

CLUSTERING - unsupervised learning

↳ similar points clustered together

- each data point assigned to a specific group where they are
 - ① similar to each other, ② dissimilar to other groups

Goal: define a process that can give us clustering

↗ to optimize this

- what does similar mean? - How do we find clustering?
- How do we know if it's good clustering?

difficulties: ambiguous

Types of Clustering

- Partitional: each object belongs to exactly 1 cluster (k clusters)
- Hierarchical: set of nested clusters organized as tree
(cut the tree @ any point \rightarrow new cluster)
- Density-Based: defined based on local density of points \rightarrow (need to define density)
- Soft Clustering: each point assigned to clusters w/ a probability

\rightarrow clustering in real world \rightarrow detecting fraudulent purchases } filter out outliers
 \rightarrow filling in NaN values } marked as spam (filter)

Partitional Clustering

need to specify k upfront = limitation

- given n data points and a k # of clusters
 - partition n data points into k clusters
- given all the ways to distribute n into k buckets
 - which is the best?

if you sum all dissimilarities (ex. using pairwise) should get a small #

\rightarrow intra-cluster distances

\rightarrow good partition = total dissimilarity is small

* want to minimize this "cost"

Cost Function

$$= \sum_{k=1}^K \sum_{x_i, x_j \in C_k} d(x_i, x_j)$$

sum of all clusters \rightarrow for each cluster C_k take distance between x_i, x_j

\rightarrow given function distance d can find centroid (center of mass)

\rightarrow when [Euclidean] \rightarrow mean/average = centroid

Turns out when d is Euclidean → just look @ all distances to the centroid → no need to look @ all pairwise distances

in a given cluster

$$\sum_k \sum_{x_i, x_j \in C_k} d(x_i, x_j)^2 = \sum_k |C_k| \sum_{x_i \in C_k} d(x_i, \mu_k)^2$$

* optimizing pairwise differences is equivalent to optimizing distance to center

→ scalar = size of C_k = # of points in C_k

→ center of the cluster

★ K-means Given $X = \{x_1, \dots, x_n\}$ — find k points that minimize the cost function $\{\mu_1, \dots, \mu_k\}$

sum over all clusters over all points at the point to centroid

minimize distance to centroid
because it's equivalent to minimizing dissimilarity

$$\sum_i \sum_{x \in C_i} \|x - \mu_i\|_2^2 = \sum_{i=1}^k \sum_{x \in C_i} L_2(x, \mu_i)$$

L_2 norm?

euclidean distance

centroid of C_i → mean (because euclidean)

recall: $p=2$ means euclidean distance

↳ when $k=1$ and $k=n$ → each point has its own little cluster
this is easy. Why?

one cluster
— no other options

limitation ↳ when x_i lives in more than 2 dim (> 2 features), this is VERY difficult (NP-hard)

Trying to optimize algorithm
K-means → Lloyd's Algorithm

- 1) Randomly pick k centers (μ_1, \dots, μ_k)
 - 2) assign each point to its closest center → calculate distance to all centers & choose smallest
 - 3) compute the new centers as means of each cluster → computing mean of points in a cluster generates new center
 - 4) repeat 2 & 3 until convergence
↳ the points assigned are the same that you get when computing new one
↳ center doesn't move
- when the mean = center
(need ALL means to stop moving)

Will this algorithm always converge? Proved by contradiction

↳ only a finite # of partitions to only a finite # of ways we can pick to set

↳ min. cost function is only needed in the limit

(infinite # of iterations)

= IMPOSSIBLE, we have finite # points

[Does it always converge?

Proof by contradiction \rightarrow 1) Suppose it does NOT converge: then...

* minimum of cost function only reached in limit (infinite iterations)

\Rightarrow IMPOSSIBLE (finite # of points) \rightarrow finite $n, \text{ and } k$
 \therefore finite # of double sums

OR

* algorithm gets stuck in the loop/cycle

\Rightarrow IMPOSSIBLE (@ each step, we lower the cost

if old \neq new, cost has improved

if old = new, cost hasn't changed

* cannot have a cost that is "smaller than itself"

[But does it converge to the OPTIMAL solution?

\rightarrow NOT always

Solution = (ex.) Run Lloyd's algorithm multiple times & choose result w/ lowest cost

(Limitation) bad results because of randomness

Solution = (ex.) Random Initialization (different initialization methods)

① Farthest First Traversal

- pick the next center to be farthest from all previous ~~pts~~

1) pick a random first one
2) \rightarrow

limitation = outliers in dataset cause issues

\rightarrow choose outliers as centers and then all the other data are identified as outliers

* Need to combine "Random" & "Farthest First Traversal"

② K-Means ++

- Initialize w/ a combo of the methods

1) start w/ random center

2) let $DC(x)$ be distance between x and centers selected so far
choose next center w/ prob. proportional to $(DC(x))^a$ \rightarrow when $a=0$ (random)
the distance

possible candidate for a center (within a cluster)

~~weight proportional to~~

- x selected based on weighted probability as a function of when
met point's distance to centers picked so far
 \rightarrow points @ greater distance = higher probability (farthest first traversal)
 $a=2$

* trying to get farthest first traversal probabilistically

(K-means ++)

- squared distances

[How do you pick points w/ probability proportional to the distance?

→ Suppose given a random number generator (black box)
- generates uniform random number between 0 and any N

↳ How do you use this to select points w/ prob. proportional to $D(x)^a$ - $a=2$

$a=2$
which of x, y, z do we pick?
(w/ prob. proportional to the distance)

x
 y
 z

$* D(x)^2 = 3^2 = 9$
 $D(y)^2 = 2^2 = 4$
 $D(z)^2 = 1^2 = 1$

* write out the distances
(note: don't have to be int distances)

* assign x to 9 buckets in random # generator
 y to 4 buckets
 z to 1 bucket

generate a # between 0 and N
↳ it necessarily falls into one of the buckets!!
the prob. of falling into x bucket is HIGHER

$$N = D(x)^2 + D(y)^2 + D(z)^2 = 14$$

* if black box can only generate numbers $(0,1)$ → divide by N

Limitations - prefers clusters of similar size

not good w/ clusters of different densities

tries to find globular shapes (fits in a circle)

→ a big ball in $N-D$ space

How to choose the right k ? (# of centers)

1) Iterate through different values of k (elbow method)

cost
↑
↓
find something in the elbow of the curve
↳ maximize distance to inner reference

← computationally expensive

2) Use empirical / domain-specific knowledge of the data

- Is there a known approx. distribution? (K-means good for →)

spherical Gaussians

K-Means Variations

1) K-medians (L_1 norm / Manhattan distance)

2) K-medoids (use any distance function + centers in the dataset)

3) Weighted K-means (each point has different weight when computing means)

→ ALWAYS