

## COMP6237 Data Mining Lecture 2: Discovering Groups

#### Zhiwu Huang

Zhiwu.Huang@soton.ac.uk Lecturer (Assistant Professor) @ VLC of ECS University of Southampton

Lecture slides available here: http://comp6237.ecs.soton.ac.uk/zh.html

(Thanks to Prof. Jonathon Hare and Dr. Jo Grundy for providing the lecture materials used to develop the slides.)



#### **Discovering Groups - Roadmap**





#### Discovering Groups – Textbook

# CHAPTER 3 Discovering Groups

Programming Collective Intelligence: Building Smart Web 2.0 Applications *T. Segaran*.

Chapter 7

## Clustering

Clustering is the process of examining a collection of "points," and grouping the points into "clusters" according to some distance measure. The goal is that points in the same cluster have a small distance from one another, while points in different clusters are at a large distance from one another. A suggestion of

Mining of Massive Datasets J. Leskovec et al

## Discovering Groups – Overview (1/5)



- Clustering algorithms group data, just using the feature vectors
  - Unsupervised: no group labels for training
  - Key idea: data with similar features grouped together
  - Can be
    - Hard (each item assigned to one group)
    - Soft (allow overlapping groups)

Feature space: color value (3D)





#### Discovering Groups – Overview (2/5)



0





Moons



https://marlabskochi.github.io



#### Discovering Groups – Overview (3/5)

#### Clustering Algorithm (1/3) - K-means



Credit: Pratik Thorat



#### Discovering Groups – Overview (4/5)

**Clustering Algorithm (2/3) - DBSCAN** 





#### Discovering Groups – Overview (5/5)

**Clustering Algorithm (3/3)** – Hierarchical/**Agglomerative** (Divisive not

learned in this lecture)





- LO1: Comprehend the key ideas and the essential mathematical formulations employed in clustering methods (exam).
  - E.g., how is sum of squared error (SSE) defined?
  - E.g., understand the pros and cons of the learned algorithms
- LO2: Compute the fundamental stages of learned clustering approaches (exam).
  - E.g., given a dataset and a distance metric, be prepared to follow the selected clustering algorithm to cluster the instances in the dataset
- LO3: Implement and evaluate the learned clustering algorithms using Python (course work)

#### Assessment hints: Multi-choice Questions (single answer: concepts, calculation etc)

- *Textbook Exercises: textbooks (Programming + Mining)*
- Other Exercises: <u>https://www-users.cse.umn.edu/~kumar001/dmbook/sol.pdf</u>
- ChatGPT or other Al-based techs

**Discovering Groups – K-means** 



- Given: data set  $\{\mathbf{x}_1, \ldots, \mathbf{x}_N\}$ , number of clusters K
- Goal: find cluster centers  $\{\mu_1, \ldots, \mu_K\}$  so that

$$J = \sum_{n=1}^{N} \sum_{k=1}^{K} r_{nk} \|\mathbf{x}_n - \boldsymbol{\mu}_k\|^2$$

is minimal, where  $r_{nk} = 1$  if  $\mathbf{x}_n$  is assigned to  $\boldsymbol{\mu}_k$ 

• Idea: compute  $r_{nk}$  and  $\mu_k$  iteratively Otherwise,  $r_{n,k} = 0$ 

# The used objective function is **Sum of Squared Error** (SSE)



```
Algorithm 1: K Means clustering
Data: X, K
initialise K centroids;
while positions of centroids change do
   for each data point do
       assign to nearest centroid
   end
   for each centroid do
       move to average of assigned data points
   end
end
return centroids, assignments;
```

A special case of Expectation Maximisation - why?

## **Discovering Groups – K-means**



```
Algorithm 2: K Means clustering
Data: X, K
initialise K centroids;
while positions of centroids change do
   for each data point do
       assign to nearest centroid ;
                                                // Expectation of
        associations E-step: Estimate the posterior probabilities...
   end
   for each centroid do
       move to average of assigned data points ;
        // Maximisation of likelihood M-step: Estimate new parameters
   end
end
return centroids, assignments;
```

Assumes spherical clusters

Initialize cluster means:  $\{ \boldsymbol{\mu}_1, \dots, \boldsymbol{\mu}_K \}$ 







Find optimal assignments:

$$r_{nk} = \begin{cases} 1 & \text{if } k = \arg\min_{j} \|\mathbf{x}_{n} - \boldsymbol{\mu}_{j}\| \\ 0 & \text{otherwise} \end{cases}$$





Find new optimal means:



Source: D. Cremers



Find new optimal assignments:

$$r_{nk} = \begin{cases} 1 & \text{if } k = \arg\min_{j} \|\mathbf{x}_{n} - \boldsymbol{\mu}_{j}\| \\ 0 & \text{otherwise} \end{cases}$$



Source: D. Cremers

Iterate these steps until means and assignments do not change any more









- Real data set
- Radom initialization

 Magenta line is 'decision boundary'



## Discovering Groups – K-means Sum of Squared Error (SSE) Curve



- After every step the cost function J is minimized
- Blue steps: update assignments
- Red steps: update means
- Convergence after 4 rounds



#### Discovering Groups – K-means

K-means can quickly and cheaply cluster data.

Problems?

- need to specify the cluster number k
- depends on good initial centroid guesses
- may converge on local minimum
- assumes spherical data (or ellipsoid-shaped clusters, or at best convex clusters)

**Gaussian Mixture models (GMM)** can work better, using a generalization of K-means (assuming each cluster is Gaussian), not discussed in this lecture.



#### Discovering Groups – DBSCAN

• Idea: uses the local density of points to determine the clusters, rather than using only the distance between points

$$N_{\epsilon}(\mathbf{x}) = B_d(\mathbf{x}, \epsilon) = \{\mathbf{y} \mid \delta(\mathbf{x}, \mathbf{y}) \le \epsilon\}$$

where  $\delta(\mathbf{x}, \mathbf{y})$  represents the distance between  $\mathbf{x}$  and  $\mathbf{y}$ ,  $\varepsilon$  indicates Max radius, and  $\mathbf{x}$  is a core point if  $|N_{\epsilon}(\mathbf{x})| \ge minpts$ , where minpts is a Min number that is user-defined local density or frequency threshold

• **x** belongs to a density-based cluster when  $|N_{\epsilon}(\mathbf{x})| \ge minpts$  or  $\mathbf{x} \in N_{\epsilon}(\mathbf{z})$ 

where **z** is another data point, *minpts* is a **Min number** that is user-defined local density or frequency threshold

Max radius is the limit on which to look for neighbours Min number is the lower limit on what can be in a cluster



### **Discovering Groups – DBSCAN**

#### Algorithm 3: DBSCAN



22/36

Now we randomly pick a **Core Point**... Next, the **Core Points** that are close to the **first cluster**, meaning they overlap the **orange circle**...



However, because this is not a **Core Point**, we do not use it to extend the **first cluster** any further.

And now we are done creating the **first cluster**.



https://www.youtube.com/watch?v=RDZUdRSDOok



#### Discovering Groups – DBSCAN

DBSCAN works well on any shape of data and is robust to outliers.

Problems?

- needs radius & minimum number specified
- needs a distance parameter
- same parameter may not work for different cluster density
- can struggle in high dimensions



#### Discovering Groups – Hierarchical Clustering

Creates a **binary tree** that **recursively groups** pairs of similar items or clusters

Can be:

- Agglomerative (bottom up)
- Divisive (top down)

We will look at Agglomerative clustering. Needs a distance measure.

#### 26/36

#### Discovering Groups – Hierarchical Clustering

Distance based linkage:

Minimum, or single-linkage clustering Distance between two closest members

 $\min d(a, b) : a \in A, b \in B$ 

Produces long, thin clusters

Maximum, or complete-linkage clustering Distance between two most distant members

 $\max d(a, b) : a \in A, b \in B$ 

Finds compact clusters, approximately equal diameter

Mean or Average Linkage Clustering (UPGMA: Unweighted Pairwise Group Method with Arithmetic Mean):

$$\frac{1}{|A||B|} \sum_{a \in A} \sum_{b \in B} d(a, b)$$







With sample data:

$$X = \begin{bmatrix} 1.5 & 1.5 \\ 2.0 & 1.0 \\ 2.0 & 0.5 \\ -1.0 & 0.5 \\ -1.5 & -0.5 \\ -1 & 0.5 \end{bmatrix}$$



Distance matrix: demo distances, not ground truth

	Α	В	С	D	Е	F
Α	0	0.7	2.7	1.8		
В	0.7	0				
С	2.7		0			
D	1.8			0		
F	:				0	































## Discovering Groups – Hierarchical Clustering

#### Pros:

- No need to pre-specify cluster numbers; cut the dendrogram at the desired level for the clusters.
- Dendrograms easily summarize data into a hierarchy, facilitating cluster examination and interpretation.

#### Cons:

- Needs a threshold to determine the number of clusters
- Non-trivial to select the best linkage method



#### **Discovering Groups – Summary**

Clustering is a key way to understand your data.

There are many different approaches

- K Means Need to chose K
- DBSCAN need to choose min points and radius
- Hierarchical Agglomerative Clustering needs a threshold or number of clusters

They are a very good way to start exploring a dataset



Source: <u>https://scikit-learn.org/stable/auto\_examples/cluster/plot\_cluster\_comparison.html#sphx-glr-auto-examples-cluster-plot-cluster-comparison-py</u>



## Discovering Groups – Appendix (1/2)





## Discovering Groups – Appendix (2/2)

#### **Spectral Clustering**

Eigenvectors of the Laplacian matrix provide an embedding of the data based on similarity.

