

Linear Regression

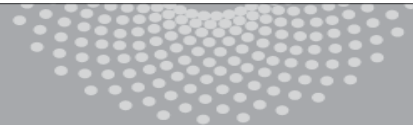
Aarti Singh

Machine Learning 10-315

Sept 23, 2020



MACHINE LEARNING DEPARTMENT



Carnegie Mellon.
School of Computer Science

Supervised Learning Tasks

Classification

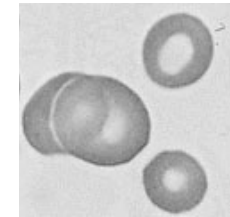
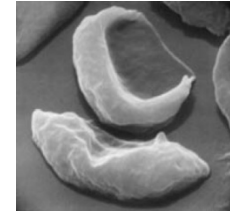


X = Document



Sports
Science
News

Y = Topic



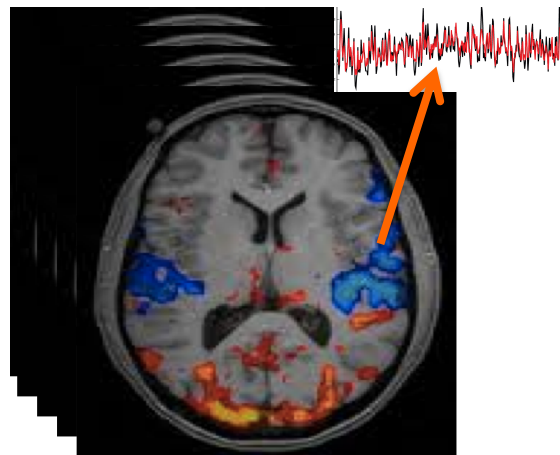
X = Cell Image



Anemic cell
Healthy cell

Y = Diagnosis

Regression



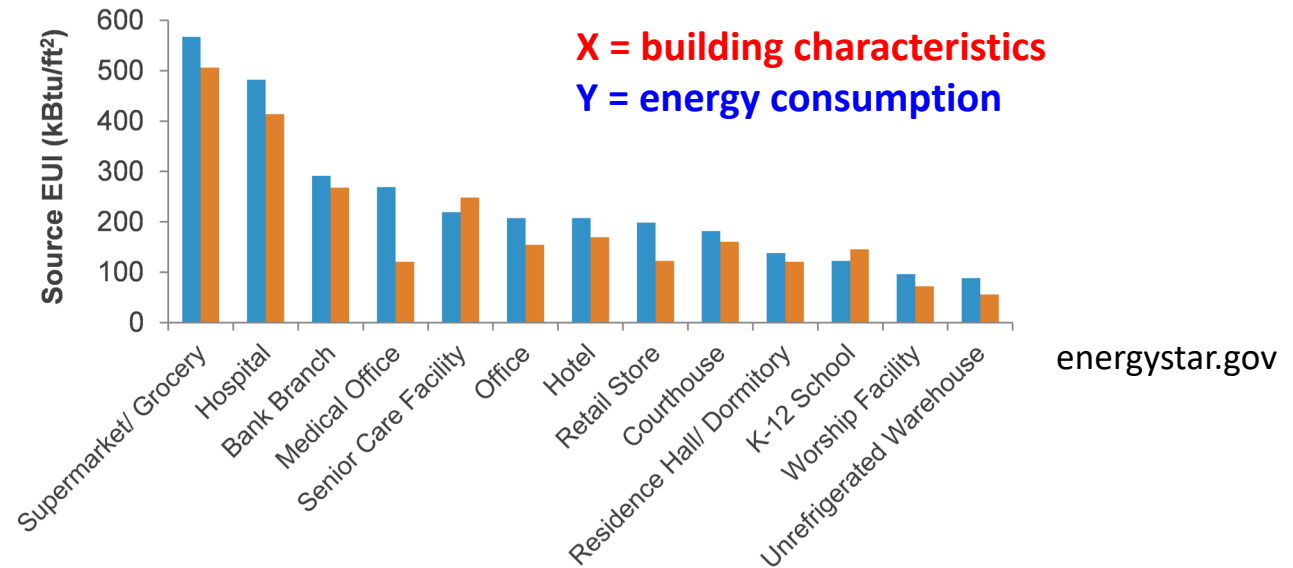
X = Brain Scan



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 (e.g. cell image), but any cell image $X \in \mathcal{X}$

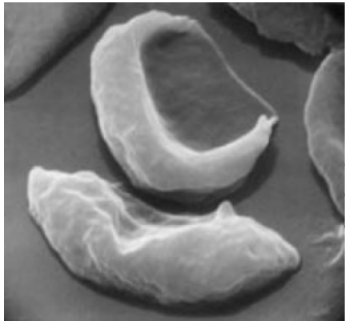
$$(X, Y) \sim P_{XY}$$

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

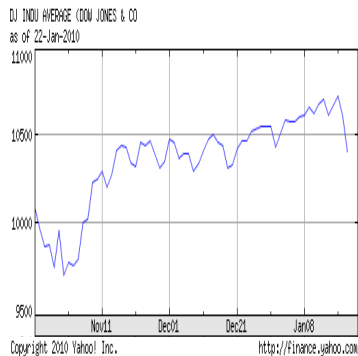
$$\text{Risk } R(f) \equiv \mathbb{E}_{XY} [\text{loss}(Y, f(X))]$$

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

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

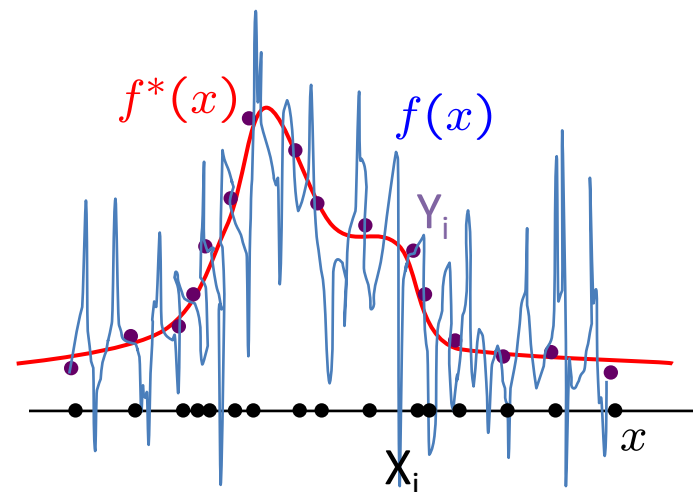
Class of predictors

➤ 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

Linear Regression

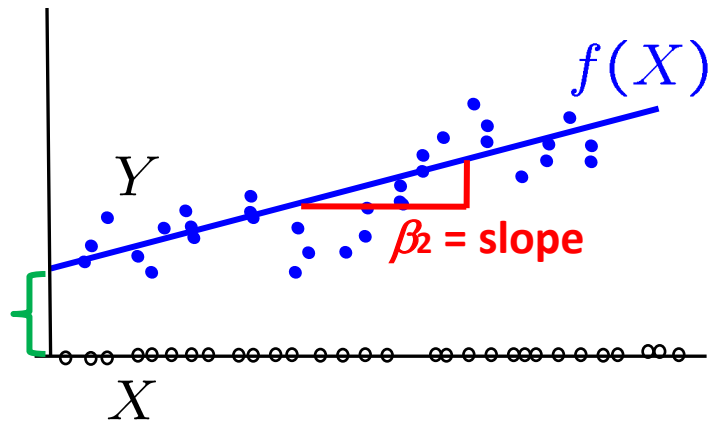
$$\hat{f}_n^L = \arg \min_{f \in \mathcal{F}_L} \frac{1}{n} \sum_{i=1}^n (f(X_i) - Y_i)^2 \quad \text{Least Squares Estimator}$$

\mathcal{F}_L - Class of Linear functions

Uni-variate case:

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

β_1 - intercept



Multi-variate case:

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

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

Linear Regression

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

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



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

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

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

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

Linear Regression

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

$$J(\beta) = (\mathbf{A}\beta - \mathbf{Y})^T (\mathbf{A}\beta - \mathbf{Y}) \quad \blacktriangleright \text{Is the objective convex?}$$

Yes!

$$= \underbrace{\beta^T \mathbf{A}^T \mathbf{A} \beta}_M - \underbrace{\beta^T \mathbf{A}^T \mathbf{Y}}_a - \underbrace{\mathbf{Y}^T \mathbf{A} \beta}_{a^T} + \mathbf{Y}^T \mathbf{Y}$$

$$\frac{\partial J(\beta)}{\partial \beta} = 2\mathbf{A}^T \mathbf{A} \beta - \mathbf{A}^T \mathbf{Y} - \mathbf{A}^T \mathbf{Y}$$

$$\frac{\partial^2 J(\beta)}{\partial \beta^2} = 2\mathbf{A}^T \mathbf{A} \geq 0$$

positive semi-definite

since $v^T (2\mathbf{A}^T \mathbf{A}) v = 2\|\mathbf{A}v\|^2 \geq 0 \quad \forall v \neq 0$

$$\frac{\partial x^T M x}{\partial x} = (M + M^T)x$$

$$\frac{\partial a^T x}{\partial x} = a = \frac{\partial x^T a}{\partial x}$$

Linear Regression

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

$$J(\beta) = (\mathbf{A}\beta - \mathbf{Y})^T (\mathbf{A}\beta - \mathbf{Y}) \quad \blacktriangleright \text{Is the objective convex?}$$

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

$$\frac{\partial J(\beta)}{\partial \beta} = 2\mathbf{A}^T\mathbf{A}\beta - 2\mathbf{A}^T\mathbf{Y}$$

$$\Rightarrow \underline{\underline{(\mathbf{A}^T\mathbf{A})\hat{\beta} = \mathbf{A}^T\mathbf{Y}}}$$

Normal equations