Clustering

Nisheeth

What we want



What regression gives us







What causality inference gives us



What we're missing



We're missing the conditional probability tables

We can get them from empirical frequencies, but empirical frequencies of what?

Words shape reality





Need succinct phenomena descriptors



Categories compress information



Today

- Clustering
- Types of clusters
 - We focus on exclusive clusters
- Types of clustering algorithms
 - Distance-based
 - Contiguity-based
 - Density-based
 - Hierarchical
- Validation
- Suggested reading: Pan, Kumar & Steinbach Ch 8
 - Lots of slides drawn from that book chapter

Clustering



Distance measures

• Vectorize the data

- Turn each attribute into a binary label

• Use any of the following measures

- Euclidean $d(\mathbf{p}, \mathbf{q}) = d(\mathbf{q}, \mathbf{p}) = \sqrt{(q_1 - p_1)^2 + (q_2 - p_2)^2 + \dots + (q_n - p_n)^2}$

$$= \sqrt{\sum_{i=1}^{n} (q_i - p_i)^2}.$$

$$- \text{Cosine} \quad \text{similarity} = \cos(\theta) = \frac{\mathbf{A} \cdot \mathbf{B}}{\|\mathbf{A}\| \|\mathbf{B}\|} = \frac{\sum_{i=1}^{n} A_i B_i}{\sqrt{\sum_{i=1}^{n} A_i^2} \sqrt{\sum_{i=1}^{n} B_i^2}}$$

$$- \text{Manhattan} \quad d_1(\mathbf{p}, \mathbf{q}) = \|\mathbf{p} - \mathbf{q}\|_1 = \sum_{i=1}^{n} |p_i - q_i|,$$

Types of clusters

- Well-separated clusters
- Center-based clusters
- Contiguous clusters
- Density-based clusters
- Conceptual clusters

Well-separated clusters

- Well-Separated Clusters:
 - A cluster is a set of points such that any point in a cluster is closer (or more similar) to every other point in the cluster than to any point not in the cluster.







3 well-separated clusters

Center-based clusters

Center-based

- A cluster is a set of objects such that an object in a cluster is closer (more similar) to the "center" of a cluster, than to the center of any other cluster
- The center of a cluster is often a centroid, the average of all the points in the cluster, or a medoid, the most "representative" point of a cluster



4 center-based clusters

Contiguity-based clusters

- Contiguous Cluster (Nearest neighbor or Transitive)
 - A cluster is a set of points such that a point in a cluster is closer (or more similar) to one or more other points in the cluster than to any point not in the cluster.



Density-based

Density-based

- A cluster is a dense region of points, which is separated by lowdensity regions, from other regions of high density.
- Used when the clusters are irregular or intertwined, and when noise and outliers are present.



6 density-based clusters

Conceptual clusters

- Shared Property or Conceptual Clusters
 - Finds clusters that share some common property or represent a particular concept.



2 Overlapping Circles

Algorithm to cluster types mapping

- K-means and its variants
 - Center-based
 - Density-based
- DBSCAN clustering
 - Density-based
 - Contiguity-based

K means clustering

- Exclusive clustering approach
- Each cluster is associated with a centroid (center point)
- Each point is assigned to the cluster with the closest centroid
- Number of clusters, K, must be specified
- The basic algorithm is very 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

Details

- Initial centroids are often chosen randomly.
 - Clusters produced vary from one run to another.
- The centroid is (typically) the mean of the points in the cluster.
- 'Closeness' is measured by Euclidean distance, cosine similarity, correlation, etc.
- K-means will converge for common similarity measures mentioned above.
- Most of the convergence happens in the first few iterations.
 - Often the stopping condition is changed to 'Until relatively few points change clusters'
- Complexity is O(n * K * I * d)
 - n = number of points, K = number of clusters,
 I = number of iterations, d = number of attributes

Stochasticity







Good initialization



Good result



Poor initialization





Bad result



How to measure bad?

- Most common measure is Sum of Squared Error (SSE)
 - For each point, the error is the distance to the nearest cluster
 - To get SSE, we square these errors and sum them.

$$SSE = \sum_{i=1}^{K} \sum_{x \in C_i} dist^2(m_i, x)$$

- x is a data point in cluster C_i and m_i is the representative point for cluster C_i
 - can show that m_i corresponds to the center (mean) of the cluster
- Given two clusters, we can choose the one with the smallest error
- One easy way to reduce SSE is to increase K, the number of clusters
 - A good clustering with smaller K can have a lower SSE than a poor clustering with higher K

How to fix?

• Multiple runs

– Helps, but probability is not on your side

- Select more than k initial centroids and then select among these initial centroids
- Select most widely separated initial centroids
- Bisecting K-means

Not as susceptible to initialization issues

Bisecting K means

- Bisecting K-means algorithm
 - Variant of K-means that can produce a partitional or a hierarchical clustering
- 1: Initialize the list of clusters to contain the cluster containing all points.
- 2: repeat
- 3: Select a cluster from the list of clusters
- 4: for i = 1 to number_of_iterations do
- 5: Bisect the selected cluster using basic K-means
- 6: end for
- 7: Add the two clusters from the bisection with the lowest SSE to the list of clusters.
- 8: until Until the list of clusters contains K clusters

Example

Iteration 10 81 6 4 2 \succ_0 -2 -4 -6 20 10 5 15 0

Х

Other problems

- K-means has problems when clusters are of differing
 - Sizes
 - Densities
 - Non-globular shapes

• K-means has problems when the data contains outliers.

Different sizes



Original Points

K-means (3 Clusters)

Different densities



Original Points

K-means (3 Clusters)

Non-globular shapes



Original Points

K-means (2 Clusters)

Can increase K



Original Points

K-means Clusters

Can increase K



Original Points

K-means Clusters

Can increase K



Original Points

K-means Clusters

WHAT TO DO?

Bigger K = Bigger CPT table = sparser observations per cell = uncertainty

DBSCAN

- DBSCAN is a density-based algorithm.
 - Density = number of points within a specified radius (Eps)
 - A point is a core point if it has more than a specified number of points (MinPts) within Eps
 - These are points that are at the interior of a cluster
 - A border point has fewer than MinPts within Eps, but is in the neighborhood of a core point
 - A noise point is any point that is not a core point or a border point.

Definitions



DBSCAN algorithm

- Eliminate noise points
- Perform clustering on the remaining points

 $current_cluster_label \leftarrow 1$

for all core points \mathbf{do}

 ${\bf if}$ the core point has no cluster label ${\bf then}$

 $current_cluster_label \leftarrow current_cluster_label + 1$

Label the current core point with cluster label $current_cluster_label$ end if

for all points in the Eps-neighborhood, except i^{th} the point itself do

 ${\bf if}$ the point does not have a cluster label ${\bf then}$

Label the point with cluster label *current_cluster_label*

end if

end for

end for

Example





Original Points

Point types: core, border and noise

Eps = 10, MinPts = 4

Strengths





Original Points

Clusters

- Resistant to Noise
- Can handle clusters of different shapes and sizes

Weaknesses



Original Points

- Varying densities
- High-dimensional data O(n²)
 - •But see (Gan & Tao, 2015)



(MinPts=4, Eps=9.75).



(MinPts=4, Eps=9.92)

Parameter fitting

- Idea is that for points in a cluster, their kth nearest neighbors are close by
- Noise points have the kth nearest neighbor far away
- So, plot sorted distance of every point to its kth nearest neighbor



Validation

- For supervised classification we have a variety of measures to evaluate how good our model is
 - Accuracy, precision, recall
- For cluster analysis, the analogous question is how to evaluate the "goodness" of the resulting clusters?
- But "clusters are in the eye of the beholder"!
- Then why do we want to evaluate them?
 - To avoid finding patterns in noise
 - To compare clustering algorithms
 - To compare two sets of clusters
 - To compare two clusters

Clusters in random data



Cluster correlation

• Two matrices

- Proximity Matrix
- "Incidence" Matrix
 - One row and one column for each data point
 - An entry is 1 if the associated pair of points belong to the same cluster
 - An entry is 0 if the associated pair of points belongs to different clusters
- Compute the correlation between the two matrices
 - Since the matrices are symmetric, only the correlation between n(n-1) / 2 entries needs to be calculated.
- High correlation indicates that points that belong to the same cluster are close to each other.
- Not a good measure for some density or contiguity based clusters.

Correlation as measure of quality

 Correlation of incidence and proximity matrices for the K-means clusterings of these two data sets.



0.58

Similarity matrix visualization



For random data



For random data





Can give you fine detail



DBSCAN

Can also use residuals

- Clusters in more complicated figures aren't well separated
- Internal Index: Used to measure the goodness of a clustering structure without respect to external information
 - SSE
- SSE is good for comparing two clusterings or two clusters (average SSE).
- Can also be used to estimate the number of clusters



What's a good SSE?

- Example
 - Compare SSE of 0.005 against three clusters in random data
 - Histogram shows SSE of three clusters in 500 sets of random data points of size 100 distributed over the range 0.2 – 0.8 for x and y values



Cohesion and separation

- Cluster Cohesion: Measures how closely related are objects in a cluster
 - Example: SSE
- Cluster Separation: Measures how distinct or well-separated a cluster is from other clusters
- Example: Squared Error
 - Cohesion is measured by the within cluster sum of squares (SSE)

$$WSS = \sum_{i} \sum_{x \in C_i} (x - m_i)^2$$

Separation is measured by the between cluster sum of squares

$$BSS = \sum_{i} |C_{i}| (m - m_{i})^{2}$$

Where |C_{i}| is the size of cluster i

Entropy and purity

Cluster	Entertainment	Financial	Foreign	Metro	National	Sports	Entropy	Purity
1	3	5	40	506	96	27	1.2270	0.7474
2	4	7	280	29	39	2	1.1472	0.7756
3	1	1	1	7	4	671	0.1813	0.9796
4	10	162	3	119	73	2	1.7487	0.4390
5	331	22	5	70	13	23	1.3976	0.7134
6	5	358	12	212	48	13	1.5523	0.5525
Total	354	555	341	943	273	738	1.1450	0.7203

Table 5.9. K-means Clustering Results for LA Document Data Set

- entropy For each cluster, the class distribution of the data is calculated first, i.e., for cluster j we compute p_{ij} , the 'probability' that a member of cluster j belongs to class i as follows: $p_{ij} = m_{ij}/m_j$, where m_j is the number of values in cluster j and m_{ij} is the number of values of class i in cluster j. Then using this class distribution, the entropy of each cluster j is calculated using the standard formula $e_j = \sum_{i=1}^{L} p_{ij} \log_2 p_{ij}$, where the L is the number of classes. The total entropy for a set of clusters is calculated as the sum of the entropies of each cluster i.e., $e = \sum_{i=1}^{K} \frac{m_i}{m} e_j$, where m_j is the size of cluster j, K is the number of clusters, and m is the total number of data points.
- **purity** Using the terminology derived for entropy, the purity of cluster j, is given by $purity_j = \max p_{ij}$ and the overall purity of a clustering by $purity = \sum_{i=1}^{K} \frac{m_i}{m} purity_j$.