

Lecture 18:
Microarray Data Analysis: clustering

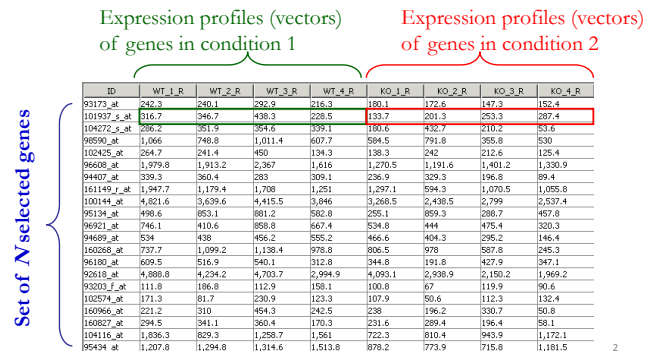
Lubica Benuskova

<http://www.cs.otago.ac.nz/cosc348/>

1

Matrix of selected genes

- We can treat rows/columns of numbers as n -dimensional vectors.



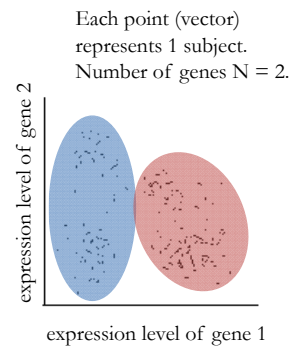
Gene expression profiles

- Each gene has its own **expression profile**, i.e. an array of values for each condition.
- The idea behind clustering is that genes with similar profiles can be grouped together.
- This can help to identify the function of “mystery” genes:
 - if a gene with an unknown function has a similar expression profile to a collection genes, then there is a good chance that the “mystery” gene also plays a role in the studied process.
- It is assumed that the genes with similar expression profiles are co-expressed, e.g. their expression may co-regulated by the same transcription factors.

3

Cluster analysis: definition

- Clustering is an exploratory procedure which provides a method for grouping objects based on some measure of similarity.
- The goal is to find groups of points that are close to each other and represent a class of objects, e.g. **cluster of all subjects which are healthy** versus a **cluster of subjects who have genetic disease**.



Cluster analysis: assumptions

- No assumptions are made about the number of groups/clusters, or the group/cluster structure.
- Grouping is done on the basis of similarities, or dissimilarities, i.e. distances between vectors (points in n -dim space).
- A quantitative scale (metric) is used to measure the closeness or similarity between objects.
- In the microarray setting, the objects (vectors) are the expression profiles of the genes in the experiment.

Cluster analysis: distance measure

- One popular measure of the closeness between two vectors is the **Euclidean distance**:

$$d(x, y) = \sqrt{(x_1 - y_1)^2 + \dots + (x_n - y_n)^2}$$

- Another popular measure is the **Pearson's correlation coefficient** (σ is the standard deviation):

$$r = \frac{1}{n-1} \sum_{i=1}^n \left(\frac{x_i - \bar{x}}{\sigma_x} \right) \left(\frac{y_i - \bar{y}}{\sigma_y} \right)$$

- There are other distance (correlation) measures. The choice depends on the measurement scale, heuristic knowledge, or we can experiment with different measures to see, which one yields better results. Note that the metrics are symmetric.

K-means clustering

- K-means clustering is a method of cluster analysis, which aims to partition N vectors into K clusters, in which each vector belongs to the cluster with the nearest mean (i.e. center of a cluster).
- It attempts to find the centers of natural clusters in the data by the iterative refinement.
- It is the most popular “bottom-up” clustering algorithm.
- If K is the number of clusters, n is the dimensionality of vectors, and N the number of vectors, the problem can be exactly solved in time $O(N^{nK+1} \log N)$.

7

K-means clustering

- Vector clustering for K centers:
 - K-means clustering partitions the vectors into K clusters $\mathbf{C}=\{C_1, C_2, \dots, C_K\}$ so as to minimize the sum of distances between all vectors belonging to the cluster and the centre of the cluster

$$\underset{\mathbf{C}}{\operatorname{argmin}} \sum_{k=1}^K \sum_{\mathbf{x}_p \in C_k} d(\mathbf{x}_p, \mathbf{c}_k)$$

- Cluster C_k has the centre \mathbf{c}_k
- Let us have P vectors in total (p is the index of input vector)
- But how do we know which vectors belong to the cluster k ?

8

K-means clustering algorithm

- Randomly choose K input vectors from the data, which will be the starting centres of clusters.
- Then for each input vector \mathbf{x}_p do:
 - Find the centre \mathbf{c}_k , which is closest to \mathbf{x}_p
 - Create a new centre \mathbf{c}_k by changing each coordinate i , so that:
$$c_{ki} \leftarrow c_{ki} + \alpha (x_{pi} - c_{ki})$$
 - Stopping criterion: either we can make α linearly decreasing to zero or we can stop when centres do not move significantly. In general, we choose initial $\alpha \in (0, 1]$.

9

K-means clustering

- K-means clustering is run for different number of clusters K as we do not know ahead how many natural clusters are in the data.
- There is no guarantee it will converge to the global optimum.
- The result depends on the initial random choice of centers of clusters.
- As the algorithm is usually very fast, it is common to run it multiple times with different starting centers.
- Particular value of $\alpha \in (0, 1]$ must be also found by experimentation.

10

K-means++ algorithm

- K-means++ is designed to improve the initial choice of cluster centers based on spreading the K initial cluster centers away from each other:
1. Choose one center uniformly at random from among the data points.
 2. For each data point x_p , compute $d(x_p, c_k)$, the distance between x_p and the nearest center that has already been chosen.
 3. Add one new data point at random as a new center, using a weighted probability distribution, where a new point x_p^* is chosen with probability proportional to the distance $d(x_p, c_k)$.
 4. Repeat Steps 2 and 3 until K centers have been chosen.
 5. Now that the initial centers have been chosen, proceed using standard K-means clustering.

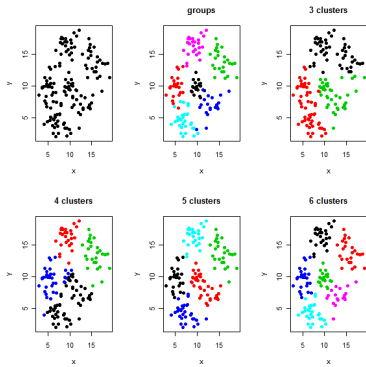
11

K-means++ clustering

- This seeding method gives out considerable improvements in the final error of K-means.
- Although the initial selection in the algorithm takes extra time, the K-means part itself converges very fast after this seeding and thus the algorithm actually lowers the computation time too.
- The authors tested their method with real and synthetic datasets and obtained typically a 2-fold improvements in speed, and for certain datasets close to 1000-fold improvements in error.
- It was proposed in 2007 by David Arthur and Sergei Vassilvitskii as a way of avoiding the sometimes poor clusterings found by the standard K-means algorithm.

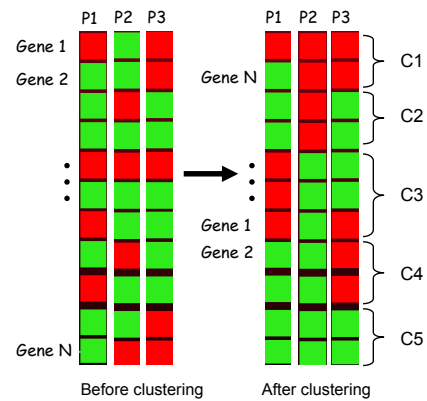
12

How to choose K ?



- How many clusters are there? We do not know.
- We have to apply some heuristics to choose K (e.g. K = number of conditions) or experiment with different K s.

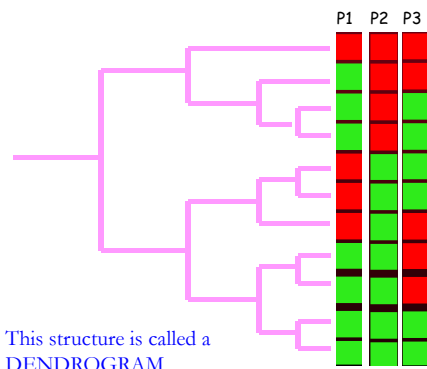
K -means clustering of gene profiles



- Let us have spotted microarrays. Red means high expression, green means low expression.
- We have clustered genes into 5 clusters according to similarity of their profiles across patients P1, P2 and P3.

14

Hierarchical clustering based on K -means



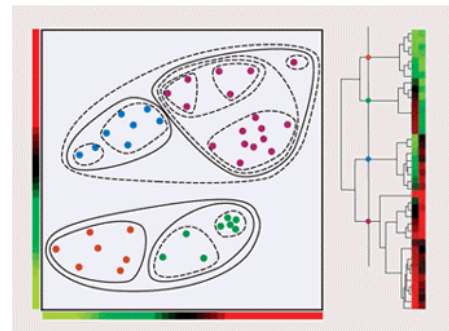
This structure is called a DENDROGRAM

- First we perform the K -means clustering
- Then we hierarchically arrange the clusters according to linkage criteria
- Result: a tree of distances between clusters of vectors

15

Hierarchical clustering

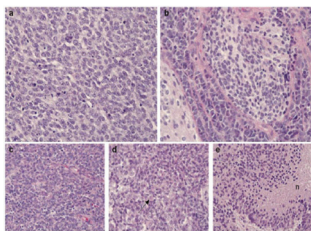
- Another illustration of hierarchical organisation of clusters produced by K -means first ($K = 9$).



16

Brain tumour data

- Brain data, Scott Pomeroy et al, Nature (vol. 415), Jan 2002
 - 99 samples, about 7000 genes, 5 classes of brain tumour
 - Authors have selected 50 genes that are most differentially expressed based on t -test.



Photomicrographs of tumours:

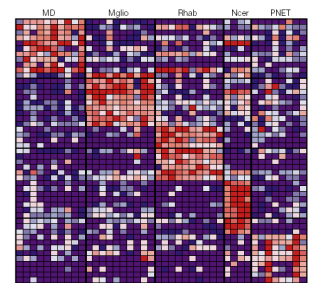
- a, MD (medulloblastoma) classis
- b, MD desmoplastic
- c, PNET
- d, rhabdoid tumour
- e, glioblastoma

Analysis also used Normal cerebella tissue (Ncer), not shown here

17

Real example of K -means clustering

- They performed gene clustering on selected genes for $K = 5$.
- As a result, each condition (disease) may involve different cluster of genes, which are differentially expressed.
- **The goal** of gene clustering: to reveal the clusters of differentially expressed genes, which characterize each condition.



Rows: selected genes (profile vectors)
Columns: patients
Classes: 4 brain tumour types + Ncer

[Pomeroy et al., Nature 415: 436-442, 2002]

18