



McCOMBS SCHOOL OF BUSINESS

Salem Center for Policy

Unsupervised learning

David Puelz

October 12, 2021



Clustering

Principal Components Analysis



1. Introduction to clustering
2. K-means
3. Implementing K-means: some practical details
4. Hierarchical clustering



You've seen a lot of models for $p(y | x)$. This is supervised learning (knowing outcomes $y = \text{"supervision"}$).

The next few topics are all about models for x alone.

- **Clustering** means dividing data points into categories which are not defined in advance.
- It's different from classification: dividing data points into categories which *are* defined in advance, and for which we have explicit labels.

Clustering: Toy example



- The horizontal and vertical locations of each point give its location $x_i = (x_{i1}, x_{i2}, \dots, x_{ip})$ in **feature** space.
- From these locations, we impute the cluster labels: $\gamma_i = k$ if point i is in cluster $k \in \{1, \dots, K\}$.

Clusters should partition the data:

- Partition = mutually exclusive, collectively exhaustive set of categories (each data point is in one and only one cluster).
- No “mixed membership.” If you encounter a **Chiweenie**, you need a new cluster!



Criteria for clustering



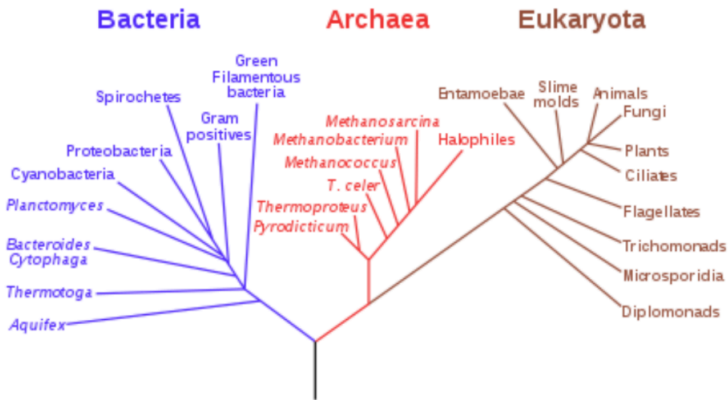
Data points within the same cluster should be close/similar, and data points in different clusters should be far apart/dissimilar.



Criteria for clustering



Clusters should (ideally) be **balanced**. A sensible clustering of living creatures:



Criteria for clustering



A less sensible clustering of living creatures:



1



2

All other living
creatures

3

How can we cluster without labels?



There are many algorithms which do this, i.e. try to find clusters that are homogenous, well separated, balanced, etc.

Key fact: we need to know about **distances** to quantify similarity (within a cluster) and difference (between clusters).

Generically, if x_i and x_j are two data points, we let $d(x_i, x_j)$ denote the distance between them.



Properties of distance functions:

1. $d(x_i, x_j) \geq 0$
2. $d(x_i, x_j) = 0 \iff x_i = x_j$.
3. $d(x_i, x_j) = d(x_j, x_i)$ (symmetry)
4. $d(x_i, x_j) \leq d(x_i, x_k) + d(x_j, x_k)$ (triangle inequality, i.e. “If you want to get from Austin to Houston, don’t go through Dallas!”)

NB: in math, distance functions are called “metrics.”



Euclidean (ℓ^2):

$$d(x_i, x_j) = \left[\sum_{d=1}^D (x_{id} - x_{jd})^2 \right]^{1/2}$$

(just the Pythagorean theorem!)

Manhattan (ℓ^1):

$$d(x_i, x_j) = \sum_{d=1}^D |x_{id} - x_{jd}|$$

(also called “taxicab” distance)



The x points can be arbitrary objects in potentially crazy spaces:

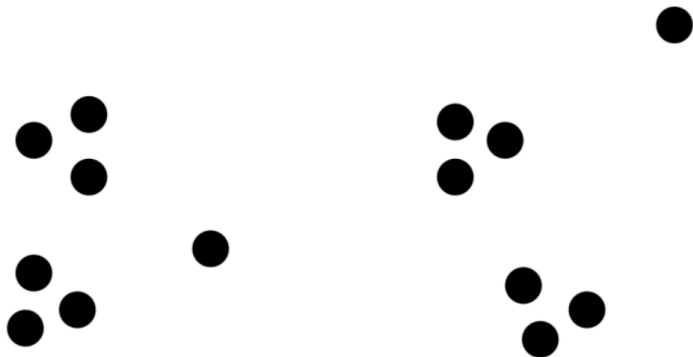
- playlists on Spotify
- phrase counts for books
- economic indicators
- DNA sequences (“Hamming distance”)
- etc.

As long as we can measure the distance between any two points, we can cluster them!

A few miscellaneous notes



Clusters can be ambiguous:



How many clusters do you see?



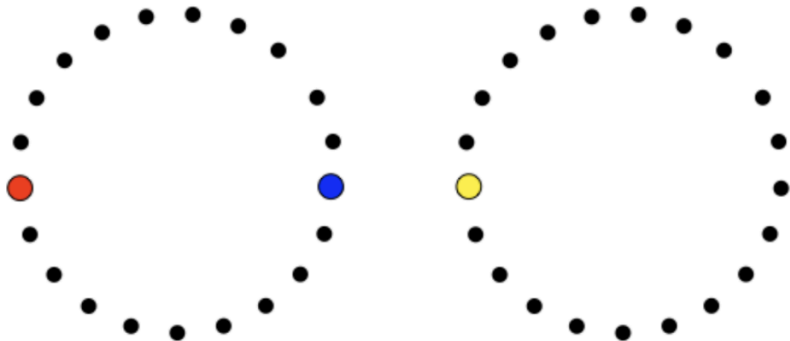
We can organize algorithms into two broad classes:

1. “**hierarchical**” clustering methods. Think the tree of life! Mammals and reptiles are both vertebrates. Vertebrates and invertebrates are both animals. Animals and plants are both eukaryotes. Etc.
2. “**partitional**” clustering methods. Here there is no tree or hierarchy, just a “flat” set of clusters.

A few miscellaneous notes



Distance-based clustering isn't magic.



Should blue cluster with red or with yellow?

(We can often deal with situations like this by redefining what distance is. The term here is “manifold learning.”)



- K-means is a partitional clustering approach. It's the "Least Squares Regression of clustering"
- Each cluster is associated with a centroid (center point). The number of clusters K must be chosen in advance
- Each point is assigned to the cluster with the closest centroid



The basic algorithm is super simple:

-
- 1: Select K points as the initial centroids.
 - 2: **repeat**
 - 3: Form K clusters by assigning all points to the closest centroid.
 - 4: Recompute the centroid of each cluster.
 - 5: **until** The centroids don't change
-

That's it! Algorithms don't get any simpler in machine learning