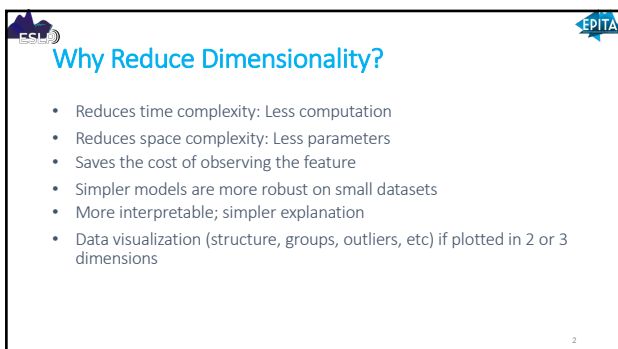
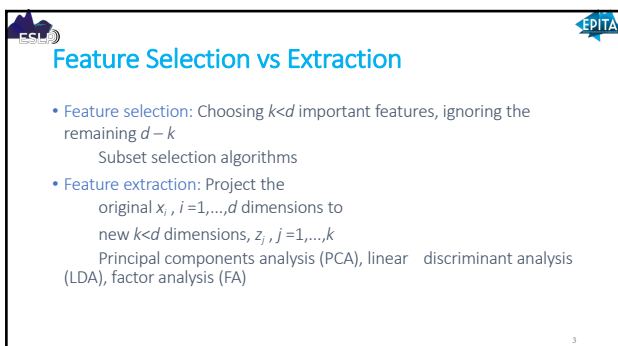




1



2



3






Subset Selection

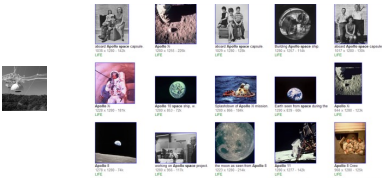
- There are 2^d subsets of d features
- Forward search: Add the best feature at each step
 - Set of features F initially \emptyset .
 - At each iteration, find the best new feature
 $j = \operatorname{argmin}_j E(F \cup x_j)$
 - Add x_j to F if $E(F \cup x_j) < E(F)$
- Hill-climbing $O(d^2)$ algorithm
- Backward search: Start with all features and remove one at a time, if possible.
- Floating search (Add k , remove l)

4

4



Recall: Representing images



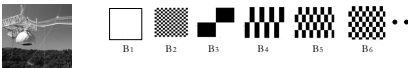
- The most common element in the image: background
 - Or rather large regions of relatively featureless shading
 - Uniform sequences of numbers

5

5

Adding more bases



- Checkerboards with different variations

$$\text{Image} \approx w_1 B_1 + w_2 B_2 + w_3 B_3 + \dots$$

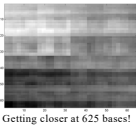
$$W = \begin{bmatrix} w_1 \\ w_2 \\ w_3 \\ \vdots \end{bmatrix}$$

$$B = [B_1 \ B_2 \ B_3]$$

$$BW \approx \text{Image}$$

$$W = \operatorname{pinv}(B) \text{ Image}$$

$$\text{PROJECTION} = BW$$



6

6

"Bases"

$image \approx w_1 B_1 + w_2 B_2 + w_3 B_3 + \dots$

- "Bases" are the "standard" units such that all instances can be expressed a weighted combinations of these units
- Ideal requirements: Bases must be orthogonal
 - Checkerboards are one choice of bases
 - Orthogonal
 - But not "smooth"
- Other choices of bases: Complex exponentials, Wavelets, etc...

7

Data specific bases?

- **Issue: The bases we have considered so far are data agnostic**
 - Checkerboards, Complex exponentials, Wavelets...
 - We use the same bases regardless of the data we analyze
 - Image of face vs. Image of a forest
 - Segment of speech vs. Seismic rumble
- How about data specific bases
 - Bases that consider the underlying data
 - E.g. is there something better than checkerboards to describe faces
 - Something better than complex exponentials to describe music?

8

Principal Components Analysis (PCA)

- Find a low-dimensional space such that when \mathbf{x} is projected there, information loss is minimized.
- The projection of \mathbf{x} on the direction of \mathbf{w} is: $z = \mathbf{w}^T \mathbf{x}$
- Find \mathbf{w} such that $\text{Var}(z)$ is maximized

$$\begin{aligned} \text{Var}(z) &= \text{Var}(\mathbf{w}^T \mathbf{x}) = E[(\mathbf{w}^T \mathbf{x} - \mathbf{w}^T \boldsymbol{\mu})^2] \\ &= E[(\mathbf{w}^T \mathbf{x} - \mathbf{w}^T \boldsymbol{\mu})(\mathbf{w}^T \mathbf{x} - \mathbf{w}^T \boldsymbol{\mu})] \\ &= E[\mathbf{w}^T (\mathbf{x} - \boldsymbol{\mu})(\mathbf{x} - \boldsymbol{\mu})^T \mathbf{w}] \\ &= \mathbf{w}^T E[(\mathbf{x} - \boldsymbol{\mu})(\mathbf{x} - \boldsymbol{\mu})^T] \mathbf{w} = \mathbf{w}^T \boldsymbol{\Sigma} \mathbf{w} \end{aligned}$$
 where $\text{Var}(\mathbf{x}) = E[(\mathbf{x} - \boldsymbol{\mu})(\mathbf{x} - \boldsymbol{\mu})^T] = \boldsymbol{\Sigma}$

9

Principal Components Analysis (PCA)

- Maximize $\text{Var}(z)$ subject to $\|w\|=1$

$$\max_{w_1} w_1^T \Sigma w_1 - \alpha (w_1^T w_1 - 1)$$

$\Sigma w_1 = \alpha w_1$ that is, w_1 is an eigenvector of Σ
Choose the one with the largest eigenvalue for $\text{Var}(z)$ to be max

- Second principal component: Max $\text{Var}(z)$, s.t., $\|w_2\|=1$ and orthogonal to w_1

$$\max_{w_2} w_2^T \Sigma w_2 - \alpha (w_2^T w_2 - 1) - \beta (w_2^T w_1 - 0)$$

$\Sigma w_2 = \alpha w_2$ that is, w_2 is another eigenvector of Σ and so on.

10

What PCA does

$$z = W^T(x - m)$$

where the columns of W are the eigenvectors of Σ and m is sample mean

Centers the data at the origin and rotates the axes

11

How to choose k ?

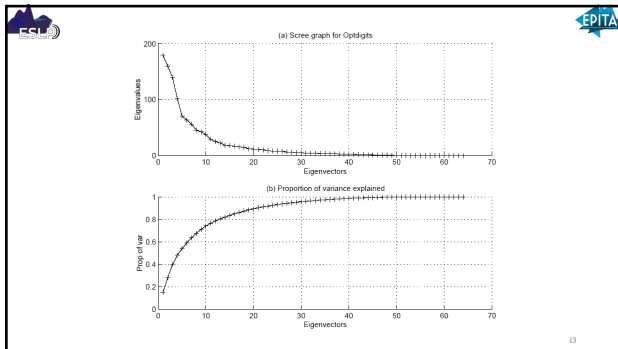
- Proportion of Variance (PoV) explained

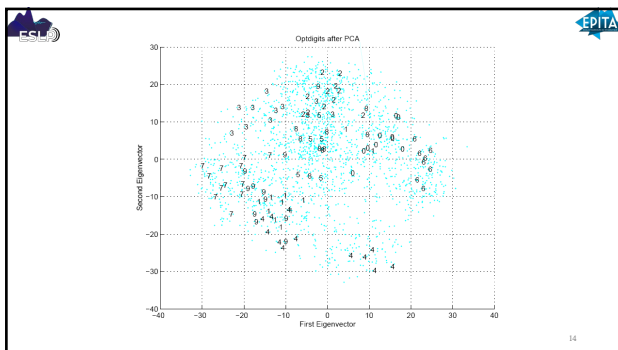
$$\frac{\lambda_1 + \lambda_2 + \dots + \lambda_k}{\lambda_1 + \lambda_2 + \dots + \lambda_k + \dots + \lambda_d}$$

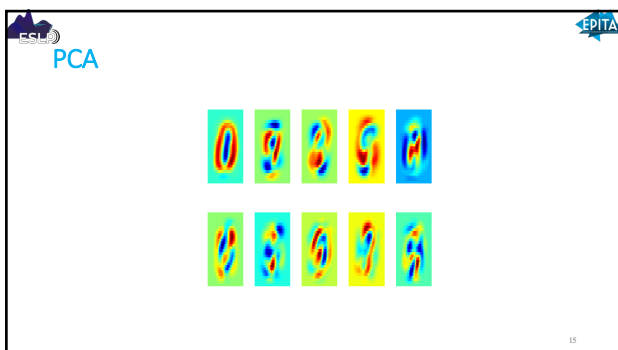
when λ_i are sorted in descending order

- Typically, stop at $\text{PoV} > 0.9$
- Scree graph plots of PoV vs k , stop at "elbow"

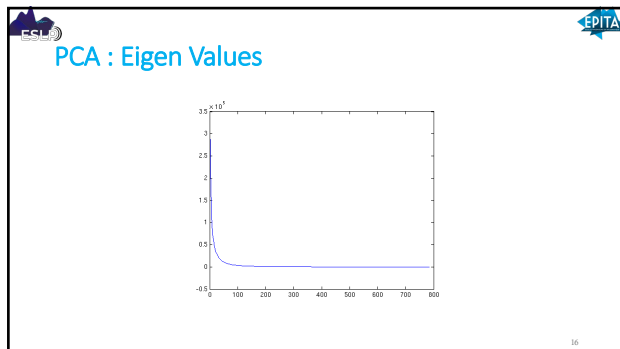
12

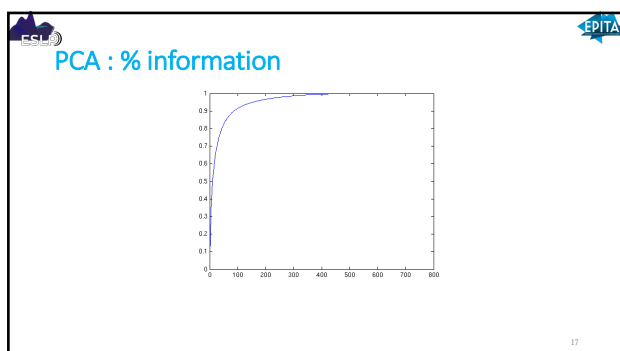


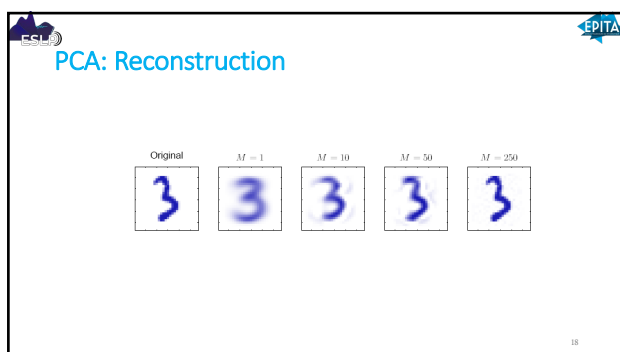




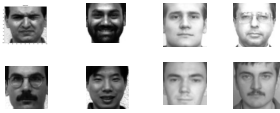
15







Data-specific description of faces




- A collection of images
 - All normalized to 100x100 pixels
- What is common among all of them?
 - Do we have a common descriptor?

19

19

Eigen Faces!




- Arrange your input data into a matrix X
- Compute the correlation $R = XX^T$
- Solve the Eigen decomposition: $RV = \Lambda V$
- The Eigen vectors corresponding to the K largest eigen values are our optimal bases
- We will refer to these as *eigen faces*.

20

20

Class specificity

- Eigen bases are class specific



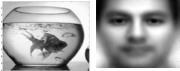
- Composing a fishbowl from Eigenfaces

21

21

Class specificity

- Eigen bases are class specific




- Composing a fishbowl from Eigenfaces
- With 1 basis

$$f = w_1 \mathbf{v}_1$$

22

Class specificity

- Eigen bases are class specific



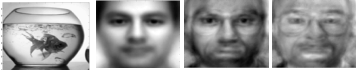
- Composing a fishbowl from Eigenfaces
- With 10 bases

$$f = w_1 \mathbf{v}_1 + w_2 \mathbf{v}_2 + \dots + w_{10} \mathbf{v}_{10}$$

23

Class specificity

- Eigen bases are class specific




- Composing a fishbowl from Eigenfaces
- With 30 bases

$$f = w_1 \mathbf{v}_1 + w_2 \mathbf{v}_2 + \dots + w_{10} \mathbf{v}_{10} + \dots + w_{30} \mathbf{v}_{30}$$

24

Class specificity

- Eigen bases are class specific



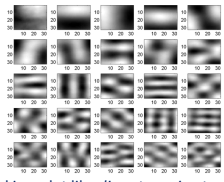
- Composing a fishbowl from Eigenfaces
- With 100 bases

$$f = w_1 \mathbf{v}_1 + w_2 \mathbf{v}_2 + \dots + w_{10} \mathbf{v}_{10} + \dots + w_{30} \mathbf{v}_{30} + \dots + w_{100} \mathbf{v}_{100}$$

25

Universal bases

- Universal bases..



- End up looking a lot like *discrete cosine transforms!!!!*
- DCTs are the best "universal" bases
 - If you don't know what your data are, use the DCT

26

Factor Analysis

- Find a small number of **factors \mathbf{z}** , which when combined generate \mathbf{x} :

$$X_i - \mu_i = V_{i1}Z_1 + V_{i2}Z_2 + \dots + V_{ik}Z_k + \epsilon_i$$

where $z_j, j=1, \dots, k$ are the **latent factors** with

$$E[Z_j] = 0, \text{Var}(Z_j) = 1, \text{Cov}(Z_i, Z_j) = 0, i \neq j,$$

ϵ_i are the **noise sources**

$$E[\epsilon_i] = \psi_i, \text{Cov}(\epsilon_i, \epsilon_j) = 0, i \neq j, \text{Cov}(\epsilon_i, Z_j) = 0,$$

and v_{ij} are the **factor loadings**

27

PCA vs FA

- PCA From \mathbf{x} to \mathbf{z} $\mathbf{z} = \mathbf{W}^T(\mathbf{x} - \boldsymbol{\mu})$
- FA From \mathbf{z} to \mathbf{x} $\mathbf{x} - \boldsymbol{\mu} = \mathbf{V}\mathbf{z} + \boldsymbol{\varepsilon}$

Diagram illustrating the relationship between variables \mathbf{x} and factors \mathbf{z} in PCA and FA. In PCA, variables x_1, x_2, x_3 are transformed into new variables z_1, z_2, z_3 using matrix \mathbf{W} . In FA, factors z_1, z_2, z_3 are transformed into variables x_1, x_2, x_3 using matrix \mathbf{V} . A curved arrow labeled \mathbf{W} indicates the transformation from factors \mathbf{z} back to variables \mathbf{x} .

28

Factor Analysis

- In FA, factors \mathbf{z} are stretched, rotated and translated to generate \mathbf{x}

Diagram illustrating Factor Analysis. The left plot shows factors z_1 and z_2 with a circle centered at the origin. The right plot shows variables x_1 and x_2 with an ellipse centered at the origin, representing the transformed factors.

29

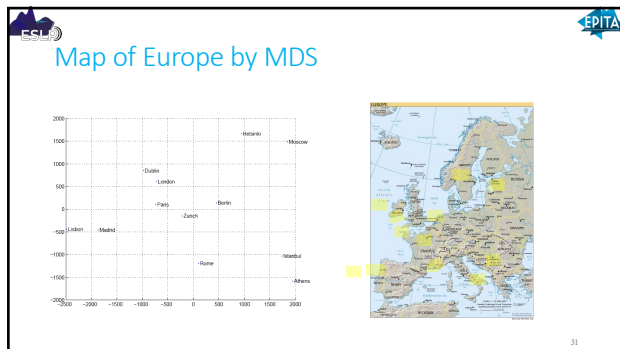
Multidimensional Scaling

- Given pairwise distances between N points, $d_{ij}, i, j = 1, \dots, N$, place on a low-dim map s.t. distances are preserved.
- $\mathbf{z} = \mathbf{g}(\mathbf{x} | \boldsymbol{\theta})$ Find $\boldsymbol{\theta}$ that min **Sammon stress**

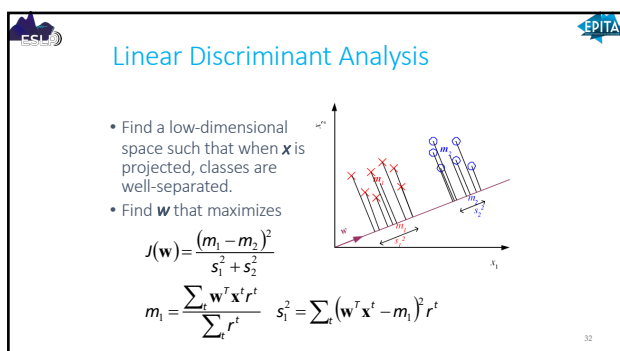
$$\varepsilon(\boldsymbol{\theta} | \mathbf{X}) = \sum_{i,j} \frac{(\|\mathbf{z}'_i - \mathbf{z}'_j\| - \|\mathbf{x}'_i - \mathbf{x}'_j\|)^2}{\|\mathbf{x}'_i - \mathbf{x}'_j\|^2}$$

$$= \sum_{i,j} \frac{(\|\mathbf{g}(\mathbf{x}'_i | \boldsymbol{\theta}) - \mathbf{g}(\mathbf{x}'_j | \boldsymbol{\theta})\| - \|\mathbf{x}'_i - \mathbf{x}'_j\|)^2}{\|\mathbf{x}'_i - \mathbf{x}'_j\|^2}$$

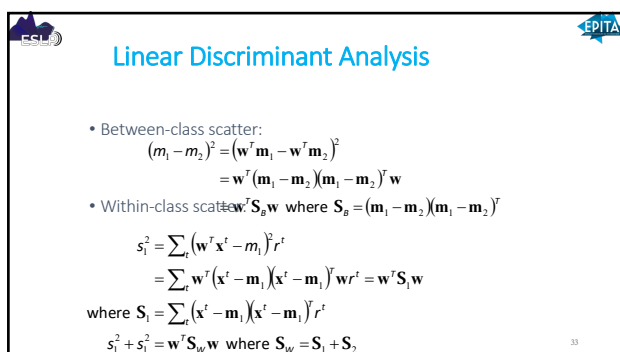
30



31



32



33

Fisher's Linear Discriminant

- Find \mathbf{w} that max

$$J(\mathbf{w}) = \frac{\mathbf{w}^T \mathbf{S}_b \mathbf{w}}{\mathbf{w}^T \mathbf{S}_w \mathbf{w}} = \frac{|\mathbf{w}^T (\mathbf{m}_1 - \mathbf{m}_2)|^2}{\mathbf{w}^T \mathbf{S}_w \mathbf{w}}$$

- LDA soln:

$$\mathbf{w} = \mathbf{c} \cdot \mathbf{S}_w^{-1} (\mathbf{m}_1 - \mathbf{m}_2)$$

- Parametric soln:

$$\mathbf{w} = \Sigma^{-1} (\mu_1 - \mu_2)$$

when $p(\mathbf{x} | C_i) \sim \mathcal{N}(\mu_i, \Sigma)$

34

K>2 Classes

- Within-class scatter:

$$\mathbf{S}_w = \sum_{i=1}^K \mathbf{S}_i \quad \mathbf{S}_i = \sum_t r_i^t (\mathbf{x}^t - \mathbf{m}_i)(\mathbf{x}^t - \mathbf{m}_i)^T$$

- Between-class scatter:

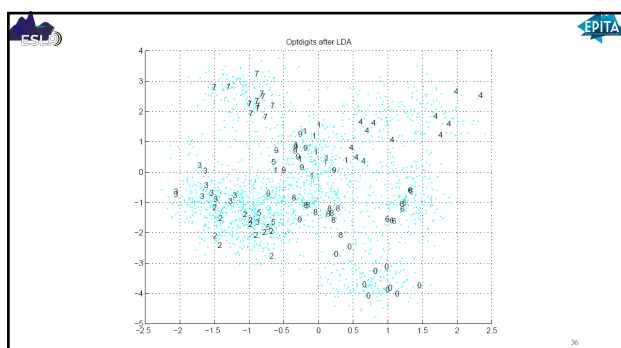
$$\mathbf{S}_b = \sum_{i=1}^K N_i (\mathbf{m}_i - \mathbf{m})(\mathbf{m}_i - \mathbf{m})^T \quad \mathbf{m} = \frac{1}{K} \sum_{i=1}^K \mathbf{m}_i$$

- Find \mathbf{W} that max

$$J(\mathbf{W}) = \frac{|\mathbf{W}^T \mathbf{S}_b \mathbf{W}|}{|\mathbf{W}^T \mathbf{S}_w \mathbf{W}|}$$

The largest eigenvectors of $\mathbf{S}_w^{-1} \mathbf{S}_b$
Maximum rank of $K-1$

35



Neighborhood Components Analysis

- NCA learns a Mahalanobis distance metric for the KNN classifier by maximizing the leave-one-out cross validation.
 - Let p_{ij} be the probability that point x_j is selected as point x_i 's neighbour.
 - The probability that points are correctly classified when x_i is used as the reference is:

$$P_i = \sum_{j \in C_i} p_{ij}$$

$$C_i = \{x_j | \text{class}(x_j) = \text{class}(x_i)\}$$
 - p_{ij} is related to the distance between x_i and x_j

$$p_{ij} = \frac{e^{-d_{ij}}}{\sum_{k \neq i} e^{-d_{ik}}}, p_{ii} = 0$$
- The expected number of correctly classification points:

$$f(A) = \sum_i P_i$$

37

Neighborhood Components Analysis

- How do we define d_{ij} ?
- Limit the distance measure within Mahalanobis (quadratic) distance:

$$d(x, y) = (x - y)^T Q (x - y)$$

$$Q = A^T A$$

$$d(x, y) = (Ax - Ay)^T (Ax - Ay)$$
- That is to say, we project the original feature vectors x into another vector space with transformation matrix A .

Distance Learning \Leftrightarrow Finding the Best matrix A

38

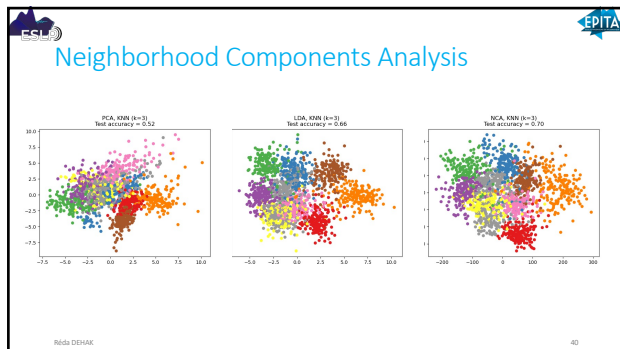
Neighborhood Components Analysis

- Substitute d_{ij} in p_{ij}

$$p_{ij} = \frac{\exp(-\|Ax_i - Ax_j\|^2)}{\sum_{k \neq i} \exp(-\|Ax_i - Ax_k\|^2)}$$
- The objective function:

$$f(A) = \sum_i P_i = \sum_i \sum_{j \in C_i} p_{ij}$$

39



40

t-SNE

- SNE - Stochastic Neighbor Embedding
- t-SNE
- KL divergence

Réda DENAE 41



41

Local Embedding

- t-SNE is an alternative dimensionality reduction algorithm.
- PCA tries to find a **global** structure
 - Low dimensional subspace
 - Can lead to local inconsistencies
 - Far away point can become nearest neighbors
- t-SNE tries to preserve **local** structure
 - Low dimensional neighborhood should be the same as original neighborhood.
- Unlike PCA almost only used for visualization
 - No easy way to embed new points

Réda DENAE 42



42

Stochastic Neighbor Embedding (SNE)

- "Encode" high dimensional neighborhood information as a distribution
- Intuition: Random walk between data points.
 - High probability to jump to a close point
- Find low dimensional points such that their neighborhood distribution is similar
- How do you measure distance between distributions?
 - Most common measure: KL divergence

43






Neighborhood Distributions

- Consider the neighborhood around an input data point $x_i \in \mathbb{R}^d$
- Imagine that we have a Gaussian distribution centered around x_i
- Then the probability that x_i chooses some other datapoint x_j as its neighbor is in proportion with the density under this Gaussian
- A point closer to x_i will be more likely than one further away
- $P_{j|i}$ probability (similar to NCA), is the probability that point x_i chooses x_j as its neighbor:

$$P_{j|i} = \frac{\exp(-\|x_i - x_j\|^2 / 2\sigma_i^2)}{\sum_{k \neq i} \exp(-\|x_i - x_k\|^2 / 2\sigma_i^2)}$$

44

Neighborhood Distributions



- $P_{j|i}$ probability (similar to NCA), is the probability that point x_i chooses x_j as its neighbor:

$$P_{j|i} = \frac{\exp(-\|x_i - x_j\|^2 / 2\sigma_i^2)}{\sum_{k \neq i} \exp(-\|x_i - x_k\|^2 / 2\sigma_i^2)}$$

With $P_{i|i} = 0$

- The parameter σ_i sets the size of the neighborhood
 - Very low σ_i : all the probability is in the nearest neighbor.
 - Very high σ_i : Uniform weights.
- We set σ_i differently for each data point
- Results depend heavily on σ_i :
 - it defines the neighborhoods we are trying to preserve.
- Final distribution over pairs is symmetrized: $P_{ij} = \frac{1}{2N} (P_{i|j} + P_{j|i})$
- Random Walk:
 - Pick i uniformly and then "jump" to j according to $P_{j|i}$

45



SNE objective

- Given $x_1, x_2, \dots, x_N \in \mathbb{R}^d$, we define the distribution P_{ij}
- Goal: Find good embedding $y_1, y_2, \dots, y_N \in \mathbb{R}^p$ for some $p \ll d$ (normally 2 or 3 for visualization)
- How do we measure an embedding quality?
- For points $y_1, y_2, \dots, y_N \in \mathbb{R}^p$, we can define distribution Q similarly the same (notice no σ_i^2 and not symmetric)

$$Q_{ij} = \frac{\exp(-\|y_i - y_j\|^2)}{\sum_{k \neq i} \exp(-\|y_i - y_k\|^2)}$$

- Optimize Q to be close to P
 - Minimize KL-divergence
- The embeddings $y_1, y_2, \dots, y_N \in \mathbb{R}^p$ are the parameters we are optimizing.
 - How do you embed a new point? **No embedding function!**

46



SNE objective

- Given $x_1, x_2, \dots, x_N \in \mathbb{R}^d$, we define the distribution P_{ij}
- Goal: Find good embedding $y_1, y_2, \dots, y_N \in \mathbb{R}^p$ for some $p \ll d$ (normally 2 or 3 for visualization)
- How do we measure an embedding quality?
- For points $y_1, y_2, \dots, y_N \in \mathbb{R}^p$, we can define distribution Q similarly the same (notice no σ_i^2 and not symmetric)

$$Q_{ij} = \frac{\exp(-\|y_i - y_j\|^2)}{\sum_{k \neq i} \exp(-\|y_i - y_k\|^2)}$$

- Optimize Q to be close to P
 - Minimize KL-divergence
- The embeddings $y_1, y_2, \dots, y_N \in \mathbb{R}^p$ are the parameters we are optimizing.
 - How do you embed a new point? **No embedding function!**

47



SNE algorithm

- We have P , and are looking for $y_1, y_2, \dots, y_N \in \mathbb{R}^p$ such that the distribution Q we infer will minimize $\mathbb{C}(Q) = \text{KL}(P\|Q)$ (notice Q on right)

$$\text{KL}(P\|Q) = \sum_{ij} P_{ij} \log\left(\frac{P_{ij}}{Q_{ij}}\right) = -\sum_{ij} P_{ij} \log(Q_{ij}) + \text{const}$$

- Not a convex problem! can use multiple restarts.
- Main issue: **crowding problem**



48

Crowding Problem: Neighborhood

- In high dimension we have more room, points can have a lot of different neighbors
- In 2D a point can have a few neighbors at distance one all far from each other
 - what happens when we embed in 1D?
- **Crowding problem:**
 - We don't have enough room to accommodate all neighbors.
 - One of the biggest problems with SNE.
- t-SNE solution: Change the Gaussian in Q to a heavy tailed distribution.
- if Q changes slower, we have more wiggle room to place points at.

49

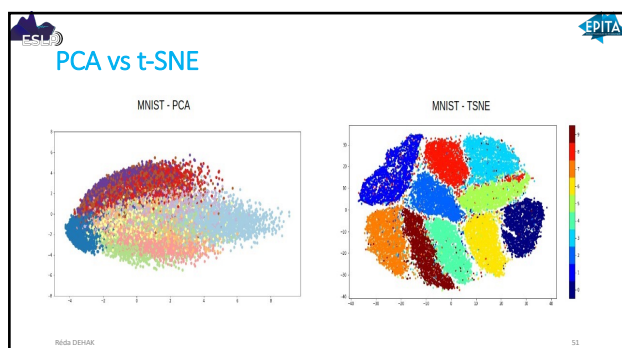



T-SNE: t-Distributed Stochastic Neighbor Embedding



- Student-t Probability density $p(x) \propto \left(1 + \frac{x^2}{v}\right)^{\frac{v+1}{2}}$
 - For $v = 1$ we get $p(x) \propto \frac{1}{1+x^2}$
- Probability goes to zero much slower than a Gaussian
- We can now redefine Q_{ij} as

$$Q_{ij} = \frac{(1 + \|y_i - y_j\|^2)^{-1}}{\sum_k \sum_{l \neq k} (1 + \|y_k - y_l\|^2)^{-1}}$$
- We leave P_{ij} as is

50



51



Conclusions

- A lot of tools can be used for dimensionality reduction linear and not linear.
- Helpful for visualization
- Helps to understand "black-box" algorithms
- t-SNE is a great way to visualize data

Réda DENAE52
