

CS725: Recap kernels

A kernel $K(x, y)$ takes x, y inputs in the original feature space and computes $\phi(x)^T \phi(y)$, which is a measure of similarity between x and y in a new (higher-dimensional) subspace defined implicitly by ϕ .

How do we determine if K is valid?

① Can we identify a ϕ function s.t. $K(x, y) = \phi(x)^T \phi(y)$

② Mercer's theorem

③ By applying known operations (sum, product, etc.) on known/valid kernels

Q. Is $K(x, y) = \|x + y\|^2$ a valid kernel?

$$K(x, y) = \|x + y\|^2 = (x + y)^T (x + y)$$

No

$$K(0, 0) = 0 = \phi(0)^T \phi(0) \Rightarrow \phi(0) = 0$$

$$K(x, 0) = 0 \rightarrow \textcircled{A} \quad \because \phi(0) = 0$$

$$K(x, 0) = \|x\|^2 \rightarrow \textcircled{B} \quad \because \text{definition of } K(x, y)$$

CONTRADICTION.

Examples of popular kernels

(A) Polynomial kernel: $K(x, y) = (x^T y)^d, d > 0$
 $= (c + x^T y)^d, d > 0$

(B) Gaussian kernel: $K(x, y) = \exp\left(-\frac{1}{2\sigma^2} \|x - y\|^2\right)$

[The implicit ϕ corresponding to a Gaussian kernel is infinite-dimensional]

\downarrow for small σ , all the pts are distinct or different
(K is a diagonal matrix)

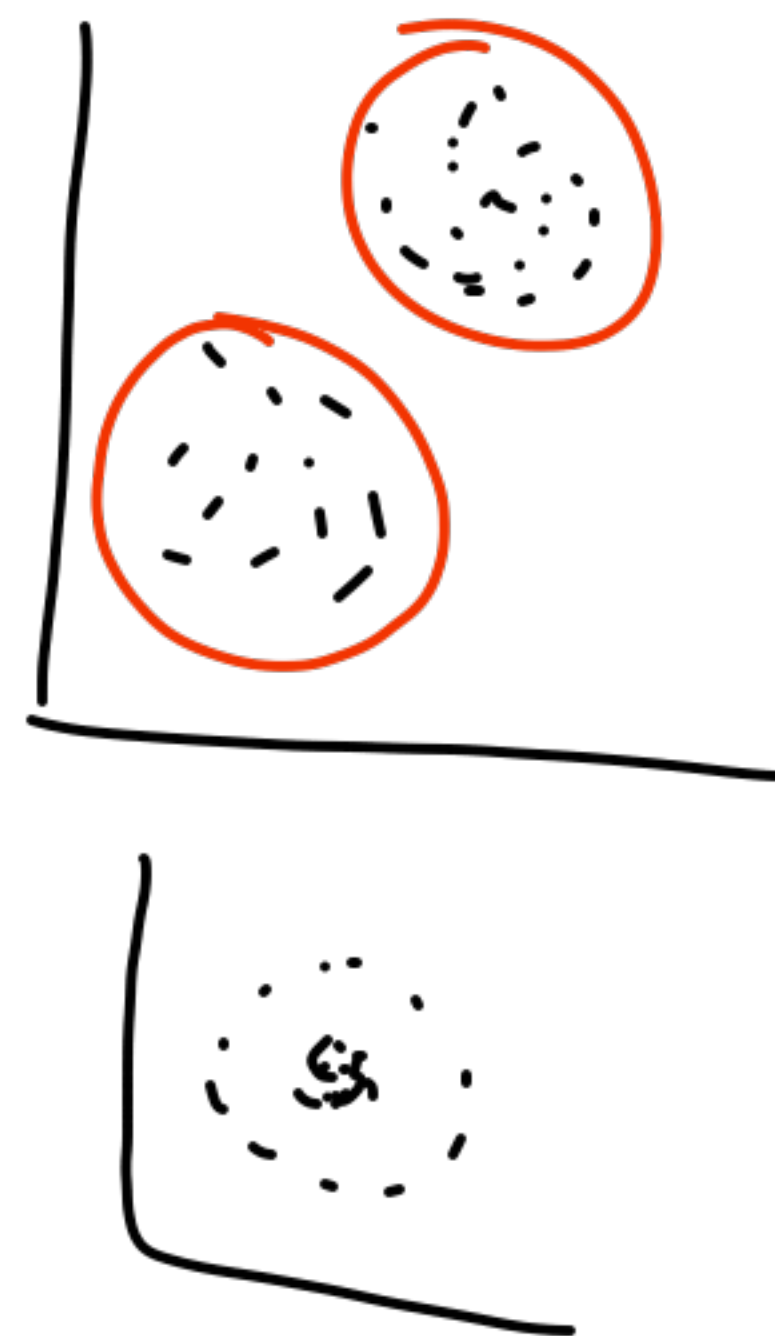
for large σ , all the pts are v. similar
(K is a uniform matrix)

CLUSTERING

Given a set of points $x_1, x_2, \dots, x_n, x_i \in \mathbb{R}^d$, we aim to find K clusters of points where the points within each cluster are similar to each other [typically in terms of Euclidean distances].

Motivation: (e.g. of clustering; image segmentation, customer aggregation, etc.)

- ① First, coarse way of labeling points
- ② Outlier detection
- ③ Lossy compression of the original data



k-means clustering algorithm

Assume we have K clusters of points; each point in a cluster is closest to its centroid (more than any other cluster centroid)

mean of all the datapoints
in a cluster

If cluster assignment is known, it is easy to compute the centroids

If cluster centroids are known, it is easy to do cluster assignment

How do we solve this chicken-egg problem? Fix one, optimize the other!
(ALTERNATING OPTIMIZATION)

k-means objective:

find $\mu_1, \mu_2, \dots, \mu_k, \mu_i \in \mathbb{R}^d$ (cluster centroids) and C^1, C^2, \dots, C^n (cluster assignments for n training points) to minimize the sum of squared distances of each point from its cluster centroid. Here, C^i is a 1-hot vector of size k .

Objective function:
$$\min_{\{\mu_k\}} \min_{\{C^i\}} \sum_{i=1}^n \sum_{k=1}^k C_k^i \underbrace{\| \mu_k - x_i \|^2}_{\text{A}}$$

where $C_k^i = 1$ [x_i is assigned to cluster index k]

Heuristic to solve (A) [to reach a local optimum] is to fix one set of variables and optimize the other

Step 1: fix μ and optimize the assignment

for the i^{th} datapoint, $\min_{C_i} \sum_{k=1}^K C_k^i \| \mu_k - x_i \|^2$

$$C_k^i = \begin{cases} 1 & \text{if } k = \underset{j}{\operatorname{argmin}} \| \mu_j - x_i \|^2 \\ 0 & \text{otherwise} \end{cases}$$

ASSIGNMENT
STEP

Step 2: fix the assignment and optimize for the cluster centroids

$$\frac{\partial}{\partial \mu_k} \sum_{i=1}^n \sum_{k=1}^K C_k^i \|\mu_k - x_i\|^2, \text{ set to 0 and solve for } \mu_k$$

$$\Rightarrow 2 \sum_{i=1}^n C_k^i (\mu_k - x_i) = 0$$

$$\Rightarrow \mu_k = \frac{\sum_i C_k^i x_i}{\sum_i C_k^i}$$

→ cluster centroids
are mean of all the
assigned datapoints

UPDATE
STEP

k-means algorithm

Given : $x_1, x_2, \dots, x_n, x_i \in \mathbb{R}^d$, number of clusters is K

Initialize $\mu_1, \mu_2, \dots, \mu_K$

Repeat until convergence (i.e. assignments do not change)

1. Assignment : Assign each point to the closest μ_k
2. Update : Recompute the cluster centroids

k-means is guaranteed to converge (to a local optimum)
which depends on μ_1, \dots, μ_k
initialization

↳ Because of the following two facts:

- ① The cost or objective either stays the same or decreases across iterations
- ② There are only a finite number of cluster assignments

k-means is very sensitive to its initialization

Common initialization strategy is to pick the centroids randomly

Another strategy is "farthest point". Pick the first centroid at random, and every succeeding centroid is the farthest point from the preceding centroids

Limitation of farthest pt \rightarrow could pick outliers

Another variant is k-means++ \rightarrow instead of farthest pt, pick a pt based on a prob distribution weighted acc. to $D(x)^2$ where $D(x)$ is the distance of point x from the preceding centroids

How do we choose k in k -means clustering?

Keep varying k and compute the squared distances of all pts from its centroids

