Unsupervised Learning: Clustering (*K*-means and other algos)

CS771: Introduction to Machine Learning Pivush Rai

Unsupervised Learning

- It's about learning interesting/useful structures in the data (unsupervisedly!)
- There is no supervision (no labels/responses), only inputs x_1, x_2, \dots, x_N
- Some examples of unsupervised learning
 - Clustering: Grouping similar inputs together (and dissimilar ones far apart)
 - Dimensionality Reduction: Reducing the dimensionality of inputs
 - Estimating the probability density of inputs (which distribution $p(x|\theta)$ "generated" the inputs)

0 0 0 - x_n belongs to cluster 3

- Unsupervised feature/representation learning
- Most unsup. learning algos also learn a new feature representation of inputs, e.g.,
 - Clustering gives a one-hot "quantized" representation \mathbf{z}_n for each input \mathbf{x}_n K = 6 clusters Assuming each input belongs

Some clustering algos learn a soft/probabilistic clustering in which \mathbf{z}_n will be be probability vector that sums to 1 (will see later)

A one-hot (quantized) rep. • Dim-red gives \boldsymbol{z}_n , a lower-dim representation of \boldsymbol{x}_n

 \boldsymbol{z}_n

- Unsup. Rep. learning can learn lower/higher dimensional representation \boldsymbol{z}_n of \boldsymbol{x}_n

deterministically to a single cluster

Clustering

In some cases, we may not know the right number of clusters in the data and may want to learn that (technique exists for doing this but beyond the scope)

- Given: N unlabeled inputs x_1, x_2, \dots, x_N ; desired no. of partitions K
- Goal: Group the examples into K "homogeneous" partitions \sim



In addition to partitioning these N inputs, we may also want to predict which partition a <u>new test point</u> belongs to

Picture courtesy: "Data Clustering: 50 Years Beyond K-Means", A.K. Jain (2008)

- Loosely speaking, it is classification without ground truth labels of training data
- A good clustering is one that achieves
 - High within-cluster similarity
 - Low inter-cluster similarity



Similarity can be Subjective

- Clustering only looks at similarities b/w inputs, since no labels are given
- Without labels, similarity can be hard to define



- Thus using the right distance/similarity is very important in clustering
- In some sense, related to asking: "Clustering based on what"?



Clustering: Some Examples

- Document/Image/Webpage Clustering
- Image Segmentation (clustering pixels)



- Clustering web-search results
- Clustering (people) nodes in (social) networks/graphs
- .. and many more..



Picture courtesy: http://people.cs.uchicago.edu/~pff/segment/

K-means Clustering

Based on the following two steps (applied in alternating fashion until not converged)

1. Assign each input to the current closest mean

2. Recompute (update) the means

 At the very beginning, the means needs to be initialized (at random locations or using other schemes that we will see later)



K-means Clustering

Good initialization matters (some methods exist to pick good initialization, e.g., *K*-means++)

1. Randomly initialize \overline{K} locations (the means)

2. Using the current means, predict the cluster id for each input (closest mean)

3. Recompute the means using the predicted cluster id of each input

4. Go to step 2 if not converged

Similar to LwP but in LwP the labels are known so the means can be computed in a single step without iterating multiple times



K-means clustering with K = 2



K-means (a.k.a. Lloyd's) Algorithm: Formally

- Input: N inputs $x_1, x_2, ..., x_N$; $x_n \in \mathbb{R}^D$; desired no. of partitions K
- Desired Output: Cluster ids of these N inputs and K cluster means $\mu_1, \mu_2, ..., \mu_K$

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

 $z_n = k$ means $z_{nk} = 1$ in one-hot representation of z_n $z_n = \arg \min_{k \in \{1,...,K\}} ||x_n - \mu_k||^2$

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 algo can also be seen as
doing a compression by "quantization":
Representing each of the *N* inputs by
one of the
$$K < N$$
 means

Can be fixed by modeling each cluster by a probability distribution, such as Gaussian (e.g., Gaussian Mixture Model; will see later)

This basic *K*-means models each cluster by a single mean μ_k . Ignores size/shape of clusters

What Loss Function is *K*-means Optimizing?

• Define the distortion or "loss" for the cluster assignment of a single input x_n

No labels here. We are
measuring how much error we
will incur if we assign
$$x_n$$
 to its
nearest cluster mean
$$\ell(x_n, \mu, z_n) = \sum_{k=1}^{K} z_{nk} ||x_n - \mu_k||^2$$

$$k_n = [z_{n1}, z_{n2}, \dots, z_{nK}]$$
encoding of x_n and $z_{nk} = 1$
if x_n is assigned to cluster k

Notation:

wi

K

m

- X is $N \times D$ matrix of inputs, Z is $N \times K$ matrix denoting cluster ids (each row is a one-hot z_n)
- μ is $K \times D$ (each row contains the mean μ_k of cluster k)
- Total distortion over all inputs defines the K-means "loss function"

$$\ell(X, \mu, Z) = \sum_{n=1}^{N} \sum_{k=1}^{K} z_{nk} ||x_n - \mu_k||^2$$
-means problem viewed as a atrix factorization problem = $||X - Z\mu||_F^2$ $\|x_n - \mu_k\|^2$
= The K-means problem is to minimize this objective wrt. μ and Z (one-hot vector)

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

CS771: Intro to ML

X approximated by the

product of **Z** and μ

K-means Loss: Several Forms, Same Meaning!

Can write the K-means loss function in several ways

$$\mathcal{L}(\mathbf{X}, \mathbf{Z}, \boldsymbol{\mu}) = \sum_{n=1}^{N} ||\boldsymbol{x}_n - \boldsymbol{\mu}_{\boldsymbol{z}_n}||^2 \qquad \mathcal{L}(\mathbf{X}, \mathbf{Z}, \boldsymbol{\mu}) = \sum_{k=1}^{K} \sum_{\substack{n:\boldsymbol{z}_n = k \\ \text{within cluster variance}}} ||\boldsymbol{x}_n - \boldsymbol{\mu}_k||^2$$

$$\mathcal{L}(\mathbf{X}, \mathbf{Z}, \boldsymbol{\mu}) = \sum_{n=1}^{N} \sum_{k=1}^{K} z_{nk} ||\boldsymbol{x}_n - \boldsymbol{\mu}_k||^2 \qquad \mathcal{L}(\mathbf{X}, \mathbf{Z}, \boldsymbol{\mu}) = \underbrace{||\mathbf{X} - \mathbf{Z}\boldsymbol{\mu}||_F^2}_{\text{as matrix factorization}}$$

Note: Not just K-means but many unsup. learning algos try to minimize a distortion or reconstruction error for X by solving a problem of the form

 ${f Z}$ denotes a new representation of the inputs and ${m \mu}$ denotes the other parameters of the model

$$\{\hat{\mathbf{Z}}, \hat{\boldsymbol{\mu}}\} = rg\min_{\mathbf{Z}, \boldsymbol{\mu}} \mathcal{L}(\mathbf{X}, \mathbf{Z}, \boldsymbol{\mu})$$

Replacing the ℓ_2 (Euclidean) squared by

 ℓ_1 distances gives the *K*-medoids

algorithm (more robust to outliers)

Optimizing the K-means Loss Function

The K-means problem is

$$\{\hat{\mathbf{Z}}, \hat{\boldsymbol{\mu}}\} = \arg\min_{\mathbf{Z}, \boldsymbol{\mu}} \mathcal{L}(\mathbf{X}, \mathbf{Z}, \boldsymbol{\mu}) = \arg\min_{\mathbf{Z}, \boldsymbol{\mu}} \sum_{n=1}^{N} \sum_{k=1}^{K} z_{nk} ||\mathbf{x}_n - \boldsymbol{\mu}_k||^2$$

- Can't optimize it jointly for Z and μ . Let's try alternating optimization for Z and μ

Alternating Optimization for K-means Problem Fix μ as μ̂ and find the optimal Z as Given the current estimates of the K means μ₁, μ₂, ..., μ_K, find the optimal cluster ids z₁, z₂, ..., z_N for all the inputs = arg min L(X, Z, μ̂) (still not easy - next slide) Fix Z as and find the optimal μ as Given the current estimates of the optimal cluster ids z₁, z₂, ..., z_N, compute the K means μ₁, μ₂, ..., μ_K Fix Z as and find the optimal μ as μ̂ arg min L(X, Â, μ) Go to step 1 if not yet converged

Solving for Z

- Solving for Z with μ fixed at $\widehat{\mu}$

$$\hat{\mathbf{Z}} = \arg\min_{\mathbf{Z}} \mathcal{L}(\mathbf{X}, \mathbf{Z}, \hat{\boldsymbol{\mu}}) = \arg\min_{\mathbf{Z}} \sum_{n=1}^{N} \sum_{k=1}^{K} z_{nk} || \mathbf{x}_n - \hat{\mu}_k ||^2$$

Still not easy. **Z** is discrete and above is an NP-hard problem

- Combinatorial optimization: K^N possibilities for Z ($N \times K$ matrix with one-hot rows)
- Greedy approach: Optimize Z one row (z_n) at a time keeping all others z_n 's (and the cluster means $\mu_1, \mu_2, \dots, \mu_K$) fixed

$$\hat{\boldsymbol{z}}_n = \arg\min_{\boldsymbol{z}_n} \sum_{k=1}^{\mathcal{K}} z_{nk} ||\boldsymbol{x}_n - \hat{\mu}_k||^2 = \arg\min_{\boldsymbol{z}_n} ||\boldsymbol{x}_n - \hat{\mu}_{\boldsymbol{z}_n}||^2$$

- Easy to see that this is minimized by assigning x_n to the closest mean
 - This is exactly what the *K*-means (Lloyd's) algo does!



Solving for μ

• Solving for μ with Z fixed at \widehat{Z}

$$\hat{\boldsymbol{\mu}} = \arg\min_{\boldsymbol{\mu}} \mathcal{L}(\mathbf{X}, \hat{\mathbf{Z}}, \boldsymbol{\mu}) = \arg\min_{\boldsymbol{\mu}} \sum_{k=1}^{K} \sum_{n:\hat{z}_n = k} ||\boldsymbol{x}_n - \boldsymbol{\mu}_k||^2$$

Not difficult to solve (each μ_k is a real-valued vector, can optimize easily)

$$\hat{\mu}_k = \arg\min_{\mu_k} \sum_{n:\hat{z}_n=k} ||\boldsymbol{x}_n - \mu_k||^2$$

- Note that each μ_k can be optimized for independently
- (Verify) This is minimized by setting $\hat{\mu}_k$ to be mean of points currently in cluster k
 - This is exactly what the *K*-means (Lloyd's) algo does!

Convergence of K-means algorithm

- Each step (updating Z or μ) can never <u>increase</u> the *K*-means loss
- When we update Z from $Z^{(t-1)}$ to $Z^{(t)}$

 $\mathcal{L}(\mathbf{X}, \mathbf{Z}^{(t)}, \boldsymbol{\mu}^{(t-1)}) \leq \mathcal{L}(\mathbf{X}, \mathbf{Z}^{(t-1)}, \boldsymbol{\mu}^{(t-1)})$

• When we update $\boldsymbol{\mu}$ from $\boldsymbol{\mu}^{(t-1)}$ to $\boldsymbol{\mu}^{(t)}$

 $\mathcal{L}(\mathbf{X}, \mathbf{Z}^{(t)}, \boldsymbol{\mu}^{(t)}) \leq \mathcal{L}(\mathbf{X}, \mathbf{Z}^{(t)}, \boldsymbol{\mu}^{(t-1)})$

because
$$\mathsf{Z}^{(t)} = rgmin_{\mathsf{Z}} \mathcal{L}(\mathsf{X},\mathsf{Z},\boldsymbol{\mu}^{(t-1)})$$

because $\boldsymbol{\mu}^{(t)} = rgmin_{\boldsymbol{\mu}} \mathcal{L}(\mathbf{X}, \mathbf{Z}^{(t)}, \boldsymbol{\mu})$

• Thus the K-means algorithm monotonically decreases the objective





14

K-means: Choosing K

One way to select K for the K-means algorithm is to try different values of K, plot the K-means objective versus K, and look at the "elbow-point"



• Can also information criterion such as AIC (Akaike Information Criterion) $AIC = 2\mathcal{L}(\hat{\mu}, \mathbf{X}, \hat{\mathbf{Z}}) + KD$

and choose K which gives the smallest AIC (small loss + large K values penalized)

 More advanced approaches, such as nonparametric Bayesian methods (Dirichlet Process mixture models also used, not within K-means but with other clustering algos)

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

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

as cluster means

most likely to be selected

K-means: Hard vs Soft Clustering

■ F.

- 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







- When clusters overlap, soft assignment is preferable (i.e., probability of being assigned to each cluster: say K = 3 and for some point x_n , $p_1 = 0.7$, $p_2 = 0.2$, $p_3 = 0.1$)
- A heuristic to get soft assignments: Transform distances from clusters into prob.

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

ach cluster mean updates changes: $\mu_k = \frac{\sum_{n=1}^{N} \gamma_{nk} \mathbf{x}_n}{\sum_{n=1}^{N} \gamma_{nk}} \quad \text{(all points contribute fractionally)}$

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

- Basic idea: Replace the Eucl. distances in *K*-means by the kernelized versions $\underset{\text{between input } \mathbf{x}_n \text{ and} \\ \underset{\text{mean of cluster } k}{\text{between input } \mathbf{x}_n \text{ and} } ||\phi(\mathbf{x}_n) - \phi(\boldsymbol{\mu}_k)||^2 = ||\phi(\mathbf{x}_n)||^2 + ||\phi(\boldsymbol{\mu}_k)||^2 - 2\phi(\mathbf{x}_n)^\top \phi(\boldsymbol{\mu}_k)$
- 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

 $\underline{\text{Not}}$ the same as the ϕ mapping of the mean of the data points assigned to cluster k

$$\begin{split} \|\phi(\mu_{k})\|^{2} &= \phi(\mu_{k})^{\mathsf{T}}\phi(\mu_{k}) \\ &= \frac{1}{|\mathcal{C}_{k}|^{2}} \sum_{n:\mathbf{z}_{n}=k} \sum_{n:\mathbf{z}_{m}=k} k(\mathbf{x}_{n},\mathbf{x}_{m}) \\ \phi(\mathbf{x}_{n})^{\mathsf{T}}\phi(\mu_{k}) &= \frac{1}{|\mathcal{C}_{k}|} \sum_{m:\mathbf{z}_{m}=k} k(\mathbf{x}_{n},\mathbf{x}_{m}) \end{split}$$

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

Kernel KMeans

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)



Hierarchical Clustering

Can be done in two ways: Agglomerative or Divisive

Agglomerative: Start t = 0 t = 1 t = 2 t = 3 t = 4 Agglomerative with each point being in Keep recursing until the desired a singleton cluster number of clusters found At each step, greedily merge At each step, break a cluster two most "similar" sub-clusters into (at least) two smaller X, homogeneous sub-clusters Stop when there is a single х, cluster containing all the points Divisive: Start with all points being in a single cluster Learns a dendrogram-like structure with inputs at the leaf Divisive nodes. Can then choose how t = 4t = 3t = 2t = 1t = 0 Tricky because no labels many clusters we want (unlike Decision Trees)

- Agglomerative is more popular and simpler than divisive (the latter usually needs complicated heuristics to decide cluster splitting).
- Neither uses any loss function

Similarity between two clusters (or two set of points) is needed in HC algos (e.g., this can be average pairwise similarity between the inputs in the two clusters)



Clustering vs Classification

- Any clustering model (prob/non-prob) typically learns two type of quantities
 - Parameters Θ of the clustering model (e.g., cluster means in K-means)
 - Cluster assignments $Z = \{z_1, z_2, ..., z_N\}$ for the points
- If cluster assignments Z were known, learning the parameters Θ is just like learning the parameters of a classifn model (typically generative classification) using labeled data
- Thus helps to think of clustering as (generative) classification with unknown labels
- Therefore many clustering problems are typically solved in the following fashion
 - 1. Initialize Θ somehow
 - 2. Predict Z given current estimate of Θ
 - 3. Use the predicted Z to improve the estimate of Θ (like learning a generative classification model)
 - 4. Go to step 2 if not converged yet

Clustering can help supervised learning, too

- Often "difficult" sup. learning problems can be seen as mixture of simpler models
- Example: Nonlinear regression or nonlinear classification as mixture of linear models



- Don't know which point should be modeled by which linear model \Rightarrow Clustering
- Can therefore solve such problems as follows
 - Initialize each linear model somehow (maybe randomly)
 - Cluster the data by assigning each point to its "closest" linear model (one that gives lower error)

Such an approach is also an example of divide and conquer

CS771: Intro to ML

and is also known as "mixture of experts" (will see it more

formally when we discuss latent variable models)

• (Re-)Learn a linear model for each cluster's data. Go to step 2 if not converged.

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, distortion or loss on test data etc



Evaluating Clustering Algorithms

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



 Rand Index (RI): Can also look at what fractions of pairs of points with same (resp. different) label are assigned to same (resp. different) cluster



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.,

 $z_n = \begin{bmatrix} 1 & 0 & 0 & 1 & 0 \end{bmatrix} \xrightarrow{\text{K=5 clusters with point } x_n \text{ belonging (in whole, not in terms of probabilities) to clusters 1 and 4}$

- 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, 2977)1: Intro to ML