

Gene Expression Analysis.

DNA arrays \Rightarrow expression levels or amount of mRNA produced in the cell of n genes, under many time points & conditions.

\Downarrow
genes switched on and switched off in the cell.

$$[I]_{n \times m} = \text{expression matrix.} \\ = [I_{ij}]$$

$I_{ij} \Rightarrow$ intensity of hybridization signal as provided by DNA array (for i^{th} gene at j^{th} time instance & conditions).

\hookrightarrow Intensities are transformed & normalized. (Not discussed here).

Expression pattern of gene $i \Leftrightarrow i^{\text{th}}$ row of I .

$n \times m$ expression matrix \rightarrow $n \times n$ distance matrix $d = (d_{ij})$.

Clustering algorithms group genes with similar expression patterns into clusters with the hopes that these clusters correspond to groups of functionally related genes.

Clustering principle:

- Homogeneity: d_{ij} small if i & j belong to same group.
- Separation: d_{ij} large if they belong to diff. cluster.

* Hierarchical Clustering (d, n)

1. Form n clusters with single element.
Let $S_c =$ set of initial clusters. $V =$ {vertex set corresponds to clusters}.
2. while (there is more than 1 cluster)
 - 2.1 Find the two closest clusters c_1 & c_2 .
 - 2.2 Merge c_1 & c_2 to C with $|c_1| + |c_2|$ elements.
 - 2.3 Compute distances of C with other clusters.
 - 2.4 Add a vertex, whose children are c_1 & c_2 .
 - 2.5 Remove c_1 & c_2 from the list & add C to the list.

Dist. distance measures:

$$d_{\min}(c^*, c) = \min_{x \in c^*, y \in c} d(x, y)$$

$$d_{\text{av}}(c^*, c) = \text{avg.} \{ d(x, y) \}_{x \in c^*, y \in c}$$

$$? d_H(c^*, c) = \frac{d_H(c^*, c_1) + d_H(c^*, c_2) - d_H(c_1, c_2)}{2}$$

* K-means clustering.

Let $\chi = \{x_1, x_2, \dots, x_k\}$ be k -centers.

Given a data point v , let the squared error distribution be:

$$d(v, \chi) = \min_{1 \leq i \leq k} d(v, x_i)$$

Let $\mathcal{V} = \{v_1, v_2, \dots, v_n\}$

$$\therefore d(\mathcal{V}, \chi) = \frac{\sum_{i=1}^n (d(v_i, \chi))^2}{n}$$

= arg. squared error distribution.

k-means clustering problem:

Given n data points, find k center points minimizing the ^{avg.} squared error distortion.

→ no efficient (polynomial) alg. known yet.
Lloyd k-means clustering algorithm (k)

1. choose arbitrary k points as cluster representatives. $\{x_1, x_2, \dots, x_k\}$
2. while (partitions converges or fluctuations are small)
 - 2.1. Assign each data point to the cluster c_i corresponding to the closest cluster representative x_i ($1 \leq i \leq k$)

- 2.2. $R_i \neq \emptyset, x_i = \frac{\sum_{u_j \in c_i} u_j}{|c_i|}$

Variation over distance measure:

$$d(V, X) = \sum_{i=1}^n d(u_i, X) \Rightarrow \text{K-Median Prob.}$$

$$d(V, X) = \max_{1 \leq i \leq n} d(u_i, X) \Rightarrow \text{K-Center Prob.}$$

While Lloyd alg. is fast, it can significantly rearrange every cluster in every iteration. A conservative approach is to move only one element between clusters in every iteration (the most suitable replacement is preferred, which decreases manually the cost fn. or distance).

* Greedy Approach for K-Means Clustering:

Let P denotes a partitioning into k partitions.

$\text{Cost}(P)$ = associated distortion measure.

$P_{i \rightarrow c}$ = Partition obtained from P by moving v_i to c .

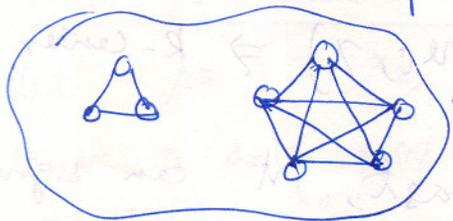
\therefore improvement $\Rightarrow \Delta(i \rightarrow c) = \text{Cost}(P) - \text{Cost}(P_{i \rightarrow c}) > 0$

Greedy method takes that move which ^{maximises this} improvement.
 or. Stop ^{iteration} iteratively after the ^{maximal} improvement becomes zero or ~~small~~.

* Clustering through clique finding / transformation.

Complete graph: K_n (Undirected) with n nodes with every two vertices connected by an edge.

Clique graph: Every connected component is a complete graph.



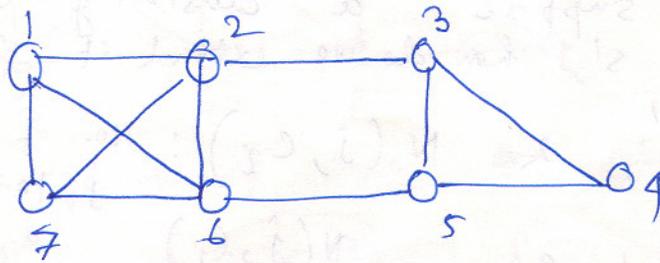
← Clique graph.

⇒ Each ~~clique~~ component is a strong cluster.

Complete Subgraph: If $V' \subset V$ in a graph $G(N, E)$ induces a subgraph which is complete, the subgraph is complete.

Clique: A clique in a subgraph graph is a maximal complete subgraph.

U.E. a complete subgraph which is not contained inside any other ~~complete~~ complete subgraph.



$\{1, 7, 6\} \Rightarrow$ Not a clique but a complete subgraph.

Cliques are: $\{1, 2, 6, 7\}$; $\{3, 4, 5\}$.

Expression matrix \rightarrow Distance Matrix \rightarrow Distance graph
 $\{E_{i,j}\}$ $\{d_{i,j}\}$ $\{E_{i,j} \mid d_{i,j} < \theta\}$
 \uparrow threshold

Errors in expression data and the absence of a "universal good threshold" θ often results in distance graphs that do not quite form clique graphs.

Corrupted clique problem:

Determine the smallest number of edges that need to be added or removed to transform a graph into a clique graph.

\rightarrow NP-hard. \rightarrow Need for approx. soln.

Two ~~approx~~ approximate algorithms:

- Parallel Classification with cores (PCC)
- Cluster Affinity Search Technique (CAST)

• Parallel Classifications with Cores (PCC) Alg. 1 (Amir Ben-Dor et al. (1999))

Let $S' \subset S$, suppose a clustering (C_1, C_2, \dots, C_k) is available on S' , how do we extend it for S ?

Let $j \in S - S'$, & $N(j, C_i)$: No. of edges from j to C_i .

$$\therefore \text{affinity } j \text{ to } C_i = \frac{N(j, C_i)}{|C_i|}$$

Assign j to that cluster which has maximum affinity.

~~PCC (S, G, k)~~ : ~~S is a distance graph,~~
 ~~k no. of clusters.~~

* PCC (S, G, k)

S : Set of n elements (say genes, forming vertices of G)

G : Dist. graph

k : no. of clusters.

1. Randomly select $S' \subset S$ & $S'' \subset S - S'$

s.t. $|S'| \approx \log \log n$ & $|S''| \approx \log n$

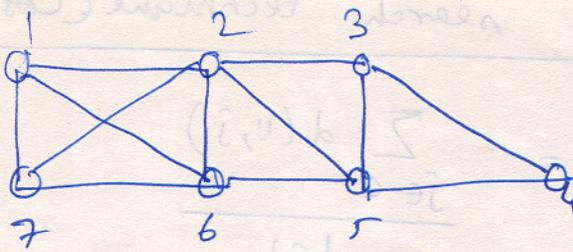
2. For all k partitions in S' ; obtain extended partitions in S , through two stages of extensions ($S' \rightarrow S'' \rightarrow S \hat{=} (S' \cup S'')$).

Choose the one which has ~~no~~ ^{minimum} ~~minimum~~ ~~cluster~~ ~~score~~ score (i.e. the

no. of edges required to add or remove ~~set~~ from G , to get a clique graph as per the partition).

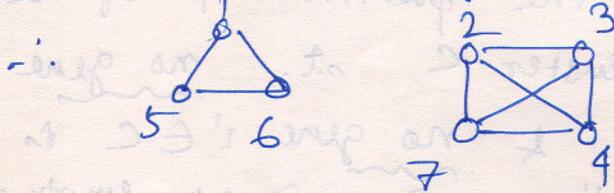
Ex 2

Let $G =$



Let a partition P into 2 ~~clique~~ partitions be
 as follows: $\{1, 5, 6\}, \{2, 3, 7, 4\}$.

Clique graph wrt P :



No. of edges to be added or removed for G :

For addition: $\{(1,5), (4,7), (2,4), (3,7)\}$

For deletion: $\{(1,2), (1,7), (7,6), (2,5), (3,5), (4,5)\}$

$= 10$

$\therefore \text{score}(P) = 10$

Time Complexity:

no. of partitions in $S' = k^{|S'|} = k^{\log \log n} = (\log n)^{\log k} = (\log n)^c$

In each iteration $O(n^2)$ operations for extension and score computation.

\therefore Total Time Complexity: $O(n^2 (\log n)^c)$

where $c = \log_2 k$

• cluster affinity search technique (CAST)

$$\text{Let, } d(u, C) = \frac{\sum_{j \in C} d(u, j)}{|C|}$$

Given a threshold θ , u is close to C , if $d(u, C) < \theta$
 else, u is distant.

CAST iteratively builds the partition P of the set S by finding a cluster C s.t. no gene $i \notin S$ is close to C & no gene $i \in C$ is distant to C . Initially P is an empty set.

CAST (S, G, θ)

1. $P = \emptyset$

2. while $(S \neq \emptyset)$

2.1. $v =$ vertex of maximal degree in G .

2.2. $C \leftarrow \{v\}$;

2.3. while $(\exists$ a close gene $i \notin C$ or \exists distant gene $j \in C)$

{ 2.3.1. Find the closest gene $i \notin C$ & add to C .

2.3.2. Find the farthest ^{distant} gene $i \in C$ & remove it.

}

2.4. Add C to P

2.5. $S = S - C$;

& update G ^{distance graph} accordingly.

}

There is no performance guarantee & it may not converge at all. But, in general it gives good performance.