

Lee 3

Disimilarity

① - Minkowski

②

Euklidian

$= 1$

triangle dist

③

Manhattan

$= 2$

black dist

params: d (dimension), $p =$ weight?

④ - Jaccard Similarity \rightarrow Manhattan for set difference.

1 when $x_i \neq y_i$ and 0 when $x_i = y_i$.

* differ all then Manhattan of 1

$$JDist(x, y) = 1 - \frac{|x \cap y|}{|x \cup y|} \quad \begin{array}{l} \text{same} \\ \text{different} \end{array}$$

⑤

Cosine similarity

- take two objects, return large if similar.

$$s(x, y) = \cos \theta$$



$= 1$



$= 0$



$= -1$

$$\hookrightarrow \text{disimilarity} = \frac{1}{s(x, y)} \quad \text{or} \quad 1 - s(x, y)$$

lec #4

Clustering - kmeans

Clustering is groups of assignment

Cost function
$$\sum_i \sum_{i \in C_i} d(x_i, \mu_i)^2$$

* minimize cost

params: d : euclidean dist
 k : center of clusters

* $k=1$ $k=n$ is easy bc.

↓ ↓
all is each k is own cluster
one cluster

- > 2 dimensions would be rly hard

Lloyd's algo

1. pick center at random, k
2. assign data point to closest k
3. compute new cluster center.
4. repeat until converge

we could pick outlier
as centers :c

but not always
optimal!

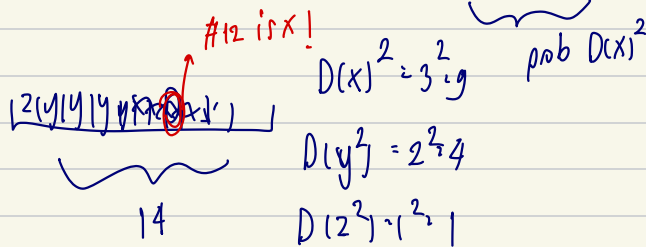
Lec #5

/ equally fine w/ K mean

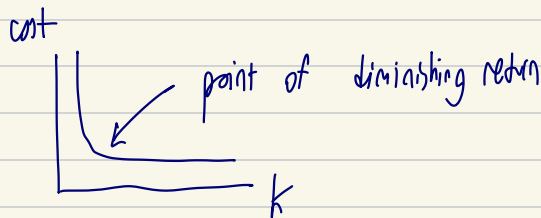
K-mean++

1. pick random center
2. let $D(x)$ be dist of x and closest center picked.
Choose next center w/ prob proportional to $D(x)^2$

- Black box and uniformly choose btw 0 and N .



♥ choose right K via elbow method



Silhouette score

- $\left\{ \begin{array}{l} - \text{avg within cluster dist} : a \\ - \text{avg into cluster dist} : b \end{array} \right\}$
 if $b-a = 0$ they identical

$b-a$ large + separate, high quality

$b-a$ small not so much

♡ Silhouette score

$$\frac{(b-a)}{\max(a,b)},$$

a : within avg dist
 b : intra cluster

-1 to 1
↓ ↓
bad good

K-means variations

1. k medians (L_1 norm / manhattan dist)
2. k-medoid (any dist + center in data)
3. weighted k-means (each point diff weight)

Hierarchical clustering

- ① Agglomerative (bottom \rightarrow top)
- ② Divisive (top \rightarrow down)

- every step record cluster to merge in Dendrogram

parameters : dist between points + dist btw clusters

① Single-Link Distance (min of pairwise dist)

pros:

cons: sensitive to noise, elongated cluster / chain-like

② Complete-Link Distance (max of pairwise)

pros: less susceptible to noise, more balanced, equal diameters

cons: split up large cluster, same, circle sphere

③ Average-Link Distance (avg. of pairwise ^{combination} dist)

pros: less susceptible to noise and outliers

cons: bias toward globular cluster

④ centroid Distance (dist b/w two centroid of cluster)

⑤ Ward's Distance (diff b/w variance of points in cluster)

Lecture #6

Density-based Clustering

↳ cluster pts that are densely banded together

params : fix radius ϵ —
min # of point min-pt —

- ① core : ϵ -neighborhood w/ at least min-pt
- ② border : in ϵ but not core.
- ③ Noise : neither core nor border

DBScan algorithm (DFS)

1. find ϵ -neighborhood for each point
— if at least min-pt \rightarrow core.
2. for each core, assign. to same cluster
all core point in same ϵ
3. if in core ϵ -neighborhood then mark border.
4. what left is define all else as noise.
5. assign border to nearby cluster

pro : identify diff sizes & shape, resistant to noise

cons : fail to identify clusters of differing density
tend to create clusters of same density
notion of density is problematic.

Lecture #7 / Normal distribution

Soft-clustering $P(s_j | x_i)$

Maximum likelihood

↳ output is weight / prob that it come from that species

$$P(s_j | x_i) = \frac{P(x_i | s_j) P(s_j)}{P(x_i)}$$

Lecture #8

Clustering Aggregation

↳ compare & combine

↳ agree or disagree

- Disagreement Dist

$$D(x,y) = \begin{cases} 1, & \text{if disagree} \\ 0, & \text{if agree} \end{cases}$$

Aggregate clustering :

pros: identify best # of clks
detect outlier
robust
privacy

cons : NP-hard