

UNSUPERVISED ML

OBJECTIVE

- Discover information from data without labeled examples
- Extract some hidden organisation, patterns, relation between elements
- Extract a (statistical ?) model of the data ?

OBJECTIVE

- Typical objectives:
 - ▶ Cluster discovery
 - ▶ Anomaly Detection
 - ▶ Latent variable discovery / Embedding / dimensionality reduction...

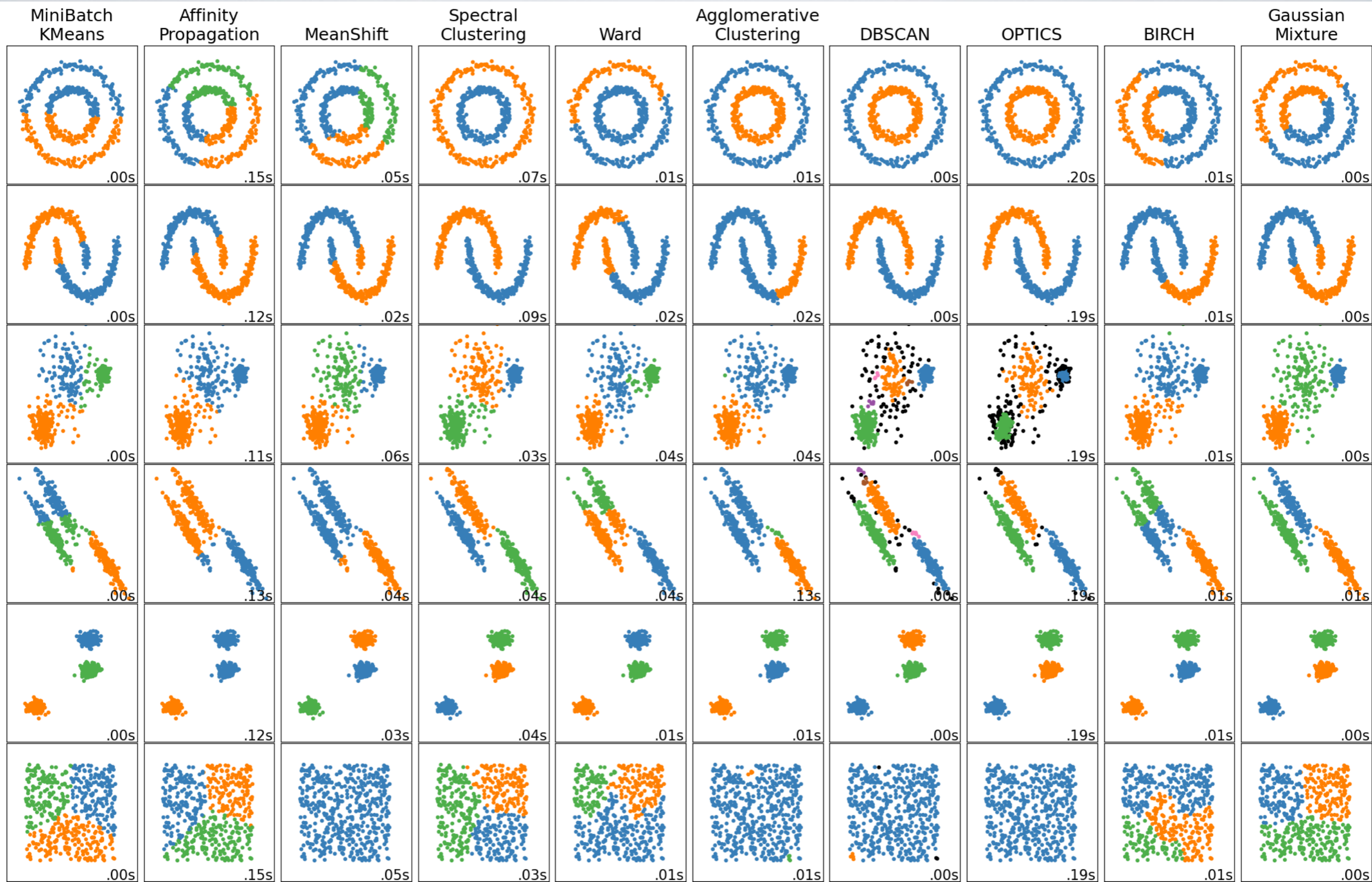
CLUSTERING

CLUSTERING

- The most famous unsupervised ML problem
- 100+ methods exist
 - Most people use “good old” methods: k-means (1967), DBSCAN (1996)
 - They are often “good enough”, well implemented, safe, ...
- Part of the problem: Clustering is not well defined
 - What is “a good cluster”?

CLUSTERING

- How would you define a good cluster ?
- A good partition in clusters ?



K-MEANS

- Definition:
 - ▶ For a target number of clusters k
 - ▶ Find the item assignment minimizing
 - The inter-cluster variance
 - Equivalently \Rightarrow The squared distance from points to their cluster center
 - Equivalently \Rightarrow The squared distance between cluster elements
- This is only one possible objective !
 - ▶ Why this one ?
 - ▶ Intuitive, good properties...

K-MEANS

$$\arg \min_{\mathbf{S}} \sum_{i=1}^k \sum_{\mathbf{x} \in S_i} \|\mathbf{x} - \boldsymbol{\mu}_i\|^2 = \arg \min_{\mathbf{S}} \sum_{i=1}^k |S_i| \text{Var}(S_i)$$

with

\mathbf{S} a cluster assignment,

x a d dimensional item, and

$\boldsymbol{\mu}_i$ the mean of items in cluster \mathbf{S}_i .

Note that without fixing k , there is a trivial solution with each item alone in its own cluster.

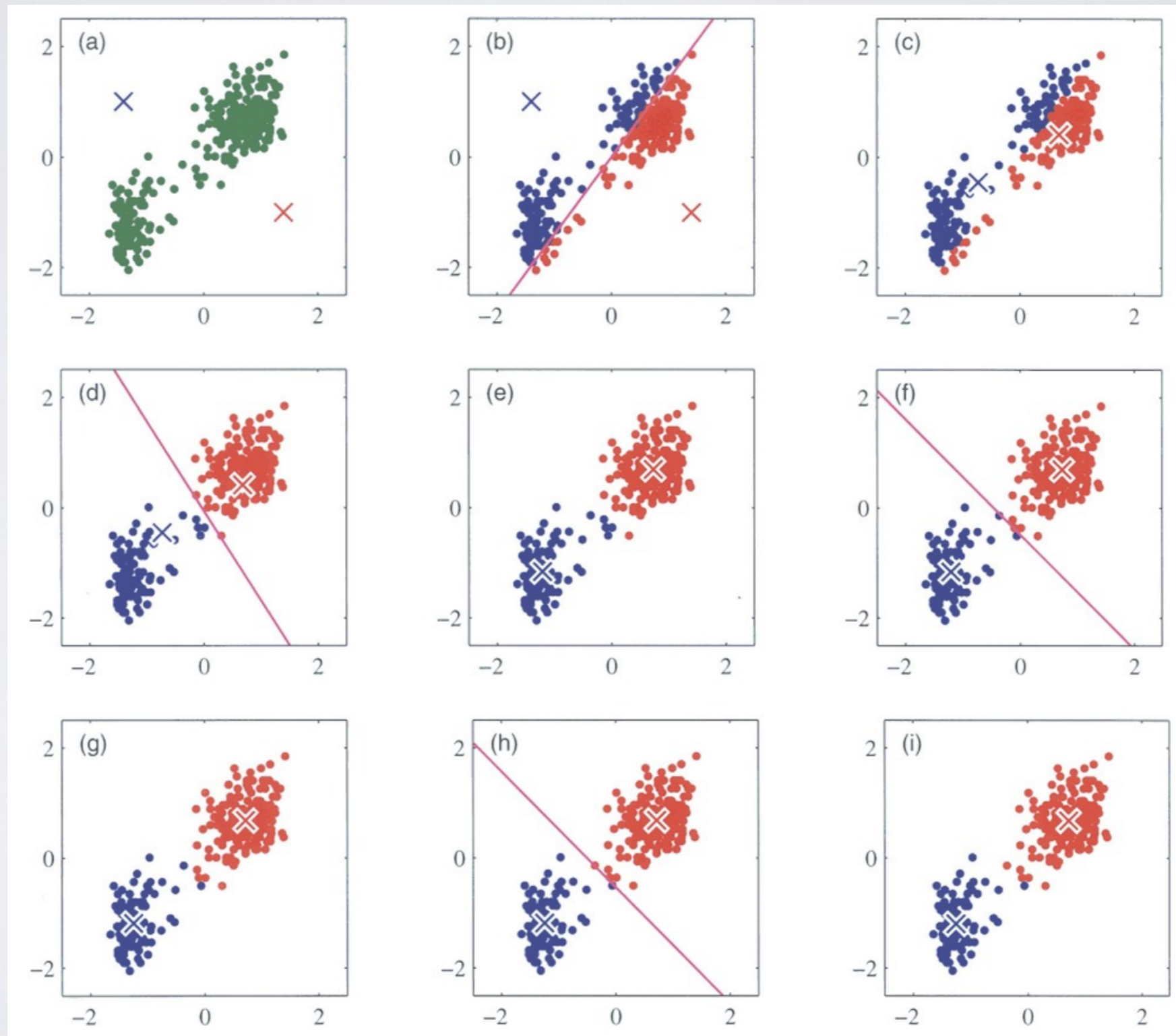
K-MEANS

- Discovering the global optimum is NP-hard
- How to find quickly a good solution ?
 - Naive k-means
 - K-means ++ (used in most current implementations)
 - Use optimized data structure (KDtrees...)

NAIVE K-MEANS

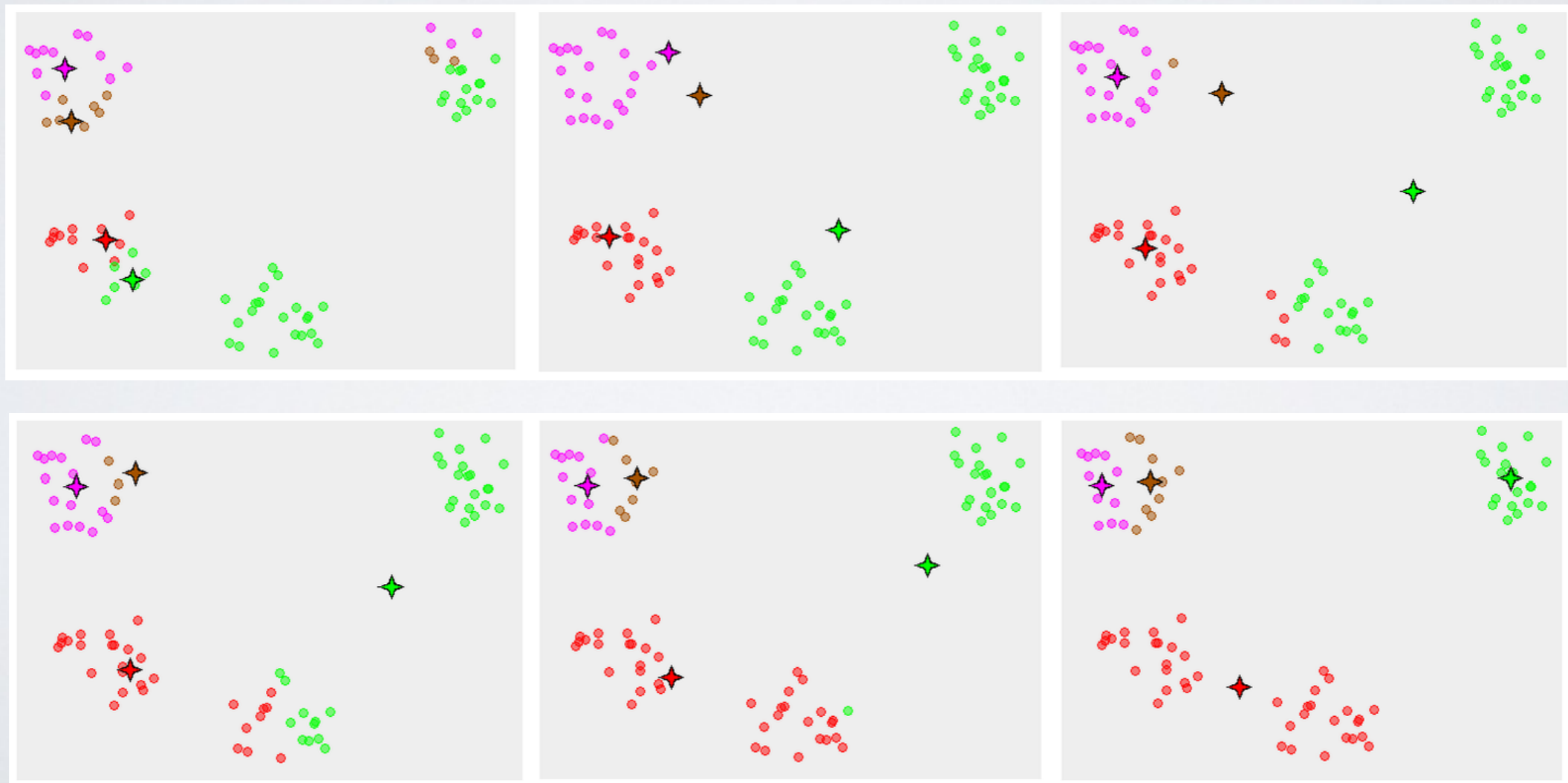
- 1) Assignment: Assign each item to its closest cluster center
- 2) Update: Recompute the center of each cluster as the mean (centroid) of items that compose that cluster
- Start with random centroids

NAIVE K-MEANS



NAIVE K-MEANS

- Known limit: convergence to poor local minimum if poor initial centroids



K-MEANS++

- Several variants to choose wisely the initial centroids
- K-means++ is proven to improve the results, statistically
 - Not always, but improves more often than deteriorate the results.

K-MEANS++

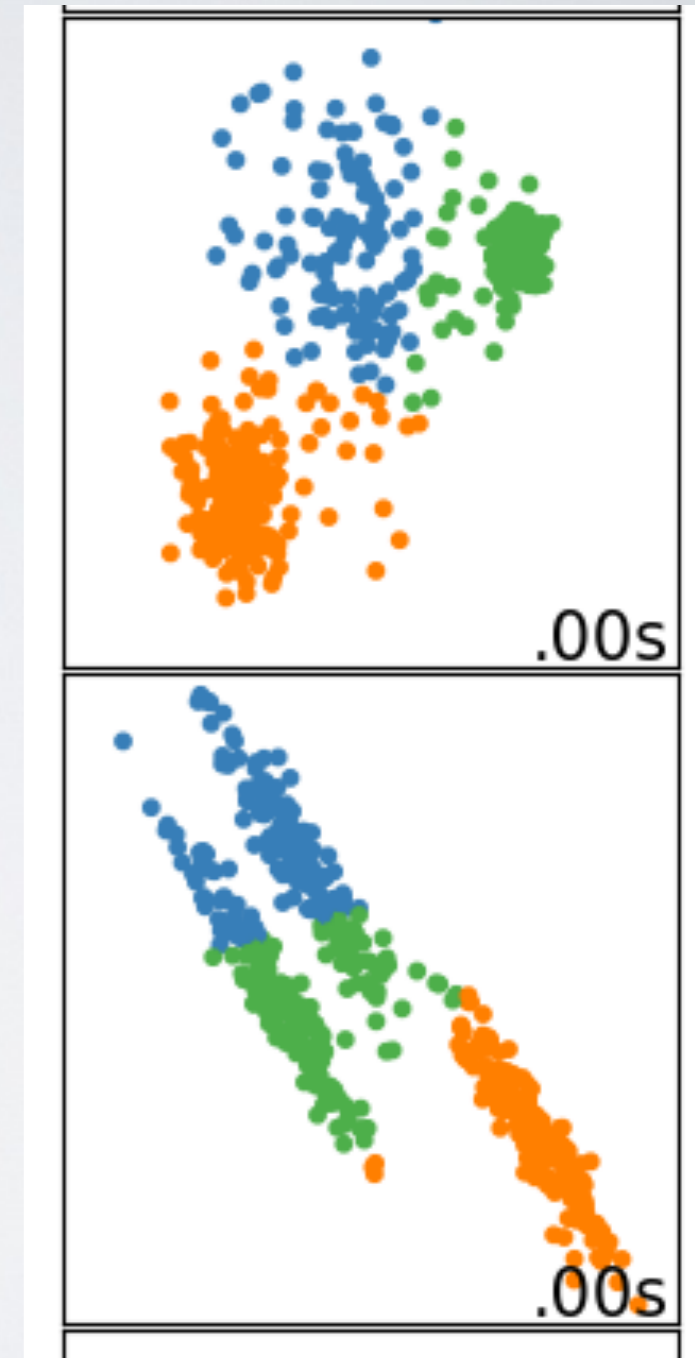
1. Choose one center uniformly at random among the data points.
2. For each data point x not chosen yet, compute $D(x)$, the distance between x and the nearest center that has already been chosen.
3. Choose one new data point at random as a new center, using a weighted probability distribution where a point x is chosen with probability proportional to $D(x)^2$.
4. Repeat Steps 2 and 3 until k centers have been chosen.

K-MEANS++



WEAKNESSES

- We can identify some clear weaknesses:
 - ▶ K-means has a tendency to search for clusters of equal sizes (minimize **overall** cluster variance)
 - ▶ Clusters tend to be **circular**, since all directions are worth the same.



NORMALIZATION

- Very important point: k-means is based on **euclidean distance**.
 - ▶ We minimize the inter-cluster euclidean distance between points
 - ▶ We could adapt the method to other distances
- Data need to be **normalized/standardized**
 - ▶ Clustering based on age in years and revenue in \$. The “distance” in \$ will dominate
 - ▶ Remember: normalization/standardization are not fixing magically problems (outliers..)
 - You need to **think**: is 1 unit in one dimension *worth* 1 unit in other dimensions ?

GAUSSIAN MIXTURES

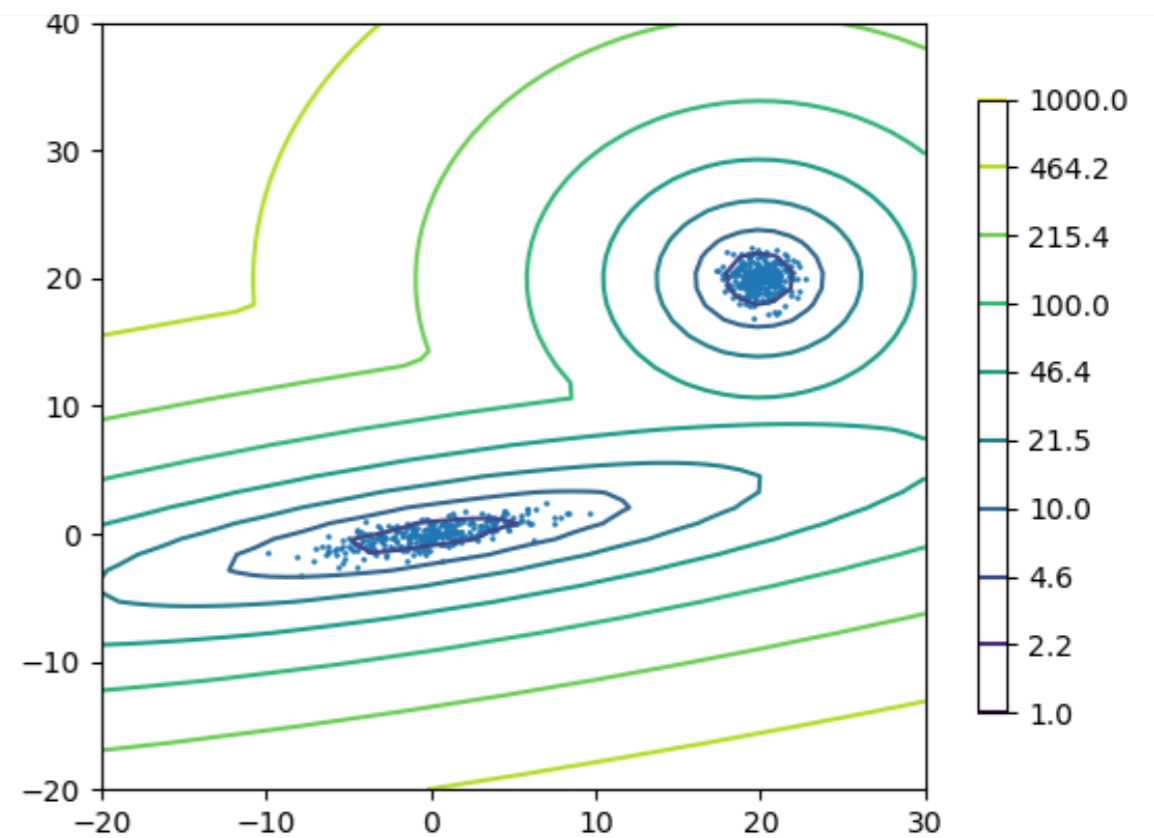
GAUSSIAN MIXTURES

- Generalize k-means concept:
 - Clusters are sets of points that are close in euclidean space
 - Different clusters tend to be far appart
- Translate it statistically:
 - Each cluster can be described using a normal distribution centered on its centroid, with the probability of observing points decreasing with the distance to the centroid.

GAUSSIAN MIXTURES

Train accuracy: 88.4

Test accuracy: 92.1



GAUSSIAN MIXTURES

- We define a **generative model** for k clusters
 - Each cluster corresponds to a gaussian distribution, defined by a center and a *variance, or covariance matrix*
 - The problem to solve is to find the parameters Θ (centers, variances) that maximize the likelihood of the corresponding model to generate the observed items X
 - More formally, we are searching for: $\arg \max_{\Theta} p(X | \Theta)$

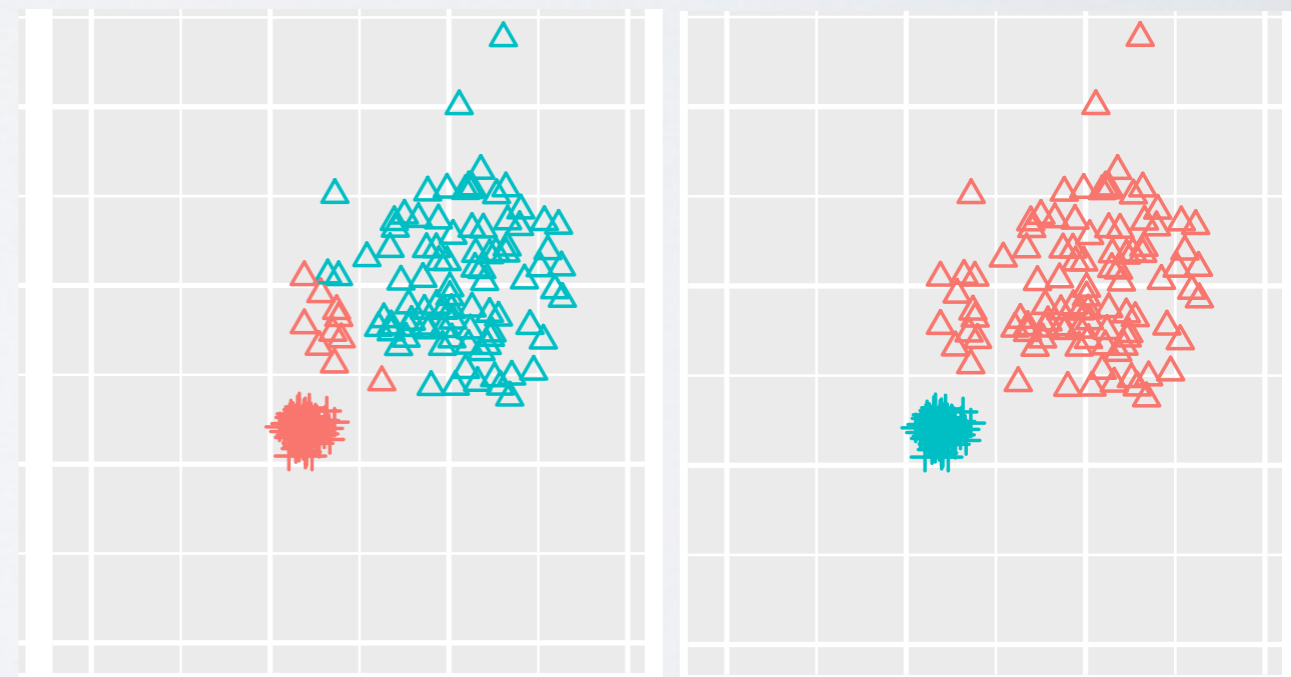
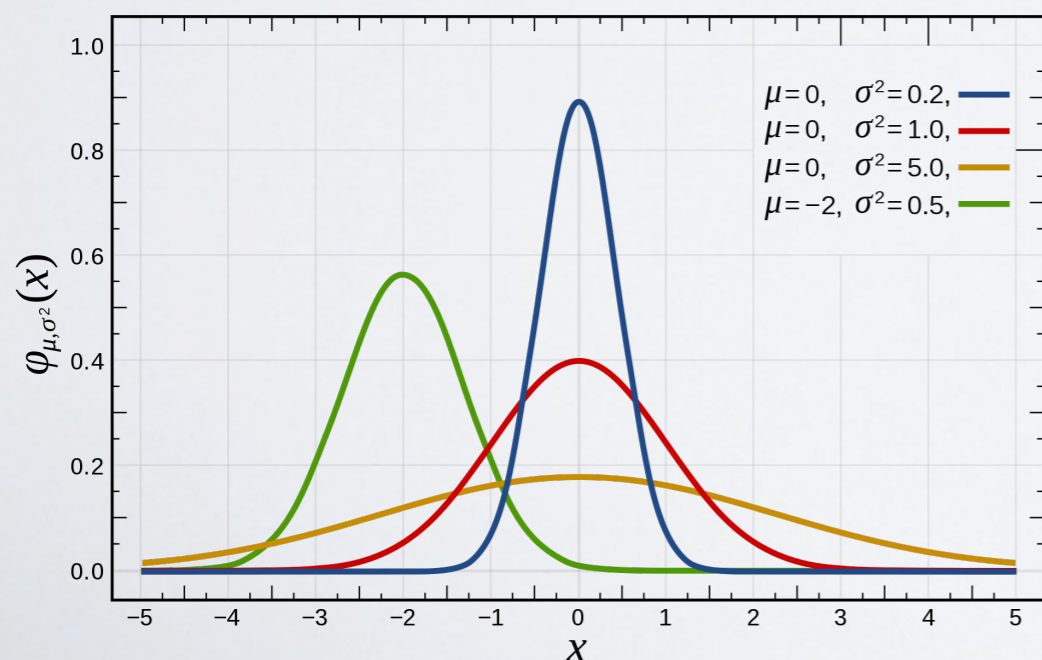
K-MEANS EQUIVALENCE

$$\begin{bmatrix} \text{Var}(x_1) & \dots & \text{Cov}(x_n, x_1) \\ \vdots & \cdot & \vdots \\ \text{Cov}(x_n, x_1) & \dots & \text{Var}(x_n) \end{bmatrix}$$

- If we assume that:
 - ▶ The gaussian distributions are defined only by their variance, not by complete covariance matrices
 - Similar in all directions, “spherical”
 - ▶ The variance value is the same for all gaussian distributions
 - Spheres of the same “size”
 - ▶ The probability for each item to be generated by each of the gaussian distribution is identical
- Then it can be shown that the objective is equivalent to the k-means objective !
 - ▶ We can relax some of those constraints to get better results

DENSITY HETEROGENEITY

- Allowing denser/sparser clusters
 - Consider the case in which gaussians are defined by a single value of variance (covariance=0)
 - If they differ for each clusters, some can be denser than others

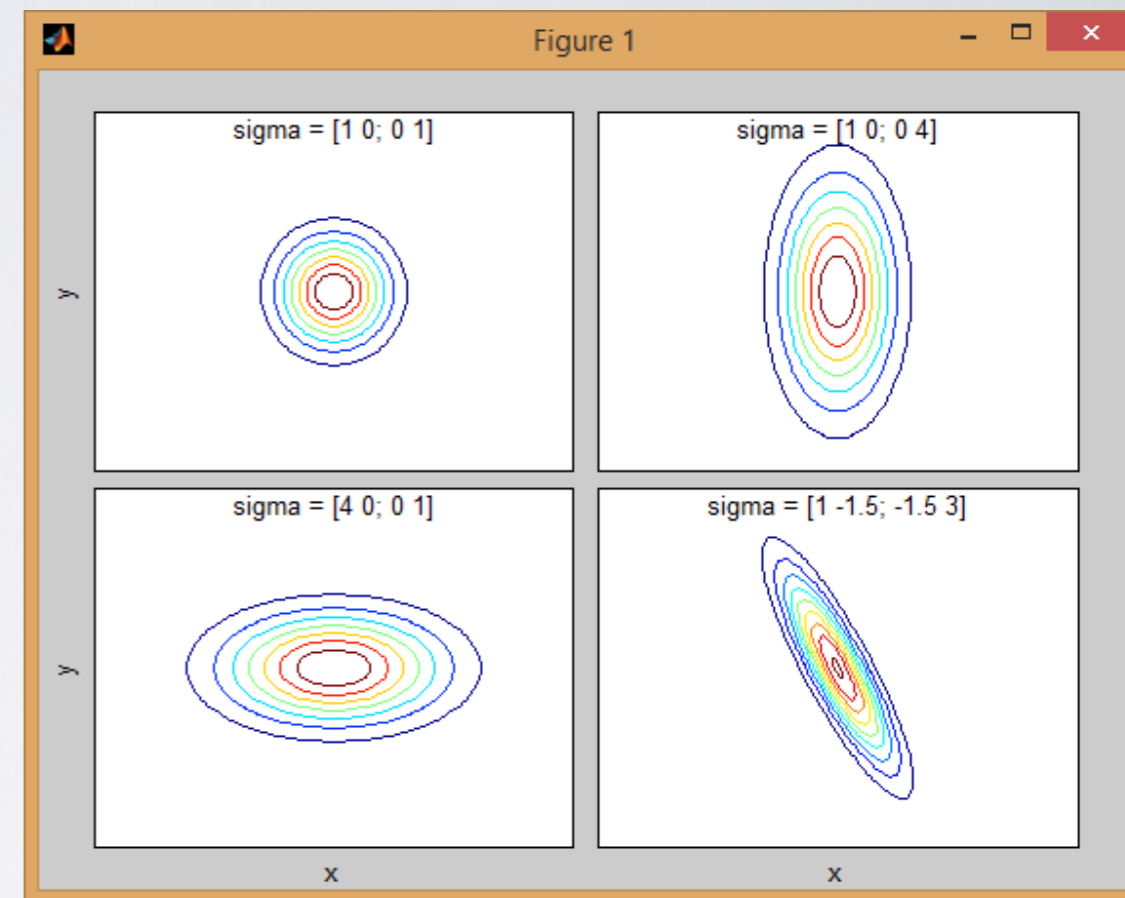


K-means

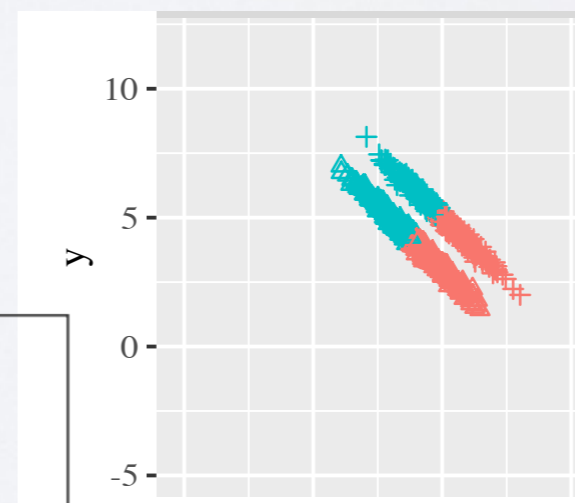
GM,
free variance

SHAPE VARIATIONS

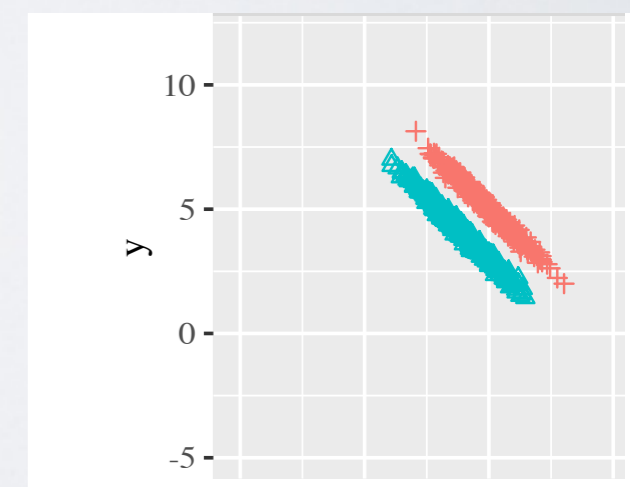
- Allowing non-circular shaped clusters
 - ▶ If values on the diagonal of the covariance matrix differs, the matrix can have ellipsoidal shape, in the direction of the axes
 - ▶ If the full covariance matrix is inferred, any ellipsoidal shape can be obtained



$$\begin{bmatrix} \text{Var}(x_1) & \dots & \text{Cov}(x_n, x_1) \\ \vdots & \ddots & \vdots \\ \text{Cov}(x_n, x_1) & \dots & \text{Var}(x_n) \end{bmatrix}$$



K-means

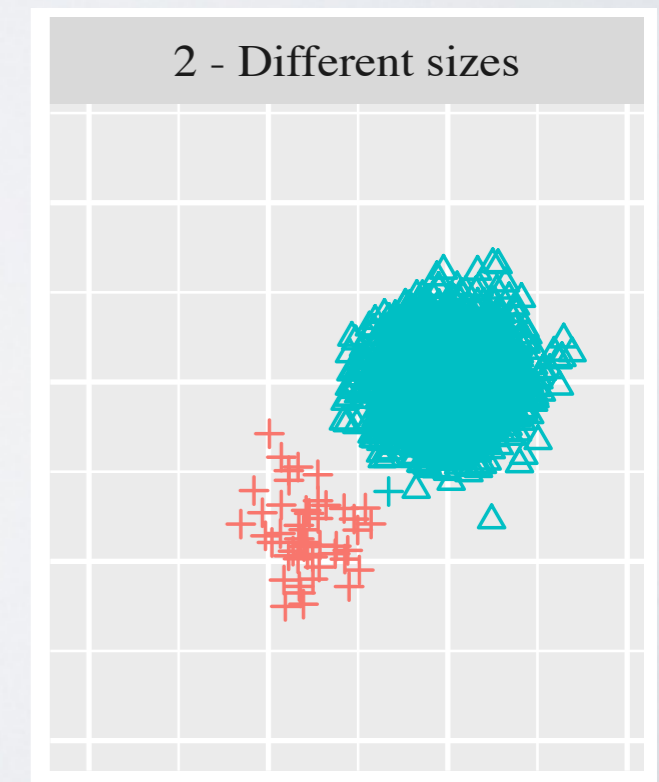
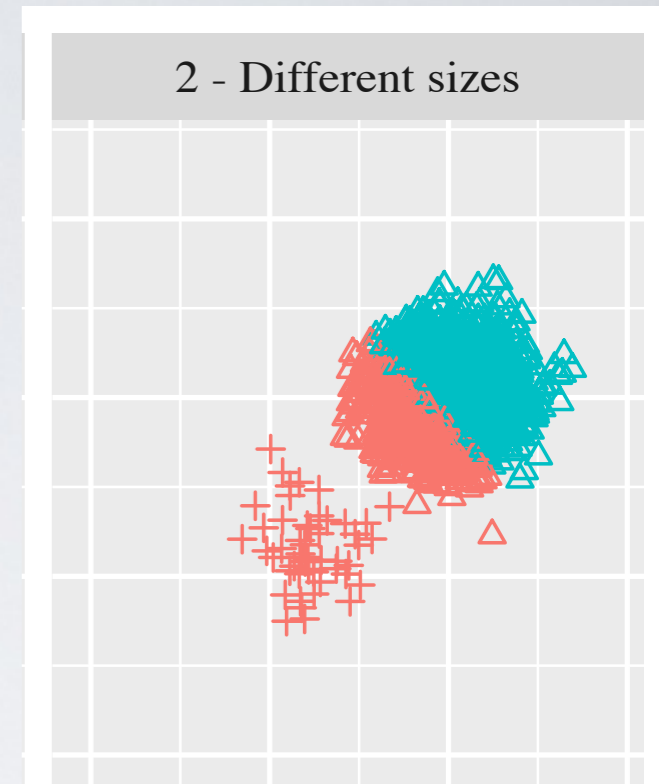


Full gaussian

SIZE HETEROGENEITY

- The fraction of all items generated by each generative gaussian (e.g., cluster) is the same.
- We usually add a *strength* parameter π to weight the fraction of items generated by each cluster

$$p(X) = \sum_{k=1}^K \pi_k G(X | \mu_k, \sigma_k)$$



ALL TOGETHER

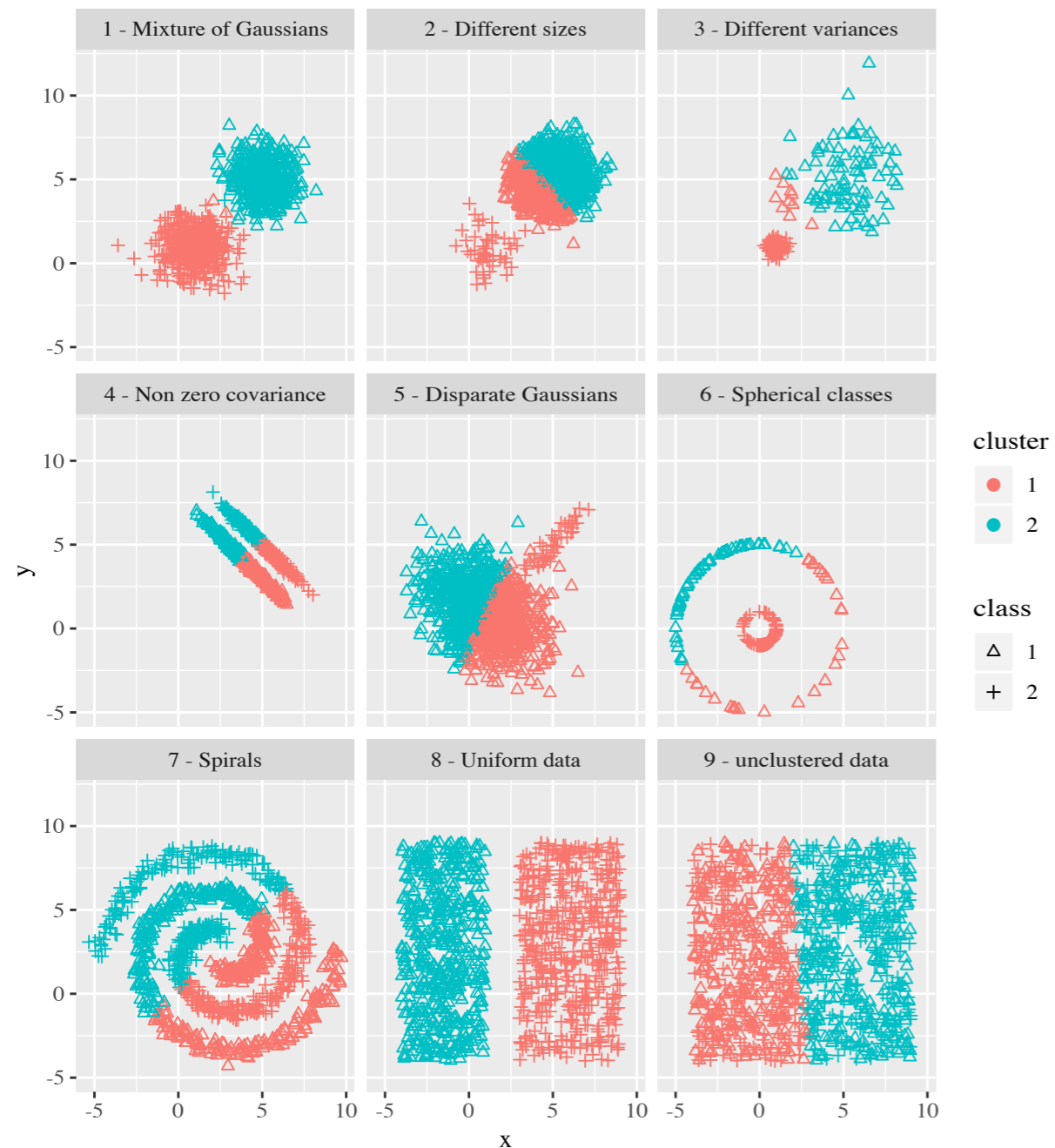
$$p(X) = \sum_{k=1}^K \pi_k G(X | \mu_k, \sigma_k)$$

$$\arg \max_{\Theta} p(X | \Theta)$$

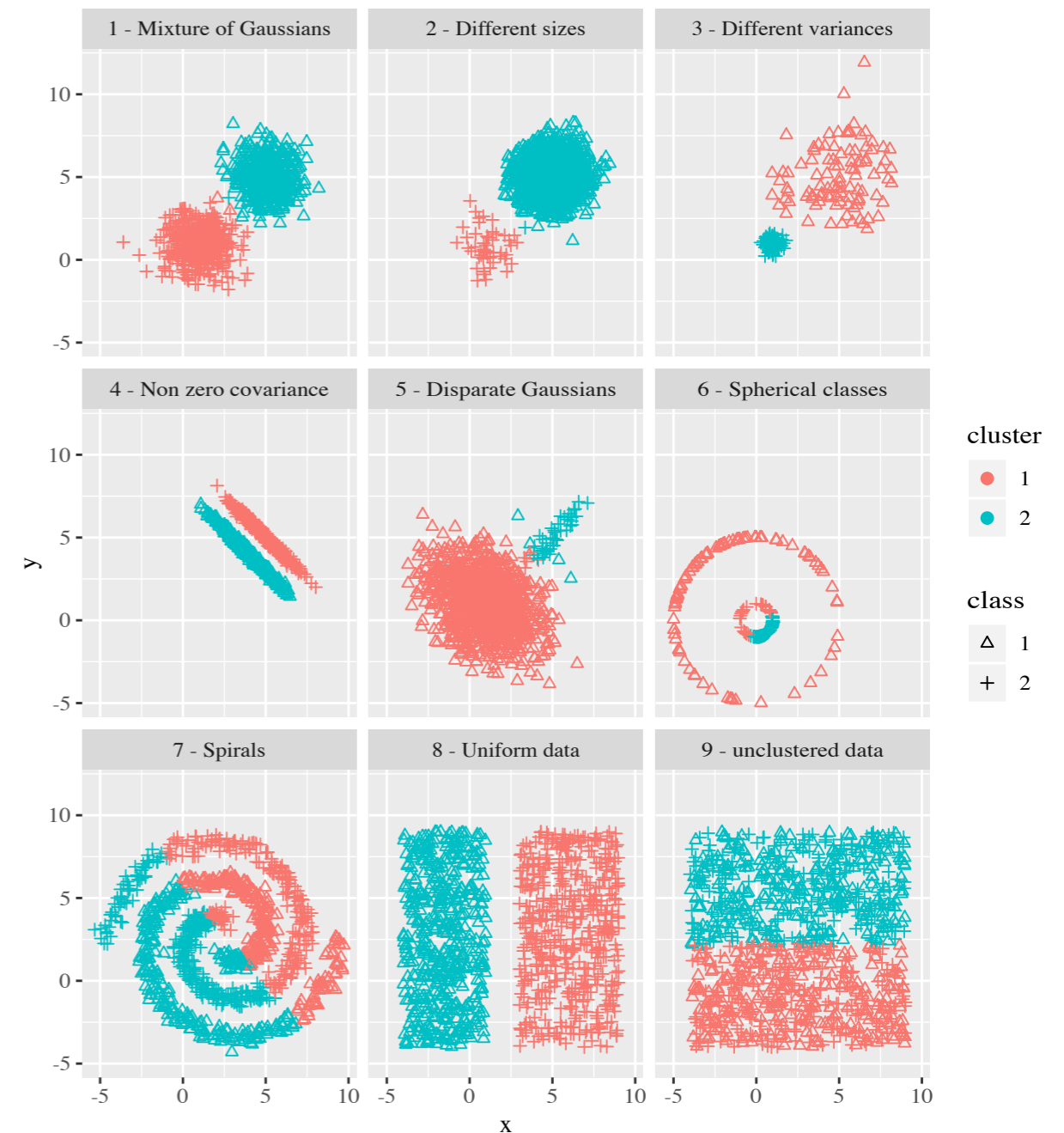
$$\Theta = \mu, \sigma, \pi$$

K-MEANS COMPARISON

K-means



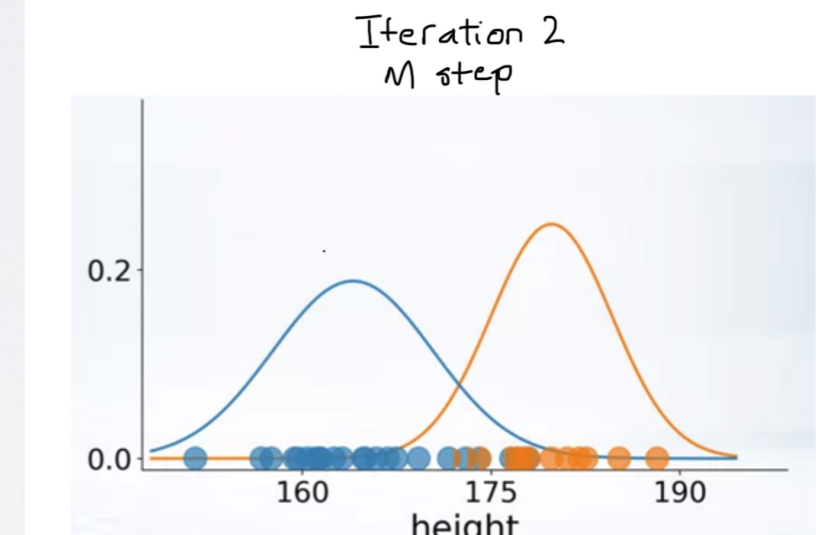
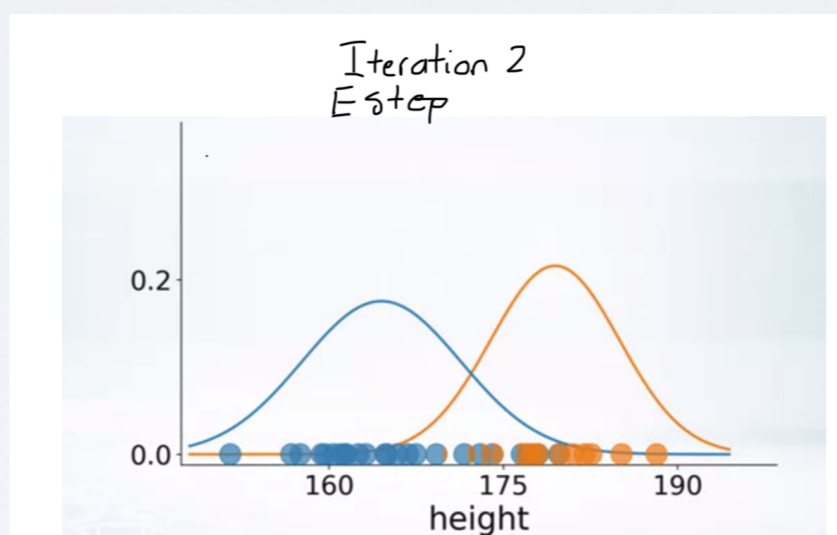
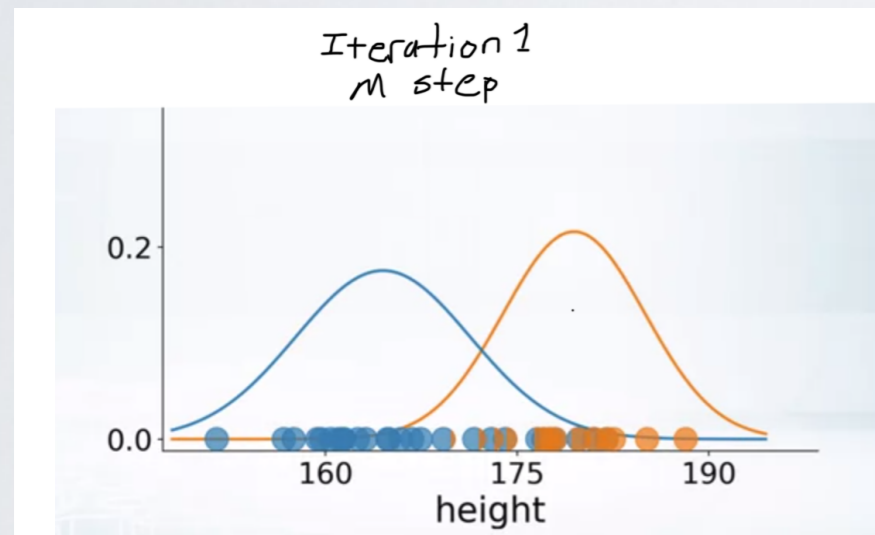
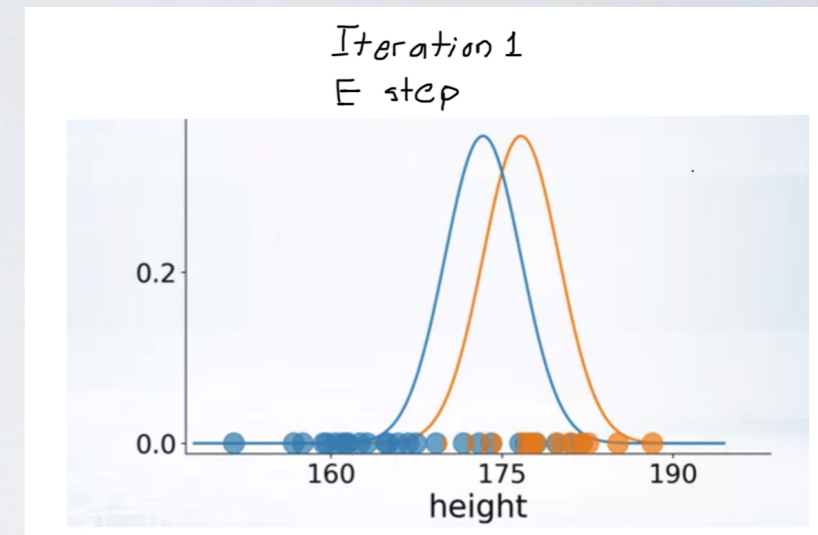
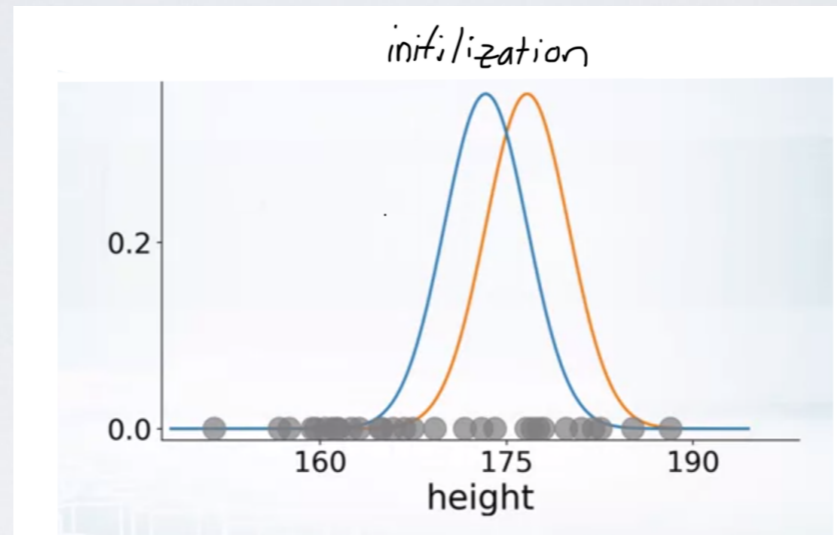
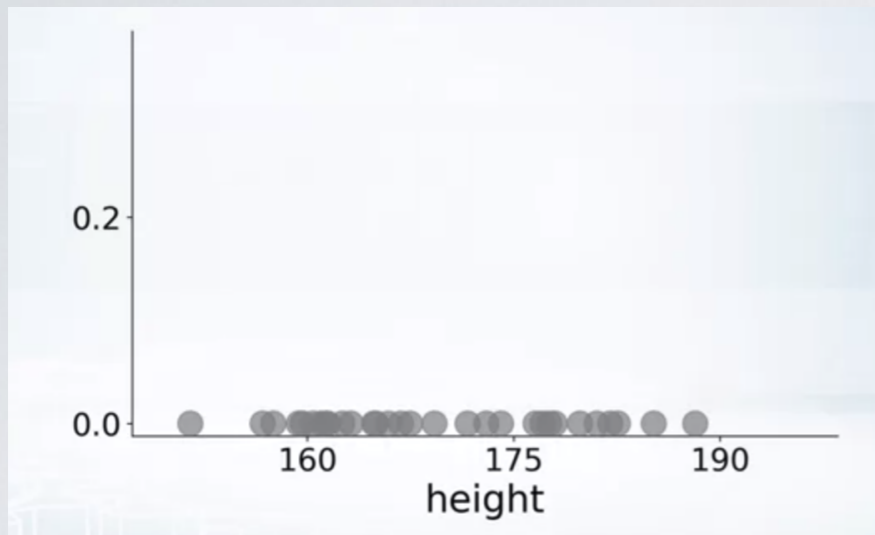
Full Gaussian Mixture



EM ALGORITHM

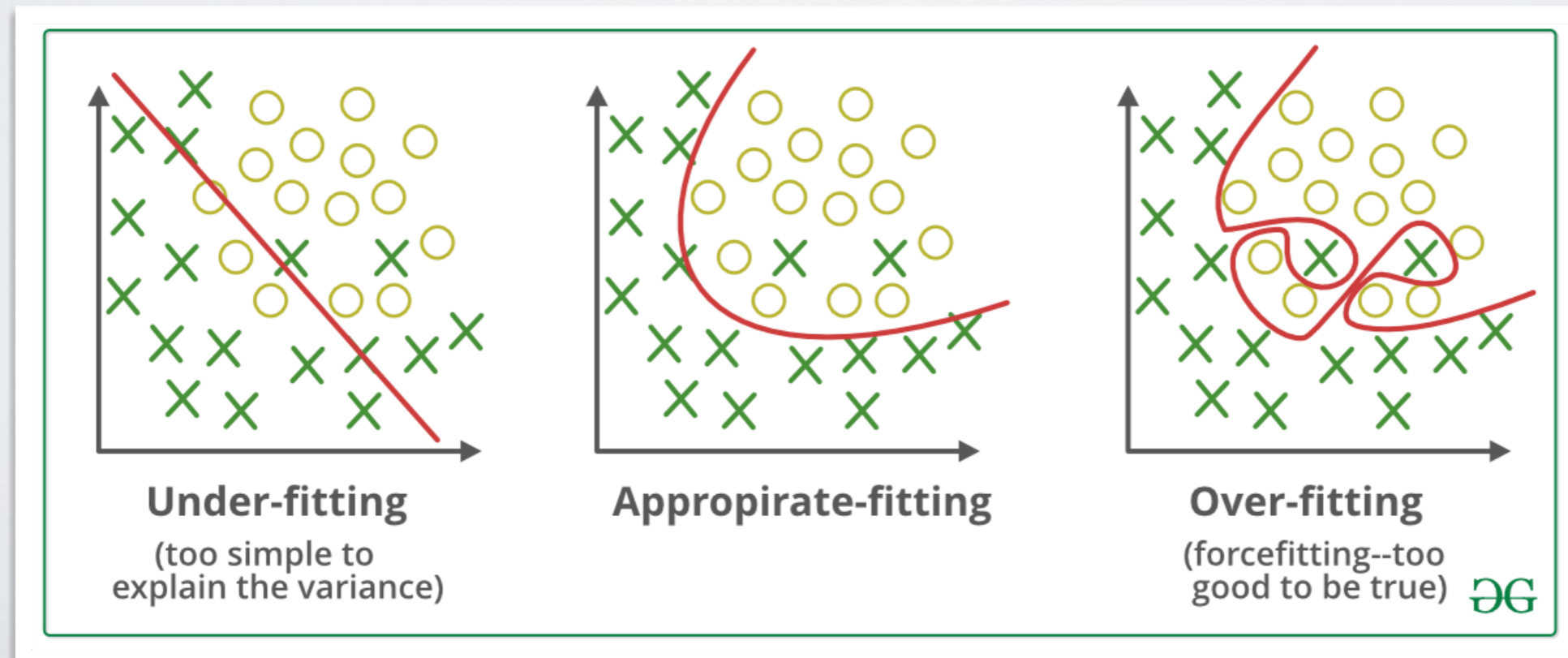
- To search for the parameters, we can use a method similar to naive k-means known as EM (Expectation Maximization)
 - Note Z the cluster assignation of items to their **most likely** clusters
 - 1) Initialize parameters Θ to random values
 - 2)(E) Compute Z , given Θ
 - 3)(M) Use assignations in Z to update values of Θ
 - 4) Iterate steps 2 and 3 until convergence

EM ALGORITHM



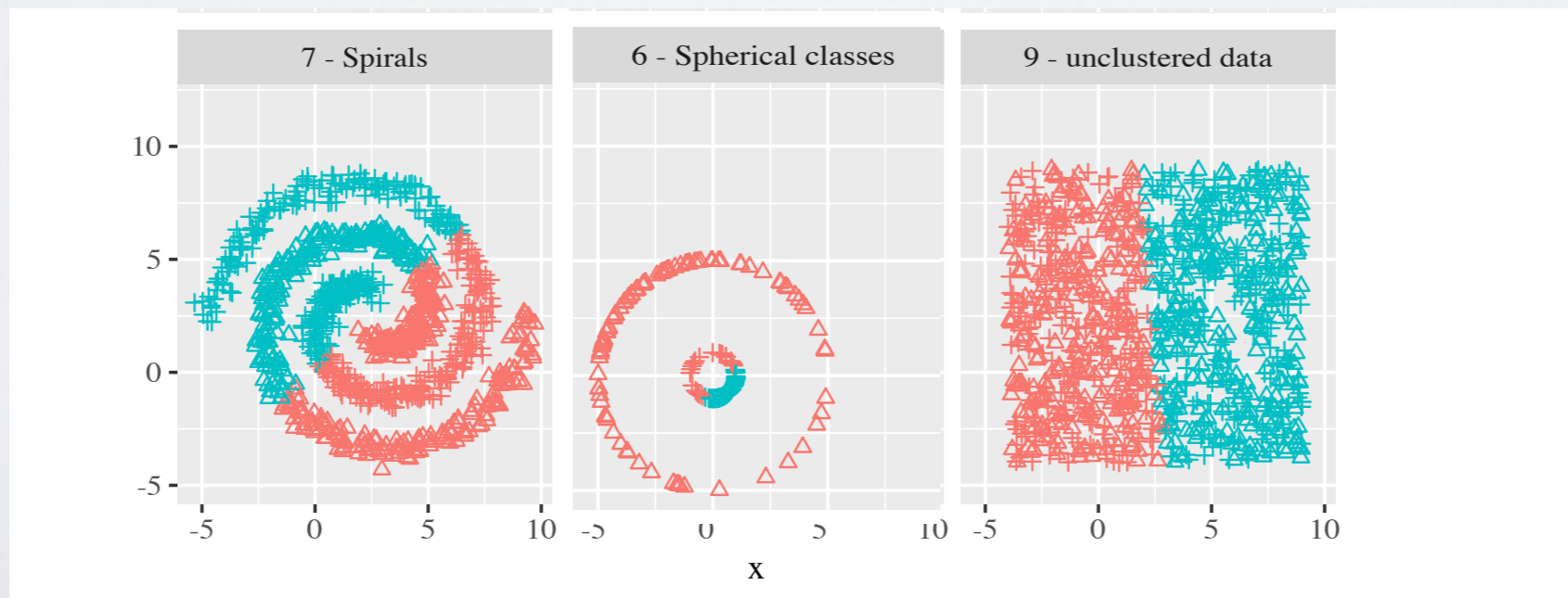
PROS AND CONS

- Gaussian mixture seems an improvement over k-means. Why not always using it?
 - Force of habits
 - Higher computational cost (More parameters => More complex problem)
 - Higher possibility of overfitting (More parameters => More overfit risk)



REMAINING PROBLEMS

- We can mention 3 problems remaining (at least)
 - ▶ The number of clusters still need to be provided.
 - If allowed to change, it will always converge to the trivial solution with each item in its own cluster
 - ▶ If the data is completely random, the method still finds clusters
 - ▶ Impossible to discover non-convex structures, such as circles or spirals

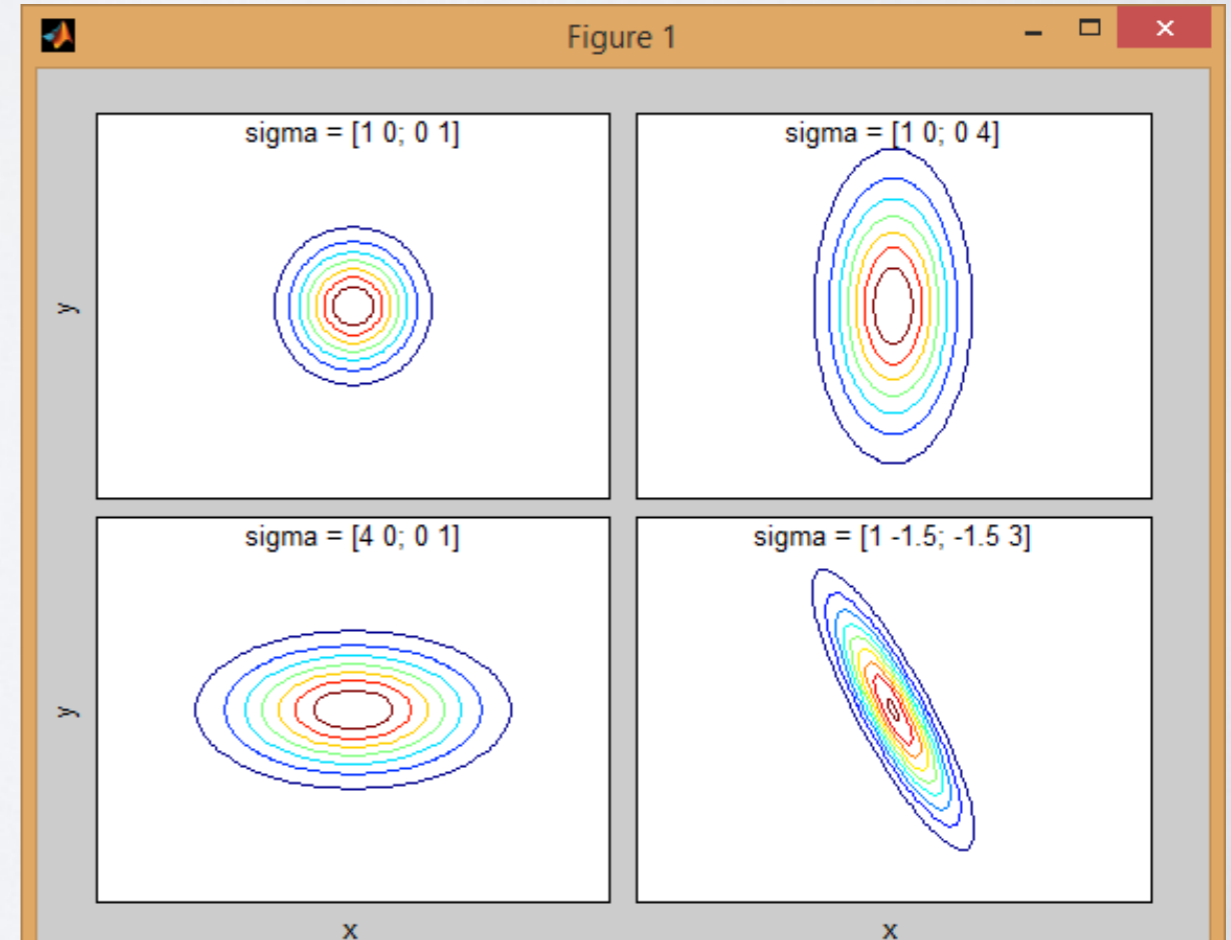


MDL

- Discovering automatically the number of clusters —and thus finding no clusters in random data— is possible using an MDL approach
- MDL = Minimum Description Length
- The principle is to search a solution maximizing the compression rate, i.e., minimizing the *cost* of the description, e.g., in bits.
- Method introduced later

NORMALIZATION

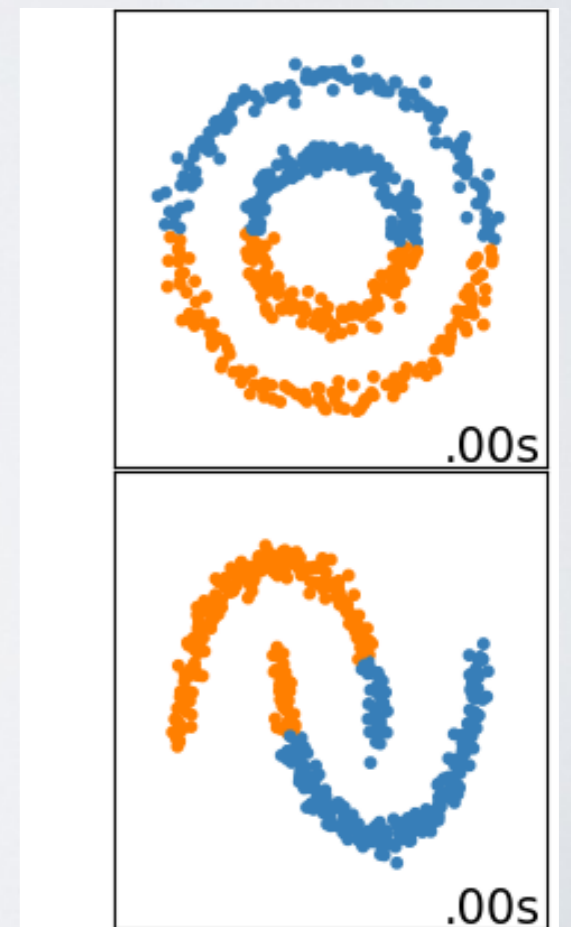
- Is normalization as important for full GM models as for k-means?



DBSCAN

K-MEANS/GM LIMITS

- The problem of spiral/Circular/weird shaped clusters comes from the assumption that items of a cluster should be “normally distributed” around their mean



LOCAL DEFINITIONS

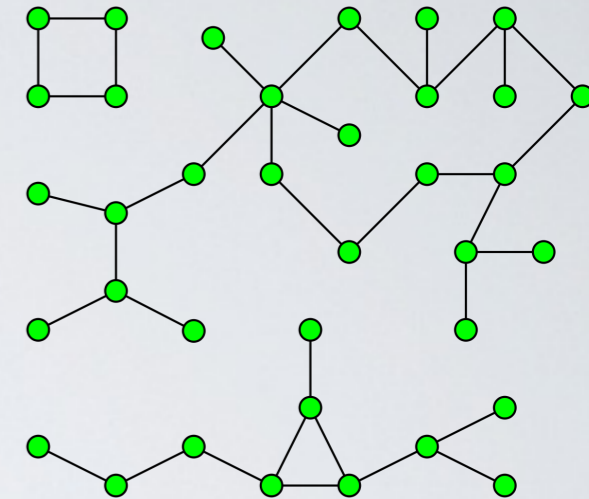
- To overcome this problem, several methods propose local definitions of clusters
 - Does not explicitly optimize a global function
 - Items belong to clusters because they are close enough, locally, to other items in that cluster
 - Clusters exist because there is continuum between all items in it, locally

DBSCAN

- Define some local parameters:
 - ▶ ϵ , the distance threshold above which items are considered “too different”
 - ▶ *minPts*, a minimal number of reachable points
 - ▶ No need to define a number of clusters !
- Define:
 - ▶ An item p is a *core point* if it has at least *minPts* items at distance less than ϵ
 - Including p itself

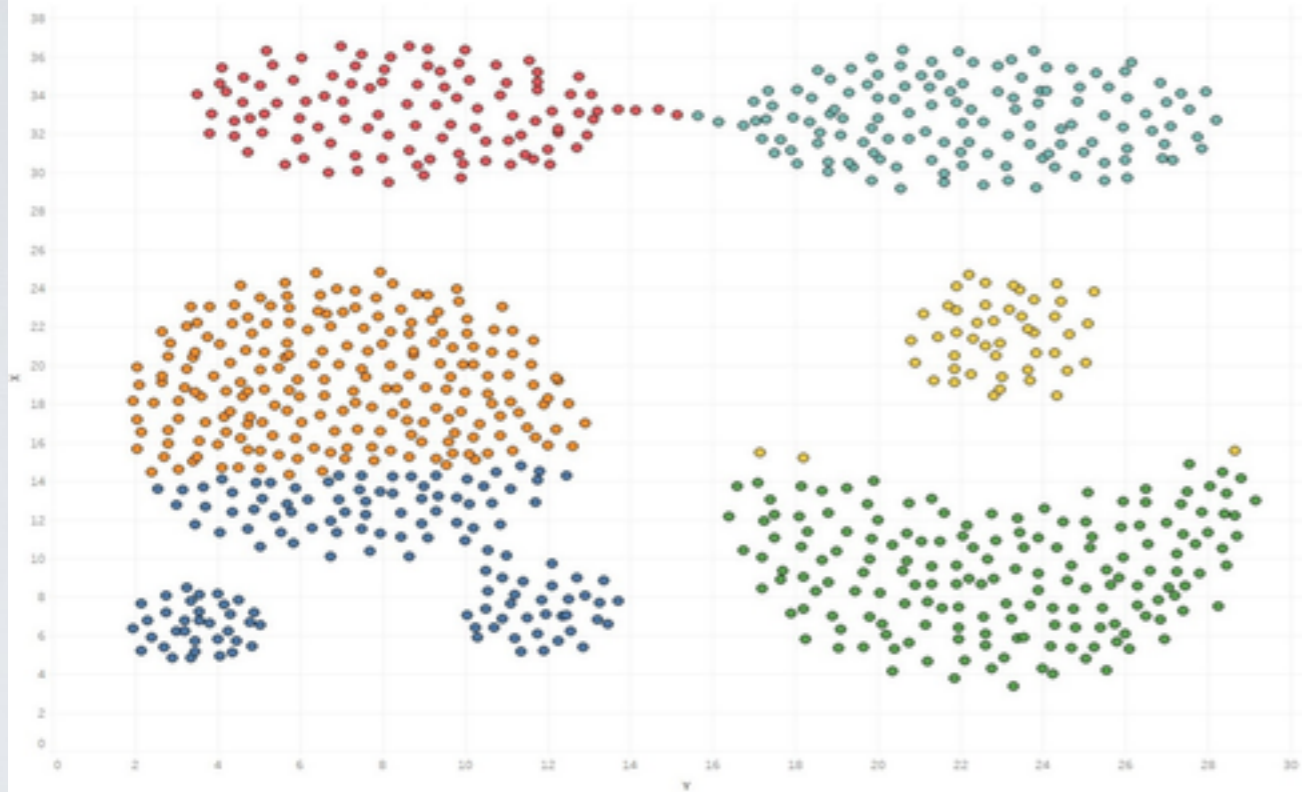
DBSCAN: GRAPH DEFINITION

- 1) Build a graph such as
 - Each core node is a node
 - A link exist between core nodes if they are at $d < \epsilon$
- 2) Detect the connected components of the graph
 - 2 nodes belong to the same connected components if there is a path between them
- 3) For all non-core nodes:
 - If they have no core points directly reachable, discard them as noise
 - Else, attribute them to (one of) the clusters for which one core point is directly reachable
 - Variant DBSCAN* => ignore those points as noise

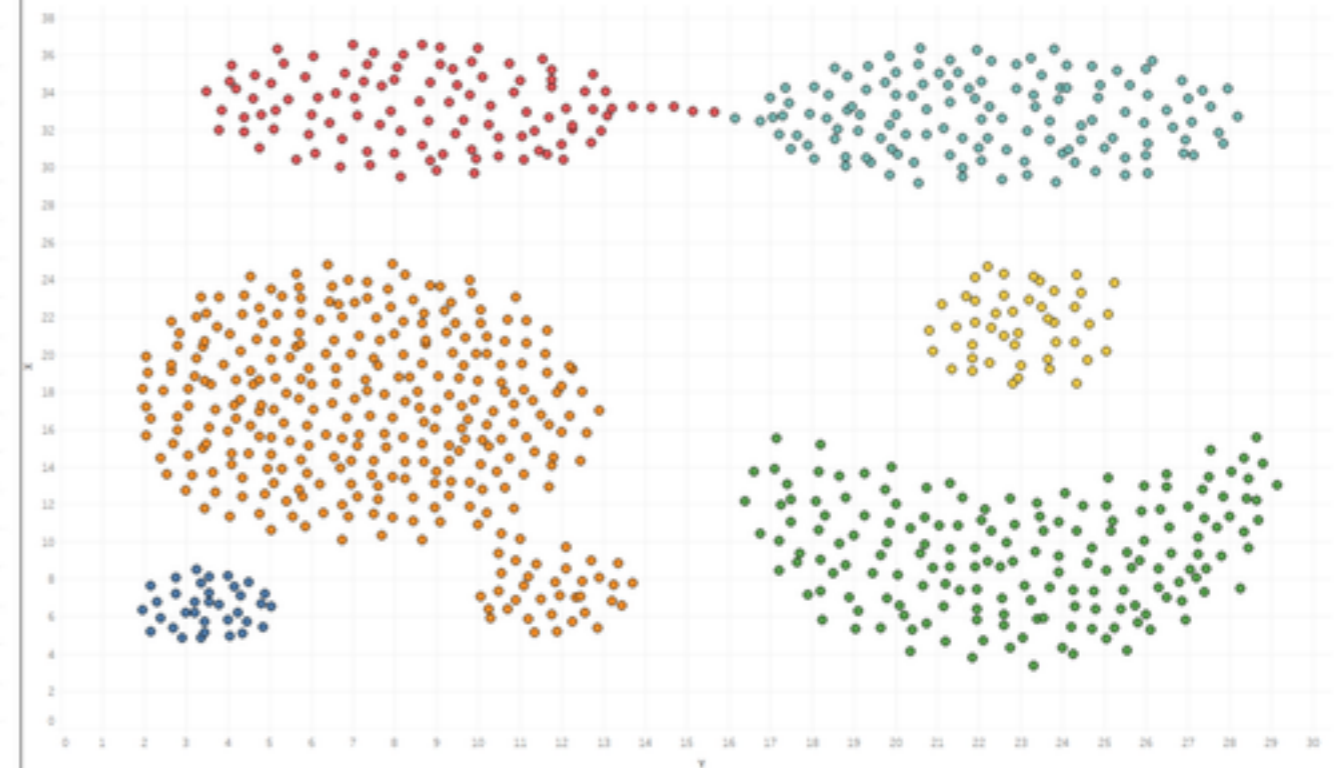


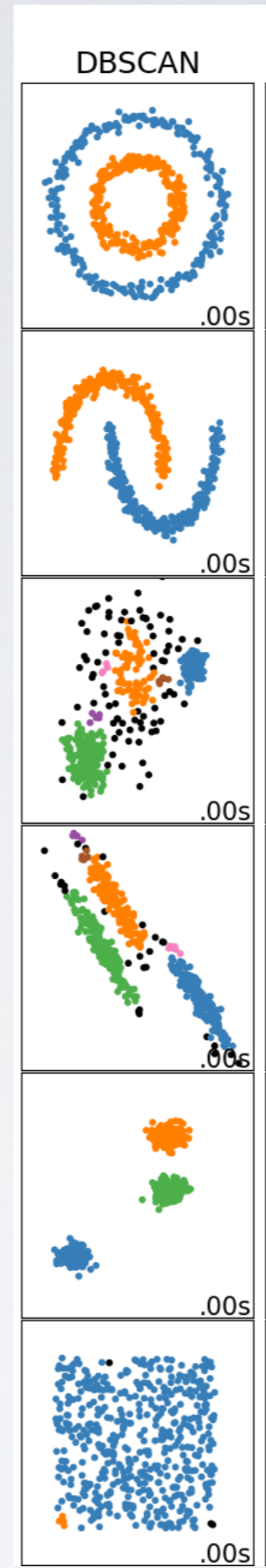
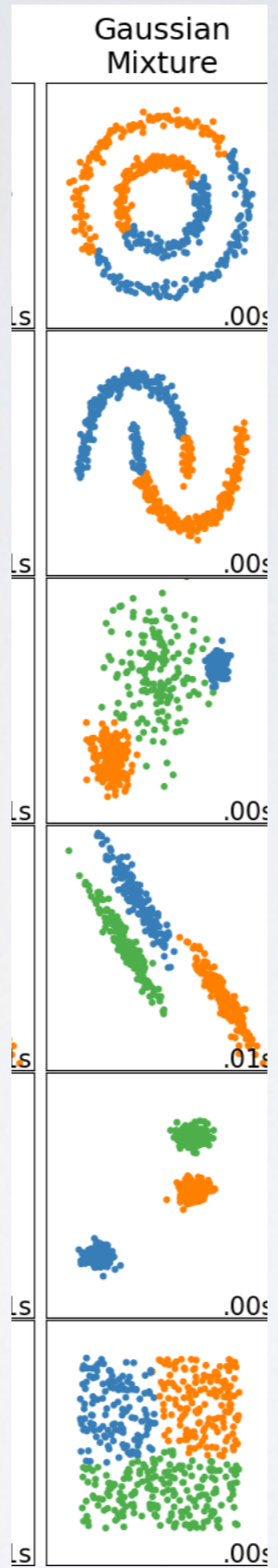
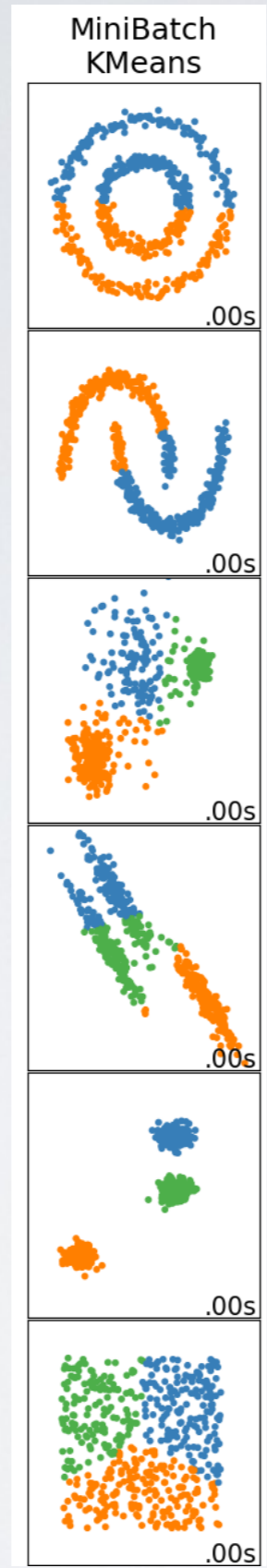
DBSCAN

Traditional Clustering (K-means)



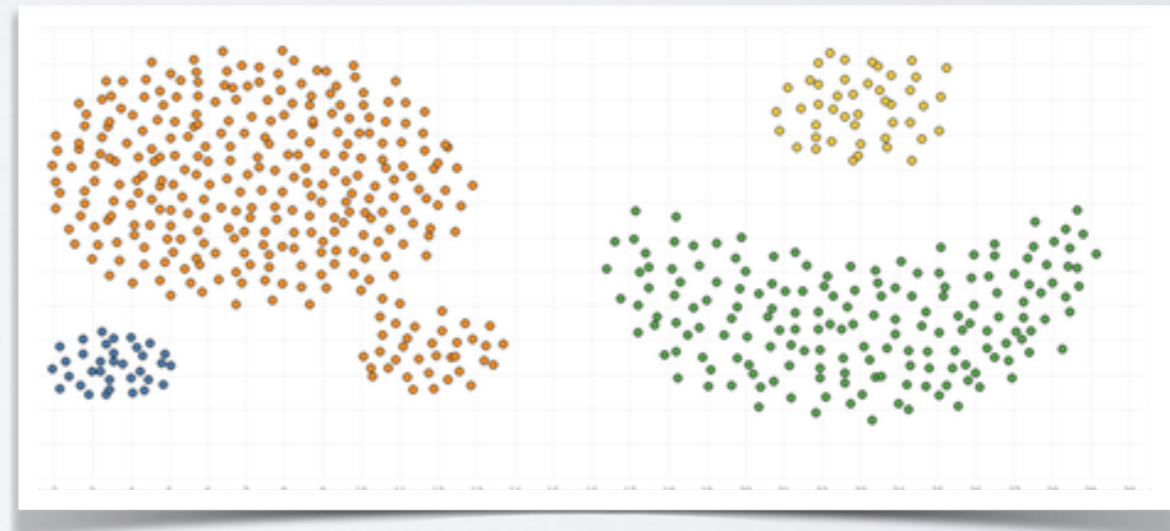
DBSCAN





DBSCAN

- Strength:
 - ▶ No need to define the number of clusters
 - ▶ Can discover arbitrarily-shaped clusters
 - ▶ A notion of noise
- Weaknesses
 - ▶ Defining ϵ is extremely difficult
 - Similar to the number of clusters.
 - In fact it determines the number of clusters...
 - ▶ Despite safeguards, risk of the stretched clusters effect



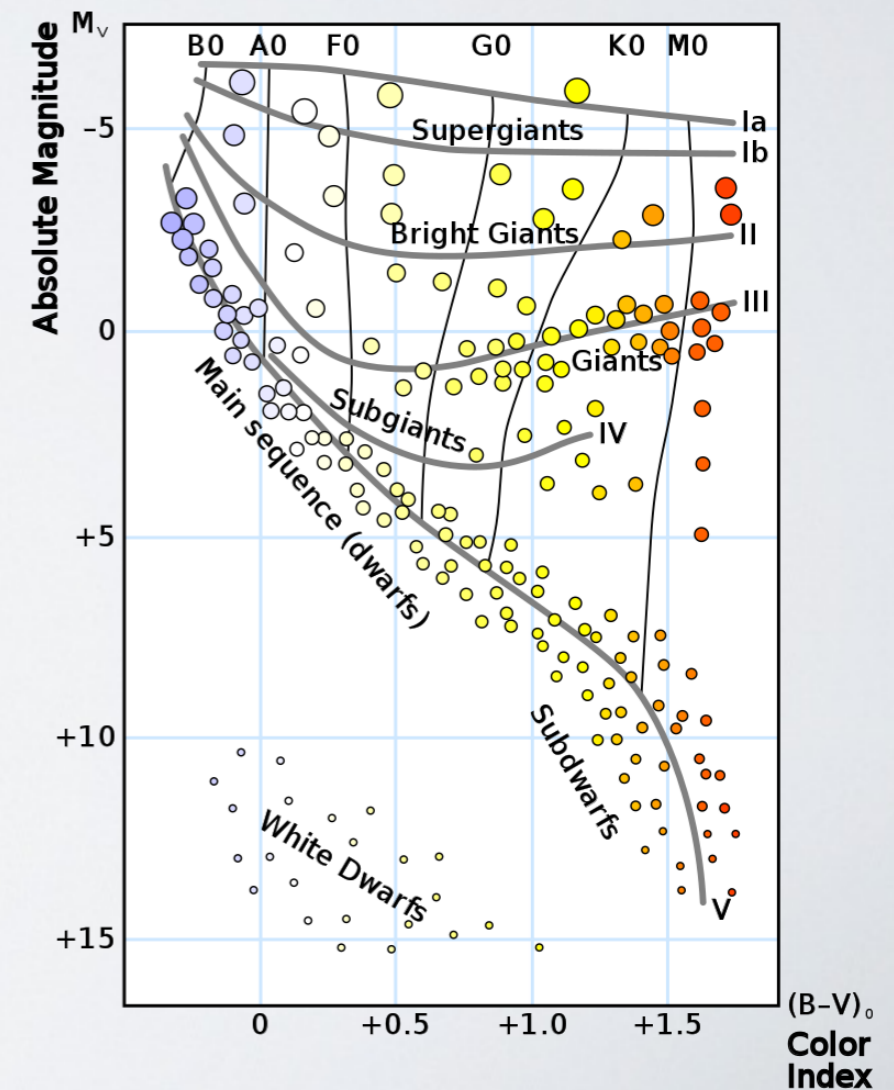
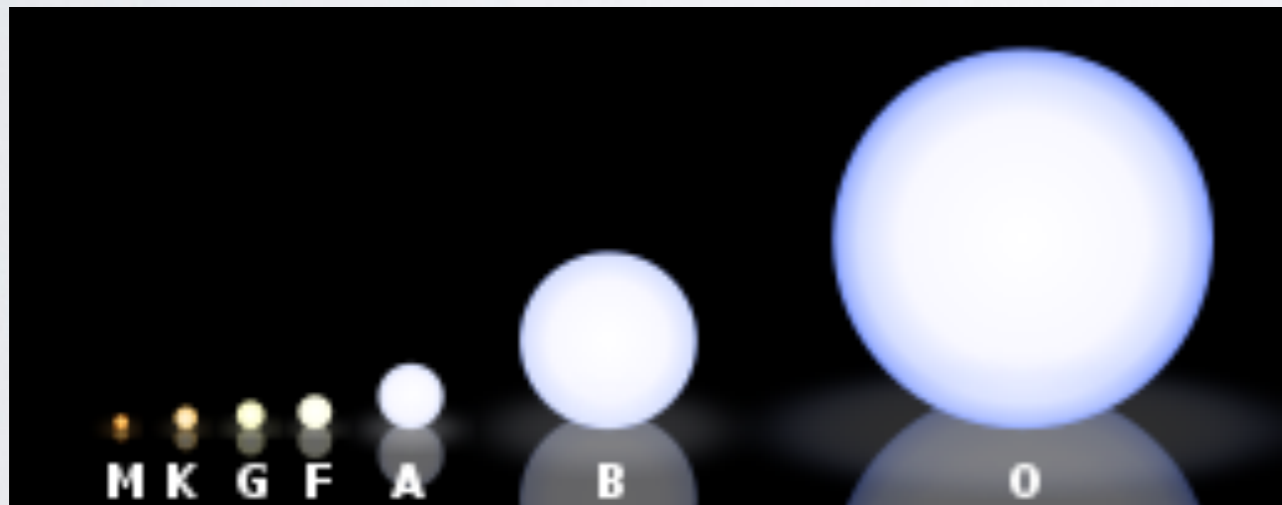
CLUSTERING EVALUATION

INTERNAL/EXTERNAL

- Two types of evaluation: internal or external
- External Evaluation (extrinsic):
 - ▶ Similarly to supervised learning, compares the clusters found with a “ground truth”
 - ▶ The ground truth can be exactly the right clustering desired
 - So we are just validating the method, since we already know the answer...
 - ▶ The ground truth can be a proxy to what we want
 - e.g., we want to cluster stars based on their characteristics (size, temperature, color...). We already have a manual historic categorization (red dwarf, Brown dwarfs, Red giants...). We assume that the new categories found should be somewhat similar

INTERNAL/EXTERNAL

Class	Effective temperature ^{[2][3]}	Vega-relative chromaticity ^{[4][5][a]}	Chromaticity (D65) ^{[6][7][4][b]}	Main-sequence mass ^{[2][8]} (solar masses)	Main-sequence radius ^{[2][8]} (solar radii)	Main-sequence luminosity ^{[2][8]} (bolometric)	Hydrogen lines	Fraction of all main-sequence stars ^[9]
O	≥ 30,000 K	blue	blue	≥ 16 M_{\odot}	≥ 6.6 R_{\odot}	≥ 30,000 L_{\odot}	Weak	~0.00003%
B	10,000–30,000 K	blue white	deep blue white	2.1–16 M_{\odot}	1.8–6.6 R_{\odot}	25–30,000 L_{\odot}	Medium	0.13%
A	7,500–10,000 K	white	blue white	1.4–2.1 M_{\odot}	1.4–1.8 R_{\odot}	5–25 L_{\odot}	Strong	0.6%
F	6,000–7,500 K	yellow white	white	1.04–1.4 M_{\odot}	1.15–1.4 R_{\odot}	1.5–5 L_{\odot}	Medium	3%
G	5,200–6,000 K	yellow	yellowish white	0.8–1.04 M_{\odot}	0.96–1.15 R_{\odot}	0.6–1.5 L_{\odot}	Weak	7.6%
K	3,700–5,200 K	light orange	pale yellow orange	0.45–0.8 M_{\odot}	0.7–0.96 R_{\odot}	0.08–0.6 L_{\odot}	Very weak	12.1%
M	2,400–3,700 K	orange red	light orange red	0.08–0.45 M_{\odot}	≤ 0.7 R_{\odot}	≤ 0.08 L_{\odot}	Very weak	76.45%



INTERNAL/EXTERNAL

- Two types of evaluation: internal or external
- Internal Evaluation (Intrinsic):
 - ▶ We have no ground truth to compare to
 - ▶ We evaluate the intrinsic properties of our clusters, typically
 - If their elements are similar
 - If clusters are far apart /if elements in different clusters are different.

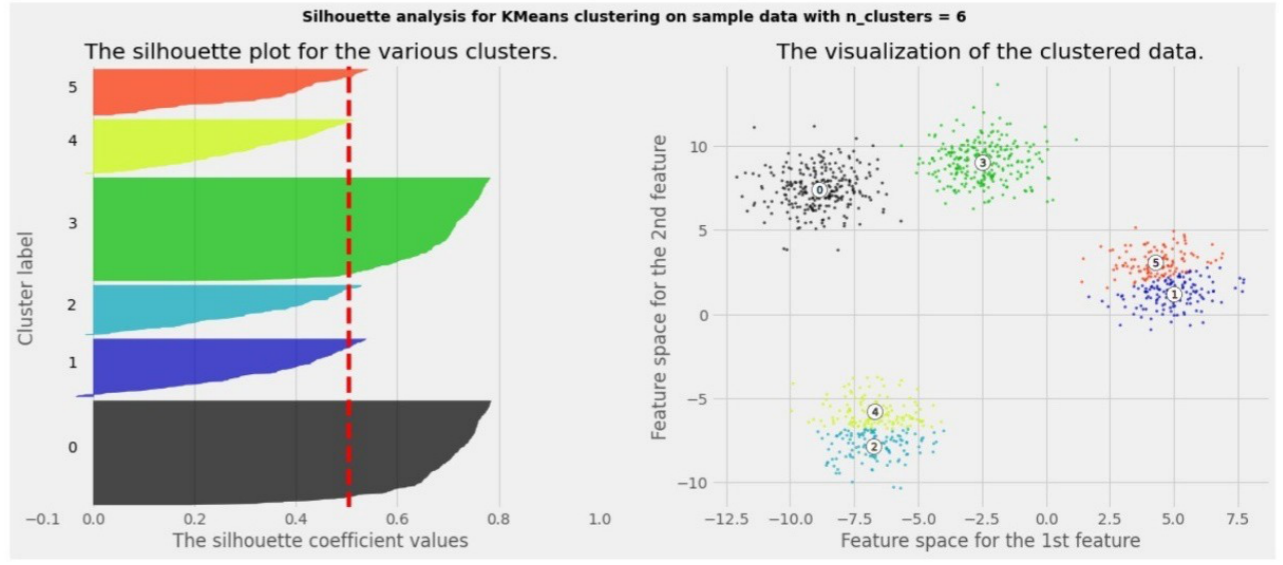
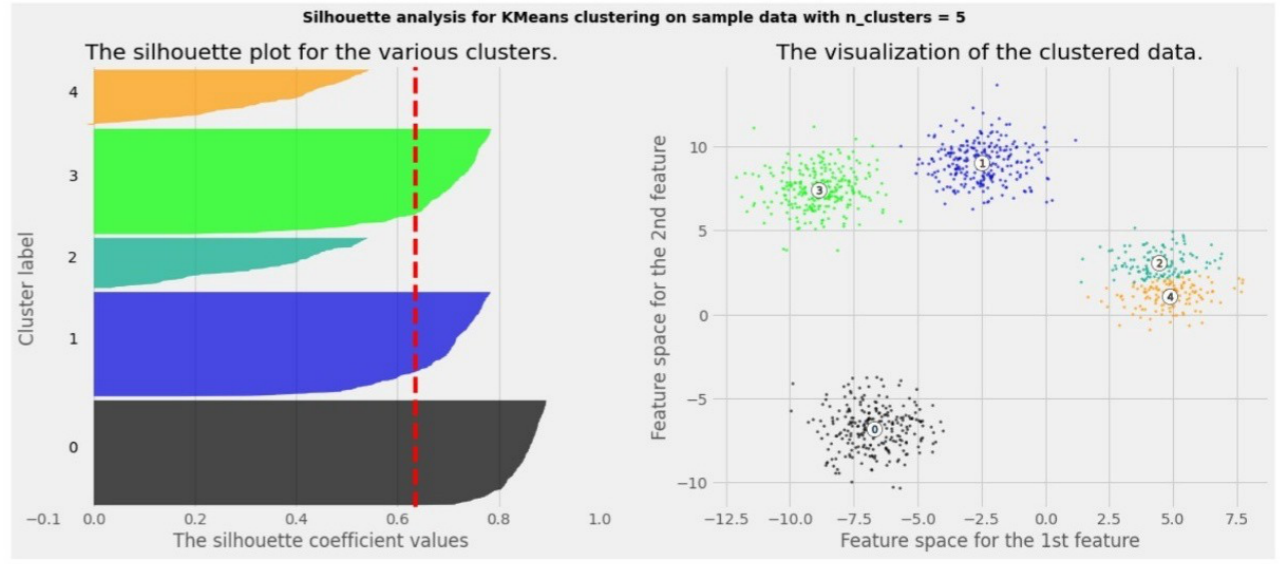
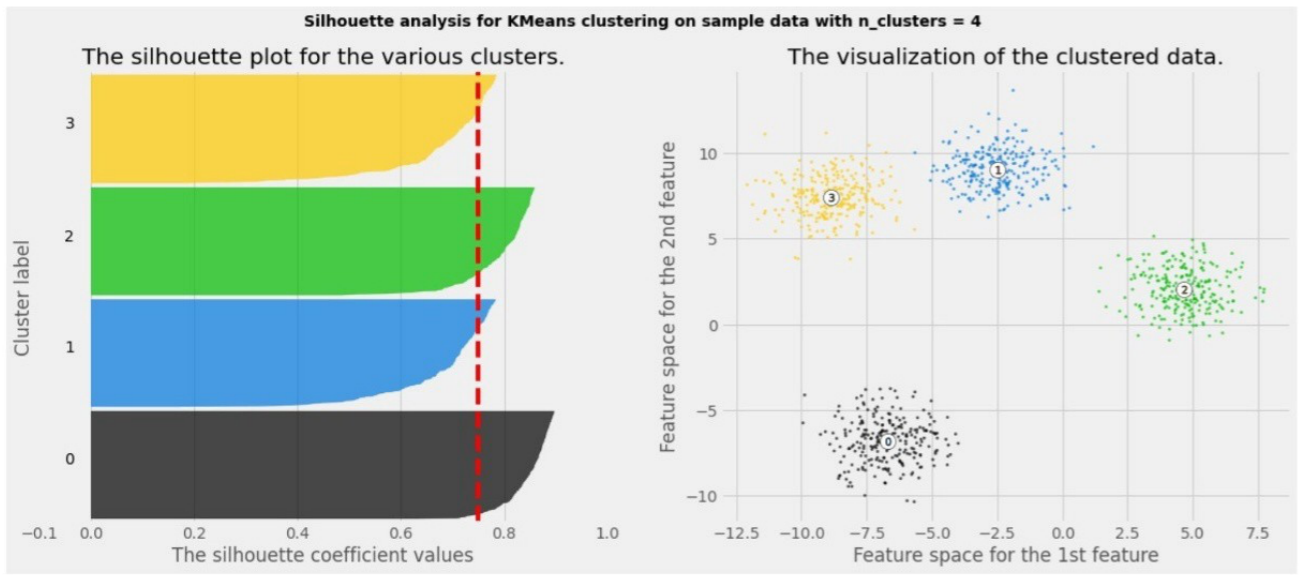
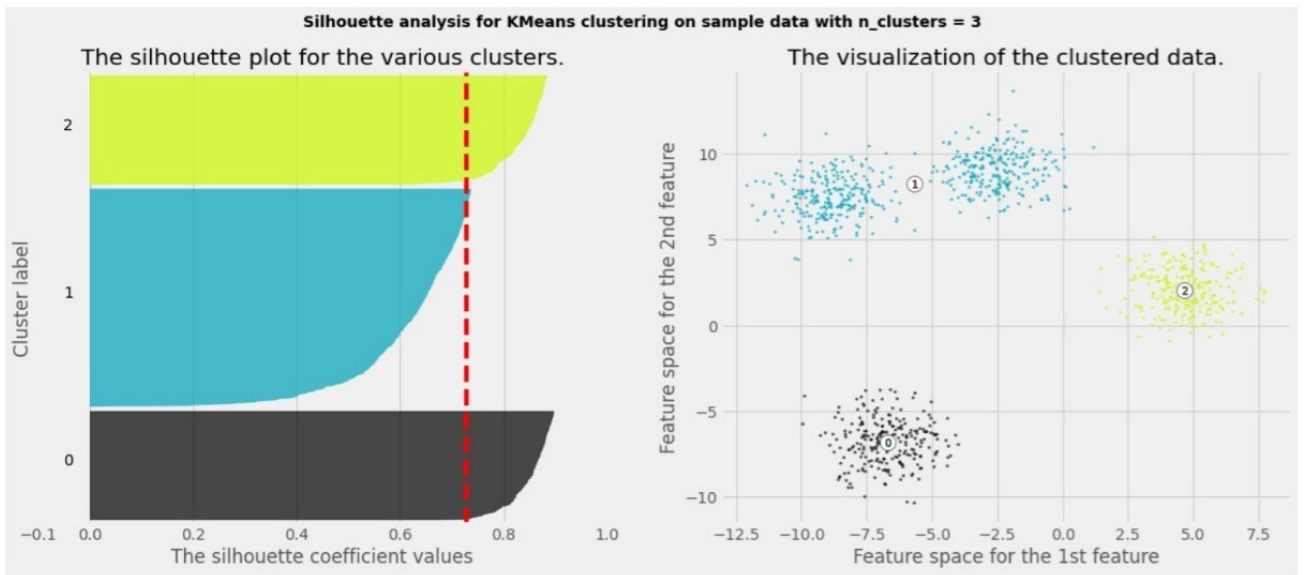
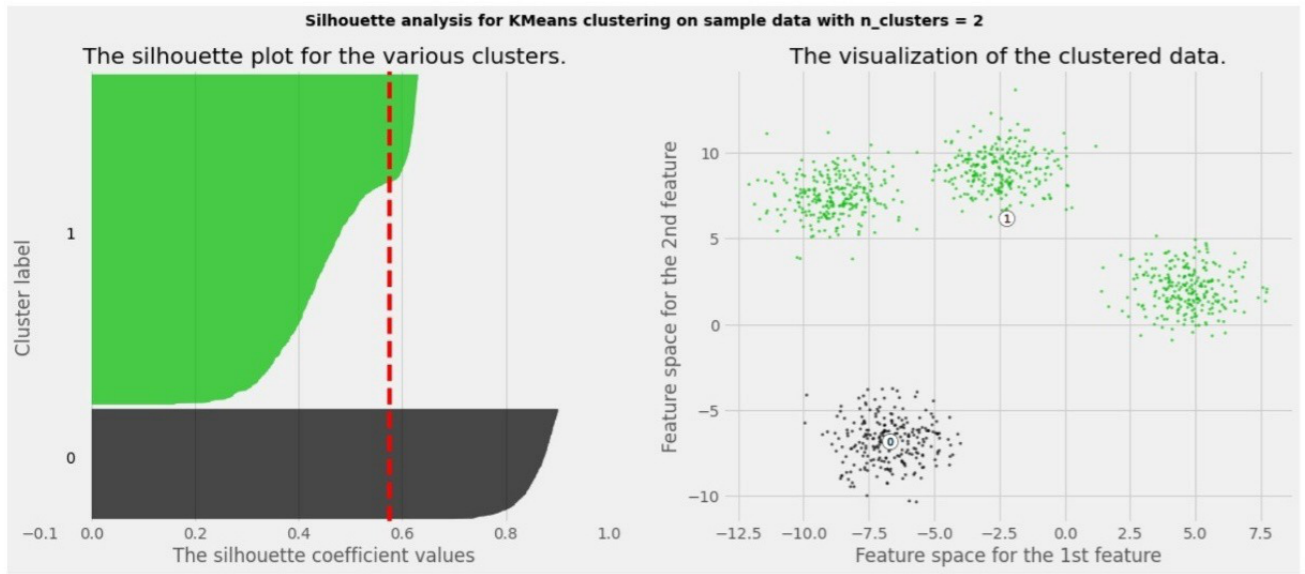
INTERNAL EVALUATION

AD-HOC SCORES

- Several clustering methods define their own objective to minimize. This objective can be used as a score for clusters obtained by this method or others
 - k-means minimizes inter-cluster variance
 - Gaussian mixture maximizes likelihood
- But can lead to unfair comparison:
 - Using inter-cluster variance to compare k-means and another method such as DBscan is unfair.
 - One explicitly minimizes this objective, the other no...
- As always, the choice of a score is equivalent to choosing a definition of cluster...

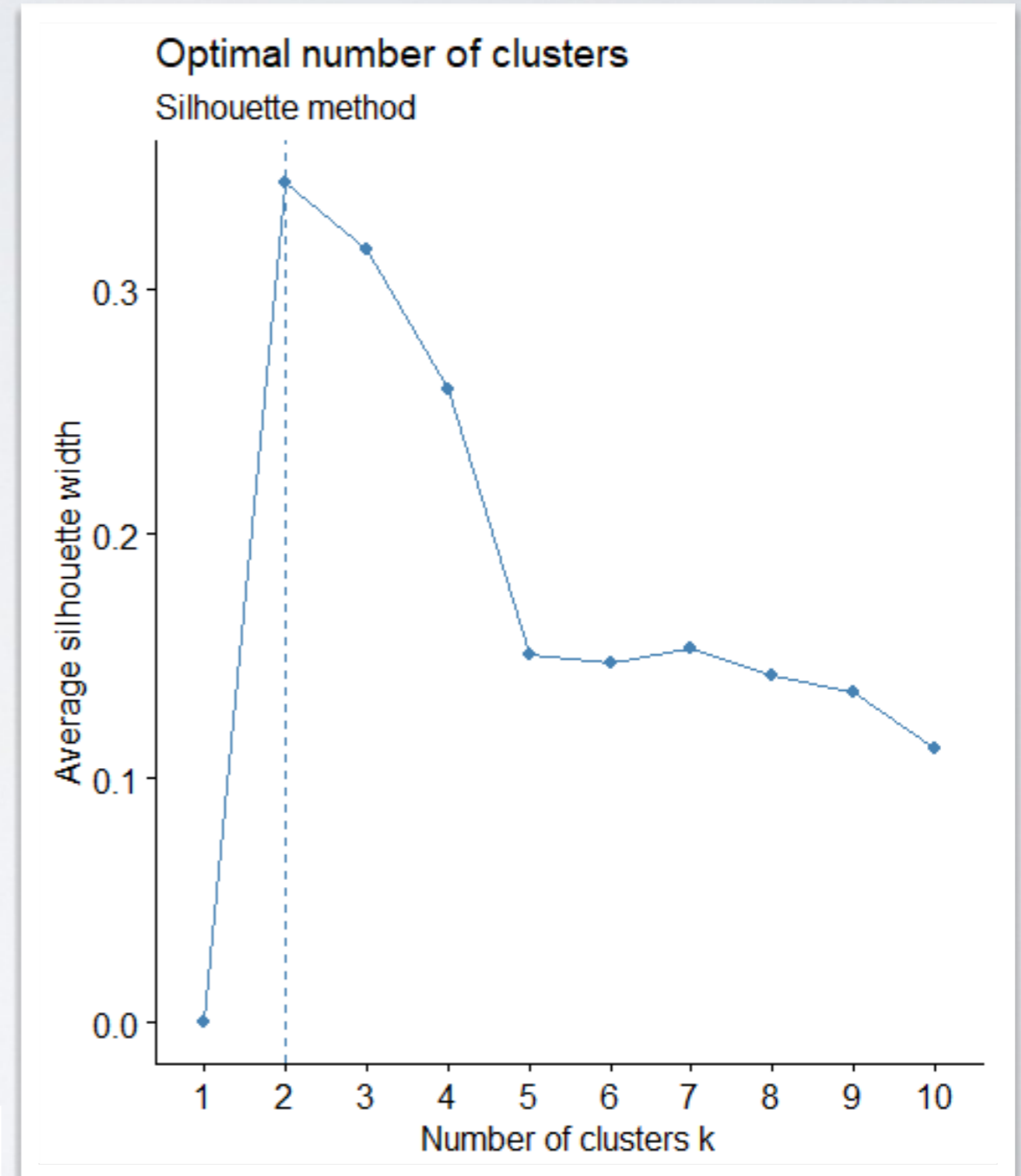
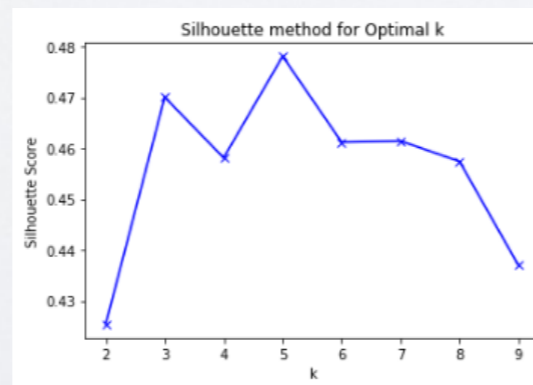
SILHOUETTE SCORE

- Intrinsic score
- Silhouette score of 1 item:
 - 1) Compute $a(i)$, average distance to all other points of the same cluster
 - 2) Compute $b(i)$, min average distance to all points of another cluster
 - 3) Silhouette: $s(i) = \begin{cases} 1 - a(i)/b(i), & \text{if } a(i) < b(i) \\ 0, & \text{if } a(i) = b(i) \\ b(i)/a(i) - 1, & \text{if } a(i) > b(i) \end{cases}$
- Silhouette coefficient:
 - Average of all individual Silhouette scores.



AUTOMATIC K SELECTION

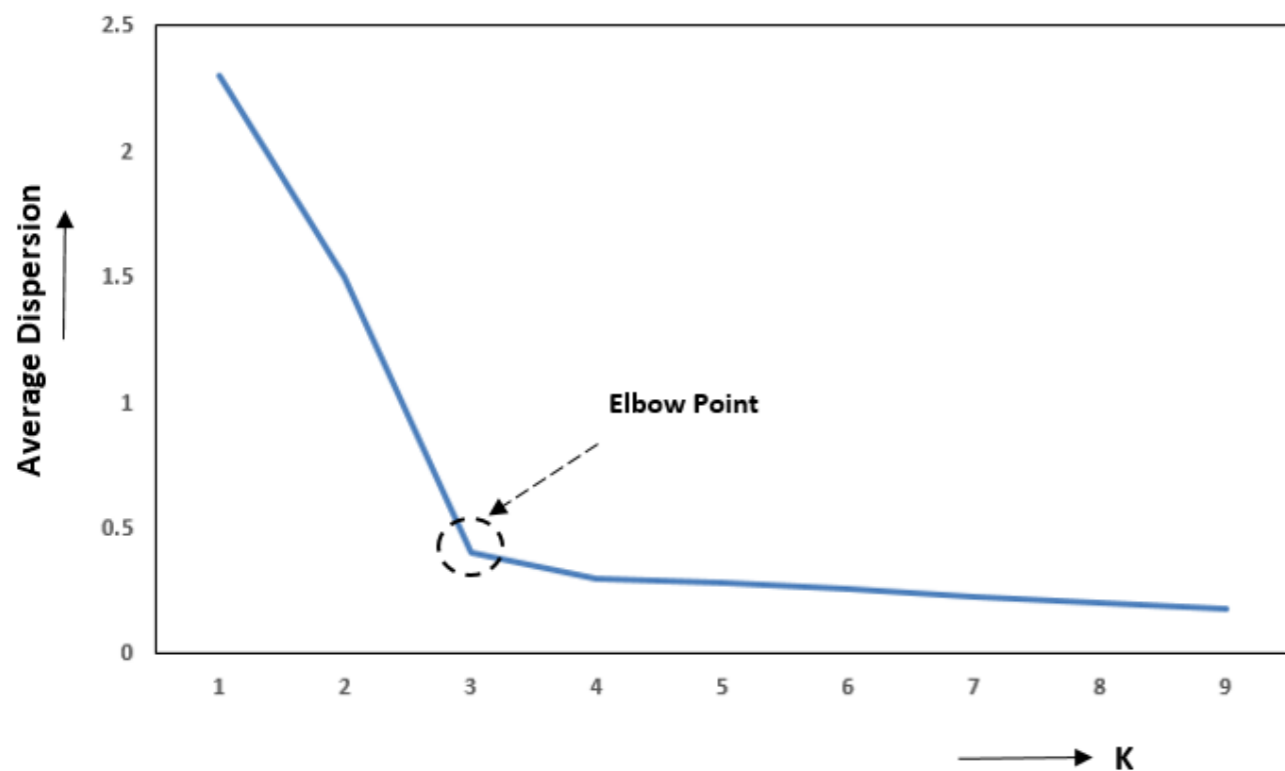
- The Silhouette score can be used to choose automatically the number of clusters:
 - We vary the number of clusters k , and search for the maximum



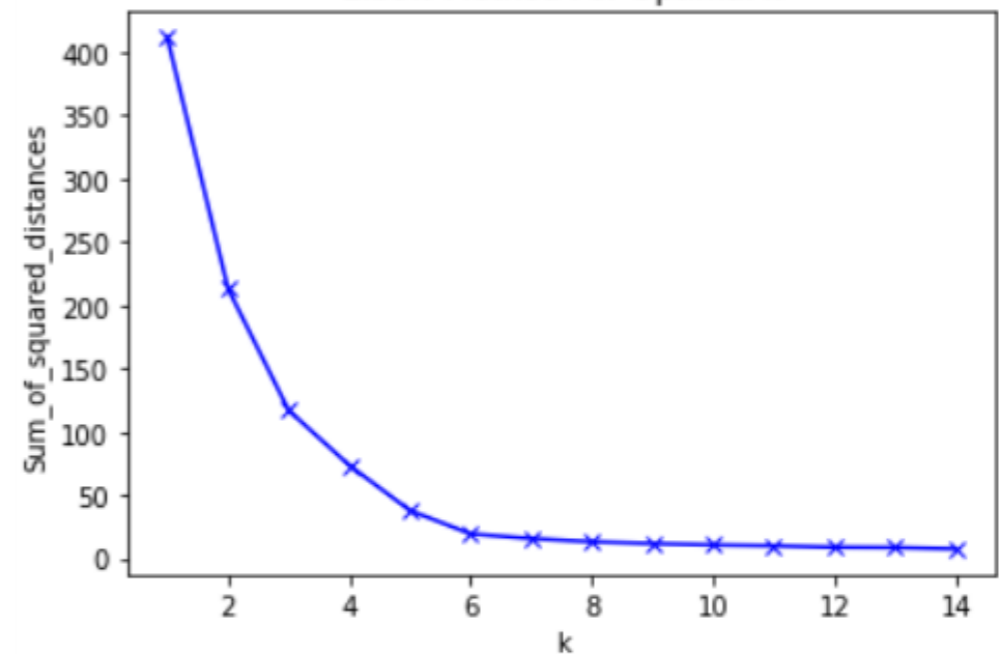
VARIANT: ELBOW METHOD

- Another well known method to find automatically the number of clusters consists in plotting a measure of quality such as the inter-cluster variance, and cut at an “elbow”
 - Diminishing returns=> less “worthy” to continue

Elbow Method for selection of optimal “K” clusters



Elbow Method For Optimal k



EXTERNAL EVALUATION

EXTERNAL EVALUATION

- Rand Index = Accuracy over node pairs

$$\bullet RI = \frac{TP + TN}{TP + FP + FN + TN}$$

- ▶ TP: two nodes in same cluster in both GT and solution
 - ▶ TN: two nodes in different clusters in both GT and solution
 - ▶ $TP + FP + FN + TN =$ all possible node pairs
- Problem: complexity. #of pairs = $\mathcal{O}(n^2)$
 - ▶ 100k items: 10 Billion pairs...

RAND INDEX

- Fast computation: Contingency table
 - Same as a confusion matrix, but not necessarily square

	S1	S2	S3	S4	S5	S6	S7	S8	S9	NA
K1	95									
K2		18								
K3			17							
K4				31	1	1	6		11	6
K5					45	1				1
K6					3	21	12	8	6	4
K7					44	1			3	1

$$\binom{x}{2} = \frac{x(x-1)}{2}$$

RAND INDEX

	x1	x2	x3		Sums
y1	1	1	0		2
y2	1	2	0		3
y3	0	0	4		4
Sums	2	3	4		

- TP: $\sum_{ij} \binom{n_{ij}}{2} = 3 \binom{1}{2} + \binom{2}{2} + \binom{4}{2} = 6 + 1 = 7$

- TN: $\binom{n}{2} - \sum_i \binom{n_i}{2} - \sum_j \binom{n_j}{2} + TP = 36 - (1+3+6) - (1+3+6) + 7 = 23$

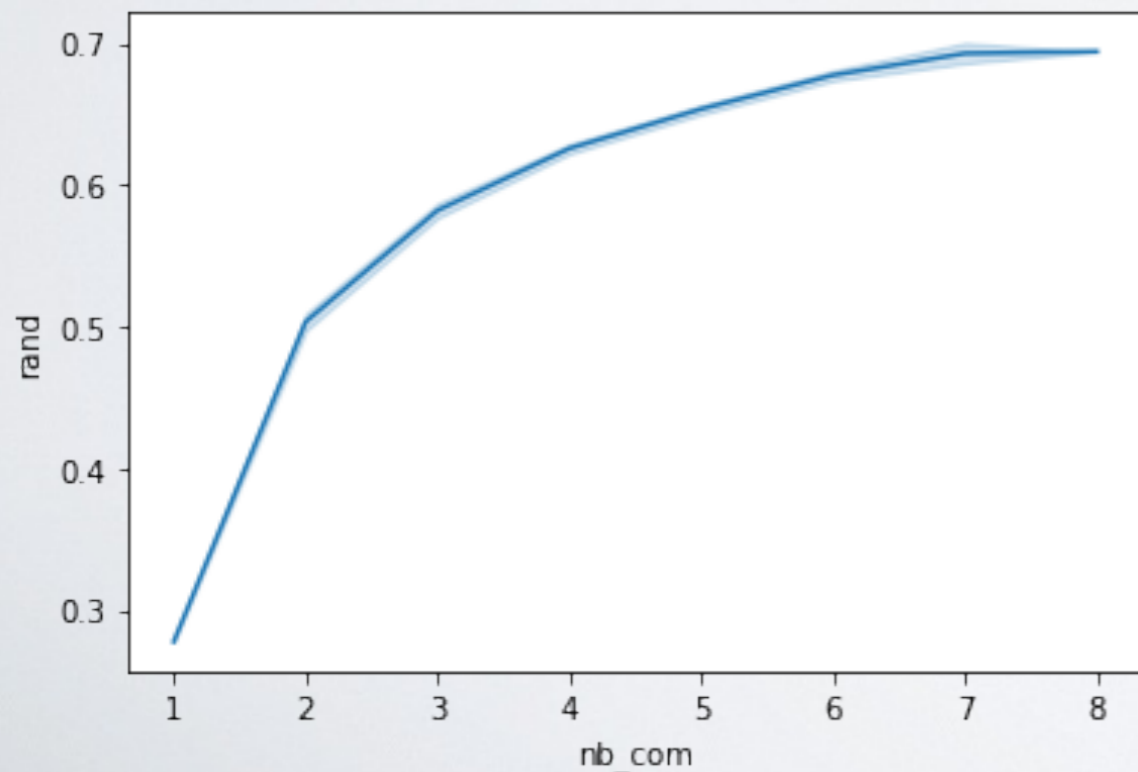
- All pairs - pairs inside clusters in our solution ($\sum_i \binom{n_i}{2}$) = N(all negatives) in our solution.

- Remove those that are positive in GT ($\sum_j \binom{n_j}{2}$) But retain among them those in common in our solution too (TP)

- Rand = $(23+7)/36 = 0.833$

RAND INDEX

- Rand Index has the same weakness as Accuracy:
 - If the classes are imbalanced, i.e., the size and number of communities vary between GT and clustering, results can be counterintuitive
- In practice:
 - Random communities have different scores depending on their size
 - => Prefer certain types of communities



Previous example,
With random communities

ARI

- Solution: Use an **adjusted for chance** score.
 - Principle: adjust (normalize) such as 0 is the score obtained with a “random” solution, 1 is the highest possible score.
 - Negative solutions are worse than random

$$\frac{\text{Index} - \text{Expected index}}{\text{Max index} - \text{Expected index}}$$

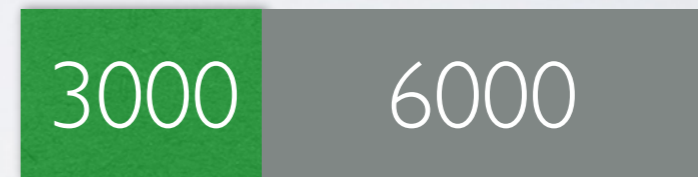
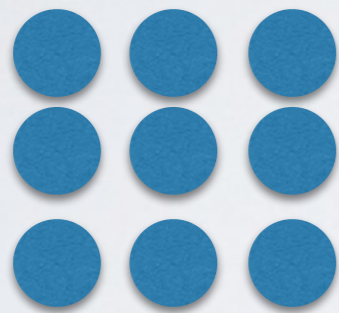
Expected nb of TP

$$E(n_{ij}) = \binom{n_i}{2} \frac{\binom{n_j}{2}}{\binom{n}{2}}$$

ARI

- Simple interpretation,

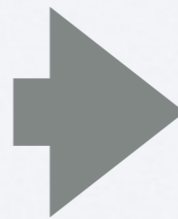
- ▶ Forget about clusters: we had X blue chocolates, we distribute at random in a box having a total of Y free spaces. The green part of the box has Z free spaces. How many blue chocolates do you expect in the green part?



Chocolate/
pairs

	x1	x2	x3	Sums
y1	1	1	0	2
y2	1	2	0	3
y3	0	0	4	4
Sums	2	3	4	

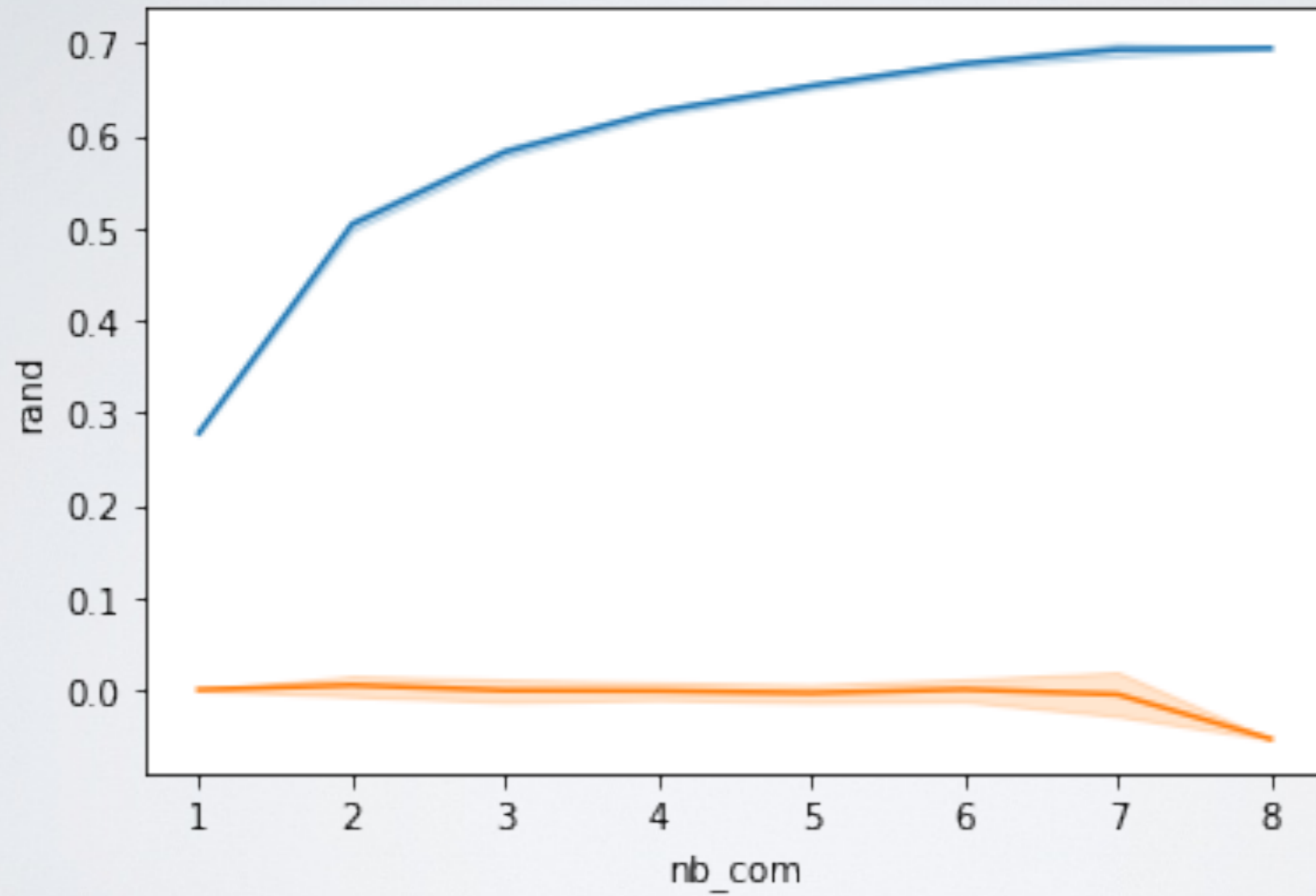
Space



	x1	x2	x3	Sums
y1	1/36	3/36	6/36	2
y2	3/36	9/36	18/36	3
y3	6/36	18/36	36/36	4
Sums	2	3	4	

$= 100/36 = 2.77$ expected TP

ARI



HOMOGENEITY/ COMPLETENESS

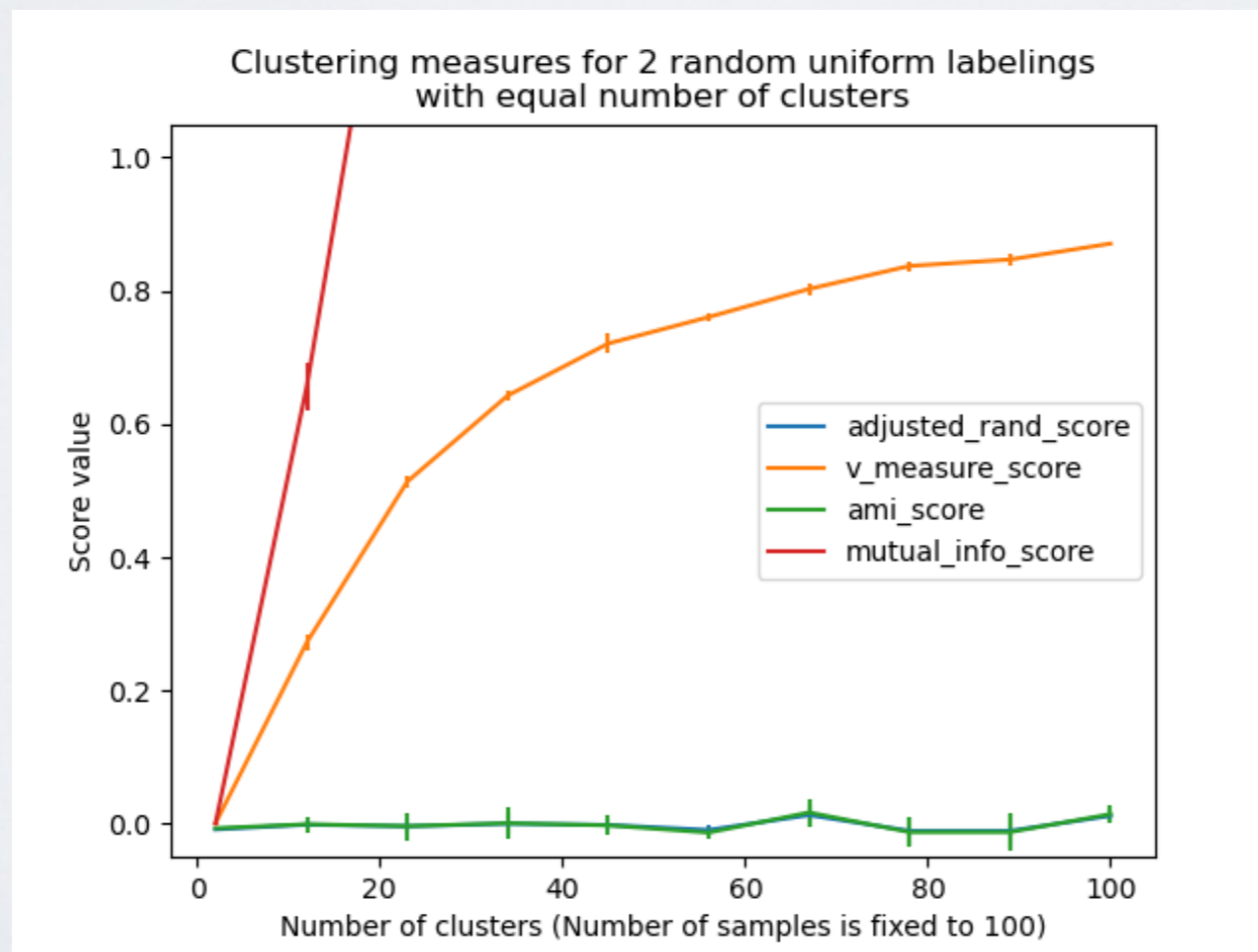
- Similar idea than Precision and Recall
 - Scores based on Information Theory
- Homogeneity
 - Items inside a same cluster belong to the same class (GT cluster)
 - $\in [0,1]$, 1: each cluster contains only items of the same class
 - But several clusters can be composed of the same class...
- Completeness
 - Items belonging to the same class (GT cluster) are found in the same cluster
 - $\in [0,1]$, 1: All items of a same class are found in the same clusters
 - But a same cluster might contain different classes...
- Trivial solutions: Each in its own cluster, a single cluster...

V-SCORE / NMI

- As with f1-score, Homogeneity and Completeness can be combined in a single score, using harmonic mean
 - Called the v-score.
 - $$v = \frac{2 \times \text{homogeneity} \times \text{completeness}}{\text{homogeneity} + \text{completeness}}$$
- The exact same score (defined differently, still from Information Theory) is also called NMI
 - NMI: Normalized Mutual Information

AMI

- NMI/v-score is known to suffer from the exact same problem as Rand Index
 - An adjusted for chance version exist, called AMI (adjusted Mutual Information)



NO FREE LUNCH THEOREM

- “Any two optimization algorithms are equivalent when their performance is averaged across all possible problems”
 - ▶ Two clustering algorithms with different objective functions are fully comparable, one is not intrinsically better than another.
 - ▶ Each is the best for the objective function it defines
 - ▶ What is “the best” cluster? Depends on your definition.
- Does not mean that some methods are not more appropriate than other for what most people consider as clusters...