

Principal Component Analysis

PCA is an unsupervised statistical technique that is used to reduce the dimensions of the dataset. ML models with higher dimensionality tend to fail when operating on a higher input dataset. PCA helps in identifying relationships among different variables & then coupling them. PCA works on some assumptions it helps developers maintain a standard.

Or

Principal Component Analysis or PCA is a simplest and very fundamental technique which is extensively used for dimensionality reduction of the large data set. Reducing the number of components or features costs some accuracy and on the other hand, it makes the large data set simpler, easy to explore and visualize. Also, it reduces the computational complexity of the model which makes machine learning algorithms run faster.

PCA transform the features into a new set of features by rotating the axes which are called PCs (Principal Components). The principal components would be equal to the number of original features in the given dataset.

The first principal component (PC1) contains the maximum variation which was present in earlier features, and this variation decreases as we move to the lower level. The final PC would have the least variation among all features and you will be able to reduce the dimensions of your feature set.

Assumptions:

1. **Linearity:** There must be linearity in the data set, i.e. the variables combine in a linear manner to form the dataset. The variables exhibit relationships among themselves.
2. **Outliers :** outliers, should be less. More number of outliers will represent experimental errors and will degrade your [ML model](#)/algorithm.
3. The feature set must be correlated and the reduced feature set after applying PCA will represent the original data set but in an effective way with fewer dimensions.
4. ***Large variance implies more structure:*** high variance axes are treated as principal components, while low variance axes are treated as noise and discarded.

Now for MNIST dataset it has 784 dimensions we can't visualize it in 2D by using pair plot or scatter plot

Principal Component analysis (PCA)

(why?)

dimensionality reduction

$$x_i \in \mathbb{R}^d$$

✓ dims

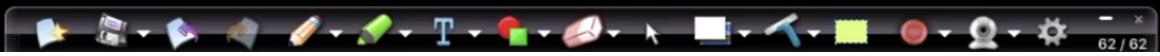
$$d' \text{-dim}$$

$$d' < d$$

Simplest
fundamental

$$\checkmark \quad ① \quad \text{MNIST} \rightarrow 184 \text{-dim} \rightarrow 2 \text{-dim} \quad , (\text{visualize})$$

$$\text{why?} \rightarrow ② \quad d \text{-dim} \rightarrow d' \text{-dim} \quad d' < d \quad d' = 10$$



Geometric interpretation of PCA:

For visualization we convert d -dim. (Higher dim.) data to lower d' -dim. Data or we converting data 2D to 1D, Let say we have two feature f_1 and f_2 we plot f_1 feature "blackness of hair" in x-axis and f_2 feature is "height" in y-axis. And our dataset x is a datamatrix of this two f_1 and f_2 features from the plot we can say variance or spread of height or f_2 feature is more compare to f_1 feature now suppose we want to convert this 2D dataset into 1D or represent the x dataset into x' data matrix which contain only one feature or one vector, in that case we only keep one feature f_2 and skip f_1 feature bcz the spread is very low in f_1 axis or we have to force to skip one feature between this two then we only skip f_1 bcz by skipping this feature we lose very less information

Geometric interpretation of PCA:

f_2
Spread
Variance

$$2D \rightarrow 1D$$

spread on f_1 is very low

Indians
gray, black

$$X = \begin{bmatrix} x_1 \\ x_2 \\ \vdots \\ x_n \end{bmatrix}$$

$2 \rightarrow 1$

$$d \text{-dim} \rightarrow d' \text{-dim}$$

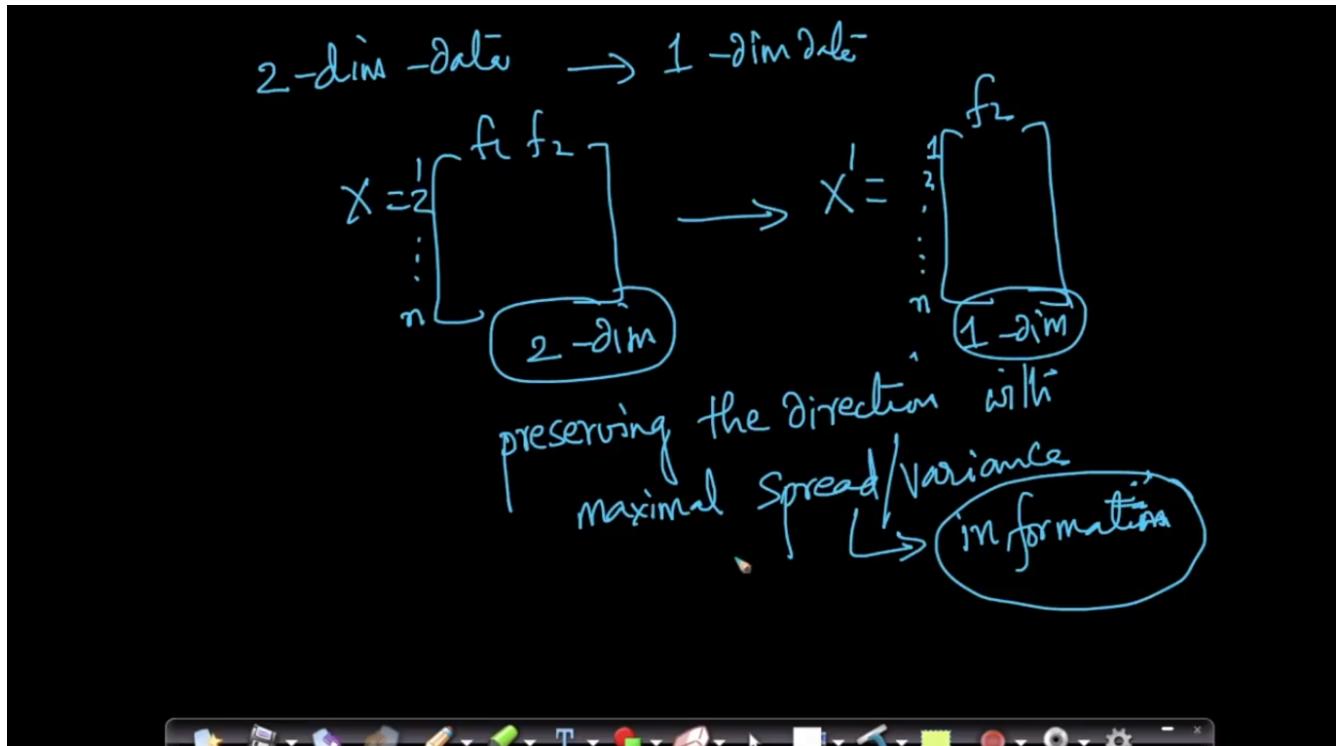
$$2 \text{-dim} \rightarrow 1 \text{-dim}$$

Variance =
Spread =
Variability

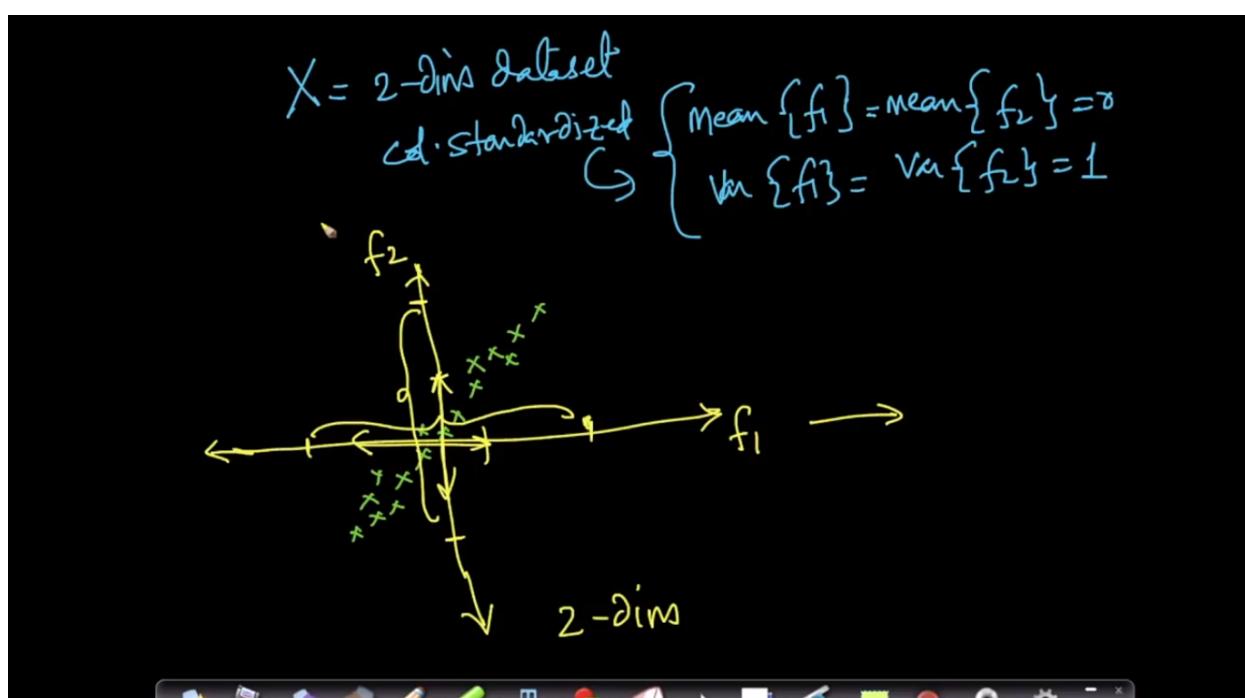
$$X' = \begin{bmatrix} f_2 \\ \vdots \\ f_1 \end{bmatrix}$$



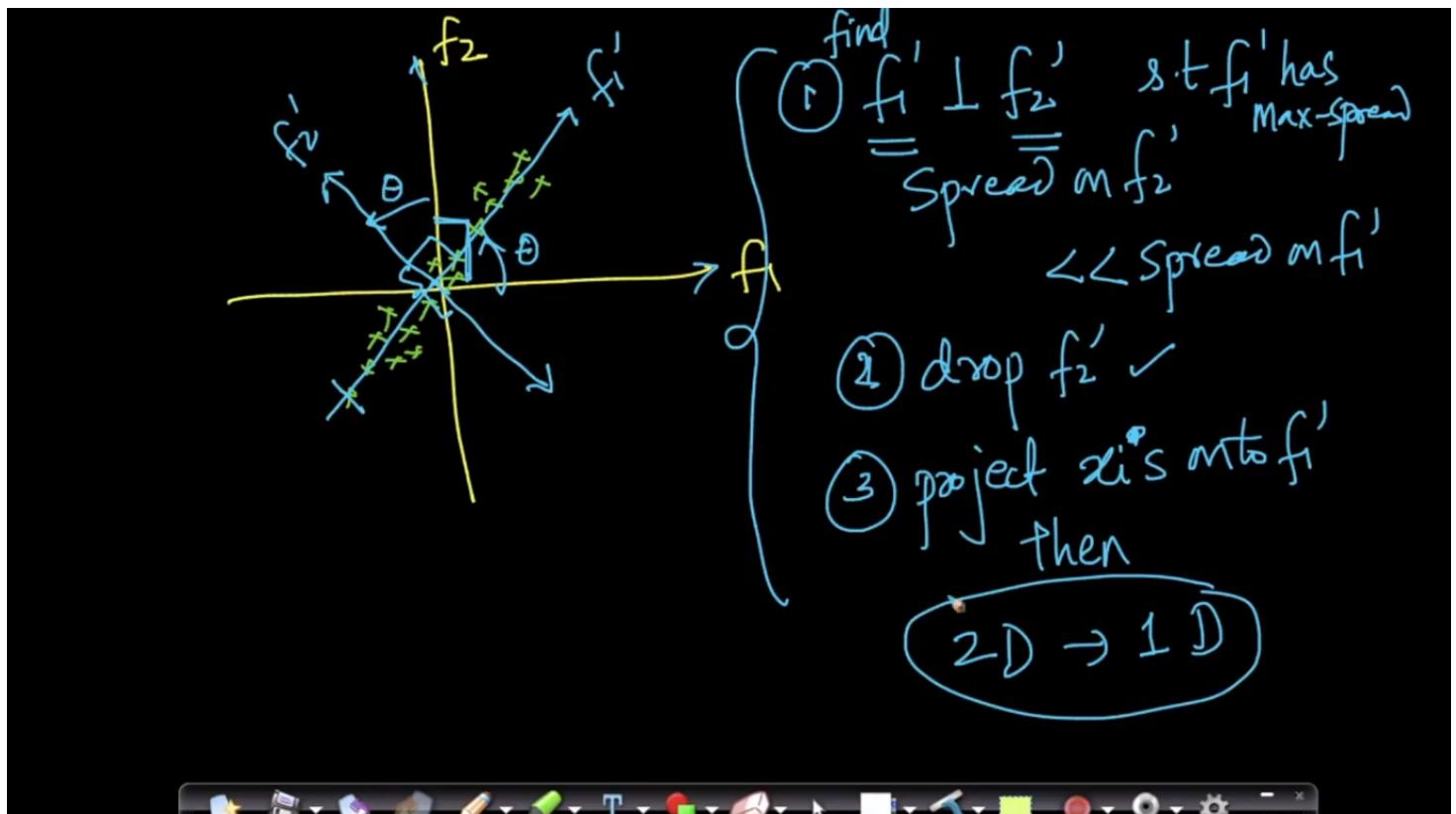
now we represent our data from 2D to 1D means convert the data x to x' , so here we preserving **the direction** with maximal spread or maximal variance so more spread means more information so at the time of converting data high dimensional to low dimensional we only preserve those features giving us more spread or more variance



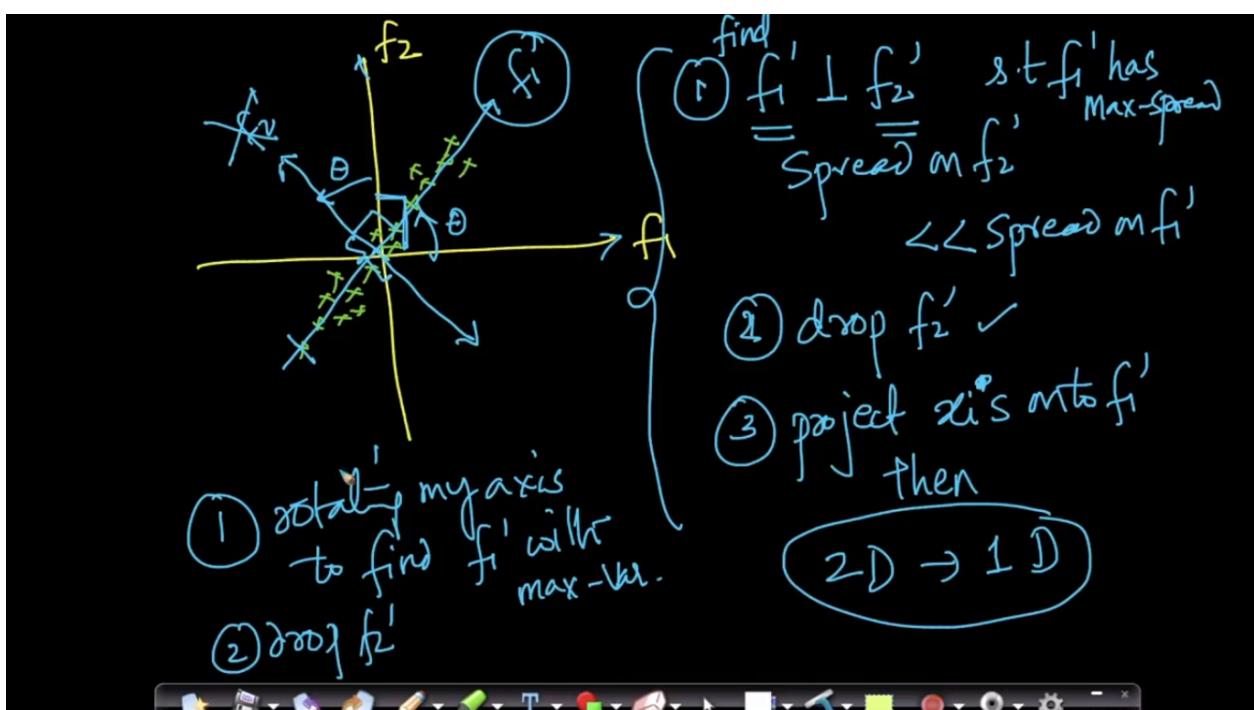
Now lets take another complex example our dataset x is 2 dimensional and dataset has been col. Standardize means mean of f_1 and f_2 is zero and variance of feature f_1 and f_2 is one and our dataset look like this and it is also a two 2d data and from the graph we can see sprade of data in f_1 and f_2 are same and also here we want to convert this 2d data into 1d



now we can see in f_1' direction spread is very high and take one direction f_2' perpendicular to f_1' and spread in f_2' is much much lesser than f_1' and if we drop f_2' and project all the point on f_1' then we can easily convert this 2d to 1d , it is look like we rotating f_1 with some angle θ and rotating f_2 with same angle θ

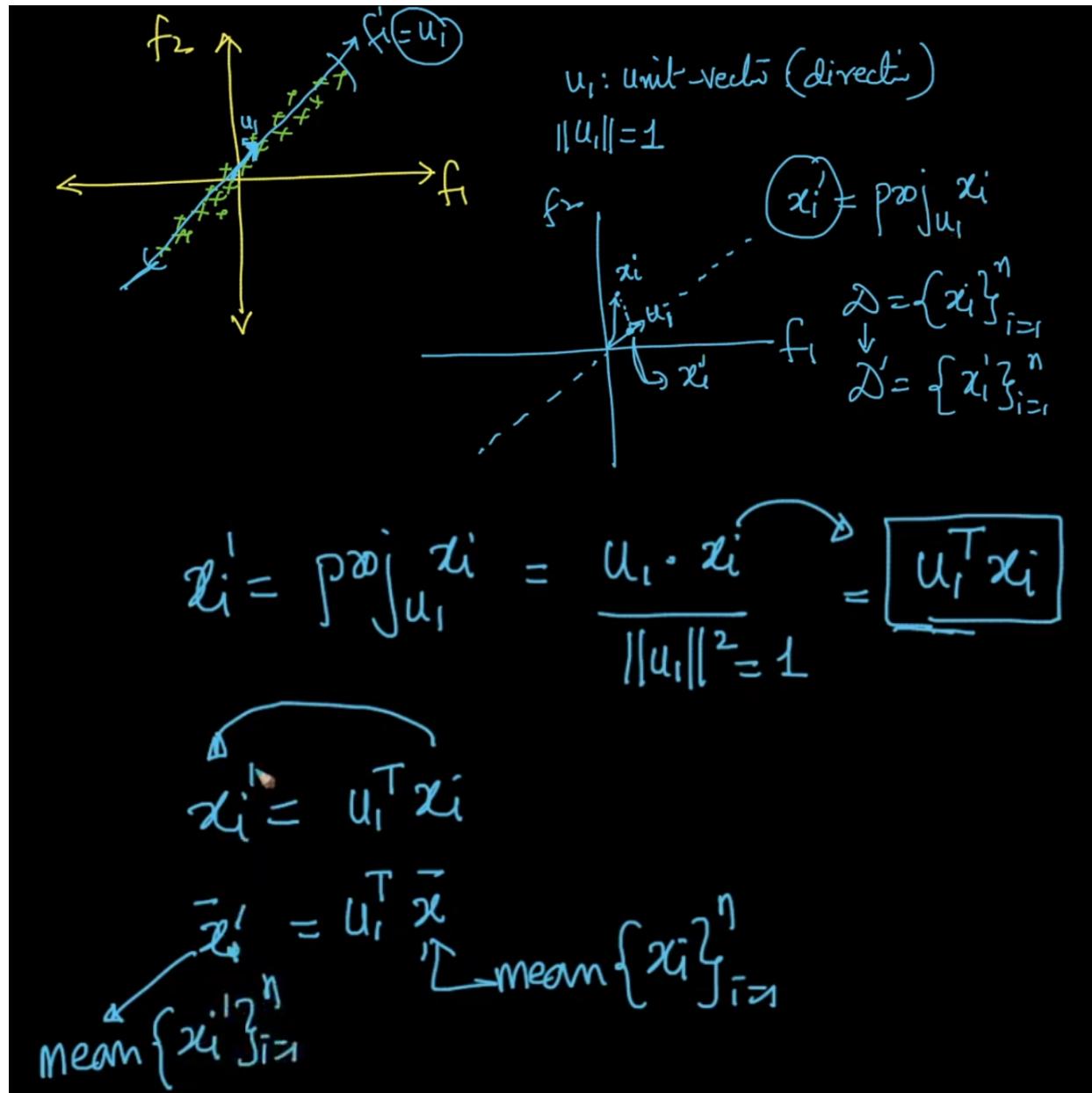


So here we want to find a direction f_1' such that the variance of x_i projected onto f_1' is maximal



Mathematically objective function of PCA

Let say we are assuming $f_1' = u_1$ & we are caring about finding direction of u_1 bcz once we know the direction we can project the points on direction and we know we represent the direction using unit vector means the length of this unit vector(u_1)=1



so given any point x_i we can convert into x_i' using u_1 transpose and we can do same operation by using mean of x .

So now our task, for a point to find u_1 such that.....

④ find u_1 so that $\text{Var}\left\{\text{Proj}_{u_1} \vec{x}_i\right\}_{i=1}^n$ is maximal.

$$\text{Var}\left\{\text{Proj}_{u_1} \vec{x}_i\right\}_{i=1}^n = \frac{1}{n} \sum_{i=1}^n \left(u_1^\top \vec{x}_i - \bar{u}_1^\top \bar{x} \right)^2$$

↓ avg ↓ \vec{x}_i ↓ mean $\{\vec{x}_i\}_{i=1}^n$

$$\text{Scale} = (u_1)^\top \vec{x}_i \quad (\vec{x} \perp)$$

\vec{x} : Col. Standardized

$$\bar{x} = [0, 0, 0, \dots, 0]$$

$$\text{Var}\left\{\vec{x}_i\right\}_{i=1}^n = \frac{1}{n} \sum_{i=1}^n (u_1^\top \vec{x}_i)^2$$

objective
of an
optimization
problem

$$\max_{u_1} \frac{1}{n} \sum_{i=1}^n (u_1^\top \vec{x}_i)^2$$

Data-matrix ✓

$\text{Var}\{\vec{x}_i\}$

Optimization problem

$$\text{so t } u_1^\top u_1 = 1 = \|u_1\|^2$$

↳ Constraint $\overbrace{u_1 \text{ is a unit vector}}$

$$u_1 = [\infty, \infty]$$



so our task is basically maximize the variance \vec{x}_i and we want to find out u_1 , u_1 is the direction that maximizes the variance which is

$$\boxed{\max_{u_1} \frac{1}{n} \sum_{i=1}^n (u_1^\top \vec{x}_i)^2}$$

Data

Alternative Formulation of PCA, distance minimization:

So in previously we saw the formula of maximizing the variance u_1 , so there is another alternative formula that is distance minimization. Means we take the distance of each points to this u_1 line lets call this distance $d_1, d_2, d_3 \dots$ so for each point x_i we can find d_i which is distance from x_i to u_1 . So alternative formulation or optimization function we can write that we want to minimize $\sum d_i^2$ we have to find u_1 which minimizes the sumation of this distances

Alternative formulation of PCA : dist. minimization

→ find u_1 which maximizes projected variance

$x_i \rightarrow (d_i)$: dist from x_i to u_1

$$\min_{u_1} \sum_{i=1}^n d_i^2$$

this is called distance minimization formula of PCA

$$\min_{u_1} \sum_{i=1}^n d_i^2$$

so now we will see how to calculate distance square

u_1 : Unit vector

$$u_1^T u_1 = 1 = \|u\|^2$$

$$d_i^2 = \|x_i\|^2 - (u_1^T x_i)^2$$

$$= x_i^T x_i - (u_1^T x_i)^2$$

so if we put d_i^2 in formula then our distance minimization formula will be

$$\min_{u_i} \sum_{i=1}^n \left(x_i^T x_i - (u_i^T x_i)^2 \right)$$

d_i^2

So we got two different optimization function for PCA one is maximization of variance and another is minimization of distance, two formulas are different but we can achieve same goal using one of them.

Steps Involved in PCA:

1. Standardize the data. (with mean = 0 and variance = 1)
2. Compute the Covariance matrix (S) of dimensions. The off-diagonal elements in the covariance matrix will represent the covariance among each pair of variables and the diagonal elements will represent the variances of each variable/dimension.
3. Obtain the Eigenvectors and Eigenvalues from the covariance matrix. The eigenvalues of the covariance matrix which represents the variability in data on an orthogonal basis in the plot. You will also have to find eigenvectors of the covariance matrix which will represent the direction in which maximum variance among the data occurs.
4. Arrange the eigenvalues in an ascending/descending order and select the higher eigenvalues. You can choose how many eigenvalues you want to proceed with. You will lose some information while ignoring the smaller eigenvalues, but those minute values will not create enough impact on the final result. The selected higher eigenvalues will become the dimensions of your updated feature set.
5. Using the feature vector, we find the principal components of the dataset under analysis. We multiply the transpose of the feature vector with the transpose of the scaled matrix (a scaled version of data after normalisation) to obtain a matrix containing principal components.

Solution to our optimization Problem and Calculate eigen values and eigen vectors :

in $d \times d$ covariance matrix S there are eigen values like $\lambda_1, \lambda_2, \lambda_3, \dots, \lambda_d$ and corresponding eigen vector is $v_1, v_2, v_3, \dots, v_d$ and let also assume that $\lambda_1 \geq \lambda_2 \geq \lambda_3 \geq \dots \geq \lambda_d$ so λ_1 is maximal value, here λ is the scalar and v is vector matrix if they satisfy this condition $\lambda_1 v_1 = S v_1$ then we can say λ_1 is eigen values of co-variance matrix S and v_1 is corresponding eigen vector of co-variance matrix S

$$\begin{array}{l}
 S_{d \times d} \quad \text{Maximal eigen-value} \\
 \downarrow \\
 \lambda_1 > \lambda_2 > \lambda_3 > \lambda_4 \dots > \lambda_d \\
 \left\{ \begin{array}{l}
 \text{eigen-values of } (S) = \lambda_1, \lambda_2, \lambda_3, \lambda_4, \dots, \lambda_d \\
 \text{eigen vectors of } (S) = v_1, v_2, v_3, v_4, \dots, v_d
 \end{array} \right. \\
 \text{def: if } (\lambda, v) = S_{d \times d} v \rightarrow \text{dx1 Vect} \\
 \text{scalar} \quad \text{dx1 Vect} \\
 \boxed{\lambda, v_1 = S v_1} \\
 \checkmark \\
 \lambda_1: \text{eigen value of } S \\
 v_1: \text{eigen vec to } S \\
 \text{corr. to } \lambda_1
 \end{array}$$

also a good properties of eigen vector if $\lambda_1 \geq \lambda_2 \geq \lambda_3 \geq \dots \geq \lambda_d$ and their corresponding vectors are $v_1, v_2, v_3, \dots, v_d$ then we can say $v_1, v_2, v_3, \dots, v_i$ perpendicular to v_j means v_1 perpendicular to v_2 so on so fourth. And then it implies

$$v_i^T v_j = 0 \Rightarrow v_i \cdot v_j = 0$$

So basically u_1 vector is nothing but v_1 , $u_1 = v_1$ and this v_1 is the eigen vector of S (co-variance matrix) corresponding to largest eigen value which is λ_1 , so here u_1 is the maximum variance of direction

$$\lambda_1 > \lambda_2 \geq \lambda_3 \dots \geq \lambda_d$$

$\downarrow \quad \downarrow \quad \downarrow \quad \dots \quad \downarrow$

$$v_1, v_2, v_3, \dots, v_d$$

$S_{d \times d}$

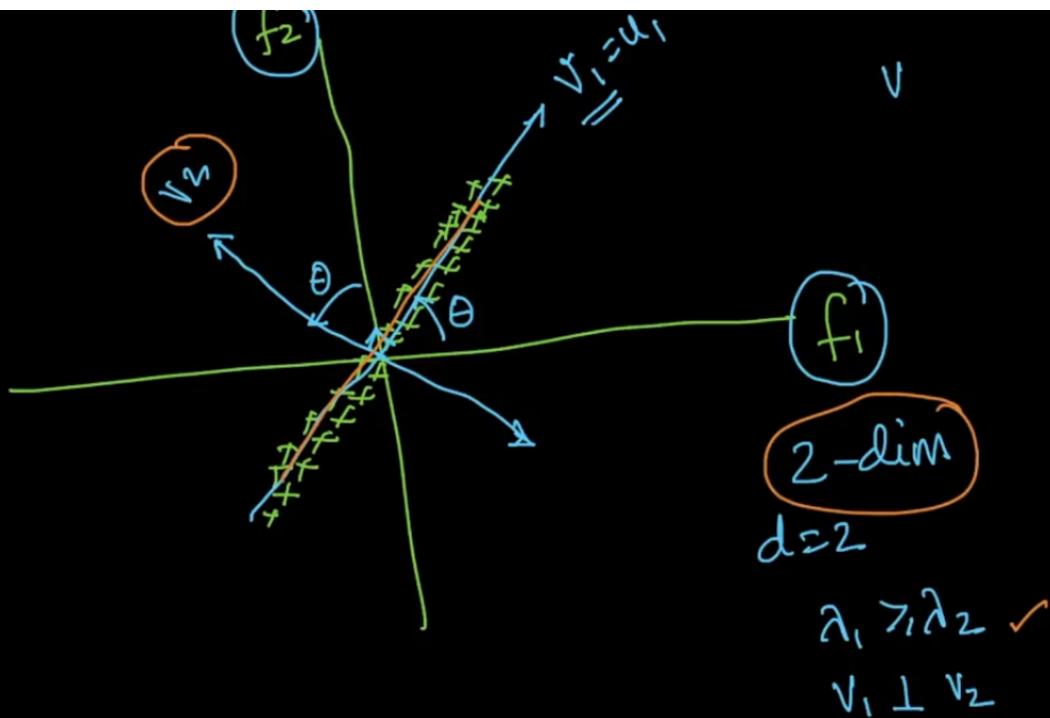
✓ $v_i \perp v_j$: $v_i^T v_j = 0 = v_i \cdot v_j = 0$

✓ $v_1 = v_1 = \text{eigen-vector of } S (= X^T X)$
 corr. to largest eigen-value ($= \lambda_1$)

Max-Variance direction

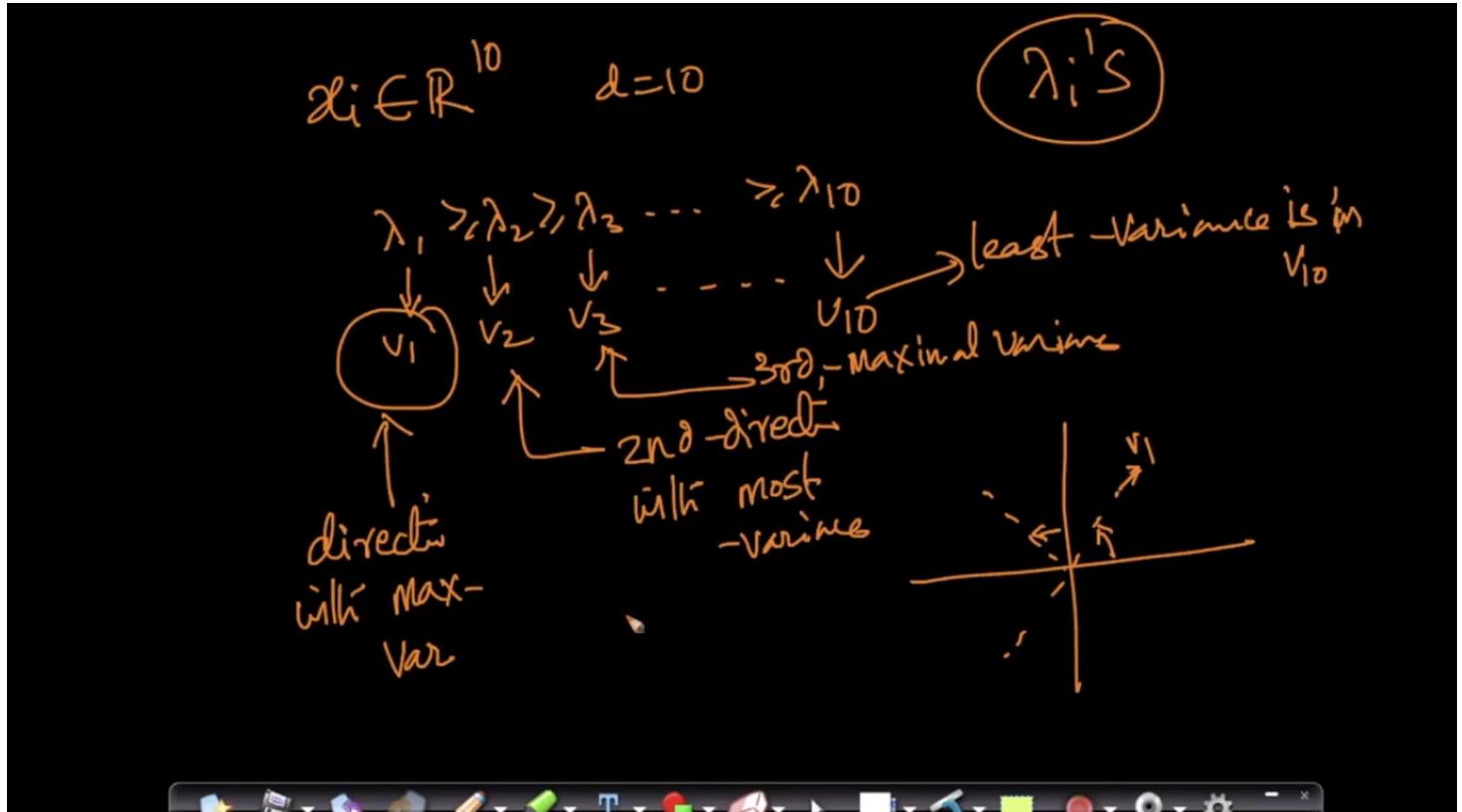
in python we can easily calculate eigen values & eigen vectors.

Let say we have two dimensional data means two eigen values & eigen vectors

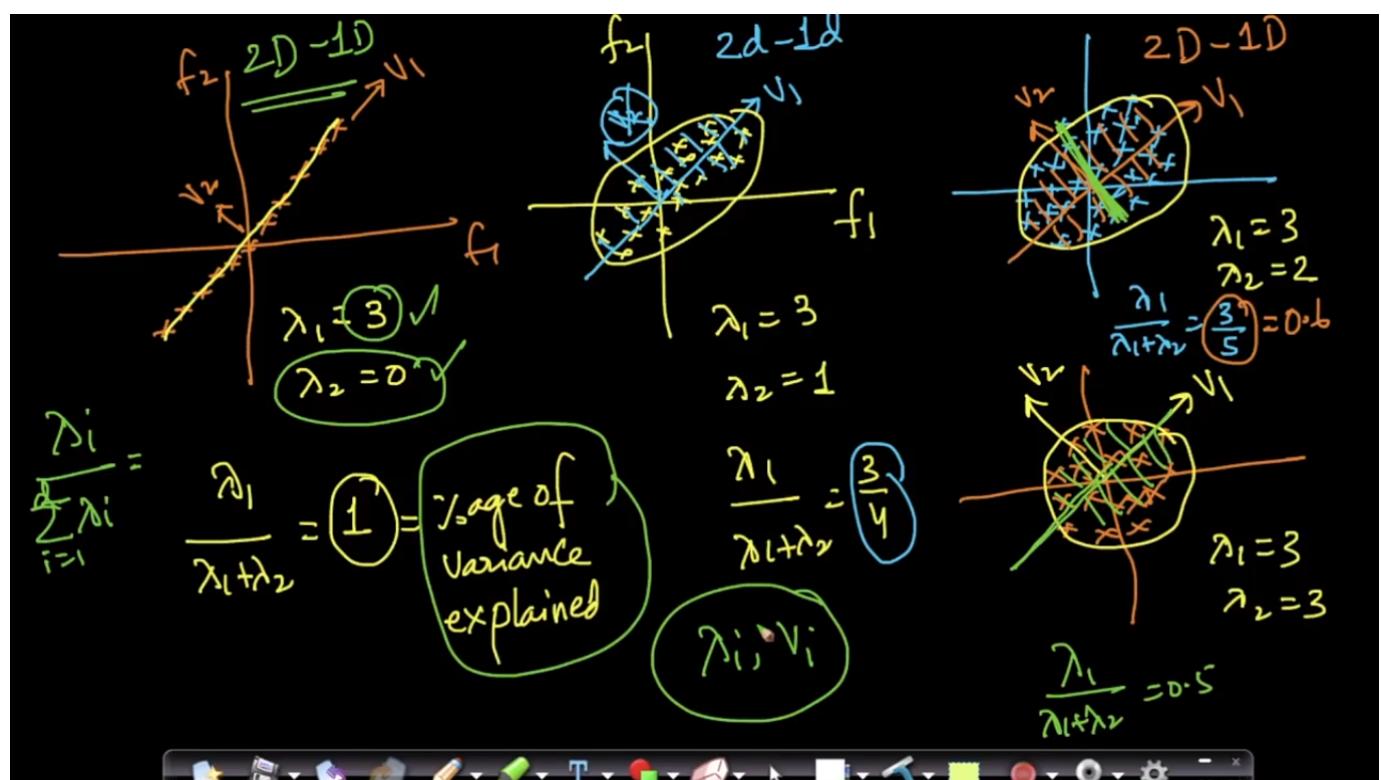


and v_1 and v_2 are perpendicular to each other .

Imagine we have a 10 dimensional data means $d=10$ so we have 10 eigen values & corresponding 10 eigen vectors like $(v_1, v_2, v_3, v_4, \dots, v_{10})$ so v_1 will be direction with maximum variance then v_2 will be 2nd direction with most variance and so on so forth. And v_{10} will be least variance

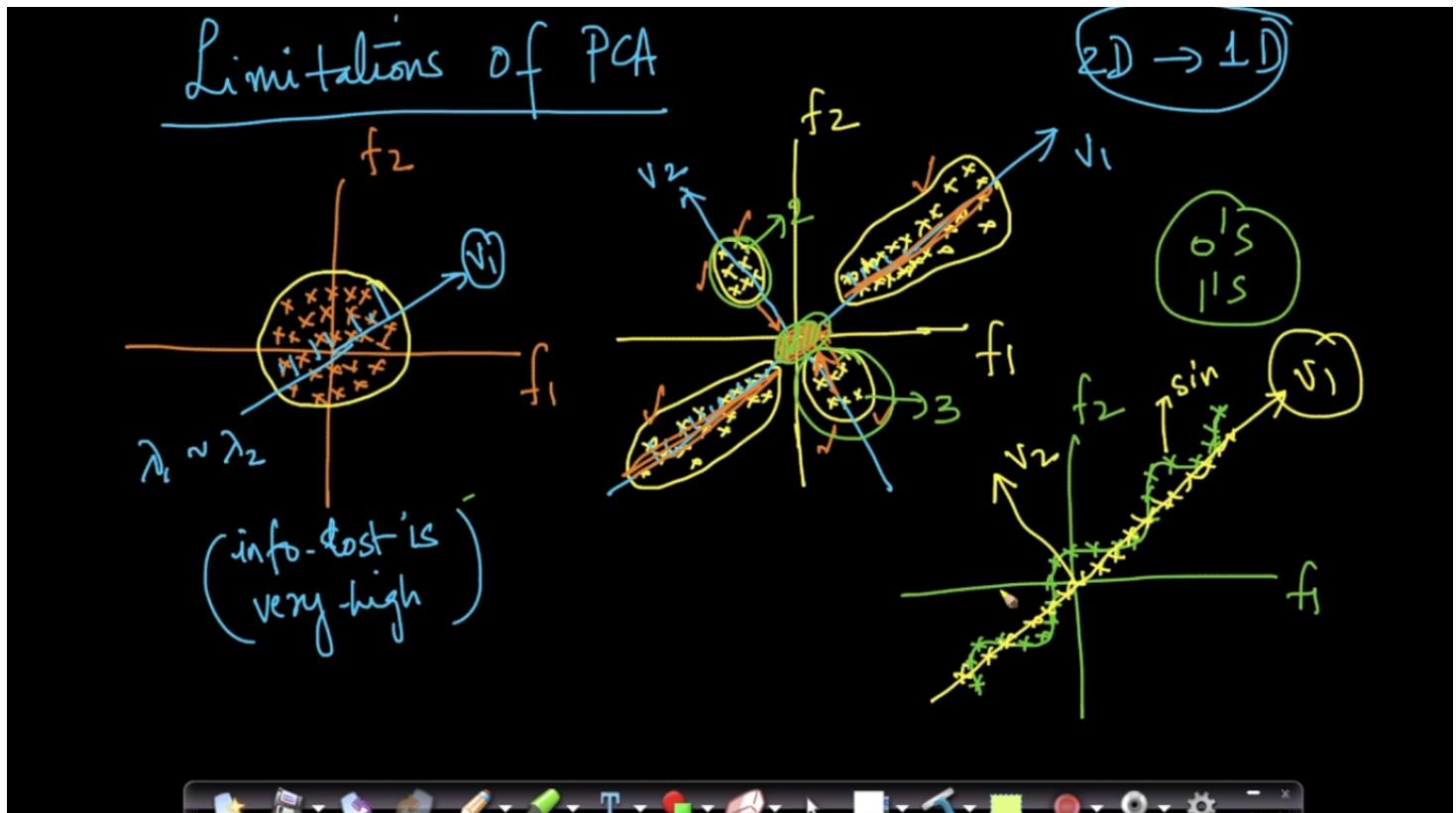


and if we plot this we will see how λ_1 & λ_2 will change for different types of data



now if we want 2d to 1d conversion in first above plot we can see data almost lie on v1 in straight line form and it covers 100% of variance its ratio says 1 means only 1 feature covers 100% variance, & next plot says ratio is 3/4 means if we skip f2 features then only 1 feature will covers 75% variance and so on so fourth.

Limitation of PCA:



- from the above first plot we can see almost same data spread in two axis so if we convert this data 2d to 1d then no doubt we will definitely lost almost 50% information.
- from the above 2nd plot we can see data is more spread on v1 axis and less spread in v2 axis so if we convert this data 2d to 1d then some percentage of data will be overlapping on one axis.
- from the above 3rd plot we can see data has nice sine structure if we convert this data 2d to 1d then we will lose this nice data structure.

Applications of PCA:

- Neuroscience – Neuroscientists use PCA to identify any neuron or to map the brain structure during phase transitions.
- Finance – PCA is used in the finance sector for reducing the dimensionality of data to create fixed income portfolios. Many other facets of the finance sector involve PCA like forecasting returns, making asset allocation algorithms or equity algorithms, etc.

- Image Technology – PCA is also used for image compression or digital image processing. Each image can be represented via a matrix by plotting the intensity values of each pixel, and then we can apply PCA on it.
- Facial Recognition– PCA in facial recognition leads to the creation of eigenfaces which makes facial recognition more accurate.