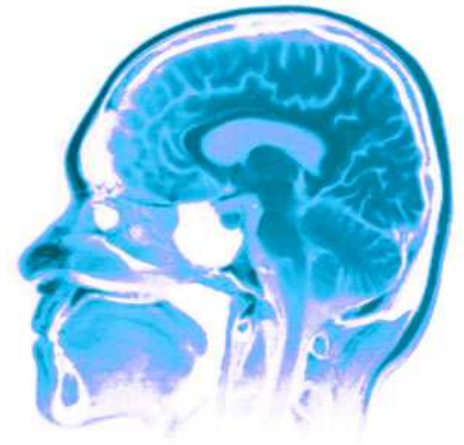




CPSC340



Linear prediction



Nando de Freitas

October, 2012

University of British Columbia

Outline of the lecture

This lecture introduces us to the topic of **supervised learning**. Here the data consists of **input-output** pairs. Inputs are also often referred to as **covariates**, **predictors** and **features**; while outputs are known as **variates** and **labels**. The goal of the lecture is for you to:

- ☐ Understand the supervised learning setting.
- ☐ Understand linear regression (aka **least squares**)
- ☐ Understand how to apply linear regression models to make predictions.
- ☐ Learn to derive the least squares estimate by optimization.

Linear supervised learning



- ❑ Many real processes can be **approximated** with linear models.
- ❑ Linear regression often appears as a **module** of larger systems.
- ❑ Linear problems can be solved **analytically**.
- ❑ Linear prediction provides an introduction to many of the **core concepts** of machine learning.

We are given a training dataset of n instances of input-output pairs $\{\mathbf{x}_{1:n}, \mathbf{y}_{1:n}\}$. Each input $\mathbf{x}_i \in \mathbb{R}^{1 \times d}$ is a vector with d attributes. The inputs are also known as predictors or covariates. The output, often referred to as the target, will be assumed to be univariate, $\mathbf{y}_i \in \mathbb{R}$, for now.



A typical dataset with $n = 4$ instances and 2 attributes would look like the following table:

Wind speed	People inside building	Energy requirement
100	2	5
50	42	25
45	31	22
60	35	18

← th. setting

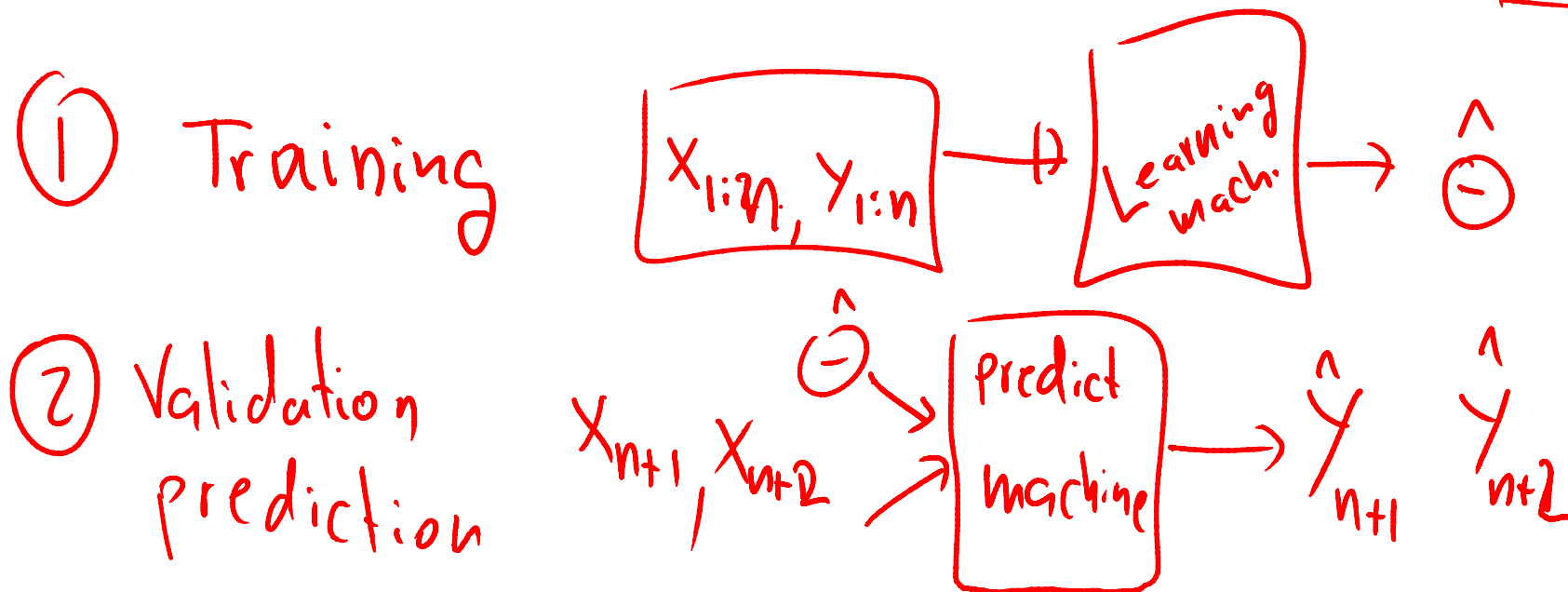
day 1 $x_1 = (100, 2)$ $y_1 = 5$

$y(x)$

Energy demand prediction



Given the training set $\{\mathbf{x}_{1:n}, \mathbf{y}_{1:n}\}$, we would like to learn a model of how the inputs affect the outputs. Given this model and a new value of the input \mathbf{x}_{n+1} , we can use the model to make a prediction $\hat{y}(\mathbf{x}_{n+1})$.



Prostate cancer example

❑ **Goal:** Predict a prostate-specific antigen (log of lpsa) from a number of clinical measures in men who are about to receive a radical prostatectomy.



❑ The **inputs** are:

- Log cancer volume (lcavol) ✓
- Log prostate weight (lweight) ✓
- Age
- Log of the amount of benign prostatic hyperplasia (lbph)
- Seminal vesicle invasion (svi) - *binary*
- Log of capsular penetration (lcp)
- Gleason score (gleason) – *ordered categorical*
- Percent of Gleason scores 4 or 5 (pgg45)

Which inputs are more important?

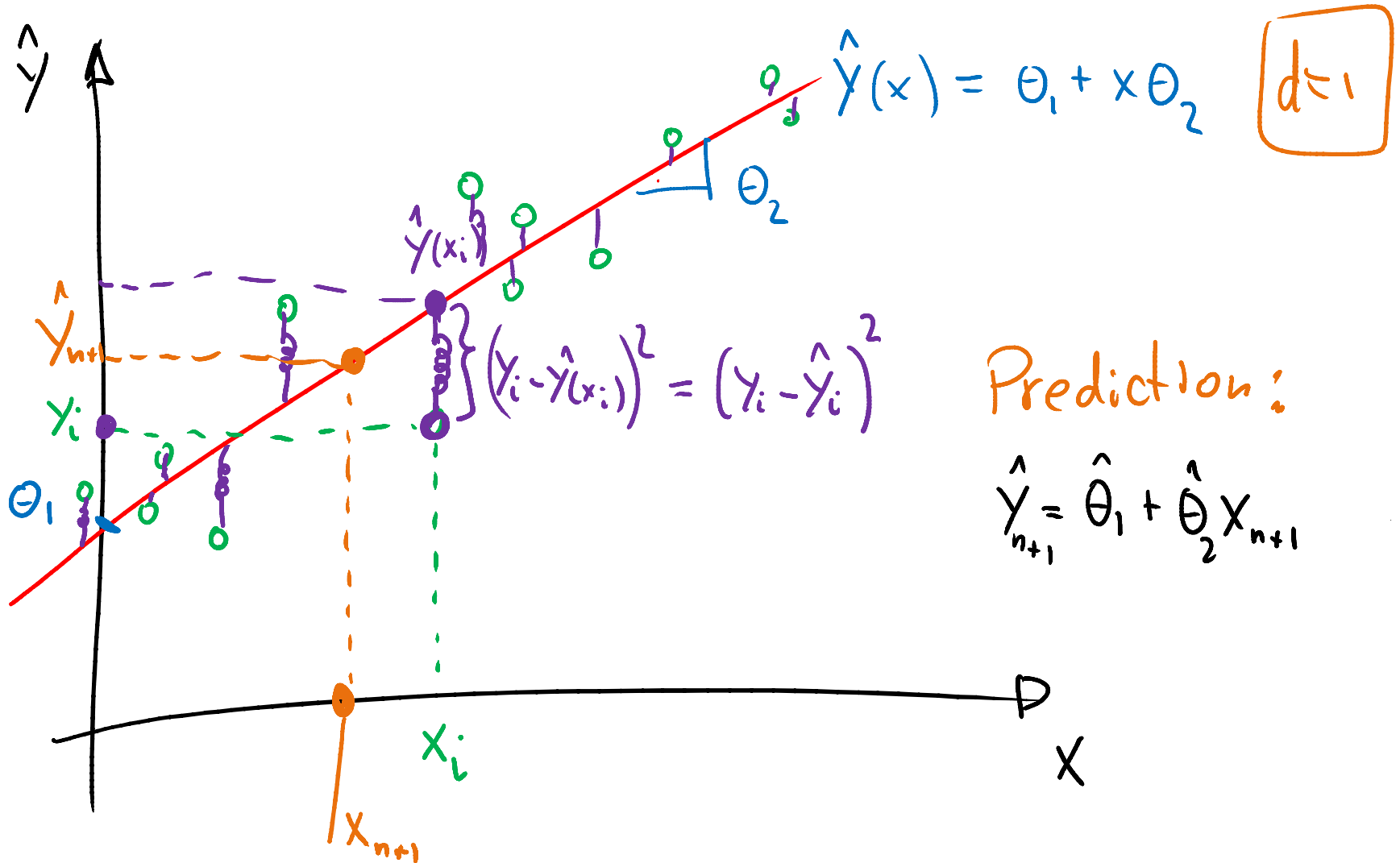
website
✓

[Hastie, Tibshirani & Friedman book]

$$\hat{y}(\mathbf{x}_i) = \theta_1 + x_i \theta_2$$

$$\underline{J(\theta)} = \sum_{i=1}^n \underbrace{(y_i - \hat{y}_i)^2}_{\text{residual squared}} = \sum_{i=1}^n (y_i - \theta_1 - x_i \theta_2)^2$$

Goal: solve for θ_1 and θ_2



x_{i2} = height
 x_{i2} = weight

Linear prediction

$$y = \theta_2 + \theta_1 x_2$$
$$y = \sum_{j=1}^d \theta_j x_j, x_i = 1$$

In general, the linear model is expressed as follows:

$$y = 1\theta_1 + x_{i2}\theta_2 + x_{i3}\theta_3$$

$$\hat{y}_i = \sum_{j=1}^d x_{ij}\theta_j = x_{i1}\theta_1 + x_{i2}\theta_2 + \dots + x_{id}\theta_d$$

where we have assumed that $x_{i1} = 1$ so that θ_1 corresponds to the intercept of the line with the vertical axis. θ_1 is known as the bias or offset.

In matrix form, the expression for the linear model is:

$$\hat{\mathbf{y}} = \mathbf{X}\boldsymbol{\theta}$$

with $\hat{\mathbf{y}} \in \mathbb{R}^{n \times 1}$, $\mathbf{X} \in \mathbb{R}^{n \times d}$ and $\boldsymbol{\theta} \in \mathbb{R}^{d \times 1}$. That is,

$$\begin{bmatrix} \hat{y}_1 \\ \vdots \\ \hat{y}_n \end{bmatrix} = \begin{bmatrix} x_{11} & \cdots & x_{1d} \\ \vdots & \vdots & \vdots \\ x_{n1} & \cdots & x_{nd} \end{bmatrix} \begin{bmatrix} \theta_1 \\ \vdots \\ \theta_d \end{bmatrix}.$$

Wind speed	People inside building	Energy requirement
100	2	5
50	42	25
45	31	22
60	35	18

For our energy prediction example, we would form the following matrices with $n = 4$ and $d = 3$:

$$\mathbf{y} = \begin{bmatrix} 5 \\ 25 \\ 22 \\ 18 \end{bmatrix}, \quad \mathbf{X} = \begin{bmatrix} 1 & 100 & 2 \\ 1 & 50 & 42 \\ 1 & 45 & 31 \\ 1 & 60 & 35 \end{bmatrix}, \quad \boldsymbol{\theta} = \begin{bmatrix} \theta_1 \\ \theta_2 \\ \theta_3 \end{bmatrix}.$$

Suppose that $\hat{\boldsymbol{\theta}} = [1 \ 0 \ 0.5]^T$. Then, by multiplying \mathbf{X} times $\hat{\boldsymbol{\theta}}$, we would get the following predictions on the training set:

$$\hat{\mathbf{y}} = \begin{bmatrix} 2 \\ 22 \\ 16.5 \\ 18.5 \end{bmatrix} = \begin{bmatrix} 1 & 100 & 2 \\ 1 & 50 & 42 \\ 1 & 45 & 31 \\ 1 & 60 & 35 \end{bmatrix} \begin{bmatrix} 1 \\ 0 \\ 0.5 \end{bmatrix}.$$

Linear prediction

Likewise, for a point that we have never seen before, say $x = [50 \ 20]$, we generate the following prediction:

$$\hat{y}(x) = [1 \ 50 \ 20] \begin{bmatrix} 1 \\ 0 \\ 0.5 \end{bmatrix} = 1 + 0 + 10 = 11.$$

Handwritten red annotations: x_{n+1} above the vector $[1 \ 50 \ 20]$, and $\hat{\theta}$ above the vector $\begin{bmatrix} 1 \\ 0 \\ 0.5 \end{bmatrix}$.

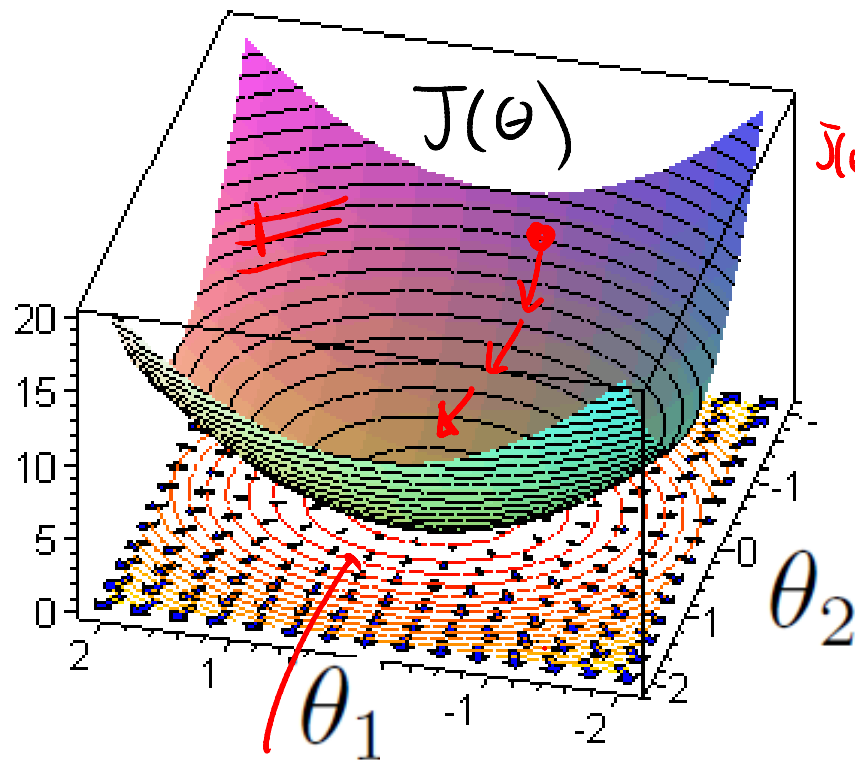
Optimization approach $\hat{y}_i = x_i^T \theta$

Our aim is to minimise the quadratic cost between the output labels and the model predictions

$$J(\theta) = (y - X\theta)^T (y - X\theta) = \sum_{i=1}^n (y_i - x_i^T \theta)^2$$

$$\hat{y}_1 = x_1^T \theta$$

$$\hat{y}_2 = x_2^T \theta$$



Contours

$$J(\theta) = \left(\begin{bmatrix} y_1 \\ y_2 \end{bmatrix} - \begin{bmatrix} x_{11} & x_{12} \\ x_{21} & x_{22} \end{bmatrix} \begin{bmatrix} \theta_1 \\ \theta_2 \end{bmatrix} \right)^T \left(\begin{bmatrix} y_1 \\ y_2 \end{bmatrix} - \begin{bmatrix} x_{11} & x_{12} \\ x_{21} & x_{22} \end{bmatrix} \begin{bmatrix} \theta_1 \\ \theta_2 \end{bmatrix} \right)$$

$$= \left(\begin{bmatrix} y_1 \\ y_2 \end{bmatrix} - \begin{bmatrix} x_1^T \theta \\ x_2^T \theta \end{bmatrix} \right)^T \left(\begin{bmatrix} y_1 \\ y_2 \end{bmatrix} - \begin{bmatrix} x_1^T \theta \\ x_2^T \theta \end{bmatrix} \right)$$

$$= \left(\begin{bmatrix} y_1 \\ y_2 \end{bmatrix} - \begin{bmatrix} \hat{y}_1 \\ \hat{y}_2 \end{bmatrix} \right)^T \left(\begin{bmatrix} y_1 \\ y_2 \end{bmatrix} - \begin{bmatrix} \hat{y}_1 \\ \hat{y}_2 \end{bmatrix} \right)$$

$$= \begin{bmatrix} (y_1 - \hat{y}_1) & (y_2 - \hat{y}_2) \end{bmatrix} \begin{bmatrix} y_1 - \hat{y}_1 \\ y_2 - \hat{y}_2 \end{bmatrix} = (y_1 - \hat{y}_1)^2 + (y_2 - \hat{y}_2)^2$$

Optimization

$$J(\theta) = \underbrace{(y - X\theta)^T}_{1 \times n} \underbrace{(y - X\theta)}_{n \times 1} = \sum_{i=1}^n \underbrace{(y_i - x_i^T \theta)}_{1 \times 1}^2$$

We will need the following results from matrix differentiation:

$$\frac{\partial A\theta}{\partial \theta} = A^T \text{ and } \frac{\partial \theta^T A \theta}{\partial \theta} = 2A^T \theta$$

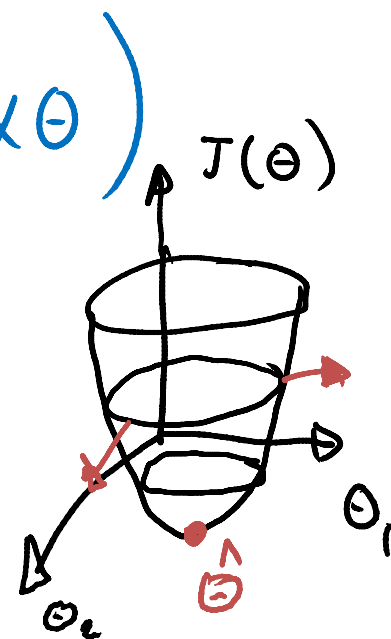
$$\frac{\partial J(\theta)}{\partial \theta} = \frac{\partial}{\partial \theta} (y - x\theta)^T (y - x\theta)$$

$$\stackrel{(2)}{=} \frac{\partial}{\partial \theta} (y^T y - \underbrace{y^T x \theta}_{1 \times 1} - \underbrace{\theta^T x^T y}_{1 \times d \times d \times n \times n \times 1} + \theta^T x^T x \theta)$$

$$\stackrel{(3)}{=} 0 - 2(y^T x)^T + 2(x^T x)^T \theta$$

$$= -2X^T y + 2X^T X \theta = \nabla J(\theta)$$

$$(X\theta)^T = \theta^T X^T$$



Least squares estimates

$$-2X^T y + 2X^T X \theta = 0$$

$$(X)^{-1} \underline{\underline{NO}}$$

$$\overset{d \times n}{X^T} \overset{n \times d}{X} \overset{d \times 1}{\theta} = \overset{d \times n}{X^T} \overset{n \times 1}{y}$$

$$\frac{2}{X} \underline{\underline{NO}}$$

$$\hat{\theta}_{ls} = (X^T X)^{-1} X^T y$$

$$X^T X \theta = X^T y$$

$$\underbrace{(X^T X)^{-1} X^T X}_{\bar{I}} \theta = \underbrace{(X^T X)^{-1} X^T y}_{\bar{I} \theta}$$

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

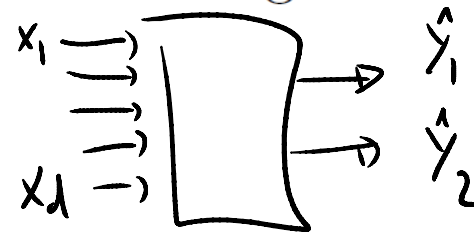
Hat

$$\hat{y} = X^* \hat{\theta}_{ls}$$

Multiple outputs

If we have several outputs $\mathbf{y}_i \in \mathbb{R}^c$, our linear regression expression becomes:

eg $c=2$



$$\begin{bmatrix} \hat{y}_{11} & \hat{y}_{12} \\ \vdots & \vdots \\ \hat{y}_{n1} & \hat{y}_{n2} \end{bmatrix} = \begin{bmatrix} 1 & \dots & 1 \\ \vdots & & \vdots \\ 1 & \dots & 1 \end{bmatrix} \begin{bmatrix} x_{12} & \dots & x_{1d} \\ \vdots & & \vdots \\ x_{nd} & & \end{bmatrix} \begin{bmatrix} \Theta_{11} & \Theta_{12} \\ \vdots & \vdots \\ \Theta_{d1} & \Theta_{d2} \end{bmatrix}$$

Next lecture

In the next lecture, we learn to derive the linear regression estimates by maximum likelihood with multivariate Gaussian distributions.