarity with SVMs

Primal problem:

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

SVM Primal problem:

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

Lagrangian:

$$\sum_{i=1}^n z_i^2 + \lambda \|\beta\|_2^2 + \sum_{i=1}^n \alpha_i (z_i - Y_i + X_i \beta)$$

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

# Ridge regression (dual)

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

Dual problem:

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

$\alpha = \{\alpha_i\}$  for  $i = 1, \dots, n$

Taking derivatives of Lagrangian wrt  $\beta$  and  $z_i$  we get:

$$\beta = -\frac{1}{2\lambda} \mathbf{A}^T \alpha \quad z_i = -\frac{\alpha_i}{2}$$

Dual problem: 
$$\max_{\alpha} -\frac{\alpha^T \alpha}{4} - \frac{1}{4\lambda} \alpha^T \mathbf{A} \mathbf{A}^T \alpha - \alpha^T \mathbf{Y}$$

n-dimensional optimization problem

# Ridge regression (dual)

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

Dual problem:

$$\max_{\alpha} -\frac{\alpha^T \alpha}{4} - \frac{1}{4\lambda} \alpha^T \mathbf{A} \mathbf{A}^T \alpha - \alpha^T \mathbf{Y} \quad \Rightarrow \hat{\alpha} = - \left( \frac{\mathbf{A} \mathbf{A}^T}{\lambda} + \mathbf{I} \right)^{-1} 2 \mathbf{Y}$$

can get back  $\hat{\beta} = -\frac{1}{2\lambda} \mathbf{A}^T \hat{\alpha} = \mathbf{A}^T (\mathbf{A} \mathbf{A}^T + \lambda \mathbf{I})^{-1} \mathbf{Y}$

Weighted average of  
training points

Weight of each training point



# Kernelized ridge regression

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

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

Using dual, can re-write solution as:

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

How does this help?

- Only need to invert  $n \times n$  matrix (instead of  $p \times p$  or  $m \times m$ )
- More importantly, kernel trick!

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

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

# Kernelized ridge regression

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

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

Examples of kernels:

Polynomials of degree exactly  $d$   $K(\mathbf{u}, \mathbf{v}) = (\mathbf{u} \cdot \mathbf{v})^d$

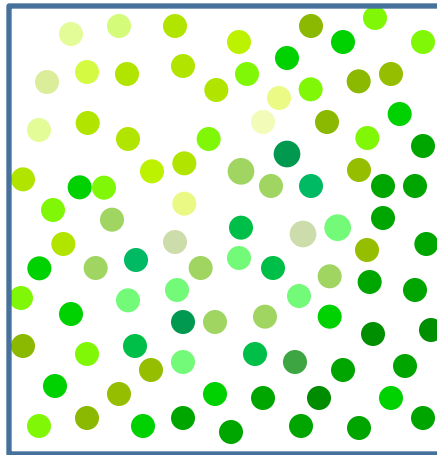
Polynomials of degree up to  $d$   $K(\mathbf{u}, \mathbf{v}) = (\mathbf{u} \cdot \mathbf{v} + 1)^d$

Gaussian/Radial kernels  $K(\mathbf{u}, \mathbf{v}) = \exp\left(-\frac{\|\mathbf{u} - \mathbf{v}\|^2}{2\sigma^2}\right)$

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

# Local Kernel Regression

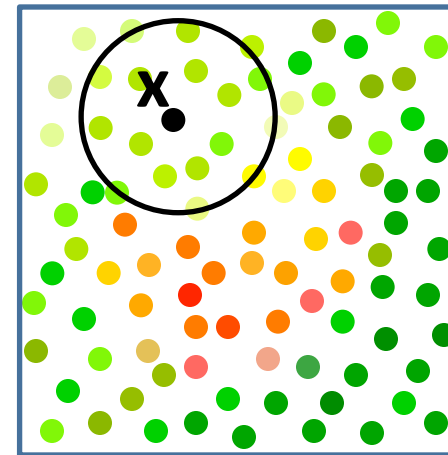
- What is the temperature in the room?



$$\hat{T} = \frac{1}{n} \sum_{i=1}^n Y_i$$

**Average**

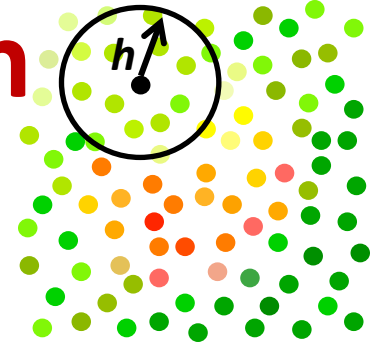
at location  $x$ ?



$$\hat{T}(x) = \frac{\sum_{i=1}^n Y_i \mathbf{1}_{\|X_i - x\| \leq h}}{\sum_{i=1}^n \mathbf{1}_{\|X_i - x\| \leq h}}$$

**"Local" Average**

# Local Kernel Regression



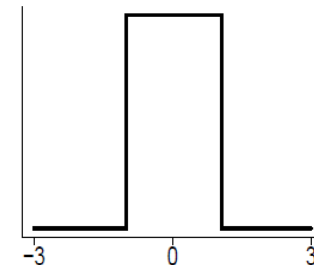
- Nonparametric estimator akin to kNN
- Nadaraya-Watson Kernel Estimator

$$\hat{f}_n(X) = \sum_{i=1}^n w_i Y_i \quad \text{Where} \quad w_i(X) = \frac{K\left(\frac{X-X_i}{h}\right)}{\sum_{i=1}^n K\left(\frac{X-X_i}{h}\right)}$$

- Weight each training point based on distance to test point
- Boxcar kernel yields local average

boxcar kernel :

$$K(x) = \frac{1}{2}I(x),$$



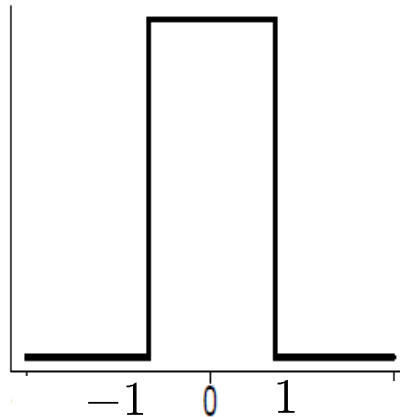
$$K(x) \geq 0,$$

$$\int K(x)dx = 1$$

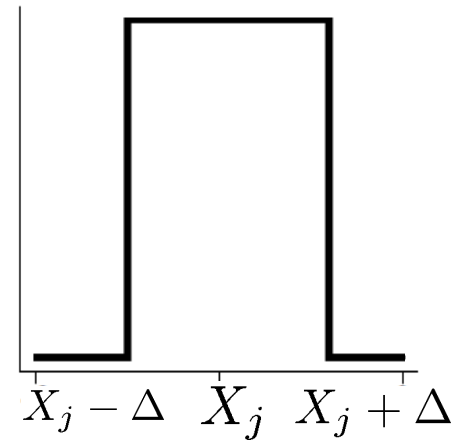
# Kernels

boxcar kernel :

$$K(x) = \frac{1}{2}I(x),$$

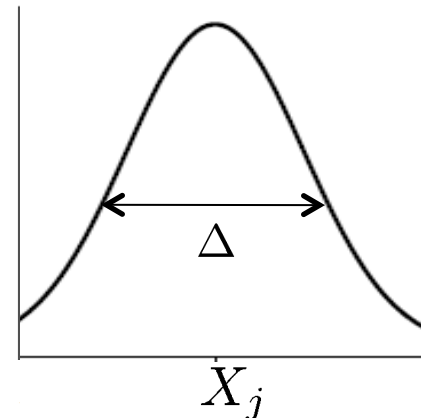
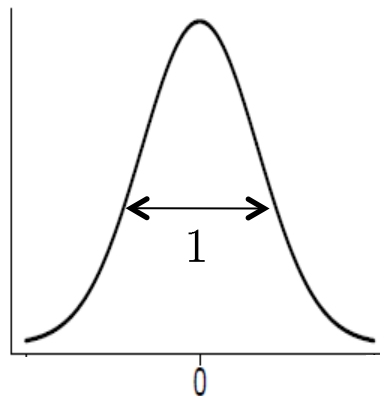


$$K\left(\frac{X_j - x}{\Delta}\right)$$



Gaussian kernel :

$$K(x) = \frac{1}{\sqrt{2\pi}}e^{-x^2/2}$$



# Choice of kernel bandwidth $h$

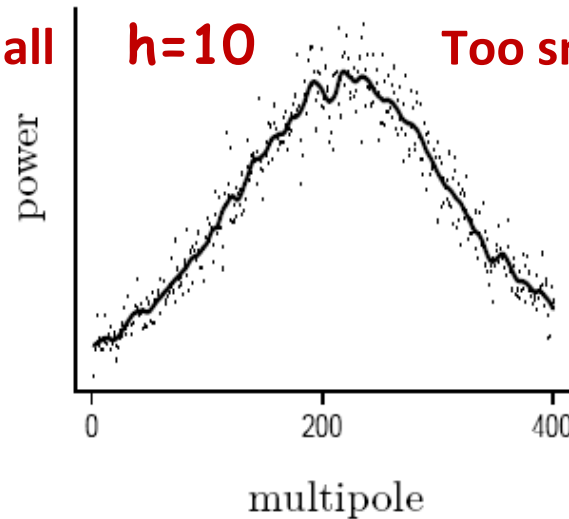
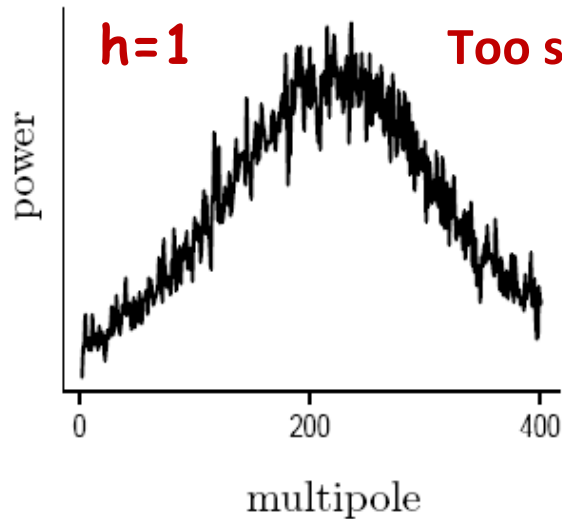
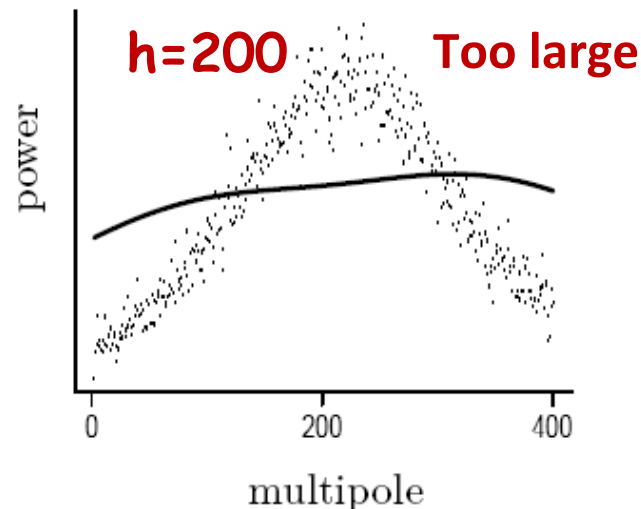
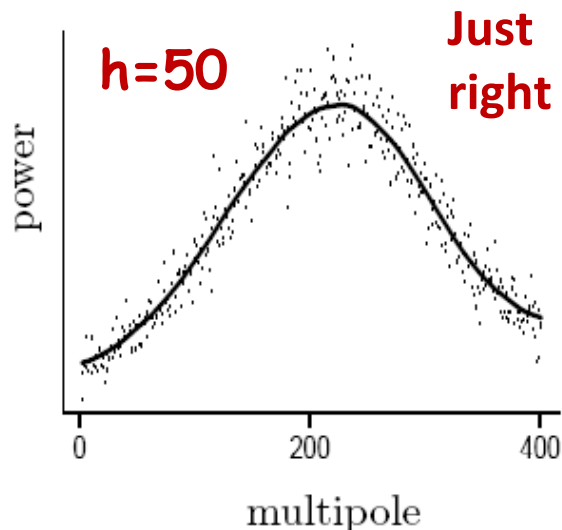


Image Source:  
Larry's book – All  
of Nonparametric  
Statistics

Choice of kernel is  
not that important



# Kernel Regression as Weighted Least Squares

$$\min_f \sum_{i=1}^n w_i (f(X_i) - Y_i)^2 \qquad w_i(X) = \frac{K\left(\frac{X - X_i}{h}\right)}{\sum_{i=1}^n K\left(\frac{X - X_i}{h}\right)}$$

Weighted Least Squares

Kernel regression corresponds to locally constant estimator obtained from (locally) weighted least squares

i.e. set  $f(X_i) = \beta$  (a constant)

# Kernel Regression as Weighted Least Squares

set  $f(X_i) = \beta$  (a constant)

$$\min_{\beta} \sum_{i=1}^n w_i (\underbrace{\beta}_{\text{constant}} - Y_i)^2$$

$$w_i(X) = \frac{K\left(\frac{X - X_i}{h}\right)}{\sum_{i=1}^n K\left(\frac{X - X_i}{h}\right)}$$

$$\frac{\partial J(\beta)}{\partial \beta} = 2 \sum_{i=1}^n w_i (\beta - Y_i) = 0$$

Notice that  $\sum_{i=1}^n w_i = 1$

$$\Rightarrow \hat{f}_n(X) = \hat{\beta} = \sum_{i=1}^n w_i Y_i$$



# Local Linear/Polynomial Regression

$$\min_f \sum_{i=1}^n w_i (f(X_i) - Y_i)^2 \qquad w_i(X) = \frac{K\left(\frac{X - X_i}{h}\right)}{\sum_{i=1}^n K\left(\frac{X - X_i}{h}\right)}$$

Weighted Least Squares

Local Polynomial regression corresponds to locally polynomial estimator obtained from (locally) weighted least squares

$$f(X_i) = \beta_0 + \beta_1(X_i - X) + \frac{\beta_2}{2!}(X_i - X)^2 + \dots + \frac{\beta_p}{p!}(X_i - X)^p$$

i.e. set

(local polynomial of degree p around X)

# What you should know

## Linear Regression

- Least Squares Estimator

- Normal Equations

- Gradient Descent

- Probabilistic Interpretation (connection to MCLE)

## Regularized Linear Regression (connection to MCAP)

- Ridge Regression, Lasso

## Beyond Linear

- Polynomial regression, Regression with Basis functions and Non-linear features, Bias-variance tradeoff, Kernelized ridge regression, Local Kernel Regression and Weighted Least Squares

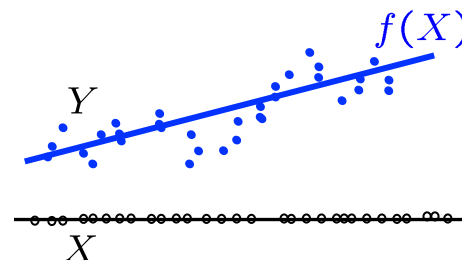


# Prior for kernelized ridge regression

- Recall ridge regression was obtained as MAP under a Gaussian prior for a linear model with Gaussian noise

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

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



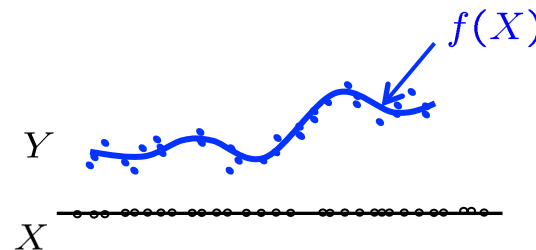
- For kernelized ridge regression, what's the prior?

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

$$f(X) = \phi(X)\beta \quad \beta \sim \mathcal{N}(0, \tau^2 \mathbf{I})$$

$$\Rightarrow f(X) \sim \mathcal{N}(0, \tau^2 \phi(X) \cdot \phi(X))$$

$$\sim \mathcal{N}(0, \tau^2 K_{X,X})$$



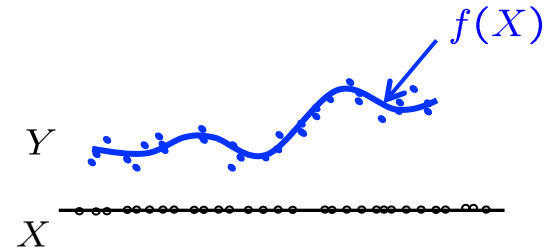
# Prior for kernelized ridge regression

- For kernelized ridge regression, what's the prior?

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

$$f(X) = \phi(X)\beta \quad \beta \sim \mathcal{N}(0, \tau^2 \mathbf{I})$$

$$\Rightarrow f(X) \sim \mathcal{N}(0, \tau^2 \phi(X) \cdot \phi(X)) \sim \mathcal{N}(0, \tau^2 K_{X,X})$$



$$\begin{bmatrix} f(X_i) \\ f(X_j) \end{bmatrix} = \begin{bmatrix} \phi(X_i) \\ \phi(X_j) \end{bmatrix} \beta$$

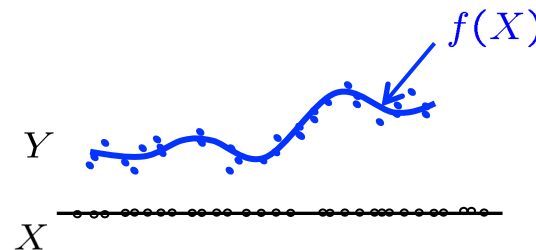
$$\begin{bmatrix} f(X_i) \\ f(X_j) \end{bmatrix} \sim \mathcal{N} \left( 0, \tau^2 \begin{bmatrix} \phi(X_i) \cdot \phi(X_i) & \phi(X_i) \cdot \phi(X_j) \\ \phi(X_j) \cdot \phi(X_i) & \phi(X_j) \cdot \phi(X_j) \end{bmatrix} \right)$$

$$\sim \mathcal{N} \left( 0, \tau^2 \begin{bmatrix} K(i,i) & K(i,j) \\ K(j,i) & K(j,j) \end{bmatrix} \right)$$

# Prior for kernelized ridge regression

- For kernelized ridge regression, what's the prior?

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



$$\begin{bmatrix} f(X_i) \\ f(X_j) \end{bmatrix} \sim \mathcal{N} \left( 0, \tau^2 \begin{bmatrix} K(i, i) & K(i, j) \\ K(j, i) & K(j, j) \end{bmatrix} \right)$$

Then  $\tau^2 K(i, j) = \Sigma(i, j)$  i.e. we can interpret Kernel function as the covariance function under a Gaussian process prior.

$$f \sim \mathcal{N}(0, \Sigma)$$

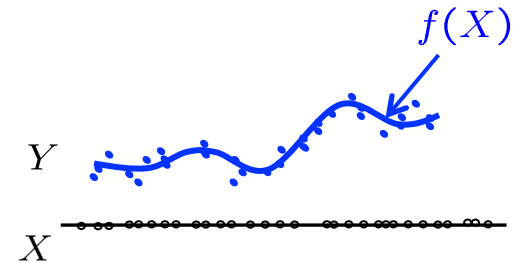
$$\text{i.e. } \mathbb{E}[f(X_i)f(X_j)] = \Sigma(i, j)$$

# Gaussian process regression

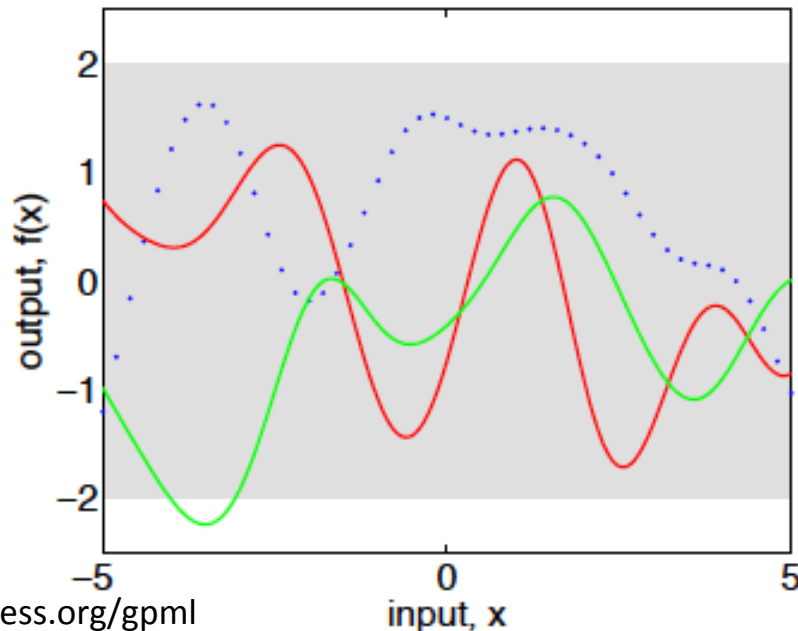
$$Y = f(X) + \epsilon, \quad \epsilon \sim \mathcal{N}(0, \sigma^2 \mathbf{I})$$

$f \sim \mathcal{N}(0, \Sigma)$       Gaussian process prior

i.e.  $\mathbb{E}[f(X_i)f(X_j)] = \Sigma(i, j)$



3 functions drawn at random from GP **prior**



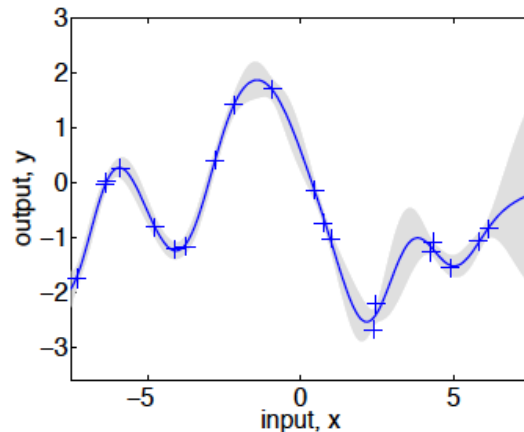
$$\Sigma(i, j) = \exp \left( -\frac{1}{2} \|X_i - X_j\|^2 \right)$$

Shaded region represents  
pointwise mean +/-  
2 standard deviation

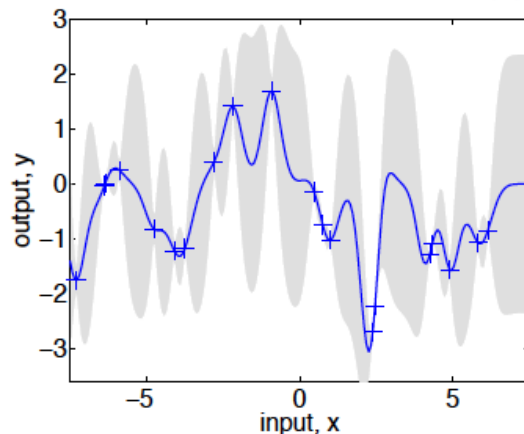
# Gaussian process regression

Controls “smoothness” of functions

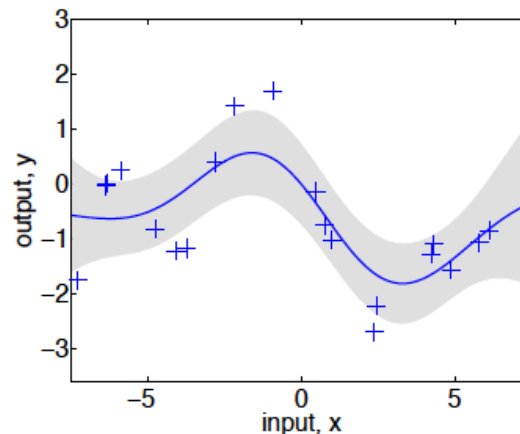
$$k_y(x_p, x_q) = \sigma_f^2 \exp\left(-\frac{1}{2\ell^2}(x_p - x_q)^2\right) + \sigma_n^2 \delta_{pq}$$



(a),  $\ell = 1$



(b),  $\ell = 0.3$



(c),  $\ell = 3$

Figure 2.5: (a) Data is generated from a GP with hyperparameters  $(\ell, \sigma_f, \sigma_n) = (1, 1, 0.1)$ , as shown by the + symbols. Using Gaussian process prediction with these hyperparameters we obtain a 95% confidence region for the underlying function  $f$  (shown in grey). Panels (b) and (c) again show the 95% confidence region, but this time for hyperparameter values  $(0.3, 1.08, 0.00005)$  and  $(3.0, 1.16, 0.89)$  respectively.

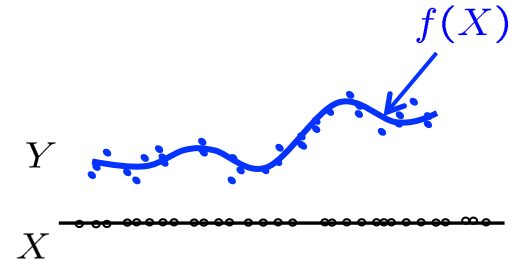


# Gaussian process regression

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

$$f \sim \mathcal{N}(0, \Sigma) \quad \text{Gaussian process prior}$$

$$\text{i.e. } \mathbb{E}[f(X_i)f(X_j)] = \Sigma(i, j)$$



Lets derive the MAP estimate under this prior:

$$\begin{bmatrix} \mathbf{Y} \\ f(X) \end{bmatrix} \sim \mathcal{N}\left(0, \begin{bmatrix} \Sigma_{nn} + \sigma^2 \mathbf{I} & \Sigma_{nX} \\ \Sigma_{Xn} & \Sigma_{XX} \end{bmatrix}\right)$$

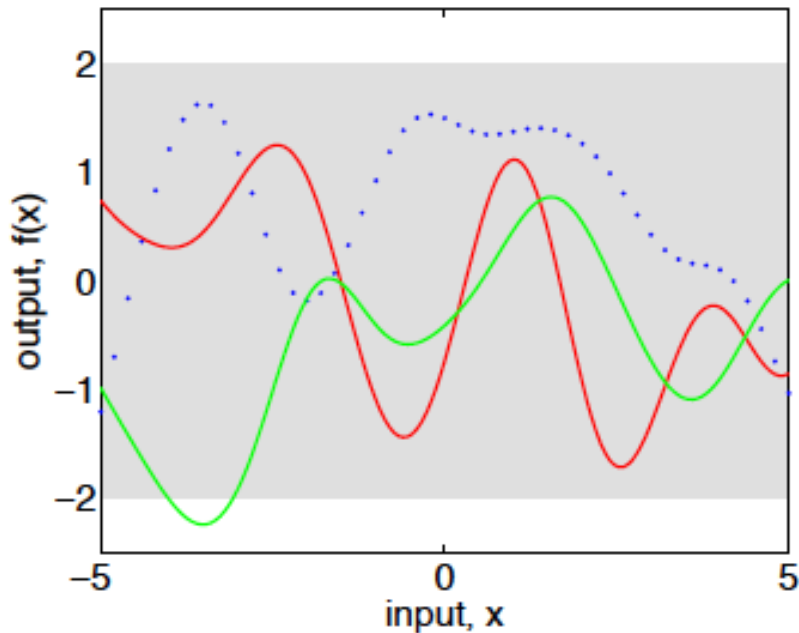
Posterior distribution

$$\Rightarrow f(X)|\mathbf{Y} \sim \mathcal{N}(\Sigma_{Xn}(\Sigma_{nn} + \sigma^2 \mathbf{I})^{-1}\mathbf{Y}, \Sigma_{XX} - \Sigma_{Xn}(\Sigma_{nn} + \sigma^2 \mathbf{I})^{-1}\Sigma_{nX})$$

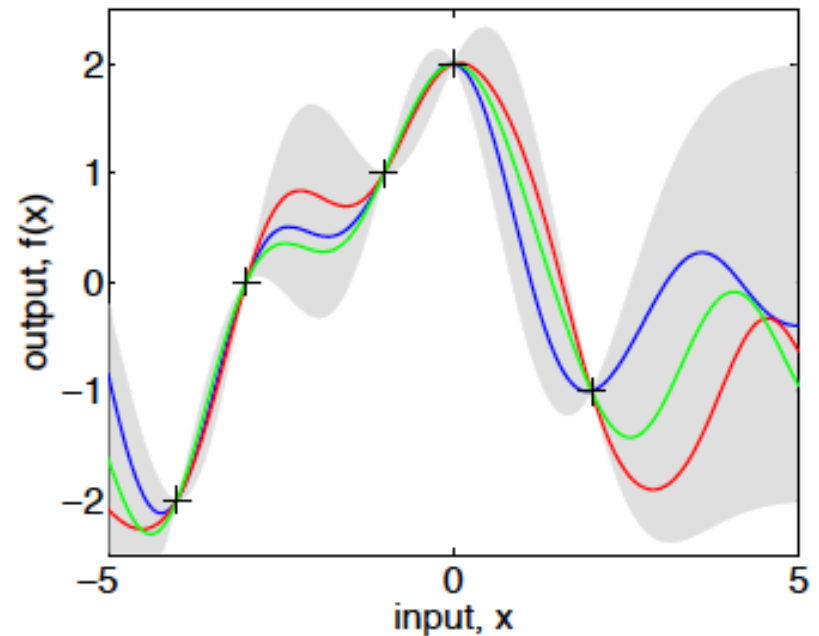
$$\text{MAP estimate } \hat{f}(X) = \Sigma_{Xn}(\Sigma_{nn} + \sigma^2 \mathbf{I})^{-1}\mathbf{Y} \equiv \mathbf{K}_X(\mathbf{K} + \lambda \mathbf{I})^{-1}\mathbf{Y}$$

# Gaussian process regression

3 functions drawn at random from GP **prior**



3 functions drawn at random from **posterior** based on five noiseless observations



Shaded region represents pointwise mean  $\pm 2$  standard deviation

# What you should know

Nonlinear regression

- Polynomial regression

- Regression with Non-linear features

- Kernelized ridge regression

- GP Regression

Bias-variance decomposition



# Ridge regression (dual)

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

$$\begin{aligned}\hat{\beta} &= (\mathbf{A}^T \mathbf{A} + \lambda \mathbf{I})^{-1} \mathbf{A}^T \mathbf{Y} \\ &= \mathbf{A}^T (\mathbf{A} \mathbf{A}^T + \lambda \mathbf{I})^{-1} \mathbf{Y}\end{aligned}$$

## Similarity with SVMs

Primal problem:

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

SVM Primal problem:

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

Dual problem:

$$\max_{\alpha} \quad -\frac{\alpha^\top \alpha}{2} - \frac{1}{2\lambda} \alpha^\top \mathbf{A} \mathbf{A}^\top \alpha - \alpha^\top \mathbf{Y}$$

$\alpha_i$  – one per training point (n-dimensional optimization problem)

# Ridge regression (dual)

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

## Alternate derivation

Dual problem:

$$\max_{\alpha} -\frac{\alpha^T \alpha}{2} - \frac{1}{2\lambda} \alpha^T \mathbf{A} \mathbf{A}^T \alpha - \alpha^T \mathbf{Y} \quad \Rightarrow \hat{\alpha} = - \left( \frac{\mathbf{A} \mathbf{A}^T}{\lambda} + \mathbf{I} \right)^{-1} \mathbf{Y}$$

$\alpha_i$  – one per training point

can get back  $\hat{\beta} = -\frac{1}{\lambda} \sum_{i=1}^n \alpha_i X_i^T = -\frac{1}{\lambda} \mathbf{A}^T \alpha = \mathbf{A}^T (\mathbf{A} \mathbf{A}^T + \lambda \mathbf{I})^{-1} \mathbf{Y}$


Weight of each training point

# Ridge regression (dual)

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

## Dual interpretation

$$\hat{\beta} = -\frac{1}{\lambda} \sum_{i=1}^n \alpha_i X_i^\top = -\frac{1}{\lambda} \mathbf{A}^\top \alpha = \mathbf{A}^\top (\mathbf{A} \mathbf{A}^\top + \lambda \mathbf{I})^{-1} \mathbf{Y}$$

 Weight of each training point

Prediction  $\hat{f}_n(X) = \mathbf{X} \hat{\beta} = -\frac{1}{\lambda} \sum_{i=1}^n \alpha_i \mathbf{X} \cdot X_i$

which can be kernelized

$$\hat{f}_n(X) = \mathbf{K}_X (\mathbf{K} + \lambda \mathbf{I})^{-1} \mathbf{Y} = -\frac{1}{\lambda} \sum_{i=1}^n \alpha_i K(\mathbf{X}, X_i)$$

Prediction is a weighted sum of “bumps” placed at each training data.

# Online ques

- If you train a linear regression estimator with only half the data, its variance is larger.
- If you train a linear regression estimator with only half the data, its bias is smaller.
- As you increase the bandwidth of a Gaussian kernel, the bias of the corresponding kernelized ridge regression decreases.
- A least square regression estimator essentially forms a weighted sum of the training labels.
- Ridge regression estimator essentially forms a weighted sum of the training labels.
- As the regularization parameter  $\lambda$  is increased, the estimated regression coefficients become sparser.
- If we have a dataset of brain scans of 20 patients, and we are trying to predict age of a patient using their recorded brain activity in 600 regions of interest in the brain as features, then we should prefer to use ridge regression over least squares linear regression.
- Polynomial regression estimator involves a nonlinear transformation of the training labels.
- Polynomial regression estimator typically involves a nonlinear transformation of the training features.
- As the number of test points increases, the average mean square error on test points goes down.
- As the number of training points increases, the average mean square error on test points goes down.
- As the number of test points increases, the average mean square error on training points goes down.
- If we set the covariance function of a Gaussian process prior to be diagonal i.e. 1 if  $X_i = X_j$  and 0 otherwise, then we can't hope to learn the regression function.
- A polynomial of degree up to 2 has a lower bias than a polynomial of degree up to 3 for fitting labels that are generated according to a quadratic signal function plus noise.
- Suppose we are solving ridge regression using gradient descent. Then increasing the regularization parameter  $\lambda$  in ridge regression causes a larger change in the regression coefficients.
- Gradient descent update for kernel ridge regression?
- The average Constant fit?

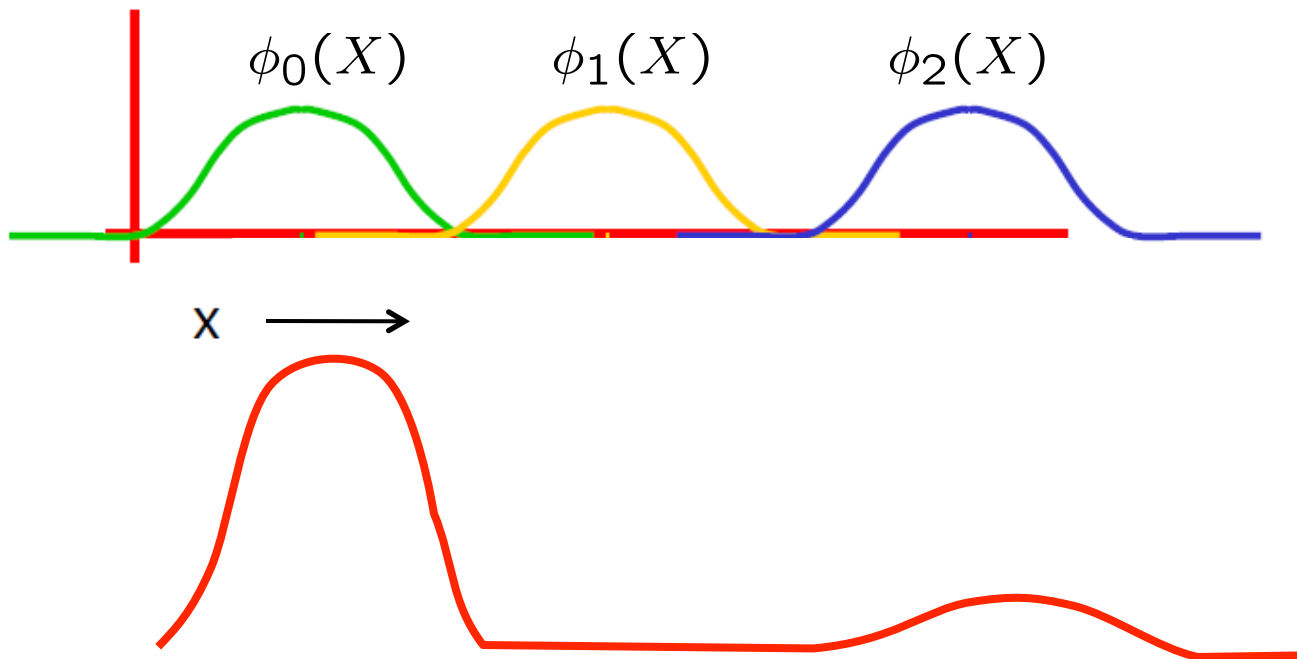




# Local Regression

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

Basis coefficients  $\leftarrow$  Nonlinear features/basis functions

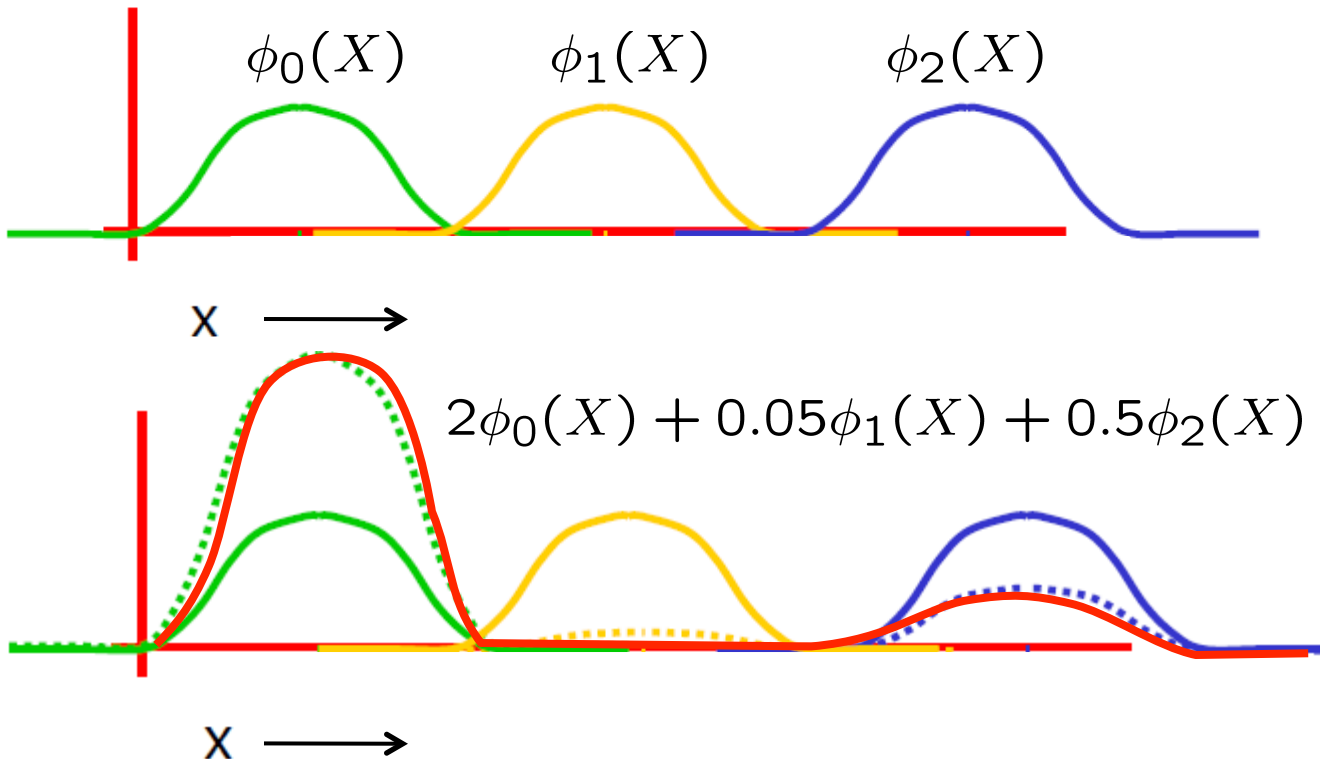


Globally supported basis functions (polynomial, fourier) will not yield a good representation

# Local Regression

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

Basis coefficients  $\leftarrow$  Nonlinear features/basis functions



Globally supported basis functions (polynomial, fourier) will not yield a good representation

# What you should know

## Linear Regression

- Least Squares Estimator

- Normal Equations

- Gradient Descent

- Geometric and Probabilistic Interpretation (connection to MLE)

## Regularized Linear Regression (connection to MAP)

- Ridge Regression, Lasso

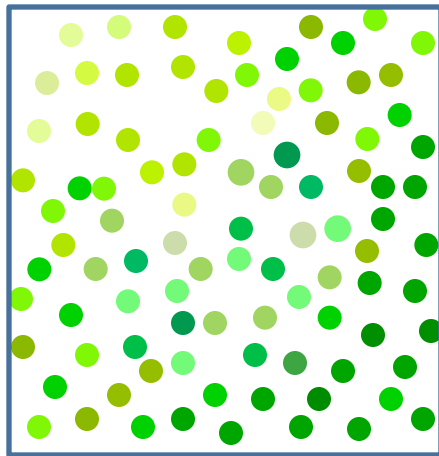
## Polynomial Regression, Basis (Fourier, Wavelet) Estimators

## Next time

- Kernel Regression (Localized)
- Regression Trees

# Temperature sensing

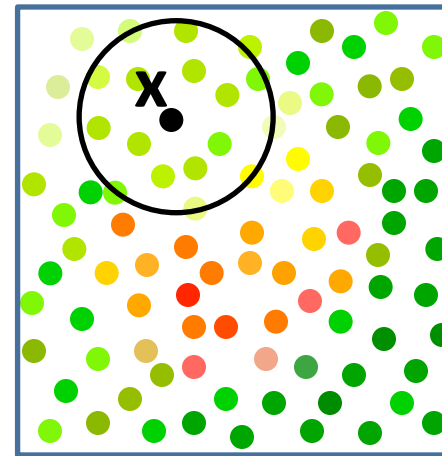
- What is the temperature in the room?



$$\hat{T} = \frac{1}{n} \sum_{i=1}^n Y_i$$

**Average**

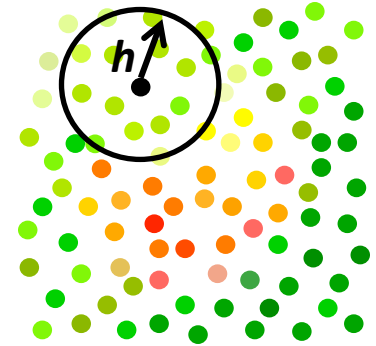
at location  $x$ ?



$$\hat{T}(x) = \frac{\sum_{i=1}^n Y_i \mathbf{1}_{\|X_i - x\| \leq h}}{\sum_{i=1}^n \mathbf{1}_{\|X_i - x\| \leq h}}$$

**"Local" Average**

# Kernel Regression



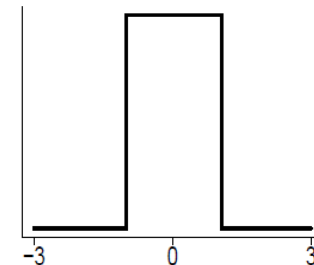
- Aka Local Regression
- Nadaraya-Watson Kernel Estimator

$$\hat{f}_n(X) = \sum_{i=1}^n w_i Y_i \quad \text{Where} \quad w_i(X) = \frac{K\left(\frac{X - X_i}{h}\right)}{\sum_{i=1}^n K\left(\frac{X - X_i}{h}\right)}$$

- Weight each training point based on distance to test point
- Boxcar kernel yields local average

boxcar kernel :

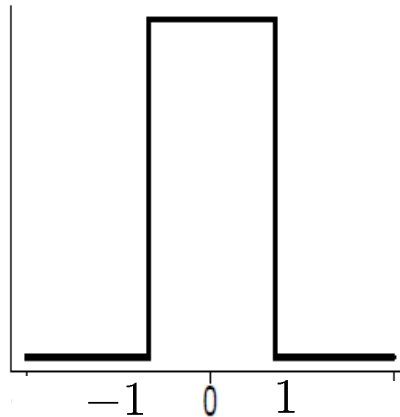
$$K(x) = \frac{1}{2}I(x),$$



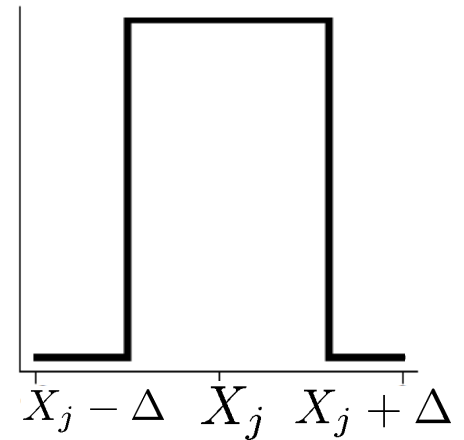
# Kernels

boxcar kernel :

$$K(x) = \frac{1}{2}I(x),$$

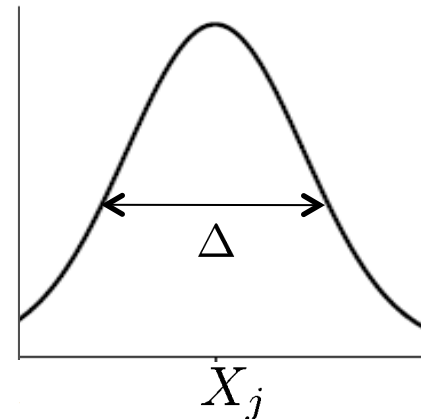
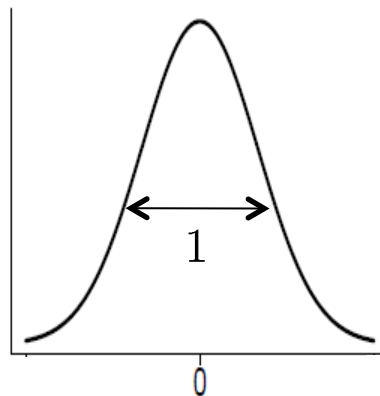


$$K\left(\frac{X_j - x}{\Delta}\right)$$



Gaussian kernel :

$$K(x) = \frac{1}{\sqrt{2\pi}}e^{-x^2/2}$$



# Choice of kernel bandwidth $h$

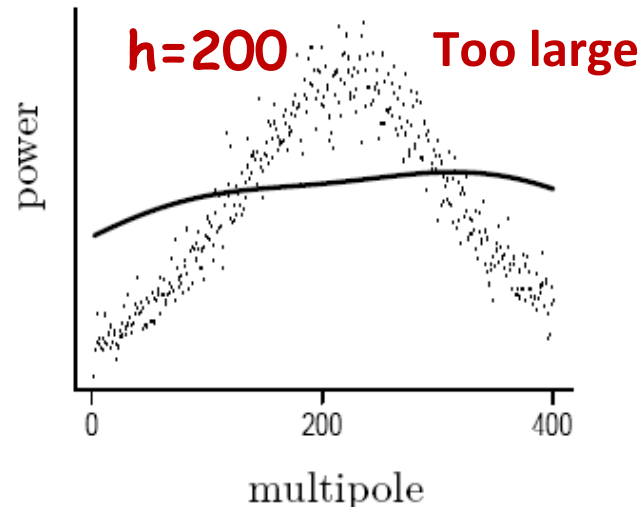
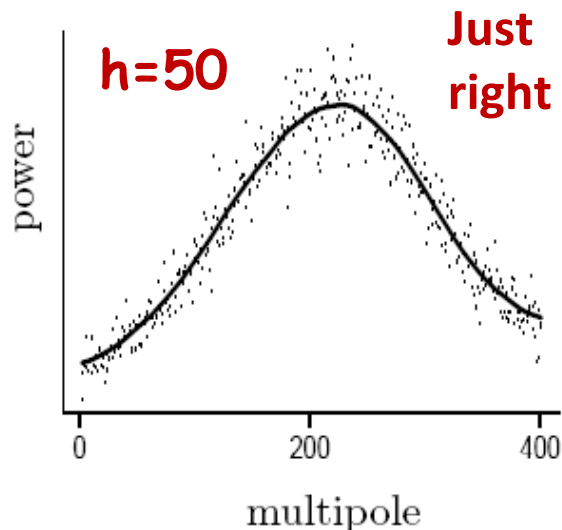
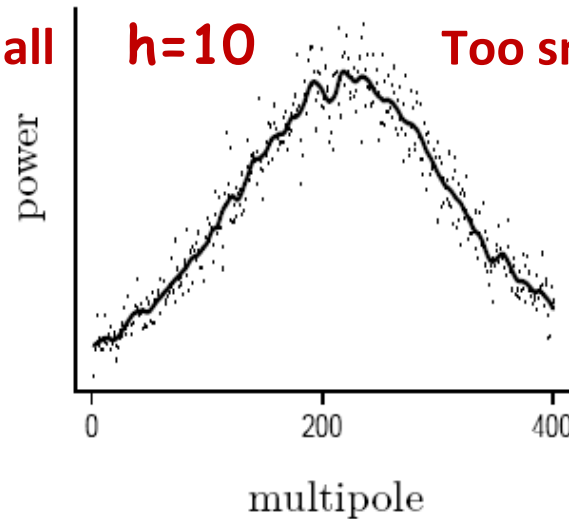
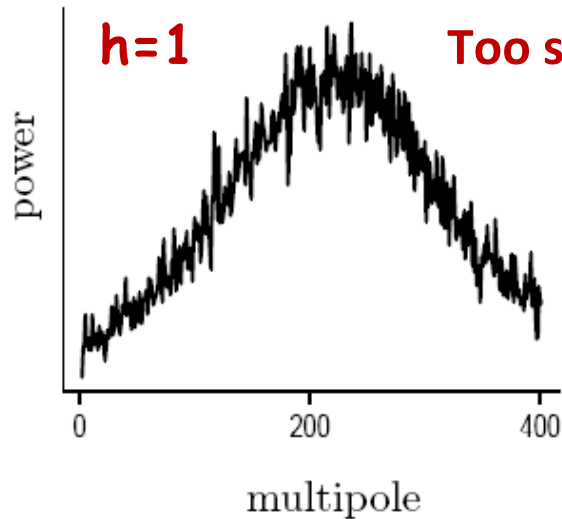


Image Source:  
Larry's book – All  
of Nonparametric  
Statistics

Choice of kernel is  
not that important



# Kernel Regression as Weighted Least Squares

$$\min_f \sum_{i=1}^n w_i (f(X_i) - Y_i)^2 \qquad w_i(X) = \frac{K\left(\frac{X - X_i}{h}\right)}{\sum_{i=1}^n K\left(\frac{X - X_i}{h}\right)}$$

Weighted Least Squares

Kernel regression corresponds to locally constant estimator obtained from (locally) weighted least squares

i.e. set  $f(X_i) = \beta$  (a constant)

# Kernel Regression as Weighted Least Squares

set  $f(X_i) = \beta$  (a constant)

$$\min_{\beta} \sum_{i=1}^n w_i (\underbrace{\beta}_{\text{constant}} - Y_i)^2$$

$$w_i(X) = \frac{K\left(\frac{X - X_i}{h}\right)}{\sum_{i=1}^n K\left(\frac{X - X_i}{h}\right)}$$

$$\frac{\partial J(\beta)}{\partial \beta} = 2 \sum_{i=1}^n w_i (\beta - Y_i) = 0$$

Notice that  $\sum_{i=1}^n w_i = 1$

$$\Rightarrow \hat{f}_n(X) = \hat{\beta} = \sum_{i=1}^n w_i Y_i$$

# Local Linear/Polynomial Regression

$$\min_f \sum_{i=1}^n w_i (f(X_i) - Y_i)^2 \qquad w_i(X) = \frac{K\left(\frac{X - X_i}{h}\right)}{\sum_{i=1}^n K\left(\frac{X - X_i}{h}\right)}$$

Weighted Least Squares

Local Polynomial regression corresponds to locally polynomial estimator obtained from (locally) weighted least squares

$$f(X_i) = \beta_0 + \beta_1(X_i - X) + \frac{\beta_2}{2!}(X_i - X)^2 + \dots + \frac{\beta_p}{p!}(X_i - X)^p$$

i.e. set

(local polynomial of degree p around X)

More in HW, 10-702 (statistical machine learning)



# Kernel Regression (Local)

$$\min_f \frac{1}{n} \sum_{i=1}^n w_i (f(X_i) - Y_i)^2$$

$$\frac{1}{n} \sum_{i=1}^n w_i = 1$$

## Weighted Least Squares

Weigh each training point based on distance to test point

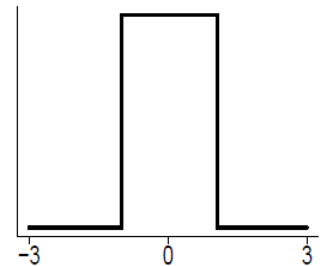
$$w_i(X) = \frac{K\left(\frac{X - X_i}{h}\right)}{\sum_{i=1}^n K\left(\frac{X - X_i}{h}\right)}$$

$K$  – Kernel

$h$  – Bandwidth of kernel

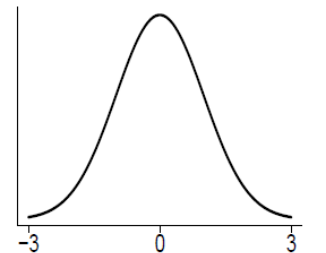
boxcar kernel :

$$K(x) = \frac{1}{2} I(x),$$




Gaussian kernel :

$$K(x) = \frac{1}{\sqrt{2\pi}} e^{-x^2/2}$$



# Nadaraya-Watson Kernel Regression

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

  
constant

$$w_i(X) = \frac{K\left(\frac{X - X_i}{h}\right)}{\sum_{i=1}^n K\left(\frac{X - X_i}{h}\right)}$$

$$\frac{\partial J(\beta)}{\partial \beta} = 2 \sum_{i=1}^n w_i (\beta - Y_i) = 0$$

$$\Rightarrow \hat{f}_n(X) = \hat{\beta} = \sum_{i=1}^n w_i Y_i$$

# Nadaraya-Watson Kernel Regression

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

$\downarrow$   
 constant

$$w_i(X) = \frac{K\left(\frac{X - X_i}{h}\right)}{\sum_{i=1}^n K\left(\frac{X - X_i}{h}\right)}$$

boxcar kernel :

$$K\left(\frac{X - X_i}{h}\right) = 1_{|X - X_i| \leq h}$$

$$\frac{\partial J(\beta)}{\partial \beta} = 2 \sum_{i=1}^n w_i (\beta - Y_i) = 0$$

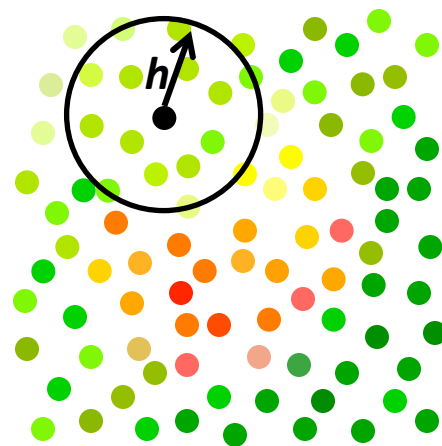
$$\Rightarrow \hat{f}_n(X) = \hat{\beta} = \sum_{i=1}^n w_i Y_i$$

with box-car kernel

$$= \frac{1}{n_h} \sum_{i=1}^n Y_i 1_{|X - X_i| \leq h}$$

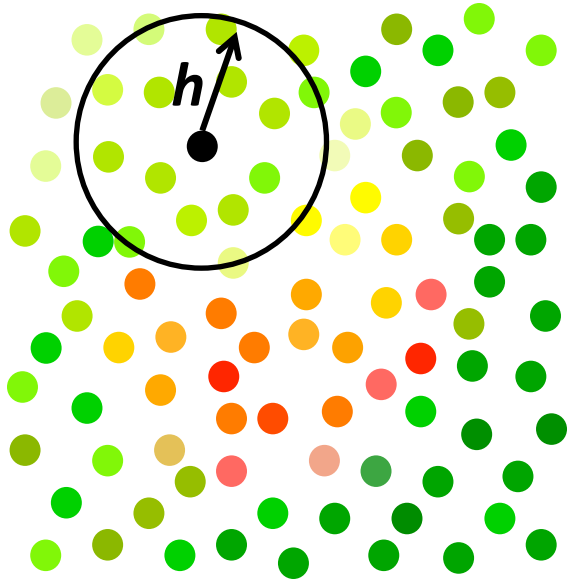
#pts in h ball around X

Sum of Ys in h ball around X



Recall: NN classifier  
Average  $\leftrightarrow$  majority vote

# Choice of Bandwidth



Should depend on  $n$ , # training data  
(determines variance)

Should depend on smoothness of  
function  
(determines bias)

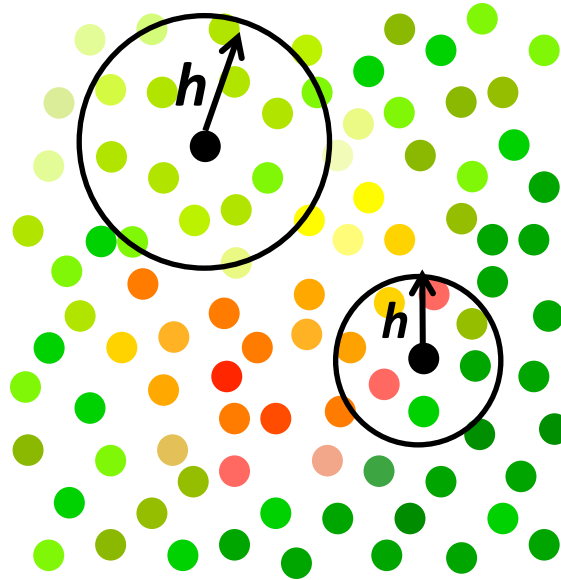
Large Bandwidth – average more data points, reduce noise (**Lower variance**)

Small Bandwidth – less smoothing, more accurate fit (**Lower bias**)

**Bias - Variance tradeoff** : More to come in later lectures



# Spatially adaptive regression



If function smoothness varies spatially, we want to allow bandwidth  $h$  to depend on  $X$

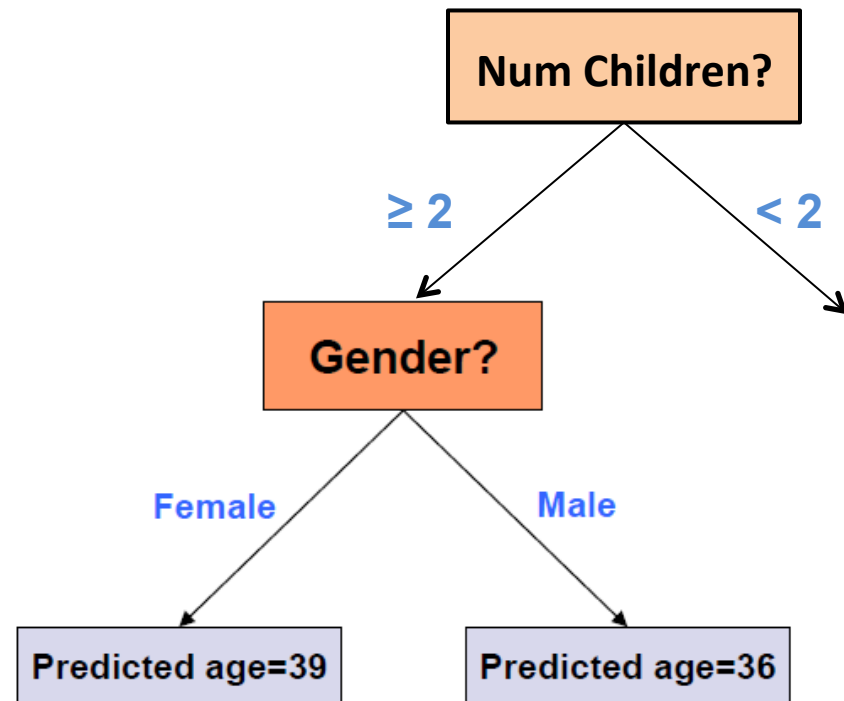
Local polynomials, splines, wavelets, regression trees ...

# Regression trees

$X^{(1)}$       ....       $X^{(p)}$        $Y$

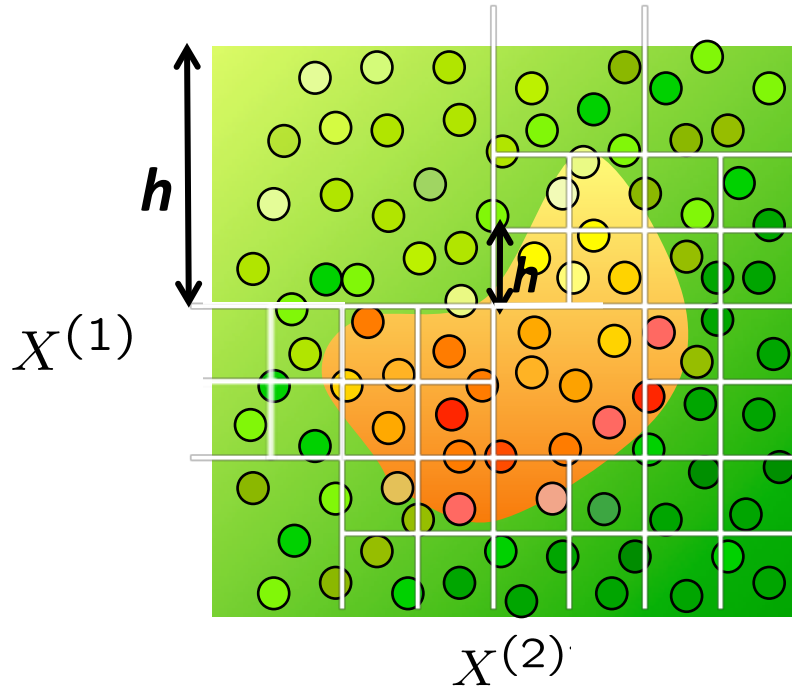
Gender	Rich?	Num. Children	# travel per yr.	Age
F	No	2	5	38
M	No	0	2	25
M	Yes	1	0	72
:	:	:	:	:

Binary Decision Tree

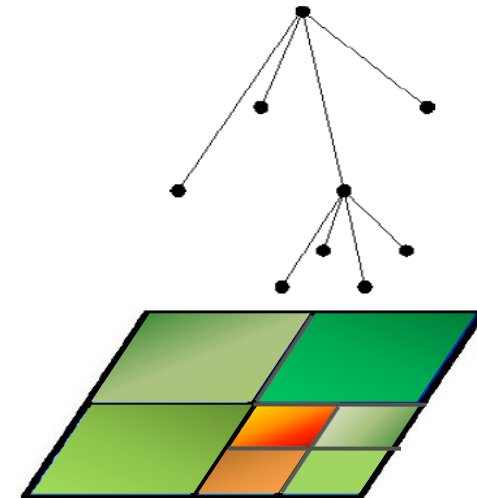


Average (fit a constant ) on the leaves

# Regression trees



Quad Decision Tree



$f$  - Polynomial fit on each leaf

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

If  $\sum_{\text{small cells}, L} \sum_{i \in L} (\hat{f}(X_i) - Y_i)^2 < \sum_{i \in \text{merged cell}} (\hat{f}(X_i) - Y_i)^2$ , then split

Else stop

Compare residual error with and without split 91

# Summary

## Discriminative vs Generative Classifiers

- Naïve Bayes vs Logistic Regression

## Regression

- Linear Regression
  - Least Squares Estimator
  - Normal Equations
  - Gradient Descent
  - Geometric Interpretation
  - Probabilistic Interpretation (connection to MLE)
- Regularized Linear Regression (connection to MAP)
  - Ridge Regression, Lasso
- Polynomial Regression, Basis (Fourier, Wavelet) Estimators
- Kernel Regression (Localized)
- Regression Trees