

- Why linear regression?
  - Basic technique w/ many extensions
  - Widely used
- Consider the following dataset which is the sales of a particular product in 200 markets w/ advertising budgets for three different media: television, newspaper and radio.

### Possible questions:

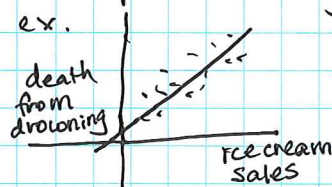
1. Are there linear relationships btwn budget vs. sales? media vs. sales?  
How strong are these relationships?

1. Which media contribute to sales?
2. How accurately can we estimate the effect of media on sales?
3. How accurately can we predict future sales from media?

• We cannot answer 2 or 3 with linear regression

(2) implies causality, and correlation/association  $\nRightarrow$  causality.

ex.



Obviously icecream sales don't cause death by drowning. They may correlate ~~as~~ if they are both affected by the same <sup>(third)</sup> variable.

ex. (good) Weather  $\rightarrow$  death by drowning  $\uparrow$   
ice cream sales  $\uparrow$

- Simple linear regression : we model linear dependence of a response variable / output / dependent variable  $Y$  on a predictor / input / feature / independent variable  $X$ .
- Stochastic linear model:



$$Y_i = \beta_0 + \beta_1 X_i + \varepsilon_i \quad \left. \vphantom{Y_i = \beta_0 + \beta_1 X_i + \varepsilon_i} \right\} \begin{array}{l} \varepsilon : \text{statistical error;} \\ \text{considered } \underline{\text{random}} \end{array}$$

- Some notes and assumptions.

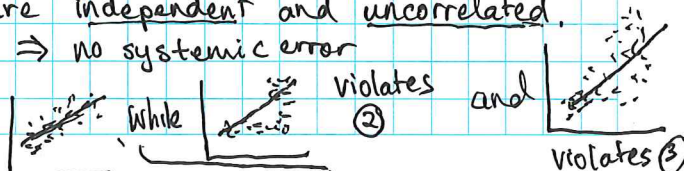
- $\varepsilon$  (statistical error) can represent ex measurement error ; unknown, non-systemic effects
- We assume  $X_i$  are fixed ie non-random for simplicity.
- $\beta_0$  and  $\beta_1$  are parameters of the model ; they are unknown and non-random
- Since we consider  $\varepsilon$  to be random  $\Rightarrow Y$  is also random  $\Rightarrow$  model is stochastic

↳ Assumptions on  $\varepsilon_i$ 's (1)  $\varepsilon_i, \dots, \varepsilon_n$  are independent and uncorrelated.

$\forall i=1, \dots, n$   $\begin{cases} \textcircled{2} E(\varepsilon_i) = 0 \Rightarrow \text{no systematic error} \\ \textcircled{3} \text{Var}(\varepsilon_i) = \sigma^2 \end{cases}$

ex.  while 

This satisfies all 3 assumptions.





(from prev page) Simple linear regression: stochastic model:

$$Y_i = \beta_0 + \beta_1 x_i + \varepsilon_i, \quad i=1, \dots, n$$

- The data are represented as realizations  $(x_1, y_1), \dots, (x_n, y_n)$  where  $y_i$  is a realization of  $Y_i$ .

- the linear model is linear in the parameters  $\beta_0, \beta_1$   
like the one above

ie.  $x_i$  may be non-linear in a linear model, ex.  $Y_i = \beta_0 + \beta_1 x_i^2 + \varepsilon_i$

Note that:

- The unknown parameters are  $\beta_0, \beta_1$  and  $\sigma^2$
- We cannot observe  $\varepsilon_i$ 's since the true line (set of parameters) is unknown

Goal:

so ...

- We want to estimate  $\beta_0, \beta_1, \sigma^2$  (true values)
- We want to find  $\hat{\beta}_0, \hat{\beta}_1$  so that ("estimates") the line  $\hat{\beta}_0 + \hat{\beta}_1 x$  is "close" to the data points.

That is ... The Game of Statistics:

- We see only one dataset

- We obtain the estimates  $\hat{\beta}_0, \hat{\beta}_1$  by understanding the precision/accuracy of these estimates.

- We want to say something about the unknown  $\beta_0, \beta_1$  (true vals) ie. We want to infer something about the true distribution from estimates  $(\hat{\beta}_0, \hat{\beta}_1)$

### How to measure "closeness"?

- often we find the params that give us the smallest residual sum of squares, ie least squares (LSQ) (RSS).

$$\text{RSS (residual sum of squares)} = \sum_{i=1}^n e_i^2$$

where

$$e_i = y_i - \hat{y}_i \text{ is the residual; and } \hat{y}_i = \hat{\beta}_0 + \hat{\beta}_1 x_i \text{ is the fitted value}$$

- Graphically we see



- For any combination  $(a, b)$  we can compute  
( $\hat{y}_i$ ) the fitted values  $a + bx_i = \hat{y}_i$ ;  
( $e_i$ ) the corresponding residuals  $e_i = y_i - \hat{y}_i$ ; and  
(RSS) the RSS  $\sum_{i=1}^n e_i^2$

$$\text{s.t. } (\hat{\beta}_0, \hat{\beta}_1) = \underset{a, b}{\operatorname{argmin}} \sum_{i=1}^n (y_i - (a + bx_i))^2$$

Note  $\hat{\beta}_0$  and  $\hat{\beta}_1$  are random!

Repetitions of regression give us

Note that regressions tend to be more accurate in the middle of the data than ends since the latter are more affected by changes in slope (small)



Lecture 1: Introduction; Linear Regression, cont'd

22 Feb 2018

Note / Tip: when examining LSQ regression lines, look at the data / line / graph interval by interval. The line segment should lie in the middle of the data in that interval

Multiple Linear Regression

$$Y_i = \beta_1 X_{i1} + \beta_2 X_{i2} + \dots + \beta_p X_{ip} + \varepsilon_i$$

Here our unknown parameters are  $\beta_1, \dots, \beta_p, \sigma^2$

• We have  $p$  predictor variables.

• We set  $X_{i1} = 1$  s.t.  $\beta_1 X_{i1} = \beta_1$  is the y-intercept of the model

• As before,  $\varepsilon_1, \dots, \varepsilon_n$  are uncorrelated (and independent)

$$\begin{cases} E(\varepsilon_i) = 0 \\ \text{Var}(\varepsilon_i) = \sigma^2 \end{cases} \forall i = 1, \dots, n.$$

We can write the above equation for our multiple linear model in

• vector form

or

$$Y_i = \mathbf{X}_i^T \boldsymbol{\beta} + \varepsilon_i, \quad \mathbf{X}_i = (X_{i1} \ X_{i2} \ \dots \ X_{ip})$$

$$\text{and } \boldsymbol{\beta} = \begin{pmatrix} \beta_1 \\ \vdots \\ \beta_p \end{pmatrix}$$

\* matrix form

$$\begin{matrix} (n \times 1) & (n \times p) & (p \times 1) & (n \times 1) \\ \mathbf{Y} & = & \mathbf{X} \boldsymbol{\beta} & + \boldsymbol{\varepsilon} \end{matrix}$$

usually we assume  $n > p$ .

$$\begin{matrix} \begin{pmatrix} Y_1 \\ \vdots \\ Y_n \end{pmatrix} & = & \begin{pmatrix} X_{11} & \dots & X_{1p} \\ \vdots & & \vdots \\ X_{n1} & \dots & X_{np} \end{pmatrix} \begin{pmatrix} \beta_1 \\ \vdots \\ \beta_p \end{pmatrix} + \begin{pmatrix} \varepsilon_1 \\ \vdots \\ \varepsilon_n \end{pmatrix} \\ n \times 1 & & n \times p & p \times 1 & n \times 1 \end{matrix}$$

$\mathbf{X}$  is called a design matrix:

the  $j^{\text{th}}$  row of  $\mathbf{X}$  is denoted by  $\mathbf{X}_j^T$ , and

the  $k^{\text{th}}$  column of  $\mathbf{X}$  is denoted by  $\mathbf{X}^{(k)}$

$$\text{ie } \mathbf{X} = (\mathbf{X}_j^{(k)}) \text{ , } j \in \{1, \dots, n\} \text{ and } k \in \{1, \dots, p\}$$