

- ③ Recap
- ④ Issues with PCA
- ⑤ t-SNE intro
- ⑥ Preserving Relative Distances
- ⑦ Inverse of Distances as PDF
- ⑧ Normalized Probability with ND
- ⑨ Perplexity as an alternative
- ⑩ Crowding problem
- ⑪ t-distribution
- ⑫ KL-Divergence for comparing Distributions

Recap of PCA

* Goal Find new Feature and remove some Not so useful ones

* How? By trying to maximize Avg projection Length

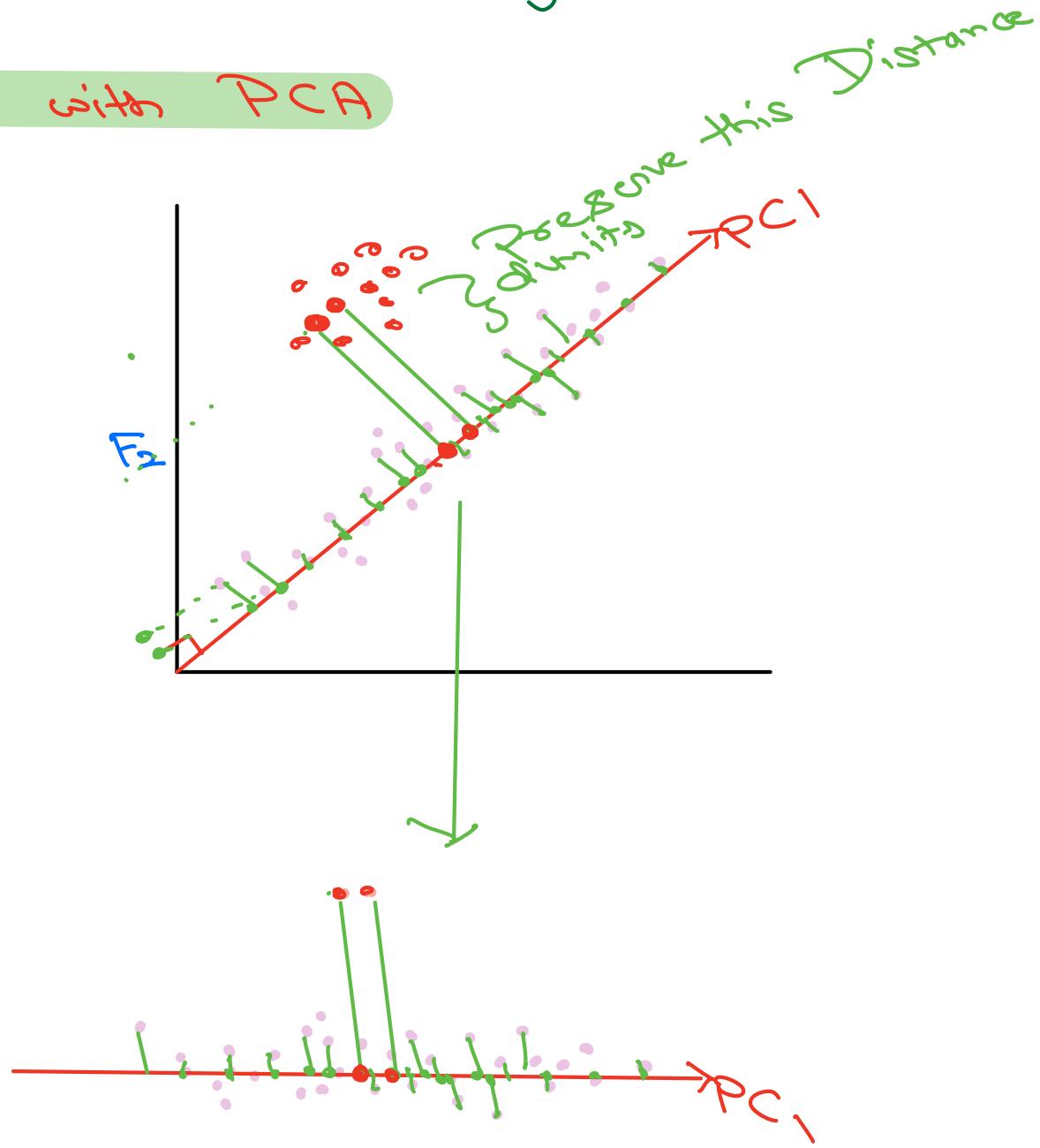
Steps for PCA

- ① Mean centering (Standard Scaling)
- ② Covariance Matrix
- ③ $\nabla \rightarrow x^T x / n$
- ④ find Eigen Values and Eigen Vectors of ∇^{\max}
- ⑤ Sort w.r.t Eigen Values and remove Eigen Vectors with small Eigen values

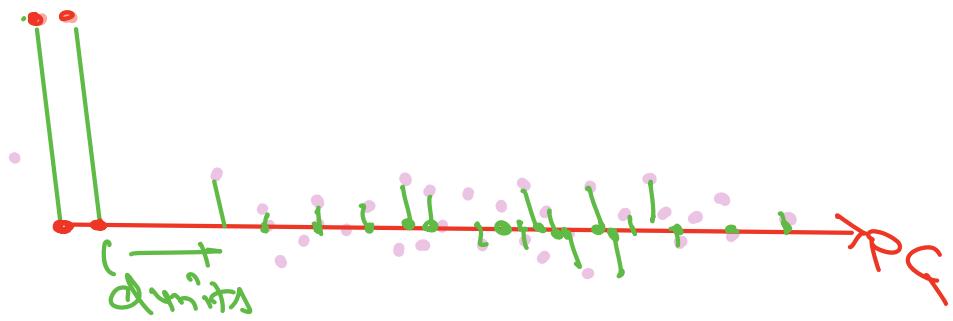
$$\nabla \rightarrow \sum_{i=1}^n \frac{x_i \cdot q_j}{\|x_i\|}$$

$\nabla \nabla = \lambda \nabla$

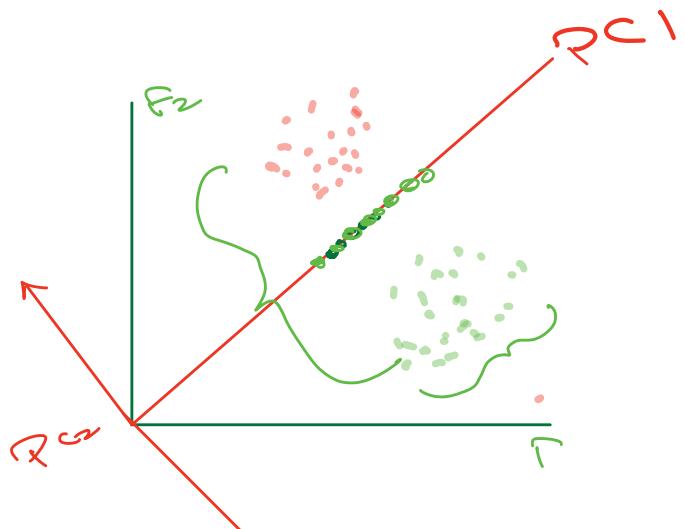
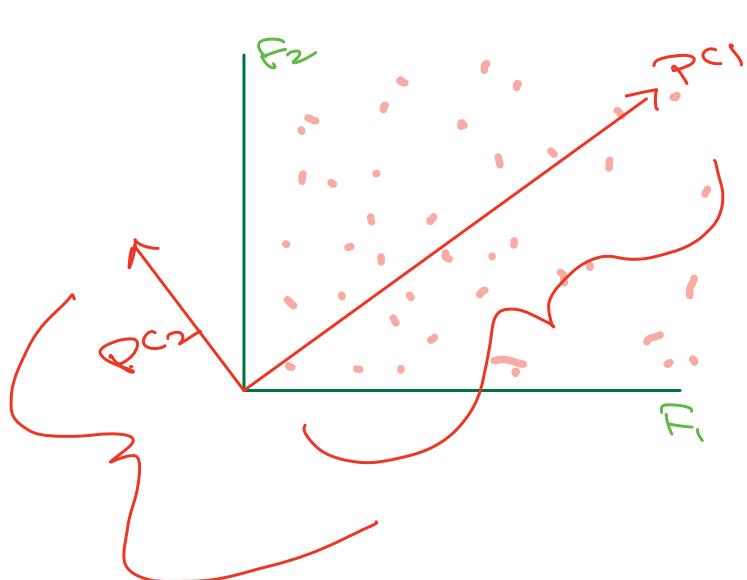
- ⑤ Project your Data to EigenVectors
 $x_i \in \mathbb{R}^n$
 - ⑥ Plot or perform ML training
- ⑦ Issue with PCA



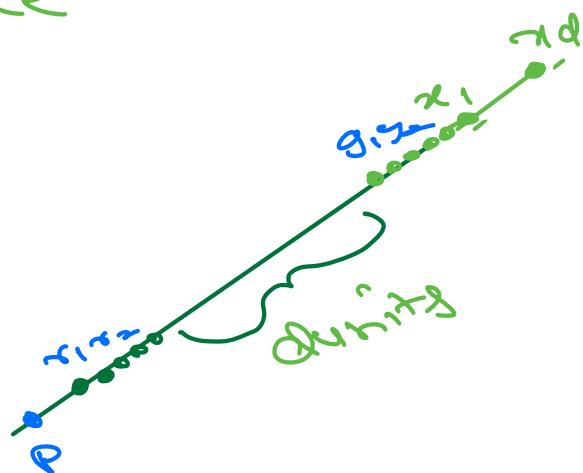
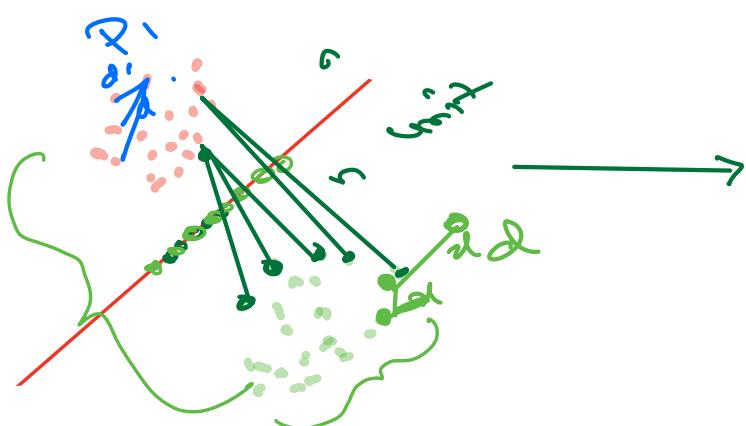
Ideal Scenario:



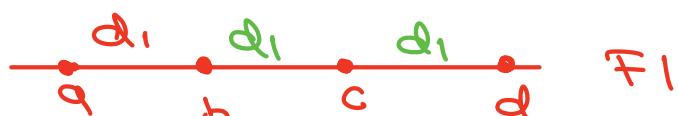
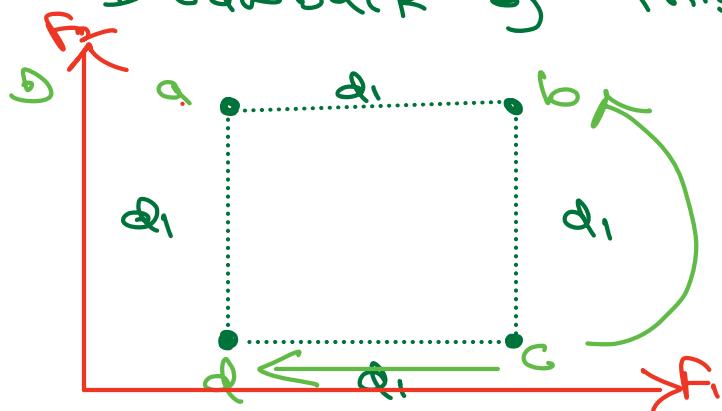
Can't drop any
of the PC's



Idea-2 : preserve the relative
Distance

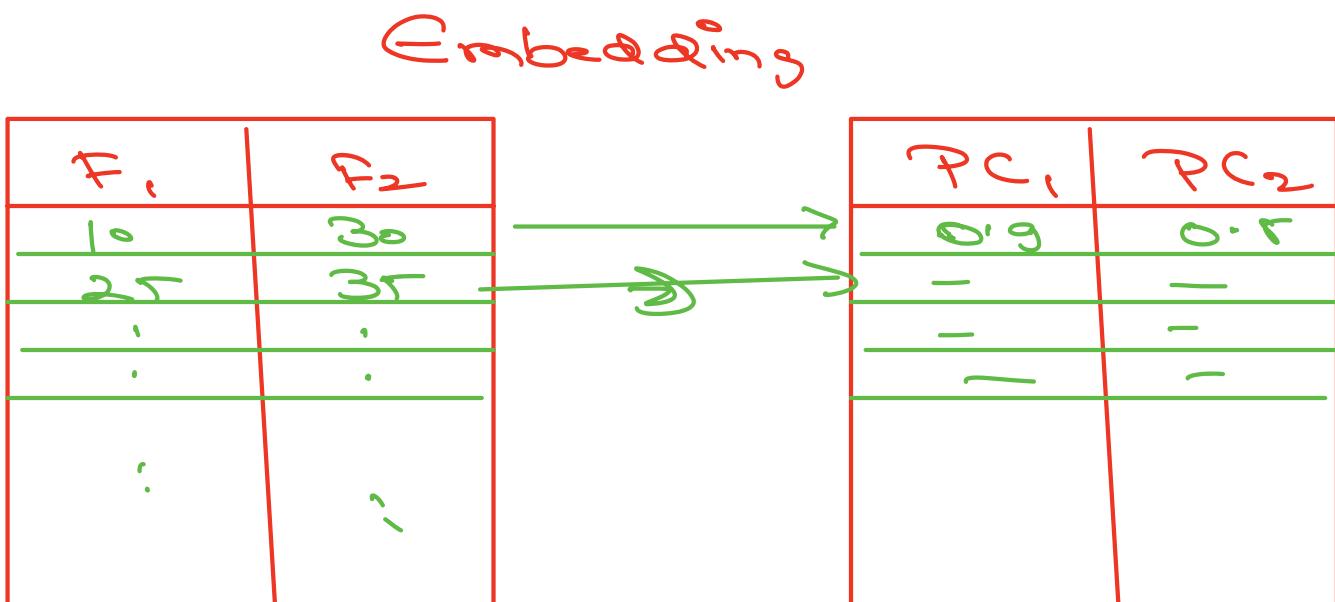
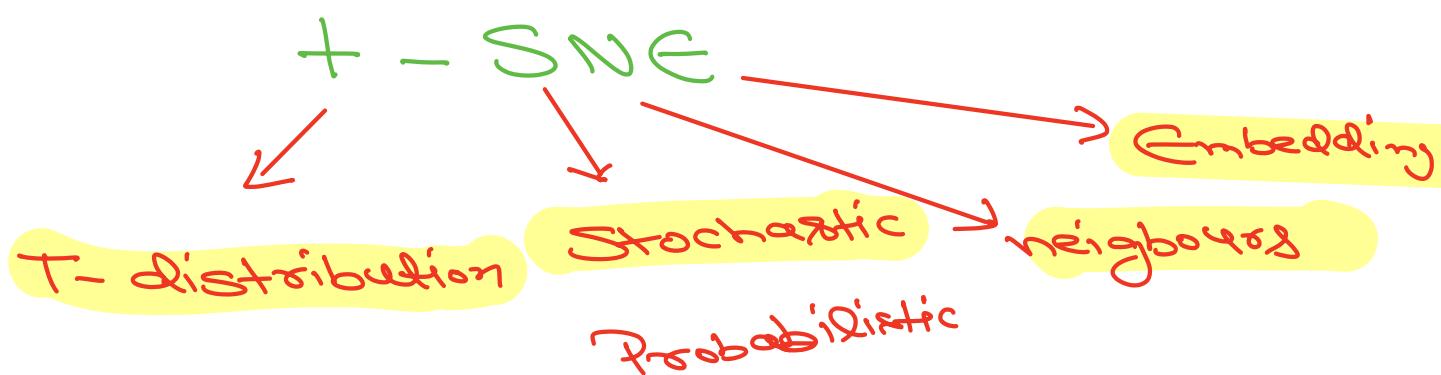


Drawback of this distance approach



* Impossible
without any
restriction

Idea 2: Preserve probability of points as Neighbours

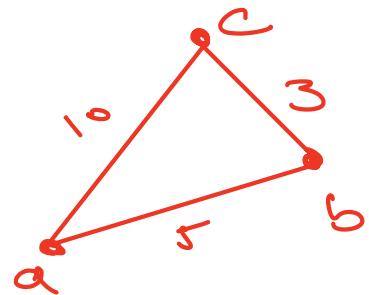


$$PC_1 \propto F_1 + \beta F_2$$

Projected feature space is also called Embedding

How can we convert distance into probability

$\Rightarrow d(A, B) \rightarrow \text{Prob score}$



$P(A, B)$ being Neighbours can be quantified as

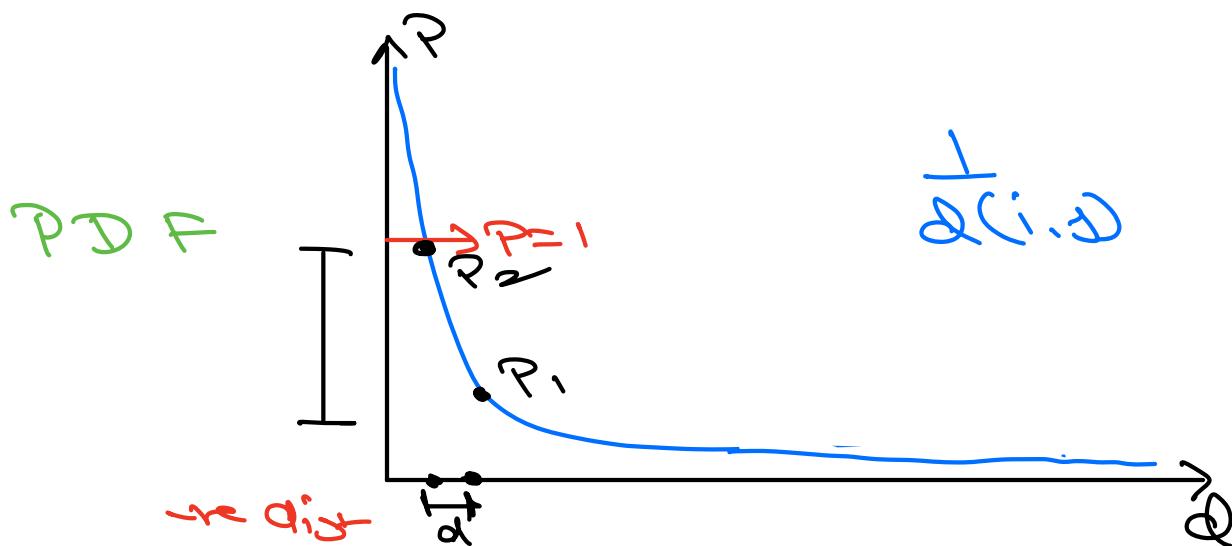
$$\frac{1}{d(A, B)}$$

$$\Rightarrow \frac{1}{\|A, B\|}$$

$$P(a, b) \Rightarrow \frac{1}{d(a, b)} \Rightarrow \frac{1}{5} \Rightarrow 0.20$$

$$P(b, c) \Rightarrow \frac{1}{d(b, c)} \Rightarrow \frac{1}{3} \Rightarrow 33.33$$

$$P(c, a) \Rightarrow \frac{1}{d(c, a)} \Rightarrow \frac{1}{10} \Rightarrow 0.1$$



Issues with PDF $\frac{1}{\delta(x_i, x)}$

① When $\delta \rightarrow 0, P \rightarrow \infty$



② when $\delta \rightarrow 0.2, P \rightarrow 5$



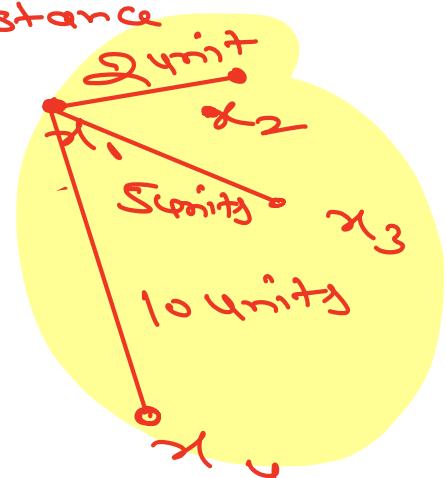
③ Very sensitive to

Small Changes in Distance

$$P(x_1, x_2) \propto \frac{1}{\delta} \propto 0.5$$

$$P(x_1, x_3) \propto \frac{1}{\delta} \propto 0.2$$

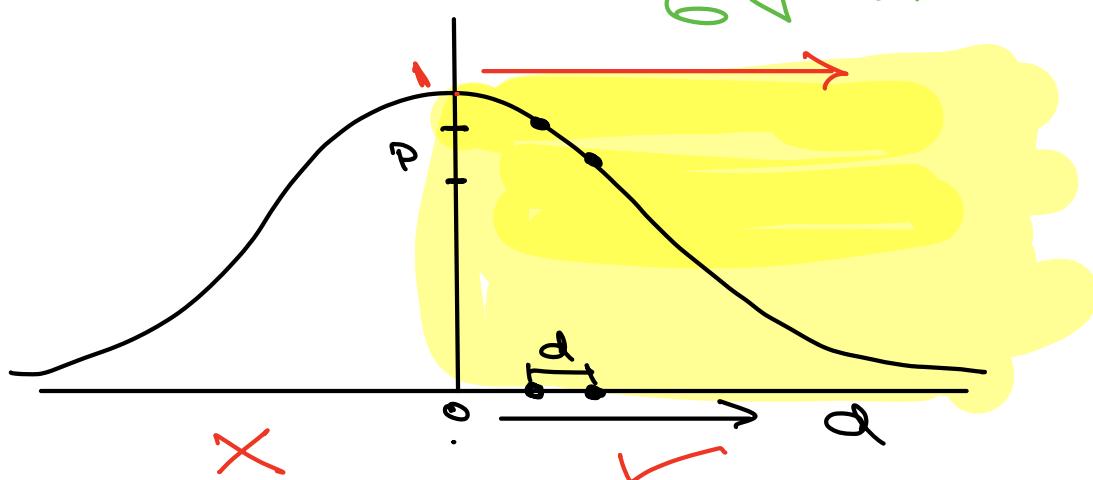
$$P(x_1, x_4) \propto \frac{1}{\delta} \propto 0.1$$



$$\frac{1}{100} \approx 0.01$$

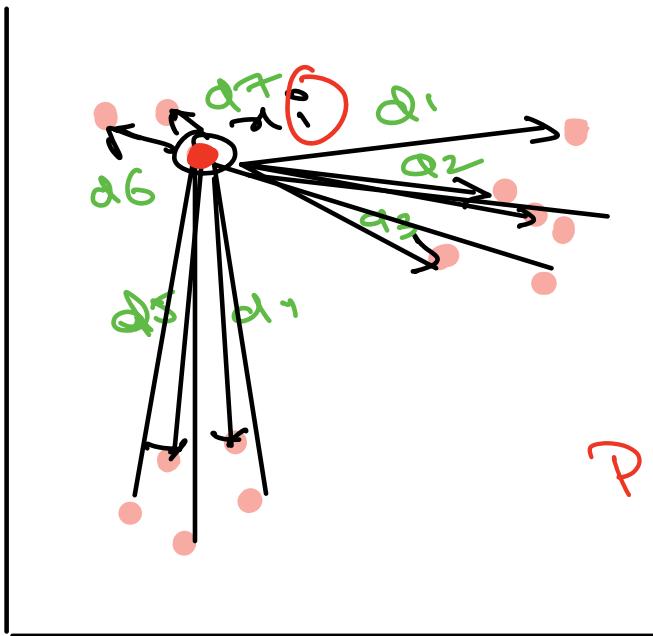
Can we use Normal Distribution instead

$$e^{-\frac{(x-\bar{x})^2}{2\sigma^2}}$$



① $d \rightarrow 0, P \rightarrow 1$

② Small 'increase' in distance does not lead to High decrease in prob
Class Sensitive than $\frac{1}{d(i,j)}$



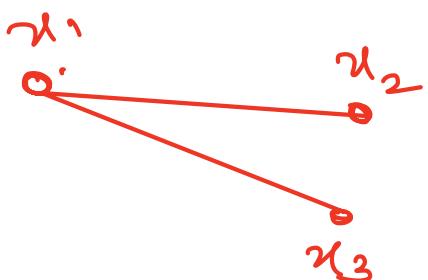
$$P \propto \frac{1}{d(i,j)}$$

with N.D as PDF

$$P(i,j) \propto \frac{1}{\sqrt{2\pi} \sigma} e^{-\frac{(x_j - x_i)^2}{2\sigma^2}}$$

$$\frac{x_j - x_i}{\sigma} \rightarrow (x_j - x_i)$$

x_i



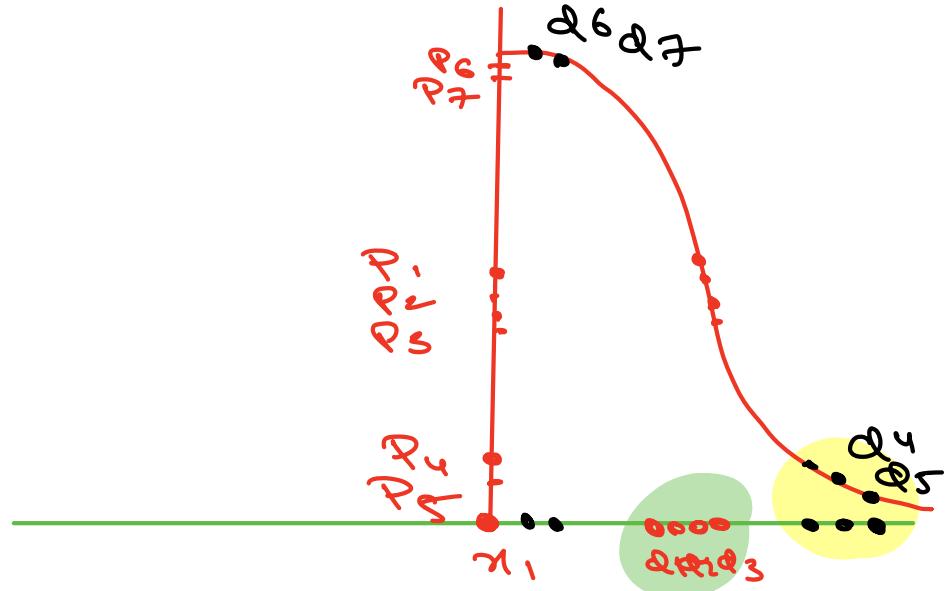
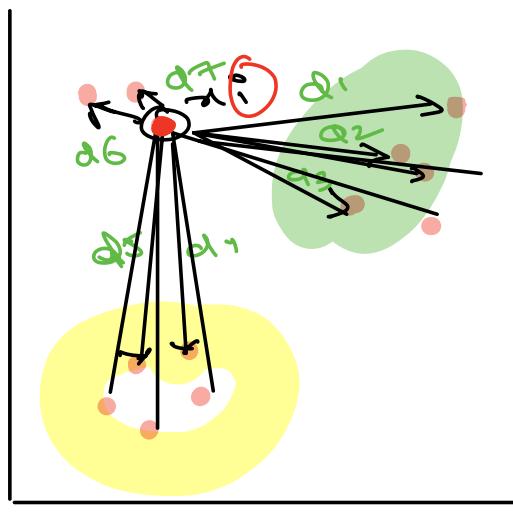
$$P(x_1, x_2)$$

$$\frac{1}{\sqrt{2\pi} \sigma} e^{-\frac{(x_1 - x_2)^2}{2\sigma^2}}$$

$$P(x_1, x_3)$$

$$\frac{1}{\sqrt{2\pi} \sigma} e^{-\frac{(x_1 - x_3)^2}{2\sigma^2}}$$

$$P(x_i, x_j) \propto \frac{1}{\sqrt{2\pi} \sigma} e^{-\frac{\|x_i - x_j\|^2}{2\sigma^2}}$$



Can we Normalize to make
 sum of $P_1, P_2, \dots, P_n \Rightarrow 1$

$$P_{(i,j)} \rightarrow \frac{\cancel{c}}{\sqrt{2\pi}^n} e^{-\frac{\|x_i - x_j\|^2}{2\sigma^2}}$$

$$\sum_{i=1}^n \sum_{j \neq i} \cancel{c} e^{-\frac{\|x_i - x_j\|^2}{2\sigma^2}}$$

Normalize P.D.F

$$P(i,j) \rightarrow$$

$$\frac{e^{-\frac{\|x_i - x_j\|^2}{2\sigma^2}}}{\sum_{i=1}^n \sum_{j \neq i} e^{-\frac{\|x_i - x_j\|^2}{2\sigma^2}}}$$

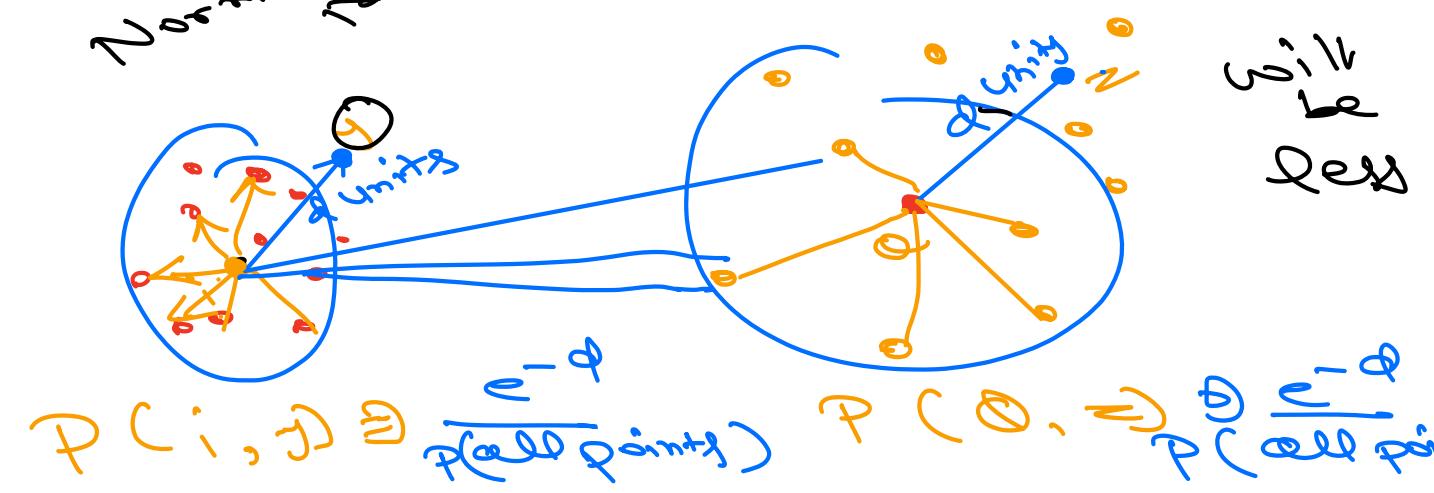
$$6. \quad P(\theta_1, \theta_2, \theta_3, \theta_4, \dots, \theta_n) = \text{std}(\theta)$$

$$P(i, j) \propto \frac{e^{-\frac{\|x_i - x_j\|}{c}}}{\sum_{k \neq i} e^{-\frac{\|x_i - x_k\|^2}{c}}}$$

$$P(i, j) \propto \frac{e^{-\frac{\|x_i - x_j\|}{c}}}{\sum_{k \neq i} e^{-\frac{\|x_i - x_k\|^2}{c}}}$$

for calculating probability
of Neighbours in Original
feature Space

Normalization
done.

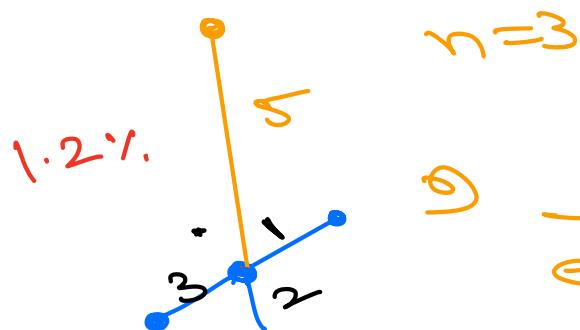


$$P(i,j) \propto \frac{e^{-\frac{\|x_i - x_j\|}{\sigma}}}{\sum_{k \neq i} e^{-\frac{\|x_i - x_k\|}{\sigma}}}$$

is all the point in
Dataset

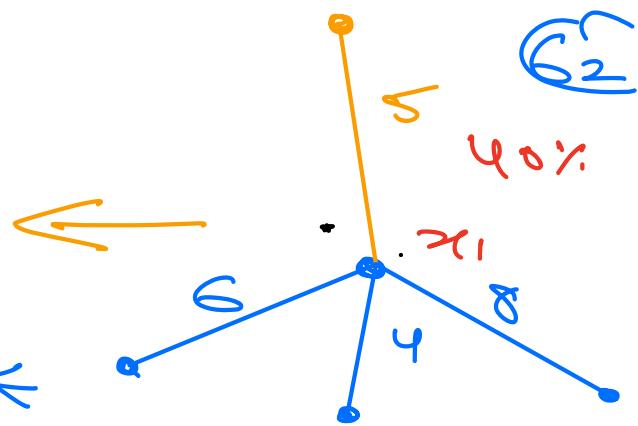
We take only K nearest neighbors for
Normalization

$$P(i,j) \propto \frac{e^{-\frac{\|x_i - x_j\|}{\sigma}}}{\sum_{k=1}^K e^{-\frac{\|x_i - x_k\|}{\sigma}}}$$



$$\text{Value} = \frac{e^{-\gamma}}{e^{-\gamma} + e^{-2} + e^{-3}} = \frac{0.005}{0.36 + 0.13 + 0.04} = 0.012$$

$$\text{Value} = \frac{e^{-4}/e^{-5}}{e^{-4} + e^{-6} + e^{-8}} = \frac{0.006}{0.01 + 0.002 + 0.003} = 0.006$$



④ So for normalize we fix a hyper-parameter n -neighbours to handle different densities

④ Perplexity is a replacement of n -neighbours

n -neighbours \Rightarrow int
Perplexity \Rightarrow float

$$\text{Perplexity} \triangleq 2^H(P_{ij})$$

where

$$H(P_{ij}) = -\sum_j P_{ij} \log_2 P_{ij}$$

Shannon's Entropy

Perplexity ≈ 0.0

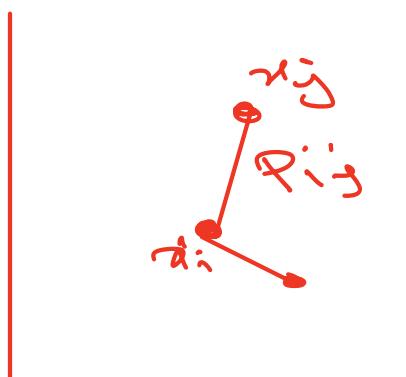
④ Dynamic n -neighbours

Probability of pairs being Neighbours
in Original / High Dimension

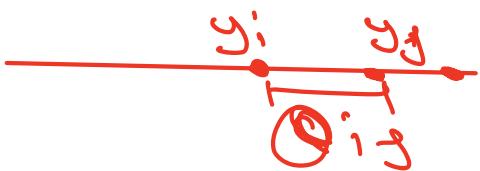
$$P_{ij} \propto \frac{e^{-\|x_i - x_j\|^2}}{\sum_{k=1}^M e^{-\|x_i - x_k\|^2}}$$

Can we use same Distribution
and calculate probability in
Lower / Transformed Dimension?

$$\Theta_{ij} \propto \frac{e^{-\|y_i - y_j\|^2}}{\sum_{k=1}^M e^{-\|y_i - y_k\|^2}}$$

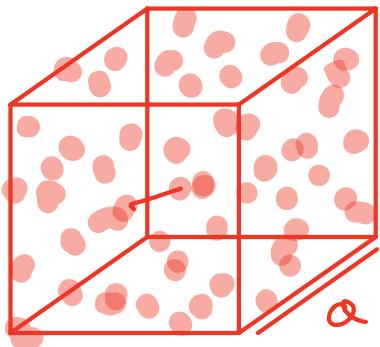


(Y)



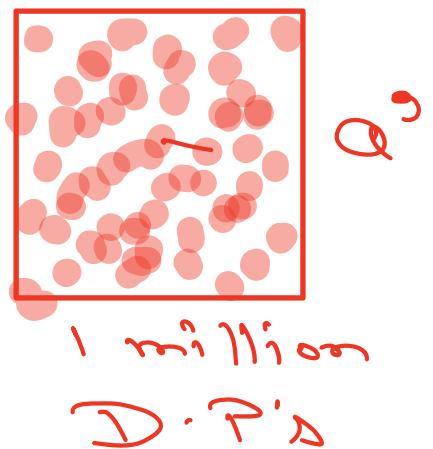
Crowding Problem

3d High



1 million D.P's

Low Dim
Qd



1 million
D.P's

① To avoid Crowding: $Q' > Q$

② $P_{ij} \leq Q_{ij}$

if i use same distribution for
represent both P and Q

$P_{ij} \leq Q_{ij}$ will not be possible

$$P_{ij} \propto \frac{1}{\sum_{j=1}^M e^{-\frac{\|x_i - x_j\|^2}{2\sigma^2}}}$$

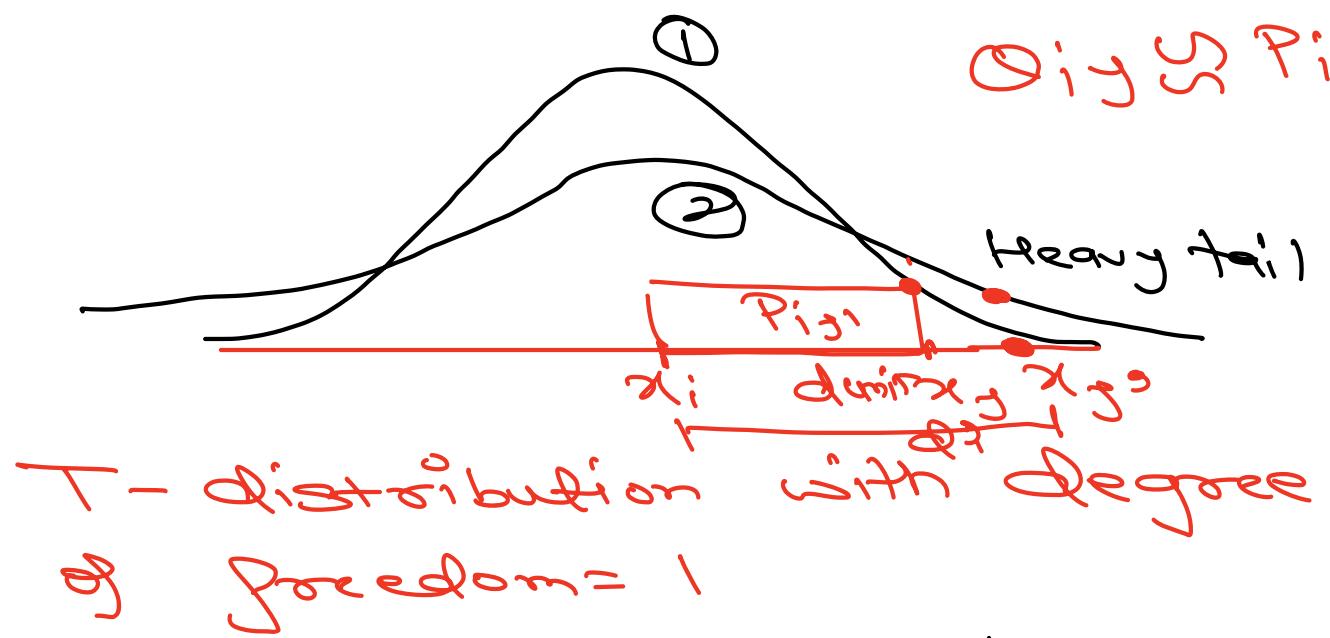
$$\sigma < \sigma'$$

$$Q_{ij} \propto \frac{1}{\sum_{j=1}^M e^{-\frac{\|y_i - y_j\|^2}{2\sigma'^2}}}$$

Hence for Lower Dimension
we can use relaxed Normal
Distribution

$$\sigma^2 > \sigma$$

$$\Omega_{ij} \leq \Omega_{ij}$$



$$dof=1$$

$$dof=10$$

$$dof=30$$

$$\Omega_{ij} \rightarrow \frac{e^{-\frac{\|y_i - y_j\|^2}{2}}}{\sum_{k=1}^M e^{-\frac{\|y_i - y_k\|^2}{2}}}$$

N-Dist

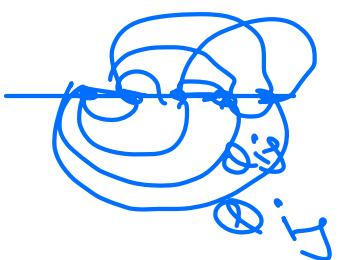
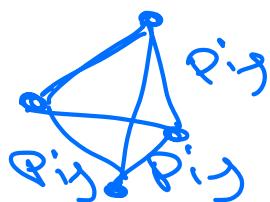
$$\frac{(1 + \|y_i - y_j\|^2)^{-\frac{1}{2}}}{\sum_{k=1}^M (1 + \|y_i - y_k\|^2)^{-\frac{1}{2}}}$$

T-Dist

We have

In high $\Rightarrow P_{ij} \rightarrow$ Normal Dist
dim

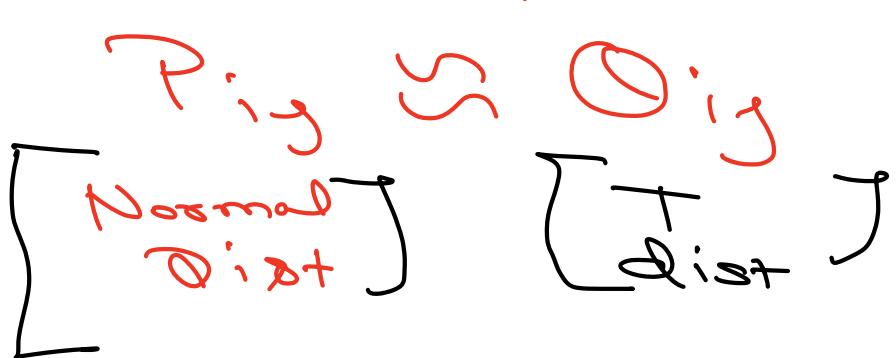
In low $\Rightarrow Q_{ij} \rightarrow T$ -Distribution
dim (Cauchy Dist)



How do we compare all the probabilities



How do we compare and say
that $P_{ij} \ll Q_{ij}$



A new Loss function is
required

③ We use KL Divergence

$$KL_{div}(P, Q) \triangleq \sum_{i=1}^M \sum_{j=1}^N P_{ij} \log \left(\frac{P_{ij}}{Q_{ij}} \right)$$

④ we use KL_{div} as Loss function

Case 1:

i and j overlap

$$KL(P, Q) \rightarrow 0 \quad \text{as } P_{ij} \rightarrow 0$$

Case 2:

$$P_{ij} > Q_{ij}$$

$$KL(P, Q) \triangleq \log \frac{P_{ij}}{Q_{ij}} \rightarrow > 1$$

$$\rightarrow (\log \cancel{Q}) \times P_{ij}$$

Case 3:

$$P_{ij} = Q_{ij}$$

$$\cancel{\log(P_{ij})} \times P_{ij} \triangleq P_{ij} \log(1) = 0$$

④ To solve and find optimal axis we can use KL-Div as loss function and apply Gradient Descent.

Drawback:

- ① T-SNE is very time intensive
- ② Hyperparameter: Perplexity Tuning

T-SNE is generally used for visualization

⇒ Hence n-components $\rightarrow 2, 3$

③ Umap and Umap

$$Q_{ij} = \frac{e^{-d_i^2}}{\sum_j e^{-d_j^2}}$$



$$P_{ij} \geq 0$$

$$Q_{ij} = P_{ij} \Rightarrow \text{Score}$$