# DATA SCIENCE
## PPGIA/PUCPR

Prof. Jean Paul Barddal

# REMINDERS

# Reminders

- Today we will discuss dimensionality reduction
- Next week we will have our test!

# DIMENSIONALITY REDUCTION

# Dimensionality reduction

- In opposition to feature selection, dimensionality reduction techniques decrease the dimensionality of a problem by **combining** features
- There are different techniques to achieve this goal
- The most famous is Principal Component Analysis (PCA)

# PRINCIPAL COMPONENT ANALYSIS

Variance and covariance

- Variance and covariance measure how "spread" a set of points are around their mean
- Variance is used for analyzing a single dimension
- Covariance measures how much each of the dimensions vary from the mean with respect to each other
- Covariance is measures between 2 dimensions to see if there is a relationship between them

# Covariance

- The covariance between 2 variables is computed by:

$$Cov(X,Y) = \frac{\sum_{i=1}^{n}(X_i - \bar{X})(Y_i - \bar{Y})}{(n-1)}$$

- If you have more than 2 variables, you need to compute a covariance matrix

# Covariance

- The exact value is not as important as its sign
- A **positive** value indicates both variables increase or decrease together
- A **negative** value indicates that while one variable increases, the other decreases, or vice-versa
- If covariance is **zero**, the two variables are independent from each other

But what about correlation?

- Correlation allowed us to infere the same thing
- Why do we need covariance?
- Covariance is used to find relationships between variables in high-dimensional scenarios, where visualization is difficult

# Principal component analysis (PCA)

- PCA is a technique used to simplify a dataset
- It is a linear transformation that chooses a new coordinate system for the dataset such that:
  - **the greatest variance by any projection lies on the first axis**: the 1st principal component (eigenvector with the largest eigenvalue)
  - **the second greatest variance lies on the y axis (2nd PC), and so forth** (eigenvector with the second largest eigenvalue, etc)
- PCA can be used for reducing dimensionality by eliminating later principal components

# Steps to use PCA

- Normalize the data
- Calculate the covariance matrix
- Calculate the eigenvalues and eigenvectors
- Choosing principal components
- Forming a feature vector
- Forming principal components

- **All but the first of these steps are covered in scikit-learn's PCA implementation**

PCA Limitations

- If the data does not follow a multidimensional normal (gaussian) distribution, the principal components extracted will be distorted

## Activity

- Let's run PCA on a dataset representing customers
- Each customer represents either a restaurant, a retail store, etc
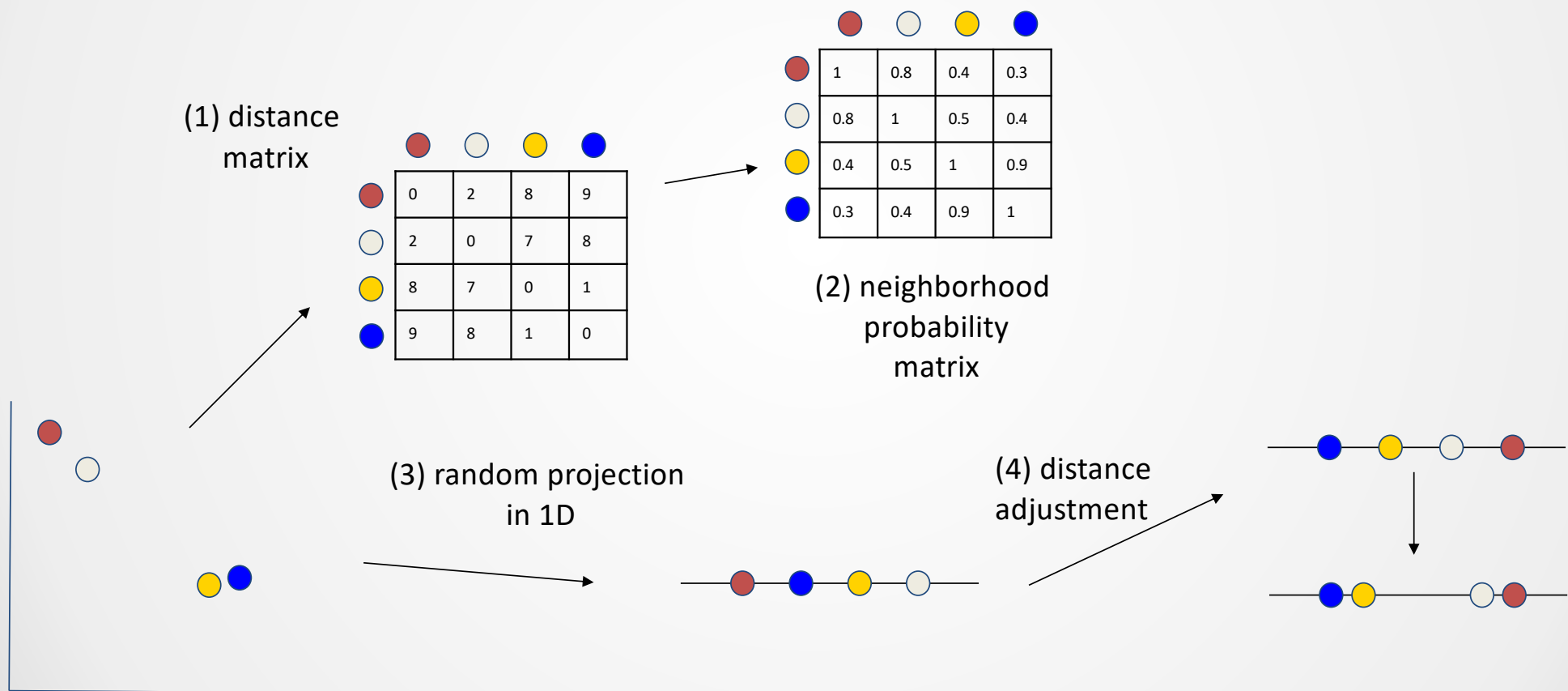- Let's analyze its principal components

# T-SNE

t-Stochastic Neighbor Embedding

- Technique tailored for visualizing high-dimensional datasets

- How do we visualize data in 2D or 3D?

- Two goals:
  - Distance preservation
  - Neighbor preservation

- Unsupervised, but it helps uncovering interesting aspects of the data

# t-SNE overall idea

- Let's say we have a 2D problem we wish to visualize in 1D

# t-SNE

---

**Algorithm 1**: Simple version of t-Distributed Stochastic Neighbor Embedding.

---

**Data**: data set $X = \{x_1, x_2, ..., x_n\}$,

cost function parameters: perplexity *Perp*,

optimization parameters: number of iterations $T$, learning rate $\eta$, momentum $\alpha(t)$.

**Result**: low-dimensional data representation $\mathcal{Y}^{(T)} = \{y_1, y_2, ..., y_n\}$.

**begin**

    compute pairwise affinities $p_{j|i}$ with perplexity *Perp* (using Equation 1)

    set $p_{ij} = \frac{p_{j|i} + p_{i|j}}{2n}$

    sample initial solution $\mathcal{Y}^{(0)} = \{y_1, y_2, ..., y_n\}$ from $\mathcal{N}(0, 10^{-4}I)$

    **for** $t=1$ **to** $T$ **do**

        compute low-dimensional affinities $q_{ij}$ (using Equation 4)

        compute gradient $\frac{\delta C}{\delta \mathcal{Y}}$ (using Equation 5)

        set $\mathcal{Y}^{(t)} = \mathcal{Y}^{(t-1)} + \eta \frac{\delta C}{\delta \mathcal{Y}} + \alpha(t)\left(\mathcal{Y}^{(t-1)} - \mathcal{Y}^{(t-2)}\right)$

    **end**

**end**

# t-SNE

---

**Algorithm 1**: Simple version of t-Distributed Stochastic Neighbor Embedding.

**Data**: data set $X = \{x_1, x_2, ..., x_n\}$,

cost function parameters: perplexity $Perp$,

optimization parameters: number of iterations $T$, learning rate $\eta$, momentum $\alpha(t)$.

**Result**: low-dimensional data representation $\mathcal{Y}^{(T)} = \{y_1, y_2, ..., y_n\}$.

**begin**

    compute pairwise affinities $p_{j|i}$ with perplexity $Perp$ (using Equation 1)

    set $p_{ij} = \frac{p_{j|i} + p_{i|j}}{2n}$

    sample initial solution $\mathcal{Y}^{(0)} = \{y_1, y_2, ..., y_n\}$ from $\mathcal{N}(0, 10^{-4}I)$

    **for** $t=1$ **to** $T$ **do**

        compute low-dimensional affinities $q_{ij}$ (using Equation 4)
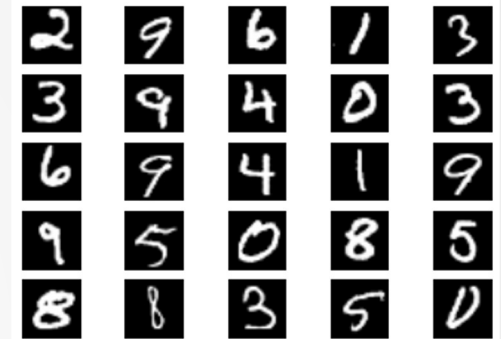
        compute gradient $\frac{\delta C}{\delta \mathcal{Y}}$ (using Equation 5)

        set $\mathcal{Y}^{(t)} = \mathcal{Y}^{(t-1)} + \eta \frac{\delta C}{\delta \mathcal{Y}} + \alpha(t)\left(\mathcal{Y}^{(t-1)} - \mathcal{Y}^{(t-2)}\right)$

    **end**

**end**



Random Sampling of MNIST

# t-SNE

**Algorithm 1**: Simple version of t-Distributed Stochastic Neighbor Embedding.

**Data**: data set $X = \{x_1, x_2, ..., x_n\}$,

cost function parameters: perplexity $Perp$,

optimization parameters: number of iterations $T$, learning rate $\eta$, momentum $\alpha(t)$.

**Result**: low-dimensional data representation $\mathcal{Y}^{(T)} = \{y_1, y_2, ..., y_n\}$.

**begin**

    compute pairwise affinities $p_{j|i}$ with perplexity $Perp$ (using Equation 1)

    set $p_{ij} = \frac{p_{j|i} + p_{i|j}}{2n}$

    sample initial solution $\mathcal{Y}^{(0)} = \{y_1, y_2, ..., y_n\}$ from $\mathcal{N}(0, 10^{-4}I)$

    **for** $t=1$ **to** $T$ **do**

        compute low-dimensional affinities $q_{ij}$ (using Equation 4)
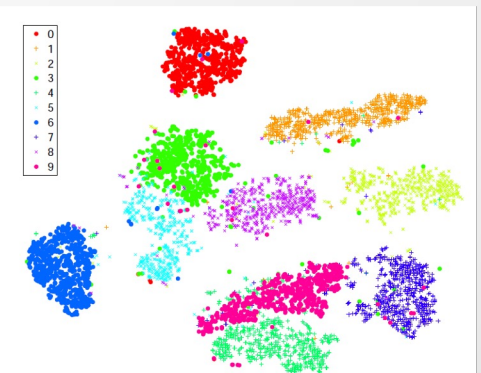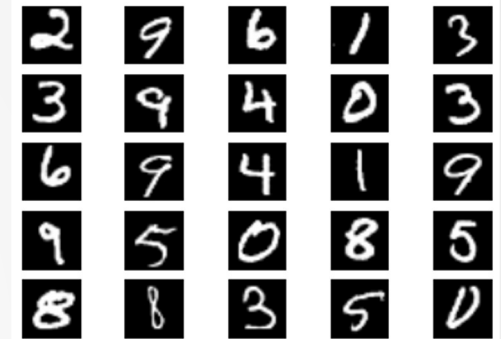
        compute gradient $\frac{\delta C}{\delta \mathcal{Y}}$ (using Equation 5)

        set $\mathcal{Y}^{(t)} = \mathcal{Y}^{(t-1)} + \eta \frac{\delta C}{\delta \mathcal{Y}} + \alpha(t) \left( \mathcal{Y}^{(t-1)} - \mathcal{Y}^{(t-2)} \right)$

    **end**

**end**



Random Sampling of MNIST



(a) Visualization by t-SNE.

# t-SNE

---

**Algorithm 1**: Simple version of t-Distributed Stochastic Neighbor Embedding.

**Data**: data set $X = \{x_1, x_2, ..., x_n\}$,

cost function parameters: perplexity $Perp$,

optimization parameters: number of iterations $T$, learning rate $\eta$, momentum $\alpha(t)$.

**Result**: low-dimensional data representation $\mathcal{Y}^{(T)} = \{y_1, y_2, ..., y_n\}$.

**begin**

    compute pairwise affinities $p_{j|i}$ with perplexity $Perp$ (using Equation 1)

    set $p_{ij} = \frac{p_{j|i} + p_{i|j}}{2n}$

    sample initial solution $\mathcal{Y}^{(0)} = \{y_1, y_2, ..., y_n\}$ from $\mathcal{N}(0, 10^{-4}I)$

    **for** $t=1$ **to** $T$ **do**

        compute low-dimensional affinities $q_{ij}$ (using Equation 4)
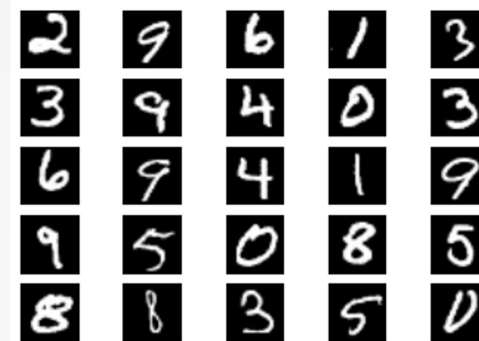
        compute gradient $\frac{\delta C}{\delta \mathcal{Y}}$ (using Equation 5)

        set $\mathcal{Y}^{(t)} = \mathcal{Y}^{(t-1)} + \eta \frac{\delta C}{\delta \mathcal{Y}} + \alpha(t)\left(\mathcal{Y}^{(t-1)} - \mathcal{Y}^{(t-2)}\right)$

    **end**

**end**



Random Sampling of MNIST

# t-SNE

**Algorithm 1**: Simple version of t-Distributed Stochastic Neighbor Embedding.

**Data**: data set $\mathcal{X} = \{x_1, x_2, ..., x_n\}$,
cost function parameters: perplexity $Perp$,
optimization parameters: number of iterations $T$, learning rate $\eta$, momentum $\alpha(t)$.
**Result**: low-dimensional data representation $\mathcal{Y}^{(T)} = \{y_1, y_2, ..., y_n\}$.
**begin**

    compute pairwise affinities $p_{j|i}$ with perplexity $Perp$ (using Equation 1)

    set $p_{ij} = \frac{p_{j|i} + p_{i|j}}{2n}$

    sample initial solution $\mathcal{Y}^{(0)} = \{y_1, y_2, ..., y_n\}$ from $\mathcal{N}(0, 10^{-4}I)$

    **for** $t=1$ **to** $T$ **do**

        compute low-dimensional affinities $q_{ij}$ (using Equation 4)

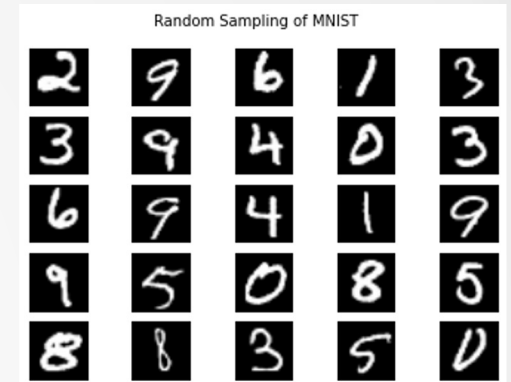        compute gradient $\frac{\delta C}{\delta \mathcal{Y}}$ (using Equation 5)

        set $\mathcal{Y}^{(t)} = \mathcal{Y}^{(t-1)} + \eta \frac{\delta C}{\delta \mathcal{Y}} + \alpha(t)\left(\mathcal{Y}^{(t-1)} - \mathcal{Y}^{(t-2)}\right)$

    **end**

**end**

Random Sampling of MNIST

**Perplexity** is the number of instances that we want to present the distances

Compute probabilities **P** that **xi** and **xj** are neighbors (based on Euclidian distance in **high-d** space)

# t-SNE

**Algorithm 1**: Simple version of t-Distributed Stochastic Neighbor Embedding.

**Data**: data set $X = \{x_1, x_2, ..., x_n\}$,

cost function parameters: perplexity $Perp$,

optimization parameters: number of iterations $T$, learning rate $\eta$, momentum $\alpha(t)$.

**Result**: low-dimensional data representation $\mathcal{Y}^{(T)} = \{y_1, y_2, ..., y_n\}$.

**begin**

    compute pairwise affinities $p_{j|i}$ with perplexity $Perp$ (using Equation 1)

    set $p_{ij} = \frac{p_{j|i} + p_{i|j}}{2n}$

    sample initial solution $\mathcal{Y}^{(0)} = \{y_1, y_2, ..., y_n\}$ from $\mathcal{N}(0, 10^{-4}I)$

    **for** $t=1$ **to** $T$ **do**

        compute low-dimensional affinities $q_{ij}$ (using Equation 4)

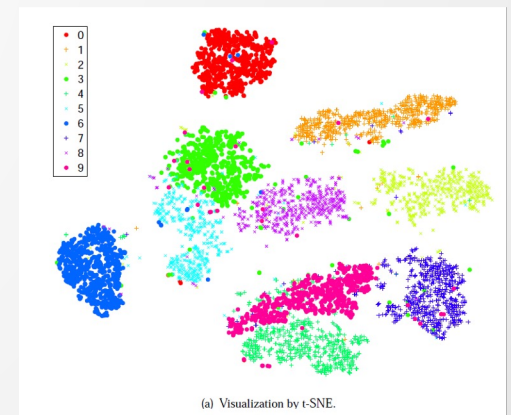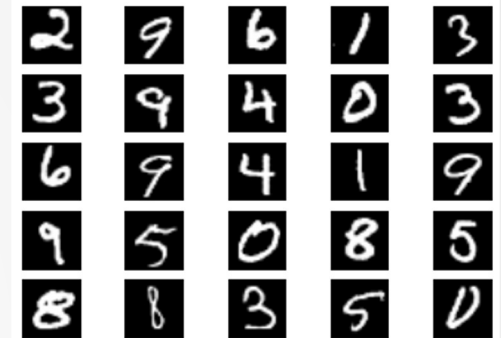        compute gradient $\frac{\delta C}{\delta \mathcal{Y}}$ (using Equation 5)

        set $\mathcal{Y}^{(t)} = \mathcal{Y}^{(t-1)} + \eta \frac{\delta C}{\delta \mathcal{Y}} + \alpha(t)\left(\mathcal{Y}^{(t-1)} - \mathcal{Y}^{(t-2)}\right)$

    **end**

**end**

Compute probabilities **Q** that **yi** and **yj** are neighbors (corresponding to *xi, xj*)
(based on Euclidian distance in **low-d** space)

Random Sampling of MNIST

(a) Visualization by t-SNE.

# t-SNE



**Algorithm 1**: Simple version of t-Distributed Stochastic Neighbor Embedding.

**Data**: data set $\mathcal{X} = \{x_1, x_2, ..., x_n\}$,
cost function parameters: perplexity $Perp$,
optimization parameters: number of iterations $T$, learning rate $\eta$, momentum $\alpha(t)$.
**Result**: low-dimensional data representation $\mathcal{Y}^{(T)} = \{y_1, y_2, ..., y_n\}$.
**begin**
    compute pairwise affinities $p_{j|i}$ with perplexity $Perp$ (using Equation 1)
    set $p_{ij} = \frac{p_{j|i} + p_{i|j}}{2n}$
    sample initial solution $\mathcal{Y}^{(0)} = \{y_1, y_2, ..., y_n\}$ from $\mathcal{N}(0, 10^{-4}I)$
    **for** $t=1$ **to** $T$ **do**
        compute low-dimensional affinities $q_{ij}$ (using Equation 4)
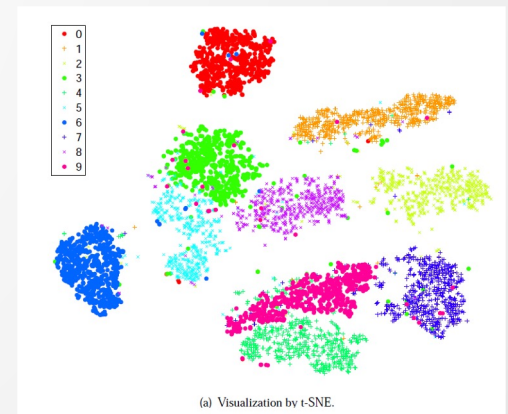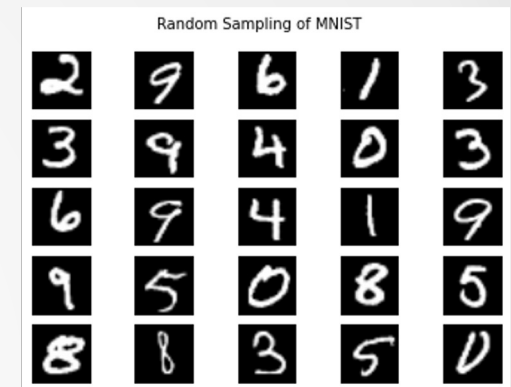        compute gradient $\frac{\delta C}{\delta \mathcal{Y}}$ (using Equation 5)
        set $\mathcal{Y}^{(t)} = \mathcal{Y}^{(t-1)} + \eta \frac{\delta C}{\delta \mathcal{Y}} + \alpha(t) \left( \mathcal{Y}^{(t-1)} - \mathcal{Y}^{(t-2)} \right)$
    **end**
**end**

Key assumption is that the **high-d** *P* and the **low-d** *Q* probability distributions should be the same


Random Sampling of MNIST


(a) Visualization by t-SNE.

# t-SNE

**Algorithm 1**: Simple version of t-Distributed Stochastic Neighbor Embedding.

**Data**: data set $X = \{x_1, x_2, ..., x_n\}$,
cost function parameters: perplexity $Perp$,
optimization parameters: number of iterations $T$, learning rate $\eta$, momentum $\alpha(t)$.
**Result**: low-dimensional data representation $\mathcal{Y}^{(T)} = \{y_1, y_2, ..., y_n\}$.
**begin**
> compute pairwise affinities $p_{j|i}$ with perplexity $Perp$ (using Equation 1)
> set $p_{ij} = \frac{p_{j|i} + p_{i|j}}{2n}$
> sample initial solution $\mathcal{Y}^{(0)} = \{y_1, y_2, ..., y_n\}$ from $\mathcal{N}(0, 10^{-4}I)$
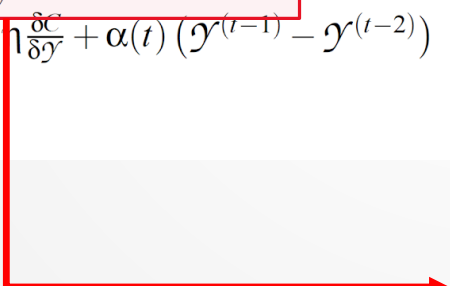> **for** $t=1$ **to** $T$ **do**
>> compute low-dimensional affinities $q_{ij}$ (using Equation 4)
>> compute gradient $\frac{\delta C}{\delta \mathcal{Y}}$ (using Equation 5)
>> set $\mathcal{Y}^{(t)} = \mathcal{Y}^{(t-1)} + \eta \frac{\delta C}{\delta \mathcal{Y}} + \alpha(t) \left(\mathcal{Y}^{(t-1)} - \mathcal{Y}^{(t-2)}\right)$
> **end**

**end**

Find a **low-d** map that minimizes the difference between the **P** (high-d) and **Q** (low-d) distributions

(if $xi, xj$ has high probability of being neighbors in **high-d**, then $yi, yj$ should have high probability in **low-d**)

# t-SNE

**Algorithm 1**: Simple version of t-Distributed Stochastic Neighbor Embedding.

**Data**: data set $X = \{x_1, x_2, ..., x_n\}$,

cost function parameters: perplexity *Perp*,

optimization parameters: number of iterations $T$, learning rate $\eta$, momentum $\alpha(t)$.

**Result**: low-dimensional data representation $\mathcal{Y}^{(T)} = \{y_1, y_2, ..., y_n\}$.

**begin**

    compute pairwise affinities $p_{j|i}$ with perplexity *Perp* (using Equation 1)

    set $p_{ij} = \frac{p_{j|i} + p_{i|j}}{2n}$

    sample initial solution $\mathcal{Y}^{(0)} = \{y_1, y_2, ..., y_n\}$ from $\mathcal{N}(0, 10^{-4}I)$

    **for** $t=1$ **to** $T$ **do**

        compute low-dimensional affinities $q_{ij}$ (using Equation 4)

        compute gradient $\frac{\delta C}{\delta \mathcal{Y}}$ (using Equation 5)

        set $\mathcal{Y}^{(t)} = \mathcal{Y}^{(t-1)} + \eta \frac{\delta C}{\delta \mathcal{Y}} + \alpha(t) \left( \mathcal{Y}^{(t-1)} - \mathcal{Y}^{(t-2)} \right)$

    **end**

**end**

We will minimize the difference between the
**high-d** and **low-d** maps using **gradient descent**

# t-SNE details

- Details on t-SNE can be found at the original paper
- [https://jmlr.org/papers/volume9/vandermaaten08a/vandermaaten08a.pdf](https://jmlr.org/papers/volume9/vandermaaten08a/vandermaaten08a.pdf)

## Visualizing Data using t-SNE

**Laurens van der Maaten**                                    LVDMAATEN@GMAIL.COM
*TiCC*
*Tilburg University*
*P.O. Box 90153, 5000 LE Tilburg, The Netherlands*

**Geoffrey Hinton**                                          HINTON@CS.TORONTO.EDU
*Department of Computer Science*
*University of Toronto*
*6 King's College Road, M5S 3G4 Toronto, ON, Canada*

# t-SNE

- In opposition to PCA, t-SNE is **not** parametric
- This means that we cannot learn a manifold representation from a dataset and apply it to another dataset
  - Therefore, this cannot be used as a dimensionality reduction technique between training and test data
- There is, however, a parametric version available at:

**Learning a Parametric Embedding by Preserving Local Structure**

**Laurens van der Maaten**
TiCC, Tilburg University
P.O. Box 90153, 5000 LE Tilburg, The Netherlands
lvdmaaten@gmail.com