

[Clustering](#)

COMP90049  
Knowledge  
Technologies

[Clustering](#)  
[Definition](#)  
[Types](#)  
[Evaluation](#)

Methods  
[Similarity](#)  
[k-Means](#)  
[k-Ms limitation](#)  
[Hierarchal](#)

[Summary](#)  
[Resources](#)

## Lecture 12: Clustering

### COMP90049 Knowledge Technologies

Hasti Samadi & Sarah Erfani & Karin Verspoor, CIS

Semester 2, 2019



THE UNIVERSITY OF  
**MELBOURNE**



## Clustering

- Clustering
- Definition
- Types
- Evaluation
- Methods
  - Similarity
  - k-Means
  - k-Ms limitation
  - Hierarchical

Outlook	Temperature	Humidity	Windy
sunny	hot	high	F FALSE
sunny	hot	high	TRUE
overcast	hot	high	F FALSE
rainy	mild	high	F FALSE
rainy	cool	normal	F FALSE
rainy	cool	normal	TRUE
⋮	⋮	⋮	⋮

# Dataset with no label

## Clustering

COMP90049  
Knowledge  
Technologies

Clustering  
Definition  
Types  
Evaluation

Methods  
Similarity  
k-Means  
k-Ms limitation  
Hierarchical

Summary  
Resources

Outlook	Temperature	Humidity	Windy
sunny	hot	high	F FALSE
sunny	hot	high	TRUE
overcast	hot	high	F FALSE
rainy	mild	high	F FALSE
rainy	cool	normal	F FALSE
rainy	cool	normal	TRUE
:	:	:	:

- What can we do with a data set without labels?

## Clustering

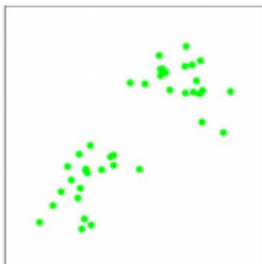
COMP90049  
Knowledge  
Technologies

Clustering  
Definition  
Types  
Evaluation

Methods  
Similarity  
k-Means  
k-Ms limitation  
Hierarchical

Summary  
Resources

- We can find groups of data in our dataset which are **similar (close)** to one another -- what we call **clusters**.



## Clustering

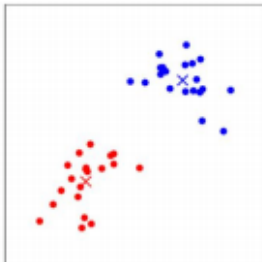
COMP90049  
Knowledge  
Technologies

Clustering  
Definition  
Types  
Evaluation

Methods  
Similarity  
k-Means  
k-Ms limitation  
Hierarchical

Summary  
Resources

- We can find groups of data in our dataset which are **similar (close)** to one another -- what we call **clusters**.



## Clustering

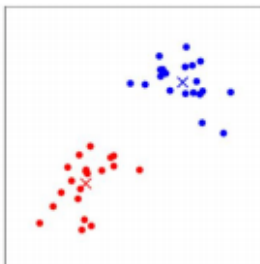
COMP90049  
Knowledge  
Technologies

Clustering  
Definition  
Types  
Evaluation

Methods  
Similarity  
k-Means  
k-Ms limitation  
Hierarchical

Summary  
Resources

- We can find groups of data in our dataset which are **similar (close)** to one another -- what we call **clusters**.



- In clustering it is also important that the objects (instances) in each cluster are **different from** (or unrelated to) the objects in other groups.

## Clustering

COMP90049  
Knowledge  
Technologies

Clustering  
Definition  
Types  
Evaluation

Methods  
Similarity  
k-Means  
k-Ms limitation  
Hierarchical

Summary  
Resources

- Clustering is an ***unsupervised learner***
- Applications in:
  - Pattern recognition
  - Spatial data analysis
  - Medical diagnosis
  - Marketing...

# A possible clustering of the weather dataset

## Clustering

COMP90049  
Knowledge  
Technologies

## Clustering

### Definition

### Types

### Evaluation

## Methods

### Similarity

### k-Means

### k-Ms limitation

### Hierarchical

## Summary

## Resources

Outlook	Temperature	Humidity	Windy	Cluster
sunny	hot	high	FALSE	?
sunny	hot	high	TRUE	?
overcast	hot	high	FALSE	?
rainy	mild	high	FALSE	?
rainy	cool	normal	FALSE	?
rainy	cool	normal	TRUE	?
overcast	cool	normal	TRUE	?
sunny	mild	high	FALSE	?
sunny	cool	normal	FALSE	?
rainy	mild	normal	FALSE	?
sunny	mild	normal	TRUE	?
overcast	mild	normal	TRUE	?
overcast	hot	high	FALSE	?
rainy	mild	high	TRUE	?



# A possible clustering of the weather dataset

## Clustering

COMP90049  
Knowledge  
Technologies

Clustering  
Definition  
Types  
Evaluation

Methods  
Similarity  
k-Means  
k-Ms limitation  
Hierarchical

Summary  
Resources

Outlook	Temperature	Humidity	Windy	Cluster
<b>sunny</b>	<b>hot</b>	<b>high</b>	FALSE	0
<b>sunny</b>	<b>hot</b>	<b>high</b>	TRUE	0
overcast	<b>hot</b>	<b>high</b>	FALSE	0
<b>rainy</b>	<b>mild</b>	high	FALSE	1
<b>rainy</b>	<b>cool</b>	<b>normal</b>	FALSE	1
<b>rainy</b>	<b>cool</b>	<b>normal</b>	TRUE	1
overcast	<b>cool</b>	<b>normal</b>	TRUE	1
<b>sunny</b>	mild	<b>high</b>	FALSE	0
sunny	<b>cool</b>	<b>normal</b>	FALSE	1
<b>rainy</b>	<b>mild</b>	<b>normal</b>	FALSE	1
sunny	<b>mild</b>	<b>normal</b>	TRUE	1
overcast	<b>mild</b>	<b>normal</b>	TRUE	1
overcast	<b>hot</b>	<b>high</b>	FALSE	0
<b>rainy</b>	<b>mild</b>	high	TRUE	1

# Clustering over the weather dataset (cf. outputs)

## Clustering

COMP90049  
Knowledge  
Technologies

Clustering  
Definition  
Types  
Evaluation

Methods  
Similarity  
k-Means  
k-Ms limitation  
Hierarchical

Summary  
Resources

Outlook	Temperature	Humidity	Windy	Cluster	Play
<b>sunny</b>	<b>hot</b>	<b>high</b>	FALSE	0	no ✓
<b>sunny</b>	<b>hot</b>	<b>high</b>	TRUE	0	no ✓
overcast	<b>hot</b>	<b>high</b>	FALSE	0	yes ✗
<b>rainy</b>	<b>mild</b>	high	FALSE	1	yes ✓
<b>rainy</b>	<b>cool</b>	<b>normal</b>	FALSE	1	yes ✓
<b>rainy</b>	<b>cool</b>	<b>normal</b>	TRUE	1	no ✗
overcast	<b>cool</b>	<b>normal</b>	TRUE	1	yes ✓
<b>sunny</b>	mild	<b>high</b>	FALSE	0	no ✓
sunny	<b>cool</b>	<b>normal</b>	FALSE	1	yes ✓
<b>rainy</b>	<b>mild</b>	<b>normal</b>	FALSE	1	yes ✓
sunny	<b>mild</b>	<b>normal</b>	TRUE	1	yes ✓
overcast	<b>mild</b>	<b>normal</b>	TRUE	1	yes ✓
overcast	<b>hot</b>	<b>high</b>	FALSE	0	yes ✗
<b>rainy</b>	<b>mild</b>	high	TRUE	1	no ✗

# Clustering, basic contrasts

## Clustering

COMP90049  
Knowledge  
Technologies

Clustering  
Definition  
Types  
Evaluation

Methods  
Similarity  
k-Means  
k-Ms limitation  
Hierarchical

Summary  
Resources

- Exclusive vs. overlapping clustering
  - Can an item be in more than one cluster?

# Clustering, basic contrasts

## Clustering

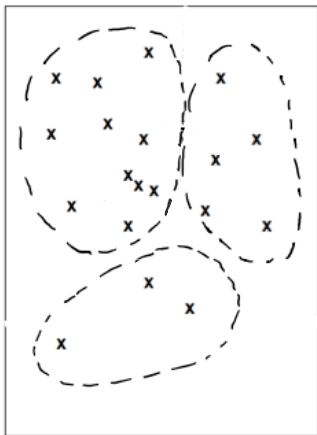
COMP90049  
Knowledge  
Technologies

Clustering  
Definition  
Types  
Evaluation

Methods  
Similarity  
k-Means  
k-Ms limitation  
Hierarchical

Summary  
Resources

- Exclusive vs. overlapping clustering
  - Can an item be in more than one cluster?



# Clustering, basic contrasts

## Clustering

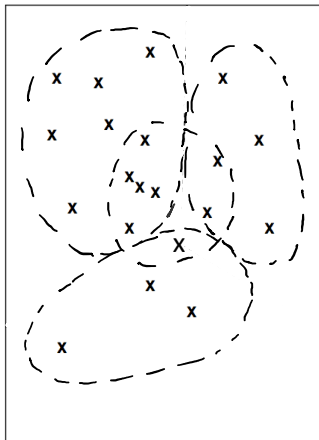
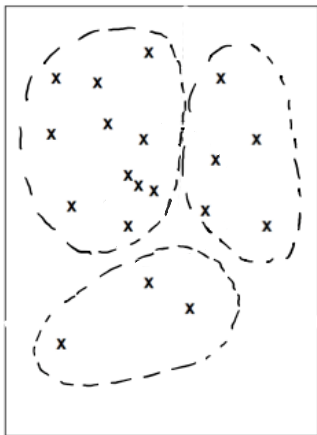
COMP90049  
Knowledge  
Technologies

Clustering  
Definition  
Types  
Evaluation

Methods  
Similarity  
k-Means  
k-Ms limitation  
Hierarchical

Summary  
Resources

- Exclusive vs. overlapping clustering
  - Can an item be in more than one cluster?



# Clustering, basic contrasts

## Clustering

COMP90049  
Knowledge  
Technologies

Clustering  
Definition  
Types  
Evaluation

Methods  
Similarity  
k-Means  
k-Ms limitation  
Hierarchical

Summary  
Resources

- Exclusive vs. overlapping clustering
  - Can an item be in more than one cluster?
- Deterministic vs. probabilistic clustering (Hard vs. soft clustering)
  - Can an item be partially or weakly in a cluster?

<i>Instance</i>	<i>Cluster</i>
1	2
2	3
3	2
4	1
5	2
6	2
7	4
:	:

<i>Instance</i>	Cluster			
	1	2	3	4
1	0.01	<b>0.87</b>	0.12	0.00
2	0.05	0.25	<b>0.67</b>	0.03
3	0.00	<b>0.98</b>	0.02	0.00
4	<b>0.45</b>	0.39	0.08	0.08
5	0.01	<b>0.99</b>	0.00	0.00
6	0.07	<b>0.75</b>	0.08	0.10
7	0.23	0.10	0.20	<b>0.47</b>
:	:	:	:	:

# Clustering, basic contrasts

## Clustering

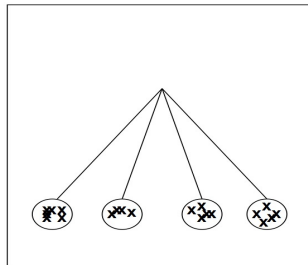
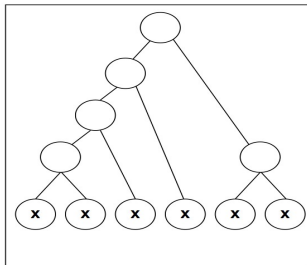
COMP90049  
Knowledge  
Technologies

[Clustering](#)  
[Definition](#)  
[Types](#)  
[Evaluation](#)

[Methods](#)  
[Similarity](#)  
[k-Means](#)  
[k-Ms limitation](#)  
[Hierarchical](#)

[Summary](#)  
[Resources](#)

- Exclusive vs. overlapping clustering
  - Can an item be in more than one cluster?
- Deterministic vs. probabilistic clustering (Hard vs. soft clustering)
  - Can an item be partially or weakly in a cluster?
- Hierarchical vs. partitioning clustering
  - Do the clusters have subset relationships between them?  
e.g. nested in a tree?



# Clustering, basic contrasts

## Clustering

COMP90049  
Knowledge  
Technologies

Clustering  
Definition  
Types  
Evaluation

Methods  
Similarity  
k-Means  
k-Ms limitation  
Hierarchical

Summary  
Resources

- Exclusive vs. overlapping clustering
  - Can an item be in more than one cluster?
- Deterministic vs. probabilistic clustering (Hard vs. soft clustering)
  - Can an item be partially or weakly in a cluster?
- Hierarchical vs. partitioning clustering
  - Do the clusters have subset relationships between them?  
e.g. nested in a tree?
- Partial vs. complete
  - In some cases, we only want to cluster some of the data



# Clustering, basic contrasts

## Clustering

COMP90049  
Knowledge  
Technologies

Clustering  
Definition  
Types  
Evaluation

Methods  
Similarity  
k-Means  
k-Ms limitation  
Hierarchical

Summary  
Resources

- Partial vs. complete
  - In some cases, we only want to cluster some of the data



# Clustering, basic contrasts

## Clustering

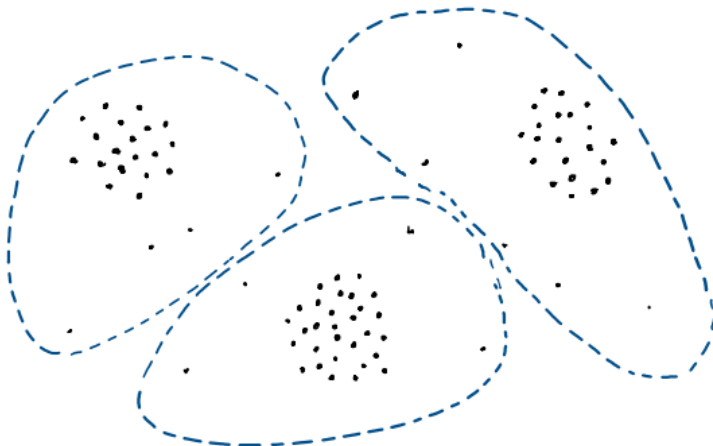
COMP90049  
Knowledge  
Technologies

Clustering  
Definition  
Types  
Evaluation

Methods  
Similarity  
k-Means  
k-Ms limitation  
Hierarchical

Summary  
Resources

- Partial vs. complete
  - In some cases, we only want to cluster some of the data



# Clustering, basic contrasts

## Clustering

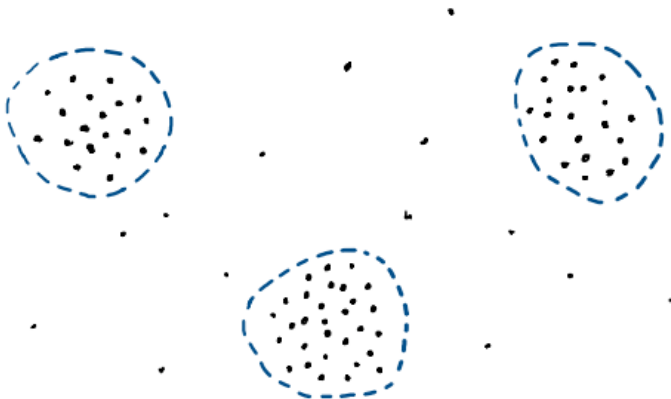
COMP90049  
Knowledge  
Technologies

Clustering  
Definition  
Types  
Evaluation

Methods  
Similarity  
k-Means  
k-Ms limitation  
Hierarchical

Summary  
Resources

- Partial vs. complete
  - In some cases, we only want to cluster some of the data



# Clustering, basic contrasts

## Clustering

COMP90049  
Knowledge  
Technologies

Clustering  
Definition  
Types  
Evaluation

Methods  
Similarity  
k-Means  
k-Ms limitation  
Hierarchical

Summary  
Resources

- Exclusive vs. overlapping clustering
  - Can an item be in more than one cluster?
- Deterministic vs. probabilistic clustering (Hard vs. soft clustering)
  - Can an item be partially or weakly in a cluster?
- Hierarchical vs. partitioning clustering
  - Do the clusters have subset relationships between them?  
e.g. nested in a tree?
- Partial vs. complete
  - In some cases, we only want to cluster some of the data
- Heterogenous vs. homogenous
  - Clusters of widely different sizes, shapes, and densities

# Heterogenous vs. homogenous

## Clustering

COMP90049  
Knowledge  
Technologies

## Clustering

### Definition

### Types

### Evaluation

## Methods

### Similarity

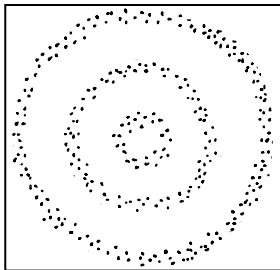
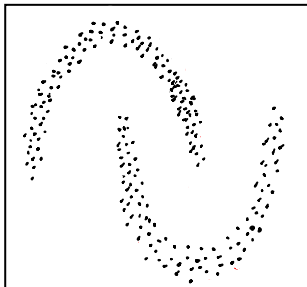
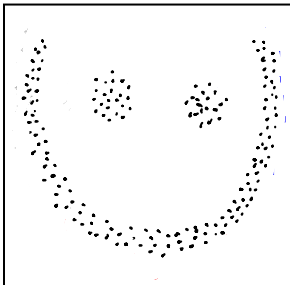
### k-Means

### k-Ms limitation

### Hierarchical

## Summary

## Resources



# Clustering, basic contrasts

## Clustering

COMP90049  
Knowledge  
Technologies

Clustering  
Definition  
Types  
Evaluation

Methods  
Similarity  
k-Means  
k-Ms limitation  
Hierarchical

Summary  
Resources

- Exclusive vs. overlapping clustering
  - Can an item be in more than one cluster?
- Deterministic vs. probabilistic clustering (Hard vs. soft clustering)
  - Can an item be partially or weakly in a cluster?
- Hierarchical vs. partitioning clustering
  - Do the clusters have subset relationships between them?  
e.g. nested in a tree?
- Partial vs. complete
  - In some cases, we only want to cluster some of the data
- Heterogenous vs. homogenous
  - Clusters of widely different sizes, shapes, and densities
- Incremental vs. batch clustering
  - Is the whole set of items clustered in one go?

# What is a good clustering?

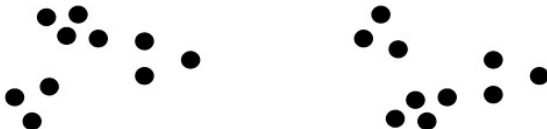
## Clustering

COMP90049  
Knowledge  
Technologies

## Clustering Definition Types Evaluation

## Methods Similarity k-Means k-Ms limitation Hierarchical

## Summary Resources



# Two clusters?

## Clustering

COMP90049  
Knowledge  
Technologies

## Clustering

### Definition

### Types

### Evaluation

## Methods

### Similarity

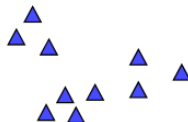
### k-Means

### k-Ms limitation

### Hierarchical

## Summary

## Resources





# Four clusters?

## Clustering

COMP90049  
Knowledge  
Technologies

## Clustering

### Definition

### Types

### Evaluation

## Methods

### Similarity

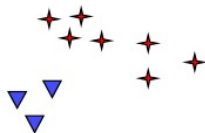
### k-Means

### k-Ms limitation

### Hierarchical

## Summary

## Resources



# Six clusters?

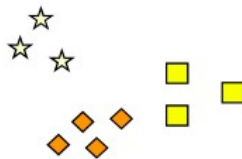
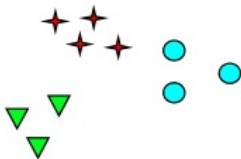
## Clustering

COMP90049  
Knowledge  
Technologies

Clustering  
Definition  
Types  
Evaluation

Methods  
Similarity  
k-Means  
k-Ms limitation  
Hierarchical

Summary  
Resources



# What number of clusters is better?

2 Clusters

4 Clusters

6 clusters

18 clusters

## Clustering

COMP90049  
Knowledge  
Technologies

Clustering  
Definition  
Types  
Evaluation

Methods  
Similarity  
k-Means  
k-Ms limitation  
Hierarchical

Summary  
Resources

Clustering evaluation measures come in two basic types:

- **Unsupervised:** Measures the goodness of a clustering structure without respect to external information.
  - How cohesive are individual clusters?
  - How separate is one cluster from other clusters?
- **Supervised:** how well do cluster labels match externally supplied class labels?

## Clustering

COMP90049  
Knowledge  
Technologies

## Clustering Definition Types Evaluation

Methods  
Similarity  
k-Means  
k-Ms limitation  
Hierarchical

Summary  
Resources

A “good” cluster analysis should have one or both of:

- **high cluster cohesion**, i.e. instances in a given cluster should be closely related to each other

$$cohesion(C_i) = \frac{1}{\sum_{x,y \in C_i} Distance(x,y)}$$

- **high cluster separation**, i.e. instances in different clusters should be distinct from each other

$$separation(C_i, C_j) = \sum_{x \in C_i, y \in C_j, i \neq j} Distance(x,y)$$

## Clustering

COMP90049  
Knowledge  
Technologies

Clustering  
Definition  
Types  
Evaluation

Methods  
Similarity  
k-Means  
k-Ms limitation  
Hierarchical

Summary  
Resources

- Most common measure evaluating the quality of clusters (esp. for k-means) is **Sum of Squared Error (SSE)** or *Scatter*

## Clustering

COMP90049  
Knowledge  
Technologies

Clustering  
Definition  
Types  
Evaluation

Methods  
Similarity  
k-Means  
k-Ms limitation  
Hierarchical

Summary  
Resources

## 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.

$$\sum_{i=1}^k \sum_{x \in C_i} \text{dist}^2(m_i, x)$$

- $x$  is a data point in cluster  $C_i$  and  $m_i$  is the representative point for cluster  $C_i$  (Centroid)
- Can show that the  $m_i$  that minimises SSE corresponds to the center (mean) of the cluster
- Given two clusters, we can choose the one with the smallest error
- One easy way to reduce SSE is to increase  $k$ , the number of clusters
- However, a good clustering with smaller  $k$  can have a lower SSE than a poor clustering with higher  $k$

# Two clusters?

## Clustering

COMP90049  
Knowledge  
Technologies

## Clustering

### Definition

### Types

### Evaluation

## Methods

### Similarity

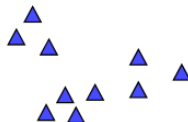
### k-Means

### k-Ms limitation

### Hierarchical

## Summary

## Resources





# Six clusters?

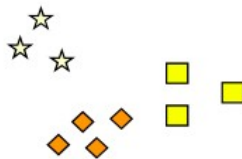
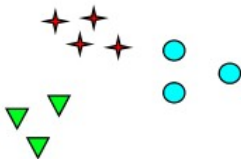
## Clustering

COMP90049  
Knowledge  
Technologies

Clustering  
Definition  
Types  
Evaluation

Methods  
Similarity  
k-Means  
k-Ms limitation  
Hierarchical

Summary  
Resources



## Clustering

COMP90049  
Knowledge  
Technologies

## Clustering

### Definition

### Types

### Evaluation

## Methods

### Similarity

### k-Means

### k-Ms limitation

### Hierarchical

## Summary

## Resources

- **Supervised** evaluation of cluster “validity” measures the degree to which predicted class labels match the actual class labels, e.g. based on the distribution of actual class labels within each cluster.

# Clustering over the weather dataset (cf. outputs)

## Clustering

COMP90049  
Knowledge  
Technologies

Clustering  
Definition  
Types  
Evaluation

Methods  
Similarity  
k-Means  
k-Ms limitation  
Hierarchical

Summary  
Resources

Outlook	Temperature	Humidity	Windy	Cluster	Play
<b>sunny</b>	<b>hot</b>	<b>high</b>	FALSE	0	no ✓
<b>sunny</b>	<b>hot</b>	<b>high</b>	TRUE	0	no ✓
overcast	<b>hot</b>	<b>high</b>	FALSE	0	yes ✗
<b>rainy</b>	<b>mild</b>	high	FALSE	1	yes ✓
<b>rainy</b>	<b>cool</b>	<b>normal</b>	FALSE	1	yes ✓
<b>rainy</b>	<b>cool</b>	<b>normal</b>	TRUE	1	no ✗
overcast	<b>cool</b>	<b>normal</b>	TRUE	1	yes ✓
<b>sunny</b>	mild	<b>high</b>	FALSE	0	no ✓
sunny	<b>cool</b>	<b>normal</b>	FALSE	1	yes ✓
<b>rainy</b>	<b>mild</b>	<b>normal</b>	FALSE	1	yes ✓
sunny	<b>mild</b>	<b>normal</b>	TRUE	1	yes ✓
overcast	<b>mild</b>	<b>normal</b>	TRUE	1	yes ✓
overcast	<b>hot</b>	<b>high</b>	FALSE	0	yes ✗
<b>rainy</b>	<b>mild</b>	high	TRUE	1	no ✗

## [Clustering](#)

COMP90049  
Knowledge  
Technologies

## [Clustering](#) [Definition](#) [Types](#) [Evaluation](#)

Methods  
[Similarity](#)  
[k-Means](#)  
[k-Ms limitation](#)  
[Hierarchal](#)

[Summary](#)  
[Resources](#)

Based our definition:

- Clustering is finding groups of data in our dataset which are **similar or close** to one another and **different or separated** from other clusters
- A key component of any clustering algorithm is a measurement of the distance between any points.

## Clustering

COMP90049  
Knowledge  
Technologies

## Clustering Definition Types Evaluation

Methods  
Similarity  
k-Means  
k-Ms limitation  
Hierarchical

Summary  
Resources

- **Data points** (in Euclidean space)
  - Euclidean (L2) distance
  - Manhattan (L1) distance

## Clustering

COMP90049  
Knowledge  
Technologies

## Clustering Definition Types Evaluation

Methods  
Similarity  
k-Means  
k-Ms limitation  
Hierarchical

Summary  
Resources

- **Data points** (in Euclidean space)
  - Euclidean (L2) distance
  - Manhattan (L1) distance
- **Discrete values**
  - Hamming distance (discrepancy between the bit strings)

Sunny            011

Overcast        101

Rainy            110

For two bit strings, the number of positions at which the corresponding symbols are different

## Clustering

COMP90049  
Knowledge  
Technologies

## Clustering Definition Types Evaluation

Methods  
Similarity  
k-Means  
k-Ms limitation  
Hierarchical

Summary  
Resources

- **Data points** (in Euclidean space)
  - Euclidean (L2) distance
  - Manhattan (L1) distance
- **Discrete values**
  - Hamming distance (discrepancy between the bit strings)

Sunny            011

Overcast        101

Rainy            110

For two bit strings, the number of positions at which the corresponding symbols are different

011

101

## Clustering

COMP90049  
Knowledge  
Technologies

## Clustering Definition Types Evaluation

Methods  
Similarity  
k-Means  
k-Ms limitation  
Hierarchical

Summary  
Resources

- **Data points** (in Euclidean space)
  - Euclidean (L2) distance
  - Manhattan (L1) distance
- **Discrete values**
  - Hamming distance (discrepancy between the bit strings)

Sunny                      011

Overcast                101

Rainy                    110

For two bit strings, the number of positions at which the corresponding symbols are different

011                      Hamming Distance = 2  
101



# Measuring the distance

## Clustering

COMP90049  
Knowledge  
Technologies

## Clustering Definition Types Evaluation

Methods  
Similarity  
k-Means  
k-Ms limitation  
Hierarchical

Summary  
Resources

- **Data points** (in Euclidean space)
  - Euclidean (L2) distance
  - Manhattan (L1) distance
- **Discrete values**
  - Hamming distance (discrepancy between the bit strings)
- **Documents**
  - Cosine similarity
  - Jaccard measure

# Measuring the distance

## Clustering

COMP90049  
Knowledge  
Technologies

## Clustering Definition Types Evaluation

Methods  
Similarity  
k-Means  
k-Ms limitation  
Hierarchical

Summary  
Resources

- **Data points** (in Euclidean space)
  - Euclidean (L2) distance
  - Manhattan (L1) distance
- **Discrete values**
  - Hamming distance (discrepancy between the bit strings)
- **Documents**
  - Cosine similarity
  - Jaccard measure
- **Other measures**
  - Correlation
  - Graph-based measures

## Clustering

COMP90049  
Knowledge  
Technologies

## Clustering Definition Types Evaluation

## Methods Similarity k-Means k-Ms limitation Hierarchical

## Summary Resources

- **K-Means** is one of the most popular "clustering" algorithms.
- Group the dataset into **K** clusters, with the use of **K centroids**

## Clustering

COMP90049  
Knowledge  
Technologies

## Clustering Definition Types Evaluation

## Methods Similarity k-Means k-Ms limitation Hierarchical

## Summary Resources

Given **k**, the k-means algorithm is implemented as follows:

1. Select K points at random to act as **seed** clusters ( $\mu_1, \dots, \mu_k$ ), the initial **centroids**
2. **repeat**
3.     Assign each instance to the cluster with **nearest** centroid
4.     Recompute the centroids of each clusters (the centroid is the center, i.e. *mean* point of the cluster)
5. **Until** the centroids don't change

# Example, Iterations

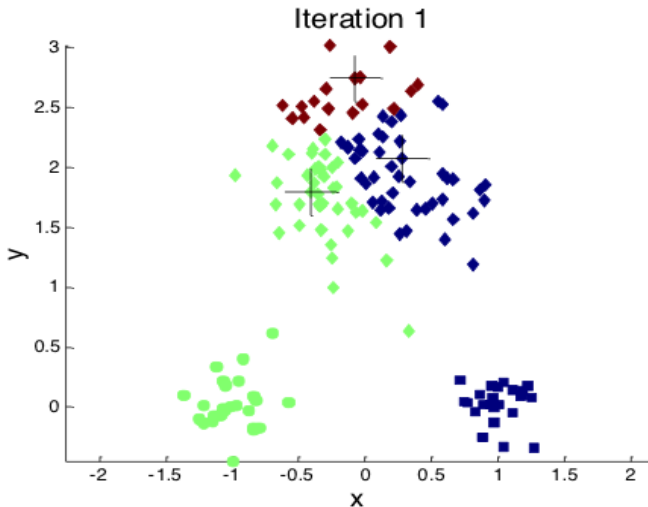
## Clustering

COMP90049  
Knowledge  
Technologies

Clustering  
Definition  
Types  
Evaluation

Methods  
Similarity  
k-Means  
k-Ms limitation  
Hierarchical

Summary  
Resources



# Example, Iterations

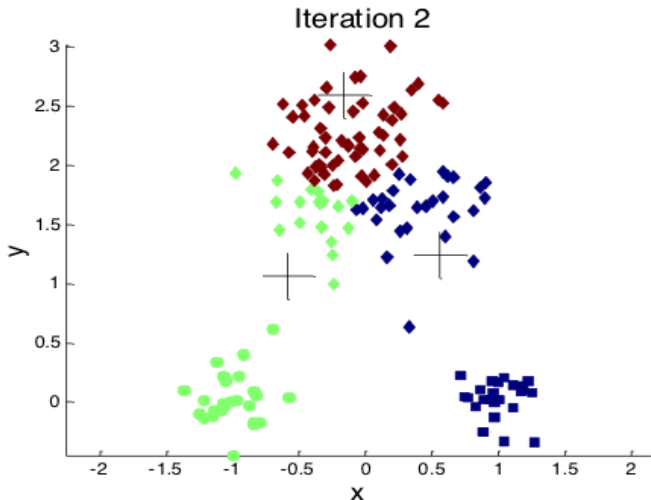
## Clustering

COMP90049  
Knowledge  
Technologies

Clustering  
Definition  
Types  
Evaluation

Methods  
Similarity  
k-Means  
k-Ms limitation  
Hierarchical

Summary  
Resources



# Example, Iterations

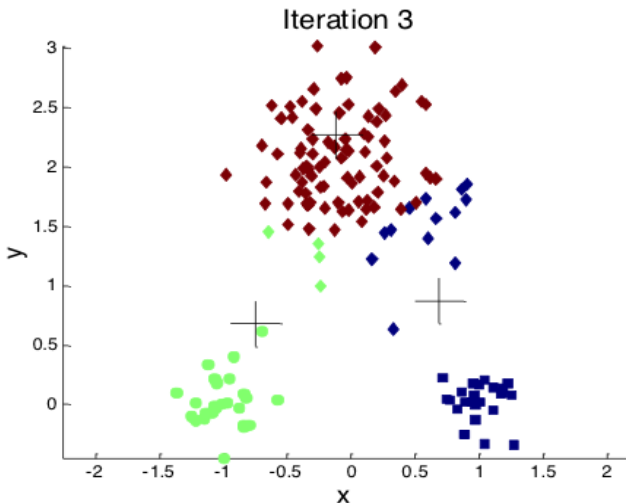
## Clustering

COMP90049  
Knowledge  
Technologies

Clustering  
Definition  
Types  
Evaluation

Methods  
Similarity  
k-Means  
k-Ms limitation  
Hierarchical

Summary  
Resources



# Example, Iterations

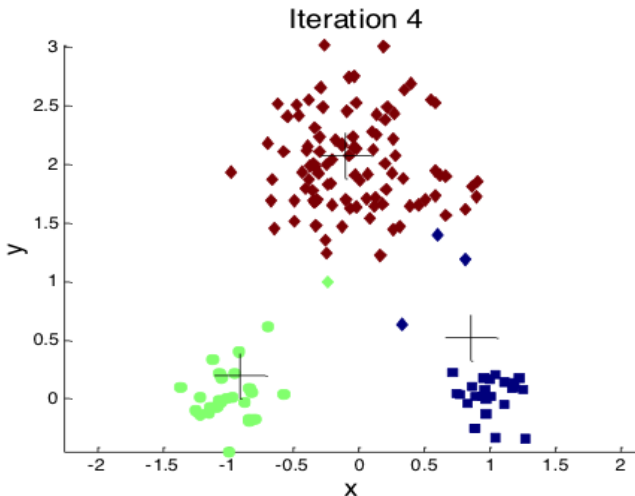
## Clustering

COMP90049  
Knowledge  
Technologies

Clustering  
Definition  
Types  
Evaluation

Methods  
Similarity  
k-Means  
k-Ms limitation  
Hierarchical

Summary  
Resources





# Example, Iterations

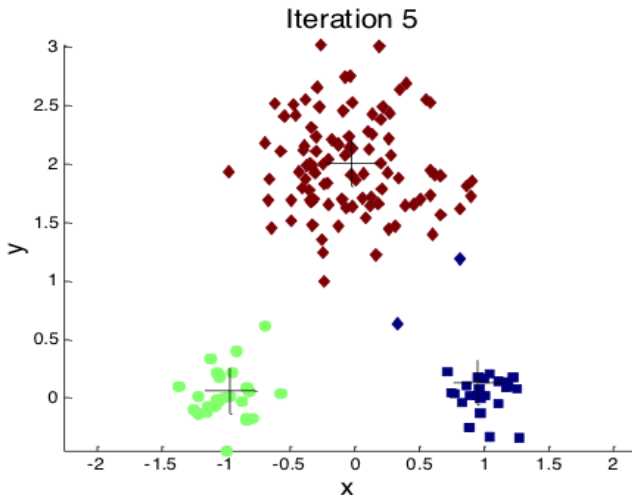
## Clustering

COMP90049  
Knowledge  
Technologies

Clustering  
Definition  
Types  
Evaluation

Methods  
Similarity  
k-Means  
k-Ms limitation  
Hierarchical

Summary  
Resources



# Example, Iterations

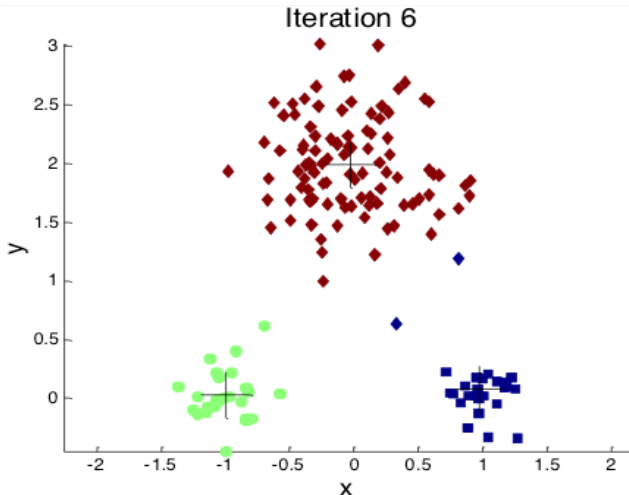
## Clustering

COMP90049  
Knowledge  
Technologies

Clustering  
Definition  
Types  
Evaluation

Methods  
Similarity  
k-Means  
k-Ms limitation  
Hierarchical

Summary  
Resources



# k-means Clustering – Details

## Clustering

COMP90049  
Knowledge  
Technologies

Clustering  
Definition  
Types  
Evaluation

Methods  
Similarity  
k-Means  
k-Ms limitation  
Hierarchical

Summary  
Resources

- Initial centroids are often chosen randomly.
  - Clusters produced vary from one run to another.

# k-means Clustering – Details

## Clustering

COMP90049  
Knowledge  
Technologies

Clustering  
Definition  
Types  
Evaluation

Methods  
Similarity  
k-Means  
k-Ms limitation  
Hierarchical

Summary  
Resources

- 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.

# k-means Clustering – Details

## Clustering

COMP90049  
Knowledge  
Technologies

Clustering  
Definition  
Types  
Evaluation

Methods  
Similarity  
k-Means  
k-Ms limitation  
Hierarchical

Summary  
Resources

- 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.
- ‘Nearest’ is based on proximity/similarity/distance metric.

# k-means Clustering – Details

## Clustering

COMP90049  
Knowledge  
Technologies

Clustering  
Definition  
Types  
Evaluation

Methods  
Similarity  
k-Means  
k-Ms limitation  
Hierarchical

Summary  
Resources

- 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.
- ‘Nearest’ is based on proximity/similarity/distance metric.
- 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’

# *k*-means, Pros and Cons

## Clustering

COMP90049  
Knowledge  
Technologies

Clustering  
Definition  
Types  
Evaluation

Methods  
Similarity  
k-Means  
k-Ms limitation  
Hierarchical

Summary  
Resources

- Exclusive, deterministic, partitioning, batch clustering method

## [Clustering](#)

COMP90049  
Knowledge  
Technologies

[Clustering](#)  
[Definition](#)  
[Types](#)  
[Evaluation](#)

Methods  
[Similarity](#)  
[k-Means](#)  
[k-Ms limitation](#)  
[Hierarchical](#)

[Summary](#)  
[Resources](#)

- Exclusive, deterministic, partitioning, batch clustering method

## Strengths:

- relatively efficient:

$O(ndki)$ , where

- $n$  is number of *instances*,
  - $d$  is number of *attributes*,
  - $k$  is number of *clusters*,
  - $i$  is number of *iterations*; normally  $k, i \ll n$
- Unfortunately we know the value of  $i$  in advance



## Clustering

COMP90049  
Knowledge  
Technologies

Clustering  
Definition  
Types  
Evaluation

Methods  
Similarity  
k-Means  
k-Ms limitation  
Hierarchical

Summary  
Resources

## Weaknesses:

- Need to specify  $k$  in advance

## Clustering

COMP90049  
Knowledge  
Technologies

Clustering  
Definition  
Types  
Evaluation

Methods  
Similarity  
k-Means  
k-Ms limitation  
Hierarchical

Summary  
Resources

## Weaknesses:

- Need to specify  $k$  in advance
- “Mean” is ill-defined for nominal or categorical attributes

## Clustering

COMP90049  
Knowledge  
Technologies

Clustering  
Definition  
Types  
Evaluation

Methods  
Similarity  
k-Means  
k-Ms limitation  
Hierarchical

Summary  
Resources

## Weaknesses:

- Need to specify  $k$  in advance
- “Mean” is ill-defined for nominal or categorical attributes
- May not work well when the data contains outliers

## Clustering

COMP90049  
Knowledge  
Technologies

Clustering  
Definition  
Types  
Evaluation

Methods  
Similarity  
k-Means  
k-Ms limitation  
Hierarchical

Summary  
Resources

## Weaknesses:

- Need to specify  $k$  in advance
- “Mean” is ill-defined for nominal or categorical attributes
- May not work well when the data contains outliers
- Tends to converge to local minimum; sensitive to seed instances

# Example, Impact of initial seeds

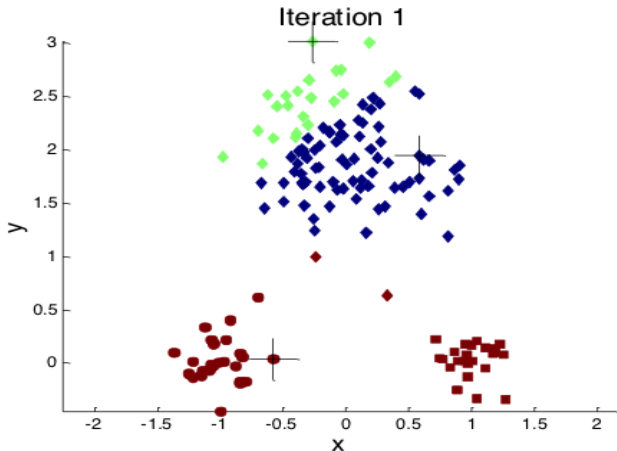
## Clustering

COMP90049  
Knowledge  
Technologies

Clustering  
Definition  
Types  
Evaluation

Methods  
Similarity  
k-Means  
k-Ms limitation  
Hierarchical

Summary  
Resources



# Example, Impact of initial seeds

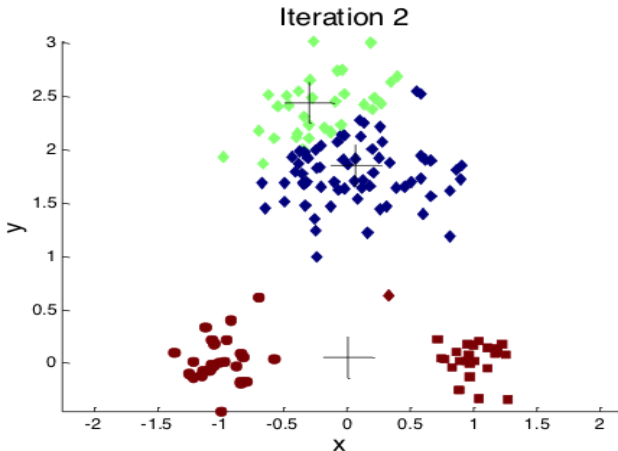
## Clustering

COMP90049  
Knowledge  
Technologies

Clustering  
Definition  
Types  
Evaluation

Methods  
Similarity  
k-Means  
k-Ms limitation  
Hierarchical

Summary  
Resources



# Example, Impact of initial seeds

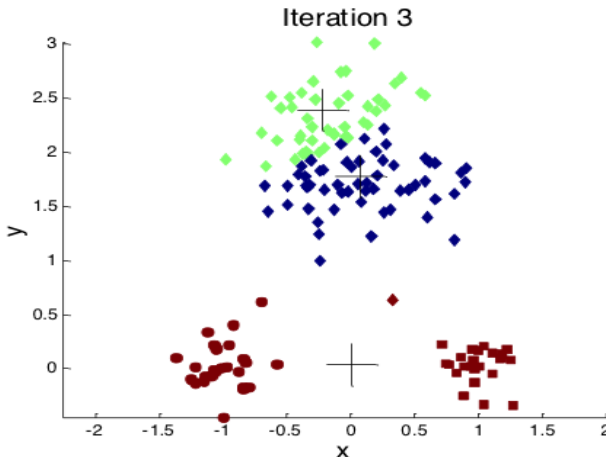
## Clustering

COMP90049  
Knowledge  
Technologies

Clustering  
Definition  
Types  
Evaluation

Methods  
Similarity  
k-Means  
k-Ms limitation  
Hierarchical

Summary  
Resources



# Example, Impact of initial seeds

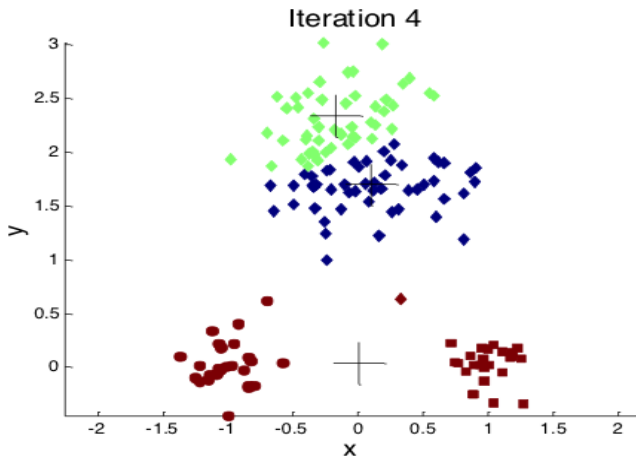
## Clustering

COMP90049  
Knowledge  
Technologies

Clustering  
Definition  
Types  
Evaluation

Methods  
Similarity  
k-Means  
k-Ms limitation  
Hierarchical

Summary  
Resources





# Example, Impact of initial seeds

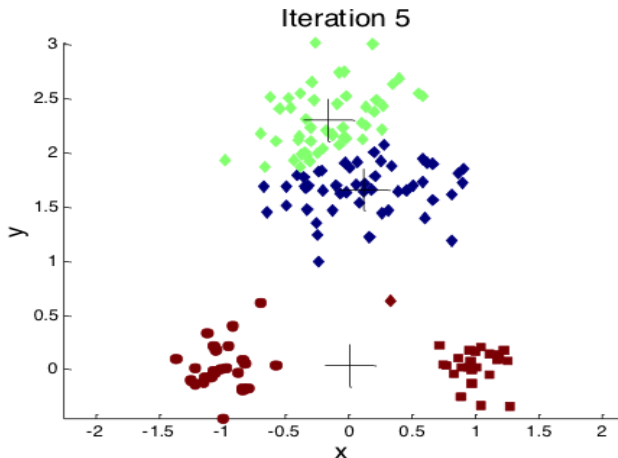
## Clustering

COMP90049  
Knowledge  
Technologies

Clustering  
Definition  
Types  
Evaluation

Methods  
Similarity  
k-Means  
k-Ms limitation  
Hierarchical

Summary  
Resources



## Clustering

COMP90049  
Knowledge  
Technologies

Clustering  
Definition  
Types  
Evaluation

Methods  
Similarity  
k-Means  
k-Ms limitation  
Hierarchical

Summary  
Resources

- K-means has problems when clusters are different in
  - Sizes
  - Densities

# Other limitations of $k$ -means

## Clustering

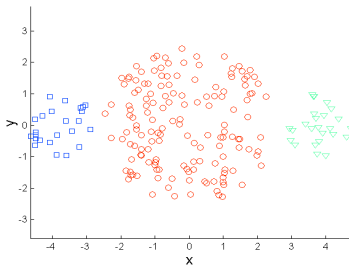
COMP90049  
Knowledge  
Technologies

Clustering  
Definition  
Types  
Evaluation

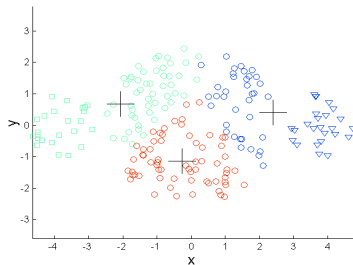
Methods  
Similarity  
k-Means  
k-Ms limitation  
Hierarchical

Summary  
Resources

## Clusters with different sizes:



Intended clusters



K-means (3 Clusters)

## Clustering

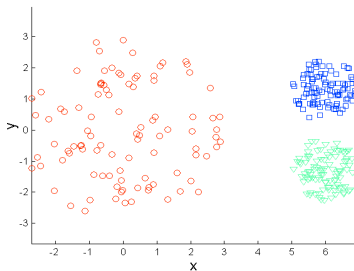
COMP90049  
Knowledge  
Technologies

Clustering  
Definition  
Types  
Evaluation

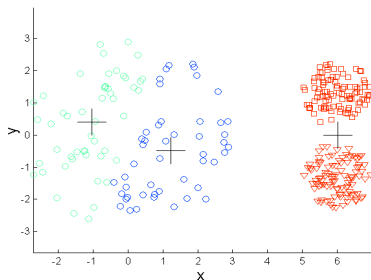
Methods  
Similarity  
 $k$ -Means  
 $k$ -Ms limitation  
Hierarchical

Summary  
Resources

## Clusters with different Densities:



Intended clusters



K-means (3 Clusters)

## Clustering

COMP90049  
Knowledge  
Technologies

Clustering  
Definition  
Types  
Evaluation

Methods  
Similarity  
k-Means  
k-Ms limitation  
Hierarchical

Summary  
Resources

- K-means has problems when clusters are different in
  - Sizes
  - Densities
- K-means does not work well on non-globular shapes

## Clustering

COMP90049  
Knowledge  
Technologies

## Clustering

### Definition

### Types

### Evaluation

## Methods

### Similarity

### k-Means

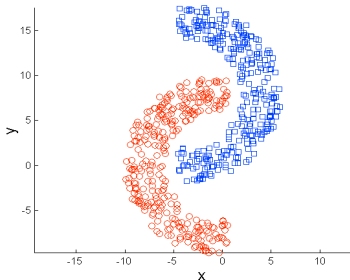
### k-Ms limitation

### Hierarchical

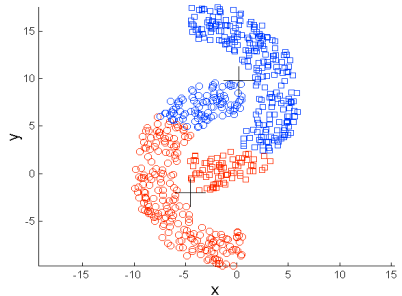
## Summary

## Resources

## Clusters with non-globular Shapes :



Intended clusters



K-means (2 Clusters)

# Other limitations of $k$ -means

## Clustering

COMP90049  
Knowledge  
Technologies

## Clustering

### Definition

### Types

### Evaluation

## Methods

### Similarity

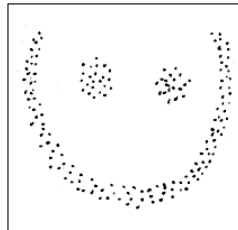
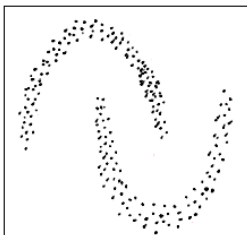
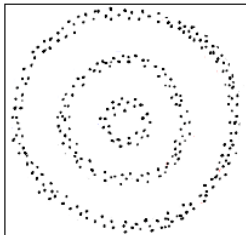
### k-Means

### k-Ms limitation

### Hierarchical

## Summary

## Resources



# Other limitations of $k$ -means

## Clustering

COMP90049  
Knowledge  
Technologies

## Clustering

Definition

Types

Evaluation

## Methods

Similarity

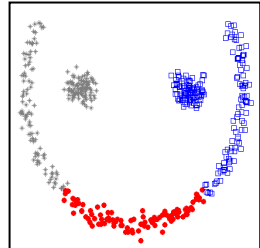
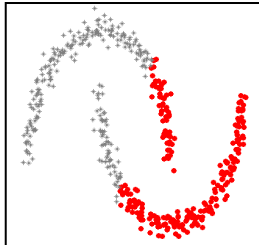
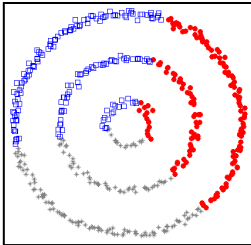
$k$ -Means

$k$ -Ms limitation

Hierarchical

## Summary

Resources





# Other limitations of $k$ -means

## Clustering

COMP90049  
Knowledge  
Technologies

## Clustering

### Definition

### Types

### Evaluation

## Methods

### Similarity

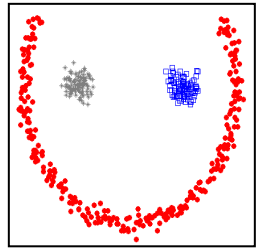
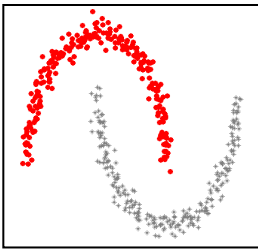
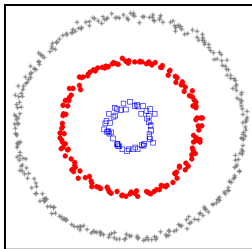
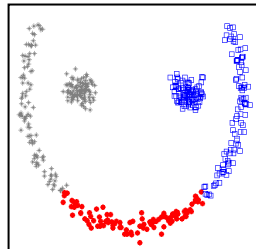
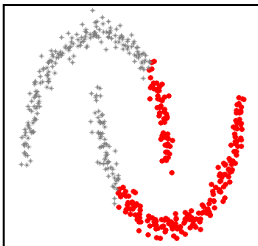
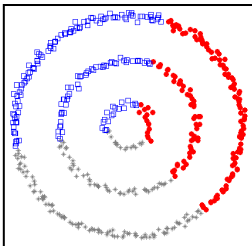
### $k$ -Means

### $k$ -Ms limitation

### Hierarchical

## Summary

## Resources



## Clustering

COMP90049  
Knowledge  
Technologies

Clustering  
Definition  
Types  
Evaluation

Methods  
Similarity  
k-Means  
k-Ms limitation  
Hierarchcal

Summary  
Resources

## Clustering

COMP90049  
Knowledge  
Technologies

Clustering  
Definition  
Types  
Evaluation

Methods  
Similarity  
k-Means  
k-Ms limitation  
Hierarchical

Summary  
Resources

## Bottom-up (= agglomerative) clustering

- Start with single-instance clusters
- At each step, join the two closest clusters (in terms of margin between clusters, distance between mean, ...)

## Clustering

COMP90049  
Knowledge  
Technologies

Clustering  
Definition  
Types  
Evaluation

Methods  
Similarity  
k-Means  
k-Ms limitation  
Hierarchical

Summary  
Resources

## Bottom-up (= agglomerative) clustering

- Start with single-instance clusters
- At each step, join the two closest clusters (in terms of margin between clusters, distance between mean, ...)

## Top-down (= divisive) clustering

- Start with one universal cluster
- Find two partitioning clusters
- Proceed recursively on each subset
- Can be very fast

## Clustering

COMP90049  
Knowledge  
Technologies

Clustering  
Definition  
Types  
Evaluation

Methods  
Similarity  
k-Means  
k-Ms limitation  
Hierarchical

Summary  
Resources

## Bottom-up (= agglomerative) clustering

- Start with single-instance clusters
- At each step, join the two closest clusters (in terms of margin between clusters, distance between mean, ...)

## Top-down (= divisive) clustering

- Start with one universal cluster
  - Find two partitioning clusters
  - Proceed recursively on each subset
  - Can be very fast
- In contrast to  $k$ -means clustering, hierarchical clustering only requires a measure of similarity between *groups* of data points (no seeds, no  $k$  value).

## Clustering

COMP90049  
Knowledge  
Technologies

Clustering  
Definition  
Types  
Evaluation

Methods  
Similarity  
k-Means  
k-Ms limitation  
Hierarchical

Summary  
Resources

Compute the proximity matrix, if necessary.

**repeat**

- Merge the closest two clusters
- Update the proximity matrix to reflect the proximity between the new cluster and the original clusters

**until** Only one cluster remains

# Example, Step 1

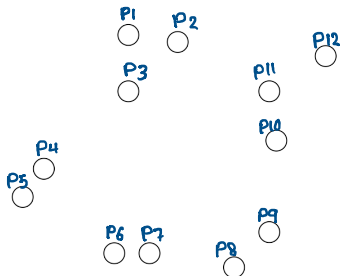
## Clustering

COMP90049  
Knowledge  
Technologies

Clustering  
Definition  
Types  
Evaluation

Methods  
Similarity  
k-Means  
k-Ms limitation  
Hierarchical

Summary  
Resources



	p1	p2	p3	p4	p5	...
p1						
p2						
p3						
p4						
p5						
...						

Proximity Matrix



# Example, Step 2

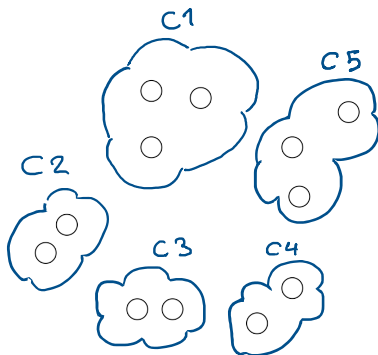
## Clustering

COMP90049  
Knowledge  
Technologies

Clustering  
Definition  
Types  
Evaluation

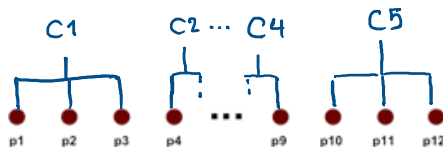
Methods  
Similarity  
k-Means  
k-Ms limitation  
Hierarchical

Summary  
Resources



	C1	C2	C3	C4	C5
C1					
C2					
C3					
C4					
C5					

Proximity Matrix





# Example, Step 3

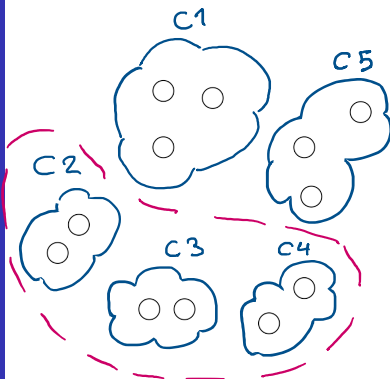
## Clustering

COMP90049  
Knowledge  
Technologies

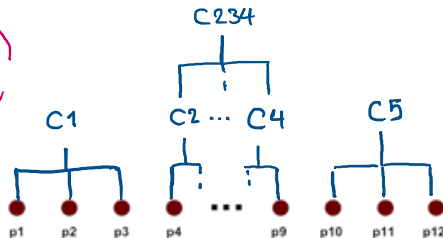
Clustering  
Definition  
Types  
Evaluation

Methods  
Similarity  
k-Means  
k-Ms limitation  
Hierarchical

Summary  
Resources



	c1	c234	c5
c1			
c234			
c5			



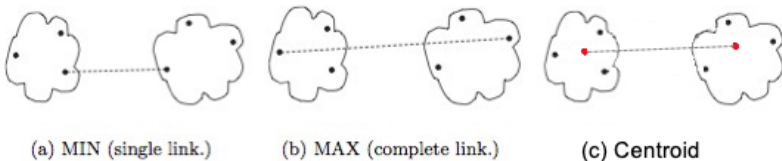
## Clustering

COMP90049  
Knowledge  
Technologies

Clustering  
Definition  
Types  
Evaluation

Methods  
Similarity  
k-Means  
k-Ms limitation  
Hierarchal

Summary  
Resources



Updating the proximity matrix:

- **Single Link:** *Minimum* distance between any two points in the two clusters. (most similar members)
- **Complete Link:** *Maximum* distance between any two points in the two clusters. (most dissimilar members)
- **Centroid:** Distance between the centroids of each cluster

## Clustering

COMP90049  
Knowledge  
Technologies

## Clustering

### Definition

### Types

### Evaluation

## Methods

### Similarity

### k-Means

### k-Ms limitation

### Hierarchical

## Summary

## Resources

	1	2	3	4	5
1	-	0.90	0.10	0.75	0.20
2	0.90	-	0.80	0.60	0.50
3	0.10	0.80	-	0.40	0.30
4	0.75	0.60	0.40	-	0.70
5	0.20	0.50	0.30	0.70	-

Based on this similarity matrix, what are the two closest points?

# Agglomerative Clustering Example

## Clustering

COMP90049  
Knowledge  
Technologies

## Clustering

### Definition

### Types

### Evaluation

## Methods

### Similarity

### k-Means

### k-Ms limitation

### Hierarchical

## Summary

## Resources

	1	2	3	4	5
1	-	0.90	0.10	0.75	0.20
2	0.90	-	0.80	0.60	0.50
3	0.10	0.80	-	0.40	0.30
4	0.75	0.60	0.40	-	0.70
5	0.20	0.50	0.30	0.70	-

Based on this similarity matrix, what are the two closest points?

# Agglomerative Clustering Example

## Clustering

COMP90049  
Knowledge  
Technologies

[Clustering](#)  
[Definition](#)  
[Types](#)  
[Evaluation](#)

Methods  
[Similarity](#)  
[k-Means](#)  
[k-Ms limitation](#)  
[Hierarchal](#)

[Summary](#)  
[Resources](#)

Merge points 1 & 2 into a new cluster: 12

**Update Single Link:** update the approximate matrix based on the most similar members

	1	2	3	4	5
1	-	0.90	0.10	0.75	0.20
2	0.90	-	0.80	0.60	0.50
3	0.10	0.80	-	0.40	0.30
4	0.75	0.60	0.40	-	0.70
5	0.20	0.50	0.30	0.70	-

	12	3	4	5
12				
3				
4				
5				

# Agglomerative Clustering Example

## Clustering

COMP90049  
Knowledge  
Technologies

[Clustering](#)  
[Definition](#)  
[Types](#)  
[Evaluation](#)

Methods  
[Similarity](#)  
[k-Means](#)  
[k-Ms limitation](#)  
[Hierarchal](#)

[Summary](#)  
[Resources](#)

Merge points 1 & 2 into a new cluster: 12

**Update Single Link:** update the approximate matrix based on the most similar members

	1	2	3	4	5
1	-	0.90	<b>0.10</b>	0.75	0.20
2	0.90	-	<b>0.80</b>	0.60	0.50
3	<b>0.10</b>	<b>0.80</b>	-	0.40	0.30
4	0.75	0.60	0.40	-	0.70
5	0.20	0.50	0.30	0.70	-

	12	3	4	5
12				
3				
4				
5				

# Agglomerative Clustering Example

## Clustering

COMP90049  
Knowledge  
Technologies

[Clustering](#)  
[Definition](#)  
[Types](#)  
[Evaluation](#)

Methods  
[Similarity](#)  
[k-Means](#)  
[k-Ms limitation](#)  
[Hierarchical](#)

[Summary](#)  
[Resources](#)

Merge points 1 & 2 into a new cluster: 12

**Update Single Link:** update the approximate matrix based on the most similar members

	1	2	3	4	5
1	-	0.90	<b>0.10</b>	0.75	0.20
2	0.90	-	<b>0.80</b>	0.60	0.50
3	<b>0.10</b>	<b>0.80</b>	-	0.40	0.30
4	0.75	0.60	0.40	-	0.70
5	0.20	0.50	0.30	0.70	-

	12	3	4	5
12	-	<b>0.80</b>		
3	<b>0.80</b>	-		
4			-	
5				-

# Agglomerative Clustering Example

## Clustering

COMP90049  
Knowledge  
Technologies

Clustering  
Definition  
Types  
Evaluation

Methods  
Similarity  
k-Means  
k-Ms limitation  
Hierarchical

Summary  
Resources

Merge points 1 & 2 into a new cluster: 12

**Update Single Link:** update the approximate matrix based on the most similar members

	1	2	3	4	5
1	-	0.90	0.10	<b>0.75</b>	0.20
2	0.90	-	0.80	<b>0.60</b>	0.50
3	0.10	0.80	-	0.40	0.30
4	<b>0.75</b>	<b>0.60</b>	0.40	-	0.70
5	0.20	0.50	0.30	0.70	-

	12	3	4	5
12	-	0.80	<b>0.75</b>	
3	0.80	-	0.40	
4	<b>0.75</b>	0.40	-	
5				-



# Agglomerative Clustering Example

## Clustering

COMP90049  
Knowledge  
Technologies

Clustering  
Definition  
Types  
Evaluation

Methods  
Similarity  
k-Means  
k-Ms limitation  
Hierarchical

Summary  
Resources

Merge points 1 & 2 into a new cluster: 12

**Update Single Link:** update the approximate matrix based on the most similar members

	1	2	3	4	5
1	-	0.90	0.10	0.75	<b>0.20</b>
2	0.90	-	0.80	0.60	<b>0.50</b>
3	0.10	0.80	-	0.40	0.30
4	0.75	0.60	0.40	-	0.70
5	<b>0.20</b>	<b>0.50</b>	0.30	0.70	-

	12	3	4	5
12	-	0.80	0.75	<b>0.50</b>
3	0.80	-	0.40	0.30
4	0.75	0.40	-	0.70
5	<b>0.50</b>	0.30	0.70	-

# Agglomerative Clustering Example

## Clustering

COMP90049  
Knowledge  
Technologies

Clustering  
Definition  
Types  
Evaluation

Methods  
Similarity  
k-Means  
k-Ms limitation  
Hierarchical

Summary  
Resources

Merge points 12 & 3 into new cluster: 123

**Using Single Link**

	12	3	4	5
12	-	<b>0.80</b>	0.75	0.50
3	<b>0.80</b>	-	0.40	0.30
4	0.75	0.40	-	0.70
5	0.50	0.30	0.70	-

# Agglomerative Clustering Example

## Clustering

COMP90049  
Knowledge  
Technologies

Clustering  
Definition  
Types  
Evaluation

Methods  
Similarity  
k-Means  
k-Ms limitation  
Hierarchical

Summary  
Resources

Merge points 12 & 3 into new cluster: 123  
**Using Single Link**

	12	3	4	5
12	-	0.80	0.75	0.50
3	0.80	-	0.40	0.30
4	0.75	0.40	-	0.70
5	0.50	0.30	0.70	-

	123	4	5
123	-	<b>0.75</b>	<b>0.50</b>
4	<b>0.75</b>	-	0.70
5	<b>0.50</b>	0.70	-

# Agglomerative Clustering Example

## Clustering

COMP90049  
Knowledge  
Technologies

Clustering  
Definition  
Types  
Evaluation

Methods  
Similarity  
k-Means  
k-Ms limitation  
Hierarchical

Summary  
Resources

Merge points 123 & 4 into new cluster: 12345  
**Using Single Link**

	1234	5
1234	-	0.70
5	0.70	-

# Agglomerative Clustering Example

## Clustering

COMP90049  
Knowledge  
Technologies

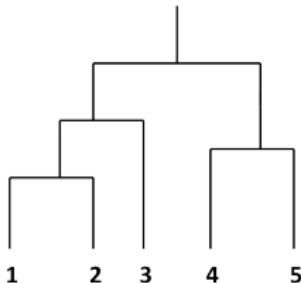
[Clustering](#)  
[Definition](#)  
[Types](#)  
[Evaluation](#)

[Methods](#)  
[Similarity](#)  
[k-Means](#)  
[k-Ms limitation](#)  
[Hierarchical](#)

[Summary](#)  
[Resources](#)

Merge points 123 & 4 into new cluster: 12345  
**Using Single Link**

	1234	5
1234	-	0.70
5	0.70	-



# Agglomerative Clustering Example

## Clustering

COMP90049  
Knowledge  
Technologies

Clustering  
Definition  
Types  
Evaluation

Methods  
Similarity  
k-Means  
k-Ms limitation  
Hierarchical

Summary  
Resources

Merge points 1 & 2 into a new cluster: 12

**Using Complete Link:** update the approximate matrix based on the most dissimilar members

	1	2	3	4	5
1	-	0.90	0.10	0.75	0.20
2	0.90	-	0.80	0.60	0.50
3	0.10	0.80	-	0.40	0.30
4	0.75	0.60	0.40	-	0.70
5	0.20	0.50	0.30	0.70	-

	12	3	4	5
12				
3				
4				
5				

# Agglomerative Clustering Example

## Clustering

COMP90049  
Knowledge  
Technologies

Clustering  
Definition  
Types  
Evaluation

Methods  
Similarity  
k-Means  
k-Ms limitation  
Hierarchical

Summary  
Resources

Merge points 1 & 2 into a new cluster: 12

**Using Complete Link:** update the approximate matrix based on the most dissimilar members

	1	2	3	4	5
1	-	0.90	<b>0.10</b>	0.75	0.20
2	0.90	-	<b>0.80</b>	0.60	0.50
3	<b>0.10</b>	<b>0.80</b>	-	0.40	0.30
4	0.75	0.60	0.40	-	0.70
5	0.20	0.50	0.30	0.70	-

	12	3	4	5
12	-	<b>0.10</b>		
3	<b>0.10</b>	-		
4			-	
5				-

# Agglomerative Clustering Example

## Clustering

COMP90049  
Knowledge  
Technologies

Clustering  
Definition  
Types  
Evaluation

Methods  
Similarity  
k-Means  
k-Ms limitation  
Hierarchical

Summary  
Resources

Merge points 1 & 2 into a new cluster: 12

**Using Complete Link:** update the approximate matrix based on the most dissimilar members

	1	2	3	4	5
1	-	0.90	0.10	<b>0.75</b>	0.20
2	0.90	-	0.80	<b>0.60</b>	0.50
3	0.10	0.80	-	0.40	0.30
4	<b>0.75</b>	<b>0.60</b>	0.40	-	0.70
5	0.20	0.50	0.30	0.70	-

	12	3	4	5
12	-	0.10	<b>0.60</b>	
3	0.10	-	0.40	
4	<b>0.60</b>	0.40	-	
5				-



# Agglomerative Clustering Example

## Clustering

COMP90049  
Knowledge  
Technologies

Clustering  
Definition  
Types  
Evaluation

Methods  
Similarity  
k-Means  
k-Ms limitation  
Hierarchical

Summary  
Resources

Merge points 1 & 2 into a new cluster: 12

**Using Complete Link:** update the approximate matrix based on the most dissimilar members

	1	2	3	4	5
1	-	0.90	0.10	0.75	<b>0.20</b>
2	0.90	-	0.80	0.60	<b>0.50</b>
3	0.10	0.80	-	0.40	0.30
4	0.75	0.60	0.40	-	0.70
5	<b>0.20</b>	<b>0.50</b>	0.30	0.70	-

	12	3	4	5
12	-	0.10	0.60	<b>0.20</b>
3	0.10	-	0.40	0.30
4	0.60	0.40	-	0.70
5	<b>0.20</b>	0.30	0.70	-

# Agglomerative Clustering Example

## Clustering

COMP90049  
Knowledge  
Technologies

Clustering  
Definition  
Types  
Evaluation

Methods  
Similarity  
k-Means  
k-Ms limitation  
Hierarchical

Summary  
Resources

Merge points 4 & 5 into new cluster: 45  
**Using Complete Link**

	12	3	4	5
12	-	0.10	0.60	0.20
3	0.10	-	0.40	0.30
4	0.60	0.40	-	<b>0.70</b>
5	0.20	0.30	<b>0.70</b>	-

# Agglomerative Clustering Example

## Clustering

COMP90049  
Knowledge  
Technologies

[Clustering](#)  
[Definition](#)  
[Types](#)  
[Evaluation](#)

Methods  
[Similarity](#)  
[k-Means](#)  
[k-Ms limitation](#)  
[Hierarchical](#)

[Summary](#)  
[Resources](#)

Merge points 4 & 5 into new cluster: 45  
**Using Complete Link**

	12	3	4	5
12	-	0.10	0.60	0.20
3	0.10	-	0.40	0.30
4	0.60	0.40	-	0.70
5	0.20	0.30	0.70	-

	12	3	45
12			
3			
45			

# Agglomerative Clustering Example

## Clustering

COMP90049  
Knowledge  
Technologies

[Clustering](#)  
[Definition](#)  
[Types](#)  
[Evaluation](#)

Methods  
[Similarity](#)  
[k-Means](#)  
[k-Ms limitation](#)  
[Hierarchical](#)

[Summary](#)  
[Resources](#)

Merge points 4 & 5 into new cluster: 45  
**Using Complete Link**

	12	3	4	5
12	-	0.10	0.60	0.20
3	0.10	-	<b>0.40</b>	<b>0.30</b>
4	0.60	<b>0.40</b>	-	0.70
5	0.20	<b>0.30</b>	0.70	-

	12	3	45
12	-	0.10	0.20
3	0.10	-	<b>0.30</b>
45	0.20	<b>0.30</b>	-

# Agglomerative Clustering Example

## Clustering

COMP90049  
Knowledge  
Technologies

## Clustering

### Definition

### Types

### Evaluation

## Methods

### Similarity

### k-Means

### k-Ms limitation

### Hierarchical

## Summary

## Resources

Merge points 45 & 3 into new cluster: 345

**Using Complete Link**

	12	345
12	-	0.10
345	0.10	-

# Agglomerative Clustering Example

## Clustering

COMP90049  
Knowledge  
Technologies

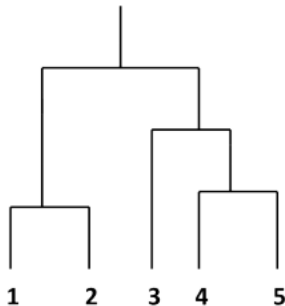
Clustering  
Definition  
Types  
Evaluation

Methods  
Similarity  
k-Means  
k-Ms limitation  
Hierarchical

Summary  
Resources

Merge points 45 & 3 into new cluster: 345  
**Using Complete Link**

	12	345
12	-	0.10
345	0.10	-



# Agglomerative Clustering Example

## Clustering

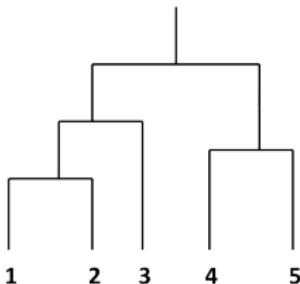
COMP90049  
Knowledge  
Technologies

Clustering  
Definition  
Types  
Evaluation

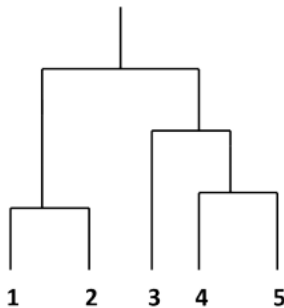
Methods  
Similarity  
k-Means  
k-Ms limitation  
Hierarchical

Summary  
Resources

### Single Link



### Complete Link



# Clustering vs Classification

## Clustering

COMP90049  
Knowledge  
Technologies

Clustering  
Definition  
Types  
Evaluation

Methods  
Similarity  
k-Means  
k-Ms limitation  
Hierarchical

Summary  
Resources

Outlook	Temperature	Humidity	Windy	Play
sunny	hot	high	FALSE	no
sunny	hot	high	TRUE	no
overcast	hot	high	FALSE	yes
rainy	mild	high	FALSE	yes
rainy	cool	normal	FALSE	yes
rainy	cool	normal	TRUE	no
⋮	⋮	⋮	⋮	⋮



# Clustering vs Classification

## Clustering

COMP90049  
Knowledge  
Technologies

Clustering  
Definition  
Types  
Evaluation

Methods  
Similarity  
k-Means  
k-Ms limitation  
Hierarchical

Summary  
Resources

Outlook	Temperature	Humidity	Windy	Play
sunny	hot	high	FALSE	no
sunny	hot	high	TRUE	no
overcast	hot	high	FALSE	yes
rainy	mild	high	FALSE	yes
rainy	cool	normal	FALSE	yes
rainy	cool	normal	TRUE	no
⋮	⋮	⋮	⋮	⋮

## Clustering

COMP90049  
Knowledge  
Technologies

Clustering  
Definition  
Types  
Evaluation

Methods  
Similarity  
k-Means  
k-Ms limitation  
Hierarchical

Summary  
Resources

- What basic contrasts are there in different clustering methods?
- How does  $k$ -means operate, and what are its strengths and weaknesses?
- How to evaluate clusters
- What are some challenges we face when clustering data?

## Clustering

COMP90049  
Knowledge  
Technologies

Clustering  
Definition  
Types  
Evaluation

Methods  
Similarity  
k-Means  
k-Ms limitation  
Hierarchical

Summary  
Resources

Tan, Steinbach, Kumar (2006) Introduction to Data Mining.  
Chapter 8, Cluster Analysis

<http://www-users.cs.umn.edu/~kumar/dmbook/ch8.pdf>

Jain, Dubes (1988) Algorithms for Clustering Data.

[http://homepages.inf.ed.ac.uk/rbf/BOOKS/JAIN/Clustering\\_Jain\\_Dubbes.pdf](http://homepages.inf.ed.ac.uk/rbf/BOOKS/JAIN/Clustering_Jain_Dubbes.pdf)

# MID-SEMESTER EXAM REVIEW

## Clustering

COMP90049  
Knowledge  
Technologies

Clustering  
Definition  
Types  
Evaluation

Methods  
Similarity  
k-Means  
k-Ms limitation  
Hierarchical

Summary  
Resources

	Q1	Q2	Q3	Q4	Q5	Q6
<b>Average</b>	1.41	1.70	0.68	1.28	1.25	0.68
<b>Median</b>	1.5	2	0.75	1.5	1.5	1
<b>SD</b>	0.52	0.55	0.30	0.85	0.69	0.41