

# Linear Regression

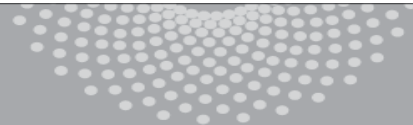
Aarti Singh

Machine Learning 10-315

Oct 23, 2019



**MACHINE LEARNING** DEPARTMENT



**Carnegie Mellon.**  
School of Computer Science

# Supervised Learning Tasks

## Classification

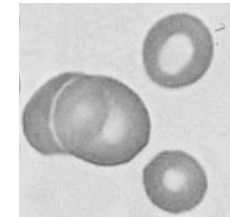
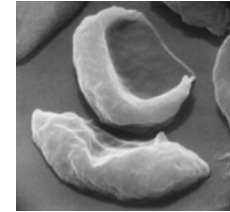


**X = Document**



Sports  
Science  
News

**Y = Topic**



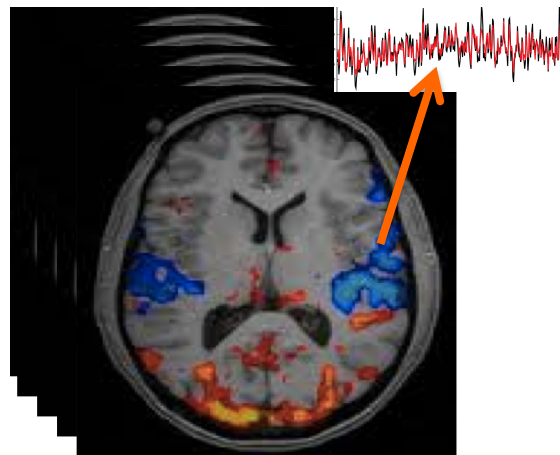
**X = Cell Image**



Anemic cell  
Healthy cell

**Y = Diagnosis**

## Regression



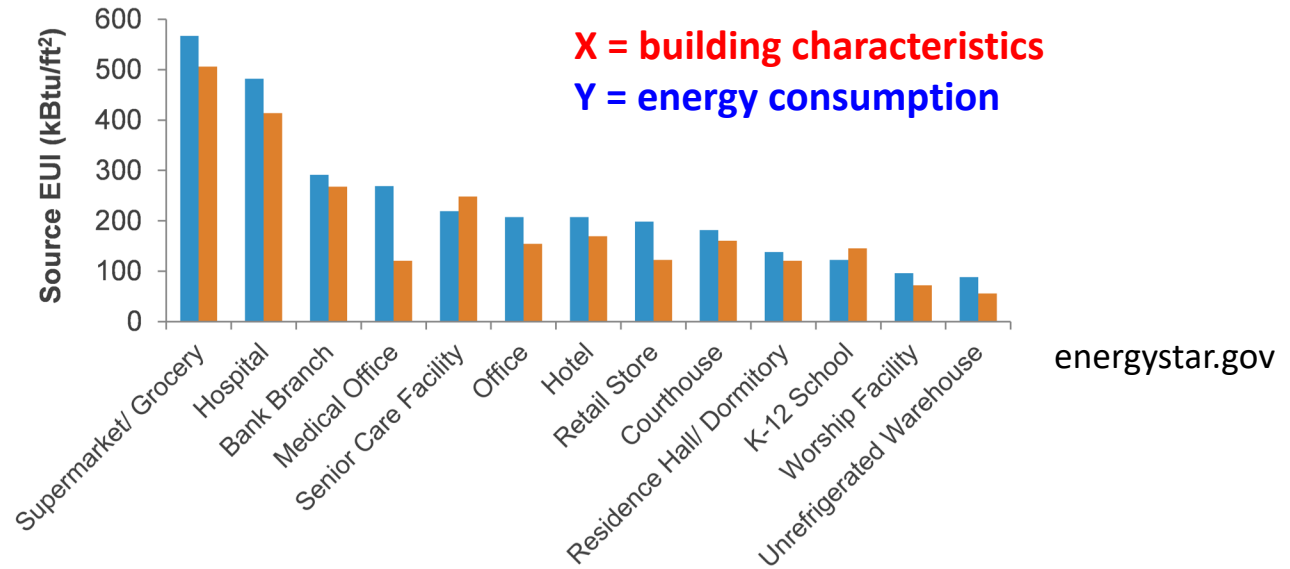
**X = Brain Scan**



**Y = Age of a subject**

# Regression Tasks

Estimating Energy Usage



Estimating Contamination



# Performance Measures

**Performance Measure:** Quantifies knowledge gained

$\text{loss}(Y, f(X))$  - Measure of closeness between true label  $Y$  and prediction  $f(X)$

Don't just want label of one test data (cell image), but any cell image  $X \in \mathcal{X}$

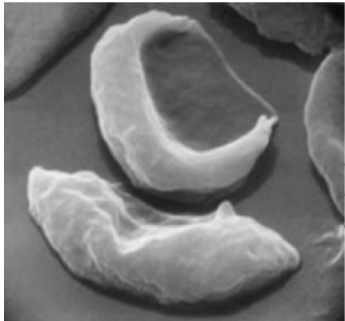
$$(X, Y) \sim P_{XY}$$

Given a cell image drawn randomly from the collection of all cell images, how well does the predictor perform on average?

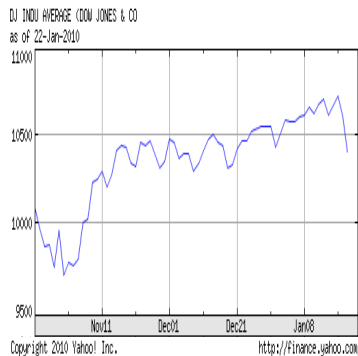
$$\text{Risk } R(f) \equiv \mathbb{E}_{XY} [\text{loss}(Y, f(X))]$$

# Performance Measures

**Performance Measure:** Risk  $R(f) \equiv \mathbb{E}_{XY} [\text{loss}(Y, f(X))]$



➔ “Anemic cell”



➔ Share Price  
“\$ 24.50”

$\text{loss}(Y, f(X))$	Risk $R(f)$
$1_{\{f(X) \neq Y\}}$ <b>0/1 loss</b>	$P(f(X) \neq Y)$ <b>Probability of Error</b>
$(f(X) - Y)^2$ <b>square loss</b>	$\mathbb{E}[(f(X) - Y)^2]$ <b>Mean Square Error</b>

# Bayes Optimal Rule

Ideal goal: Construct **prediction rule**  $f^* : \mathcal{X} \rightarrow \mathcal{Y}$

$$f^* = \arg \min_f \mathbb{E}_{XY} [\text{loss}(Y, f(X))]$$

**Bayes optimal rule**

What's the rule  
for Mean  
Square Error?

[HW3](#)

Best possible performance:

**Bayes Risk**  $R(f^*) \leq R(f)$  for all  $f$

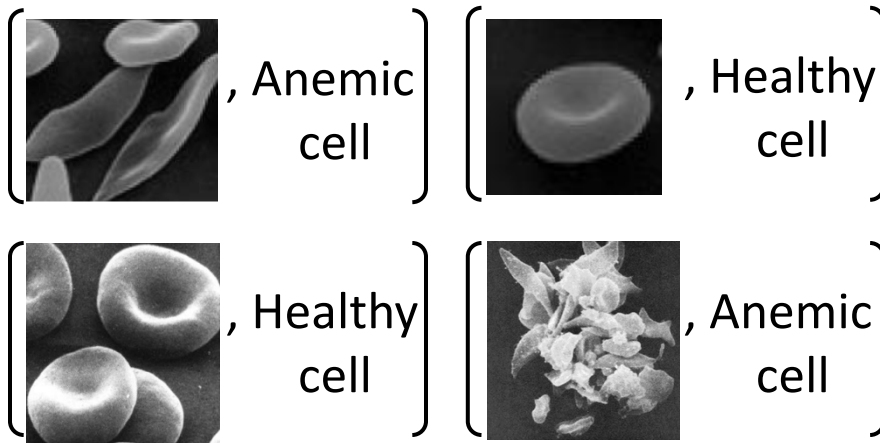
**BUT... Optimal rule is not computable - depends on unknown  $P_{XY}$  !**

# Experience - Training Data

Can't minimize risk since  $P_{XY}$  unknown!

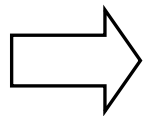
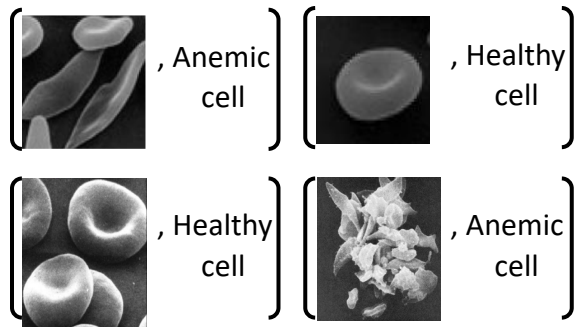
Training data (experience) provides a glimpse of  $P_{XY}$

**(observed)**  $\{(X_i, Y_i)\}_{i=1}^n \stackrel{i.i.d.}{\sim} P_{XY}$  **(unknown)**  
↓ independent, identically distributed

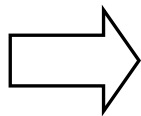


Provided by expert,  
measuring device,  
some experiment, ...

# Machine Learning Algorithm



Learning algorithm



$\hat{f}_n$

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

$\hat{f}_n$  is a mapping from  $\mathcal{X} \rightarrow \mathcal{Y}$

$\hat{f}_n \left[ \begin{array}{c} \text{Microscopic image of anemic cells} \end{array} \right] = \text{"Anemic cell"}$

Test data  $X$



# Empirical Risk Minimization

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

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

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

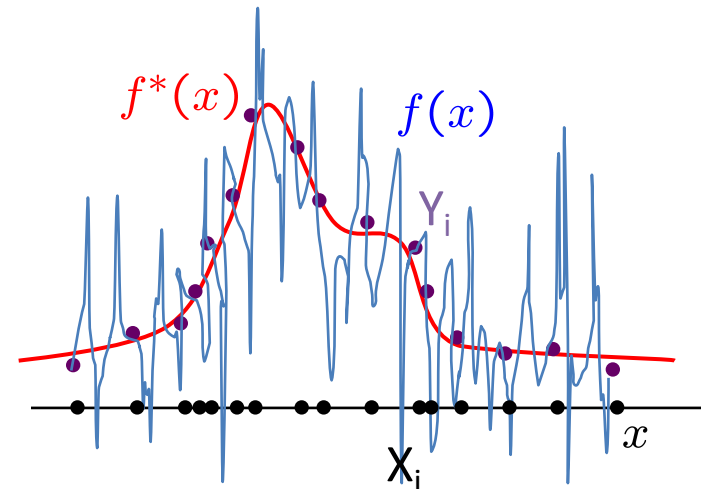
Class of predictors

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

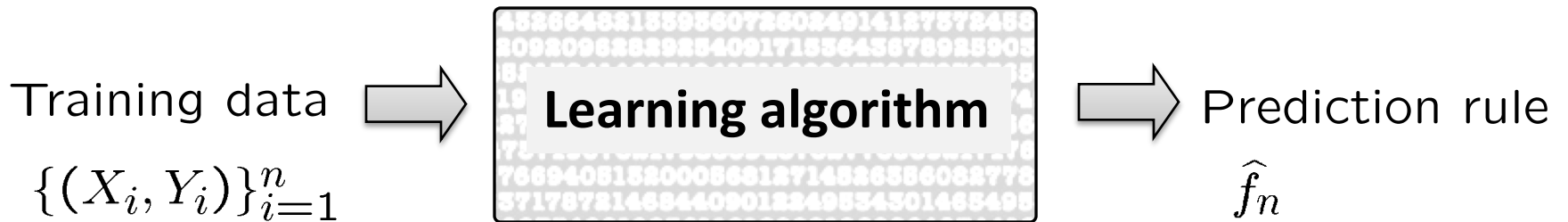
Optimal predictor:  $f^* = \arg \min_f \mathbb{E}[(f(X) - Y)^2]$

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

# 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, ...

# Linear Regression

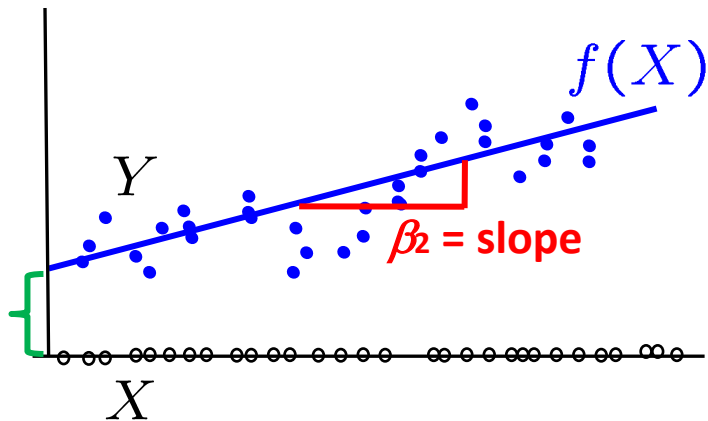
$$\hat{f}_n^L = \arg \min_{f \in \mathcal{F}_L} \frac{1}{n} \sum_{i=1}^n (f(X_i) - Y_i)^2 \quad \text{Least Squares Estimator}$$

$\mathcal{F}_L$  - 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,

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

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

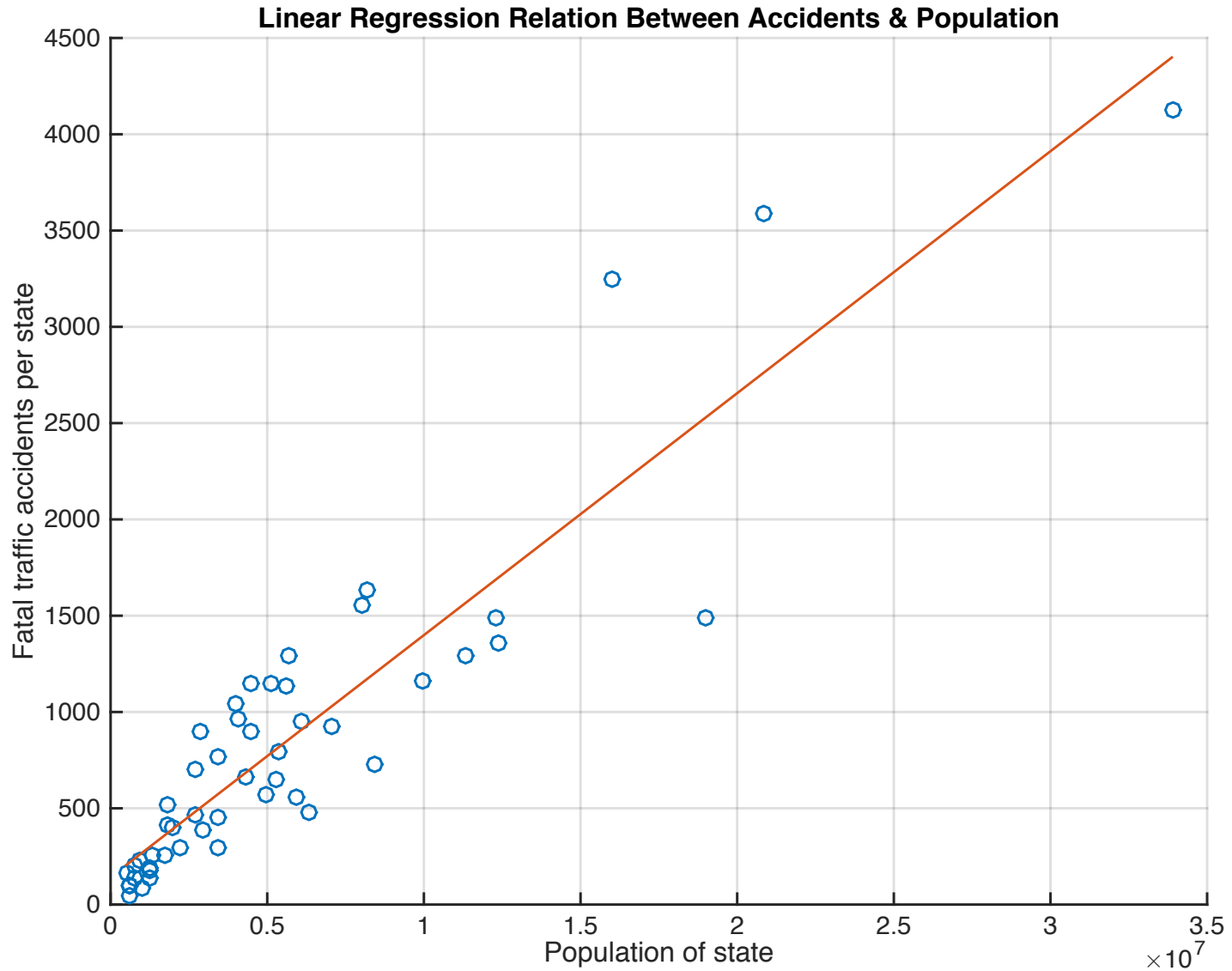


# Matlab example – linear regression

```
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,

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

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

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

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

# Gradient Descent

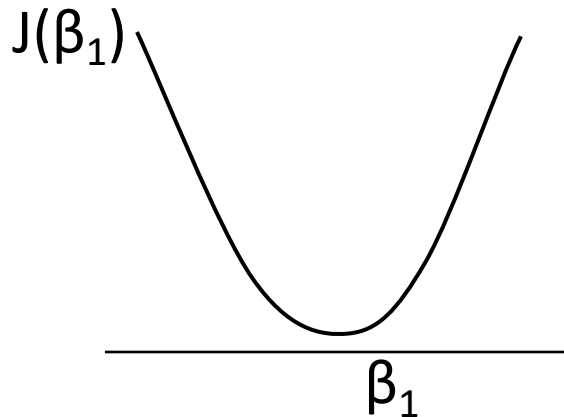
Even when  $(\mathbf{A}^T \mathbf{A})$  is invertible, might be computationally expensive if  $\mathbf{A}$  is huge.

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

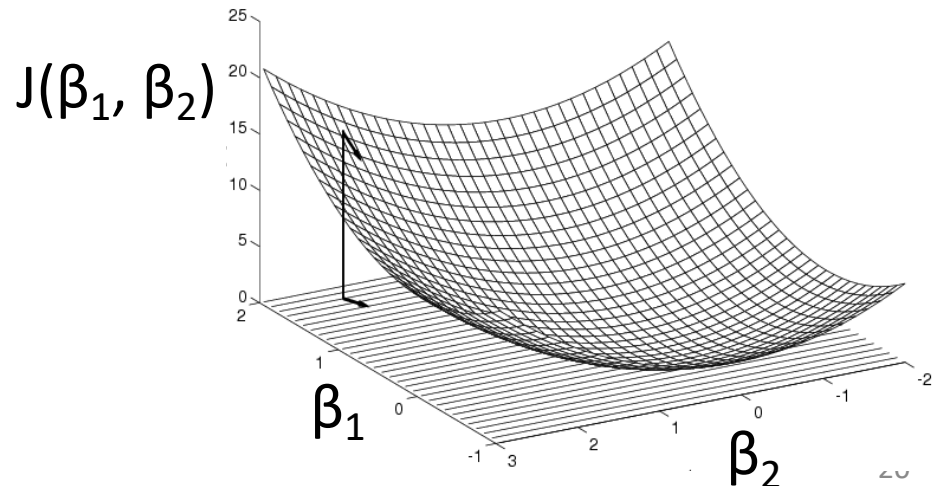
Treat as optimization problem

Why? HW3

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.

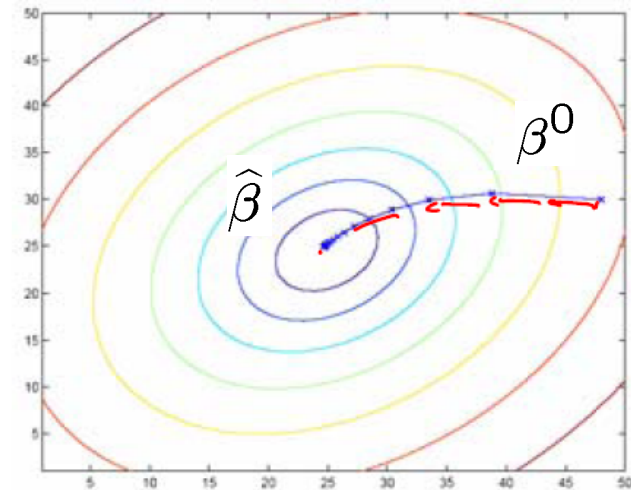
$$\hat{\beta} = \arg \min_{\beta} \frac{1}{n} (\mathbf{A}\beta - \mathbf{Y})^T (\mathbf{A}\beta - \mathbf{Y}) = \arg \min_{\beta} J(\beta)$$

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

Initialize:  $\beta^0$

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

step size



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