

Data Science Course (ATI) Notes on 1st Part

S. MARCIAL-DAMER

2019

Chapter ① Linear Algebra

- Main Concepts:
- what is a vector / vector space?
 - what does the scalar product represent?
 - why is projection an important concept?
 - what is critical about eigen-things?

I. Vector spaces and vectors

Given a field k (\mathbb{R} or \mathbb{C}) and a set V , a vector space is using the addition on V and the external product (s-multiplication) to transform vectors (slide 4). The axioms to be satisfied by the addition and the s-multiplication are generally enumerated as CANI-ADDU.

The important point is that one can compute linear combinations of vectors E.g:

$$\vec{w} = \sum_{i=1}^n \alpha_i \vec{v}_i \quad \alpha_i \in k \quad \vec{v}_i \in V$$

In particular

$$\vec{g} = \sum_{i=1}^n \frac{1}{n} \vec{v}_i = \frac{1}{n} \sum_{i=1}^n \vec{v}_i$$

In turn this defines the idea of a set of linearly independent vectors (slide 6)

If $\vec{\alpha}_i$'s such that $\sum_{i=1}^n \alpha_i \vec{v}_i = 0$ and $\exists j / \alpha_j \neq 0$
 $\Rightarrow \sum_{\substack{i=1 \\ i \neq j}}^n \vec{v}_i + \alpha_j \vec{v}_j = 0 \Rightarrow \vec{v}_j = \frac{1}{\alpha_j} \sum_{\substack{i=1 \\ i \neq j}}^n \alpha_i \vec{v}_i$ dependent

→ linear independent vectors cannot be made up from others

Bases: this notion of linear independence induces that
of a (Hamel) basis:

A set of linearly independent vectors $B = \{e_1, \dots, e_n\}$
is a basis of V iff any vector $v \in V$ can be
constructed from the vectors in B . i.e.:

$$\exists \alpha_1, \dots, \alpha_n \in k \quad v = \sum_{i=1}^n \alpha_i e_i$$

The number $n \in \mathbb{N}$ of basis vectors e_i is the dimension
of the vector space. n is both the minimum number
of vectors necessary to represent any $v \in V$ and it is
also a maximal number of linearly independent
vectors in V .

In turn this defines the notion of components as the
values $\alpha_i \in k$, which completely characterize v in B
this is why any vector in V can be identified
to its components in B , the "vector" of k^n
 $(\alpha_1, \dots, \alpha_n)^T$ as a representation of v .

Hence, clearly the same vector v may have different
representations, depending on the choice of the basis.

Given 2 bases $B = \{e_1, \dots, e_n\}$ and $B' = \{e'_1, \dots, e'_n\}$,
every e_j can be written in B' as

$$e_j = \sum_{i=1}^n p_{ij} e'_i$$

Hence if we have a vector $v \in V$,

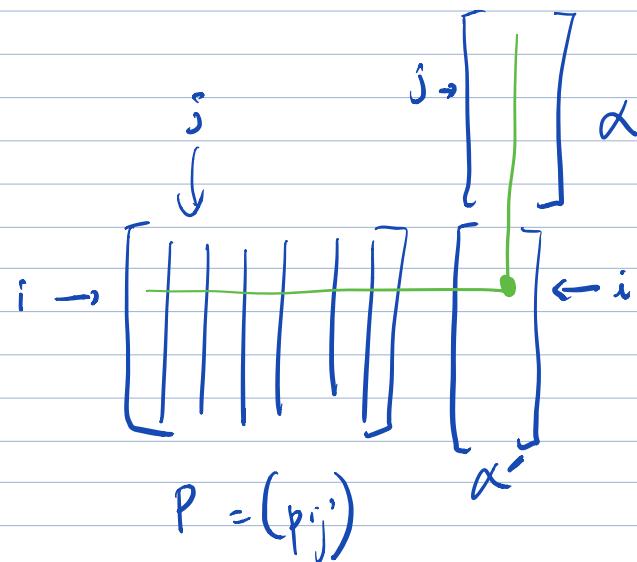
in B :

$$v = \sum_{j=1}^n \alpha_j e_j = \sum_{j=1}^n \alpha_j \sum_{i=1}^n p_{ij} e'_i = \sum_{i=1}^n \left[\sum_{j=1}^n \alpha_j p_{ij} \right] e'_i$$

in B' :

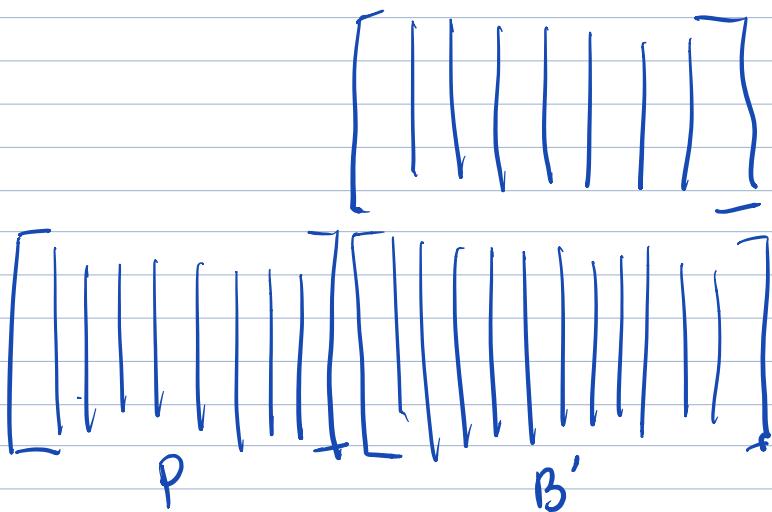
$$v = \sum_{i=1}^n \alpha'_i e'_i$$

$$\Rightarrow \alpha'_i = \sum_{j=1}^n p_{ij} \alpha_j$$



\Rightarrow la matrice de passage P qui transforme les coordonnées de v dans B (α) en coordonnées de v dans B' (α')
à comme colonnes
 p_{ij} les coordonnées de e_j dans B'

$$\Rightarrow \alpha' = P \alpha$$



More generally $B' = PB$ (matrix form)

2 - Scalar product

This is a structure added to the vector space.

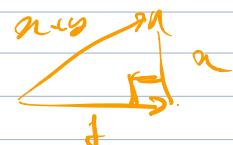
It defines :

- (Axioms 8-9)
- orthogonality : $F^\perp := \{x \in E / \langle x, y \rangle = 0 \forall y \in F\}$
 - a norm : $\|x\|^2 = \langle x, x \rangle$

This is related by the Pythagoras theorem.

If x and y are orthogonal then

$$\|x+y\|^2 = \|x\|^2 + \|y\|^2$$



Since

$$\|x+y\|^2 = \langle x+y, x+y \rangle = \langle x, x \rangle + \langle y, y \rangle + 2 \langle x, y \rangle$$

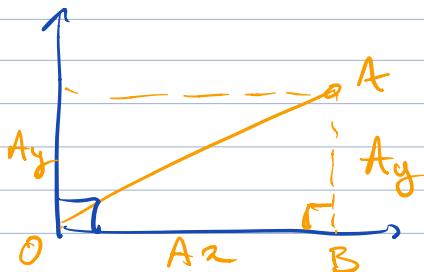
$$\|x\|^2$$

$$\|y\|^2$$

$$0$$

In a Euclidean space with Cartesian coordinates

this helps deriving



$$\|\vec{OA}\|^2 = \langle \vec{OA}, \vec{OA} \rangle$$

Here $\|\vec{OA}\|^2 = Ax^2 + Ay^2$ by Pythagoras theorem
for triangle $\triangle OAB$

thus $\|\vec{OA}\|^2 = \langle \vec{OA}, \vec{OA} \rangle = Ax \cdot Ax + Ay \cdot Ay$

which generalizes in more dimensions

and for different vectors as $\langle x, y \rangle = \sum_{i=1}^n x_i y_i$:

Note: You may also relate this to (in fact start from) duality via $\langle x, y \rangle = x^T y$ with transpose a duality

A way to think of the scalar product is via the formula:

$$\langle \alpha, y \rangle = \|\alpha\| \cdot \|y\| \cdot \cos(\hat{\alpha}, \hat{y})$$

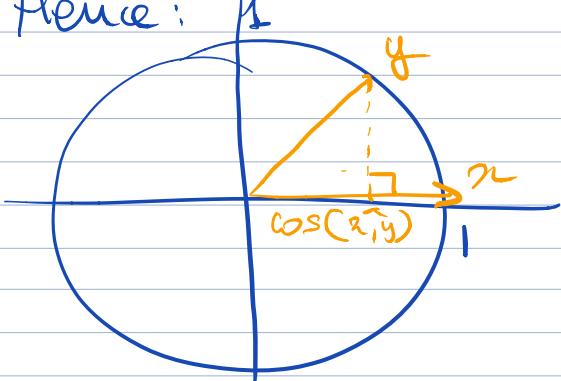
hence:

$$\langle \alpha, \alpha \rangle = \|\alpha\| \cdot \|\alpha\| \cdot \cos(0) = \|\alpha\|^2$$

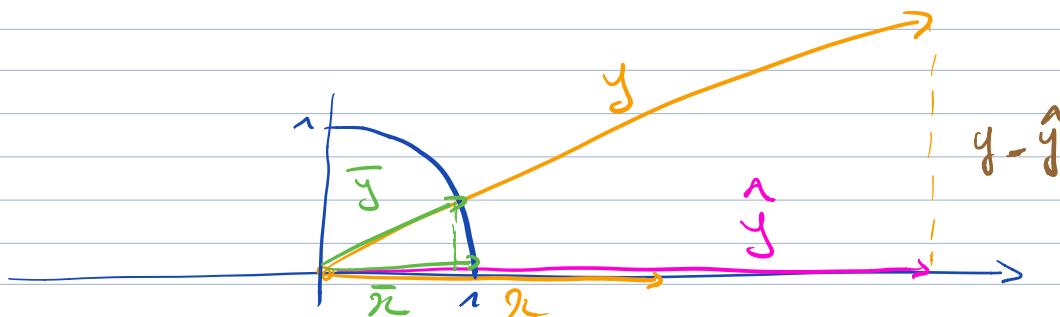
α, y orthogonal: $\langle \alpha, y \rangle = \|\alpha\| \cdot \|y\| \cdot \cos\left(\frac{\pi}{2}\right) = 0$

α, y of norm 1: $\langle \alpha, y \rangle = \cos(\hat{\alpha}, \hat{y})$

Hence:



leading to projection of a vector y of any norm onto a vector α of any norm:



$$\langle \alpha, y \rangle = \|\alpha\| \cdot \|y\| \cos(\hat{\alpha}, \hat{y})$$

that's: $\|\hat{y}\| = \|y\| \cos(\hat{\alpha}, \hat{y}) = \langle \frac{\alpha}{\|\alpha\|}, y \rangle$ and $\hat{y} = \alpha \frac{\langle \alpha, y \rangle}{\|\alpha\|}$

$$\hat{y} = \langle \frac{\alpha}{\|\alpha\|}, y \rangle \alpha = \frac{\alpha}{\|\alpha\|} \langle \alpha, y \rangle$$

Note that the norm of the projection vector, i.e. the distance between the point and its projection is: $\|y - \hat{y}\|$ with

$$y - \hat{y} = y - \frac{\alpha}{\|\alpha\|} \langle \alpha, y \rangle$$

3. Orthogonal basis - Canonical basis

Combining the scalar product and linear independence, we can construct orthogonal bases by choosing

$B = \{e_1, \dots, e_n\}$ with e_i such that

$$\langle e_i, e_j \rangle = 0 \quad \forall i \neq j$$

$$\text{and } \langle e_i, e_i \rangle = \|e_i\|^2 = 1 \quad (\text{for orthonormal basis})$$

in this case, the transition matrix P containing as columns the coordinates of e_j in B is orthogonal since

$$P^T = \begin{pmatrix} & & & & \\ & & & & \\ & & & & \\ & & & & \\ & & & & \end{pmatrix} \quad P = \begin{pmatrix} & & & & \\ & & & & \\ & & & & \\ & & & & \\ & & & & \end{pmatrix}$$

$$\begin{aligned} P^T P &= I_d \\ \Rightarrow (P^T P) &= P P^T = (I_d)^T = I_d \\ \Rightarrow P^T &= P^{-1} \end{aligned}$$

As particular (canonical) basis we can choose:

$B = \{e_1, \dots, e_n\}$ such that

$$e_i = \begin{bmatrix} 0 \\ \vdots \\ 0 \\ 1 \\ 0 \\ \vdots \\ 0 \end{bmatrix} \leftarrow_i \quad \text{which is orthogonal.}$$

4. Representation of linear maps-

Given $f: V \rightarrow W$ a map from K -vector space V to K -vector space W , f is said to be linear iff:

$$f(\alpha u + \beta v) = \alpha f(u) + \beta f(v) \in W \quad \forall u, v \in V \quad \alpha, \beta \in K$$

Given $B = \{e_1, \dots, e_n\}$ a basis on V

and $B' = \{e'_1, \dots, e'_m\}$ a basis on W

We have

$$\forall v \in V \quad \exists \alpha_1, \dots, \alpha_n \text{ s.t. } v = \sum_{j=1}^n \alpha_j e_j$$

$$\text{by linearity: } f(v) = f\left(\sum_{j=1}^n \alpha_j e_j\right) = \sum_{j=1}^n \alpha_j f(e_j)$$

similarly since $f(e_j) \in W \quad \exists \alpha_{j1}, \dots, \alpha_{jn} \text{ s.t. } f(e_j) = \sum_{i=1}^n \alpha_{ij} e'_i$

$$\text{Together: } f(v) = \sum_{j=1}^n \alpha_j \sum_{i=1}^m \alpha_{ij} e'_i = \sum_{i=1}^m \sum_{j=1}^n \alpha_{ij} \alpha_j e'_i$$

$$\text{Denoting } w = f(v) = \sum_{i=1}^m \alpha'_i e'_i$$

Hence in matrix form:

$$\begin{bmatrix} \alpha'_1 \\ \vdots \\ \alpha'_m \end{bmatrix}_{m \times 1} = \begin{bmatrix} & & & \\ & a_{11} & \dots & \\ & \vdots & & \\ & a_{m1} & \dots & \end{bmatrix}_{m \times n} \begin{bmatrix} \alpha_1 \\ \vdots \\ \alpha_n \end{bmatrix}_{n \times 1} \quad m \begin{pmatrix} & & \\ & A & \\ & & \end{pmatrix}(v)$$

Hence the $m \times n$ coefficients of matrix $A = (a_{ij})_{i=1, \dots, m; j=1, \dots, n}$

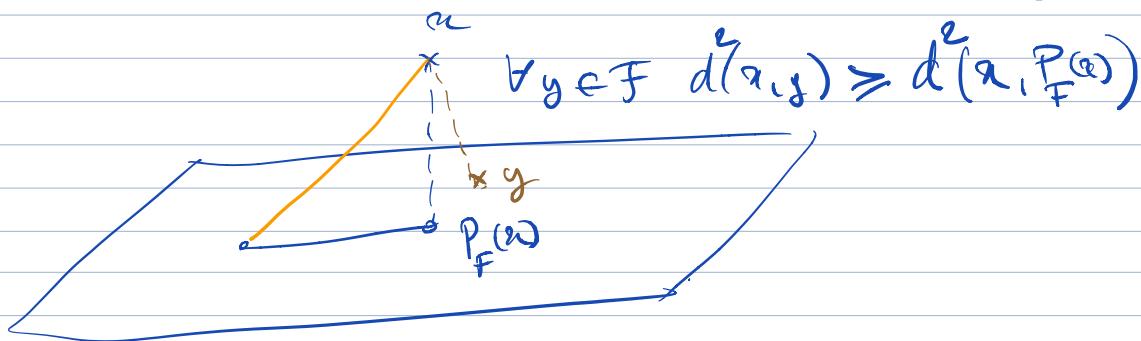
Completely determine the linear map f

S - Projection:

Projection is a basic concept in data analysis and geometry in general.

The projection defines a notion of distance between a point x and a subspace F .

$$d(x, F) = d(x, P_F(x)) \text{ where } P_F(x) = \underset{y \in F}{\operatorname{arg\min}} d^2(x, y)$$



If we define $d^2(x, y) = \|x - y\|^2 = \langle x - y, x - y \rangle$

We can use scalar product to compute the orthogonal projection.

Basically: $P_F(x)$ is the best representative (approximation) of x on F . This relates to operations that use the distance between x and $P_F(x)$.

- (slide 3a) - MSE criterion (linear systems, regression)
(lecture practice) - Computation of variance (F "mean subspace")

This definition of projection is equivalent to always considering an orthogonal projection as shown next.

Note: $d(a, F) = d(a, P_F(a))$ avec $P_F(a) = \underset{y \in F}{\operatorname{argmin}} d^2(a, y)$

$\Rightarrow P_F(a)$ is such that

$$d^2(P_F(a), a) = \min_{y \in F} d^2(a, y)$$

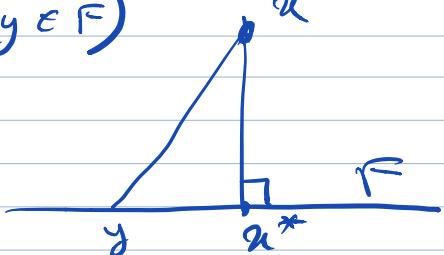
If F is flat then by Pythagoras: $\forall y \in F$

$$d^2(a, y) = d^2(a, a^*) + d^2(y, a^*) \text{ where } a^* \text{ is such that } (a - a^*) \perp F \quad (\Rightarrow \langle a - a^*, y \rangle = 0 \quad \forall y \in F)$$

Hence since $d^2(y, a^*) \geq 0$

$$\min_{y \in F} d^2(a, y) = d^2(a, a^*)$$

$$\Rightarrow P_F(a) = a^*$$



More generally consider the curve $Y: [0, 1] \rightarrow F$
 $t \mapsto y(t) \in F$

such that $\exists t^* \in [0, 1] / y(t^*) = P_F(a)$

$$\Rightarrow P_F(a) = \underset{y \in Y}{\operatorname{argmin}} d^2(a, y) \Rightarrow t^* = \underset{t \in [0, 1]}{\operatorname{argmin}} d^2(a, y(t))$$

$$\text{let } v(t) = y(t) - a \Rightarrow d^2(a, y(t)) = \|v(t)\|^2 = \langle v(t), v(t) \rangle$$

$\Rightarrow a$ being fixed $\frac{\partial v(t)}{\partial t} \Big|_{t=t^*}$ is the tangent vector to y at $y(t^*)$

Now $P_F(a)$ is such that $\frac{\partial \langle v(t), v(t) \rangle}{\partial t} \Big|_{t=t^*} = 0$

$$\begin{aligned} \text{By linearity: } \frac{\partial \langle v(t), v(t) \rangle}{\partial t} &= \langle \frac{\partial v(t)}{\partial t}, v(t) \rangle + \langle v(t), \frac{\partial v(t)}{\partial t} \rangle \\ &= 2 \langle \frac{\partial v(t)}{\partial t}, v(t) \rangle \end{aligned}$$

$$\Rightarrow \text{at } t^* \cdot \langle \frac{\partial v(t)}{\partial t}, v(t) \rangle \Big|_{t=t^*} = 0 \Rightarrow (a - a^*) \perp F$$

The projection happens perpendicularly to F

5. Eigen-things

The eigenvector of a linear map f is defined as a vector that is only rescaled by the map f

$$\text{An eigenvector} \Rightarrow \exists \lambda > 0 \quad / \quad f(u) = \lambda u$$

λ is the eigenvalue associated to eigenvector u

In matrix terms this translates to

$$A u = \lambda u \quad (\text{slide 19})$$

The set of all eigen-pairs can be gathered into one system $Au = u\Lambda$ where

- a) Columns of u are eigenvectors u_i
- b) $\Lambda = \text{diag}(\lambda_1, \dots, \lambda_n)$

$$\begin{aligned} a) \Rightarrow \quad & U^T u = \text{Id} \quad \text{since } (U^T u)_{ij} = u_i^T u_j = \delta_{ij} \\ & \Rightarrow U^{-1} = U^T \\ \Rightarrow \quad & A = u \Lambda u^{-1} = u \Lambda U^T \end{aligned}$$

One immediate consequence is that

$$\text{Tr}(A) = \text{Tr}(u \Lambda u^{-1}) = \text{Tr}(\Lambda (u^{-1} u)) = \text{Tr}(\Lambda) = \sum_{i=1}^n \lambda_i$$

i.e. the Trace is invariant to a change of Basis

$$\text{Also } A \cdot A = (u \Lambda u^{-1})(u \Lambda u^{-1}) = u \Lambda \Lambda u^{-1} = u \Lambda^2 u^{-1}$$

$$\text{Hence } A^k = u \Lambda^k u^{-1} \text{ with } \Lambda^k = \text{diag}(\lambda_1^k, \dots, \lambda_n^k)$$

Note that $A^k v$ is attached by eigenvectors

$$(A^k v \xrightarrow[k \rightarrow \infty]{} u_i \text{ for some } i)$$

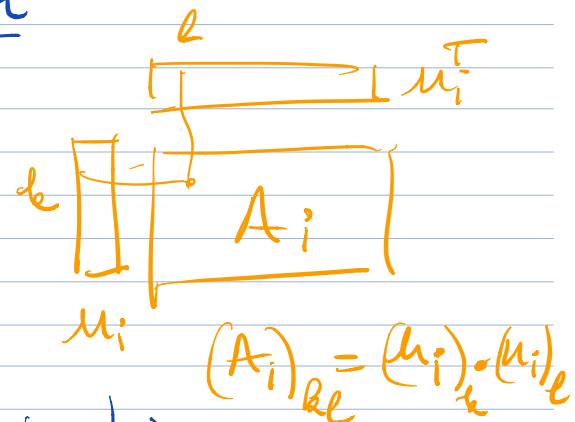
Reconstruction and Approximation

From

$$A = U \Lambda U^T$$

we can write

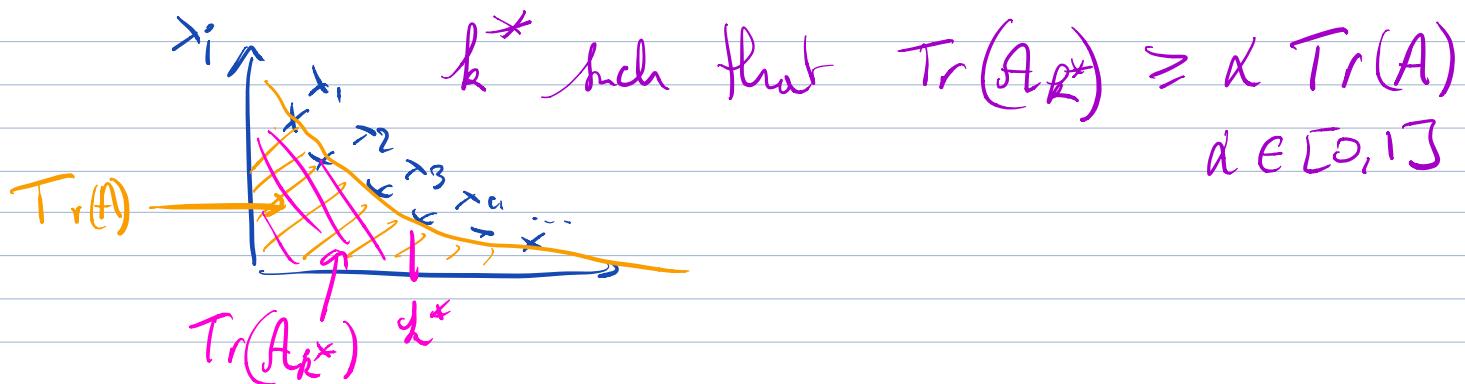
$$A = \sum_{i=1}^n \lambda_i \underbrace{u_i}_{\text{u}_i} \underbrace{u_i^T}_{\text{u}_i^T}$$



This is the basis for approximation:

$$A_k = \sum_{i=1}^k \lambda_i u_i u_i^T = \sum_{i=1}^k \lambda_i A_i \quad (\text{or any other combination})$$

The choice of approximation may be motivated by conservation of the trace:



Interpretation:

A corresponds to the representation of a linear map $f: \mathbb{R}^n \rightarrow \mathbb{R}^m$ in the original coordinate system. Λ is the view of that map along coordinate vectors u_i 's.

A_k is therefore a "simplification" of this map

in some optimal sense (eg such that it operates only within the \mathbb{R}^{k^*} subspace

of basis (u_1, \dots, u_{k^*}) (\rightarrow cf applications for PCA, LSI)

More generally, given a linear endomorphism $f: \mathbb{R}^n \rightarrow \mathbb{R}^n$, given a basis $B = \{e_1, \dots, e_n\}$ one can write

$$f(e_i) = \sum_{j=1}^n a_{ij} e_j$$

As we know, the n^2 numbers a_{ij} completely describe the action of a linear endomorphism

If we happen to choose $B = \{u_1, \dots, u_n\}$ we have

$$f(u_i) = \lambda_i u_i \quad (\text{by definition of a eigenvector})$$

Hence

$$f(v) = f\left(\sum_{i=1}^n v_i u_i\right) = \sum_{i=1}^n v_i f(u_i) = \sum_{i=1}^n v_i \lambda_i u_i$$

Therefore, in the eigenspace (the space whose basis is eigenvectors) only n numbers (λ_i) are needed to represent the action of the linear map.

Clearly the interpretation of $A = P^{-1} \Lambda P$ is $f = g \circ \lambda \circ g^{-1}$ where g is a rotation (represented by orthogonal matrix P) taking a vector in basis $\{e_1, \dots, e_n\}$ and placing it into basis $\{u_1, \dots, u_n\}$

Then λ is a scaling transform of matrix $\Lambda = \begin{pmatrix} \lambda_1 & & \\ & \ddots & \\ & & \lambda_n \end{pmatrix}$ applying the map and g^{-1} then returns this result into basis $\{e_1, \dots, e_n\}$.

Chapter ② Probabilistic Space

Data Analysis often starts with a set of documents (images, texts, objects, ...) to which we must associate a representation to immerse them into a mathematical space within which the analysis is performed. In general it is assumed that the representation is so that the space is a topological space (defining a notion of neighborhood) and a vector space (so we can sum and scale elements).

The topology is in general constructed as the standard topology where open sets are open balls for some given distance function.

(slide 8)

The distance function can be used as a notion of similarity. This view also leads to the idea of a probabilistic space where the occurrence of a "document" or a group of "similar document" may occur with more or less "chance".

A probability space can be thought of as a space of neighbourhoods (with certain properties - of a σ -algebra) with each associated with a positive measure P summing to 1

(slide 14)

1. Probability functions:



$P(A)$ can be thought of as the area of A if the total area is 1

Eg

shewix counted twice

$$P(A \cup B) = P(A) + P(B) - P(A \cap B)$$

2. Conditional probability

The information an event (A) brings about another (B) $P(B|A)$ (slide 15)

$$P(A, B) = P(B|A) P(A) = P(A|B) P(B)$$

↑ ↑ ↑ ↑ ↑
 Prob of A and B Prob of B readie A Prob of A readie B
 knowing A knowing first knowing B . first

Exple: Two tokens, one red, one blue in envelopes and given to X and Y

$$A = \{X \text{ gets red}\}$$

$$B = \{Y \text{ gets blue}\}$$

$$P(A) = \frac{1}{2} \quad P(B) = \frac{1}{2} \quad P(B|A) = 1$$

uncertain \Rightarrow uncertain \Rightarrow the knowledge of A

$$P(A|B) = \frac{1}{2} \quad (\text{other possibility } \Rightarrow X \rightarrow \text{blue}, Y \rightarrow \text{red})$$

$$P(B|A) = \frac{P(A|B)}{P(B)} = \frac{\frac{1}{2}}{\frac{1}{2}} = 1$$

3 - Random Variable.

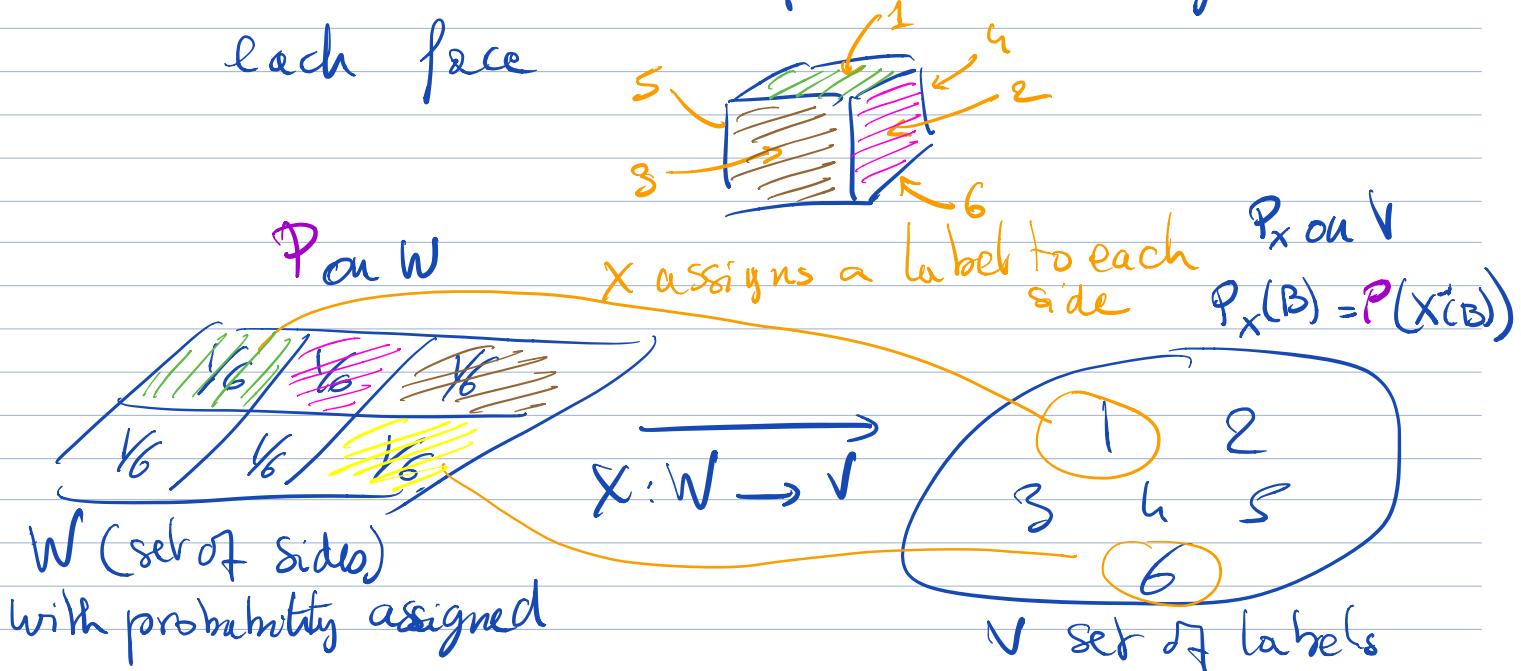
Creating a r.v is to give each event a label (on which comparisons can be made) and transferring the probability function on the space of labels.

Exle: Dice. (Exle of a discrete case)

Formally a dice is a cube. The raw event is $\omega = \text{"the cube will fall on the side"}$.

$\omega \in \Omega$ where a probability function is defined giving $\frac{1}{6}$ to each side.

A random variable X is a function from Ω to \mathbb{R} a set of labels assigned to each face



Hence I can ask $X(\omega) \leq 3$ (or $X \leq 3$) this answer "sides green ($X=1$), pink ($X=2$), brown ($X=3$)"

and use this to assign a probability : (slide 17)

$$P_X(X \leq 3) := P(X^{-1}(\{1, 2, 3\})) = P(\text{green, pink, brown}) = \frac{3}{6}$$

Note: writing " $X \leq 3$ " is not possible before assigning labels to "sides"

This is why \mathbb{V} is generally chosen as a space where an order or comparison can be made (\mathbb{N}, \mathbb{Z} for discrete, \mathbb{R} for continuous)

Continuous example: Temperature

X assigns a scale (degrees) to the height of mercury so I can mention " $X < 0$ " and compute

$$P_X(X < 0) = P(\text{"all height below that labeled 0"})$$

$$X^{-1}(0)$$

Note: for integration consistency

$$P_X(x) = 0 \quad (\text{Lebesgue notion "a.e.")}$$

Hence we extend the notion by defining

$$f(a) = P_X([a, a+da])$$

infinitesimal

Summary: Defining a random variable allows to install a probability function P_X on a space (\mathbb{V}) where an order may be defined. and then define entities like:

$$\sum_{a \in \mathbb{V}} a \cdot P_X(a)$$

discrete

$$\int_{a \in \mathbb{V}} a \cdot f(a) da$$

continuous

4 - Expectation, Variance.

The Expectation $E(X)$ of a random variable is the "center of mass" of the probability P_X

(slide 60) $E(X) = \sum_{x \in V} x P_X(x)$ $E(X) = \int a f(a) da$

Note: $E(X)$ is not the "most likely event".

Eg: Fair dice: $E(X) = \sum_{x=1}^6 x \cdot \frac{1}{6} = \frac{21}{6} = 3,5$
not an event value

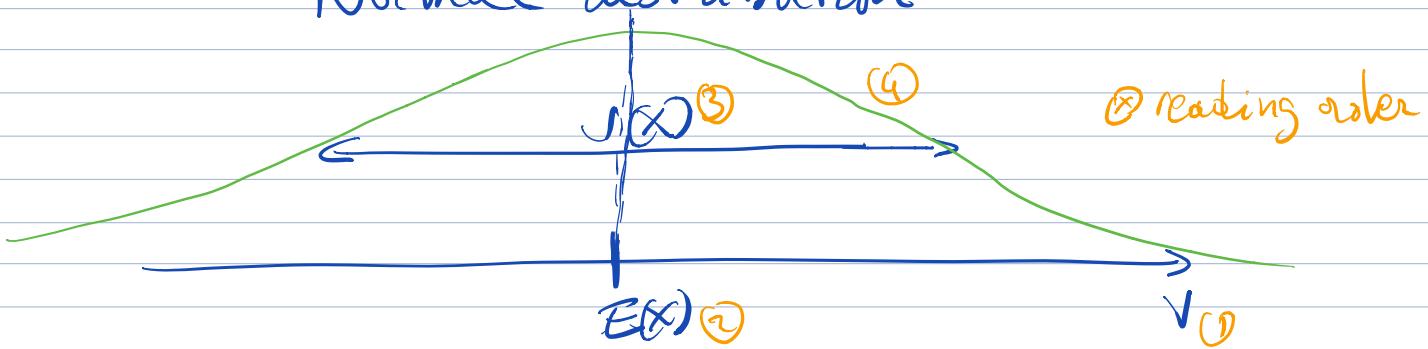
The variance is related to the expectation of deviation from this center of mass:

$$V(X) = E((X - E(X))^2)$$

↑
Expectation from $E(X)$ deviation

Note: Considering a r.v X with only non-zero expectation and variance (all other moments taken to 0) leads naturally to the

Normal distribution

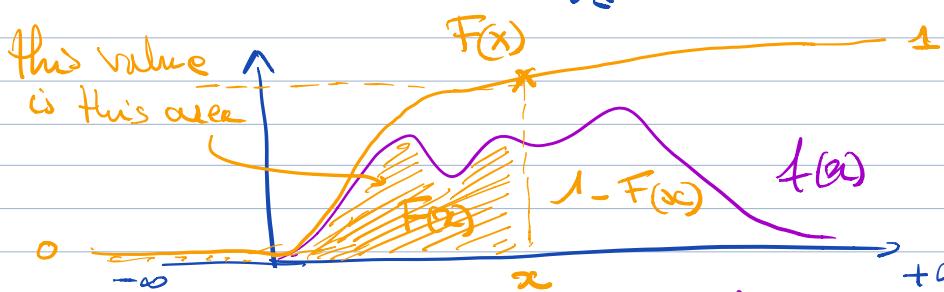


This can model a phenomenon whose "normal" (regular) realization is $E(X)$ and samples deviate in average of $V(X)$.

5 - Simulating a Probability Density Function

$X: W \rightarrow V$ Thanks to the order on V we can define the cumulative function F that tells us the cumulative chance of realization of an event w whose label $y = X(w)$ is less than $x \in V$.

$$F(x) = \sum_{y \in V} P_X(y) \leq F(x) = \int_{-\infty}^x f(a) da$$

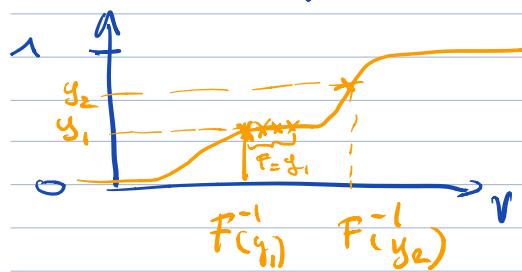


$$F: V \rightarrow [0, 1]$$

Total area under f sums to 1

(slides 30-31)

We define the inverse of F as



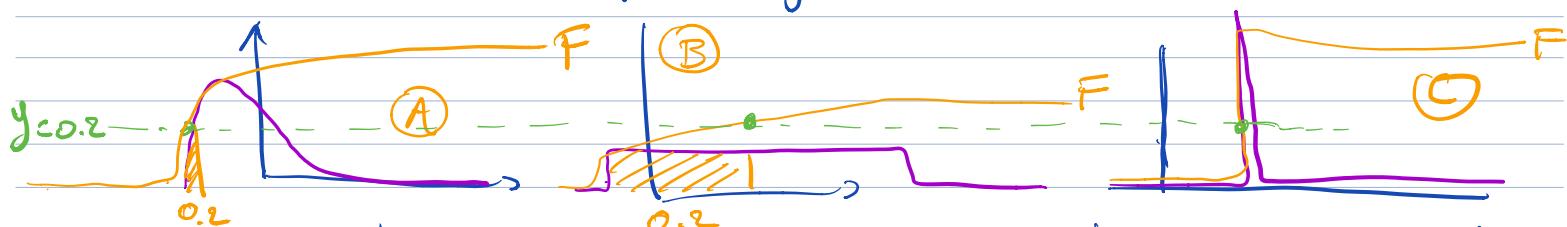
$$F^{-1}: [0, 1] \rightarrow V$$

$$y \mapsto F^{-1}(y) = \inf\{x \in V / F(x) = y\}$$

if F is not strictly monotone

Since $F(x) = y$ is the 'number' of realizations of "labels" less than x , F (and F^{-1}) "encode" f (or P_X)

Take $y = 0.2$ depending on P_X $F^{-1}(y)$ will be different



F^{-1} "remaps" uniform values to the original distribution. In extreme case, all goes to one value (C)

chapter ③ High dimensional data

This is to motivate the idea that one should be careful when choosing the representation space for the documents since the structure of the resulting (vector) space, including the dimension will influence the statistics over the data -

This is an artifact of:

- The way we compute distances:
bhaskovsky norms sum over dimensions
hence add up this effect
- the fact we only access the space (and the probability density function) via a sample and the dimension has an effect on how large this sample should be.

(slides 40-56)

Chapter ④ Principal Component Analysis

We look at a data sample \mathcal{X} whose representation is $\mathbf{X} = \{\mathbf{x}_i\}_{i=1 \dots n}$ in vector space \mathbb{R}^N .

\mathbb{R}^N is equipped with scalar product

$$\langle \mathbf{x}, \mathbf{y} \rangle = \mathbf{x}^T \mathbf{y} = \sum_{i=1}^N x_i y_i.$$

1- Data statistics

By way of representation, each sample is associated to a vector $\mathbf{x}_i \in \mathbb{R}^n$. The given of many samples allows to compute statistics over each dimension. The more samples, the better the statistics.

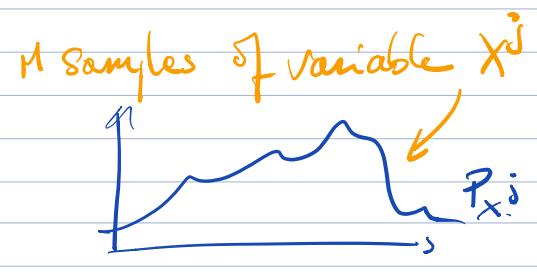
If we consider the additive noise model

$$\hat{\mathbf{x}} = \mathbf{x} + \boldsymbol{\varepsilon} \leftarrow \begin{array}{l} \text{noise} \\ \text{---} \end{array} \quad \left\{ \begin{array}{l} \text{- measure} \\ \text{- formation} \\ \text{---} \end{array} \right.$$

↑
True, never observed
Sample observed

and associate a random variable x_i to each dimension, and consider these N random variables independent, then, we can compute the statistics of the data sample as the statistics of the N -dimensional random variable whose sample is \mathbf{X} (the random variable will then also be called \mathbf{X}).

(slide 30)



2 - Notion of inertia (physics)

Identify samples with points of unit mass
in a N -dimensional space

$$\text{Center of mass} : g = \frac{1}{\sum_{i=1}^n m_i} \sum_{i=1}^n m_i \mathbf{x}_i = \frac{1}{n} \sum_{i=1}^n \mathbf{a}_i$$

Define the inertia of the sample wrt a point

$$a \in \mathbb{R}^n \text{ as } I_a(\mathbf{x}) = \sum_{i=1}^n d^2(a, \mathbf{x}_i) \quad (\text{slide 42})$$

We can extend the definition to the
inertia wrt a subspace $F \subseteq \mathbb{R}^N$

$$I_F(\mathbf{x}) = \sum_{i=1}^n d^2(a, \text{Proj}_F(\mathbf{x}_i))$$

(cf properties of projection above)

Note: Given \mathbf{x} , point a is the point minimizing
the inertia of \mathbf{x} if

$$\hat{a}^* = \underset{a}{\operatorname{argmin}} I_a(\mathbf{x}) \iff \frac{\partial I_a(\mathbf{x})}{\partial a} \Big|_{a=\hat{a}^*} = 0$$

$$\begin{aligned} \frac{\partial I_a(\mathbf{x})}{\partial a} &= \frac{\partial}{\partial a} \left(\sum_i (a - \mathbf{x}_i)^T (a - \mathbf{x}_i) \right) = \sum_i \frac{\partial}{\partial a} (a - \mathbf{x}_i)^T (a - \mathbf{x}_i) \\ &= \sum_i \frac{\partial}{\partial a} (a^T a - 2a^T \mathbf{x}_i + \mathbf{x}_i^T \mathbf{x}_i) = \sum_i -2\mathbf{x}_i + 2a \end{aligned}$$

$$\Rightarrow \sum_i -2\mathbf{x}_i + 2\hat{a}^* = 0 \Rightarrow \sum_i \mathbf{x}_i = \sum_i \hat{a}^* \Leftrightarrow \hat{a}^* = \frac{1}{M} \sum_i \mathbf{x}_i$$

→ The center of mass g is the point of minimum inertia

→ it makes sense to center our space wrt g ($\mathbf{x}_i \leftarrow \mathbf{x}_i - g$)

→ similarly we will be interested in subspaces

F_g containing g

→ relation to the grassmann theorem - slide 43

As we know, the space can be decomposed into two orthogonal components: a subspace F_g and its orthogonal F_g^\perp : $V = F_g \oplus F_g^\perp$

We can always choose (by rotation) a basis such that F is the subspace spanned by the first N_F components and F^\perp by the rest.

$$d(x, y) = (x - y)^T (x - y) = \sum_{i=1}^{N_F} (x_i - y_i)(x_i - y_i)$$

decomposes into $\sum_{i=1}^{N_F} (x_i - y_i)^2 + \sum_{i=N_F+1}^N (x_i - y_i)^2$ so that $I_g(x)$ decomposes into:

$$I_g(x) = I_F(x) + I_{F^\perp}(x)$$

I_F is the explained inertia (the inertia within the considered subspace F)

I_{F^\perp} is the remain, called the residual inertia

3. relation to statistics

Clearly there is an analogy between a physical system and statistics. Considering the mass of a particle as related to its probability, one can associate:

g center of mass $\rightarrow E(x) = \mu$ Expectation

$x - g$ centered data $\rightarrow x - \mu$ centered data

$I_g(x) = \sum_{i=1}^n d^2(x_i, g)$ inertia $\leftrightarrow \text{Var}[E((x - \mu)^2)]$ variance

In matrix form, with centered data

$$I_g = \text{Tr}(XX^T) = \sum_{i=1}^n x_i^T x_i$$

4- PCA

PCA seeks a new coordinate system into which the data will be better legible. Here "legible" means "incrementally decomposed".

The central notion is Variance / Inertia.

PCA seeks to incrementally construct subspaces which will retain the maximum of the original variance.

⇒ PCA seeks axis maximizing the variance of the data projected onto that axis

From the above analogy, because:

- Variance of projected data ~ Explained inertia
- (Explained + residual = total) inertia
- Inertia ~ Trace (invariant)

⇒ line \vec{m} maximizing variance of $\text{Proj}_{\vec{m}}(x)$

= line \vec{m} minimizing the residual inertia

$$\Rightarrow \vec{m}^* = \underset{\|\vec{m}\|=1}{\operatorname{argmin}} \sum_{i=1}^n \|x_i - \underbrace{\langle x_i, \vec{m} \rangle \vec{m}}_{\text{projection}}\|^2$$

$$\|x_i - \langle x_i, \vec{m} \rangle \vec{m}\|^2 = \langle x_i - \langle x_i, \vec{m} \rangle \vec{m}, x_i - \langle x_i, \vec{m} \rangle \vec{m} \rangle$$

$$= \langle x_i, x_i \rangle - 2 \langle \langle x_i, \vec{m} \rangle \vec{m}, x_i \rangle + \langle \vec{m}, \vec{m} \rangle \quad \leftarrow 1$$

$$= -2 \langle x_i, \vec{m} \rangle \langle \vec{m}, x_i \rangle = -2 \langle \vec{m}, x_i \rangle \langle x_i, \vec{m} \rangle \quad \text{since scalars}$$

$$= -2 \vec{m}^T x_i x_i^T \vec{m}$$

$$\Rightarrow \underset{\|\vec{m}\|=1}{\operatorname{argmin}} \sum_i \text{Tr}(\underline{u^T X X^T u}) \Rightarrow J = \vec{u}^T X X^T \vec{u} - \lambda (\vec{u}^T \vec{u})$$

$$\frac{\partial J}{\partial \vec{u}} = 0 \Rightarrow X X^T \vec{u} = \lambda \vec{u} \Rightarrow \text{an eigenvector of } \underline{X^T X}$$

In Summary:

on Subspace Δ_1 , (line of direction it's going thru g)
the projected dataset retains the maximum
of the variance of the original X .

The rest of the variance is on the orthogonal
Subspace Δ_1^\perp of dimension $N-1$
The procedure is iterated over $\text{Proj}_{\Delta_1}(X)$ to find
 $\Delta_2 \dots$ etc

In effect, this corresponds to finding all
eigenvectors u_i of the covariance matrix of
the centered data $\Sigma = XXT$ (slide 69)

u_i = Principal Component associated to $\lambda_i = \text{I}_{\Delta_i}$

5 - Properties of PCA

One can summarize PCA as:

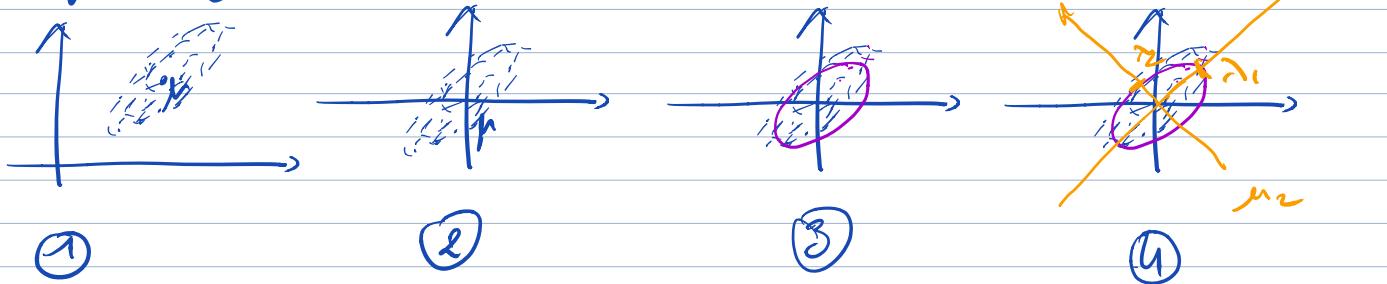
① Get data $X \in \mathbb{R}^{N \times n}$

② $X \leftarrow X - \mu$ Centering

③ $\Sigma \leftarrow XXT$ Covariance

④ $\Sigma = U \Lambda U^T$ eigen vectors

Graphically:



we have $XX^T = U \Lambda U^T \rightarrow \Lambda = U^T X X^T U$

$$\Rightarrow \Lambda = (U^T X)(U^T X)^T$$

Define new representation $Y = U^T X$ ($y_i = U^T x_i$)

the Covariance Matrix in that new representation

$$\Rightarrow Y Y^T = \Lambda \text{ (diagonal)}$$

\Rightarrow Data components are decorrelated

given a new point $x_{\text{new}} \Rightarrow$ its representation is

$$y_{\text{new}} = U^T(x_{\text{new}} - \mu) \Rightarrow x_{\text{new}} = U y_{\text{new}} + \mu.$$

Note: By considering mean (μ) and variance (Σ), PCA considers the data distributed with an underlying Gaussian (Normal) distribution.

Principal directions u_i are in the directions of maximally decreasing variance. u_i make it simple and optimal to determine the dimensions of the space

Exle: Take data distributed on a 2D ellipse add one component and rotate and translate

