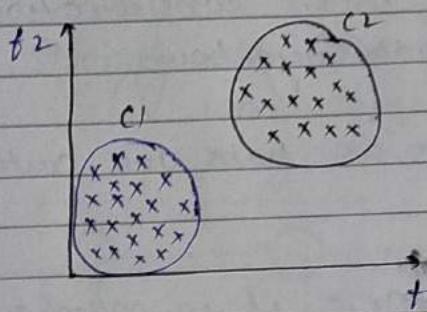


UNSUPERVISED LEARNING

(*) Clustering :-

$D = \{x\}$: no y_i 's
Task: Group/cluster "Similar" data points.



: All points are clustered in 2 diff groups.

: The process of grouping similar points into diff groups/clusters is known as clustering.

clustering :- K-Means, Hierarchical clustering, DBSCAN.

(*) Applications of clustering :-

eg: ① e-commerce - Amazon, Alibaba, eBay, Flipkart.

task:- ↗ Group "Similar" customers based on their purchasing behavior → money
↪ credit
↪ products

Youtube: Programming Cradle

② Image Segmentation (Computer Vision and image processing)
↪ grouping / clustering similar pixels

③ Amazon food reviews :- Manually label all the reviews as +ve or -ve to create a dataset! (very time consuming and expensive)

⇒ $D = 1M$ reviews no y_i 's

↪① cluster them in 10k groups of similar words (Bow, tf-idf, w2v)

↪② review & label each cluster.

↪
C_i → pick one point randomly
+ve
-ve
if +ve
if -ve

✳ Metrics for clustering :-

Geometrically :- What is a good clustering result?

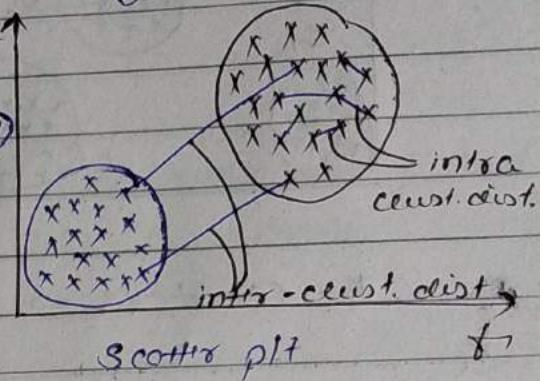
Given :- $D = \{x_i\} \quad x_i \in \mathbb{R}^2$

→ Grouped in clusters such that intra-^{is small} dist.

intra-cluster :- with in a cluster

inter-cluster :- across or b/w
clusters

(large)

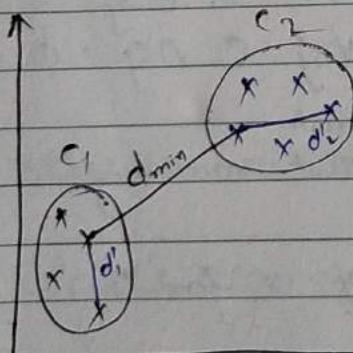


✳ Density :- inter cluster dist. → v. high
intra cluster dist. → v. low.

✳ Dunn Index :-

$$D = \frac{\min_{i,j} d(i,j)}{\max_k d'(k)}$$

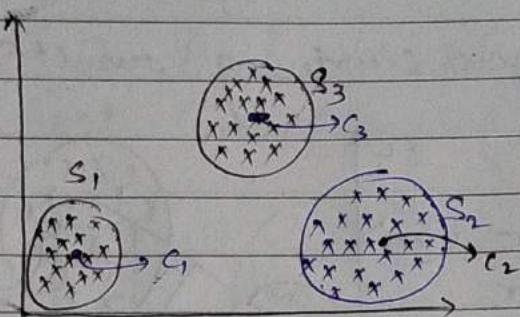
dist. b/w c_i, c_j
 K clusters
 $= \{c_1, c_2, c_3, \dots, c_n\}$



Youtube: Programming Cradle

- for every pt in c_1 , find dist. to all pts in c_2 : $\{d_1, d_2, d_3, \dots, d_n\}$
take min among others
- find all intra cluster dists and take max
 $\max\{d'_1, d'_2\}$

④ K-Means: Geometric Intuition, Centroids

w/ $K=3$

↳ no. of clusters (hyperparam)

$$S_1 \cap S_2 = \emptyset$$

 $c_1, c_2, c_3 \rightarrow$ Centroids

$$S_1 \cap S_3 = \emptyset$$

$$S_2 \cup S_3 = \Omega$$

$$S_2 \cap S_3 = \emptyset$$

K-clusters \Rightarrow K-centroids and K-sets.

$$c_i = \frac{1}{n} \sum_{x_j \in S_i} x_j \leftarrow \text{mean point to } S_i$$

K-means:- Centroid-based clustering.

⑤ K-Means:- Mathematical formulation of Objectivefn.

$$D = \{x_1, x_2, \dots, x_n\}$$

Youtube: Programming Cradle

task:- K-centroids :- $c_1, c_2, c_3, \dots, c_K$ $\left\{ \forall i, x_i \in S_j \right\}$
 sets:- $S_1, S_2, S_3, \dots, S_K$ $\left\{ \forall i, j, S_i \cap S_j = \emptyset \right\}$

$$\Rightarrow \underset{\{c_1, c_2, \dots, c_K\}}{\text{augmin}} \sum_{i=1}^K \sum_{x \in S_i} \|x - c_i\|^2 \leftarrow \text{intra dist. minimiz.} \quad \text{cost.}$$

such that $x \in S_i$ } $\left\{ \begin{array}{l} \text{SD. dist. of } x \text{ for } c_i \\ \text{S.t. } S_i \neq \emptyset \end{array} \right\} \rightarrow \text{constraints.}$

Very hard to solve as its NP-hard problem
 it takes exponential time complexity.

↓
 Solved using approximation (Lloyd's algo.)

* K-Means Algorithm :- (Lloyd's algo)

① Initialization :-

↳ randomly pick k pts from D and call them c_1, c_2, c_3, \dots

② Assignment :-

for each point in D select the nearest c_j and
find $\text{dist.}(x_i, c_j)$ $\forall j = 1, 2, \dots, k$ and
add x_i to set S_j

③ Recompute Centroid :-

↳ calculate / update c_j 's as follows

$$c_j = \frac{1}{|S_j|} \sum_{x_i \in S_j} x_i$$

no. of pts in set S_j mean.pt

④ Repeat step ② and ③ until convergence.

↳ when centroids doesn't change much.

* How to initialize : K-means++

Lloyd's algo :- initialization step is done using randomly picking points. BUT there is a problem with it.

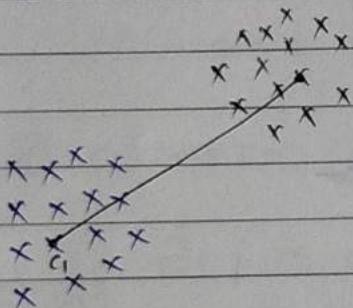
Problem :- initialization sensitivity.

final cluster is dependent on how points are picked.

i) How to deal with initialization sensitivity problem?
 Ans: Repeat K-means multiple times with diff initialization
 ↳ pick the best clustering based on:
 smaller intra clust. dist. and
 larger inter clust. dist.

ii) K-means++ :- Replace random initialization with
 smart init. { c_1, c_2, \dots, c_k }

probabilistic approach



x_i	d_i	$\ x_i - c_j\ ^2$
x_1	d_1	
x_2	d_2	
x_3	d_3	
x_4	d_4	
:	:	
x_n	d_n	

Youtube: Programming Cradle

- ① Pick the 1st centroid randomly $\rightarrow c_1$ for i^*
- ② $\forall x_i \in D$ create a distribution as below.
 $x_i \rightarrow \text{dist}^2(x_i, \text{nearest centroid})$
- ③ pick opt. from $D - \{c_i\}$ with a prob. proportional to d_i

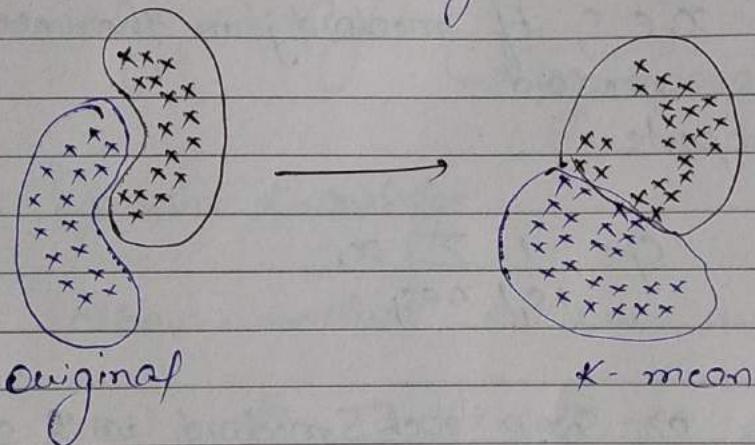
Q) Why select centroid probabilistically? Why not deterministically?

- A) Deterministic approach :- Select point with max distance from already selected nearest centroid.
- ↳ this approach is heavily affected by the outliers.

Probabilistic approach :- less affected by outliers.

* Failure cases / Limitations :-

- ① Different size of clusters :- K-means tries to make clusters of some size, and hence if the original clusters are of diff size K-means may not work well.
- ② Different density of clusters.
- ③ Clusters with Non globular shape. (non convex)



Youtube: Programming Cradle

④ Outliers

NOTE :- To solve the above limitations we can ↑ 'k'
but it doesn't work very well in all the cases.

* K-Medoids :-

problem with K-mean :- it is not interpretable
↓
clusters

one of the date points.
So, instead of giving me c_j 's computed using means,
if $x_0 \in \mathcal{N}$ as a first centroid was given it is more
interpretable

Big idea :- what if each centroid is a datapoint in \mathcal{D} ?
 ↓
 it is called as k -medoids

* Portioning around medoids (PAM) :- algorithm to implement k -medoids.

Step 1:- Initialization :- k -Means++ \rightarrow probabilistic method
 pick K pts from \mathcal{D}

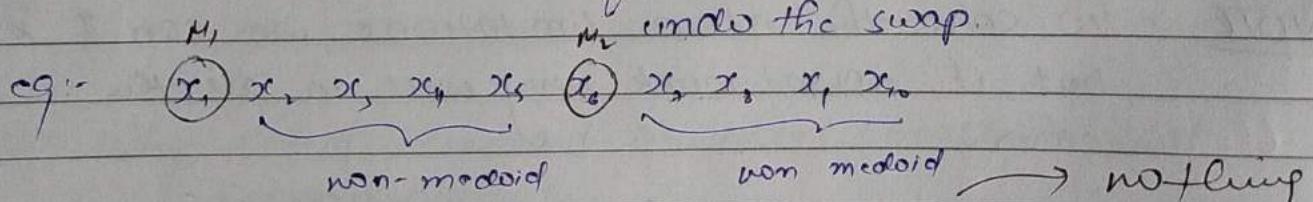
Step 2:- Assignment :- closest medoid
 $x_i \in S_j$ if medoid j is the closest medoid

Step 3:- Update / Recompute :-

$$k\text{-means: } c_j = \frac{1}{|S_j|} \sum_{x_i \in S_j} x_i$$

k -medoids: a) Swap each medoid with a non-medoid pt.
 (PAM)

b) if loss \downarrow swap else undo the swap.



$$\therefore \text{loss in } k\text{-means: } \sum_{j=1}^k \sum_{x_i \in S_j} \|x_i - m_j\|^2 \text{ but distance}$$

$$\text{eg: } (x_1, x_2, x_3, x_4, x_5, x_6, x_7, x_8, x_9, x_{10}) \xrightarrow{\text{medoid}} (x_1, x_2, x_3, x_4, x_5, x_6, x_7, x_8, x_9, x_{10})$$

a) loss value $x_1 = m_1, x_6 = m_2 \rightarrow l_1$

b) swap $M_1 = x_2, M_2 = x_6$: x_1 as non medoid pt.

$$\hookrightarrow \text{loss value } \left. \begin{array}{l} M_1 = x_2 \\ M_2 = x_6 \end{array} \right\} \rightarrow l_2$$

if $b_2 < b_1$,

$$m_1 = x_2$$

$$m_2 = x_0$$

else

$$m_1 = x_1$$

$$m_2 = x_0$$

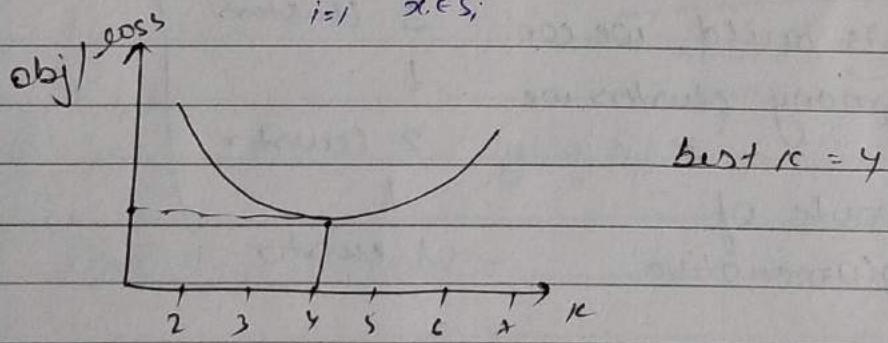
NOTE :- If similarity or distance matrix is given
 k -medoid can be implemented easily.

* Determining the Right k :-
 ↗ hyperparameter.

① Domain Knowledge Youtube: Programming Cradle

② elbow - method or sum method.

$$\text{cost} : \sum_{i=1}^k \sum_{x \in S_i} \|x - c_i\|^2 \rightarrow \text{minimize}$$



* Time Complexity :-

$$K\text{means} = O(nKd) \approx O(nd)$$

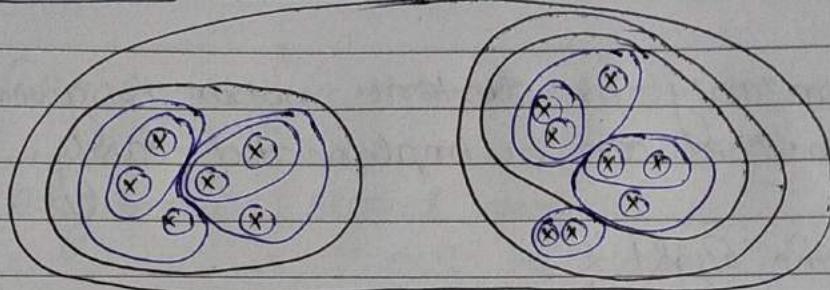
↓ ↓ ↓
 no. of pts iteration dim
 no. of centers

$$\text{Space Complexity} :: O(nd + kd) = O(nd)$$

* Hierarchical clustering :-

Two types :- Agglomerative and Divisive
more popular.

•> Agglomerative



Youtube: Programming Cradle

Initially each point is a cluster : 14 clusters

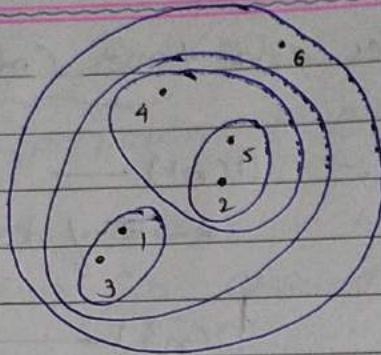
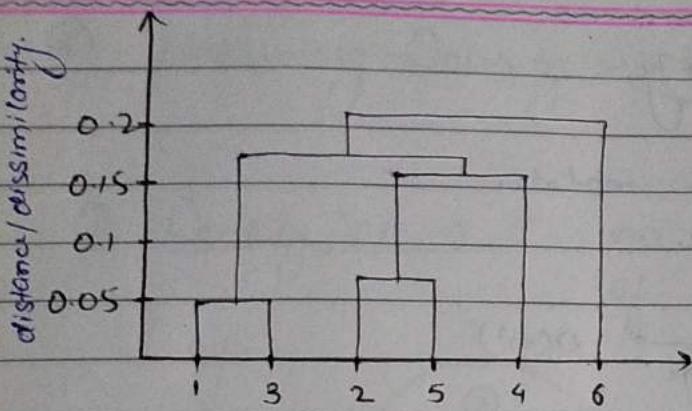
Now it groups close-by clusters : 9 clusters

NOTE :- NO need to specify no. of clusters (K). Once the hierarchy is build, we can select how many clusters we want.

•> Divisive :- Opposite of Agglomerative

Initially whole dataset / all points are part of 1 big cluster. Then in each step points will be divided into various clusters.

1 cluster \rightarrow 2 clusters \rightarrow 3 clusters \rightarrow 5 clusters \rightarrow 9 clusters
 \downarrow
14 clusters.



→ Dendrogram :- A tree like diagram that records the sequences of merges or splits.

• How to define intra-cluster similarity?

$$\left\{ \begin{array}{l}
 \text{Min} := \text{sim}(c_1, c_2) = \min_{p_i \in c_1, p_j \in c_2} \text{sim}(p_i, p_j) \text{ such that} \\
 \text{compute link} \\
 \text{Max} := \text{sim}(c_1, c_2) = \max_{p_i \in c_1, p_j \in c_2} \text{sim}(p_i, p_j) \text{ such that} \\
 \text{Group avg} := \text{sim}(c_1, c_2) = \frac{\sum_{p_i \in c_1, p_j \in c_2} \text{sim}(p_i, p_j)}{\text{size of } c_1 * \text{size of } c_2}
 \end{array} \right.$$

⇒ Dist. B/w Centroid :-

• Limitations :-

Min :- Sensitive to outliers/noise.

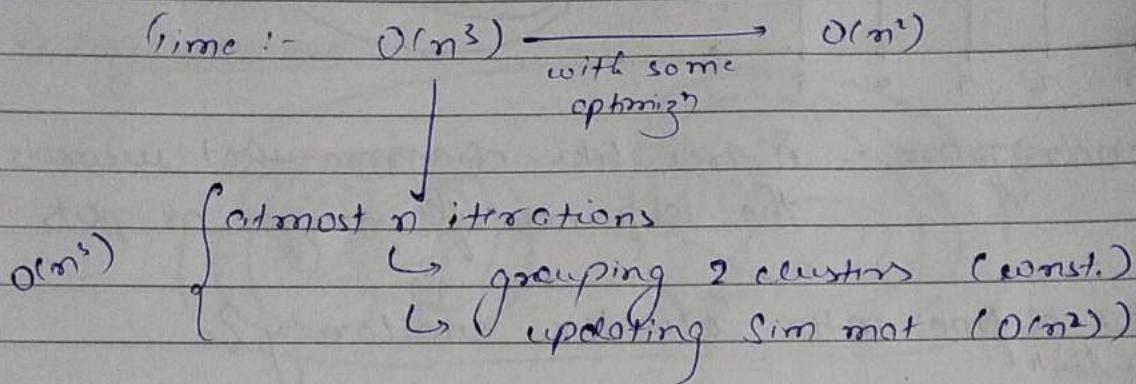
Max :- • Biased towards globular clusters

• Tends to break large clusters into cliff parts.

Group avg :- Biased towards globular clusters.

* Space and time Complexity:-

Space :- $O(n^2) \rightarrow$ Sim. matrix.
 $n \rightarrow$ # of datapts.



* Limitations of Hierarchical clustering.

- ① No mathematical objective function that we are directly solving. Hence we can not link it with other concepts.
- ② Difficult to select one from Min, Max, avg, clust.
- ③ Space and time complexity. (can't be used for large dataset)

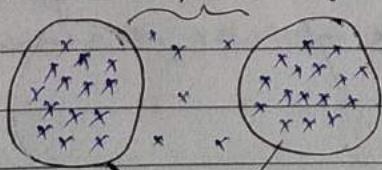
* DBSCAN :- Density based clustering

↳ Density-based spatial clustering of applications with noise (DBSCAN)

Dense-region \rightarrow clusters (valid data points)

Sparse-region \rightarrow noise (outliers)

Sparse region (noise)

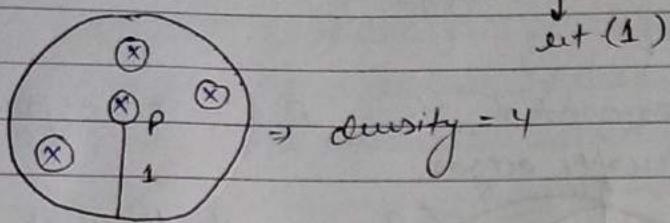


Dense-region (clusters)

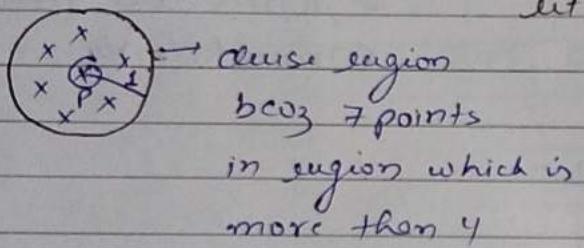
* Measuring Density :- $\underbrace{\text{Min Pts.}}_{\text{hyperparam of DBSCAN}} \cdot \text{Eps} (\epsilon)$

① Density at a P : no of pts within a hypersphere of radius Eps around p.

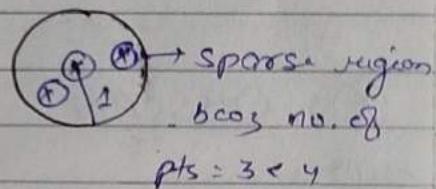
eg:-



② Dense region :- a hypersphere/circle of radius Eps that contains atleast $\underbrace{\text{Min Pts.}}_{\text{at } 4}$ points



Youtube: Programming Cradle



* Core point; border pt and noise pt.

① Core point :- if p has $\geq \text{Min Pts}$ points in an Eps radius around it then p is core pt.

② Border pt :- i) p is not a core pt \Rightarrow p has $< \text{Min Pts}$ pts in Eps radius
ii) $p \in \text{Neighbourhood}(\delta)$ where δ is core pt
 \downarrow
 $\text{dist}(p, \delta) \leq \text{eps}$

③ Noise pts :- point which is neither core pt nor Border pt.

* Density edge & density connected pts.

① Density edge :- $p, q \rightarrow$ core point and

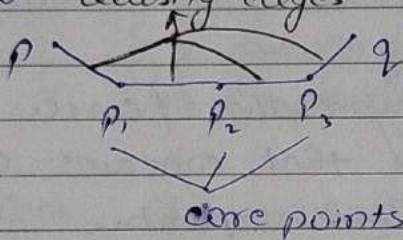
$$d(p, q) \leq \epsilon_{ps}$$

\curvearrowleft Density edge connecting these pts is density edge.

p q

② Density connected pts:- $p, q \rightarrow$ core point.

density edges



if there is a path b/w two core points and intermediate points are also core points then 2 points are density connected pts.

* DBSCAN Algorithm - Youtube: Programming Cradle

Steps

① $\forall x_i \in D \rightarrow$ Label them as core pt, border pt, noise pt.

set of pts. labeled as core, border, noise

e.g. $x_1 \rightarrow$ core pt $x_4 \rightarrow$ border pt.

$x_2 \rightarrow$ noise pt $x_5 \rightarrow$ core pt.

$x_3 \rightarrow$ core pt $x_6 \rightarrow$ border pt.

$S_i = \text{rangeQuery}(x_i, D, \epsilon_{ps})$

SKLearn uses kd-tree to implement range query in DBSCAN

② Remove noise pts.

③ For each core pt. 'p' not assigned to a cluster

a) create a new cluster in the p

b) add all pts that are density connected to p into this new cluster.

④ Each border pt. \rightarrow assign it to the nearest core pts cluster.

* Hyperparam :- MinPts, Eps

① Min Pts :-

a) Rule of thumb :- $\text{Min Pts} > d + 1$

↳ dimensionality

Typically $\text{min Pts} \approx 2 * d$

b) If dataset is more noisy :- choose larger value of Min Pts to remove noise.

c) Min Pts often chosen by domain expert.

② Eps : cut $\text{Min Pts} = 4$

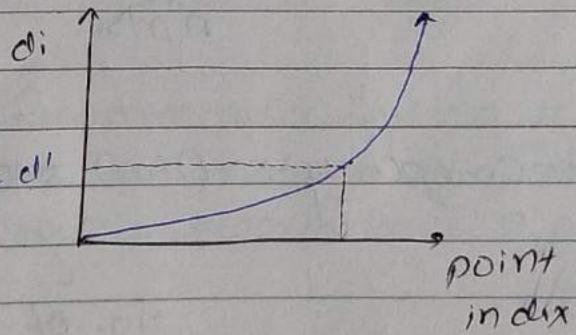
a) $x_i \rightarrow d_i$: $d_i \rightarrow$ dist from x_i to the 4th NN of x_i

b) $x_1 \quad d_1$
 $x_2 \quad d_2$
 $x_3 \quad d_3$
 $\vdots \quad d_y$
 $x_n \quad d_n$

sort d_i 's in increasing order.

Youtube: Programming Cradle

c) plot point index vs d_i and use elbow method.



if $d_i \rightarrow$ high:
 chances that $Eps = d_i$
 x_i is noisy is
 also high

- ### * Advantages:-
- i) Resistant to Noise
 - ii) can handle clusters of diff shapes and sizes
 - iii) No need to specify no. of clusters
 - iv) DBSCAN works really well with standard DB's such as MySQL.

- Disadvantages :-
- Very sensitive to hyperparameters
 - High-dimensional data
(may not work very well with text data)
 - using Euclidean distance hence can get in problem of curse of dimensionality.
 - DBSCAN cannot cluster datasets well with large diff' in density. Since the minpts-e combination cannot then be chosen appropriately for all clusters.
 - If domain expert could not understand the data and its scale suitable hyperparam is difficult.

* Time and Space Complexity :-

Time complexity $\approx O(n \log n)$

n pts. range Query
 Euclidean Query

Space Complexity : $O(n) \approx O(nd)$

no. of data pts. often d is small as DBSCAN doesn't work very well with high dim.

Problem formulation: Recommender Systems

Recommender System:- Recommends relevant items to a user based on historical data.

Dataset: A

	Q_1	Q_2	Q_3	\dots	Q_j	\dots	Q_m	
$A =$	u_1							$n: \text{users}$
	u_2							$m: \text{items}$
	u_i				A_{ij}			
	u_n							

Youtube: Programming Cradle

A : matrix A of ratings is very sparse.

$A_{ij} = \text{rating given by } u_i \text{ on } Q_j \text{ (1 to 5)}$

(or)

u_i watched Q_j or not
(0 or 1)

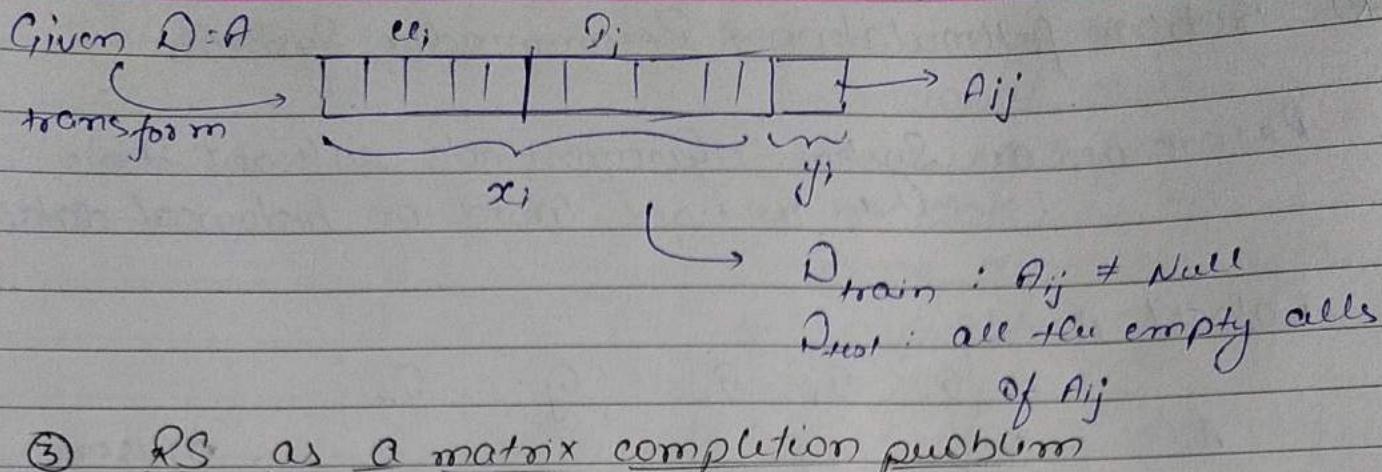
Sparisty of A : $\frac{\text{number of non-empty cells}}{\text{total number of cells}} : \frac{5 \times 10^6}{10^9} = 5 \times 10^{-4}$

→ Big task of building Recommender System is :-

Given a u_i and his/her ratings of previously watched movies / shows, recommend a new item Q_j

① RS can be considered as regression problem if rating is given. (1 to 5)

② RS can be considered as classification problem if watched or not is given (0 or 1)



$A = \begin{bmatrix} \checkmark & \checkmark \\ \checkmark & \checkmark \\ \checkmark & \checkmark \\ \checkmark & \end{bmatrix} \rightarrow \text{Some values of } A_{ij} \text{ are given}$
 $\rightarrow \text{many } A_{ij}'s \text{ are empty.}$

this is known as matrix completion

Youtube: Programming Cradle

Fill up the empty cells with reasonable values based on values in non empty cells.

④ Content based RS vs Collaborative Filtering RS

⇒ Collab. filt. :-

We see collab. b/w users / objects.
 have collab. filt.

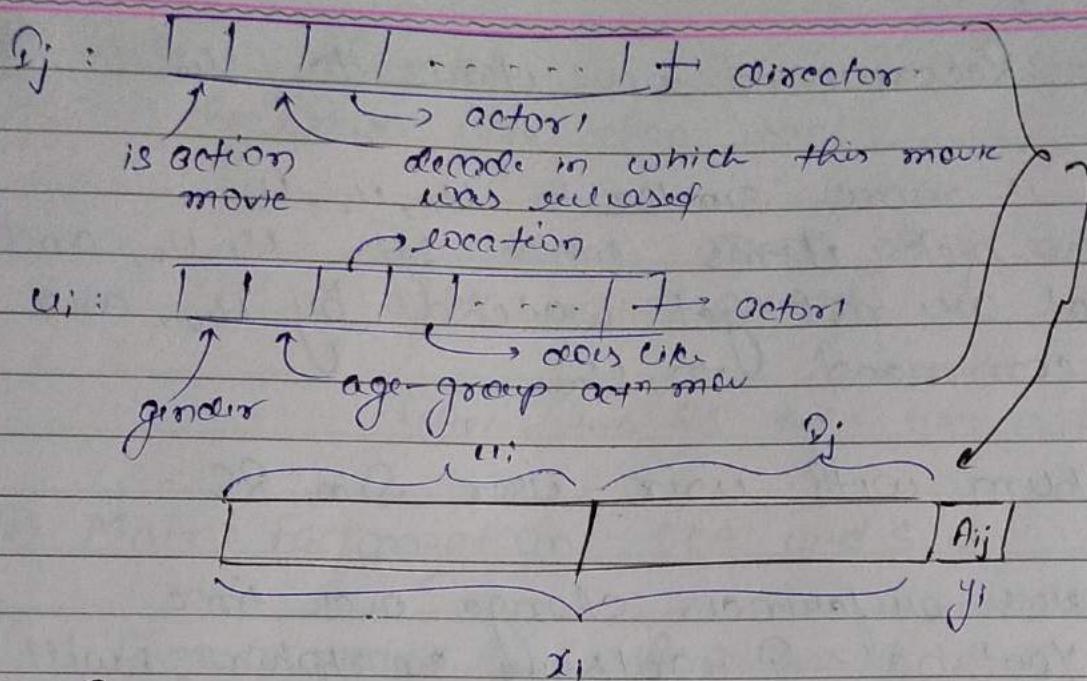
$\left\{ \begin{array}{l} U_1 : M_1, M_2, M_3 \\ U_2 : M_1, M_2, M_3 \\ U_3 : M_1, ? , M_3 \end{array} \right\} \text{ mat A}$

core-idea :- users who agreed in the past (assumption) tend to also agree in the future.

⇒ Content based :-

- doesn't use the A_{ij} info as label
- represent v_j & u_i into various features.

features v_j / u_i : action, drama, 90's, beach, octer, director, ...
 features u_i : like action movie, likes Rom-com, male, ...



* Similarity based Algorithms :-

- ↳ ① item-item Similarity
- ↳ ② user-user Similarity.

* User-user Similarity :-

u_i :- user-vector $A = \begin{bmatrix} u_1 \\ u_2 \\ u_i \\ u_n \end{bmatrix}$

$$u_i = \begin{bmatrix} u_1 \\ u_2 \\ \vdots \\ u_n \end{bmatrix}$$

$$\text{Sim}(u_i, u_j) = \text{Cosine}(u_i, u_j) = \frac{u_i^T u_j}{\|u_i\| \|u_j\|}$$

$\rightarrow S_{ij}$

$S^u = S_{ij} = \text{user Similarity matrix } (S^u)_{n \times n}$

$$\begin{array}{c|c} u_1 & u_1 \\ u_2 & u_2 \\ \vdots & \vdots \\ u_i & u_i \\ \vdots & \vdots \\ u_n & u_n \end{array} \quad T \rightarrow u_i \text{ Similar to } u_j ?$$

Task :- Recommend new items to U_{10}

Let U_{10} is more similar to U_1, U_2, U_3 .

Now pick items liked by U_1, U_2 , and U_3 that are not yet watched by U_{10} , and recommend to U_{10} .

★ Problem with user-user Sim RS :-

user's preferences change over time

e.g. - YouTube :- I want to buy smartphone, i will watch more smartphone reviews.

After some day I want to buy laptop. i will watch laptop reviews and so on...

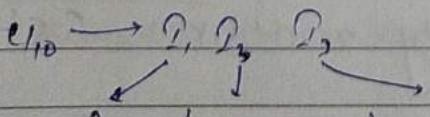
Youtube: Programming Cradle

NOTE :- This problem can be resolved by limiting the data to recent 3 months, but this will lead to more sparse data ^{set}.

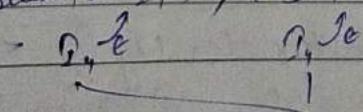
★ Item-Item based RS :-

$$\text{Sim}(I_i, I_j) = \cos(\theta_i, \theta_j)$$

★ Ratings on a given product/item do not change significantly over time after the initial period.



{prod Sim to $I_1\}$, {sim to $I_2\}$, {sim to $I_3\}$ }



Recom. I_2 to U_{10}

✳ Rule of thumb:-

When there're more user than items



Item ratings do not change much over time after initial period.



Item-Item RS over user-user RS

✳ Matrix Factorization: PCA and SVD

Given a matrix A , the process of decomposition of A into several matrix is called as matrix fact. multiplicative

$$A = \underbrace{BCD}_{\text{product of other matrix}} \therefore B, C, D \text{ are factors of } A$$

PCA :- $X_{n \times d} \rightarrow$ Data Matrix

PCA functionality :-

① dim ↓

② mat.fact.

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

Youtube: Programming Cradle

eigen-values :- $\lambda_1, \lambda_2, \dots, \lambda_d$

eigen-vectors :- w_1, w_2, \dots, w_d

$$S_{d \times d} = w_{d \times d} \Lambda_{d \times d} w_{d \times d}^T \quad \left. \right\} \rightarrow \text{eigen-decomposition of}$$

covar mat

$$w = \begin{bmatrix} w_1 & w_2 & \dots & w_d \end{bmatrix} ; \Lambda_{d \times d} = \begin{bmatrix} \lambda_1 & & & \\ & \ddots & & \\ & & \lambda_d & \end{bmatrix}$$

$d \times d$

eigen vectors of S

diagonal mat.

(*) Singular value Decomposition (SVD)

↳ Mat. fact. technique related to PCA

→ PCA is applied on co-var(s) mat. which is sgr. mat.

→ SVD can be applied on any rectangular mat.

$$\rightarrow \text{PCA} \leftarrow S_{d \times d} \leftarrow (W_{n \times d})$$

$$\rightarrow \text{SVD} \leftarrow (D_{n \times d})$$

$$(*) X_{n \times d} = U \sum_{i=1}^r \sqrt{\lambda_i} V_i^T$$

cols of U , are called
 left singular vectors of X

$\sqrt{\lambda_i}$ called right singular
 vectors of X

Singular values (s_i)
 $\frac{s_i^2}{n-1} = \lambda_i$

$$Z = \begin{bmatrix} s_1 & 0 & 0 & 0 \\ 0 & s_2 & 0 & 0 \\ 0 & 0 & s_3 & 0 \\ 0 & 0 & 0 & s_d \end{bmatrix}_{n \times d} \rightarrow \text{diagonal matrix.}$$

Youtube: Programming Cradle

U_i = ith col of U = eig. vect of $(X_{n \times d} X_{d \times n}^T)$

V_i = ith col of V = ith eig. vect of $(X_{d \times n}^T X_{n \times d})$
 S (cov. mat.)

$$\text{cov}(x) = S = W \Lambda W^T \rightarrow \text{PCA}$$

eig-val of S (λ_i)

eig-val of S (λ_i)

eig-vect of $(X^T X)$

$$X = U \sum_{i=1}^r \sqrt{\lambda_i} V_i^T$$

eig-vect of $(X^T X)$

$$\text{Singular val } s_i : \frac{s_i^2}{n-1} = \lambda_i$$

★ Non-negative Matrix Factorization :- (NMF or NNMF)

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

$B : n \times d$

Such that $B_{ij} \geq 0 \quad \forall i, j$ $C : m \times d$
 $C_{ij} \geq 0 \quad \forall i, j$

∴ All the values of mat B and C (factors of A) are non-negative

★ Matrix factorization for collaborative filtering :-

$$A = \begin{bmatrix} & & \\ & \ddots & \\ u_i & & \boxed{a_{ij}} \end{bmatrix}$$

→ A_{ij} : rating on i_j by u_i

→ sparse matrix : many empty cells

Youtube: Programming Cradle

Let $A = B * C^T$

$$\begin{array}{c} \text{m rows} \\ \downarrow \\ \text{users} \end{array} \quad \begin{array}{c} \text{n columns} \\ \downarrow \\ \text{items} \end{array} \quad \begin{array}{c} n \times d \\ \text{d} > 0 \\ \text{d} \leq \min(m, n) \end{array}$$

$$\begin{cases} B : n \times d \\ C : m \times d \end{cases}$$

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

$$A = \begin{bmatrix} & & \\ & \ddots & \\ & & a_{ij} \end{bmatrix}_{n \times m}$$

$$B = \begin{bmatrix} u_1 \\ u_2 \\ \vdots \\ u_n \end{bmatrix}_{n \times d}$$

$$C = \begin{bmatrix} a_1 \\ a_2 \\ \vdots \\ a_m \end{bmatrix}_{m \times d}$$

→ Task : find 'B' and 'C' :- (MF)

optimization eqn squared loss.

$$\underset{B, C}{\text{augmin}} \sum_{ij} \overbrace{(A_{ij} - B_i^T C_j)^2}^{\text{such that}} \rightarrow ①$$

A_{ij} is not empty

• After solving eqn ① we will get

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

$$C = \begin{bmatrix} c_1 \\ c_2 \\ \vdots \\ c_j \\ \vdots \\ c_n \end{bmatrix}_{d \times n}$$

① Solve optimization problem using SGD

$$\underset{B, C}{\text{augmin}} \sum_{ij} \overbrace{(A_{ij} - B_i^T C_j)^2}^{\text{valid nonempty } A_{ij}}$$

$$② B = \begin{bmatrix} B_1 \\ \vdots \\ B_i \\ \vdots \\ B_n \end{bmatrix}, C = \begin{bmatrix} c_1 \\ \vdots \\ c_j \\ \vdots \\ c_n \end{bmatrix}$$

construct B and C

③ Matrix completion :-

$$A = \begin{bmatrix} & & 5 \\ 10 & \begin{bmatrix} \square & 1 \end{bmatrix} & \text{empty} \end{bmatrix}$$

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

we have computed this hence we can fill empty cell with this value

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

actual rating: $\nwarrow A$

predicted rating: $\nwarrow A = B * C^T$

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

without empty cell

* Matrix factorization for future engineering.

want to make u_i and v_j by using

$$A = \begin{bmatrix} v_j \\ u_i \\ \square \end{bmatrix}$$

$\left\{ \begin{array}{l} u_i \quad \boxed{\square \square \square} \rightarrow d \text{ dim} \\ v_j \quad \boxed{\square \square \square} \rightarrow d \text{ dim} \end{array} \right.$

using the data in A_{ij}

$$A = B C^T$$

$n \times m$ $n \times d$ $d \times m$

users \nwarrow items

$$B = \begin{bmatrix} 1 & 2 & 3 & \dots & d \\ \vdots & & & & \\ i & & & & \\ \vdots & & & & \\ n & & & & \end{bmatrix}$$

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

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

* If $u_i, v_j \rightarrow$ very similar in their ratings of products/items
 $\Rightarrow \text{dist}(u_i, v_j) \rightarrow$ small

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

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

$c_i = v_j \in \mathbb{R}^d$

* $u_i, v_j \rightarrow$ very similar in their ratings $\therefore \text{dist}(u_i, v_j) \rightarrow$ small

* clustering as MF :-

K-Means $\therefore D = \{x_i\}$

$$\min_{c_i} \sum_{i=1}^k \sum_{x \in s_i} \|x - c_i\|^2 \quad x \in \mathbb{R}^d$$

lets define Z s.t $Z_{ij} = \begin{cases} 1 & \text{if } x_j \in s_i \\ 0 & \text{otherwise} \end{cases}$

$$Z = \begin{matrix} & 1 & 2 & 3 & \dots & j & \dots & n \\ \begin{matrix} 1 \\ 2 \\ \vdots \\ K \end{matrix} & \left| \begin{matrix} 0 \\ 0 \\ 1 \\ 0 \\ 0 \end{matrix} \right| \end{matrix}$$

$Z_{ij} = 1 \Rightarrow$ this means point j belongs to cluster i

$\xrightarrow{\text{no. of clusters}} \xrightarrow{\text{no. of datapts/rows in original dataset}}$

$$\min_{c_i} \sum_{i=1}^K \left(\sum_{x \in s_i} \|x - c_i\|^2 \right) \rightarrow ①$$

$$\min_{c_i, z_{ij}} \sum_{i=1}^K \sum_{j=1}^n z_{ij} \|x_j - c_i\|^2 \rightarrow ②$$

$$z_{ij} = \begin{cases} 1 & \text{if } x_j \in s_i \\ 0 & \text{otherwise} \end{cases}$$

let's define some matrices

$$X = \begin{bmatrix} & \uparrow & \uparrow \\ x_1 & x_2 & \dots & x_n \\ \downarrow & \downarrow & \dots & \downarrow \\ \underbrace{\quad \quad \quad}_{\text{each col is a data point}} \end{bmatrix}_{d \times n}$$

$$C = \begin{bmatrix} \uparrow & \uparrow & \uparrow \\ c_1 & c_2 & \dots & c_K \\ \downarrow & \downarrow & \dots & \downarrow \\ \underbrace{\quad \quad \quad}_{\text{each col is centroid}} \end{bmatrix}_{d \times K}$$

$$Z \rightarrow z_{ij} = \begin{cases} 1 & \text{if } x_j \in s_i \\ 0 & \text{otherwise} \end{cases}$$

using above matrices we can write eq ② as follows.

$$\min_{c, z} \|x - cz\|^2 \quad \because \text{find } c, z \text{ such that } x - cz \text{ is minimized } (x - cz \text{ ideally})$$

$$\left\{ \begin{array}{l} \text{Such that } z_{ij} = 0 \text{ or } 1 \quad ① \\ \therefore z_j: \text{col of } Z \quad ② \end{array} \right.$$

\hookrightarrow sum of elements in a particular col is 1

NOTE:- \rightarrow If we don't have above 2 constraints K -means \approx MF
 \rightarrow If we have constraints K -means = MF + col. constraints + 0/1 constraints.

MF :- c and z

NNF :- $z_{ij} \geq 0$ and $c_{ij} \geq 0$

K-means :- $\rightarrow z_{ij} = 0$ or 1

\rightarrow col-sum of $z_j = 1$

* Hyperparameter tuning :- d

$$A_{n \times m} = \begin{bmatrix} 1 & 2 & 3 & \dots & m \\ i \\ n \end{bmatrix}_{n \times m}$$

$$A_{n \times m} = B_{n \times d} C_{m \times d}^T$$

$$d > 0$$

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

Methods

① problem specific :-

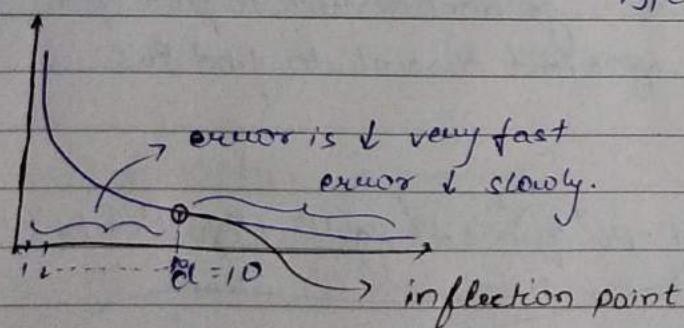
eg:- Given ratings data where $(v_i, \text{rates } r_j) \mapsto (A_{ij})$
 then for each v_i and r_j we want 100-dim
 we don't want dim to be too small or too large

② Systematic way

Youtube: Programming Cradle

$$\text{optimization} : \min_{B, C} \sum_{i, j} (A_{ij} - B_i^T C_j)^2$$

loss / error



* Netflix Prize :- (Wikipedia)

Performance-metric \rightarrow RMSE

$$\sqrt{\frac{1}{n} \sum_{ij} (x_{ij} - \hat{x}_{ij})^2}$$

↑ ↑
rating pred. val.

data & rating,
given u, m, r, d
user movies rating

Read paper [datajobs.com \rightarrow recommendation-systems-
(Netflix).pdf]

* Optimization :-

$$\min_{p, q} \left\{ \sum_{ui} (x_{ui} - q_i^T p_u)^2 + \lambda (||q_i||^2 + ||p_u||^2) \right\}$$

Solve:

Youtube: Programming Cradle

① SGD :- $\frac{\partial L}{\partial p_u}$ and $\frac{\partial L}{\partial q_i}$ \rightarrow takes more time.

② Alternating Least Squares (ALS)
 a) fix p_u ; gradient descent to find q_i 's
 b) fix q_i ; gradient descent to find p_u 's

$$\min_{q, p, b} \sum_{ui} (x_{ui} - \mu - b_u - b_i - p_u^T q_i)^2 + \lambda (||q_i||^2 + ||p_u||^2 + b_u^2 + b_i^2)$$

avg rating across
all users and
items

$$\mu = \frac{1}{n} \sum u_{ui}$$

const.
 user bias
 +ve if user gives
 higher rating than
 other users
 -ve if user usually
 gives lower rating
 than other users
 item/more bias
 interaction effect

→ ut Joe rates Titanic as 3.9

3.9 can be decomposed as below:-

$$\text{ut} = \underbrace{(3.0 + 0.5 + (-0.3) + 0.7)}_{\text{avg rating}} + \underbrace{\text{item/movie bias}}_{\downarrow} + \underbrace{\text{user bias}}_{\rightarrow \text{user}} + \underbrace{\text{put } g_i}_{\text{interaction effect}}$$

→ In Research paper:-

ratings by users and items are time dependent

$$\rightarrow bu(t) \rightarrow p_u$$

$$\rightarrow bi(t) \rightarrow x_{ui}(t)$$

* Cold Start Problem :- Youtube: Programming Cradle

→ new user (what item to recommend to this user?)

→ new Item (to whom this item should be recommended?)

① New user :- Top items based on user's meta-data such as

content based RS $\begin{cases} \rightarrow \text{geo-location} \\ \rightarrow \text{browser (chrome, mozilla etc)} \\ \rightarrow \text{device (android, ios)} \end{cases}$

② New Item :- meta data

↳ category, price, description

* Word-vectors using MF

→ W2V using SVD

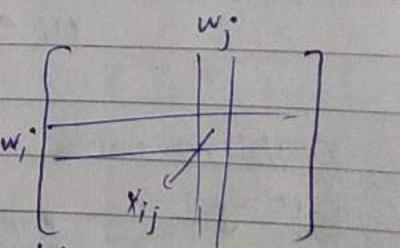
① find co-occurrence matrix (X)

x_{ij} = no. of times w_j occurs in context of w_i

w_j in context with w_i if :- they are in neighbour.

e.g. reviews :- $w_1 w_2 w_3 w_4 w_5 w_6 w_7 w_8 w_9 w_{10}$ $\therefore \text{neighbour} = 5$

w_j is in context with w_i if w_j is at a distance less than 5



② Decompose $X_{n \times n}$ using SVD

$$X_{n \times n} = U_{n \times n} \sum_{n \times n} V^T_{n \times n}$$

$$\frac{s^2}{n-1} = \sigma_i^2 \text{ eig. value of } S(\text{cov. mat})$$

$$\begin{matrix} & u \\ \begin{matrix} 1 \\ 2 \\ \vdots \\ n \end{matrix} & \left[\begin{matrix} 1 & 2 & \dots & n \end{matrix} \right] \end{matrix} \left[\begin{matrix} 1 & 2 & \dots & n \\ s_1 & s_2 & \dots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ 0 & & & s_n \end{matrix} \right] \left[\begin{matrix} 1 & 2 & \dots & n \\ V^T & & & \end{matrix} \right]$$

$$s_1 > s_2 > \dots > s_n$$

Truncated SVD :- $1 < K < n$ \rightarrow hyperparameter

$$\begin{matrix} & u_1 & u_2 & \dots & u_k \\ \begin{matrix} 1 \\ 2 \\ \vdots \\ n \end{matrix} & \left[\begin{matrix} 1 & 2 & \dots & n \end{matrix} \right] & \left[\begin{matrix} 1 & 2 & \dots & n \end{matrix} \right] & \dots & \left[\begin{matrix} 1 & 2 & \dots & n \end{matrix} \right] \end{matrix}$$

$\boxed{\dots}$:- Truncated / deleted part.

$$X \approx \hat{X} = \hat{U} \hat{\Sigma} \hat{V}^T$$

u_i = word-vector corresponding to word;

③ \hat{U} : each row of \hat{U} as u_1, u_2, \dots, u_n

$$u_i \in \mathbb{R}^K$$

\uparrow word, word, word

problem with :- co-occurrence matrix + Truncated SVD

Soln \rightarrow (Top k words)
use most imp. top m
words instead to all
the words.

\downarrow
v. large $n =$ no. of words

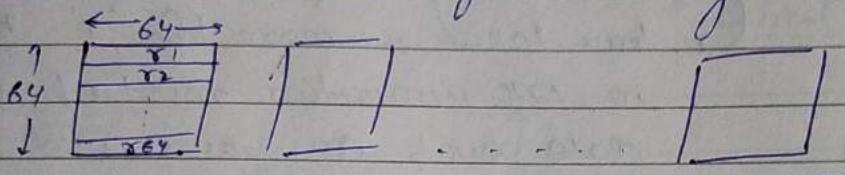
\uparrow because of this computation of
Truncated SVD is very expensive

uk.research.att.com / facerecognition.html

* Eigen faces : PCA on images

→ early technique to recognize faces in images.

dataset consists of 400 imgs. 64×64 px.



$$64 \times 64 = 4096 \quad \text{Step ①} \quad 1 \ 2 \ 3 \dots 4096$$

$$\begin{matrix} & 1 & 2 & \dots & n \\ \begin{matrix} 64 & 64 & \dots & 64 \end{matrix} & \boxed{x_1} & \boxed{x_2} & \dots & \boxed{x_{4096}} \end{matrix} \Rightarrow A = \frac{1}{\sqrt{n}} \begin{bmatrix} x_1 & x_2 & \dots & x_{4096} \end{bmatrix}$$

Youtube: Programming Cradle

Now from $A_{n \times d}$ we can compute 400

$S_{\text{PCA}} = S_{\text{SVD}}$ and :-

$$② S_{\text{SVD}} = U_{d \times d} V_{d \times d}^T$$

$\begin{bmatrix} \downarrow & \downarrow & \dots & \downarrow \\ \downarrow & \downarrow & \dots & \downarrow \\ \downarrow & \downarrow & \dots & \downarrow \end{bmatrix} \quad \begin{bmatrix} 1 & 0 & \dots & 0 \\ 0 & 1 & \dots & 0 \\ \vdots & \vdots & \ddots & \vdots \end{bmatrix}$

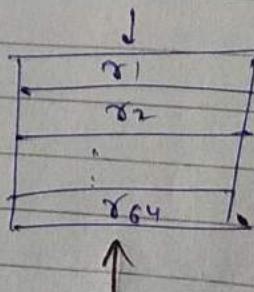
 $d_1 > d_2 > \dots > d_d$

from d-dim to
 k-dim such that $k < d$
 using dimensionality reduction
 (Crop x eigen-val \rightarrow eigen vect.)

$$③ w_k = \begin{bmatrix} v_1 & v_2 - v_1 & \dots & v_k \end{bmatrix}_{d \times k}$$

$v_i \in \mathbb{R}^d \quad d = 4096 \rightarrow \boxed{x_1 \ x_2 \ \dots \ x_{4096}}$

$$④ A_{n \times d} w_k = \tilde{A}_{n \times k} \rightarrow \begin{bmatrix} 1 & 2 & \dots & k \end{bmatrix}_{400}$$



NOTE :- Eigen faces is nothing but taking face dataset and applying PCA on it. eigen faces.

* Important point:-

$u_i \in \mathbb{R}^{17K}$ (is a sparse data) \rightarrow Similarity mat
 \downarrow PCA / SVD

$u_i \in \mathbb{R}^{500}$ \rightarrow Similarity mat

time taken is ' T '

time taken is more than ' T ' because:
 in $17K$ dimension most of the
 data points are zero and hence
 less multiplication, whereas in
 500 dim. it is a dense data (most
 pts are non-zero)