# Linear Regression

Aarti Singh

Machine Learning 10-701/15-781
Sept 27, 2010

**MACHINE LEARNING** DEPARTMENT

**Carnegie Mellon.**
**School of Computer Science**

# Discrete to Continuous Labels

**Classification**



Sports
Science
News



Anemic cell
Healthy cell

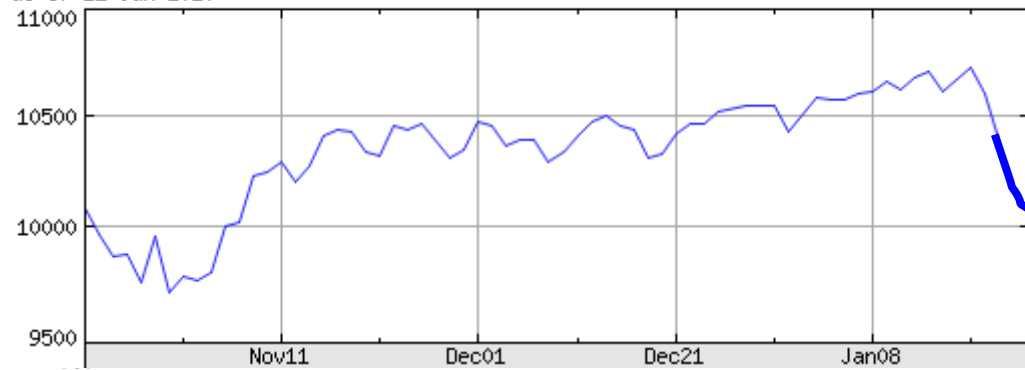**X = Document**     **Y = Topic**     **X = Cell Image**     **Y = Diagnosis**

**Regression**

Stock Market
Prediction



DJ INDU AVERAGE (DOW JONES & CO
as of 22-Jan-2010

**Y = ?**

**X = Feb01**

Copyright 2010 Yahoo! Inc.     http://finance.yahoo.com/
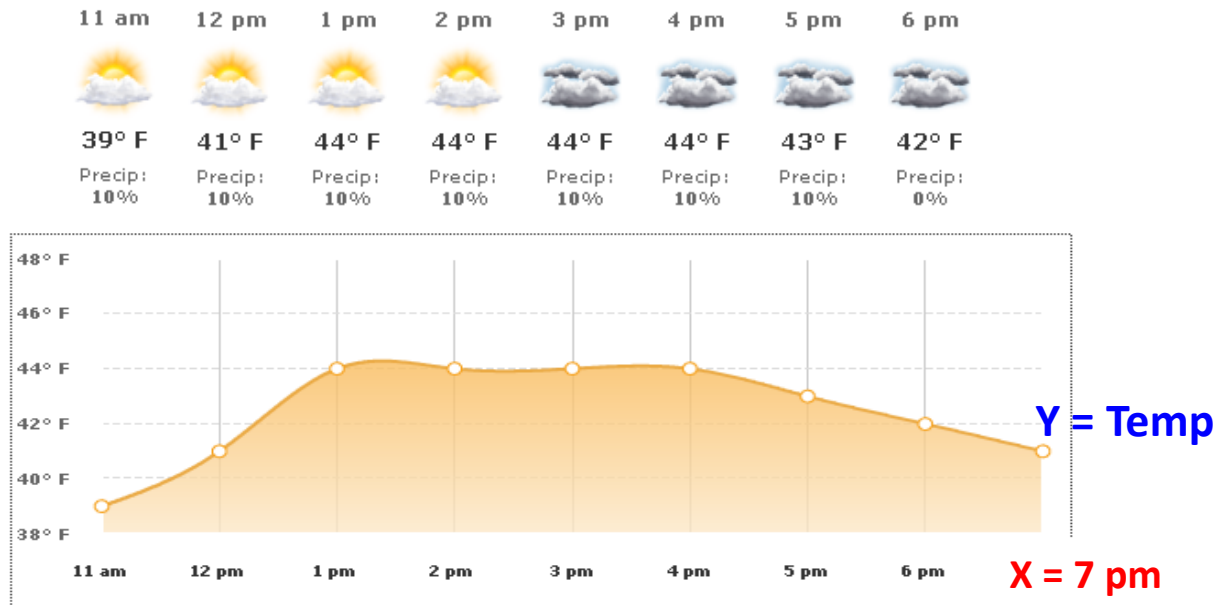
2

# Regression Tasks

Weather Prediction



Estimating Contamination



3

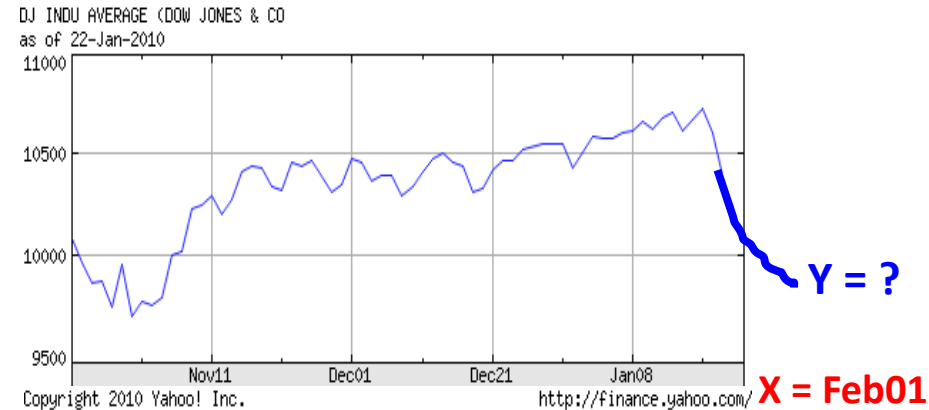# Supervised Learning

**Goal:** Construct a **predictor** $f : X \to Y$ to minimize a risk (performance measure) $R(f)$



Sports
Science
News

Y = ?

X = Feb01

**Classification:**

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

**Probability of Error**

**Regression:**

$$R(f) = \mathbb{E}[(f(X) - Y)^2]$$
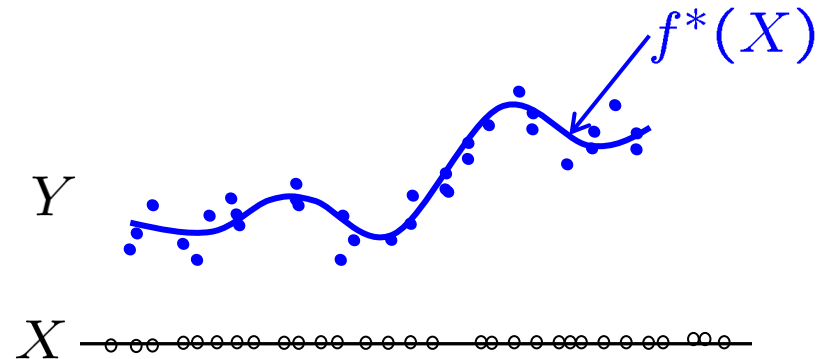
**Mean Squared Error**

4

# Regression

Optimal predictor:

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

$$= \mathbb{E}[Y|X] \qquad \text{(Conditional Mean)}$$

Intuition: Signal plus (zero-mean) Noise model

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



$f^*(X)$

$Y$

$X$

# Regression

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

Proof Strategy: $R(f) \geq R(f^*)$ for any prediction rule $f$

$$R(f) = \mathbb{E}_{XY}[(f(X) - Y)^2] = \mathbb{E}_X[\mathbb{E}_{Y|X}[(f(X) - Y)^2|X]]$$

**Dropping subscripts for notational convenience**

$$= E\left[E\left[(\underline{f(X) - E[Y|X]} + \underline{E[Y|X] - Y})^2|X\right]\right]$$

$$= E\left[\begin{array}{l} E[(f(X) - E[Y|X])^2|X] \\ +2E\left[(f(X) - E[Y|X])(E[Y|X] - Y)|X\right] \\ +E[(E[Y|X] - Y)^2|X]\end{array}\right]$$

$$= E\left[\begin{array}{l} E[(f(X) - E[Y|X])^2|X] \\ +2(f(X) - E[Y|X]) \times 0 \\ +E[(E[Y|X] - Y)^2|X]\end{array}\right]$$

$$= \underline{E\left[(f(X) - E[Y|X])^2\right]} + R(f^*).$$

**≥ 0**

# Regression

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

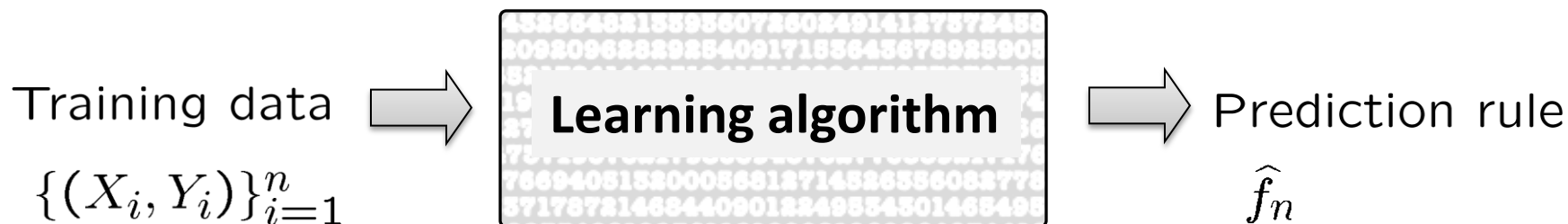$$= \mathbb{E}[Y|X] \qquad \text{(Conditional Mean)}$$

Intuition: Signal plus (zero-mean) Noise model

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

$f^*(X)$

$Y$

$X$

Depends on **unknown** distribution $P_{XY}$

# Regression algorithms

Training data

$\{(X_i, Y_i)\}_{i=1}^n$

$\Longrightarrow$

**Learning algorithm**

$\Longrightarrow$

Prediction rule

$\widehat{f_n}$

Linear Regression

Lasso, Ridge regression (Regularized Linear Regression)

Nonlinear Regression

Kernel Regression

Regression Trees, Splines, Wavelet estimators, …

# Empirical Risk Minimization (ERM)

Optimal predictor:

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

Empirical Risk 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**          **Empirical mean**

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

More later…

# ERM – you saw it before!

- Learning Distributions

  Max likelihood = Min -ve log likelihood empirical risk

$$\max_{\theta} P(D|\theta) = \min_{\theta} \frac{1}{n} \sum_{i=1}^{n} \underbrace{- \log P(X_i|\theta)}_{\text{loss}(X_i, \theta)}$$

**Negative log Likelihood loss**

What is the class $\mathcal{F}$ ?

Class of parametric distributions

Bernoulli ($\theta$)
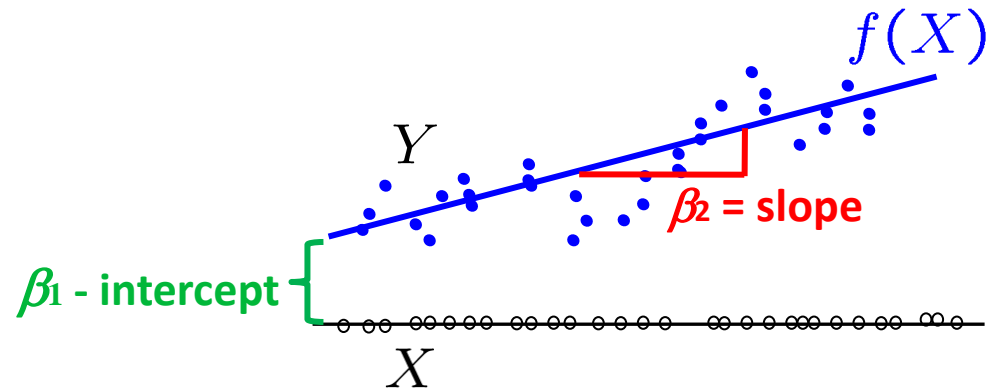
Gaussian ($\mu$, $\sigma^2$)

# Linear Regression

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

$\mathcal{F}_L$ - 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:

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

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

11

# Least Squares Estimator

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

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

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

$$\mathbf{A} = \left[\begin{array}{c} X_1 \\ \vdots \\ X_n \end{array}\right] = \left[\begin{array}{ccc} X_1^{(1)} & \ldots & X_1^{(p)} \\ \vdots & \ddots & \vdots \\ X_n^{(1)} & \ldots & X_n^{(p)} \end{array}\right] \qquad \mathbf{Y} = \left[\begin{array}{c} \mathbf{Y}_1 \\ \vdots \\ \mathbf{Y}_n \end{array}\right]$$

# Least Squares Estimator

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

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

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

# Normal Equations

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

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

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

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

When is $(\mathbf{A}^T\mathbf{A})$ invertible ?
Recall: <span style="color:red">Full rank matrices are invertible. What is rank of</span> $(\mathbf{A}^T\mathbf{A})$ <span style="color:red">?</span>

What if $(\mathbf{A}^T\mathbf{A})$ is not invertible ?
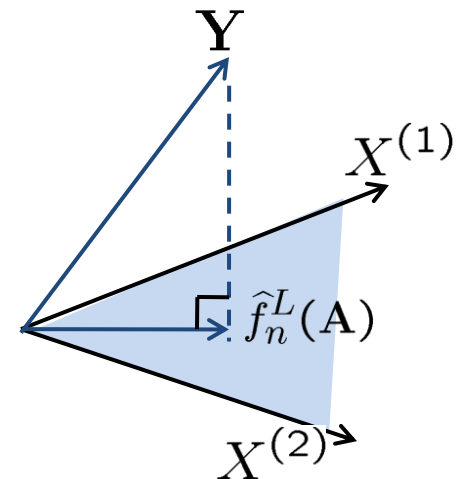<span style="color:red">Regularization (later)</span>

# Geometric Interpretation

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

Difference in prediction on training set:

$$\widehat{f}_n^L(\mathbf{A}) - \mathbf{Y} =$$

$$\mathbf{A}^T(\widehat{f}_n^L(\mathbf{A}) - \mathbf{Y}) = 0$$

$\widehat{f}_n^L(\mathbf{A})$ is the orthogonal projection of $\mathbf{Y}$ onto the linear subspace spanned by the columns of $\mathbf{A}$
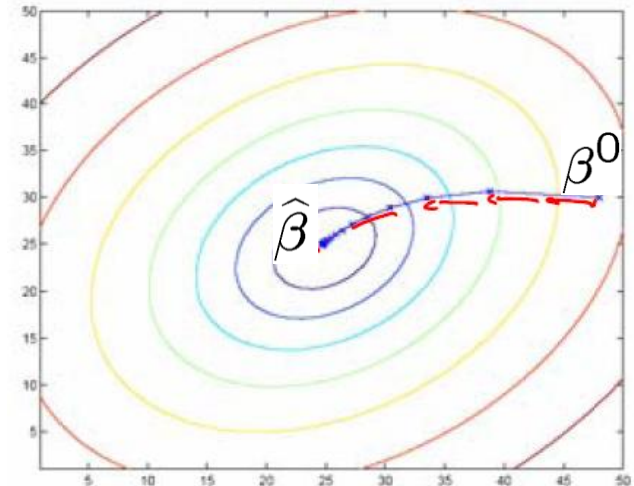
# Revisiting Gradient Descent

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

$$\widehat{\beta} = \arg\min_{\beta} \frac{1}{n}(\mathbf{A}\beta - \mathbf{Y})^T(\mathbf{A}\beta - \mathbf{Y}) = \arg\min_{\beta} J(\beta)$$
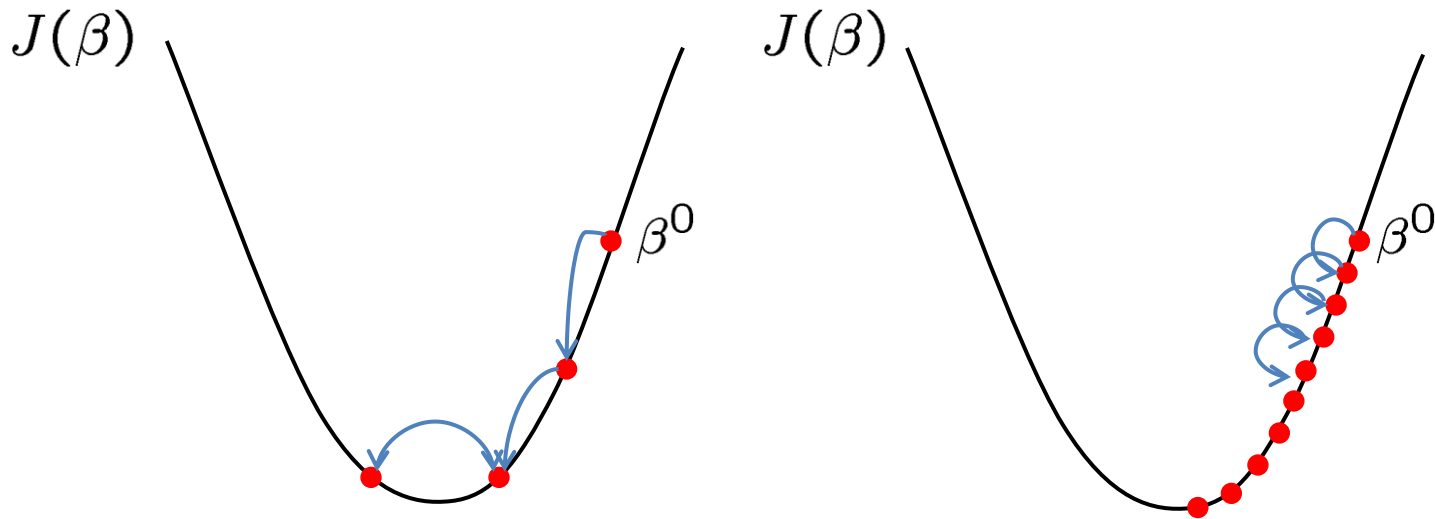
**Gradient Descent since J(β) is convex**

Initialize: $\beta^0$

Update: $\beta^{t+1} = \beta^t - \dfrac{\alpha}{2}\dfrac{\partial J(\beta)}{\partial \beta}\Big|_t$

$\phantom{\text{Update: } \beta^{t+1}} = \beta^t - \alpha\,\mathbf{A}^T\underbrace{(\mathbf{A}\beta^t - Y)}$

0 if $\beta^t = \widehat{\beta}$



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

# Effect of step-size α



Large α  => Fast convergence but larger residual error
          Also possible oscillations

Small α  => Slow convergence but small residual error

# Least Squares and MLE

Intuition: Signal plus (zero-mean) Noise model

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

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

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

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

**Least Square Estimate is same as Maximum Likelihood Estimate under a Gaussian model !**
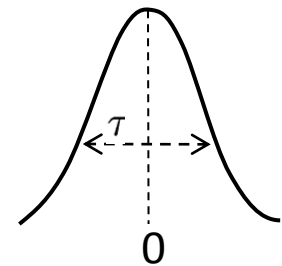
# Regularized Least Squares and MAP

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

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

I) Gaussian Prior

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

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

$$\underset{\downarrow}{\phantom{x}}$$

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

**Closed form: HW**

**Ridge Regression**

Prior belief that β is Gaussian with zero-mean biases solution to "small" β
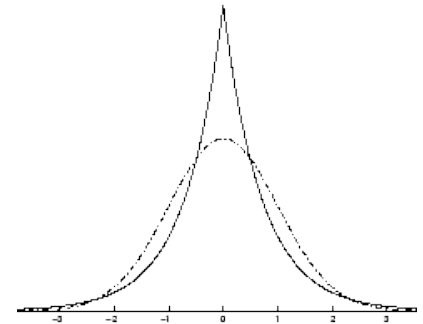
# Regularized Least Squares and MAP

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

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

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 "small" β

# Ridge Regression vs Lasso

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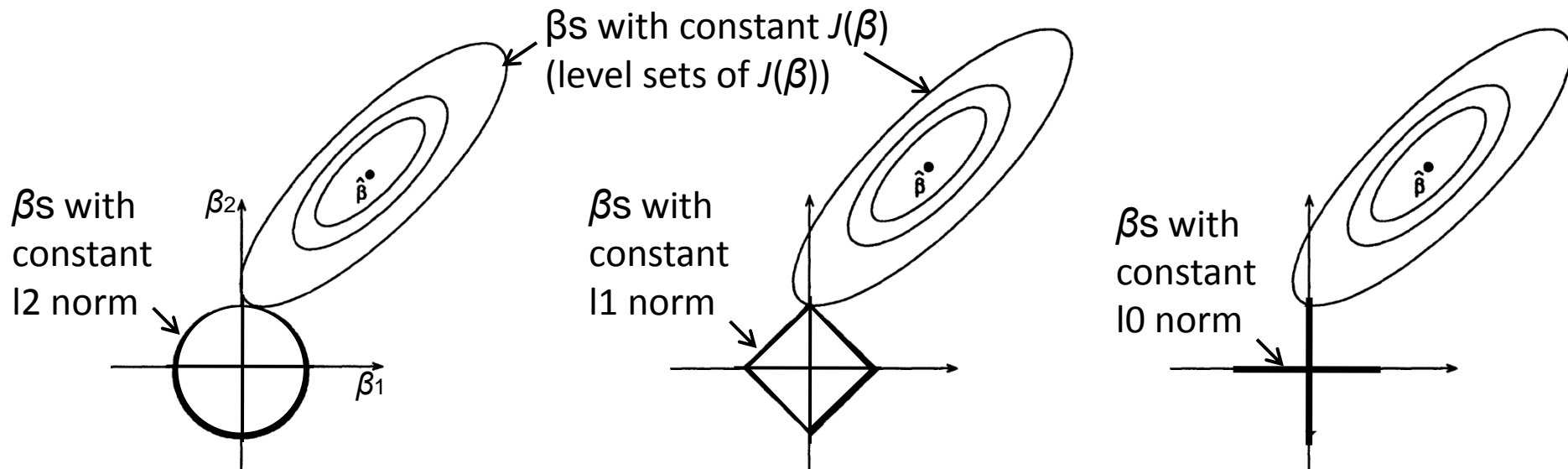
Ridge Regression:
$$\text{pen}(\beta) = \|\beta\|_2^2$$

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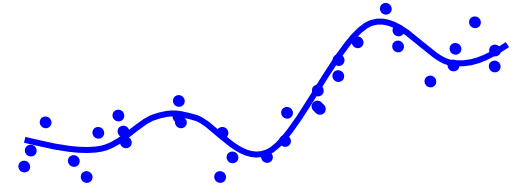$$\text{pen}(\beta) = \|\beta\|_1$$

**HOT!**

Ideally l0 penalty, but optimization becomes non-convex



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

βs with constant l2 norm

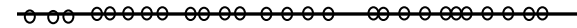βs with constant l1 norm

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!**
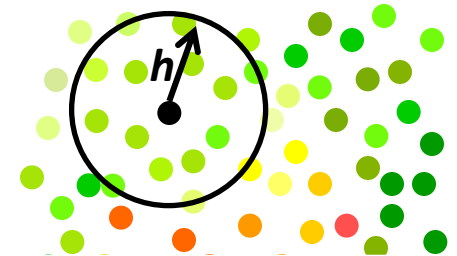
22

# Beyond Linear Regression

Polynomial regression
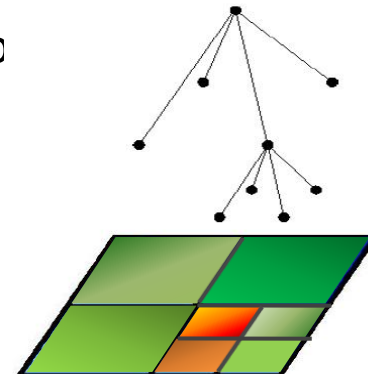
Regression with nonlinear features/basis functions

Kernel regression - Local/Weighted regression

*h*

Regression trees – Spatially adaptive regressio

# Polynomial Regression

Univariate (1-d) case:

$$f(X) = \beta_0 + \beta_1 X + \beta_2 X^2 + \cdots + \beta_m X^m = \mathbf{X}\beta$$

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

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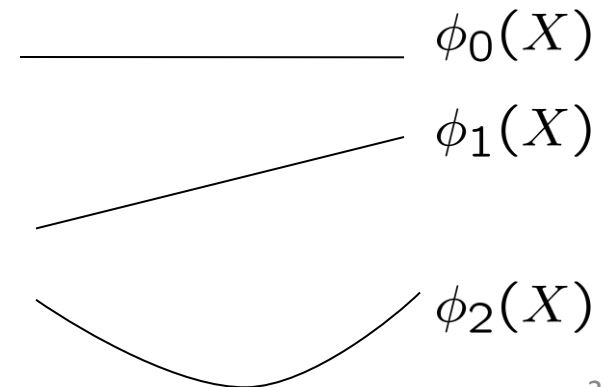
$$\widehat{f_n}(X) = \mathbf{X}\widehat{\beta}$$

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

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

Weight of each feature

Nonlinear features

$\phi_0(X)$

$\phi_1(X)$

$\phi_2(X)$

# Polynomial Regression
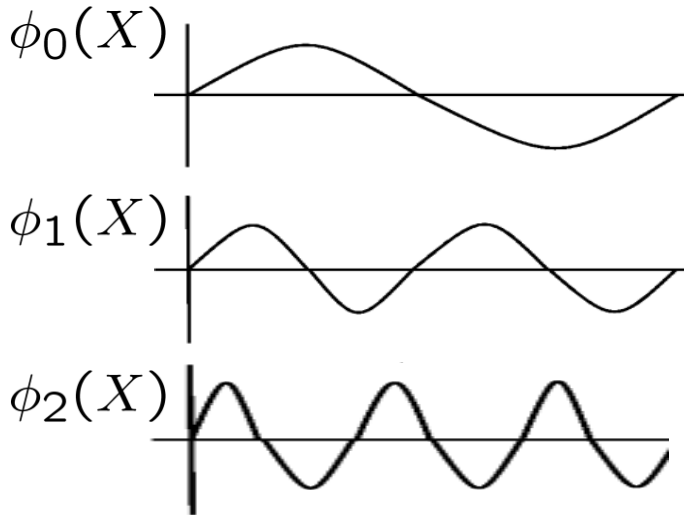
http://mste.illinois.edu/users/exner/java.f/leastsquares/

# Nonlinear Regression

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

Basis coefficients ← Nonlinear features/basis functions

Fourier Basis

$\phi_0(X)$

$\phi_1(X)$

$\phi_2(X)$

Wavelet Basis

$\phi_0(X)$

$\phi_1(X)$

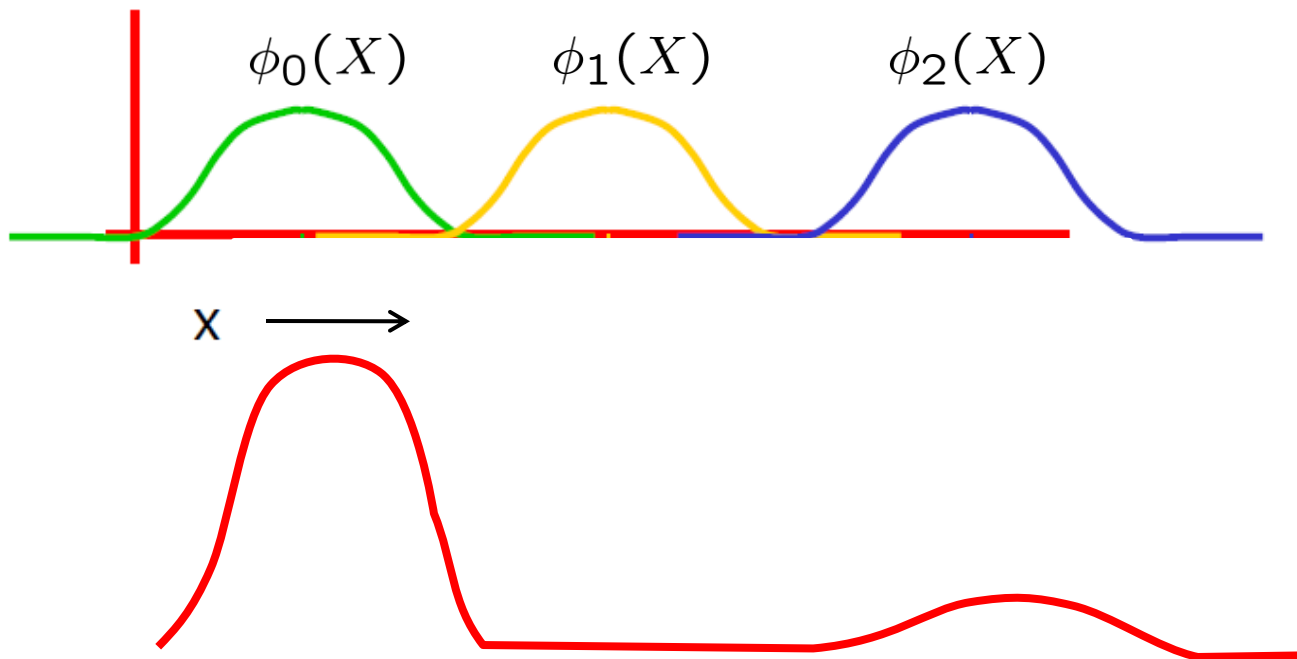$\phi_2(X)$

Good representation for oscillatory functions

Good representation for functions localized at multiple scales

# Local Regression

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

Basis coefficients ← Nonlinear features/basis functions



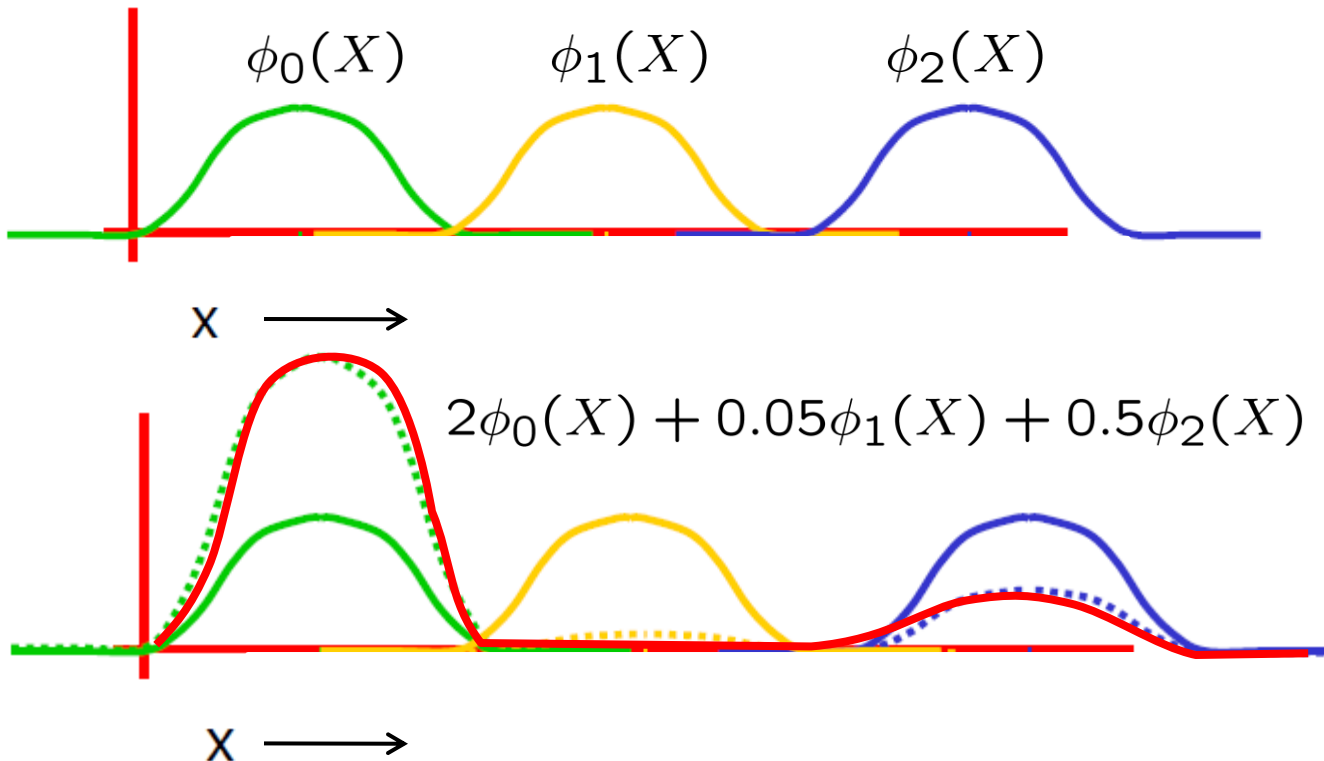$\phi_0(X)$ $\phi_1(X)$ $\phi_2(X)$

X →

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



$\phi_0(X)$      $\phi_1(X)$      $\phi_2(X)$

X $\longrightarrow$

$2\phi_0(X) + 0.05\phi_1(X) + 0.5\phi_2(X)$

X $\longrightarrow$

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