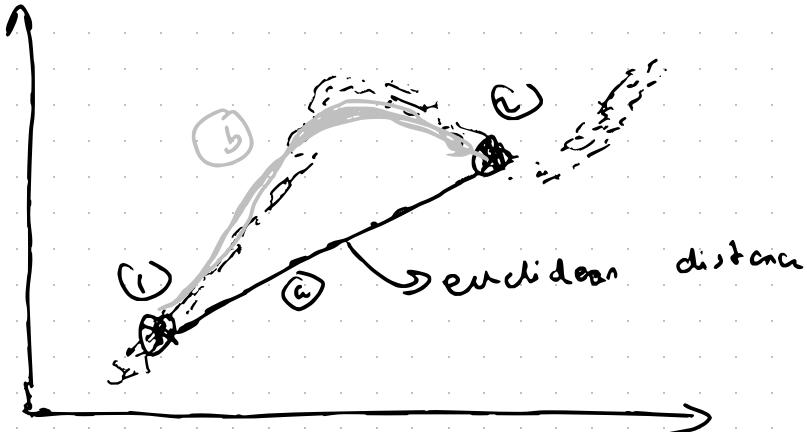


MDS: Multi-Dimensional Scaling
For eg,

(when the embedding manifold is NOT on a hyperplane)



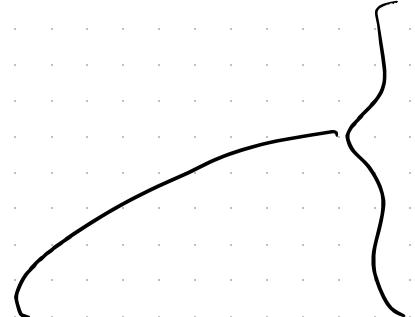
- data is efficiently 1-D, but manifold containing the data is NOT on a hyperplane.
- In many practical scenarios this occurs.
- We also introduce a def^c of a "DISTANCE BETWEEN THE DATA POINTS WHICH IS CONTAINED IN THE EMBEDDING MANIFOLD"
- ⇒ To measure the dist. b/w pts. (1) & (2), we can consider the euclidean distance (black line → ③).

→ However, this dist. clearly lies outside the manifold. So an alternative is the gray line which measures the dist along the manifold.

→ Think of dist. b/w = NYC + Delhi

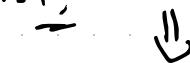
 (1) dist walking along surface / curvature of earth

 (2) dist walking through the earth's core



obviously option (1) is more natural.

→ Many, the natural way of measuring distances b/w 2 st. points on a curved manifold is the geodesic distance? & not the Euclidean dist.



"Geodesic Distance"

→ Informally geodesic dist. can be understood as min. dist. b/w^o the points when the manifold is curved ← joined in a variety of way.
 More formal treatment will be done in future lectures.

⇒ Fundamental idea:

If we are able to find a parameterization γ , in which the geodesic dist. b/w^o points coincides with the cartesian dist., then we have unwrapped our manifold.

⇒ Different manner of defining Dimensional Reduction :-

$$x^i \in \mathbb{R}^d \rightarrow y^i \in \mathbb{R}^d, \text{ originally.}$$



$$\Delta^{ij} \equiv \text{suitable dist. b/w } i \text{ & } j \text{ using original coord.}$$

$$\equiv \text{geodesic dist. } \Delta^{ij} \quad (\text{for example, could also be something else})$$

→ Now we try to find γ^i s.t.

$$\|\gamma^i - \gamma^j\| \approx \delta^{ij} + \epsilon_{ij}$$

⇒ Cartesian dist. \approx if

\approx geodesic dist. \approx if;

for all i, j

⇒ if we can solve this problem, we will have mapped the manifold & found a global parameterization.

$$\|\gamma^i - \gamma^j\| = \text{Cartesian dist.}$$

$$\delta^{ij} = \text{geodesic dist.}$$

→ How can we proceed with finding a set of coordinates $\{\gamma^i\}$, s.t.

$$\gamma: \|\gamma^i - \gamma^j\| \sim \delta^{ij} + \epsilon_{ij} ?$$

\rightarrow method 1 :- Metric MDS (metric multi-dimensional scaling)

$$\sum_{ij} \{ (\|\gamma^i - \gamma^j\|) - d^{ij} \}^2$$

minimize this

nonlinear optimization, $\therefore \|\gamma^i - \gamma^j\|$ is obtained from a square root.

To make it linear, we instead optimize,

$$\sum_{ij} (\|\gamma^i - \gamma^j\|^2 - d^{ij})^2$$

\rightarrow method 2 :- Classical MDS

$$\sum_{ij} (\|\gamma^i - \gamma^j\|^2 - d^{ij})^2$$

\rightarrow (Torgersen, 1916)

$\therefore d^{ij}$ is an arbitrary quantity so we can write it as d^{2ij} or d^{ij} without loss of generality, $\therefore d^{ij}$ is still not determined.

→ Method 3 :- Sketch map :-

$$\sum_{ij} \left(\gamma (||\mathbf{y}^i - \mathbf{y}^j||) - \gamma(\mathbf{z}^{ij}) \right)^2$$

→ γ & γ one nonlinear func^{ns}
cartesian & target distance metric
respectively.

→ Because nonlinear func^{ns} so that neural networks can be used to learn them as well.

→ In method (1) & method (2),

\sum_{ij} double sums are dominated by high-density distribution of points, & global property of the embedding manifold is not captured

→ method (3) :- Was introduced in order to take care of strong density variations, so that all points are not equally weighted.

Classical MDS :-

$$\rightarrow \Delta^{ij} = \|x^i - x^j\|^2$$

→ we want to prove this

$$\rightarrow \text{goal} \leftarrow \text{Find } y^i \in \mathbb{R}^d \text{ s.t., } \|y^i - y^j\|^2 \approx \|\Delta^{ij}\|^2$$

→ This procedure is formally totally equivalent to PCA.

→ But, it can be carried on using ONLY Δ^{ij} ,

→ i.e., we can use only

the distances instead of the coordinates.

→ ∵ we can eventually throw away the coordinates $\{x^i\}$ & keep ONLY Δ^{ij} .

⇒ This mathematical equivalence will hold regardless of data distribution.

→ We start by looking for Y coordinates that are linear transformations of original X coordinates.

$$\therefore \mathbf{y} = \mathbf{A} \mathbf{x}$$

→ rows of \mathbf{A} are mutually orthogonal,
i.e., $\sum_k A_{k1} A_{m1} = \delta_{km}$

Some assumptions as in PCA.

$$\Rightarrow \text{If } \|x^i - x^j\|^2 = \|y^i - y^j\|^2$$

equivalent to requiring :-

$$y^i \cdot y^j = x^i \cdot x^j + t_{ij} \quad [\cdot \text{ is scalar product}]$$

$$\begin{aligned} & \because \|x^i\|^2 + \|x^j\|^2 - 2x^i \cdot x^j \\ &= \|y^i\|^2 + \|y^j\|^2 - 2y^i \cdot y^j \end{aligned}$$

$$\underbrace{\mathbf{y} = \mathbf{A} \mathbf{x}}_{\text{s.t.}}, \quad \sum_k A_{k1} A_{m1} = \delta_{km}$$

then \mathbf{A} is a matrix of that conserves the norm of the vector.

$$\therefore \|x^i\|^2 = \|y^i\|^2$$

$$\therefore x^i, x^j = \mu^i, \mu^j + \epsilon_{i,j}$$

→ Gram matrix of the data :-

dataset: x^i $i = 1, 2, \dots, N$
components of dataset: x^i_λ (no. of features)
 $\lambda = 1, \dots, D$

\therefore Gram matrix $\equiv g$

$$g^{ij} = x^i, x^j = \sum_l x^i_l, x^j_l$$

→ dual of covariance matrix

$$C_{lm} = \frac{1}{N} \sum_i x^i_l x^i_m$$

in Cov. Mat. \equiv sum over upper index
 in Gram. mat. \equiv sum over lower index
 (feature labels)

$$\Rightarrow g^{ij} = (N \times N)$$

$$C_{lm} = (D \times D)$$

$$\Rightarrow \sum_i g^{ij} = 0 = \sum_j g^{ij}$$

$\sum_{i=1}^N x^i = 0$ \leftarrow ($\because x^i$ is centered)
we can always obtain this by subtracting (\bar{x})
from each row of the original data.

$$\Rightarrow g^{ij} = g^{ji} \quad (\text{symmetric})$$

\Rightarrow data we try to construct a matrix which
satisfies the dimensionality + the 2 above
properties of the gram matrix.

$$y \quad N = 10 \text{ K}, \quad D = 10$$

$$C \geq 10 \times 10$$

$$y \geq 10K \times 10K$$

$$\therefore g^{ij} = \sum_{l=1}^D x_l^i x_l^j$$

\therefore max. possible rank of $g^{ij} = 10$

det's say, $M = \sum_{\alpha=1}^n \lambda_\alpha v_\alpha v_\alpha^t$

(Spectral decomposition of the matrix)

$$\sum_{\alpha} \lambda_{\alpha} = 0$$

det's say, $\lambda_1 \neq 0, \lambda_{\alpha > 1} = 0$

rank is 1

$$\therefore M = \lambda_1 v_1 v_1^t$$

det's say, $\lambda_1, \lambda_2 \neq 0, \lambda_{\alpha > 2} = 0$

$$\therefore M = \lambda_1 v_1 v_1^t + \lambda_2 v_2 v_2^t$$

rank = 2

$$\therefore g^{ij} = \sum_{i=1}^D x_i^i \cdot x_j^j$$

Max rank = D

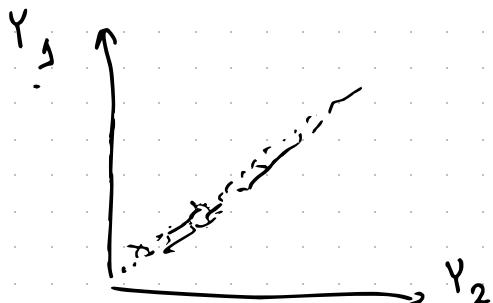
if all x^i are linearly independent

if some x^i are not " then rank < D"

i.e. The rank of g is at most = D

→ Meaning of linearly indep. x^i

embedding manifold in a hyperplane



If the data manifold is contained in a hyperplane of dimension S, then the rank of g = S

\therefore we have shown that

$$\text{rank}(C) = \text{rank}(S)$$

$$(= r)$$

despite $C \equiv D \times D \quad (10 \times 10)$

$$S \equiv N \times N \quad (10K \times 10K)$$

\rightarrow Next, we want to show that

eigenvalues of C = eigenval. of S .

\rightarrow This is imp., "to compute C we need all features, for (S) we need only distances, so proving some equivalence b/w them is vital"

→ We want to solve the following eq².

$$\sum_i g^{ij} \tilde{y}_i = \tilde{\lambda} \tilde{y}_j$$

(eigenvalue eq²)

Let's assume →

$$G \tilde{y}_m = \sum_k A_m x_k$$

$$A: CA = \lambda A$$

(i.e. A_m = eigenvectors of covariance matrix)

$$\sum_k C_{hk} A_{kl} = \lambda_h A_{hl}$$

(similar to in PCA)

$$\sum_i g^{ij} \tilde{v}_m^i = \sum_i g^{ij} \sum_k A_{mk} x_k^i$$

$$= \sum_{l \neq k} x_l^i x_k^j A_{ml} x_l^i$$

$$\left(\because g^{ij} = \sum_{l \neq i} x_l^i x_l^j \right)$$

projection sum over i .

$$= N \sum_{lk} \left(\frac{1}{N} \sum_i x_k^i x_l^i \right) A_{ml} x_k^j$$

$\rightarrow (C_{kl} \equiv \text{cov. matrix})$

$$= N \sum_{lk} C_{kl} A_{ml} x_k^j$$

$$\therefore \sum_i g^{ij} \tilde{y}_m^i = N \sum_k x_k^j C_{kl} A_{ml}$$

(Sum over l)

$$= N \sum_k x_k^j x_m A_{mk}$$

$$\left(\therefore \sum_k C_{hk} A_{kl} = \lambda_l A_{hl} \right)$$

(OR λ_h ?)

$$= N \lambda_m \sum_k x_k^j A_{mk}$$

$$\left(\therefore \tilde{y}_m^i = \sum_k A_{mi} x_k^i \right)$$

$$> N \lambda_m \tilde{y}_m$$

$$\therefore \sum_i g^{ij} \tilde{y}_m^i = N \lambda_m \tilde{y}_m$$

→ From we see that eigenvalues of S are equal to λ_i^N times N .

→ In practice, the eigenvectors of S are given by AX , where A are the eigenvectors of C .

⇒ i. eigenvectors of S are just the Principal Components of our dataset.



→ MDS (contd.) :-

→ focus on the distances.

→ We define in $d^{ij} + t_{ij}$
 $s.t. \gamma^i \in \mathbb{R}^{\delta}, (\delta \ll D)$
 $(x^i \in \mathbb{R}^D)$

$$\|\gamma^i - \gamma^j\| \approx d^{ij} + t_{ij}$$

→ In desired MDS we are solving the following
 least square eq² to minimize the quantity

$$\min_{Y} \left(\sum_{ij} (\|y^i - y^j\|^2 - d^{ij})^2 \right) \quad \text{eq.(1)}$$

$$\text{Def} \quad d^{ij} = \|x^i - x^j\|^2$$

\equiv Square of L^2 norm in original
 coordinate space.

→ We introduced the Gram matrix in the previous lecture.

$$g_{ij} = \mathbf{x}_i^T \mathbf{x}_j$$

(assuming that \mathbf{X} is centered)

$$\Rightarrow \mathbf{y}^i = A \mathbf{x}^i \quad (\text{linear combination})$$

$$\Rightarrow \mathbf{y}_n^i = \sum_m a_{im} \mathbf{x}_m^i \quad (\text{Ansatz } (1))$$

⇒ rows of A are orthonormal
($\because A$ is an orthonorm)

Solving eq (1) using these 2 cond's :-

\mathbf{Y} are chosen such that

Gram matrix in \mathbf{Y} space

$$= \begin{pmatrix} 1 & \dots & \dots & \dots & \dots & \dots \\ \vdots & \ddots & \ddots & \ddots & \ddots & \ddots \\ 1 & \dots & \dots & \dots & \dots & \dots \end{pmatrix} \in \mathbb{R}^{n \times n}$$

i.e., $\mathbf{y}_i^T \mathbf{y}_j = g_{ij} + \epsilon_{ij}$

; i.e., $\sum_j (A\mathbf{x}^i - \mathbf{y}_j)^T (A\mathbf{x}^i - \mathbf{y}_j) = \epsilon_i$

\Rightarrow In previous lecture, we have proved that eigenvalues + eigenvectors of S^{ij} are trivially related to " " " " " C (covariance matrix).

$$\Rightarrow C = \langle x \cdot x^T \rangle$$

$$C_{lm} = \frac{1}{N} \sum_i x_l^i x_m^i \quad (D \times D)$$

$$(S^{ij} = N \times N)$$

$$\Rightarrow \text{rank}(C) = \text{rank}(S^{ij})$$

$$CA_l = \lambda_l A_l$$

$$S \tilde{Y}_l = \tilde{\lambda}_N \tilde{Y}_N$$

$$\left. \begin{array}{l} \\ \\ \end{array} \right\}$$

$$\tilde{\lambda}_N = N \lambda_l$$

$$\tilde{Y}_N^i = A_l \cdot x^i$$

$\rightarrow [l = \text{eigenvector index}]$

Vector in \mathbb{R}^5

Vector in \mathbb{R}^D

\rightarrow if To do $P(A)$, G^{ij} is sufficient, C is NOT necessary. We can compute eigenvalues & evecs of G instead of C .

$$D = 100$$

$$N = 10,000$$

$$C \\ G$$

$$100 \times 100$$

$$(10K \times 10K)$$

However, $\text{rank}(G) \leq 100$

\rightarrow diagonalizing & finding eval & evec of C & G will give us same result. However, it's obviously much easier to do it for C than for G , due to massive difference in matrix sizes.

\rightarrow This makes it possible to compute eval & evec of G : many $\lambda \approx 0$. Otherwise, computing eval & evec of G would be impossible. This computation is done iteratively.

~~→ Ref in Lanczos Method~~

one of the methods that allows us to find
the dominant eigenvector of a matrix.

↓
(~longest)

→ then, we project out the dominant evec by doing
a projection on the " " & therefore
removing it, + then we find the 2nd longest,
also with Lanczos + so on.

Cov, Power method / LOBPCG method

→ So, we keep finding the eigenvalues in decreasing
order, in an iterative fashion, instead of
computing them all at once.



→ How to estimate \mathbf{g} , if we know only the distances?

◻ "double centering" allows us to estimate \mathbf{g} from the dist. matrix, under weak assumptions.

→ Result :-

$$g^{ij} = -\frac{1}{2} (\Delta^{ij}) + \frac{1}{N^2} \sum_{kh} \Delta^{kh}$$
$$-\frac{1}{N} \sum_k (\Delta^{ik} + \Delta^{jk})$$

(symm. matrix)

avg. over both indices.

① Proof :-

$$\sum_i g^{ij} = 0$$
$$\sum_j g^{ij} = g^{ji} \rightarrow \text{key properties of Gram Matrix}$$
$$\sum_{ij} g^{ij} = 0$$

→ Double - Centering :- i.e., the Gram Matrix must be centered w.r.t. both its indices.

Let assume, $\Delta^{ij} = \|x^i - x^j\|$

$$= \|x^i\|^2 + \|x^j\|^2 - 2 \sum_{\ell} g^{\ell j} \rightarrow @$$

\rightarrow summing over j :- where, $\sum_j g^{ij} = x^i \cdot x^j$

$$\sum_j \Delta^{ij} = N \|x^i\|^2 + \underbrace{\sum_j \|x^j\|^2}_{\text{as } x \text{ is centered}} + 0$$

↓
b)

$\therefore \sum_j g^{ij} = 0$
as x is centered

\rightarrow summing over i :-

$$\sum_{ij} \Delta^{ij} = N \sum_i \|x^i\|^2 + N \underbrace{\sum_j \|x^j\|^2}_{\text{but } i \neq j}$$

(over dummy indices)

$$\therefore \sum_{ij} \Delta^{ij} = 2N \sum_i \|x^i\|^2 \rightarrow c)$$

Putting eq @, b, c together, we have

$$g^{ij} = -\frac{1}{2} \Delta^{ij} + \frac{1}{2} (\|x^i\|^2 + \|x^j\|^2)$$

$$\therefore g^{ij} = -\frac{1}{2} \Delta^{ij}$$

$$+ \frac{1}{2} \left(\sum_k \Delta^{ik} - \sum_i \Delta^{ik} \|x^i\|^2 \right)$$

eq. (3) for index

from eq. (6)

$$+ \frac{1}{2} \sum_k \Delta^{ik} - \frac{1}{n} \sum_i \|x^i\|^2)$$

$$= -\frac{1}{2} (\Delta^{ij} - \frac{1}{n} \sum_k (\Delta^{ik} + \Delta^{jk}) + \frac{1}{n} \sum_i \|x^i\|^2)$$

$$= -\frac{1}{2} (\Delta^{ij} - \frac{1}{n} \sum_k (\Delta^{ik} + \Delta^{jk}) + \frac{2}{n} \sum_i \|x^i\|^2)$$

(dropping indices)

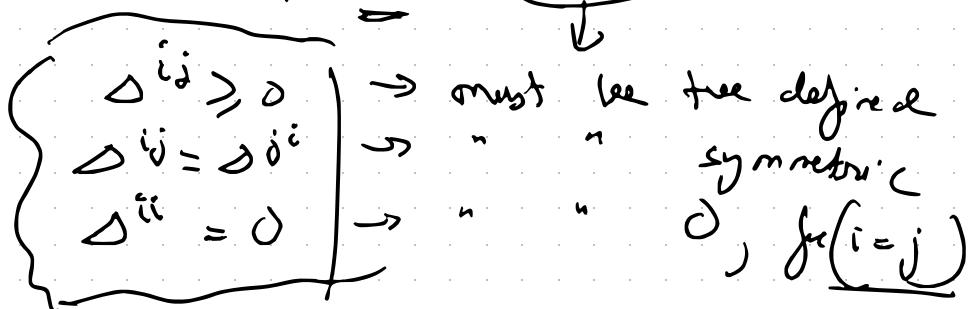
$$+ \frac{2}{n} \cdot \frac{1}{n} \sum_k \Delta^{kk}) \rightarrow \text{from eq. (2)}$$



$$\therefore g^{ij} = -\frac{1}{2} \left(\Delta^{ij} + \frac{1}{n} \sum_k \Delta^{kk} - \frac{1}{n} \sum_k (\Delta^{ik} + \Delta^{jk}) \right)$$

Classical MDS Algorithms

- Given a set of data, estimate distances b/w each pair of data points :-



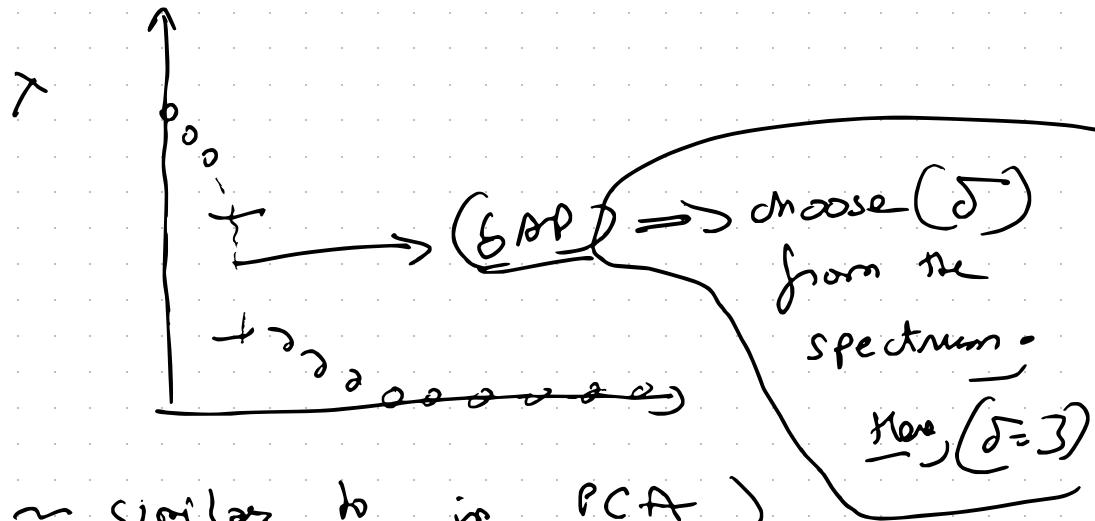
- Estimate δ_{ij} from eq² (j previous page)

Verify that $\delta_{ij} = \delta_{ji}$

$$+ \sum_k \delta_{ik} = 0 \rightarrow \text{addition to define from matrix}$$

$$4 \sum_{k \neq i} \delta_{ik} = 0 \rightarrow \text{from matrix}$$

- Find spectrum of (δ)



• $y^i \in \mathbb{R}^\delta$ are defined by

$$y^i_n = \sqrt{\lambda_i} A_i x^i_n$$

(eigenvectors of the Gram matrix)

if we have δ -“dominant” eigenvectors, we will have

$$\begin{aligned} y^i_n &= A_i x^i_n \\ y^i &= \sqrt{\lambda_i} A_i x^i \end{aligned}$$

“ δ ” \mathbb{R} -coordinates.

$$y^i_n = \sqrt{\lambda_i} A_i x^i_n$$

Scaling factor

→ Why is this procedure justified?

- ① in the specific case that

$$\Delta^{ij} = \|x^i - x^j\|^2 \quad (\text{sum of } L^2\text{-norm})$$

This procedure reduces exactly to PCA.

- ② in the case where $\Delta^{ij} \neq \|x^i - x^j\|^2$,
this procedure can be considered as
a generalization of PCA.

$$y_i \Delta^{ij} = \|x^i - x^j\|^2$$

then $\bar{Y}_1 \propto Y_1 \rightarrow$ (PCA obtained components)
(most retained words) $(\bar{Y}_1 = \sqrt{\lambda_1} Y_1)$

→ We can exploit this procedure, even

$$\text{when } \Delta^{ij} = f(x^i, x^j, X)$$

(some complex func²)

→ However when we have euclidean distances, it
doesn't make sense to use MDS, PCA is more efficient.

$\therefore \text{PCA} \rightarrow d^2 \text{ operation } (CC)$

$\text{MDS} \rightarrow n^2 \rightarrow (G^{ij})$

\rightarrow Be careful MDS often we don't explicitly have the words

\Rightarrow i.e., often we have euclidean distances
PCA is way more efficient than MDS