

Linear Regression

Maria-Florina Balcan

10/10/2018

From Discrete to Continuous Labels

Classification

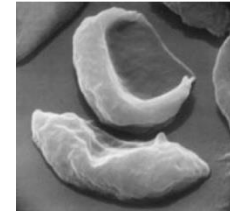


X = Document



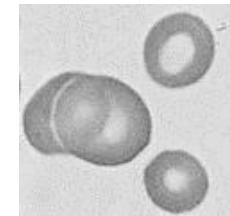
Sports
Science
News

Y = Topic



Anemic cell
Healthy cell

Y = Diagnosis

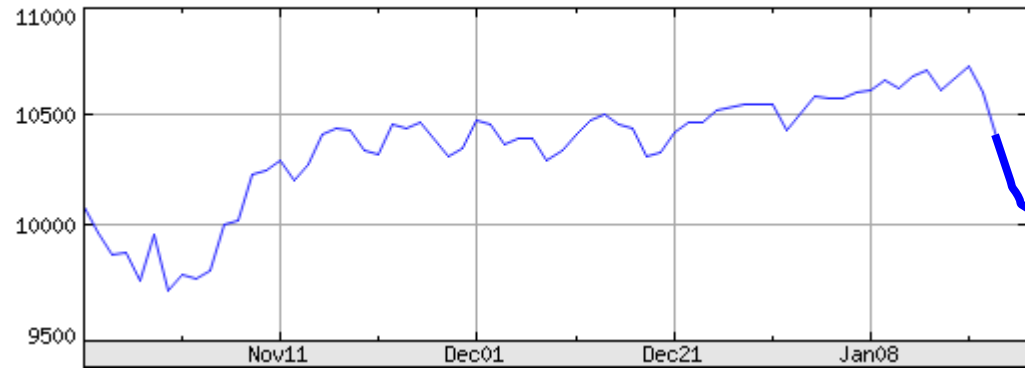


X = Cell Image

Regression

Stock Market
Prediction

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



Y = ?

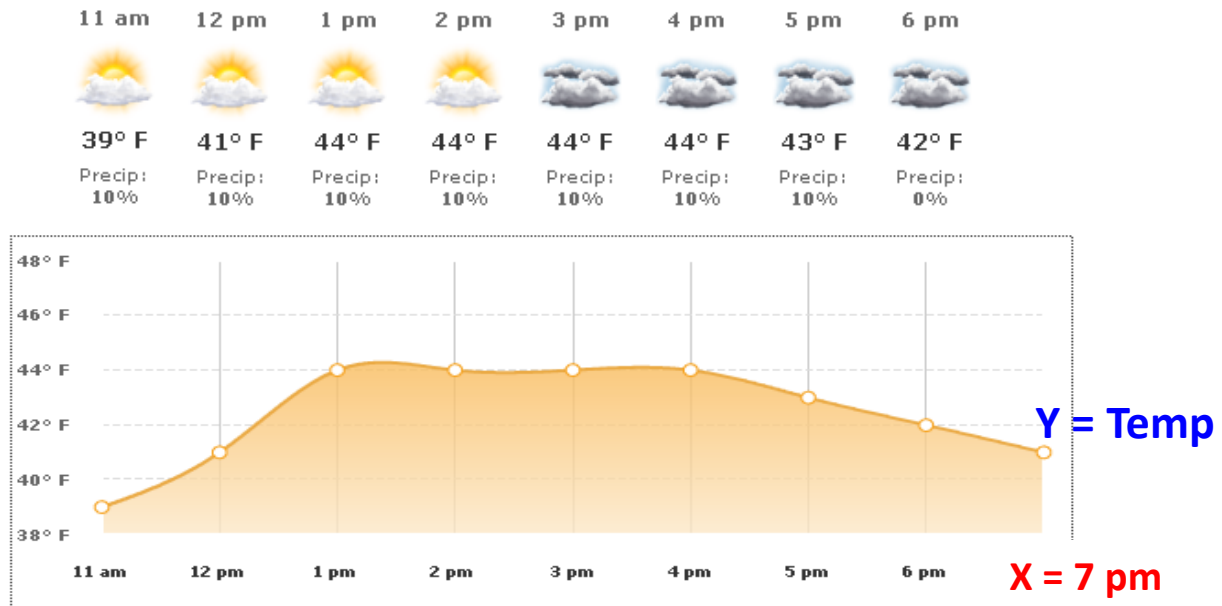
X = Feb01

Copyright 2010 Yahoo! Inc.

<http://finance.yahoo.com/>

Regression Tasks

Weather Prediction



Estimating Contamination



Supervised Learning

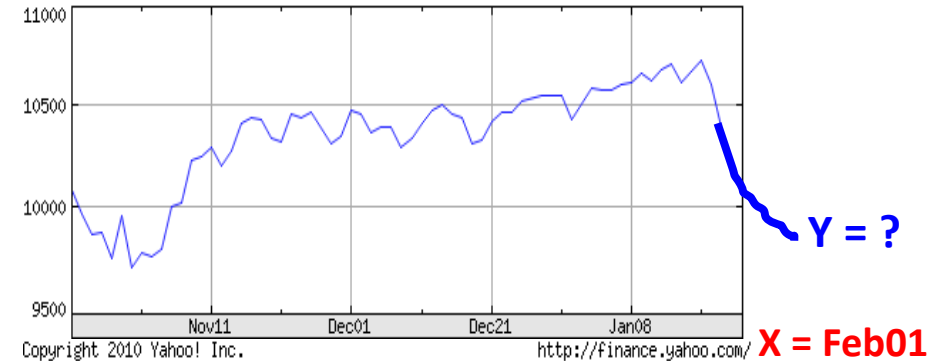
Goal: Construct a **predictor** $f: X \rightarrow Y$ to minimize a risk (error measure) $\text{err}(f)$.

Typical Error Measures



Sports
Science
News

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



Classification:

$$\text{err}(f) = P(f(X) \neq Y)$$

Probability of Error

Regression:

$$\text{err}(f) = E[(f(X) - Y)^2]$$

Mean Squared Error

Linear Regression

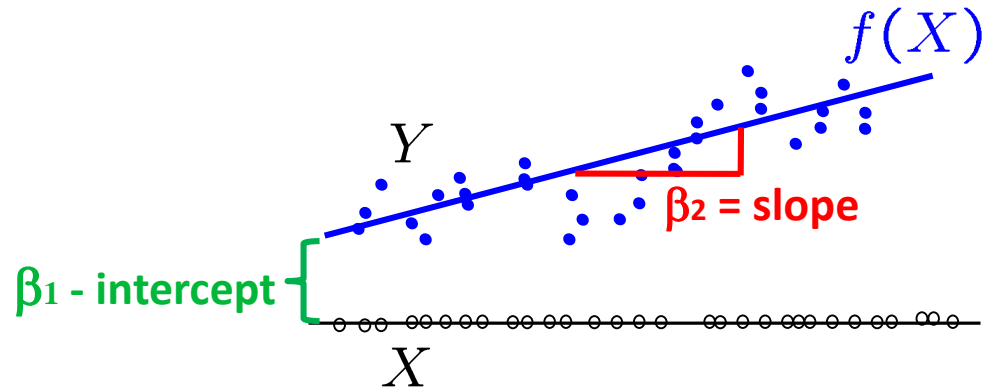
$$\hat{f}_n^L = \arg \min_{f \in F_L} \frac{1}{n} \sum_{i=1}^n (f(X_i) - Y_i)^2$$

Least Squares Estimator

F_L - Class of Linear functions

Uni-variate case:

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



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

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



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

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

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

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

Least Squares Estimator

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

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

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

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,

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

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

Geometric Interpretation

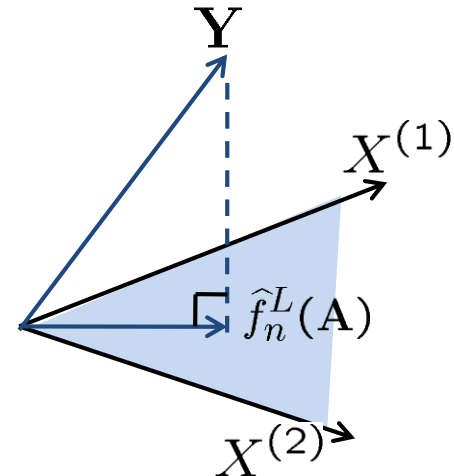
$$\hat{f}_n^L(X) = X\hat{\beta} = X(A^T A)^{-1} A^T Y$$

Difference in prediction on training set:

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

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

$\hat{f}_n^L(A)$ is the orthogonal projection of Y onto the linear subspace spanned by the columns of A .



Revisiting Gradient Descent

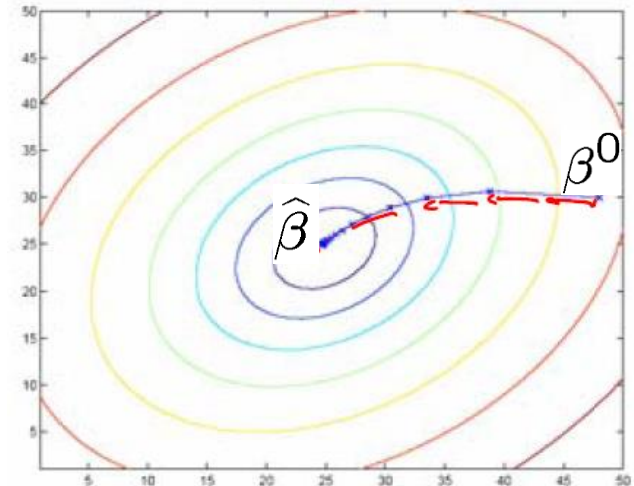
Even when $(\mathbf{A}^T \mathbf{A})$ is invertible, might be computationally expensive if \mathbf{A} is huge.

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

Gradient Descent since $J(\beta)$ is convex

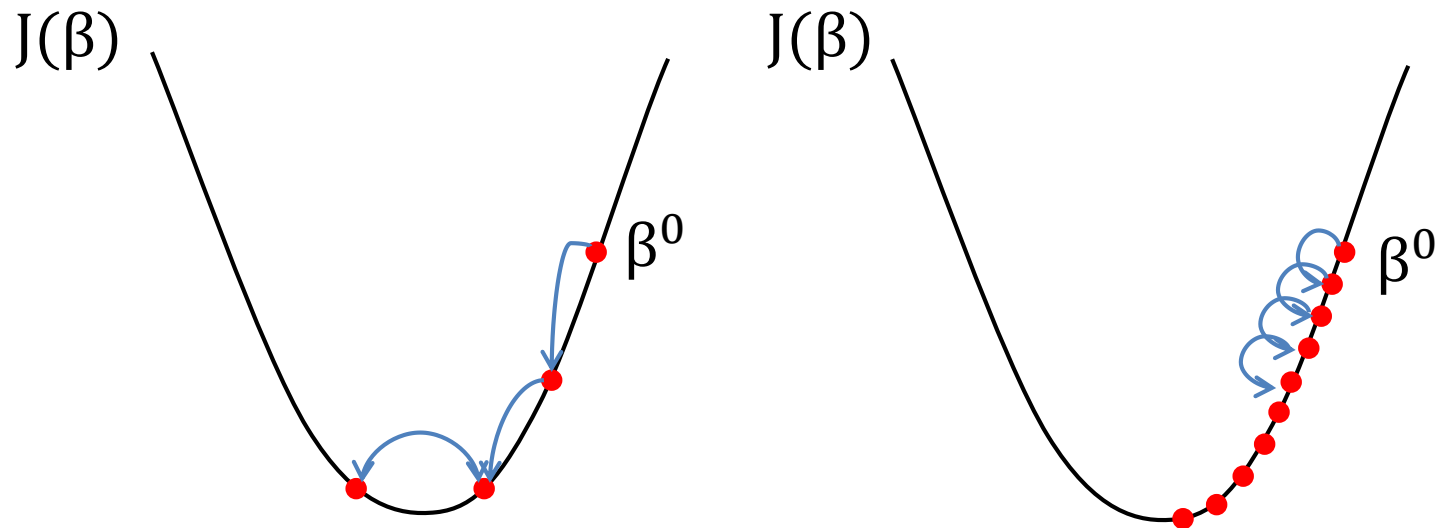
Initialize: β^0

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



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

Effect of step-size α



Large $\alpha \Rightarrow$ Fast convergence but larger residual error
Also possible oscillations

Small $\alpha \Rightarrow$ Slow convergence but small residual error

Least Squares and MLE

Intuition: Signal plus (zero-mean) Noise model

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

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

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

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

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

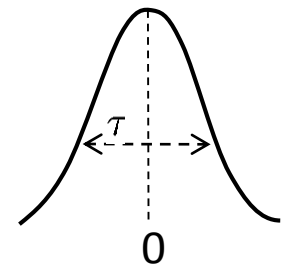
Regularized Least Squares and MAP

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

I) Gaussian Prior

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

$$\hat{\beta}_{\text{MAP}} = \arg \min_{\beta} \sum_{i=1}^n (Y_i - X_i \beta)^2 + \lambda \underbrace{\|\beta\|_2^2}_{\text{constant}(\sigma^2, \tau^2)}$$



Ridge Regression

Prior belief that β is Gaussian with zero-mean biases solution to “small” β

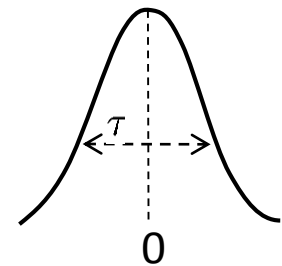
Regularized Least Squares and MAP

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

I) Gaussian Prior

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

$$\hat{\beta}_{\text{MAP}} = \arg \min_{\beta} \sum_{i=1}^n (Y_i - X_i \beta)^2 + \lambda \underbrace{\|\beta\|_2^2}_{\text{constant}(\sigma^2, \tau^2)}$$



Ridge Regression

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

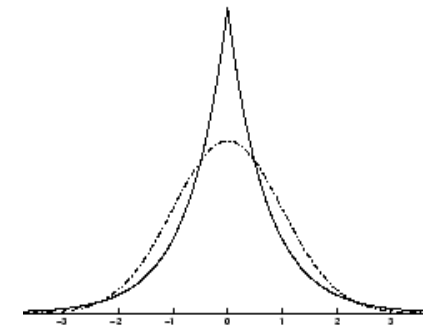
Regularized Least Squares and MAP

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

II) Laplace Prior

$$\beta_i \sim \text{Laplace}(0, t) \quad [\text{iid}] \quad p(\beta_i) \propto e^{-|\beta_i|/t}$$

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



Lasso

Prior belief that β is Laplace with zero-mean biases solution to “small” β