Clustering

Rich Caruana Cornell University

Supervised Learning

- Linear regression
- Logistic regression
- Artificial neural nets
- Decision trees
- K-nearest neighbor
- Support vectors
- ...

Supervised Learning

Train Set:



Test Set:

Supervised Learning

Well Defined Goal:

Learn function y = f(x) that is a good approximation from the training sample D

Know How to Measure Error:

Accuracy, Squared Error, ...

Clustering ≠ Supervised Learning

Clustering

Unsupervised Learning

Supervised Learning

Train Set:

57, M, 195, 0, 125, 95, 39, 25, 0, 1, 0, 0, 0, 1, 0, 0, 0, 0, 0, 0, 1, 1, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0,	
78,M,160,1,130,100,37,40,1,0,0,0,1,0,1,1,1,0,0,0,0,0,0,0,0,0,	
69,F,180,0,115,85,40,22,0,0,0,0,0,1,0,0,0,0,1,0,0,0,0,0,0,0	
18,M,165,0,110,80,41,30,0,0,0,1,0,0,0,0,0,0,0,0,0,0,0,0,0,0	
54,F,135,0,115,95,39,35,1,1,0,0,0,1,0,0,0,1,0,0,0,1,0,0,0,1,0,0,0,0	
84,F,210,1,135,105,39,24,0,0,0,0,0,0,0,0,0,0,0,0,0,0,0,0,0,0,	
89,F,135,0,120,95,36,28,0,0,0,0,0,0,0,0,0,0,0,0,1,1,0,0,0,0,0,	
49,M,195,0,115,85,39,32,0,0,0,1,1,0,0,0,0,0,0,1,0,0,0,0,1,0,0,0,0	
40,M,205,0,115,90,37,18,0,0,0,0,0,0,0,0,0,0,0,0,0,0,0,0,0,0,	
74,M,250,1,130,100,38,26,1,1,0,0,0,1,1,0,0,0,0,0,0,0,0,0,0,0,0,	
77, F, 140, 0, 125, 100, 40, 30, 1, 1, 0, 0, 0, 0, 0, 0, 0, 1, 0, 0, 0, 0, 0, 0, 0, 0, 1, 1	

Test Set:

Un-Supervised Learning

Train Set:

Test Set:

Un-Supervised Learning

Train Set:

Test Set:

Un-Supervised Learning

Data Set:

Clustering

• Goal:

- Find groups in data
- Group similar items together

71,M,160,1,130,105,20,20,10,0,0,0,0

- Separate different kinds of objects
- Find patterns in data
- Organize/understand complex data

Clustering

• Given:

- Data Set D (training set)
- Similarity/distance metric/information
- Find:
 - Partitioning of data
 - Groups of similar/close items

Similarity?

- Groups of similar customers
 - Similar demographics
 - Similar buying behavior
 - Similar health
- Similar products
 - Similar cost
 - Similar function
 - Similar store
 - -...
- Similarity usually is domain/problem specific

Types of Clustering

- Hierarchical
 - Divisive clustering (top down)
 - Agglomerative clustering (bottom up)
- Partitioning
 - K-means clustering
 - K-medoids clustering
 - EM (expectation maximization) clustering
- Density-Based Methods
 - Regions of dense points separated by sparser regions of relatively low density







Two Types of Data/Distance Info• N-dim vector space representation and distance metricDimension of the state of

Agglomerative Clustering

- Put each item in its own cluster (N singletons)
- Calculate pairwise distances between clusters
- Merge the two *closest* clusters
- Repeat until all points in one cluster
- Hierarchical clustering
- Yields a clustering with each possible # of clusters
- Greedy clustering: not optimal for any cluster size

Do Agglomerative Clustering Demo

Agglomerative Clustering of Proteins



Merging: Closest Clusters

- Nearest neighbors (shortest link)
- Nearest average distance (average link)
- Smallest greatest distance (maximum link)
- Domain specific similarity measure – word frequency, TFIDF, KL-divergence, ...
- Merge clusters that optimize criterion after merge – minimum mean_internal_distance



Minimum Distance Between Clusters

$Min_Dist\ (c_1,c_2) = MIN_{i \in c_1, j \in c_2}(Dist\ (i, j))$



			"mc.5.1000.plot"	+
1	\$\$ \$ \$	#辈 +*幕	\$* \$\$, \$\$	
	## 4 <u>*</u>	業集 井+	4# ## ## ## ##	
0.8	老道 耕什	2년 분호		
	30 M 144	***	44 ¥1 44 44	
0.6				
0.4				
	14 14	御物 活性	花井 花井 古井 井井	
0.2	Þ¥ 4 *	74 44	魏 48 《韓 48	
	## ##	₩£ \$¢	#4 #2 ** ** **	
	19.# ±#	t.1 .++	15 38 *** ***	

				"mc.5.1000.plot" +
0.35				. + .
0.3	+ * * + + + +	4. 4. 1. 1 <u>1</u>	+ + +++ +++ + ±	+ + 4 + ≠
0.25	7 + + 4 + + + +	** *** **	······································	#++ + * *++ +#
0.2	+ +			
0.15				
0.1	+ ++ + +	*	++ +	+ +
	\$4 \$ <i>#</i>	ŧ ‡	*+ +	***
0.05	+12 4	++ ‡‡	· ·	-
	\$\$ <u>\$</u> ‡	+ +	+ + +	+ + +

Recursive Clusters



Recursive Clusters



Weighted Mean Internal Distance



Weighted Mean Internal Distance





Distance Between Helices

- Vector representation of protein data in 3-D space that gives x,y,z coordinates of each atom in helix
- Use a program developed by chemists (fortran) to convert 3-D atom coordinates into average atomic distances in angstroms between aligned helices
- 641 helices = 641 * 640 / 2
 - = 205,120 pairwise distances



Agglomerative Clustering of Proteins



Agglomerative Clustering of Proteins





Agglomerative Clustering of Proteins



Agglomerative Clustering of Proteins













Agglomerative Clustering

- Greedy clustering
 - once points are merged, never separated
 - suboptimal w.r.t. clustering criterion
- Combine greedy with iterative refinement
 - post processing
 - interleaved refinement







How do you tell if clustering is finding *REAL* structure in data?



			"ma	.5.1000.plot" +
1	## ## ## 21	## ## ## #+	북학 출시 2월 3월	## ## ## ##
0.8	蒜 耕 森 姓	114 #4 12 51	标 44 19 14	4.結 \$## #15_###
0.6				
0.4	14 1 4	難得	拉許 幣款	14 ##
0.2	₽¥ 4*	74 44	学花 非 英	₩1 # 2
	***	\$\$\$ \$\$ [†]	# #	11 H

Recursive Clusters



Recursive Clusters

Weighted Mean Internal Distance





Recursive Clusters + Random Noise



Bergmark Data

• 35.000 web documents

- results from 26 focused web crawls
- have crawl label of each document (know "class")
- Bags-of-words
 - how many times each word occurs in each document
- Cosine distance
 - good metric for sparse vectors
 - computationally efficient



14



Agglomerative on Bergmark





Agglomerative on Bergmark

K-Means Clustering

- Inputs: data set and k (number of clusters)
- Output: each point assigned to one of k clusters

• K-Means Algorithm:

- -Initialize the k-means
 - +assign from randomly selected points
- +randomly or equally distributed in space
- -Assign each point to nearest mean
- -Update means from assigned points
- -Repeat until convergence

Agglomerative Clustering

Computational Cost

- O(N²) just to read/calculate pairwise distances
- N-1 merges to build complete hierarchy
 - + scan pairwise distances to find closest
- + calculate pairwise distances between clusters
 + fewer clusters to scan as clusters get larger
- Overall O(N³) for simple implementations
- Improvements
 - sampling
 - dynamic sampling: add new points while merging
 - Kruskal's algorithm for minimum spanning tree (shortest link)
 - tricks for updating pairwise distances

K-Means Clustering: Convergence

• Squared-Error Criterion

Squared _Error = $\sum_{c} \sum_{i \in c} (Dist(i, mean(c)))^2$

- Converged when SE criterion stops changing
- Increasing K reduces SE can't determine K by finding minimum SE
- Instead, plot SE as function of K



K-Means Clustering

- Efficient
 - K << N, so assigning points is O(K*N) < O(N²)
 - updating means can be done during assignment
 - usually # of iterations << N</p>
 - Overall O(N*K*iterations) closer to O(N) than O(N²)
- Gets stuck in local minima
 - Sensitive to initialization
- Number of clusters must be pre-specified
- Requires vector space date to calculate means



Soft K-Means Clustering

- Instance of EM (Expectation Maximization)
- Like K-Means, except each point is assigned to each cluster with a probability
- Cluster means updated using weighted average
- Generalizes to Standard_Deviation/Covariance
- Works well if cluster models are known

Soft K-Means Clustering (EM)

- -Initialize model parameters:
 - +means
 - $+ std_devs$
 - +••
- -Assign points probabilistically to each cluster
- -Update cluster parameters from weighted points
- -Repeat until convergence to local minimum

What do we do if we can't calculate cluster means?



K-Medoids Clustering

- Inputs: data set and k (number of clusters)
- Output: each point assigned to one of k clusters
- •
- Initialize k medoids
 - pick points randomly
- Pick medoid and non-medoid point at random
- Evaluate quality of swap
 - Mean point happiness
- Accept random swap if it improves cluster quality

Comparison: Cost of K-Means

- n cases; d dimensions; k centers; i iterations
- compute distance each point to each center: $O(n^*d^*k)$
- assign each of n cases to closest center: O(n*k)
- update centers (means) from assigned points: O(n*d*k)
- repeat i times until convergence
- overall: O(n*d*k*i)
- much better than $O(n^2)$ - $O(n^3)$ for HAC
- sensitive to initialization run many times
- usually don't know k run many times with different k
- requires many passes through data set

Clustering Summary

- Goal
 - Find groups of similar objects in complex data
- Agglomerative clustering
 - Finds a hierarchy
 - Can cut clustering tree at any number of clusters
 - Works with both vector and similarity data
 - Expensive: O(N2) cost
- K-Means clustering
 - Must be told how many clusters to find
 - Usually requires vector data
 - Faster with large data sets
 - Gets stuck in local minima --- run multiple times and pick best

Thank You!