# UMAP in context of other dim. reduction methods

Anh Vu Le
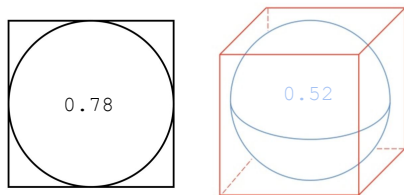
# Outline

1. Introduction
2. PCA
3. Methods related to UMAP
   a. ISOMAP
   b. t-SNE
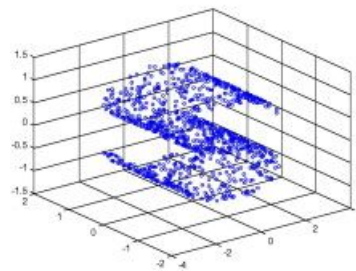4. Taxonomy of dimensionality reduction methods
5. UMAP

# What is dimensionality reduction about?

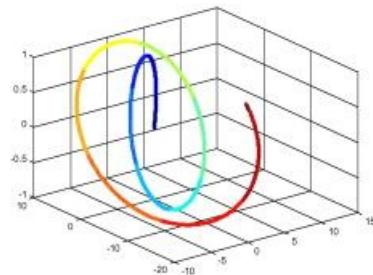*" ...transformation of **high-dimensional data** into a **meaningful representation of reduced dimensionality"***

- **intrinsic** dimensions of data (~ **manifold**)
  - smaller dimension, own geometry
- **mitigates the curse of dimensionality**



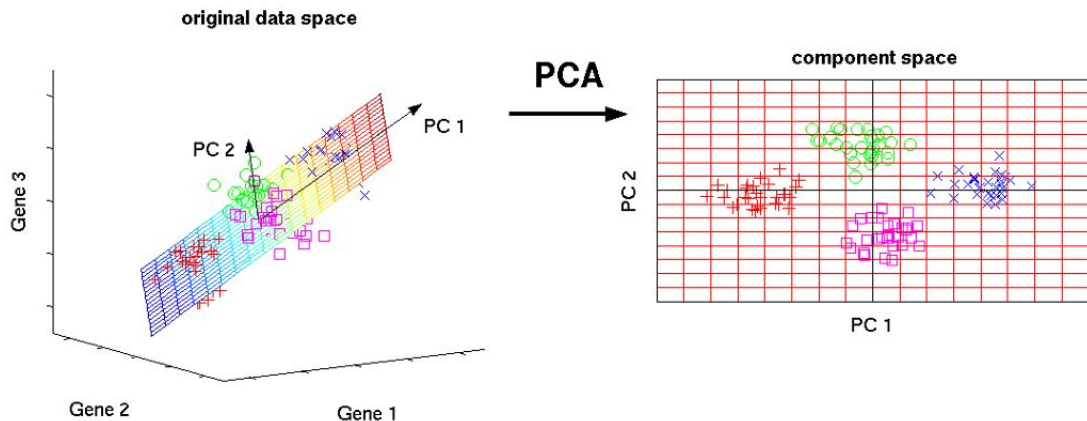- **classification**, **visualization**, and **compression**



(d)



(g)

# PCA - the queen of dimensionality reduction

*"PCA transforms (possibly correlated) data **linearly into new properties** that are **not correlated**"*

Properties of transformation:

- **Preserves variance**
  - Axis are directions of greatest variance
- properties are **uncorrelated**

# The elegance of PCA

"PCA transforms *possibly correlated* data **linearly into new properties** that are **not correlated**"

$cov(\mathbf{X}) = \mathbf{C_X} = \mathbf{X}^T\mathbf{X}$ is a covariance matrix

We want to find the projection $\mathbf{P}$ to uncorralated space, s.t.:
$$\mathbf{XP} = \mathbf{T}$$
where the uncorrelateness is expressed in the fact, that $\mathbf{C_T}$ is **diagonal**

# The elegance of PCA

PCA transforms *possibly correlated* **data linearly into new properties** that are **not correlated** with each other

$cov(\mathbf{X}) = \mathbf{C_X} = \mathbf{X}^T \mathbf{X}$ is a covariance matrix

We want to find the projection $\mathbf{P}$ to uncorralated space, s.t.:
$$\mathbf{XP} = \mathbf{T}$$
where the uncorrelateness is expressed in the fact, that $\mathbf{C_T}$ is **diagonal**

The covariance matrix of $\mathbf{T}$ can be expressed as:
$$\boxed{\mathbf{C_T}} = \mathbf{T}^T\mathbf{T} = (\mathbf{XP})^T\mathbf{XP} = \mathbf{P}^T(\mathbf{X}^T\mathbf{X})\mathbf{P} = \boxed{\mathbf{P}^T\mathbf{C_X}\mathbf{P}}$$

We know, that $\mathbf{C_X}$ is a **symetric** matrix.
The fine property of symetrix matrices, that their eigenvectors are **orthogonal** and:
$$\boxed{\mathbf{C_X} = \mathbf{EDE}^T} \text{ where D is diagonal matrix}$$

Now the magic: what if $\mathbf{P}$ is said to be $\mathbf{E}$?

# The elegance of PCA

$$\mathbf{C_T} = \mathbf{P}^T \mathbf{C_X} \mathbf{P}$$
$$\mathbf{C_X} = \mathbf{EDE}^T$$

If $\mathbf{P}$ is said to be $\mathbf{E}$, then
$$\mathbf{C_T} = \mathbf{P}^T \mathbf{C_X} \mathbf{P} = \mathbf{P}^T (\mathbf{EDE}^T) \mathbf{P} = \mathbf{P}^T (\mathbf{PDP}^T) \mathbf{P} = \mathbf{D}$$
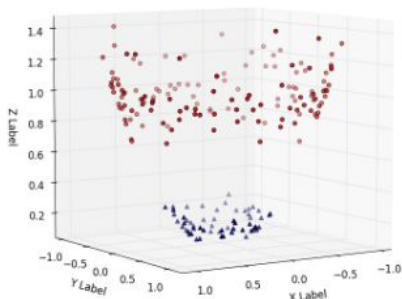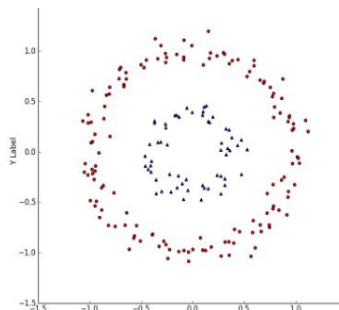
Choosing the **projection matrix** to be
**eigenvectors** of the **covariance matrix of *X***
gives us the projection we want:
- uncorrelated properties (diagonal $C_T$)
- no loss of information (P is orthogonal)

# Kernel PCA

- Some data linearly inseparable



$$cov(\mathbf{X}) \rightarrow K(\mathbf{X})$$
$$k_{i,j} = \kappa(\mathbf{x}_i, \mathbf{x}_j)$$

Transformation matrix is again the eigenmatrix of $K$

- Vapnik–Chervonenkis theory - projection into a **higher dimensional space** may provide us with **better classification power**.

- **Kernel trick**, a method to project original data into higher dimension **without sacrificing too much computational time**
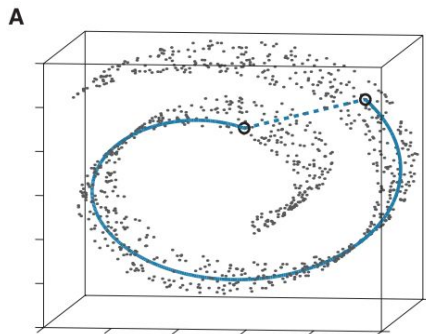
# Focusing on local patches of manifold

**ISOMAP**
- Close points in original space don't need to be close on manifold
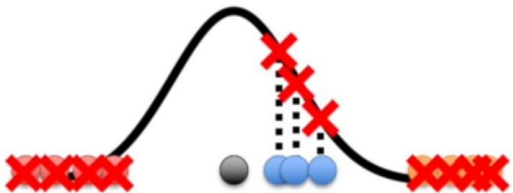
**Steps**:
1. Construct a **neighbourhood** graph (B)
2. Compute **shortest paths** between points in the graph (B)
3. Embed points with knowledge shortest paths in **lower dimension**

# Focusing on local patches of manifold (**t-SNE**)

**Key koncept:**   t-SNE maps **distances** to **probabilities**

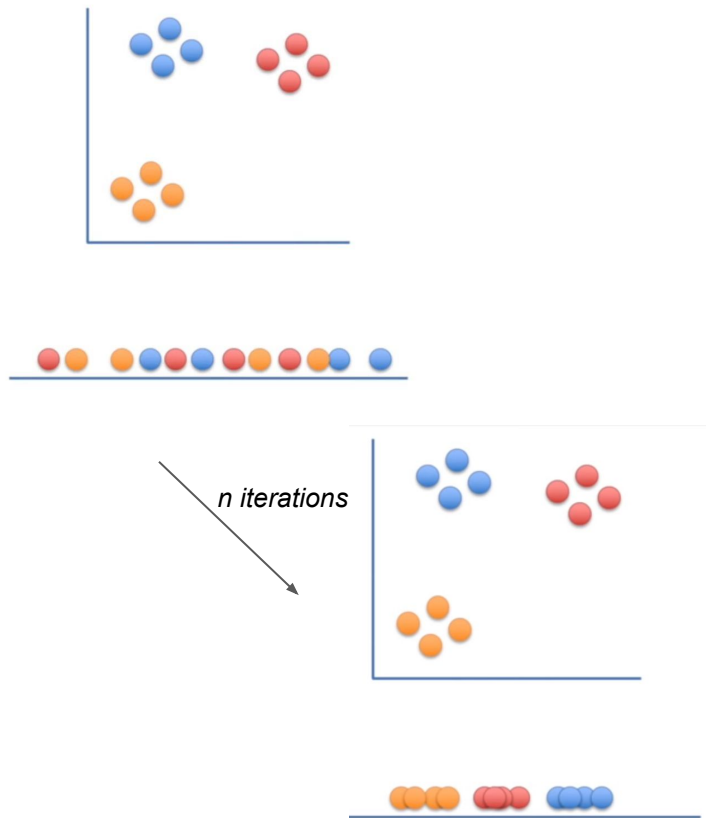**Each point** in original space forms a **_Gaussian_** around itself

On **lower dimension,** place points **randomly** at first

**Each point** in lower dim. forms a **_t-distribution_** around itself

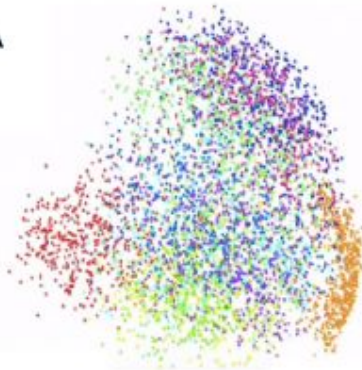**Minimize entropy:**    $C = \sum_{i \neq j} p_{ij} \log \frac{p_{ij}}{q_{ij}}$

$p_{ij}$...probability distance in original space
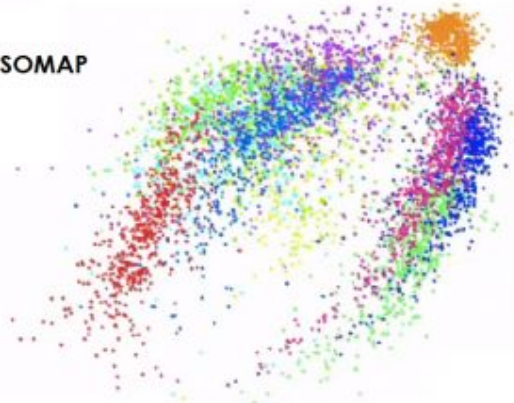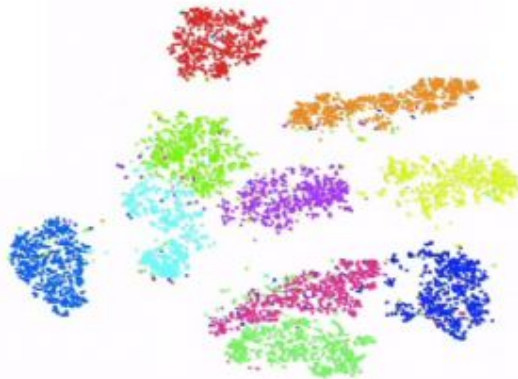$q_{ij}$...probability distance in reduced space
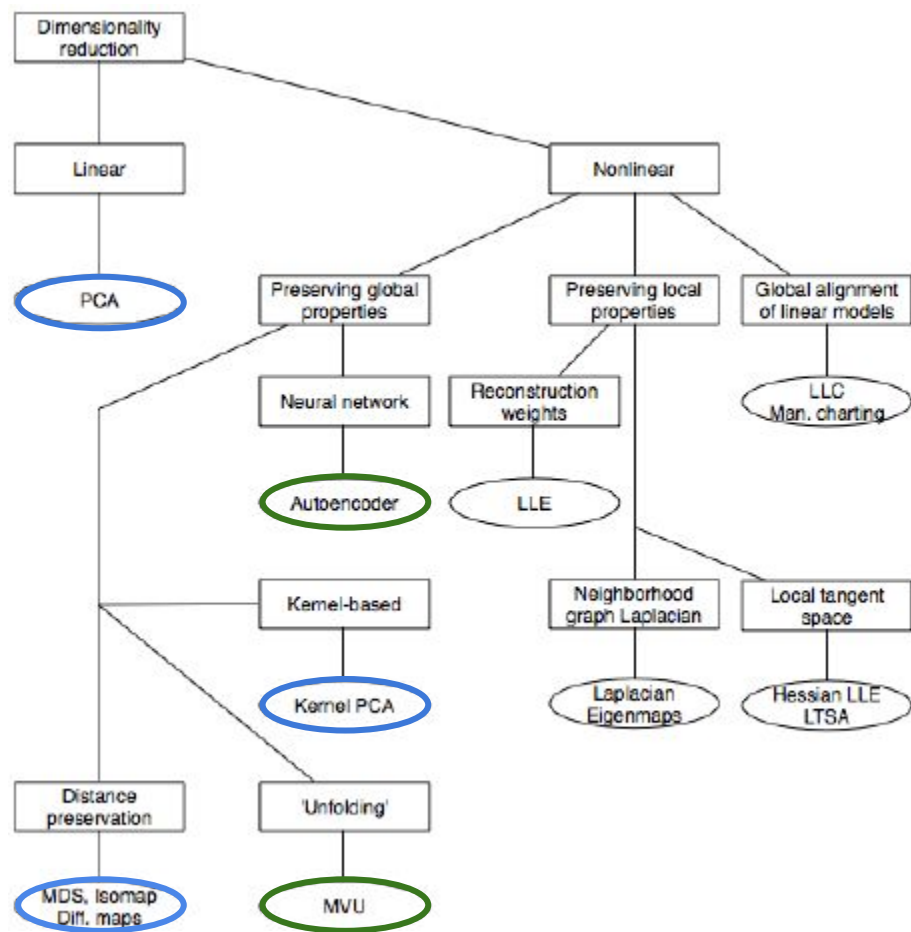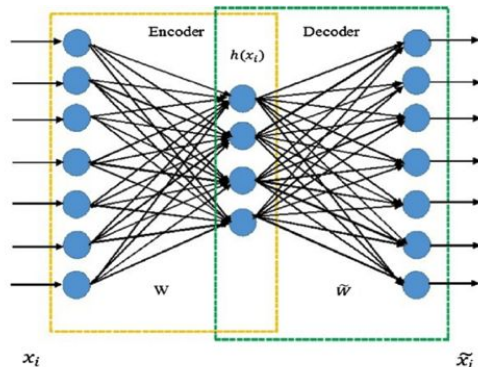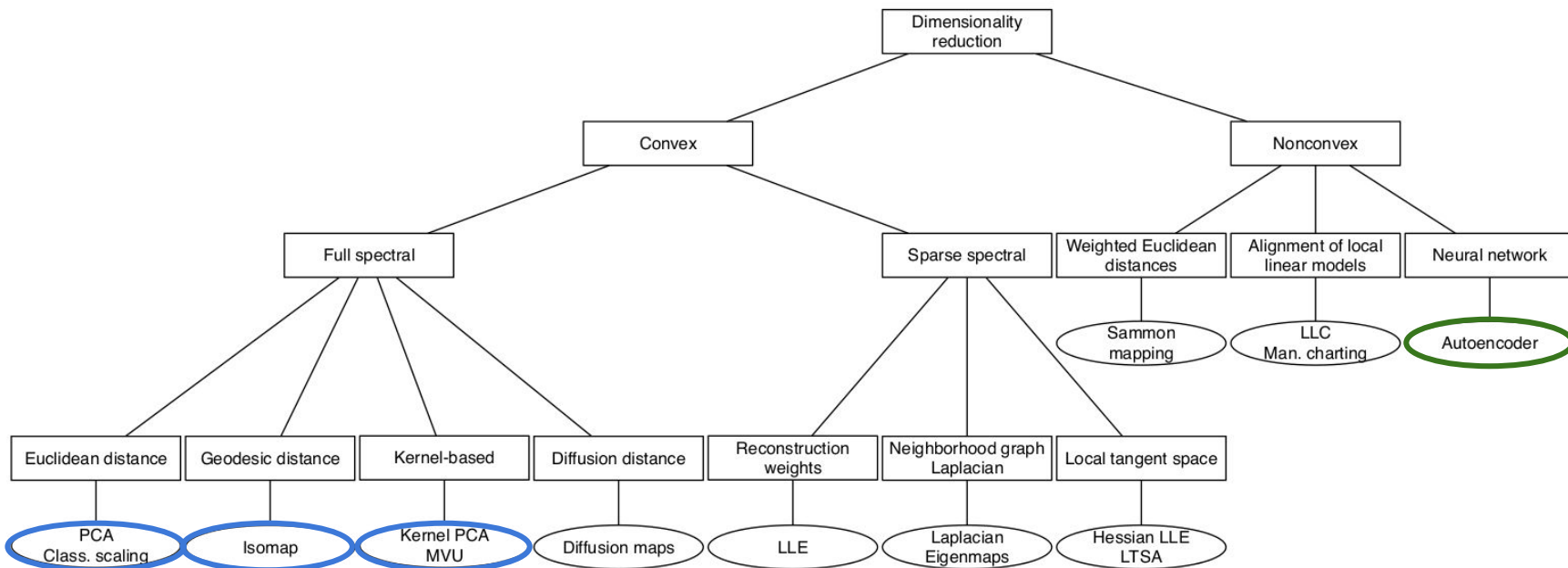
*n iterations*

# t-SNE in practice

# Taxonomy of techniques

- **MDS** (Multidimensional scaling)
  - "Fabricating" coordinate system based on similarity/distances
  - 3rd step of ISOMAP

- **MVU** (Maximum Variance Unfolding)
  - Learning kernel function for kernel PCA

- **Autoencoder**

# Taxonomy of techniques

# UMAP

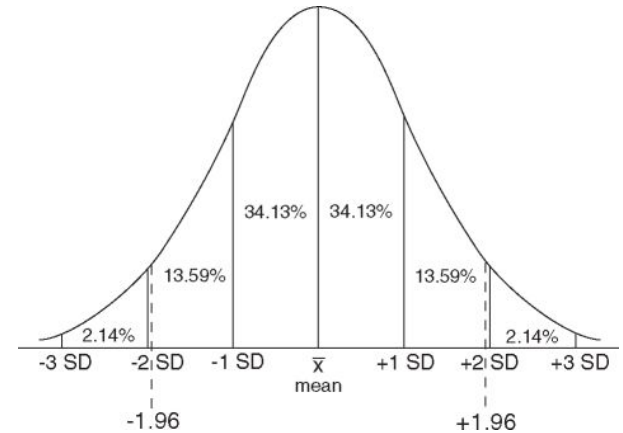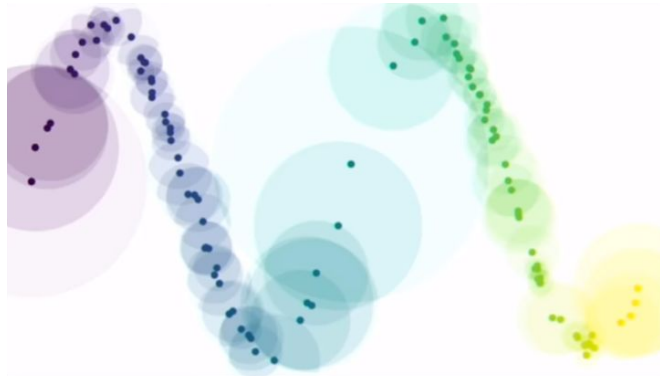| | t-SNE | UMAP |
|---|---|---|
| COIL20 | 20 seconds | 7 seconds |
| MNIST | 22 minutes | 98 seconds |
| Fashion MNIST | 15 minutes | 78 seconds |
| GoogleNews | 4.5 hours | 14 minutes |

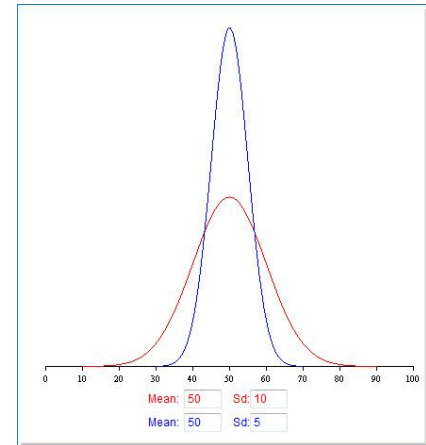**Algorithm 2** Constructing a local fuzzy simplicial set

**function** LocalFuzzySimplicialSet($X$, $x$, $n$)
    knn, knn-dists ← ApproxNearestNeighbors($X$, $x$, $n$)
    $\rho$ ← knn-dists[1]                   ▷ Distance to nearest neighbor
    $\sigma$ ← SmoothKNNDist(knn-dists, $n$, $\rho$)      ▷ Smooth approximator to
knn-distance
    fs-set$_0$ ← $X$
    fs-set$_1$ ← $\{([x, y], 0) \mid y \in X\}$
    **for all** $y \in$ knn **do**
        $d_{x,y}$ ← $\max\{0, \text{dist}(x, y) - \rho\}/\sigma$
        fs-set$_1$ ← fs-set$_1 \cup ([x, y], \exp(-d_{x,y}))$
    **return** fs-set

**Algorithm 3** Compute the normalizing factor for distances $\sigma$

**function** SmoothKNNDist(knn-dists, $n$, $\rho$)
    Binary search for $\sigma$ such that $\sum_{i=1}^{n} \exp(-(\text{knn-dists}_i - \rho)/\sigma) = \log_2(n)$
    **return** $\sigma$





https://ebrary.net/74165/environment/sampling_theory
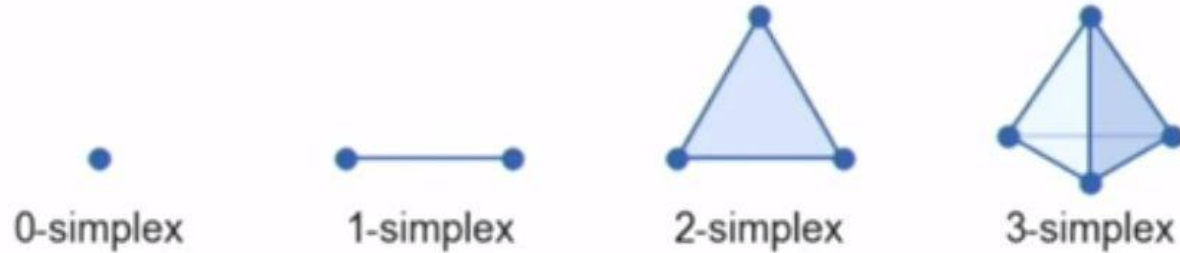


http://onlinestatbook.com/2/summarizing_distributions/spread_sim.html

*"...converting the [set of] **metric spaces** into **fuzzy simplicial sets**"*



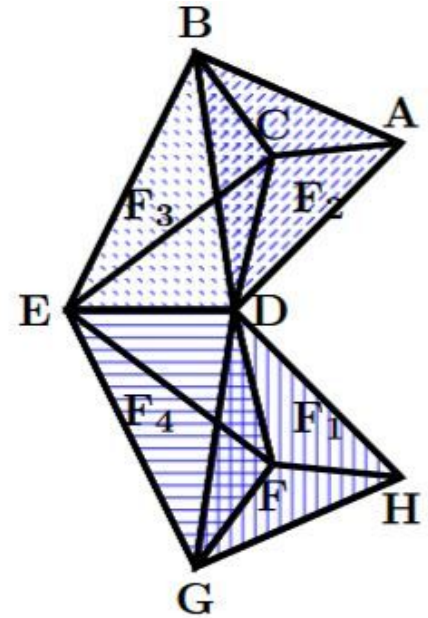0-simplex  1-simplex  2-simplex  3-simplex



**Fuzzy set:**
- Carrier set *A*
- Membership function $\mu : A \rightarrow [0, 1]$

Example:
$\mu(x)$ for $x \in A$ to be the **membership strength** of $x$ to the set *A*.

# UMAP - cost function

Optimizes fuzzy set cross entropy:

$$C_{UMAP} = \sum_{i \neq j} \boxed{v_{ij} \log \left( \frac{v_{ij}}{w_{ij}} \right)} + \boxed{(1 - v_{ij}) \log \left( \frac{1 - v_{ij}}{1 - w_{ij}} \right)}$$

$v_{ij}$...membership strenght of $j$ to fuzzy set of $i$ in orig. space

$w_{ij}$...membership strenght of $j$ to fuzzy set of $i$ in red. space

$$C = \sum_{i \neq j} p_{ij} \log \frac{p_{ij}}{q_{ij}}$$

$p_{ij}$...probability distance in original space

$q_{ij}$...probability distance in reduced space

---
**Algorithm 5** Optimizing the embedding
---
**function** OptimizeEmbedding(top-rep, $Y$, min-dist, n-epochs)

$\quad \alpha \leftarrow 1.0$

$\quad$ Fit $\Phi$ from $\Psi$ defined by min-dist

$\quad$ **for** $e \leftarrow 1, \ldots,$ n-epochs **do**

$\quad\quad$ **for all** $([a, b], p) \in$ top-rep$_1$ **do**

$\quad\quad\quad$ **if** Random( ) $\leq p$ **then** $\qquad \triangleright$ Sample simplex with probability $p$

$\quad\quad\quad\quad \boxed{y_a \leftarrow y_a + \alpha \cdot \nabla(\log(\Phi))(y_a, y_b)}$

$\quad\quad\quad\quad$ **for** $\imath \leftarrow 1, \ldots,$ n-neg-samples **do**

$\quad\quad\quad\quad\quad c \leftarrow$ random sample from Y

$\quad\quad\quad\quad\quad \boxed{y_a \leftarrow y_a + \alpha \cdot \nabla(\log(1 - \Phi))(y_a, y_c)}$

$\quad \alpha \leftarrow 1.0 - e/$n-epochs

$\quad$ **return** $Y$
---

# Sources

McInnes, Leland, John Healy, and James Melville. "Umap: Uniform manifold approximation and projection for dimension reduction." *arXiv preprint arXiv:1802.03426* (2018).

He, Jinrong, et al. "Intrinsic dimensionality estimation based on manifold assumption." *Journal of Visual Communication and Image Representation* 25.5 (2014): 740-747.

van der Maaten, Laurens & Postma, Eric & Herik, H.. (2007). Dimensionality Reduction: A Comparative Review. Journal of Machine Learning Research - JMLR. 10.

https://blog.datawow.io/the-curse-of-dimensionality-c99409eb58e9

https://towardsdatascience.com/kernel-pca-vs-pca-vs-ica-in-tensorflow-sklearn-60e17eb15a64

https://towardsdatascience.com/understanding-pca-autoencoders-algorithms-everyone-can-understand-28ee89b570e2

https://medium.com/bluekiri/understanding-principal-component-analysis-once-and-for-all-9f75e7b33635

http://blog.thegrandlocus.com/img/Gauss_vs_Cauchy.png