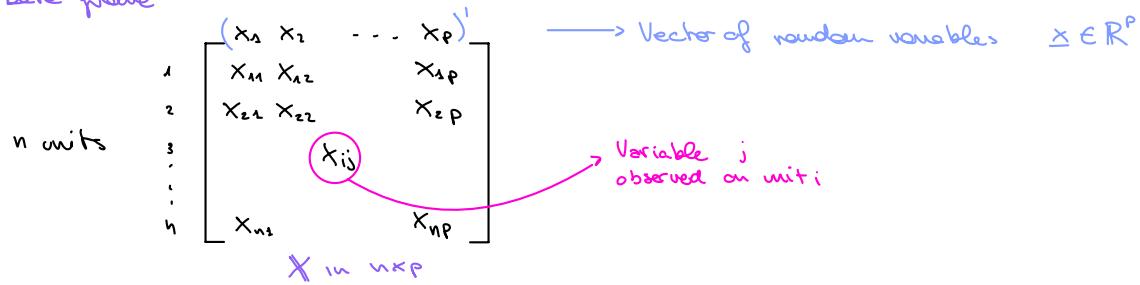


Data frame



Row (unit) perspective

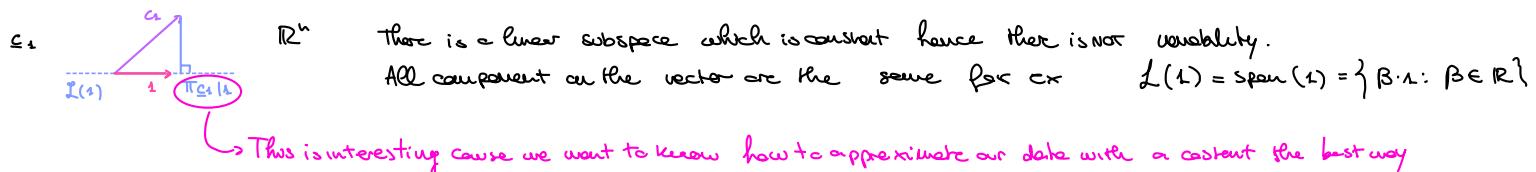
$$\begin{aligned} \underline{x}_1' &= (x_{11}, \dots, x_{1p}) \\ \underline{x}_2' &= (x_{21}, \dots, x_{2p}) \\ &\vdots \\ \underline{x}_n' &= (x_{n1}, \dots, x_{np}) \end{aligned} \quad \left. \right\} \in R^p$$

Column perspective

$$\underline{c}_z = \begin{pmatrix} x_{11} \\ x_{21} \\ \vdots \\ x_{n1} \end{pmatrix} \in R^n$$

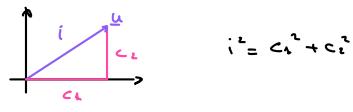
These numbers are all talking about the same information
1 variable for n statistical units (we have information about mutual variability)

Having dependence between 2 variables we can create a prediction.



Exclusion spaces

$$\begin{aligned} u, v &\in R^n \\ \|u\| &= \sqrt{\sum_{i=1}^n u_i^2} \\ u \cdot v &= \langle u, v \rangle \end{aligned}$$



$$\cos \theta = \frac{\langle u, v \rangle}{\|u\| \|v\|} = \frac{u \cdot v}{\sqrt{u \cdot u} \sqrt{v \cdot v}}$$

$$\pi_{u \perp v} = \|u\| \frac{u}{\|u\| \|v\|} \cdot \frac{v}{\|v\|} = \frac{u \cdot v}{\|u\| \|v\|} v = \boxed{\frac{u \cdot v}{\|u\| \|v\|} v}$$

orth proj operator on $L(u)$

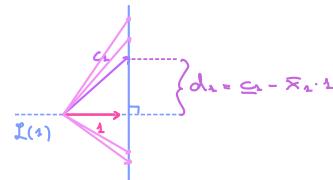
$$\pi_{L(u)} = \frac{1}{\|u\|} u \cdot 1 = \frac{1}{\|u\|} c_u \cdot 1 = \frac{1}{n} \sum_{i=1}^n x_{i1} = \bar{x}_1 \cdot 1$$

The mean

Different sample distributions can have the same mean because they have the same orthogonal projection, but they have different distance from it.

$$d_n = \begin{bmatrix} c_{11} - \bar{x}_1 \\ c_{21} - \bar{x}_1 \\ \vdots \\ c_{n1} - \bar{x}_1 \end{bmatrix}$$

$$\|d_n\| = \sqrt{\sum_{i=1}^n (x_{i1} - \bar{x}_1)^2}$$



SAMPLE STANDARD DEVIATION

$$\sqrt{s_{11}} = \sqrt{\frac{1}{n-1} \sum_{i=1}^n (x_{i1} - \bar{x}_1)^2}$$

$$\|d_n\| \propto \sqrt{s_{11}} = \sqrt{n-1} \sqrt{s_{11}}$$

$$s_{11} = \frac{1}{n-1} \sum_{i=1}^n (x_{i1} - \bar{x}_1)^2$$

SAMPLE VARIANCE

Tchebychev's inequality

$$\text{Freq} [\bar{x}_1 - k\sqrt{s_{11}} < x_1 < \bar{x}_1 + k\sqrt{s_{11}}] \geq 1 - \frac{1}{k^2} \quad k > 0$$

$$\text{Freq} [x_1 \in [\bar{x}_1 - 2\sqrt{s_{11}} < x_1 < \bar{x}_1 + 2\sqrt{s_{11}}] \geq 75\%$$



$$\underline{x}_1 = \begin{pmatrix} x_{11} \\ x_{12} \\ \vdots \\ x_{1n} \end{pmatrix} \in \mathbb{R}^n \quad \rightarrow \quad \bar{x}_{1 \cdot 1}, \quad d_1 = \underline{x}_1 - \bar{x}_{1 \cdot 1}$$

$$\underline{x}_2 = \begin{pmatrix} x_{21} \\ x_{22} \\ \vdots \\ x_{2n} \end{pmatrix} \in \mathbb{R}^n \quad \rightarrow \quad \bar{x}_{2 \cdot 1}, \quad d_2 = \underline{x}_2 - \bar{x}_{2 \cdot 1}$$



Learning cases

$$\theta_{12} = 0 \quad \xrightarrow{\underline{d}_1 \parallel \underline{d}_2} \quad \underline{d}_2 \in \text{Span}(\underline{d}_1) \quad \underline{d}_2 = \beta \underline{d}_1$$

$$x_{i2} - \bar{x}_2 = \beta(x_{i1} - \bar{x}_1) \quad i=1, 2, \dots, n \quad \text{all information of one in the other}$$

$$x_{i2} = \beta x_{i1} + (\bar{x}_2 - \beta \bar{x}_1)$$

$$\theta_{12} = \frac{\pi}{2} \quad \xrightarrow{\underline{d}_1 \perp \underline{d}_2} \quad \text{no information one in the other}$$

$$\cos \theta_{12} = \frac{\langle \underline{d}_1, \underline{d}_2 \rangle}{\|\underline{d}_1\| \|\underline{d}_2\|} = \frac{\underline{d}_1^\top \underline{d}_2}{\sqrt{n-1} \sqrt{n-1} \sqrt{s_{11}} \sqrt{s_{22}}} = \frac{\sum_{i=1}^{n-1} (x_{i1} - \bar{x}_1)(x_{i2} - \bar{x}_2)}{(n-1) \sqrt{s_{11} s_{22}}} = \frac{s_{12}}{\sqrt{s_{11}} \sqrt{s_{22}}} = \text{Corr}(x_1, x_2) = \rho_{12}$$

$$0 < \theta_{12} < \frac{\pi}{2} \quad \xrightarrow{\underline{d}_1 \text{ and } \underline{d}_2 \text{ not perpendicular}} \quad \text{Corr}(x_1, x_2) = \frac{1}{n-1} \sum_{i=1}^{n-1} (x_{i1} - \bar{x}_1)(x_{i2} - \bar{x}_2) = s_{12}$$

$0 \leq \rho_{12} \leq 1$ (because we have a cosine, hence the correlation is 1 because we have a angle)

Data frame

$$\begin{pmatrix} \bar{x}_1 \\ \vdots \\ \bar{x}_p \end{pmatrix} = \bar{x} \quad \text{vector of means}$$

$$S_{jk} = \frac{1}{n-1} \sum_{i=1}^{n-1} (x_{ij} - \bar{x}_j)(x_{ik} - \bar{x}_k) \quad \text{sample covariance if } j \neq k$$

$$S_{jj} = \text{Var}(x_j) \quad \text{sample variance}$$

$$S = \begin{bmatrix} S_{11} & S_{12} & \dots & S_{1p} \\ \ddots & \ddots & & \vdots \\ & & S_{pp} & \end{bmatrix} \quad p \times p$$

$$\rho_{jk} = \frac{S_{jk}}{\sqrt{S_{jj} S_{kk}}} \quad j, k = 1, \dots, p \quad \rho = \begin{bmatrix} 1 & \rho_{12} & \dots & \rho_{1p} \\ \ddots & \ddots & & \vdots \\ & & 1 & \end{bmatrix}$$

Rank (standardization)

$$x_1 \xrightarrow{\text{std}} \frac{x_1 - \bar{x}_1}{\sqrt{s_{11}}} = z_1 \rightarrow \text{how far you are from the mean}$$

$$z_1 \quad \text{Var}(z_1) = 1 \quad \mathbb{E}[z_1] = 0$$

Stat Learning (?)

$$X = X \cdot Y \rightarrow \text{Target variable}$$

1. Supervised learning
2. Unsupervised learning

Fundamental ideas about Stat Learning

Supervised setup $\underline{x} = (x_1, \dots, x_p)^\top \in \mathbb{R}^p \quad Y \in \mathbb{R}$

Observed: on n units $\Rightarrow X$ data frame

Goal: explain the variability of Y as a function of X

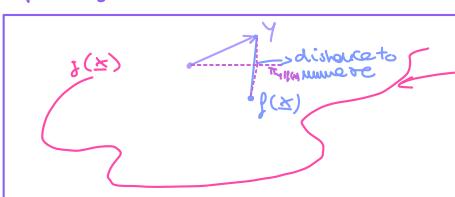
$Y \rightarrow$ Random variable

$X \rightarrow$ Random vector

? $f: \mathbb{R}^p \rightarrow \mathbb{R}$ s.t. $f(\underline{x})$ is a "prediction" of Y

What f ? "The best"

Space of all r.v. in \mathbb{R} ($w.s.t. \mathbb{E}[W^2] < \infty$)



2 random vectors, hence we have to minimize it for each unit that we will ever possibly observe

$$\text{find } f: \mathbb{R}^p \rightarrow \mathbb{R} \text{ s.t. } \underbrace{|Y - f(\underline{x})|^2}_{\text{measuring not this BUT}}$$

$\mathbb{E}[|Y - f(\underline{x})|^2]$

$$\pi_{Y|x} = E[Y|x]$$

REGRESSION FUNCTION

(can have any form, don't confuse it with linear)

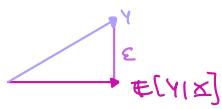
$$E[(Y-f(x))^2] = E[(Y - E[Y|x] + E[Y|x] - f(x))^2] = \\ = E[(Y - E[Y|x])^2] + E[(E[Y|x] - f(x))^2] + 2E[(Y - E[Y|x])(E[Y|x] - f(x))]$$

$$w, z \text{ r.v. } E[w] = E[E[w|z]]$$

$$\textcircled{*} = E[(\cdot)(\cdot)] = E[E[(Y - E[Y|x])(E[Y|x] - f(x))|x]] = \\ = E[(E[Y|x] - f(x)) \cdot \underbrace{E[Y - E[Y|x]|x]}_{=0}]$$

$$E[(Y-f(x))^2] = E[(Y - E[Y|x])^2] + E[(E[Y|x] - f(x))^2]$$

$$\text{Hence } f(x) = \arg \min_{g(x)} E[(Y - g(x))^2] = E[Y|x] \quad \text{This holds for everything in } \mathcal{C}$$



General model:

$$Y = E[Y|x] + \epsilon$$

$$Y = f(x) + \epsilon \quad \text{where } f(x) = E[Y|x]$$

$$E[Y] = E[f(x) + \epsilon] = E[E[Y|x]] + E[\epsilon] = E[Y] + E[\epsilon]$$

$$\rightarrow E[\epsilon] = 0$$

General model

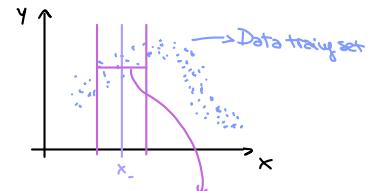
$$Y = f(x) + \epsilon$$

$$f(x) = E[Y|x]$$

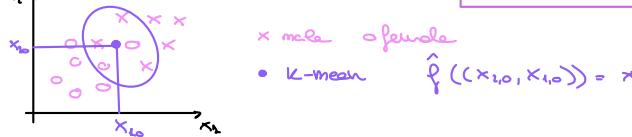
$$\epsilon \text{ r.v. s.t. } \epsilon \perp \sigma(x) \quad E[\epsilon] = 0$$

Learn f from data (and from the engineering knowledge)
 $\Rightarrow X \xrightarrow{\text{data}} \hat{f} (= E[Y|x])$

Do we really need a model for f ?



$$\text{For } x \in \mathbb{R} \\ \text{let } N_x = \{x_i : i \in \text{train set} \text{ and close to } x\} \\ \hat{f}(x) = \frac{1}{k} \sum_{x_i \in N_x} y_i$$



Local methods for obtaining \hat{f} are not (always) good!

Curse of dimensionality

An idea

$S_r(p)$ universe of radius r in \mathbb{R}^p

i.e. sphere centered in 0, with radius r in \mathbb{R}^p

Friends are uniformly distributed in $S_r(p)$. How far to go to meet 10% of friends.

$$p=1 \quad \text{---} \quad S_1 \quad ?r \rightarrow \frac{2r}{2} = 0, 1 \quad r=0, 1$$

$$p=2 \quad \text{---} \quad S_2 \quad \frac{\pi r^2}{\pi 1^2} \rightarrow r^2 = 0, 1 \quad r = 0, \sqrt{2}$$

$$p=3 \quad \text{---} \quad S_3 \quad \frac{4\pi r^3}{4\pi 1^3} \rightarrow r^3 = 0, 1 \quad r = 0, \sqrt[3]{3}$$

$$p=100 \quad r^{100} = 0, 1 \quad r = 0, \sqrt[100]{1}$$

To extend the concept we have to face higher data to get information on higher dimension. Basically the space is too big and near me there are not people like me, because of dimension.



We have to reduce the dimensionality of the problem:

1. Reducing p (in a smart way, we have to still have information) \rightarrow PCA, ICA ...
2. Reduce the dim of f \rightarrow Structured Models

Struct model: Find $\hat{f} \in \{f_\theta \mid \theta \in \mathbb{R}^k\}$ if θ is known $\Rightarrow f_\theta$ is known
 $x \rightarrow \hat{\theta} \Rightarrow \hat{f} = f_\theta$

EXAMPLE

$$f_\theta(x) = \beta_0 + \beta_1 x_1 + \dots + \beta_p x_p \quad x = (x_1, \dots, x_p)'$$

$$\theta = (\beta_0, \beta_1, \dots, \beta_p)' \quad \text{this has not to be linear, we linearize it in this form}$$

Suppose that:

- obtained the data
 - estimated f (approx. of $E[Y|x]$)
- \Rightarrow Ready to solve my problem

x_0 new features for a new unit \Rightarrow predict y_0 by means of $\hat{f}(x)$

$$(y_0 - \hat{f}(x_0))^2 \quad Y = f(x) + \varepsilon \quad \text{But given the same input the outputs haven't to be the same value}$$

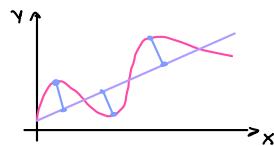
$$y_0 = f(x_0) + \varepsilon$$

$$\mathbb{E}_{\mathbb{X}}[(y_0 - \hat{f}(x_0))^2] = \mathbb{E}_{\mathbb{X}}[(f(x_0) + \varepsilon_0 + \hat{f}(x_0))^2] = \mathbb{E}_{\mathbb{X}}[(f(x_0) - \hat{f}(x_0))^2] + \mathbb{E}_{\mathbb{X}}[\varepsilon_0^2] + \mathbb{E}_{\mathbb{X}}[\varepsilon_0(\hat{f}(x_0) - f(x_0))]$$

$$= (f(x_0) - \hat{f}(x_0))^2 + \text{Var}(\varepsilon_0)$$

Generalization problem:

$$\mathbb{E}_{\mathbb{X}}[(y_0 - \hat{f}(x_0))^2] = \underbrace{(f(x_0) - \hat{f}(x_0))^2}_{\text{reducible error}} + \underbrace{\text{Var}(\varepsilon_0)}_{\text{irreducible}}$$



If we try to reduce too much we get overfitting \rightarrow basically overfit set, hence we don't have any new information nor anything useful.

Expected General. Error:

$$\mathbb{E}[\mathbb{E}_{\mathbb{X}}[(y_0 - \hat{f}(x_0))^2]] = \mathbb{E}[(f(x_0) - \hat{f}(x_0))^2] + \text{Var}(\varepsilon_0) = (f(x_0) - \mathbb{E}[\hat{f}(x_0)])^2 + \mathbb{E}[(\hat{f}(x_0) - \mathbb{E}[\hat{f}(x_0)])^2] + \text{Var}(\varepsilon_0)$$

$$= \text{bias}^2 + \text{Var}(\hat{f}(x_0)) + \text{Var}(\varepsilon_0)$$

$$\mathbb{E}[\mathbb{E}_{\mathbb{X}}[(y_0 - \hat{f}(x_0))^2]] = \text{bias}^2 + \text{Var}(\hat{f}(x_0)) + \text{Var}(\varepsilon_0)$$

