

# Multivariate Clustering

Ranga **Raju** Vatsavai, Ph.D.  
 Chancellors Faculty Excellence Associate Professor in Geospatial Analytics  
 Department of Computer Science, North Carolina State University (NCSU)

March. 11-13, 2019

## K-means Clustering

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

3/11/19

Raju Vatsavai

## K-means Clustering

- 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
  
```

---

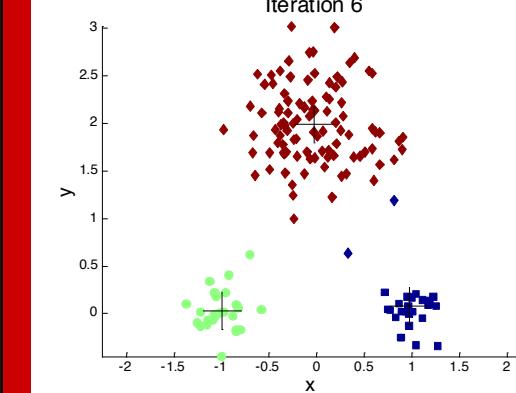
3/11/19

Raju Vatsavai

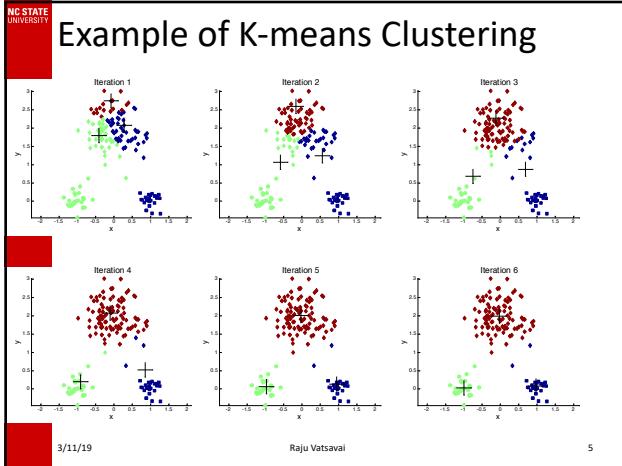
3

## Example of K-means Clustering

Iteration 6



4



**K-means Clustering – 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

3/11/19      Raju Vatsavai      6

**K-Means: Detailed Algorithm**

```

K-MEANS ( $\mathbf{D}, k, \epsilon$ ):
1  $t = 0$ 
2 Randomly initialize  $k$  centroids:  $\mu_1^t, \mu_2^t, \dots, \mu_k^t \in \mathbb{R}^d$ 
3 repeat
4    $t \leftarrow t + 1$ 
5    $C_j \leftarrow \emptyset$  for all  $j = 1, \dots, k$ 
6   // Cluster Assignment Step
7   foreach  $\mathbf{x}_j \in \mathbf{D}$  do
8      $j^* \leftarrow \operatorname{argmin}_i \left\{ \|\mathbf{x}_j - \mu_i^t\|^2 \right\}$  // Assign  $\mathbf{x}_j$  to closest centroid
9      $C_{j^*} \leftarrow C_{j^*} \cup \{\mathbf{x}_j\}$ 
10  // Centroid Update Step
11  foreach  $i = 1$  to  $k$  do
12     $\mu_i^t \leftarrow \frac{1}{|C_i|} \sum_{\mathbf{x}_j \in C_i} \mathbf{x}_j$ 
13 until  $\sum_{i=1}^k \|\mu_i^t - \mu_i^{t-1}\|^2 \leq \epsilon$ 
  
```

M. Zaki      7

3/11/19      Raju Vatsavai

**K-Means: Basic Computations**

- Compute new cluster centroid (mean)
 
$$m_k = \frac{\sum_{i \in C(i)} x_i}{N_k}, k = 1, \dots, K.$$
- Cluster assignment
 
$$C(i) = \arg \min_{1 \leq k \leq K} \|x_i - m_k\|^2, i = 1, \dots, N$$

3/11/19      Raju Vatsavai      8

## Evaluating K-means Clusters

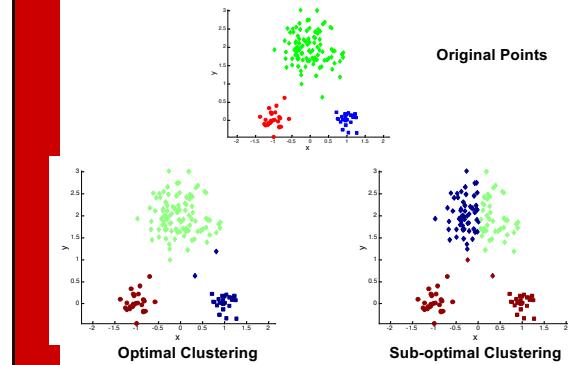
- Most common measure is Sum of Squared Error (SSE)
    - For each point, the error is the distance to the nearest cluster
    - To get SSE, we square these errors and sum them.
- $$SSE = \sum_{i=1}^K \sum_{x \in C_i} dist(m_i, x)^2$$
- $x$  is a data point in cluster  $C_i$  and  $m_i$  is the representative point for cluster  $C_i$ 
    - show that  $m_i$  corresponds to the center (mean) of the cluster
  - Given two sets of clusters, we prefer the one with the smallest error
  - One easy way to reduce SSE is to increase K, the number of clusters
    - But may not lead to better clustering!

3/11/19

Raju Vatsavai

9

## Two different K-means Clusterings



3/11/19

Raju Vatsavai

10

## Limitations of K-means

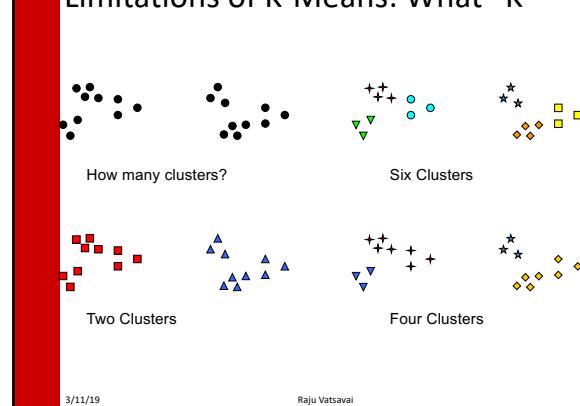
- What is appropriate "K"
- K-means has problems when clusters are of differing
  - Sizes
  - Densities
  - Non-globular shapes
- K-means has problems when the data contains outliers.

3/11/19

Raju Vatsavai

11

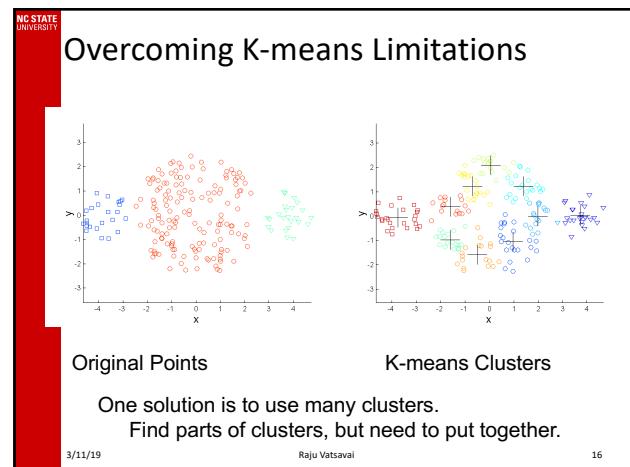
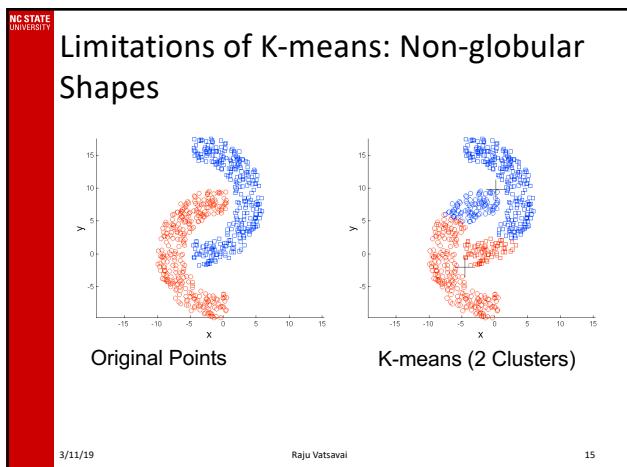
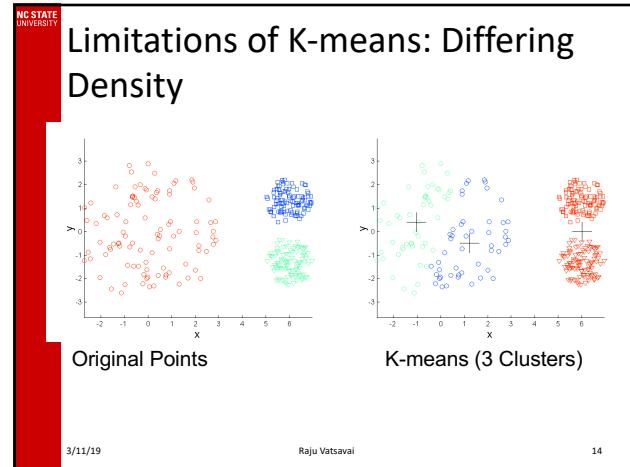
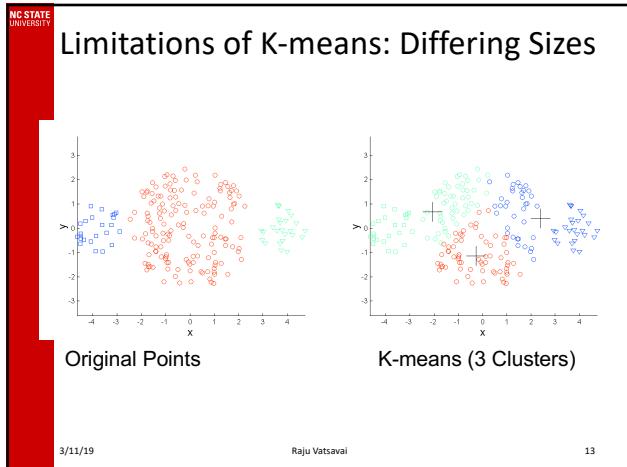
## Limitations of K-Means: What "K"

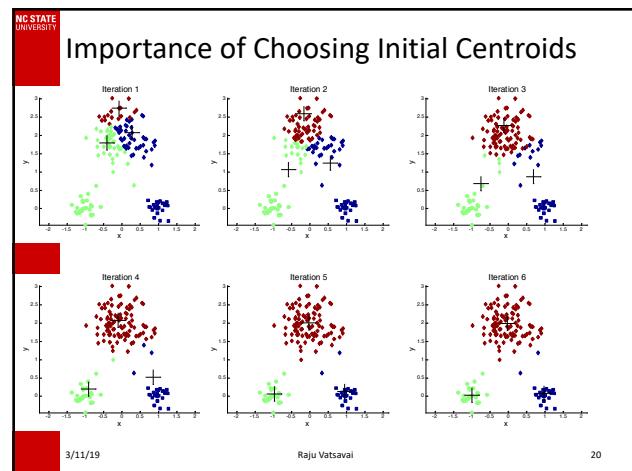
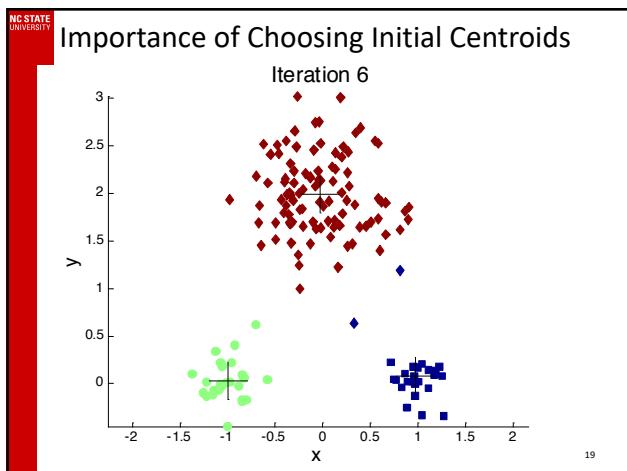
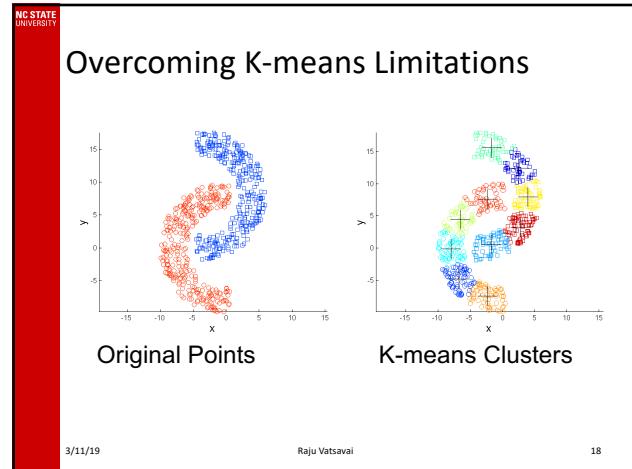
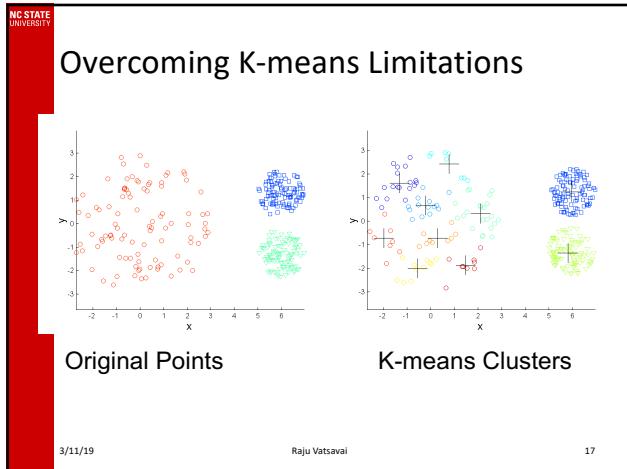


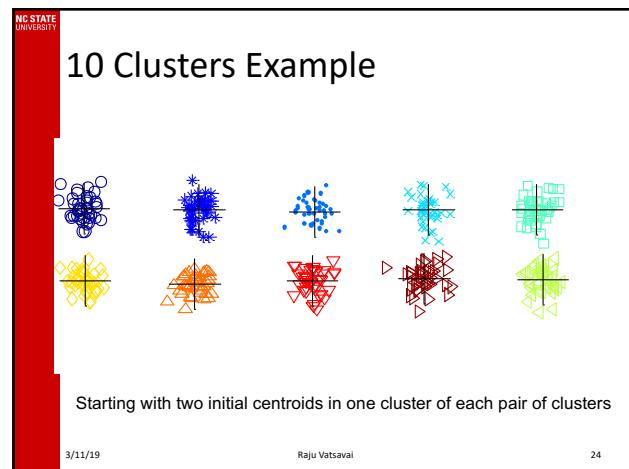
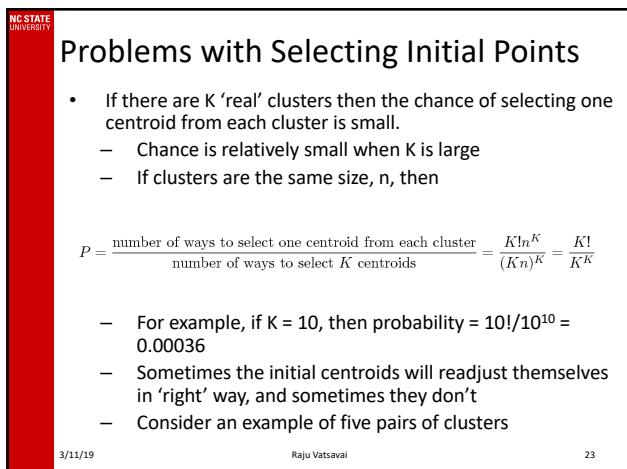
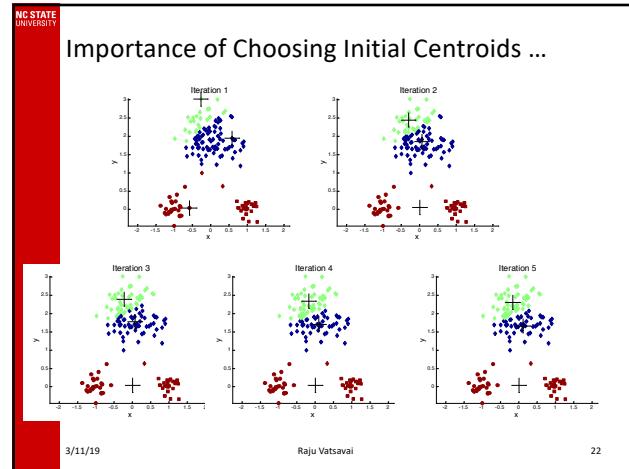
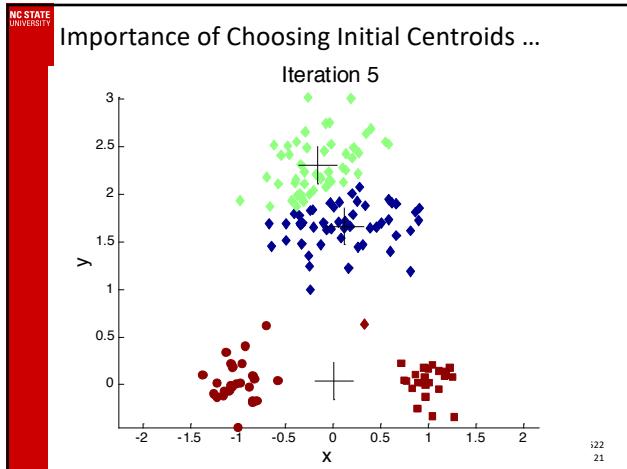
3/11/19

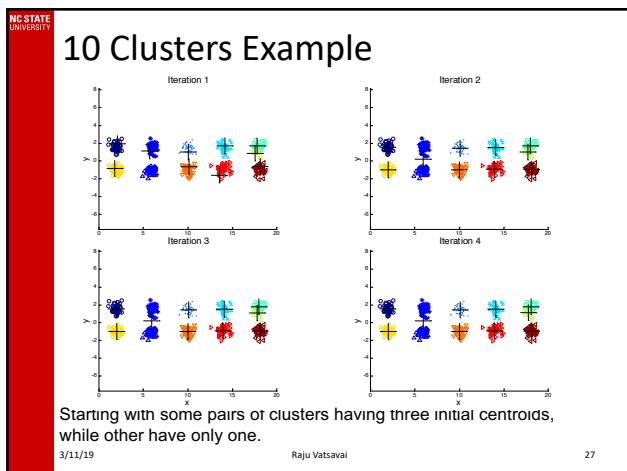
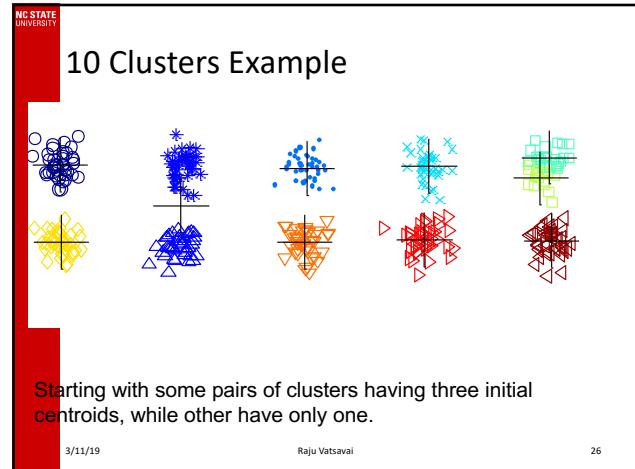
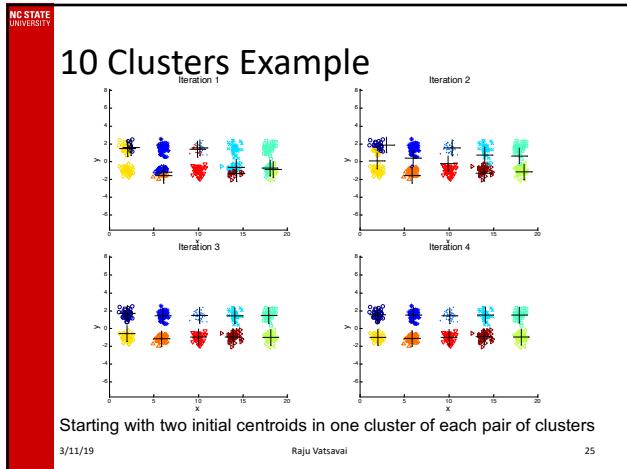
Raju Vatsavai

12









- Solutions to Initial Centroids Problem**
- Multiple runs
    - Helps, but probability is not on your side
  - Sample and use hierarchical clustering to determine initial centroids
  - Select more than K initial centroids and then select among these initial centroids
    - Select most widely separated
  - Post processing
  - Generate a larger number of clusters and then perform a hierarchical clustering
  - Bisection K-means
    - Not as susceptible to initialization issues
- 3/11/19 Raju Vatsavai 28

**Empty Clusters**

- K-means can yield empty clusters

3/11/19 Raju Vatsavai 29

## Handling Empty Clusters

- Basic K-means algorithm can yield empty clusters
- Several strategies
  - Choose the point that contributes most to SSE
  - Choose a replacement centroid from the cluster with the highest SSE
  - If there are several empty clusters, the above can be repeated several times.

3/11/19 Raju Vatsavai 30

**Updating Centers Incrementally**

- In the basic K-means algorithm, centroids are updated after all points are assigned to a centroid
- An alternative is to update the centroids after each assignment (incremental approach)
  - Each assignment updates zero or two centroids
  - More expensive
  - Introduces an order dependency
  - Never get an empty cluster
  - Relative weight of the point being added may be adjusted (decreased as clustering proceeds)

3/11/19 Raju Vatsavai 31

## Pre-processing and Post-processing

- Pre-processing
  - Normalize the data
  - Eliminate outliers
- Post-processing
  - Eliminate small clusters that may represent outliers
  - Split 'loose' clusters, i.e., clusters with relatively high SSE
  - Merge clusters that are 'close' and that have relatively low SSE
  - Can use these steps during the clustering process
    - ISODATA

3/11/19 Raju Vatsavai 32

## Bisecting K-means

- Bisecting K-means algorithm

- Variant of K-means that can produce a partitional or a hierarchical clustering

**Algorithm 8.2** Bisecting K-means algorithm.

```

1: Initialize the list of clusters to contain the cluster consisting of all points.
2: repeat
3:   Remove a cluster from the list of clusters.
4:   {Perform several “trial” bisections of the chosen cluster.}
5:   for  $i = 1$  to number of trials do
6:     Bisect the selected cluster using basic K-means.
7:   end for
8:   Select the two clusters from the bisection with the lowest total SSE.
9:   Add these two clusters to the list of clusters.
10: until Until the list of clusters contains  $K$  clusters.

```

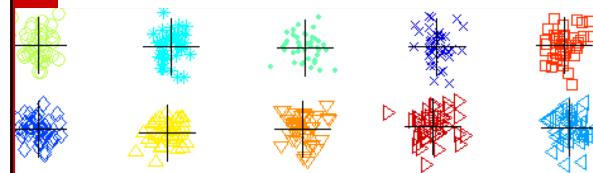
CLUTO: <http://glaros.dtc.umn.edu/gkhome/cluto/cluto/overview>

3/11/19

Raju Vatsavai

33

## Bisecting K-means Example



3/11/19

Raju Vatsavai

34

## Acknowledgements

- Slides are customized from: Introduction to Data Mining by Tan, Steinbach, Kumar

3/11/19

Raju Vatsavai

35