

# Linear Regression

Aarti Singh

Co-instructor: Barnabas Póczos

Machine Learning 10-401

Mar 24, 2016



**MACHINE LEARNING** DEPARTMENT



# Discrete to Continuous Labels

## Classification

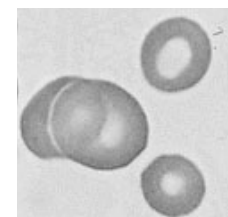
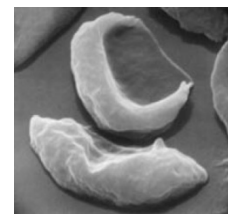


**X = Document**



Sports  
Science  
News

**Y = Topic**



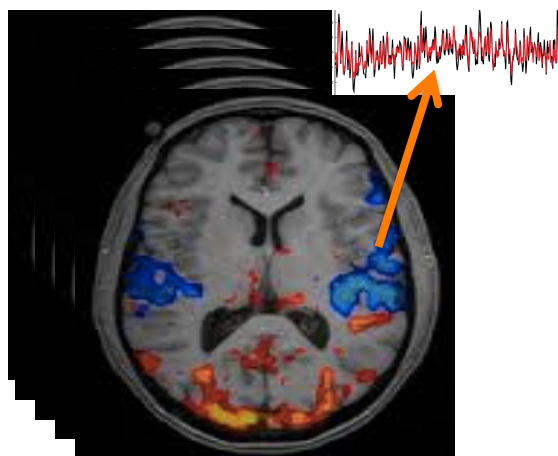
**X = Cell Image**



Anemic cell  
Healthy cell

**Y = Diagnosis**

## Regression



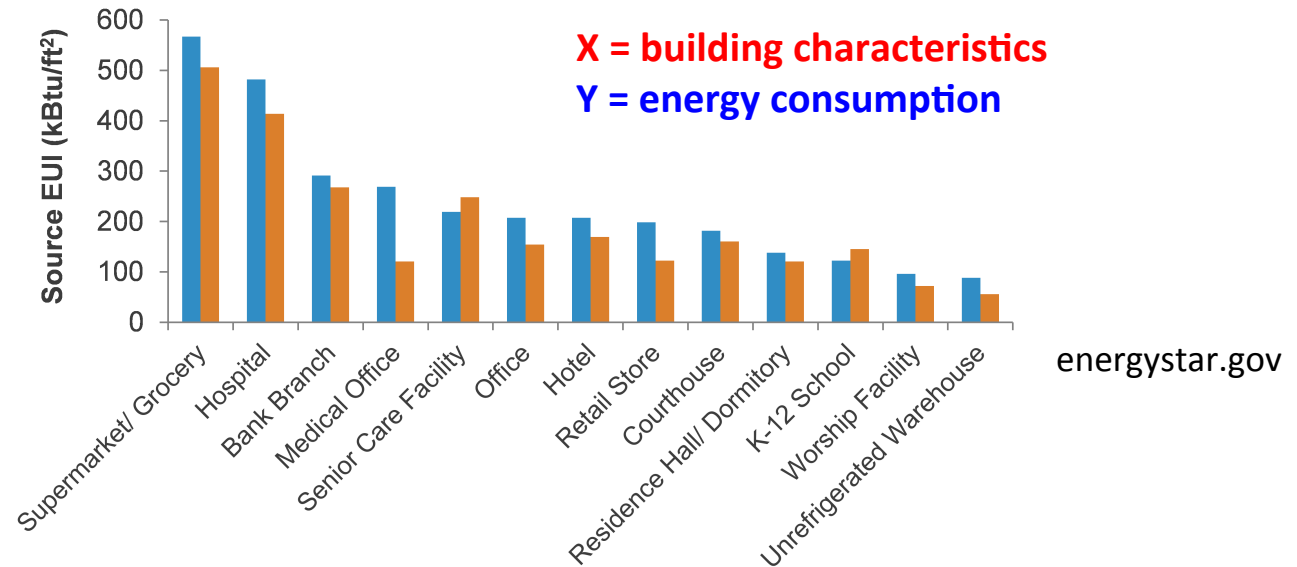
**X = Brain Scan**



**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)

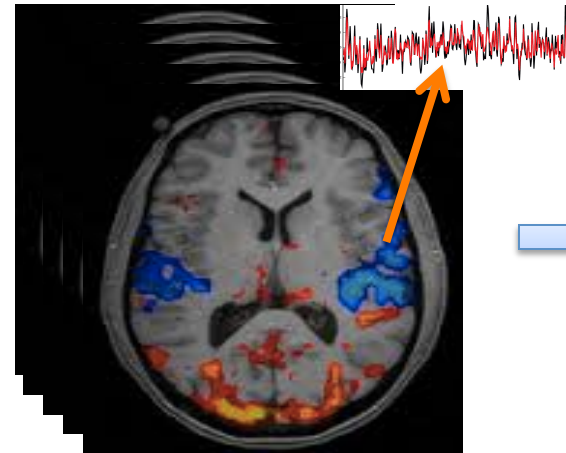


Sports  
Science  
News

**Classification:**

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

**Probability of Error**



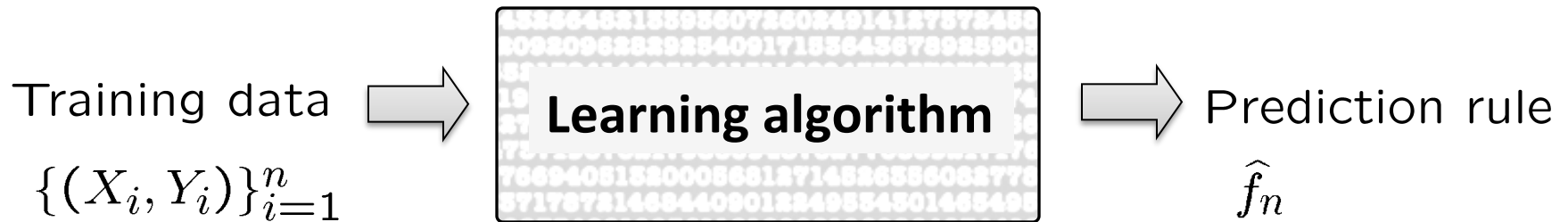
**Y = Age of  
a subject**

**Regression:**

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

**Mean Squared Error**

# Regression algorithms



Linear Regression

Regularized Linear Regression – Ridge regression, Lasso

Polynomial Regression

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

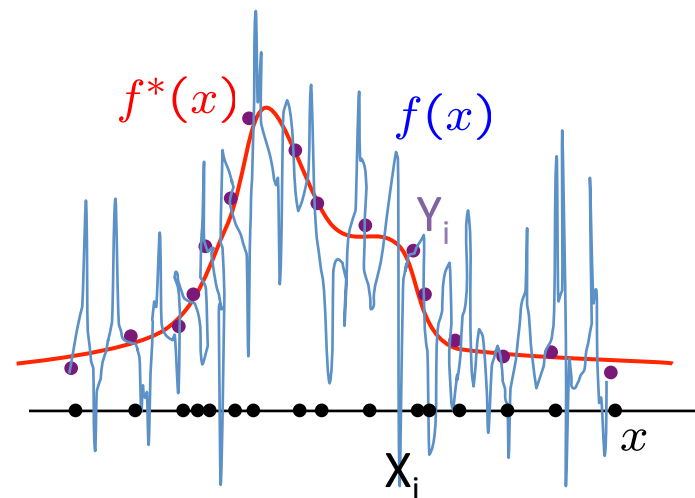
Class of predictors

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

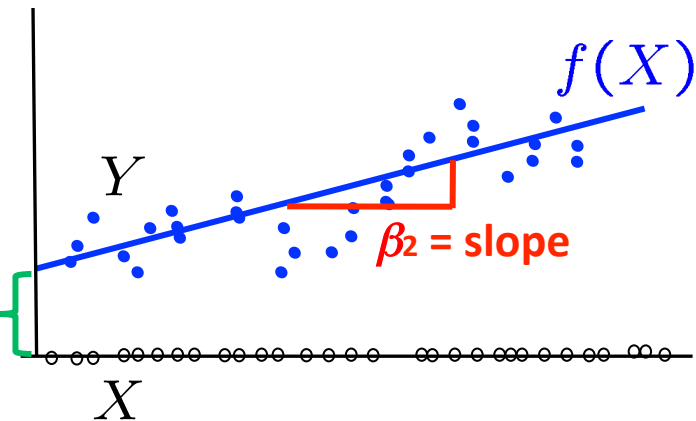
$$\hat{f}_n^L = \arg \min_{f \in \mathcal{F}_L} \frac{1}{n} \sum_{i=1}^n (f(X_i) - Y_i)^2 \quad \text{Least Squares Estimator}$$

$\mathcal{F}_L$  - 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,

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

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