



# CS145 Discussion: Week 5 Classification evaluation, Clustering

Junheng Hao Friday, 11/06/2020







- Announcement
- Similarity Measurement
- Classification evaluation
- Clustering: K-means





- Homework 3 due Nov. 9 (Monday, Week 6) 11:59 PT
  - Submit through GradeScope of 1 PDF (2 python file and 2 jupyter notebooks into 1 PDF file)
  - Assign pages to the questions on GradeScope
- Midterm project due on Nov. 11 (Wednesday, Week 6)
  - 3-page midterm project report
  - At least one submission to Kaggle

- Approximately 3 pages
- Current progress about project, including
  - $\circ$   $\quad$  Data processing and transformation
  - $\circ \quad \text{Designed & tested models / methods}$
- Discussion and future project plan
  - Some conclusions and findings
  - Analysis of current models and techniques
  - $\circ$  ~ Timeline of future project plan (around the next 4 weeks)
- Midterm exam on Nov. 16 (Monday, Week 7) on CCLE (with proper browser setting)





- Similarity and dissimilarity matrix
  - Pairwise measures how alike/different two data points are
- Example of numerical attribute
  - Three 2-dim input, x1,x2,x3
  - We write the dissimilarity matrix as a 3x3 lower-triangular matrix

$$\begin{bmatrix} 3 & 5 \\ 6 & 9 \\ 11 & 21 \end{bmatrix} \longrightarrow \begin{bmatrix} 0 & 0 & 0 \\ d(2,1) & 0 & 0 \\ d(3,1) & d(3,2) & 0 \end{bmatrix}$$

• Dissimilarity under Euclidean distance

$$\begin{bmatrix} 0 & 0 & 0 \\ \sqrt{(3-6)^2 + (5-9)^2} & 0 & 0 \\ \sqrt{(3-11)^2 + (5-21)^2} & \sqrt{(11-6)^2 + (21-9)^2} & 0 \end{bmatrix}$$





- Student 1: likes Jazz, eats pizza, roots for the cubs, wears socks
- Student 2: likes Rock, eats pizza, roots for the cubs, goes barefoot
- d(Student 1, Student 2):
  - $\circ \quad m: \text{ \# of matches} \to 2$
  - $\circ \quad p: \text{ total \# of variables} \to 4$
  - d(Student 1, Student 2) = (4-2)/4 = 0.5







- Symmetric binary attributes:
  - $\circ$  Gender
- Asymmetric attributes:
  - Preference, diagnosis, etc.
- Can be manually defined





- Dissimilarity of Binary Attributes:
  - Define 0 and 1 and calculate q,s,r,t,p
- Symmetric binary variables:

 $d(i, j) = \frac{r+s}{q+r+s+t}$ 

• Asymmetric binary variables (0 is less important):

$$d(i,j) = \frac{r+s}{q+r+s}$$

• Jaccard coefficient (similarity measure for asymmetric binary variables):

$$sim_{Jaccard}(i, j) = \frac{q}{q+r+s}$$







Name	Gender	Fever	Cough	Test-1	Test-2	Test-3	Test-4
Jack	Μ	Y	N	Р	N	N	N
Mary	F	Y	N	Р	N	Р	N
Jim	Μ	Y	Р	N	N	N	N

		0	bject <i>j</i>	
		1	0	sum
Object	t / <sup>1</sup>	q	r	q + r
	0	8	t	s+t
	sum	q + s	r+t	p

- Define M,Y,P as 1; Define F,N as 0
- Assume symmetric for Gender, asymmetric for other attributes
- i = Jack, j = Mary, what is d(i, j) in terms asymmetric attributes?
  - $\circ$  r = 0, s = 1, q = 2
  - o d(i, j) = 0.32

**Engineer Change.** 





- Order is important
  - Freshman, Sophomore, Junior, Senior
- Replace attribute by rank

 $r_{if} \in \{1, ..., M_f\}$ 

• Convert rank to numeric values

$$z_{if} = \frac{r_{if} - 1}{M_f - 1}$$





- For vector data
- d1: I like to go to the store
- d2: I like the cubs, go cubs go

 $\cos(d_1, d_2) = (d_1 \bullet d_2) / ||d_1|| ||d_2||$ ,

Document	I.	like	to	go	the	store	cubs
d1	1	1	2	1	1	1	0
d2	1	1	0	2	1	0	2

• cos(d1, d2)?

 $1 \cdot 1 + 1 \cdot 1 + 2 \cdot 0 + 1 \cdot 2 + 1 \cdot 1 + 1 \cdot 0 + 0 \cdot 2$ 

 $\overline{\sqrt{1^2+1^2+2^2+1^2+1^2+1^2+0^2}} \cdot \sqrt{1^2+1^2+0^2+2^2+1^2+0^2+2^2}$ 







• d(3,1)?  
• 
$$\frac{1(1)+1(0.5)+1(0.45)}{3}$$





#### Confusion Matrix: How your model got confused

Actual class\Predicted class	C <sub>1</sub>	- C <sub>1</sub>		
C <sub>1</sub>	True Positives (TP)	False Negatives (FN)		
- C <sub>1</sub>	False Positives (FP)	True Negatives (TN)		







Accuracy = (TP + TN)/All Error rate = (FP + FN)/All Sensitivity = TP/P Specificity = TN/N  $precision = \frac{TP}{TP + FP}$  $recall = \frac{TP}{TP + FN}$ 



## **Evaluation Metric**





Why do we need all these measures? Isn't accuracy telling us how good the model is?

- Imbalanced data
- Imbalanced importance of positive and negative

Accuracy = (TP + TN)/All Error rate = (FP + FN)/All Sensitivity = TP/P Specificity = TN/N  $precision = \frac{TP}{TP + FP}$  $recall = \frac{TP}{TP + FN}$ 





Single number metric

• F measure (F1 score)

 $F = \frac{2 \times precision \times recall}{precision + recall}$ 

- Area under the curve (AUC)
  - ROC curve (true positive against false positive)
  - PRC (precision against recall)





TPR



True positive rate: TPR = TP/P (sensitivity) False positive rate: FPR = FP/N (1-specificity)

#### http://mlwiki.org/index.php/ROC\_Analysis







- Input / Output / Goal of clustering analysis
  - Large amount of unlabeled data in real life
- Supervised learning v.s. unsupervised learning
- Typical clustering algorithm examples:
  - K-means
  - Hierarchical clustering
  - DB-SCAN
  - Mixture models







- Demo 1: <u>http://stanford.edu/class/ee103/visualizations/kmeans/kmeans.html</u>
- Demo 2: <u>https://www.naftaliharris.com/blog/visualizing-k-means-clustering/</u>









Until no change







- Key idea of K-means algorithms:
  - Step 1: Partition into k non-empty subsets (select K points as initial centroids)
  - Step 2: Iteration: Update mean point and assign object to cluster again
  - Step 3: Stop when converge
- Partition-based clustering methods
- Can be considered as a special case of GMM







- Q1: Will K-means converge?
- Q2: Will different initialization of K-means generate different clustering results?



### K-means



- Q1: Will K-means converge?
- A1: Yes.  $J = \sum_{j=1}^{k} \sum_{C(i)=j} d(x_i, c_j)^2$
- Q2: Will different initialization of K-means generate different clustering results?
- A2: Yes. Initialization matters!







- Efficiency: O(tkn) normally k,t are much smaller than n → efficient
- Can terminate at a local optimum
- Need to specify k (or take time to find best k)
- Sensitive to noisy data and outliers
- Different sizes and variances
- Not suitable to discover clusters with non-convex shapes
  - Can K-medoids help?
- Many variants of K-means:
  - K-means++, Genetics K-means, etc.







#### **Hierarchical Clustering**



- Method
  - Divisive (Top-down)
  - Agglomerative (Bottom-up)

- Distance metrics
  - Single linkage
  - Complete linkage
  - Average linkage
  - Centroid
  - Medoid





### **Hierarchical Clustering**



• Single Linkage

• Complete Linkage

• Average Linkage



 $L(r,s) = \min(D(x_{ri}, x_{sj}))$ 



 $L(r,s) = \max(D(x_{ri}, x_{sj}))$ 



 $L(r,s) = \frac{1}{n_r n_s} \sum_{i=1}^{n_r} \sum_{j=1}^{n_s} D(x_{ri}, x_{sj})$ 







- Density-based clustering method
- Discover clusters of arbitrary shape
- Handle noise
- <u>Demo</u>





### DB-SCAN



- Two parameters:
  - Eps: Maximum radius of the neighborhood
  - *MinPts*: Minimum number of points in an Epsneighborhood of that point
- N<sub>Eps</sub>(q): {p belongs to D | dist(p,q) ≤ Eps}
- Directly density-reachable: A point *p* is directly densityreachable from a point *q* w.r.t. *Eps*, *MinPts* if
  - p belongs to  $N_{Eps}(q)$ • q is a core point, core point condition:  $|N_{Eps}(q)| \ge MinPts$  Min Eps

MinPts = 5 Eps = 1 cm







#### Density-reachable:

A point *p* is density-reachable from a point *q* w.r.t. *Eps*, *MinPts* if there is a chain of points *p*<sub>1</sub>, ..., *p*<sub>n</sub>, *p*<sub>1</sub> = *q*, *p*<sub>n</sub> = *p* such that *p*<sub>i+1</sub> is directly density-reachable from *p*<sub>i</sub>



- Density-connected
  - A point *p* is density-connected to a point *q* w.r.t. *Eps*, *MinPts* if there is a point *o* such that both, *p* and *q* are density-reachable from *o* w.r.t. *Eps* and *MinPts*







- Relies on a *density-based* notion of cluster: A *cluster* is defined as a maximal set of density-connected points
- Noise: object not contained in any cluster is noise
- Discovers clusters of arbitrary shape in spatial databases with noise









- Density-based clustering method
- <u>Demo</u>
- Pros and Cons
  - It allows noise, so it is robust to outliers
  - It can figure out number of clusters automatically, as opposed to k-means
  - It can find arbitrarily shaped clusters
  - It is not deterministic, depending on the data processing order
  - It is not a good choice for clustering data sets with large differences in densities









# Thank you!

Q & A