

# Principal Component Analysis

lec 1 To convert high dimensional dataset to smaller dimension.

for eg- MNIST dataset was 784 dimensions and we convert it to 2 dimensions.

The Principal Component Analysis is used there.

lec 2 Geometric Intuition of PCA.

We have to convert d-dimensions to d' dimensions, where  $d' < d$ .

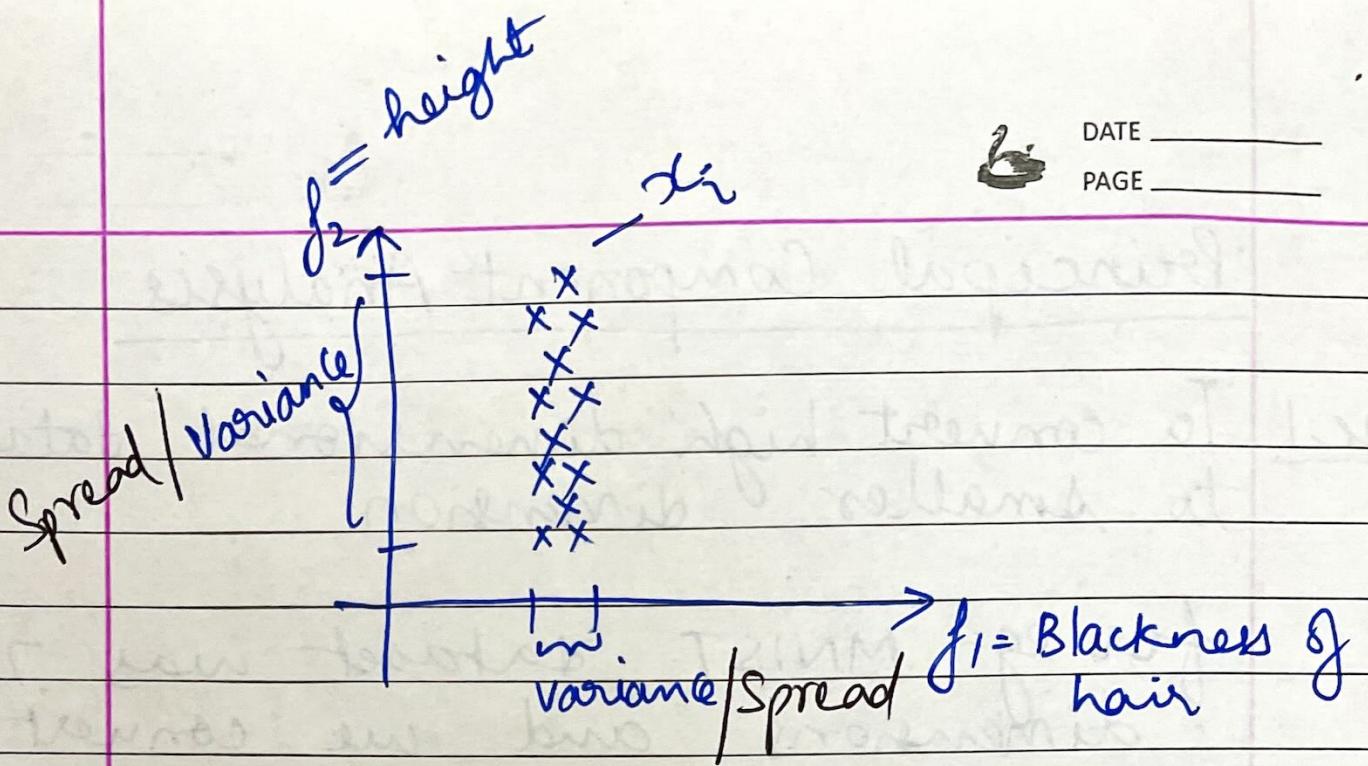
Say  $d = 2$ -dim and  $d'$  is 1dim.

Let's assume two features  $f_1$  and  $f_2$ .

where  $f_1$  = Blackness of hair  $\in \mathbb{R}$  and  $f_2$  is Height of people.

The data matrix can be represented as

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



To reduce this 2 dim. data to 1 dim we can skip  $f_1$  and  $f_2$  as it has high variance and more information.

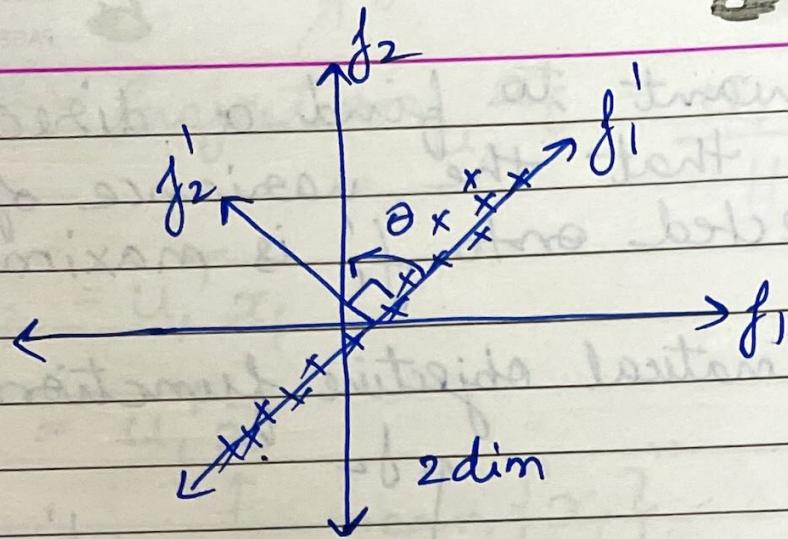
$$x' = \begin{bmatrix} f_2 \\ ? \\ ? \\ \vdots \\ n \end{bmatrix}$$

You lose less info in  $f_1$  so skip  $f_1$ .  
You lose more info in  $f_2$  so don't skip  $f_2$ .

Preserving the direction with maximal spread/variance. More variance means more information.

A different dataset now:-  
Let  $X = 2$  dim dataset.

It is column standardized which means  
So,  $\text{mean } \{f_1\} = \text{mean } \{f_2\} = 0$   
 $\text{var } \{f_1\} = \text{var } \{f_2\} = 1$



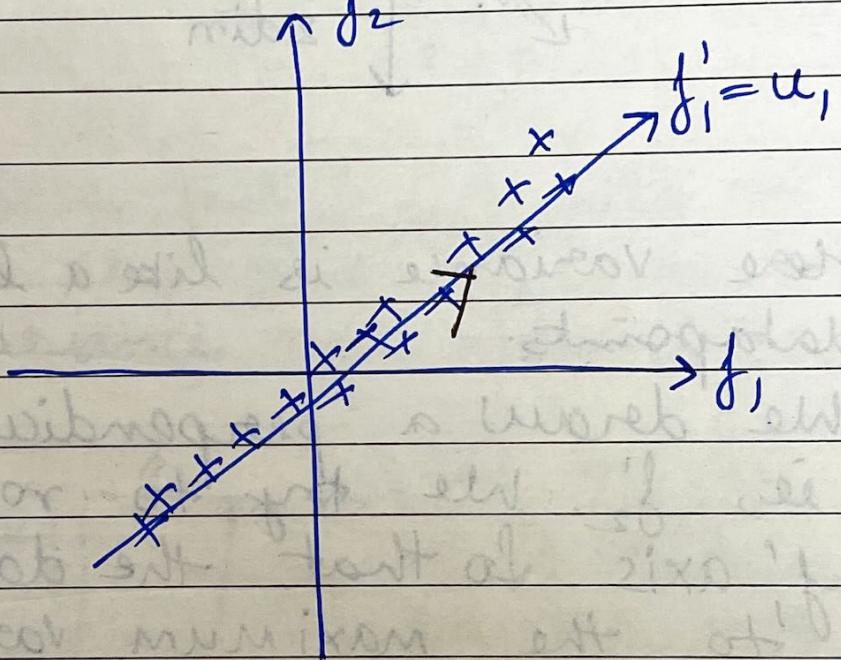
Here variance is like a line with datapoints.

- ① We draw a perpendicular to  $f_1'$ , ie,  $f_2'$ . We try to rotate the  $f_1'$  axis so that the datapoints get to the maximum variance.  
 $\therefore$  Spread of mean of  $f_2' \ll$  spread of mean  $f_1'$
- ② Drop  $f_2'$ .
- ③ Project  $x_i$ 's mean onto  $f_1'$  then 2D to 1D.  
such that  $f_1'$  has max variance.

- ① we want to find a direction  $f'_i$  such that the variance of  ~~$x_i$~~   $x_i$  projected onto  $f'_i$  is maximum.

Lec 3

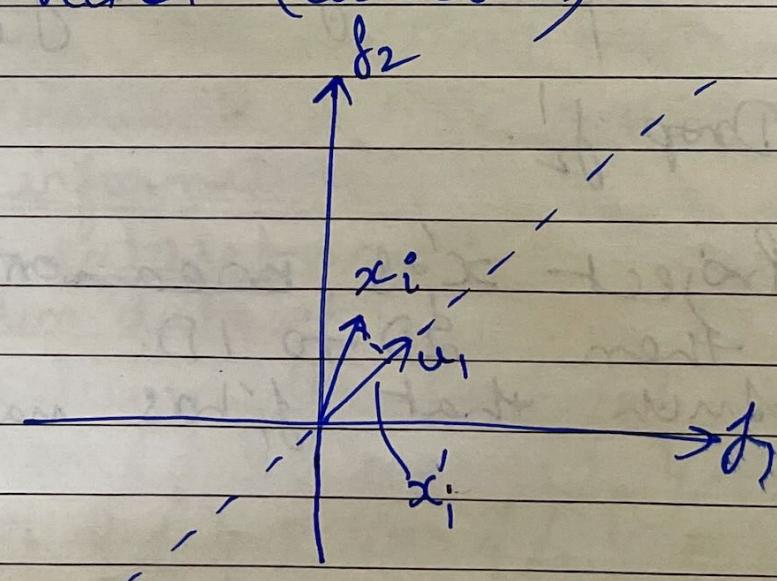
Mathematical objective function of PCA.



Find direction of  $u_1$

$u_1$  = unit vector (direction)

$$\|u_1\|=1$$



projection of  $x_i$  onto  $u_1$

DATE \_\_\_\_\_  
PAGE \_\_\_\_\_

$$x'_i = \text{projection}_{u_1} x_i = \frac{u_1 \cdot x_i}{\|u_1\|} = u_1^T x_i$$

$$x'_i = u_1^T x_i$$

$$\bar{x}' = u_1^T \bar{x}$$

$$\text{mean}\{x'_i\}_{i=1}^n$$

$u_1$  is the

problem: Find  $u_1$  st.  $\text{Var}\{\text{proj}_{u_1} x_i\}_{i=1}^n$  is maximal

$$\text{variance}\left\{u_1^T x_i\right\}_{i=1}^n = \frac{1}{n} \sum_{i=1}^n \underbrace{\left(u_1^T x_i - \bar{u}_1^T \bar{x}\right)^2}_{\text{aug. } x'_i - \text{mean of } x'_i}$$

$$\text{Scalar} = (u_1)^T x_i \underset{(1 \times n)}{\underset{(1 \times 1)}{\underset{\diagdown}{}}}$$

$$\text{Variance of } \{x_i'\}_{i=1}^n = \frac{1}{n} \sum_{i=1}^n (u^T x_i)^2$$

We want to maximize Variance, ie

$$\max_{u_1} \frac{1}{n} \sum_{i=1}^n (u^T x_i)^2$$

↑  
Data matrix

objective  
of an  
optimization  
problem

such that  $u^T u = 1 = \|u\|^2$

constraint

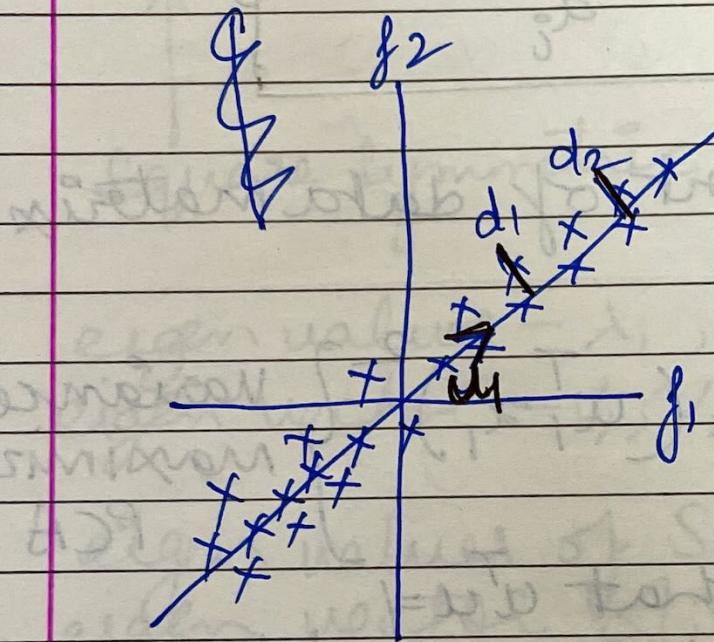
$u_1$  is a unit vector  
 $u_1 = [\infty, \infty]$ .

Find  $u_1$  which maximizes projected variance.

lec9

## Alternative Formulation of PCA: minimization

→ Let's assume our points



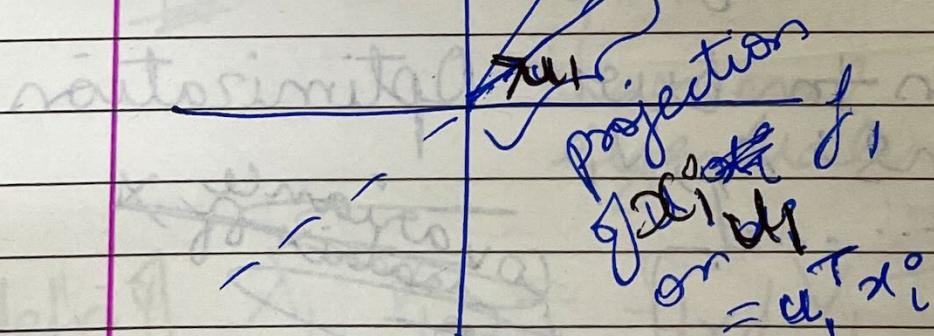
for each point  $x_i^0$  we get  $d_i^0$  which is distance from  $x_i^0$  to  $u_1$

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

Assumptions -

$u_1$  = unit vector

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



$$\begin{aligned} d_i^2 &= \|x_i^0\|^2 - (u_1^T x_i^0)^2 \\ &= x_i^T x_i - (u_1^T x_i)^2 \end{aligned}$$

DATE \_\_\_\_\_  
PAGE \_\_\_\_\_

optimization problem for distance.

objective function.

$\min_{u_1} \sum_{i=1}^n \left( x_i^T x_i - (u^T x_i)^2 \right)$

$\text{constraint: } u^T u = 1$

$d_i^2$

distance min PCA

$x_i$  are part of data matrix

We have done this

$\max_{u_1} \frac{1}{n} \sum_{i=1}^n (u^T x_i)^2$

Variance maximization PCA

such that  $u^T u = 1$

Constraint  $u^T u$  is same in both the optimization of variance maximization & distance minimization

lec 5 Eigen Values and Eigen Vectors (PCA)  
Dimensionality Reduction.

Solutions to our Optimization problems:  $\lambda_1, v_1, \dots, \lambda_d, v_d$

Data Matrix  $X \in \mathbb{R}^{n \times d}$

Column Standardized  $\mu = 0, \sigma^2 = 1$

$n \times d$

~~Covariance matrix~~  $\Sigma \in \mathbb{R}^{d \times d}$

$d$  features &  $n$  rows of  $X$

# Covariance matrix of $X = S$

$$S_{d \times d} = X_{d \times n}^T X_{n \times d}$$

↑  
Square Symmetric matrix

eigen values -  $\lambda_1, \lambda_2, \lambda_3, \dots, \lambda_d$   
eigen vector -  $v_1, v_2, v_3, \dots, v_d$

eigen values of  $S = \lambda_1, \lambda_2, \dots, \lambda_d$   
eigen vector of  $S = v_1, v_2, \dots, v_d$

Every eigen value has an eigen vector  
where  $\lambda_1 > \lambda_2 > \lambda_3 > \lambda_4 > \dots > \lambda_d$

~~Maximal eigen values~~

def'n of  $\lambda_1, v_1 = S_{d \times d} v_1 \rightarrow d \times 1 \text{ vector}$

eigen values & scalar column vectors  
and  $v_1$  is corresponding to eigen vector

then  $\lambda_1$  = eigen value

and  $v_1$  is corresponding to eigen vector

numpy gives eigen values of an Matrix  $X$

## Property of Eigen vectors

$$\rightarrow v_i \perp v_j$$

which,  $v_i^T v_j = 0$  or  $\underbrace{v_i, v_j}_\text{dot product} = 0$

dot product

$u_i = v_i$  = eigen vector of  $S (= X^T X)$   
maximal variance direction  
Corresponding to largest eigen value ( $\lambda_i$ )

$$X = \begin{bmatrix} & & \\ & & \\ & & \end{bmatrix}$$

Steps to find  $u_i$        $n \times d$

1. Do column standardization of  $X$

$$S = \frac{1}{n} (X^T X)$$

2. Compute Eigen Values & Eigen Vectors of  $S$ . Transpose the vector.  
DO  $(\sqrt{\text{natural}} X^T)$

$$u_i = v_i$$

eigen vector

You will get 2-D Data.

Use for dimensionality reduction als

Blue points  $\rightarrow$  To find max variance (max)

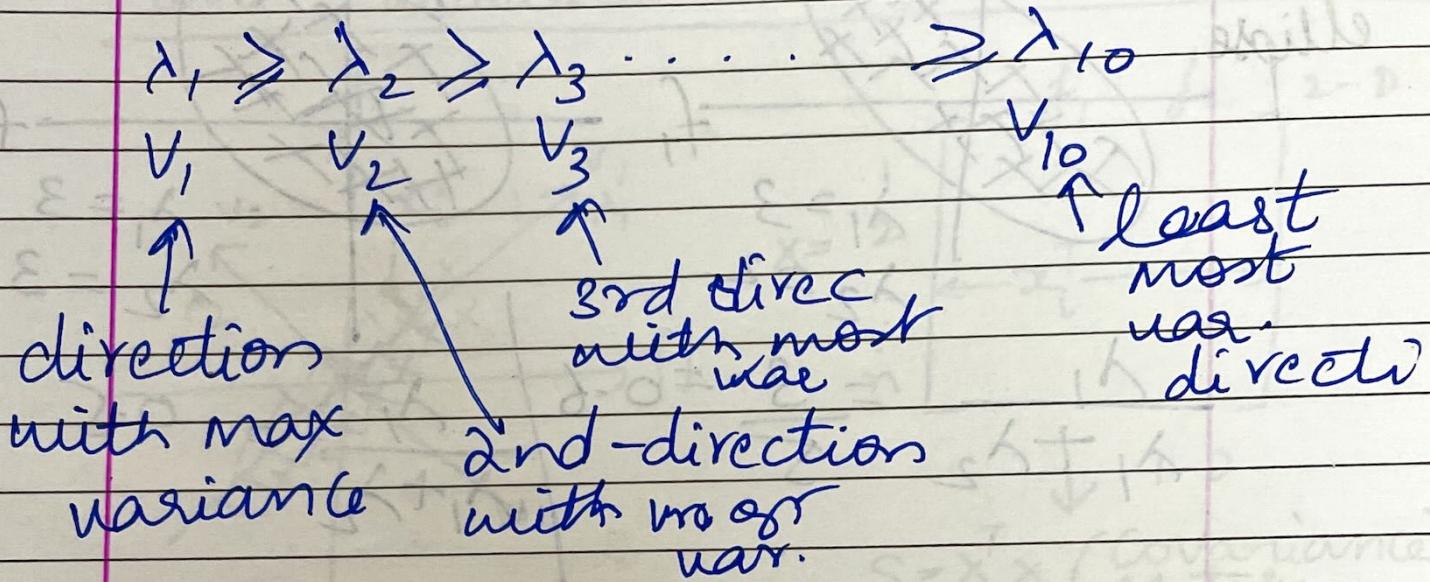
Black points  $\rightarrow$  Dim. Reduction.

For  $d$  dim data, we have  $d$  eigen values and corresponding  $d$  eigenvectors

DATE \_\_\_\_\_

PAGE \_\_\_\_\_

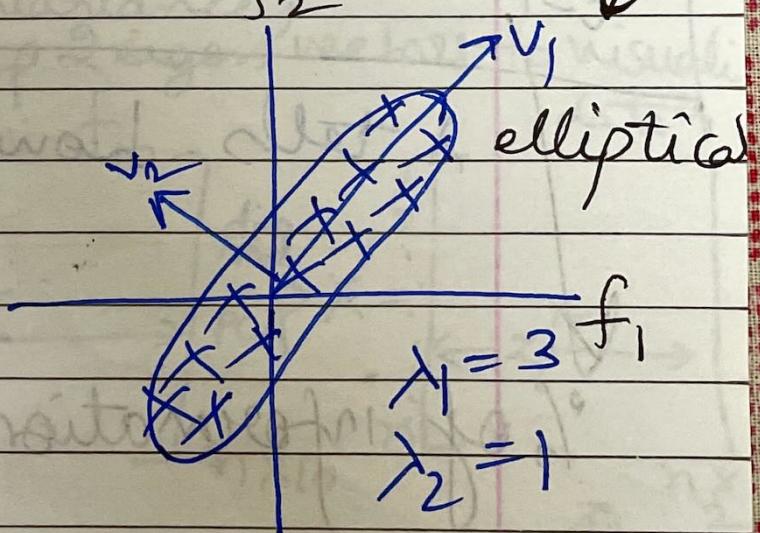
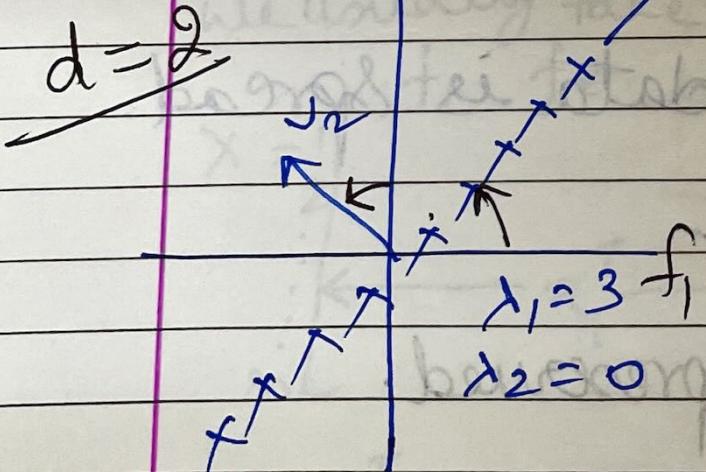
$$x_i \in \mathbb{R}^d \text{ If, } d=10$$



Geometric explanations of Eigenvalues & Eigen Vectors

$\sum_{i=1}^d \lambda_i$  Linear spread  $f_2$

Some spread  $f_2$

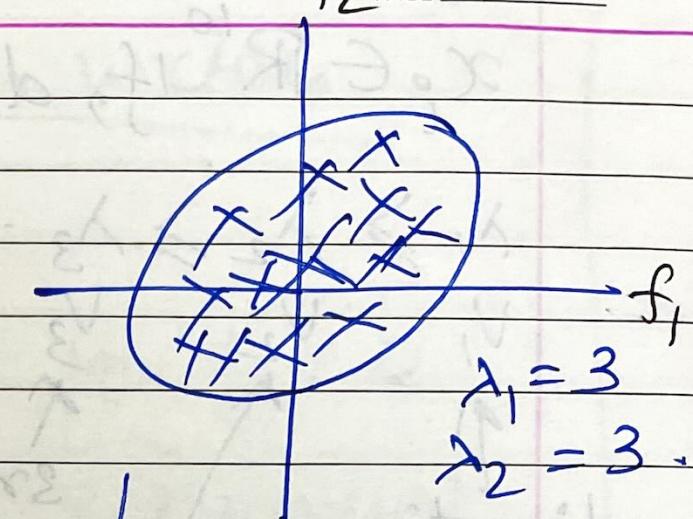
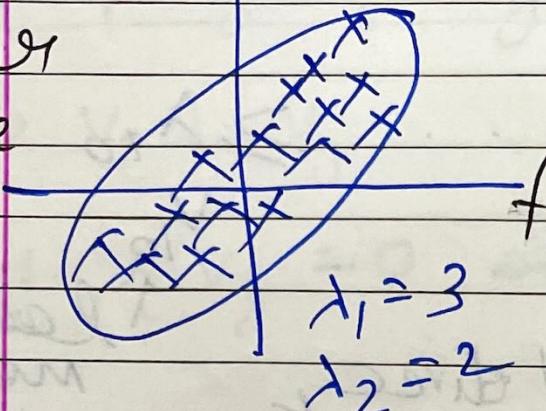


$$\frac{\lambda_1}{\lambda_1 + \lambda_2} = \frac{1}{4} \text{ of variance explained}$$

$$\frac{\lambda_1}{\lambda_1 + \lambda_2} = \frac{3}{4} \text{ of variance explained}$$

$f_2$ 

Thicker ellipse



$$\frac{\lambda_1}{\lambda_1 + \lambda_2} = \frac{3}{5} = 0.6 \quad \left| \quad \frac{\lambda_1}{\lambda_1 + \lambda_2} = 0.5$$

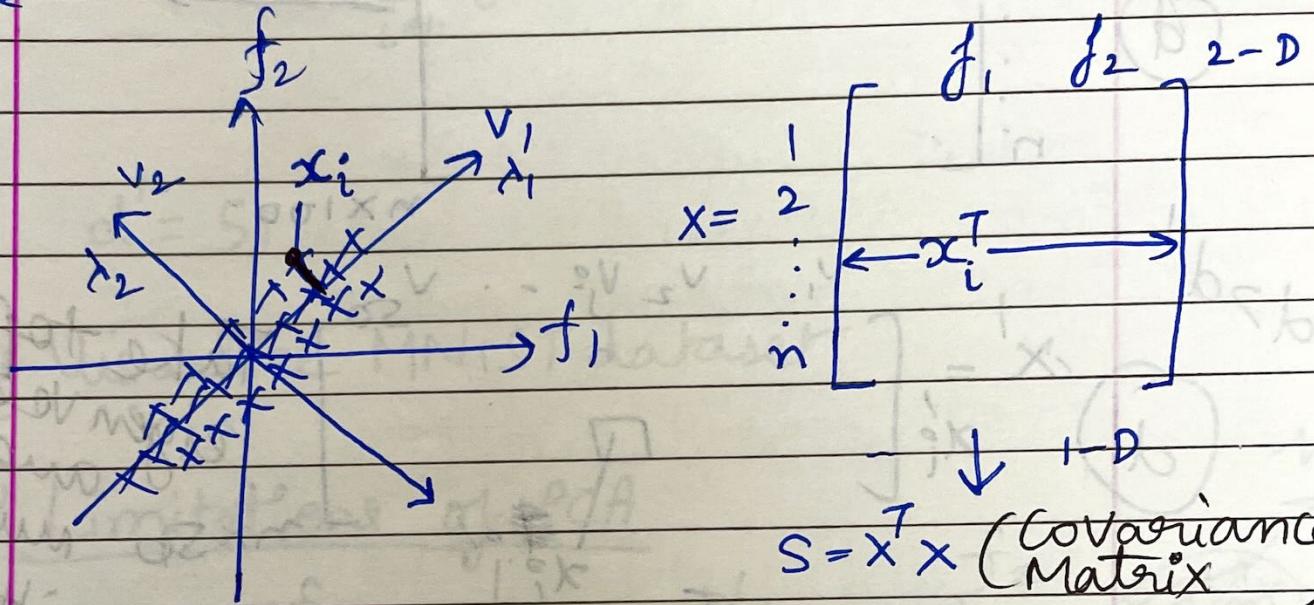
$$\frac{\lambda_i}{\sum_{i=1}^n \lambda_i} = \text{Percentage of variance explained.}$$

$\lambda$  tells how data is spread.

% of information preserved.

$\lambda$  tells whether variance is spread on eigen value  $\lambda$  no. of axes.

# lec6 PCA for dimensionality reduction and visualization



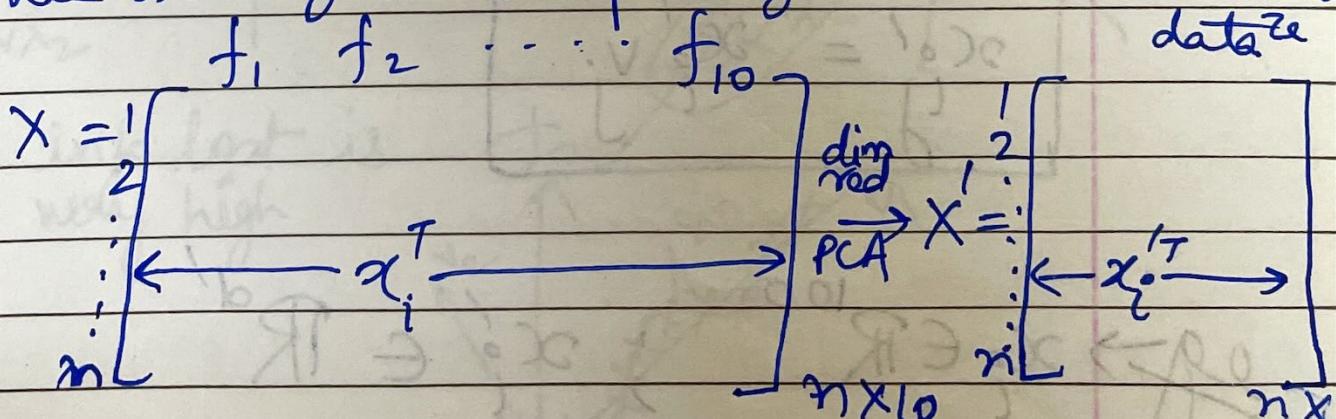
$$S = X^T X \text{ (Covariance Matrix)}$$

$$x'_i = x_i^T v_1 \quad \text{For } v_1 \text{ to maximize variance, we take top 2 eigenvectors.}$$

By projecting  $x'_i$  on  $v_1$

(which is maximal variance)

we usually take top 2 eigenvectors to visualize data.



$$S = X^T X$$

$$\text{eigen}(S) = \lambda_1 \geq \lambda_2 \geq \lambda_3 \dots \geq \lambda_{10} \quad \left| \begin{array}{l} x'_i = [x_i^T v_1 \\ x_i^T v_2] \\ \vdots \\ f_1 \\ f_2 \end{array} \right.$$

①  
②

③  
④

$f_1, f_2, \dots, f_{100}$

DATE \_\_\_\_\_  
PAGE \_\_\_\_\_

(a)

$$X = \begin{bmatrix} 1 \\ 2 \\ \vdots \\ n \end{bmatrix} \quad \xrightarrow{\text{PCA}} \quad n \times 100$$

(d)

$$d > d' \quad X' = X'_o \begin{bmatrix} v_1 & v_2 & v_j & \dots & v_{50} \end{bmatrix}_{n \times 50} \quad \text{Take top eigen vectors and multiply}$$

$S_{100 \times 100} = X^T X$

$\lambda_1 \geq \lambda_2 \geq \dots \geq \lambda_{100}$

$v_1, v_2, \dots, v_{100}$

$x_o' = x_o^T v_j$

with  $X^T$  you get top features

$\Rightarrow x_o \in \mathbb{R}^{100}; x'_o \in \mathbb{R}^{d'}$

$d' < 100$

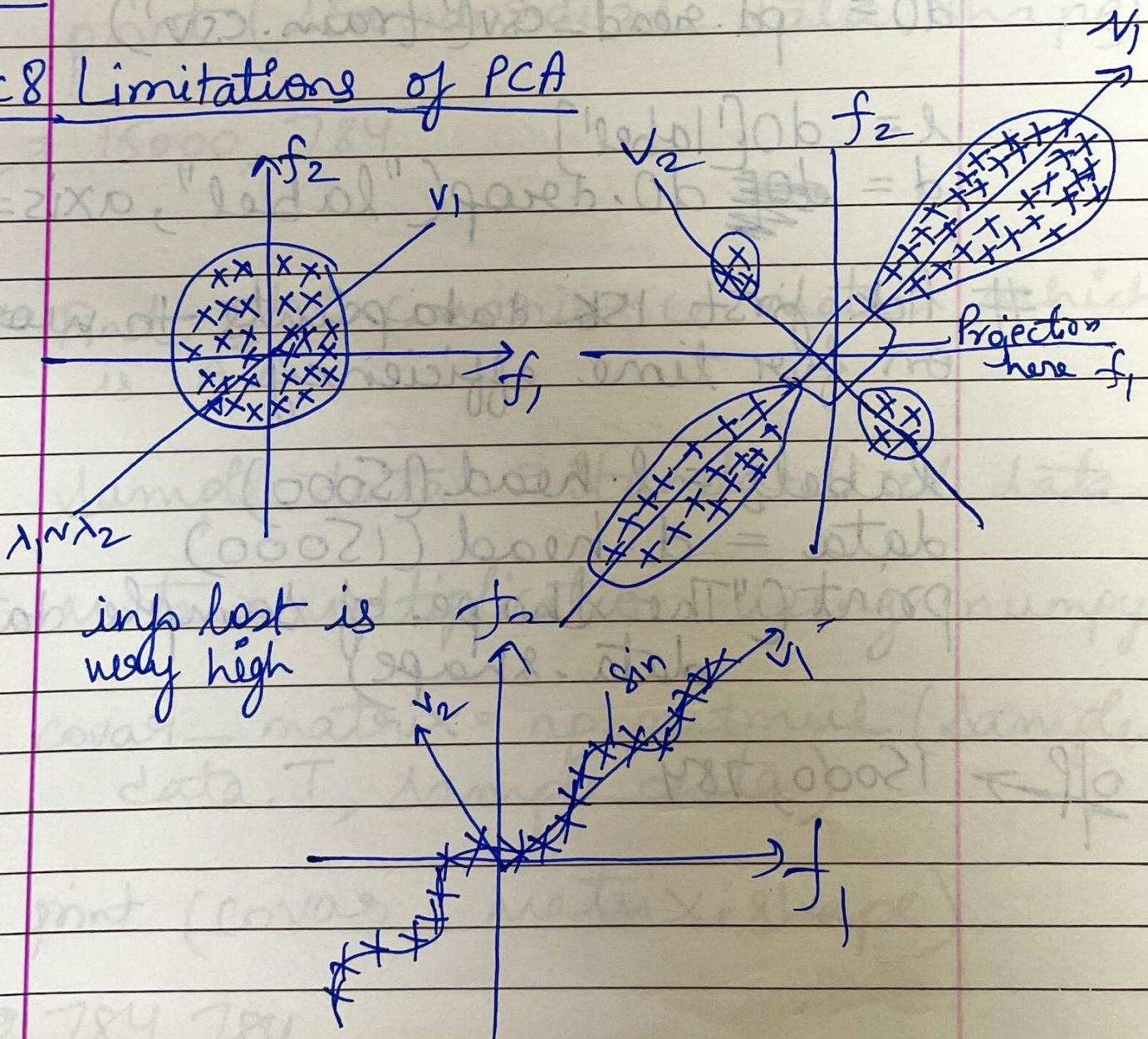
Preserve 99% of the variance

Let  $\frac{\lambda_1 + \lambda_2 + \lambda_3 + \dots + \lambda_{57}}{\sum_{i=1}^{100} \lambda_i} = 0.99$

$$d' = 51$$

## lec 7 Visualize MNIST dataset

## lec 8 Limitations of PCA





## lec9 PCA code example

# imports.

```
import numpy as np  
import pandas as pd  
import matplotlib.pyplot as plt
```

```
d0 = pd.read_csv('train.csv')
```

```
l = d0['label']
```

```
d = d0.drop("label", axis=1)
```

# pick first 15K datapoints to work  
on for time efficiency.

```
labels = l.head(15000)
```

```
data = d.head(15000)
```

```
print("The shape of sample data",  
      data.shape)
```

~~off~~ → 15000, 784



DATE \_\_\_\_\_

PAGE \_\_\_\_\_

# Data preprocessing: Standardizing  
the data

from sklearn.preprocessing import  
StandardScaler

standardized\_data = StandardScaler.  
fit\_transform(data)

print(standardized\_data.shape)

~~off~~ → 15000, 784

# Find the covariance matrix which  
is  $A^T A$

sample\_data = standardized\_data

→ Matrix multiplication using numpy

covar\_matrix = np.matmul(sample  
data.T, sample\_data)

print(covar\_matrix.shape)

~~off~~ → 784, 784

# Finding the top 2 eigen values  
and corresponding eigen vectors  
for projecting onto 2-Dim  
Space.

# The parameter 'eigvals' is defined  
(low value to high value) eigh  
function will return the eigen  
values in ascending order.

# This code generates only top  
2 (782, 783) eigenvalues.

from scipy.linalg import eigh  
values, vectors = eigh (covar\_matrix  
eigvals = (782, 783)  
point(vectors.shape)

op> 784, 2

vectors = vectors.T  
point(vectors.shape)

op> 2, 784

# Projecting the original data sample  
on the plane formed by two  
principal eigenvectors by vector-vector  
multiplication

import matplotlib.pyplot as plt  
new\_coordinates = np.matmul(eigenvectors,  
sample\_data.T)

of  $2,784 \times 784,15000 = 2,15000$

import pandas as pd

# appending label to the 2D projected  
data

new\_coordinates = np.vstack((new\_coordinates,  
labels)).T

# Creating a new data frame for  
plotting the labelled points.

dataframe = pd.DataFrame(data=  
new\_coordinates, columns= ("1st principal",  
"2nd principal", "label"))  
print(dataframe.head())

# plotting the 2d. data points with Seaborn.

```
import seaborn as sns  
sns.FacetGrid(dataframe, hue="label",  
size=6).map(plt.scatter, "1st principal",  
"2nd principal").add_legend()  
plt.show()
```

## PCA using Scikit Learn.

# initializing the PCA

```
from sklearn import decomposition  
pca = decomposition.PCA()
```

# Configuring the parameters the no. of parameters = 2

```
pca.n_components = 2
```

```
pca_data = pca.fit_transform(  
sample_data)
```

# pca reduced will contain the  
2-d projects of simple data

print(pca\_data.shape)

~~'op> 15000, 2~~

# attaching the label for each 2-D  
data point

pca\_data = np.vstack((pca\_data.T,  
labels)).T

# Creating a new data frame which  
help us in plotting the result data

pca\_df = pd.DataFrame(data=pca\_data,  
columns=["1st principal", "2nd principal",  
"label"])

sns.FacetGrid(pca\_df, hue="label",  
size=6).map(plt.scatter, "1st principal",  
"2nd principal").add\_legend()  
plt.show()

lec10

PCA for dimensionality reduction  
not for ~~code~~ visualization.

# initializing the PCA

```
from sklearn import decomposition  
pca = decomposition.PCA()
```

# Configuring the no. of parameters = 784.

```
pca.n_components = 784
```

```
pca_data = pca.fit_transform(  
sample_data)
```

percentage\_var\_explained =

```
pca.explained_variance_ /
```

```
np.sum(pca.explained_variance_)
```

cum\_var\_explained = np.cumsum  
(percentage\_var\_explained)

# Plot the PCA spectrum

DATE \_\_\_\_\_

PAGE \_\_\_\_\_

