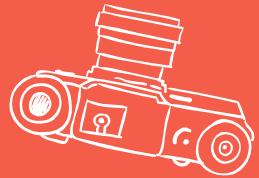


# JSC 270 - LECTURE 9

## EMBEDDING

<https://jsc270.github.io/>



# ANNOUNCEMENTS

- Perusall paper is posted
- Reflection quiz will be graded by the end of today
- Next week lecture: data science in healthcare  
talk by Devin Singh, MD, MSc (SickKids/HeroAI)
- Paper I mentioned last class by Zack Lipton on NLP:  
<https://arxiv.org/pdf/1909.12434.pdf>

# DIMENSIONALITY REDUCTION

How would you define it?

# DIMENSIONALITY REDUCTION

*Dimensionality Reduction* is the transformation of data from a *high-dimensional space* into a *low-dimensional space* so that the low-dimensional representation retains some meaningful properties of the original data

Wikipedia

# DIMENSIONALITY REDUCTION

What applications can you think of?

# DIMENSIONALITY REDUCTION

In any AI and stats applications where

- a) There is a reason to believe that the data can be well captured by fewer dimensions
- b) There is a need to reduce data, for e.g. for computational reasons

Primary Uses:

- Exploratory data analysis
- Visualization

Applications: ubiquitous across fields (healthcare, industry; text, speech, imaging, etc)

# DIMENSIONALITY REDUCTION

Which techniques do you know?

# DIMENSIONALITY REDUCTION

- Matrix factorization based
  - Principle Component Analysis (PCA)
    - Graph-based PCA
    - Kernel-based PCA
  - Non-negative Matrix Factorization (NMF)
  - Linear Discriminant Analysis (LDA)
  - Generalized Discriminant Analysis (GDA)
  - Random projections
- Neighborhood based
  - t-SNE
  - UMAP
  - Locally Linear Embedding
  - ISOMAP
- Autoencoders (reconstruction based)

And many more

[https://en.wikipedia.org/wiki/Nonlinear\\_dimensionality\\_reduction](https://en.wikipedia.org/wiki/Nonlinear_dimensionality_reduction)

## DIMENSIONALITY REDUCTION VS EMBEDDING

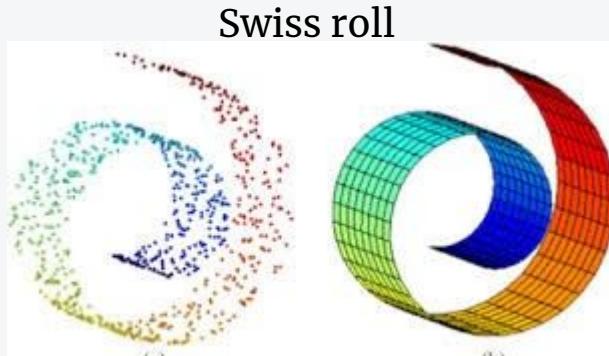
High-dimensional data often lies on or near a much lower dimensional manifold.

We reduce dimension or *embed* data into the manifold to get the *embedding*, in other words a low dimensional representation



## USEFUL CONCEPTS: MANIFOLD

**Manifold** - a collection of points forming a certain kind of set, such as those of a topologically closed surface in three or more dimensions (such topological space that locally resembles Euclidean space near each point)



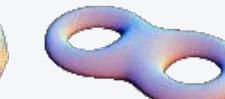
*sphere*



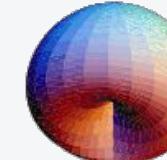
*torus*



*double torus*



*cross surface*

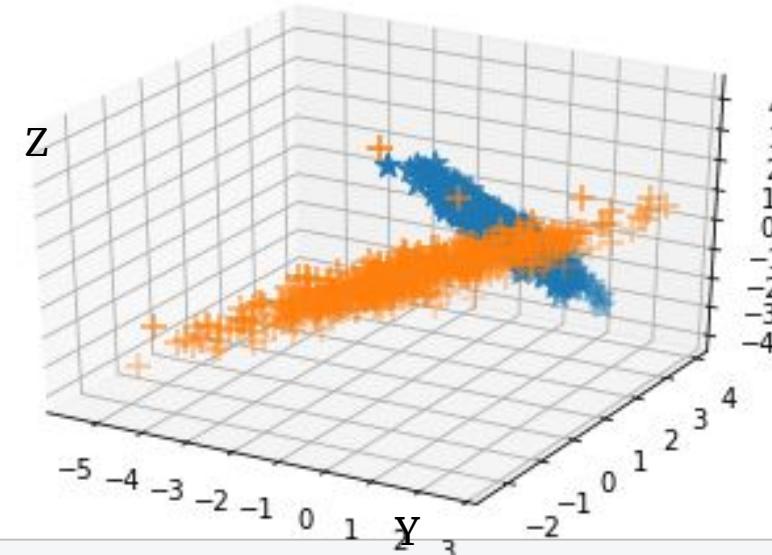


*Klein bottle*



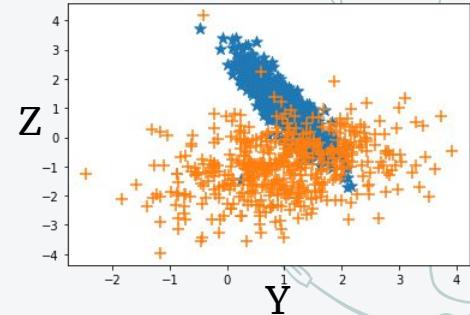
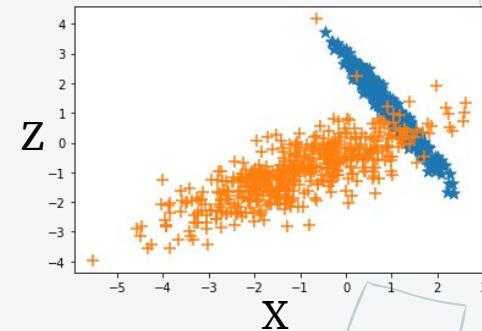
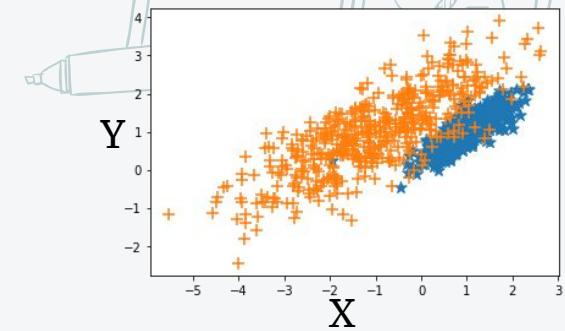
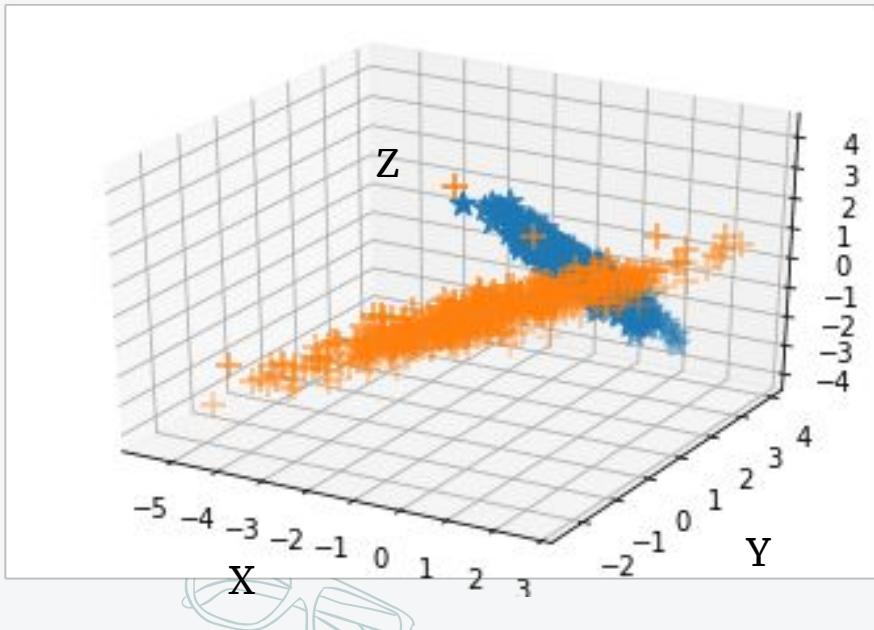
# DATA

Number of samples  $N = 1000$   
Number of features  $P = 3$



# DATA

Number of samples  $N = 1000$   
Number of features  $P = 3$



# PRINCIPAL COMPONENT ANALYSIS

Principle components = axes of greatest variability

When we talk about  $X_1, X_2, \dots, X_p$  - these features are in cartesian coordinate system

Principle components are a *new* set of  $p$  axes in the direction of greatest variability, that are a linear combination of the original

# SELECTING PRINCIPAL COMPONENTS

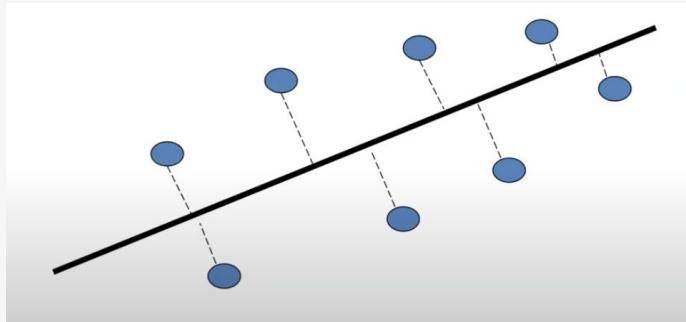
In one dimension, we would projecting points to a line.

How do we do that?

# SELECTING PRINCIPAL COMPONENTS

In one dimension, we would project points to a line.

How do we do that? Regression (squared loss)



## SELECTING PRINCIPAL COMPONENTS

In one dimension, we would projecting points to a line.

How do we do that? Regression (squared loss)

How do we do it in many dimensions?

# LINEAR ALGEBRA (REVIEW)

- Eigenvectors for a square  $m \times m$  matrix  $\mathbf{S}$

$$\mathbf{S}\mathbf{v} = \lambda\mathbf{v}$$

Eigenvector  $\mathbf{v} \in \mathbb{R}^m, \mathbf{v} \neq \mathbf{0}$

Eigenvalue  $\lambda \in \mathbb{R}$

E.g

$$\begin{pmatrix} 6 & -2 \\ 4 & 0 \end{pmatrix} \begin{pmatrix} 1 \\ 2 \end{pmatrix} = \begin{pmatrix} 2 \\ 4 \end{pmatrix} = 2 \begin{pmatrix} 1 \\ 2 \end{pmatrix}$$

$\uparrow$                              $\uparrow$   
eigenvector                    eigenvalue

# LINEAR ALGEBRA (REVIEW)

How many eigenvalues are there at most?

## LINEAR ALGEBRA (REVIEW)

How many eigenvalues are there at most?

$$S\mathbf{v} = \lambda\mathbf{v} \equiv (S - \lambda I)\mathbf{v} = 0$$

only has non-zero solutions if  $|S - \lambda I| = 0$

Can have at most  $m$  distinct solutions

If  $S = S^T \Rightarrow \lambda \in \mathbb{R}$

## LINEAR ALGEBRA (REVIEW)

For *symmetric matrices*, eigenvectors for distinct eigenvalues are *orthogonal*

All eigenvalues for a *positive semidefinite matrix* are *non-negative*

Note: S is positive semidefinite if  $\forall w \in \mathbb{R}^m, w^T S w \geq 0$

# EIGEN DECOMPOSITION

Eigen decomposition transforms a matrix into its principal vectors of variation

$$S = \begin{matrix} u_1, u_2, u_3, \dots, u_m \\ U \\ \lambda_1 \\ \lambda_2 \\ \lambda_m \\ U^{-1} \end{matrix}$$

$$S = U \Lambda U^{-1}$$

U are eigenvectors of S, diagonal values of Lambda are eigenvalues.

# EIGEN DECOMPOSITION

Eigen decomposition transforms a matrix into its principal vectors of variation

$$S = \begin{matrix} u_1, u_2, u_3, \dots, u_m \\ U \end{matrix} \begin{matrix} \lambda_1 & & & \\ & \ddots & & \\ & & \lambda_2 & \\ & & & \lambda_m \end{matrix} \begin{matrix} & O & \\ & & O \end{matrix} U^{-1}$$

Decomposition is unique if eigenvalues are unique

$$S = U \Lambda U^{-1}$$

Columns of U are **eigenvectors** of S

Diagonal values of Lambda are **eigenvalues** of S

# SINGULAR VALUE DECOMPOSITION (SVD)

For an  $m \times n$  matrix  $A$  of rank  $r$  there exists a factorization (SVD) as follows:

$$A = U \Sigma V^T$$

$m \times m$     $m \times n$     $n \times n$

The columns of  $U$  are orthogonal eigenvectors of  $AA^T$

The columns of  $V$  are orthogonal eigenvectors of  $A^TA$

Eigenvalues  $(\sigma_1, \sigma_2, \dots, \sigma_r)$  of  $AA^T$  are eigenvalues of  $A^TA$

$$\Sigma = \begin{matrix} & \sigma_1 & & 0 \\ & \sigma_2 & & \\ 0 & & & \\ & & & \sigma_r \end{matrix}$$

Note: rank  $r$  means there are max of  $r$  linearly independent rows and columns in  $A$

# SINGULAR VALUE DECOMPOSITION (SVD)

For an  $m \times n$  matrix  $A$  of rank  $r$  there exists a factorization (SVD) as follows:

$$A = U\Sigma V^T$$

$m \times m$     $m \times n$     $n \times n$

The columns of  $U$  are orthogonal eigenvectors of  $AA^T$

The columns of  $V$  are orthogonal eigenvectors of  $A^TA$

Eigenvalues  $(\sigma_1, \sigma_2, \dots, \sigma_r)$  of  $AA^T$  are eigenvalues of  $A^TA$

$$\Sigma = \begin{matrix} & \sigma_1 & & 0 \\ & \sigma_2 & & \\ 0 & & & \\ & & & \sigma_r \end{matrix}$$

←  
Singular  
Values

Note: rank  $r$  means there are max of  $r$  linearly independent rows and columns in  $A$

# SINGULAR VALUE DECOMPOSITION (SVD)

For an  $m \times n$  matrix  $A$  of rank  $r$  there exists a factorization (SVD) as follows:

$$A = U \Sigma V^T$$

$m \times m$     $m \times n$     $n \times n$

The columns of  $U$  are orthogonal eigenvectors of  $AA^T$

The columns of  $V$  are orthogonal eigenvectors of  $A^TA$

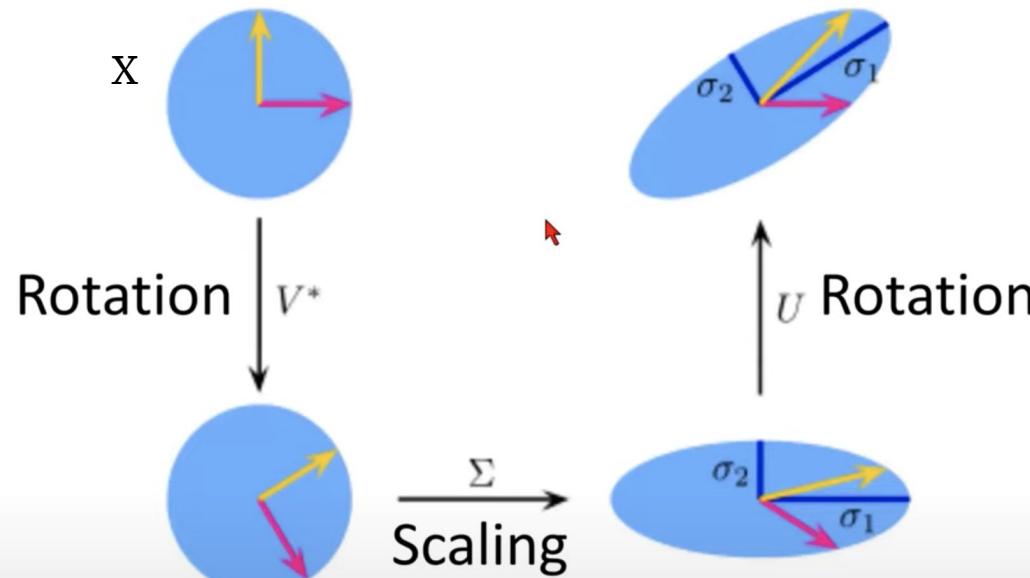
Eigenvalues  $(\sigma_1, \sigma_2, \dots, \sigma_r)$  of  $AA^T$  are eigenvalues of  $A^TA$

$$\Sigma = \begin{matrix} & \sigma_1 & & 0 \\ & \sigma_2 & & \\ 0 & & & \sigma_r \end{matrix}$$

$\lambda_i = \sqrt{\sigma_i}$

Note: rank  $r$  means there are max of  $r$  linearly independent rows and columns in  $A$

# GEOMETRIC INTERPRETATION OF SVD



## LOW RANK APPROXIMATION

Rank r of matrix A might not be small, but we can find  $A_k$  with rank  $k \ll r$ :

$$A_k = \min_{Z: \text{rank}(Z)=k} \|A - Z\|_F$$

Frobenius  
(Euclidean) norm

$$\|X\|_F = \sqrt{\sum_{i=1}^m \sum_{j=1}^n |X_{ij}|^2}$$

A and Z are both  $m \times n$  matrices

# LOW RANK APPROXIMATION

Solution via SVD

$$A_k = U \underbrace{diag(\sigma_1, \dots, \sigma_k, 0, \dots, 0)}_{\text{Set smallest r-k singular values to 0}} V^T$$

Error:  $\|A - A_k\|_F = \sigma_{k+1}$

## BACK TO PCA

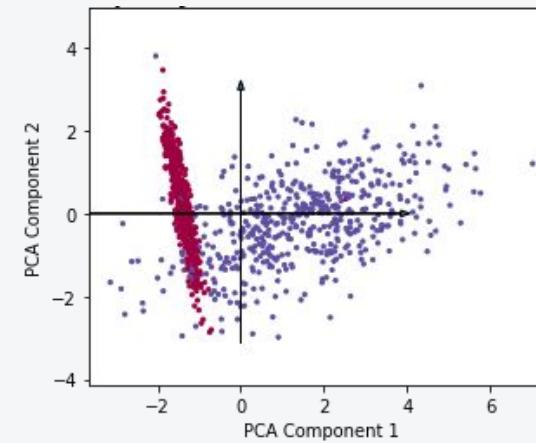
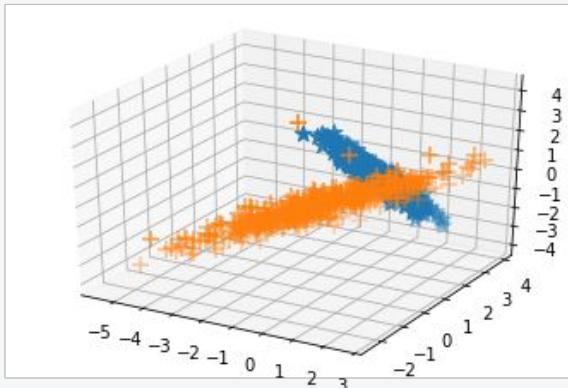
- Find eigenvectors and eigenvalues using SVD
- Examine how much variance is explained using  $k$  components ( $k$  is varied)
- For the purposes of explaining (reconstructing the data), pick the number of components based on variance explained
- For plotting purposes can plot the first 2 PCA or all pairs of interest

PCA

Variance explained:

PC1 PC2 PC3  
[0.73 0.24 0.03]

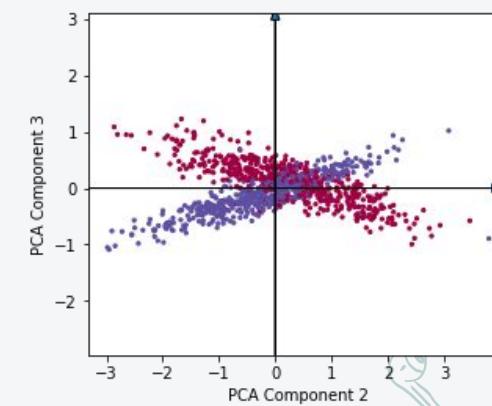
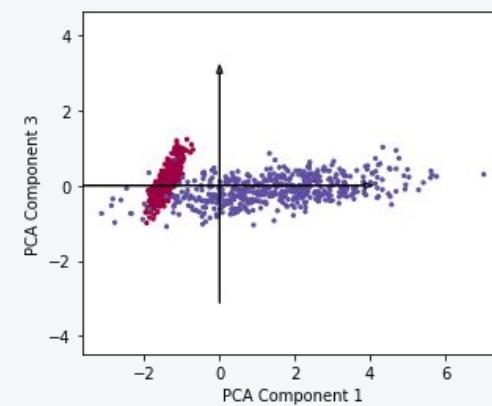
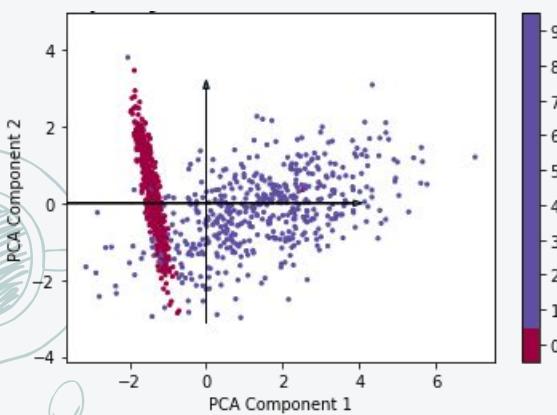
DATA



# PCA

Variance explained:

PC1    PC2    PC3  
[0.73 0.24 0.03]



## ANOTHER DATASET: DIGITS

0	1	2	3	4	5	6	7	8	9	0	1	2	3	4	5	6	7	8	9
0	1	2	3	4	5	6	7	8	9	0	1	2	3	4	5	6	7	8	9
3	9	1	7	7	5	5	1	0	0	2	2	7	8	2	0	1	2	6	3
3	3	3	3	3	4	6	6	9	1	5	0	5	5	2	2	2	0	0	0
1	7	4	2	1	3	9	6	1	1	1	3	1	7	2	1	4	3	1	1
4	0	5	7	6	1	7	5	4	9	2	2	8	2	2	5	7	5	9	3
5	9	1	1	9	0	7	9	0	1	2	3	4	5	6	7	8	3	5	0
0	1	2	3	4	5	6	3	8	9	0	1	2	3	4	5	6	7	8	9
0	5	5	6	5	0	9	8	5	3	4	1	7	7	3	5	1	0	0	0
1	2	7	8	1	0	1	2	6	3	3	3	4	6	6	6	4	5	5	5
1	5	0	5	5	2	1	1	0	0	1	7	6	3	1	4	3	1	3	1
5	4	7	6	4	7	4	4	0	5	7	6	3	5	6	7	5	4	4	4
2	1	8	1	5	5	4	8	5	4	9	0	8	9	8	0	1	2	3	3
4	5	6	7	8	7	0	4	2	3	9	1	6	7	8	7	0	4	2	3
4	5	6	7	8	9	0	3	5	5	6	5	0	9	8	9	8	4	1	7
7	3	5	1	0	0	2	2	7	8	2	0	4	2	6	3	3	7	3	3
5	6	6	6	4	9	1	5	0	9	5	2	8	2	0	0	4	7	6	3
2	4	2	6	6	3	1	3	9	1	2	6	8	4	3	1	5	0	5	3
6	7	6	1	7	5	4	4	7	2	8	2	2	5	7	7	5	4	3	8
6	9	0	1	9	3	0	1	1	3	4	5	6	7	8	9	0	1	1	3

Optical recognition of handwritten digits dataset

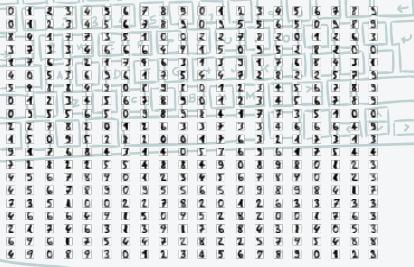
Number of Instances: 5620

Number of Attributes: 64

Attribute Information: 8x8 image of integer pixels in the range 0..16

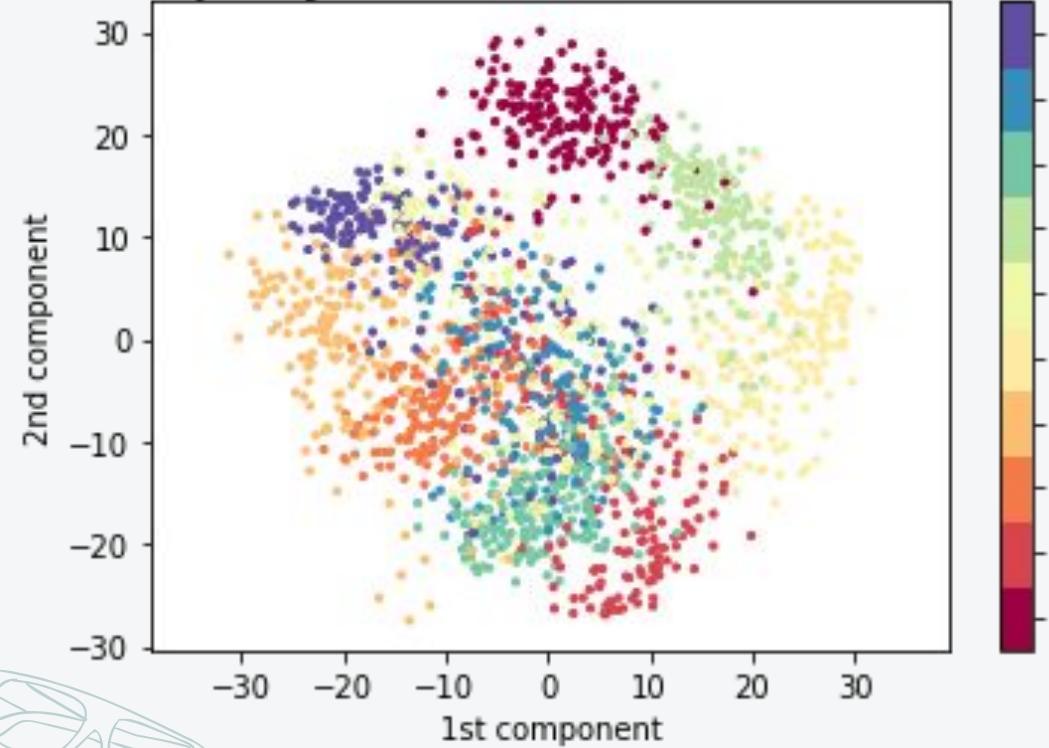
Creator: E. Alpaydin

Date: July, 1998



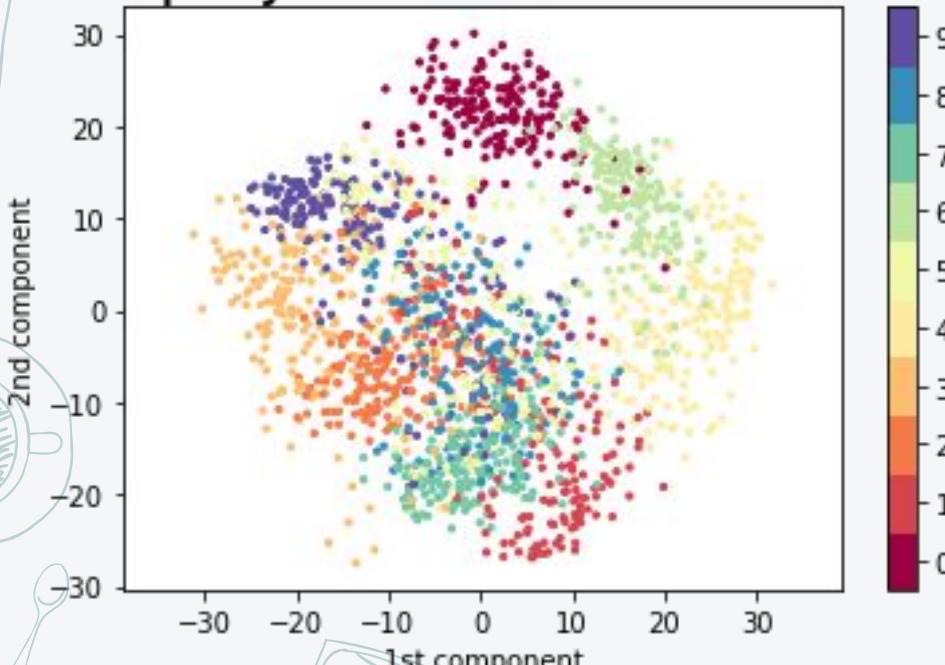
## PCA: DIGITS

### PCA projection of the dataset

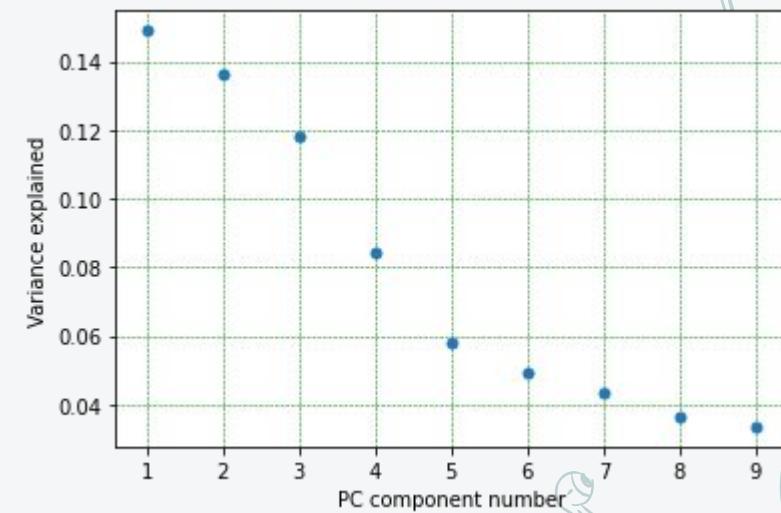


# PCA: DIGITS

PCA projection of the dataset



Variance explained



# PCA

## Advantages:

- Relatively Fast
- Gives some sense of intrinsic dimensionality
- Can compute variation explained
- Axes are interpretable

## Disadvantages:

- Linear reduction in dimensionality
- Orthogonal components (hard to model non-linear data)

# NON-LINEAR LOWER DIMENSIONAL EMBEDDING

Neighborhood based

- T-SNE
- UMAP
- ISOMAP
- LLE

Autoencoder-based

And many more...

# T-SNE : T-DISTRIBUTED STOCHASTIC NETWORK EMBEDDING

- *Distance preservation*: when 2 points are close in high dimensional space, they should be close in 2 dimensional space
- *Neighbor preservation*: neighborhoods of a given point in high-d should be neighborhoods of a given point in low-d
  - Note1: effective number of neighbors is a parameter called *perplexity*

Note2: t-SNE focuses on the order of things, not necessarily on the distance itself

# T-SNE

Similarity of data points in High Dimensional space

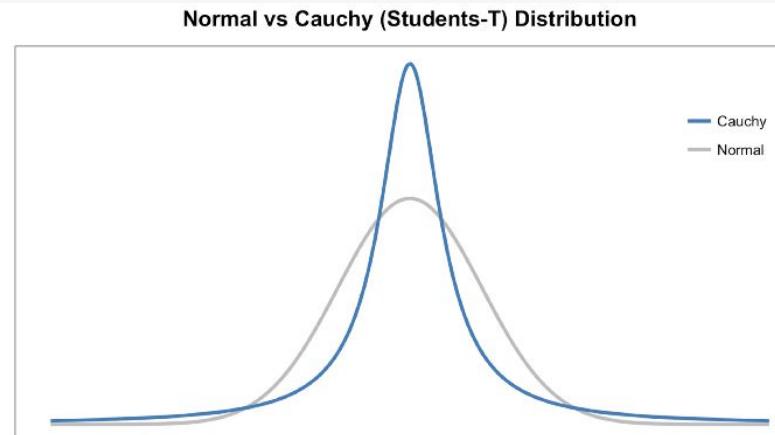
$$p_{ij} = \frac{\exp(-\|x_i - x_j\|^2/2\sigma^2)}{\sum_{k \neq l} \exp(-\|x_k - x_l\|^2/2\sigma^2)}$$

Similarity of points in Low Dimensional space

$$q_{ij} = \frac{(1 + \|y_i - y_j\|^2)^{-1}}{\sum_{k \neq l} (1 + \|y_k - y_l\|^2)^{-1}}$$

# T-SNE

Reason for t-distribution as opposed to Gaussian as in high d is crowding - need fatter tails to accommodate distant points



# T-SNE COST FUNCTION

Cost function: Kullback Leibler divergence (relative entropy)  
- measures distance between two distributions

$$C = KL(P||Q) = \sum_i \sum_j p_{ij} \log \frac{p_{ij}}{q_{ij}}$$

(Kullback and Leibler, 1951)

Properties?

# T-SNE COST FUNCTION

Cost function: Kullback Leibler divergence (relative entropy)  
- measures distance between two distributions

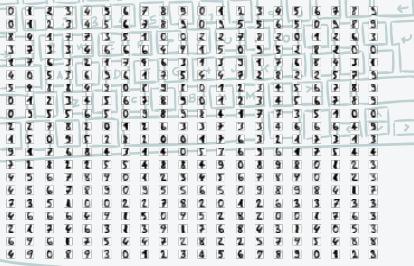
$$C = KL(P||Q) = \sum_i \sum_j p_{ij} \log \frac{p_{ij}}{q_{ij}}$$

(Kullback and Leibler, 1951)

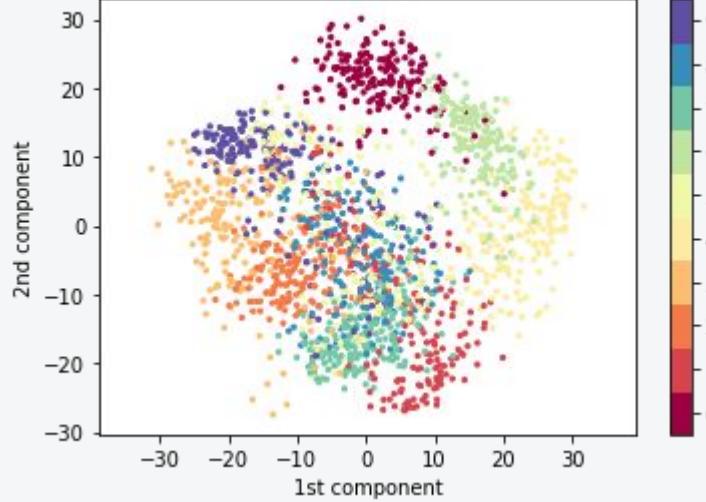
Properties:

- Large  $p_{ij}$  modeled by small  $q_{ij}$  - large penalty (distant points shouldn't be close in low d)
- Small  $p_{ij}$  modeled by large  $q_{ij}$  - small penalty (ok to separate nearby points -most are beyond neighborhood)

Optimize the cost using gradient descent

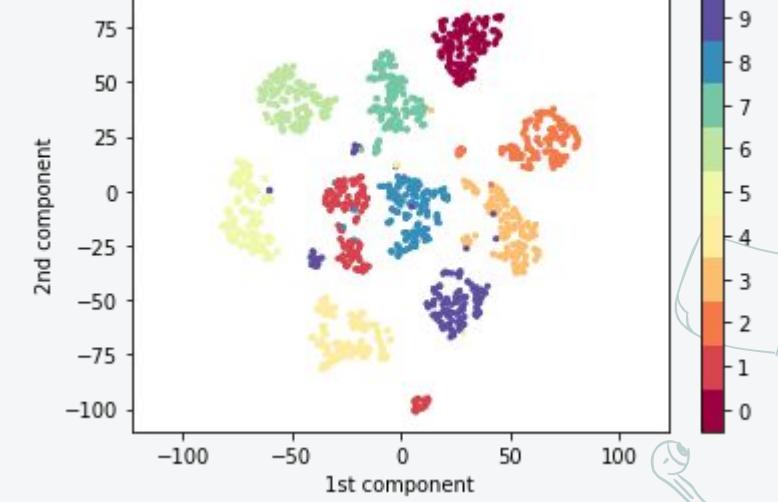


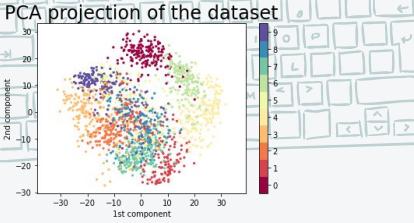
PCA projection of the dataset



# T-SNE ON DIGITS

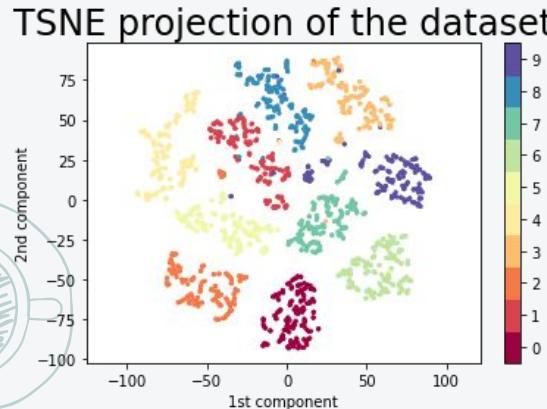
TSNE projection of the dataset



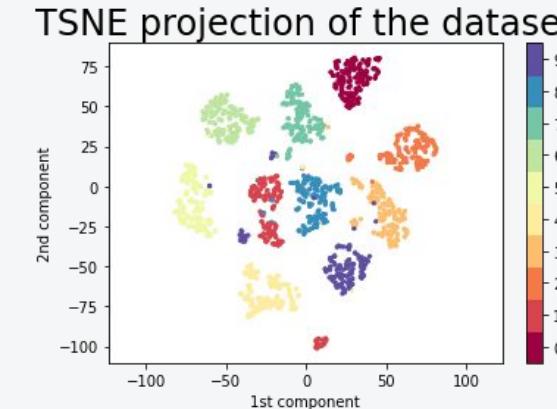


# T-SNE SENSITIVITY TO PARAMETERS: PERPLEXITY

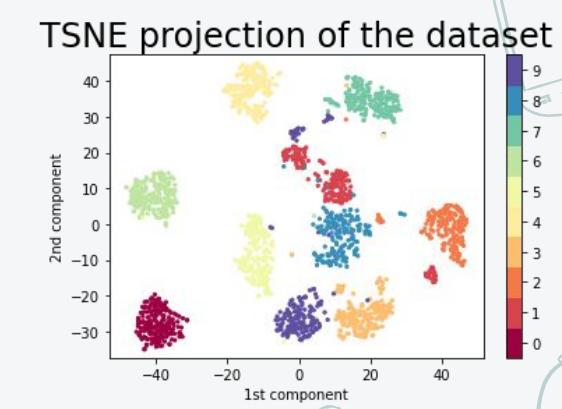
Perplexity = 5



Perplexity = 10

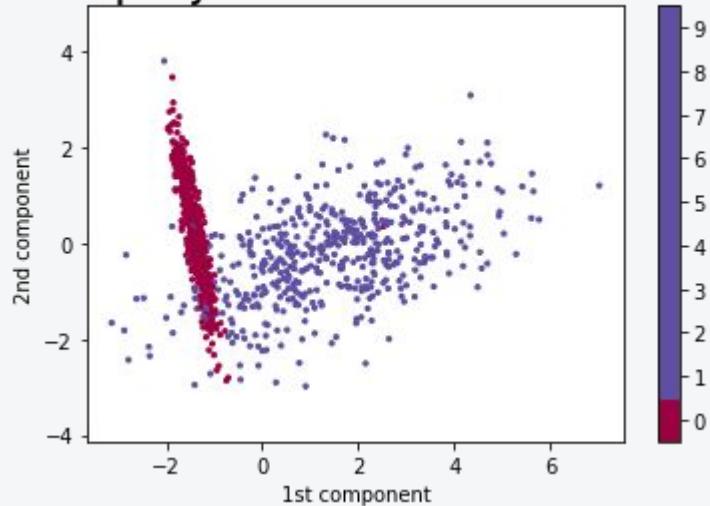


Perplexity = 50

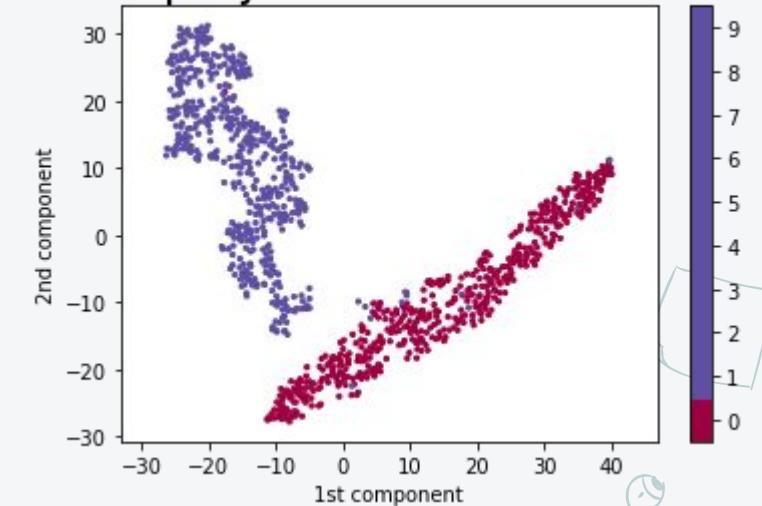


# T-SNE ON THE ORIGINAL DATASET

PCA projection of the dataset



TSNE projection of the dataset



# T-SNE

## Advantages

- Retains local structure
- Excels at visualizing complex high-d data in 2d

## Disadvantages

- Slower
- Doesn't really capture global structure
- Sensitive to params
- Clustering and distance depend on perplexity (shouldn't be interpreted off the plot)
- Stochasticity leads to different results
- Not guaranteed to converge to the global optimum of the cost function

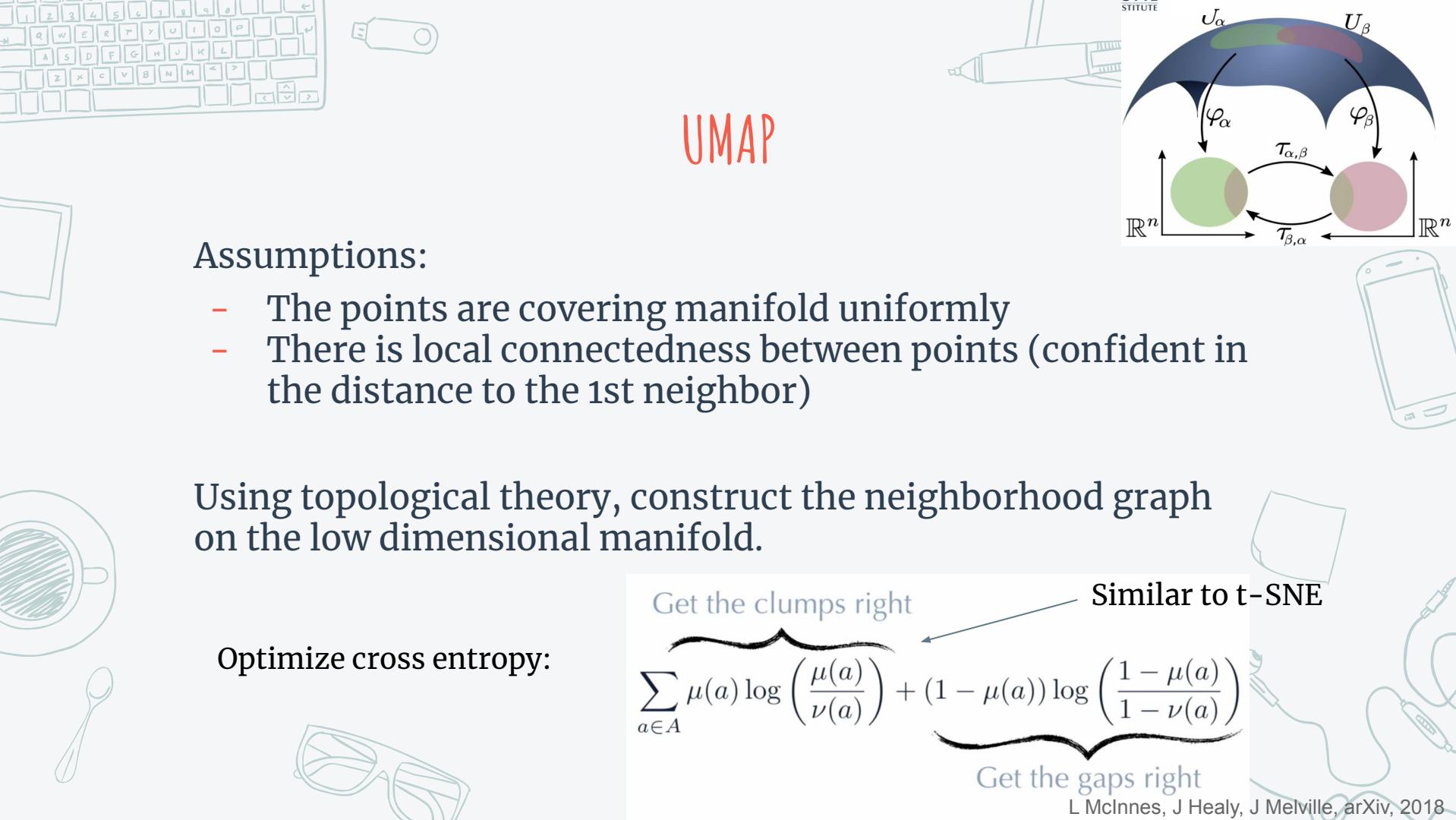
# UMAP: UNIFORM MANIFOLD APPROXIMATION AND PROJECTION

UMAP – preserves both local and global structure by

1. Learning the manifold of the data
2. Embedding that manifold into a low dimensional space
  - Builds on neighborhood based approaches but adds mathematical foundations

Key intuition:

While the assumptions may not be true in the original space, they can be true on the manifold



## Assumptions:

- The points are covering manifold uniformly
- There is local connectedness between points (confident in the distance to the 1st neighbor)

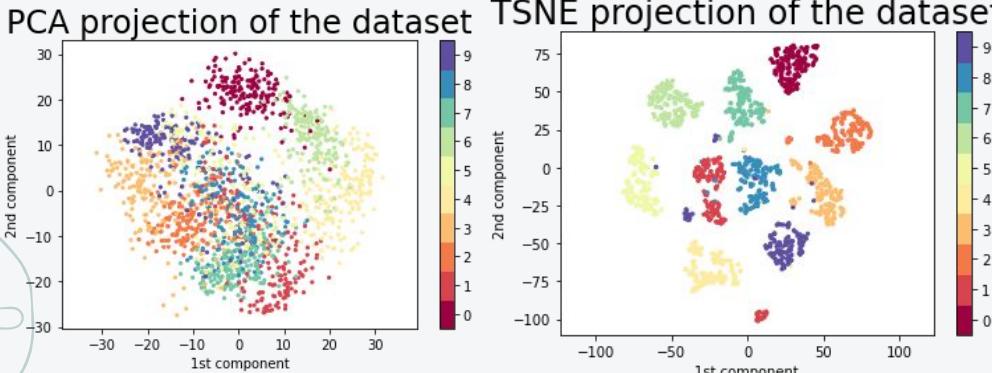
Using topological theory, construct the neighborhood graph on the low dimensional manifold.

Optimize cross entropy:

$$\sum_{a \in A} \mu(a) \log \left( \frac{\mu(a)}{\nu(a)} \right) + (1 - \mu(a)) \log \left( \frac{1 - \mu(a)}{1 - \nu(a)} \right)$$

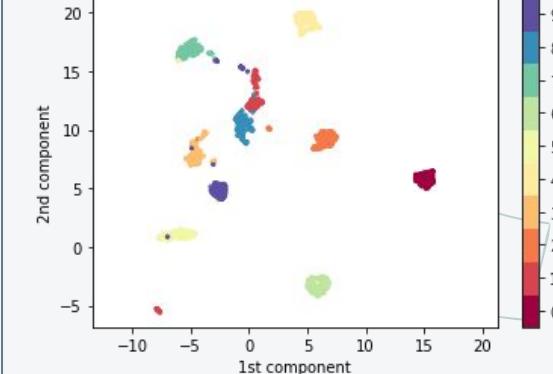
Get the clumps right      Similar to t-SNE  
Get the gaps right

L McInnes, J Healy, J Melville, arXiv, 2018



# UMAP ON DIGITS

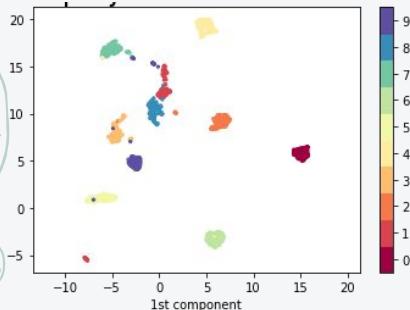
UMAP projection of the dataset



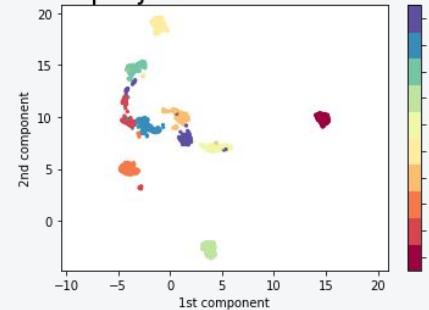
# UMAP HYPERPARAMETERS

- Number of neighbors - local vs global structure
- Minimum distance - how tightly points can be packed together
- Number of components (we set 2 for all)
- Distance metric

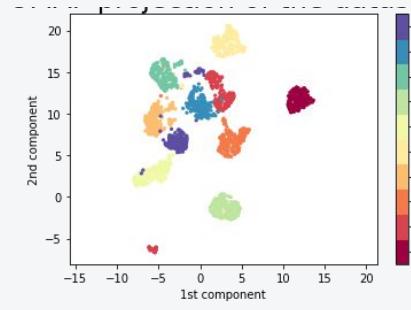
neighbors = 15  
metric = euclidian  
min\_dist = 0.1



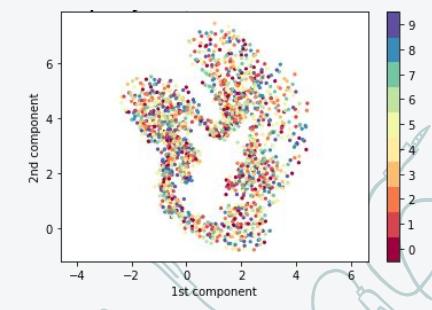
neighbors = 50  
metric = euclidian  
min\_dist = 0.1



neighbors = 5  
metric = euclidian  
min\_dist = 0.5

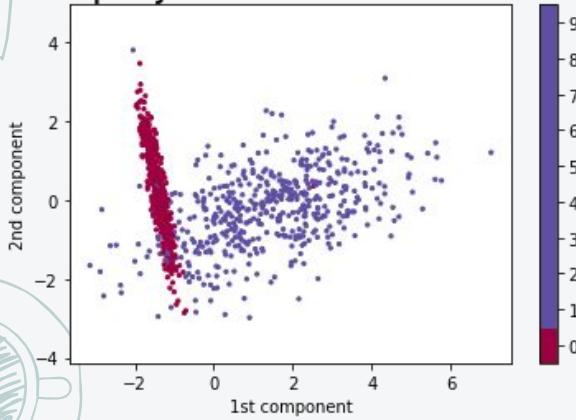


neighbors = 5  
metric = mahalanobis  
min\_dist = 0.1

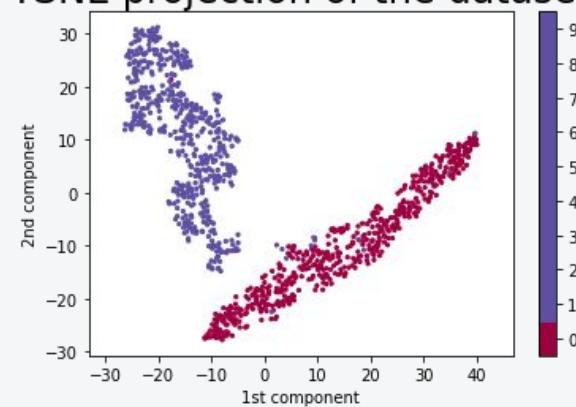


# UMAP ON THE ORIGINAL DATA

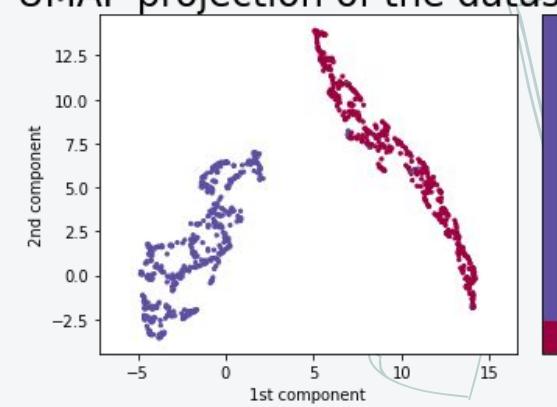
PCA projection of the dataset



TSNE projection of the dataset



UMAP projection of the dataset



# UMAP

## Advantages

- Preserves global as well as local structure
- Has mathematical guarantees
- More robust to hyperparameter changes (e.g. number of neighbors)

## Disadvantages

- Depends on hyperparameters
- Slow-ish

# READING

## PCA

- In depth PCA with examples from the Python handbook:  
<https://jakevdp.github.io/PythonDataScienceHandbook/05.09-principal-component-analysis.html>

## t-SNE

- Playing with t-sne: <https://distill.pub/2016/misread-tsne/>
- Intro to t-sne with python:  
<https://towardsdatascience.com/an-introduction-to-t-sne-with-python-example-5a3a293108d1>

## UMAP

- Talk (26min): <https://www.youtube.com/watch?v=nq6iPZVUxZU>
- Examples with code in python:  
[https://umap-learn.readthedocs.io/en/latest/basic\\_usage.html](https://umap-learn.readthedocs.io/en/latest/basic_usage.html)

## t-SNE vs UMAP:

<https://towardsdatascience.com/tsne-vs-umap-global-structure-4d8045acba17>