Welcome to

DS3010: DS-III: Computational Data Intelligence Unsupervised Learning Prof. Yanhua Li

Time: 11:00am – 12:50pm M & R Location: HL 114 D-term 2022

Data pipeline



Urban Computing: concepts, methodologies, and applications. Zheng, Y., et al. *ACM transactions on Intelligent Systems and Technology*.





Unsupervised Learning: Clustering



Bishop: Chapter 9.1

Unsupervised Learning

Unsupervised Learning

 Clustering & dimension reduction (our focus)



Generation



- K-means
 - Clustering $X = \{x^1, \dots, x^n, \dots, x^N\}$ into K clusters
 - Initialize cluster center c^i , i=1,2, ... K (K random x^n from X)
 - Repeat
 - For all x^n in X: $b_i^n \begin{cases} 1 & x^n \text{ is most "close" to } c^i \\ 0 & 0 \end{cases}$ Otherwise

• Updating all
$$c^i$$
: $c^i = \sum_{x^n} b^n_i x^n / \sum_{x^n} b^n_i$

Bishop Chapter 9.1

Example: Assigning Clusters





Clusters after round 1

J. Leskovec, A. Rajaraman, J. Ullman: Mining of Massive Datasets, http://www.mmds.org

Example: Assigning Clusters



x ... data point ... centroid

Clusters after round 2

J. Leskovec, A. Rajaraman, J. Ullman: Mining of Massive Datasets, http://www.mmds.org

Example: Assigning Clusters



x ... data point ... centroid

Clusters at the end

J. Leskovec, A. Rajaraman, J. Ullman: Mining of Massive Datasets, http://www.mmds.org

Getting the k right

How to select k?

- Try different k, looking at the change in the average distance to centroid as k increases
- Average falls rapidly until right k, then changes little



J. Leskovec, A. Rajaraman, J. Ullman: Mining of Massiv<u>e1</u> Datasets, http://www.mmds.org

Example: Picking *k=2*

Too few; many long distances to centroid.



Example: Picking *k=3*

Just right; distances rather short.



Example: Picking k



Clustering

• Hierarchical Agglomerative Clustering (HAC)



Using single linkage similarity

• Hierarchical Agglomerative Clustering (HAC)



Using complete linkage similarity

• Hierarchical Agglomerative Clustering (HAC)



Questions