

Unsupervised Learning and K-Means Clustering

Data Science Dojo

Unsupervised Learning (1/5)

- Trying to find hidden structure in unlabeled data
- No error or reward signal to evaluate a potential solution. *No need to pick a response class.*
- Common techniques: **K-Means clustering**, hierarchical clustering, hidden Markov models, etc.
 - It has a long history, and used in almost every field, e.g., medicine, psychology, botany, sociology, biology, archeology, marketing, insurance, libraries, etc.

Unsupervised Learning (2/5)

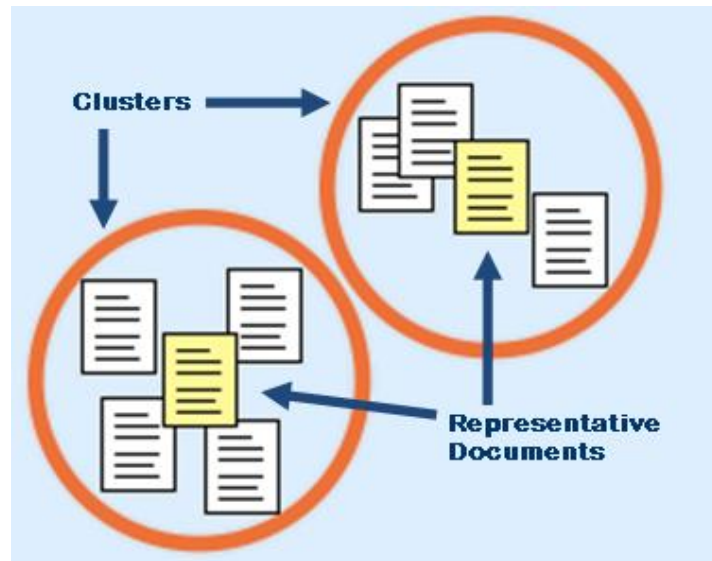
Example 1: Clothing size

- Tailor-made for each person is too expensive
- One-size-fits-all: does not work!
- Groups people of similar sizes together to make "small", "medium", and "large" t-shirts

Unsupervised Learning (3/5)

Example 2: Text document organization

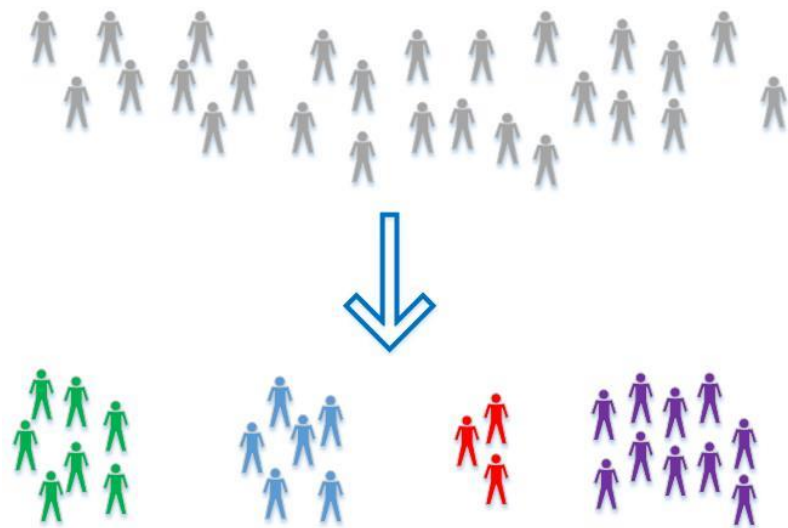
- To find groups of documents that are similar to each other based on the important terms appearing in them



Unsupervised Learning (4/5)

Example 3: Target Marketing

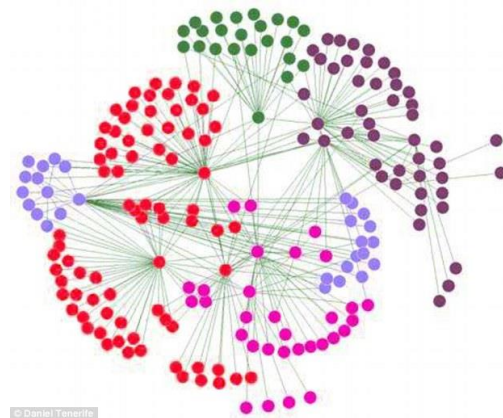
- Subdivide market into distinct subsets of customers
- where any subset may conceivably be selected as a segment to be reached with a particular offer



Unsupervised Learning (5/5)

Example 4: Social network graphs

- Subdivide social network into distinct subsets of user groups (Facebook friends, LinkedIn contacts...)
- Group users by their similar characteristics



K-Means Clustering

- Partitions data points into similarity clusters
- Unsupervised technique: there is no partitioning into a learning or a test set in unsupervised learning
- Useful in grouping observations
- Only works for numeric data

K-Means Clustering

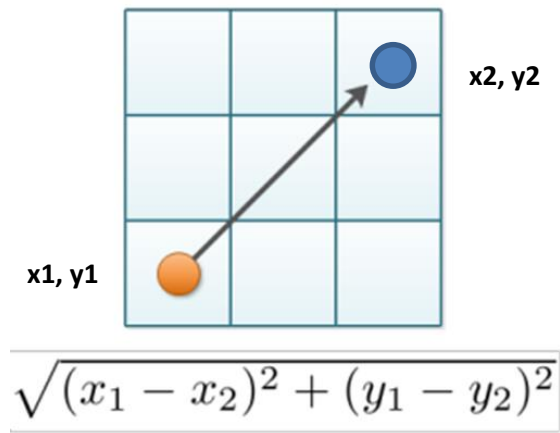
- Transform categorical variables into numeric
- Datasets will become wide quickly
- Needed to compute similarity

Often called
“dummy variables” or
“one-hot encoding”

Age	Pclass.1	Pclass.2	Pclass.3	Sex.female	Sex.male
19	0	1	0	0	1
28	1	0	0	1	0
64	0	0	1	0	1

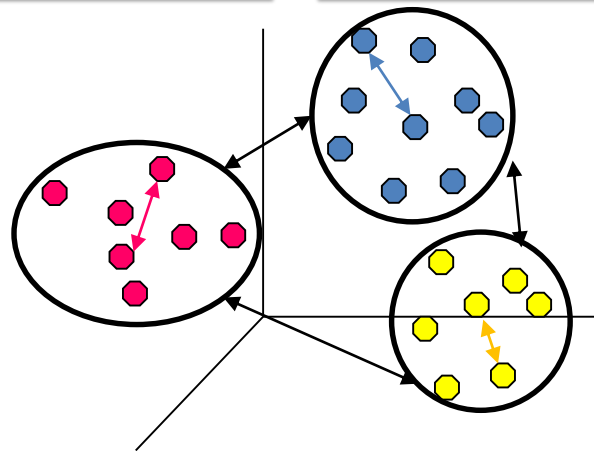
Euclidean Distance

points in a two-dimensional space to determine intra- and inter-cluster similarity

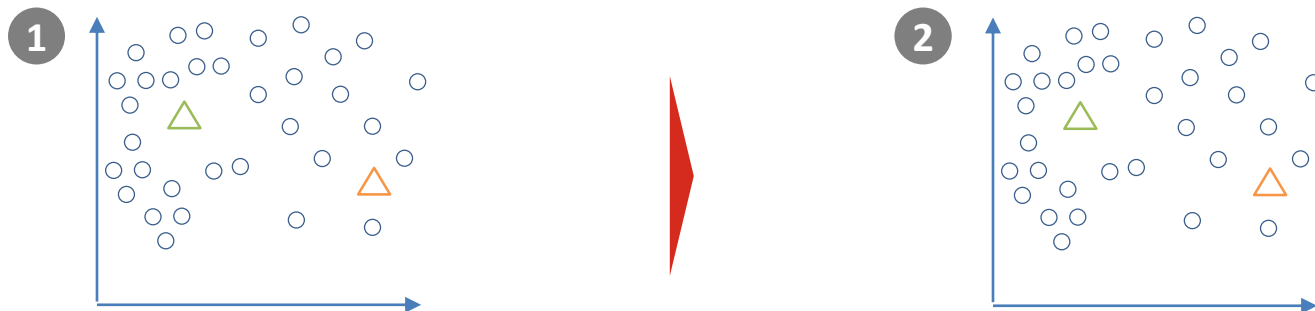


Intra-cluster distances
are minimized

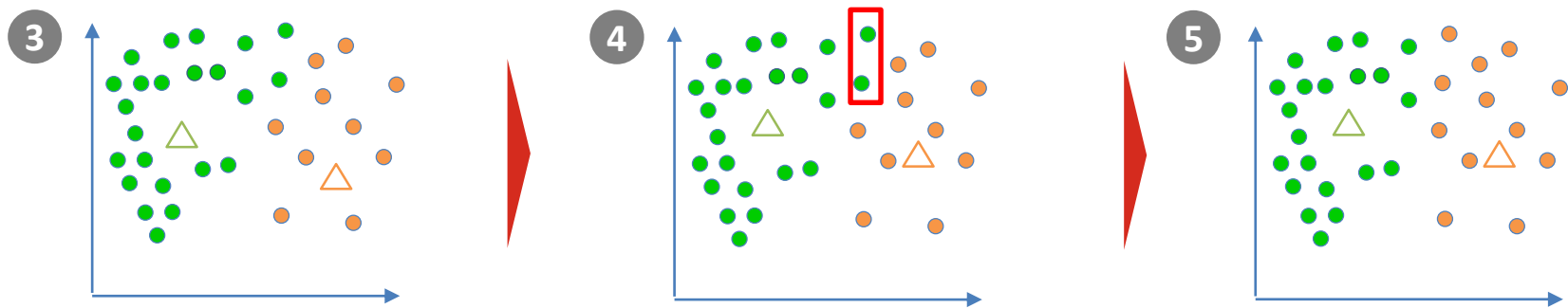
Inter-cluster distances
are maximized



K-Means Clustering (1/2)



K-Means Clustering (2/2)



The positions of the cluster centers are determined by the mean of all the points in the cluster.

K-Means Clustering Algorithm

Suppose set of data points: $\{x_1, x_2, x_3, \dots, x_n\}$

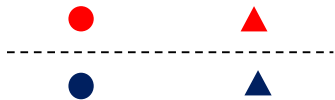
- **Step 1:** Decide the number of clusters, $K=1,2,\dots,k$.
- **Step 2:** Place centroids at random locations
 - c_1, c_2, \dots, c_k
- **Step 3:** Repeat until convergence:
 - { for each point x_i → find nearest centroid c_j (eg. Euclidean distance)
→ assign the point x_i to cluster j
 - for each cluster $j = 1 \dots k$ → calculate new centroid c_j
 $c_j = \text{mean of all points } x_i \text{ assigned to cluster } j \text{ in previous step}$
 - }
- **Step 4:** Stop when none of the cluster assignments change

K-Means Clustering

- Minimizes aggregate intra-cluster distance
 - Measure squared distance from point to center of its cluster.

$$\sum_{j=1}^K \sum_{x \in g_j} D(c_j, x)^2$$

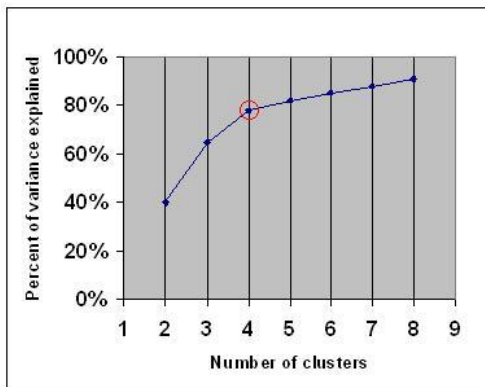
- Could converge to local minimum
 - Different starting points → very different results
 - Run many times with random starting points
- Nearby points may not be assigned to the same cluster



K-Means Clustering

- Strengths
 - Simple: easy to understand and to implement
 - Efficient: linear time, minimal storage
- Weaknesses
 - Mean must be well defined
 - The user needs to specify k
 - Algorithm is sensitive to outliers

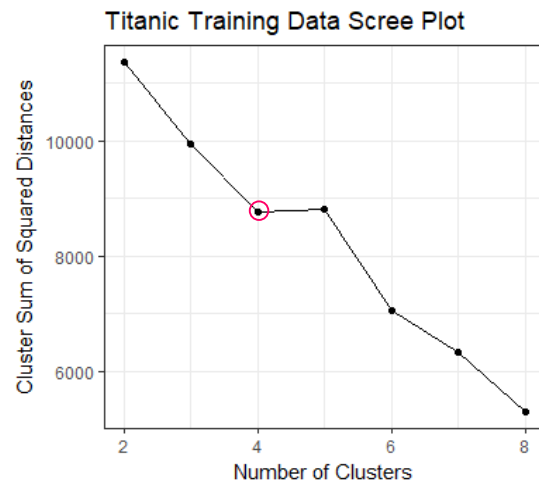
Finding K with Elbow Method



Option 1 - Percentage of variance explained as a function of the number of clusters.

Option 2 - Total of the squared distances of cluster point to center.

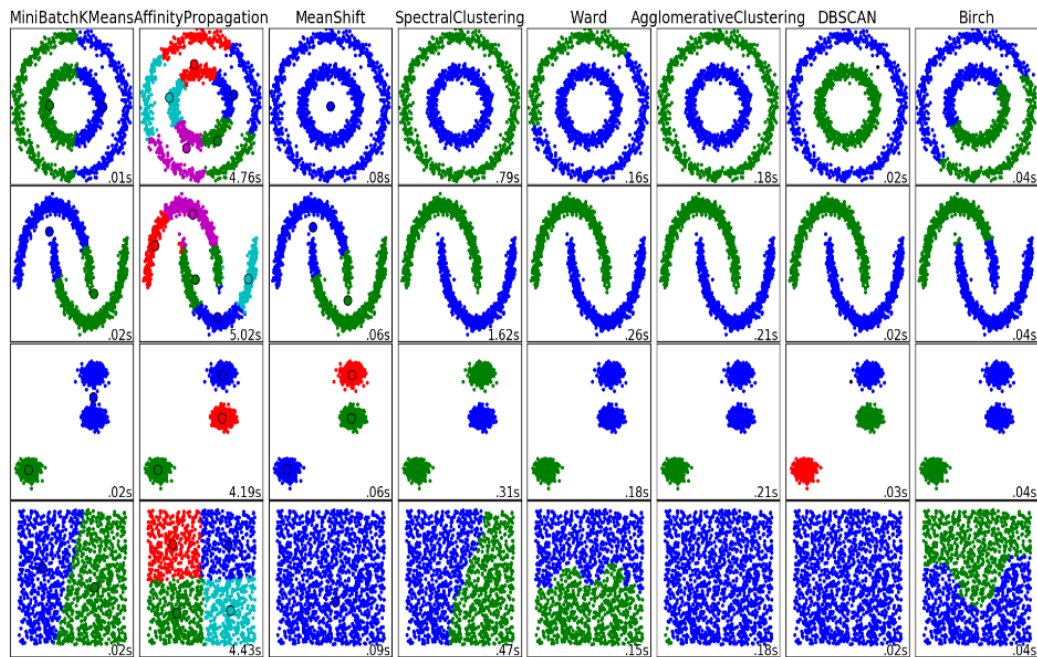
Goal - Choose a number of clusters so that adding another cluster doesn't give much better modeling of the data.



Other K Optimization Techniques

- Silhouette
- Calinski criterion
- Bayesian Information Criterion
- Affinity propagation (AP) clustering
- Gap statistic

Comparing Clustering Algorithms



A comparison of the clustering algorithms in scikit-learn

QUESTIONS