

Techniques in Dimensionality Reduction

COD

Curse of Dimensionality describes the explosive nature of increasing data dimensions and its resulting exponential increase in computational efforts required for its processing and/or analysis.

Techniques in Dimensionality Reduction

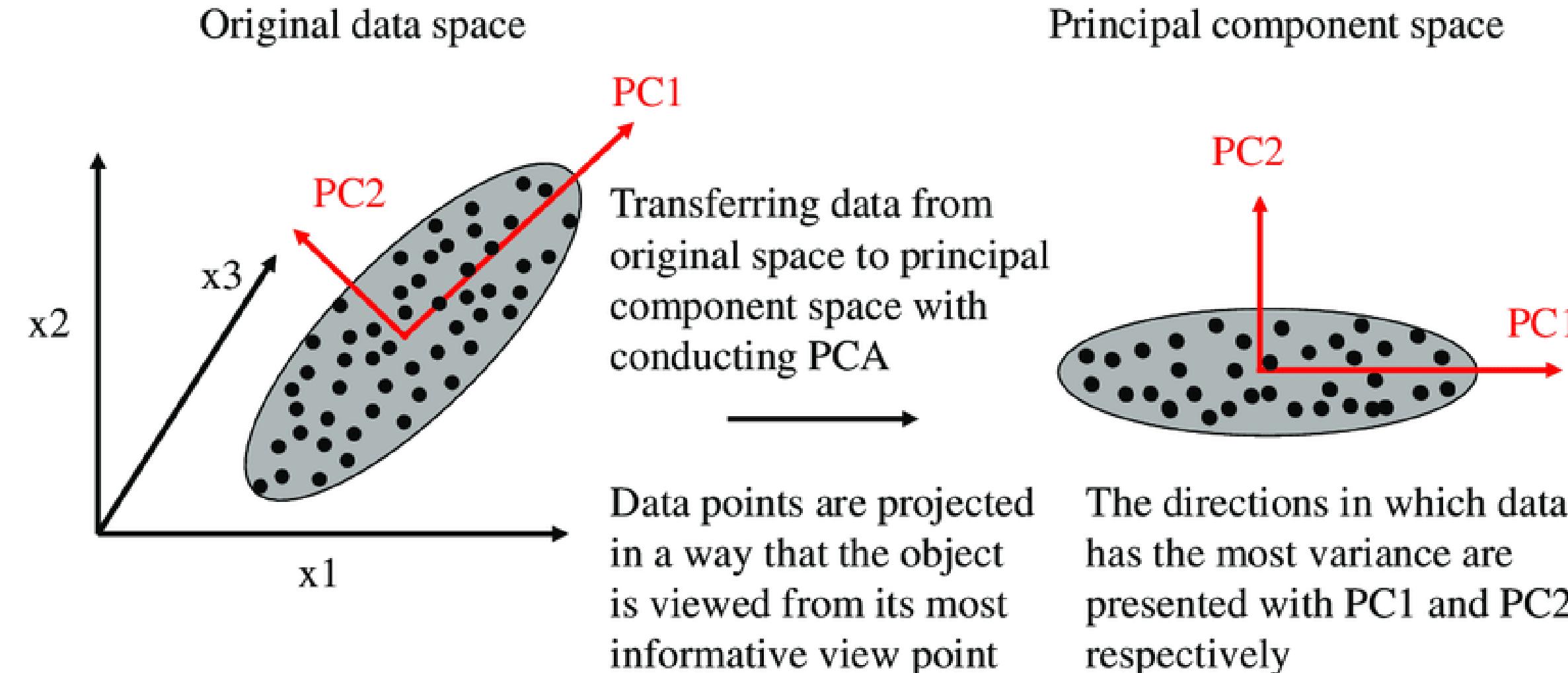
1. Feature Selection Methods.

2. Manifold Learning (e.g. t-sne, MDS, etc..)

3. Matrix Factorization (e.g PCA, Kernal PCA, etc..)

PCA

Principal Component Analysis (PCA) is a widely used dimensionality reduction technique and it comes under an unsupervised machine learning algorithm because we don't need to provide a label for dimension reduction. We can use PCA for dimensionality reduction or we can use PCA for analysis of higher dimension data in a lower dimension. PCA algorithm's, task is to find new axes or basis vectors that preserve a higher variance for data in lower dimensions. In PCA, new axes are known as PC's.



Why & when PCA:-

- 1- Better Perspective and less Complexity: When we need a more realistic perspective and we have many features on a given data set and specifically when we have this intuitive knowledge that we don't need this much number of features.
Similarly, in many other practices modelling is easier in 2D than 3D , right?
- 2 – Better visualization: When we cannot get a good visualization due to high number of dimensions we use PCA to reduce it into a shadow of 2D or 3D features (or even more but convenient enough for better parallel coordinates or Andrew Curve, e.g. when you transfer 100 features into 10 features you cannot still depict it as 2D or 3D but you can get a much better Andrew Curve)
- 3- Reduce size: When we have too much data and we are going to use process-intensive algorithms (like many supervised algorithms) on the data so we need to get rid of redundancy . Sometimes change of perspective matters more than Dimension reduction and we want to exploit dimensionality :
- 4- Different perspective: Maybe you don't have any of these motivations but you merely need to improve your knowledge on your data. PCA can give you the best linearly independent and different combinations of features so you can use to describe your data differently.

Steps in Principal Component Analysis and Mathematical Proof

Step1: Standardization of the continuous features in the dataset.

Step2: Computing the covariance matrix.

Step3: Computing eigenvalues and eigenvectors for a covariance matrix.

Step 4: Sort the eigenvalues of the covariance matrix in descending order and select the n number of eigenvalues from the top. where n is the number of axes you need.

Step 5: Select the eigenvectors corresponding to the above-selected eigenvalues.

Step 6: To get the data set for the new feature space, just use the formula of Projection of data points onto new axes. $U_{\text{transpose}} * X_i_{\text{vectors}}$. Where U is the basis vector(axes) and X_i (data points) . .

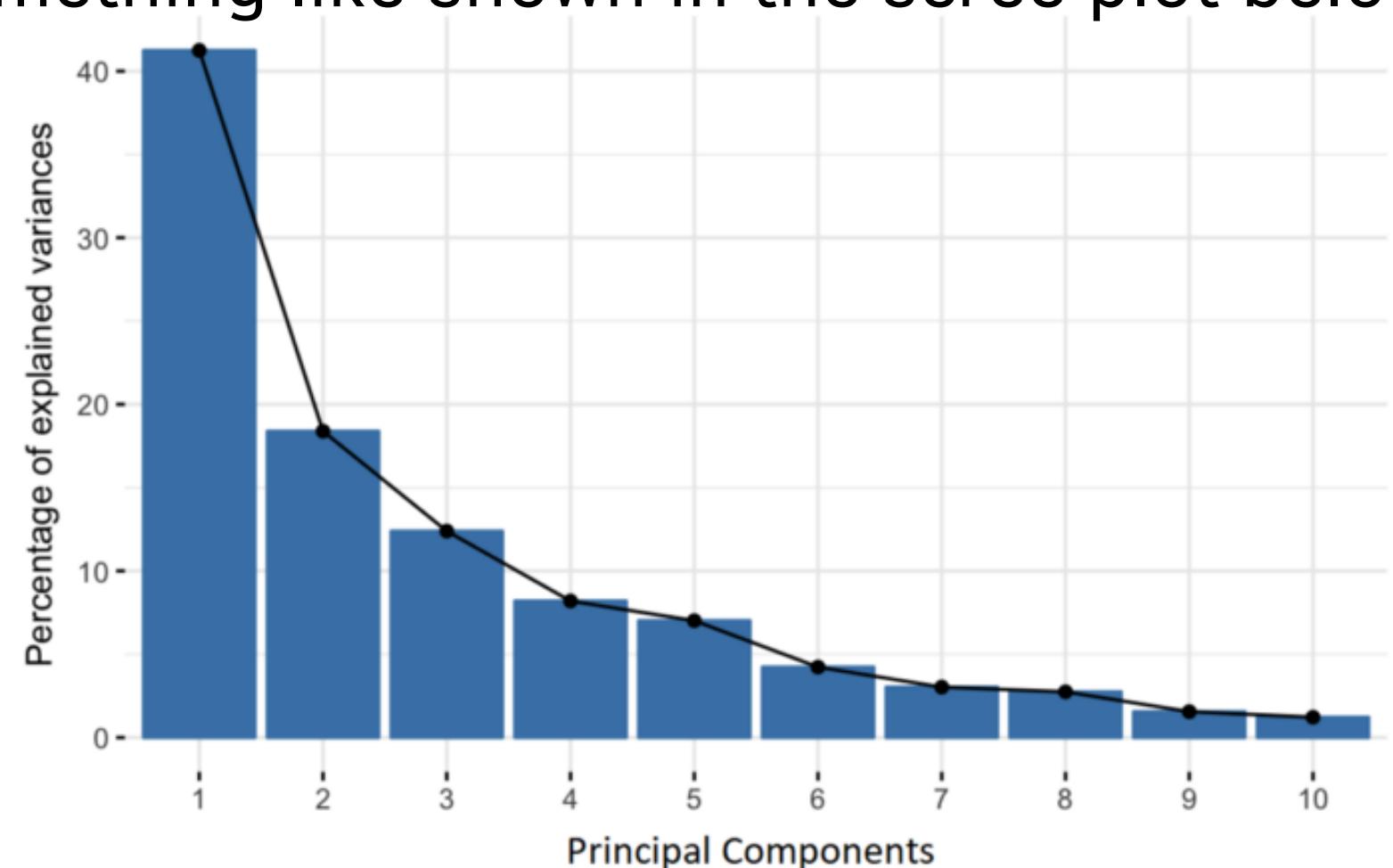
Normalize your Data before performing PCA :-

The reason you should think about Normalizing/ standardizing your data before performing PCA is that PCA is a variance maximizing exercise . PCA calculates a new projection of your data set. It projects your original data onto directions which maximize the variance. So the new axis are based on the standard deviation of your variables. So a variable with a high standard deviation will have a higher weight for the calculation of axis than a variable with a low standard deviation. If you normalize your data, all variables have the same standard deviation, thus all variables have the same weight and your PCA calculates relevant axis.

More specifically, the reason why it is critical to perform standardization prior to PCA, is that the latter is quite sensitive regarding the variances of the initial variables. That is, if there are large differences between the ranges of initial variables, those variables with larger ranges will dominate over those with small ranges (For example, a variable that ranges between 0 and 100 will dominate over a variable that ranges between 0 and 1), which will lead to biased results. So, transforming the data to comparable scales can prevent this problem.

What are Principal Components

Principal components are new variables that are constructed as linear combinations or mixtures of the initial variables. These combinations are done in such a way that the new variables (i.e., principal components) are uncorrelated and most of the information within the initial variables is squeezed or compressed into the first components. So, the idea is 10-dimensional data gives you 10 principal components, but PCA tries to put maximum possible information in the first component, then maximum remaining information in the second and so on, until having something like shown in the scree plot below.



Organizing information in principal components this way, will allow you to reduce dimensionality without losing much information, and this by discarding the components with low information and considering the remaining components as your new variables.

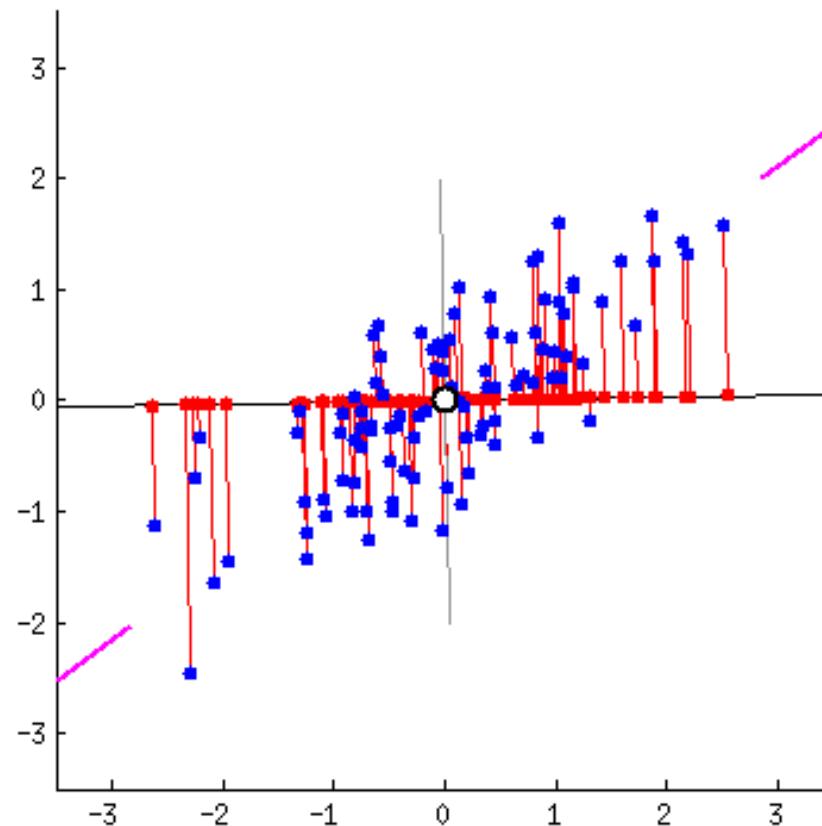
An important thing to realize here is that, the principal components are less interpretable and don't have any real meaning since they are constructed as linear combinations of the initial variables

What are Principal Components

Geometrically speaking, principal components represent the directions of the data that explain a **maximal amount of variance**, that is to say, the lines that capture most information of the data. The relationship between variance and information here, is that, the larger the variance carried by a line, the larger the dispersion of the data points along it, and the larger the dispersion along a line, the more the information it has. To put all this simply, just think of principal components as new axes that provide the best angle to see and evaluate the data, so that the differences between the observations are better visible.

How PCA Constructs the Principal Components

As there are as many principal components as there are variables in the data, principal components are constructed in such a manner that the first principal component accounts for the largest possible variance in the data set. For example, let's assume that the scatter plot of our data set is as shown below, can we guess the first principal component ? Yes, it's approximately the line that matches the purple marks because it goes through the origin and it's the line in which the projection of the points (red dots) is the most spread out. Or mathematically speaking, it's the line that maximizes the variance (the average of the squared distances from the projected points (red dots) to the origin).



The second principal component is calculated in the same way, with the condition that it is uncorrelated with (i.e., perpendicular to) the first principal component and that it accounts for the next highest variance.

This continues until a total of p principal components have been calculated, equal to the original number of variables.

STEP 3: COMPUTE THE EIGENVECTORS AND EIGENVALUES OF THE COVARIANCE MATRIX TO IDENTIFY THE PRINCIPAL COMPONENTS

it is eigenvectors and eigenvalues who are behind all the magic explained above, because the eigenvectors of the Covariance matrix are actually the directions of the axes where there is the most variance(most information) and that we call Principal Components. And eigenvalues are simply the coefficients attached to eigenvectors, which give the amount of variance carried in each Principal Component.

Why do we need a covariance matrix

The aim of this step is to understand how the variables of the input data set are varying from the mean with respect to each other, or in other words, to see if there is any relationship between them. Because sometimes, variables are highly correlated in such a way that they contain redundant information. So, in order to identify these correlations, we compute the covariance matrix.

Why do we need a covariance matrix and eigenvalues and eigenvectors of it?

We know PCA will try to find new PCs (axes), but how is it gonna do it. it tries a bunch of vectors(x_i) and tries to project all of the data points onto that vectors(x_i) and it will calculate the variance of the vectors(x_i) and select the vectors(x_i) which preserve a higher variance of the datapoint.

Think of this like a linear regression without gradient descent, where our model will try a bunch of lines and calculate the error between the actual and prediction values. For the best fit line, we considered the cost function in linear regression. But the PCA will select the axes based on the Eigenvalues. And the axes are nothing but an Eigenvector corresponding to that eigenvalues.

New Axes = Eigenvectors

But why eigenvectors and eigenvalues ?

Why do we need a covariance matrix and eigenvalues and eigenvectors of it?

Before the proof, we need to understand some math formulas involved in Principal Component Analysis.

1. We need a formula for the Variance of the projection vector.
2. Need to know about a Closed-Form of the covariance.

* Formula for Closed form of Covariance Matrix

Covariance

Closed form

$$\text{Covariance} = \frac{1}{N-1} \sum_{i=1}^{N-1} (x_i - \mu)(x_i - \mu)^T$$

$N \Rightarrow$ No of Records

Eqn① and Eqn② are equal (Same).

②

* Formula for Covariance Matrix

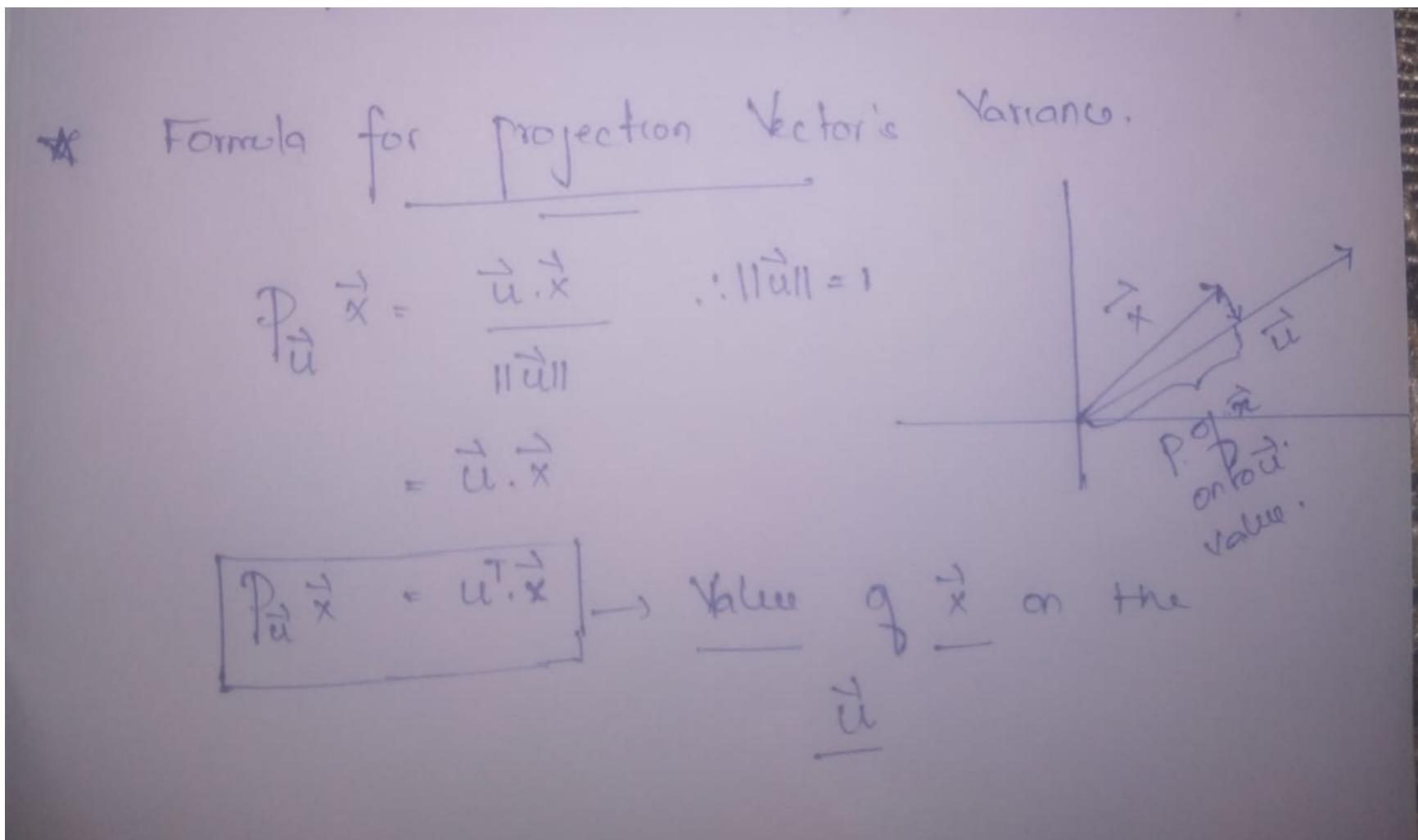
Covariance

$$\text{Covariance} = \frac{1}{N-1} \sum_{i=1}^{N-1} (x_i - \mu)^2$$

$N \Rightarrow$ No of Records

①

Why do we need a covariance matrix and eigenvalues and eigenvectors of it?



From the equation of the Variance of projection vector, we can see some magic that both closed covariance matrix and variance of projection vector are the same. So we already know that the closed covariance matrix and covariance matrix are the same. So, now variance of projection vector and covariance matrix are the same

Variance of projection (P) $\left\{ \frac{1}{N-1} \sum_{i=1}^{N-1} (\vec{u}^T \vec{x}_i - \bar{u} \cdot \vec{x})^2 \right\} \quad \bar{u} \rightarrow \text{Mean value}$

$$(\vec{u} \cdot \vec{x}) = \frac{1}{N-1} \sum_{i=1}^{N-1} (\vec{u}^T (\vec{x}_i - \bar{x}))^2$$

$$= \frac{1}{N-1} \sum_{i=1}^{N-1} (\vec{u}^T (\vec{x}_i - \bar{x}) \cdot (\vec{x}_i - \bar{x})^T \cdot \vec{u})$$

part of variance $\vec{u}^T \vec{C}_{uu} \vec{u} = \frac{1}{N-1} \sum_{i=1}^{N-1} [(\vec{x}_i - \bar{x}) \cdot (\vec{x}_i - \bar{x})^T]$

Variance of projection vector $\left\{ \frac{1}{N-1} \sum_{i=1}^{N-1} [(\vec{x}_i - \bar{x}) \cdot (\vec{x}_i - \bar{x})^T] \right\}$

LAST STEP: RECAST THE DATA ALONG THE PRINCIPAL COMPONENTS AXES

In the previous steps, apart from standardization, you do not make any changes on the data, you just select the principal components and form the feature vector, but the input data set remains always in terms of the original axes (i.e., in terms of the initial variables).

In this step, which is the last one, the aim is to use the feature vector formed using the eigenvectors of the covariance matrix, to reorient the data from the original axes to the ones represented by the principal components (hence the name Principal Components Analysis). This can be done by multiplying the transpose of the original data set by the transpose of the feature vector.

$$\text{FinalDataSet} = \text{FeatureVector}^T * \text{StandardizedOriginalDataSet}^T$$

Pros and Cons :

Pros :

- 1) Reduces Overfitting : PCA helps the model not to overfit by reducing unwanted features from the dataset.
- 2) Helps in Visualization : It is difficult to visualize data which are high in dimension (4D or more), PCA helps to visualize it by reducing the dimensionality.
- 3) Improves Model Performance : Having too many features will lead the model not to give best and accurate results, but thanks to PCA , It speeds up the machine learning algorithm by getting rid of correlated variables.

Cons :

- 1) Standardization is a must on the dataset before implementing PCA, else PCA will not be able to optimal principal Components.
- 2) Chances of information loss if we don't choose the principal components with care.
- 3) Input variables tend to become less interpretable.

DATE

Feature Selection

Curse of Dimensionality (COD)

मध्ये Feature ही more dimension

as No of feature ↑ Bay limit GTD accuracy ↓ वैद्यता

Complex GT आपेक्षा acc. May decrese
ही अंग Feature add करते चाल
किंवा मध्ये नाही नोटी,

COD

Sparsity द्वारा कृत वा in (higher dim) eg wallet जुळा
on road, playground,
building करते ती किंवा
essential ways to find wallet

As dim ↑.

- Performance decreases

- Computation ↑ in 1-D (constant)

How to

How to solve

After Reducing

Reduced

As dim \uparrow ,

- Performance decreases
- Computations \uparrow (less easy ways to find wallet is 1-D (on road))

How to solve

Dimensionality Reduction

How to

Code

feature selection

- forward selection
- backward elimination
- ...

feature extraction

- PCA
- LDA
- tSNE

Feature Extraction

PCA Principle Component Analysis

- unsupervised learning prob.
- Complex technique (old)

PCA is a technique which can transform higher dim data

to lower dim data while keeping essence of data

(highlight data on behavior of
that are lower dim & reduce noise)

(e.g. Photos जैसे 3-4 करोड़ 2D फोटो)

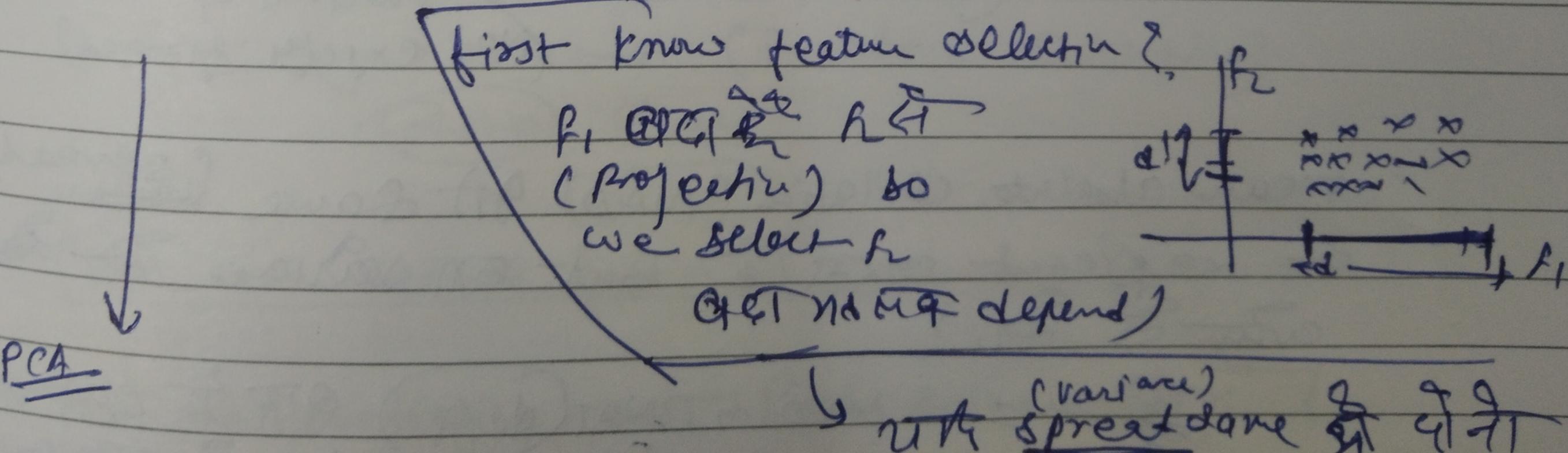
why we use PCA in ML & (Benefits)

1 → higher dim lower dim → algo size कम हो जाता है और algo run fast

(faster execution of algo.)

2 → Visualization
Higher dim data not easy to understand
it reduce data dimensions to visualize the data

Geometric Intuition:-



we select r

אַלְפָה אֲלָמָה

PCA

Get rid of depend)

उक्त ^(Variance) spread दाने की विधि
पर तो decide करी पर प्रयोग
कर करता select (Remove)
ही feature selection की
so solves by feature extraction

Has feature extraction
Solve this problem

<u>PCA</u>	eg.	room	washroom	price
------------	-----	------	----------	-------

Price depends on room / washroom etc let & बिनाे f₁ & f₂

equally as agreed in the PCA and onto the fit

Chy ~~to~~ remove only new dry got ५ मी ८३१

Cream. Sust. (g per)

So PCA rotate the axis

$$\text{Now } d > d'$$

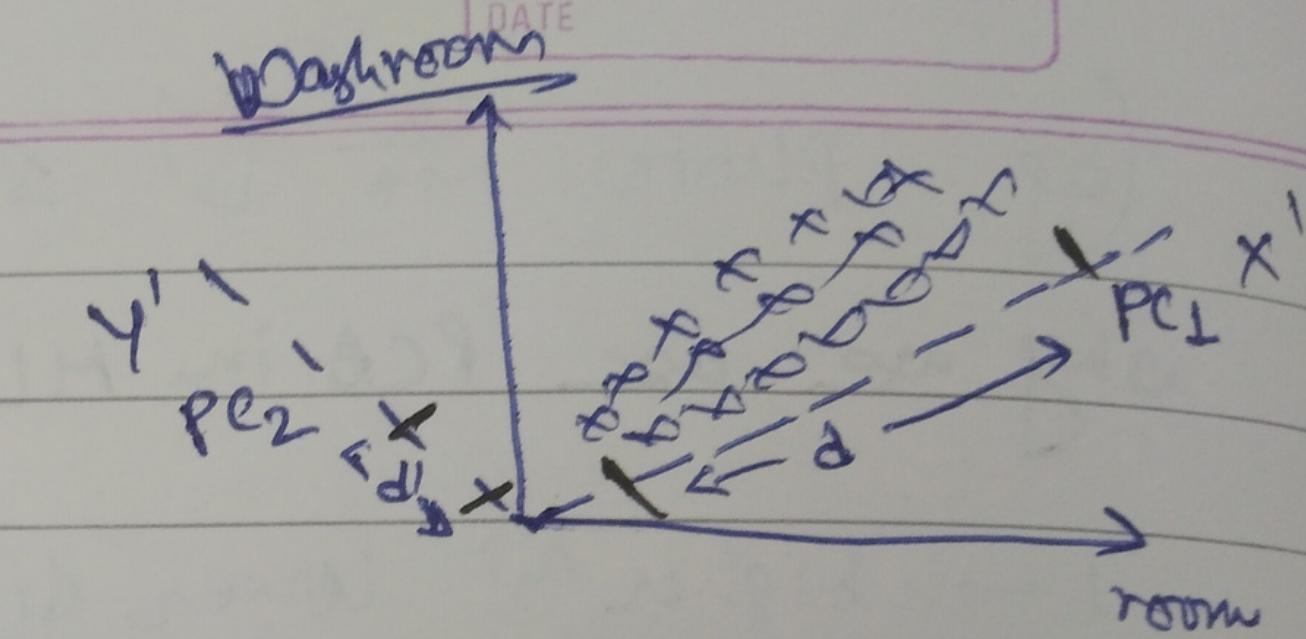
So transform the data acc to. PC1

↳ why variance (~~spread~~) is important & in (Stat.)

Variance ~~variance~~ variance

(spread & variance) variance is
(Not exactly spread)

mean absolute deviation (MAD) get some ~~same~~ spread
to represent out ~~out~~ But not variance out ~~out~~
out or ~~out~~ ?



answ at 28 %

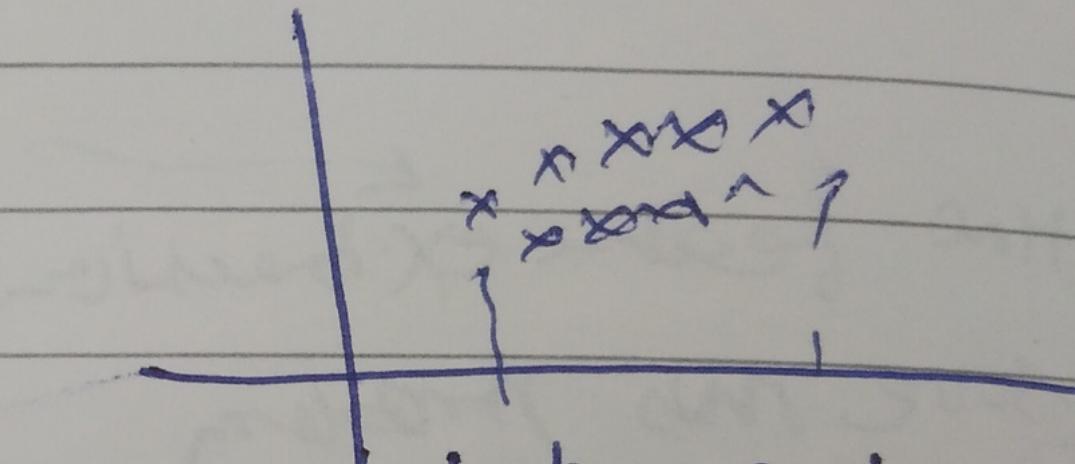
Because MAD is not diff

ATD is optimiser
tech diff w/ AT
AT is hard

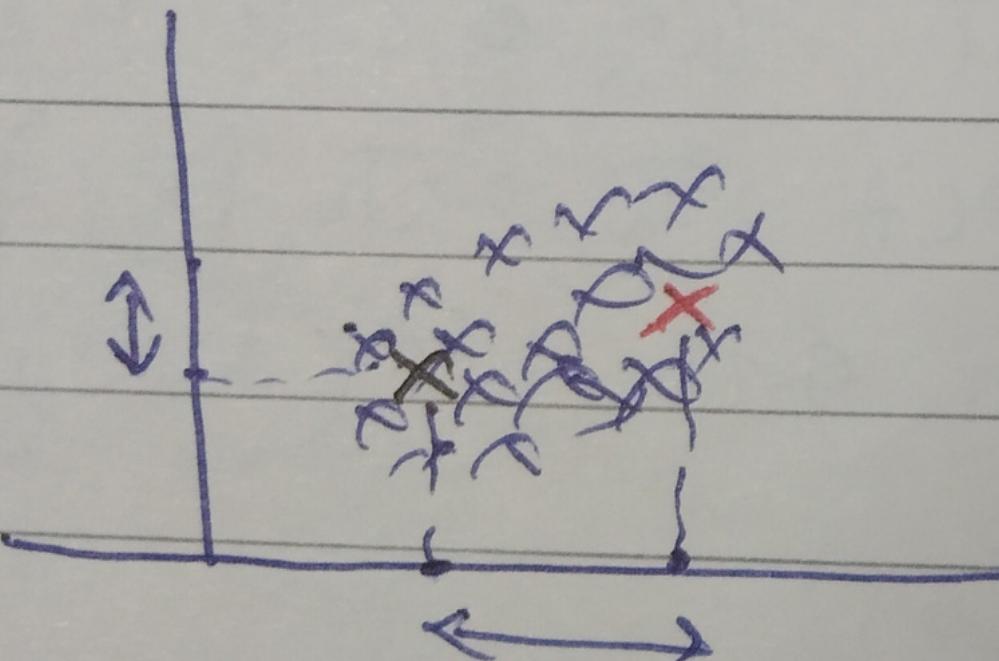
$$MAD = \frac{\sum |x_i - \bar{x}|}{n}$$

$$\text{Variance} = \frac{\sum (x_i - \bar{x})^2}{n}$$

Variance imp for PCA ?



Take projection then
find variance.



let PCA help AT fed/black
point after checke answ

Proj on X-axis & acc Variance w.r.t. Σ
as compare to Y

It take distance based algo uses \sqrt{d} to rainy
to acc distinguish off wt. point's activity list
variance acc to so variance w.r.t. Σ

e.g. Photographer taking photo's of Mayas (Marry)
to take angle & (area spread (or) &
(or) spread) best angle & (or)
say X-axis (in above example)
photo clear

(PCA dim, reduce Σ)

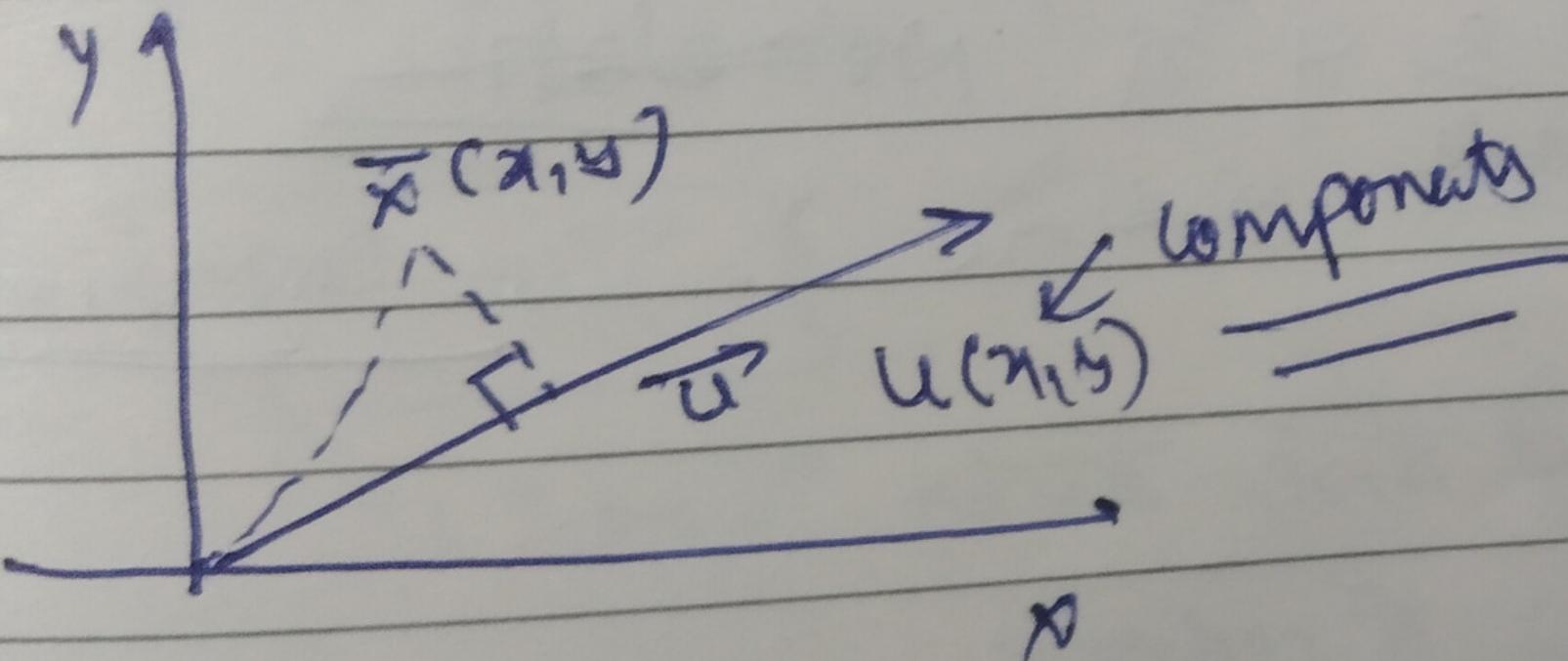
~~# PCA mathematically~~ ~~out~~ Prob solve ~~on~~ ~~in~~ } problem formulation

In PCA, we have to find a direction (^{unit} Vector) such that the variance of all points (after projection) is Maximum

Project (unit vector)

$$\frac{\vec{u} \cdot \vec{x}}{\|\vec{u}\|}$$
 x on \vec{u} vector

$$\vec{u} \cdot \vec{x} = \boxed{u^T x}$$



We know Variance $\sigma^2 = \frac{1}{n} \sum (x_i - \bar{x})^2$

$\Leftrightarrow u^T \bar{x} \leftarrow$ mean at point \bar{x} projected

Objectives $f(u)$

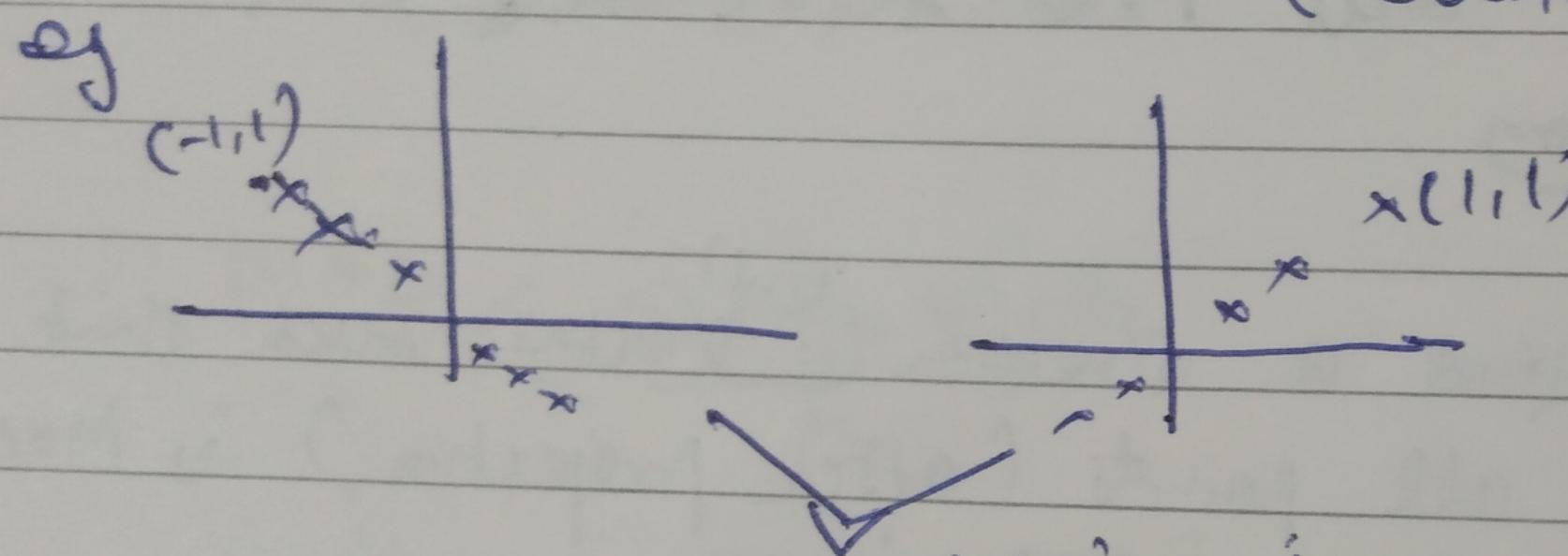
$$\boxed{\frac{\sum (u^T x_i - \bar{u}^T \bar{x})^2}{n} = \text{Variance}}$$

find unit vector such that variance is maximizing

Some concepts:

Covariance & Covariance Matrix

variance only say ~~any~~ ~~at~~ spread बताता है
But Multiple ~~to~~ compare ~~at~~ are पड़ता
(relationship)



variance is same (~~say on~~
हमें दो डेटेसेट समान हैं लेकिन
भिन्न ओर हैं)

so ~~if~~ x & y ~~is~~ ~~no relation~~
relation इसके बहाव क्या है? covariance

1) Correlation & Covariance are nearly same but correlation is $-1 \text{ to } 1$ But no restriction on covariance

Raishree
PAGE NO.
DATE

Covariance Matrix

Let two feature

$$\begin{matrix} x_1 & \left[\begin{matrix} \text{cov}(x_1, x_1) & \text{cov}(x_2, x_1) \\ \text{cov}(x_1, x_2) & \text{cov}(x_2, x_2) \end{matrix} \right] \\ x_2 & \quad \quad \quad x_1 \quad \quad \quad x_2 \end{matrix}$$

$$\text{cov}(x_1, x_2) = \text{var}(x_1)$$

$$\text{cov}(a, b) = \text{cov}(b, a)$$

$$\begin{bmatrix} \text{var } x_1 & \text{cov}(x_1, x_2) \\ \text{cov}(x_1, x_2) & \text{var } x_2 \end{bmatrix}$$

\rightarrow Symmetric

Ans

Cov. matrix to give full details about data.

Diag \leftarrow give all site var by data $\frac{\partial}{\partial t} \frac{\partial}{\partial t} \frac{\partial}{\partial t}$ विलम्बी

Other \leftarrow cov term \rightarrow उन्हें तु देते हैं विवरण अंकों में

etc

Linear Transformation:

Matrices

→ vector transformation

→ are basically LT. b/w co-ordination system to change in dir.

Eigenvectors: - \vec{v} all vectors b/w sys. b/w LT change in dir. doesn't change b/w sys. But scale change in sys. may P.↓, 2x -ve dir. b/w sys. change in eigenvalue.

A

Let perple eigenvector on mey I am \vec{v} is

ϕ is it am the factor by which \vec{v} has changed

so λ is eigenvalue

$$A\vec{v} = \lambda\vec{v}$$

LT on \vec{v} makes \vec{v} a multiple (mey v, 1 or but direction n't change)

after solving last optimisation f^n (objective β^n) we find that,

The largest eigenvector of the covariance matrix always points into the direction of the largest variance of the data & magnitude of this vector equals the corresponding eigenvalue.

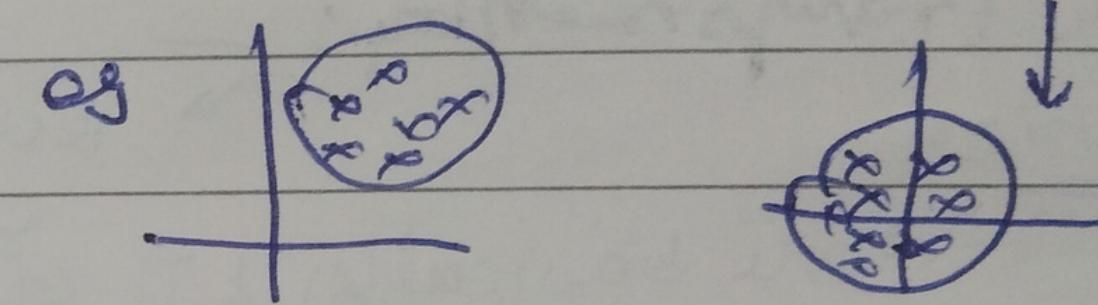
Do ~~simplify~~ covariance matrix w.r.t eigenvectors
eigenvalue ~~rotation~~ ~~to~~ soft

larger $e^T \rightarrow$ largest spread

Step by step for $\mathbf{A}^T \mathbf{A}$:-

$$f_1 \left| \begin{array}{c} \mathbf{B}_2 \\ \mathbf{B}_3 \end{array} \right|$$

1) mean center of all column.



2) Go find the Cov. matrix \mathbf{P}

3) Find the eigen value / vector

4) $v_1, v_2, \dots, d_1, d_2, \dots$

$\max(d_1, d_2, \dots)$ is PC1

~~sort~~

eg $3 \geq d_1, d_2, d_3$ but $d_1 > d_2 > d_3$

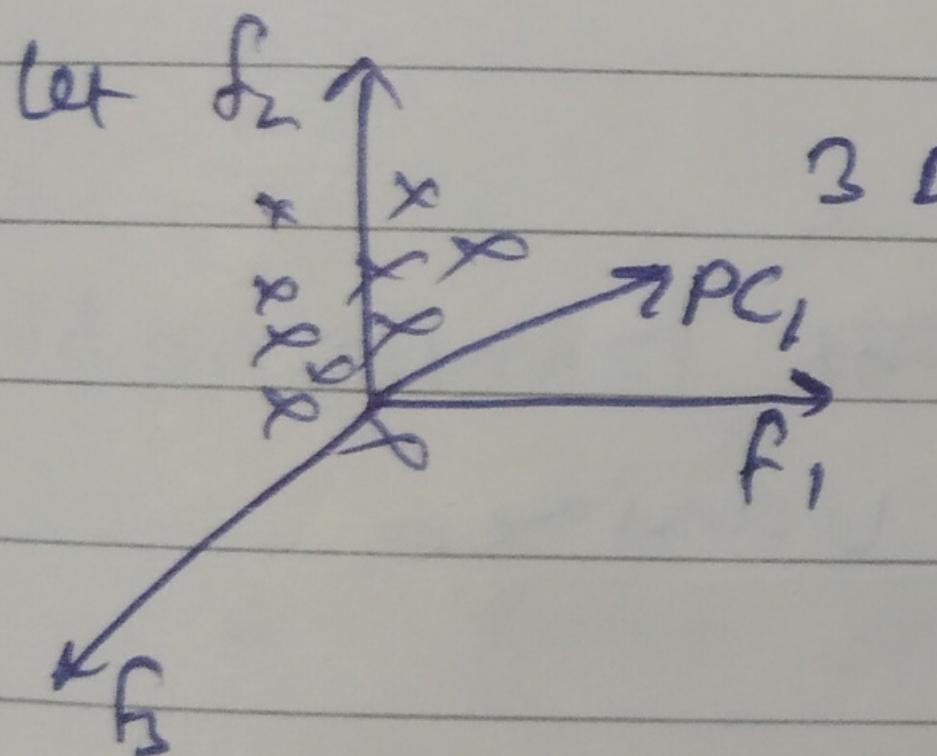
\downarrow \downarrow \downarrow
PC1 PC2 PC3

④

$2 \geq$

How to transform Points?

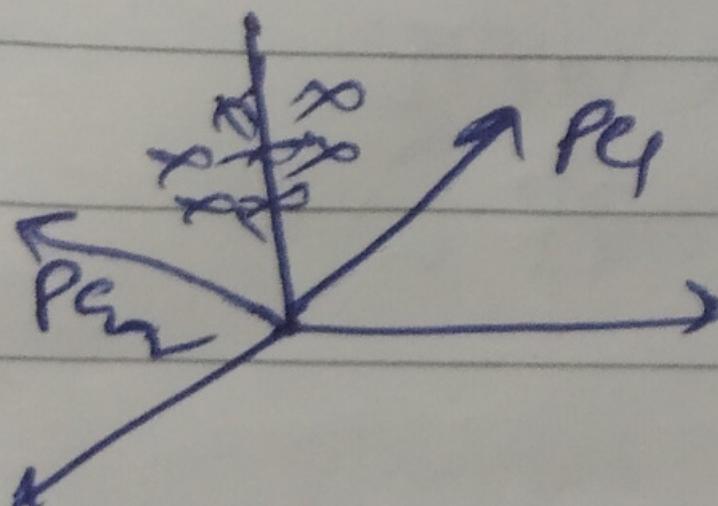
(Now we get PC's. → Has to map)



2 D

By dot Product
(projective) concept
on P_C1

P_C1 | output



2 D

on P_{C1}, P_{C2}

P_{C1} | P_{C2} | output

Finding optimum number of Principle Components:-

$d_1 \dots d_n$

$$d_i' = \left(\frac{d_i}{d_1 + d_2 + \dots + d_n} \right) \times 100 \quad \leftarrow \text{Percentage Value}$$

→ ये अंकों का गुणांक - % Variance
explain करते हैं ये data की

we need 90% variance approx.

add -- $d_1' \dots d_n'$

let d_{15}' comp give $\approx 90\%$. ✓
so LS (explained variance)

pca-explained-variance-ratio → eigen values

pca.

np.cumsum()) so we get explained variance

(or)

plt.plot()

$\approx 20 \quad 90\%$.

what PCA will produce a set of Principal components which

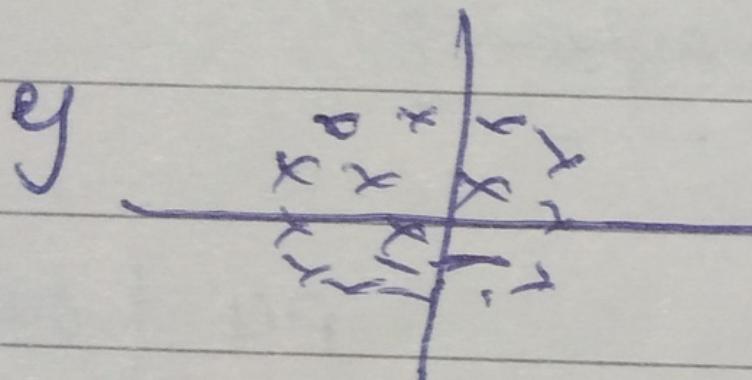
Rajshree

PAGE NO.

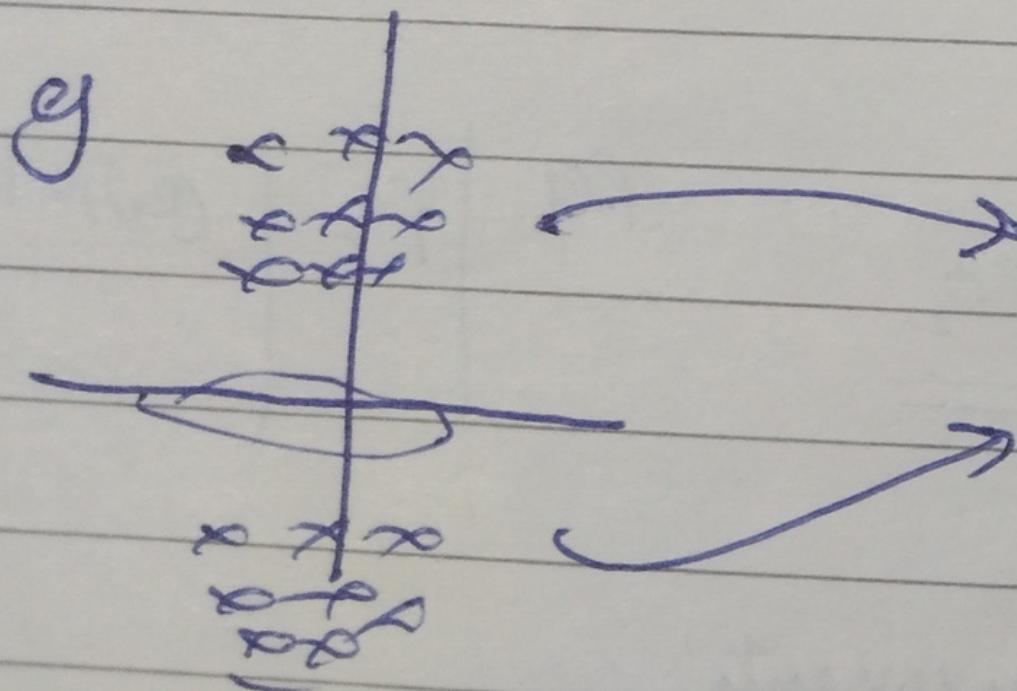
- Reduce noise, by maximizing feature variance

- Reduce redundancy, by minimizing the covariance b/w pairs of features.

when PCA does not work.



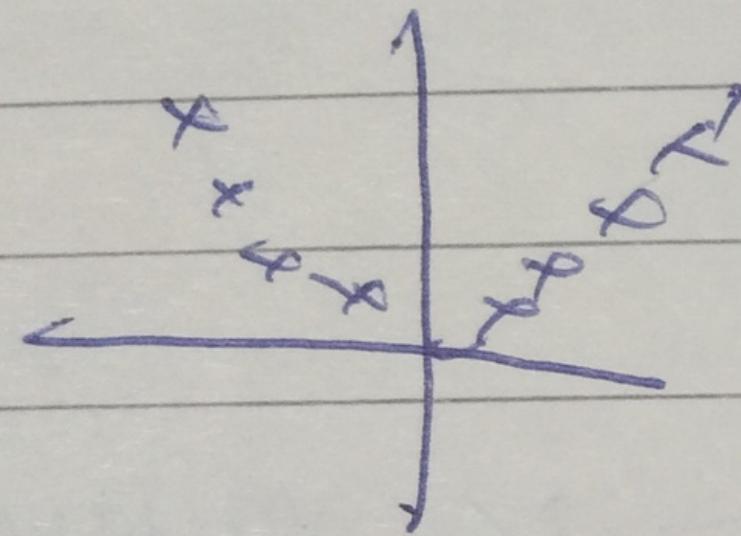
Let data is circle तो लेट ही \overline{var} on transformation
मात्र ही दो आसे ही समीक्षा वारिए



दोनों Project same वारि

ही differentiate ही नहीं

og



data fit specify pattern

→ if L fit must fit pattern loosely
but not so well enough
L fit against other orthogonal

To reduce dimensionality, autoencoder is another commonly used method. But, the latent space of the autoencoder is not necessarily uncorrelated. While PCA guarantees that all features are uncorrelated with each other,

PCA. Effect of normalization → common rep., all feature & principal components are → linearly indep, Orthogonal. Normalization also guarantees that

- Standard St.

e.g. from sklearn.preprocessing import StandardScaler

scaler = StandardScaler()

X_train = scaler.fit_transform(X_train)

X_test = scaler.transform(X_test)

- from sklearn.decomposition import PCA

pca = PCA(n_components = 2)

No of Comp.
Pca \leq n \leq total

2 \leq 2

By default **None** \leftarrow 2

X_train = pca.fit_transform(X_train)

X_test = pca.transform(X-test)

Now

fit in any Model

e.g. knn

~~from~~

from sklearn.neighbors import
KNeighborsClassifier

knn = KNeighborsClassifier()

knn.fit(X_train, y_train)

y_pred = knn.predict(X_test)

accuracy_score(y_test, y_pred)

~~plot~~

import plotly.express as px

y_train ~~df~~ = y_train.astype('int')

for 3d

px.scatter_3d

fig = px.scatter(x=X_train[:, 0], y=X_train[:, 1],

color = y_train,

color_discrete_sequence = px.colors.qualitative.G10

G10

fig.show()