

☆ Dimensionality Reduction (Usually used for unsupervised Learning)

2D, 3D : Scatterplots
4D, 5D, 6D : pair plot

but what if
10-D, 100-D, 1000-D plots will not work
so to visualize data we try to reduce the dimensionality.

☆ Row vector and column vector

flower [SL, PL, SW, PW]
 └──┬──┘
 real values

$\mathbb{R} \rightarrow$ Real space | Real numbers

i^{th} point : $x_i \in \mathbb{R}^d$
 $x_i = \begin{bmatrix} x_{i,1} \\ x_{i,2} \\ \vdots \\ x_{i,d} \end{bmatrix}_d$: column vector

$x_i = [2.1, 3.2, 4.1, 1.2] \rightarrow$ data point (row vector)

default representation is column vector

☆ How to represent a dataset?

$\mathcal{D} = \{x_i, y_i\}_{i=1}^n$ $n =$ number of datapoint

$x_i \in \mathbb{R}^d$
 $y \in \{ \text{\textcolor{red}{\#1}}, \text{\textcolor{red}{\#2}}, \text{\textcolor{red}{\#3}} \}$
 └──┬──┘
 classifications / labels

Data matrix

$$X = \begin{matrix} & f_1 & f_2 & f_3 & f_4 & \dots & f_d \\ \begin{matrix} 1 \\ 2 \\ 3 \\ \vdots \\ n \end{matrix} & \left[\begin{array}{cccccc} & & & & & \\ & & & & & \\ & & & & & \\ & & & & & \\ & & & & & \\ & & & & & \end{array} \right] \end{matrix}$$

$\xleftarrow{\quad} x^T_i \xrightarrow{\quad}$

$n \times d$

each datapoint : row
each column : feature (specific)

☆ Data Preprocessing : Feature Normalisation :

Column $a_1, a_2, \dots, a_n \rightarrow n$ -values of f_1

$$\max(a_i) = a_{\max} \geq a_i$$

$$\min(a_i) = a_{\min} \leq a_i$$

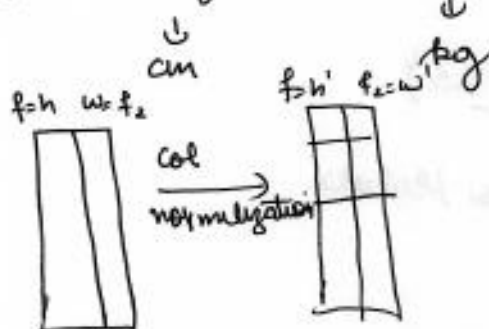
$$a'_i = \frac{a_i - a_{\min}}{a_{\max} - a_{\min}} \quad a'_i \in [0, 1]$$

$$a'_{\min} = \frac{a'_{\min} - a'_{\min}}{a_{\max} - a_{\min}} = 0$$

$$a'_{\max} = \frac{a_{\max} - a_{\min}}{a_{\max} - a_{\min}} = 1$$

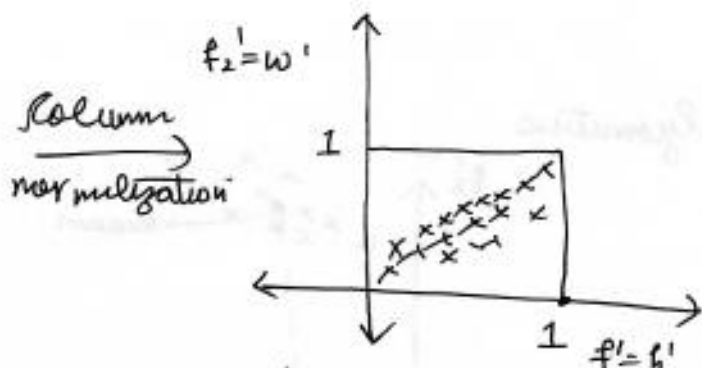
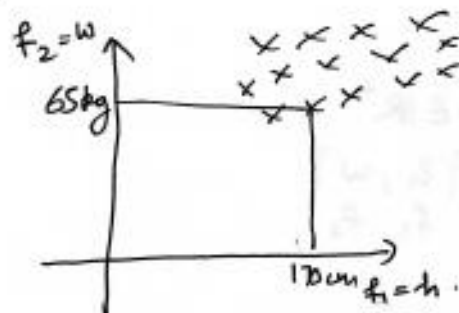
Transform the given data into a' .

why? height and weight as 2 features

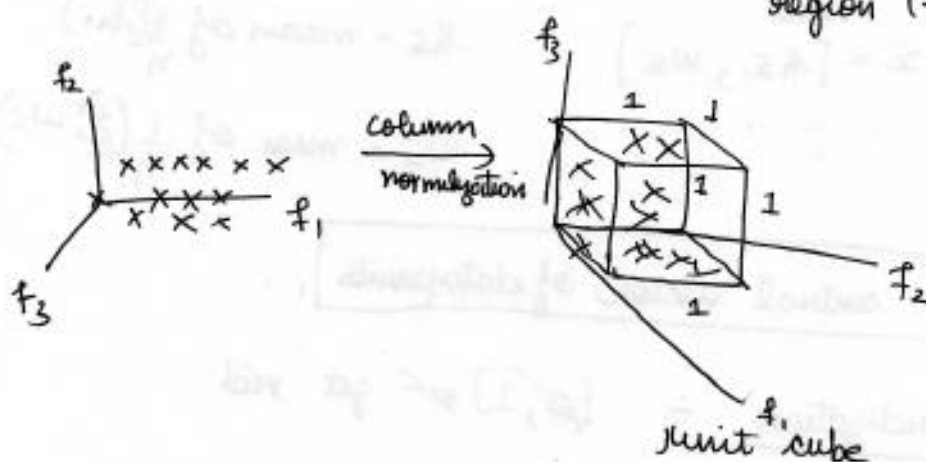


getting rid of scale.

Geometric Intuition



All points brought into region (unit square)



unit cube

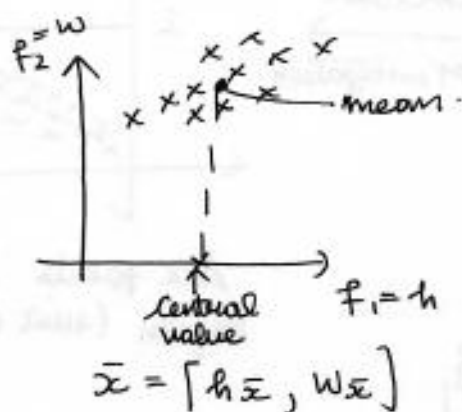
Note: Getting rid of scales

Data mean

$$X = \begin{matrix} & f_1 & f_2 & & f_d \\ \begin{matrix} 1 \\ 2 \\ 3 \\ \vdots \\ n \end{matrix} & \left[\begin{array}{c} \\ \\ \\ \\ \end{array} \right] & \left[\begin{array}{c} \\ \\ \\ \\ \end{array} \right] & \left[\begin{array}{c} \\ \\ \\ \\ \end{array} \right] & \left[\begin{array}{c} \\ \\ \\ \\ \end{array} \right] \\ & & x_i^T & & \end{matrix} \quad n \times d$$

mean vector $\rightarrow \bar{x} = \frac{1}{n} \sum_{i=1}^n x_i = \frac{1}{n} (x_1 + x_2 + \dots + x_n)$

Geometric



$$x_i \in \mathbb{R}^2$$

$$[h, w]$$

$$f_1 \quad f_2$$

$$h\bar{x} = \text{mean of } \frac{1}{n} \left(\sum_{i=1}^n h_i \right)$$

$$w\bar{x} = \text{mean of } \frac{1}{n} \left(\sum_{i=1}^n w_i \right)$$

mean vector central vector of datapoints

→ Column Standardization $\div [\phi, 1]$ get rid

Column normalization $\div [0, 1] \leftarrow$ get rid of scales of each features

→ Column standardization \div It is more often used

$$\text{Column standardization} \begin{cases} \text{mean} = 0 \\ \text{Standard dev} = 1 \end{cases}$$

$$\bar{a} = \text{mean} (a_i)_{i=1}^n \leftarrow \text{sample mean.}$$

$$s = \text{std. dev. } \{a_i\}_{i=1}^n \leftarrow \text{sample std. dev}$$

$$a_i' = \frac{a_i - \bar{a}}{s}$$

$$\text{mean } \{a_i\}_{i=1}^n = 0$$

$$\text{std. dev } \{a_i\}_{i=1}^n = 1$$

$$\text{Cov}(X, Y) = \frac{1}{n} \sum_{i=1}^n (x_i - \mu_x)(y_i - \mu_y)$$

$$S_{ij} = \text{cov}(f_i, f_j)$$

$$\text{cov}(f_i, f_j) = \text{var}(f_i)$$

$$\left\{ \begin{array}{l} \text{Cov}(X, Y) = \text{var}(X) \quad \text{--- (1)} \\ \text{cov}(f_i, f_j) = \text{cov}(f_j, f_i) \quad \text{--- (2)} \end{array} \right.$$

$$S = \begin{bmatrix} S_{11} & S_{12} & S_{13} \\ S_{21} & S_{22} & S_{23} \\ S_{31} & S_{32} & S_{33} \end{bmatrix}_{d \times d}$$

Symmetric matrix

$$A_{ij} = A_{ji} \quad \forall i, j$$

avg feature

$$\text{cov}(f_1, f_2) = \frac{1}{n} \sum_{i=1}^n (x_{i1} - \mu_1)(x_{i2} - \mu_2)$$

\uparrow mean(f_1) \uparrow mean(f_2)

$$\begin{array}{l} \mu_1 = 0 \\ \mu_2 = 0 \end{array}$$

* S has been column standardised

$$\text{cov}(f_1, f_2) = \frac{1}{n} \sum_{i=1}^n (x_{i1}) * (x_{i2})$$

$$\text{cov}(f_1, f_2) = (f_1^T f_2) * \frac{1}{n}$$

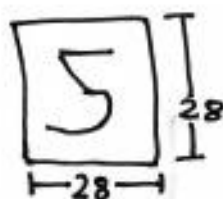
MNIST dataset

data visualization of high dimensions

Basic \rightarrow database of hand written digits

It is basic simple computer vision dataset.

It consists of 28×28 pixel images of handwritten digits (0-9)



28×28 pixels.

60K training datapoints

10K Test datapoints

$$D = \{x_i, y_i\}_{i=1}^{60K}$$

$$x_i : \begin{array}{|c|} \hline \square \\ \hline \end{array} \begin{array}{l} 28 \\ 28 \end{array}$$

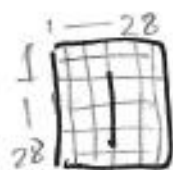
$$y_i \in \{0, 1, 2, 3, 4, 5, 6, 7, 8, 9\}$$

Objective: Given Image determine the number user has written.

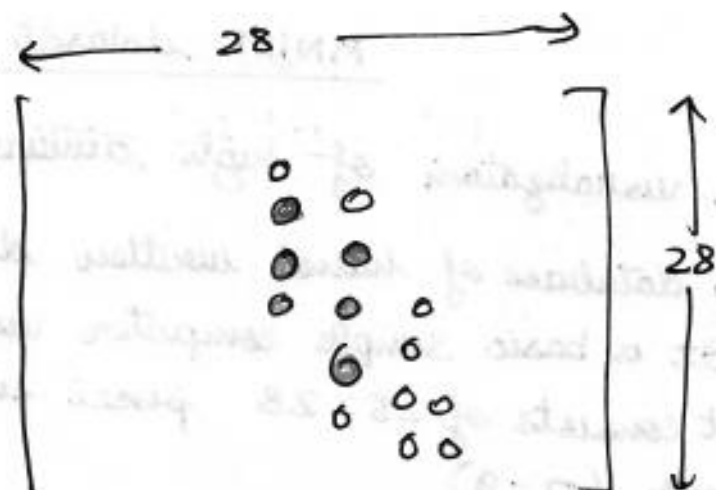
x_i is image of 28×28 so usually x_i are vectors (column)

$$x_i = \begin{bmatrix} \vdots \end{bmatrix} \quad x_i \in \mathbb{R}^{28}$$

for eg



=



where fully dark
put

0 denotes 1

0 → light gray
↓
fully dark

$x_i = \text{image}$



numerical / real matrix

NOTE:
matrix representation of
image

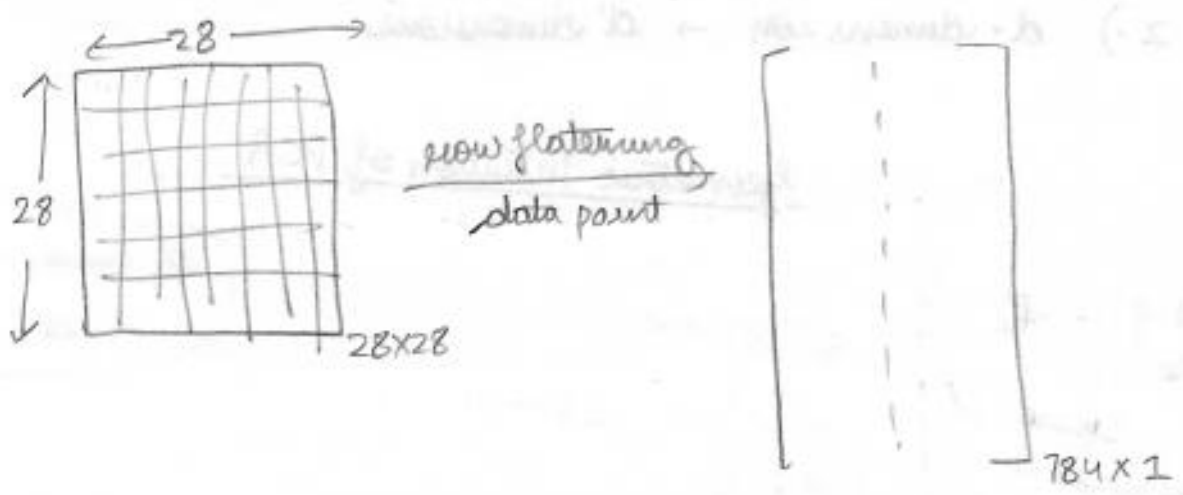
FLATENING

	1	2	3	4	5
1	1	2	4	8	8
2	3	2	1	8	2
3	2	1	5	8	4
4	3	2	1	8	2
5	4	2	6	8	1

FLATENING

1
2
4
6
8
3
2
1
8
2
1
8
1

25x1



$$X = \begin{bmatrix} f_1 & f_2 & f_3 & \dots & f_{784} \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ \vdots & \vdots & \vdots & \ddots & \vdots \end{bmatrix} \quad \begin{bmatrix} y_i \in \{0, \dots, 9\} \\ \vdots \\ \vdots \end{bmatrix}$$

← 784 →

$n \times d$ $n \times 1$

$n = 60K$
 $d = 784$

784-dimensional dataset

Principal Component Analysis (PCA)

Why? (dimensionality reduction)

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

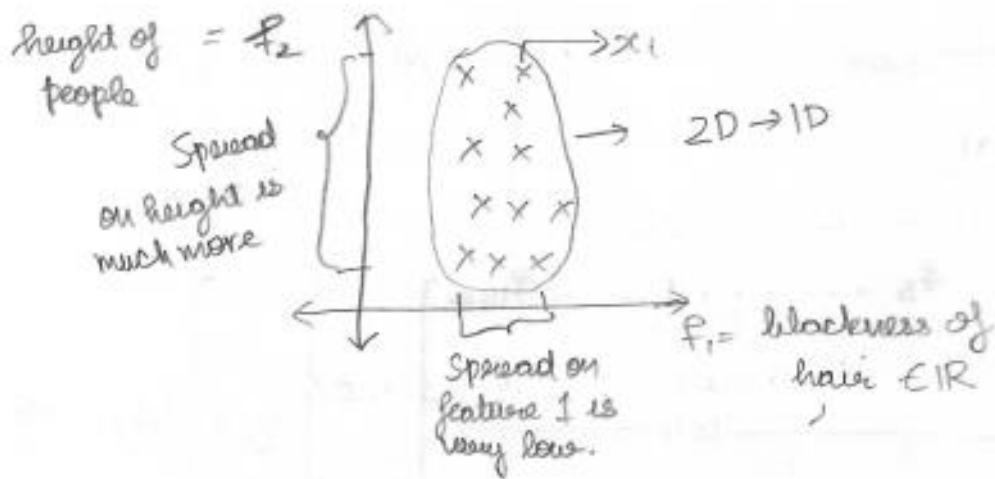
$$x_i \in \mathbb{R}^d \quad (d' < d)$$

MNIST \rightarrow 784 dimensions \rightarrow 2-dimension (visualize)

2.) d -dimension $\rightarrow d'$ dimensions

Geometric Intuition of PCA

d -dimension $\rightarrow d'$ dimensions
 $d' < d$



Indians \leftarrow

$$X = \begin{bmatrix} 1 & f_1 & f_2 \\ 2 & & \\ 3 & & \\ \vdots & & \\ n & & \end{bmatrix}$$

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

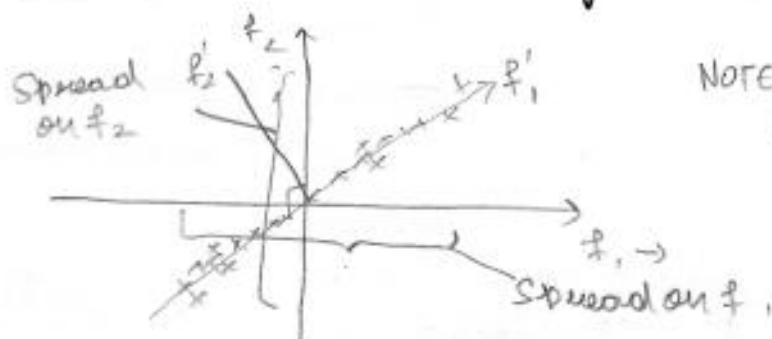
As spread on f_1 is low (low variance)
on the other hand f_2 variance is large

If forced to convert from $2D \rightarrow 1D$ we would consider f_2 as spread is more.

NOTE: preserving the direction with maximal spread (more information)
(spread is measure of information)

Example: 2 $X = 2$ dimensional dataset

column standardized $\left\{ \begin{array}{l} \text{mean } \{f_1\} = \text{mean } \{f_2\} = 0 \\ \text{var } \{f_1\} = \text{var } \{f_2\} = 1 \end{array} \right.$



NOTE: You can't drop f_1 or f_2 spread are large

2 Dimensional \rightarrow 1 Dimensional

f_1' has lot of spread

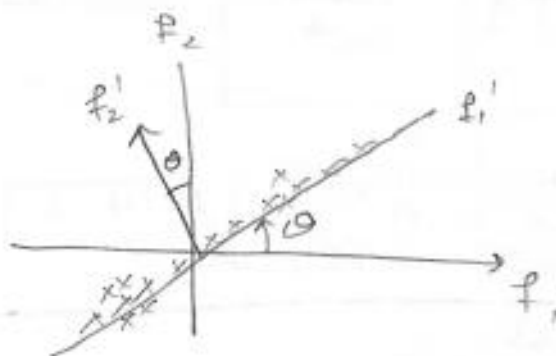
s.t f_1' has max spread.

① On the other hand $f_2' \perp f_1'$ has less spread

Spread on $f_2' \ll$ Spread on f_1'

② drop f_2'

③ project x_i 's on f_1' then 2D \rightarrow 1D



Objective

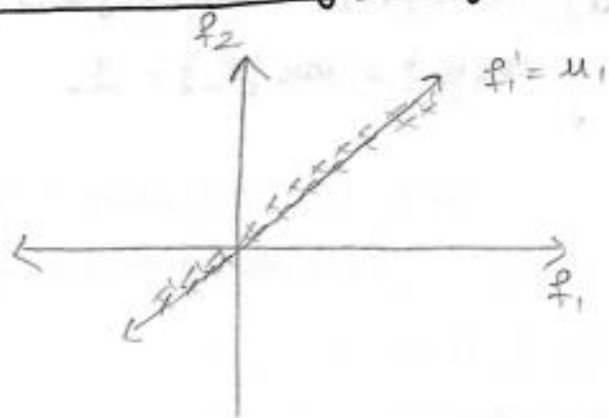
① we want to find a direction f_1' s.t. the variance of x_i 's projected onto f_1' is maximized.

① Rotating my axes to find f_1' with max-var.

② drop by f_2'

represent direction

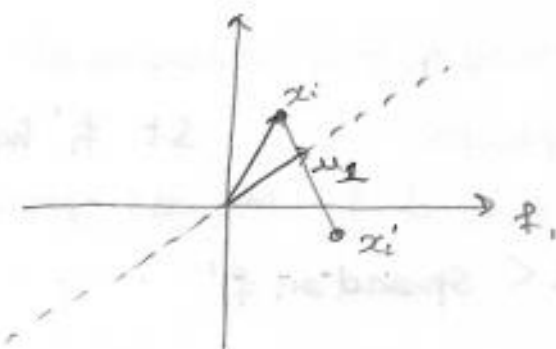
Mathematical objective function of PCA.



u_1 : unit vector

$$\|u_1\| = 1$$

Objective: To find direction by using u_1 , where spread is maximum



$$x_i' = \text{proj}_{u_1} x_i$$

$$\text{Dataset} = \{x_i\}_{i=1}^n$$

$$\text{Dataset}' = \{x_i'\}_{i=1}^n$$

$$x_i' = \frac{u_1 \cdot x_i}{\|u_1\|}$$

from linear algebra

$$PC_i = U_1^T x_i$$

$$\text{As } \|u_1\| = 1$$

If given

$$\overline{x_i'} = U_1^T \overline{x_i}$$

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

$$\text{mean } \{x_i\}_{i=1}^n$$

Important

Objective: u_1 s.t. $\text{var. } \{ \text{proj}_{u_1} x_i \}_{i=1}^n$ is maximum
(Spread)

$$\text{var} \{ u_1^T x_i \}_{i=1}^n = \frac{1}{n} \sum (\underbrace{u_1^T x_i}_{x'_i} - \underbrace{u_1^T \bar{x}}_{\text{mean} \{ x_i \}_{i=1}^n})^2$$

Remember

$$= (u_1^T)_{1 \times n} \underbrace{x_i}_{n \times 1}$$

1 scalar value

X : Column standardized {mean=0}.

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

$$\text{var} \{ x'_i \}_{i=1}^n = \frac{1}{n} \sum_{i=1}^n (u_1^T x_i)^2$$

variance $\{ x'_i \}$

NOTE

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

Objective
of optimization
problem

data matrix (given)

s.t $u_1^T u_1 = 1 = \|u_1\|^2$ Find u_1 to be maximized

constraint u_1 is unit vector as we want to represent direction

→ Eigen values & Eigen vectors

Objective → As we have very large dataset

(for eg Image of something, ~~set~~ weather forecast)

↓
we extract lot of data from it in numerical values. (lot of numerical values are 0 or sparse)

To reduce the dimensionality we predict eigen values and eigen vector

$$\boxed{A \vec{x} = \lambda \vec{x}}$$

eigenvalue

eigen vector

Singular Value Decomposition

$$A = U \Sigma V^T$$

orthogonal

diagonal

$$AA^T = U \Sigma V^T V \Sigma^T U^T$$

Solution to our optimization problems: λ, V ,

cov. $S_{d \times d}$

$$X = \begin{bmatrix} 1 & 2 & 3 & \dots & d \\ 2 & & & & \\ \vdots & & & & \\ n & & & & \end{bmatrix}_{n \times d}$$

$d \rightarrow$ features
 $n \rightarrow$ data points

Covariance matrix of $X = S$

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

\rightarrow symmetric matrix

NOTE:

$\lambda_n \rightarrow$ here represents features

Eigen values & Eigen vectors

$$\downarrow \quad \downarrow$$

$$(\lambda_1, \lambda_2, \lambda_3, \dots, \lambda_d) \quad (v_1, v_2, \dots, v_d)$$

$S_{d \times d}$

eigen values of $S = \lambda_1, \lambda_2, \lambda_3, \dots, \lambda_d$

eigen vector of $S = v_1, v_2, v_3, \dots, v_d$

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

maximal
eigen value

$d \times 1$ vector

def: $\lambda_1 v_1 = S_{d \times d} v_1$

$\downarrow \quad \downarrow$

Scalar $d \times 1$ vector

λ_1 : eigen value of S

v_1 : eigen vec. to S corr. to λ_1

Property $\Rightarrow \boxed{v_i \perp v_j}$

$$v_i^T v_j = 0$$

$$v_i \cdot v_j = 0$$

$\mu_1 = v_1$ = eigen-vector of $S (= X^T X)$
corresponding to largest eigen-value ($= \lambda_1$)

max-variance direction

Eigen values & vectors

Steps

stretching the value of X not changing

$$\text{Given } X = \begin{bmatrix} \dots \\ \dots \\ \dots \end{bmatrix}$$

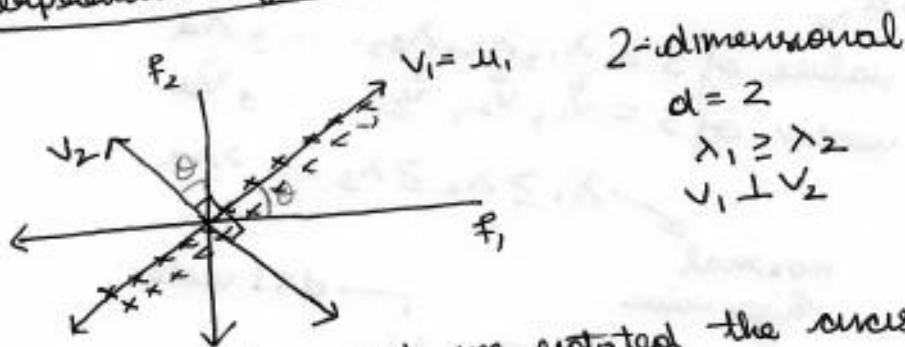
2.) Column Standardization done $\mu=0, \sigma=1$

$$3.) S_{\text{and}} = X^T X$$

4.) Compute values ~~eigen~~ eigenvalues & vectors
 $\lambda_1 \geq \lambda_2 \geq \dots \geq \lambda_d$
 v_1, v_2, \dots, v_d

5.) $\mu_1 = v_1$ (why)

→ Geometric Interpretation of λ_i & v_i



We just take the data and we rotated the axes with θ such that where variance is maximum.

Example: If we have 10 dimensional data
 $d=10$

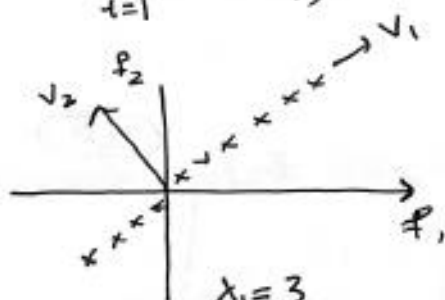
$$\lambda_1 \geq \lambda_2 \geq \lambda_3 \geq \dots \geq \lambda_{10}$$

$(v_1), v_2, \dots, v_{10}$
 \downarrow direction with max variance
 \downarrow 2nd most variance
 \downarrow least variance

λ_i significance

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

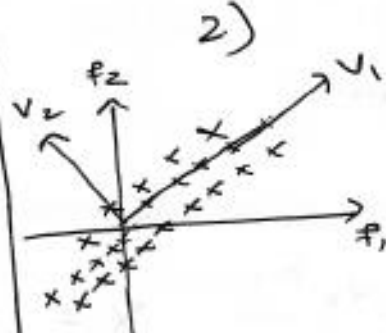
$$\sum_{i=1}^d \lambda_i$$



$$\lambda_1 = 3$$

$$\lambda_2 = 0$$

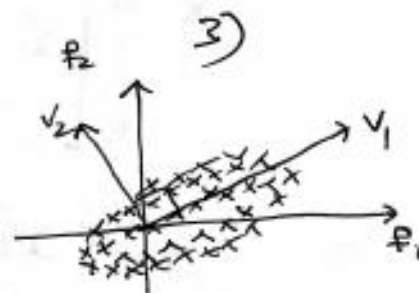
$$\frac{\lambda_1}{\lambda_1 + \lambda_2} = \frac{3}{3} = 1$$



$$\lambda_1 = 3$$

$$\lambda_2 = 1$$

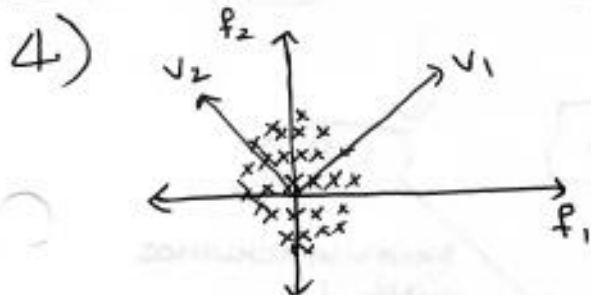
$$\frac{\lambda_1}{\lambda_1 + \lambda_2} = \frac{3}{4}$$



$$\lambda_1 = 3$$

$$\lambda_2 = 2$$

$$\frac{\lambda_1}{\lambda_1 + \lambda_2} = \frac{3}{5}$$



$$\lambda_1 = 3$$

$$\lambda_2 = 3$$

$$\frac{\lambda_1}{\lambda_1 + \lambda_2} = \frac{3}{6} = \frac{1}{2}$$

NOTE

$\frac{\lambda_i}{\lambda_1 + \lambda_2}$ %age of variance or information preserved

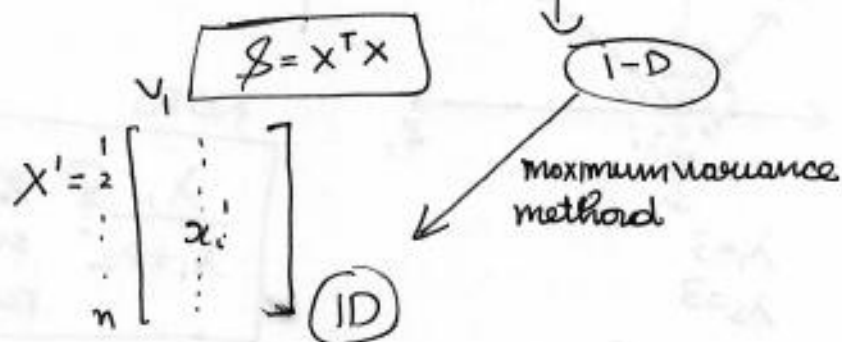
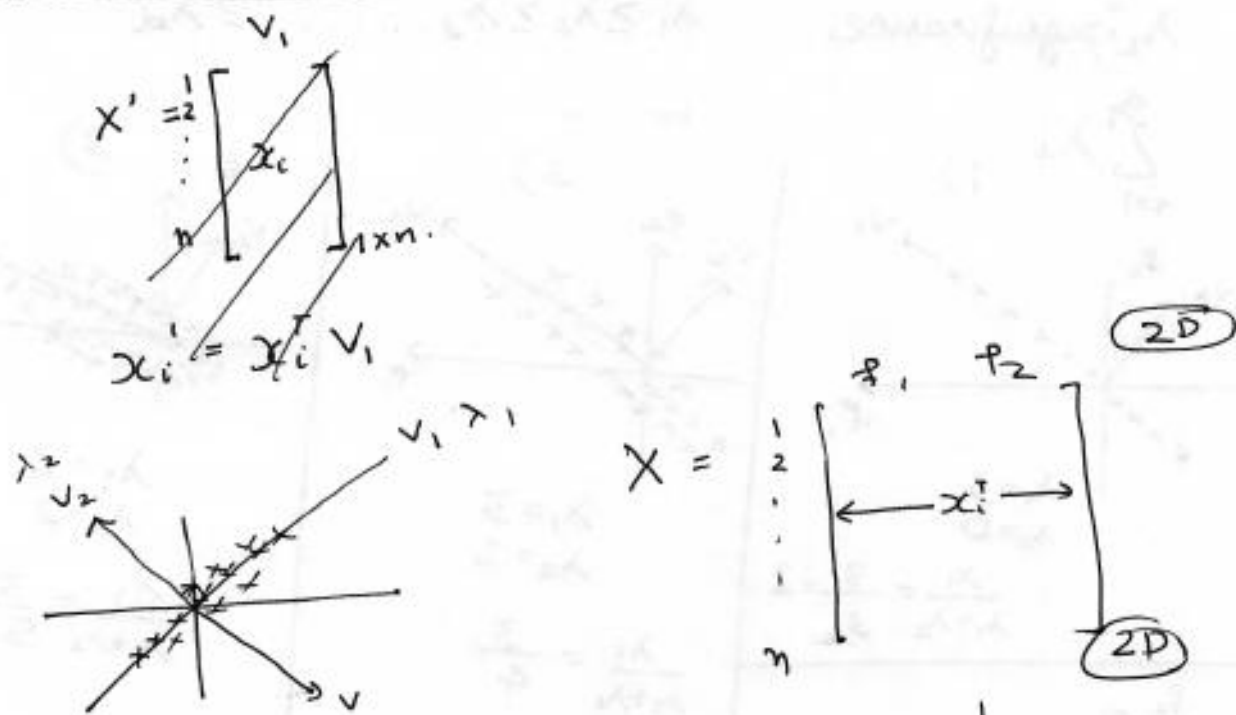
Project all data to v_1 →
we are going to preserve
100% of information for 1
75% of information for 2
60% _____ 3
50% _____ 4

$\frac{\lambda_i}{\sum \lambda_i}$ what % age of variance

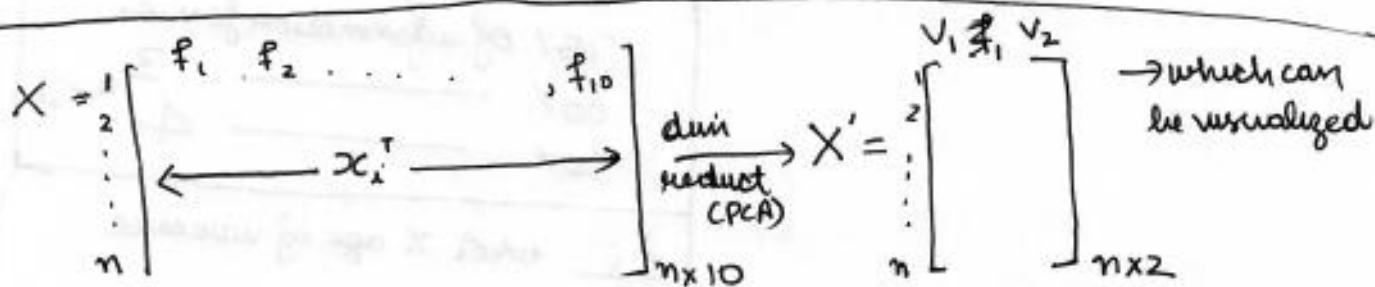
v_i tells direction where variance is maximum



PCA for dimensionality Reduction and Normalization



$x_i' = x_i^T v_1$ → v_1 choose as it has maximum variance.



$S = X^T X$
eigen value & eigen vector
 $\lambda_1 \geq \lambda_2 \geq \lambda_3 \dots \geq \lambda_{10}$
 $\downarrow \quad \downarrow \quad \downarrow \quad \downarrow$
 $v_1 \quad v_2 \quad v_3 \quad v_{10}$

$x_i' = [x_i^T v_1, x_i^T v_2]$
 $\downarrow \quad \downarrow$
 $v_1 \quad v_2$
Top 2 features

Suppose want to do 50 dimensions instead of 2

$$X = \begin{bmatrix} f_1 & f_2 & \dots & f_{100} \end{bmatrix}_{n \times 100} \xrightarrow[\text{convert 50 values}]{\text{PCA}} X' = \begin{bmatrix} v_1 & v_2 & \dots & v_{50} \end{bmatrix}_{n \times 50}$$

-10-
v₁ v₂ ... v₅₀

x_{ij}

$\lambda_1 \geq \lambda_2 \geq \lambda_3 \dots \geq \lambda_{50}$
v₁ v₂ v₃ ... v₅₀

$$x_i \in \mathbb{R}^{100} ; x_i' \in \mathbb{R}^{d'} \\ d' < 100$$

preserve 99% of the variance

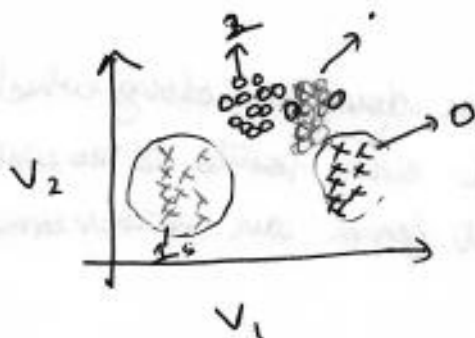
$$\text{let } \frac{\lambda_1 + \lambda_2 + \dots + \lambda_{51}}{\sum_{i=1}^{100} \lambda_i} = 0.99$$

99% variance preserved

$$d' = 51$$

✧ Visualize MNIST dataset

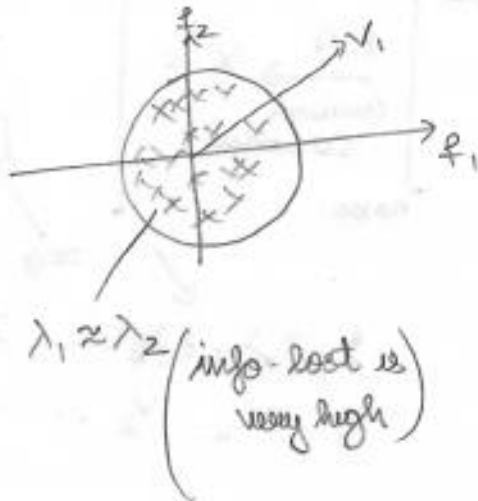
MNIST \rightarrow 784 $\xrightarrow[\text{dimension reduction}]{\text{PCA}}$ 2 (top 2 eigen vectors v_1 & v_2)



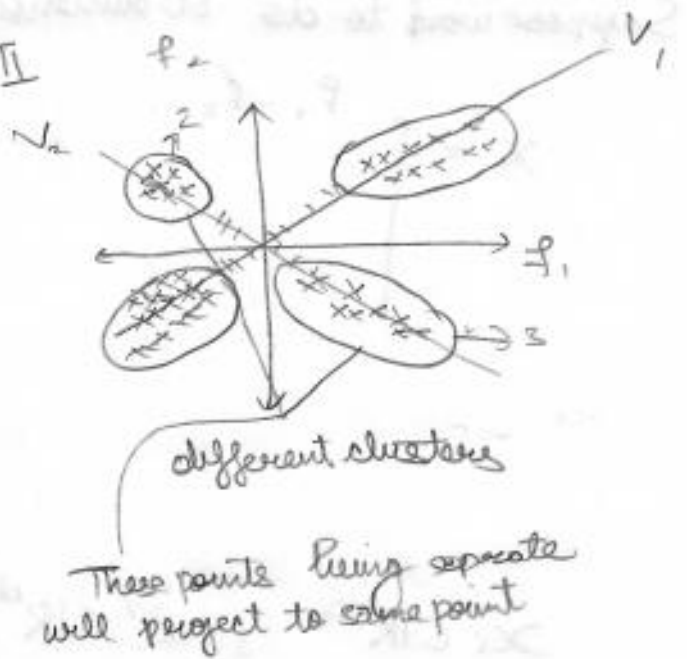
All 0 and 1's grouped together
2's are not well separated

★ Limitations of PCA

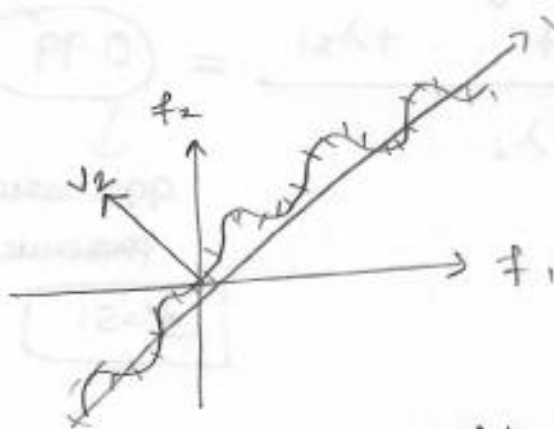
Case I



Case II



Case III



- visualizing in 1D all will be same
- but in 2D dimensional, we can be easily distinguished

Goal : PCA is to find projection directions along which the variance of the projected data points is maximum, and projection vectors must form an orthonormal basis

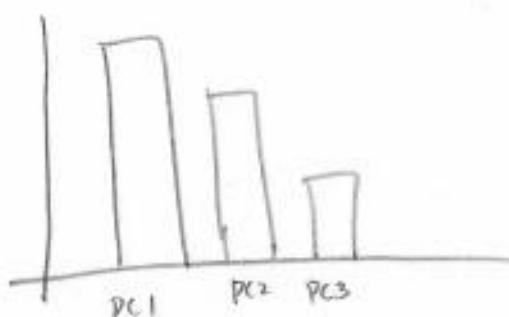
PCA: constrained optimization:

$$\begin{aligned}
 & \left\{ \begin{array}{l} \max_u u^T S u \\ \text{s.t. } u^T u = 1 \end{array} \right\} \rightarrow \mathcal{L}(u, \lambda) = u^T S u - \lambda(u^T u - 1) \\
 & \frac{\partial \mathcal{L}}{\partial u} = 0 \Rightarrow S u - \lambda u = 0 \\
 & u \text{ is eig. vec of } S \\
 & \lambda \text{ is eig. val of } S \\
 & S u = \lambda u \\
 & \lambda_1 \geq \lambda_2 \geq \lambda_3 \geq \lambda_4 \geq \dots \\
 & u_1, u_2, u_3, u_4, \dots, u_i \text{'s} \\
 & u_i^T u_i = 1 \quad | \text{ As unit vector} \\
 & \text{eig. vec.}
 \end{aligned}$$

$$\begin{aligned}
 & u^* = \max_u u^T S u \\
 & \text{s.t. } u^T u = 1 \\
 & \lambda_1 = \max_u \lambda \Rightarrow \text{From } \lambda
 \end{aligned}$$

$$\begin{aligned}
 & u^T S u \\
 & \Rightarrow \text{We know } S u = \lambda u \\
 & u^T (\lambda u) \\
 & \downarrow \text{scalar} \\
 & \lambda (u^T u) \\
 & \rightarrow \lambda * 1 = \lambda
 \end{aligned}$$

Plot for PCA visualization is called **scree plot**



Practical Tips for PCA:

- 1) Make sure that data is ~~scaled~~ normalized.
- 2) Make sure data is centered that is $\text{mean} = 0$
- 3) How many principal component you expect to find?
- 4) In summary, technically there is a PC for each variable in the data set. However, if there are fewer sample than variables, then the number of samples put an upper bound on the number of PCs with eigenvalues greater than 0.

(t-SNE) T-distributed Stochastic Neighbourhood Embedding

- State of the best dimensional reduction → visualization
- PCA → basic, old } → 2dim visualization not very good



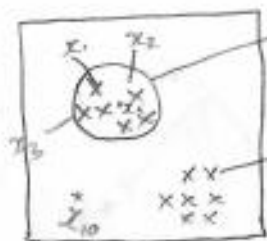
PCA: tries to preserve global structure of data.

t-SNE: preserves local-structure of data

Neighbourhood; Embedding

1) Neighbourhood

d-dim



Neighbourhood of point $N(x_1)$

MNIST: 784-dim

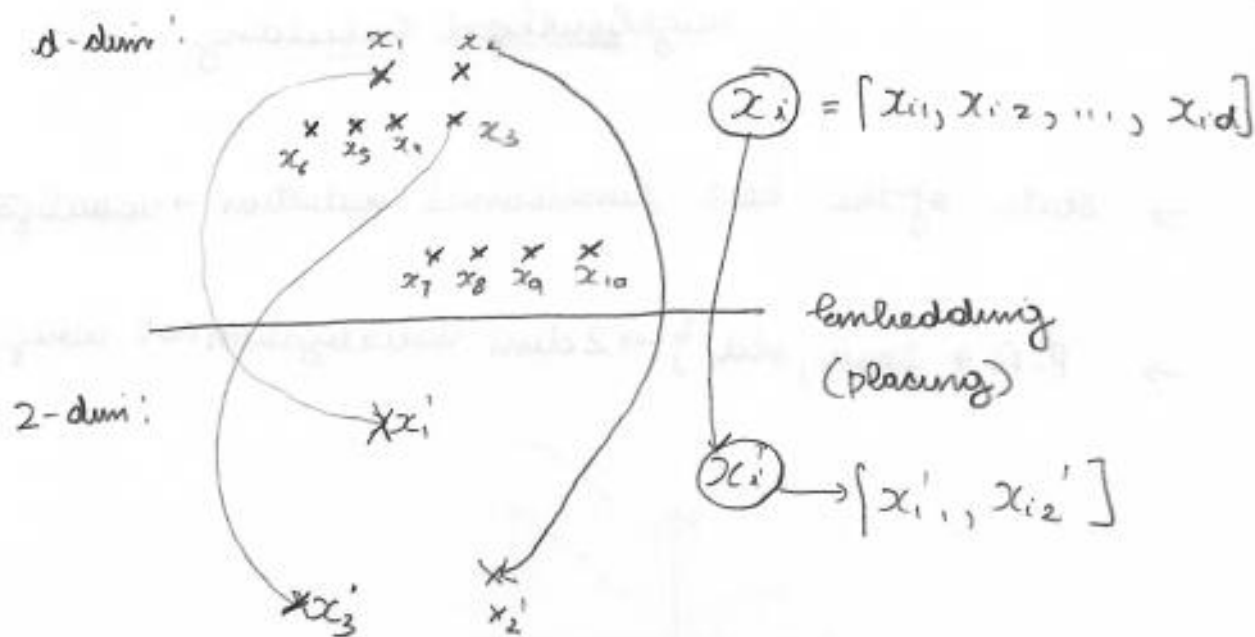
$$N(x_1) = \{x_i \mid x_i \text{ is geometrically close to } x_1\}$$

$$\|x_1 - x_2\|^2 = \text{dist}$$

$$N(x_1) = \{x_2, x_3, x_4\}$$

does Not contain x_{10} & x_{20}

Embedding

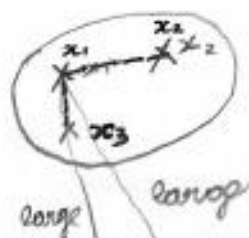


Every point in high dimensional space have corresponding in low dimension space such a thing is called embedding.

Geometric Intuition of t-SNE

d-dimension

$N(x_1) = \{x_2, x_3\}$



$N(x_4) = \{x_5\}$

\Rightarrow Embedding

2-dimension

$x_i' \in \mathbb{R}^2$



preserving distances of point in the neighborhood

$d(x_1, x_4) \neq d(x_1', x_4')$
But $d(x_4, x_5) \approx d(x_4', x_5')$

$$\frac{d(x_1, x_2) \approx d(x_1', x_2')}{d(x_1, x_3) \approx d(x_1', x_3')}$$

distance between the neighbourhood are preserved i.e

$$d(x_1, x_2) \approx d(x'_1, x'_2)$$

$$d(x_1, x_3) \approx d(x'_1, x'_3)$$

$$d(x_4, x_5) \approx d(x'_4, x'_5)$$

on the other hand

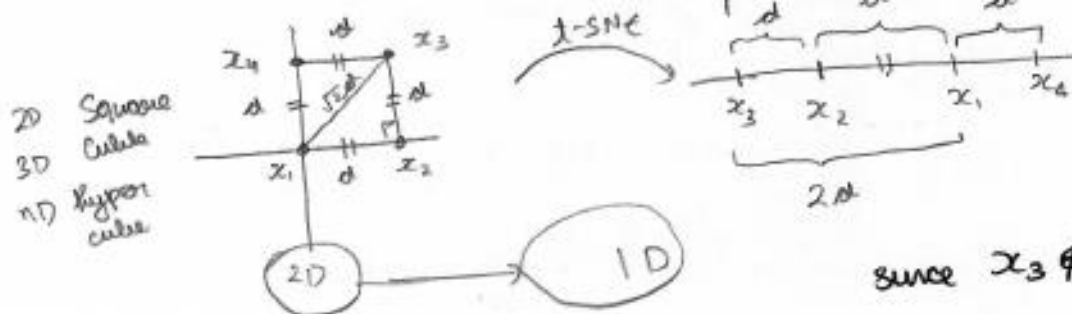
$$d(x_1, x_4) \neq d(x'_1, x'_4)$$

$$d(x_1, x_5) \neq d(x'_1, x'_5)$$

Crowding Problem

t-SNE: preserve dist in a N

$$N(x_3) = \{x_4, x_2\}$$



since $x_3 \notin N_d(x_1)$

some times it is impossible to preserve distance in all the N .

such problem is called crowding problem

It can be resolved using

t-distribution

Stochastic neighbourhood embedding

Student's t-distribution

distill. pub → go to the web link.

How to use t-SNE effectively

$(x_1, x_2) \rightarrow (x_1, x_2)$
 $(x_1, x_2) \rightarrow (x_1, x_2)$
 $(x_1, x_2) \rightarrow (x_1, x_2)$
 $(x_1, x_2) \rightarrow (x_1, x_2)$
 $(x_1, x_2) \rightarrow (x_1, x_2)$
 $(x_1, x_2) \rightarrow (x_1, x_2)$
 $(x_1, x_2) \rightarrow (x_1, x_2)$
 $(x_1, x_2) \rightarrow (x_1, x_2)$
 $(x_1, x_2) \rightarrow (x_1, x_2)$
 $(x_1, x_2) \rightarrow (x_1, x_2)$

distill. pub

$(x_1, x_2) \rightarrow (x_1, x_2)$

It is a lot more interesting



$(x_1, x_2) \rightarrow (x_1, x_2)$

It is a lot more interesting

It is

distill. pub

It is a lot more interesting



$(x_1, x_2) \rightarrow (x_1, x_2)$

Matrix Factorization : PCA & SVD

Matrix Factorization : $A = BCD$
 matrix product of other matrices

$A = BC$ → multiplicative decomposition

$6 \rightarrow 2 \times 3 \rightarrow (2, 3) \text{ factors of } 6$

Similarly → we say B, C, D are factors of A .

MF : → PCA (Principal component Analysis)
 ↳ Dimensionality reduction (high dim → low dim.)

PCA : → $X_{n \times d}$ Data Matrix (Centered)

$$S_{d \times d} = \frac{X^T X}{n-1} \Rightarrow \text{Covariance matrix}$$

eigen values : $\lambda_1, \lambda_2, \dots, \lambda_d$
 eigen vectors : w_1, w_2, \dots, w_d

$d \rightarrow d' \rightarrow$ basically take top λ 's → w 's
 ↓ eigen value ↓ eigen vectors

$$S_{d \times d} = \underbrace{W}_{\substack{\text{Covariance} \\ \text{matrix}}} \underbrace{\Lambda}_{d \times d} \underbrace{W^T}_{d \times d}$$

$w^T = \text{Transpose of } w$.

$$W = \begin{bmatrix} \uparrow & \uparrow & \uparrow & & \uparrow \\ w_1 & w_2 & w_3 & \dots & w_d \\ \downarrow & \downarrow & \downarrow & & \downarrow \end{bmatrix}_{d \times d}$$

eigen vector of \mathcal{S}

$$\lambda_{d \times d} = \begin{bmatrix} \lambda_1 & 0 & 0 & \dots & 0 \\ 0 & \lambda_2 & 0 & \dots & 0 \\ 0 & 0 & \lambda_3 & \dots & 0 \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & 0 & \dots & \lambda_d \end{bmatrix}_{d \times d}$$

These are eigenvalues of \mathcal{S} .

This whole decomposition is called eigen decomposition of \mathcal{S} which is purely matrix factorization

★ Singular value decomposition (Singular VD)

It is related to PCA

PCA: performed on co-variance matrix (S) \rightarrow square, symmetric.

SVD: performed on rectangular matrix.

$X_{n \times d}$

$$X_{n \times d} = U_{n \times n} \Sigma_{n \times d} V_{d \times d}^T \rightarrow \text{SVD}$$

for PCA $\text{Cov}(X) = S = W \lambda W^T \rightarrow \text{PCA}$

$$\Sigma = \begin{bmatrix} \lambda_1 & 0 & 0 & \dots & 0 \\ 0 & \lambda_2 & 0 & \dots & 0 \\ 0 & 0 & \lambda_3 & \dots & 0 \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & 0 & \dots & \lambda_d \end{bmatrix}_{n \times d}$$

Singular values of X \rightarrow diagonal matrix

$$\frac{s_i^2}{n-1} = \lambda_i$$

Relationship between eig-val of S and singular value.

$$X_{n \times d} = U_{n \times n} \Sigma_{n \times d} V_{d \times d}^T$$

\downarrow left singular vector of X \downarrow right singular vectors of X

u_i :- i^{th} column of $u =$ eig. vector of $(X X^T)$ not S

v_i :- i^{th} column of $v =$ eig. vect of $(X^T X) = S$

$\text{cov}(X) \quad S = W \Lambda W^T \quad \longrightarrow \text{PCA}$

\uparrow eig. vector of S \uparrow eig. values of S (w_i)
 \uparrow eig. vect of $(X^T X)$

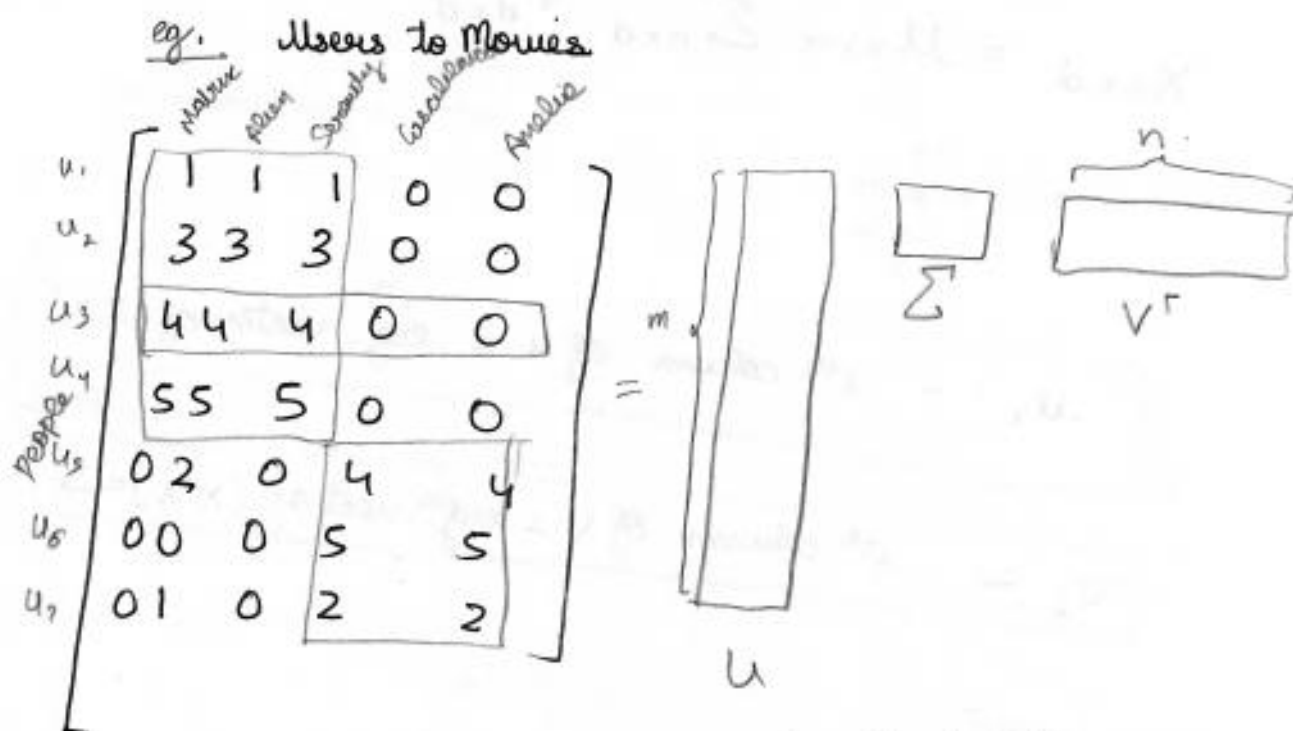
$X = U \Sigma V^T \longrightarrow \text{SVD}$
 \downarrow singular values
 $s_i ; \frac{s_i^2}{n-1} = \lambda_i$

or //

SVD → Stanford University Video

$$A_{[m \times n]} = U_{[m \times r]} \Sigma_{[r \times r]} (V_{[n \times r]})^T$$

o A: Input data matrix



movies and user both break into 2 groups.

if we supply SVD to the matrix

Self concept

← Romance concept

u_1	0.13	0.02	-0.01
u_2	0.41	0.07	-0.03
u_3	0.55	0.09	-0.04
u_4	0.68	0.11	-0.05
u_5	0.15	-0.59	0.65
u_6	0.07	-0.73	-0.67
u_7	0.07	-0.29	0.32

Strength of Self

Strength of Romance concept

low strength ignore it

Σ

V^T

12.4	0	0
0	9.5	0
0	0	1.3

Self concept

← movies to concept matrix

Romance concept

0.56	0.59	0.56	0.09	0.09
0.12	-0.02	0.12	-0.69	-0.69
0.40	-0.80	0.40	0.09	0.09

Note:

From above 'movies', 'users' & 'concepts'

-16-

U: user to concept similarity matrix

V: movies to concept similarity matrix

Σ : its diagonal element: 'strength of each concept'

Non-negative Matrix Factorization (NMF)

$$A_{n \times m} = B_{n \times d} (C^T)_{d \times m}$$

$$B: n \times d$$

$$C: m \times d$$

$$s.t. \quad B_{ij} \geq 0 \quad \forall i, j$$

$$C_{ij} \geq 0 \quad \forall i, j$$

$$\begin{bmatrix} & & & \\ & & & \\ & & & \\ & & & \end{bmatrix} = \begin{bmatrix} 0 & 0 & 0 \\ 0 & 0 & 0 \\ 0 & 0 & 0 \end{bmatrix} \begin{bmatrix} 0 & 0 \\ 0 & 0 \\ 0 & 0 \end{bmatrix}$$

A B C^T

Factors here are non negative

MF for Collaborative Filtering

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

u_i r_{ij}

r_{ij} - rating on I_j by U_i :
Sparse matrix: many empty cells.

$$A = B * C^T$$

\swarrow \searrow
 users items

$$B: n \times d$$

$$d \geq 0$$

$$d \leq \min(m, n)$$

$$d \leq m, d \leq n$$

$$C^T: d \times m$$

A can be decomposed as a product of 2 matrices B & C.

Q lot of empty values in matrix

$$A_{ij} = B_i \cdot C_j$$

\nwarrow jth row of matrix C.

\downarrow
ith row of matrix B

$$A = \begin{bmatrix} u_1 \\ u_2 \\ u_3 \\ \vdots \\ u_n \end{bmatrix} \begin{bmatrix} \\ \\ \\ \\ \end{bmatrix}_{n \times m}$$

I_j

$$B = \begin{matrix} & \text{user} \\ \begin{bmatrix} u_1 \\ u_2 \\ u_3 \\ \vdots \\ u_n \end{bmatrix} & \begin{bmatrix} \\ \\ \\ \\ \end{bmatrix} \end{matrix}_{n \times d}$$

$$C = \begin{matrix} \begin{bmatrix} I_1 \\ I_2 \\ \vdots \\ I_m \end{bmatrix} & \begin{bmatrix} \\ \\ \\ \end{bmatrix} \end{matrix}_{m \times d}$$

\downarrow
item

$$A_{ij} = B_i^T * C_j$$

\downarrow
(many empty cells)

$A = BC^T$

\downarrow
optimization problem
 \downarrow
SGD

Find B & $C \Rightarrow MF$

argmin $\sum_{i,j} (A_{ij} - B_i^T C_j)^2$ can't use empty $A_{i,j}$

$\{i,j \text{ s.t. } A_{ij} \text{ is not empty}\}$

B, C Full matrices

minimizing squared loss

$$B^T = \begin{bmatrix} \uparrow B_1 \downarrow & \uparrow B_2 \downarrow & \uparrow B_3 \downarrow & \dots & \uparrow B_n \downarrow \end{bmatrix}_{d \times n}$$

$$B = \begin{bmatrix} B_1 \\ B_2 \\ \vdots \\ B_n \end{bmatrix}_{n \times d}$$

$n = \text{number of users}$

$$C = \begin{bmatrix} c_1 \\ c_2 \\ \vdots \\ c_m \end{bmatrix}_{m \times d}$$

$m = \text{number of items}$

① solve the optimization problem using SGD

nonempty $A_{i,j}$'s to find B, C s.t. $A_{i,j} \neq 0$

argmin $\sum_{i,j \text{ valid}} (A_{ij} - B_i^T C_j)^2$

② $B = \begin{bmatrix} \leftarrow B_i \rightarrow \end{bmatrix}$; $C = \begin{bmatrix} \leftarrow c_j \rightarrow \end{bmatrix}$

③ matrix completion

$$A = \begin{bmatrix} \square & \square & \square & \square \\ 0 & \square & \square & \square \\ \vdots & \vdots & \vdots & \vdots \end{bmatrix}$$

empty $(A_{10,5}) = \text{empty}$

$$A_{10,5} = B_{10}^T * C_5$$

Fill my empty cells in A using $B_{10}^T * C_5$.

↓
it is ~~to~~ not using $A_{10,5}$
as it was zero it is prediction
based on some other user
whose output $\neq 0$

$$A = \begin{bmatrix} \quad \end{bmatrix} \rightarrow \text{mostly empty}$$

no
empty
cells

$$\hat{A} = B * C^T$$

computed using A_{ij} that are not empty
by solving optimization problem

$$\begin{bmatrix} \quad \end{bmatrix}$$

$$A - \hat{A} \rightarrow \text{v. small}$$

↓
it is fully filled / completed form of A using
matrix factorization

$$= \begin{bmatrix} 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 \end{bmatrix} - u_i \begin{bmatrix} \quad \quad \quad \quad \\ \quad \quad \quad \quad \\ \quad \quad \quad \quad \\ \quad \quad \quad \quad \end{bmatrix}$$

A A_1

↓
 u_i & $\Sigma_1 \rightarrow$ get a rating
↓
predicted rating

Matrix Factorization for feature engineering

$$A = u_i \begin{bmatrix} I_r \\ \square \end{bmatrix} \quad \left\{ \begin{array}{l} \mu_i \rightarrow \text{[box]} \rightarrow d-d_u \\ I_r \rightarrow \text{[box]} \rightarrow d-d_u \end{array} \right.$$

→ using the state in A_{ij}

$$A_{n \times m} = B C^T$$

$\swarrow \quad \searrow$
 $n \times d \quad d \times m$

$$d \geq 0, d \leq \min(m, n)$$

$$B = \begin{bmatrix} 1 & 2 & 3 & \dots & d \\ 2 & & & & \\ \vdots & & & & \\ n & & & & \end{bmatrix} \quad \begin{array}{c} \text{[row of boxes]} \rightarrow u_i \\ n \times d \end{array}$$

\swarrow
 number of users

$$C = \begin{bmatrix} 1 & 2 & \dots & d \\ 2 & & & \\ \vdots & & & \\ m & & & \end{bmatrix} \quad m \times d$$

\swarrow
 number of items

$$B_i = i^{\text{th}} \text{ row of } B$$

$$= u_i \in \mathbb{R}^d$$

$$C_i = i^{\text{th}} \text{ row of } C$$

$$= I_i \in \mathbb{R}^d$$

$I_i, I_j \rightarrow v$ similar based on rating data

$\text{dist}(I_i, I_j) \rightarrow \text{small}$

MF for feature