



CSC380: Principles of Data Science

Clustering

Credit:

- Jason Pacheco,
- Kwang-Sung Jun,
- Chicheng Zhang
- Xincheng yu

1

1

Announcements

2

- Fill out SCS (<https://scsonline.oia.arizona.edu/>) – if 80% responses, will add 5 points to the homework with lowest grade (33% right now).
- The final project due date is Friday, Dec 8.
- No lecture next Thursday, Dec 7

2

Announcements

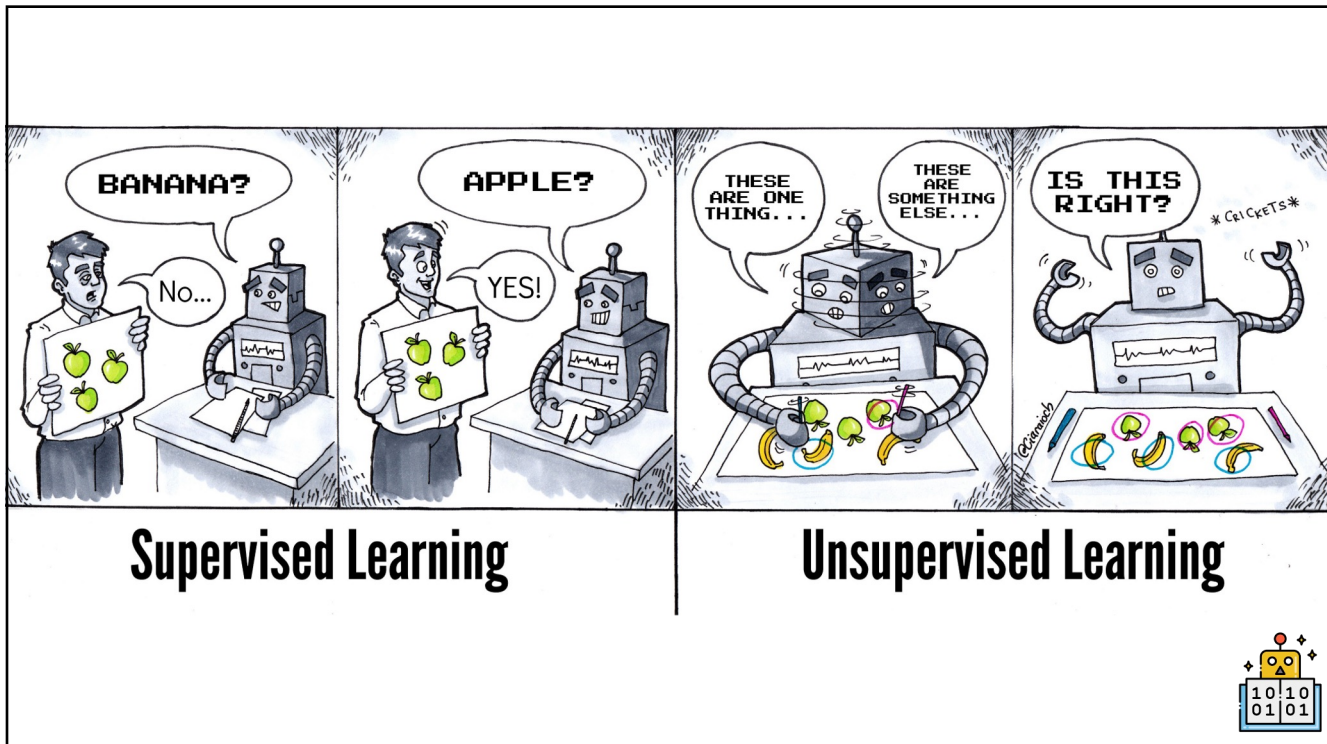
- Final exam
 - Time: Dec 13, 3:30 - 5:30pm
 - Location: C E Chavez Bldg, Rm 111 (same room)
 - What you can bring:
 - one letter size cheat sheet, you can use double sides
 - calculator (not necessary)

3

Announcements

- ~20 questions and 50% questions will be before midterm.
- Office hours next week will be announced this Thursday.
- Practice questions will be out by next Monday Dec 4.
- No coding questions.
- How to prepare
 - **Slides**
 - Practice problems (helpful but do not only rely on it!)
 - HW questions

4



5

Task 1 : Group These Set of Document into 3 Groups based on meaning

- Doc1 : Health , Medicine, Doctor
- Doc 2 : Machine Learning, Computer
- Doc 3 : Environment, Planet
- Doc 4 : Pollution, Climate Crisis
- Doc 5 : Covid, Health , Doctor



6

Task 1 : Group These Set of Document into 3 Groups based on meaning

Doc1 : Health , Medicine, Doctor

Doc 2 : Machine Learning, Computer

Doc 3 : Environment, Planet

Doc 4 : Pollution, Climate Crisis

Doc 5 : Covid, Health , Doctor



7

Task 1 : Group These Set of Document into 3 Groups based on meaning

Doc1 : Health , Medicine, Doctor

Doc 5 : Covid, Health , Doctor

Doc 3 : Environment,
Planet

Doc 4 : Pollution, Climate
Crisis

Doc 2 : Machine
Learning, Computer

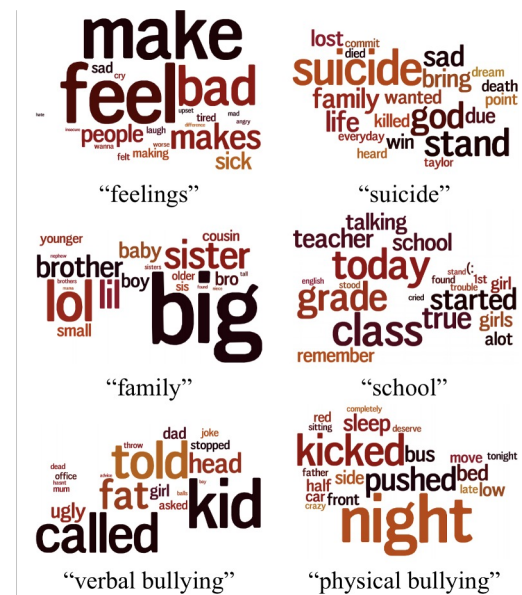


8

Task 2: Topic modeling

9

- Provides a summary of a corpus.
- Collected n tweets containing the keyword “bullying”, “bullied”, etc.
- Extracts k topics: each topic is a list of words with importance weights.
 - A set of words that co-occurs frequently throughout.



Learning from bullying traces in social media.

Jun-Ming Xu, Kwang-Sung Jun, Xiaojin Zhu, and Amy Bellmore.

In the Conference of North American Chapter of the Association for Computational Linguistics: Human Language Technologies (NAACL HLT), 2012. [pdf]

9

What is unsupervised learning?

10

- Learning with unlabeled data
- What can we expect to learn?
 - **Clustering**: obtain partition of the data that are well-separated.
 - a preliminary classification without predefined class labels.
 - **Components**: extract common components
 - e.g., topic modeling given a set of articles: each article talks about a few topics => extract the topics that appear frequently.
- How can we use?
 - As a summary of the data
 - **Exploratory data analysis**: what are the **patterns** even without labels?
 - As a ‘preprocessing techniques’
 - e.g., extract useful **features** using soft clustering assignments

10

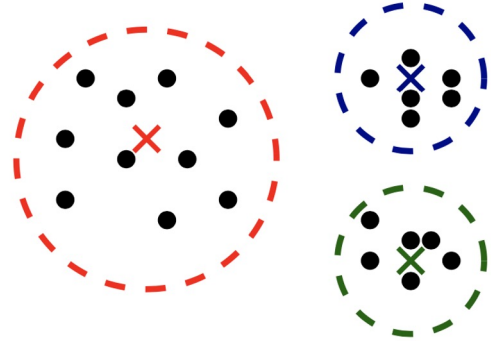
Clustering

11

- Input: k : the number of clusters (hyperparameter)

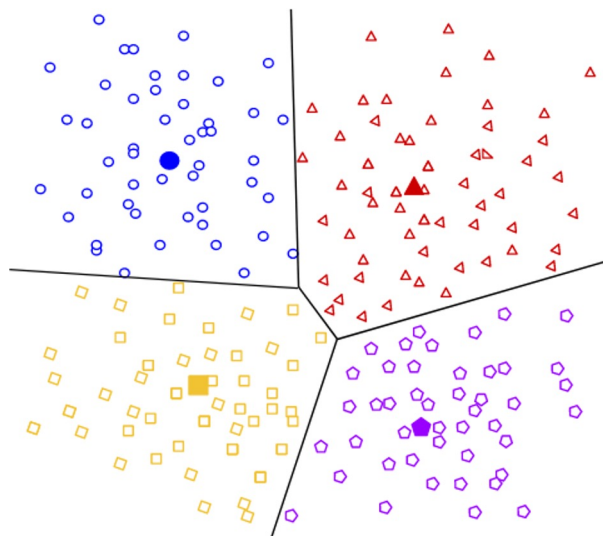
$$S = \{x_1, \dots, x_n\}$$

- Output
 - partition $\{G_i\}_{i=1}^k$ s.t. $S = \cup_i G_i$ (disjoint union).
 - often, we also obtain 'centroids'



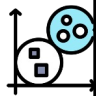
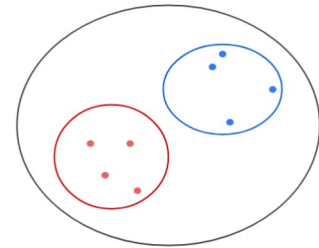
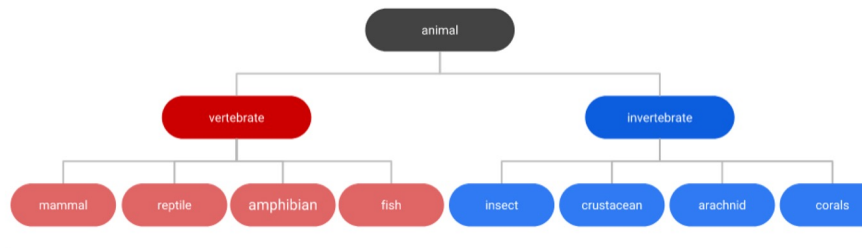
11

Centroid-based Clustering



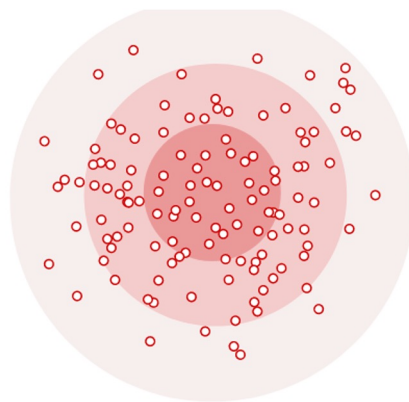
12

Hierarchical Clustering

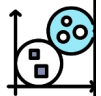
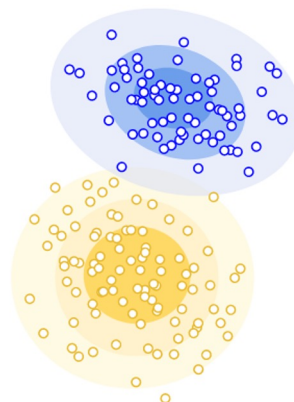


13

Distribution-based Clustering



(probabilistic treatment)



14

Clustering

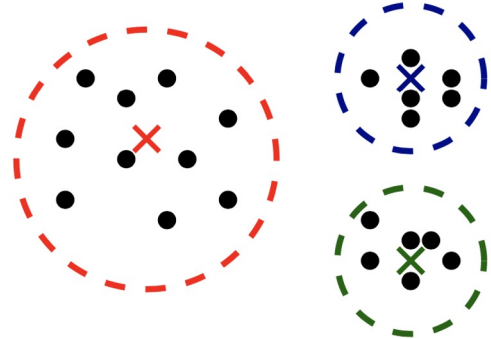
15

- Input: k : the number of clusters (hyperparameter)

$$S = \{x_1, \dots, x_n\}$$

- Output

- partition $\{G_i\}_{i=1}^k$ s.t. $S = \cup_i G_i$ (disjoint union).
- often, we also obtain 'centroids'



- Q: if we are given the groups, what would be a reasonable definition of centroids?

- The point that has the minimum average distance to the datapoints?
- The datapoint that has the minimum average distance to the datapoints?
- The point that has the minimum average squared distance to the datapoints?

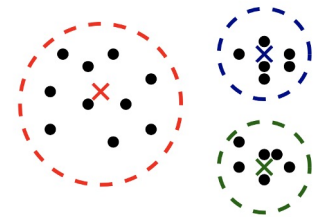
=> Turns out, the last one corresponds to the average point!

15

k-means Clustering

16

Lloyd's algorithm: solve it approximately (heuristic)



Observation: The chicken-and-egg problem.

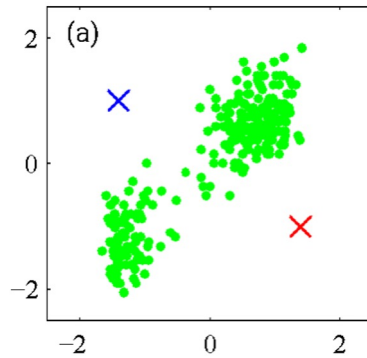
- If you knew the cluster assignments... just find the centroids as the average
- If you knew the centroids... make cluster assignments by the closest centroid.

Why not: start from some centroids and then alternate between the two?

16

Initialization

17

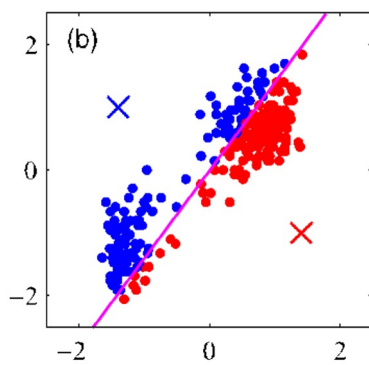


Arbitrary/random initialization of c_1 and c_2

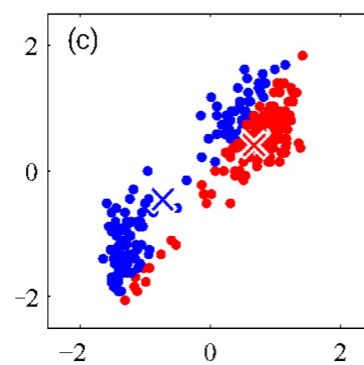
17

Iteration 1

18



(A) update the cluster assignments.

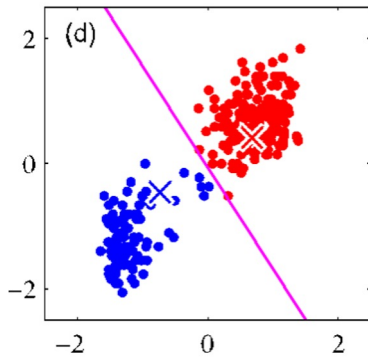


(B) Update the centroids $\{c_j\}$

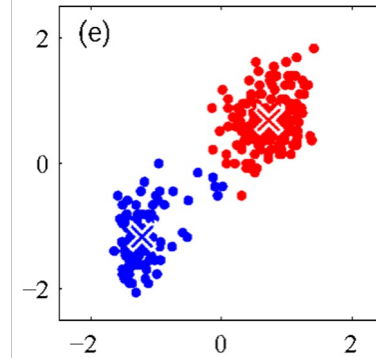
18

Iteration 2

19



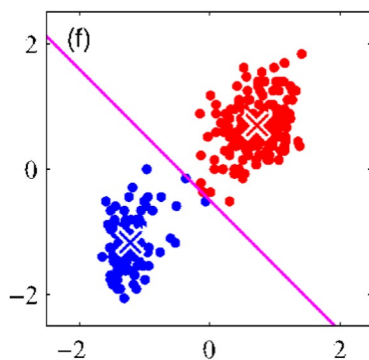
(A) update the cluster assignments.

(B) Update the centroids $\{c_j\}$

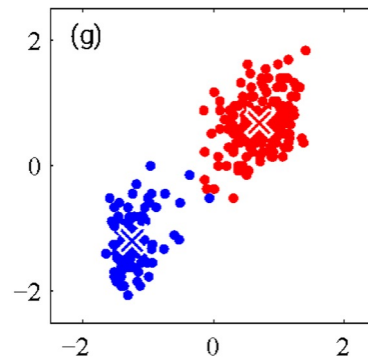
19

Iteration 3

20



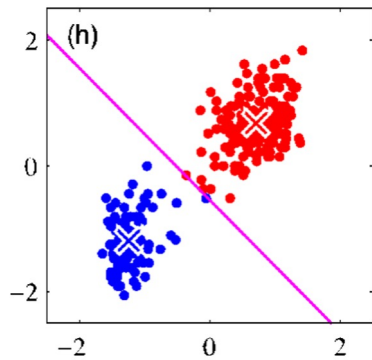
(A) update the cluster assignments.

(B) Update the centroids $\{c_j\}$

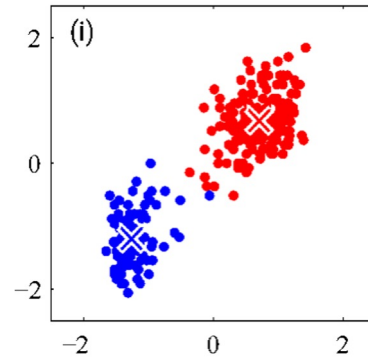
20

Iteration 4

21

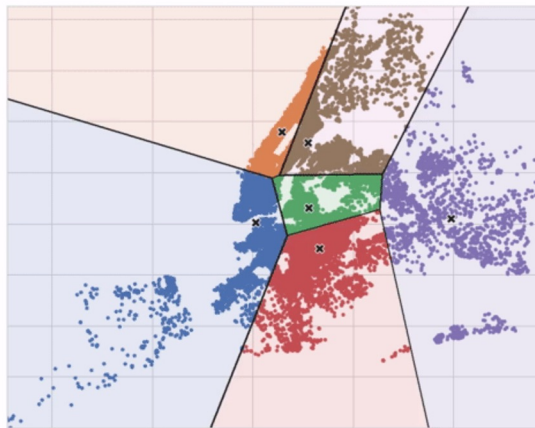


(A) update the cluster assignments.

(B) Update the centroids $\{c_j\}$

21

Iterating until Convergence

Animation from [Kaggle](https://www.kaggle.com)

22

k-means clustering

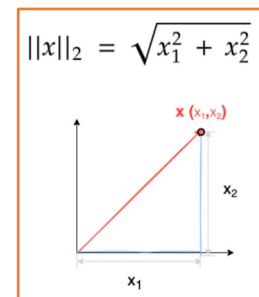
23

Input: k : num. of clusters, $S = \{x_1, \dots, x_n\}$

[Initialize] Pick c_1, \dots, c_k as randomly selected points from S (see next slides for alternatives)

For $t=1,2,\dots,\text{max_iter}$

- **[Assignments]** $\forall x \in S, \quad a_t(x) = \arg \min_{j \in [k]} \|x - c_j\|_2^2$
- If $t \neq 1$ AND $a_t(x) = a_{t-1}(x), \forall x \in S$
 - break



23

k-means clustering

24

Input: k : num. of clusters, $S = \{x_1, \dots, x_n\}$

[Initialize] Pick c_1, \dots, c_k as randomly selected points from S (see next slides for alternatives)

For $t=1,2,\dots,\text{max_iter}$

- **[Assignments]** $\forall x \in S, \quad a_t(x) = \arg \min_{j \in [k]} \|x - c_j\|_2^2$
- If $t \neq 1$ AND $a_t(x) = a_{t-1}(x), \forall x \in S$
 - break
- **[Centroids]** $\forall j \in [k], \quad c_j \leftarrow \text{average}(\{x \in S: a_t(x) = j\})$

Output: c_1, \dots, c_k and $\{a_t(x_i)\}_{i \in [n]}$

24

But,

It may converge to a local rather than global minimum.

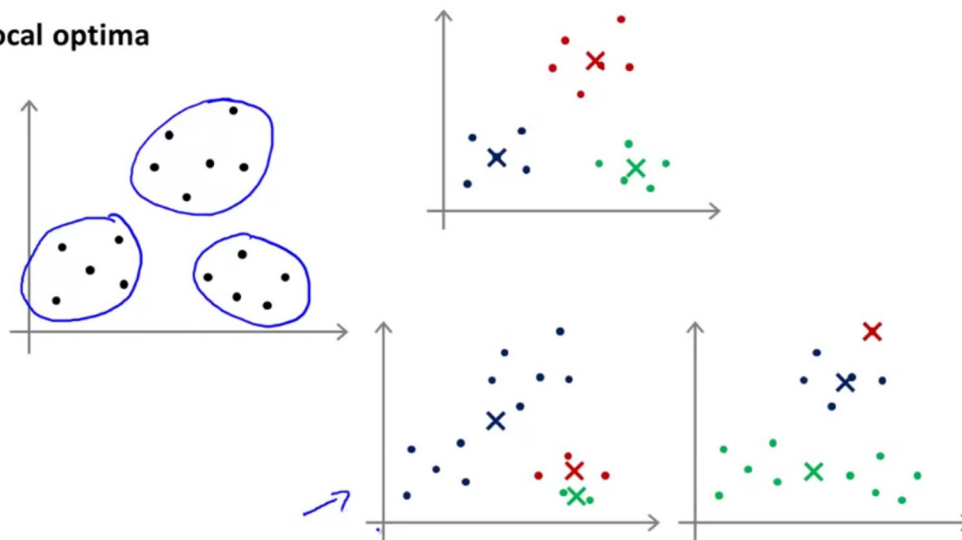
$$\text{objective function} \leftarrow J = \sum_{j=1}^k \sum_{i=1}^n \underbrace{\|x_i^{(j)} - c_j\|^2}_{\text{Distance function}}$$

number of clusters $\rightarrow k$
 number of cases $\rightarrow n$
 case i $\rightarrow x_i^{(j)}$
 centroid for cluster j $\rightarrow c_j$



25

Local optima



Andrew Ng

Image from Andrew NG Coursera Machine Learning Course



26

Issue 1: Unreliable solution

27

- You usually get suboptimal solutions
- You usually get different solutions every time you run.
- **Standard practice:** Run it 50 times and take the one that achieves the smallest objective function

- Recall: $\min_{c_1, \dots, c_k} \sum_{i=1}^n \min_{j \in [k]} \|x_i - c_j\|_2^2$ Each run of algorithm outputs c_1, \dots, c_k . Compute this to evaluate the quality!

- And/or, change the initialization (next slide)
 - Idea: ensure that we pick a widespread c_1, \dots, c_k

27

Alternative initialization

28

- **k-means++**
 - Pick $c_1 \in \{x_1, \dots, x_n\}$ uniformly at random
 - For $j = 2, \dots, k$
 - Define a distribution $\forall i \in [n], \mathbb{P}(c_j = x_i) \propto \min_{j'=1, \dots, j-1} \|x_i - c_{j'}\|_2^2$
 - Draw c_j from the distribution above.

More likely to choose x_i
that is farthest from
already-chosen centroids.

=> has a mathematical guarantee that it will be better than an arbitrary starting point!

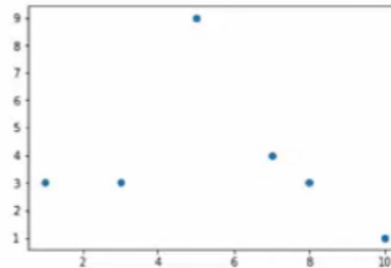
28

Suppose we have the small dataset

☞ $[(7,4),(8,3),(5,9),(3,3),(1,3),(10,1)]$ to which we wish to assign 3 clusters.

We begin by randomly selecting $(7,4)$ to be a cluster center.

| x | $\min(d(x, z_i)^2)$ |
|----------|---------------------|
| $(7,4)$ | |
| $(8,3)$ | |
| $(5,9)$ | |
| $(3,3)$ | |
| $(1,3)$ | |
| $(10,1)$ | |



[From Sara Jensen's Youtube Channel](#)



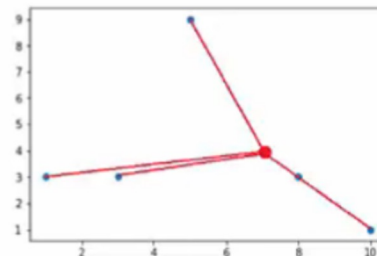
29

Suppose we have the small dataset

$[(7,4),(8,3),(5,9),(3,3),(1,3),(10,1)]$ to which we wish to assign 3 clusters.

We begin by randomly selecting $(7,4)$ to be a cluster center.

| x | $\min(d(x, z_i)^2)$ |
|----------|---------------------|
| $(7,4)$ | - |
| $(8,3)$ | 2 |
| $(5,9)$ | 29 |
| $(3,3)$ | 17 |
| $(1,3)$ | 37 |
| $(10,1)$ | 18 |

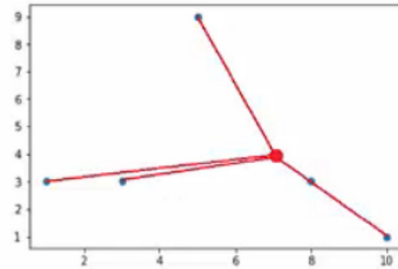


30

Suppose we have the small dataset
 $[(7,4), (8,3), (5,9), (3,3), (1,3), (10,1)]$ to which we wish to assign 3 clusters.

We begin by randomly selecting $(7,4)$ to be a cluster center.

| x | prob |
|----------|----------|
| $(7,4)$ | - |
| $(8,3)$ | $2/103$ |
| $(5,9)$ | $29/103$ |
| $(3,3)$ | $17/103$ |
| $(1,3)$ | $37/103$ |
| $(10,1)$ | $18/103$ |

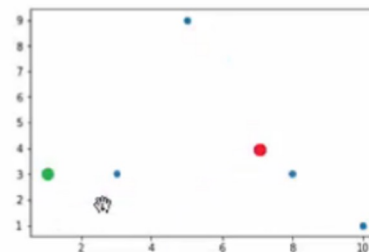


31

Suppose we have the small dataset
 $[(7,4), (8,3), (5,9), (3,3), (1,3), (10,1)]$ to which we wish to assign 3 clusters.

We add $(1,3)$ to the list of cluster centers.

| x | $\min(d(x, z_i)^2)$ |
|----------|---------------------|
| $(7,4)$ | - |
| $(8,3)$ | |
| $(5,9)$ | |
| $(3,3)$ | |
| $(1,3)$ | - |
| $(10,1)$ | |



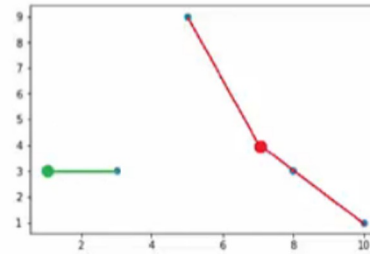
32

Suppose we have the small dataset

$[(7,4),(8,3),(5,9),(3,3),(1,3),(10,1)]$ to which we wish to assign 3 clusters.

We add $(1,3)$ to the list of cluster centers.

| x | $\min(d(x, z_i)^2)$ |
|----------|---------------------|
| $(7,4)$ | - |
| $(8,3)$ | 2 |
| $(5,9)$ | 29 |
| $(3,3)$ | 4 |
| $(1,3)$ | - |
| $(10,1)$ | 18 |



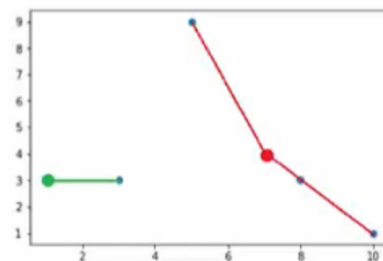
33

Suppose we have the small dataset

$[(7,4),(8,3),(5,9),(3,3),(1,3),(10,1)]$ to which we wish to assign 3 clusters.

We add $(1,3)$ to the list of cluster centers.

| x | prob |
|----------|---------|
| $(7,4)$ | - |
| $(8,3)$ | $2/53$ |
| $(5,9)$ | $29/53$ |
| $(3,3)$ | $4/53$ |
| $(1,3)$ | - |
| $(10,1)$ | $18/53$ |

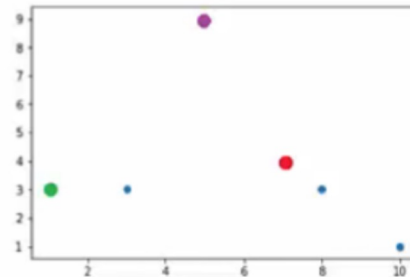


34

Suppose we have the small dataset
 $[(7,4),(8,3),(5,9),(3,3),(1,3),(10,1)]$ to which we wish to assign 3 clusters.

We add (5,9) to the list of cluster centers.

| x | prob |
|--------|------|
| (7,4) | - |
| (8,3) | |
| (5,9) | - |
| (3,3) | |
| (1,3) | - |
| (10,1) | |



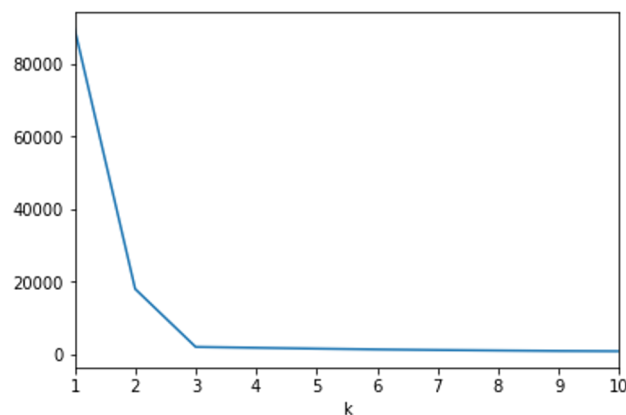
35

Issue 2: Choose k

36

- No principled way.
- Elbow method: see where you get saturation.

Objective function



<https://medium.com/analytics-vidhya/how-to-determine-the-optimal-k-for-k-means-708505d204eb>

36