

# STK 2100

Riccardo De Bin

debin@math.dio.no

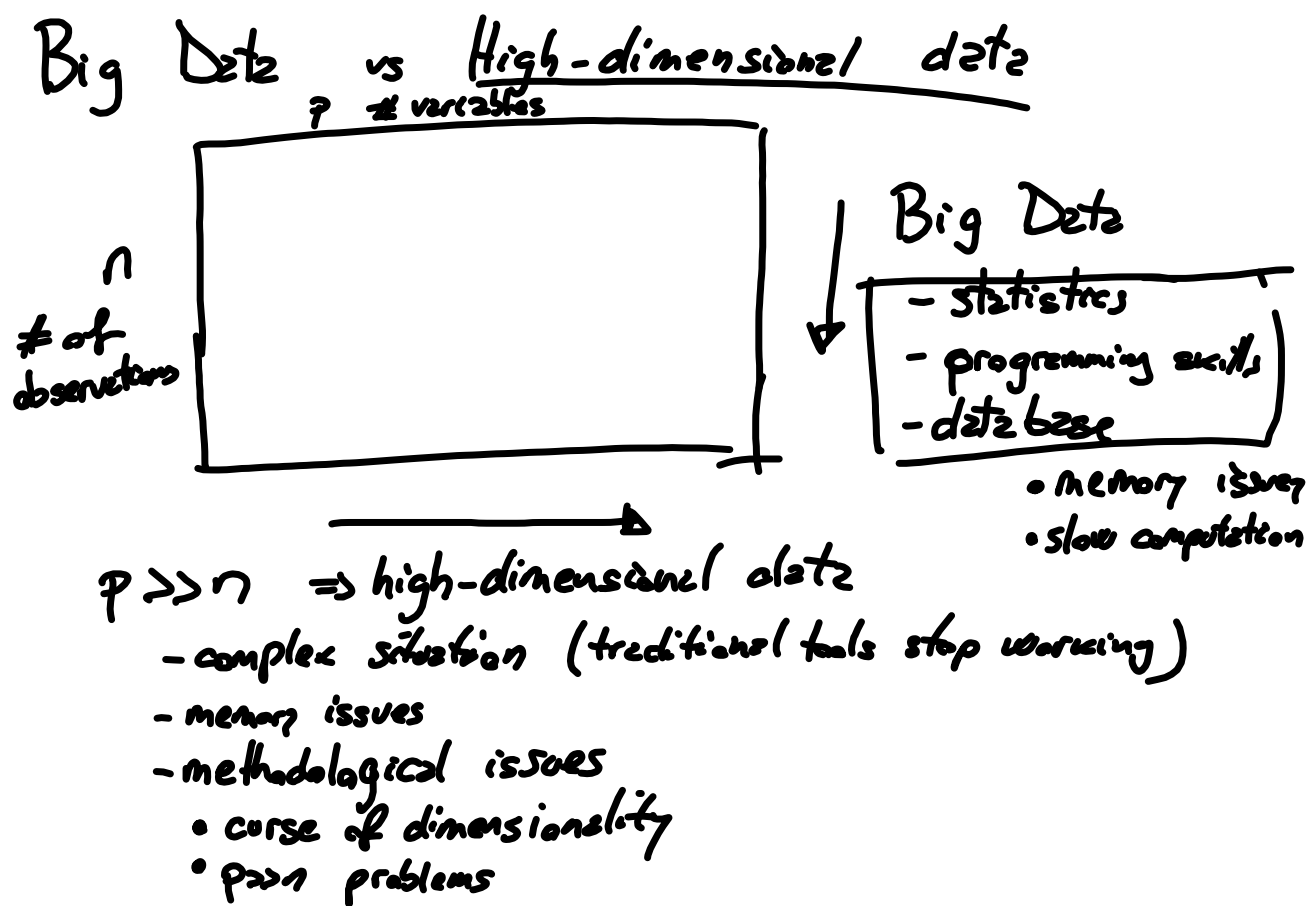
- what is statistical learning (data mining)
- linear models
- variable transformations



- black-box vs interpretable model
- prediction vs explanation

Huge amount of data: examples

- receipts
- credit card
- telephone companies
- web
- genetic data
- academic
- freud (class imbalance)



# What is a model?

"All models are wrong, but some are useful" (G.E.P. Box)

$$y = \underbrace{f(x)}_{\text{model}} + \varepsilon \rightarrow \text{error (everything we cannot explain)}$$

↖ systematic part

Azzolini & Scarpz (2012)'s definition:

'A model is a simplified representation of the phenomenon of interest, functional for a specific objective'

- simplified representation
  - portability/usability
  - focus on the important aspect
  - eliminate everything not essential
  - degree of complexity
    - trade off (e.g. bias/variance trade-off)
- functional for a specific objective
  - same data, we can have different models based on the goal (description/explanation vs prediction)

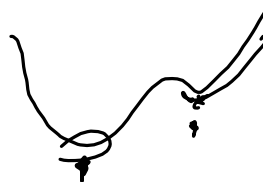
simplicity

Is there a true model?

Experimental studies vs observational studies

Against  
Pressing a button

- understand the theory behind → what are strengths and weaknesses & possible approaches
- interpret the results → good analyst behind
- reliability → find <sup>to be sure that we find</sup> a global minimum and not a local minimum



Software

R

→ CRAN

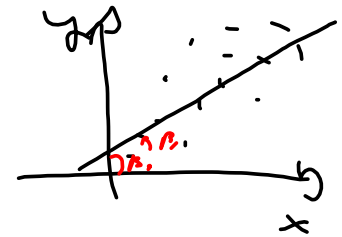
→ RStudio

→ Google's R Style Guide

↗ Manuals

# Basic: linear model

simplest way to relate two variables  
simple linear regression model



$$y = \underbrace{f(x; \beta)} + \varepsilon = \beta_0 + \beta_1 x + \varepsilon$$

$y$  = response, dependent variable, outcome

$x$  = independent variable, input, predictor, explanatory variable

$\beta_0, \beta_1$  = regression parameters / coefficients

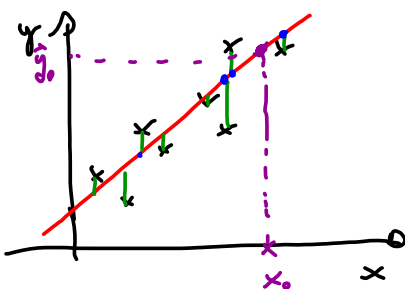
# regr. coeff = 2  $\rightarrow \beta_0$  = intercept (always)  
 $\beta_1$  = slope

$\varepsilon$  = error  $E[\varepsilon_i] = 0$   
 $Var(\varepsilon_i) = \sigma^2$  homoskedasticity  
 $Cor(\varepsilon_i, \varepsilon_j) = 0 \quad \forall i \neq j$

GOAL: estimate the values of  $\beta_0, \beta_1$  using the information in the data  
 $(x_i, y_i), i=1, \dots, n \rightarrow$  sample size

HOW: by minimizing a loss function (objective function)  
most often the sum of squares (square loss)

$$\hat{\beta} = \underset{\beta}{\operatorname{argmin}} \left\{ \underbrace{\sum_{i=1}^n (y_i - f(x_i, \beta))^2}_{\text{RSS}(\beta) = S(\beta)} \right\} \quad \|y - f(x; \beta)\|^2$$



$$\hat{y}_i = f(x_i; \hat{\beta}) = \hat{\beta}_0 + \hat{\beta}_1 x_i \quad \text{fitted value}$$

$$\hat{y}_0 = f(x_0; \hat{\beta}) = \hat{\beta}_0 + \hat{\beta}_1 x_0 \quad \text{predicted value}$$

$f(x; \beta)$  is not a line  $\rightarrow$  polynomial form

$$f(x; \beta) = \beta_0 + \beta_1 x + \beta_2 x^2 + \dots + \beta_{p-1} x^{p-1}$$

$\beta$  is a  $p$ -dimensional vector

$$f(x, \beta) = \beta_0 + \beta_1 x + \beta_2 x^2 + \beta_3 x^3$$

The model is linear in the parameters:

- conceptually and mathematically simple
- easy to compute

$$f(x; \beta) = X\beta$$

(n x p) (p x 1)

$X$  is a matrix  $X = (1, x, \dots, x^{p-1})$

$X$  is also called design matrix

In broad generality, the linear model has form

$$y = X\beta + \epsilon$$

Minimizing the sum of squares

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

$$\frac{\partial D(\beta)}{\partial \beta} = X^T (y - X\beta)$$

$$\frac{\partial D(\beta)}{\partial \beta} = 0$$

$$X^T y - X^T X \beta = 0$$

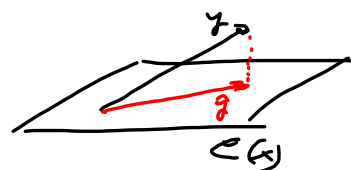
$$X^T X \beta = X^T y$$

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

LEAST SQUARES ESTIMATOR

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

HAT MATRIX  $\leftarrow$  PROJECTION MATRIX



space spanned by the columns of  $X$

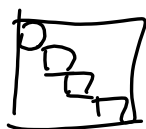
$$D(\hat{\beta}) = \|y - \hat{y}\|^2$$

We can use the deviance to estimate  $\sigma^2$ :  $s^2 = \frac{D(\hat{\beta})}{n-p}$

and the variance of  $\hat{\beta}$

$$\text{Var}(\hat{\beta}) = s^2 (X^T X)^{-1}$$

$p \times p$  matrix in which the diagonal terms are the variance of  $\hat{\beta}_j$



Include information about fuel type  
(a new variable that can help explaining the variability of  $y$ )

create a dummy variable  $I_A = \begin{cases} 0 & \text{if fuel type = 'gas'} \\ 1 & \text{if fuel type = 'diesel'} \end{cases}$

simplest way to include the variable  $\rightarrow$  additive effect

$$y = \beta_0 + \beta_1 x + \beta_2 x^2 + \beta_3 x^3 + \beta_4 I_A + \varepsilon$$

in matrix form

$$y = X\beta + \varepsilon$$

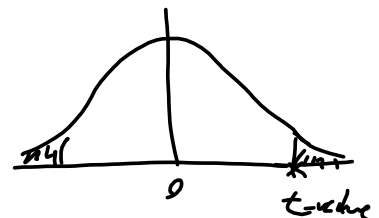
where

$$X = (1, x, x^2, x^3, I_A)$$

$$\beta = (\beta_0, \beta_1, \beta_2, \beta_3, \beta_4)^T$$

Assuming  $\varepsilon \sim N(0, \sigma^2) \rightarrow$  linear Gaussian regression

$$t\text{-values: } \frac{\hat{\beta}_j}{\hat{se}(\hat{\beta}_j)} = \frac{\hat{\beta}_j}{\sqrt{S^2(x^T x)^{-1}_{j,j}}}$$



$p$ -values =  $\Pr(\text{obtaining a } t\text{-value larger in absolute value than that we actually obtained under the null hypothesis } \beta_j = 0)$