

# What is Clustering?

- Given a set of points, with a notion of distance between points, group the points into some number of *clusters*, so that
  - Members of the same cluster are close/similar to each other
  - Members of different clusters are dissimilar
- Usually:
  - Points are in a high-dimensional space
  - Similarity is defined using a distance measure
  - Euclidean, Cosine, Jaccard, edit distance, ...



# A Different Example

- How would you group
  - 'The price of crude oil has increased significantly'
  - 'Demand for crude oil outstrips supply'
  - 'Some people do not like the flavor of olive oil'
  - 'The food was very oily'
  - 'Crude oil is in short supply'
  - 'Oil platforms extract oil'
  - 'Canola oil is supposed to be healthy'
  - 'Iraq has significant oil reserves'
  - 'There are different types of cooking oil'

A note: you might or might not know how many clusters to look for.

# A Different Example

- How would you group
  - 'The price of crude oil has increased significantly'
  - 'Demand for crude oil outstrips supply'
  - 'Some people do not like the flavor of olive oil'
  - 'The food was very oily'
  - 'Crude oil is in short supply'
  - 'Oil platforms extract oil'
  - 'Canola oil is supposed to be healthy'
  - 'Iraq has significant oil reserves'
  - 'There are different types of cooking oil'











# Sample problem: Music

- Intuitively: Music can be divided into categories, and customers prefer a few categories
  - But what are categories really?
- Represent a CD by a set of customers who bought it
- Similar CDs have similar sets of customers, and vice-versa
- But what are the characteristics of CDs that we care about?
  - Genre?
  - Popularity?
  - Cover color?

# Music example

- Think of a space with one dimension for each customer
  - Values in a dimension may be 0 or 1 only
  - A CD is a "point" in this space (x<sub>1</sub>, x<sub>2</sub>,..., x<sub>d</sub>), where x<sub>i</sub> = 1 iff the *i* th customer bought the CD
- For Amazon, the dimension is tens of millions
- Task: Find clusters of similar CDs

# **Clustering Basics**

- Collect examples
- Compute **similarity** among examples according to some metric
- Group examples together such that:
  - Examples within a cluster are similar
  - Examples in different clusters are different
- Sometimes: Summarize each cluster
- Sometimes: Assign new instances to the most similar cluster

# Measures of Similarity

- To do clustering we need some measure of similarity.
- This is basically our "critic"
- Computed over a vector of values representing instances
- Types of values depend on domain:
  - Documents: bag of words, linguistic features
  - Purchases: cost, purchaser data, item data
  - Census data: most of what is collected
- Multiple different measures exist



# Measures of Similarity

- Semantic similarity (but that's hard)
  - For example, olive oil vs. crude oil—these are semantically different
- Similar attribute counts
  - Number of attributes with the same value
  - Appropriate for large, sparse vectors
  - Bag-of-Words: BoW
- (Slightly) more complex vector comparisons:
  - Euclidean Distance
  - Cosine Similarity

#### **Euclidean Distance**

• Euclidean distance: distance between two measures summed across each feature

```
dist(x_i, x_j) = sqrt((x_{i1}-x_{j1})^2 + (x_{i2}-x_{j2})^2 + .. + (x_{in}-x_{jn})^2)
```

• Squared differences give more weight to larger differences

• dist([1,2],[3,8]) = sqrt((1-3)<sup>2</sup>+(2-8)<sup>2</sup>) =  
sqrt((-2)<sup>2</sup>+(-6)<sup>2</sup>) =  
sqrt(4+36) =  
sqrt(40) = 
$$\sim 6.3$$







# Euclidean Distance vs Cosine Similarity vs Other

- Cosine Similarity:
  - Measures relative proportions of various features
  - Ignores magnitude
  - When all the correlated dimensions between two vectors are in proportion, you get maximum similarity
- Euclidean Distance:
  - Measures actual distance between two points
  - More concerned with absolutes
- Either can deal with many dimensions
- Often similar in practice, especially on high dimensional data
- Consider meaning of features/feature vectors for your domain

Justin Washtell @ semanticvoid.com/blog/2007/02/23/similarity-measure-cosine-similarity-or-euclidean-distance-or-both/

19

#### What about non-Euclidean space?

- Is the space Euclidean or non-Euclidean?
- In Euclidean:
  - Points are vectors of real numbers, i.e. coordinates
  - It is possible to summarize a collection of points as their average ("centroid")
- Non-Euclidean:
  - There is no notion of location or centroid
  - We summarize a collection of points differently
  - Distance measures: Jaccard, Hamming, cosine

# **Clustering Algorithms**

- Flat:
  - K means
- Hierarchical:
  - Bottom up
  - Top down (not common)
- Probabilistic:
  - Expectation Maximization (E-M)

#### 21

#### Partitioning (Flat) Algorithms

- Partitioning method
  - Construct a **partition** of *n* instances into a set of *k* clusters
- Given: a set of documents and the number k
- Find: a partition of k clusters that optimizes the chosen partitioning criterion
  - Globally optimal: exhaustively enumerate all partitions.
  - Usually too expensive.
  - Effective heuristic methods: k-means algorithm.



# k-means Algorithm

- Choose k (the number of clusters)
- Randomly choose k instances to center clusters on
- Assign each point to the centroid it's closest to, forming clusters +
- Recalculate centroids of new clusters
- Reassign points based on new centroids
- Iterate until...
- Convergence (no point is reassigned) or after a fixed number of iterations.



# <section-header><section-header><list-item><list-item><list-item><list-item><list-item><list-item>









# Properties of k-means

- Guaranteed to converge in a finite number of iterations
- Running time per iteration:
- Assign data points to closest centroid: O(KN)
- Change the cluster centroid to the average of points: O(N)







#### Evaluation of k-means

- Advantages:
  - Easy to understand, implement
  - Most popular clustering algorithm
  - Efficient, almost linear
    - Time complexity: O(tkn)
      - n = number of data points
      - k = number of clusters
      - t = number of iterations.
  - In practice, performs well (especially on text)

- Disadvantages:
  - Must choose *k* beforehand
    - Bad  $k \rightarrow$  bad clusters
    - Sometimes we don't know k
  - Sensitive to initialization
    - One fix: run several times with different random centers and look for agreement
  - Sensitive to outliers, irrelevant features



# EM Clustering Algorithm

- Goal: maximize overall probability of data
- Iterate between:
  - Expectation: estimate probability that each instance belongs to each cluster
  - Maximization: recalculate parameters of probability distribution for each cluster
- Until convergence or iteration limit.

# Expectation Maximization (EM)

- Probabilistic method for soft clustering
- Idea: learn k classifications from unlabeled data
- Assumes k clusters: $\{c_1, c_2, \dots c_k\}$
- "Soft" version of k-means
- Assumes a probabilistic model of categories (such as Naive Bayes)
- Allows computing  $P(c\mathrm{i} \mid I)$  for each category,  $c\mathrm{i}$  , for a given instance I

# (Slightly) More Formally

- Iteratively learn probabilistic categorization model from unsupervised data
- Initially assume random assignment of examples to categories
  - "Randomly label" data
- Learn initial probabilistic model by estimating model parameters θ from randomly labeled data
- Iterate until convergence:
  - Expectation (E-step):
    - Compute  $P(c_i \mid I)$  for each instance (example) given the current model
    - Probabilistically re-label the examples based on these posterior probability estimates
  - **Maximization (M-step):** Re-estimate model parameters, θ, from re-labeled data











#### **EM** Summary

- Basically a probabilistic k-means.
- Has many of same advantages and disadvantages
  - Results are easy to understand
  - Have to choose k ahead of time
- Useful in domains when we want likelihood that an instance belongs to more than one cluster
  - Natural language processing for instance