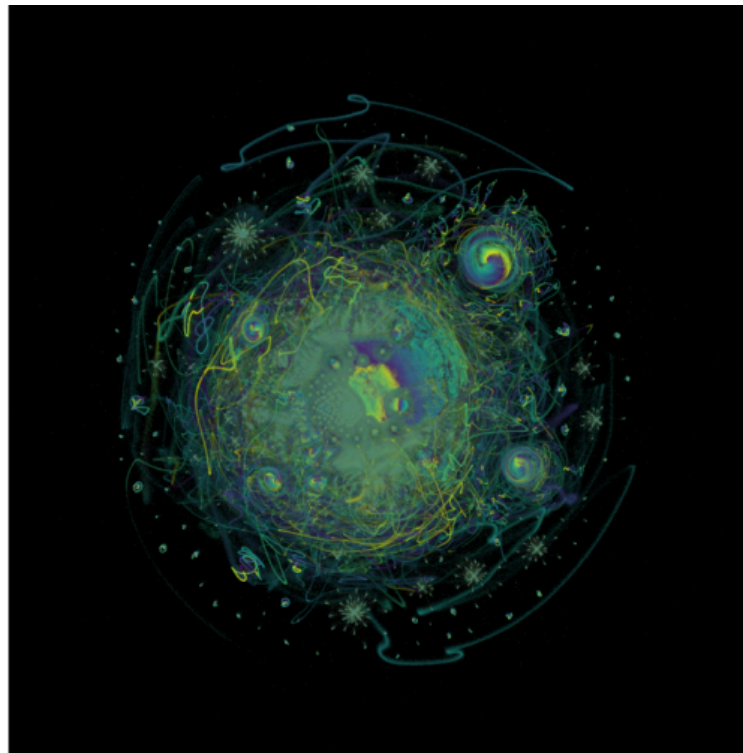
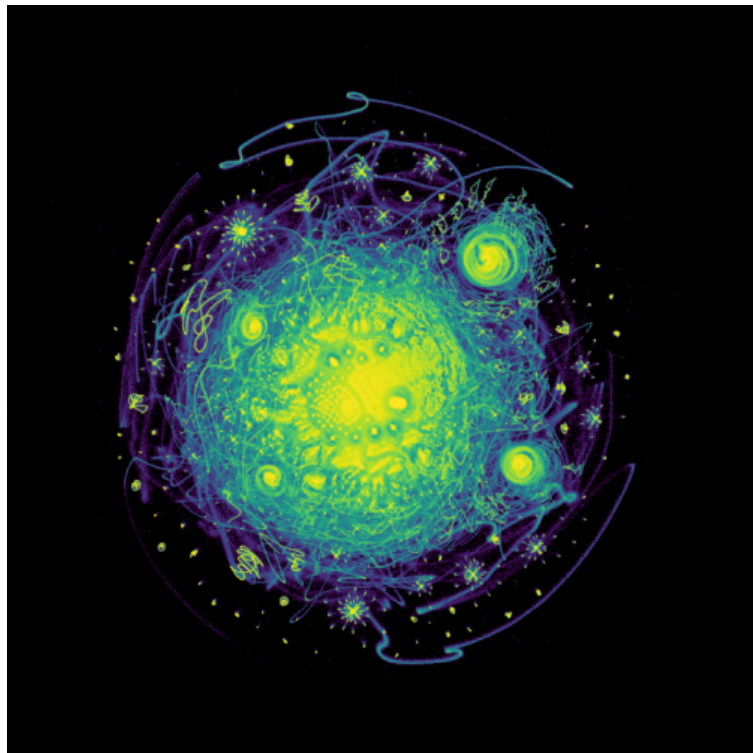


# Dimension Reduction for Single Cell Data Analysis

Nikolay Oskolkov, Lund University, NBIS SciLifeLab, Sweden  
scRNAseq course, 13.02.2024



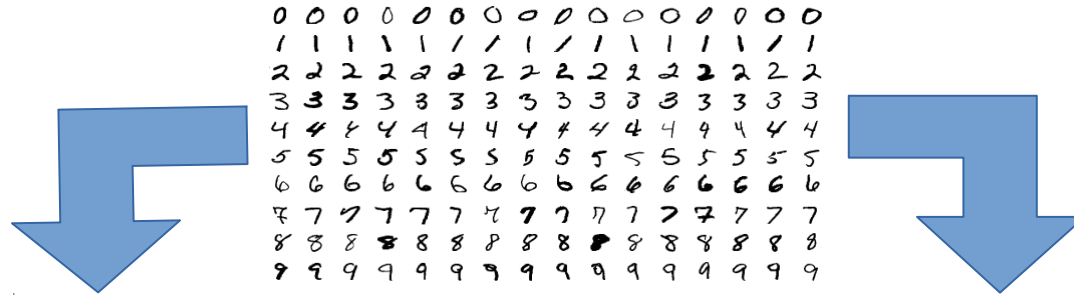
@NikolayOskolkov



**GitHub**

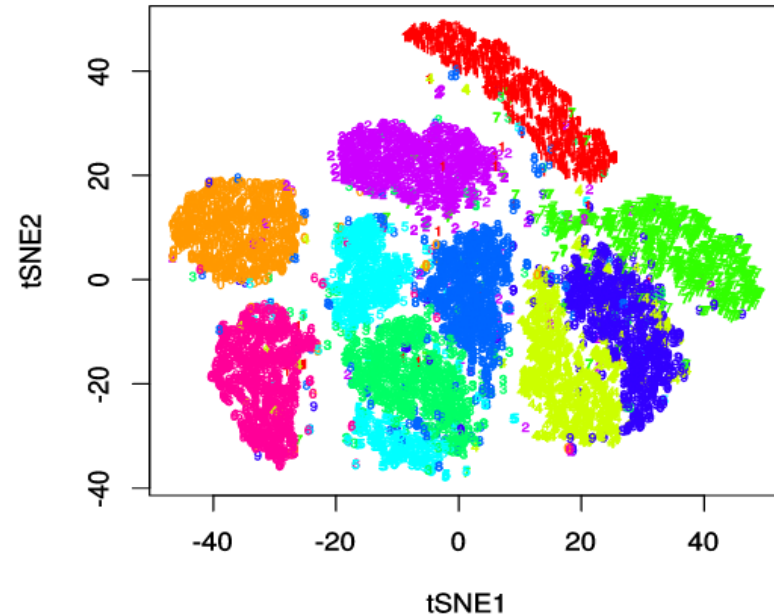
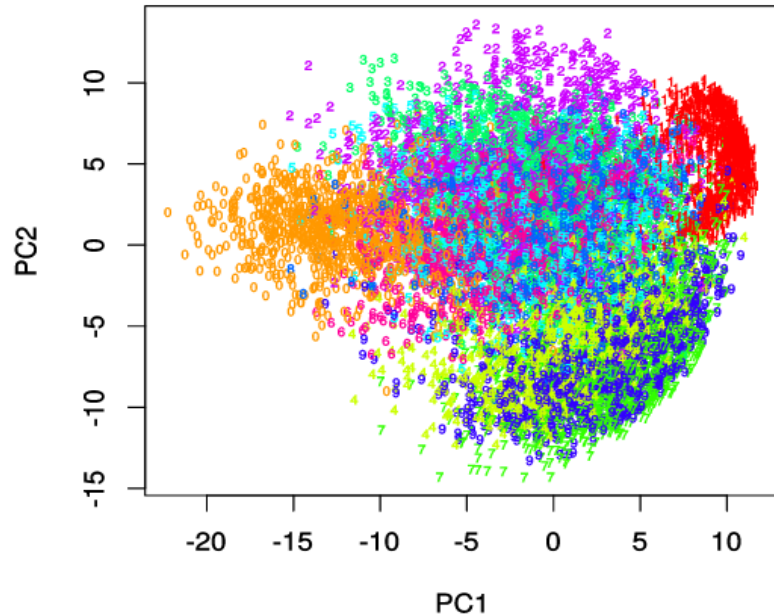
[github.com/NikolayOskolkov](https://github.com/NikolayOskolkov)

**Dimensionality reduction  
is also supposed to ... reduce dimensions**



**PCA PLOT WITH PRCOMP**

**tSNE MNIST**



The goal of dimension reduction is not only visualization but also reducing dimensions

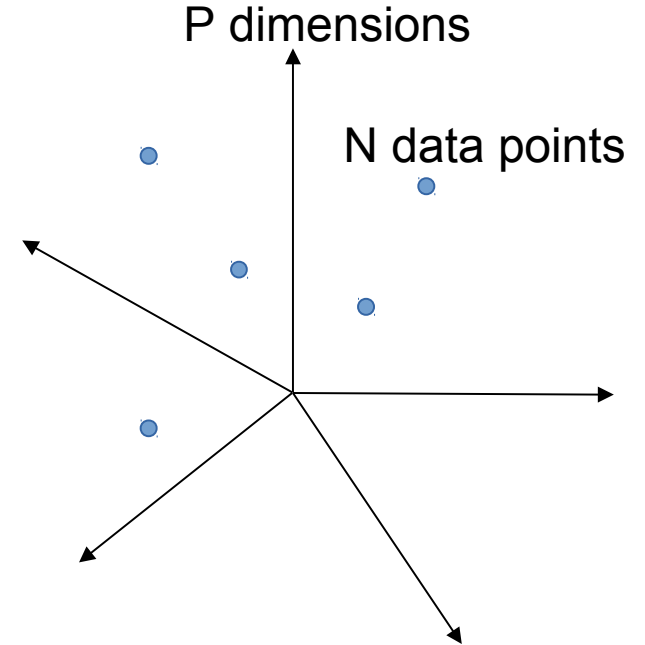
Statistical observations:  
e.g. samples, cells etc.

Features: genes, proteins,  
microbes, metabolites etc.

**N** →

0	3	1	0	2	3	8	1	1	3
1	1	0	0	7	1	2	2	3	3
1	2	2	0	0	6	7	1	2	2
1	2	3	10	0	4	6	1	0	5
3	2	2	1	4	3	2	1	6	0
7	4	4	5	3	9	6	1	6	1
7	1	1	5	2	8	9	1	3	6
5	0	1	6	2	0	0	0	1	5
1	6	3	3	4	6	2	0	1	1
1	2	2	4	1	1	3	0	8	2

← **P**



**High Dimensional Data:**  
 **$P \gg N$**

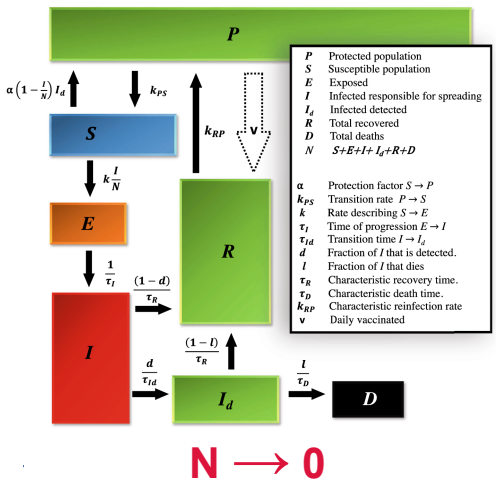
For a robust statistical analysis, one should properly “sample” the P-dimensional space, hence large sample size is required,  $N \gg P$



**P** is the number of features (genes, proteins, genetic variants etc.)  
**N** is the number of observations (samples, cells, nucleotides etc.)

Biology / Biomedicine

## Mathematical modeling



## Bayesianism

$N \ll P$

## Frequentism

$N \approx P$

## Machine Learning

$N \gg P$



## Amount of Data

The Curse of Dimensionality

$$Y = \alpha + \beta X$$

$$\beta = (X^T X)^{-1} X^T Y$$

$$(X^T X)^{-1} \sim \frac{1}{\det(X^T X)} \dots \rightarrow \infty, \quad n \ll p$$

We need to reduce dimensions to overcome the Curse of Dimensionality!

## The curse(s) of dimensionality

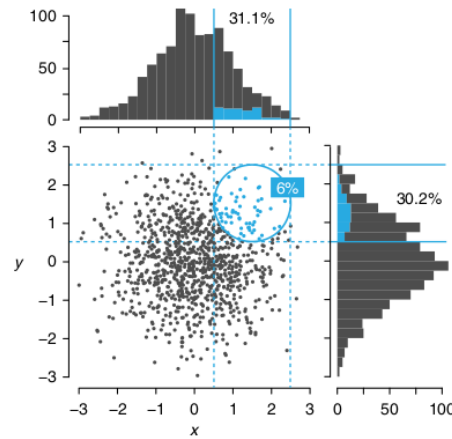
There is such a thing as too much of a good thing.

Naomi Altman and Martin Krzywinski

We generally think that more information is better than less. However, in the 'big data' era, the sheer number of variables that can be collected from a single sample can be problematic. This embarrassment of riches is called the 'curse of dimensionality'<sup>1</sup> (CoD) and manifests itself in a variety of ways. This month, we discuss four important problems of dimensionality as it applies to data sparsity<sup>1,2</sup>, multicollinearity<sup>3</sup>, multiple testing<sup>4</sup> and overfitting<sup>5</sup>. These effects are amplified by poor data quality, which may increase with the number of variables.

Throughout, we use  $n$  to indicate the sample size from the population of interest and  $p$  to indicate the number of observed variables, some of which may have missing values for some samples. For example, we may have  $n = 1,000$  subjects and  $p = 200,000$  single-nucleotide polymorphisms (SNPs).

First, as the dimensionality  $p$  increases, the 'volume' that the samples may occupy

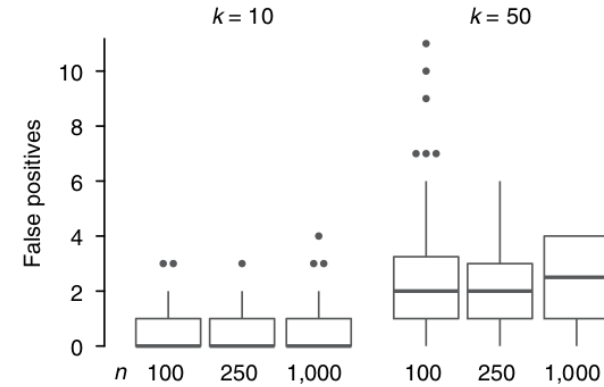


**Fig. 1 | Data tend to be sparse in higher dimensions.** Among 1,000  $(x, y)$  points in which both  $x$  and  $y$  are normally distributed with a mean of 0 and s.d.  $\sigma = 1$ , only 6% fall within  $\sigma$  of  $(x, y) = (1.5, 1.5)$  (blue circle). However, when the data are projected into a lower dimension—shown by histograms—about 30% of the points (all bins

A and 100 to have the minor allele  $a$ . If we tabulate on two SNPs, A and B, we will expect only ten samples to exhibit both minor alleles with genotype  $ab$ . With SNPs A, B and C, we expect only one sample to have genotype  $abc$ , and with four or more SNPs, we expect empty cells in our table. We need a much larger sample size to observe samples with all the possible genotypes. As  $p$  increases, we may quickly find that there are no samples with similar values of a predictor.

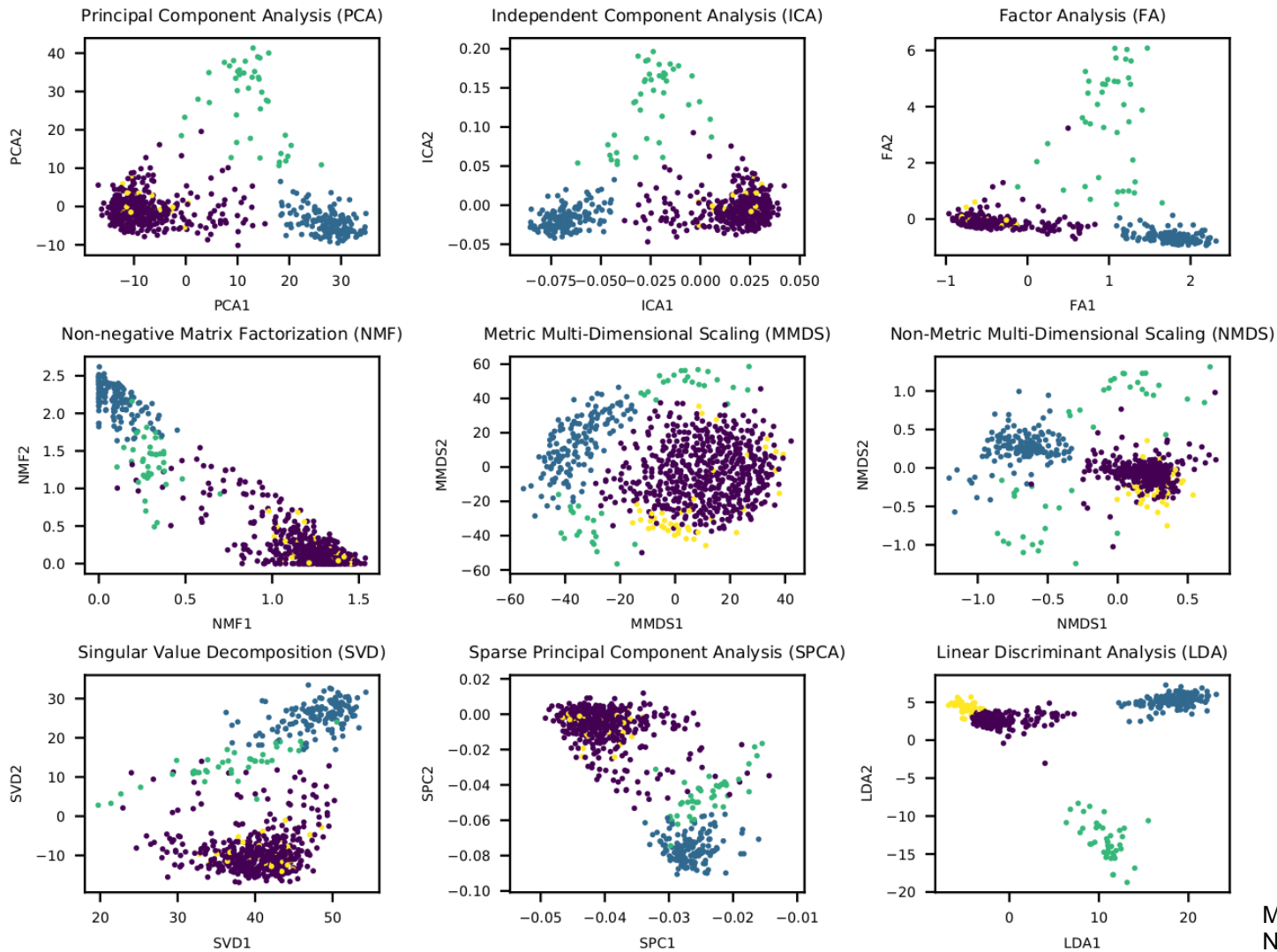
Even with just five SNPs, our ability to predict and classify the samples is impeded because of the small number of subjects that have similar genotypes. In situations where there are many gene variants, this effect is exacerbated, and it may be very difficult to find affected subjects with similar genotypes and hence to predict or classify on the basis of genetic similarity.

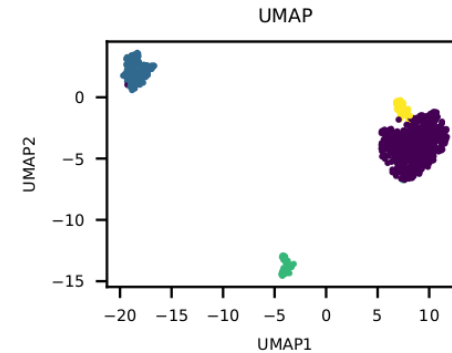
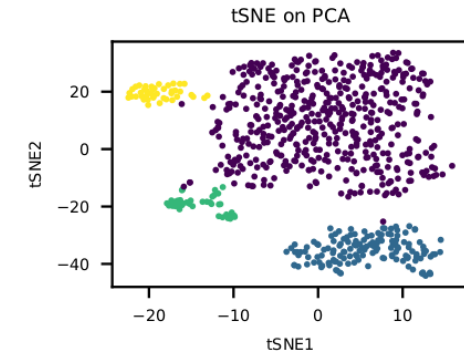
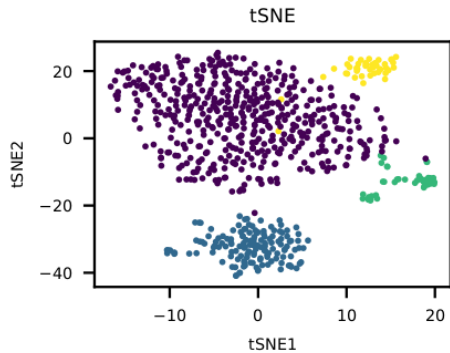
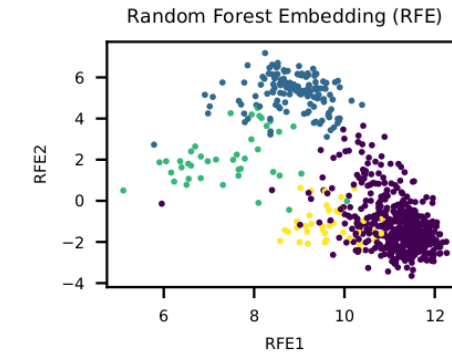
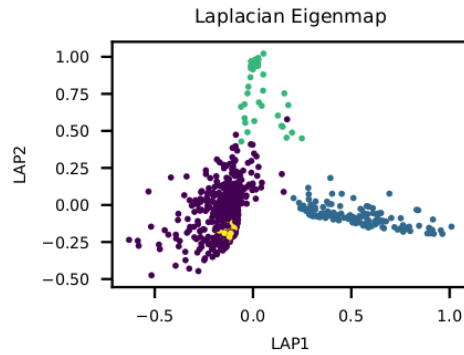
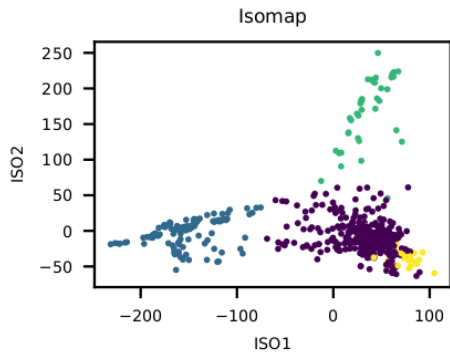
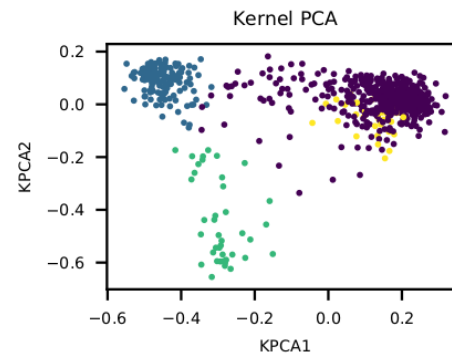
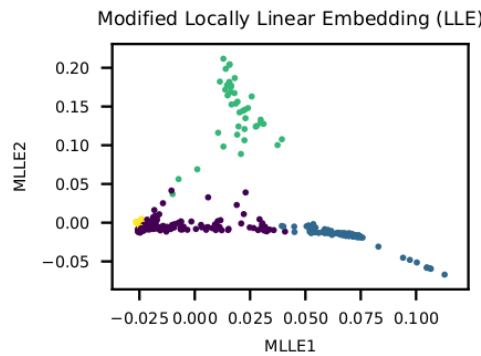
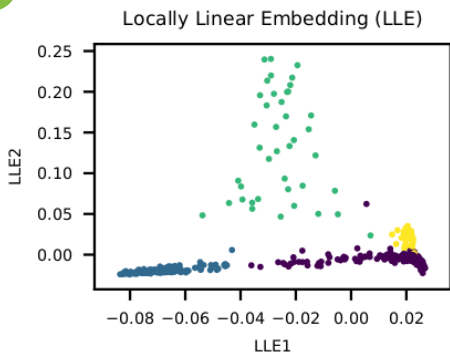
If we treat the distance between points (e.g., Euclidian distance) as a measure of similarity, then we interpret greater distance



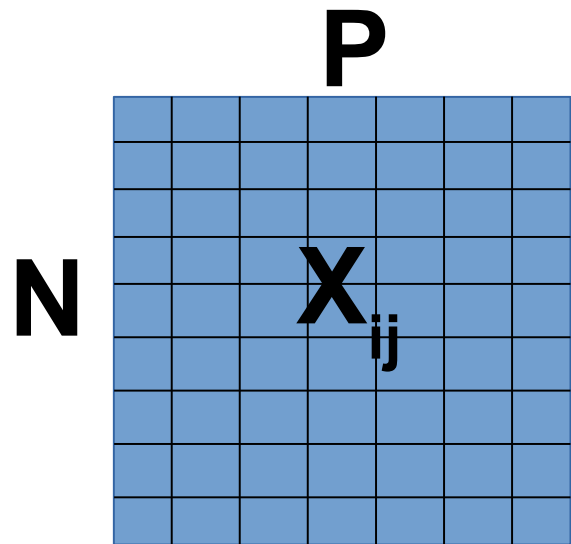
**Fig. 3 | The number of false positives increases with each additional predictor.** The box plots show the number of false positive regression-fit  $P$  values (tested at  $\alpha = 0.05$ ) of 100 simulated multiple regression fits on various numbers of samples ( $n = 100, 250$  and  $1,000$ ) in the presence of one true predictor and  $k = 10$  and  $50$  extraneous uncorrelated predictors. Box plots show means (black center lines), 25th and 75th percentiles (box edges), and minimum and maximum values (whiskers). Outliers (dots) are jittered.

# Dimension reduction techniques: linear vs. non-linear



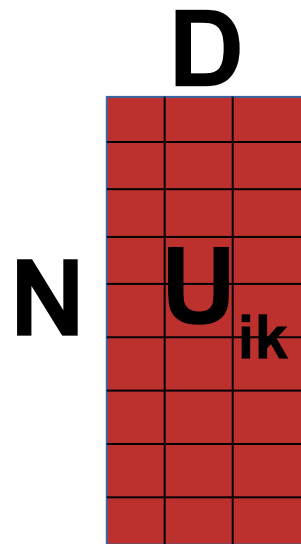


$$X_{ij} \approx U_{ik} V_{kj}$$

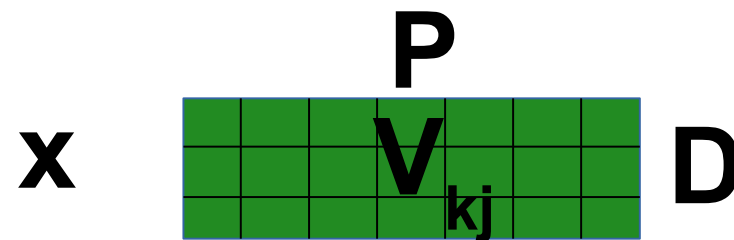


Data

≈



Low-dimensional  
data representation  
(embeddings)



Loadings

$$\text{Loss} = \sum_{i=1}^N \sum_{j=1}^P (X_{ij} - U_{ik} V_{kj})^2$$

## Coding in R:

```
data_centered <- scale(data, center = TRUE, scale = FALSE)
```

```
covariance <- t(data_centered) %*% data_centered
```

```
eig <- eigen(covariance)
```

```
plot(eig$vectors[,1:2]);
```

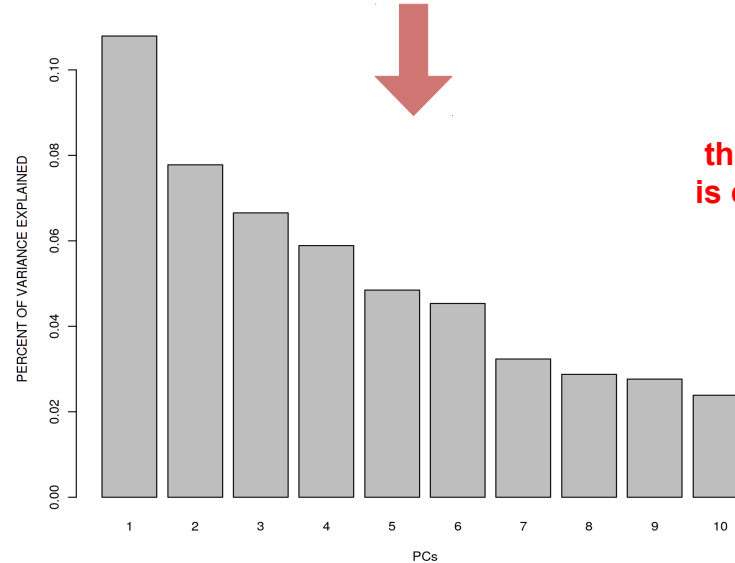
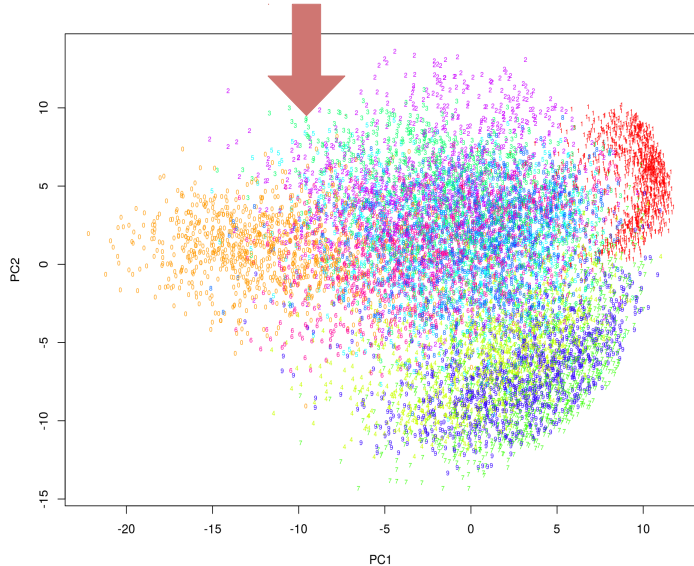
```
barplot(eig$values / sum(eig$values))
```

## Mathematically:

$$M_{ij} = X_{ij} - \mu_j$$

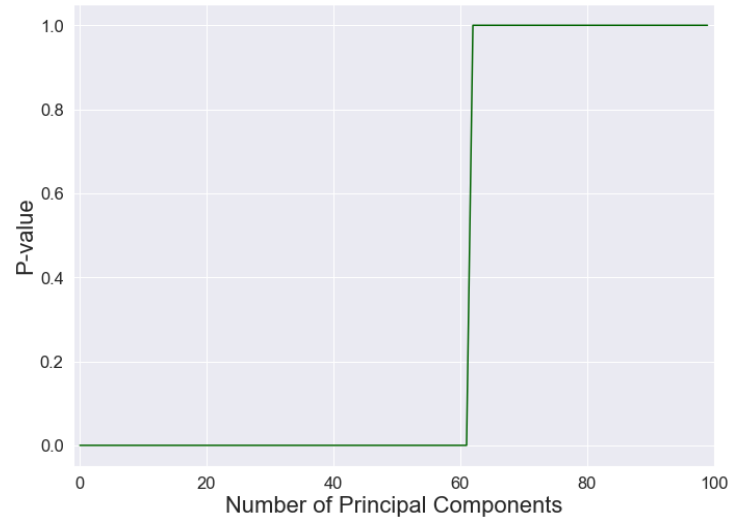
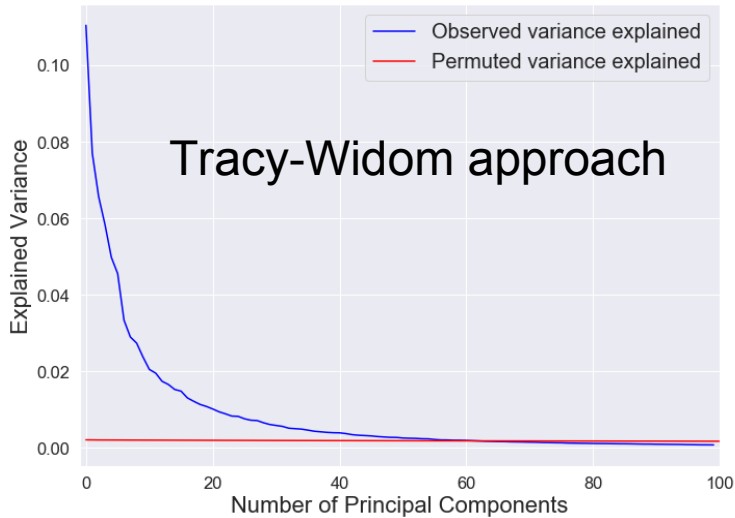
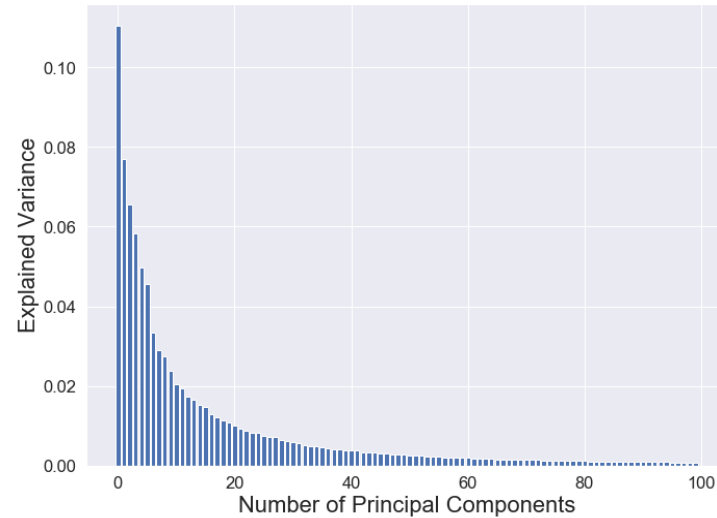
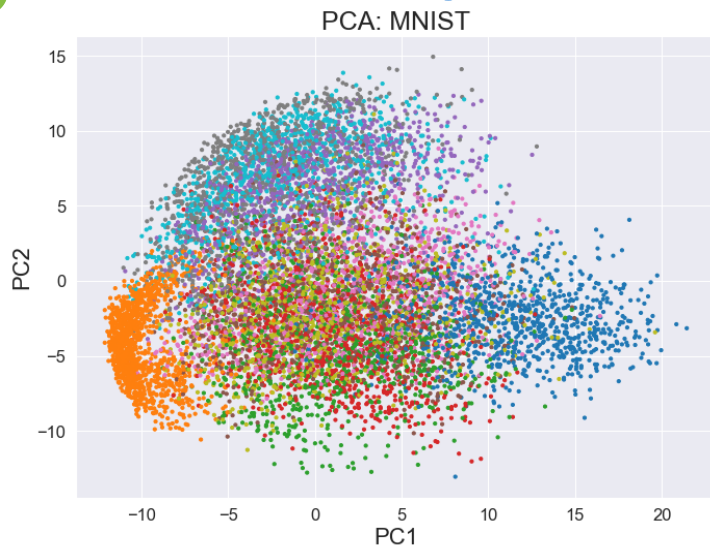
$$A = (1/N) * M^T * M$$

$$A * u = \lambda * u$$



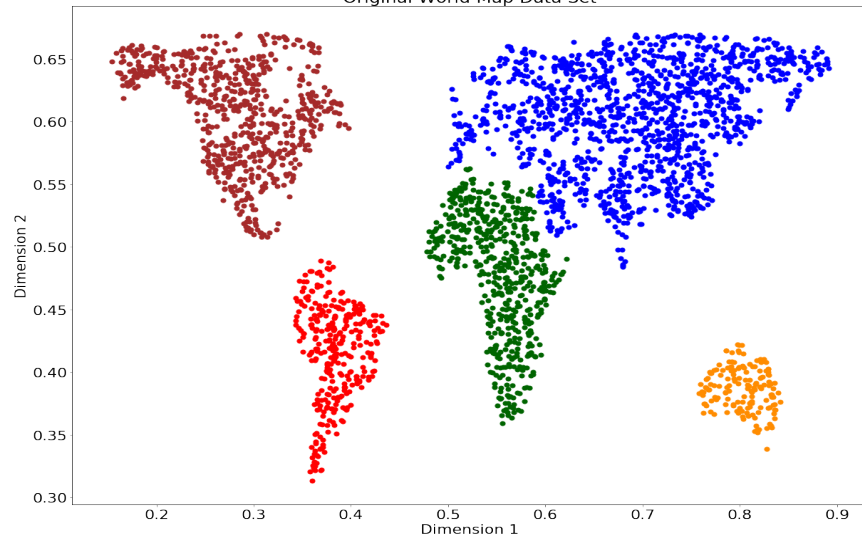
**It can be analytically derived that the eigen value decomposition in PCA is equivalent to projecting data on axes of maximal variation in the data**



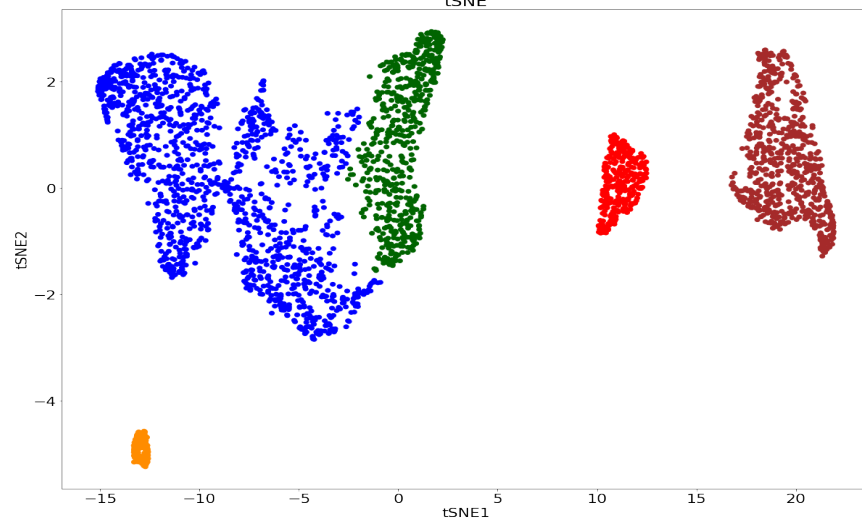
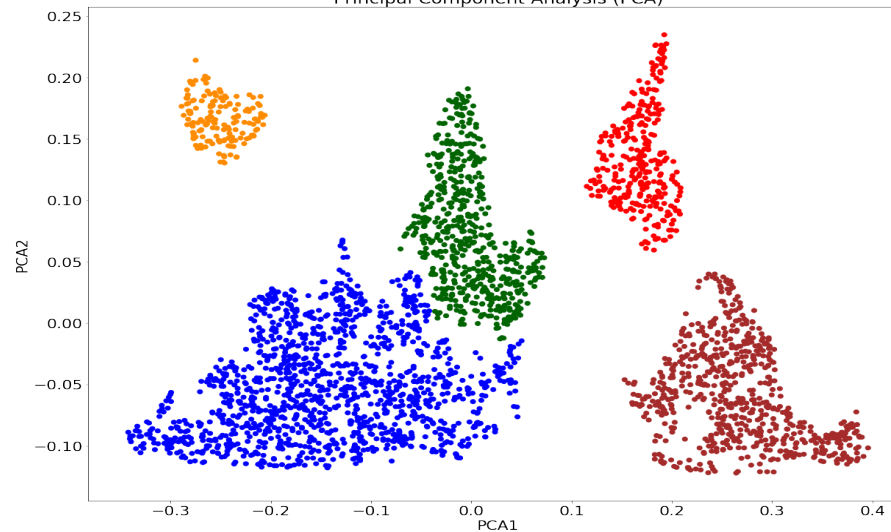


**In Seurat:  
JackStraw**

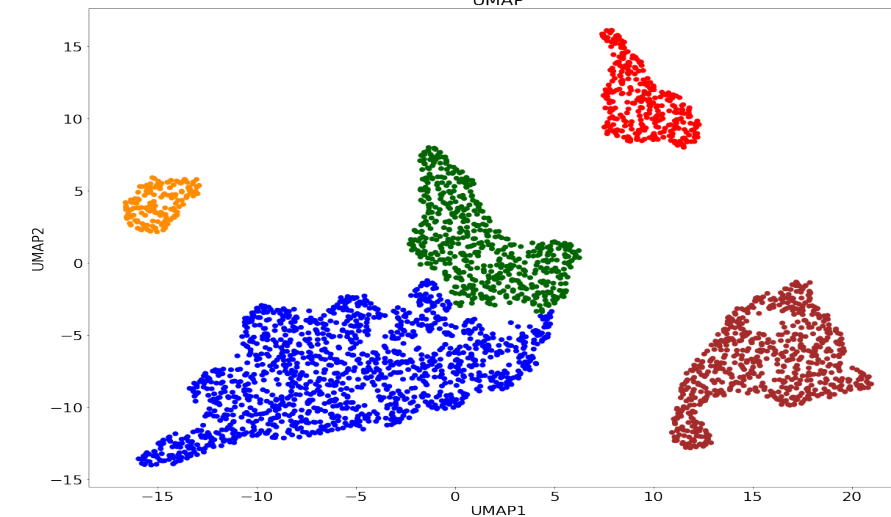
Original World Map Data Set

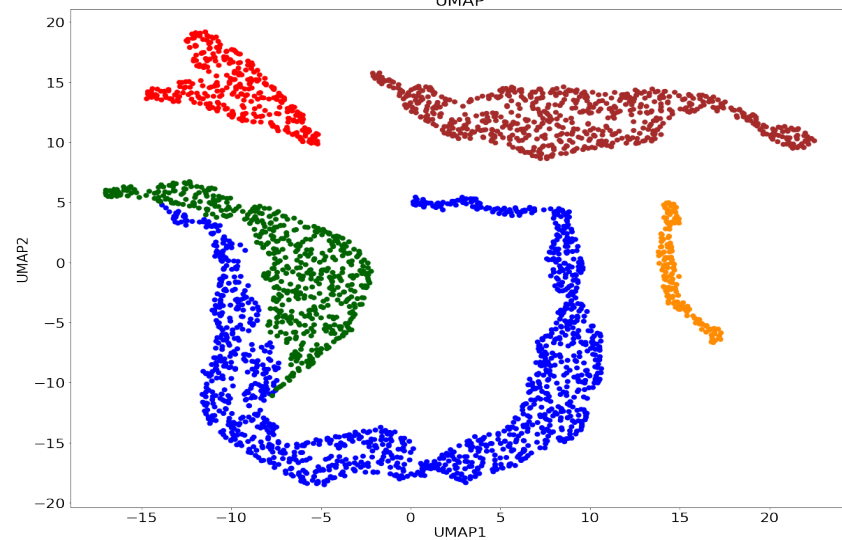
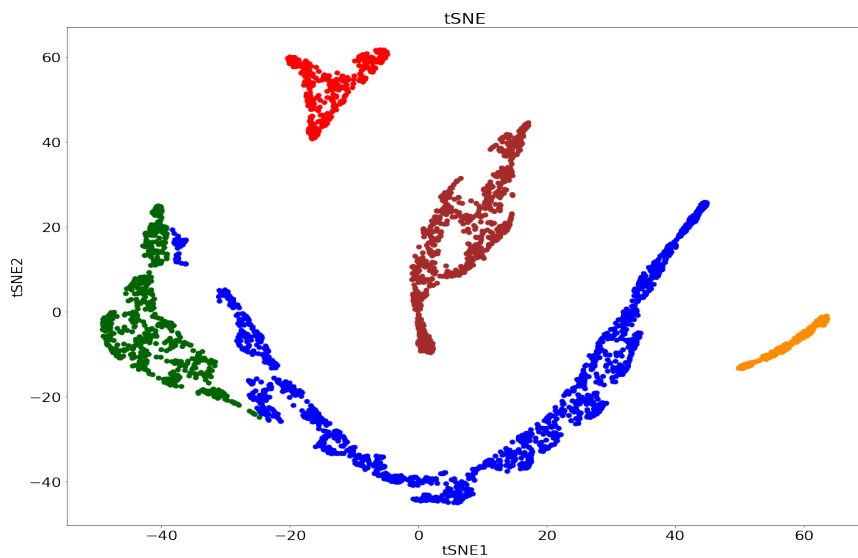
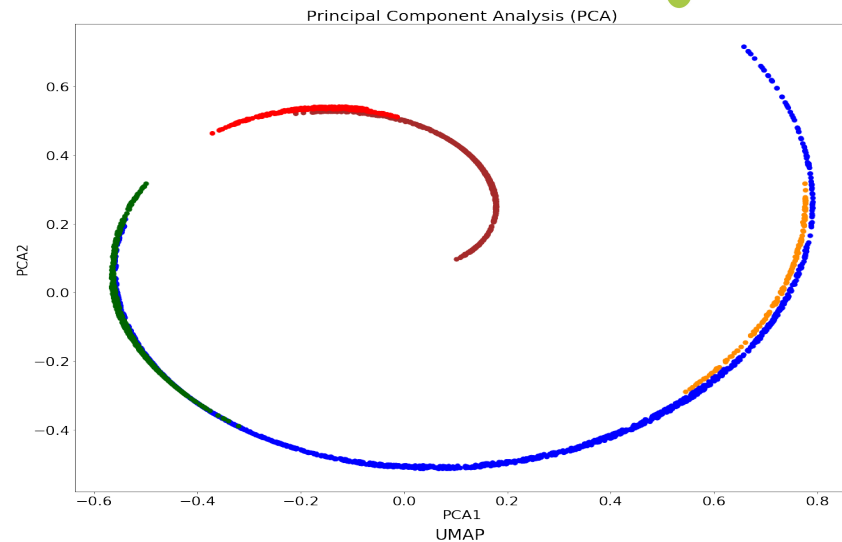
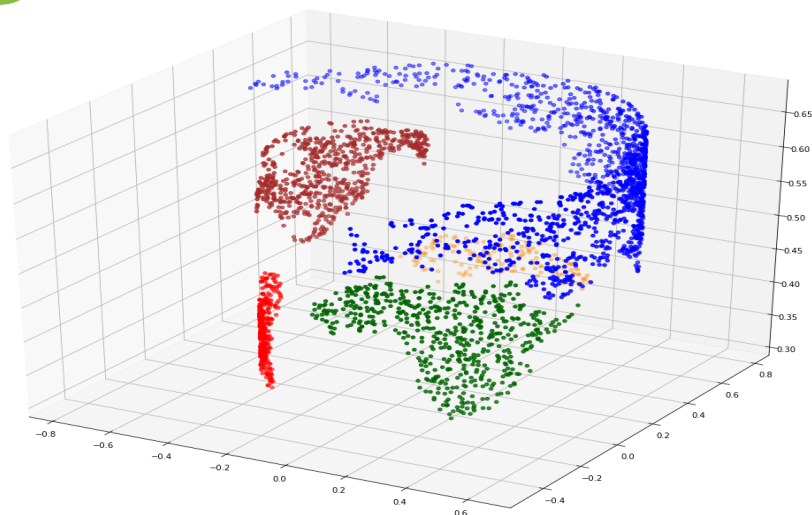


Principal Component Analysis (PCA)

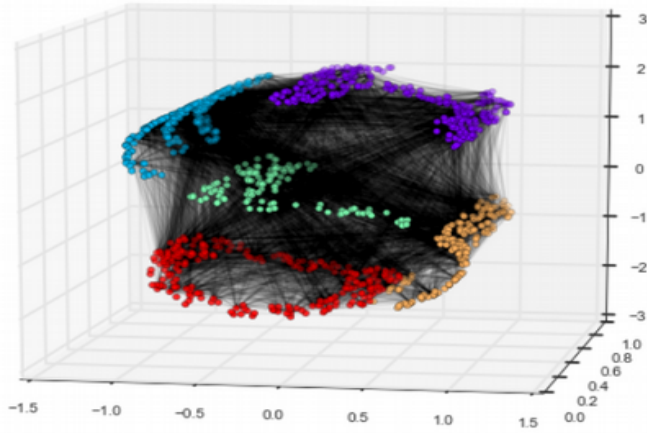


UMAP

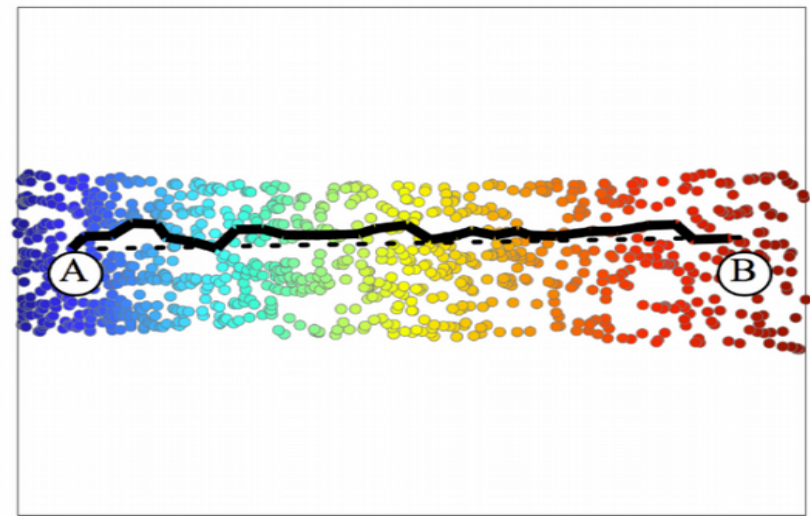
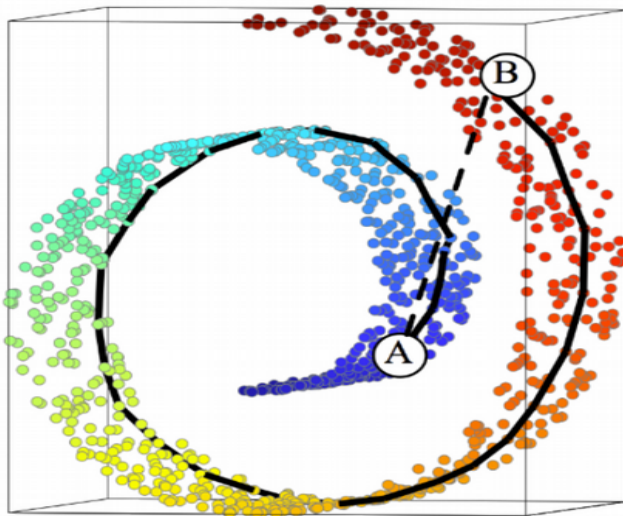
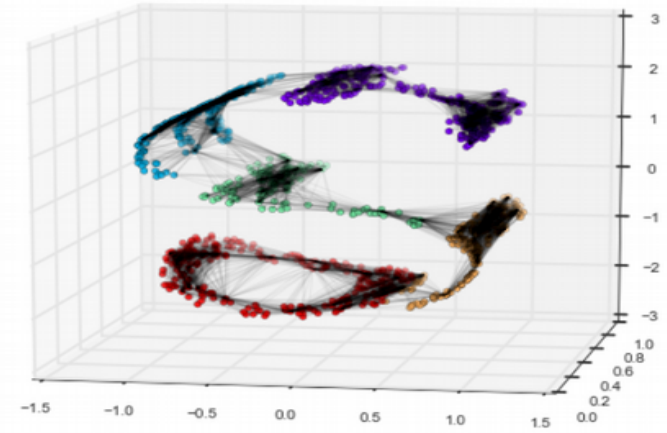




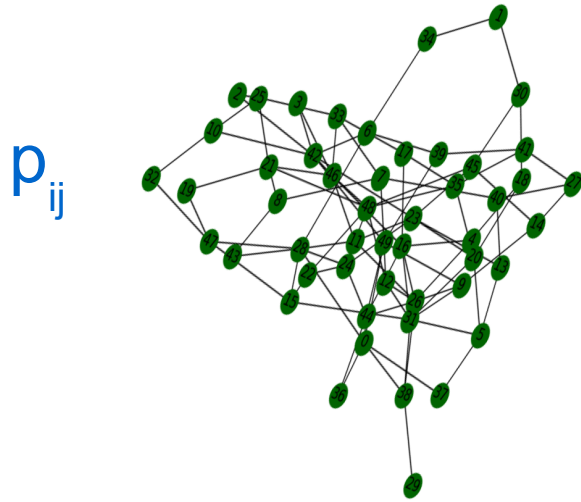
MDS Linkages



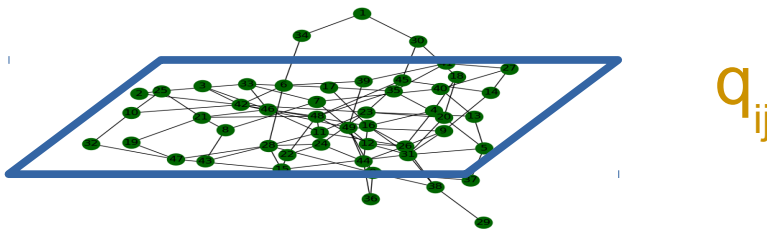
LLE Linkages (100 NN)



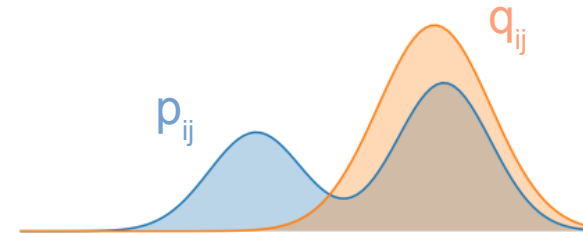
1) Construct high-dimensional graph



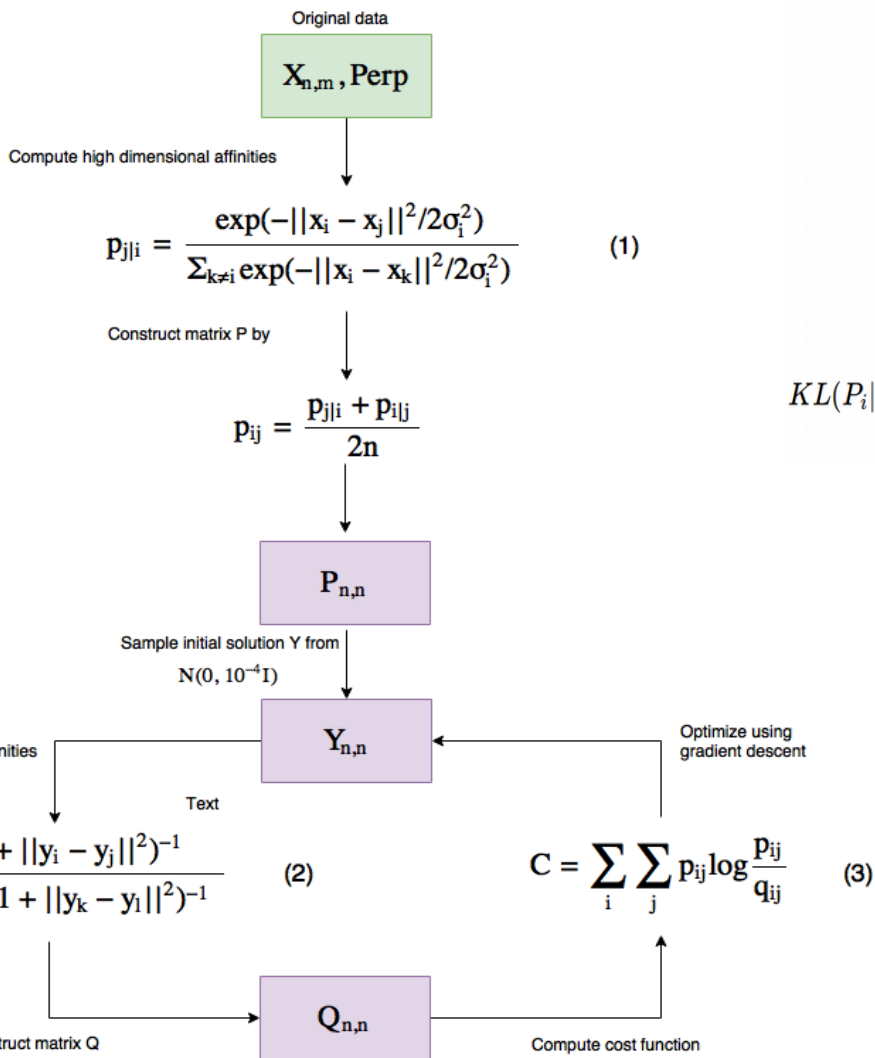
2) Construct low-dimensional graph



3) Collapse the graphs together



Kullback-Leibler divergence

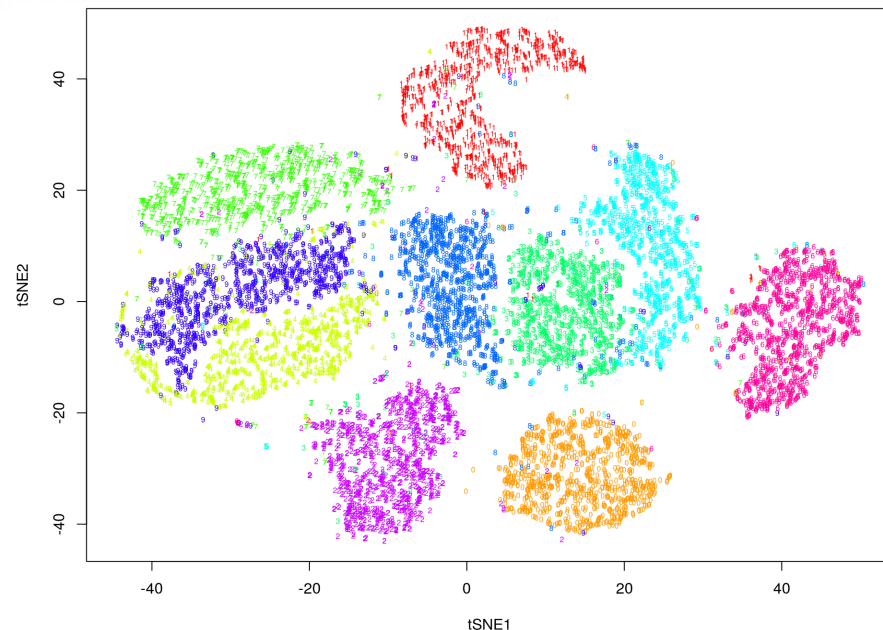


$$p_{j|i} = \frac{\exp(-\|x_i - x_j\|^2 / 2\sigma_i^2)}{\sum_{k \neq i} \exp(-\|x_i - x_k\|^2 / 2\sigma_i^2)}, \quad p_{ij} = \frac{p_{i|j} + p_{j|i}}{2N} \quad (1)$$

$$\text{Perplexity} = 2^{-\sum_j p_{j|i} \log_2 p_{j|i}} \quad (2)$$

$$q_{ij} = \frac{(1 + \|y_i - y_j\|^2)^{-1}}{\sum_{k \neq i} (1 + \|y_k - y_i\|^2)^{-1}} \quad (3)$$

$$KL(P_i || Q_i) = \sum_i \sum_j p_{j|i} \log \frac{p_{j|i}}{q_{j|i}}, \quad \frac{\partial KL}{\partial y_i} = 4 \sum_j (p_{ij} - q_{ij})(y_i - y_j) (1 + \|y_i - y_j\|^2)^{-1} \quad (4)$$

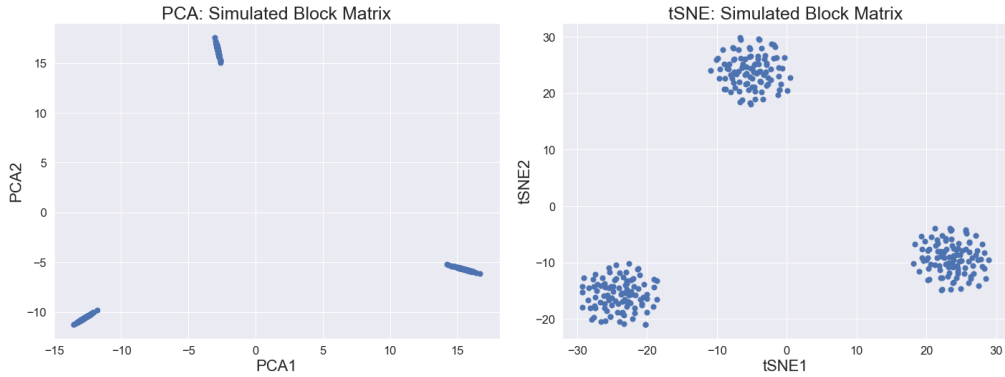
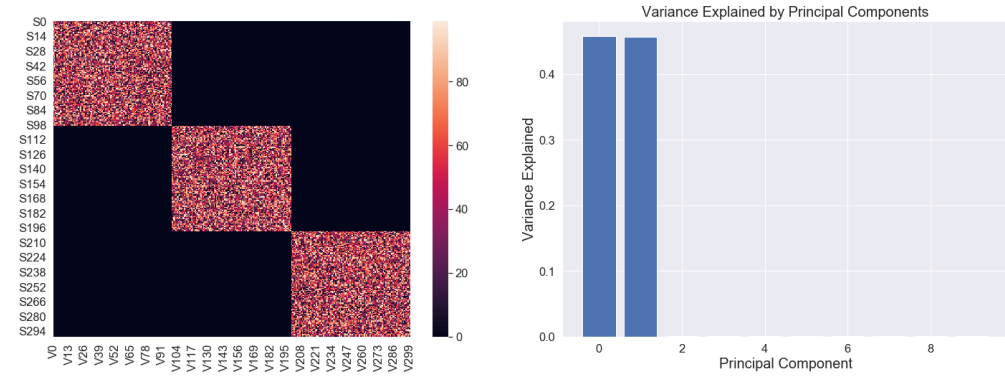




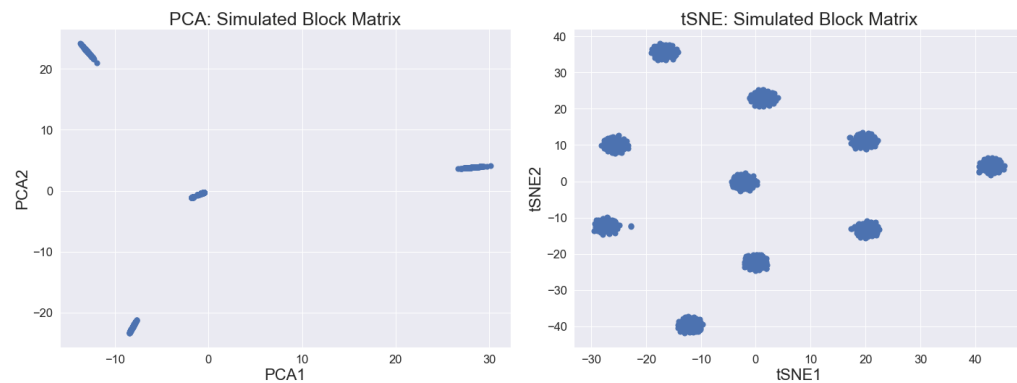
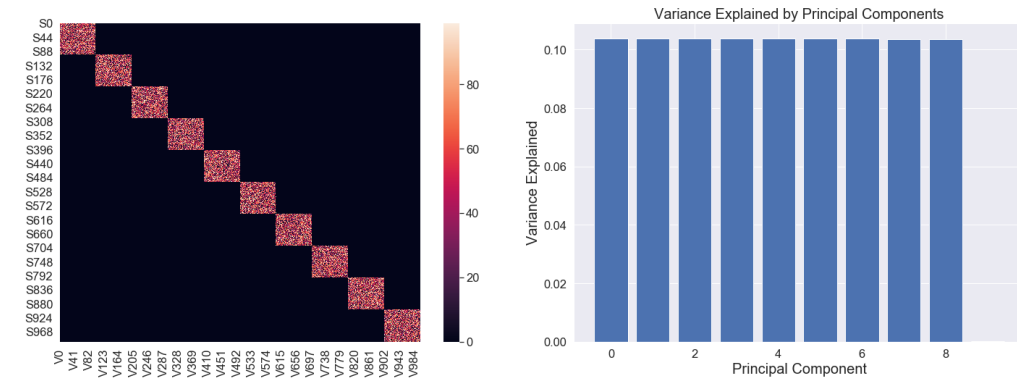
# PCA vs. tSNE

## when number of populations increases

### Three classes of data points

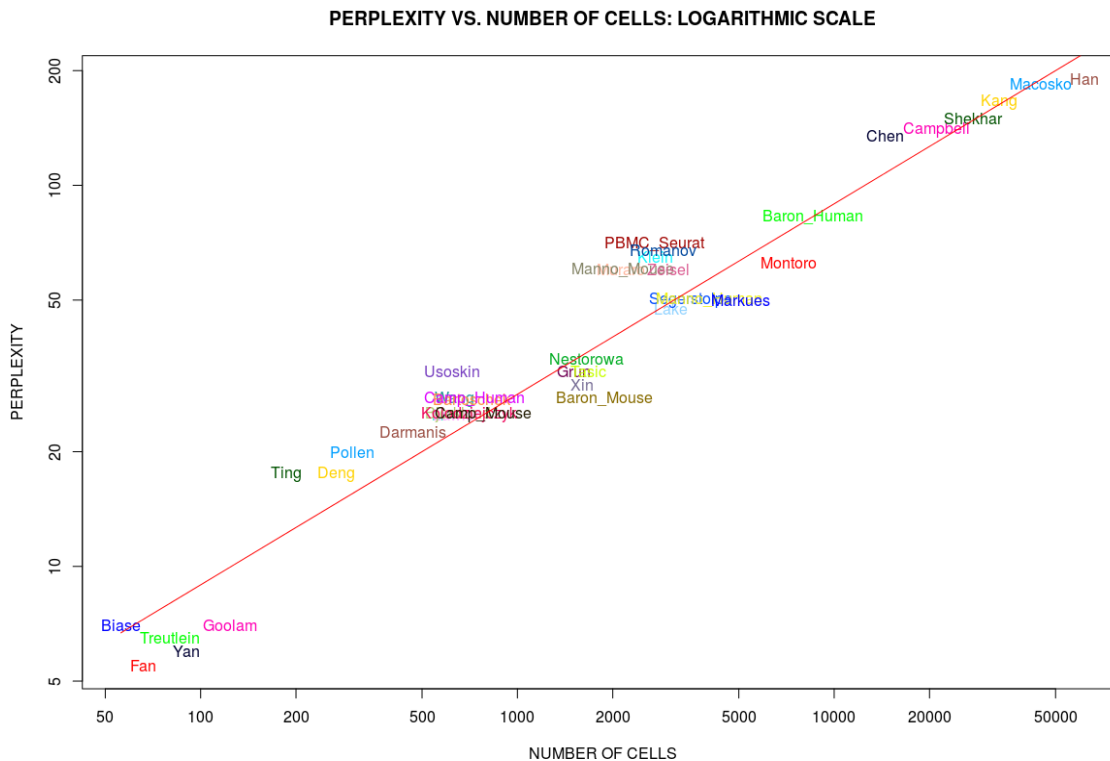


### Ten classes of data points





Van der Maaten: “Loosely speaking, one could say that a larger / denser dataset requires a larger perplexity.”



$$\log(\text{Perp}) = -0.179 + 0.51 \cdot \log(N)$$

$$\text{Perp} \sim N^{(1/2)}$$

tSNE does not scale for large data sets?

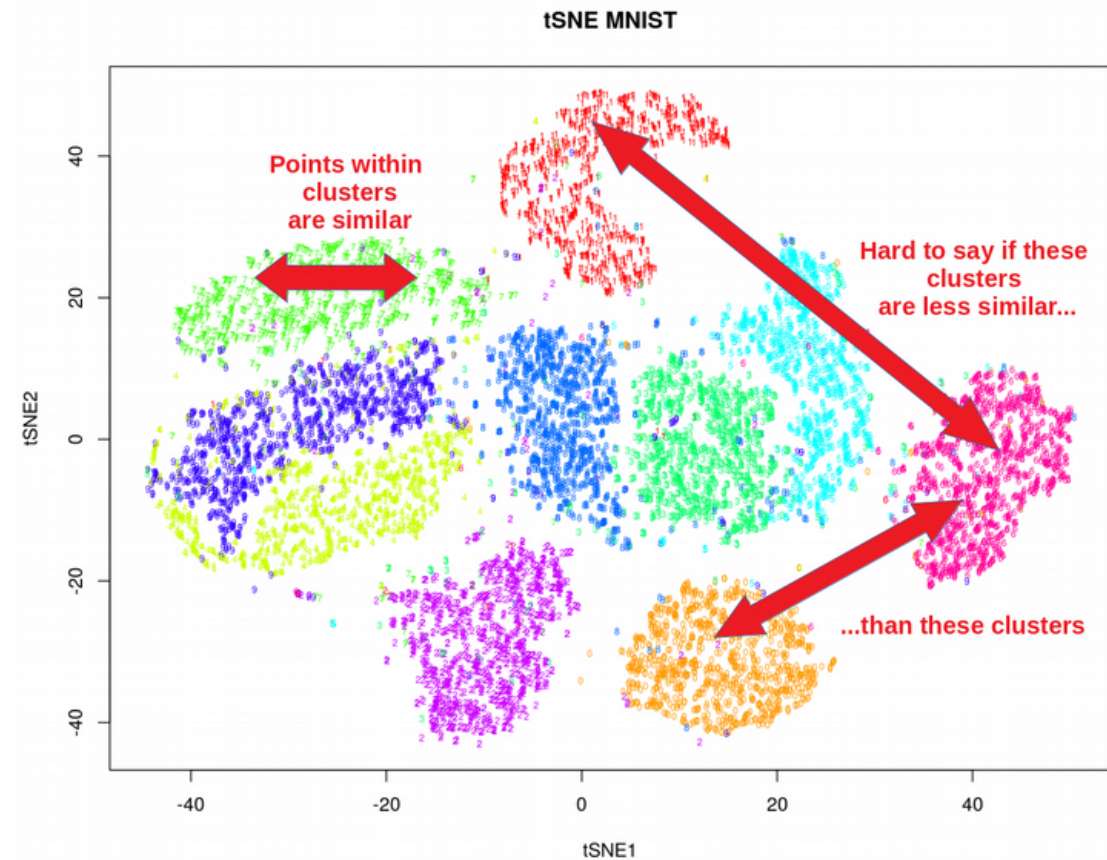
tSNE does not preserve global structure?

tSNE can only embed into 2-3 dims?

tSNE performs non-parametric mapping (no variance explained statistics)?

tSNE can not work with high-dimensional data directly (PCA needed)?

tSNE uses too much RAM at large perp?



UMAP uses local connectivity for high-dim probabilities

$$p_{i|j} = e^{-\frac{d(x_i, x_j) - \rho_i}{\sigma_i}}$$

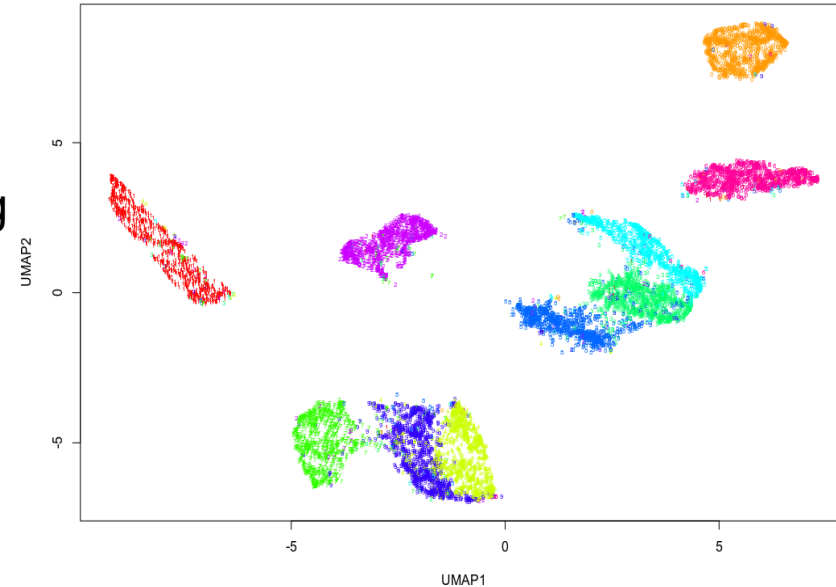
UMAP MNIST

UMAP does not normalize probabilities (speed-up)

UMAP can deliver a number of components for clustering

UMAP uses Laplacian Eigenmap for initialization

UMAP uses Cross-Entropy (not KL) as cost function

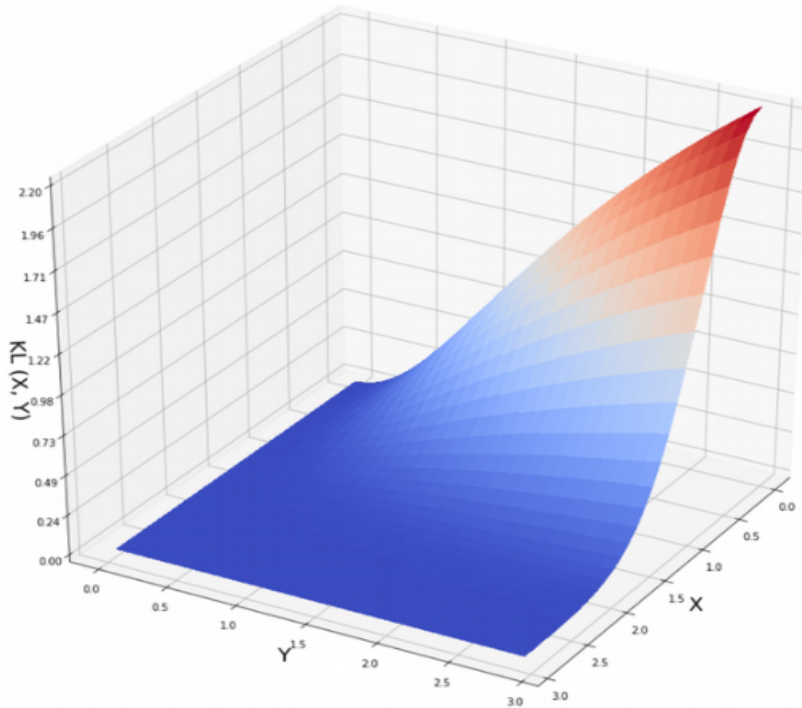


$$CE(X, Y) = \sum_i \sum_j \left[ p_{ij}(X) \log \left( \frac{p_{ij}(X)}{q_{ij}(Y)} \right) + (1 - p_{ij}(X)) \log \left( \frac{1 - p_{ij}(X)}{1 - q_{ij}(Y)} \right) \right]$$

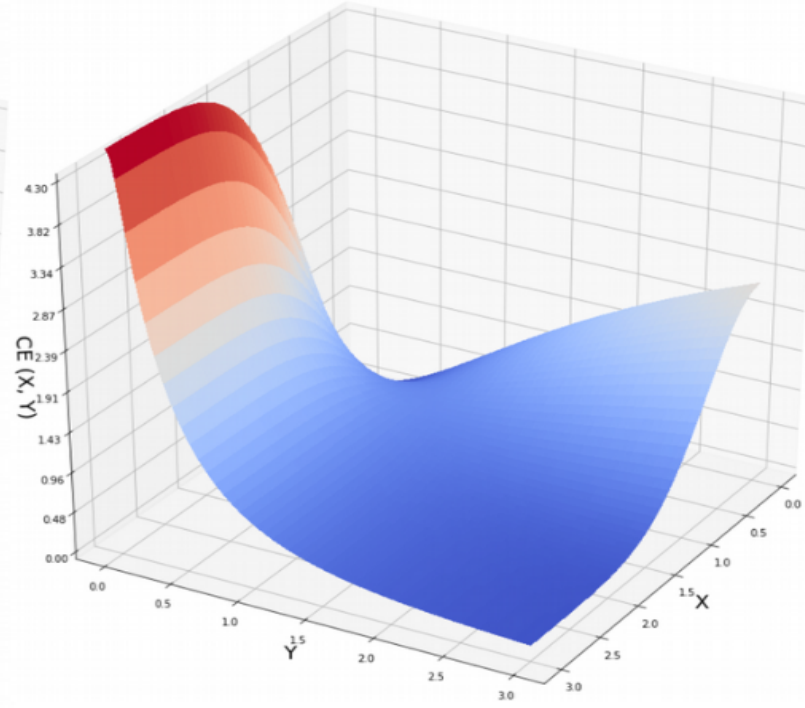
This is similar to tSNE cost function

This term is UMAP specific

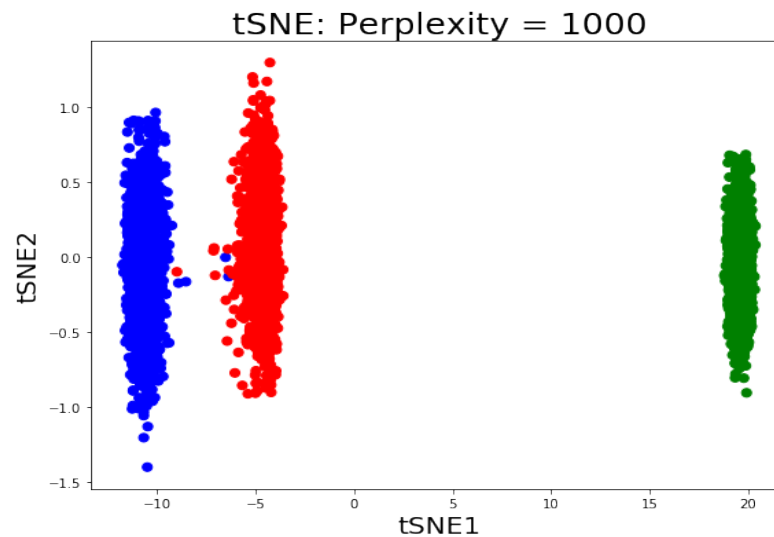
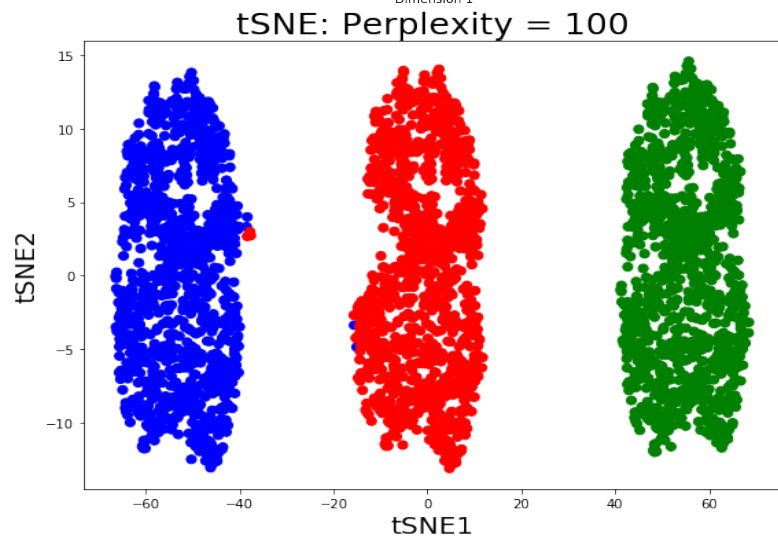
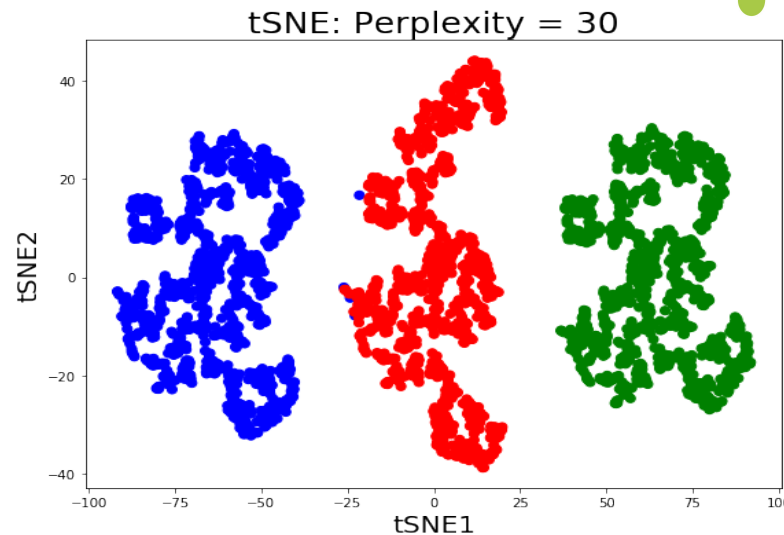
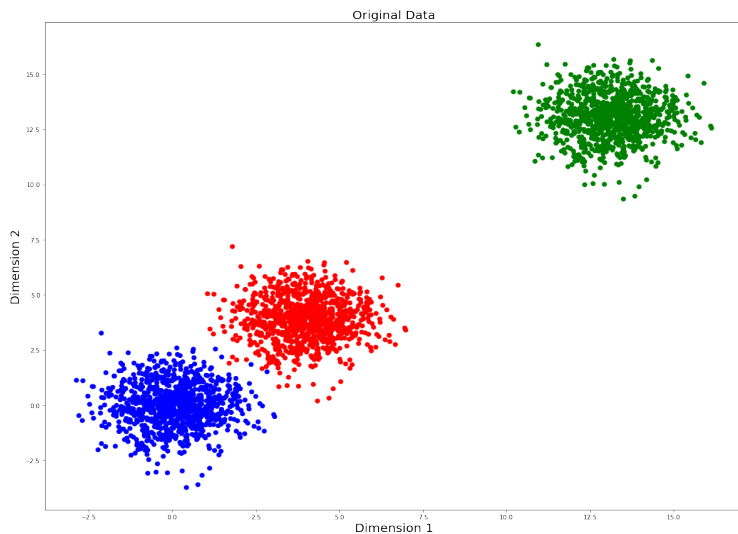
# **tSNE vs. UMAP: global structure preservation**



$X \rightarrow \text{infinity}, Y \text{ can be any}$

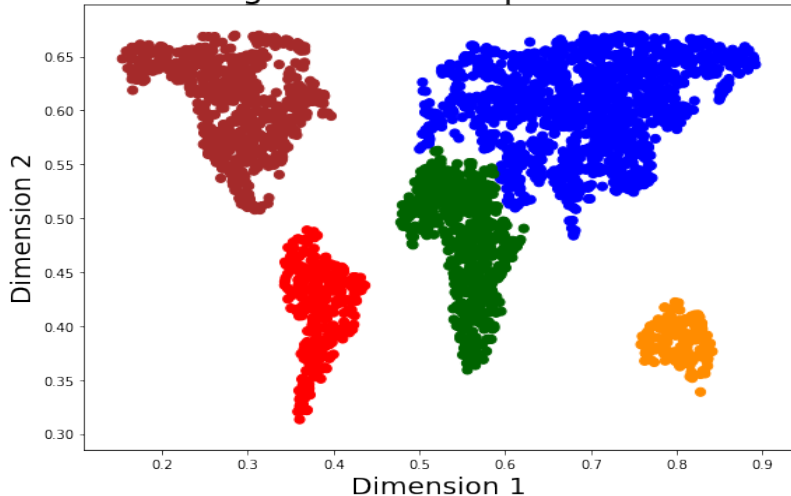


$X \rightarrow \text{infinity}, Y \rightarrow \text{infinity}$

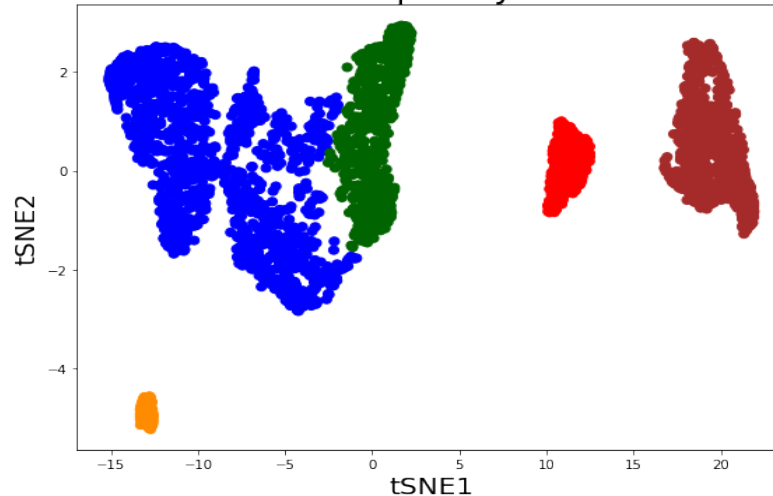


Can large perplexity solve the problem of global structure for tSNE?

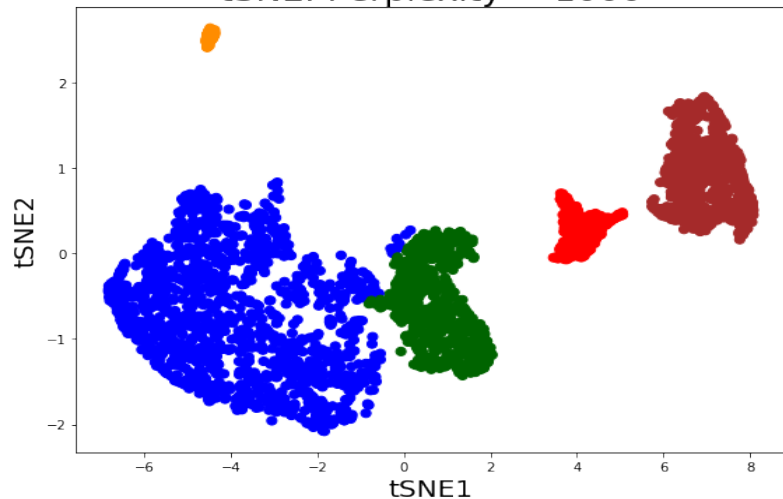
Original World Map Data Set



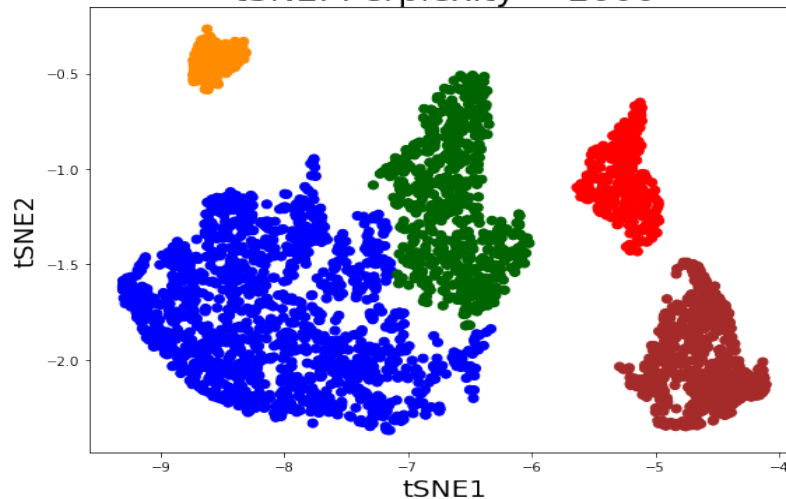
tSNE: Perplexity = 500



tSNE: Perplexity = 1000

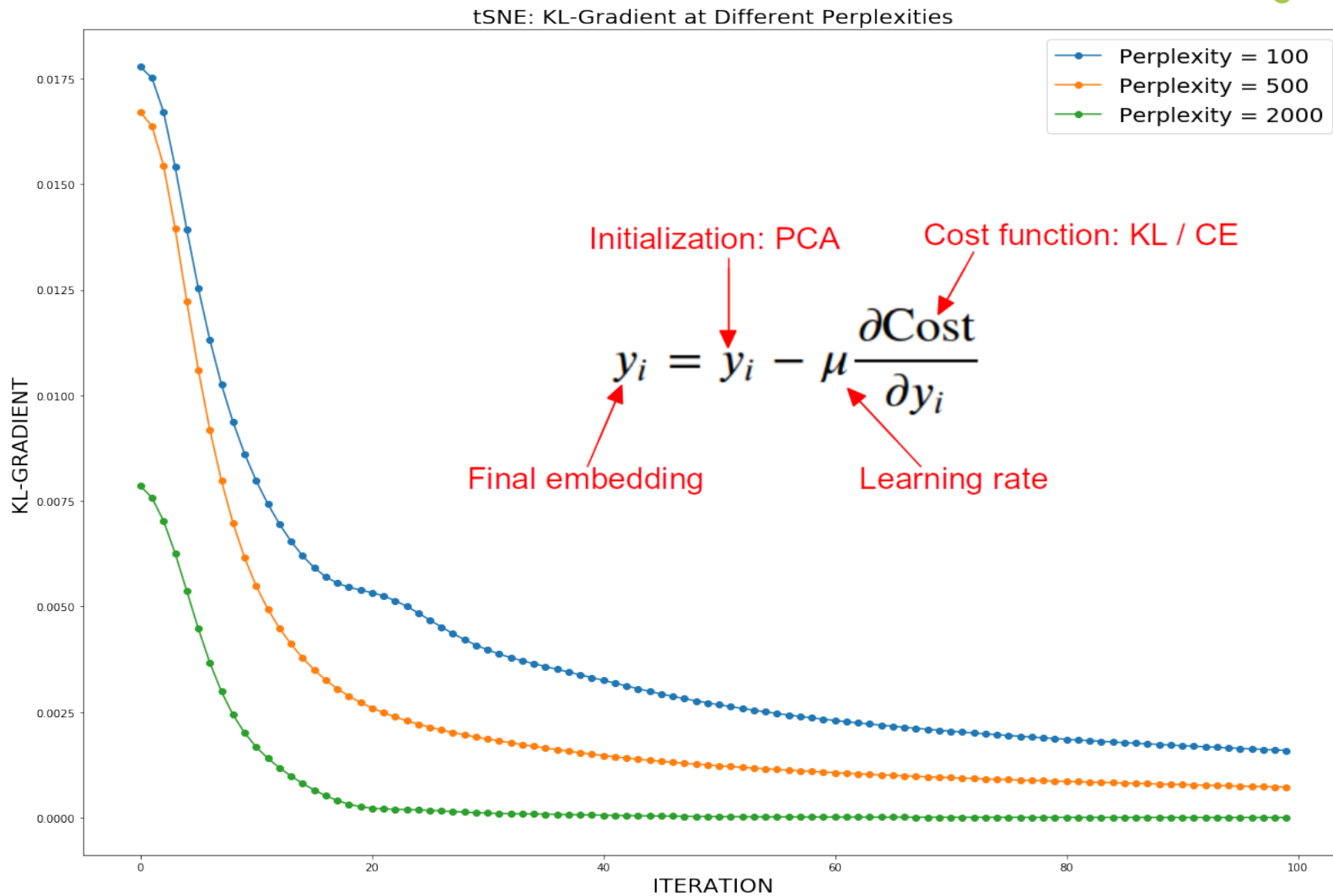


tSNE: Perplexity = 2000

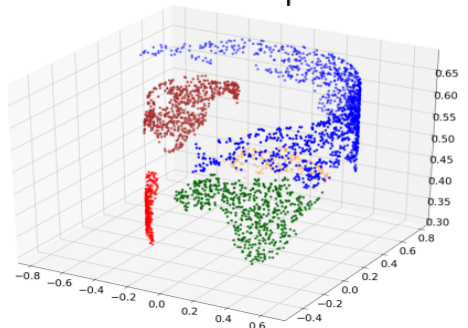


Can large perplexity solve the problem of global structure for tSNE?

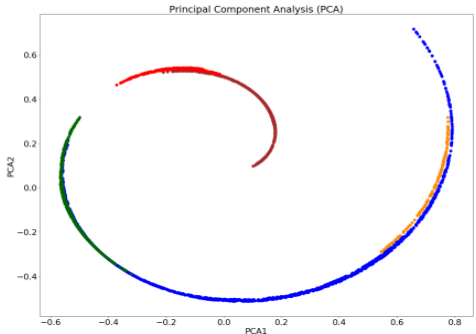




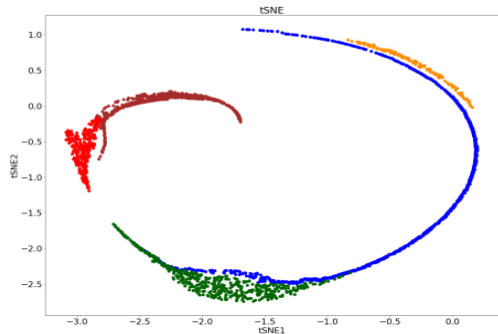
Swiss Roll: 3023 points



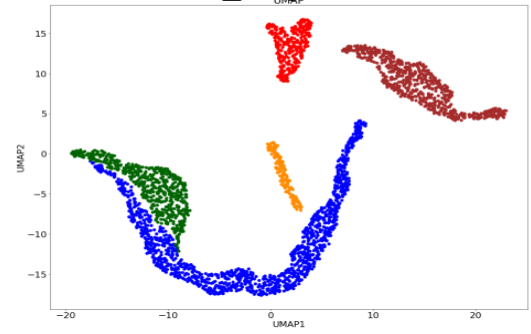
Principal Component Analysis



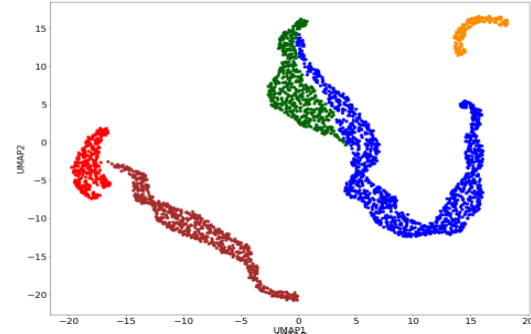
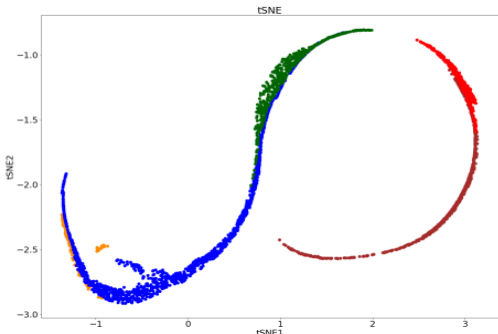
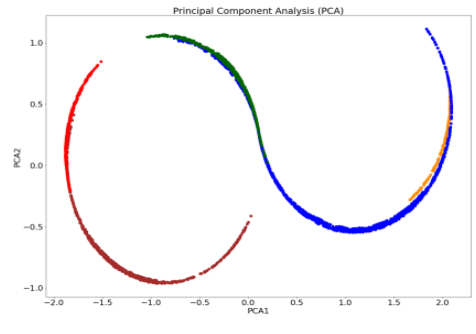
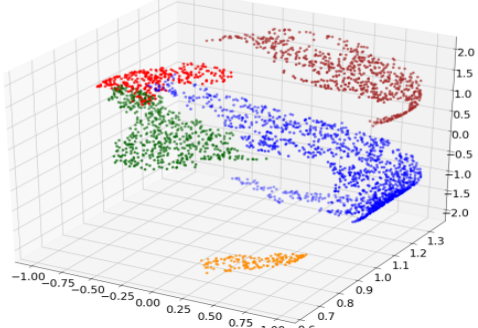
tsNE: perplexity = 2000



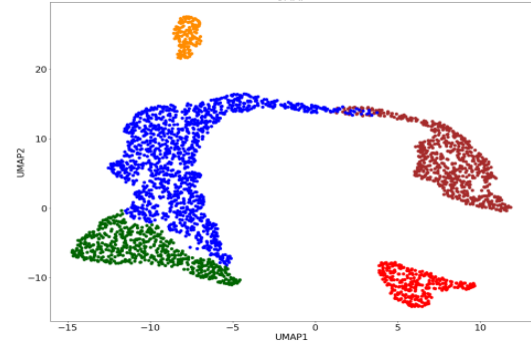
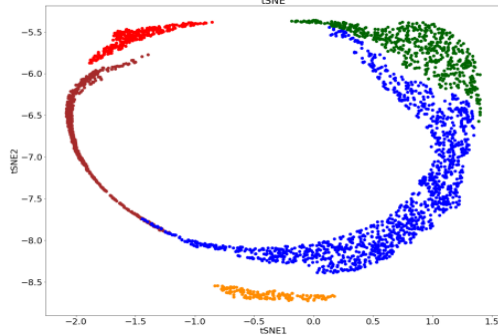
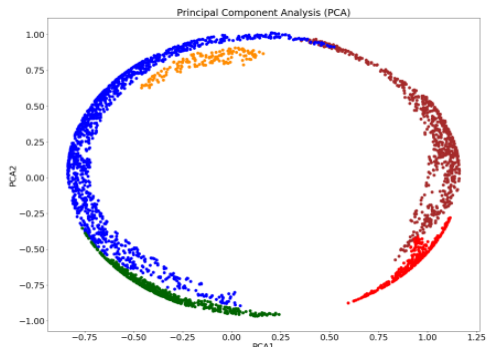
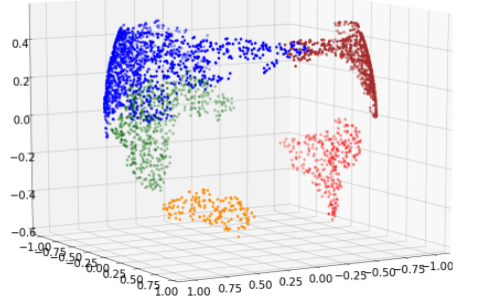
UMAP: n\_neighbor = 2000



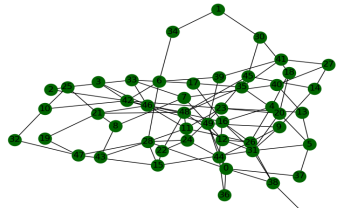
S-shape: 3023 points



Sphere: 3023 points



Graph Laplacian, Laplacian Eigenmap, spectral clustering, diffusion maps, spectral dimension reduction methods etc.



$$s(x_i, x_j) = \exp(-\alpha \|x_i - x_j\|^2)$$



$$L = I - D^{-1} * S$$

$$L * u = \lambda * u$$

Laplacian Eigenmap

$$P = D^{-1} * S$$

$$P^t * u = \lambda^t * u$$

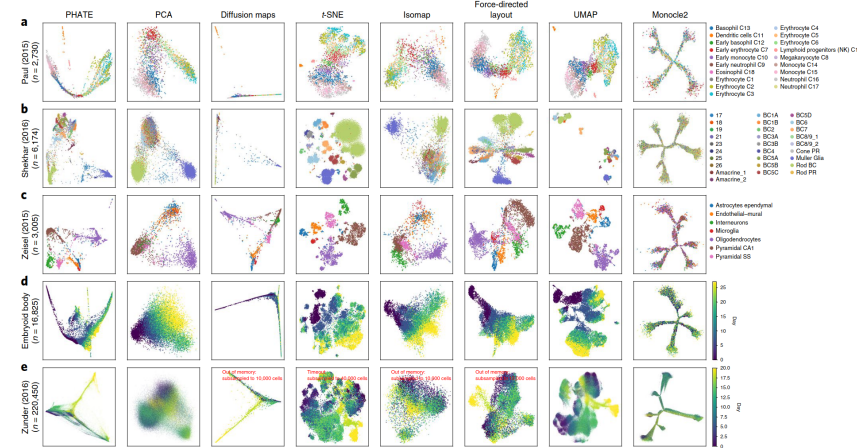
Diffusion Maps

**S =**

	[,1]	[,2]	[,3]	[,4]	[,5]	[,6]	[,7]
[1,]	1.00000000	0.7429016	0.6319343	0.0000000	0.0000000	0.0000000	0.0000000
[2,]	0.00000000	1.00000000	0.0000000	0.0000000	0.0000000	0.0000000	0.0000000
[3,]	0.7429016	0.0000000	1.0000000	0.0000000	0.0000000	0.0000000	0.6657756
[4,]	0.6319343	0.0000000	0.0000000	1.0000000	0.0000000	0.0000000	0.7195922
[5,]	0.00000000	0.00000000	0.00000000	0.00000000	1.00000000	0.7765565	0.0000000
[6,]	0.00000000	0.00000000	0.00000000	0.00000000	0.7765565	1.0000000	0.0000000
[7,]	0.00000000	0.00000000	0.6657756	0.7195922	0.0000000	0.0000000	1.0000000
[8,]	0.00000000	0.00000000	0.00000000	0.00000000	0.00000000	0.00000000	0.00000000

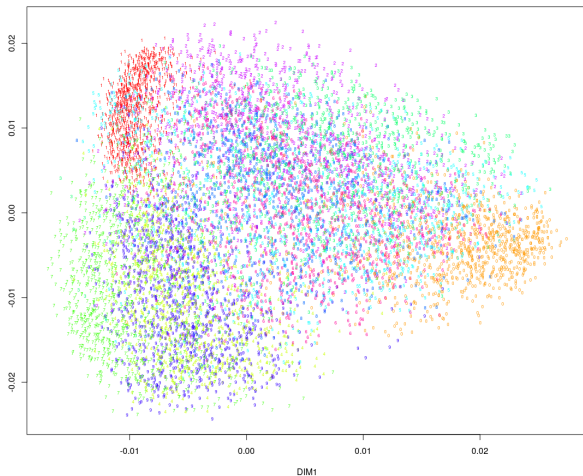
**D =**

	[,1]	[,2]	[,3]	[,4]	[,5]	[,6]	[,7]
[1,]	2.374836	0.000000	0.000000	0.000000	0.000000	0.000000	0.000000
[2,]	0.000000	2.597451	0.000000	0.000000	0.000000	0.000000	0.000000
[3,]	0.000000	0.000000	2.408677	0.000000	0.000000	0.000000	0.000000
[4,]	0.000000	0.000000	0.000000	2.351526	0.000000	0.000000	0.000000
[5,]	0.000000	0.000000	0.000000	0.000000	2.523175	0.000000	0.000000
[6,]	0.000000	0.000000	0.000000	0.000000	0.000000	2.519936	0.000000
[7,]	0.000000	0.000000	0.000000	0.000000	0.000000	0.000000	3.170424
[8,]	0.000000	0.000000	0.000000	0.000000	0.000000	0.000000	0.000000

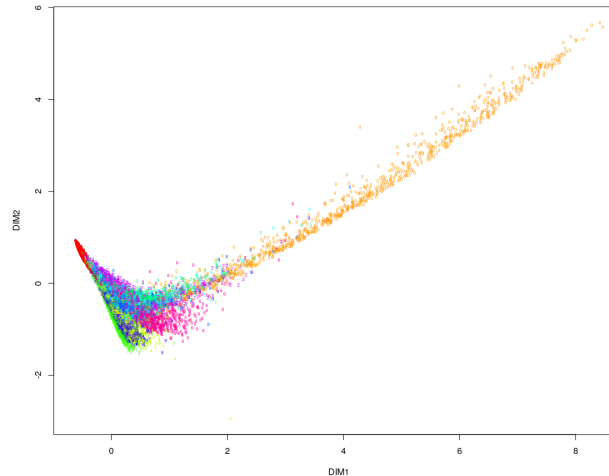


Moon et al., Nat Biotechnol. 2019; 37(12):1482-1492

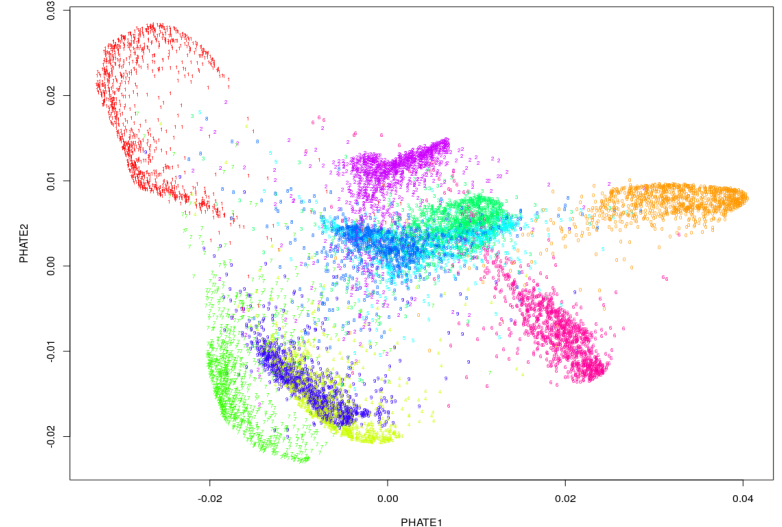
MNIST LAPLACIAN EIGENMAP



MNIST DIFFUSION MAP



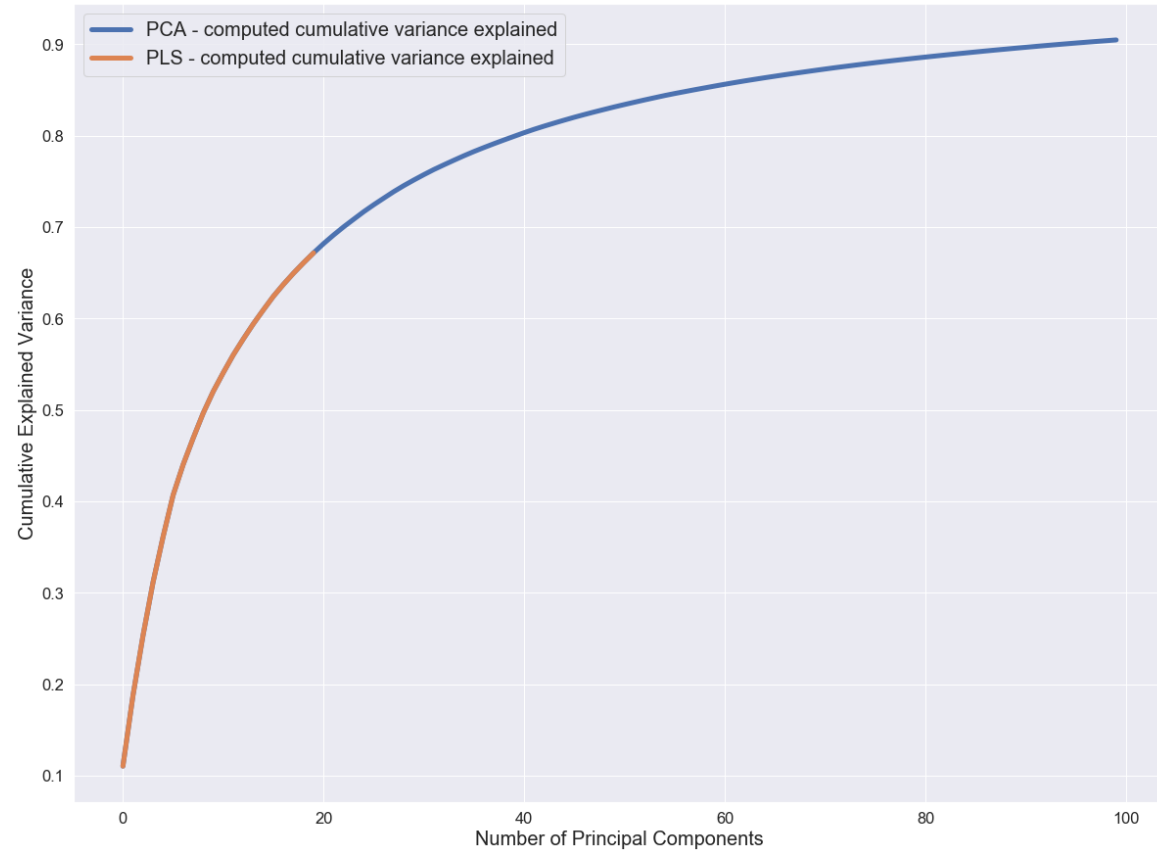
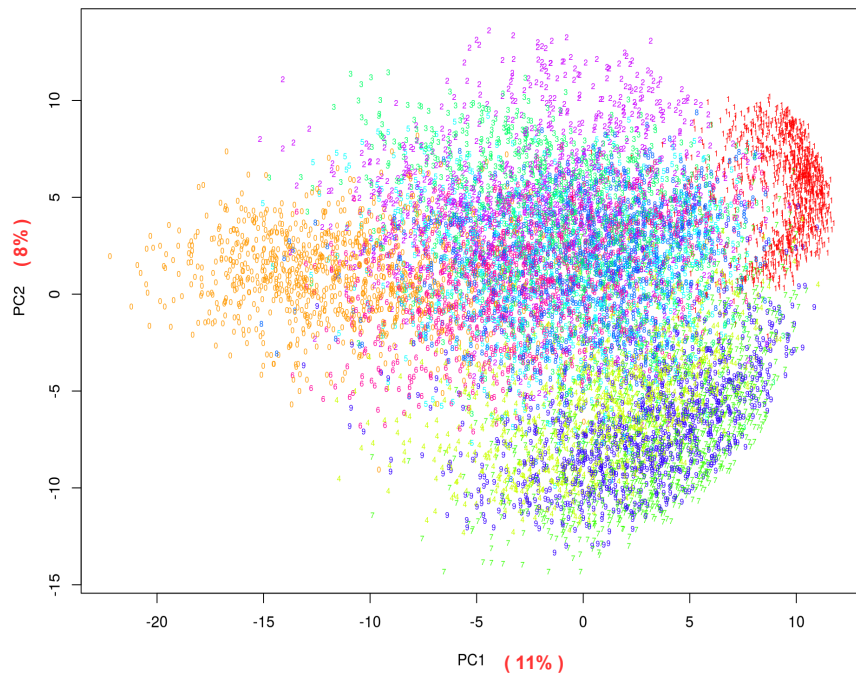
PHATE PLOT



# Variance explained by PCA, tSNE and UMAP

$$\mathbf{X} = \alpha + \beta \text{PCA}_{\text{matrix}} + \epsilon$$

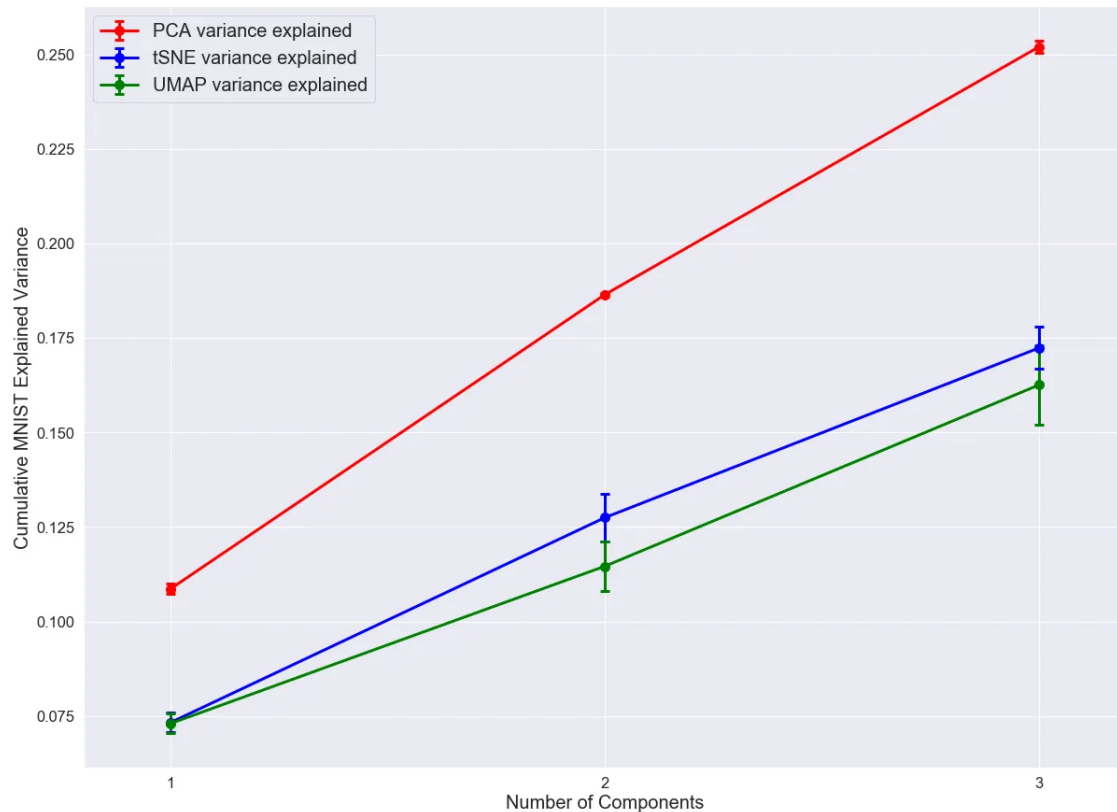
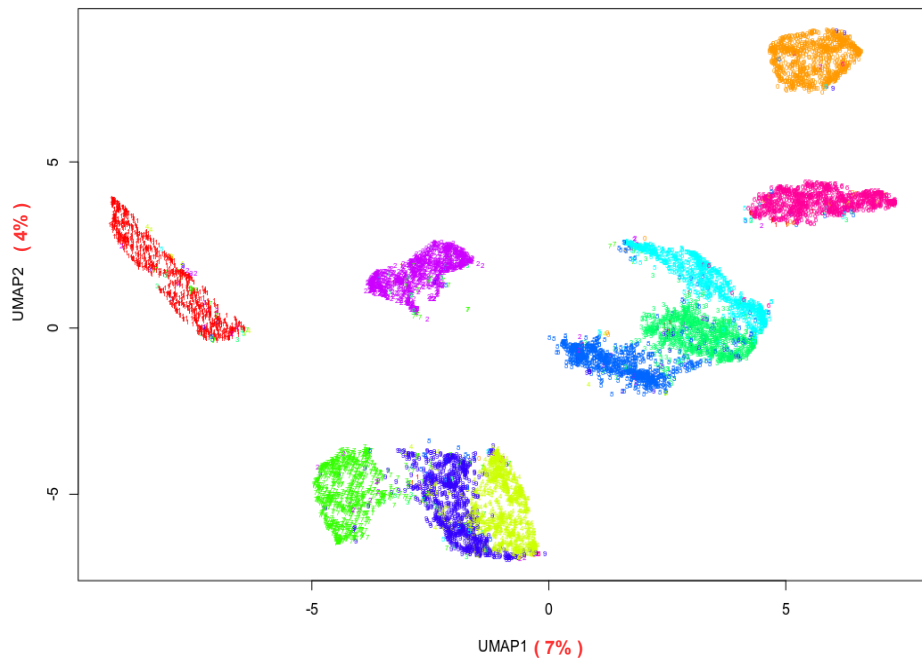
$$R^2 = 1 - \frac{\|\mathbf{X} - \mathbf{B} * \text{PCA}_{\text{matrix}}\|^2}{\|\mathbf{X}\|^2}$$



$$X = \alpha + \beta \text{UMAP}_{\text{matrix}} + \epsilon$$

$$R^2 = 1 - \frac{\|X - B * \text{UMAP}_{\text{matrix}}\|^2}{\|X\|^2}$$

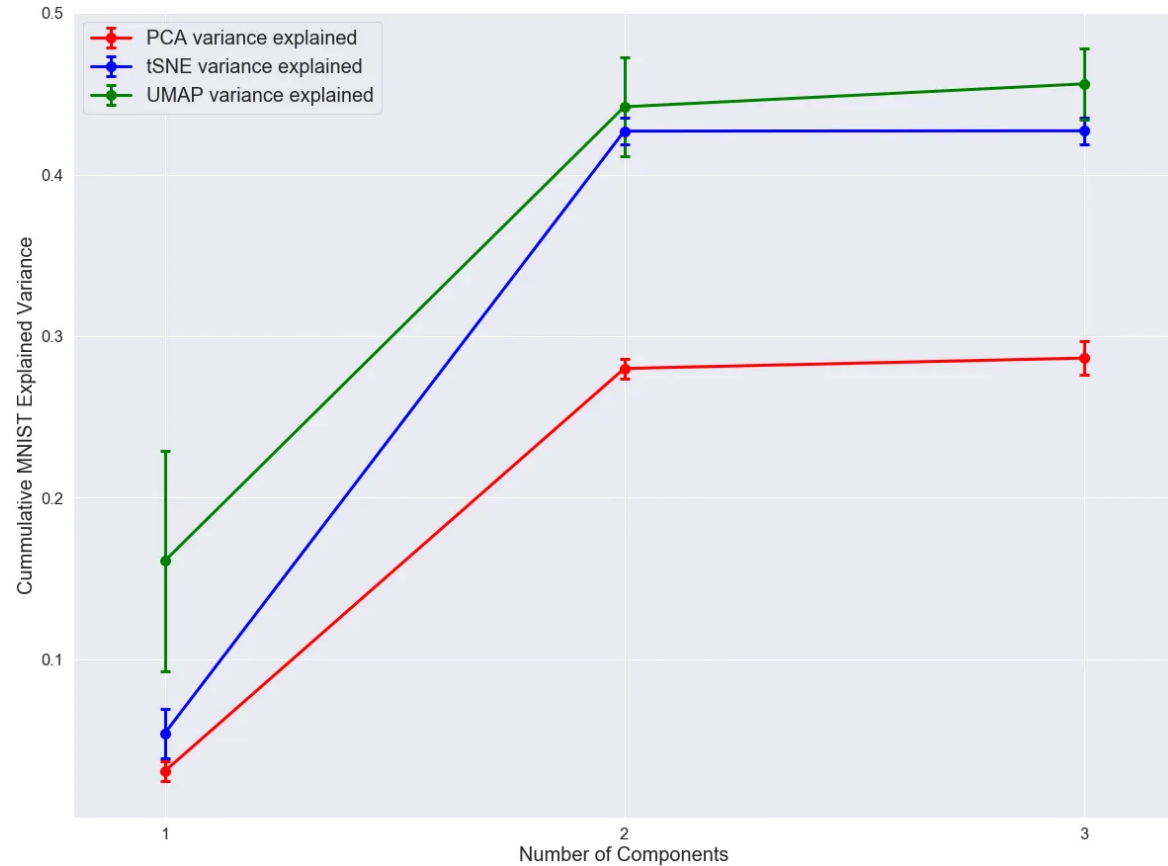
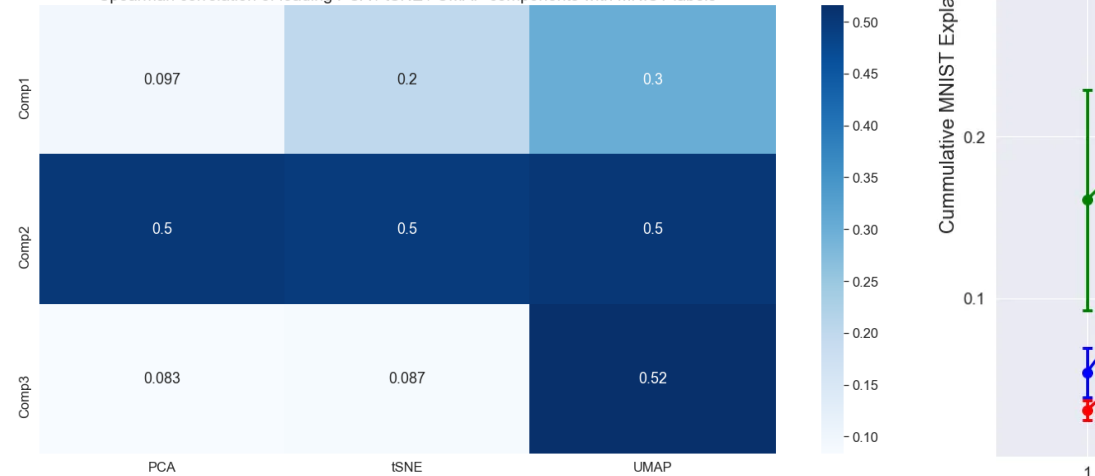
UMAP MNIST



$$\text{labels} = \alpha + \beta \text{UMAP}_{\text{matrix}} + \epsilon$$

$$R^2 = 1 - \frac{\|\text{labels} - B * \text{UMAP}_{\text{matrix}}\|^2}{\|\text{labels}\|^2}$$

Spearman correlation of leading PCA / tSNE / UMAP components with MNIST labels







*Knut och Alice  
Wallenbergs  
Stiftelse*



**LUNDS**  
UNIVERSITET