



Machine learning principles and communications for material scientists

Lecture 2: Unsupervised learning

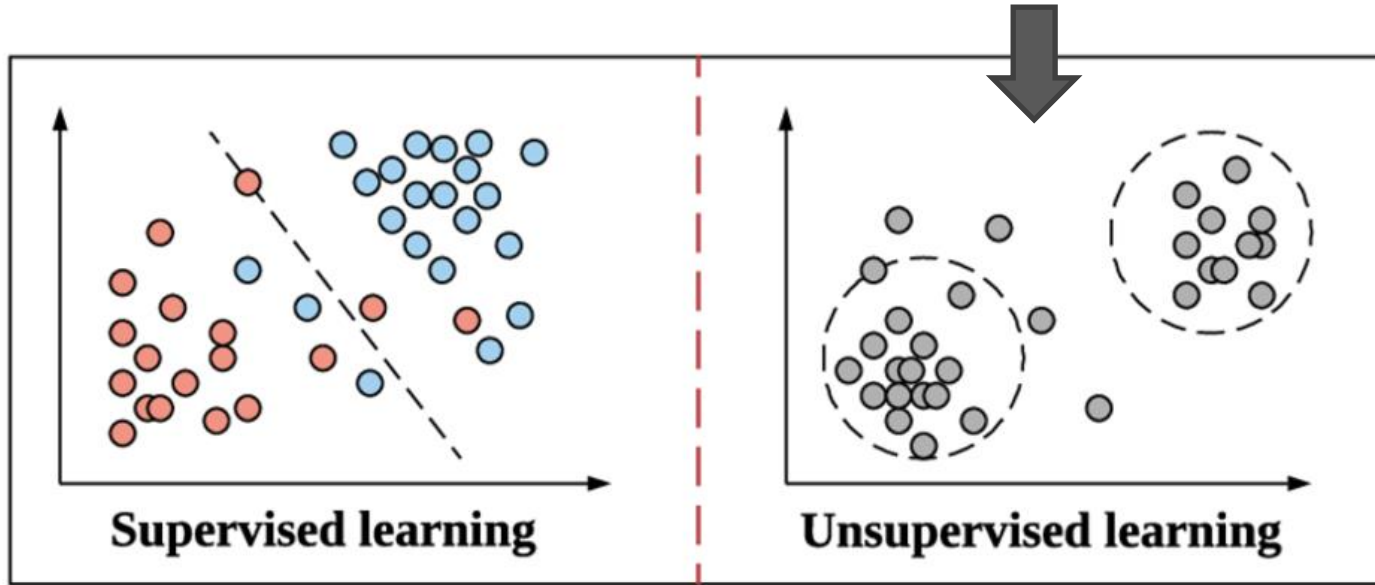
September 12th, 2022



Sira Sriswasdi, PhD

- Research Affairs
- Center of Excellence in Computational Molecular Biology (CMB)
- Center for Artificial Intelligence in Medicine (CU-AIM)

Machine learning paradigms

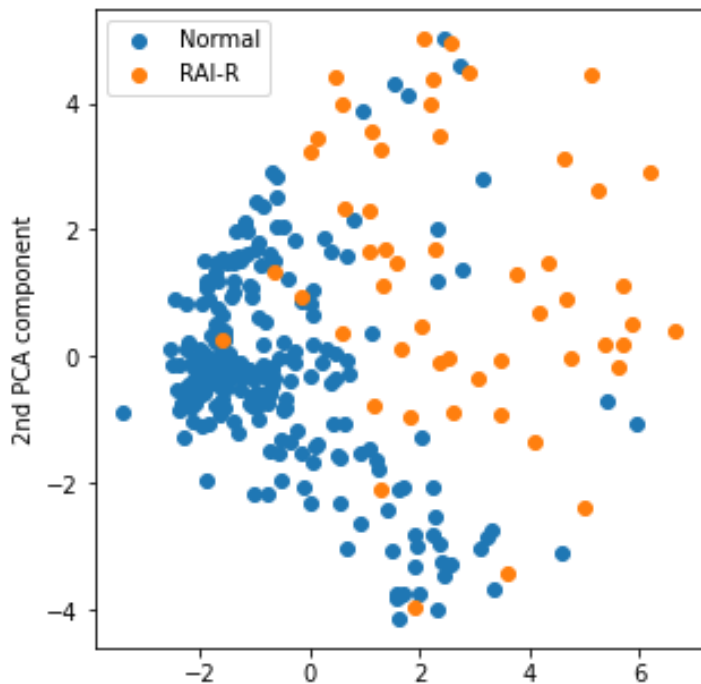


Qian, B. et al. "Orchestrating the Development Lifecycle of Machine Learning-Based IoT Applications: A Taxonomy and Survey"

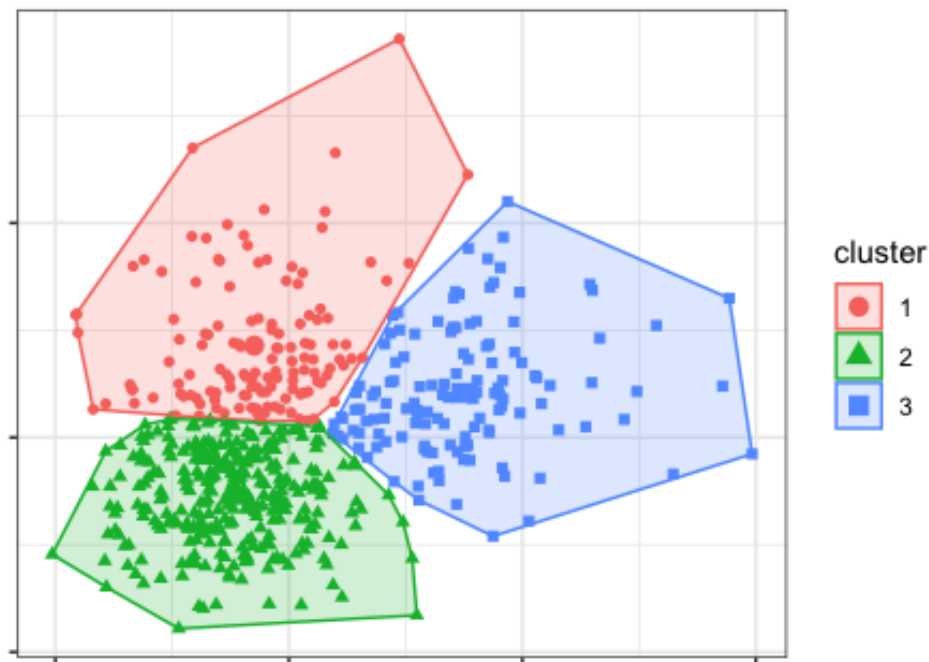
- Identify **robust patterns** that can be **generalized to new data**

Two primary branches of unsupervised learning

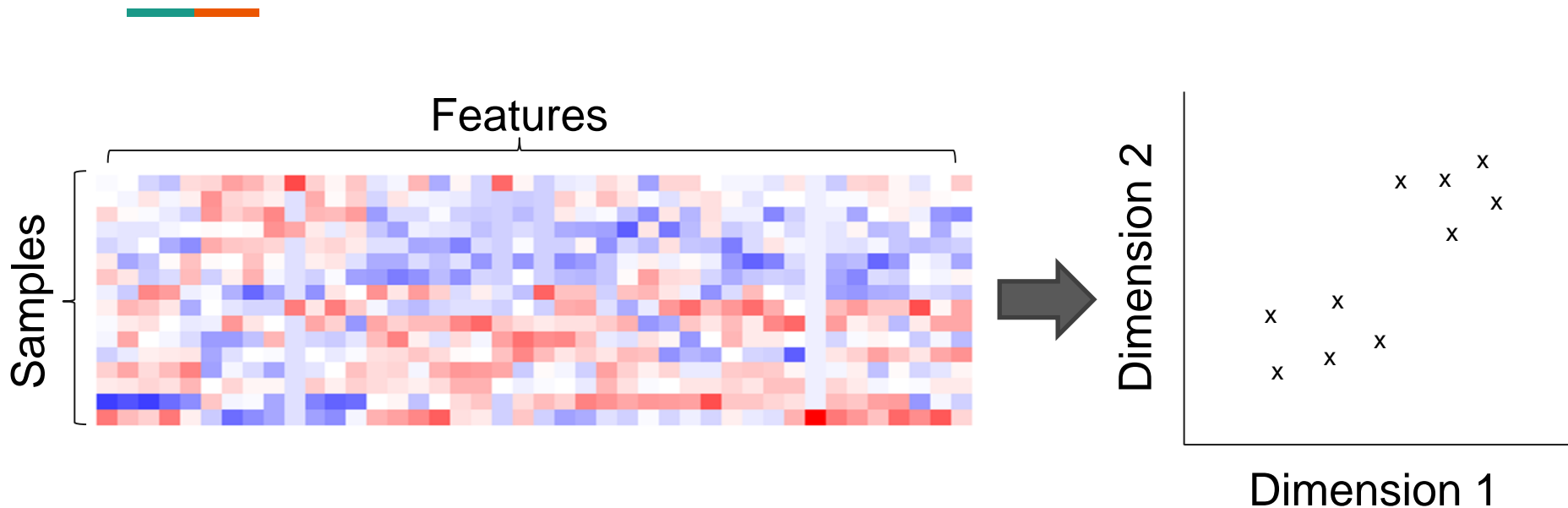
Dimensionality Reduction



Clustering

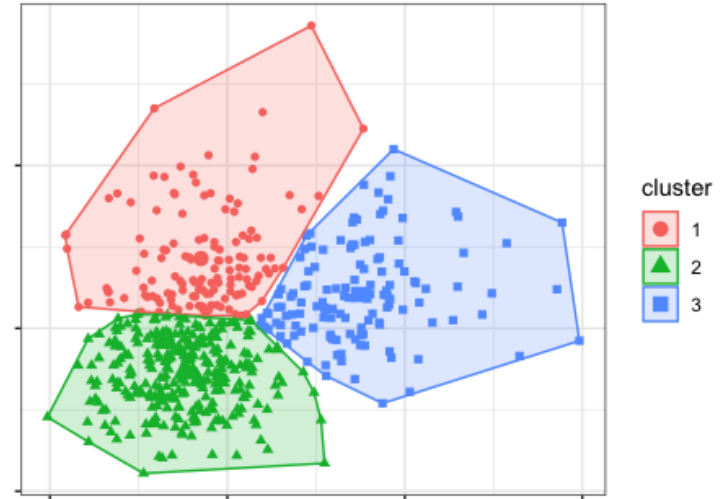
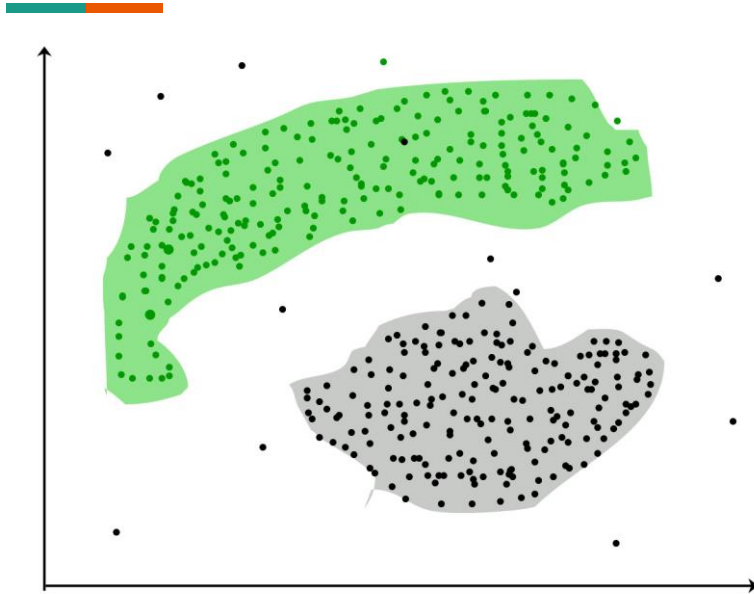


Dimensionality reduction



- Understand data distribution & gauge the difficulty of supervised learning
- Visualize on high-dimensional data on 2D-3D

Clustering

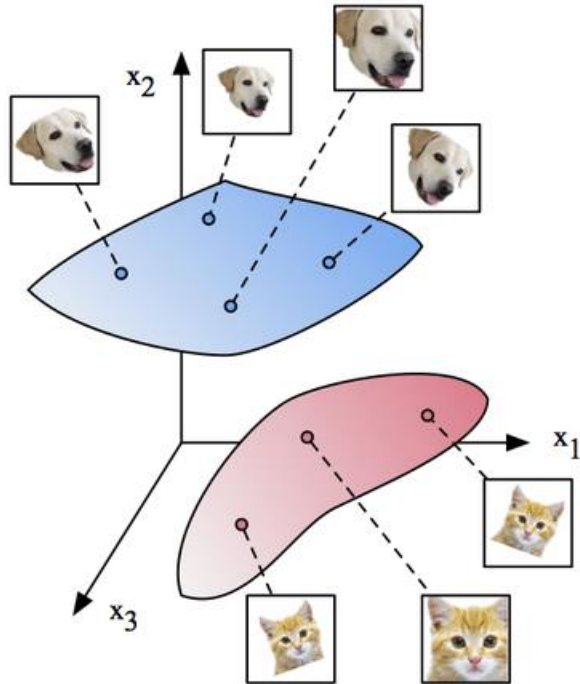


- Identify data subgroups
 - Predict shared characteristics
 - Generate hypothesis

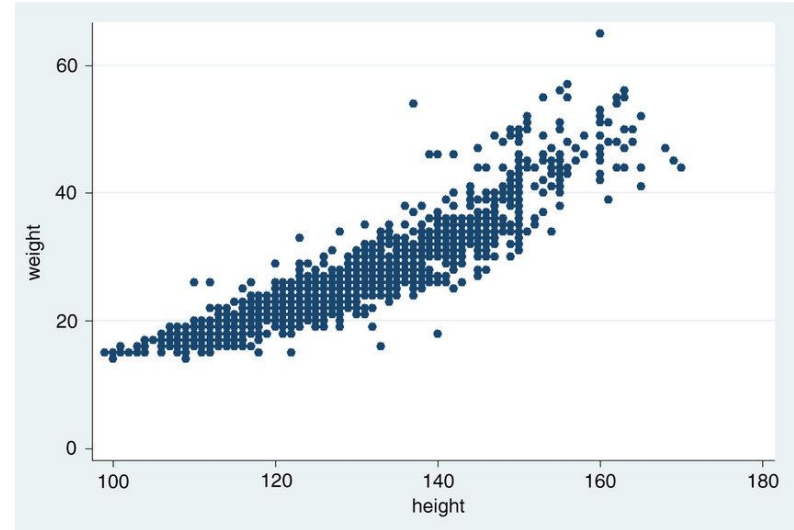


Dimensionality reduction

Manifold hypothesis



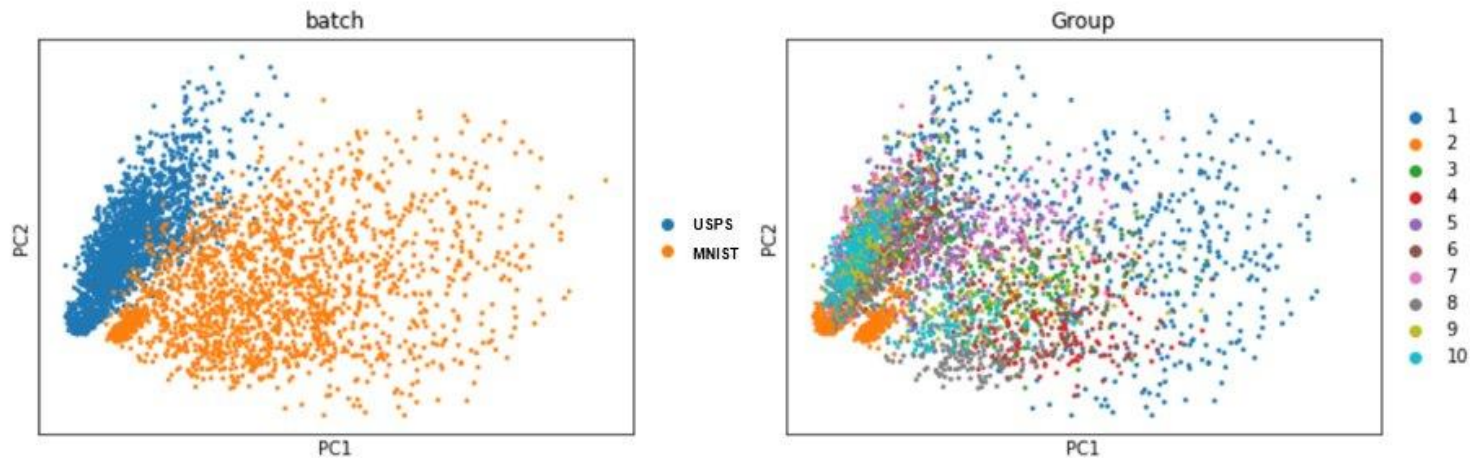
“Real-world, high-dimensional data lie on some low-dimensional manifolds”



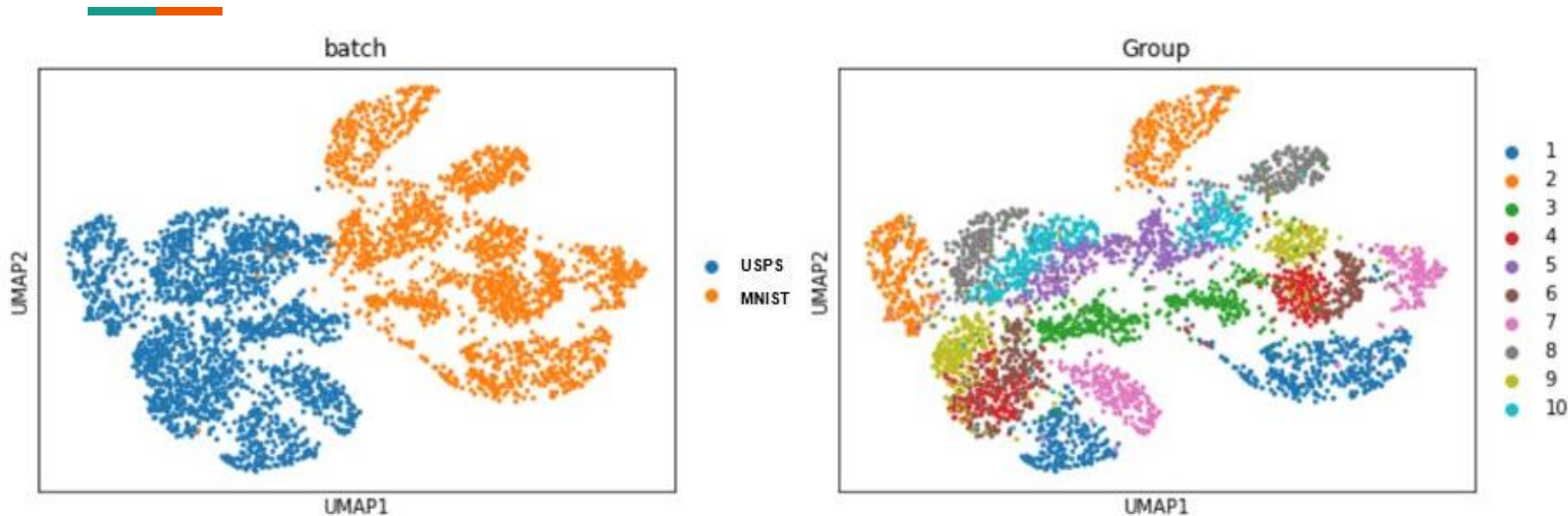
An example: Digit datasets



*adapted from doi:10.1109/TKDE.2017.2669193



A powerful 2D visualization

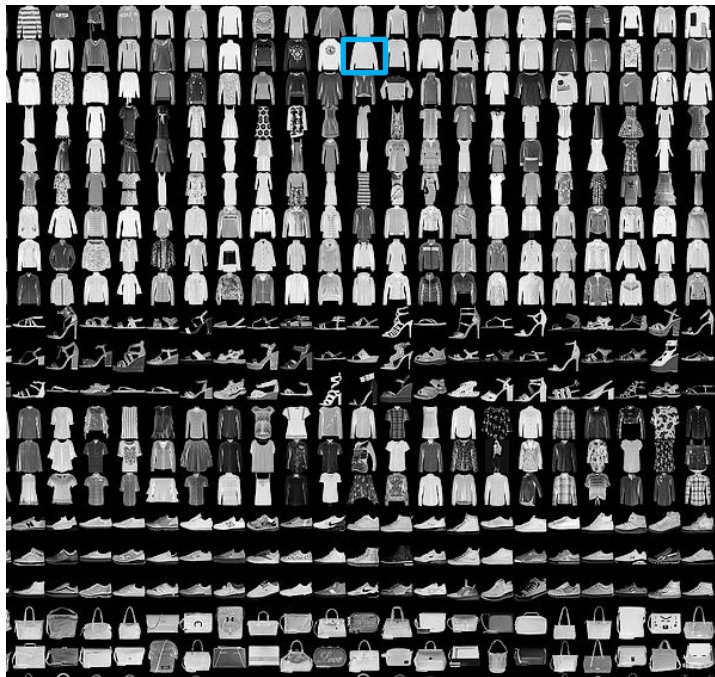


<https://twitter.com/lkmklsmn/status/1436357177887895555>

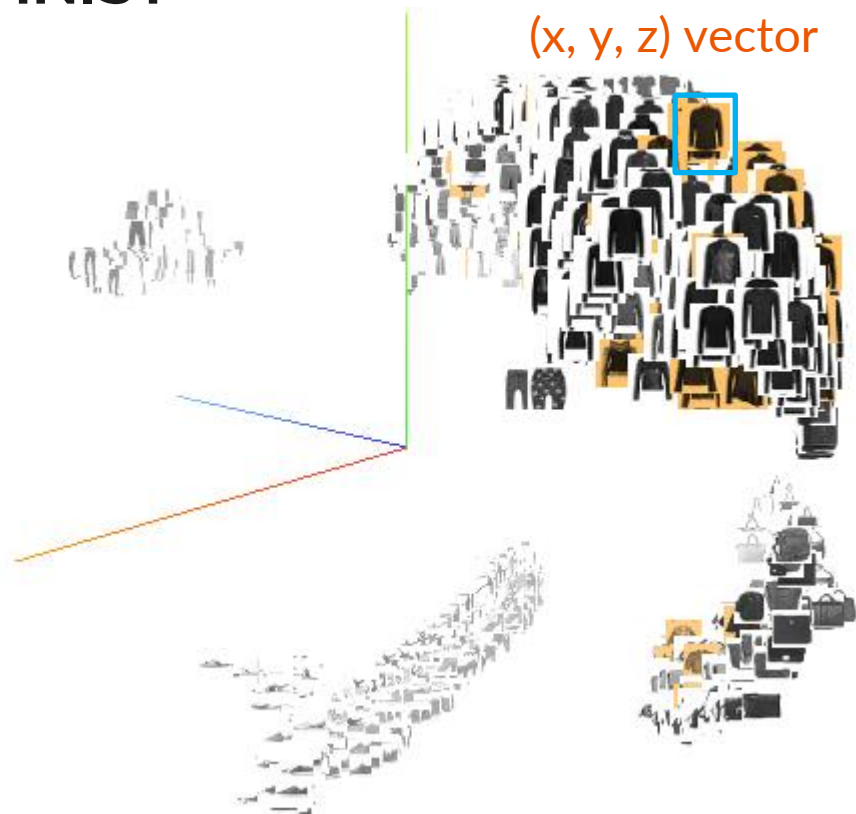
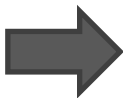
- Both data source and digit identity can be distinguished

Another example: Fashion MNIST

28x28 pixel image

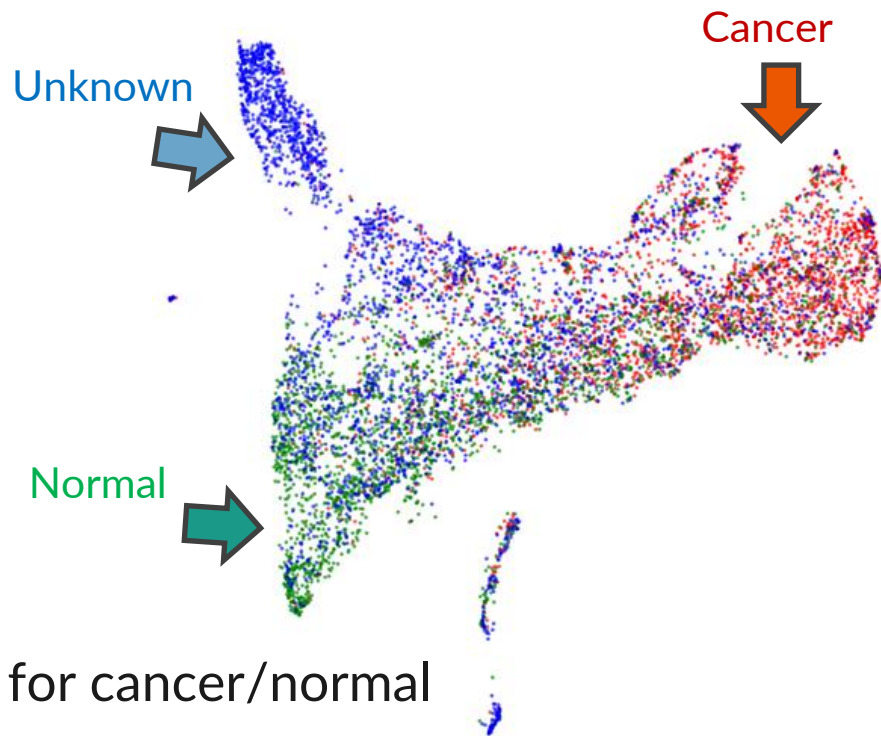
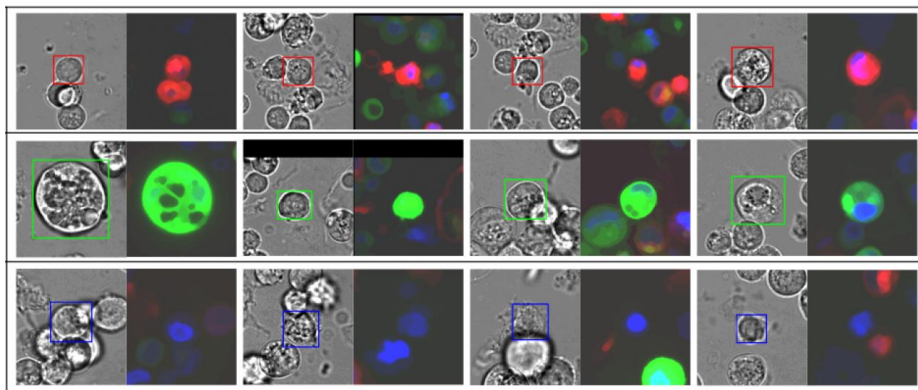


<https://github.com/zalando-research/fashion-mnist>



<https://pair-code.github.io/understanding-umap/>

2D visualization for cell images

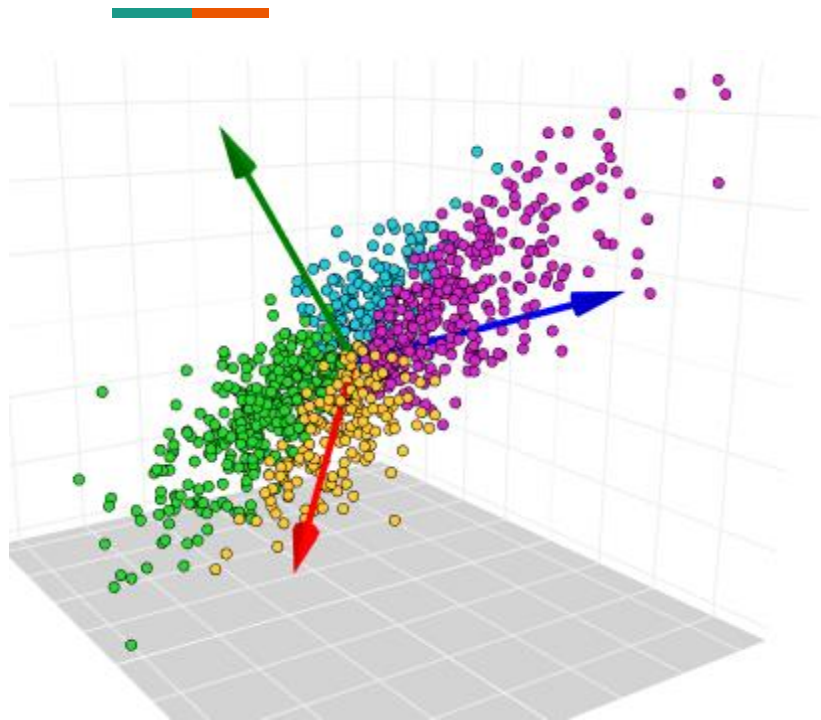


- Extracted from **deep learning model** for cancer/normal cell type classification

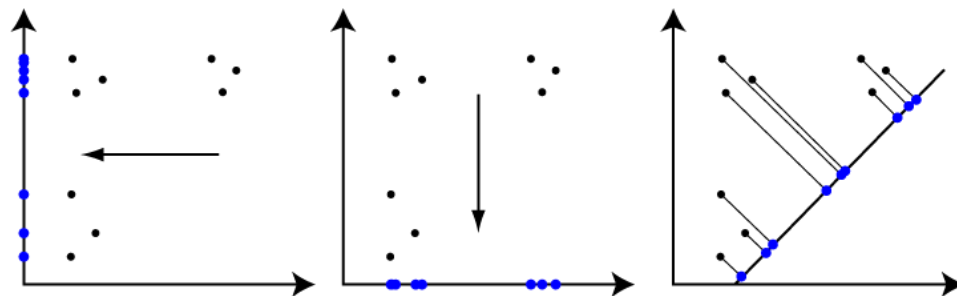


Principal component analysis (PCA)

Variance is information



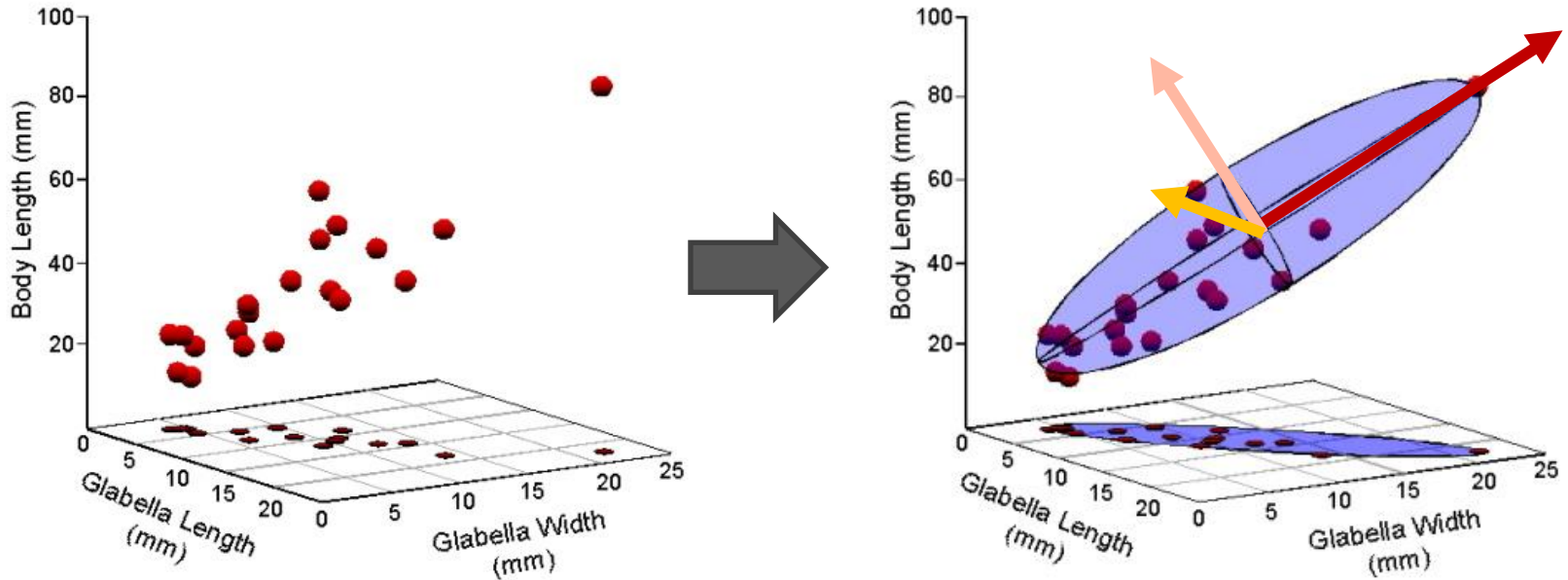
Projection



<https://shapeofdata.wordpress.com/2013/04/16/visualization-and-projection/>

- High variances = more power to distinguish groups of data points

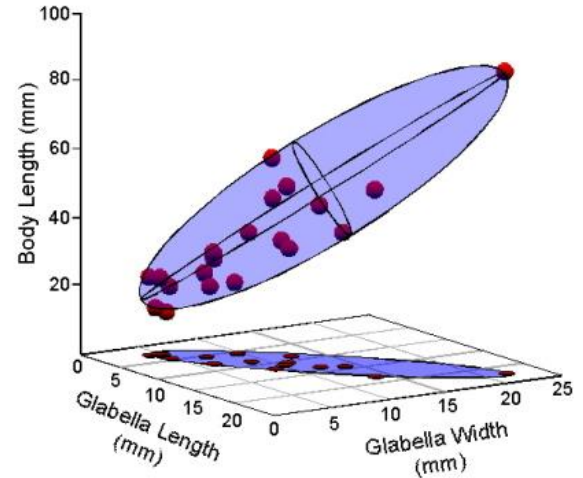
PCA identifies directions with high variances



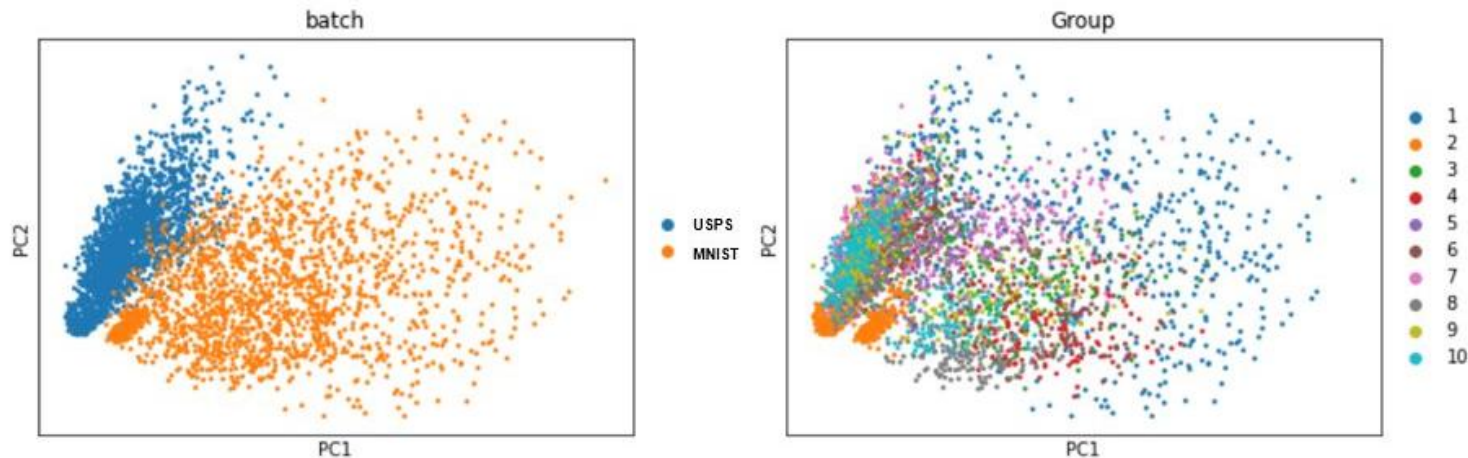
Source: the paleontological association

Interpretations of PCA algorithm

- Fitting the data with an n -dimensional ellipsoid
 - Axes = **principal component (PC)** direction
 - Axis lengths = magnitude of variances captured
- Orthogonal linear transformation
 - New axes are **rotations** of the original axes
 - $PC1 = w_1x_1 + w_2x_2 + \dots + w_nx_n$
 - $PC2 = v_1x_1 + v_2x_2 + \dots + v_nx_n$
 - w_i 's and v_i 's are called **loadings**



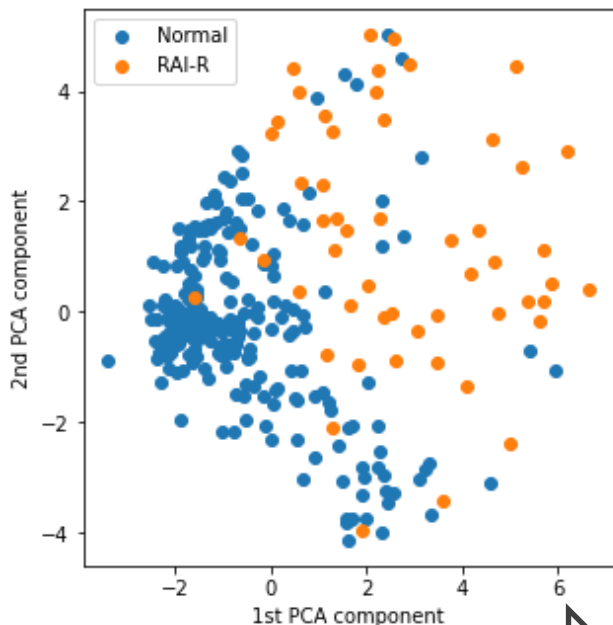
Interpretation of PCA result



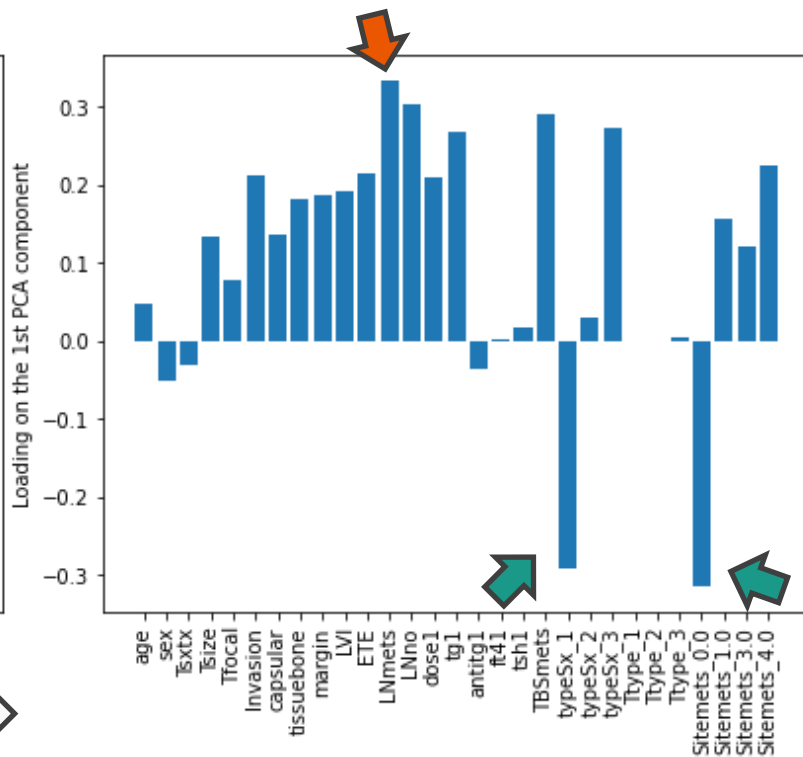
<https://twitter.com/lkmklsmn/status/1436357177887895555>

- PC1 captures the variance between data sources
- PC2 somewhat captures the variance between digit identity

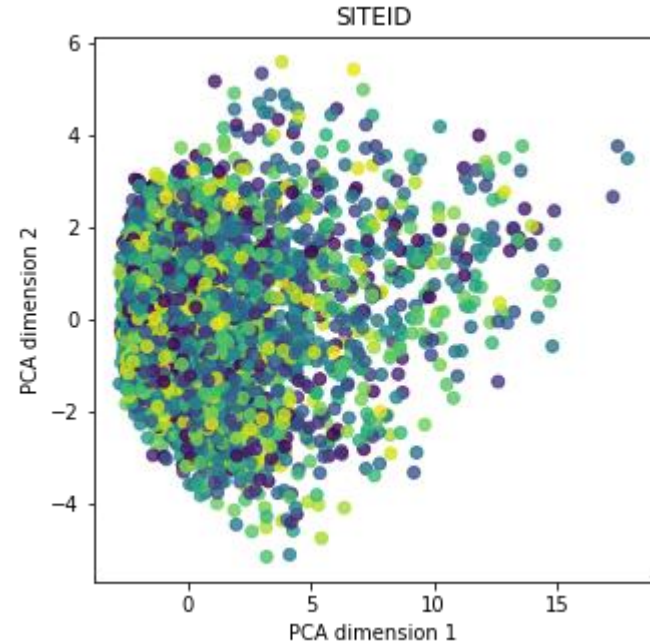
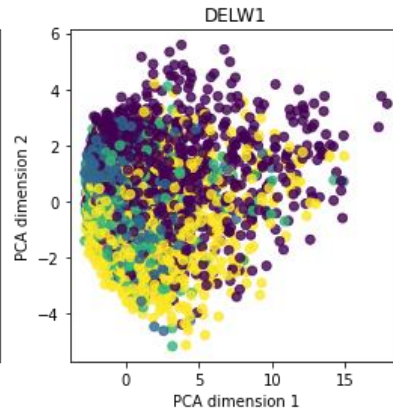
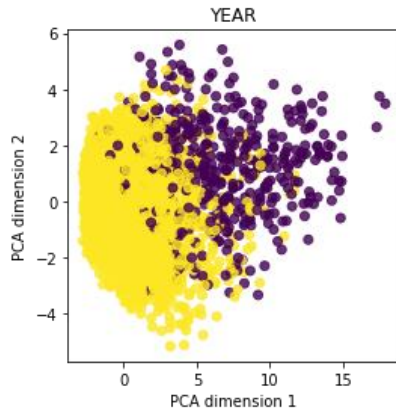
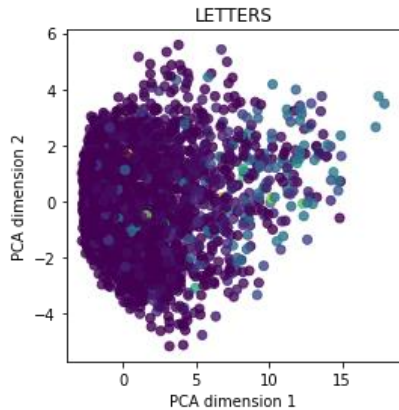
Interpreting loadings on individual PC



Resistance to treatment

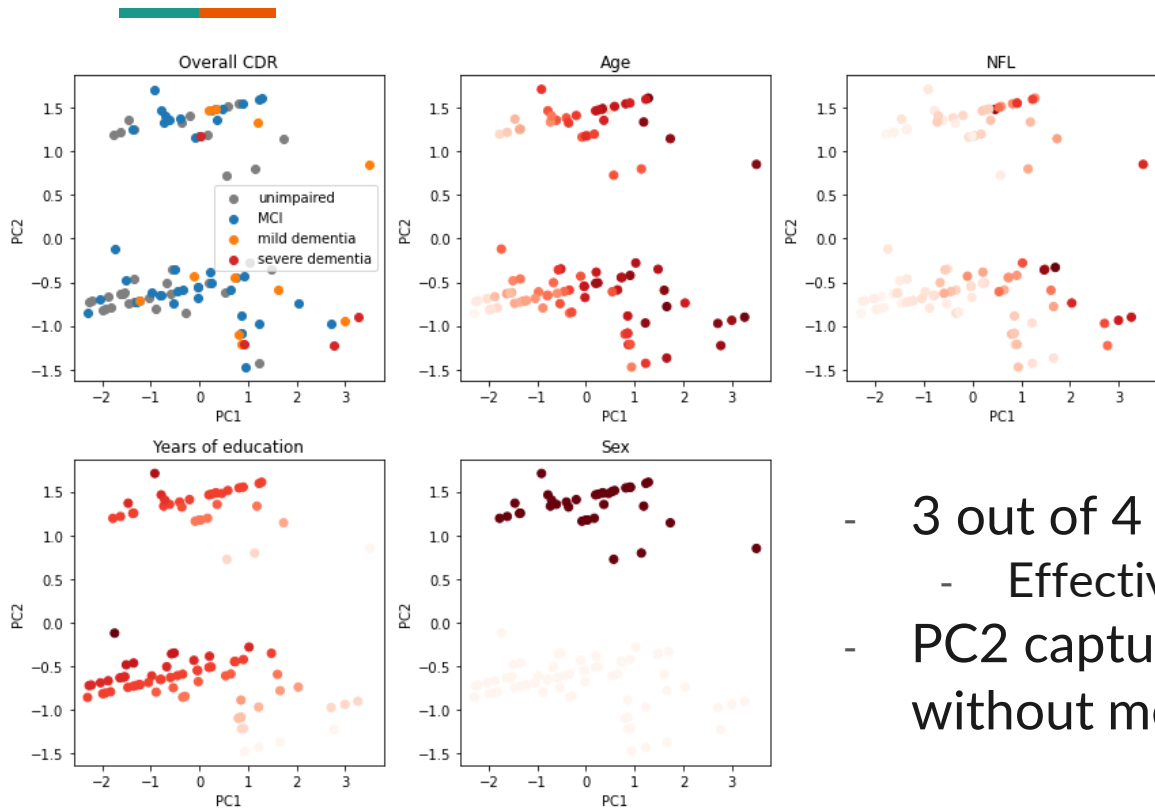


Exploring PCA results



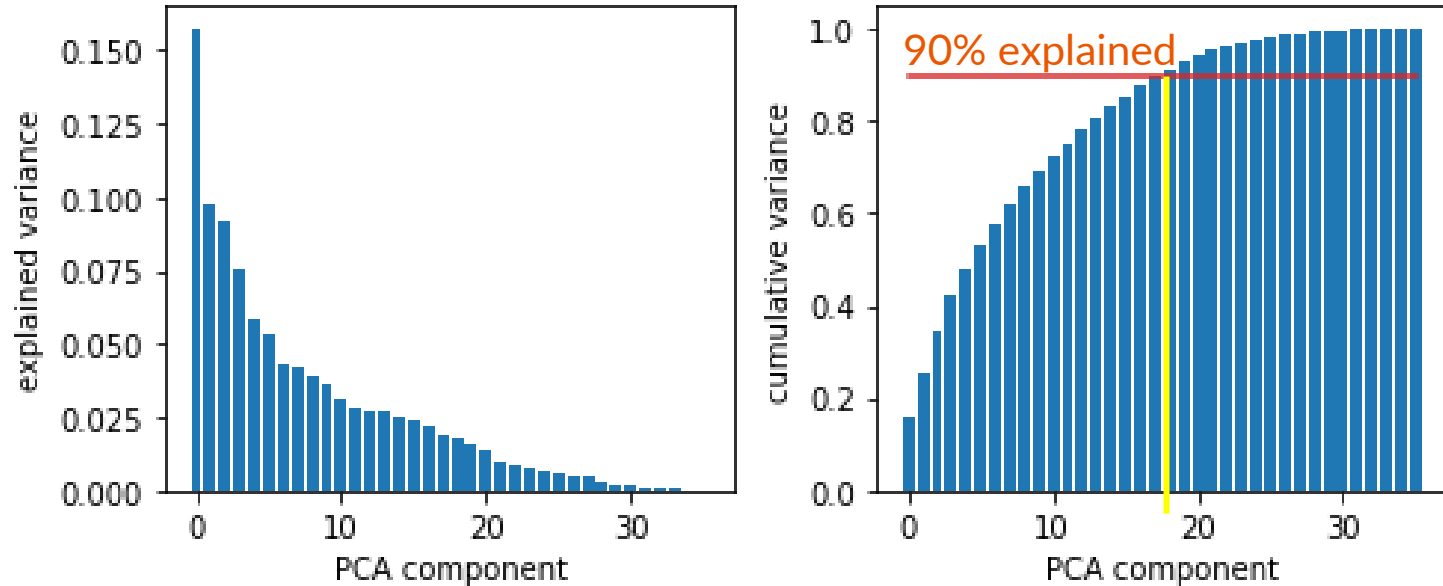
- Color by feature values to understand how PCA group data points
- Color by potential confounding factors

Be careful of correlated features



- 3 out of 4 features are correlated
 - Effective dimension is 2
- PC2 captures the gender difference without meaningful interpretation

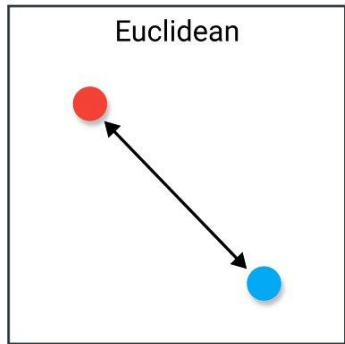
PCA for dimensionality reduction



- By default, PCA retains the number of dimensions
- We can **select only the first k PC** for downstream analyses

Pros and cons of PCA

- Each PC can be interpreted from the **loadings**
- **Highly correlated features tend** to be grouped into the same PC
- PCA is a good initial dimensionality reduction step
- PCA strictly **preserves Euclidean distance**
 - But some datasets require different distance metric!

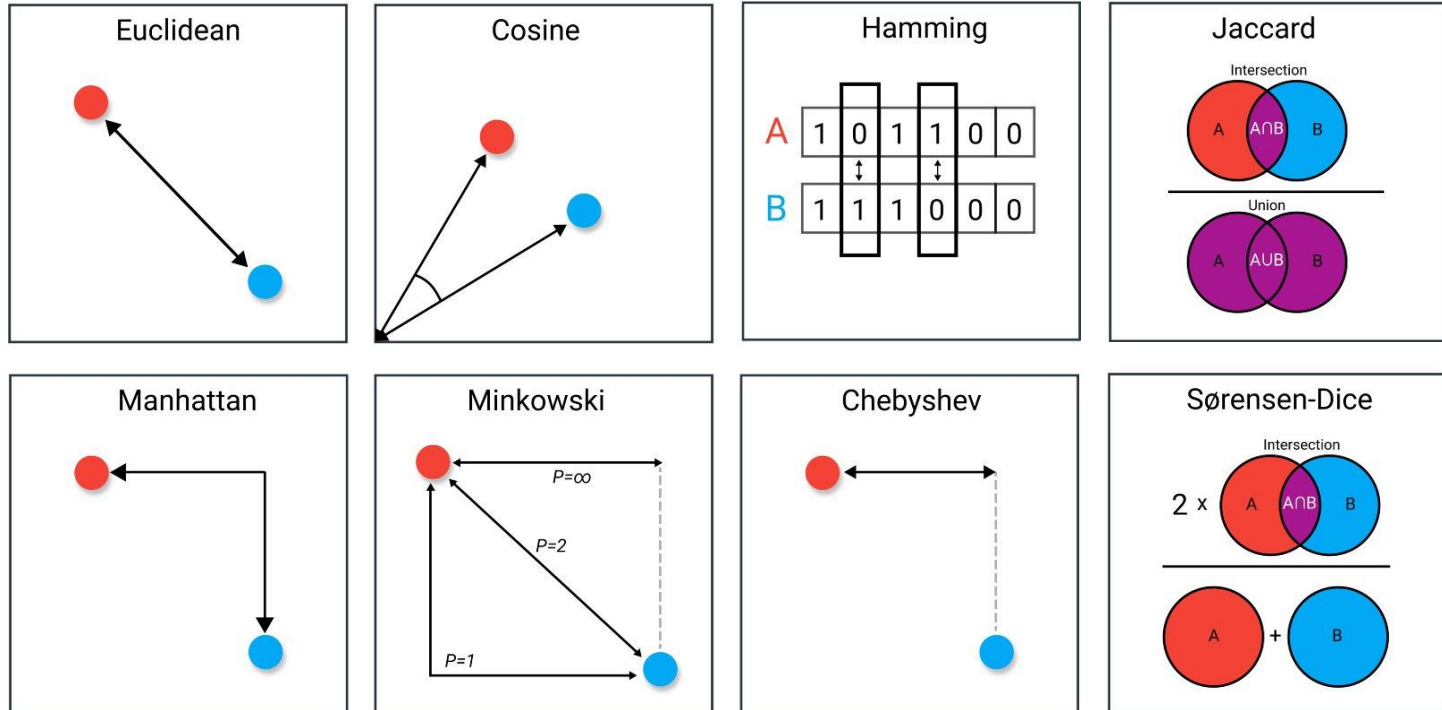


<https://towardsdatascience.com/9-distance-measures-in-data-science-918109d069fa>

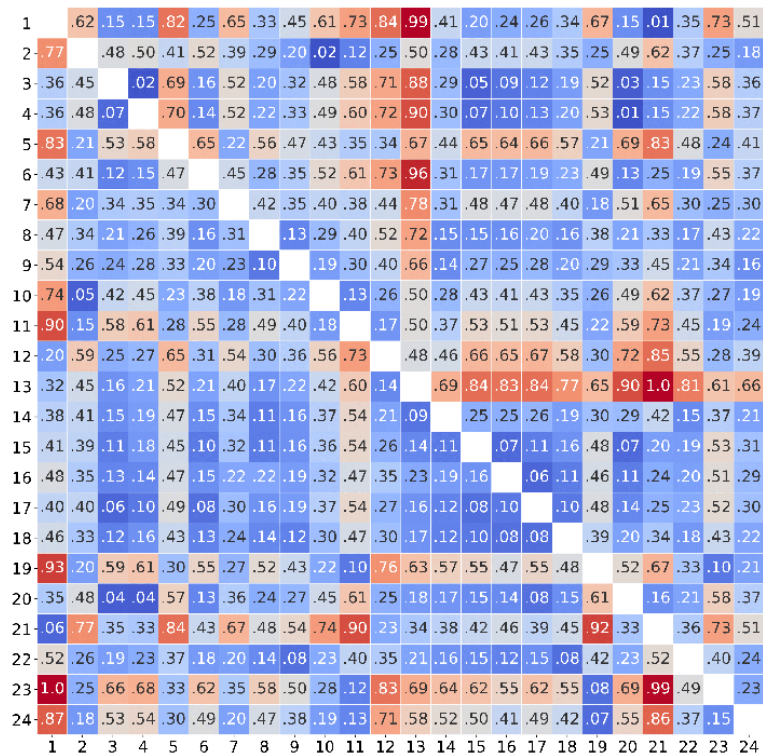


Multidimensional Scaling (MDS)

Distances

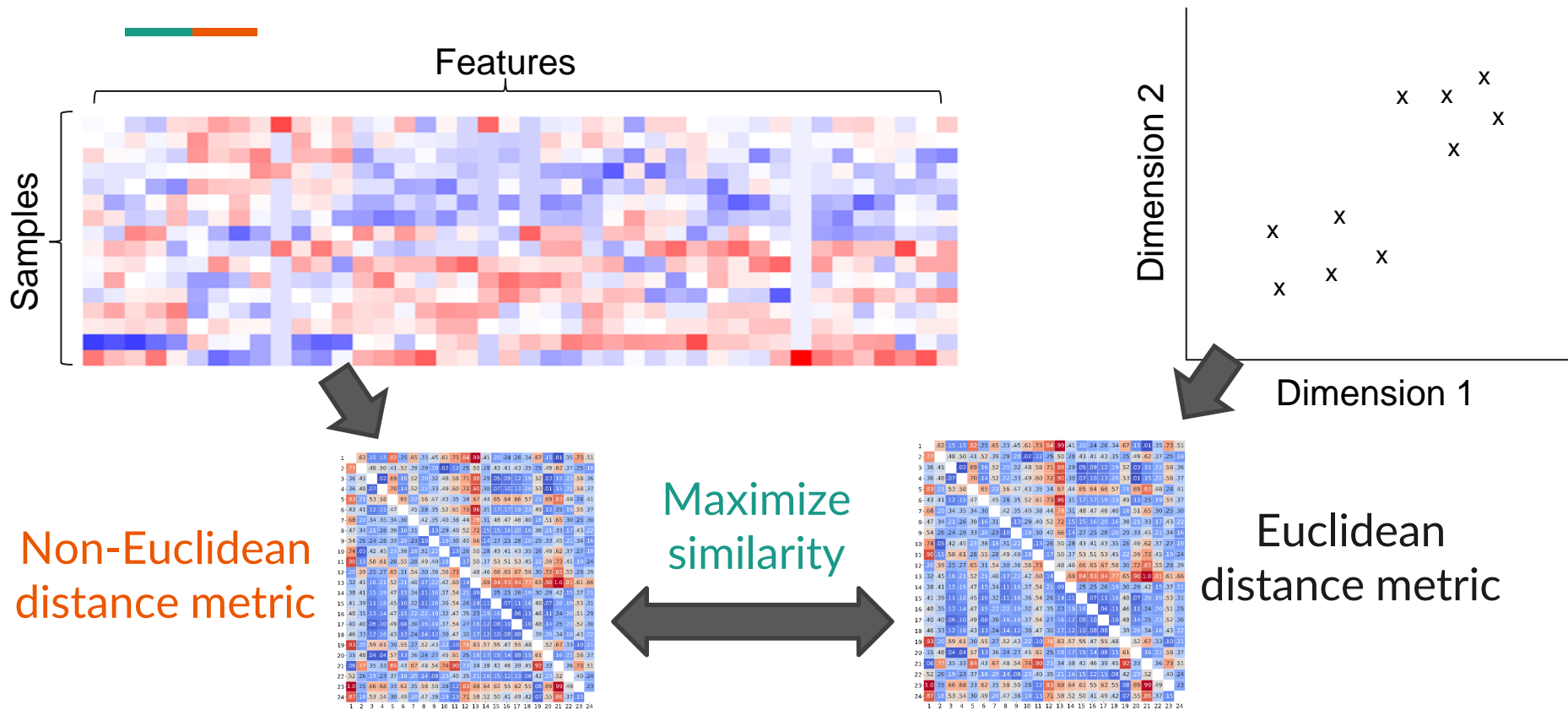


Pairwise distance matrix



- $D(i, j)$ = distance between sample i and sample j
- $D(i, i) = 0$
- $D(i, j) = D(j, i)$
- User-defined distance metric

Principal Coordinate Analysis (PCoA)

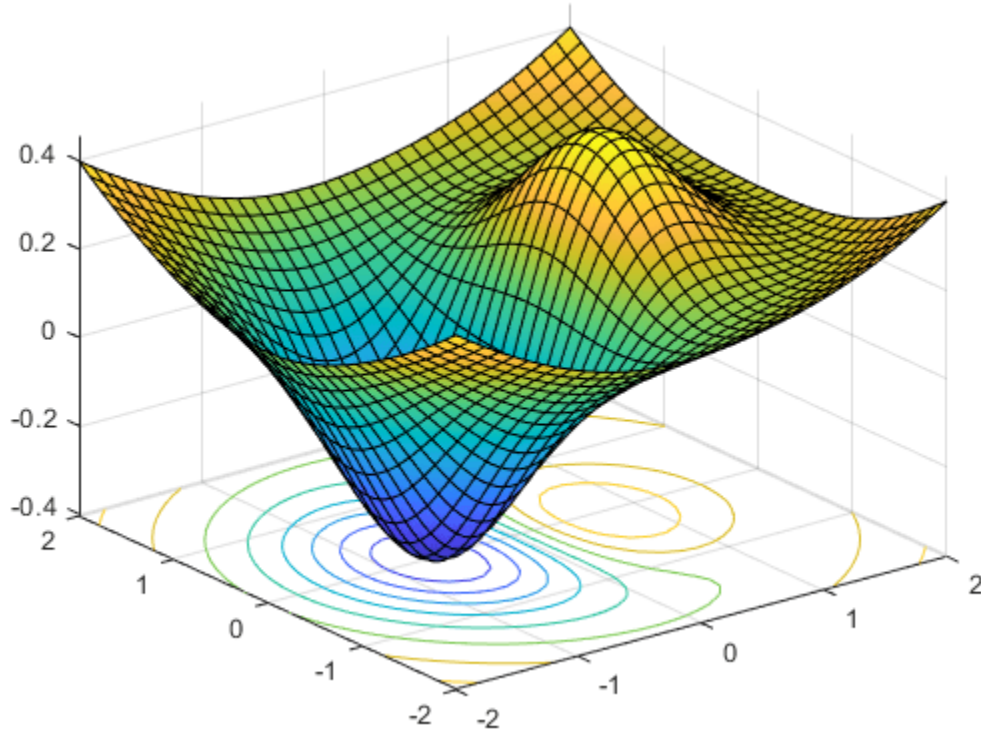


PCoA algorithm sketch



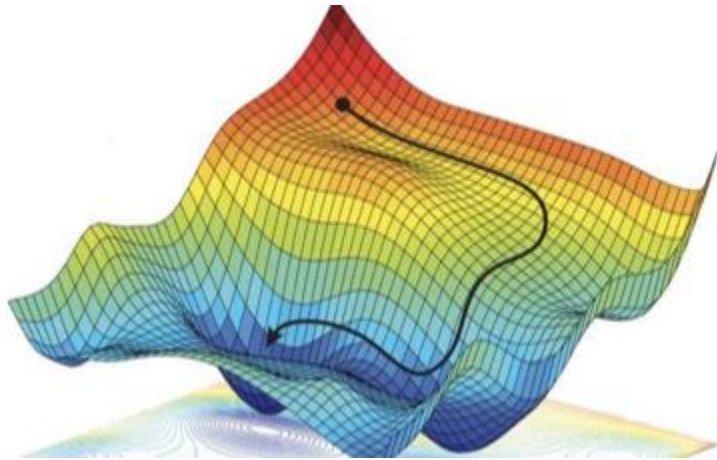
- Calculate pairwise distance matrix D for the original data
- For a 2D projection of the data: sample i projected onto (x_i, y_i)
 - Calculate pairwise distance matrix D' for the projection
 - Each element of D' is a function of (x_i, y_i) 's
- Calculate a similarity score between D and D'
 - Such as Person's correlation
 - This similarity score is a function of (x_i, y_i) 's
- Find (x_i, y_i) 's that maximize this score!

How to optimize a function?

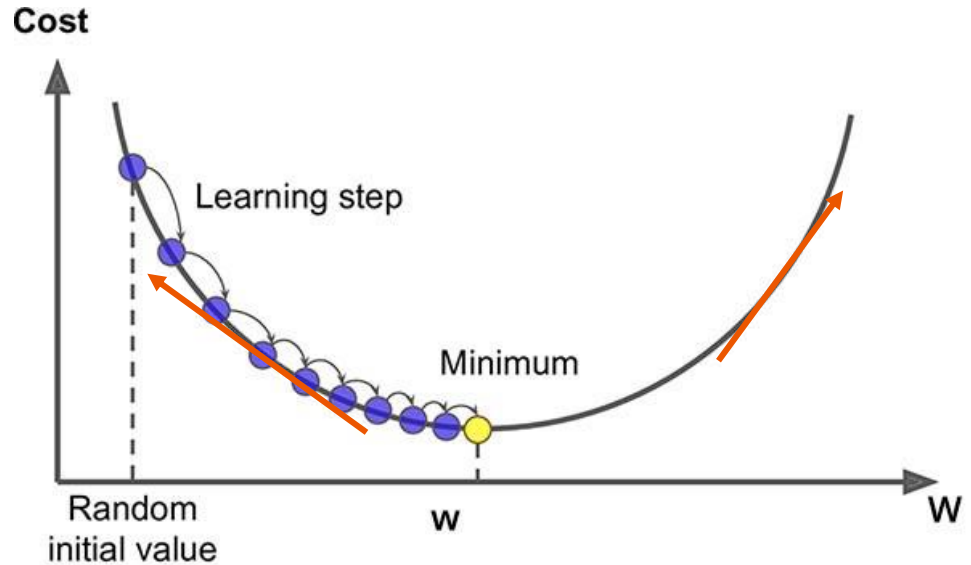


- Find (x_1, x_2, \dots, x_n) that minimize $f(x_1, x_2, \dots, x_n)$
- At minimum, the slope is zero in all directions
- Take derivative of each variable and set to zero
 - $\frac{\delta f}{\delta x_1} = 0$
 - $\frac{\delta f}{\delta x_2} = 0$
 - n equations with n variables

Gradient descent

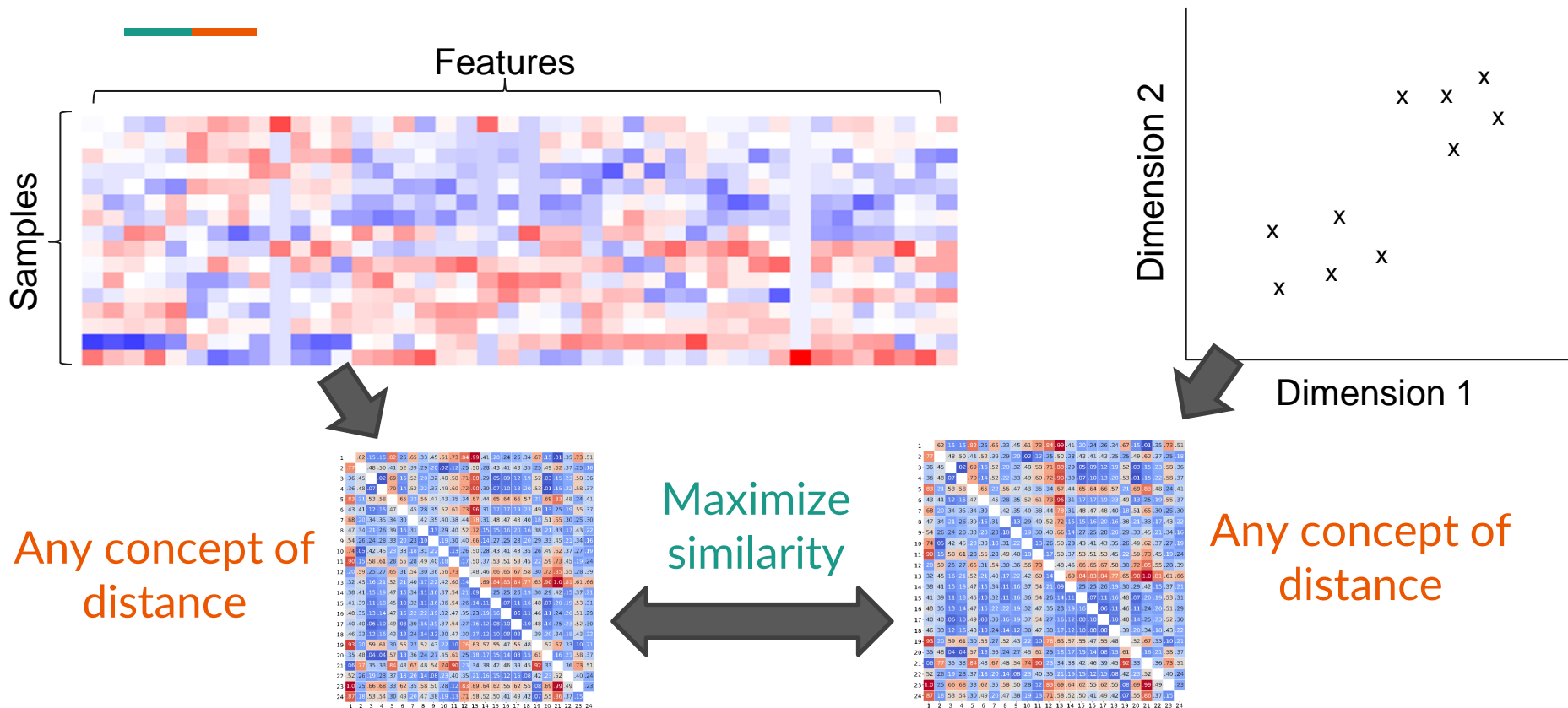


<https://medium.com/analytics-vidhya/gradient-descent-b0dc1af33517>



- Slope tells us if the function is increasing or decreasing if we increase x_i
 - So, we can update x_i accordingly

Generalized MDS



Limitation of PCA and MDS

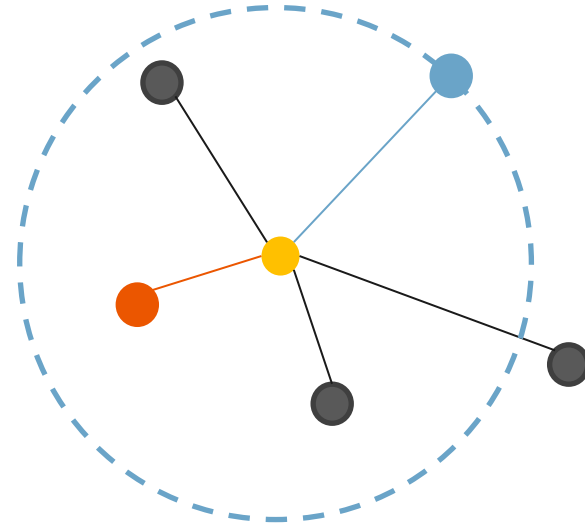
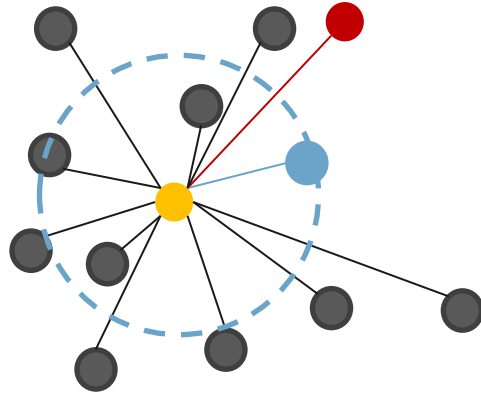


- A single definition of distance metric is used throughout the data space
- What if some data groups are noisier than the others?
 - Difference in data density



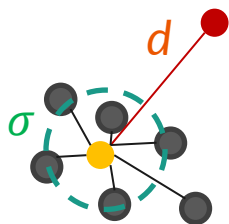
***t*-distributed stochastic neighbor embedding (*t*-SNE)**

Measuring data density



- Distance to the k -th nearest neighbor reflects data density
 - Small distance in dense area
 - Large distance in sparse area

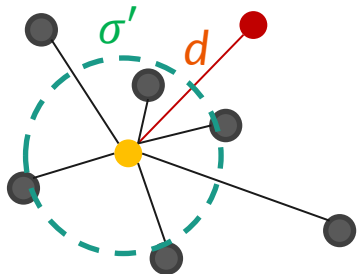
Probability of being a neighbor



$\text{score}(\text{red} \mid \text{yellow})$ = probability that yellow would pick red as neighbor under a **normal distribution** center at yellow

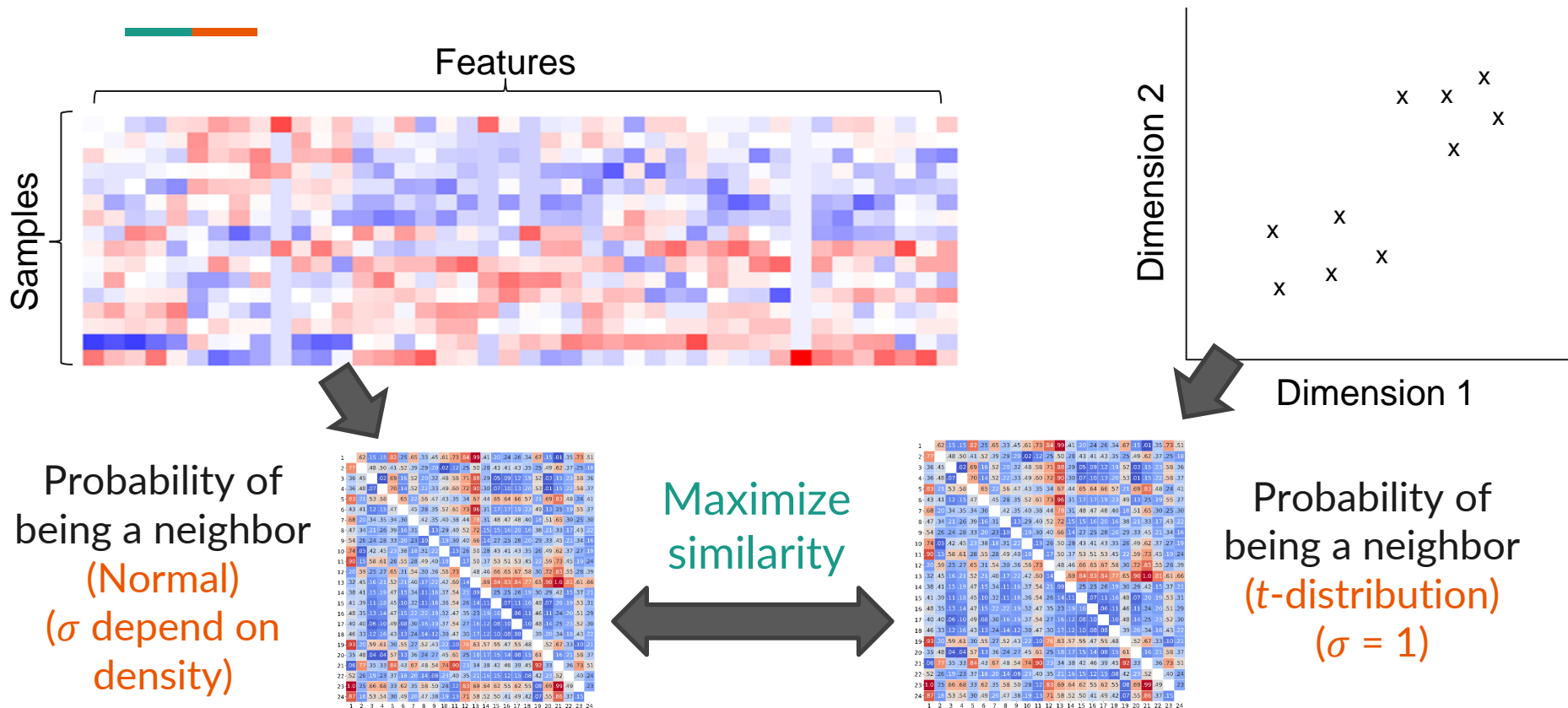
$$= \frac{e^{-\frac{d^2}{2\sigma^2}}/\sigma}{\sum e^{-\frac{(\text{dist}(\text{blue}, \text{yellow}))^2}{2\sigma^2}}/\sigma}$$

blue = other data points



- Same distance d normalized against density σ and distances to other nearby data points blue

Finding the optimal projection for t -SNE



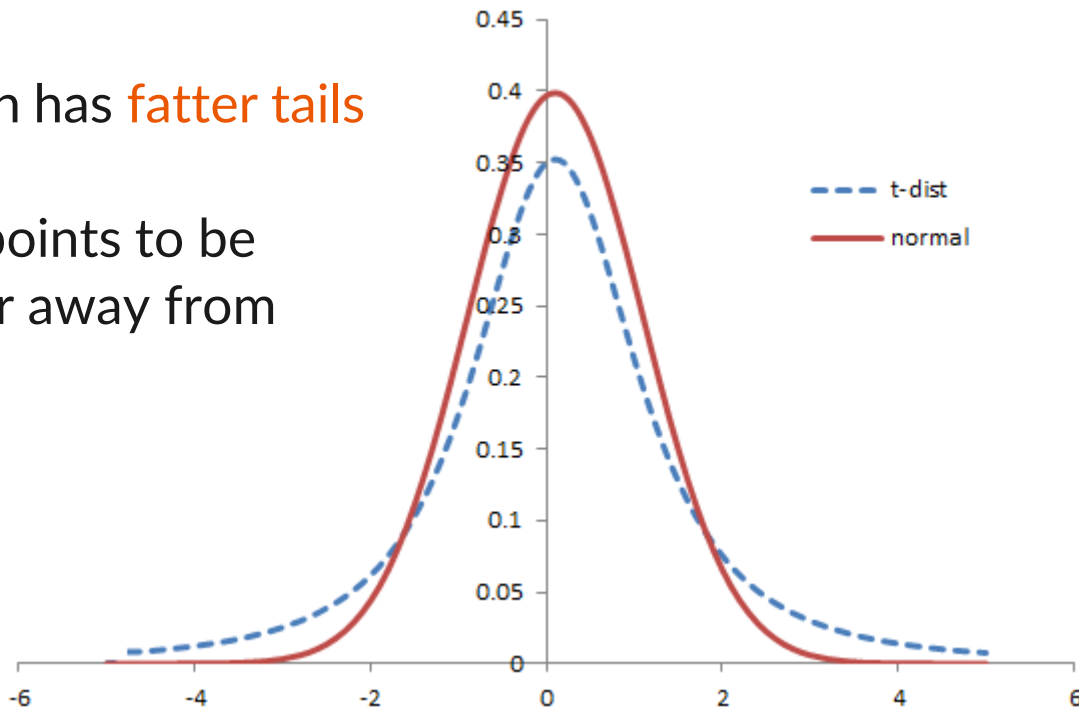
A similarity score for probability distribution



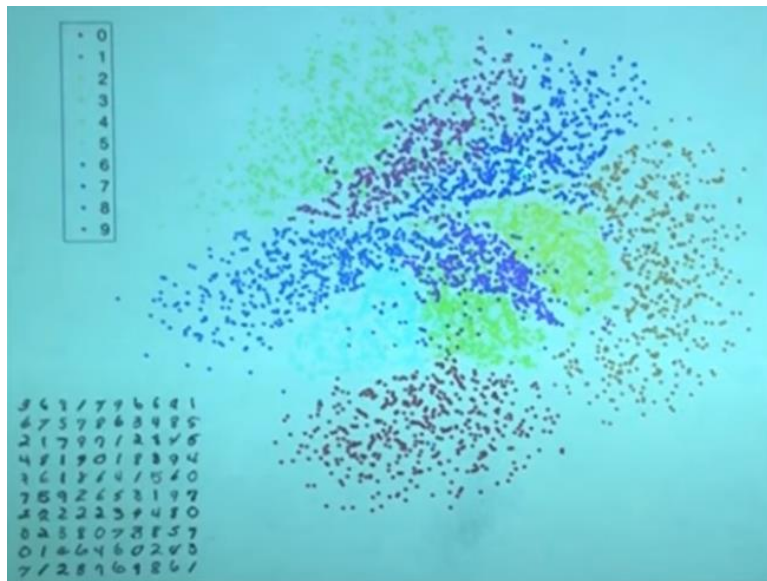
- Kullback-Leibler (KL) divergence
- Measure how distribution P differs from another distribution Q
 - $D_{KL}(P \parallel Q) = -\sum_i \sum_j p_{j|i} \log_2 \frac{p_{j|i}}{q_{j|i}}$
- P = Probability of neighbor from the original data
- Q = Probability of neighbor from the projection
- Solve for the best $q_{j|i}$'s using gradient descent

Why t -distribution for the projection?

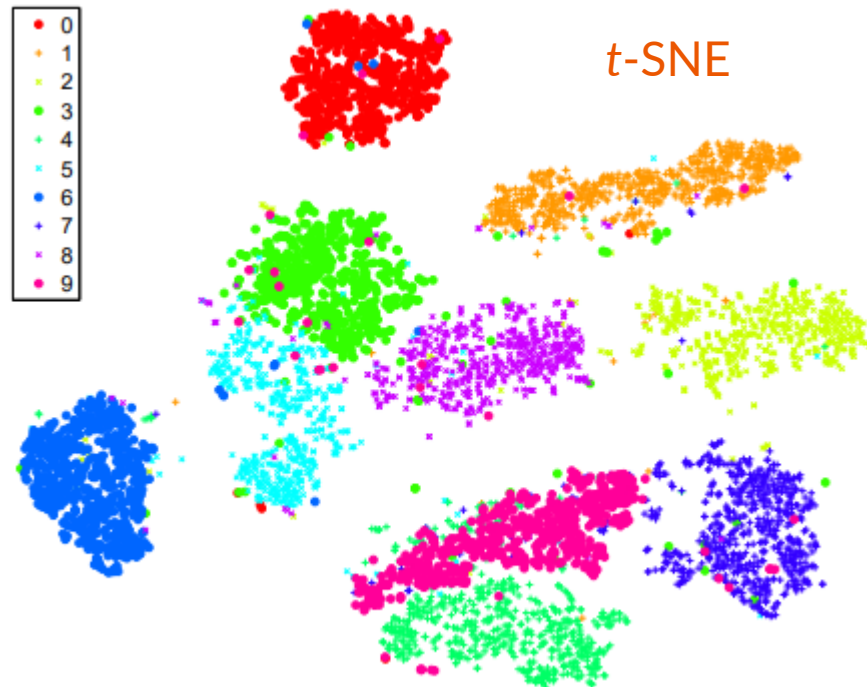
- t -distribution has **fatter tails**
- Allow data points to be projected far away from each other



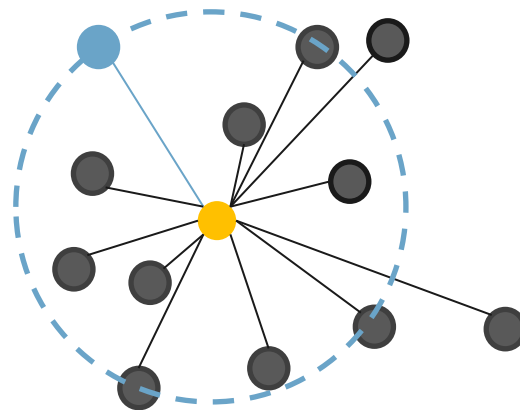
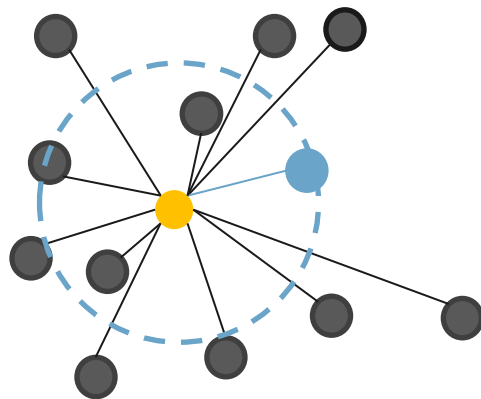
Impact of t -distribution



SNE (Normal \rightarrow Normal)

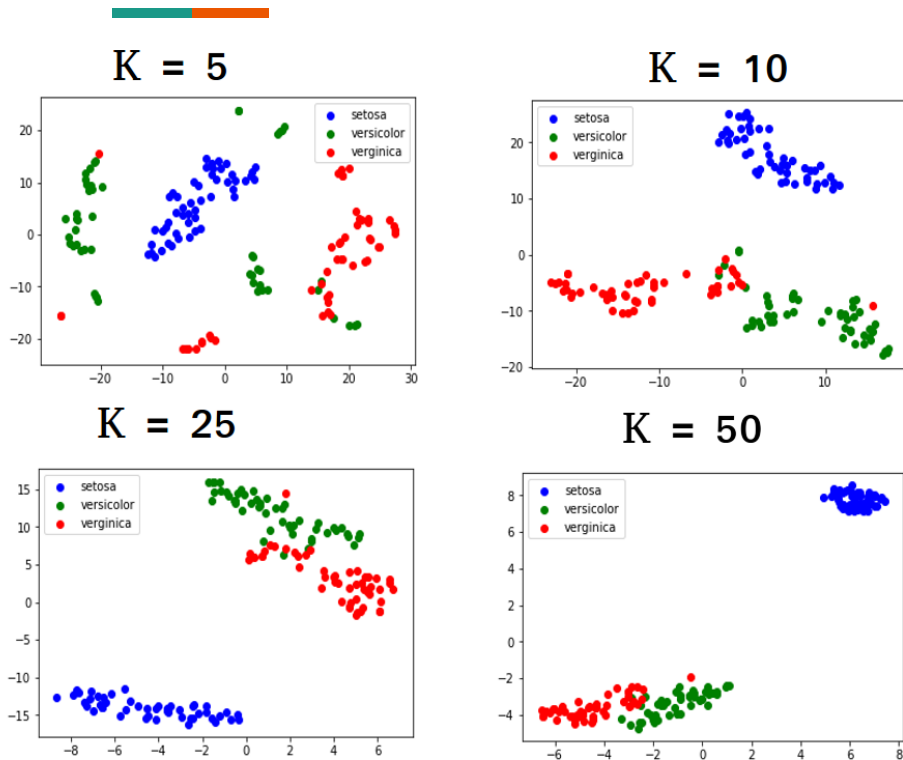


Perplexity



- How many nearest neighbors to consider to normalize data density?
 - Perplexity parameters

Impact of perplexity

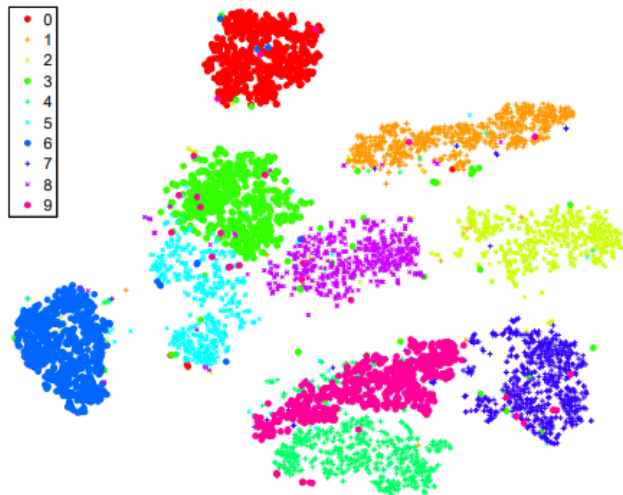


- Too small perplexity = a lot of scatted data groups
- Try varying the perplexity and identify patterns that **consistently** appear

Pros and cons of t -SNE



- Capture qualitative neighbor relationship
- Normalize data density
- Recompute every time new data is added
- Lose **long-range** relationship
- Axes of the resulting projection have no meaning
 - **Don't use t -SNE coordinates for clustering or interpretation**

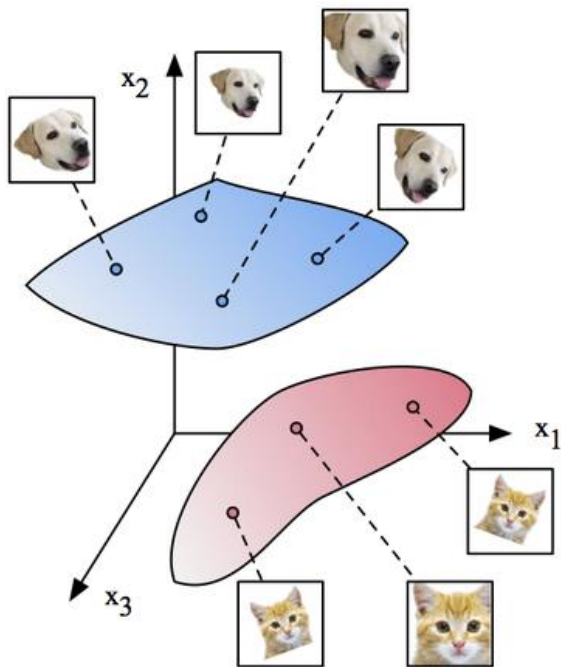


Maaten, L. and Hinton, G. J of Machine Learning Research 9:2579-2605 (2008)

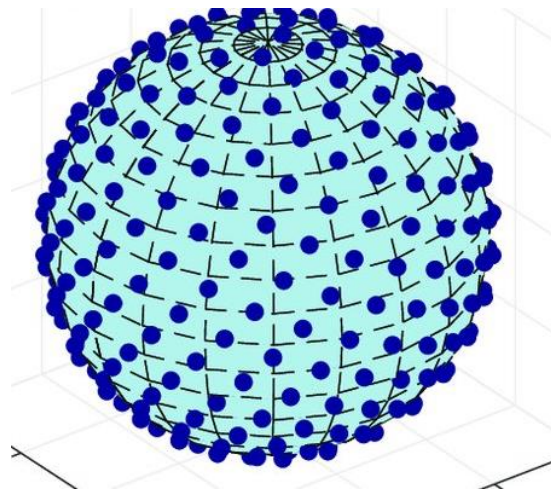


Uniform manifold approximation and projection (UMAP)

Two key assumptions



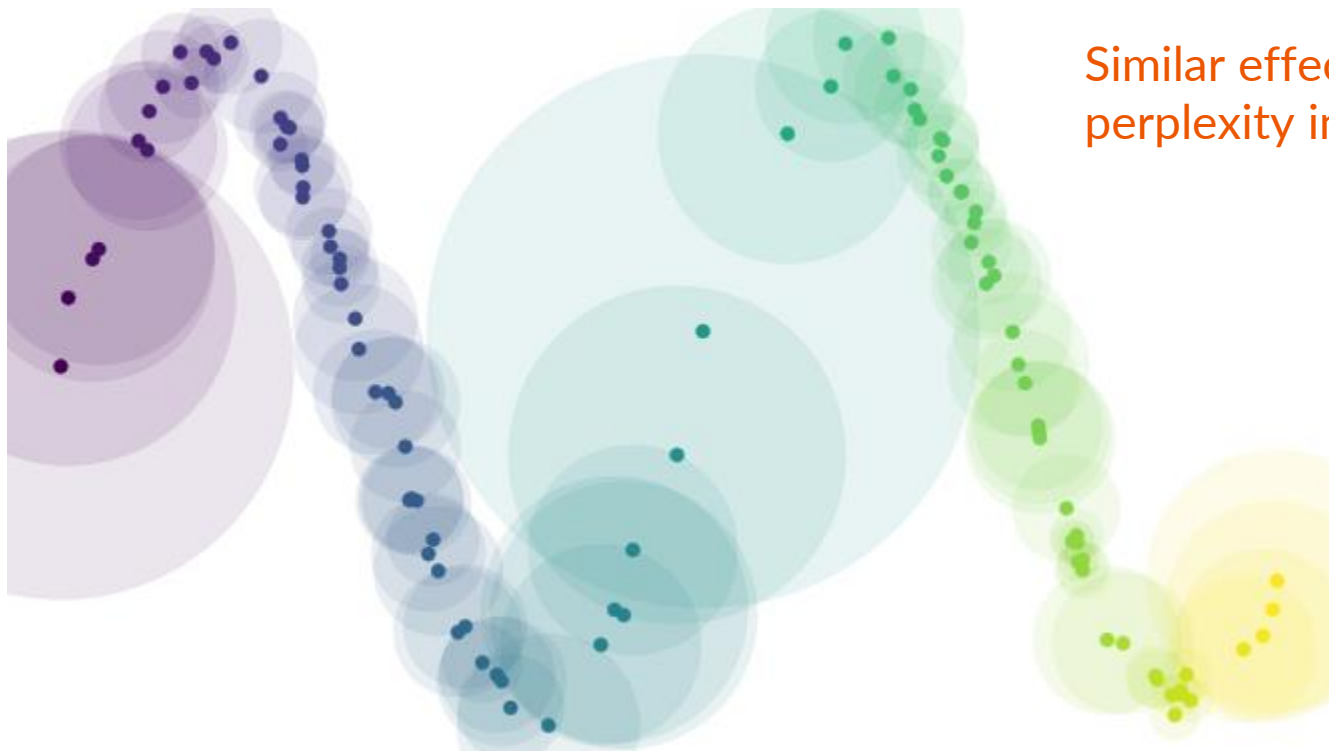
Chung, S. et al. "Classification and Geometry of General Perceptual Manifolds"



Ali, A. et al. IEEE Access PP(99):1 (2021)

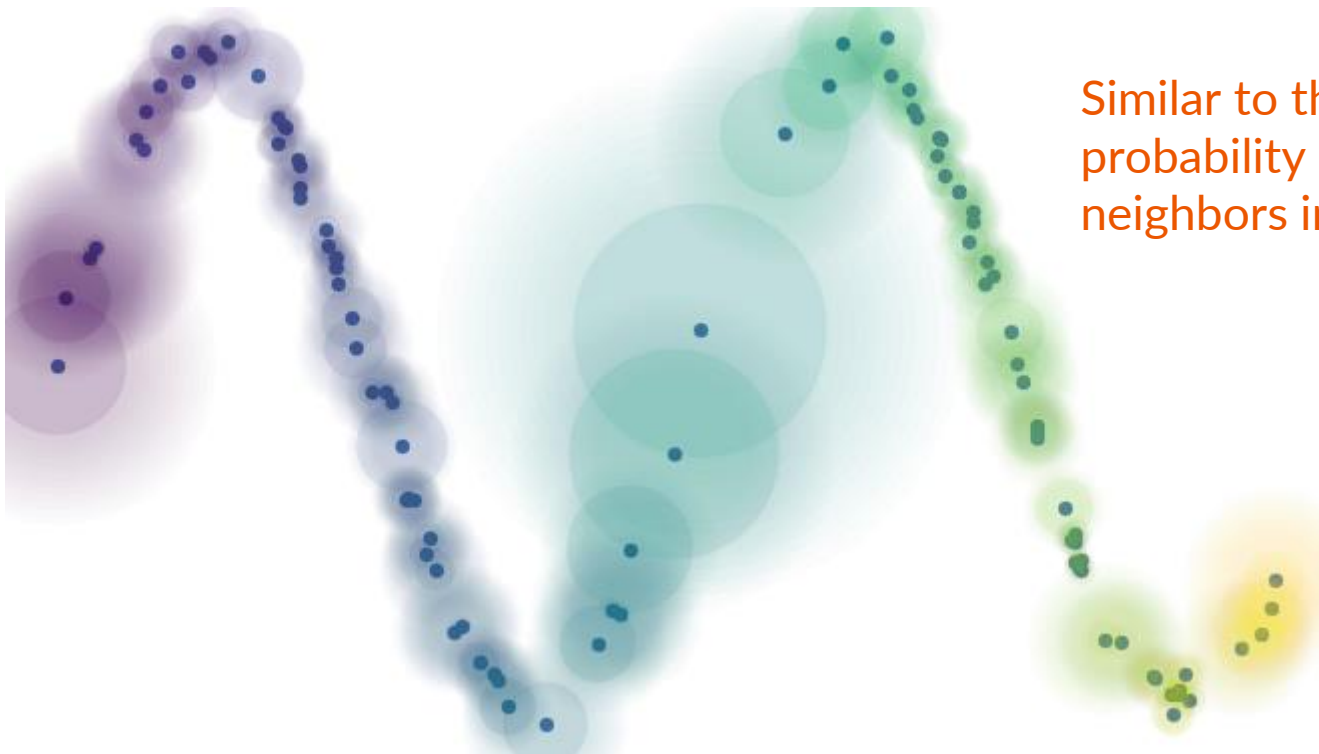
- Data came from **multiple manifolds**
- Data points were **sampled uniformly**

Uniform sampling = similar distance to k -th neighbor

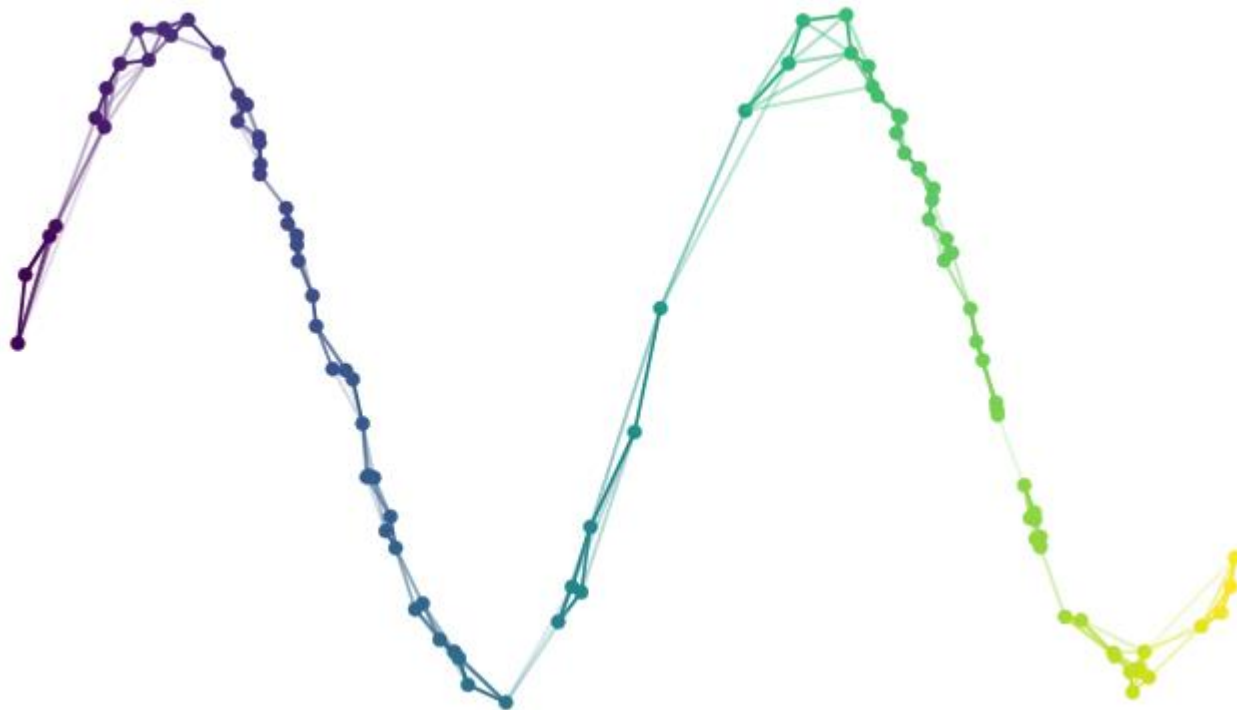


Similar effect as
perplexity in t-SNE

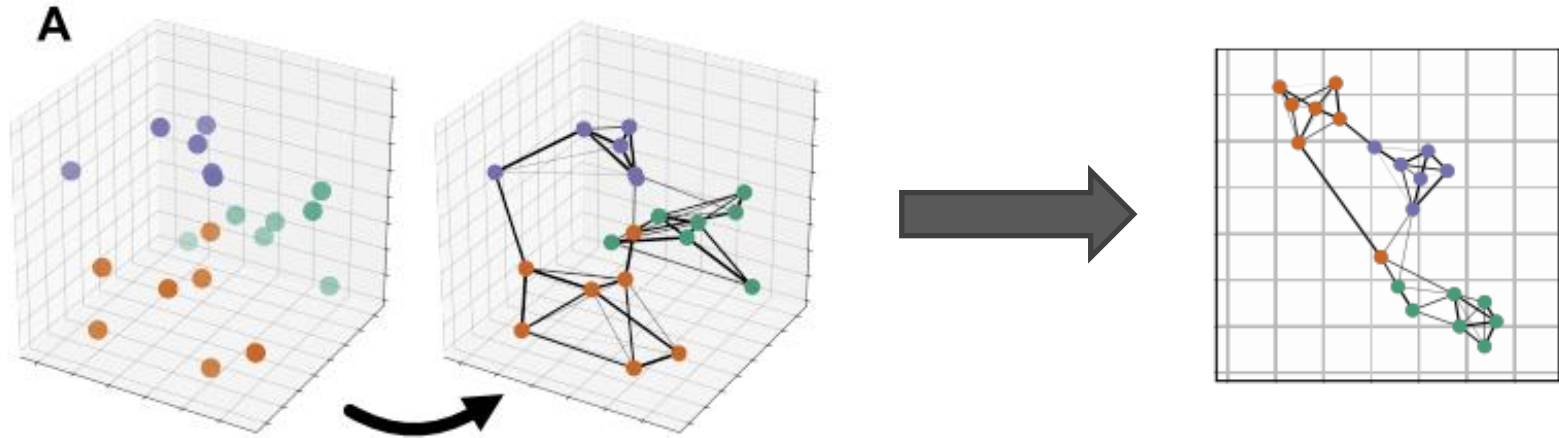
Adding uncertainty between distant data points



Network representation of neighbor relationship



Projecting network representation



Sainburg, T. *et al.*, Neural Comput 33(11):2881-2907 (2021)

- Preserve scores on edges: probability of being neighbors

Another similarity score for probability distribution

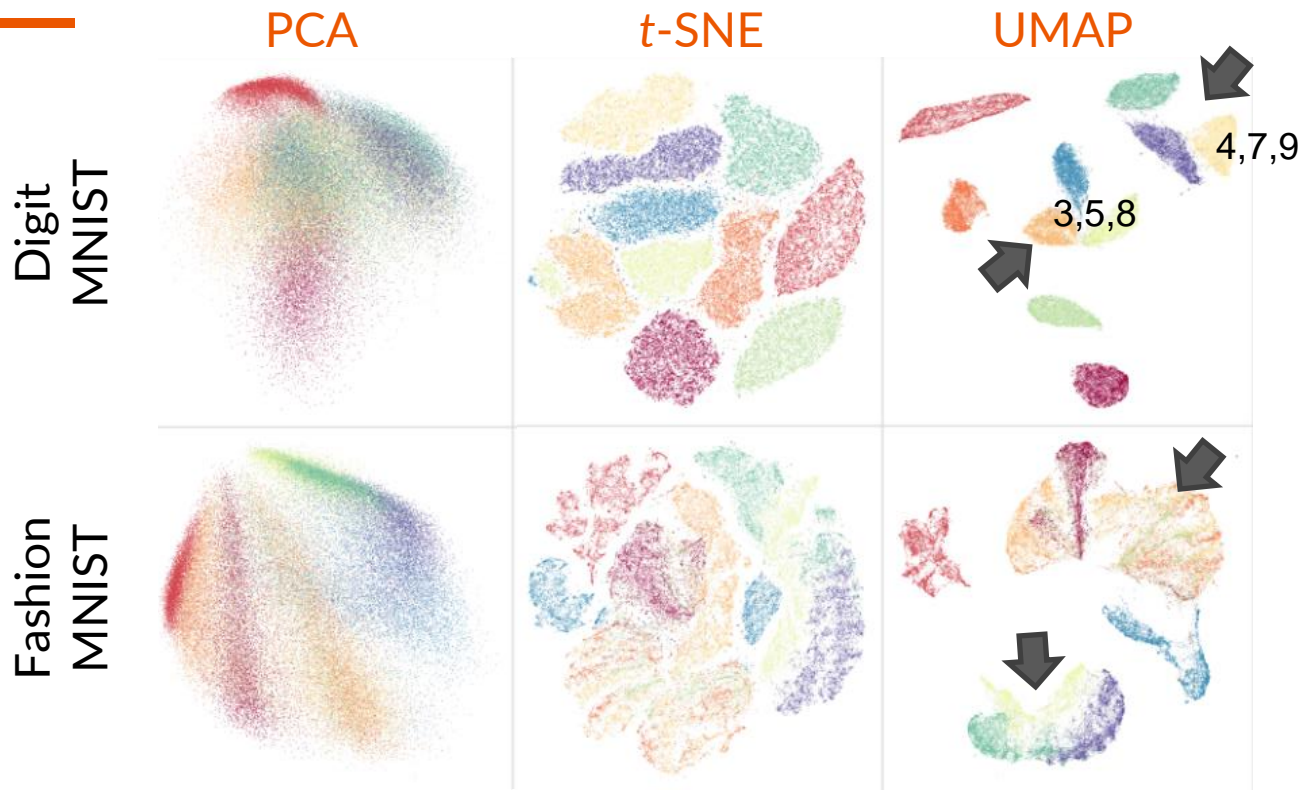


- Cross-entropy

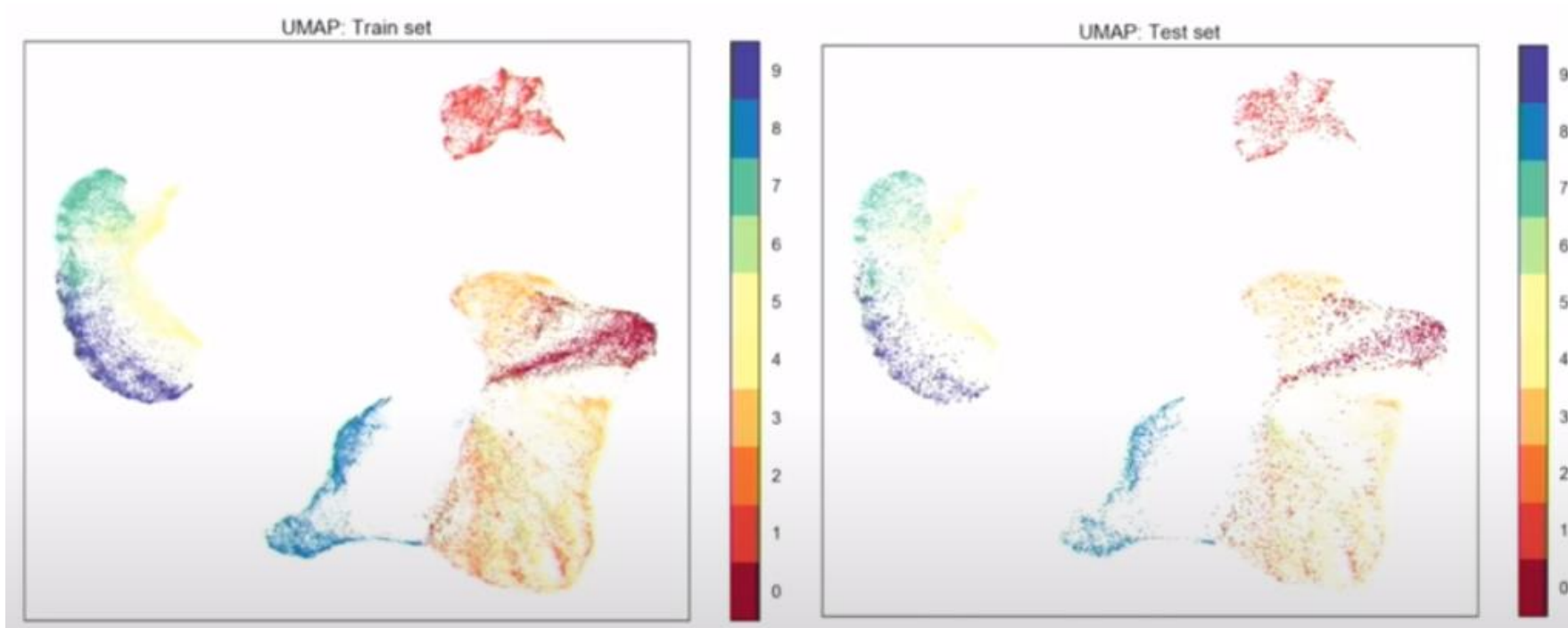
- $$\text{CE}(P \parallel Q) = \sum_i \sum_j p_{i,j} \log_2 \frac{p_{i,j}}{q_{i,j}} + \sum_i \sum_j (1 - p_{i,j}) \log_2 \frac{(1-p_{i,j})}{(1-q_{i,j})}$$

- KL divergence only has the first term
- Cross-entropy considers both when $p_{i,j}$ is high (similar data points) and when $p_{i,j}$ is low (distant data points)

Power of UMAP



UMAP can transform new data points

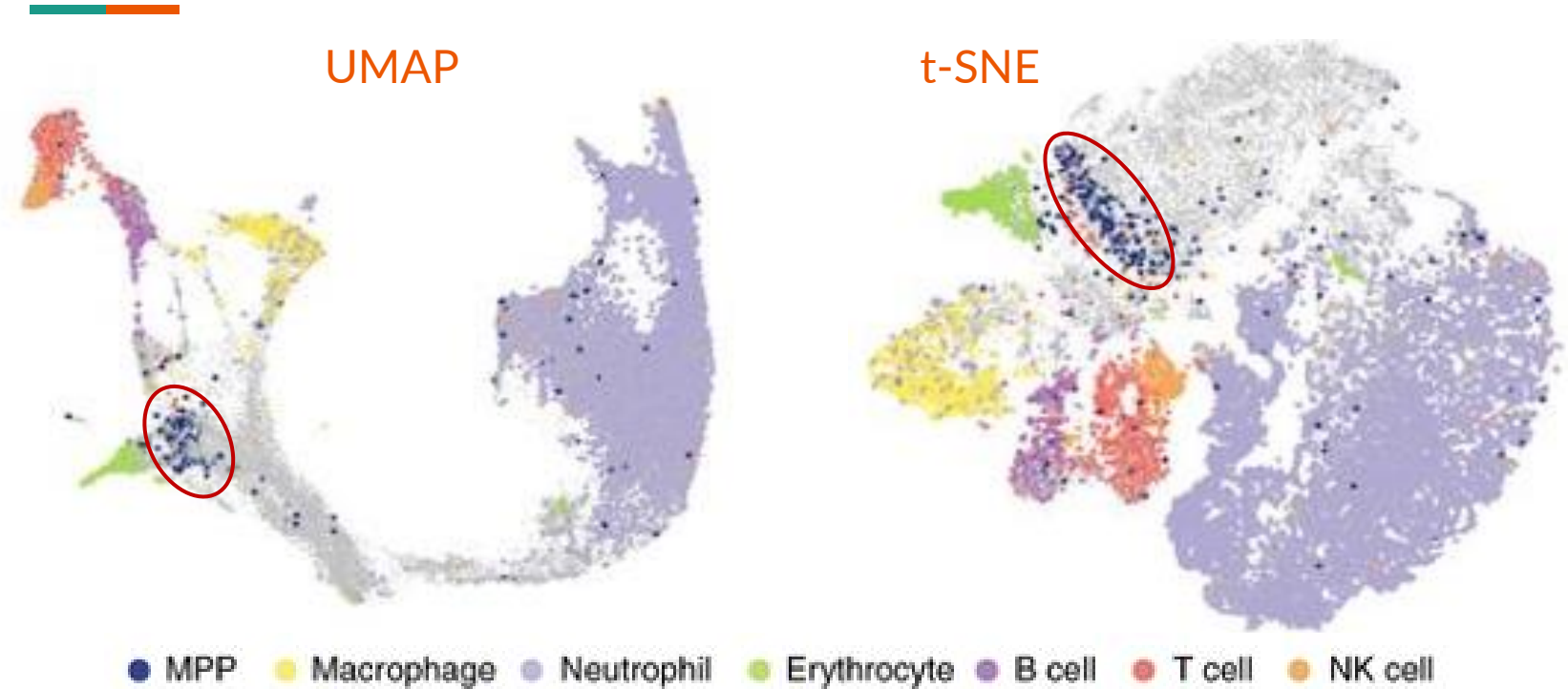


Pros and cons of UMAP

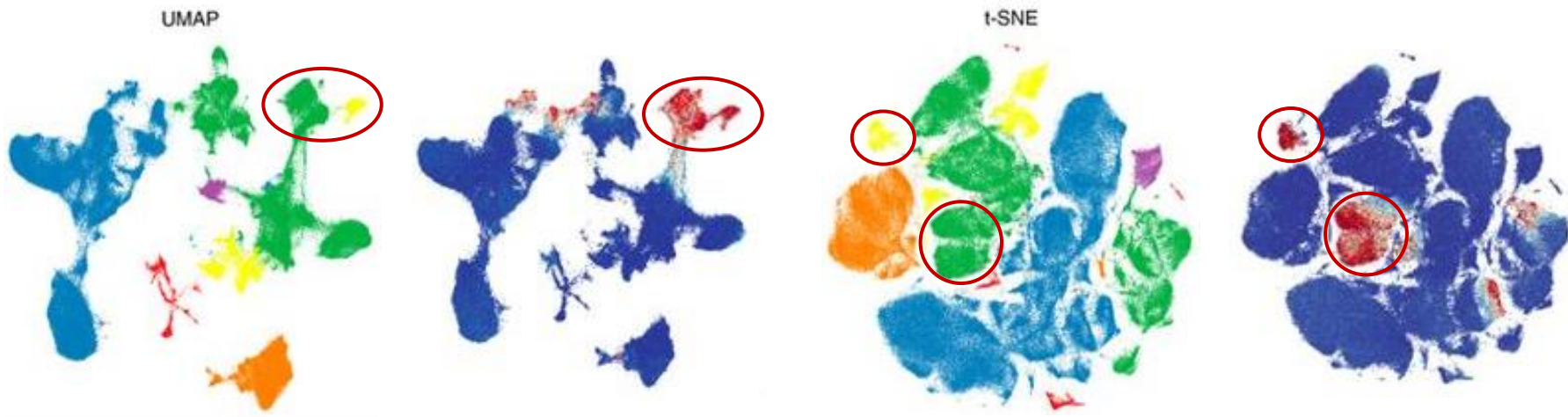


- Can capture long-range relationship
- Can be applied to new data points without recomputing
- Require a **strong assumption** of uniform sampling

t-SNE vs UMAP on biological data



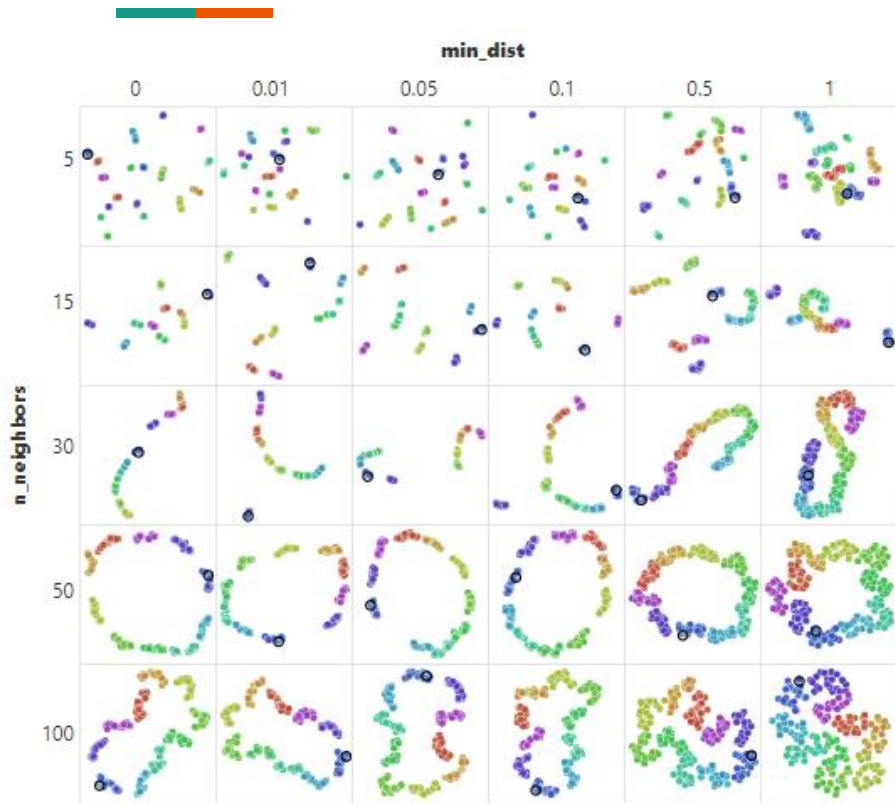
t -SNE vs UMAP on biological data



Becht, E. et al. Nature Biotechnology 37:38-44 (2019)

- Both are equally good at detecting individual data groups
- But UMAP is better at capturing transitions across data groups

Customizing UMAP outputs



- Number of neighbors (**n_neighbors**) is perplexity
- Minimum distant for placing similar data point (**min_dist**) is for adjusting the scale of visualization



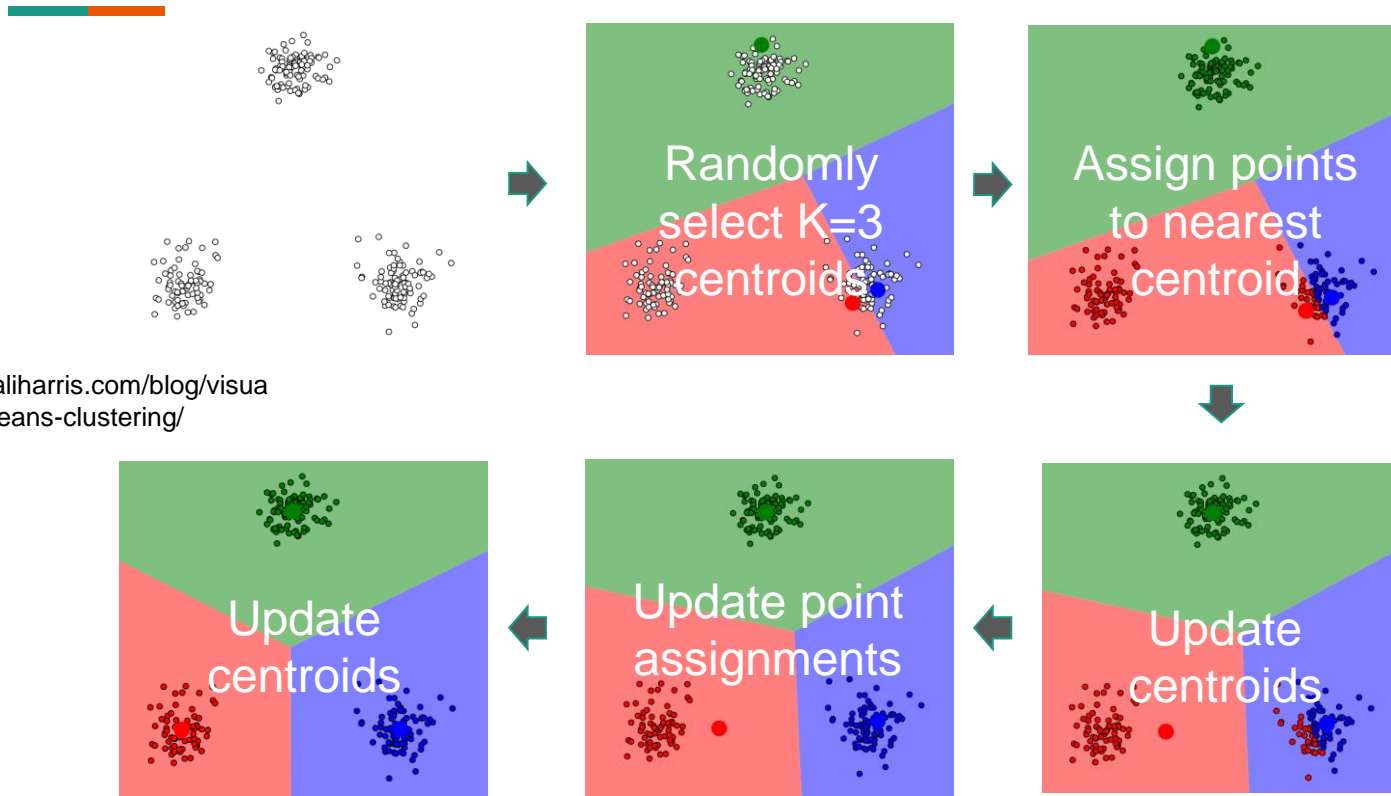
Clustering

The heart of clustering



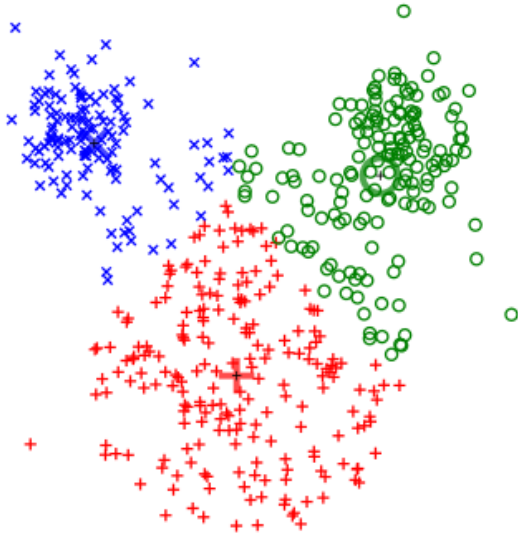
- **Goal:** Group similar data point together
- How to define similarity?
 - **Distance:** Between two data points
 - **Linkage:** Between groups of data points
- How many clusters is appropriate?
 - Within-cluster (small) versus between-cluster (large) distance

An example: k -mean clustering



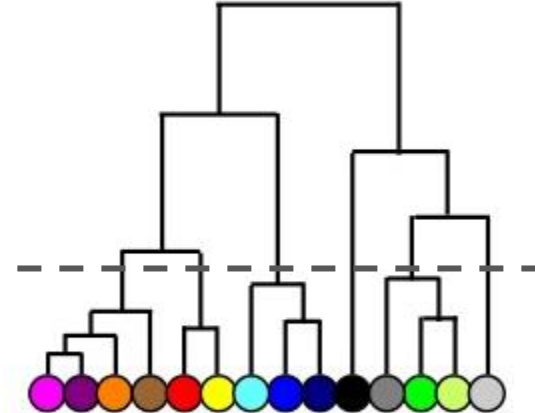
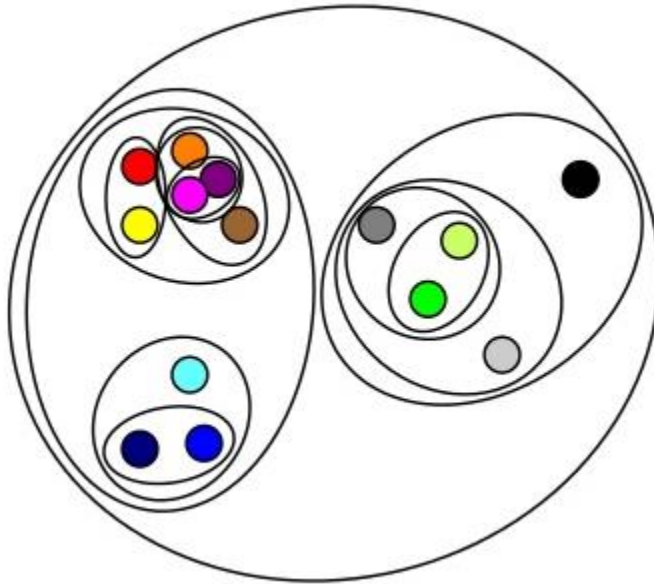
www.naftaliharris.com/blog/visualizing-k-means-clustering/

Limitation of k -mean



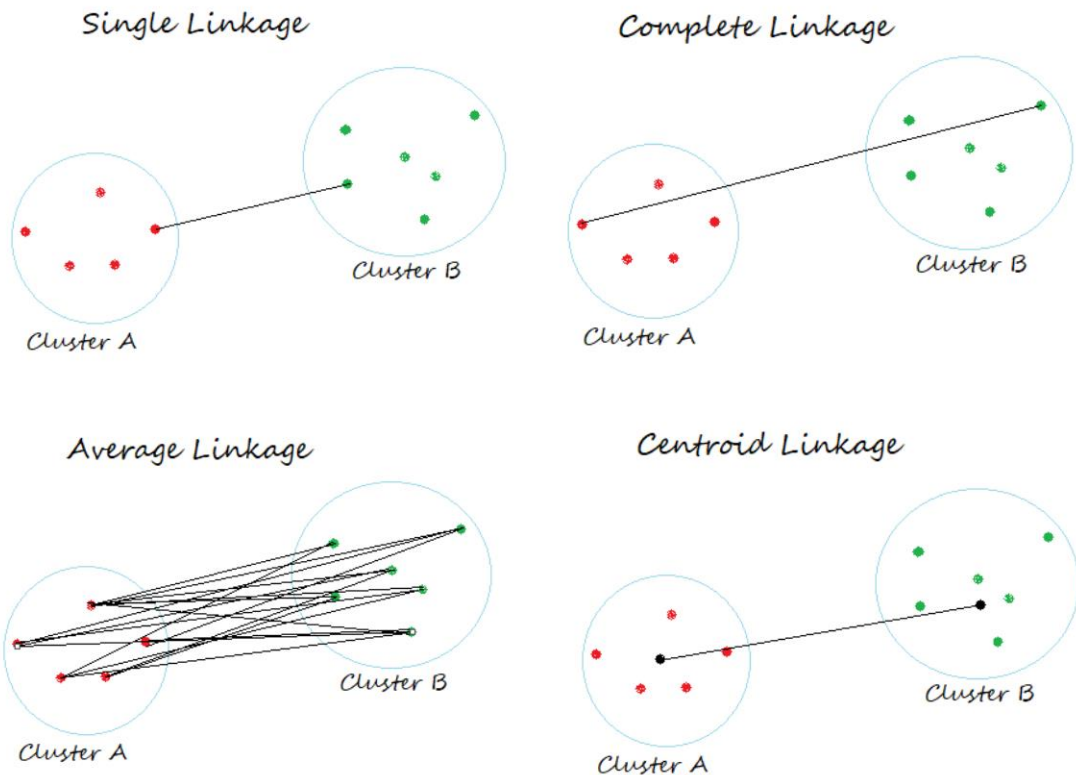
- Assume Euclidean distance
- Assume that clusters are of equal radius
- The initial guess of the locations of k means can affect the final clusters
 - Repeat multiple times

Agglomerative / Hierarchical clustering

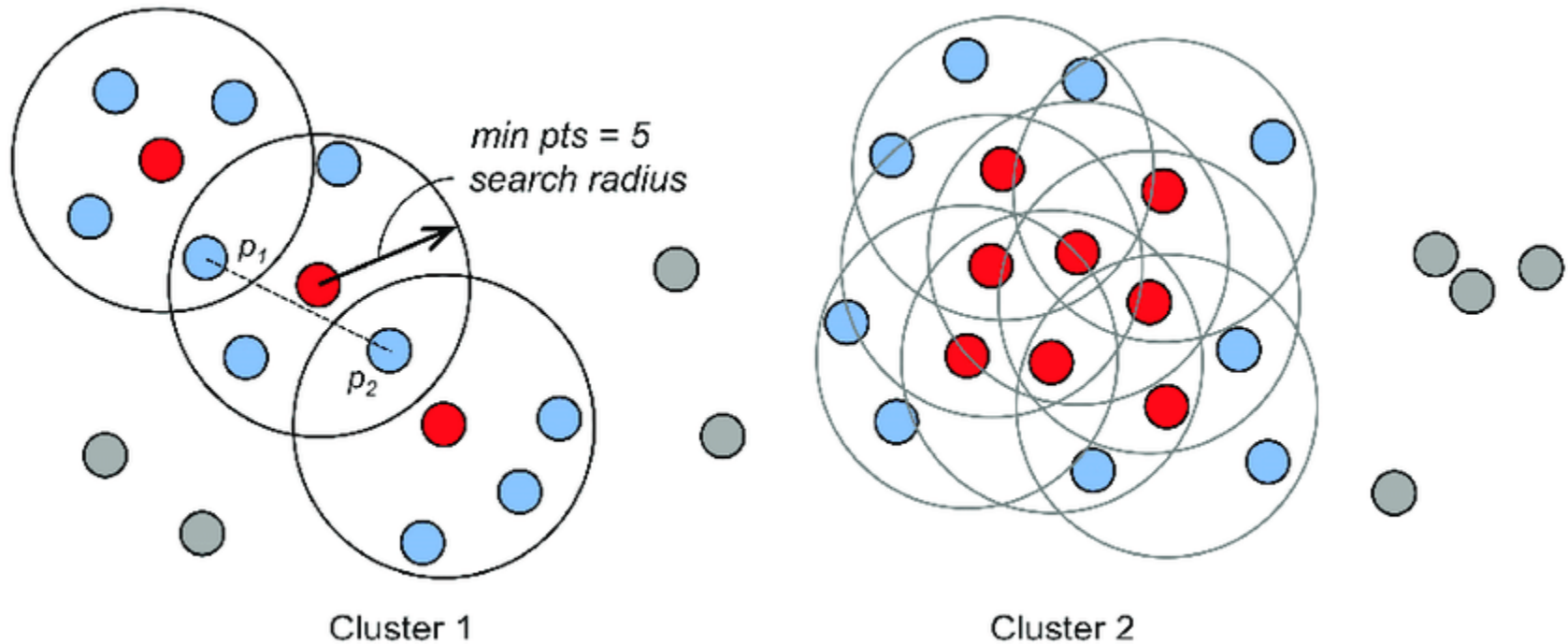


Source: www.slideshare.net/ElenaSgis/data-preprocessing-and-unsupervised-learning-methods-in-bioinformatics

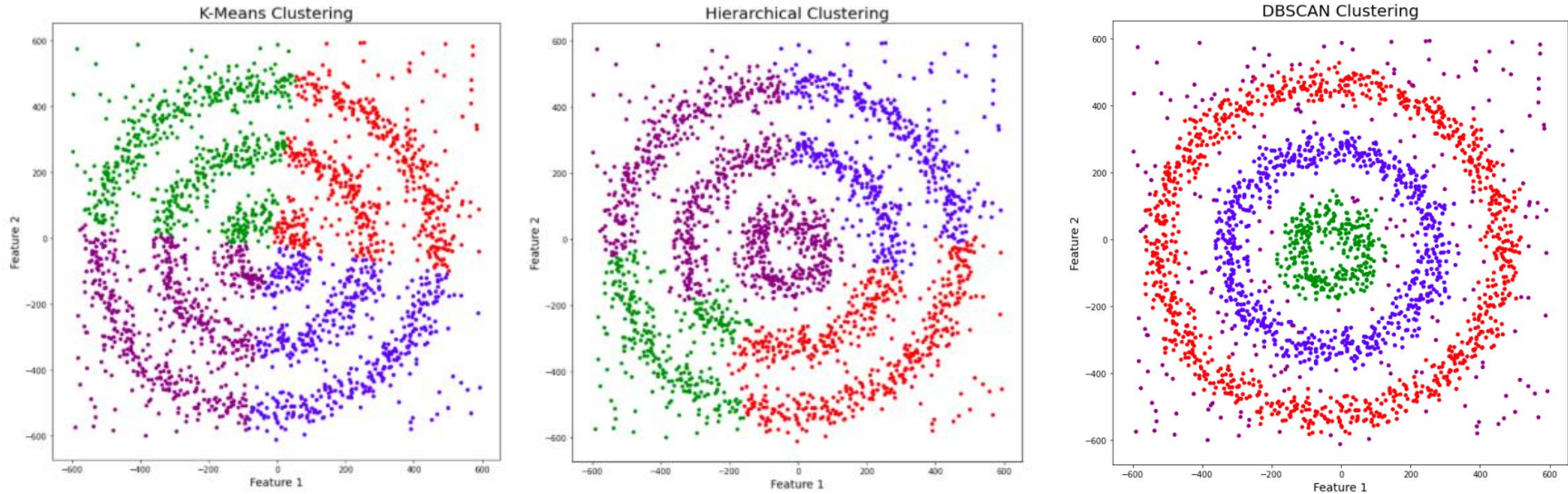
Linkage = distance metric for groups of data points



DBSCAN: A density-based technique



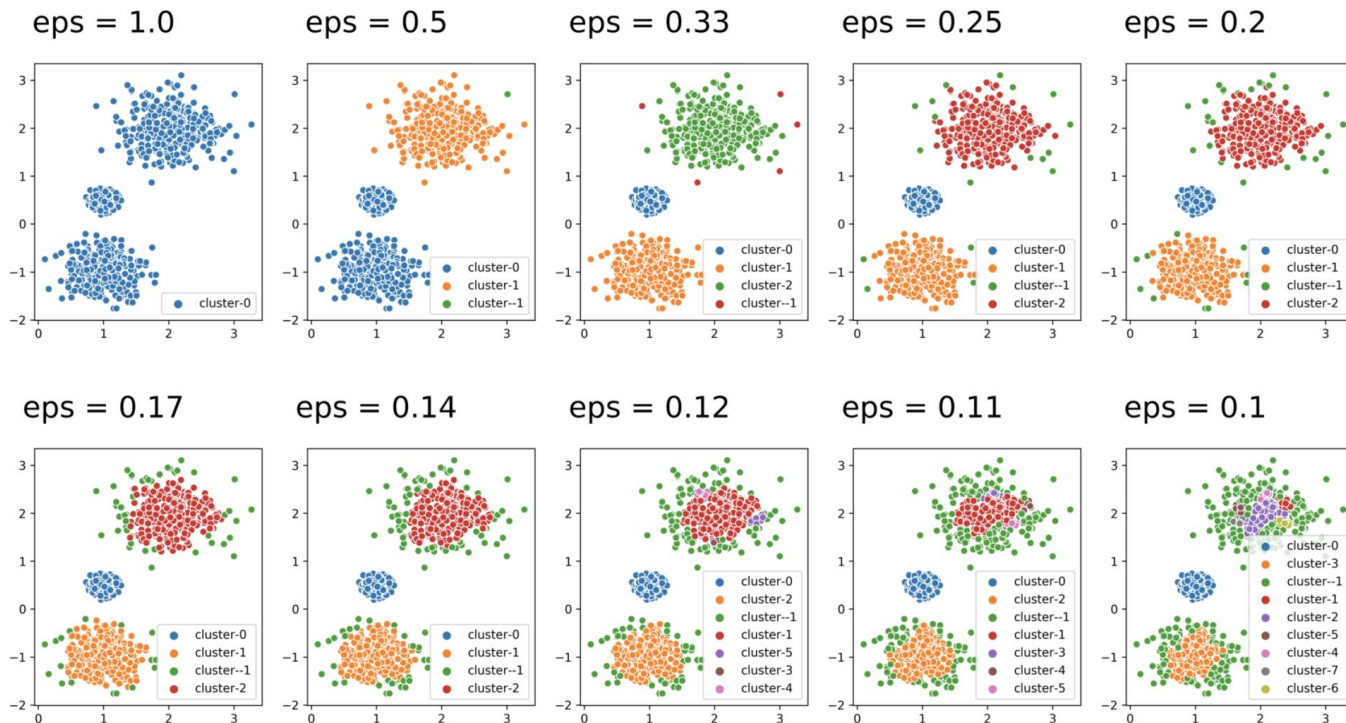
DBSCAN can handle complex cluster shape



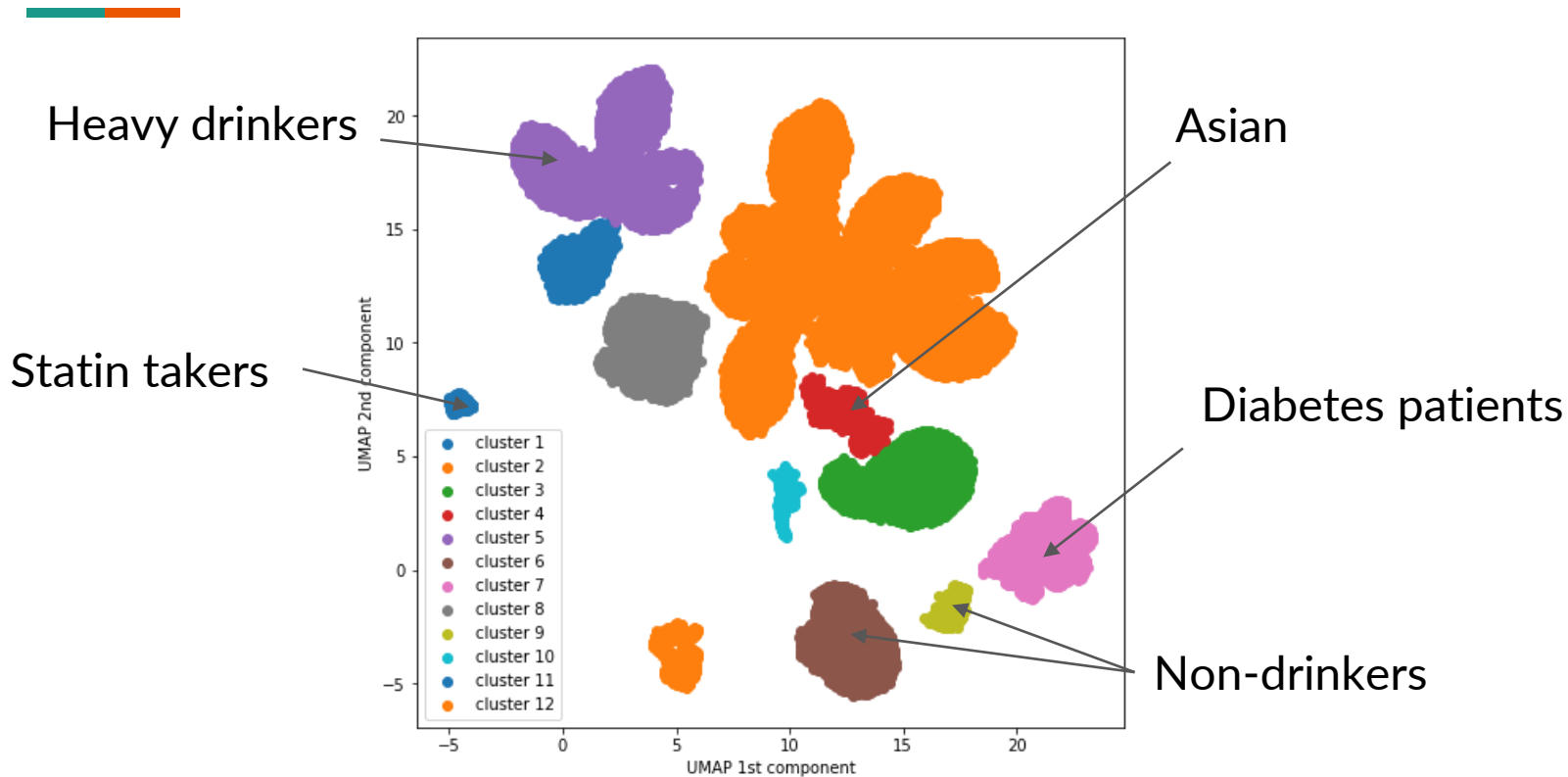
<https://www.analyticsvidhya.com/blog/2020/09/how-dbscan-clustering-works/>

- Distance-based techniques assume that data are spread in all directions

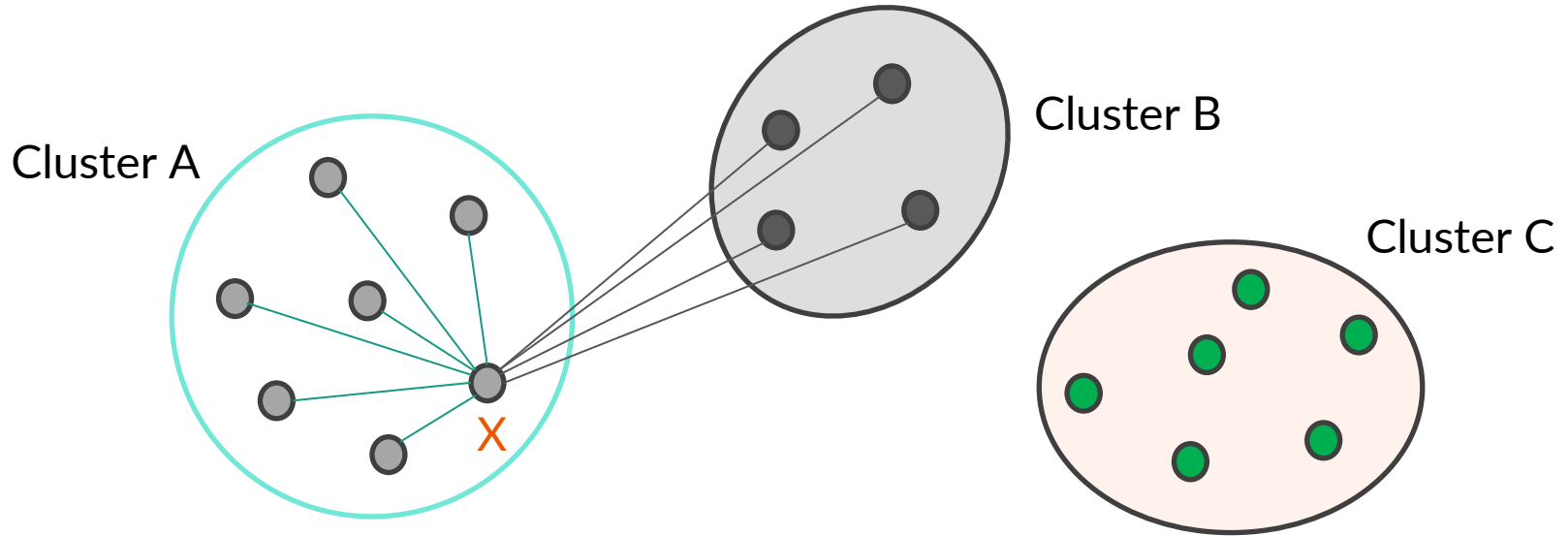
Tuning DBSCAN



An example: DBSCAN on UK Biobank data

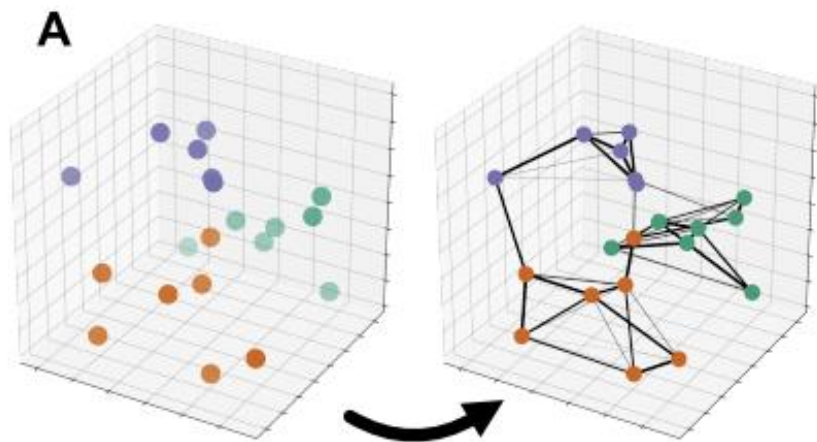


Silhouette score

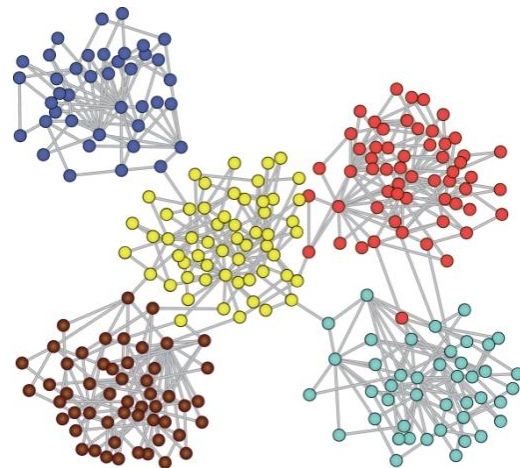


- Compare distances from **X** to other members of cluster A versus distances from **X** to members of cluster B (the closest cluster from A)

Spectral clustering and network clustering



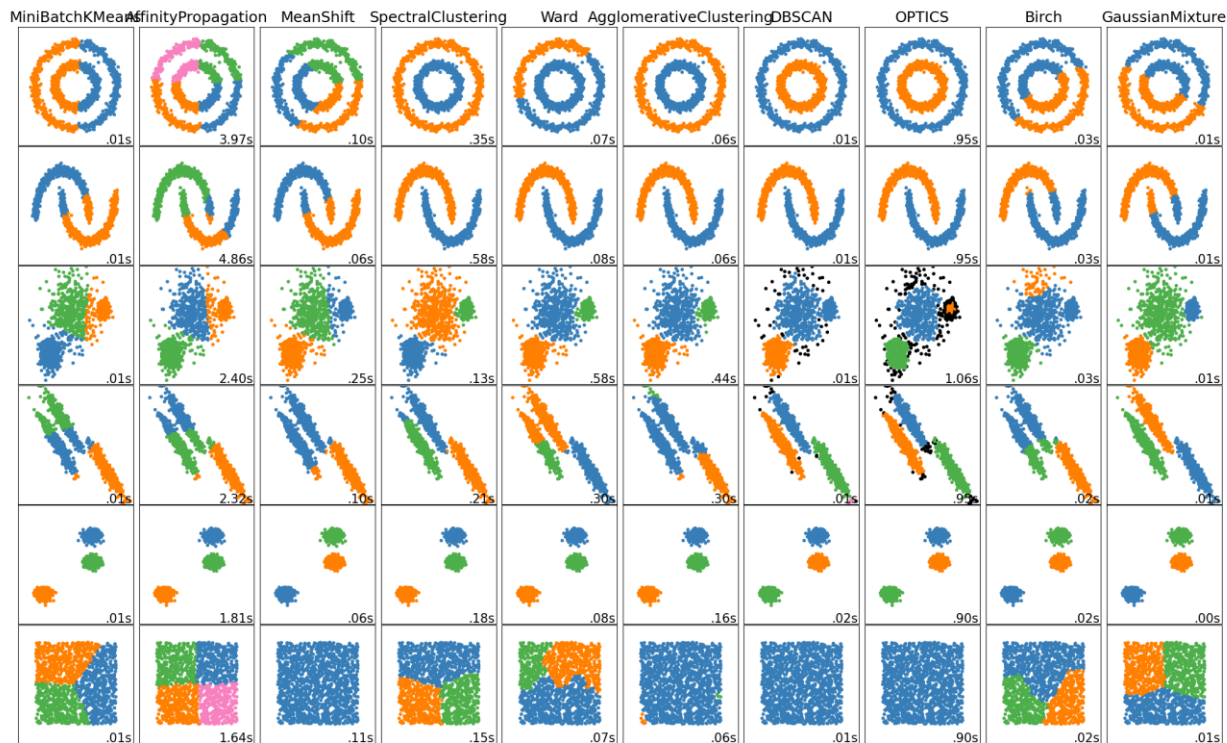
Sainburg, T. *et al.*, Neural Comput 33(11):2881-2907 (2021)



<https://github.com/topics/graph-clustering>

- View distance matrix as network
- Apply some threshold on the distance to create sparse network
- Split network into **modules with dense edges**

No one-size fits all



https://scikit-learn.org/0.23/auto_examples/cluster/plot_cluster_comparison.html

Summary



- High-dimensional data can often be simplified onto a 2D/3D visualization
- Explore distribution of feature values on the 2D/3D plot
- Picking appropriate distance metric is the key!
 - Using all features vs informative features
 - Euclidean(x, y) = $\sqrt{(x_1 - y_1)^2 + (x_2 - y_2)^2 + \dots + (x_n - y_n)^2}$
- Unsupervised learning needs domain knowledge for soft validation

Any question?



- See you next week on Sep 19th 9-10:30am