POLITECNICO DI MILANO

# Clustering Text
## Natural Language Processing

ALICE was beginning to get very tired of sitting by her sister on the bank, and of having nothing to do: once or twice she had peeped into the book her sister was reading, but it had no pictures or conversations in it, " and what is the use of a book," thought Alice, " without pictures or conversations?" So she was considering in her own mind (as well as she could, for the hot day made her feel very sleepy and stupid,) whether the pleasure of making a daisy-chain would be worth the trouble of getting up and picking the daisies, when suddenly a white rabbit with pink eyes ran close by her. There was nothing so very remarkable in that; nor did Alice think it so very much out of the way to hear the Rabbit say to itself, " Oh dear ! Oh dear ! I shall be too late !" (when she thought it over afterwards, it occurred to her that she ought to have wondered at this, but at the time it all seemed quite natural); but when the Rabbit actually took a watch out of its waistcoat-pocket, and looked at it, and then hurried on, Alice started to her feet, for it flashed across her mind that she had never before seen a rabbit with either a waistcoat-pocket or a watch to take out of it, and

Miner Image source: https://freesvg.org/miner-157442484

# Lecture content:

- Why cluster text?
- Similarities between documents
- Basic clustering algorithms
- Evaluating clustering
- Topic Modelling
- Visualisation via dimensionality reduction
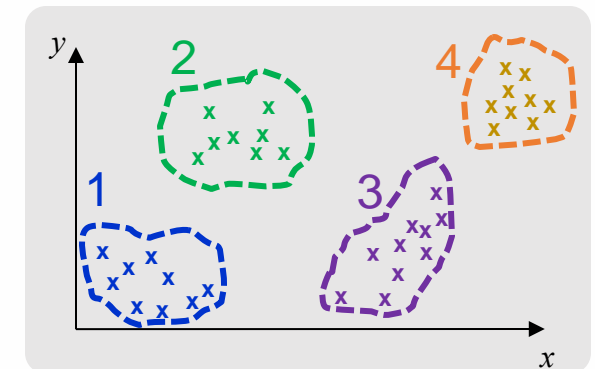
POLITECNICO DI MILANO
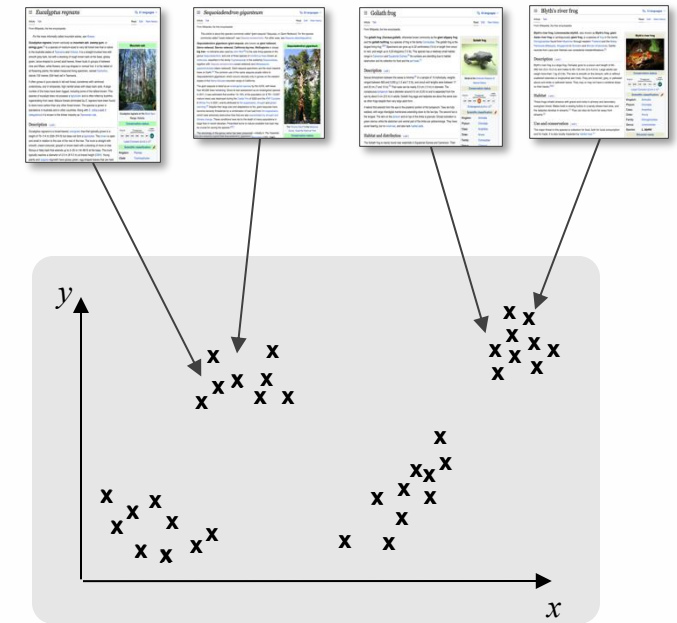
# Why cluster text?

# What is text clustering?

Process of **grouping documents** into **coherent subgroups**

- based on some **distance function** or **similarity measure**

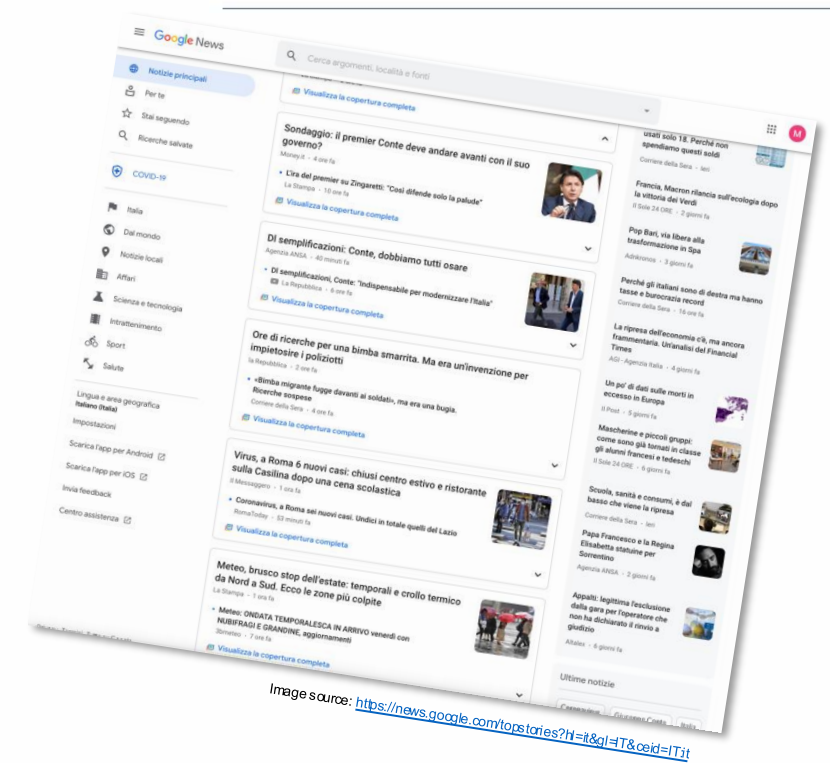For **text documents**, usually means grouping documents into:

- same **topic**: e.g. *sports news vs politics vs technology vs …*
- same **emotion**: e.g. *angry, disappointed, elated, hopeful, …*
- same **problem** to be fixed: e.g. *no heating (faulty thermostat), unable to print (missing drivers), no hot water, etc.*
- same **complaint/praise**: e.g. *hotel has beautiful view, motel on noisy road, breakfast was excellent, …*
- same **author**: e.g. *Shakespeare, Dylan, Eminem, Trump, Satoshi Nakamoto*
- same **text generator**: e.g. *used to generate spam emails*
- same **source document**: e.g. *for de-duplicating content during a web-crawl*

# Example Clustering Applications



Image source: https://news.google.com/topstories?hl=it&gl=IT&ceid=IT:it

News aggregators like Google News (https://news.google.com/)

use clustering to **summarise/characterise** a collection:

- what **topics** are being discussed?
- how **prevalent** is each topic?
- what topics are **trending**?
- which articles are most **exemplary** (central to topic)?

In security applications, such as phishing and malware detection,

- determine which groups might be malicious (based on a few identified instances)
- use clusters to discover useful **features** for describing documents
- what aspects distinguish the groups?

POLITECNICO DI MILANO

# Representing documents & measuring similarity

# Document similarity

To cluster documents need to **measure similarity** between them

- based on the similarity of their content
- which is subjective of course: *is Stephen King closer to Agatha Christie or J K Rowling?*

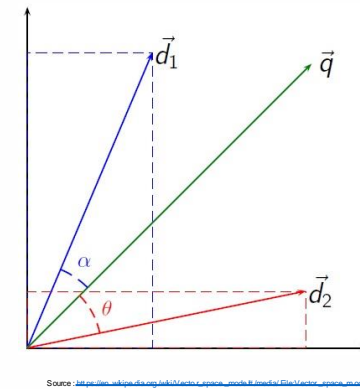Traditionally represent documents as **tf-idf vectors** (scaled **bag-of-words**)

- sparse representation: *most values in vector are zero*
- so measure similarity based on **number of shared words**
- **scaled by the importance** of the word: *how rare is it?*

Usually assume **similarity doesn't depend on length** of documents

- so compute **cosine similarity** between vectors
- = dot product between length normalised vectors

Source : http://en.wikipedia.org/wiki/Vector_space_model#/media/File:Vector_space_model.jpg

$$sim(\mathbf{d}_1, \mathbf{d}_2) = \frac{\mathbf{d}_1 \cdot \mathbf{d}_2}{||\mathbf{d}_1|| ||\mathbf{d}_2||} \qquad ||\mathbf{x}|| = \sqrt{\sum_{i=1}^{n} x_i^2}$$

# Clustering techniques

# Types of clustering algorithms

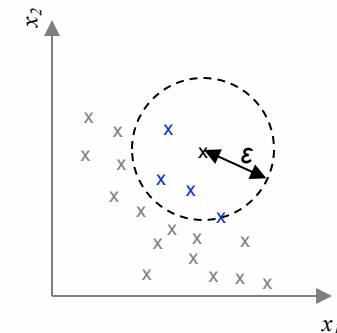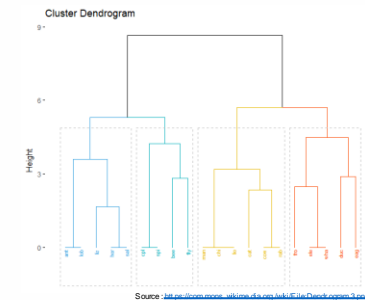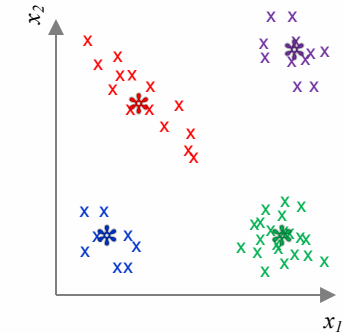Many different clustering algorithms in common use

- k-Means/k-Medoids
- Agglomerative hierarchical clustering
- Density-based clustering
- Spectral clustering
- Topic Models

All these models are implemented in sklearn:

- https://scikit-learn.org/stable/modules/clustering.html

We'll concentrate on those **most used with text**:

- primarily *k-Means* & **Topic Models**





Cluster Dendrogram

Source: https://commons.wikimedia.org/wiki/File:Dendrogram.png

# Clustering Algorithms:
# *k*-Means

# *k*-Means clustering

"Goto" clustering method
- simple, fast and relatively robust

How does it work?
- searches for exactly $k$ clusters
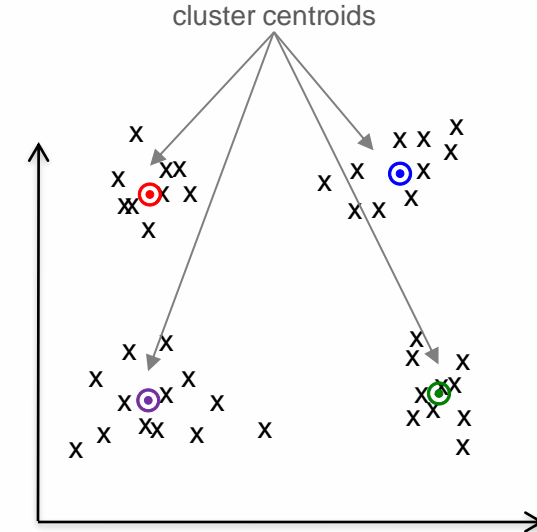- represents each cluster by its centroid

👍 Advantages:
- scales well to large collections
- (requires no pairwise distance calculations)

👎 Disadvantages:
- searches for globular clusters
- implicitly assumes Euclidean distance metric (not ideal for text)
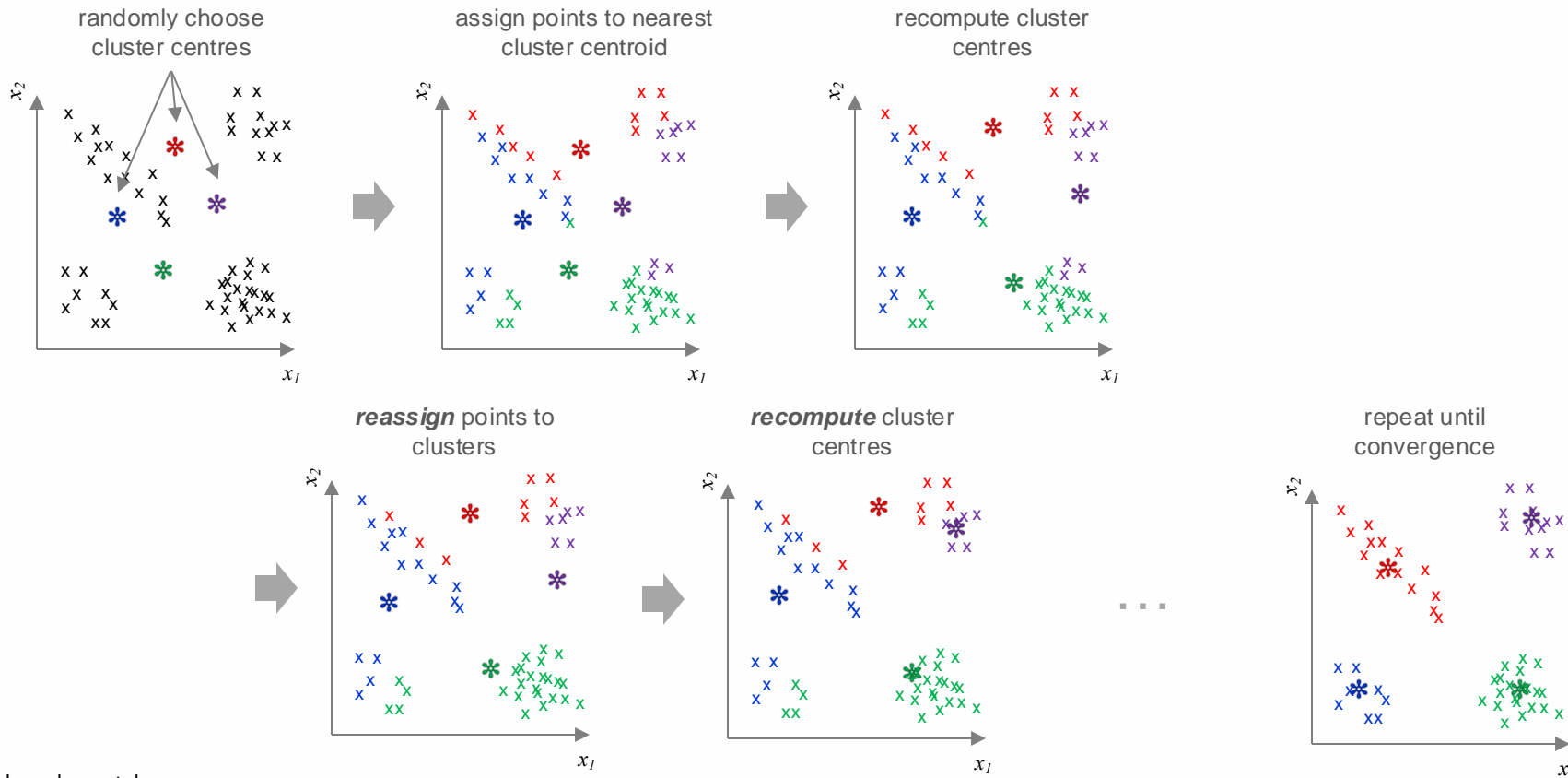- number of clusters must be specified in advance



cluster centroids

Centroid:

$$\vec{\mu}_C = \frac{1}{|C|} \sum_{\vec{x} \in C} \vec{x}$$

Euclidean distance:

$$d(x, y) = \sqrt{\sum_{i=1}^{m} (x_i - y_i)^2}$$

# *k*-Means clustering: the process



randomly choose cluster centres

assign points to nearest cluster centroid

recompute cluster centres

*reassign* points to clusters

*recompute* cluster centres

repeat until convergence

Simple algorithm:

1. initialize $k$ centroids randomly
2. (re)assign each data point to closest centroid
3. (re)compute centroids by averaging data points in cluster
4. repeat steps (2) & (3) until convergence, when centroids stop moving

algorithm **minimises variance of clusters** (squared distances between points & cluster centres)

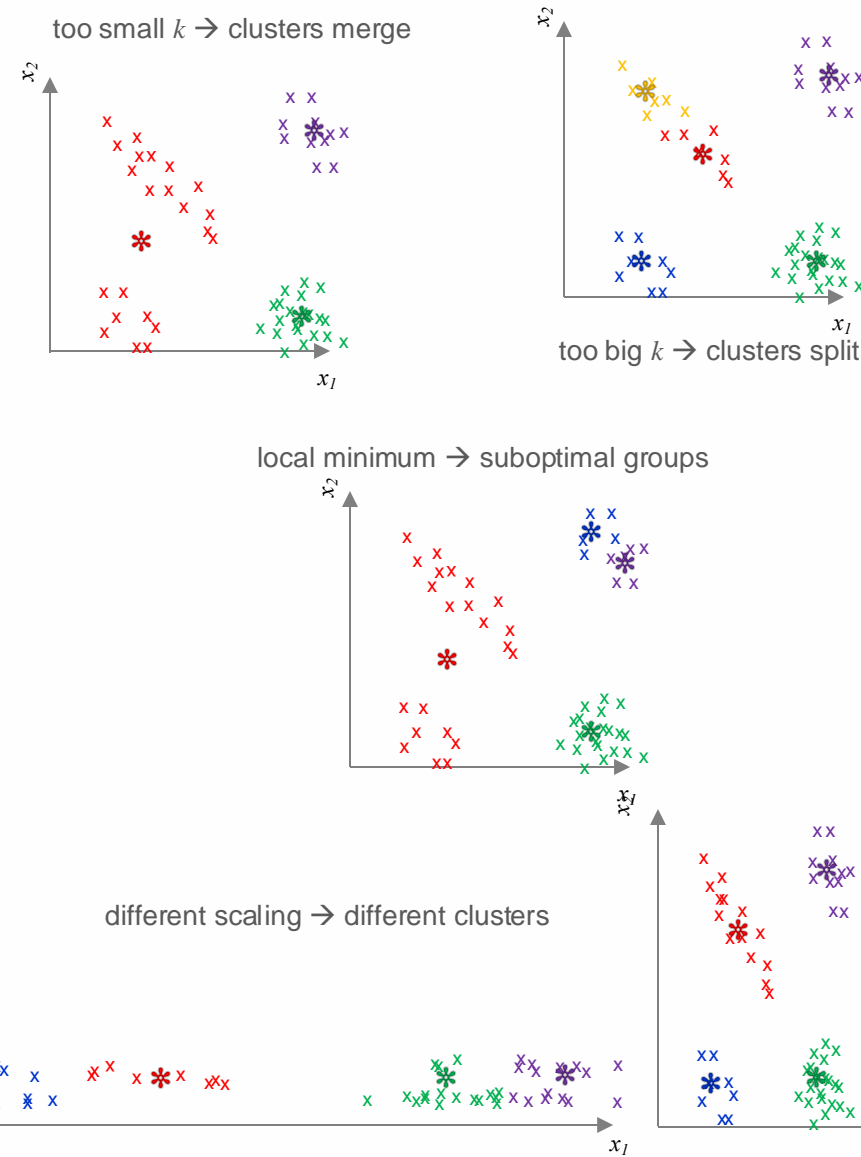$$\min \sum_{i=1}^{k} \sum_{\vec{x} \in C_i} d(\vec{x}, \vec{\mu}_{C_i})^2$$

# *k*-Means clustering – potential problems

Potential problems with k-means:

- choosing "right" value of $k$ is critical
  - can use elbow method to choose $k$
- algorithm can converge on **local minimum**
  - run algorithm multiple times algorithm

As for all clustering algorithms, **scaling effects** similarity measure
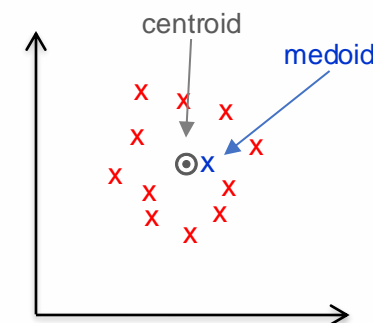
- and thus clusters found
- but for text have both tf-idf weighting and document length normalisaton

too small $k$ → clusters merge

too big $k$ → clusters split

local minimum → suboptimal groups

different scaling → different clusters

# Clustering Algorithms:
## *k*-medoids

# *k*-medoids

centroid

medoid

Represents each cluster by its **medoid** rather than its centroid

- **medoid** is one of the datapoints from dataset
  - point that is closest (*smallest average distance*) to all other points
- at each iteration, algorithm must:
  - reassign datapoints to cluster with closest medoid, and then recompute medoids

👍 Advantage:
  - clustering algorithm can be used with other metrics, not just Euclidean
  - medoid is a real document so provides realistic representation of cluster
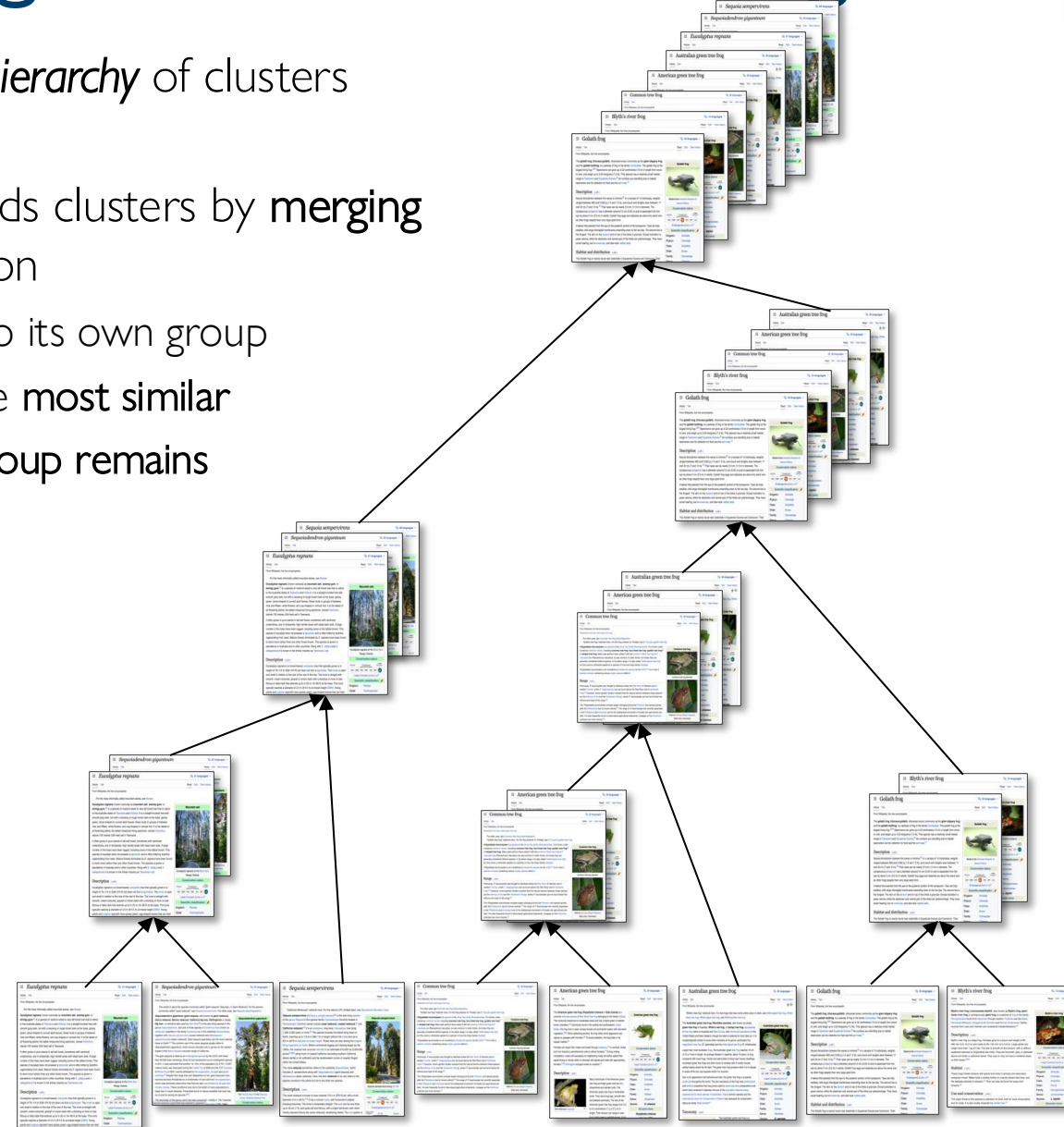
👎 Disadvantage:
  - algorithm has much **higher computational complexity** with respect to k-means
  - needs to calculate distances for pairs of datapoints, which is an $O(n^2)$ operation

# Clustering Algorithms:
# Agglomerative Hierarchical Clustering

# Hierarchical Agglomerative Clustering

Hierarchical clustering builds *hierarchy* of clusters called a *dendrogram*
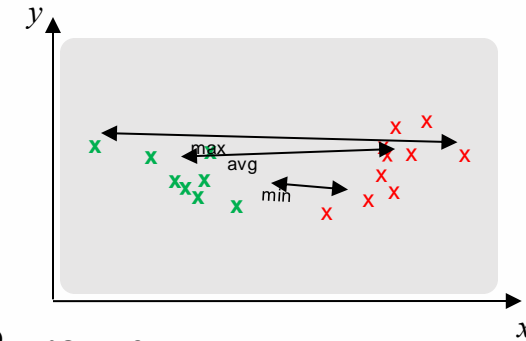
- *agglomerative clustering* builds clusters by **merging groups** in bottom-up fashion

    1. assign each document to its own group
    2. **merge** 2 groups that are **most similar**
    3. repeat **until** only **one group remains**

# Agglomerative clustering: linking

Hierarchical clustering needs to compute **distances between groups** of documents

- uses a *linkage criteria*:
  - complete-linkage: **maximum** distance between points in 2 groups
  - single-linkage: **minimum** distance across groups
  - average-linkage: **average** distance across groups
- note: linking criteria used influences shape of clusters found
  - complete- or average-linkage will restrict to tight, globular clusters
  - single-linkage clustering will allow for finding long, thin clusters

👍 Advantage:
  - works with *any distance metric / similarity function*
  - *dendrogram* provides information about underlying structure of data

👎 Disadvantage:
  - high time complexity makes it unsuitable for large datasets

# Clustering Algorithms: DBScan

# DBScan: what is it?

*Density-Based Spatial Clustering of Applications with Noise*
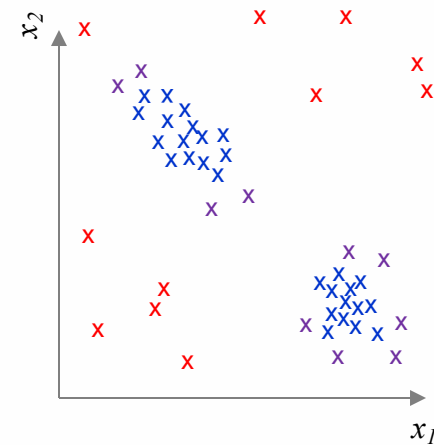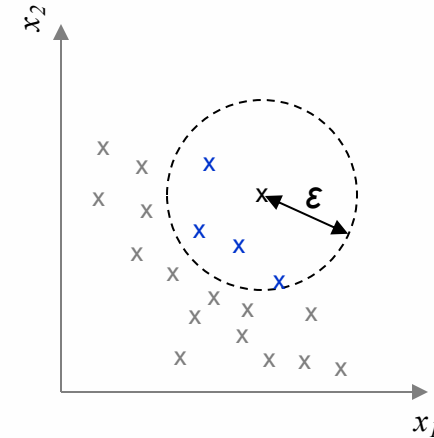
- **#clusters** not defined by user but inferred from data
- density-based algorithm: clusters found in high-density regions
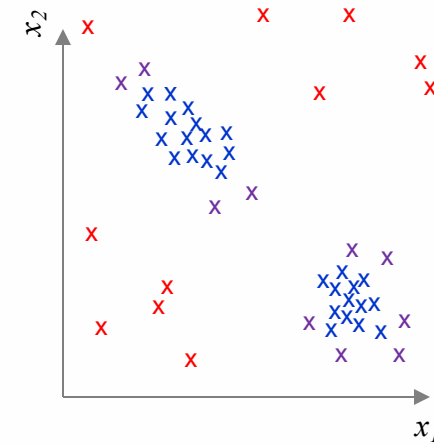
Algorithm takes 2 parameters:

- $\varepsilon$ = radius of **neighbourhood** around each point
- *minPoints* = minimum # of points required to form cluster

How does it work?

- # of datapoints in $\varepsilon$-*neighbourhood* is **density estimate**
- and *minPoints* is just a **threshold on density**
- algorithm classifies each point as either:
    - *core point*: has ≥*minPoints* neighbours inside radius $\varepsilon$
    - *border point*: not core, but within $\varepsilon$ of core point
    - *noise point*: neither core nor border point

POLITECNICO DI MILANO

# DBScan: pros and cons



👍 Advantages:

- no need to choose number of clusters in advance
- identifies and ignores noise/outlier instances
- find *arbitrarily (oddly) shaped clusters*

👎 Disadvantages:

- performance depends critically on chosen parameters
  (usually set $\varepsilon$ as a function of average distance to *MinPoints* nearest neighbour)
- may not work well if clusters have *different densities*

# Evaluating Clustering

# Evaluating clustering

Can be difficult to assess clustering results
- usually don't have access to ground-truth labelling
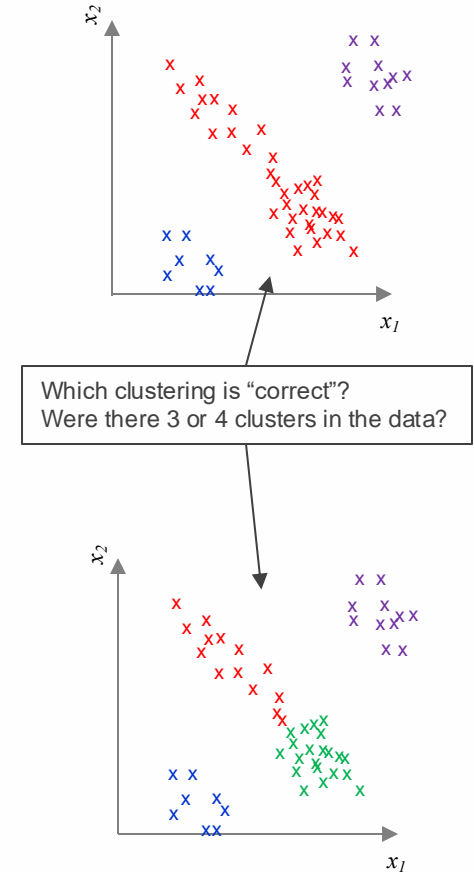- without it, evaluation can be **subjective** (people disagree on what "right" clusters are)

Two ways to evaluate results of clustering algorithm:

## 1. intrinsic evaluation:
- check **how coherent** and **well separated** clusters are
- typical evaluation measures:
    - **Within-cluster Sum of Squares** (cluster variance)
    - **Silhouette Coefficient**: compares average distance to point in current cluster with points in next closest cluster

## 2. extrinsic evaluation:
- check how well clusters **align with known document labels**
- typical evaluation measures:
    - **Homogeneity, Completeness and V-Measure**: effectively Precision, Recall and F-measure applied to clustering problems

Which clustering is "correct"?
Were there 3 or 4 clusters in the data?
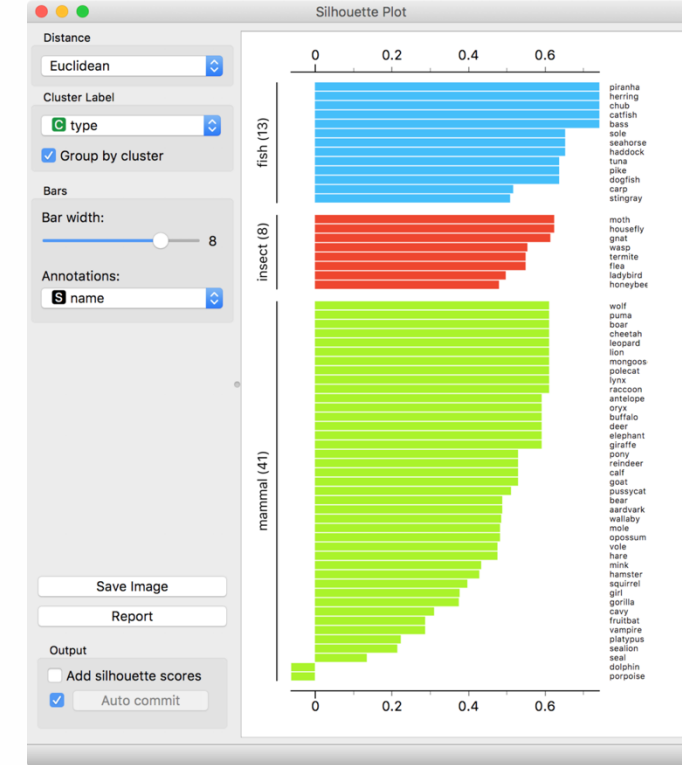
# Evaluation: silhouette

*Silhouette coefficient*

- score calculated **separately** for each sample
- uses distance metric (e.g., Euclidean)
- find average distance from sample $x_i$ to all samples:
  - in *same cluster* (denote $a_i$)
  - in *next closest cluster* (denote $b_i$)
- compute coefficient as:

$$s_i = \frac{b_i - a_i}{\max(a_i, b_i)}$$

  - ranges from −1 to +1 (really 0 to 1), higher = better
- average silhouette score over all datapoints to get single value representing quality of clustering

- often use silhouette plots to **visualise quality of clusters**
- issue: expensive to calculate, since done independently for each datapoint



Silhouette-plot-orange:
https://commons.wikimedia.org/wiki/File:Silhouette-plot-orange.png

POLITECNICO DI MILANO

# Evaluation: sum of squares

Very common to compute sum of squares values:

- *Within-cluster Sum of Squares*
  - sum over clusters the squared distances between cluster centroid and each point
- *Between-cluster Sum of Squares*
  - sum over clusters the squared distances between data centroid and each cluster centroid (weighted by data)
- **Variance Ratio Criterion** (aka Calinski-Harabasz index)
  - combines both measures into single metric
  - measures whether clusters are dense and well separated
  - faster to compute than Silhouette index, but really just objective function of *k*-Means

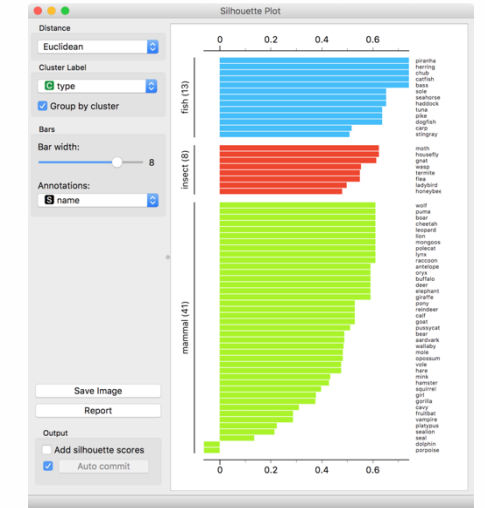$$\text{WSS}(C) = \sum_{i=1}^{k} \sum_{x_j \in C_i} d(x_j, \mu_i)^2$$

$$\text{BSS}(C) = \sum_{i=1}^{k} |C_i| d(\mu, \mu_i)^2$$

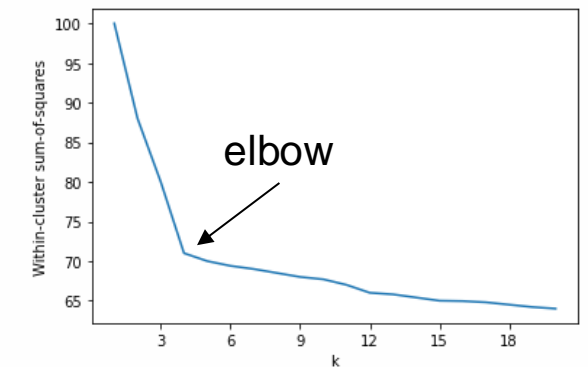$$s = \frac{BSS(C)}{WSS(C)} \times \frac{N-k}{k-1}$$

# Evaluating clustering – cont.



Silhouette-plot-orange: https://commons.wikimedia.org/wiki/File:Silhouette-plot-orange.png

Intrinsic evaluation measures can be useful to:

- **visualise** cluster quality in *Silhouette plots*

- **select a value for** *k* using the *Elbow method*
  - plot within cluster sum of squares versus number of clusters
  - as k gets larger, intrinsic measure always improves
  - so don't look for best value, but rather biggest jump:
    - point where performance increases substantially, but then slows is called **elbow**

# Topic Modelling

# What are topic models?

**Soft clustering** of documents
- allows documents to belong to multiple clusters

Provides **low-dimensional representation** of documents
- compared to high dimensional vocabulary space

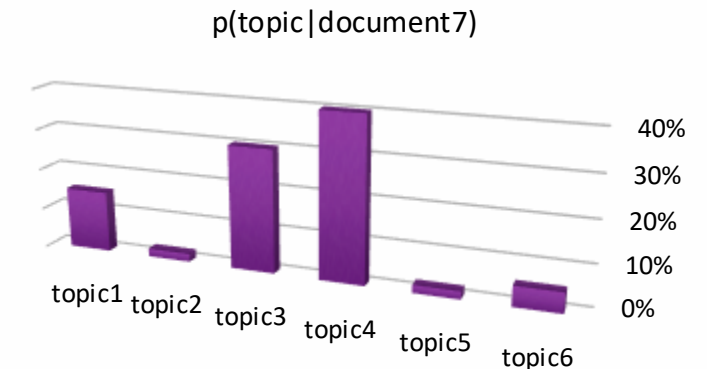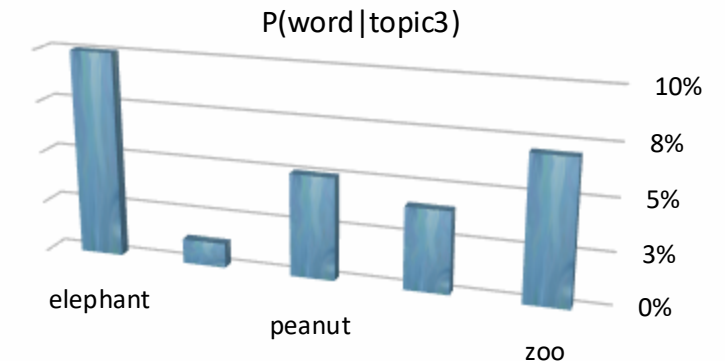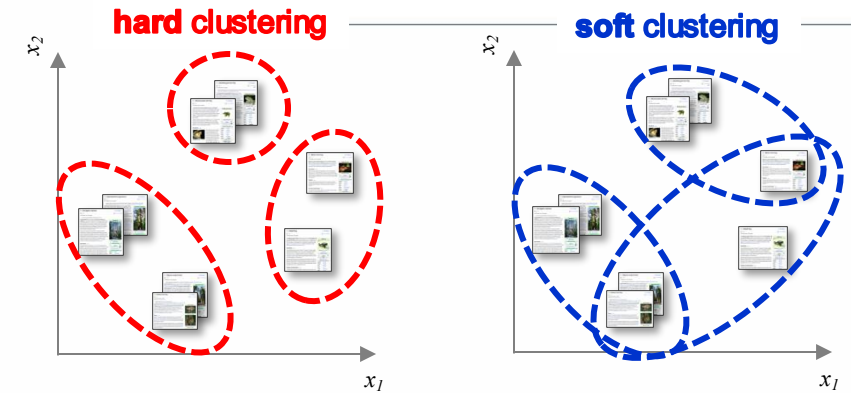Each cluster referred to as a "**topic**"
- described by distribution over words
- terms should be coherent
  (discuss single theme)
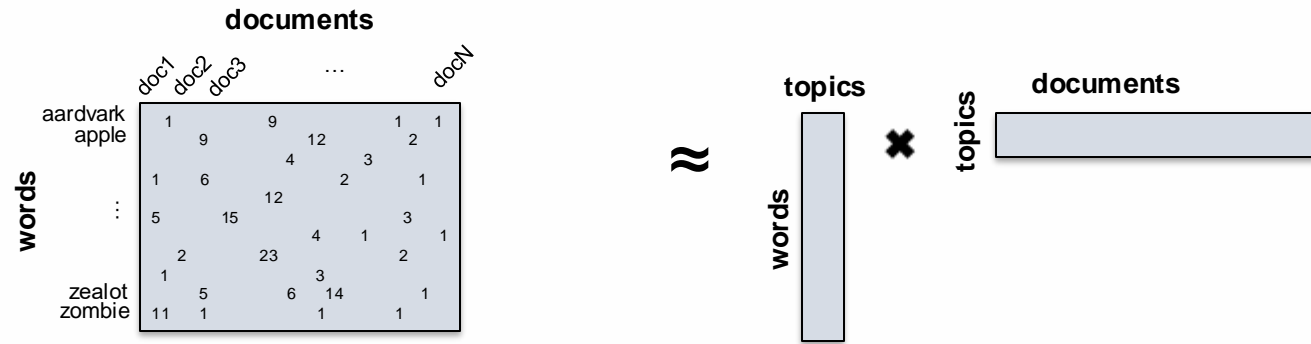
Topic vector = probability distribution over words

$$P(\mathbf{elephant}|\mathbf{topic3}) = 0.09$$

Document vector = probability distribution over topics

$$P(\mathbf{topic3}|\mathbf{doc7}) = 0.1$$

# Matrix Decomposition



Topic Modeling is a form of matrix decomposition
- decomposes *terms\*documents* count matrix
- into 2 smaller matrices: *terms\*topics* and *topics\*documents*

Much smaller number of parameters need to be estimated: $V*T+T*D <<< V*D$
- for vocabulary size $V$, document count $D$ and topic count $T$

Most famous technique is **Latent Dirichlet Allocation** (LDA)
- so named because it uses a Dirichlet prior when estimating the parameters
- related to Non-negative Matrix Factorization (NMF)
- also Latent Semantic Indexing (LSI), which applies Singular Value Decomposition (SVD) to TFIDF matrix

# What is topic modelling useful for?

Factorization/clustering **generalizes** observed term counts:
- making document representations more dense & useful
- allows calculating **more meaningful distances between documents**

Deals with problems of:
- **polysemy**: bank (financial institution or edge of river?)
- **synonymy**: cancer = oncology ?
- **short documents**: limited vocabulary

Sometimes useful for **visualizing** collections
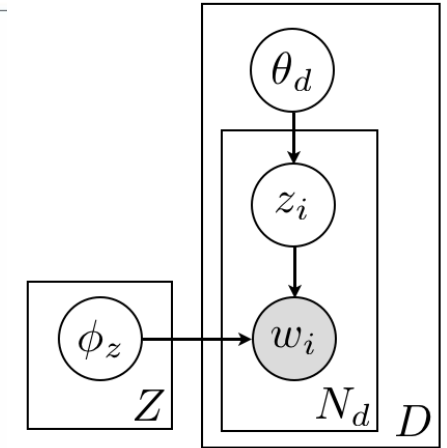- after further dimensionality reduction (using say t-SNE)



Image source: https://www.flickr.com/photos/wwarby/4915969081

# Details of LDA generative model



$$\phi_z \sim Dirichlet(\beta)$$
$$\theta_d \sim Dirichlet(\alpha)$$
$$z_i \sim Discrete(\theta_d)$$
$$w_i \sim Discrete(\phi_{z_i})$$

Generative model for latent Dirichlet Allocation (LDA)
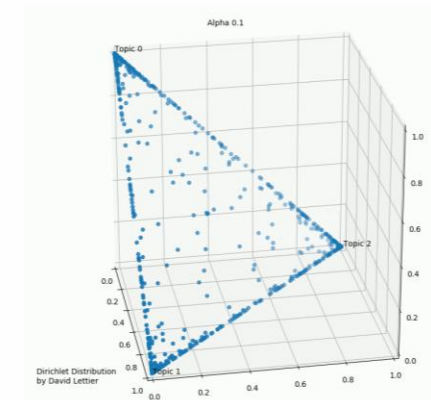
1. choose **word proportions** for each topic
2. choose **topic proportions** for each document
3. for each position in document
    1. choose **topic** based on document proportions
    2. choose **word** based on topic proportions

Estimating parameters for topic model:

- requires iteratively updating parameters
    1. topic probabilities for each document
    2. word probabilities for each topic
- place Bayesian priors on parameters to avoid over-fitting
- use (Gibbs) sampling techniques to avoid local maxima
- hyperparameters **α** and **β** determine the prior on the topic/document distributions (how concentrated they are)
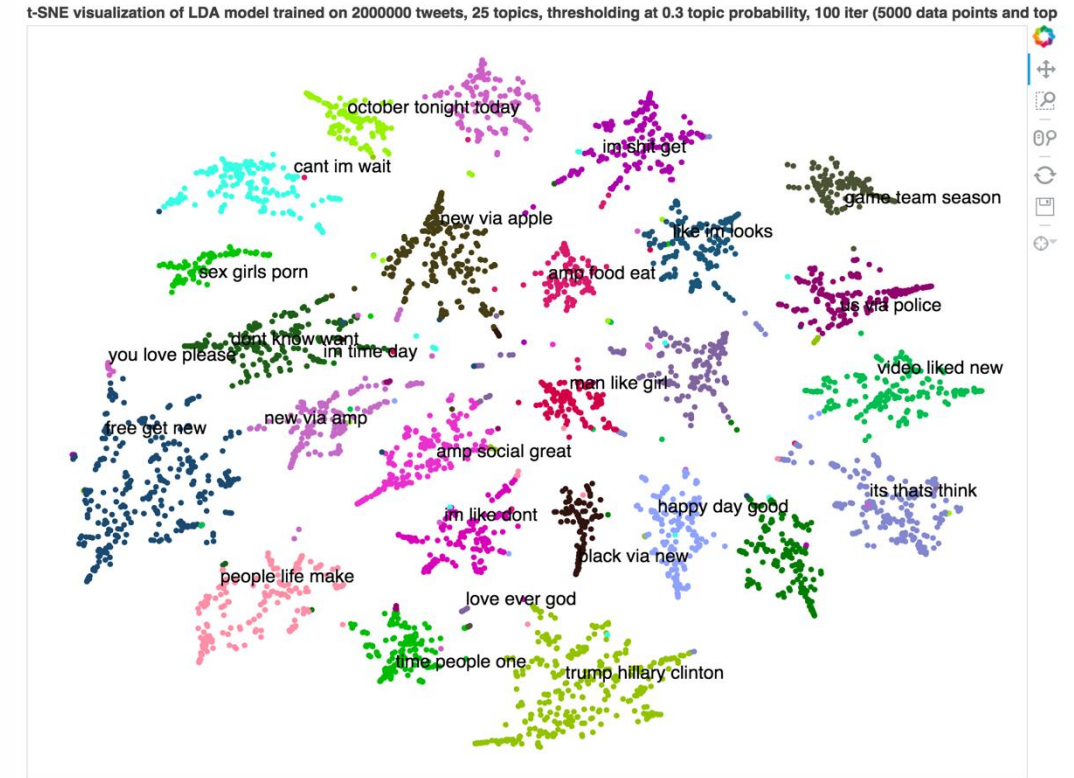


Source: David Lettier
https://medium.com/@lettier/how-does-lda-work-ill-explain-using-emoji-108abf40fa7d

# Applications of Topic Models

Useful for visualising collections
- represent topics by most frequent terms
- common also to show changes to collection over time



t-SNE visualization of LDA model trained on 2000000 tweets, 25 topics, thresholding at 0.3 topic probability, 100 iter (5000 data points and top

Source: https://shuaiw.github.io/2016/12/22/topic-modeling-and-tsne-visualization.html

POLITECNICO DI MILANO

# Visualising document vectors: dimensionality reduction techniques

Can project high-dimensional data into lower dimensions using:

1. a **linear projection**
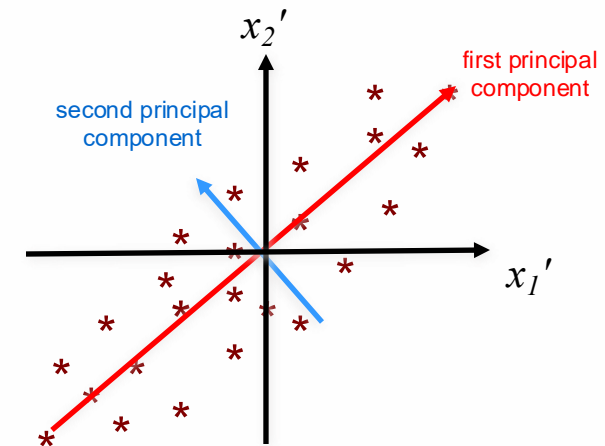
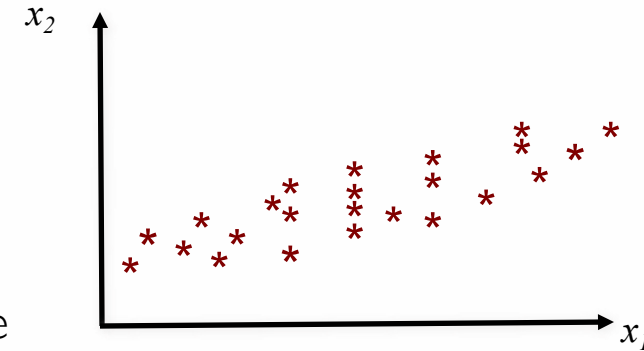   with Principle Component Analysis

2. a **non-linear projection**

   e.g. using t-distributed Stochastic Neighbor Embeddings
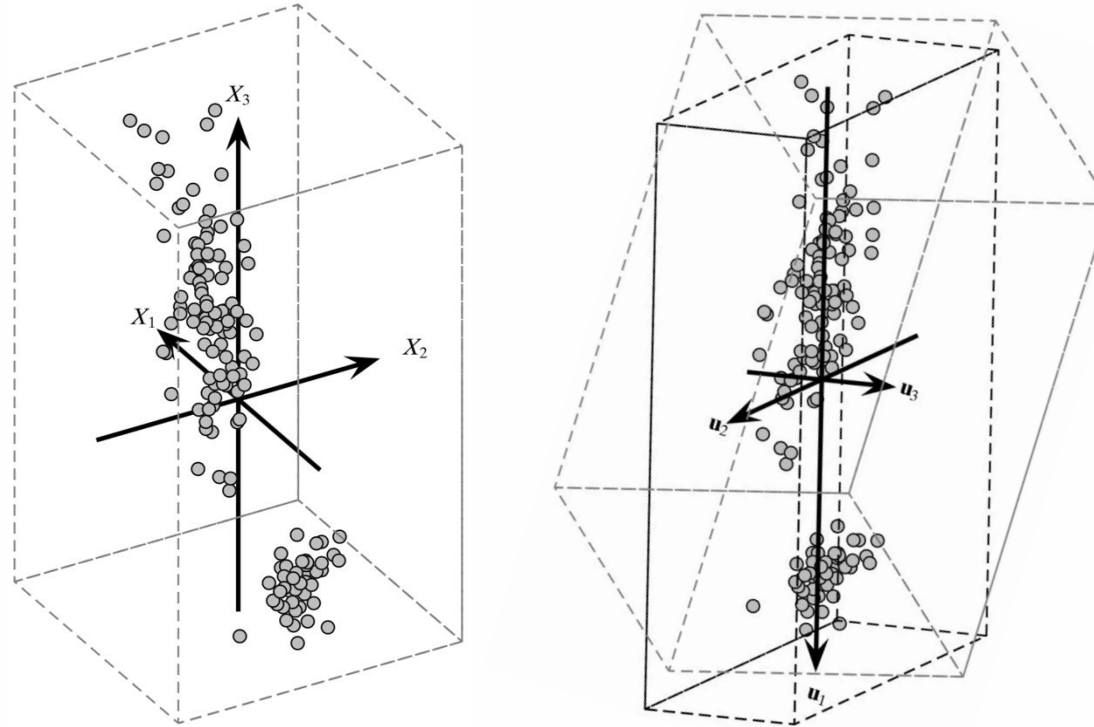
# Principal Component Analysis (PCA)

Investigates main sources of variation in dataset

- by finding **projection** that captures **largest amount of variation** in data

- method:
    1. scale features to have zero mean, unit variance
    2. find **direction** in **d**-dimensions that represent largest variation in data ($k<d$)
    3. remove variation from data
    4. repeat until $k$ dimensions found

- use new dimensions to approximate original data using only largest components

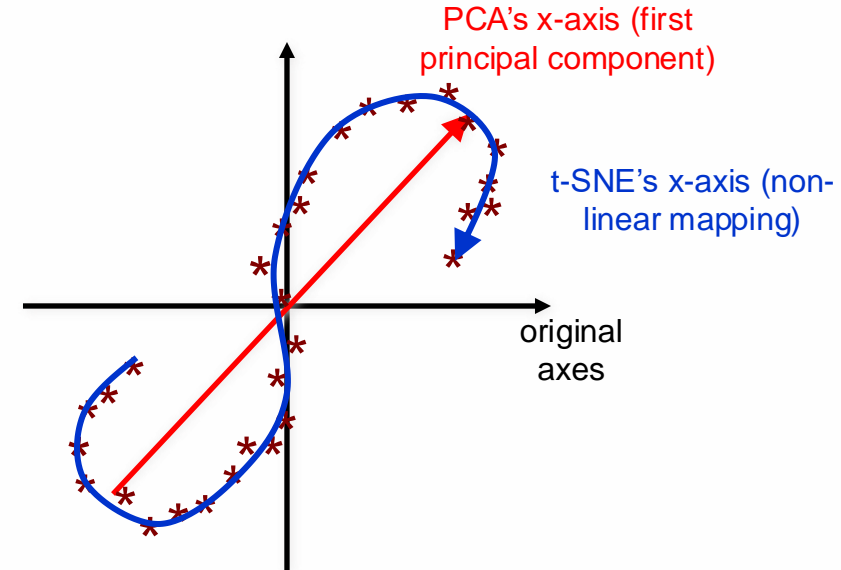# PCA: example in 3d

For more information on PCA

- see Chapter 7 of the textbook "Data Mining and Analysis"
- and: http://sebastianraschka.com/Articles/2015_pca_in_3_steps.html

POLITECNICO DI MILANO

# t-distributed Stochastic Neighbor Embedding

# Non-linear Dimensionality Reduction

PCA's x-axis (first principal component)

t-SNE's x-axis (non-linear mapping)

original axes

Data in high dimensions never fills the entire space

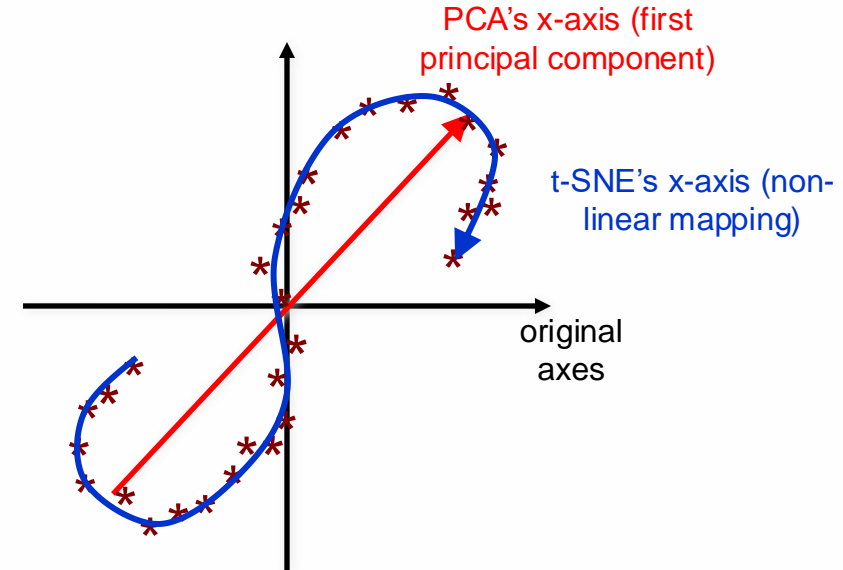- always lives on some lower-dimensional manifold

t-SNE is a non-linear dimensionality reduction technique

- used to map high-dimensional data into 2 or 3 dimensions
  - points from original space mapped onto "map points" in 2D/3D
- unlike PCA:
  - mapped point is **not linear combination** of **original attribute values**
  - axis of mapped space is not linear combination (rotation) of original axes

# t-distributed Stochastic Neighbor Embedding

PCA's x-axis (first principal component)

t-SNE's x-axis (non-linear mapping)

original axes

t-SNE:

- tries to preserve **local distances** to **nearby points**
  - unlike PCA which tries to preserve global (long range) distances between points
- converts distances between data points to joint probabilities
  - then maps original points to lower dimensional points such that position of mapped points **preserves structure of the data**
    - i.e. similar data points are assigned nearby map points
    - while dissimilar data points are assigned distant map points

# t-SNE: algorithm

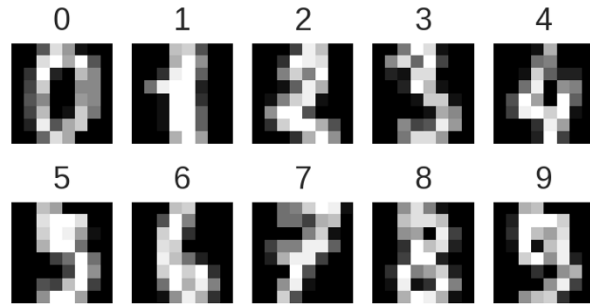The t-SNE algorithm has two main steps:

1. define **probability distribution over pairs** of high-dimensional points:
   - such that pairs of **similar data points have high probability** of being picked
   - and dissimilar points have extremely low probability of being picked
2. define a similar distribution over the points in the map space
   - and move around the mapped points to **minimize Kullback–Leibler divergence** between the two distributions
   - gradient descent is used to minimize the score

Caveats of t-SNE:
   - **perplexity** parameter controls balance between local and global aspects of the data
   - different initializations will lead to different results
   - applies only to data with "reasonable" number of dimensions (e.g. 30-50)
     - if data has more dimensions, computational time will be prohibitive

POLITECNICO DI MILANO

# t-SNE example: handwritten digit dataset

- Handwritten Digits Data Set from the UCI machine learning repository
  http://archive.ics.uci.edu/ml/datasets/Optical+Recognition+of+Handwritten+Digits
- Contains 1797 images with 8x8 pixels each



9

https://github.com/oreillymedia/t-SNE-tutorial

# Conclusions

# Conclusions

- Many applications of text clustering
- Many clustering techniques can be applied to tf-idf vectors
- Evaluating clustering can be subjective
- Topic modeling builds lower-dimensional representation of documents by performing soft clustering
- Dimensionality reduction techniques can be used to visualize document vectors