

Weekly meeting 4

Dr. Doina Bein

Friday, June 16, 10:30am-12pm

Surveys to be completed

To be done today, before starting research:

CIC-PCUBED Pre-event survey:

https://fullerton.qualtrics.com/jfe/form/SV_6YIVSkC6hLxbunA

Project 1: Data Science

What you need to do: topics & objectives

Objective 1: Learn Python using some textbook or some online courses such as (<https://www.codecademy.com/learn/learn-python>). Shared by Stephanie Pocchi: Learn Python in a couple hours. This YouTuber does a very beginner-friendly crash course about the capabilities of Python and its uses. Here is the link:

<https://www.youtube.com/watch?v=rfscVS0vtbw>

Objective 2: Learn how to use Jupyter Notebook. Start here http://jupyter-notebook-beginner-guide.readthedocs.io/en/latest/what_is_jupyter.html

Objective 3: For data science, find a suitable dataset and start training some neural network using with Google tensorflow.

Logistics for all students

- Who is participating: [list of current research students](#) and their availability
- Research will be conducted virtually during the week with in-person meetings throughout the week
- Zoom meetings for me to teach new topics and for you to participate in open discussions
- Support:
 - If needed, you can meet me
Zoom: Mon, Tu, Wed from 8:30-10:25 am
IN PERSON: Mon, Tu, Wed from 8:30-9:30 am, Thursday 8:30-10am or by email
 - CIC-PCUBED peer mentor: (tentative) [availability](#)

Logistics for all students (contd.)

- Make a copy of this GDoc [Work schedule](#), share the Gdoc copy with me, and maintain it weekly and daily; due at the end of Week 2
- Before the end of week 3, make a copy and maintain your [Proposed work](#) by individual or teams of up to three; due by the end of Week 3
- Complete your [availability here](#); try to have it consistent over the 7 weeks such that it will be easy to partner in the project
- Group projects: to be decided; sample list [here](#)
- Oral or poster presentations: tentatively scheduled for Friday, July 28, from 8:30am-12:30 pm and if needed, from 1:30-4 pm

Please checkout:

- [Other websites and ebooks](#)
- [Websites with free datasets](#)
- If you find good, free resources, please share it by email or during weekly meetings
- Next meeting: I will lecture on ZOOM on Data Science: Thursday, June 22, from 10:30am-12pm

Progress on Learning Python

- Free course: <https://www.codecademy.com/learn/learn-python>
- Free course: <https://www.kaggle.com/learn/python>
- Youtube video (about 4 hours):
<https://www.youtube.com/watch?v=rfscVS0vtbw>

Data Science

k-means Clustering

- The *k-means cost function* asks to minimize the sum of squared Euclidean distances of data points to their closest prototype centers:

$$e_k(X; C) = \sum_{i=1}^n \min_{j \in \{1, \dots, k\}} \|x_i - c_j\|^2$$

- *k*-means cost function seeks compact globular clusters of small variances: the cost of a single cluster $e_1(G) = e_1(G, c)$ is minimized when we choose for the cluster prototype its center of mass c , called the *centroid*:

$$c(G) = \operatorname{argmin}_c \sum_{x \in G} \|x - c\|^2 = \frac{1}{|G|} \sum_{x \in G} x$$

where $|G|$ denotes the cardinality of G , that is the number of elements contained in group G and $\operatorname{argmin}_x f(x)$ to denote the argument that yields the minimum in case this minimum value is unique

- If instead of choosing the squared Euclidean distance, we had chosen the ordinary Euclidean distance, one obtains the so-called *Fermat-Weber point* that generalizes the notion of median. It is thus also called the *geometric median*.
- Although the Fermat-Weber point is unique and often used in operations research for facility location problems, it does not admit a closed-form solution, but can be arbitrarily finely approximated
- *k-median clustering* is the clustering obtained by minimizing the cost function

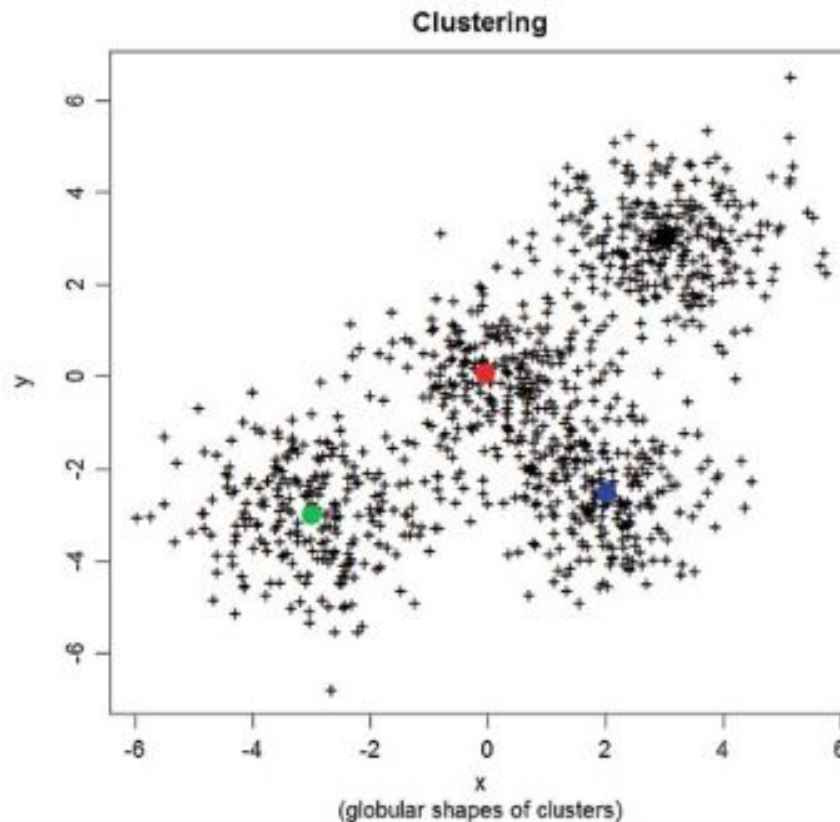
$$\min_C \sum_{i=1}^n \min_{j \in \{1, \dots, k\}} \|x_i - c_j\|$$

where $\|\cdots\|$ means the regular Euclidean distance

- Partitions from *k-means* and *k-medians* can be very different from each other: the centroid location can be different to the median for a single cluster

- Centroids can be easily corrupted by adding a single outlier point
- We say that the breakdown point of the centroid is 0: A single outlier p_0 diverging to ∞ will impact the centroid to be diverging to ∞ too.
- But the median is more robust since it requires $\left\lfloor \frac{n}{2} \right\rfloor$ outliers (that is, about 50% of outliers) to steer the median point to ∞ .
- Therefore k-median clustering is often preferred when there are many outliers in data-sets.

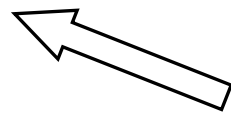
Fig. 7.3 The k -means cost function tend to find globular-shaped clusters that minimize the weighted sum of the cluster variances. k -Means clustering is a model-based clustering where each cluster is associated to a prototype: its center of mass, or centroid. Here, we have choosen $k = 4$ groups for the k -means: Cluster prototypes, centroids, are illustrated with large disks



The k-means Problem

(<https://www.math.uwaterloo.ca/~cswamy/talks/kmeans-short.ppt>)

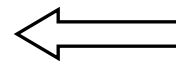
Given: n points in d -dimensional space



$X \subseteq \mathbb{R}^d$: point set with $|X| = n$

- partition X into k clusters X_1, \dots, X_k
- assign each point in X_i to a common **center**

$c_i \in \mathbb{R}^d$

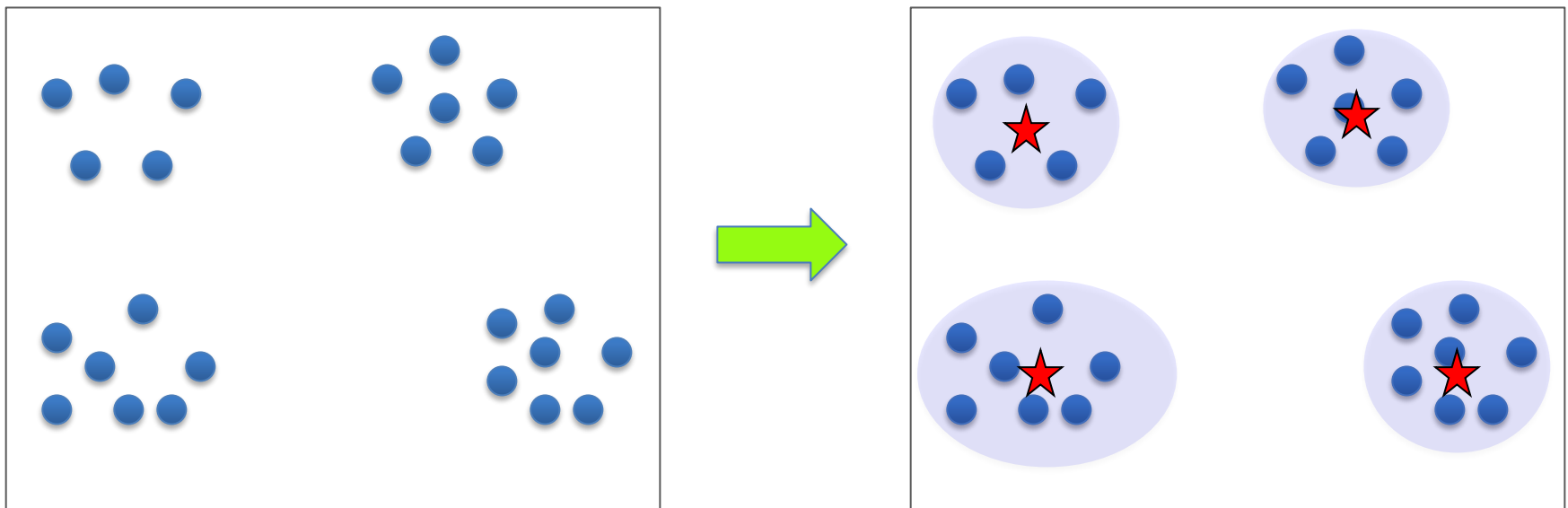


d : L_2 distance

Goal: Minimize $\sum_i \sum_{x \in X_i} d(x, c_i)^2$

(web.stanford.edu/group/mmds/slides2012/s-speaker-ppt-2012-temp/Bahmani.pptx)

- K-means clustering is a fundamental problem in data analysis and machine learning
- “By far the most popular clustering algorithm used in scientific and industrial applications” [Berkhin '02]
- Identified as one of the top 10 algorithms in data mining [Wu et al '07]



$K = 4$

k -Means Optimization Problem

- Finding the center of a single cluster is a particular case of clustering with $k = 1$ cluster.
- With the squared Euclidean distance cost, the center is the mean of the attributes, hence its naming *k-means*.
- Finding the minimum of a k -means cost function is a NP-hard problem when the dimension $d > 1$ and the number of clusters $k > 1$.
- When $k = 1$, we have shown that we can compute the optimal solution (the centroid) in linear time (computing the mean of the group).
- When $d = 1$, we can compute an optimal k -means solution using dynamic programming: Using $O(nk)$ memory, we can solve the k -means for n scalar values in time $O(n^2k)$

- For NP-hard problems, we seek efficient heuristics to approximate the cost function. We distinguish two classes of such heuristics:
 1. Global heuristics that do not depend on initialization, and
 2. Local heuristics that iteratively starts from a solution (a partition) and iteratively improves this partition using “pivot rules.”

Lloyd's Batched k-Means Local Heuristic

- Lloyd's heuristic (1957): from a given input and initialization:

Input: $\{x_1, x_2, \dots, x_n\}$, each $x_i \in \mathbb{R}^d$, the number of partitions k

Initialization: Start with k cluster centers $\{c_1, c_2, \dots, c_k\}$ (typically chosen uniformly at random from data points)

iteratively repeat until convergence the following two steps:

1. Assign points to clusters: (each $x_i \in X$ is assigned to its closest cluster center) for each $x_i \in X$, let $l_i = \operatorname{argmin}_l \|x_i - c_l\|^2$, and define the k cluster groups as $G_j = \{x_i : l_i = j\}$ with $n_j = |G_j|$, the number of elements of X falling into the j th cluster.

2. Update centers (perform an Expected Maximization-type local search):

For all $j \in \{1, \dots, k\}$, update the centers to their cluster centroids :

$$c_j = \frac{1}{n} \sum_{x \in G_j} x \text{ (or the barycenters } c_j = \frac{1}{\sum_{x \in G_j} w(x)} \sum_{x \in G_j} w(x)x \text{ for}$$

weighted data-sets).

- A possible convergence criteria: cluster centers do not change anymore

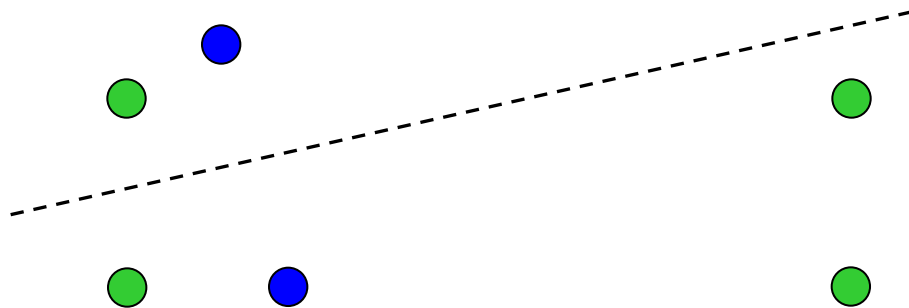
b) **Heuristics: Lloyd's method** invented in 1957 and remains an extremely popular heuristic even today

1) Start with k initial / “seed” centers c_1, \dots, c_k .

2) Iterate the following **Lloyd step**

a) Assign each point to nearest center c_i to obtain clustering X_1, \dots, X_k .

b) Update $c_i \leftarrow \text{ctr}(X_i) = \sum_{x \in X_i} x / |X_i|$.



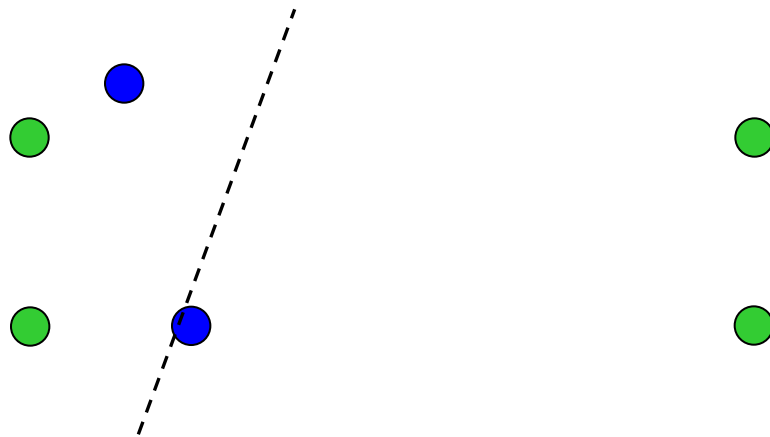
b) **Heuristics: Lloyd's method** invented in 1957 and remains an extremely popular heuristic even today

1) Start with k initial / “seed” centers c_1, \dots, c_k .

2) Iterate the following **Lloyd step**

a) Assign each point to nearest center c_i to obtain clustering X_1, \dots, X_k .

b) Update $c_i \leftarrow \text{ctr}(X_i) = \sum_{x \in X_i} x / |X_i|$.



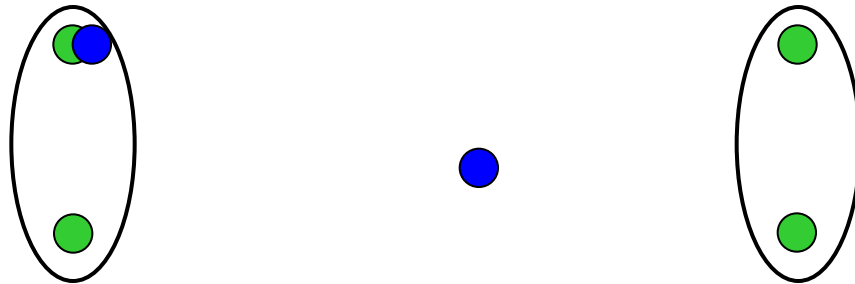
b) **Heuristics: Lloyd's method** invented in 1957 and remains an extremely popular heuristic even today

1) Start with k initial / “seed” centers c_1, \dots, c_k .

2) Iterate the following **Lloyd step**

a) Assign each point to nearest center c_i to obtain clustering X_1, \dots, X_k .

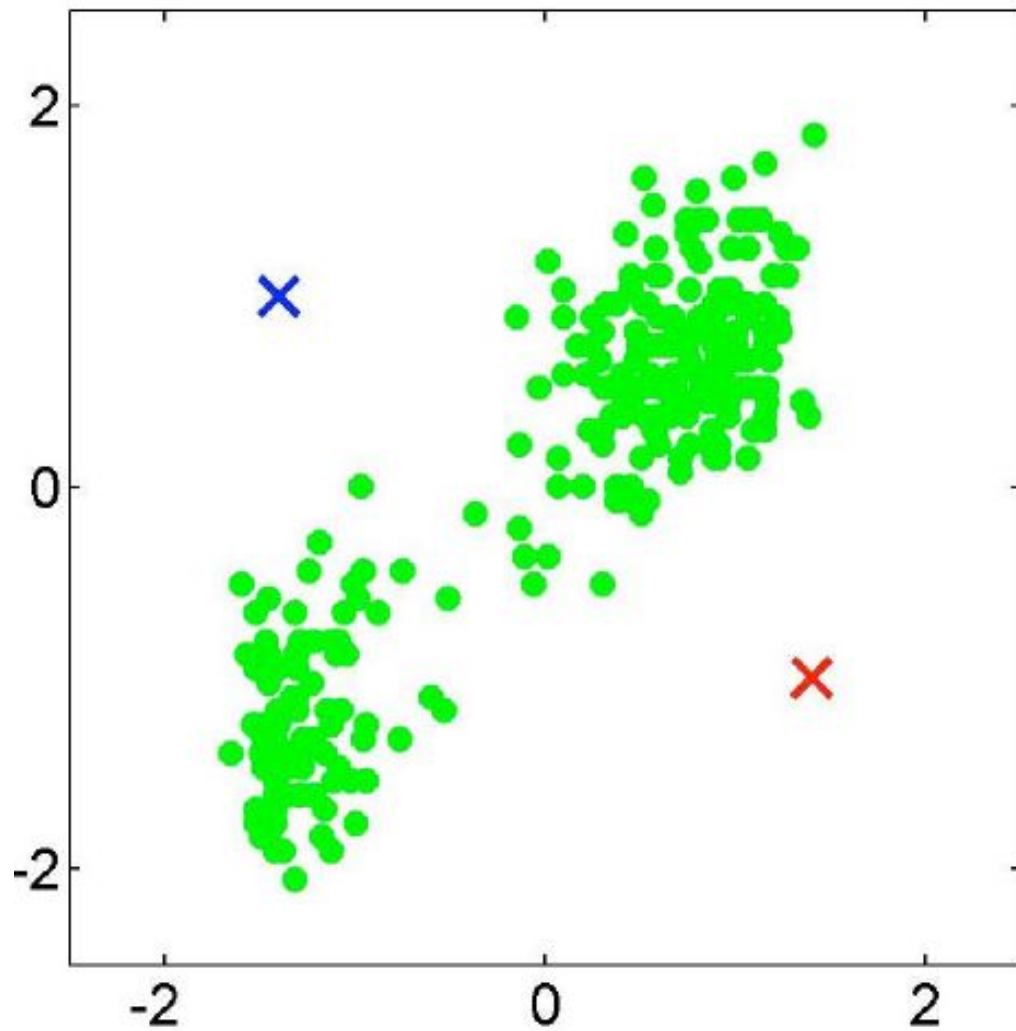
b) Update $c_i \leftarrow \text{ctr}(X_i) = \sum_{x \in X_i} x / |X_i|$.



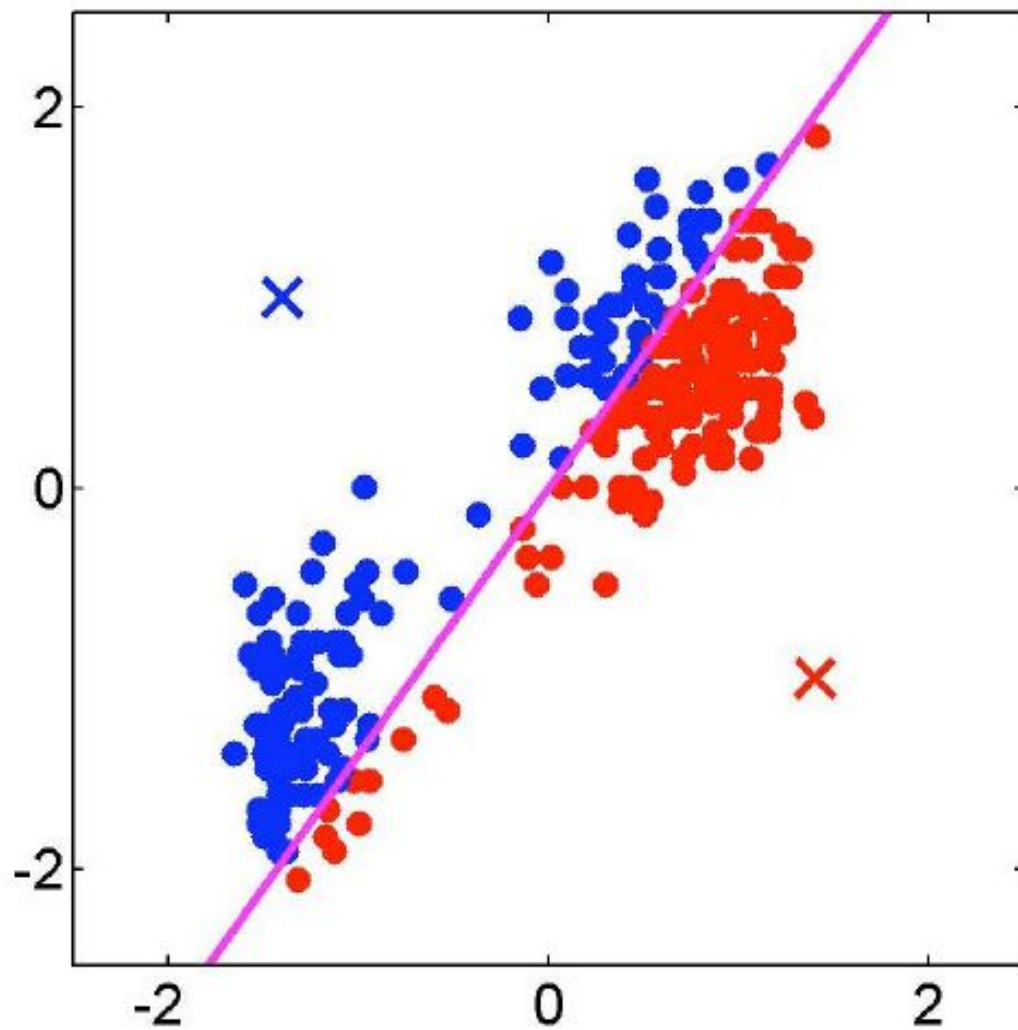
Other example

(<https://www.cs.utah.edu/~piyush/teaching/4-10-print.pdf>)

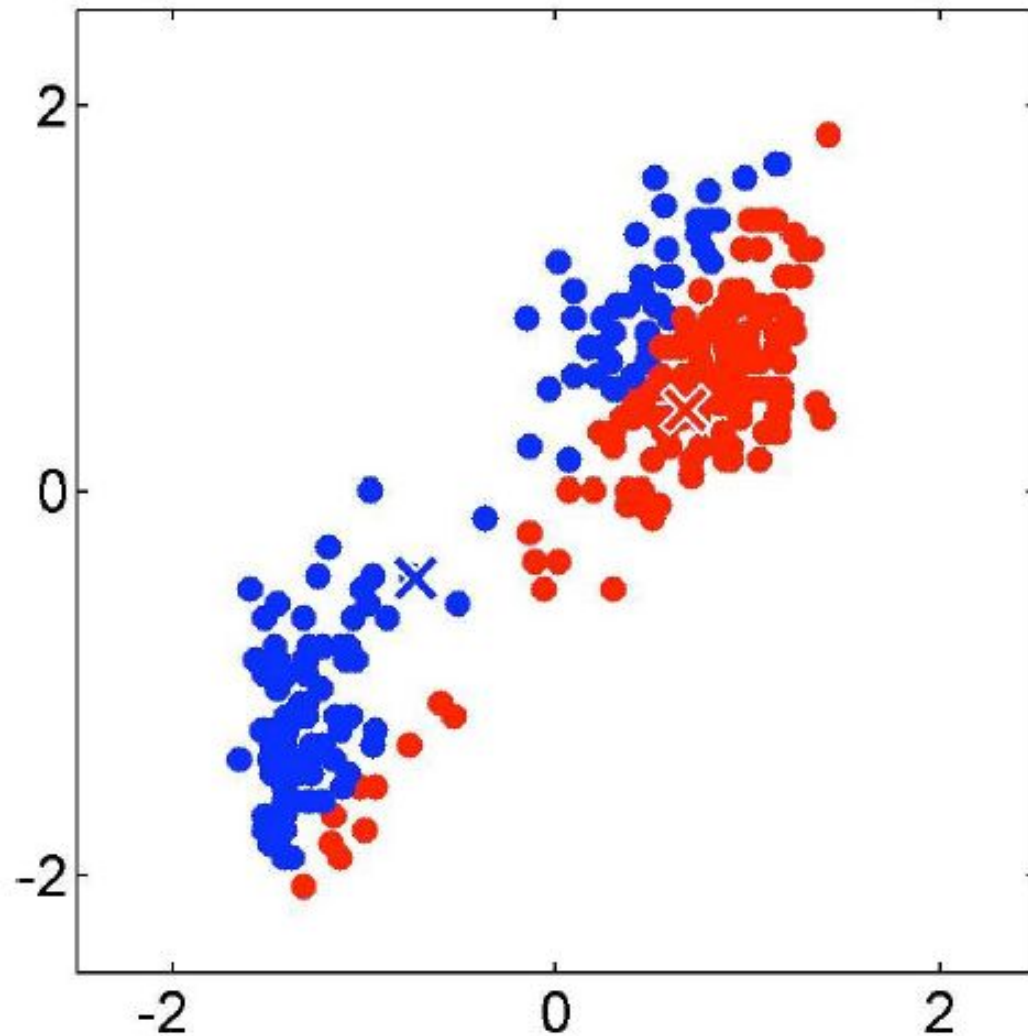
Initialization (assume $K = 2$)



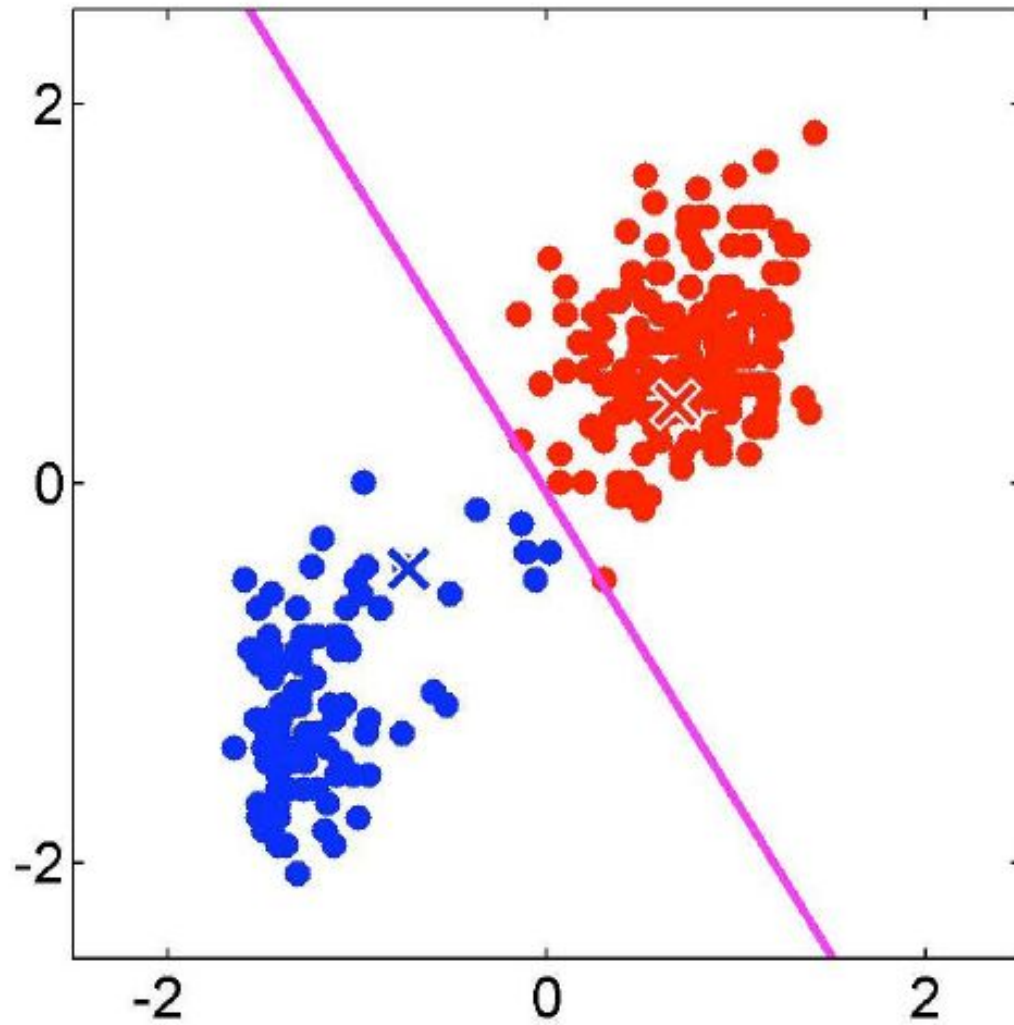
Iteration 1: Assign points to clusters



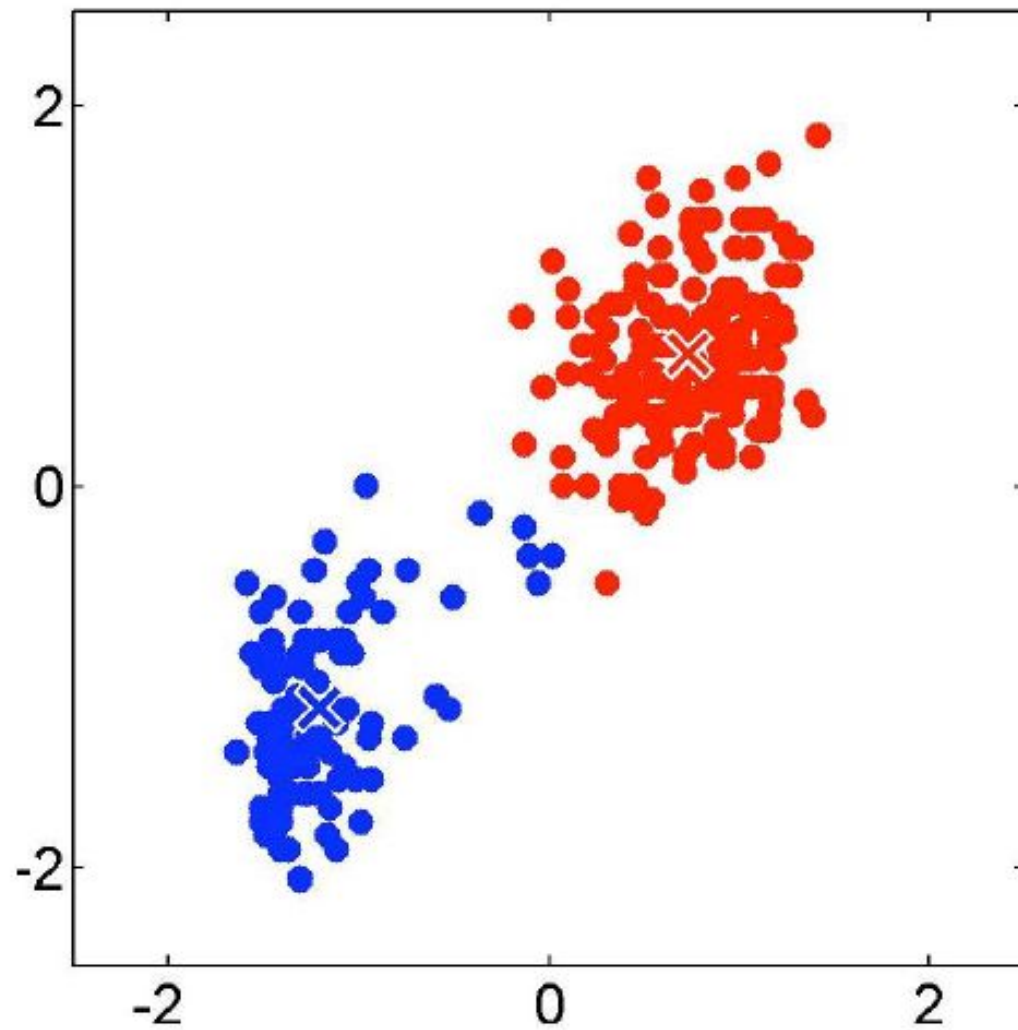
Iteration 1: Update centers



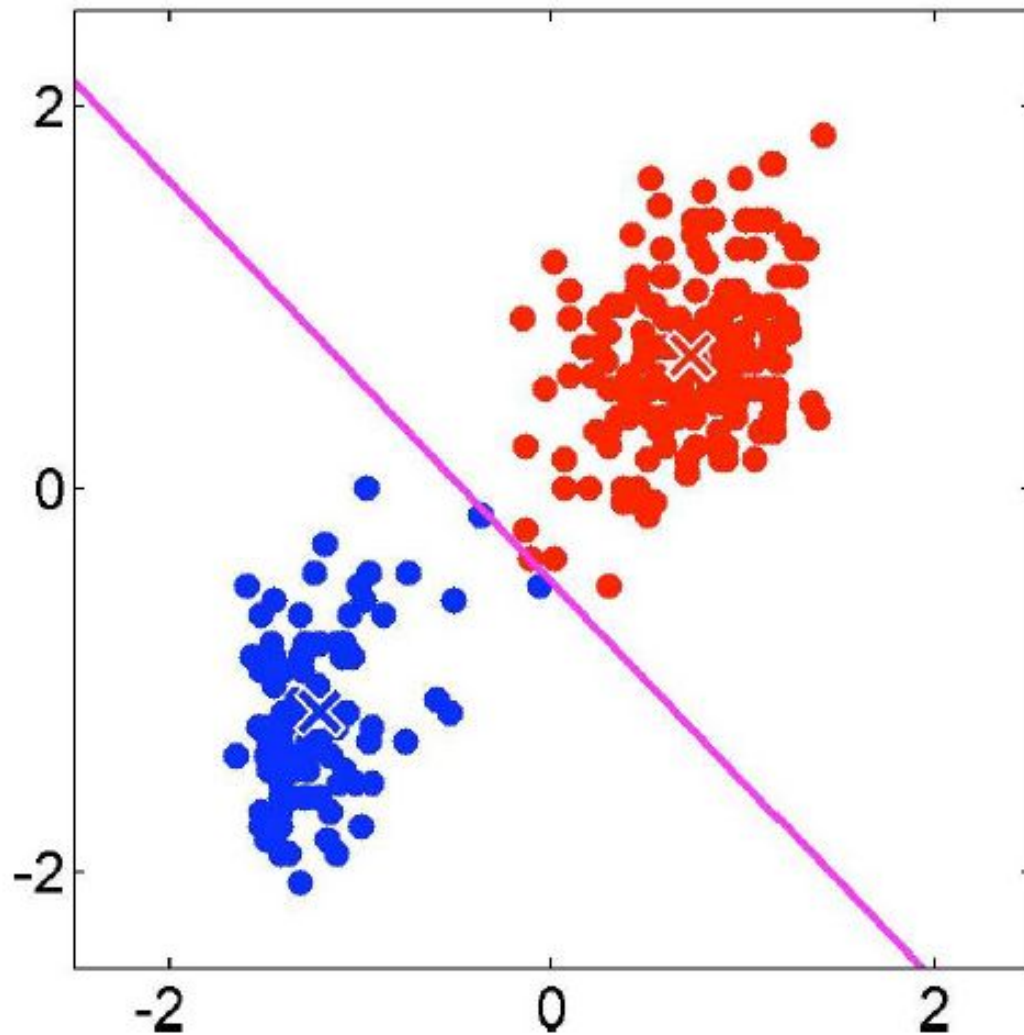
Iteration 2: Assign points to clusters



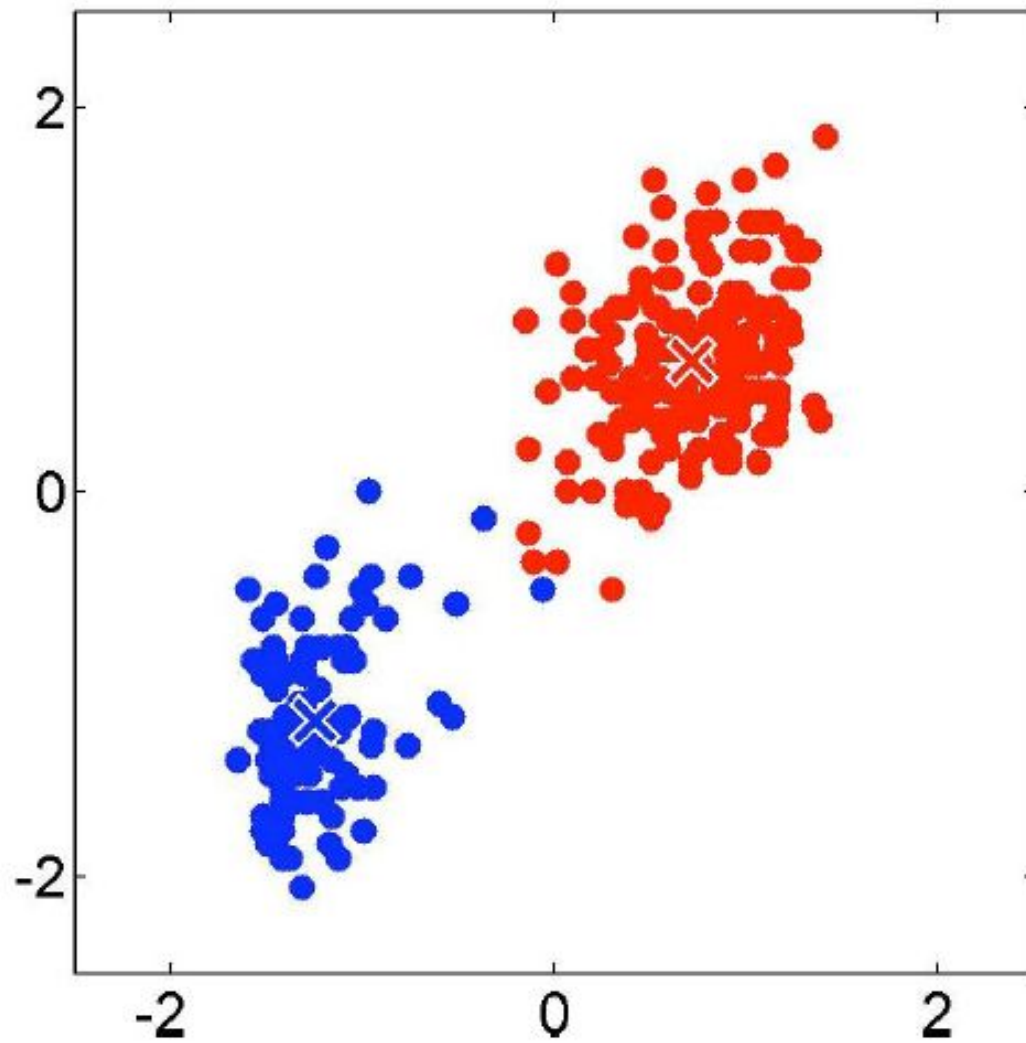
Iteration 2: Update centers



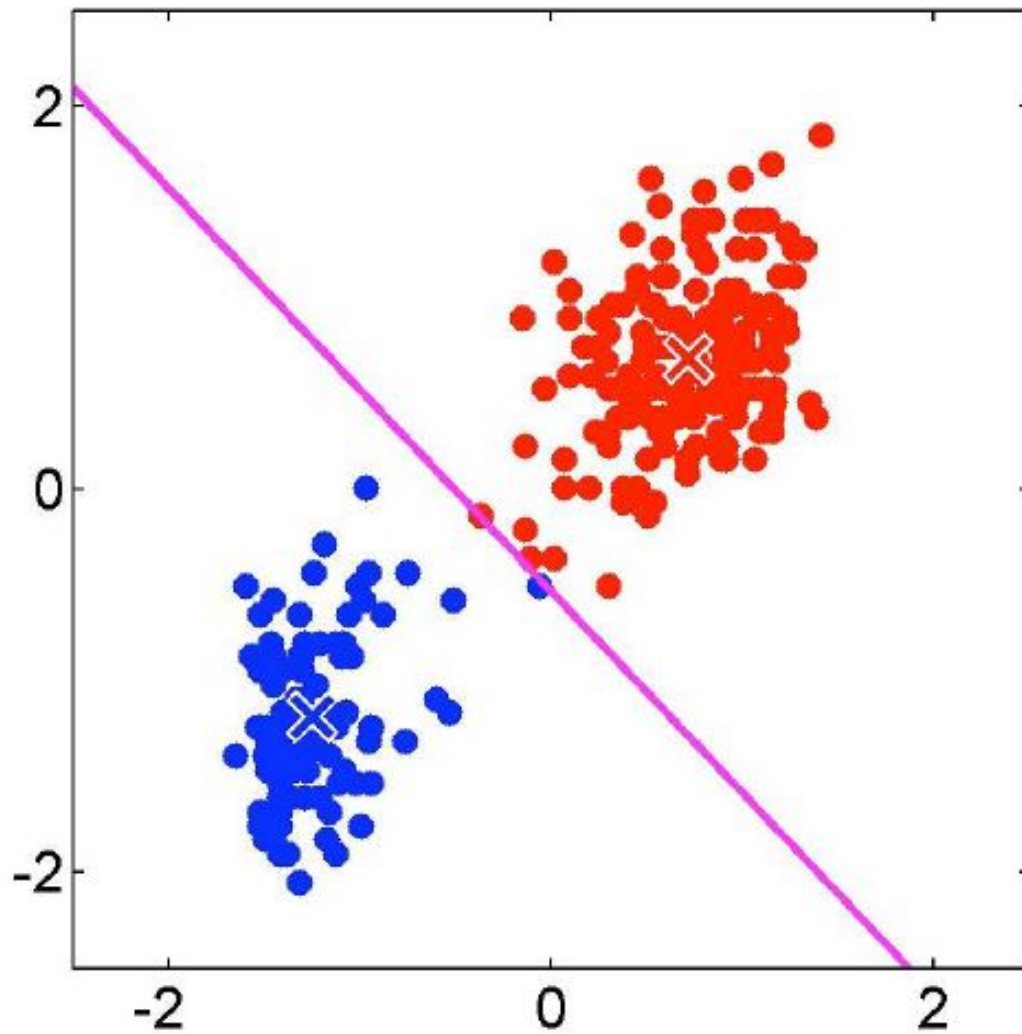
Iteration 3: Assign points to clusters



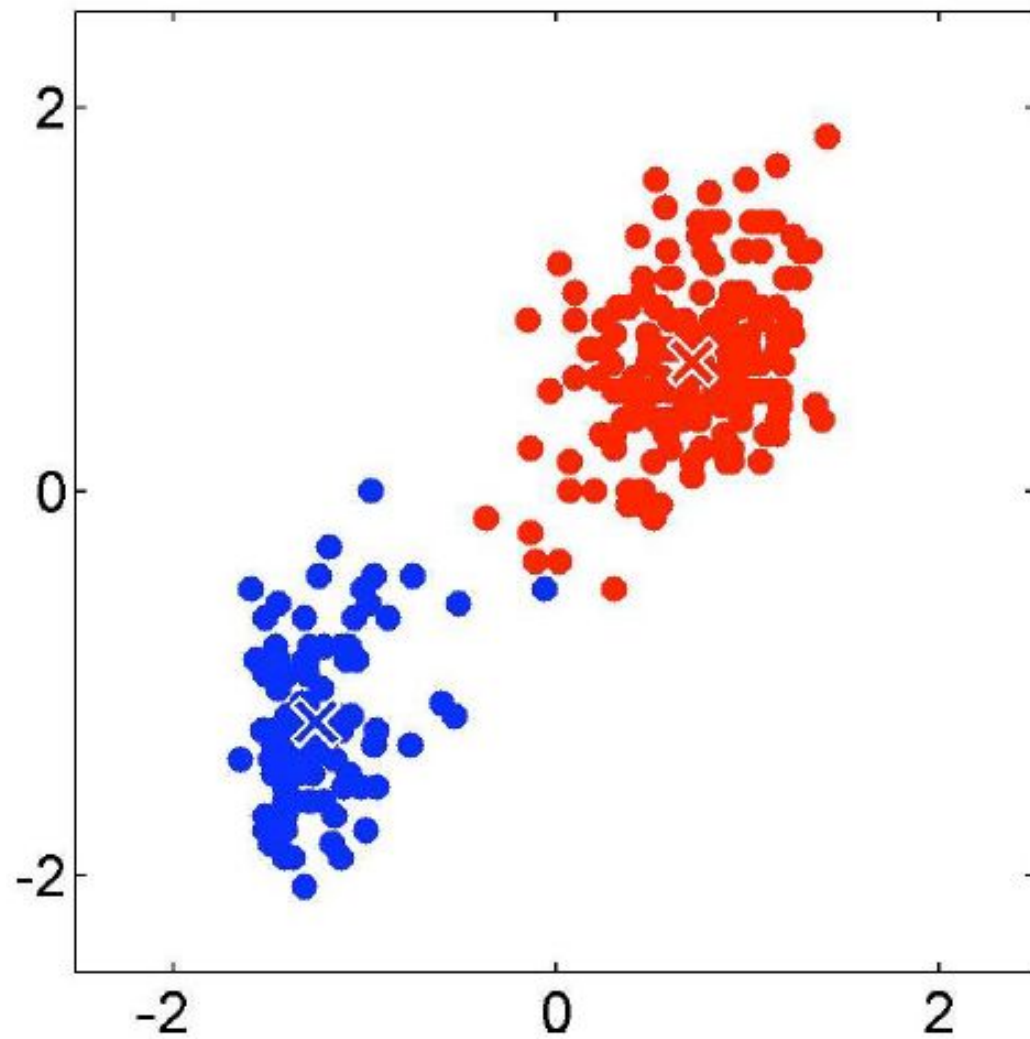
Iteration 3: Update centers



Iteration 4: Assign points to clusters



Iteration 4: Update centers



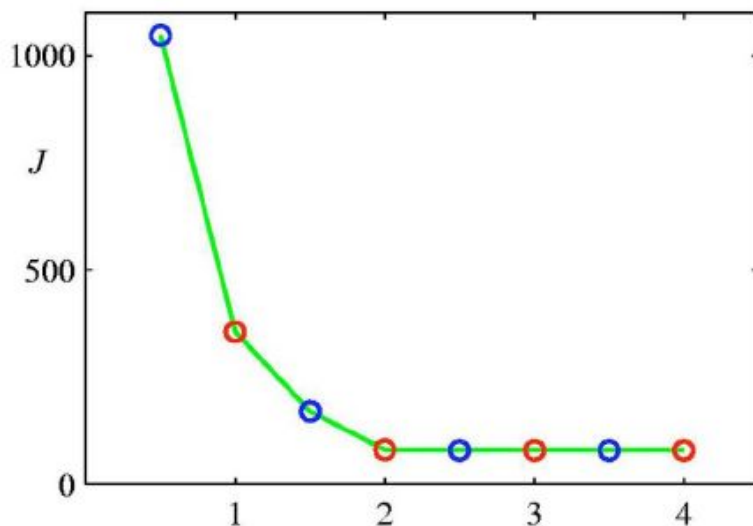
Choosing the number of clusters K

(<https://www.cs.utah.edu/~piyush/teaching/4-10-print.pdf>)

- Recall that k -means objective is to minimize the cost function, aka the sum of distances of points from their cluster centers

$$J(C, r) = \sum_{i=1}^n \sum_{j=1}^k r_{nk} \|x_i - c_j\|^2$$

- One way to select the number of clusters k is to try different values of k , plot the k -means objective versus k , and choose the “elbow-point”
- One drawback of this method is that it is computationally very expensive, and sometimes (depending on the data-sets), the inflexion point between the sharp decrease and the plateau is not so well defined



For the plot to the left, $k = 2$ is the elbow point

Initialization Issues:

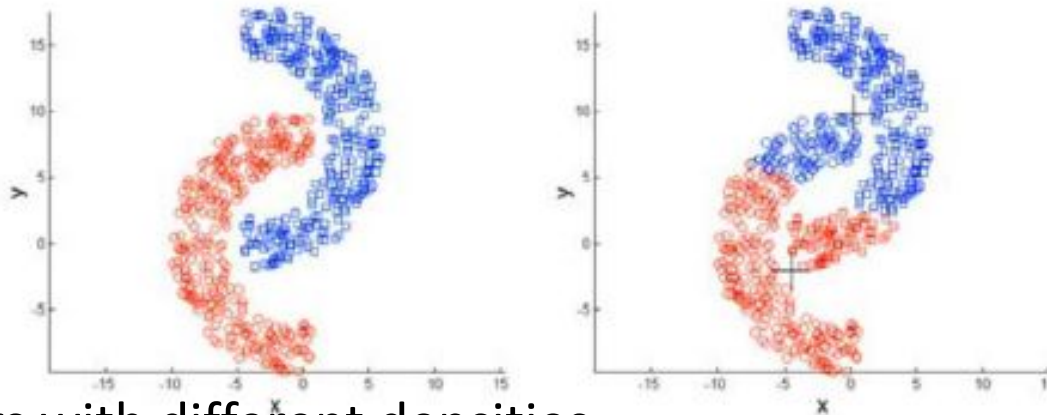
(<https://www.cs.utah.edu/~piyush/teaching/4-10-print.pdf>)

- K-means is extremely sensitive to cluster center initialization
- Bad initialization can lead to
 - Poor convergence speed
 - Bad overall clustering
- Safeguarding measures:
 - Choose first center as one of the examples, second which is the farthest from the first, third which is the farthest from both, and so on.
 - Try multiple initializations and choose the best result
- Other smarter initialization schemes (e.g., look at the K-means++ algorithm by Arthur and Vassilvitskii)

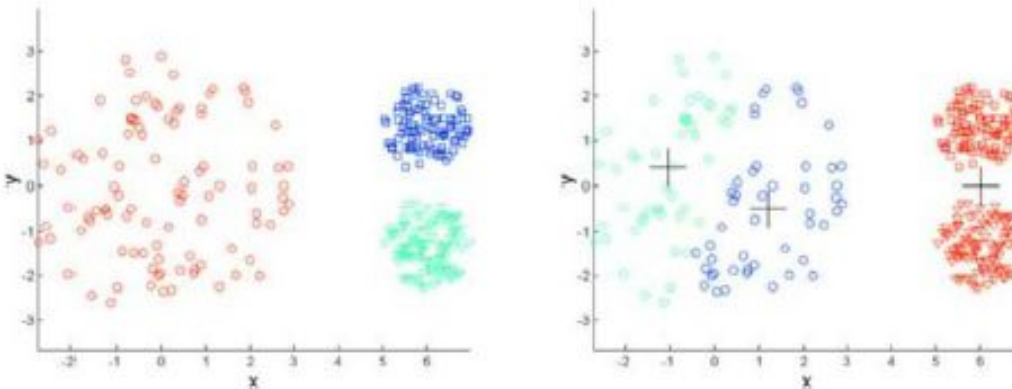
Limitations Illustrated

(<https://www.cs.utah.edu/~piyush/teaching/4-10-print.pdf>)

- K-Means might not find the best possible assignments and centers.
- Non-convex/non-round-shaped clusters: Standard K-means fails!



- Clusters with different densities



k-means++

(cs.haifa.ac.il/hagit/courses/seminars/.../Presentations/Lecture01_K-means+EM.pptx)

- *k*-Means with smart initial seeding
- Choose one center uniformly at random from among the data points.
- For each data point x , compute $D(x)$, the distance between x and the nearest center that has already been chosen.
- 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 $D(x)^2$.
- Repeat Steps 2 and 3 until k centers have been chosen.
- Now that the initial centers have been chosen, proceed using standard *k*-means.

Advantages

- This seeding method yields considerable improvement in the final error of k-means
- Takes more time to initialize
- Once initialized, K-Means converges quickly
- Usually faster than K-Means
- 1000 times less prone to error than K-Means

Computing the Similarity of Two Partitions

- We need ground-truth datasets that tell us for each data element is true cluster membership
- Without these ground-truth data sets, we could only perform a subjective or qualitative evaluation of the clustering methods.
- Human eyes can sometimes evaluate whether the obtained clustering is good or not for $d < 4$
- When ground-truth data-sets are available (say, data-sets annotated by experts), we can compute various metrics that are quantitative values that measure the similarity of two partitions: the one induced by the labels in the ground-truth data-set (assumed to be the optimal clustering by definition!) and the one reported by an automatic clustering algorithm
- The Rand index (1971) computes the similarity of two partitions

Rand Index (1971)

- computes the similarity of two partitions
- Let $G = \cup G_i$ and $G' = \cup G'_i$ cluster decomposition of the k-means heuristic and of the ground-truth data-set, respectively.
- Compare all the $n*(n-1)/2$ (n choose 2) pairs (x_i, x_j) of points, and count those that belong to the same cluster (a) from those that are found to belong to different clusters (b).
- The Rand index

$$\text{Rand}(G, G') = \frac{a + b}{\binom{n}{2}}$$

with

a: $\#\{(i, j) : l(x_i) = l(x_j) \wedge l'(x_i) = l'(x_j)\}$

b: $\#\{(i, j) : l(x_i) \neq l(x_j) \wedge l'(x_i) \neq l'(x_j)\},$

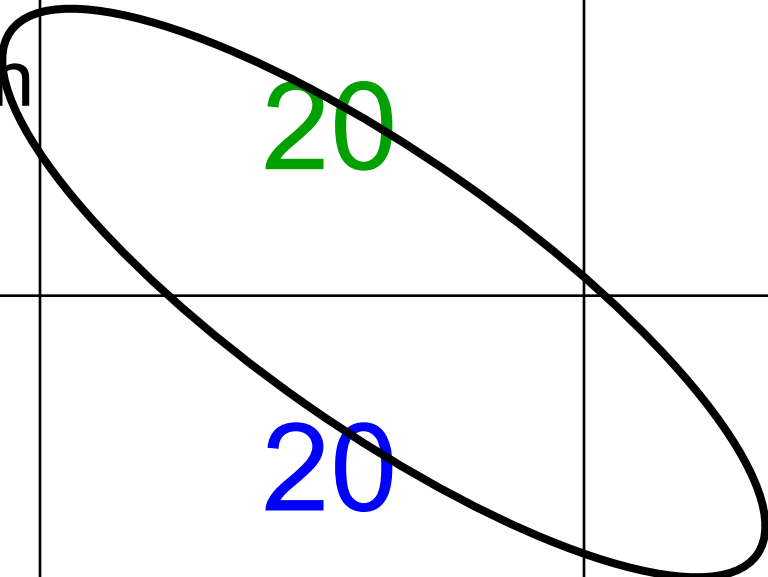
l and l' are the two cluster labeling functions of the ground-truth clustering and the automatic clustering

- Rand index belongs to the interval $[0, 1]$

Rand Index measures between pair decisions. Here $RI = 0.68$

(<https://web.stanford.edu/class/cs276/handouts/lecture14-clustering.ppt>)

| Number of point pairs | Same Cluster in clustering | Different Clusters in clustering |
|-----------------------------------|----------------------------|----------------------------------|
| Same class in ground truth | 20 | 24 |
| Different classes in ground truth | 20 | 72 |



Supervised Learning

- In supervised learning, we are given a labeled *training set* $Z = \{(x_i, y_i)\}$ with $y_i \in \pm 1$ (the ground truth labeled data) and the task is to learn a *classifier* so that we can classify new unlabeled observations of a *testing set* $Q = \{x_i\}_i$.
- When the training set has only two classes, we deal with *binary classification*, otherwise it is a multi-class classification problem.
- Statistical learning assumes that both the training set and the testing set are independently and identically sampled for an arbitrary but fixed unknown distribution.
- Our focus: learn a target function that can be used to predict the values of a discrete class attribute
- The task is commonly called: Supervised learning, classification, or inductive learning.