

Dimension Reduction →

PCA

— Feature Extraction → $[d = d']$

spread
or
variance

Principal Component Analysis

* **Feature Selection** works in general by selecting features / columns with high variance or high spread over the axis of graph.

↳ Variance threshold ← $[d > d']$

* **PCA** → Reduces high Dimension data which are imp to low Dimension Data for better understanding and visualization. $[2D \rightarrow 1D]$ $[3D \rightarrow 2D \text{ or } 1D]$

No. of PCA's = or $< n$ Shifting of Axis
Transform the data acc to Based on Variance of each PC line
max value

PCA

Find unit vector

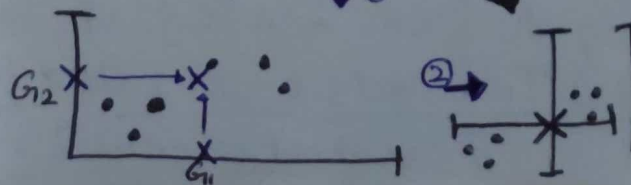
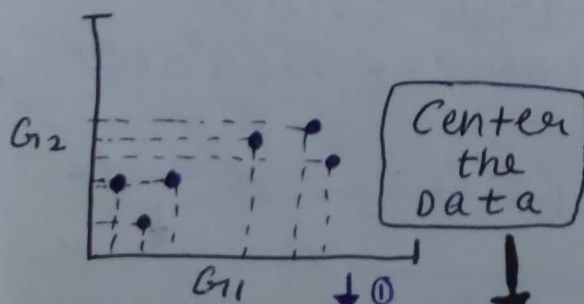
• Project of one vector onto another (one point (dist from $(0,0)$ to (n,y)) vector (\vec{x}) and (\vec{u}) → $\frac{\vec{u} \cdot \vec{x}}{|\vec{u}|}$ as \vec{u} is unit vector so
Projected vector = $\vec{u} \cdot \vec{x} = \vec{u}^T \vec{x}$

PCA retains the valuable insights after Reducing the Dimensions
PCA Also Tells which feature is more important for clustering
PCA Also Tells us about the accuracies of the Reduced Graph

Graphical understanding $[2D - 1D]$

TYPICAL DATA

$\begin{bmatrix} G_1 \\ G_2 \end{bmatrix}$ Mouses



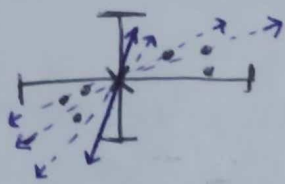
- ① Avg measurement for Gene 1
- ② avg measurement for Gene 2

• With avg values — calculate the center of the data

→ Shift the center of the data To the origins

→ Shifting the data don't change the position of datapoints relatively

→ Try to fit a line that goes through Origin → Rotate the line until it fits the data as well as goes through Origin



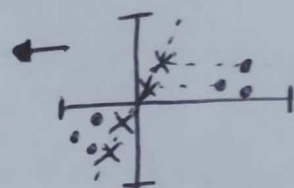
But how to decide the Best fit line?

→ 2 Factors to decide Best fitness of line

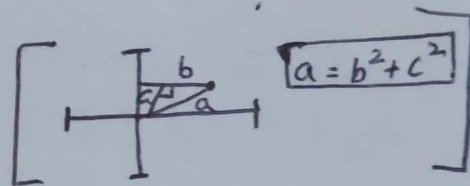
a] measure the Distances from the data to the line that minimizes those distances — (a)

b] maximizes the Distance from The Projected Points to the Origin. — (b)

PCA projects the data onto it



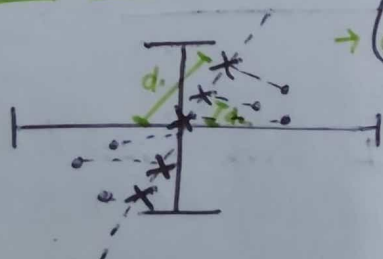
* Best fit line → (a) minimizes
(b) maximizes



What PCA Finally do?

as it is easier to calculate c (or here b) so PCA Finds the Best fitting line by → Maximizing The sum of The squared Distances from the projected points to the origin.

$$\rightarrow (d_1^2 + d_2^2 + d_3^2 + d_4^2 + d_5^2 + d_6^2) = \max \rightarrow SS(\text{distance})$$



∴ decide the line with Largest SS (Distance)

The line is known as PC1. (let's say slope of PC1 = 0.25)

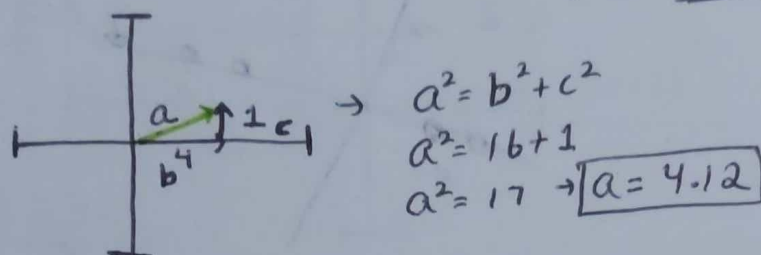
for every 4 units that we go out along the Gene 1, we go up 1 unit along Gene 2. Data is more spreaded Towards Gene 1 and less Towards Gene 2.

Linear Combination of Gene 1 & 2

To make PC 1
 mix 4 Parts Gene 1
 with 1 Part Gene 2

Gene 1 is more significant when it is about describing the spread of data

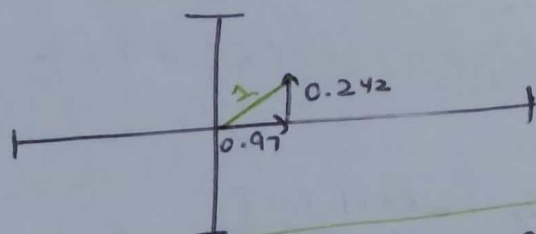
4-1 → Linear Combⁿ



length 'a' which forms as Hypotenous with the help of Linear Combⁿ is scaled to = 1 → unit vector

$\frac{a}{4.12}, \frac{b}{4.12}, \frac{c}{4.12} = 1, 0.97, 0.242$

ratio of G_1 to G_2 Remain same



This unit vector consisting of 0.97 parts of Gene 1 and 0.242 parts Gene 2 / Linear combⁿ of Gene 1 and Gene 2 is known as eigen vector for PC 1

and proportion of each gene/feature is known as

LOADING SCORES

SS (Distance) = sum of square of Projected Points to origin of PC 1

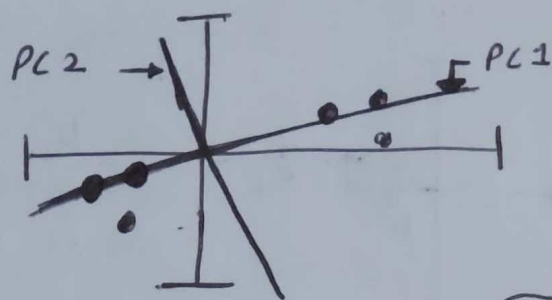
$= d_1^2 + d_2^2 + d_3^2 + d_4^2 + d_5^2 + d_6^2 = \text{Eigenvalue for PC 1}$

SS (Distance for PC 1) = Eigenvalue for PC 1

$\sqrt{\text{Eigenvalue for PC 1}} = \text{Singular value for PC 1}$

PC-2 →

PC 2 is simply the line through the origin that is \perp to PC 1, without any further optimization that has to be done.



loading scores of PC2

PC2 \rightarrow -1 parts of Gene 1
4 parts of Gene 2

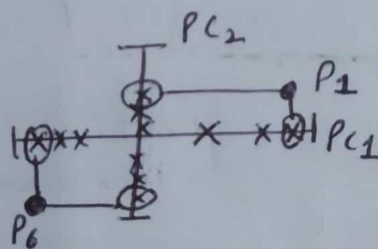
-0.242 parts of Gene 1
0.96 parts of Gene 2

unit \rightarrow eigenvector of PC2

SS (Distance of PC2) = Eigenvalues for PC2

FINAL PCA PLOT

- ① SIMPLY Rotate everything so that PC 1 is horizontal
- ② Use Projected points to find where the samples go in the PCA plot.



That's the general way of doing PCA using Singular Value Decomposition (SVD).

$$\therefore \frac{SS(\text{Distance for PC1})}{n-1} = \text{Variation for PC1} \quad \text{--- X}$$

$$\therefore \frac{SS(\text{Distance for PC2})}{n-1} = \text{Variation for PC2} \quad \text{--- Y}$$

$$X + Y = \text{Total } \frac{X}{2} \quad \left| \quad \left(\frac{X}{2} \right) 100\% \text{ accounts for PC1} \right| \quad \left| \quad \left(\frac{Y}{2} \right) 100\% \text{ accounts for PC2} \right| \quad \left| \quad \text{Variation as } 100\% \text{ PCs} \right|$$

A graphical representation of the % of variation that each PC accounts for = SCREE PLOT

$$\left[\left(\frac{x}{z} \right) \times 100 \right] \left(\frac{y}{z} \right) \times 100$$

X ————— X ————— X

PCA with 3 variables

Follows General Steps

① Center the Data

② Find the Best fitting line that goes through the origin - PC1
 ↳ Reape or loading Score = no. of variables

③ PC2 → Best fitting line that gives through the origin and is perp to PC1
 ↳ Reape or loading Score = no. of variables

④ PC3 → Best fitting line that goes through the origin and is perp to PC1 and PC2.

∴ If we have more genes, we would just keep on finding more and more principal components by adding 1st lines and rotating them.

→ unit vector derived from loading scores normalized or scaled = eigenvector of that PC.

⑤ USE eigen values (SS(Distance)) to determine the proportion of variance (x, y, z) (Divide by n-1)

$$\frac{(x+y+z)}{x} \left(\left(\frac{x}{x+y+z} \right) 100 \right) = \text{proportion of variance.}$$

Eg → PC1 (79.1%) } majority → 94.1% so a 2D graph would be better as 94.1% variation.
 PC2 (15.1%)
 PC3 (4.1%)

IMP on the basis of proportion of Variation the
Decision is made to reduce Complexity of Graph.

Here 3D-2D (Remove everything except PC1, PC2 and data)

- Project the samples on PC1 and PC2
- Rotate so that PC1 is horizontal and PC2 is vertical
- Place the points on new plot.

4D-2D decided on the Basis of score plot or proportion of Variance accountable by Each PCs

If bars in score plot are similar and not have high Difference then Plotting 4D to 2D is noisy But can be used to form CLUSTERS of Data

IMP

Dimension

Reduction

U-MAP

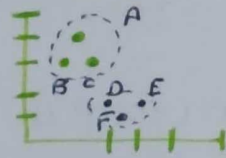
Uniform Manifold Approximation & Projection

• Create a low-dimension graph of this data that preserves the high Dimensional clusters and their rel^{nth} to each other.

• Initialize the low dim points and then move the low-dimension points around until they form clusters that have the same rel^{nth} we saw in the high-dimensional data.

* Projecting points will create mishmash of points irrespective of the ans
Increasing dimensions will make things worse.

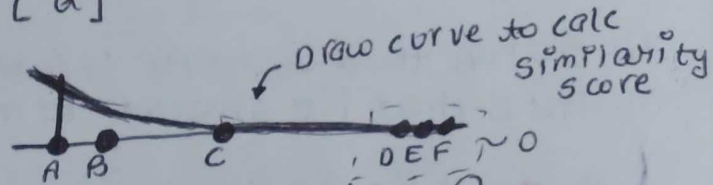
① Calculate the SIMILARITY SCORES that help identify clustered points so it can try to preserve that clustering in the low-dimensional graph.



First thing → calculate the Distance B/w each pair of high dimensional points... [d]

* Similarity Scores for A

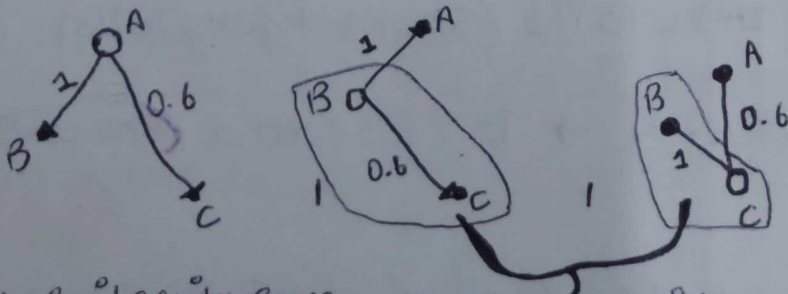
- Plot A on graph
- Plot Point B 'd' distance (eg 0.5) away from A.
- Plot Point C 'd' distance (eg 2.4) away from A.
- Plot Points 'D', 'E' and 'F' far away from A, as they are far away from A in high dimension.



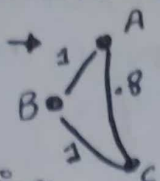
NOTE The shape of the curve depends on the number of high-dimensional neighbors that you want each point to have
 $k \rightarrow$ Default value = 15 But for small it can be 3 (it includes point itself) ($k=3$ so 2 neighbors) ←

- $\log_2(\text{neighbors}) \rightarrow \log_2(3) = 1.6$ (defines shape of curve)
- Curve shape in such a way that y-axis coordinate for the nearest neighbour (B, C) add up to the $\log_2(\text{neighbors}) = 1.6$
 $\frac{1}{1} + \frac{1}{0.6} = 1.6$
- These values of B & C (neighbors) are the similarity score rel to A.

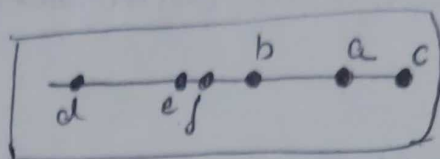
* Similarly curves for B, C, D, E and F are created



similarity score are not symmetrical } They are different } Because Curve for each point are different
 They relative score of B & C

* UMAP makes them symmetrical using a Method similar to taking the average \rightarrow  same for DEF

* THEN \rightarrow UMAP initializes a low-dimensional graph



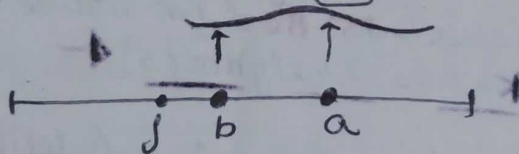
\rightarrow Not ideal as clusters are not similar

To solve this \rightarrow UMAP picks 2 low-dimensional points that it should move closer together

\hookrightarrow Does this by randomly selecting a pair of points in a cluster \propto to their high dimensional score. [A and B]

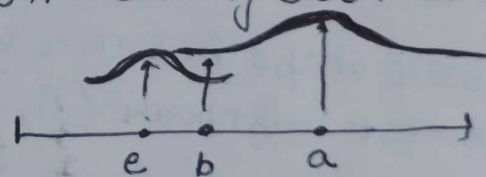
\hookrightarrow a and b closer, b far from f. \rightarrow Randomly picks points that are not in B's high dimensional cluster. [F]

Now \rightarrow How much closer to a and far from f



\hookrightarrow calculate low dimensional similarity scores y-axis coordinate on a curve

fixed Bell shaped curve
same size



t-distribution

* t-distribution curve is like a Gaussian or normal dist but has shorter and fatter tails (short + fatty tails)

* As a and b are in same cluster \rightarrow b close to a (maximize the low dim score)

* B and F \rightarrow Diff cluster \rightarrow b far from F (minimize the low dim score).
or
F

similarly we do it for all remaining pair.

Imp \rightarrow In case where $\left\{ \begin{array}{c} \text{d} \quad \text{e} \quad \text{c} \end{array} \right\}$ The score we want to minimize will still be pretty small + score - we want to minimize would be much larger

t-SNE similar to U-MAP

① T-SNE always starts with a random init of the low-dim graph

① BUT U-MAP uses Spectral Embedding to init the low-dim graph

② t-SNE moves every single point a little bit each iteration

② U-MAP just (move) one point - or a small subset of points, each time - helps it score well with

Number of neighbors (K)

• A relatively low value for the number of neighbors (K) results in \rightarrow small - independent clusters. (details but not the big picture)

• A relatively large value for the no. of neighbors gives you more of the big picture and less of details.

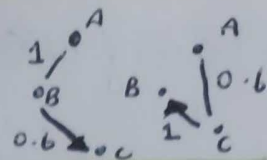
THE MATH

The general eq for the similarity scores (y-axis coordinate)

$$\text{similarity score} = e^{-(\text{raw dist} - \text{dist to nearest neighbor})/K}$$

\therefore gets diff curve for diff point \Rightarrow Change shape of curve \rightarrow make $SS = \log_2(K)$.

Asymmetrical scores B/w points \rightarrow



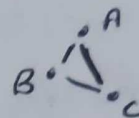
\rightarrow A verging to make them symmetrical $\rightarrow (S_1 + S_2) - S_1 S_2$

$S_1 (B \rightarrow C), S_2 (C \rightarrow B)$

$$\rightarrow (0.6 + 1.0) - 0.6 \times 1.0 = 1.0$$

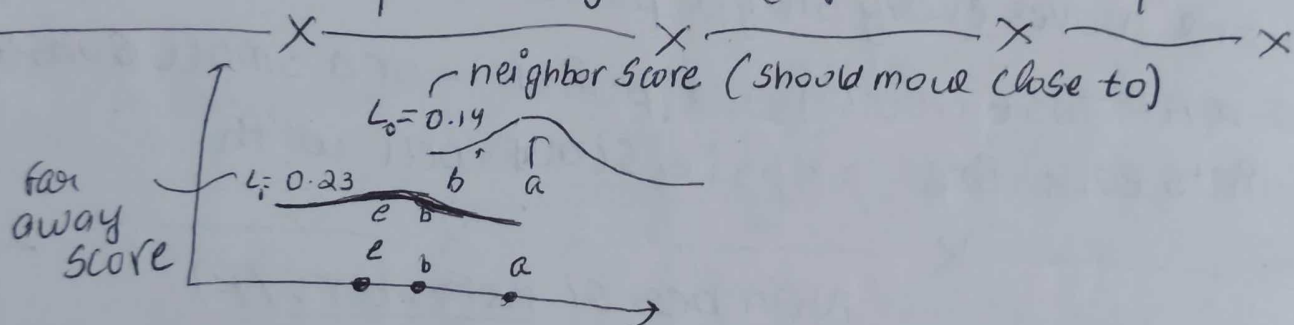
update the score to 1.0
(B and C)

Symmetrical scores $\rightarrow (S_1 + S_2) - S_1 S_2$



low dimensional similarity scores $= \frac{1}{1 + \alpha d^2 \beta} = L$

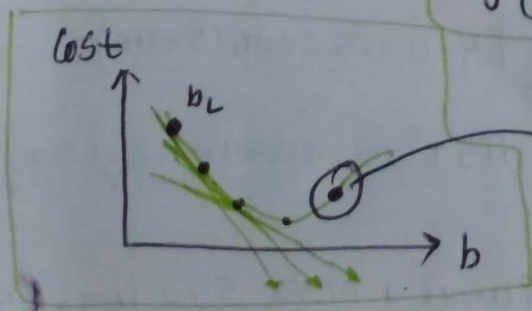
$L \propto \alpha \& \beta \rightarrow$ How tightly the low-dimensional points can be packed Together (default $\rightarrow \alpha = 1.577 / \beta = 0.8951$)



\rightarrow As we move b to a & b from the L_1 gets smaller & L_0 gets higher

\rightarrow Uses L_1 & L_0 to evaluate if b is at Right Place or not \rightarrow Cost funcⁿ = $\log\left(\frac{1}{\text{neighbor}}\right) + \log\left(\frac{1}{1 - \text{not neighbor}}\right)$

$$L \log\left(\frac{1}{L_0}\right) + \log\left(\frac{1}{1 - L_1}\right) = |b_L|$$



b crosses a

Stochastic GD

\rightarrow Concept of gradient descent to achieve minima