K-means extensions and evaluating clusterings

CS771: Introduction to Machine Learning Nisheeth

K-means algorithm: recap

- Notation: $z_n \in \{1, 2, ..., K\}$ or z_n is a K-dim one-hot vector
 - $(z_{nk} = 1 \text{ and } z_n = k \text{ mean the same})$

K-means Algorithm

- Initialize K cluster means μ_1, \ldots, μ_K
- **2** For n = 1, ..., N, assign each point x_n to the closest cluster

$$z_n = \arg\min_{k \in \{1,\ldots,K\}} ||\boldsymbol{x}_n - \boldsymbol{\mu}_k||^2$$

Suppose $C_k = \{x_n : z_n = k\}$. Re-compute the means

$$\mu_k = \operatorname{mean}(\mathcal{C}_k), \quad k = 1, \ldots, K$$

Go to step 2 if not yet converged



K-means loss function: recap

• Let $\mu_1, \mu_2, \dots, \mu_K$ be the K cluster centroids/means



Κ

(one-hot vector)

CS771: Intro to ML

- Let $z_{nk} \in \{0, 1\}$ be s.t. $z_{nk} = 1$ if x_n belongs to cluster k, and 0 otherwise
- Define the distortion or "loss" for the cluster assignment of $oldsymbol{x}_n$

$$\ell(\boldsymbol{\mu}, \boldsymbol{x}_n, \boldsymbol{z}_n) = \sum_{k=1}^{n} z_{nk} ||\boldsymbol{x}_n - \boldsymbol{\mu}_k||^2$$

• Total distortion over all points defines the K-means "loss function" $_{\rm K}$

$$L(\mu, \mathbf{X}, \mathbf{Z}) = \sum_{n=1}^{N} \sum_{k=1}^{K} z_{nk} ||\mathbf{x}_n - \mu_k||^2 = \underbrace{||\mathbf{X} - \mathbf{Z}\mu||_F^2}_{\text{matrix factorization view}} \mathbb{N} \mathbf{X} \approx \mathbf{Z} \approx \mathbf{Z} \operatorname{Row} k \text{ is } \mu_k$$
Row k is μ_k

- The K-means problem is to minimize this objective w.r.t. μ and Z
 - Alternating optimization on this loss would give the K-means (Lloyd's) algorithm we saw earlier!

K-means++

K-means results can be sensitive to initialization



- K-means++ (Arthur and Vassilvitskii, 2007) an improvement over K-means
 - Only difference is the way we initialize the cluster centers (rest of it is just K-means)
 - Basic idea: Initialize cluster centers such that they are reasonably far from each other
 - Note: In K-means++, the cluster centers are chosen to be K of the data points themselves

CS771: Intro to ML

K-means++

- K-means++ works as follows
 - Choose the first cluster mean uniformly randomly to be one of the data points
 - The subsequent K 1 cluster means are chosen as follows
 - 1. For each unselected point \boldsymbol{x} , compute its smallest distance $D(\boldsymbol{x})$ from already initialized means
 - 2. Select the next cluster mean unif. rand. to be one of the unselected points based on probability prop. to $D(x)^2$
 - 3. Repeat 1 and 2 until the K 1 cluster means are initialized
 - Now run standard K-means with these initial cluster means
- K-means++ initialization scheme sort of ensures that the initial cluster means are located in different clusters

Thus farthest points are most

likely to be selected as cluster

CS771: Intro to ML

means

K-means: Soft Clustering

- K-means makes hard assignments of points to clusters
 - Hard assignment: A point either completely belongs to a cluster or doesn't belong at all





• A heuristic to get soft assignments: Transform distances from clusters into prob.

$$\sum_{k=1}^{K} \gamma_{nk} = 1 \qquad \gamma_{nk} = \frac{\exp(-||\boldsymbol{x}_n - \boldsymbol{\mu}_k||^2)}{\sum_{\ell=1}^{K} \exp(-||\boldsymbol{x}_n - \boldsymbol{\mu}_\ell||^2)} \quad \text{(prob. that } \boldsymbol{x}_n \text{ belongs to cluster } k\text{)}$$

• Cluster mean updates also change: $\mu_k = \frac{\sum_{n=1}^{N} \gamma_{nk} \mathbf{x}_n}{\sum_{n=1}^{N} \gamma_{nk}}$ (all points contribute, fractionally)

CS771: Intro to ML

K-means: Decision Boundaries and Cluster Sizes/Shapes'

- K-mean assumes that the decision boundary between any two clusters is linear
- Reason: The K-means loss function implies assumes equal-sized, spherical clusters





May do badly if clusters are not roughly equi-sized and convex-shaped





Kernel K-means

Helps learn non-spherical clusters and nonlinear cluster boundaries

Basic idea: Replace the Eucl. distances in K-means by the kernelized versions

$$\begin{aligned} ||\phi(\boldsymbol{x}_n) - \phi(\boldsymbol{\mu}_k)||^2 &= ||\phi(\boldsymbol{x}_n)||^2 + ||\phi(\boldsymbol{\mu}_k)||^2 - 2\phi(\boldsymbol{x}_n)^\top \phi(\boldsymbol{\mu}_k) \\ &= k(\boldsymbol{x}_n, \boldsymbol{x}_n) + k(\boldsymbol{\mu}_k, \boldsymbol{\mu}_k) - 2k(\boldsymbol{x}_n, \boldsymbol{\mu}_k) \end{aligned}$$

- Here k(.,.) denotes the kernel function and ϕ is its (implicit) feature map
- Note: $\phi(\mu_k)$ is the mean of ϕ mappings of the data points assigned to cluster k

<u>Not</u> the same as the ϕ mapping of the mean of the data points assigned to cluster k

input

$$\phi(\mu_k) = \frac{1}{|\mathcal{C}_k|} \sum_{n: z_n = k} \phi(\boldsymbol{x}_n)$$



Can also used landmarks or kernel random features idea to get new features and run standard k-means on those



Note: Apart from kernels, it is also possible to use other distance functions in K-means. Bregman Divergence* is such a family of distances (Euclidean and Mahalanobis are special cases)

Overlapping Clustering

- Have seen hard clustering and soft clustering
- In hard clustering, z_n is a one-hot vector
- In soft clustering, z_n is a vector of probabilities

Kind of unsupervised version of multi-label classification (just like standard clustering is like unsupervised multi-class classification)

Example: Clustering people based on the interests they may have (a person may have multiple interests; thus may belong to more than one cluster simultaneously)

- Overlapping Clustering: A point can <u>simultaneously</u> belong to multiple clusters
 - This is different from soft-clustering
 - z_n would be a binary vector, rather than a one hot or probability vector, e.g.,

 $z_n = [1 \ 0 \ 0 \ 1 \ 0]$ K=5 clusters with point x_n belonging (<u>in whole</u>, not in terms of probabilities) to clusters 1 and 4

- In general, more difficult than hard/soft clustering (for N data points and K clusters, the size of the space of possible solutions is not K^N but 2^{NK} exp in both N and K)
- K-means has extensions* for doing overlapping clustering. There also exist latent variable models for doing overlapping clustering

*An extended version of the k-means method for overlapping clustering (Cleuziou, 2008); Non-exhaustive, Overlapping k-means (Whang et al, 2975)1: Intro to ML

Evaluating Clustering Algorithms

- Clustering algos are in general harder to evaluate since we rarely know the ground truth clustering (since clustering is unsupervised)
- If ground truth labels not available, use output of clustering for some other task
 - For example, use cluster assignment z_n (hard or soft) as a new feature representation
 - Performance on some task using this new rep. is a measure of goodness of clustering
- If ground truth labels are available, can compare them with clustering based labels
 - Not straightforward to compute accuracy since the label identities may not be the same, e.g.,

Ground truth = $[1 \ 1 \ 1 \ 0 \ 0 \ 0]$ Clustering = $[0 \ 0 \ 0 \ 1 \ 1 \ 1]$

(Perfect clustering but zero "accuracy" if we just do a direct match)

- There are various metrics that take into account the above fact
 - Purity, Rand Index, F-score, Normalized Mutual Information, etc



Evaluating Clustering Algorithms

Purity: Looks at how many points in each cluster belong to the majority class in that cluster

