

## Analysis of Variance (ANOVA)

### Introduction

- Generalization for  $n$  samples (arbitrary number).
- ANOVA tests if 'at least one' of samples is significantly different ~~than~~ <sup>from</sup> the others.
- Some examples:
  - + Testing 3 different designs of a website, sample for every design has its own:  $\bar{X}, s, n$

web version 1	web version 2	web version 3
$\bar{X}_1$	$\bar{X}_2$	$\bar{X}_3$
$s_1$	$s_2$	$s_3$
$n_1$	$n_2$	$n_3$

+ Testing 5 different products in olive trees.

- $H_0 \Rightarrow \mu_1 = \mu_2 = \dots = \mu_k$ ;  $H_1 \Rightarrow$  At least one is significantly different, don't know which one in advance.
- Hand-made example:

one way anova (one 'dimension')

Sample 1	Sample 2	Sample 3
3	5	5
2	3	6
1	4	7

$m = 3$        $n = 3$

① Compute global mean

$$\bar{X} = \frac{1+2+3+3+4+5+5+6+7}{9} = 4$$

$$d.f. = 9 - 1 = 8$$

③ Compute sample means

$$\bar{X}_1 = \frac{1+2+3}{3} = 2 \quad \bar{X}_2 = \frac{3+4+5}{3} = 4 \quad \bar{X}_3 = \frac{5+6+7}{3} = 6$$

② Compute SST

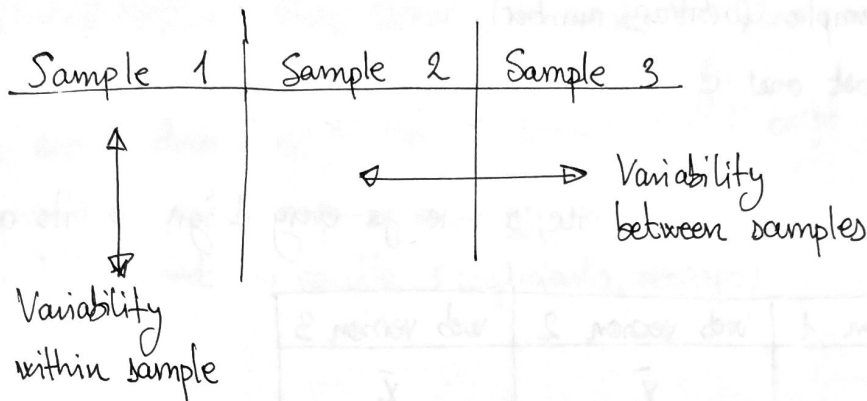
$$\begin{aligned} SST &= (3-4)^2 + (2-4)^2 + (1-4)^2 + \\ &\quad + (5-4)^2 + (3-4)^2 + (4-4)^2 + \\ &\quad + (5-4)^2 + (6-4)^2 + (7-4)^2 = 30 \end{aligned}$$

④ Compute SSW (Sum of Squares Within Samples)  $= \sum_{i=1}^m \sum_{j=1}^n (x_{ij} - \bar{x}_i)^2$

$$SSW = (3-2)^2 + (2-2)^2 + (1-2)^2 + (5-4)^2 + (3-4)^2 + (4-4)^2 + (5-6)^2 + (6-6)^2 + (7-6)^2 = 6$$

$$d.f = m \cdot (n-1) = 3 \cdot (3-1) = 6$$

↑                      ↑  
diff. samples      sample mean



⑤ Compute SSB (Sum of Squares Between Samples)

$$SSB = \left[ (\bar{x}_1 - \bar{\bar{x}})^2 \cdot n_1 + (\bar{x}_2 - \bar{\bar{x}})^2 \cdot n_2 + (\bar{x}_3 - \bar{\bar{x}})^2 \cdot n_3 \right] = (2-4)^2 \cdot 3 + (4-4)^2 \cdot 3 + (6-4)^2 \cdot 3 = 24$$

$d.f = m - 1 = 2$

⑥ Compute F statistic (F distribution) Fischer - Snedecor

$$MSW = \frac{SSW}{d.f_{SSW}} = \frac{6}{6} = 1$$

$$MSB = \frac{SSB}{d.f_{SSB}} = \frac{24}{2} = 12$$

$$F = \frac{MSB}{MSW}$$

$$F \sim \frac{X_1^2}{X_2^2} = \frac{12}{1} = 12$$

$$p = 0.008 \quad \alpha = 0.05$$

⑦ Perform test  $p < \alpha \rightarrow$  Reject null hypothesis!!

# Linear Regression

## Introduction

- One of the simplest 'Machine Learning' algorithms.
- Very good for modeling linear relationships in data. There is no free lunch in Data Science.

It has some advantages:

- Simple and not too prone to overfit → explain overfit a little
- Easy to interpret its results. (coefficients, intercept)
- Analytical solution

And drawbacks:

- Too simple to capture non-linear complex relationships out of the box.
- Analytical solution not suitable for big data problems.

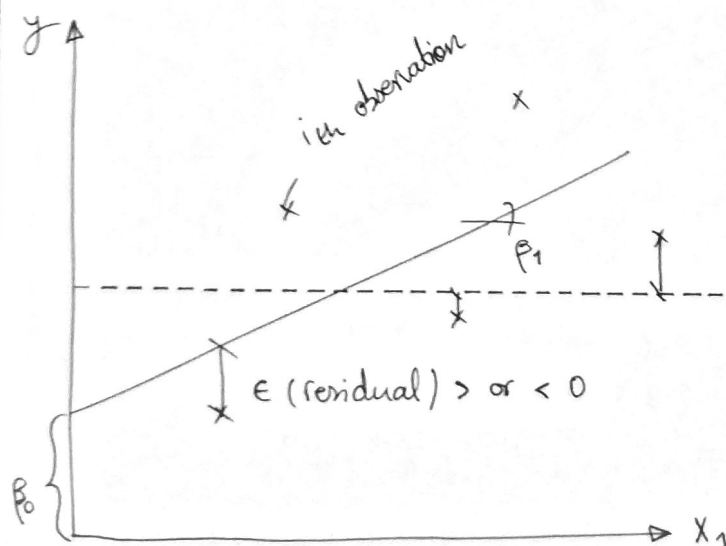
$$y = \beta_0 + \beta_1 \cdot X_1 + \beta_2 \cdot X_2 + \dots + \beta_n \cdot X_n$$

Annotations:  $\beta_0$  is the intercept,  $X_1, X_2, \dots, X_n$  are independent variables/features, and  $\beta_1, \beta_2, \dots, \beta_n$  are coefficients.  $y$  is the dependent variable.

Some examples

- + salary vs years of experience
- + price of an apartment vs  $m^2$

## Explanation (Single feature version)



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

$$\begin{bmatrix} y_1 \\ y_2 \\ \vdots \\ y_m \end{bmatrix} = \begin{bmatrix} 1 & X_1^1 & \dots & 1 \\ \beta_0 & \beta_1 & \dots & \beta_n \\ X_1^2 & \dots & X_1^m & \dots & X_n^2 & \dots & X_n^m \end{bmatrix} \cdot \begin{bmatrix} \beta_0 \\ \vdots \\ \beta_n \end{bmatrix}$$

Annotations: The first column of the matrix (all 1s) is for the intercept. The columns  $X_1^1, \dots, X_n^1$  are the first row of features.

$$y_i = \beta_0 + \beta_1 \cdot X_1^i + \epsilon^i$$

① Compute residuals " $\epsilon^i$ " for every point.

② Compute Sum of Squared Error:

$$SSE = (\epsilon^1)^2 + (\epsilon^2)^2 + \dots + (\epsilon^m)^2$$

eg. SSW in ANOVA! | Problems? Units<sup>2</sup>  
Grow with  $m$

③ Compute RMSE: How good is my model?

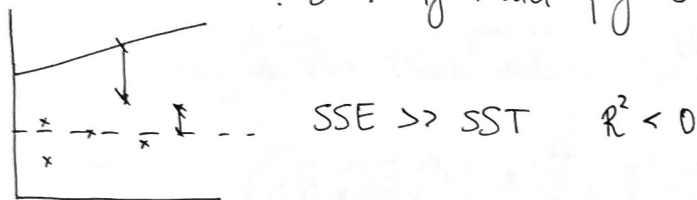
SSE is hard to interpret and growing with samples! Let's normalize:

$$RMSE = \sqrt{SSE / N}$$

④ How to compare models?

$$\text{Compute } SST = \sum_i^m (X_i - \bar{X})^2 \quad \text{and} \quad R^2 = 1 - \frac{SSE}{SST}$$

Can be  $< 0$ ? Yes! If model performs worse than just taking the mean!



### Features significance

$$t \text{ statistic} = \frac{\text{coef}}{\text{std. error}} \quad \uparrow \uparrow \text{ (want it high or low, extreme)}$$

p (probability of getting  $\frac{s}{\sqrt{n}}$  such an extreme result by chance)  $\downarrow \downarrow$

Colinearity can exist, compute Pearson's R and remove correlated features.

$$[0, \pm 1]$$