

# clustering and k-means

Week 7 -- Monday

# Clustering = Unsupervised ML Technique

- The difference between supervised and unsupervised machine learning: There are no labels (outputs; target) in the dataset.
- The goal for unsupervised learning is to model the underlying structure or distribution in the data in order to learn more about the data.
- There are several clustering algorithms, including k-means, mean-shift, hierarchical clustering, and DBSCAN
- The various algorithms differ significantly in their notion of what constitutes a cluster and how to efficiently find them (as stated in Wikipedia article here: [https://en.wikipedia.org/wiki/Cluster\\_analysis](https://en.wikipedia.org/wiki/Cluster_analysis))

## A comparison of the clustering algorithms in scikit-learn

- <http://scikit-learn.org/stable/modules/clustering.html>

Method name	Parameters	Scalability	Usecase	Geometry (metric used)
<b>K-Means</b>	number of clusters	Very large <code>n_samples</code> , medium <code>n_clusters</code> with <b>MiniBatch</b> code	General-purpose, even cluster size, flat geometry, not too many clusters	Distances between points
<b>Affinity propagation</b>	damping, sample preference	Not scalable with <code>n_samples</code>	Many clusters, uneven cluster size, non-flat geometry	Graph distance (e.g. nearest-neighbor graph)
<b>Mean-shift</b>	bandwidth	Not scalable with <code>n_samples</code>	Many clusters, uneven cluster size, non-flat geometry	Distances between points
<b>Spectral clustering</b>	number of clusters	Medium <code>n_samples</code> , small <code>n_clusters</code>	Few clusters, even cluster size, non-flat geometry	Graph distance (e.g. nearest-neighbor graph)
<b>Ward hierarchical clustering</b>	number of clusters	Large <code>n_samples</code> and <code>n_clusters</code>	Many clusters, possibly connectivity constraints	Distances between points
<b>Agglomerative clustering</b>	number of clusters, linkage type, distance	Large <code>n_samples</code> and <code>n_clusters</code>	Many clusters, possibly connectivity constraints, non Euclidean distances	Any pairwise distance
<b>DBSCAN</b>	neighborhood size	Very large <code>n_samples</code> , medium <code>n_clusters</code>	Non-flat geometry, uneven cluster sizes	Distances between nearest points
<b>Gaussian mixtures</b>	many	Not scalable	Flat geometry, good for density estimation	Mahalanobis distances to centers
<b>Birch</b>	branching factor, threshold, optional global clusterer.	Large <code>n_clusters</code> and <code>n_samples</code>	Large dataset, outlier removal, data reduction.	Euclidean distance between points

# K-Means clustering

- Works with numeric data only!
- Algorithm:
  1. Pick a number  $k$  of random cluster centers
  2. Assign every item to its nearest cluster center using a distance metric
  3. Move each cluster center to the mean of its assigned items
  4. Repeat 2-3 until convergence (change in cluster assignment less than a threshold)

# K-Means clustering: Pros and Cons

- Pros: very efficient (even if multiple runs are performed), can be used for a large variety of data types, easy to understand
- Cons: Susceptible to initialization problems and outliers, restricted to data in which there is a notion of a center, have to specify a k (different starting points may give you different clusters--you won't necessarily get an optimal cluster), requires prior knowledge or assumption about number of clusters

# 1) Format and preprocess data for clustering

- Handle missing values (either impute or remove any NaNs)
- Convert any categorical data to numeric (nominal variables need to be dummified)
- Scale data
- Check out this great article on preprocessing data for k-means clustering:  
<http://www.edupristine.com/blog/k-means-algorithm>

## 2) Perform a K-Means clustering analysis

### --Selecting k

- **Random initialization** – k-Means clustering is prone to initial seeding i.e. random initialization of centroids which is required to kick-off iterative clustering process. Bad initialization may end up getting bad clusters.

So how do we choose the right k?

## --Selecting k (cont'd)

- Visualize the data for a clue
- Trial and Error (using performance evaluation)
- Domain Knowledge or Literature Review
- Guess

Check out this great article about various methods to select k:

<http://www.edupristine.com/blog/beyond-k-means>



## 2) Perform a K-Means clustering analysis

- Set k
- Call instance of algorithm with k
  - `kmeans = cluster.KMeans(n_clusters=k)`
- Fit model
  - `kmeans.fit(dfn)`
- Get labels and centroids:
  - `labels = kmeans.labels_`
  - `centroids = kmeans.cluster_centers_`

### 3) Evaluate clusters for fit

- Accuracy Score
- Silhouette Score
- Confusion Matrix
- Classification Report
- Visualization of Results