

EXPLORING A MULTIVARIATE DATASET

Descriptive statistics and graphical display. The geometry of multivariate sample. Sample mean, covariance and correlation. Generalized variance and total variance. The metric induced by the covariance matrix.

What is statistical learning?

Consider for instance a dataset of sales:



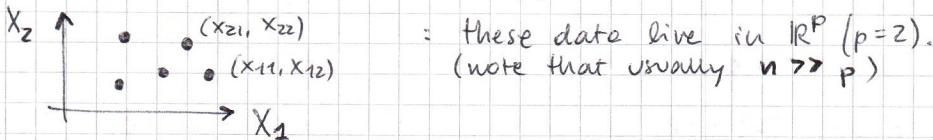
A linear regression line fit separately each one. Can we predict sales using all the informations? (All the features together?)

We'll always assume that we observed n statistical units, each one with p variables ($\equiv p$ features).

$$\begin{aligned} \underline{x}_1 &= [x_{11}, x_{12}, \dots, x_{1p}]^T \in \mathbb{R}^p \\ \underline{x}_2 &= [x_{21}, x_{22}, \dots, x_{2p}]^T \\ &\vdots \\ \underline{x}_n &= [x_{n1}, x_{n2}, \dots, x_{np}]^T \end{aligned} \quad \left\{ \begin{array}{l} \text{n-th statistical unit} \\ \text{DATA MATRIX} \end{array} \right. \quad \begin{matrix} X_1 & X_2 & \dots & X_p \end{matrix} \leftarrow \text{Features}$$

$$X = \begin{bmatrix} x_{11} & x_{12} & \dots & x_{1p} \\ x_{21} & x_{22} & \dots & x_{2p} \\ \vdots & \vdots & \ddots & \vdots \\ x_{n1} & x_{n2} & \dots & x_{np} \end{bmatrix}$$

We can read this data in two ways: by columns and by rows. Each row talks about a unit, each column about a variable. In the past statistical courses we saw how to determine the distribution of one feature, so we developed tools to analyze data by columns (mean, variance, quantile, ...), but the usual way to look at the data is by rows (by units). In this case we obtain a set of points in a \mathbb{R}^p space.



Often it happens that we get a "special" variable:

$$\begin{matrix} X_1 & X_2 & \dots & X_p & Y \end{matrix} \rightarrow \begin{array}{l} \text{what we want to predict/explain} \\ \text{in terms of the other features} \end{array}$$

$$\begin{bmatrix} x_{11} & x_{12} & \dots & x_{1p} & y_1 \\ x_{21} & x_{22} & \dots & x_{2p} & y_2 \\ \vdots & \vdots & \ddots & \vdots & \vdots \\ x_{n1} & x_{n2} & \dots & x_{np} & y_p \end{bmatrix} \quad (Y \text{ can also be categorical (a label): in this case we use information of } X_1, \dots, X_p \text{ to say something about membership of a group} \Rightarrow \text{classification problem})$$

So the problem is: observing p features we have to predict Y . How can we translate this problem mathematically?

- $\underline{X} \in \mathbb{R}^p$: random vector of features (\equiv what we observe)
- $Y \in \mathbb{R}$: what we want to predict

\Rightarrow Problem: use X to predict Y . What is the best function $f: \mathbb{R}^p \rightarrow \mathbb{R}$ to predict Y ? (f in terms of X)

What do we mean for "best"? We need a criterium for optimization. We want a function f such that:

$$\mathbb{E}[(Y - f(X))^2] \text{ is minimum.}$$

MEAN SQUARE ERROR:

average of all the possible realizations of X and Y . We use $(-)^2$ because we don't want the negative errors to compensate the positive ones.

If f was a constant we would have: $k = \arg \min_k \mathbb{E}[(Y - k)^2] = \mathbb{E}[Y]$, so we have sort of an intuition of what the function should be, except for the

fact that we don't want f to be a constant because it's useless for accurate predictions (in fact we wouldn't even be using \underline{x}).
The important fact is: whatever f will be, it should be sort of a mean: $f(\underline{x}) = E[\dots]$.

$$E[(Y - f(\underline{x}))^2] = E[a + b]$$

$E[Y|\underline{x}]$ = predictional expectation of Y given \underline{x}
(\equiv the best guess of Y once you know \underline{x}).

From a geometrical point of view this is simply
a projection of Y on the σ -field generated by \underline{x}
(it's a Radon-Nikodym derivative)

$$= E[(Y - E[Y|\underline{x}])^2] + E[(E[Y|\underline{x}] - f(\underline{x}))^2] + 2 \underbrace{E[(Y - E[Y|\underline{x}])(E[Y|\underline{x}] - f(\underline{x}))]}_{(*)}$$

We remember that: $E[W] = E[E[W|Z]]$

$$\Rightarrow (*) = 2 E[E[(\dots)(\dots) | \underline{x}]]$$

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

$$E[g(\underline{x}) | \underline{x}] = E[g(\underline{x})]$$

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

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

$$= 0$$

$$\Rightarrow E[(Y - f(\underline{x}))^2] = \underbrace{E[(Y - E[Y|\underline{x}])^2]}_{(1)} + \underbrace{E[(E[Y|\underline{x}] - f(\underline{x}))^2]}_{(2)}$$

So the problem is to find f s.t. $(1) + (2)$ is minimum. For doing that, (1) is useless, it's like a constant, and for minimizing (2) we take (looking at the const case):

$$f(\underline{x}) = E[Y|\underline{x}] \quad \leftarrow \text{solution of the optimization problem}$$

Of course there always be an error ϵ :

$$\epsilon = Y - f(\underline{x}) \quad (\text{what we cannot capture from } \underline{x} \text{ to predict } Y, \text{ it's IRREDUCIBLE using } \underline{x})$$

\Rightarrow We want to explain Y through a function f of \underline{x} and a residual ϵ .
What can we say (statistically) about ϵ ?

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

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

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

(this mean that if f is the best predictor
then f has the same average that Y)

$$\Rightarrow \text{The model is: } \begin{cases} Y = f(\underline{x}) + \epsilon \\ f(\underline{x}) = E[Y|\underline{x}] \\ E[\epsilon] = 0 \end{cases}$$

What do we know about f ? We should know the joint distribution $Y|\underline{x}$ but usually we don't know it
 \Rightarrow we need to use the data to estimate f

Suppose now that we have an estimation of f : \hat{f} .
 How good is this estimation? Note that we need \hat{f} not only to fitting the data, but for predictions too. We want to know what are the uncertainties of \hat{f} . We want a criteria not only to check if the model is good for the data, we want to see how good it is when we try to predict something new.

\hat{f} : estimation of f through data \mathbb{X}

$x_0 \in \mathbb{R}^p$: new observation ($\notin \mathbb{X}$) for which we want to predict y_0

We want to know about the error in predicting y_0 :

$$\begin{aligned} \mathbb{E}_{\mathbb{X}}[(y_0 - \hat{f}(x_0))^2] &= \mathbb{E}_{\mathbb{X}}[(y_0 - \hat{f}(x_0))^2 | \mathbb{X} = x] \\ &= \mathbb{E}_{\mathbb{X}}[\underbrace{(f(x_0) + \varepsilon - \hat{f}(x_0))^2}_{\text{they're constant}}] \quad \Rightarrow y_0 = f(x_0) + \varepsilon_0 \quad \text{error referring to } y_0 \\ &\quad (\text{we're conditioning on } \mathbb{X}, \text{ the only random is } \varepsilon) \\ &= \mathbb{E}_{\mathbb{X}}[(f(x_0) - \hat{f}(x_0))^2] + \mathbb{E}_{\mathbb{X}}[\varepsilon^2] + 2 \mathbb{E}_{\mathbb{X}}[(f(x_0) - \hat{f}(x_0)) \varepsilon] \\ &= (f(x_0) - \hat{f}(x_0))^2 + \text{Var}(\varepsilon_0) + 2(f(x_0) - \hat{f}(x_0)) \mathbb{E}_{\mathbb{X}}[\varepsilon_0] \\ &\quad \text{NOTE: } \mathbb{E}_{\mathbb{X}}[\varepsilon_0] = \mathbb{E}[\varepsilon_0] = 0 \quad \text{since } \mathbb{E}_{\mathbb{X}}[\varepsilon_0] = \mathbb{E}[\varepsilon_0] = 0 \\ &\Rightarrow \mathbb{E}_{\mathbb{X}}[\varepsilon_0] = \mathbb{E}[\varepsilon_0] \end{aligned}$$

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

we can try to make this error as small as possible.

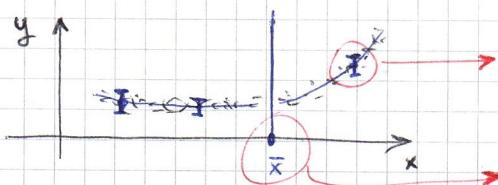
NOTE: x_0 is a new observation, $x_0 \notin \mathbb{X}$, it didn't influence the model, this is the right way to evaluate the model

we're not able to reduce this.

this is what we are not able to catch with any \hat{f} (or any f), it's the part of Y which is not explainable by X
It's **IRREDUCIBLE!**

How can we get an estimation \hat{f} ?

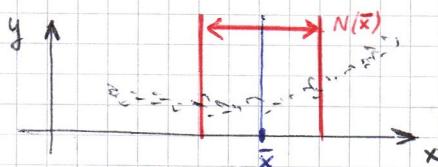
We wanted the **regression function** $\hat{f}(x) = \mathbb{E}[Y | X = x]$ but there is a problem: sometimes we have few data for a particular x so we can't estimate $\hat{f}(x)$.



in each point x we do the mean of all the observations and we state $\hat{f}(x)$ equal to that mean.

in this point we don't have any observation so we can't find an estimation

A possible solution is to relax the definition: $\hat{f}(x) = \text{Ave}(Y | X \in N(x))$. Now the estimation of $\hat{f}(x)$ is the average in a neighborhood ($N(x)$):

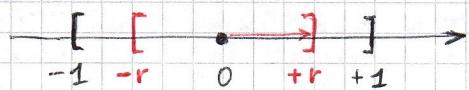


for the estimation of $\hat{f}(x)$ we use all the observations in a neighborhood.

But there is a problem: **CURSE OF DIMENSIONALITY**.

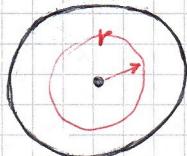
This method fails when p becomes large. Why? To be efficient, the average must be taken in a neighborhood, more precisely in a local neighborhood. When p becomes large the points result to be far away from each other, so we lose the local property.

- p=1 : let's consider $X \sim U([-1, 1]) = U(S^1(1))$ (sphere of radius 1)
How far do we have to go from 0 to capture 10% of the data?



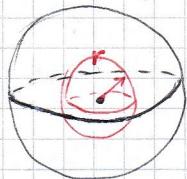
$$0,1 = \frac{S^1(r)}{S^1(1)} = \frac{\pi r}{\pi} = r \Rightarrow \text{to capture 10% of the data we have to move } r = 0,1$$

- p=2 : $X \sim U(S^2(1))$



$$0,1 = \frac{S^2(r)}{S^2(1)} = \frac{\pi r^2}{\pi} = r^2 \Rightarrow \text{to capture 10% of the data we have to move } r = \sqrt{0,1} = 0,31$$

- p=3 : $X \sim U(S^3(1))$



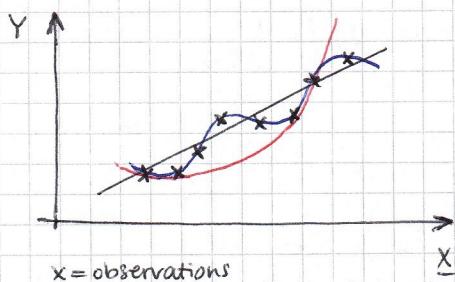
$$0,1 = \frac{S^3(r)}{S^3(1)} = \frac{\frac{4}{3}\pi r^3}{\frac{4}{3}\pi} = r^3 \Rightarrow \text{to capture 10% of the data we need to move } r = \sqrt[3]{0,1} = 0,46$$

- p=1000 : $r=0,99$: we have to travel the 99% of the universe to capture the 10% of the data. In this way we completely lose the sense of "local"
(every time we're solving: $r^p = 0,1$)

so, with p large we have the curse of dimensionality. How do we react?
We have two possibilities:

1. We reduce the dimensionality (we reduce p).
How? We forget part of the information. That's ok but we need a lot of knowledge of the phenomena. We want a DATA DRIVEN REDUCTION so we project the data in a linear subspace (of dim < p) where the variability will be preserved.
(PCA: PRINCIPAL COMPONENT ANALYSIS)
2. We can use a PARAMETRIC MODEL: $Y = f(X) + \epsilon \rightarrow Y = \beta_0 + \beta_1 X_1 + \dots + \beta_p X_p + \epsilon$.
Everything is reduced to estimate $p+1$ parameters: we're choosing the form of the model. By parametrizing f we reduce the dimensionality of the problem. It's not Data Driven (we don't know if it is better to consider, for instance, highs or highs²), that's on us.

Problem of the error: BIAS-VARIANCE TRADE OFF



$$\underline{f(x) = \beta_0 + \beta_1 x}$$

$$\underline{f(x) = \beta_0 + \beta_1 x + \beta_2 x^2}$$

$$\underline{f(x) = \beta_0 + \beta_1 x + \dots + \beta_k x^k}$$

model that goes through all the points

FITTING EXACTLY the data is the worst thing for a predictor since given a x there is variability of y not captured by f . The more we use the data for the model (for the estimate of the parameters) the worst is the model in terms of variability.
(BIAS-VARIANCE trade-off)

We want to judge how good is a model for predictions: $x_0 \xrightarrow{\hat{f}} y_0$
 $(\hat{f} \text{ based on } \mathbb{X})$:

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

this estimate the prediction error given the data \mathbb{X} but with another \mathbb{X} this will change

We want our model to be that good that if we have \hat{f}_1 based on \mathbb{X}_1 and \hat{f}_2 based on \mathbb{X}_2 (where \mathbb{X}_1 and \mathbb{X}_2 are different observations of the same phenomenon)
 $\Rightarrow \hat{f}_1 = \hat{f}_2$.

We need to do an average through all the possible \mathbb{X} :

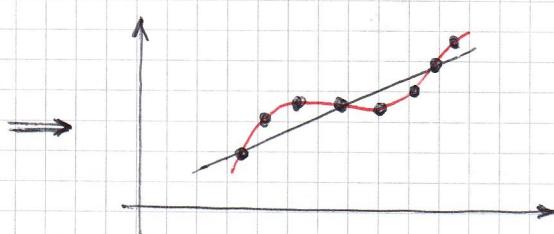
$$\begin{aligned} \mathbb{E}[\mathbb{E}_{\mathbb{X}}[(y_0 - \hat{f}(x_0))^2]] &= \mathbb{E}[(f(x_0) - \hat{f}(x_0))^2 + \text{Var}(\varepsilon_0)] \\ &\stackrel{\substack{\text{depends} \\ \text{on the dataset}}}{=} \mathbb{E}[(f(x_0) - \hat{f}(x_0)) \pm \mathbb{E}[\hat{f}(x_0)]]^2 + \text{Var}(\varepsilon_0) \\ &= \boxed{\mathbb{E}[(f(x_0) - \mathbb{E}[\hat{f}(x_0)])^2]} + \boxed{\mathbb{E}[(\mathbb{E}[\hat{f}(x_0)] - \hat{f}(x_0))^2]} + \boxed{2 \mathbb{E}[(f(x_0) - \mathbb{E}[\hat{f}(x_0)])(\mathbb{E}[\hat{f}(x_0)] - \hat{f}(x_0))]} + \text{Var}(\varepsilon_0) \end{aligned}$$

$$\begin{aligned} (*) &= 2 \mathbb{E}[(f(x_0) - \mathbb{E}[\hat{f}(x_0)]) \underbrace{(\mathbb{E}[\hat{f}(x_0)] - \hat{f}(x_0))}_{\text{number}}] \\ &= 2 (f(x_0) - \mathbb{E}[\hat{f}(x_0)]) (\underbrace{\mathbb{E}[\hat{f}(x_0)] - \mathbb{E}[\hat{f}(x_0)]}_{=0}) \\ &\stackrel{(*)}{=} 0 \end{aligned}$$

$$\begin{aligned} (*) &= \mathbb{E}[(f(x_0) - \mathbb{E}[\hat{f}(x_0)])^2] + \mathbb{E}[(\mathbb{E}[\hat{f}(x_0)] - \hat{f}(x_0))^2] \\ &= (f(x_0) - \mathbb{E}[\hat{f}(x_0)])^2 + \text{Var}(\hat{f}(x_0)) \end{aligned}$$

$$\Rightarrow \mathbb{E}[\mathbb{E}_{\mathbb{X}}[(y_0 - \hat{f}(x_0))^2]] = \underbrace{\text{Var}(\varepsilon_0)}_{\text{variability of the model}} + \underbrace{\text{Var}(\hat{f}(x_0))}_{\text{IRREDUCIBLE variability}} + \underbrace{(f(x_0) - \mathbb{E}[\hat{f}(x_0)])^2}_{:= \text{Bias}^2}$$

(how far is the model from what we want to estimate (on average))

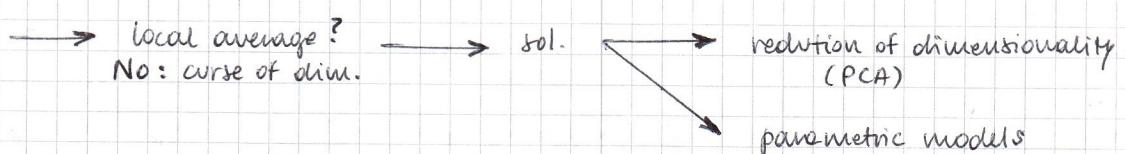


— : fits all the points, reduce the bias but the variability is pretty high (OVERFITTING). If we change 1 observ. then the model changes all

— : linear model : the bias is high but the variability is very low

We want to minimize all \Rightarrow bias-variance trade off.

RECAP: Data: units (\downarrow) \rightarrow predict data: regression funct. \rightarrow estimate of the regression function



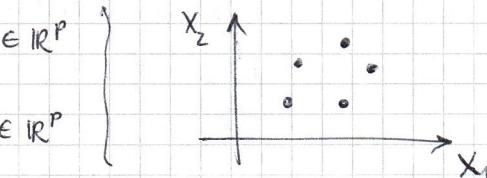
GEOMETRY OF THE DATA

We'll always assume to have:

- n statistical units
- p features / variables

$$\begin{matrix} \mathbf{X} = \\ \text{Data matrix} \\ \text{(or Data frame)} \end{matrix} \quad \begin{bmatrix} x_1 & x_2 & \dots & x_p \\ 1 & [x_{11} & x_{12} & \dots & x_{1p}] \\ 2 & [x_{21} & x_{22} & \dots & x_{2p}] \\ \vdots & \ddots & \ddots & \ddots \\ n & [x_{n1} & x_{n2} & \dots & x_{np}] \end{bmatrix} \in \mathbb{R}^{n \times p}$$

We can look at the data in two ways:

- (1) BY ROWS: $\underline{x}_i^T = [x_{i1} \ x_{i2} \ \dots \ x_{ip}] \in \mathbb{R}^p$
- \vdots
 $\underline{x}_i^T = [x_{i1} \ x_{i2} \ \dots \ x_{ip}] \in \mathbb{R}^p$
- i-th row
- (2) BY COLUMNS: $y_j = [x_{1j} \ x_{2j} \ \dots \ x_{nj}]^T \in \mathbb{R}^n$ = sample from \underline{x}_j
- \vdots
 $y_j = [x_{1j} \ x_{2j} \ \dots \ x_{nj}]^T \in \mathbb{R}^n$ = sample from \underline{x}_j
- j-th column
- (We can represent them as vectors, we'll see why.)
- 
- in \mathbb{R}^n)

1. BY COLUMN

The way of summarizing the information coming from a sample is through the distribution of the sample, and we can summarize the distribution with mean and variance. Let's formally define it. Let's look at \mathbf{X} by columns:

$$\begin{aligned} y_j &= [x_{1j} \ x_{2j} \ \dots \ x_{nj}]^T \\ \Rightarrow \bar{x}_j &= \frac{1}{n} \sum_{i=1}^n x_{ij} \longrightarrow \begin{bmatrix} \bar{x}_1 \\ \bar{x}_2 \\ \vdots \\ \bar{x}_p \end{bmatrix} = \text{mean vector} \end{aligned}$$

MEAN for variable j

The mean is the barycentrum of the sample:

This is only a partial information.

We want to know about the variability too (the dispersion around the mean), so we introduce the variance.

$$\Rightarrow S_{jj} = \frac{1}{n} \sum_{i=1}^n (x_{ij} - \bar{x}_j)^2 \longrightarrow \sqrt{S_{jj}} = \text{STANDARD DEVIATION of } \underline{x}_j$$

VARIANCE for the variable \underline{x}_j

(For the moment we want to divide by n and we will consider dividing by $n-1$ later. When we want the variability of the sample we divide by n . When we consider S_{jj} as an estimator of the variance of the entire population then we divide by $n-1$.)

What about the covariability between two variables?

$$\Rightarrow \text{Cov}(x_j, x_k) = \frac{1}{n} \sum_{i=1}^n (x_{ik} - \bar{x}_k)(x_{ij} - \bar{x}_j) := s_{kj} \quad k, j = 1, \dots, p$$

COVARIANCE between x_j and x_k

$$\text{Cov}(x_i, x_j) = S_{ij} = \frac{1}{n} \sum_{i=1}^n (x_{ij} - \bar{x}_j)^2$$

→ we can introduce a matrix: $S = \begin{bmatrix} S_{11} & S_{12} & \dots & S_{1p} \\ S_{21} & S_{22} & \dots & S_{2p} \\ \vdots & \vdots & \ddots & \vdots \\ S_{p1} & S_{p2} & \dots & S_{pp} \end{bmatrix} \in \mathbb{R}^{p \times p}$, symmetric

COVARIANCE MATRIX
of the sample

$$\Rightarrow r_{kj} = \frac{s_{kj}}{\sqrt{s_{kk} s_{jj}}} = \text{Cor}(x_k, x_j)$$

CORRELATION
between x_i and x_k
($r_{kj} \in [-1, 1]$)

$$\rightarrow r = \begin{bmatrix} r_{11} & \dots & r_{1p} \\ \vdots & \ddots & \vdots \\ r_{p1} & \dots & r_{pp} \end{bmatrix}, \text{symmetric}$$

CORRELATION
MATRIX

(covariance matrix normalized
by the product of std deviation)

These are the basic tools for a dataset.

Example: $X_1 = \text{height of italians}$, $\bar{X}_1 = 1.80 \text{ m}$, $\sqrt{s_{11}} = 0.03 \text{ m}$

$$\Rightarrow [\bar{X}_1 - 3\sqrt{s_{11}}, \bar{X}_1 + 3\sqrt{s_{11}}] = [1.71, 1.86]$$

Whatever is the distribution of X_1 , we can say that (almost) 90% of the population is in the interval (thanks to Chebyshew).

$$\Rightarrow \text{Fr}(\bar{X}_1 - k\sqrt{s_{11}} \leq X_1 \leq \bar{X}_1 + k\sqrt{s_{11}}) \geq 1 - \frac{1}{k^2}$$

CHEBY SHEV
alternatively: $\text{Pr}(|X - \mu| \leq k\sigma) \geq 1 - \frac{1}{k^2}$

→ Thanks to Chebyshew, whenever we have the mean and the standard deviation we can construct an interval. Note that, if we know that X is gaussian then the Chebyshew inequality changes a little (s.t. the frequency of $(|X - \mu| \leq k\sqrt{s_{11}})$ is higher than $1 - \frac{1}{k^2}$).

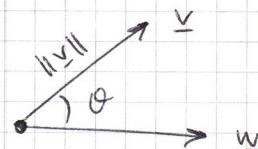
GEOMETRICAL INTERPRETATION

$$\underline{v}, \underline{w} \in \mathbb{R}^n$$

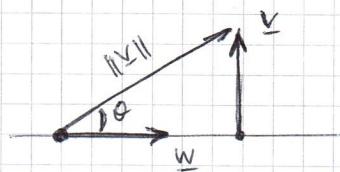
$$\langle \underline{v}, \underline{w} \rangle = \underline{v}^T \underline{w}$$

$$\|\underline{v}\| = \sqrt{\langle \underline{v}, \underline{v} \rangle} = \sqrt{\sum_i v_i^2}$$

$$\cos(\theta) = \frac{\langle \underline{v}, \underline{w} \rangle}{\|\underline{v}\| \cdot \|\underline{w}\|} = \frac{\underline{v}^T \underline{w}}{\sqrt{\underline{v}^T \underline{v}} \sqrt{\underline{w}^T \underline{w}}} = \frac{\underline{v}^T \underline{w}}{\sqrt{(\underline{v}^T \underline{v})(\underline{w}^T \underline{w})}} = \frac{\sum_i v_i w_i}{\sqrt{(\sum_i v_i^2)(\sum_i w_i^2)}}$$



Orthogonal projector:



$P(\underline{v}| \underline{z}(\underline{w})) = P(\underline{v}| \underline{z}(\underline{w})) = \text{projection of } \underline{v} \text{ on } \underline{w}$
 $\underline{z}(\underline{w}) = \{ \underline{z} : \underline{z} = c\underline{w}, c \in \mathbb{R} \}$
 linear space generated by \underline{w}
 (on the linear space generated by \underline{w})

$$P(\underline{v}| \underline{z}(\underline{w})) = \underbrace{\|\underline{v}\| \cos(\theta)}_{\text{length}} \cdot \underbrace{\frac{\underline{w}}{\|\underline{w}\|}}_{\text{direction}} = \|\underline{v}\| \cdot \frac{\underline{v}^T \underline{w}}{\|\underline{v}\| \cdot \|\underline{w}\|} \cdot \frac{\underline{w}}{\|\underline{w}\|}$$

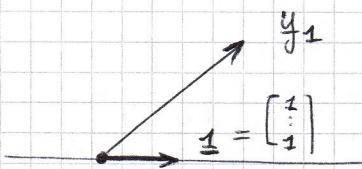
$$= \frac{\underline{v}^T \underline{w}}{\|\underline{w}\|^2} \cdot \underline{w} = \frac{\underline{w}^T \underline{v}}{\underline{w}^T \underline{w}} \cdot \underline{w} = \frac{\underline{w} \underline{w}^T}{\underline{w}^T \underline{w}} \cdot \underline{v}$$

ORTHOGONAL PROJECTOR :
matrix that project any \underline{v} on $\underline{z}(\underline{w})$

→ going back to $\mathbb{X} = [y_1 \dots y_p]$, $y_i \in \mathbb{R}^n$ sample from X ;

CONSIDER 1 VARIABLE:

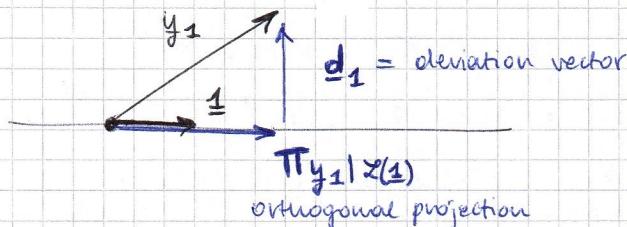
(For instance) $y_1 \in \mathbb{R}^n$: y_1 = sample of heights



$\mathbb{X}(1)$ = space with "no statistic"
(because there is no variability,

all the components are the same)
if $v \in \mathbb{X}(1) \Rightarrow v = c \cdot 1$, $c \in \mathbb{R} \Rightarrow v = \begin{bmatrix} c \\ c \\ \vdots \\ c \end{bmatrix}$

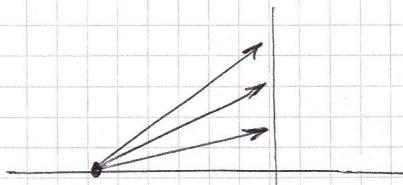
What is the closest vector in $\mathbb{X}(1)$ to y_1 ? (The closer approximation where there is no statistics) → the orthogonal projection:



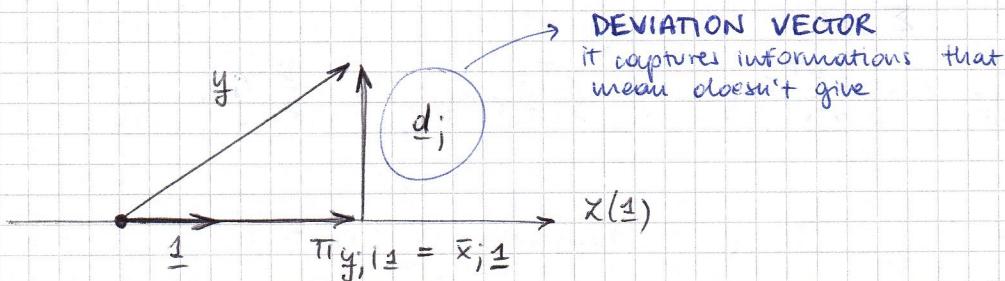
$$\Pi_{y_1 | 1} = \frac{1 \cdot 1^T}{1^T \cdot 1} \cdot y_1 = \frac{\sum_{i=1}^n (y_1)_i}{n} \cdot 1 = \frac{\sum_{i=1}^n x_{i1}}{n} \cdot 1 = \bar{x}_n \cdot 1 = \begin{bmatrix} \bar{x}_n \\ \vdots \\ \bar{x}_n \end{bmatrix}$$

→ the best approximation we can get for the vector of observations when we require that there's no variability = $\bar{x}_n \cdot 1$

But we're not capturing the deviations:



Many vectors can have the same mean (the same projection on $\mathbb{X}(1)$) but all these vectors have different deviation → the mean is not a sufficient summary



$$d_j = y_j - \bar{x}_j \cdot 1 = \begin{bmatrix} x_{1j} - \bar{x}_j \\ x_{2j} - \bar{x}_j \\ \vdots \\ x_{nj} - \bar{x}_j \end{bmatrix}$$

The longer the vector d_j is the worse the approximation is:

$$\|d_j\| = \sqrt{d_j^T d_j} = \sqrt{\sum_{i=1}^n (x_{ij} - \bar{x}_j)^2} = \sqrt{n} \sqrt{\frac{1}{n} \sum_{i=1}^n (x_{ij} - \bar{x}_j)^2} = \sqrt{n} \sqrt{s_{jj}}$$

$$\rightarrow \|d_j\| = \sqrt{n} \sqrt{s_{jj}}$$

this is why s_{jj} talks about variability

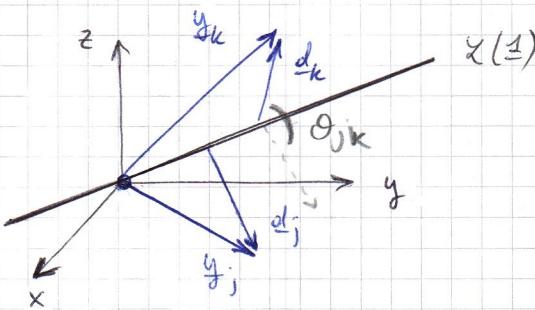
CONSIDER 2 VARIABLES:

$$y_j, y_k$$

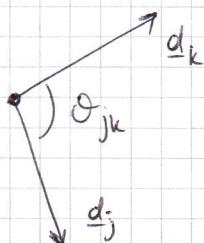
$$y_j = \bar{x}_j \perp + d_j$$

$$y_k = \bar{x}_k \perp + d_k$$

$$\in \mathbb{X}(\perp) \quad \in \mathbb{X}(\perp)^\perp$$



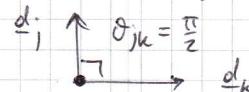
Consider the angle θ_{jk} :



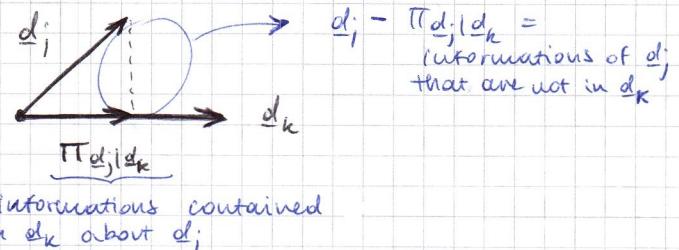
- ① $\theta_{jk} = 0$: $d_j \in \mathbb{X}(d_k) \Rightarrow d_j = \beta d_k, \beta \in \mathbb{R}$
 $\Rightarrow y_j - \bar{x}_j \perp = \beta(y_k - \bar{x}_k \perp)$
 $\Rightarrow y_j = \bar{x}_j \perp + \beta y_k - \beta \bar{x}_k \perp$
 \Rightarrow there is a linear relation between y_j and y_k (\nexists new informations in y_j that is not contained in y_k)

$d_k \rightarrow d_j$: if we know one variable then we know them both

- ② $\theta_{jk} = \frac{\pi}{2}$: \nexists informations that are contained in d_k related to d_j ;



- ③ $0 < \theta_{jk} < \frac{\pi}{2}$: there is some information of d_j in d_k (more exactly $\text{Proj}_{d_k} d_j$):



informations contained in d_k about d_j

We need to compute θ_{jk} to know how much one variable talks about the other:

$$\cos \theta_{jk} = \frac{d_j \cdot d_k}{\|d_j\| \cdot \|d_k\|} = \frac{\sum_{i=1}^n (x_{ij} - \bar{x}_j)(x_{ik} - \bar{x}_k)}{\sqrt{(\sum_{i=1}^n (x_{ij} - \bar{x}_j)^2)(\sum_{i=1}^n (x_{ik} - \bar{x}_k)^2)}} \cdot \frac{1}{\sqrt{n}} \frac{1}{\sqrt{n}}$$

$$= \frac{S_{jk}}{\sqrt{S_{jj} S_{kk}}}$$

$$= \text{Corr}(x_j, x_k) = r_{jk}$$

that's why $r_{jk} \in [-1, 1]$, it's a cosine!

$$\Rightarrow \boxed{\cos(\theta_{jk}) = r_{jk}}$$

- $\theta_{jk} = 0 \Rightarrow \cos(\theta_{jk}) = 1 \Rightarrow r_{jk} = 1 \Rightarrow d_j \in \mathbb{X}(d_k)$
- $\theta_{jk} = \frac{\pi}{2} \Rightarrow \cos(\theta_{jk}) = 0 \Rightarrow r_{jk} = 0 \Rightarrow d_j \perp d_k$

Conclusion: mean and variance (standard deviation)

2. BY ROWS

Now let's look at \underline{X} by rows : $\underline{X} = \begin{bmatrix} \underline{x}_1^T \\ \vdots \\ \underline{x}_n^T \end{bmatrix} \quad \underline{x}_i \in \mathbb{R}^p \quad i=1, \dots, n$

We need some models for these data:

We'll assume that \underline{x}_i are the realizations of a random vector \underline{X}_i and the basic assumption is that :

$$\underline{x}_1, \dots, \underline{x}_n \stackrel{iid}{\sim} \underline{X} \in \mathbb{R}^p$$

We need tools to manage with random vectors:

$$\underline{X} = \begin{bmatrix} X_1 \\ \vdots \\ X_p \end{bmatrix} \rightarrow \text{random variable}$$

When we work with random variables we need to introduce the probability law (\equiv distribution) of \underline{X} (it gives all the informations about the random variable):

$$J_{\underline{X}} : \mathcal{B}(\mathbb{R}^p) \rightarrow [0, 1] : J_{\underline{X}}(B) = P(X \in B) \quad \forall B \in \mathcal{B}(\mathbb{R}^p)$$

law of \underline{X}
Borel sets of \mathbb{R}^p

$$\text{We'll always be in the special case: } J_{\underline{X}}(B) = \int_B f_{\underline{X}}(t) dt \quad \forall B \in \mathcal{B}(\mathbb{R}^p)$$

DENSITY OF \underline{X}

Besides the density we want something that summarize the entire distribution:

$$\rightarrow \mathbb{E}[\underline{X}] = \begin{bmatrix} \mathbb{E}[X_1] \\ \vdots \\ \mathbb{E}[X_p] \end{bmatrix} = \underline{\mu} = \begin{bmatrix} \mu_1 \\ \vdots \\ \mu_p \end{bmatrix}$$

MEAN of \underline{X}
Expected value of \underline{X}

$$\Rightarrow \sigma_{jk} = \mathbb{E}[(X_j - \mu_j)(X_k - \mu_k)] = [\Sigma]_{jk} \quad j, k = 1, \dots, p$$

COVARIANCE
between X_j and X_k

$$\Rightarrow \Sigma = [\sigma_{jk}] \in \mathbb{R}^{p \times p} \text{ matrix} := \text{COVARIANCE MATRIX OF } \underline{X}$$

Notice that:

$$\sigma_{jj} = \mathbb{E}[(X_j - \mu_j)^2] = \text{Var}(X_j)$$

$$\Sigma = \mathbb{E}[(\underline{X} - \underline{\mu})(\underline{X} - \underline{\mu})^T]$$

$$\Sigma = \begin{bmatrix} & & \end{bmatrix} = \begin{bmatrix} & & \end{bmatrix} \cdot \begin{bmatrix} & & \end{bmatrix}$$

$p \times p$ $p \times 1$ $1 \times p$

Define now:

$$V = \begin{bmatrix} \sigma_{11} & & \\ \vdots & \ddots & \\ \sigma_{pp} & & \end{bmatrix} = \text{diag}(\sigma_{11}, \dots, \sigma_{pp}) \quad \sigma_{ii} \geq 0 \quad \forall i$$

$$V^{\frac{1}{2}} = \begin{bmatrix} \sqrt{\sigma_{11}} & & \\ \vdots & \ddots & \\ \sqrt{\sigma_{pp}} & & \end{bmatrix} \xrightarrow{\sigma_{ii} > 0 \quad \forall i} V^{-\frac{1}{2}} = \begin{bmatrix} \frac{1}{\sqrt{\sigma_{11}}} & & \\ \vdots & \ddots & \\ \frac{1}{\sqrt{\sigma_{pp}}} & & \end{bmatrix} ; \quad \underline{\gamma} = V^{-\frac{1}{2}} \Sigma V^{\frac{1}{2}}$$

CORRELATION
MATRIX
 $\in \mathbb{R}^{p \times p}$, symmetric

We want to work with LINEAR COMBINATIONS of the components of \underline{X} :

$$\underline{c} \in \mathbb{R}^p$$

$$\underline{c}^T \underline{X} = c_1 X_1 + c_2 X_2 + \dots + c_p X_p \in \mathbb{R}$$

- $\mathbb{E}[\underline{c}^T \underline{X}] = c_1 \mathbb{E}[X_1] + \dots + c_p \mathbb{E}[X_p]$

$$\stackrel{\downarrow}{=} c_1 \mu_1 + \dots + c_p \mu_p$$

$$\stackrel{\downarrow}{=} \underline{c}^T \underline{\mu}$$

- $\text{Var}(\underline{c}^T \underline{X}) = \underline{c}^T \Sigma \underline{c}$

(Note: $\text{Var}(c_1 X_1 + c_2 X_2) = c_1^2 \text{Var}(X_1) + c_2^2 \text{Var}(X_2) + 2c_1 c_2 \text{Cov}(X_1, X_2)$)

Not only linear combinations, we want to consider K LINEAR COMBINATIONS of \underline{X} :

$C \in \mathbb{R}^{k \times p}$ matrix

$$C = \begin{bmatrix} \underline{c}_1^T \\ \vdots \\ \underline{c}_p^T \end{bmatrix} \quad \underline{c}_i \in \mathbb{R}^p \quad \Rightarrow \quad C\underline{X} = \begin{bmatrix} \underline{c}_1^T \underline{X} \\ \vdots \\ \underline{c}_p^T \underline{X} \end{bmatrix}$$

k -linear combinations
of the components of \underline{X}

- $\mathbb{E}[C\underline{X}] = C\underline{\mu}$
- $\text{Cov}(C\underline{X}) = C \Sigma C^T$ \rightarrow Particular case: $C = \underline{c}^T = [c_1 \dots c_p]$ ($1 \times p$ matrix)
 $\Rightarrow \text{Cov}(C\underline{X}) = \text{Var}(\underline{c}^T \underline{X})$
 $= C \Sigma C^T$
 $= \underline{c}^T \Sigma \underline{c}$
 \Rightarrow coherent with the case before

ESTIMATORS

We have data \underline{X} ($\in \mathbb{R}^{n \times p}$) and the model \underline{X} random vector.
 Can we do inference with this? Can we use the data to estimate $\underline{\mu}$, Σ ?
 Can we say something from the sample that is true for the entire population?

Estimator of $\underline{\mu}$:

$$\underline{\mu} = \begin{bmatrix} \mu_1 \\ \vdots \\ \mu_p \end{bmatrix} \xleftarrow{\text{estimated}} \begin{bmatrix} \bar{x}_1 \\ \vdots \\ \bar{x}_p \end{bmatrix} \xleftarrow{\text{Data}}$$

Why does it work?
 We again consider $X_1, \dots, X_n \stackrel{\text{iid}}{\sim} \underline{X}$ that generate x_1, \dots, x_n (Data):

$$\underline{\bar{x}} = \begin{bmatrix} \bar{x}_1 \\ \vdots \\ \bar{x}_p \end{bmatrix} \xleftarrow{\text{realization of}} \underline{\bar{X}} = \frac{1}{n} \sum_{i=1}^n \underline{X}_i$$

Note: there is a difference between the ESTIMATOR (algorithm) and the ESTIMATE (product (output) of the algorithm). We want to know how good is our estimator, not an estimate (an estimate is good if it's produced by a good estimator).

Prop. Assuming $X_1, \dots, X_n \stackrel{\text{iid}}{\sim} \underline{X}$, $\mathbb{E}[\underline{X}] = \underline{\mu}$, $\text{Cov}(\underline{X}) = \Sigma$:

- $\mathbb{E}[\underline{\bar{X}}] = \underline{\mu}$ (unbiased estimator: on average this estimator is good)
- $\text{Cov}(\underline{\bar{X}}) = \frac{1}{n} \Sigma$

proof:

$$(1) \quad \mathbb{E}[\underline{\bar{X}}] = \mathbb{E} \left[\frac{1}{n} \sum_i \underline{X}_i \right] = \left[\frac{1}{n} \sum_i \mu_1 \right] = \begin{bmatrix} \mu_1 \\ \vdots \\ \mu_p \end{bmatrix}$$

$$\begin{aligned}
 (2) \text{ Cov}(\bar{\mathbf{x}}) &= \mathbb{E}[(\bar{\mathbf{x}} - \mu)(\bar{\mathbf{x}} - \mu)^T] \\
 &= \mathbb{E}\left[\left(\frac{1}{n} \sum_i (\underline{x}_i) - \mu\right)\left(\frac{1}{n} \sum_i (\underline{x}_i) - \mu\right)^T\right] \\
 &= \mathbb{E}\left[\frac{1}{n} \left(\sum_i (\underline{x}_i - \mu)\right) \frac{1}{n} \left(\sum_i (\underline{x}_i - \mu)\right)^T\right] \\
 &= \frac{1}{n^2} \mathbb{E}\left[\sum_{i=1}^n \sum_{k=1}^n (\underline{x}_i - \mu)(\underline{x}_k - \mu)^T\right] \\
 &= \frac{1}{n^2} \sum_{i=1}^n \underbrace{\mathbb{E}[(\underline{x}_i - \mu)(\underline{x}_i - \mu)^T]}_{\mathbb{E}[(\underline{x}_i - \mu)(\underline{x}_k - \mu)^T] = \begin{cases} 0 & i \neq k \\ \Sigma & i = k \end{cases}} \\
 &= \frac{1}{n^2} \sum_{i=1}^n \Sigma = \frac{n \Sigma}{n^2} = \frac{1}{n} \Sigma
 \end{aligned}$$

Estimator of Σ :

$$\mathbf{S} = \frac{1}{n} \sum_{i=1}^n (\underline{x}_i - \bar{\mathbf{x}})(\underline{x}_i - \bar{\mathbf{x}})^T \in \mathbb{R}^{p \times p} \text{ (random matrix)}$$

$$\mathbf{S} = \begin{bmatrix} s_{11} & \dots & s_{1p} \\ \vdots & & \vdots \\ s_{p1} & \dots & s_{pp} \end{bmatrix} \xleftarrow{\text{realization of}} \mathbf{S}$$

Sample covariance matrix

$$\text{Prop. } \mathbb{E}[\mathbf{S}] = \frac{n-1}{n} \Sigma$$

Corollary. $\mathbb{E}\left[\frac{n}{n-1} \mathbf{S}\right] = \Sigma \Rightarrow \frac{n}{n-1} \mathbf{S}$ is an unbiased estimator of Σ

$$\frac{n-1}{n} \mathbf{S} = \frac{n-1}{n} \frac{1}{n} \sum_{i=1}^n (\underline{x}_i - \bar{\mathbf{x}})(\underline{x}_i - \bar{\mathbf{x}})^T = \frac{1}{n-1} \sum_{i=1}^n (\underline{x}_i - \bar{\mathbf{x}})(\underline{x}_i - \bar{\mathbf{x}})^T$$

$$\Rightarrow \text{from now on: } \mathbf{S} = \frac{1}{n-1} \sum_{i=1}^n (\underline{x}_i - \bar{\mathbf{x}})(\underline{x}_i - \bar{\mathbf{x}})^T$$

(we want the covariance of the population, not of the sample,
when we'll need \mathbf{S} (the old)
we'll write \mathbf{S}_n . Note that BIG DATA
uses an n so large that there is no difference)

Focus on \mathbf{S} :

$$\underline{d}_j = \underline{y}_j - \mathbb{P}_{\mathbf{1}} \underline{y}_j = \underline{y}_j - \frac{\mathbf{1} \cdot \underline{1}^T}{\underline{1}^T \underline{1}} \underline{y}_j = \underbrace{\left(\mathbf{I} - \frac{\mathbf{1} \cdot \underline{1}^T}{\underline{1}^T \underline{1}} \right)}_{\text{orthogonal projector on } \mathcal{X}(\underline{1})} \underline{y}_j$$

$\underline{d} = [\underline{d}_1 \dots \underline{d}_p]$ \rightarrow $\underline{d} = \left(\mathbf{I} - \frac{\mathbf{1} \cdot \underline{1}^T}{\underline{1}^T \underline{1}} \right) \mathbf{X}$

deviation of variable j from its mean

DEVIATION MATRIX
($n \times p$ matrix)

$$S = \begin{bmatrix} s_{11} & s_{12} & \dots & s_{1p} \\ \vdots & & & \\ s_{n1} & \dots & s_{np} \end{bmatrix} = \frac{1}{n-1} \begin{bmatrix} \sum_i (x_{i1} - \bar{x}_1) & \sum_i (x_{i1} - \bar{x}_1)(x_{i2} - \bar{x}_2) & \dots \\ \vdots & & \\ \underline{d}_1^T \underline{d}_1 & \underline{d}_2^T \underline{d}_2 & \dots & \underline{d}_1^T \underline{d}_p \\ \vdots & & & \vdots \\ \underline{d}_p^T \underline{d}_p & & & \end{bmatrix}$$

$$\Rightarrow S = \frac{1}{n} \underline{d}^T \underline{d} = \frac{1}{n} \underline{\mathbb{X}}^T \underbrace{\left(I - \frac{1 \ 1^T}{1+1} \right)^T \left(I - \frac{1 \ 1^T}{1+1} \right)}_{\text{since it's an orthogonal projector}} \underline{\mathbb{X}}$$

then it's symmetric ($T^T = T$) and idempotent ($T^2 = T$)

$$\Rightarrow S = \frac{1}{n-1} \underline{\mathbb{X}}^T \left(I - \frac{1 \ 1^T}{1+1} \right) \underline{\mathbb{X}}$$

VARIABILITY IN A MULTIVARIATE JENSE

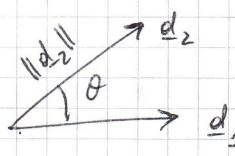
There are two possibilities:

- Generalized variance : $\text{Det}(S)$
- Total variance : $\text{Tr}(S)$

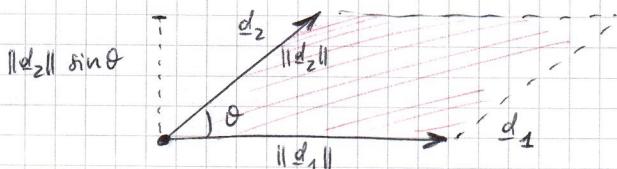
In IR the variability is simply captured by variance.

• Det(S)

Consider the case $p=2$:

$$S = \frac{1}{n-1} \underline{d}^T \underline{d} = \frac{1}{n-1} \begin{bmatrix} \underline{d}_1^T \underline{d}_1 & \underline{d}_1^T \underline{d}_2 \\ \underline{d}_2^T \underline{d}_1 & \underline{d}_2^T \underline{d}_2 \end{bmatrix} = \frac{1}{n-1} \begin{bmatrix} \|\underline{d}_1\|^2 & \|\underline{d}_1\| \cdot \|\underline{d}_2\| \cos \theta \\ \|\underline{d}_1\| \cdot \|\underline{d}_2\| \cos \theta & \|\underline{d}_2\|^2 \end{bmatrix}$$


$$\begin{aligned} \text{Det}(S) &= [\|\underline{d}_1\|^2 \|\underline{d}_2\|^2 - \|\underline{d}_1\|^2 \|\underline{d}_2\|^2 \cos^2(\theta)] \left(\frac{1}{n-1}\right)^2 \\ &\stackrel{!}{=} \|\underline{d}_1\|^2 \|\underline{d}_2\|^2 (1 - \cos^2 \theta) \left(\frac{1}{n-1}\right)^2 \\ &\stackrel{!}{=} \|\underline{d}_1\|^2 \|\underline{d}_2\|^2 \sin^2 \theta \left(\frac{1}{n-1}\right)^2 \end{aligned}$$



$\Rightarrow \text{Det}(S) \propto (\text{Area of the parallelogram } (\underline{d}_1, \underline{d}_2))^2$

$\Rightarrow [\text{Det}(S) \text{ increase} \Rightarrow \text{Area of the parallelogram increase}]$

can increase for two reasons : ① because of the angle (max area with $\theta = \frac{\pi}{2}$) or
② because of the length of $\underline{d}_1, \underline{d}_2$

$\Rightarrow \text{If } \text{Det}(S) = 0 \Rightarrow (\|\underline{d}_1\| \vee \|\underline{d}_2\|) = 0$
 $\qquad \qquad \qquad \theta = 0 \quad (\Rightarrow \text{correlation} = 1)$

$\Rightarrow \text{If } \text{Det}(S) = 0 \text{ it doesn't mean that there is no variability}$

• Tr(S)

$$\text{Tr}(S) = (\|\underline{d}_1\|^2 + \|\underline{d}_2\|^2 + \dots + \|\underline{d}_{n-1}\|^2)$$

The total variance is capturing the sum of the marginal variabilities of the variables.

In general case (general p) : • $\text{Det}(S) \propto \text{Vol}^2(\text{parallelotope } (\underline{d}_1, \dots, \underline{d}_p))$
• $\text{Tr}(S) = \|\underline{d}_1\|^2 + \|\underline{d}_2\|^2 + \dots + \|\underline{d}_p\|^2$

Note: both (generalized variance and total variance) are interesting. One is talking about volume, the other about perimeter. We can't, for instance, decide to use only generalized variance (since it depends also on correlation).

Prop. $\text{Det}(S) = 0 \iff \underline{d}_1, \dots, \underline{d}_p$ are linearly dependent
(i.e. $\exists \underline{c} \neq \underline{0}$ s.t. $\underline{d} \underline{c} = \underline{0} = c_1 \underline{d}_1 + c_2 \underline{d}_2 + \dots + c_p \underline{d}_p$)

proof.

(\Leftarrow) suppose $\underline{d}_1, \dots, \underline{d}_p$ are linearly dependent ($\exists \underline{c} \in \mathbb{R}^p$ s.t. $\underline{d} \underline{c} = \underline{0}$)

$$S = \frac{1}{n-1} \underline{d}^T \underline{d} \implies S \underline{c} = \frac{1}{n-1} \underline{d}^T (\underline{d} \underline{c}) = \underline{0}$$

\implies the columns of S are linearly dependent $\implies \det(S) = 0$

(\Rightarrow) suppose $\det(S) = 0 \implies \exists \underline{c} \neq \underline{0}$ s.t. $S \underline{c} = \underline{0}$
(column must be linearly dependent)

$$\implies \underline{c}^T S \underline{c} = 0 \implies \frac{1}{n-1} \underline{c}^T \underline{d}^T \underline{d} \underline{c} = 0 \implies \frac{1}{n-1} \|\underline{d} \underline{c}\|^2 = 0$$

$\implies \underline{d} \underline{c} = \underline{0} \implies \underline{d}_1, \dots, \underline{d}_p$ are linearly dependent ■

So, $\text{Det}(S) = 0 \implies \underline{d}_1, \dots, \underline{d}_p$ are linearly dependent ($c_1 \underline{d}_1 + \dots + c_p \underline{d}_p = \underline{0} \quad \underline{c} \neq \underline{0}$),
without loss of generality let's assume that $c_1 \neq 0$:

$$\implies \underline{d}_1 = - \sum_{i=2}^p \frac{c_i}{c_1} \underline{d}_i \quad (\underline{d}_k = y_k - \bar{x}_k \underline{1})$$

$$\implies \boxed{y_1 = \bar{x}_1 \underline{1} - \sum_{i=2}^p \frac{c_i}{c_1} (y_i - \bar{x}_i \underline{1})}$$

\implies there is a perfect linear relationship between the var 1 and all the other variables,
so the informations coming from y_1 are useless if we have the ones from y_j , $j=2, \dots, p$

Prop. $\mathbb{X} \in \mathbb{R}^{n \times p}$. If $p \geq n \implies \det(S) = 0$

(If we have too many variables wrt the sample size then it's sure that we'll find a perfect linear relationship between them)

proof.

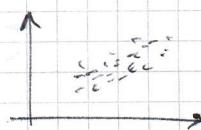
$$\underline{d} = [\underline{d}_1 \dots \underline{d}_p] \quad \underline{d}_i \in \mathbb{R}^n \quad \text{and} \quad \underline{d}_i \in \mathbb{Z}^{\perp}(\underline{1})$$

$\dim(\mathbb{Z}^{\perp}(\underline{1})) = n-1 \implies$ we can have at most $n-1$ linearly independent variables

If $p \geq n \implies \underline{d}_1, \dots, \underline{d}_p$ are linearly dependent $\implies \det(S) = 0$ ■

Example: Raw data :
(received data
≠ data matrix)

x	y
x_1	y_1
\vdots	\vdots
x_n	y_n



We can create the data matrix that we want
(we are free to take what variables we want : x, x^2, \dots)

$$\text{Data matrix } \mathbb{X} = \left[\begin{array}{cccc|c} x_1 & x_1^2 & \dots & x_1^p & y_1 \\ \vdots & \vdots & & \vdots & \vdots \\ x_n & x_n^2 & \dots & x_n^p & y_n \end{array} \right] \quad (*)$$

Consider for instance $p=n-1 \rightarrow \# \text{columns of } \mathbb{X} = p+1 = n$
 $\Rightarrow \det(S) = 0$.

Consider for instance $(*)$ to be linearly independent

$$\Rightarrow y_i = \beta_0 + \beta_1 x_i + \dots + \beta_p x_i^p$$

\Rightarrow perfect fit! Are we happy with that? No, this is an overfitting
 (we always have to think at the bias-variance trade-off)

\Rightarrow this model is terrible for prediction

Necessary condition for $\det(S) > 0 : n \geq p+1$

SPECTRAL DECOMPOSITION OF S

Since S is a $p \times p$ symmetric and real ($s_{ij} \in \mathbb{R}$) matrix

$$\Rightarrow \begin{cases} \exists \lambda_1, \dots, \lambda_p \in \mathbb{R} \\ \exists e_1, \dots, e_p \in \mathbb{R}^p \text{ s.t. } e_j^T e_i = \begin{cases} 0 & i \neq j \\ 1 & i = j \end{cases} \end{cases} \quad (\text{orthonormal system})$$

$$\text{s.t. } S = \sum_{i=1}^p \lambda_i e_i e_i^T \quad (\text{spectral decomposition of a symmetric real matrix : in fact } (\lambda_i, e_i) \text{ are an eigenvalue and eigenvector couple for } S : S e_i = \lambda_i e_i \text{ for } i=1, \dots, p)$$

$$P := [e_1, \dots, e_p], \Lambda = \begin{bmatrix} \lambda_1 & & 0 \\ 0 & \ddots & \\ & & \lambda_p \end{bmatrix} \Rightarrow S = P \Lambda P^T$$

$$\Rightarrow \begin{cases} \det(S) = \prod_{i=1}^p \lambda_i & (\text{generalized variance}) \\ \text{Tr}(S) = \sum_{i=1}^p \lambda_i & (\text{total variance}) \end{cases}$$

Prop. S is positive semi-definite.
 If $\det(S) \neq 0 \Rightarrow S$ is positive definite. $(\lambda_i \geq 0 \forall i)$
 $(\lambda_i > 0 \forall i)$

proof.

We need to prove that $\underline{e}^T S \underline{e} \geq 0 \quad \forall \underline{e} \in \mathbb{R}^p$.

$$\underline{e}^T S \underline{e} = \frac{1}{n-1} \underline{e}^T d d^T \underline{e} = \frac{1}{n-1} \|d\underline{e}\|^2 \geq 0$$

Suppose $\exists \underline{e} \neq 0$ s.t. $\underline{e}^T S \underline{e} = 0$ (semi-def.)

$$\Rightarrow \|d\underline{e}\| = 0 \Rightarrow d\underline{e} = 0 \Rightarrow c_1 d_1 + \dots + c_p d_p = 0$$

$\Rightarrow d_1, \dots, d_p$ are linearly dependent $\Rightarrow \det(S) = 0$

So if $\det(S) \neq 0 \Rightarrow \forall \underline{e} \Rightarrow S$ positive definite $\Rightarrow \det(S) > 0$
 (not only $\neq 0$)

Notations: from now on $S = \sum_{i=1}^p \lambda_i e_i e_i^T, \lambda_1 \geq \lambda_2 \geq \dots \geq \lambda_p \geq 0$

INDUCED DISTANCE (induced from S)

Assume that $\det(S) > 0 \quad (\lambda_1 \geq \lambda_2 \geq \dots \geq \lambda_p > 0)$

$$S = \sum_{i=1}^p \lambda_i e_i e_i^T \Rightarrow S^{-1} = \sum_{i=1}^p \frac{1}{\lambda_i} e_i e_i^T \quad \leftarrow \text{induces a metric on } \mathbb{R}^p$$

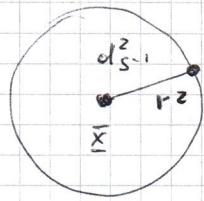
$$x, y \in \mathbb{R}^p : d_{S^{-1}}^2(x, y) = (x - y)^T S^{-1} (x - y)$$

MAHALANOBIS DISTANCE

(It's a distance) \geq in \mathbb{R}^p ,
 we can prove it with spectral decomposition

Let's consider the neighborhood with this distance :

$$\mathcal{E}_{r^2, S^{-1}}(\bar{x}) := \{x \in \mathbb{R}^p : d_{S^{-1}}^2(x, \bar{x}) \leq r^2\}$$



$\mathcal{E}_{r^2, S^{-1}}$ is a circle when we're using Mahalanobis' glasses, but what do we see if we use the euclidean glasses?

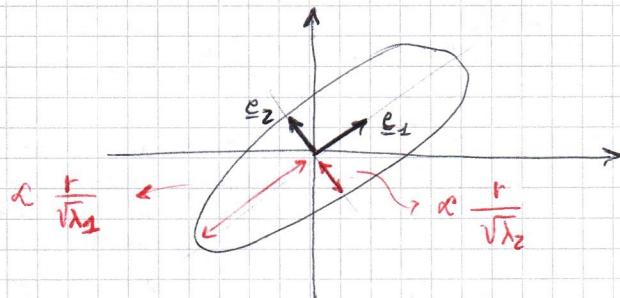
$$\mathcal{E}_{r^2, S^{-1}}(\bar{x}) = \{x \in \mathbb{R}^p : \underbrace{(x - \bar{x})^T S^{-1} (x - \bar{x}) \leq r^2}\}_{\text{in the euclidean sense this is an ellipse centered in } \bar{x}}$$

in the euclidean sense this is an ellipse centered in \bar{x}

Generic case : if we consider a generic positive definite matrix B :

$$B = \sum_{i=1}^p \lambda_i e_i e_i^T$$

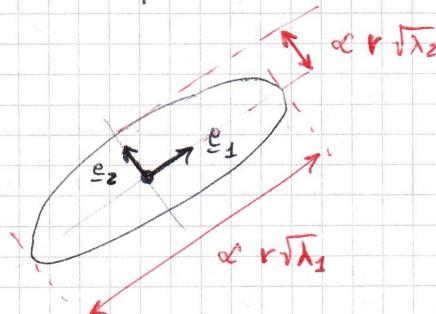
$$x^T B x = r^2$$



Our case :

$$S^{-1} = \sum_{i=1}^p \frac{1}{\lambda_i} e_i e_i^T$$

$$r^2, S^{-1}(x) \longrightarrow$$



$$\text{Volume } (\mathcal{E}_{r^2, S^{-1}}(\bar{x})) = K_p r^p \sqrt{\prod_{i=1}^p \lambda_i} = K_p r^p \sqrt{\det(S)}$$

(in 2D is the area)

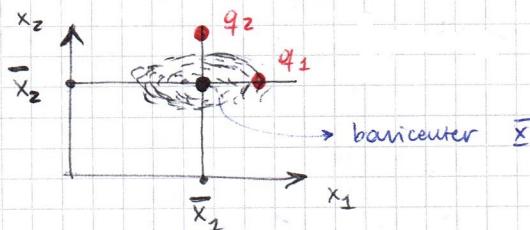
constant depending on the dimension p

the larger the generalized variance, the larger is the neighborhood (given the same r) around the mean

→ the neighborhood depends on the generalized variance.

But, why the Mahalanobis distance is the right one for capturing the distance between points considering the variability of the data?

Example : $p=2$, $S = \begin{bmatrix} s_{11} & 0 \\ 0 & s_{22} \end{bmatrix}$ $s_{11} > s_{22}$



q_1 and q_2 are s.t. $d(\bar{x}, q_1) = d(\bar{x}, q_2)$, so they're equivalently far from \bar{x} considering the euclidean distance. But are they the same from a statistical point of view? The individual q_1 has the same characteristics wrt (x_1, x_2) as individual q_2 ? No, q_1 is in the cloud, q_2 not → the euclidean distance is not a good distance

→ Standardize data first (why? Because we notice that the variance of x_1 and x_2 is not the same. We want to measure distances in terms of statistical units : std deviations, not in other units (meter, kg...))

$$q_1 = \begin{bmatrix} q_1 \\ \bar{x}_2 \end{bmatrix} \rightarrow d(\text{std}(q_1), \text{std}(\bar{x})) = \frac{|q_1 - \bar{x}_1|}{\sqrt{s_{11}}} \quad ①$$

$$q_2 = \begin{bmatrix} \bar{x}_1 \\ q_2 \end{bmatrix} \rightarrow d(\text{std}(q_2), \text{std}(\bar{x})) = \frac{|q_2 - \bar{x}_2|}{\sqrt{s_{22}}} \quad ②$$

→ since $s_{11} > s_{22} \rightarrow ① < ②$

In the general case (general point) :

$$q = \begin{bmatrix} q_1 \\ q_2 \end{bmatrix} \rightarrow d(\text{std}(q), \text{std}(\bar{x})) = \sqrt{\frac{(q_1 - \bar{x}_1)^2}{s_{11}} + \frac{(q_2 - \bar{x}_2)^2}{s_{22}}} \\ = \sqrt{(q - \bar{x})^T S^{-1} (q - \bar{x})}$$

(MAHALANOBIS DISTANCE)

→ the Mahalanobis' distance is just the euclidean distance after we standardize the data. So we'll always consider :

$$\mathcal{E}_{r^2, S^{-1}}(\bar{x}) = \{x \in \mathbb{R}^p : (x - \bar{x})^T S^{-1} (x - \bar{x}) \leq r^2\}$$

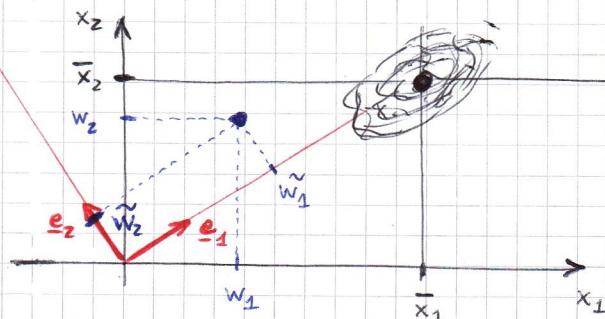
$$\rightarrow \text{Area}(\mathcal{E}_{r^2, S^{-1}}(\bar{x})) \propto r^2 \sqrt{\text{Det}(S)} = r^2 \sqrt{s_{11} s_{22}}$$

So, if $\underbrace{\text{Det}(S)}$ or $\underbrace{\text{Tr}(S)}$ increase
longer the variability → longer the neighborhood (same r)

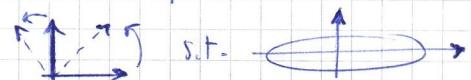
Is this example a special case?

- $p=2 \rightarrow$ easy to extend to general $p \geq 1$
- S diagonal is special, but ...

Let $p \geq 1$ and let S be s.t. $\text{Det}(S) > 0$



We can always find a reference system s.t. S is diagonal w.r.t. that system



We want to introduce the system identified by the eigenvalues of S ($= \sum_{i=1}^p \lambda_i e_i e_i^T$)

OLD SYSTEM	NEW SYSTEM
$w = [x_1 \ x_2 \ \dots \ x_p]^T$	$\tilde{w} = [e_1^T w \ e_2^T w \ \dots \ e_p^T w]^T = p^T w$
$X = \begin{bmatrix} x_1^T \\ \vdots \\ x_n^T \end{bmatrix}$	$\tilde{X} = \begin{bmatrix} (p^T x_1)^T \\ (p^T x_n)^T \end{bmatrix} = \begin{bmatrix} x_1^T p \\ \vdots \\ x_n^T p \end{bmatrix} = X p$
$S = \frac{1}{n-1} X (I - \frac{1}{n} \frac{1^T}{1^T 1}) X^T$	$\hat{S} = \frac{1}{n-1} \tilde{X}^T (I - \frac{1}{n} \frac{1^T}{1^T 1}) \tilde{X}$ $\downarrow \frac{1}{n-1} p^T \tilde{X}^T (I - \frac{1}{n} \frac{1^T}{1^T 1}) \tilde{X} p$ $\downarrow p^T S p = p^T p \Lambda p^T p = \Lambda \quad (\star)$ $= [\lambda_1 \ \dots \ \lambda_p]$

(*) $P^T P = I$, why?

I. P is an orthogonal matrix
II. $P^T P = \begin{bmatrix} e_1^T \\ e_2^T \\ \vdots \\ e_p^T \end{bmatrix} [e_1, \dots, e_p] = \begin{bmatrix} e_1^T e_1 & 1 & & \\ e_1^T e_2 & & 1 & \\ \vdots & & \ddots & \ddots \\ e_1^T e_p & & & 1 \end{bmatrix} = I$

→ In this new system (reference system of eigenvectors) the covariances are simply 0, the only thing left are only the variances.

Take home message: there is always a reference system for which the coordinates are uncorrelated → correlation is not a property of the data.

Do we lose some information using this system? No.

- Generalized variance: $\text{Det}(S) = \prod_{i=1}^p \lambda_i = \text{Det}(\tilde{S}) = \text{Det}(\Lambda) = \prod_{i=1}^p \lambda_i$
- Total variance : $\text{Tr}(S) = \sum_{i=1}^p \lambda_i = \text{Tr}(\tilde{S}) = \text{Tr}(\Lambda) = \sum_{i=1}^p \lambda_i$

Observation: $\|\underline{w}\|^2 = \underline{w}^T \underline{w} = \tilde{\underline{w}}^T \tilde{\underline{w}} = \|\tilde{\underline{w}}\|^2 = \underline{w}^T P P^T \underline{w}$

since we're just changing reference,
the length of the vectors remain the same

$$\Rightarrow \underline{w}^T \underline{w} = \underline{w}^T P P^T \underline{w} \quad (!)$$

PRINCIPAL COMPONENT ANALYSIS (PCA)

Consider the milky way :

~~the linear structure that we see is not the correlation of the stars, it's just the first principal component.~~

Let's refer to the model (not the data) :

$$\underline{X} = \text{random vector in } \mathbb{R}^P : \quad E[\underline{X}] = \underline{\mu}, \quad \text{Cov}(\underline{X}) = \Sigma$$

If we take $\underline{a} \in \mathbb{R}^P$ we can do a linear comb. :

$$\underline{a}^T \underline{X} = a_1 X_1 + \dots + a_P X_P$$

Problem: find \underline{a} s.t. $\text{Var}(\underline{a}^T \underline{X})$ is maximum.

(We want to find the linear combination of the components of the vector s.t. the variability is maximized)

Is the problem well posed?

Suppose we find the solution : $\underline{a} = [10, 0, \dots, 0]^T \Rightarrow \max \text{Var}(\underline{a}^T \underline{X})$ is reached by : $\underline{a}^T \underline{X} = 10 X_1$

But : $10 \underline{a} = [100, 0, \dots, 0]^T \Rightarrow \text{Var}(10 \underline{a}^T \underline{X}) = 100 \text{Var}(\underline{a}^T \underline{X})$

→ \underline{a} cannot be the max → the problem is bad posed

→ Problem: find \underline{a} s.t. $\|\underline{a}\|=1$ and $\text{Var}(\underline{a}^T \underline{X})$ is maximized

Note that :

$$\max_{\underline{a} \in \mathbb{R}^P: \|\underline{a}\|=1} \text{Var}(\underline{a}^T \underline{X}) = \max_{\underline{a} \in \mathbb{R}^P: \|\underline{a}\|=1} \underline{a}^T \Sigma \underline{a} = \max_{\underline{a} \in \mathbb{R}^P} \frac{\underline{a}^T \Sigma \underline{a}}{\|\underline{a}\|^2}$$

$$\text{Var}(\underline{a}^T \underline{X}) = \underline{a}^T \Sigma \underline{a}$$

equivalent since $\left\| \frac{\underline{a}}{\|\underline{a}\|} \right\| = 1$
(we transfer $\|\underline{a}\|=1$ in the objective function)

Lemma: Let $B \in \mathbb{R}^{P \times P}$ positive semidefinite and $B = \sum_{i=1}^P \lambda_i e_i e_i^\top$ its spectral decomposition.
Then:

$$1. \max_{\underline{x} \in \mathbb{R}^P} \frac{\underline{x}^\top B \underline{x}}{\underline{x}^\top \underline{x}} = \lambda_1, \quad \text{argmax}(\dots) = \underline{e}_1$$

$$2. \max_{\substack{\underline{x} \in \mathbb{R}^P \\ \underline{x} \perp \underline{e}_1}} \frac{\underline{x}^\top B \underline{x}}{\underline{x}^\top \underline{x}} = \lambda_2, \quad \text{argmax}(\dots) = \underline{e}_2$$

:

$$P. \max_{\substack{\underline{x} \in \mathbb{R}^P \\ \underline{x} \perp \underline{e}_1 \perp \dots \perp \underline{e}_{p-1}}} \frac{\underline{x}^\top B \underline{x}}{\underline{x}^\top \underline{x}} = \lambda_p = \min_{\underline{x} \in \mathbb{R}^P} \frac{\underline{x}^\top B \underline{x}}{\underline{x}^\top \underline{x}}$$

proof.

$$1. B = \sum_{i=1}^P \lambda_i \underline{e}_i \underline{e}_i^\top = P \Lambda P^\top$$

$$\rightarrow \frac{\underline{x}^\top B \underline{x}}{\underline{x}^\top \underline{x}} = \frac{\underline{x}^\top P \Lambda P^\top \underline{x}}{\underline{x}^\top P P^\top \underline{x}} \quad \forall \underline{x} \in \mathbb{R}^P \quad (*)$$

$$\text{we proved: } \underline{w}^\top \underline{w} = \underline{w}^\top P P^\top \underline{w}$$

$$y := P^\top \underline{x} \implies (*) = \frac{\underline{y}^\top \Lambda \underline{y}}{\underline{y}^\top \underline{y}} = \frac{\sum_{i=1}^P (\lambda_i y_i)^2}{\sum_{i=1}^P (y_i)^2} \leq \lambda_1 \downarrow \frac{\sum_{i=1}^P (y_i)^2}{\sum_{i=1}^P (y_i)^2}$$

since they're ordered

$$\text{If } \underline{x} = \underline{e}_1 : \frac{\underline{x}^\top B \underline{x}}{\underline{x}^\top \underline{x}} = \frac{\underline{e}_1^\top B \underline{e}_1}{\underline{e}_1^\top \underline{e}_1} = \lambda_1 \underline{e}_1^\top \underline{e}_1 = \lambda_1$$

$$\rightarrow \frac{\underline{x}^\top B \underline{x}}{\underline{x}^\top \underline{x}} \leq \lambda_1 \quad \text{but for } \underline{x} = \underline{e}_1 : \frac{\underline{x}^\top B \underline{x}}{\underline{x}^\top \underline{x}} = \lambda_1$$

$$\rightarrow \max_{\underline{x} \in \mathbb{R}^P} \frac{\underline{x}^\top B \underline{x}}{\underline{x}^\top \underline{x}} = \lambda_1 \quad \text{and } \text{argmax}(\dots) = \underline{e}_1$$

$$2. \frac{\underline{x}^\top B \underline{x}}{\underline{x}^\top \underline{x}} = \frac{\underline{x}^\top P \Lambda P^\top \underline{x}}{\underline{x}^\top P P^\top \underline{x}} = \frac{\underline{y}^\top \Lambda \underline{y}}{\underline{y}^\top \underline{y}} = \frac{\sum_{i=2}^P \lambda_i (y_i)^2}{\sum_{i=2}^P (y_i)^2} \leq \lambda_2 \quad \forall \underline{x} \perp \underline{e}_1$$

$$y := P^\top \underline{x} = \begin{bmatrix} \underline{e}_1^\top \\ \vdots \\ \underline{e}_p^\top \end{bmatrix} \underline{x} = \begin{bmatrix} 0 \\ \underline{e}_2^\top \underline{x} \\ \vdots \\ \underline{e}_p^\top \underline{x} \end{bmatrix} \quad \text{because we're looking for } \underline{x} \text{ s.t. } \underline{x} \perp \underline{e}_1$$

$$\underline{x} = \underline{e}_2 \implies \frac{\underline{x}^\top B \underline{x}}{\underline{x}^\top \underline{x}} = \frac{\underline{e}_2^\top B \underline{e}_2}{\underline{e}_2^\top \underline{e}_2} = \lambda_2 \underline{e}_2^\top \underline{e}_2 = \lambda_2^2$$

$$\rightarrow \max_{\substack{\underline{x} \in \mathbb{R}^P, \\ \underline{x} \perp \underline{e}_1}} \frac{\underline{x}^\top B \underline{x}}{\underline{x}^\top \underline{x}} = \lambda_2, \quad \text{argmax}(\dots) = \underline{e}_2$$

[...]

p. (proof of the min)

$$\frac{\underline{x}^\top B \underline{x}}{\underline{x}^\top \underline{x}} = \frac{\underline{y}^\top \Lambda \underline{y}}{\underline{y}^\top \underline{y}} = \frac{\sum_{i=1}^P \lambda_i (y_i)^2}{\sum_{i=1}^P (y_i)^2} \geq \lambda_p$$

$$y := P^\top \underline{x}$$

$$\underline{x} = \underline{e}_p \implies \frac{\underline{x}^\top B \underline{x}}{\underline{x}^\top \underline{x}} = \lambda_p \implies \min_{\underline{x} \in \mathbb{R}^P} \frac{\underline{x}^\top B \underline{x}}{\underline{x}^\top \underline{x}} = \lambda_p$$

Back to PCA :

$$\max_{\underline{a} \in \mathbb{R}^P : \|\underline{a}\|=1} \text{Var}(\underline{a}^T \underline{X}) = \max_{\underline{a} \in \mathbb{R}^P} \frac{\underline{a}^T \Sigma \underline{a}}{\underline{a}^T \underline{a}} = \lambda_1 \quad (\Sigma = \sum_{i=1}^P \lambda_i \underline{e}_i \underline{e}_i^T)$$

$$\arg \max_{\underline{a} \in \mathbb{R}^P : \|\underline{a}\|=1} \text{Var}(\underline{a}^T \underline{X}) = \underline{e}_1$$

Def. PC1 : $\underline{Y}_1 = \underline{e}_1^T \underline{X}$ is the first principal component ($/ \underline{Y}_1 = \underline{e}_1^T (\underline{X} - \mu)$)
 (projection of \underline{X} on the first eigenvector: along this direction we have max variability)

We want to find the second principal component:

Problem : $\boxed{\begin{array}{l} \max_{\underline{a} \in \mathbb{R}^P : \|\underline{a}\|=1} \text{Var}(\underline{a}^T \underline{X}) \\ \text{Cov}(\underline{a}^T \underline{X}, \underline{e}_1^T \underline{X}) = 0 \end{array}} = \boxed{\begin{array}{l} \max_{\underline{a} \in \mathbb{R}^P} \frac{\underline{a}^T \Sigma \underline{a}}{\underline{a}^T \underline{a}} \\ \text{Cov}(\underline{a}^T \underline{X}, \underline{e}_1^T \underline{X}) = 0 \end{array}} (*)$

we want to find another direction
 of max variability but we want the
 projection on this direction to be
 uncorrelated with those we have
 already found

$$\text{Cov}(\underline{a}^T \underline{X}, \underline{b}^T \underline{X}) = ?$$

$$C \in \mathbb{R}^{k \times P} : \text{Cov}(C \underline{X}) = C \Sigma C^T$$

$$C = \begin{bmatrix} \underline{a}^T \\ \underline{b}^T \end{bmatrix} \in \mathbb{R}^{2 \times P} ; \text{Cov}(C \underline{X}) = \text{Cov}\left(\begin{pmatrix} \underline{a}^T \underline{X} \\ \underline{b}^T \underline{X} \end{pmatrix}\right) = \begin{bmatrix} \underline{a}^T \\ \underline{b}^T \end{bmatrix} \Sigma \begin{bmatrix} \underline{a} & \underline{b} \end{bmatrix} = \begin{bmatrix} \underline{a}^T \Sigma \underline{a} & \underline{a}^T \Sigma \underline{b} \\ \underline{b}^T \Sigma \underline{a} & \underline{b}^T \Sigma \underline{b} \end{bmatrix}$$

$$\Rightarrow \text{Cov}(\underline{a}^T \underline{X}, \underline{b}^T \underline{X}) = \underline{a}^T \Sigma \underline{b}$$

Hence:

$$0 = \text{Cov}(\underline{a}^T \underline{X}, \underline{e}_1^T \underline{X}) = \underline{a}^T \Sigma \underline{e}_1 = \lambda_1 \underline{a}^T \underline{e}_1 \Rightarrow \underline{a} \perp \underline{e}_1$$

$$\Rightarrow (*) = \max_{\underline{a} \in \mathbb{R}^P : \underline{a} \perp \underline{e}_1} \frac{\underline{a}^T \Sigma \underline{a}}{\underline{a}^T \underline{a}} = \lambda_2 , \quad \arg \max(\cdot) = \underline{e}_2$$

Def. PC2 : $\underline{Y}_2 = \underline{e}_2^T \underline{X}$ is the second principal component ($/ \underline{Y}_2 = \underline{e}_2^T (\underline{X} - \mu)$)

So, for the j-th principal component:

Problem : $\max_{\underline{a} \in \mathbb{R}^P : \|\underline{a}\|=1} \text{Var}(\underline{a}^T \underline{X}) = \lambda_j , \quad \arg \max(\cdot) = \underline{e}_j$
 $\text{Cov}(\underline{a}^T \underline{X}, \underline{e}_i^T \underline{X}) = 0 \quad \forall i = 1, \dots, j-1$

Def. PCj : $\underline{Y}_j = \underline{e}_j^T \underline{X}$ ($/ \underline{Y}_j = \underline{e}_j^T (\underline{X} - \mu)$) j-th principal component

$$\Rightarrow \underline{Y} = \begin{bmatrix} \underline{Y}_1 \\ \vdots \\ \underline{Y}_p \end{bmatrix} \quad \text{vector of PC's} = \text{P}^T \underline{X} \quad \text{we simply project } \underline{X} \text{ on the eigenvectors of } \Sigma \text{ to have PCs}$$

Characteristics of \underline{Y} :

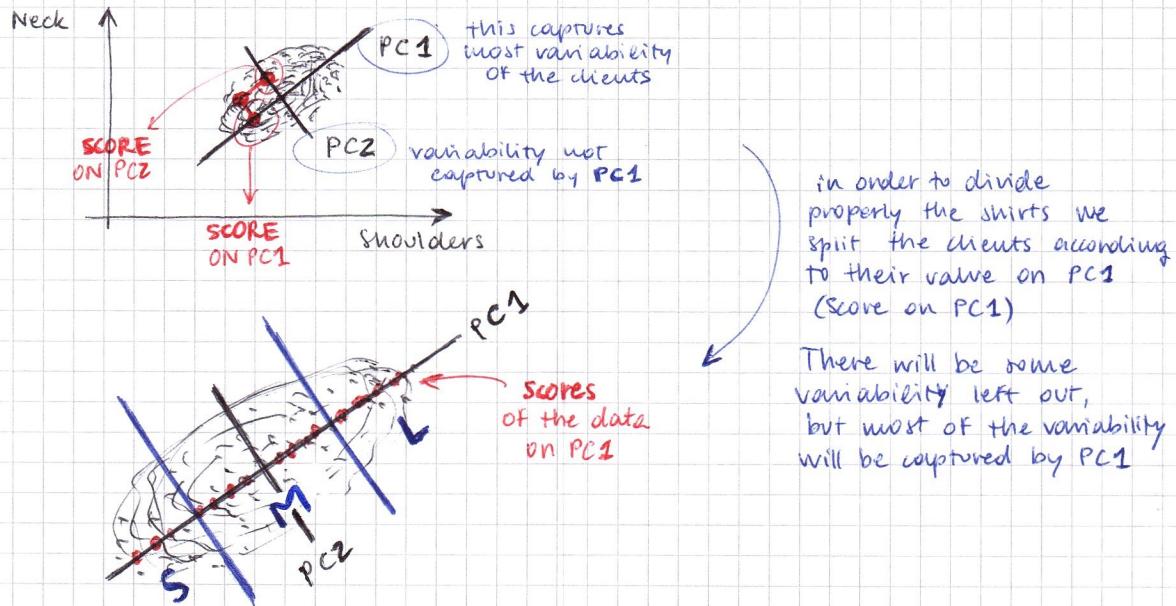
$$\text{Prop.} \bullet \mathbb{E}[\underline{Y}] = \mathbb{E}[P^T \underline{X}] = P^T \mu \quad (\text{If } \underline{Y} = P^T (\underline{X} - \mu) \Rightarrow \mathbb{E}[\underline{Y}] = 0)$$

$$\bullet \text{Cov}(\underline{Y}) = \text{Cov}(P^T \underline{X}) = P^T \Sigma P = P^T P \Lambda P^T P = \Lambda$$

$$\Rightarrow \begin{cases} \text{Cov}(\underline{Y}_i, \underline{Y}_j) = 0 & i \neq j \\ \text{Var}(\underline{Y}_i) = \lambda_i & \forall i \end{cases} \quad i, j = 1, \dots, p$$

- Observations :
- there is no correlation between the coordinates of the r.v. \mathbf{Y}
 - we have ordered the components : the first component is the one with the larger variability, the second is the second larger and so on. We can capture most of the variability with the first components (wrt this system), we can forget about the last components because they express small variability.

Example : We want to industrialize the process of shirts making, so we want to categorize them in groups (S, M, L). How do we do this?



So, what is Y_1 ? It's a linear combination of Neck and Shoulders :

$$Y_1 = X_1 e_{11} + X_2 e_{21} \quad (\text{if we want to know in which category we are, we compute } Y_1)$$

Observations :

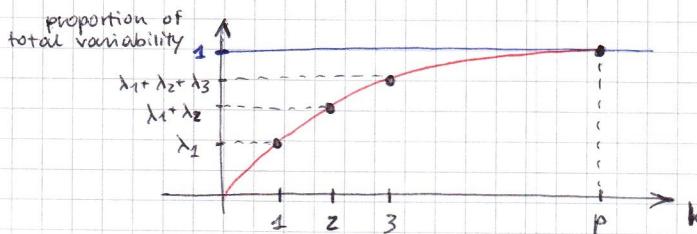
- Are we losing some variability looking at data in this system?

$$\text{Generalized variance } (\mathbf{Y}) = \text{Det}(\Lambda) = \prod_{i=1}^p \lambda_i \quad (= \text{Det}(\Sigma)) \\ = \text{generalized var}(\mathbf{X})$$

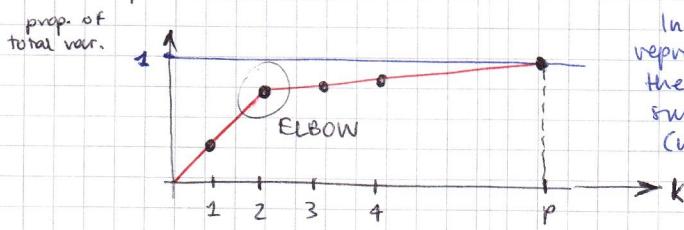
$$\text{Total variance } (\mathbf{Y}) = \text{Tr}(\Lambda) = \sum_{i=1}^p \lambda_i \quad (= \text{Tr}(\Sigma) = \text{total variance } (\mathbf{X}))$$

\Rightarrow we're not losing informations.

- $\text{Var}(Y_1) = \lambda_1 \Rightarrow Y_1 \text{ is capturing } \frac{\lambda_1}{\sum_i \lambda_i} \text{ of total variance}$
- $\text{Var}(Y_2) = \lambda_2 \Rightarrow Y_1 \text{ and } Y_2 \text{ are capturing } \frac{\lambda_1 + \lambda_2}{\sum_i \lambda_i} \text{ of total var.}$



Usually it looks like :



In this case we chose $k=2$ to represent the data. Why? Because the marginal gain get smaller and smaller (not worthy).

When we should stop? 80% / after an elbow (thumb rules)

PCA (continued)

\underline{X} random vector : $E[\underline{X}] = \mu$, $Cov(\underline{X}) = \Sigma$

$\Sigma = \sum_{i=1}^p \lambda_i \underline{e}_i \underline{e}_i^T$ spectral decomposition (s.t. $\lambda_1 \geq \lambda_2 \geq \dots \geq \lambda_p \geq 0$)

$$Y_i = \underline{e}_i^T \underline{X}$$

i-th principal component

$$(Y_i = \underline{e}_i^T (\underline{X} - \mu))$$

$$Y_i = e_{1i} X_1 + e_{2i} X_2 + \dots + e_{pi} X_p$$

↓
:= LOADINGS
(WEIGHTS)

: How to interpret them?

Prop. $\text{Corr}(Y_i, X_k) = \frac{e_{ki} \sqrt{\lambda_i}}{\sqrt{\sigma_{kk}}} \quad i, k = 1, \dots, p$

proof-

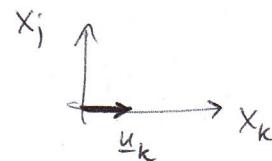
$$\text{Corr}(Y_i, X_k) = \frac{\text{Cov}(Y_i, X_k)}{\sqrt{\lambda_i \sigma_{kk}}}$$

$$\text{Cov}(Y_i, X_k) = \text{Cov}(\underline{e}_i^T \underline{X}, \underline{u}_k^T \underline{X})$$

$$\underline{u}_k = [0, 0, \dots, 0, \underset{k}{1}, 0, \dots, 0]^T \quad (\underline{u}_k^T \underline{X} = X_k)$$

$$\text{Cov}(Y_i, X_k) = \text{Cov}(\underline{e}_i^T \underline{X}, \underline{u}_k^T \underline{X}) = \underline{e}_i^T \sum \underline{u}_k = \underline{u}_k^T \sum \underline{e}_i = \lambda_i \underline{u}_k^T \underline{e}_i$$

$$\text{Corr}(Y_i, X_k) = \frac{\lambda_i e_{ki}}{\sqrt{\lambda_i \sigma_{kk}}} = \frac{\sqrt{\lambda_i}}{\sqrt{\sigma_{kk}}} e_{ki} \quad \text{k-th component of the i-th eigenvector}$$



Maybe we don't want to compute the PCA wrt the original variables, but first we standardize the variables.

$$\underline{X}$$

$$\underline{Z} = V^{-\frac{1}{2}} (\underline{X} - \mu) = \left[\frac{X_1 - \mu_1}{\sqrt{\sigma_{11}}}, \dots, \frac{X_p - \mu_p}{\sqrt{\sigma_{pp}}} \right]^T$$

$$\left(\begin{array}{l} V = \begin{bmatrix} \sigma_{11} & \cdots & \phi \\ \phi & \cdots & \sigma_{pp} \end{bmatrix} \\ \Sigma = [\sigma_{ij}] \end{array} \right)$$

$$E[\underline{Z}] = 0$$

$$\text{Cov}(\underline{Z}) = V^{-\frac{1}{2}} \Sigma V^{-\frac{1}{2}} = P$$

CORRELATION MATRIX (covariance matrix of the standardized variables)

$$\rightarrow \rho = \sum_{i=1}^p \lambda_i \underline{e}_i^T \underline{e}_i$$

λ_i eigenvalue of P
 \underline{e}_i eigenvectors of P

why?
so we have a way to compare variables

$$Y_i = \underline{e}_i^T \underline{Z} = \underline{e}_i^T V^{-\frac{1}{2}} (\underline{X} - \mu)$$

Observations:

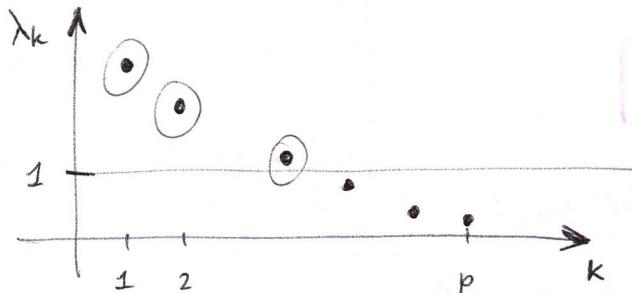
$$1. \sum_i^p \text{Var}(Y_i) = \text{Tr}(P) = \sum_{i=1}^p \text{Var}(Z_i) = p$$

$$2. \text{Corr}(Y_i, Z_k) = \frac{e_{ki} \sqrt{\lambda_i}}{\sqrt{1}} = e_{ki} \sqrt{\lambda_i}$$

3. Proportion of total variability explained by the first k -components:

$$\frac{\sum_{i=1}^k \lambda_i}{p}$$

4. $\bar{\lambda} = \frac{\sum_{i=1}^p \lambda_i}{p} = \frac{\text{Tr}(\Sigma)}{p} = \frac{p}{p} = 1$ (average value for the eigenvalues for Σ)
 (Rule of thumb: Chose Y_i if $\lambda_i > 1$)



If we know that the average variability is 1 (and we're looking for the max variability λ_1) we select all the λ_i above the average.

Example:

Why do PCA on Σ is different from PCA on ρ ?

$$\Sigma = \begin{bmatrix} 1 & 4 \\ 4 & 100 \end{bmatrix} \quad \rho = \begin{bmatrix} 1 & 0.4 \\ 0.4 & 1 \end{bmatrix}$$

Why doing PCA on Σ is not the right thing? ($100 \gg 1$)

$$\Sigma = \text{Cov}([X_1 \ X_2])$$

Suppose that X_1 is in cm, X_2 is in mm, even if the range of variability maybe is the same, the $\text{Var}(X_2)$ is 100 higher than $\text{Var}(X_1)$ simply for the units of measure!

PCA of Σ :

eigenvalues: $\lambda_1 = 100.16$
 & vectors $\lambda_2 = 0.84$

$$e_1 = [0.04, 0.999]^T$$

$$e_2 = [0.999, -0.04]^T$$

$$Y_1 = 0.04 X_1 + 0.999 X_2$$

basically X_2

$$Y_2 = 0.999 X_1 - 0.04 X_2$$

basically X_1

PCA of ρ :

$$\lambda_1 = 1.4 \quad e_1 = [\frac{1}{\sqrt{2}}, \frac{1}{\sqrt{2}}]^T$$

$$\lambda_2 = 0.6 \quad e_2 = [\frac{1}{\sqrt{2}}, -\frac{1}{\sqrt{2}}]^T$$

$$Y_1 = \frac{1}{\sqrt{2}} Z_1 + \frac{1}{\sqrt{2}} Z_2$$

$$= \frac{1}{\sqrt{2}} \left(\frac{X_1 - \mu_1}{\sqrt{\sigma_{11}}} + \frac{X_2 - \mu_2}{\sqrt{\sigma_{22}}} \right) \quad \neq Y_1 \text{ of } \Sigma \text{'s PCA}$$

$$Y_2 = \frac{1}{\sqrt{2}} Z_1 - \frac{1}{\sqrt{2}} Z_2$$

$$= \frac{1}{\sqrt{2}} \left(\frac{X_1 - \mu_1}{\sqrt{\sigma_{11}}} - \frac{X_2 - \mu_2}{\sqrt{\sigma_{22}}} \right) \quad \neq Y_2 \text{ of } \Sigma \text{'s PCA}$$

and in this case the $\text{Var}(X_2)$ completely masks $\text{Var}(X_1)$!

$\text{Var}(X_2)$ is so big that if we are looking for the direction where there is maximum variability the algorithm tells us that it is the direction of X_2 , same for the direction of minimum variability (X_1)

Now we move on Data. How do we use data to perform PCA?

Usually μ and Σ are unknown but we have data:

$$\bar{X} = \begin{bmatrix} \bar{x}_1 \\ \vdots \\ \bar{x}_n \end{bmatrix} \quad \bar{x}_i \text{ realizations of } \underline{X},$$

$$\underline{X}_1, \dots, \underline{X}_n \stackrel{iid}{\sim} \underline{X}$$

- $\text{PCA}(\Sigma) \neq \text{PCA}(\rho)$
- if we think that there are variables with very different variabilities (maybe because of the units of measure or maybe because of the phenomena) then it will be better $\text{PCA}(\rho)$

We want to use data to estimate μ and Σ

Σ estimated by S

μ estimated by \bar{X} .

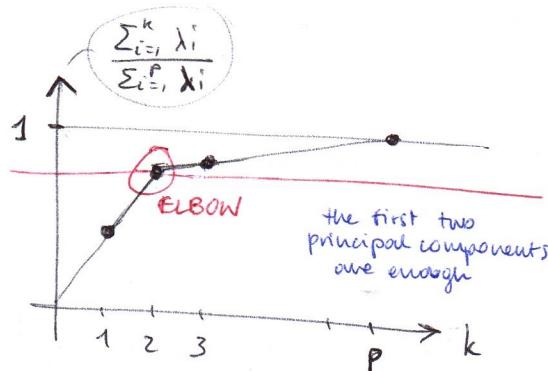
Perform PCA on S . (and not on Σ , since we don't know it)

$$S = \sum_{i=1}^p \lambda_i e_i e_i^T e_i \implies \text{PCA: } Y_i \text{ projection on } e_i$$

$$\underline{x}_i \xrightarrow{\text{PCA}} y_i = \begin{pmatrix} e_1^T x_i \\ \vdots \\ e_p^T x_i \end{pmatrix} \xrightarrow{\text{scores on the first principal component}}$$

$$\underline{X}_{\text{data}} = \begin{pmatrix} x_1^T \\ \vdots \\ x_n^T \end{pmatrix} \xrightarrow{\text{PCA}} \begin{pmatrix} y_1^T \\ \vdots \\ y_n^T \end{pmatrix} = \begin{pmatrix} Y_1 & Y_2 & \cdots & Y_p \\ y_{11} & y_{12} & \cdots & y_{1p} \\ y_{21} & y_{22} & \cdots & y_{2p} \\ \vdots & \vdots & \ddots & \vdots \\ y_{n1} & y_{n2} & \cdots & y_{np} \end{pmatrix}$$

$$\frac{\sum_{i=1}^k \lambda_i}{\sum_{i=1}^p \lambda_i} \geq \text{threshold (0.8)} \implies$$



We can get rid of part of the matrix (and we know that by doing this we're just missing 20% of the variability of the data)

SCREE PLOTS

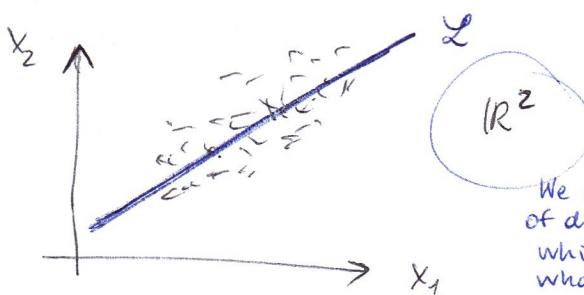
Observation 1: PCA for categorical variables is called CORRESPONDENCE ANALYSIS and it's performed on the table of joint frequencies among variables. (CONTINGENCY TABLES)

Observation 2: look also at the smaller λ_i : if $\lambda_p \approx 0$ it means that \exists linear relationship between x_1, \dots, x_k

OPTIMAL ORTHONORMAL BASIS

$$x_1, \dots, x_n \in \mathbb{R}^p \text{ (data)}$$

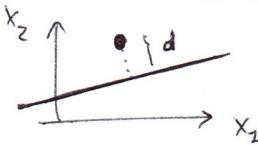
Problem: find a linear space V of dimension $k \leq p$ which best approximate the data, i.e. s.t. V is "closest" to the data



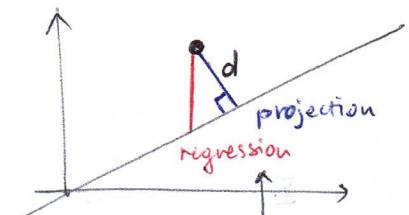
Regression line? NO

We want to find a linear space of dimension 1 (since we're in \mathbb{R}^2) which is the closest to the data. What is this linear space? The regression line? NO

What is "close"?



the regression line minimize the distance between what we're trying to predict and the observation (vertical distance). the projection is minimizing the distance between points and line



that's why they're different
($d \neq \text{regression}$)

uses the axes (the regression),
the projection doesn't
depend on the system

(on the way we're looking on the data,
it's a property of the data)

We want to identify γ :

consider $\gamma_1, \dots, \gamma_k :=$ orthonormal basis spanning \mathcal{L}

$$\mathcal{L} = \text{span}(\gamma_1, \dots, \gamma_k)$$

Problem: find $\gamma_1, \dots, \gamma_k$ orthonormal basis s.t.

$$(d^2 :=) \|(\underline{x}_i - \bar{\underline{x}}) - \sum_{j=1}^k \gamma_j \gamma_j^T (\underline{x}_i - \bar{\underline{x}})\|^2 = \left(\text{distance between the rigid datum (after centering) and the projection on the linear space generated by } \gamma_1, \dots, \gamma_k \right)^2$$

projecting on γ_j

$$\rightarrow \sum_{i=1}^n \|(\underline{x}_i - \bar{\underline{x}}) - \sum_{j=1}^k \gamma_j \gamma_j^T (\underline{x}_i - \bar{\underline{x}})\|^2 \text{ is minimum } (*)$$

Let's define: $\underline{v}_i = \underline{x}_i - \bar{\underline{x}}$

$$\begin{aligned} (*) &= \|\underline{v}_i - \sum_{j=1}^k \gamma_j \gamma_j^T \underline{v}_i\|^2 = \left(\underline{v}_i - \sum_{j=1}^k \gamma_j \gamma_j^T \underline{v}_i \right)^T \left(\underline{v}_i - \sum_{j=1}^k \gamma_j \gamma_j^T \underline{v}_i \right) \\ &= \underline{v}_i^T \underline{v}_i - \sum_{j=1}^k \underline{v}_i^T \gamma_j \gamma_j^T \underline{v}_i - \sum_{j=1}^k \underline{v}_i^T \gamma_j \gamma_j^T \underline{v}_i + \left(\sum_{j=1}^k \gamma_j \gamma_j^T \underline{v}_i \right)^T \left(\sum_{t=1}^n \gamma_t \gamma_t^T \underline{v}_i \right) \\ &= \underline{v}_i^T \underline{v}_i - \sum_{j=1}^k (\gamma_j^T \underline{v}_i)^2 - \sum_{j=1}^k (\gamma_j^T \underline{v}_i)^2 + \sum_{j=1}^k (\gamma_j^T \underline{v}_i)^2 \quad \gamma_j \cdot \gamma_t = \begin{cases} 0 & j \neq t \\ 1 & j = t \end{cases} \\ &= \underline{v}_i^T \underline{v}_i - \sum_{j=1}^k (\gamma_j^T \underline{v}_i)^2 \end{aligned}$$

Problem:

Find $\gamma_1, \dots, \gamma_k$ o.b. s.t.

$$\sum_{i=1}^n (\underline{v}_i^T \underline{v}_i - \sum_{j=1}^k (\gamma_j^T \underline{v}_i)^2) \text{ is min}$$

\Leftrightarrow

$$\left[\sum_{i=1}^n \sum_{j=1}^k (\gamma_j^T \underline{v}_i)^2 \right] \text{ is max } (*)$$

$$\begin{aligned}
 (\star) &= \sum_{j=1}^k \left(\sum_{i=1}^n \underline{\gamma}_j^\top \underline{x}_i \underline{x}_i^\top \underline{\gamma}_j \right) = \sum_{j=1}^k \underline{\gamma}_j^\top \underbrace{\left(\sum_{i=1}^n (\underline{x}_i - \bar{\underline{x}})(\underline{x}_i - \bar{\underline{x}})^\top \right)}_{= (n-1)S} \underline{\gamma}_j \\
 &= (n-1) \sum_{j=1}^k \underline{\gamma}_j^\top S \underline{\gamma}_j
 \end{aligned}$$

$$k=1 : \max_{\underline{\gamma}: \|\underline{\gamma}\|=1} \underline{\gamma}^\top S \underline{\gamma} = \lambda_1 \quad , \quad \text{arg max} = \underline{e}_1$$

by induction on k : $\Rightarrow \underline{\gamma}_1 = \underline{e}_1, \dots, \underline{\gamma}_k = \underline{e}_k$
s.t. $S = \sum \lambda_i \underline{e}_i \underline{e}_i^\top$

CONCLUSION

so if we're looking for a linear space that is close to the data (in terms of minimizing the dist.) this linear space of dimension k is given by the first k principal components (direction)

Observation: What is the error of approximation?

$$\begin{aligned}
 \sum_{i=1}^n d_{\text{euclidean}}^2(\underline{x}_i - \bar{\underline{x}}, \sum_{j=1}^k \underline{e}_j \underline{e}_j^\top (\underline{x}_i - \bar{\underline{x}})) &= \\
 &= \boxed{\sum_{i=1}^n (\underline{x}_i - \bar{\underline{x}})^\top (\underline{x}_i - \bar{\underline{x}})} - \boxed{(n-1) \sum_{j=1}^k \underline{e}_j^\top S \underline{e}_j}
 \end{aligned}$$

using what we did before with \underline{y}_i & co.

$\text{Tr}(a) = a$ if $a \in \mathbb{R}$

$$\begin{aligned}
 (\star) &\stackrel{\downarrow}{=} \text{Tr} \left(\sum_{i=1}^n (\underline{x}_i - \bar{\underline{x}})^\top (\underline{x}_i - \bar{\underline{x}}) \right) = \sum_{i=1}^n \text{Tr} ((\underline{x}_i - \bar{\underline{x}})^\top (\underline{x}_i - \bar{\underline{x}})) \\
 &\stackrel{\text{Tr}(ABC) = \text{Tr}(CAB) = \text{Tr}(BCA)}{=} \sum_{i=1}^n \text{Tr} ((\underline{x}_i - \bar{\underline{x}})(\underline{x}_i - \bar{\underline{x}})^\top) \\
 &\stackrel{\downarrow}{=} \text{Tr} \left(\sum_{i=1}^n (\underline{x}_i - \bar{\underline{x}})(\underline{x}_i - \bar{\underline{x}})^\top \right) \\
 &\stackrel{\downarrow}{=} \text{Tr}((n-1)S) \\
 &\stackrel{\downarrow}{=} (n-1) \sum_{i=1}^k \lambda_i
 \end{aligned}$$

$$\Rightarrow ((\star)) = (n-1) \sum_{j=1}^k \underline{e}_j^\top S \underline{e}_j = (n-1) \sum_{j=1}^k \lambda_j \underline{e}_j^\top \underline{e}_j = (n-1) \sum_{j=1}^k \lambda_j$$

$$\begin{aligned}
 \Rightarrow \sum_{i=1}^n d_{\text{eucl.}}^2(\underline{x}_i - \bar{\underline{x}}, \sum_{j=1}^k \underline{\gamma}_j \underline{\gamma}_j^\top (\underline{x}_i - \bar{\underline{x}})) &= \\
 &= (n-1) \sum_{j=1}^p \lambda_j - (n-1) \sum_{j=1}^k \lambda_j = (n-1) \sum_{j=k+1}^p \lambda_j
 \end{aligned}$$

sum of the λ_j that we're leaving out

error of approximation

Extensions:

- **ICA** : INDEPENDENT COMPONENT ANALYSIS
- **NON-LINEAR** dimensional reduction

the directions are not \perp , they're stochastically independent