

Lecture 5 – Multivariate Linear Regression

Dan Sheldon

September 19, 2012

Question From Last Time

Can you find matrices A and B such that $BA = I$, but $AB \neq I$?

$$A = \begin{bmatrix} 1 & 0 \\ 0 & 1 \\ 0 & 0 \end{bmatrix}, \quad B = A^T = \begin{bmatrix} 1 & 0 & 0 \\ 0 & 1 & 0 \end{bmatrix}.$$

Check that

$$BA = \begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix} = I, \quad AB = \begin{bmatrix} 1 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & 0 & 0 \end{bmatrix}$$

Question From Last Time

Vectors \mathbf{x} and \mathbf{y} are *orthogonal* if their dot product is zero:

$$\mathbf{x}^T \mathbf{y} = \sum_{i=1}^N x_i y_i = 0.$$

E.g.,

$$\begin{bmatrix} 1 & 1 & 1 \end{bmatrix} \begin{bmatrix} 1 \\ -0.5 \\ -0.5 \end{bmatrix} = 0$$

Question From Last Time

Let $\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_k \in \mathbb{R}^n$ be *mutually orthogonal unit vectors*

$$\mathbf{x}_i^T \mathbf{x}_j = 0, \quad i \neq j$$

$$\mathbf{x}_i^T \mathbf{x}_i = 1$$

Consider the product

$$A^T A = \begin{bmatrix} - & \mathbf{x}_1^T & - \\ \vdots & & \\ - & \mathbf{x}_2^T & - \\ - & \mathbf{x}_k^T & - \end{bmatrix} \begin{bmatrix} | & | & & | \\ \mathbf{x}_1 & \mathbf{x}_2 & \dots & \mathbf{x}_k \\ | & | & & | \end{bmatrix} = I_k$$

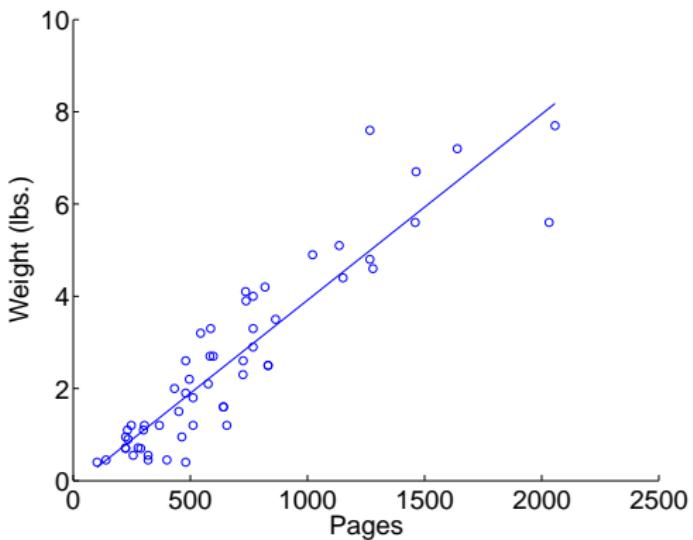
Today's Topics

Multivariate linear regression

- ▶ Model
- ▶ Cost function
- ▶ Normal equations
- ▶ Gradient descent
- ▶ Features

MATLAB (time permitting)

Book Data



$$y = 0.004036x - .122291$$

Book Data

Can we predict better with multiple features?

Width	Thickness	Height	# Pages	Hardcover	Weight
8	1.8	10	1152	1	4.4
8	0.9	9	584	1	2.7
7	1.8	9.2	738	1	3.9
6.4	1.5	9.5	512	1	1.8

Training data

$$(\mathbf{x}_1, y_1), (\mathbf{x}_2, y_2), \dots, (\mathbf{x}_N, y_N)$$

Multivariate Linear Regression

- ▶ Input: $\mathbf{x} \in \mathbb{R}^d$
- ▶ Output: $y \in \mathbb{R}$
- ▶ Model (hypothesis class): ?
- ▶ Cost function: ?

Model

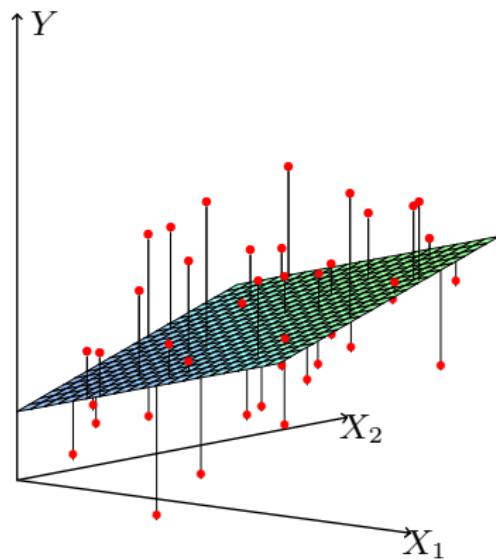
$$h_{\mathbf{w}}(\mathbf{x}) = w_0 + w_1x_1 + w_2x_2 + \dots + w_dx_d$$

$$h_{\mathbf{w}}(\mathbf{x}) = [w_0 \quad w_1 \quad \dots \quad w_d] \begin{bmatrix} 1 \\ x_1 \\ \vdots \\ x_d \end{bmatrix}$$

$$h_{\mathbf{w}}(\mathbf{x}) = \mathbf{w}^T \mathbf{x} = \mathbf{x}^T \mathbf{w}$$

(Augment feature vector with 1)

Illustration



The Problem

Find \mathbf{w} such that

$$y_i \approx h_{\mathbf{w}}(\mathbf{x}_i), \quad i = 1, \dots, N$$

$$\begin{bmatrix} y_1 \\ y_2 \\ \dots \\ y_N \end{bmatrix} \approx \begin{bmatrix} 1 & x_{1,1} & x_{1,2} & \dots & x_{1,d} \\ 1 & x_{2,1} & x_{2,2} & \dots & x_{2,d} \\ \dots \\ 1 & x_{N,1} & x_{N,2} & \dots & x_{N,d} \end{bmatrix} \begin{bmatrix} w_0 \\ w_1 \\ \dots \\ w_d \end{bmatrix}$$

$$\mathbf{y} \approx X\mathbf{w}$$

Inputs

Data matrix, label vector

$$X = \begin{bmatrix} 1 & x_{1,1} & x_{1,2} & \dots & x_{1,d} \\ 1 & x_{2,1} & x_{2,2} & \dots & x_{2,d} \\ & \vdots & & & \\ 1 & x_{N,1} & x_{N,2} & \dots & x_{N,d} \end{bmatrix} \quad \mathbf{y} = \begin{bmatrix} y_1 \\ y_2 \\ \vdots \\ y_N \end{bmatrix}.$$

Width	Thickness	Height	# Pages	Hardcover	Weight
8	1.8	10	1152	1	4.4
8	0.9	9	584	1	2.7
7	1.8	9.2	738	1	3.9
6.4	1.5	9.5	512	1	1.8

Cost Function

$$J(\mathbf{w}) = \sum_{i=1}^N (h_{\mathbf{w}}(\mathbf{x}_i) - y_i)^2$$

Exercise: write this succinctly in matrix-vector notation

Cost Function

Answer:

$$J(\mathbf{w}) = (\mathbf{X}\mathbf{w} - \mathbf{y})^T(\mathbf{X}\mathbf{w} - \mathbf{y})$$

Solution 1: Normal Equations

Normal equations

$$\mathbf{w} = (X^T X)^{-1} X^T \mathbf{y}$$

Heuristic derivation:

Proper Approach

- ▶ Set all partial derivatives to zero

$$0 = \frac{\partial}{\partial w_j} J(\mathbf{w})$$

- ▶ Solve a system of $d + 1$ linear equations for w_0, \dots, w_d
- ▶ Tedious, but leads to normal equations

Aside: Matrix Calculus

Succinct (and cool!) way to solve for normal equations:

$$0 = \frac{d}{d\mathbf{w}} (\mathbf{X}\mathbf{w} - \mathbf{y})^T (\mathbf{X}\mathbf{w} - \mathbf{y})$$

$$0 = 2(\mathbf{X}\mathbf{w} - \mathbf{y})^T \mathbf{X}$$

$$0 = \mathbf{X}^T (\mathbf{X}\mathbf{w} - \mathbf{y})$$

$$\mathbf{X}^T \mathbf{X}\mathbf{w} = \mathbf{X}^T \mathbf{y}$$

$$\mathbf{w} = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{y}$$

Solution 2: Gradient Descent

1. Initialize w_0, w_1, \dots, w_d arbitrarily
2. Repeat until convergence

$$w_j = w_j - \alpha \frac{\partial}{\partial w_j} J(\mathbf{w}), \quad j = 0, \dots, d.$$

Partial derivatives:

$$\frac{\partial}{\partial w_j} J(\mathbf{w}) = 2 \sum_{i=1}^N (h_{\mathbf{w}}(\mathbf{x}_i) - y_i) x_{i,j}$$

Feature Normalization

- ▶ Features may have very different numeric ranges

Width	Thickness	Height	# Pages	Hardcover	Weight
8	1.8	10	1152	1	4.4
8	0.9	9	584	1	2.7
7	1.8	9.2	738	1	3.9
6.4	1.5	9.5	512	1	1.8

- ▶ Advice: normalize your features!
 - ▶ Subtract mean (center)
 - ▶ Divide by standard deviation (scale)

Feature Normalization

For each feature j , compute

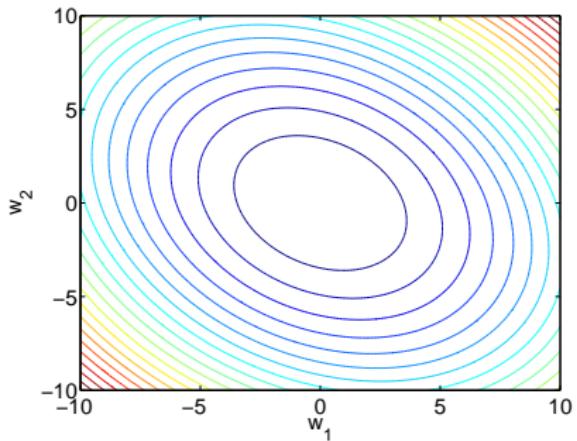
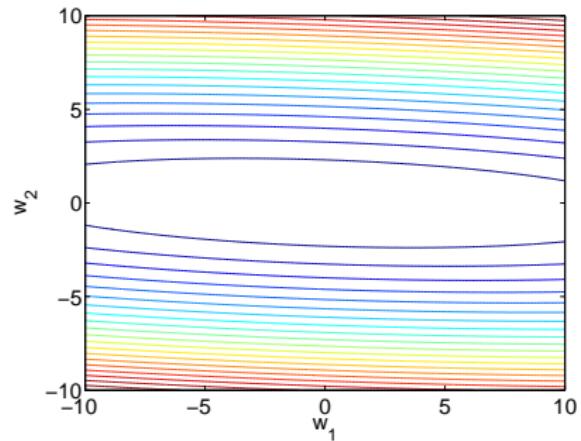
$$\mu_j = \frac{1}{N} \sum_{i=1}^N x_{i,j}, \quad \sigma_j^2 = \frac{1}{N} \sum_{i=1}^N (x_{i,j} - \mu_j)^2$$

Then, subtract μ_j and divide by σ_j :

$$x_{i,j} \leftarrow (x_{i,j} - \mu_j) / \sigma_j$$

Feature Normalization

Example: cost function contours before and after normalization



Feature Design

It is possible to fit *nonlinear* functions using linear regression:

$$(x_1, x_2, x_3) \mapsto (x_1, x_2, x_3, x_1^2, \log(x_2), x_1 + x_3)$$

Approaches

- ▶ Try standard transformations
- ▶ Design features you think will work

Polynomial Regression

$$x \mapsto (1, x, x^2, x^3, \dots)$$

