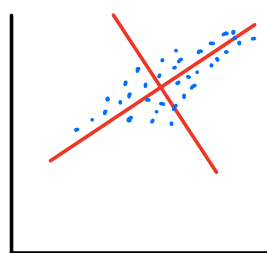


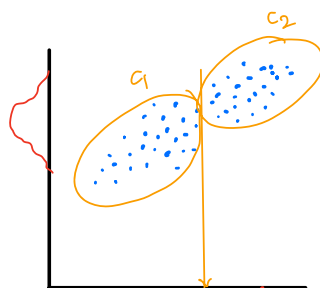
# Linear Discriminant Analysis



PCA

(operates on data,  
linear method,  
unsupervised)

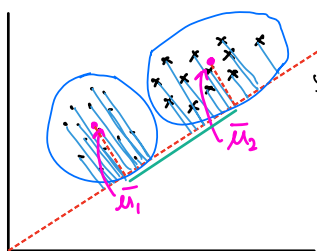
\* give me top 3 principal components



LDA

(operates on feature subspace,  
linear method,  
supervised)

- Data  $\langle x_1, x_2 \dots x_N \rangle$
- $N_1$  samples belong to class  $C_1$
- $N_2$  samples belong to class  $C_2$
- Find a line that maximize the class separation.



$y = w^T x$  find weights

only  $w^T$  change to rotate the line

\* maximize difference between averages

- Define a good separation measure

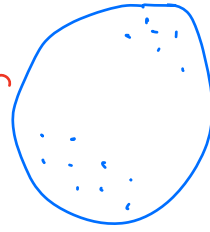
- Measure 
$$\mu_i = \frac{1}{N_i} \sum_{x \in C_i} x$$

$$\mu_i = \frac{1}{N_i} \sum_{y \in C_i} y = \frac{1}{N_i} \sum_{x \in C_i} w^T x = w^T \mu_i$$

- Driving force for separation

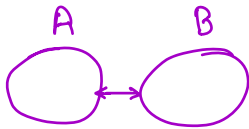
$\text{argmax}_w$  objective  $f(w) = |\tilde{\mu}_1 - \tilde{\mu}_2|$  L1 norm

$$= |w^T(\mu_1 - \mu_2)|$$

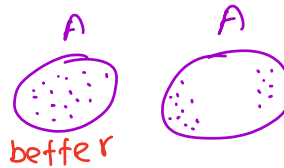


better class  $\longrightarrow$  variance inside class has to be minimum  
Distance to other classes has to be maximum

But we are ignoring variability inside classes.



Fisher Approach



Normalize the distance (difference)

between the means by intra-class

scatter = variance

variance inside class  $\tilde{S}_i^2 = \sum_{y \in c_i} (y - \tilde{\mu}_i)^2$

Intra class scatter  $= \hat{S}_1^2 + \hat{S}_2^2$  (sum should be minimum)

Objective(w) = maximum inter class separation  
minimum intra class variability

Fisher Discriminant

$$\text{Objective}(w) = \frac{|\bar{\mu}_1 - \bar{\mu}_2|^2}{\bar{\Sigma}_1^2 + \bar{\Sigma}_2^2} \leftarrow L2 \text{ norm}$$

t-SNE (t-Distributed stochastic neighbor embeddings)

\* bring anything to 2D

→ non linear data visualizer

→ t-test | t-distribution (normal distribution)

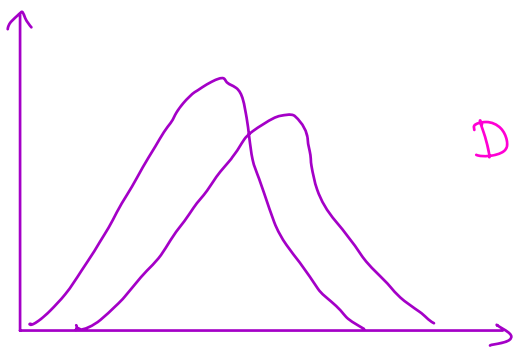
t-SNE doesn't use any norm (distance metric)

~~L1~~, ~~L2~~ .. ~~LP~~

Kullback Leibler Divergence > distance metric

Given 2 probability distributions  $p, q$  the KL divergence measures the distance

$D(p||q)$  How much <sup>does</sup>  $p$  distribution diverges from  $q$



$$D(p||q) = \sum_{x \in \mathcal{X}} p(x) \log \frac{p(x)}{q(x)}$$

universe of discourse  
(set of main events)

$$D(p||q) \neq D(q||p)$$

→ KL divergence is not a metric

ex:- Divergence  $D(\text{observed} || \text{normal})$

How far data is deviate from normal Distribution

$D(\text{observed} \parallel \text{binomial})$

Relationship to entropy  $H(x)$

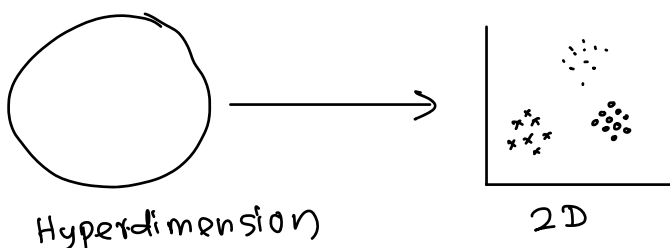
$$H(x) = \sum_{x \in X} p(x) \log \frac{1}{p(x)}$$

$$= \log N - D(p(x) \parallel p_u(x))$$

$\downarrow$  true distribution       $\downarrow$  Uniform distribution

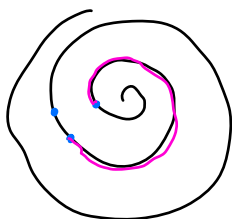
The shannon entropy is the number of bits necessary to identify  $x$  from  $N$  equally 'likely' possibilities less the KL divergence of the uniform distribution from the true distribution.

t-SNE idea



classifier X  
visualizer ✓

Similarity in high dimensions corresponds to short distance in low dimensions.



perplexing !

t-SNE minimizes the sum of KL divergence over all data points using gradient decent method.

$$\begin{aligned}\text{Objective} &= \sum_i D(p_i || q_i) \\ &= \sum_i \sum_j p_{ji} \log \frac{p_{ji}}{q_{ji}}\end{aligned}$$

$$p_{ji} = \frac{\exp\left(-\frac{\|x_i - x_j\|^2}{2\sigma_i^2}\right)}{\sum_{k \neq i} \exp\left(-\frac{\|x_i - x_k\|^2}{2\sigma_i^2}\right)} \quad \left. \vphantom{\frac{\exp\left(-\frac{\|x_i - x_j\|^2}{2\sigma_i^2}\right)}{\sum_{k \neq i} \exp\left(-\frac{\|x_i - x_k\|^2}{2\sigma_i^2}\right)}} \right\} \text{high dimension}$$

$$q_{ji} = \frac{\exp\left(-\|y_i - y_j\|^2\right)}{\sum_{k \neq i} \exp\left(-\|y_i - y_k\|^2\right)} \quad \left. \vphantom{\frac{\exp\left(-\|y_i - y_j\|^2\right)}{\sum_{k \neq i} \exp\left(-\|y_i - y_k\|^2\right)}} \right\} \begin{matrix} \delta_i = \frac{1}{\sqrt{2}} \\ \text{low dimension} \end{matrix}$$

Dim reduction and visualization

\* PCA

\* LDA

\* t-SNE