Introduction to Machine Learning



Goals of Supervised Learning

Learn a hypothesis h from labeled datase $D = \{x^{(i)}, y^{(i)}\}_{i=1}^{n}$ that has low error on unseen data.





Goals of Reinforcement Learning

Find a "policy" $\pi: S \to A$ that maximizes reward in an environment.



Goals of Unsupervised Learning?

Clustering



Cluster Dendrogram



https://www.datanovia.com/en/lessons/agglomerative-hierarchical-clustering/

Goals of Unsupervised Learning?

Representation Learning





(Lab 7)

The ML Landscape

- **Supervised Learning**: Train a model that performs well (e.g., high accuracy) on unseen data.
- **Reinforcement Learning:** Learn a policy that maximizes expected reward in some environment.
- Unsupervised Learning:
 - Extract useful insights from data
 - Learn features for downstream tasks

Clustering

- Find a mapping from each data point to a cluster.
- Modeling choices:
 - How many clusters?
 - How do we define "close"?
 - How do we know if we have succeeded?



K-Means Clustering Objective

$$\arg\min_{\mu,y} \sum_{i=1}^{n} \sum_{j=1}^{k} \mathbf{1}\{y^{(i)} = j\} \|x^{(i)} - \mu^{(j)}\|_{2}^{2}$$

Find K cluster assignments and cluster means such that across all data points, the squared Euclidean distance between the data point and the cluster mean of its assigned cluster is minimized.



Input:

- 1. Data points $\{x^{(i)}\}_{i=1}^n$
- 2. Number of clusters k
- 3. Number of iterations

Output:

- 1. An assignment of each point to a cluster.
- 2. A "centroid" of each cluster with which to assign new points.

Input:

- 1. Data points $\{x^{(i)}\}_{i=1}^n$
- 2. Number of clusters k
- 3. Number of iterations

Output:

- 1. An assignment of each point to a cluster.
- 2. A "centroid" of each cluster with which to assign new points.

k-means (k, τ) Init $\{\mu^{(j)}\}_{i=1}^k, \{y^{(i)}\}_{i=1}^n$ for t = 1 to τ $y_{\text{old}} = y$ **for** i = 1 to n $y^{(i)} = \sup_{j} ||x^{(i)} - \mu^{(j)}||_2^2$ for j = 1 to k $\mu^{(j)} =$ $\frac{\sum_{i=1}^{n} \mathbf{1}\{y^{(i)} = j\}x^{(i)}}{\sum_{i=1}^{n} \mathbf{1}\{y^{(i)} = j\}}$ if $y = y_{\text{old}}$ break return $\{\mu^{(j)}\}_{j=1}^k, \{y^{(i)}\}_{i=1}^n$



```
k-means (k, \tau)
  Init \{\mu^{(j)}\}_{j=1}^k, \{y^{(i)}\}_{i=1}^n
 for t = 1 to \tau
      y_{\text{old}} = y
     for i = 1 to n
       y_{\arg\min_{i}}^{(i)} = \|x^{(i)} - \mu^{(j)}\|_{2}^{2}
     for j = 1 to k
        \mu^{(j)} =
           \frac{\sum_{i=1}^{n} \mathbf{1}\{y^{(i)} = j\}x^{(i)}}{\sum_{i=1}^{n} \mathbf{1}\{y^{(i)} = j\}}
     if y = y_{\text{old}}
         break
 \texttt{return} \{\mu^{(j)}\}_{i=1}^k, \{y^{(i)}\}_{i=1}^n
```



k-means(k, τ) Init $\{\mu^{(j)}\}_{j=1}^k, \{y^{(i)}\}_{i=1}^n$ for t = 1 to τ $y_{\text{old}} = y$ K = 5**for** i = 1 to n $y^{(i)} = rg\min_{i} \|x^{(i)} - \mu^{(j)}\|_2^2$ for j = 1 to k $\mu^{(j)} =$ $\frac{\sum_{i=1}^{n} \mathbf{1}\{y^{(i)} = j\}x^{(i)}}{\sum_{i=1}^{n} \mathbf{1}\{y^{(i)} = j\}}$ if $y = y_{\text{old}}$ break $\texttt{return} \{\mu^{(j)}\}_{i=1}^k, \{y^{(i)}\}_{i=1}^n$







k-means (k, τ) Init $\{\mu^{(j)}\}_{j=1}^k, \{y^{(i)}\}_{i=1}^n$ for t = 1 to τ $y_{\text{old}} = y$ **for** i = 1 to n $y_{\arg\min_{i}}^{(i)} = \|x^{(i)} - \mu^{(j)}\|_{2}^{2}$ for j = 1 to k $\mu^{(j)} =$ $\frac{\sum_{i=1}^{n} \mathbf{1}\{y^{(i)} = j\}x^{(i)}}{\sum_{i=1}^{n} \mathbf{1}\{y^{(i)} = j\}}$ if $y = y_{\text{old}}$ break $\texttt{return} \{\mu^{(j)}\}_{i=1}^k, \{y^{(i)}\}_{i=1}^n$



k-means(k,
$$\tau$$
)
Init $\{\mu^{(j)}\}_{j=1}^k, \{y^{(i)}\}_{i=1}^n$
for t = 1 to τ
 $y_{\text{old}} = y$
for i = 1 to n
 $y^{(i)} =$
 $\arg\min_j ||x^{(i)} - \mu^{(j)}||_2^2$
for j = 1 to k
 $\mu^{(j)} =$
 $\frac{\sum_{i=1}^n \mathbf{1}\{y^{(i)} = j\}x^{(i)}}{\sum_{i=1}^n \mathbf{1}\{y^{(i)} = j\}}$
if $y = y_{\text{old}}$
break
return $\{\mu^{(j)}\}_{j=1}^k, \{y^{(i)}\}_{i=1}^n$



k-means (k, τ) Init $\{\mu^{(j)}\}_{j=1}^k, \{y^{(i)}\}_{i=1}^n$ for t = 1 to τ $y_{\text{old}} = y$ **for** i = 1 to n $y_{rgmin_{j}}^{(i)} = \sup_{i} \|x^{(i)} - \mu^{(j)}\|_{2}^{2}$ for j = 1 to k Assign each data $\mu^{(j)} =$ point to closest $\sum_{i=1}^{n} \mathbf{1}\{y^{(i)} = j\}x^{(i)}$ "centroid" $\sum_{i=1}^{n} \mathbf{1}\{y^{(i)} = j\}$ if $y = y_{\text{old}}$ break $\texttt{return} \{\mu^{(j)}\}_{i=1}^k, \{y^{(i)}\}_{i=1}^n$





k-means (k, τ) Init $\{\mu^{(j)}\}_{i=1}^k, \{y^{(i)}\}_{i=1}^n$ for t = 1 to τ $y_{\text{old}} = y$ **for** i = 1 to n $y^{(i)} = rg\min_{i} \|x^{(i)} - \mu^{(j)}\|_2^2$ for j = 1 to k $\mu^{(j)} =$ $\frac{\sum_{i=1}^{n} \mathbf{1}\{y^{(i)} = j\}x^{(i)}}{\sum_{i=1}^{n} \mathbf{1}\{y^{(i)} = j\}}$ if $y = y_{\text{old}}$ break $\texttt{return} \{\mu^{(j)}\}_{i=1}^k, \{y^{(i)}\}_{i=1}^n$







k-means (k,
$$\tau$$
)
Init $\{\mu^{(j)}\}_{j=1}^{k}, \{y^{(i)}\}_{i=1}^{n}$
for t = 1 to τ
 $y_{old} = y$
for i = 1 to n
 $y^{(i)} =$
 $\arg\min_{j} ||x^{(i)} - \mu^{(j)}||_{2}^{2}$
for j = 1 to k
 $\mu^{(j)} =$
 $\frac{\sum_{i=1}^{n} \mathbf{1}\{y^{(i)} = j\}x^{(i)}}{\sum_{i=1}^{n} \mathbf{1}\{y^{(i)} = j\}}$
if $y = y_{old}$
break
return $\{\mu^{(j)}\}_{j=1}^{k}, \{y^{(i)}\}_{i=1}^{n}$



k-means (k, τ) Init $\{\mu^{(j)}\}_{j=1}^k, \{y^{(i)}\}_{i=1}^n$ for t = 1 to τ $y_{\text{old}} = y$ **for** i = 1 to n $y_{rgmin_{i}}^{(i)} = \sup_{i} ||x^{(i)} - \mu^{(j)}||_{2}^{2}$ for j = 1 to k $\mu^{(j)} =$ $\frac{\sum_{i=1}^{n} \mathbf{1}\{y^{(i)} = j\}x^{(i)}}{\sum_{i=1}^{n} \mathbf{1}\{y^{(i)} = j\}}$ if $y = y_{\text{old}}$ break $\texttt{return} \{\mu^{(j)}\}_{i=1}^k, \{y^{(i)}\}_{i=1}^n$



data



```
k-means (k, \tau)
                        Init \{\mu^{(j)}\}_{j=1}^k, \{y^{(i)}\}_{i=1}^n
                        for t = 1 to \tau
                            y_{\text{old}} = y
                            for i = 1 to n
                             y^{(i)} = rg\min_{i} \|x^{(i)} - \mu^{(j)}\|_2^2
                            for j = 1 to k
                               \mu^{(j)} =
Use centroid to \sum_{i=1}^{n} \mathbf{1}\{y^{(i)} = j\}x^{(i)}
cluster new \sum_{i=1}^{n} \mathbf{1}\{y^{(i)} = j\}
                           if y = y_{\text{old}}
                                break
                        return \{\mu^{(j)}\}_{j=1}^k, \{y^{(i)}\}_{i=1}^n
```

K-means is sensitive to initialization!



K-means is guaranteed to converge with iterations, but not necessarily to the global minimum.

K-means is sensitive to initialization! And number of clusters Why can't we just increase K?



K-means is sensitive to initialization! And number of clusters k And choice of distance metric

What are some issues with this "distance"?

k-means (k, τ) Init $\{\mu^{(j)}\}_{j=1}^k, \{y^{(i)}\}_{i=1}^n$ for t = 1 to τ $y_{\text{old}} = y$ **for** i = 1 to n $u^{(i)} =$ $\arg\min_{i} \|x^{(i)} - \mu^{(j)}\|_{2}^{2}$ for j = 1 to k $\mu^{(j)} =$ $\frac{\sum_{i=1}^{n} \mathbf{1}\{y^{(i)} = j\}x^{(i)}}{\sum_{i=1}^{n} \mathbf{1}\{y^{(i)} = j\}}$ if $y = y_{\text{old}}$ break $\texttt{return} \{\mu^{(j)}\}_{i=1}^k, \{y^{(i)}\}_{i=1}^n$





When to use K-Means?

K-means works well when:

- Data "circular"



When to use K-Means?

K-means works well when:

- Data "circular"
- Clusters have roughly the same size





When to use K-Means?

K-means works well when:

- Data "circular"
- Clusters have roughly the same size
- Clusters are well separated



