



Clustering

1-1



RNA-seq: What is it good for?

High-throughput RNA sequencing experiments (RNA-seq) offer the ability to measure simultaneously the expression level of thousands of genes in a single experiment!

1-2



Data... And Lots of It!

	Experiment 1	Experiment 2	Experiment 3	Experiment 4	...	Experiment $m-1$	Experiment m
Gene 1	0.6	4.4	1.3	1.0	...	3.1	2.2
Gene 2	1.5	2.6	5.2	0.8	...	2.8	2.9
Gene 3	0.7	3.7	2.4	1.9	...	1.5	1.6
Gene 4	0.3	0.7	0.2	1.3	...	4.9	3.0
Gene 5	3.1	3.0	2.1	1.4	...	4.2	0.9
...
Gene $n-1$	1.8	2.5	1.8	0.7	...	2.7	3.1
Gene n	0.5	3.4	3.0	0.5	...	1.8	2.5

1-3



Finding Similarly Expressed Genes

- It may be useful to partition the n genes into groups of similarly expressed genes
- Clustering is the art of finding groups of genes, such that genes in the same group are as similar to each other as possible and as dissimilar to genes in other groups as possible

1-4



Clustering Algorithms

- Hierarchical clustering
- *CAST*
- *k*-means clustering
- Model-based clustering

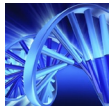
1-5



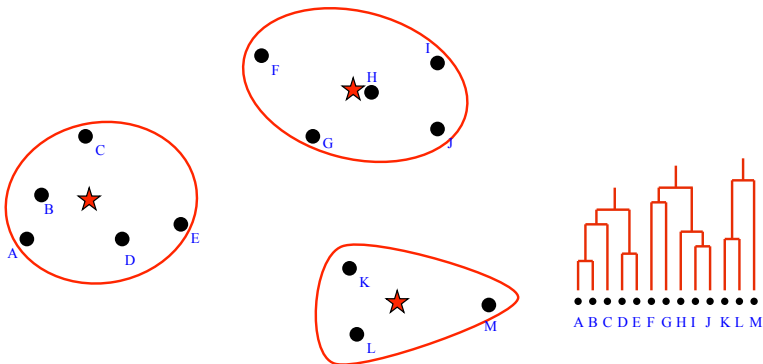
Hierarchical Clustering Algorithm

- Assign each point to its own cluster
- Repeat the following step until the desired number of clusters is reached
 - Merge together the two closest clusters

1-6



Hierarchical Clustering



1-7



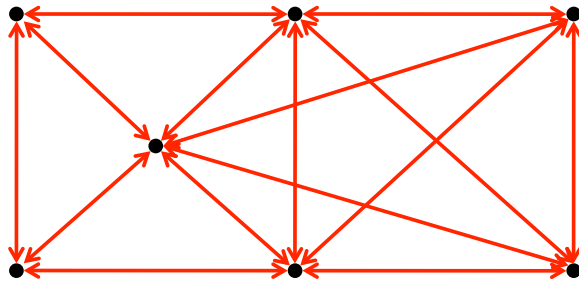
Clustering Algorithms

- Hierarchical clustering
- *CAST*
- *k*-means clustering
- Model-based clustering

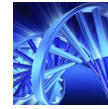
1-8



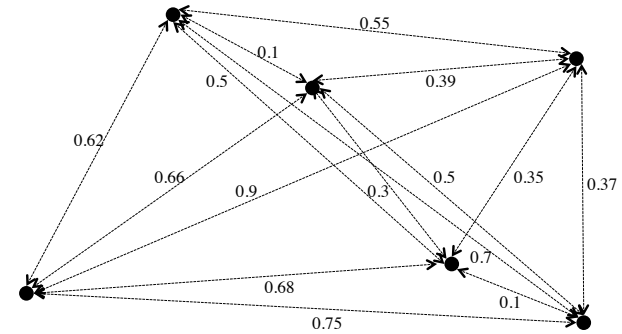
Cliques in Graphs



1-9



Cluster Affinity Search Technique (CAST)

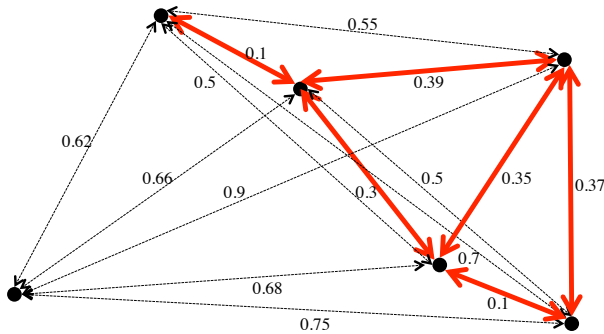


threshold = 0.0

1-10



Cluster Affinity Search Technique (CAST)

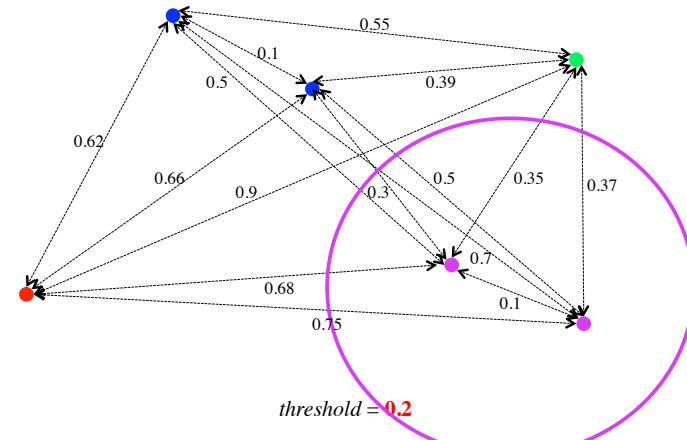


threshold = 0.4

1-11



Cluster Affinity Search Technique (CAST)



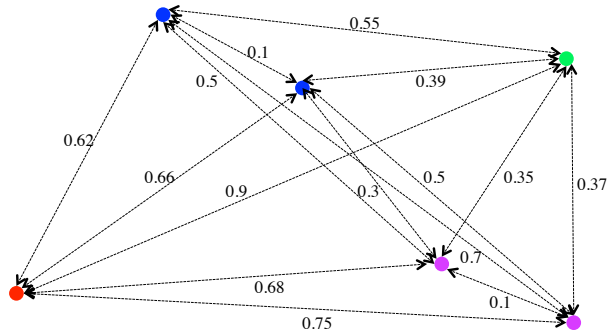
threshold = 0.2

1-12



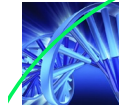
Cluster Affinity Search Technique (CAST)

4 Clusters

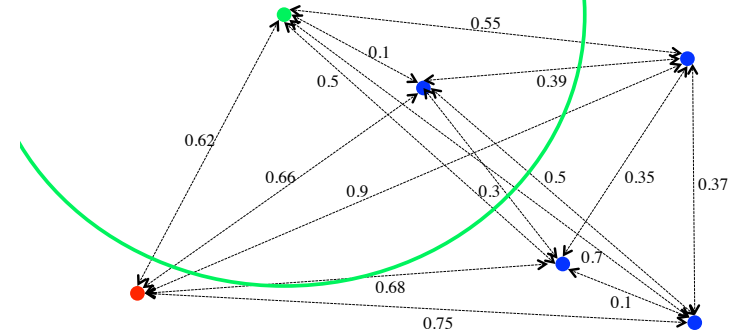


threshold = 0.2

1-13



Cluster Affinity Search Technique (CAST)



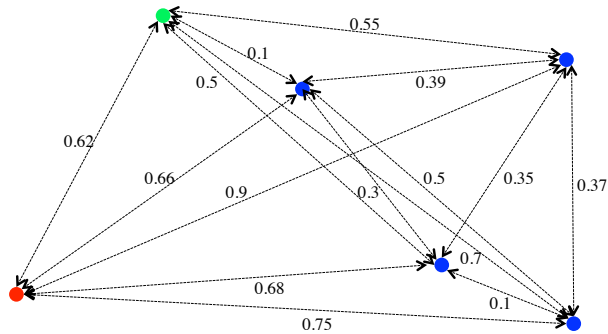
threshold = 0.45

1-14



Cluster Affinity Search Technique (CAST)

3 Clusters

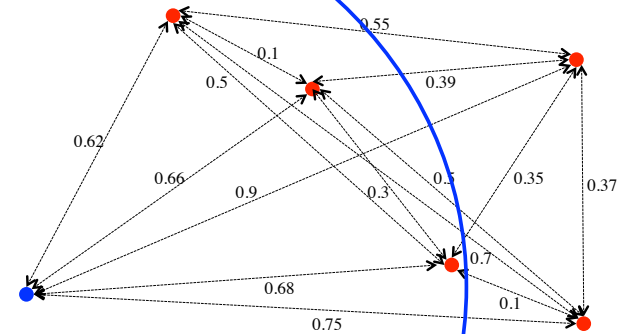


threshold = 0.45

1-15



Cluster Affinity Search Technique (CAST)



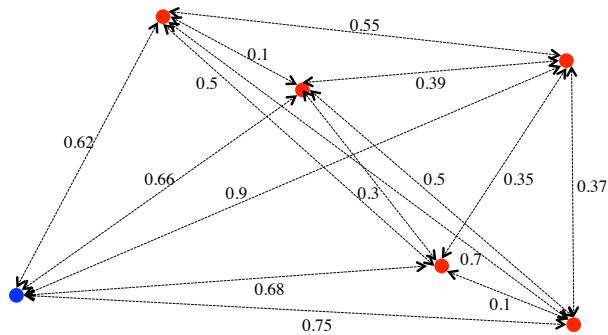
threshold = 0.7

1-16



Cluster Affinity Search Technique (CAST)

2 Clusters



threshold = 0.7

1-17



CAST Clustering Algorithm

- Repeat until all points (genes) are assigned to a cluster
 - Choose a point (gene) not already assigned to a cluster and assign it to a new cluster \mathcal{C}
 - Repeat until the cluster \mathcal{C} converges
 - Add to \mathcal{C} any unassigned points (genes) that are closer to the points (genes) in \mathcal{C} , on average, than some *threshold*
 - Remove from \mathcal{C} any points (genes) that are farther from the other points (genes) in \mathcal{C} , on average, than some *threshold*

1-18



Clustering Algorithms

- Hierarchical clustering
- CAST
- k -means clustering
- Model-based clustering

1-19



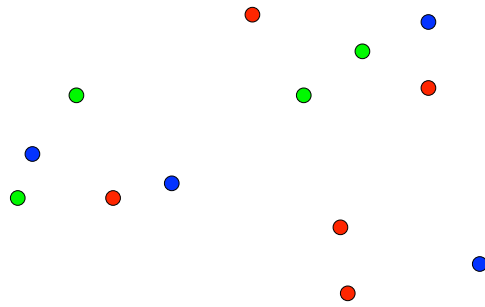
k -means Clustering Algorithm

- Randomly assign each point (gene) to one of k clusters
- Repeat until convergence
 - Calculate *mean* of each of the k clusters
 - Assign each point (gene) to the cluster with the closest *mean*

1-20



k-means Clustering Example

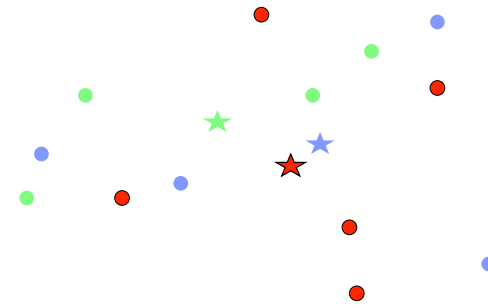


Randomly assign each point to one of k clusters

1-21



k-means Clustering Example

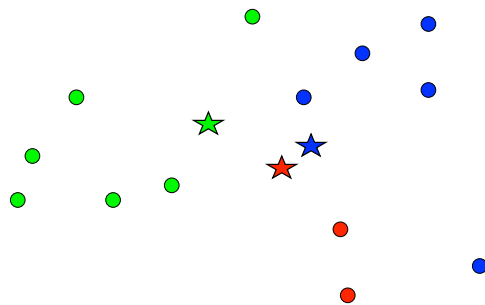


Calculate mean of each cluster

1-22

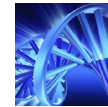


k-means Clustering Example

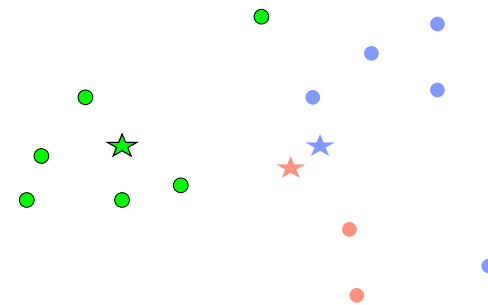


Assign each point to closest cluster mean

1-23



k-means Clustering Example

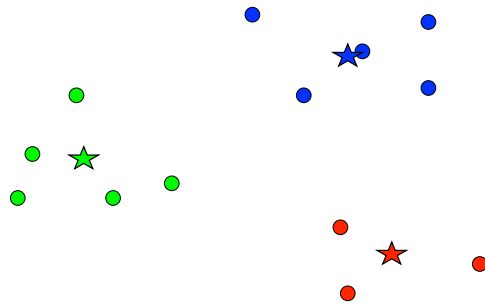


Calculate mean of each cluster

1-24



k-means Clustering Example



Convergence

1 - 25



How Good Is Our Clustering?

For a given number of clusters, k , one measure of a clustering's quality is the sum of the distances between each point and the mean of the point's cluster

1 - 26



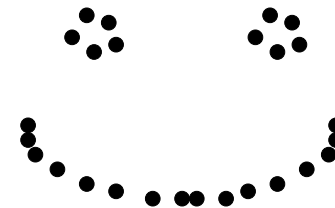
Clustering Problem

- *Clustering Problem*: Partition n data points into k clusters such that the total distance from each point to its cluster center is minimized.
- Clustering is an NP-complete problem

1 - 27



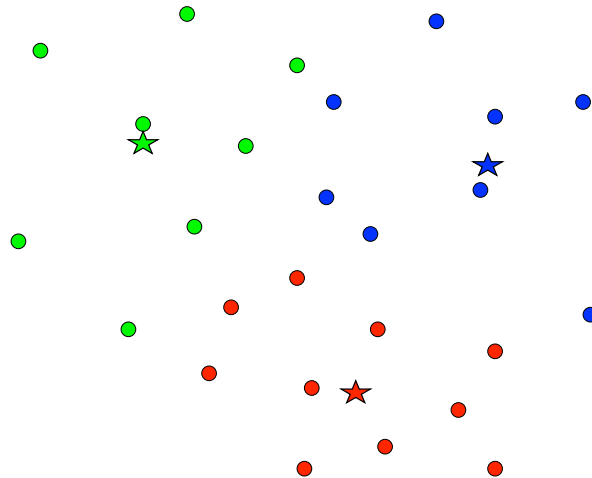
Does k -means Always Work?



1 - 28



Does k -means Always Work?



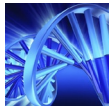
1 - 29



Clustering Algorithms

- Hierarchical clustering
- CAST
- k -means clustering
- Model-based clustering

1 - 30



Model-Based Clustering

- Randomly assign each point (*gene*) to one of k clusters
- Repeat until convergence
 - Calculate *model* of each of the k clusters
 - Assign each point (*gene*) to the cluster with the closest *model*

1 - 31



Model-Based Clustering

Randomly assign each point to one of k clusters (models)

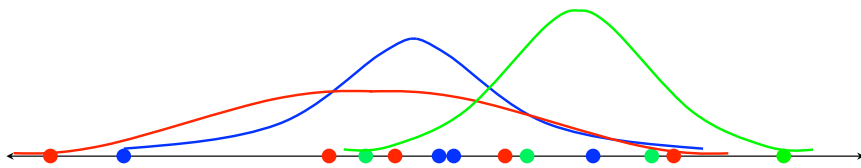


1 - 32



Model-Based Clustering

Calculate model for each of the k clusters

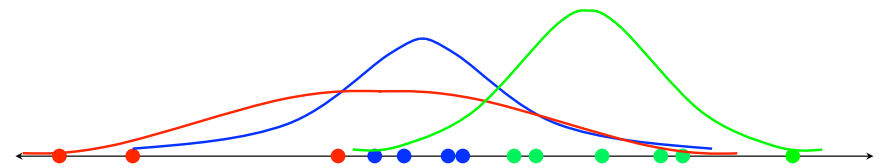


1-33



Model-Based Clustering

Assign each point to the most likely model

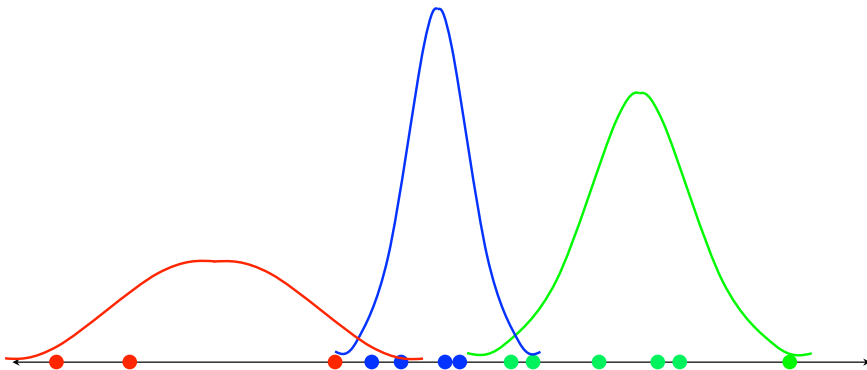


1-34



Model-Based Clustering

Calculate model for each of the k clusters

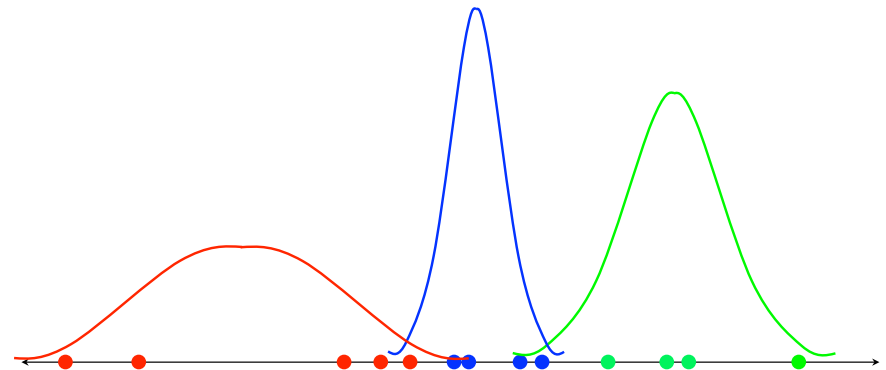


1-35



Model-Based Clustering

Assign each point to the most likely model

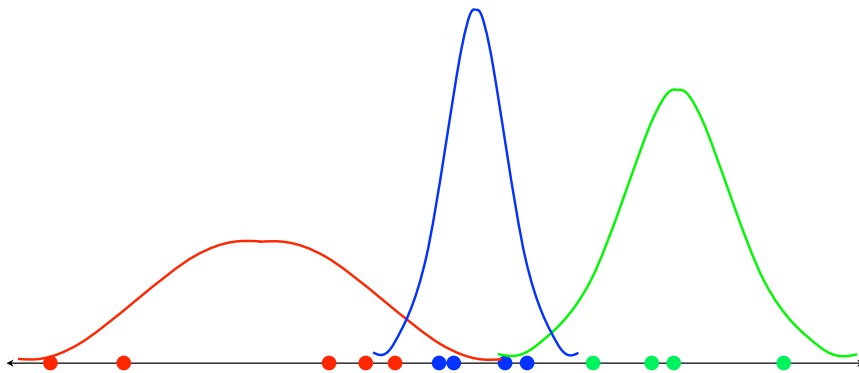


1-36



Model-Based Clustering

Calculate model for each of the k clusters



1-37



Clustering Genes vs. Clustering Experiments

	Experiment 1	Experiment 2	Experiment 3	Experiment 4	...	Experiment $m-1$	Experiment m
Gene 1	0.6	4.4	1.3	1.0	...	3.1	2.2
Gene 2	1.5	2.6	5.2	0.8	...	2.8	2.9
Gene 3	0.7	3.7	2.4	1.9	...	1.5	1.6
Gene 4	0.3	0.7	0.2	1.3	...	4.9	3.0
Gene 5	3.1	3.0	2.1	1.4	...	4.2	0.9
...
Gene $n-1$	1.8	2.5	1.8	0.7	...	2.7	3.1
Gene n	0.5	3.4	3.0	0.5	...	1.8	2.5

1-38



Assessing Clustering

• Figure of Merit (FOM)

Apply clustering algorithm to all but one experimental condition and use the left-out condition to assess the predictive power of the clustering algorithm

1-39



Assessing Clustering

• Rand Index

Measures similarity of two clusterings, X and Y

$$R = \frac{a + b}{a + b + c + d}$$

a = number of pairs of points in the same cluster in X and in the same cluster in Y

b = number of pairs of points in different clusters in X and in different clusters in Y

c = number of pairs of points in the same cluster in X and in different clusters in Y

d = number of pairs of points in different clusters in X and in the same cluster in Y

1-40