

# Unsupervised Learning

---

Pierre Gaillard – ENS Paris

September 28, 2018

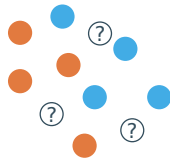
# Supervised vs unsupervised learning

Two main categories of machine learning algorithms:

- **Supervised learning:** predict output  $Y$  from some input data  $X$ . The training data has a known label  $Y$ .

Examples:

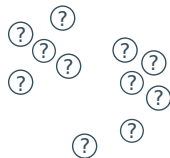
- $X$  is a picture, and  $Y$  is a cat or a dog
- $X$  is a picture, and  $Y \in \{0, \dots, 9\}$  is a digit
- $X$  is are videos captured by a robot playing table tennis, and  $Y$  are the parameters of the robots to return the ball correctly



- **Unsupervised learning:** training data is not labeled and does not have a known result

Examples:

- detect change points in a non-stationary time-series
- detect outliers
- **clustering:** group data in homogeneous groups
- **principal component analysis:** compress data without losing much information
- density estimation
- dictionary learning



- **Others:** reinforcement learning, semi-supervised learning, online learning, ...

**Reminder of last class on supervised  
learning: OLS and logistic regression**

---

# Least square Linear regression

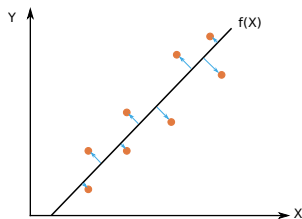
Given training data  $(X_i; Y_i)$  for  $i = 1; \dots; n$ , with  $X_i \in \mathbb{R}^d$  and  $Y_i \in \mathbb{R}$  learn a predictor  $f$  such that our expected square loss

$$E (f(X) - Y)^2$$

is small.

We assume here that  $f$  is a **linear combination** of the input  $x = (x_1; \dots; x_d)$

$$f(x) = \sum_{i=1}^d w_i x_i = w^T x$$



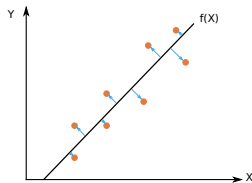
# Ordinary Least Square

Input  $X \in \mathbb{R}^d$ , output  $Y \in \mathbb{R}$ , and  $\ell$  is the square loss:  $\ell(a; y) = (a - y)^2$ .

The Ordinary Least Square regression (OLS) minimizes the **empirical risk**

$$R_n(w) = \frac{1}{n} \sum_{i=1}^n (Y_i - w^T X_i)^2$$

This is minimized in  $w \in \mathbb{R}^d$  when  $X^T X w - X^T Y = 0$ ; where  $X = [X_1; \dots; X_n]^T \in \mathbb{R}^{n \times d}$  and  $Y = [Y_1; \dots; Y_n]^T \in \mathbb{R}^n$ .



Assuming  $X$  is injective (i.e.,  $X^T X$  is invertible) and there is an exact solution

$$w = (X^T X)^{-1} X^T Y;$$



What happens if  $d > n$ ?

# Ordinary Least Square: how to compute $\hat{\theta}_n$ ?

If the design matrix  $X^T X$  is invertible, the OLS has the closed form:

$$\hat{\beta}_n \triangleq \arg \min_{\beta} \mathbf{r}_n(\beta) = (X^T X)^{-1} X^T Y;$$

**Question:** how to compute it?

- **inversion of  $(X^T X)$**  can be prohibitive (the cost is  $O(d^3)$ !)
- **QR-decomposition:** we write  $X = QR$ , with  $Q$  an orthogonal matrix and  $R$  an upper-triangular matrix. One needs to solve the linear system:

$$R\hat{\beta} = Q^T Y; \quad \text{with} \quad R = \begin{matrix} & 0 & & & & & & & 1 \\ & x & & x & \dots & \dots & \dots & & x \\ \textcircled{Q} & & \cdot & & & & & & \\ & & & \cdot & & & & & \\ & & & & \cdot & & & & \\ & & & & & \cdot & & & \\ & & & & & & \cdot & & \\ & & & & & & & & \cdot \\ & & & & & & & & x \\ & & & & & & & & \end{matrix}$$

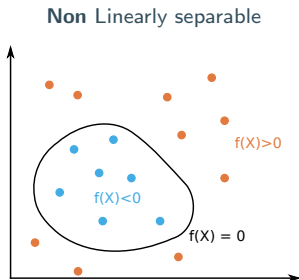
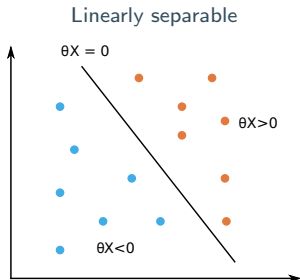
- **iterative approximation with convex optimization algorithms** [Bottou, Curtis, and Nocedal 2016]: (stochastic)-gradient descent, Newton,...

$$\hat{\beta}_{i+1} = \hat{\beta}_i - \gamma \mathbf{r}_n(\hat{\beta}_i)$$

# Classification

Given training data  $(X_i; Y_i)$  for  $i = 1; \dots; n$ , with  $X_i \in \mathbb{R}^d$  and  $Y_i \in \{-1, +1\}$  learn a classifier  $f(x)$  such that

$$f(X_i) \begin{cases} > 0 \\ < 0 \end{cases} \quad \begin{matrix} ) \\ ) \end{matrix} \quad \begin{matrix} Y_i = +1 \\ Y_i = 0 \end{matrix}$$



# Linear classification

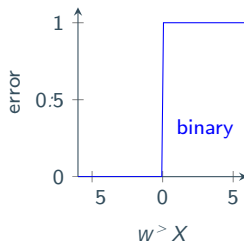
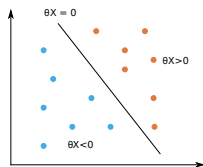
We would like to find the best linear classifier such that

$$f(X) = \begin{cases} > X^T w > 0 & Y = +1 \\ < X^T w < 0 & Y = 0 \end{cases}$$

**Empirical risk minimization** with the binary loss?

$$b_n = \arg \min_{w \in \mathbb{R}^d} \frac{1}{n} \sum_{i=1}^n \mathbb{1}_{Y_i \neq \mathbb{1}_{w^T x_i \geq 0}}$$

👎 This is **not convex** in  $w$ . Very hard to compute!

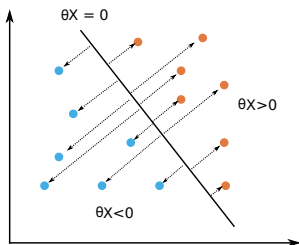
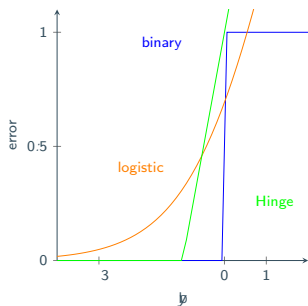




# Logistic regression

**Idea:** replace the loss with a convex loss

$$\ell(\theta; X; y) = y \log \frac{1}{1 + e^{-\theta^T X}} + (1 - y) \log \frac{1}{1 + e^{\theta^T X}}$$



**Probabilistic interpretation:** based on likelihood maximization of the model:

$$P(Y = 1|X) = \frac{1}{1 + e^{-\theta^T X}} \in [0; 1]$$

Satisfied for many distributions of  $X|Y$ : Bernoulli, Gaussian, Exponential, ...

**Computation of the minimizer of the empirical risk** (No closed form of the solution)

- Use a convex optimization algorithm (Newton, gradient descent, ...)

# Clustering

---

Reminder of last class on supervised learning: OLS and logistic regression

## Clustering

- What is clustering?

- Clustering algorithms

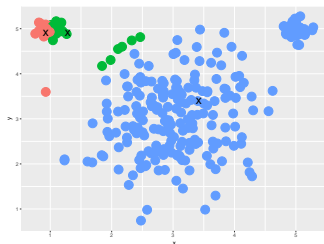
Dimensionality Reduction Algorithms

# What is clustering?

**Goal:** group together similar instances

Requires data but no labels

Useful when you don't know what you are looking for



## Types of clustering algorithms:

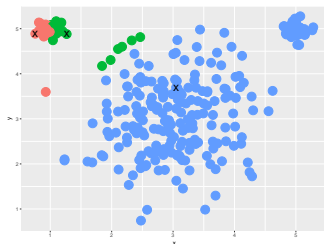
- model based clustering: **mixture of Gaussian**
- **hierarchical clustering**: a hierarchy of nested clusters is build using divisive or agglomerative approach
- Flat clustering: no hierarchy: **k-means**, spectral clustering

# What is clustering?

**Goal:** group together similar instances

Requires data but no labels

Useful when you don't know what you are looking for



## Types of clustering algorithms:

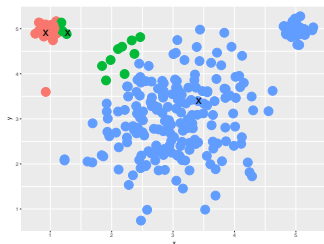
- model based clustering: **mixture of Gaussian**
- **hierarchical clustering**: a hierarchy of nested clusters is build using divisive or agglomerative approach
- Flat clustering: no hierarchy: **k-means**, spectral clustering

# What is clustering?

**Goal:** group together similar instances

Requires data but no labels

Useful when you don't know what you are looking for



**Types of clustering algorithms:**

- model based clustering: **mixture of Gaussian**
- **hierarchical clustering**: a hierarchy of nested clusters is build using divisive or agglomerative approach
- Flat clustering: no hierarchy: **k-means**, spectral clustering

## Examples of applications

Clustering has a variety of goals which all relates to grouping or segmenting a collection of objects into homogeneous objects.

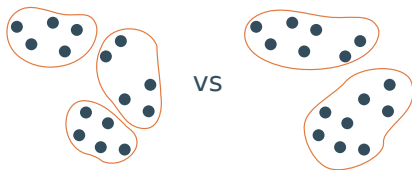
- **Biology:** group species and generate phylogenies
- **Marketing:** group similar shopping items and similar customers / market outreach, recommender systems
- **Computer science:** image segmentation



# What do we need for clustering?

Clustering does not need any labeled data. But:

1. **A proximity measure:** it can be
  - **similarity measure**  $s(X_i; X_j)$  which is large if  $X_i$  and  $X_j$  are similar. Example: correlation
  - **dissimilarity measure**  $d(X_i; X_j)$  which is large if  $X_i$  and  $X_j$  are different. Example: distance from a metric  $d(X_i; X_j) = \|X_i - X_j\|$
2. **A criterion to evaluate the clustering:**



3. **A clustering algorithm** optimizing the criterion



**Proximity matrix:** the data may be represented directly in terms of the proximity between pairs of objects.

This is the case for categorical variables.

**Example:** experiment where participants are asked to judge how much objects differ from one another.

The proximity matrix is assumed to be symmetric:  $D = (D + D^T)/2$

**Distance measure:** the similarity can be measured by a **metric**

- Euclidean distance:  $d(X_i; X_j) = \|X_i - X_j\|_2$
- Correlation
- Manhattan distance:  $d(X_i; X_j) = \|X_i - X_j\|_1$
- Minkowski distance:  $d(X_i; X_j) = \|X_i - X_j\|_p$

The results crucially depends on the metric choice: depends on data.

This is often a hard problem. Much harder than for supervised learning!

A trade-off **Intra-cluster homogeneity**: objects within the same cluster should be similar from one another

**Inter-clusters heterogeneity**: clusters should be well-separated, their centers should be far from one another

In many application, we ask experts

## Examples of performance criterion

**Known ground truth: mutual information:** if  $U; V : f_1; \dots; n_g ! P(f_1; \dots; n_g)$  are two clustering of  $n$  objects, the mutual information measures their agreement as:

$$MI(U; V) = \sum_{i=1}^n \sum_{j=1}^n P(i; j) \log \frac{P(i; j)}{P(i)P^0(j)}$$

where  $P(i) = jU_{ij=j=n}$  (resp.  $P^0(j) = jV_{jj=n}$  and  $P(i; j) = jV_{jj} \setminus jU_{ij=j=n}$ ) is the probability that an object picked at random falls into class  $U_i$  (resp.  $V_j$  and into both classes). **Drawback:** the ground truth is rarely known

**Unknown ground truth: Silhouette score** defined as

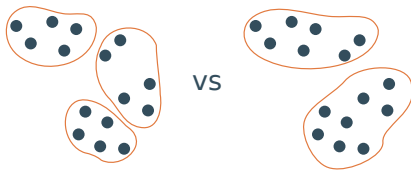
$$Score = \frac{b - a}{\max(\hat{a}; bg)}$$

where

- $a$ : the mean distance between a sample and all other points in the same class.
- $b$ : the mean distance between a sample and all other points in the next nearest cluster.

The score is higher when clusters are dense and well separated, which relates to a standard concept of a cluster.

## How many clusters?



Possible approaches:

- Fix the number of clusters  $k$
- Choose the number of clusters by optimizing the performance criterion such as the silhouette score

Reminder of last class on supervised learning: OLS and logistic regression

## Clustering

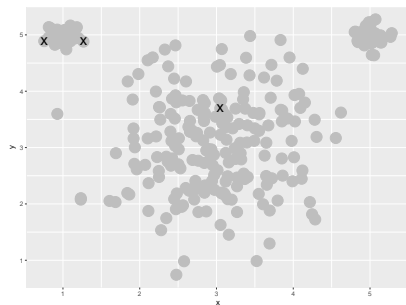
- What is clustering?

- Clustering algorithms

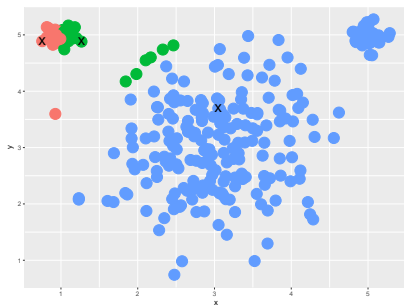
Dimensionality Reduction Algorithms

## Types of clustering algorithms:

- **Kmeans**: flat clustering
- **Spectral clustering**
- **Combinatorial clustering**
- **Gaussian mixture**: model based clustering
- **Hierarchical clustering**: a hierarchy of nested clusters is build using divisive or agglomerative approach
- **Affinity propagation**: based on message passing in a graph

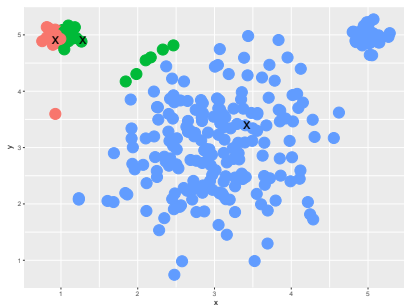


- **Initialization:** sample  $K$  points as cluster centers
- **Alternate:**
  1. Assign points to closest center
  2. Update cluster to the averaged of its assigned points
- **Stop** when no point's assignment change.

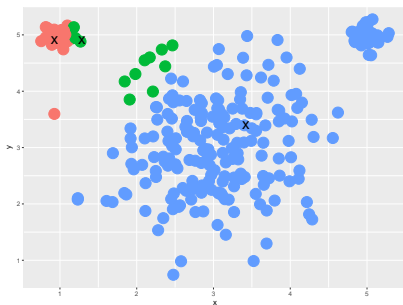


- **Initialization:** sample  $K$  points as cluster centers
- **Alternate:**
  1. Assign points to closest center
  2. Update cluster to the averaged of its assigned points
- **Stop** when no point's assignment change.

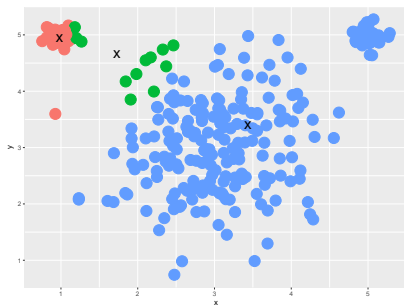




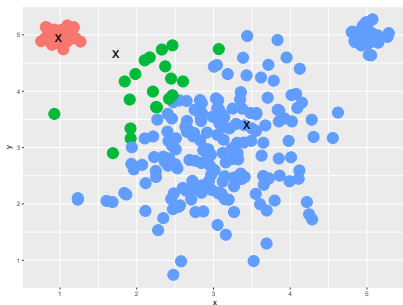
- **Initialization:** sample  $K$  points as cluster centers
- **Alternate:**
  1. Assign points to closest center
  2. Update cluster to the averaged of its assigned points
- **Stop** when no point's assignment change.



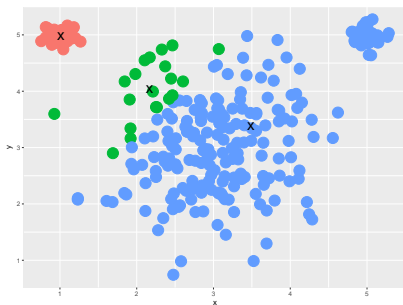
- **Initialization:** sample  $K$  points as cluster centers
- **Alternate:**
  1. Assign points to closest center
  2. Update cluster to the averaged of its assigned points
- **Stop** when no point's assignment change.



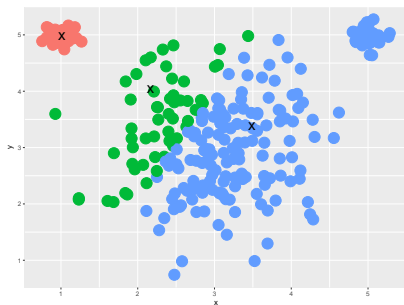
- **Initialization:** sample  $K$  points as cluster centers
- **Alternate:**
  1. Assign points to closest center
  2. Update cluster to the averaged of its assigned points
- **Stop** when no point's assignment change.



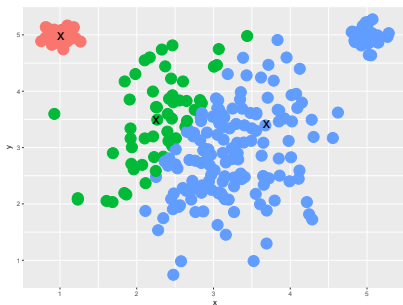
- **Initialization:** sample  $K$  points as cluster centers
- **Alternate:**
  1. Assign points to closest center
  2. Update cluster to the averaged of its assigned points
- **Stop** when no point's assignment change.



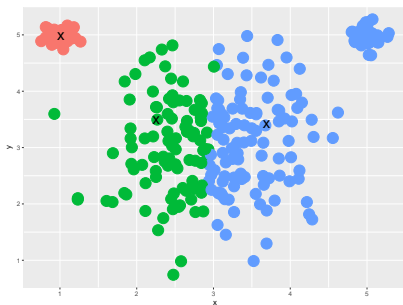
- **Initialization:** sample  $K$  points as cluster centers
- **Alternate:**
  1. Assign points to closest center
  2. Update cluster to the averaged of its assigned points
- **Stop** when no point's assignment change.



- **Initialization:** sample  $K$  points as cluster centers
- **Alternate:**
  1. Assign points to closest center
  2. Update cluster to the averaged of its assigned points
- **Stop** when no point's assignment change.

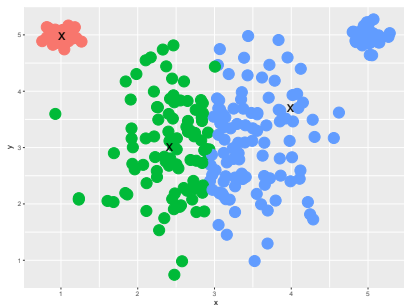


- **Initialization:** sample  $K$  points as cluster centers
- **Alternate:**
  1. Assign points to closest center
  2. Update cluster to the averaged of its assigned points
- **Stop** when no point's assignment change.

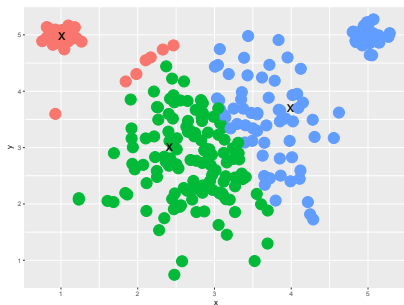


- **Initialization:** sample  $K$  points as cluster centers
- **Alternate:**
  1. Assign points to closest center
  2. Update cluster to the averaged of its assigned points
- **Stop** when no point's assignment change.

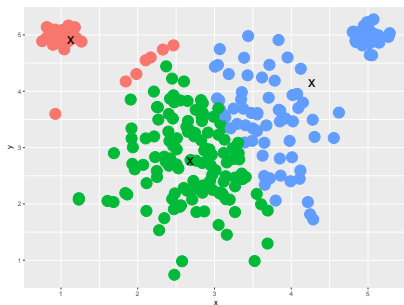




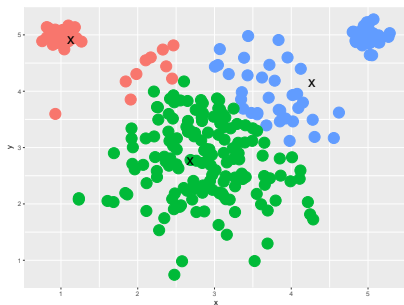
- **Initialization:** sample  $K$  points as cluster centers
- **Alternate:**
  1. Assign points to closest center
  2. Update cluster to the averaged of its assigned points
- **Stop** when no point's assignment change.



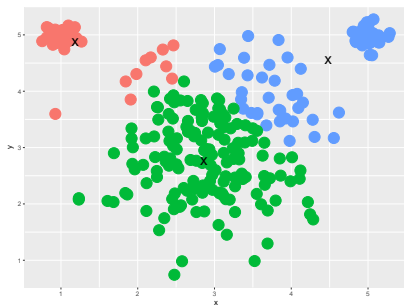
- **Initialization:** sample  $K$  points as cluster centers
- **Alternate:**
  1. Assign points to closest center
  2. Update cluster to the averaged of its assigned points
- **Stop** when no point's assignment change.



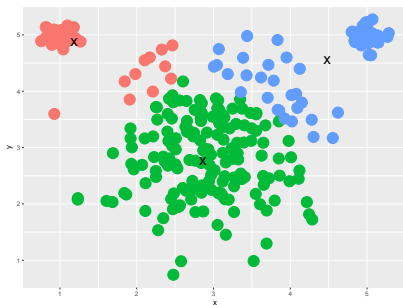
- **Initialization:** sample  $K$  points as cluster centers
- **Alternate:**
  1. Assign points to closest center
  2. Update cluster to the averaged of its assigned points
- **Stop** when no point's assignment change.



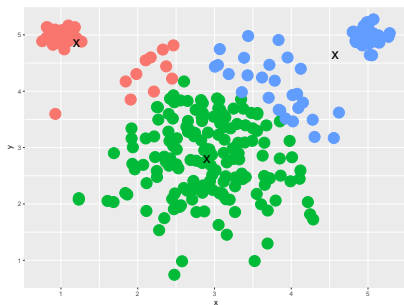
- **Initialization:** sample  $K$  points as cluster centers
- **Alternate:**
  1. Assign points to closest center
  2. Update cluster to the averaged of its assigned points
- **Stop** when no point's assignment change.



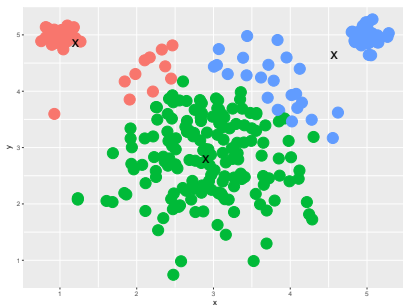
- **Initialization:** sample  $K$  points as cluster centers
- **Alternate:**
  1. Assign points to closest center
  2. Update cluster to the averaged of its assigned points
- **Stop** when no point's assignment change.



- **Initialization:** sample  $K$  points as cluster centers
- **Alternate:**
  1. Assign points to closest center
  2. Update cluster to the averaged of its assigned points
- **Stop** when no point's assignment change.

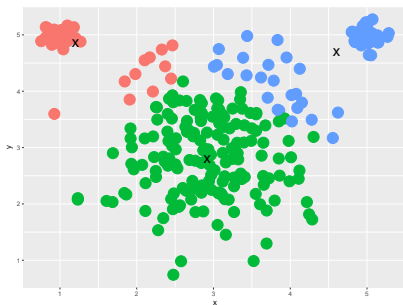


- **Initialization:** sample  $K$  points as cluster centers
- **Alternate:**
  1. Assign points to closest center
  2. Update cluster to the averaged of its assigned points
- **Stop** when no point's assignment change.

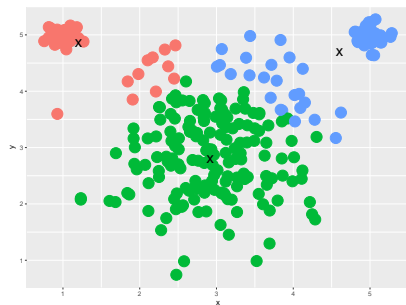


- **Initialization:** sample  $K$  points as cluster centers
- **Alternate:**
  1. Assign points to closest center
  2. Update cluster to the averaged of its assigned points
- **Stop** when no point's assignment change.





- **Initialization:** sample  $K$  points as cluster centers
- **Alternate:**
  1. Assign points to closest center
  2. Update cluster to the averaged of its assigned points
- **Stop** when no point's assignment change.



- **Initialization:** sample  $K$  points as cluster centers
- **Alternate:**
  1. Assign points to closest center
  2. Update cluster to the averaged of its assigned points
- **Stop** when no point's assignment change.

Guaranteed to converge in a finite number of iterations.

Initialization is crucial.

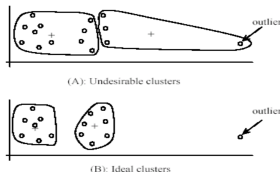
# Kmeans: Pros and cons

## Pros

- **simple** easy to implement
- **efficient**: guaranteed to converge in finite number of iteration. Complexity:  $O(tkn)$ ,  $k$ : number of clusters (small),  $n$ : number of data,  $t$  number of iteration (small)
- **popular**

## Cons

- the means need to be defined
- the number of cluster  $k$  need to be specified
- sensitivity to outliers
  - ! perform subsampling or outlier detection
- sensitivity to initial seeds: often get stuck in local minima
  - ! initialize with different random seeds
- not suited for special data structure
- fails if clusters are not round



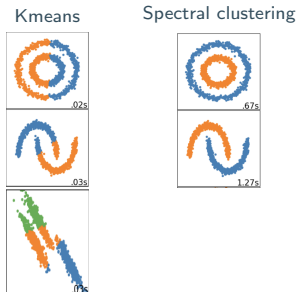
# Kmeans: Pros and cons

## Pros

- **simple** easy to implement
- **efficient**: guaranteed to converge in finite number of iteration. Complexity:  $O(tkn)$ ,  $k$ : number of clusters (small),  $n$ : number of data,  $t$  number of iteration (small)
- **popular**

## Cons

- the means need to be defined
- the number of cluster  $k$  need to be specified
- sensitivity to outliers
  - ! perform subsampling or outlier detection
- sensitivity to initial seeds: often get stuck in local minima
  - ! initialize with different random seeds
- not suited for special data structure
- fails if clusters are not round



Changing the feature  
(or distance) may help

## What is a Gaussian mixture?

- Parameter  $\theta = (w; \mu; \Sigma)$
- There are  $K$  groups.
- Group  $k$  is associated a mean  $\mu_k$  and a variance  $\Sigma_k$
- The data of each group is generated from a Gaussian  $\mathcal{N}(\mu_k; \Sigma_k)$
- A data point  $X^{(i)}$  is generated as follows:
  - randomly sample the group  $Z^{(i)} \in \{0, 1\}^K$  with  $\sum_k Z_k^{(i)} = 1$  the group  $\mathbb{P}(Z_k^{(i)} = 1) = w_k$
  - If  $Z_k^{(i)} = 1$ , sample  $X^{(i)} \sim \mathcal{N}(\mu_k; \Sigma_k)$
- The points  $X^{(i)}$  are revealed but **the groups  $Z^{(i)}$  are hidden**

The density of  $X^{(i)}$  is

$$f(x; \theta) = \sum_{k=1}^K w_k \mathcal{N}(x; \mu_k; \Sigma_k)$$

The density of  $X_i; Z_i$  is

$$f(x; z; \theta) = \sum_{k=1}^K w_k \mathcal{N}(x; \mu_k; \Sigma_k)^{z_k}$$

## Gaussian Mixture – EM algorithm

We want to maximize the log-likelihood of  $f(X; \theta) = \prod_{k=1}^K w_k \mathcal{N}(X; \mu_k; \Sigma_k)$  but the expression of the log-likelihood  $L(X; \theta)$  is complicated.

Instead, we consider the complete density of  $X^{(i)}; Z^{(i)}$

$$f(X; Z; \theta) = \prod_{k=1}^K w_k \mathcal{N}(X; \mu_k; \Sigma_k)^{z_k}$$

whose log-likelihood is

$$L(X; Z; \theta) = \log \prod_{i=1}^n f(X^{(i)}; Z^{(i)}; w; \mu; \Sigma) = \sum_{i=1}^n \sum_{k=1}^K z_k^{(i)} \log \mathcal{N}(X^{(i)}; \mu_k; \Sigma_k) + \sum_{i=1}^n \sum_{k=1}^K z_k^{(i)} \log w_k$$

### How to maximize $L$ ?

1. If we knew the  $Z^{(i)}$  we could maximize  $L$  in  $\theta = (w; \mu; \Sigma)$
2. If we knew  $\theta$  we could choose the best  $Z^{(i)}$  maximizing their probability

$$P \left( Z_k^{(i)} = 1 | X^{(i)}; \theta \right) = \frac{w_k \mathcal{N}(X^{(i)}; \mu_k; \Sigma_k)}{\sum_{j=1}^K w_j \mathcal{N}(X^{(i)}; \mu_j; \Sigma_j)}$$

The expectation-maximization algorithm alternates between those two.

## Gaussian Mixture – EM algorithm

We want to maximize the log-likelihood of  $f(x; \theta) = \sum_{k=1}^K w_k \mathcal{N}(x; \mu_k; \Sigma_k)$  but the expression of the log-likelihood  $L(X; \theta)$  is complicated.

Instead, we consider the complete density of  $X^{(i)}; Z^{(i)}$

$$f(x; z; \theta) = \sum_{k=1}^K w_k \mathcal{N}(x; \mu_k; \Sigma_k)^{z_k}$$

whose log-likelihood is

$$L(X; Z; \theta) = \log \prod_{i=1}^n f(X^{(i)}; Z^{(i)}; w; \mu; \Sigma) = \sum_{i=1}^n \sum_{k=1}^K Z_k^{(i)} \log \mathcal{N}(X^{(i)}; \mu_k; \Sigma_k) + \sum_{i=1}^n \sum_{k=1}^K Z_k^{(i)} \log w_k$$

### Expectation-Maximization algorithm

Initialize:  $\theta = \theta_0$

Until convergence do

- **Expectation step:**  $q(Z) = P f(Z|X; \theta)$
- **Maximization step:**

$$\theta_t = \arg \max_{\theta} E_Z q L(X; Z; \theta) :$$

There are closed form formulas!

Produce a **hierarchical representation** in which the clusters at each level of the hierarchy are created by

- **merging** clusters at the lower level: **agglomerative** (bottom-up) strategy  
! we group the pair of clusters with the smallest intergroup dissimilarity
- **splitting** clusters at the higher level: **divisive** (top-down) strategy  
! we split the cluster to produce two groups with the largest intergroup dissimilarity

Does not require the number  $k$  of clusters to be pre-specified.

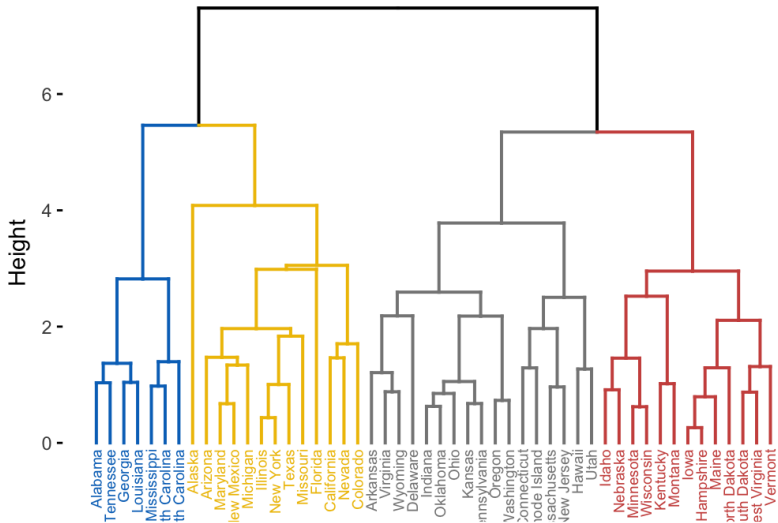
Each level of the hierarchy represents a particular clustering of the data

For some data (like biological) this is better suited than flat clustering with a hierarchy of groups: animal, mammal,...



# Hierarchical clustering dendrogram

## Cluster Dendrogram



Initialize with a single cluster

Divide each group recursively into two childs using for example kmeans with  $k = 2$ .

Stop when each cluster has only a single element.

**Algorithm:** Initialize with  $n$  cluster each containing a single data point.

While there is more than one cluster:

- find the two nearest clusters
- merge them

**How to measure the distance between clusters?** Four common ways

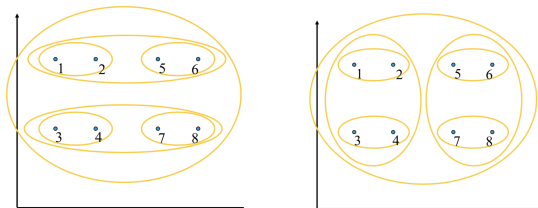
- minimum distance:  $d_{\min}(U; V) = \min_{u \in U, v \in V} \|u - v\|$
- maximum distance:  $d_{\max}(U; V) = \max_{u \in U, v \in V} \|u - v\|$
- average distance
- distance between means

Different choices create different clustering behaviors!

**Algorithm:** Initialize with  $n$  cluster each containing a single data point.

While there is more than one cluster:

- find the two nearest clusters
- merge them



[Pictures from Thorsten Joachims]

Agglomerative is faster (no kmeans computation at each iteration)

Divisive is more global: first step find the best split among the all data, while agglomerative only look at pairwise comparison.

# Dimensionality Reduction Algorithms

---

# Principal component analysis

Assume that you have **centered** observations  $x_1, \dots, x_n \in \mathbb{R}^p$  represented as a data matrix (with column-wise zero empirical mean  $\bar{x} = \frac{1}{n} \sum_{i=1}^n x_i = 0$ )

$$X := [x_1, \dots, x_n]^T \quad \text{of dimension } n \times p:$$

If  $p$  is large, some columns (i.e., explanatory variables) may be **linearly correlated**.

**bad statistical property:** risk minimization not identifiable, the covariance matrix  $(X^T X)$  is not invertible  
/ unstable estimators

**bad computational property:** we need to store  $p - 1$  columns with redundant information

Is it possible to reduce the dimension of the data from  $\mathbb{R}^p$  to  $\mathbb{R}^q$  through **linear transformations** without losing much information?

The principal components of a set of data in  $\mathbb{R}^p$  provide a **sequence of best linear approximations** of that data.

# Principal component analysis

Assume that you have **centered** observations  $x_1; \dots; x_n \in \mathbb{R}^p$  represented as a data matrix (with column-wise zero empirical mean  $\bar{x} = \frac{1}{n} \sum_{i=1}^n x_i = 0$ )

$$X := [x_1; \dots; x_n]^T \quad \text{of dimension } n \times p:$$

**Goal:** find  $q$  good directions  $u_k$  that preserve important aspects of the data

- two equivalent ways to define what are good directions
  - find directions of maximum variation
  - find projections that minimize the reconstruction error
- the directions  $u_1; \dots; u_q$  turn out to be the top  $q$ -eigenvalues of the covariance matrix  $X^T X$



# Principal Component Analysis (Maximum Variance)

1. Find the vector  $u_1$  such that the projection of the data on  $u$  has the greatest variance.

$$\begin{aligned} u_1 &\geq \arg \max_{\|u\|=1} \frac{1}{n} \sum_{i=1}^n (u_1^T x_i - u_1^T \bar{x}_i)^2 \\ &= \arg \max_{\|u\|=1} u^T X^T X u \end{aligned}$$

- ) this is the principal eigenvector of  $X^T X$ .
2. More generally, if we wish a  $k$ -dimensional subspace we choose  $u_1; \dots; u_k$  the top  $k$  eigenvectors of  $X^T X$ .
  3. The  $u_j$  form a new orthogonal basis of the data

## Principal component analysis (Minimum reconstruction Error)

We can also think of PCA as minimizing the reconstruction error of the compressed data.

Consider the linear transformation from  $\mathbb{R}^q$  to  $\mathbb{R}^p$  for  $q \leq p$ :

$$f : z \in \mathbb{R}^q \mapsto \mu + U_q z$$

where  $\mu \in \mathbb{R}^p$  is a location vector and  $U_q$  is a  $p \times q$  matrix with  $q$  orthogonal unit vectors.

The goal of PCA is to fit such a linear model to the data by **minimizing the reconstruction error**

$$\min_{\{z_i\}; U_q} \sum_{i=1}^n \|x_i - U_q z_i\|^2$$

The solution turns out to be the same as maximal variation.

## How many components should we choose?

If we use principal components as a summary of our data, how many components are sufficient?

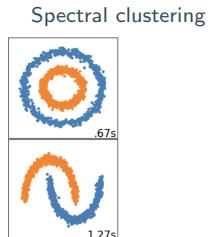
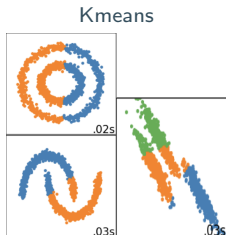
- No simple answer to this question. Cross validation as we would use for supervised learning is not possible.
- We may choose by looking visually at the reconstruction error or the percentage of variance explained.

**PCA** looks for a low-dimensional representation of the observations that explains a good fraction of the variance.

**Clustering** looks for homogeneous subgroups among the observations.

They can be combined in methods such as **Spectral Clustering**:

1. apply PCA like techniques to find a better representation of the data
2. apply Kmeans



Unsupervised learning is important for **understanding the variation and grouping structure of a set of unlabeled data**, and can be a useful pre-processor for supervised learning.

It is intrinsically **more difficult than supervised learning** because there is no gold standard (like an outcome variable) and no single objective (like test set accuracy).

It is an **active field of research**, with many recently developed tools such as self-organizing maps, independent components analysis and spectral clustering.

## References

---



L. Bottou, F. E. Curtis, and J. Nocedal. "Optimization methods for large-scale machine learning". In: *arXiv preprint arXiv:1606.04838* (2016).



N. Cesa-Bianchi and G. Lugosi. *Prediction, learning, and games*. Cambridge University Press, 2006.



T. J. Hastie, R. J. Tibshirani, and J. H. Friedman. *The Elements of Statistical Learning*. Springer Series in Statistics. New York, NY, USA: Springer New York Inc., 2001.



Wikipedia. The free encyclopedia.