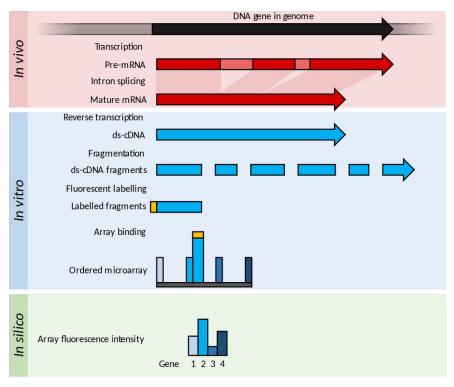
Week 07: Functional genomics

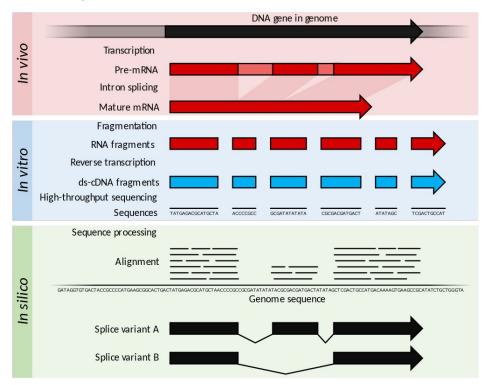
- Measuring gene-expression
- Distance/Similarity measures
- Clustering genes/samples
- Differential expression
- Functional enrichment

Measuring gene-expression on a large-scale

DNA microarrays



RNA-seq

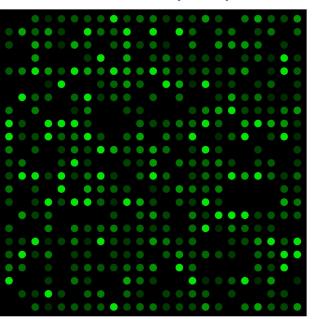


Measuring gene-expression on a large-scale

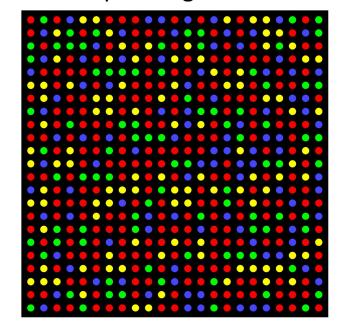
DNA microarrays

RNA-seq

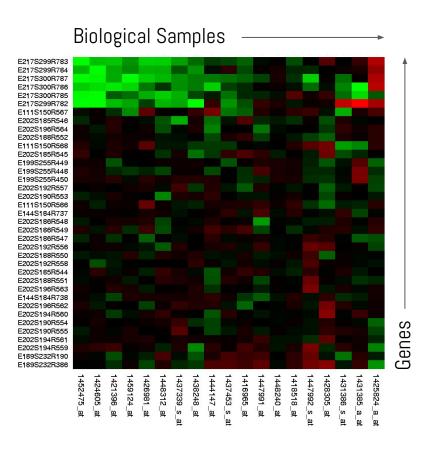
Microarray chip



Sequencing flow cell



Measuring gene-expression on a large-scale



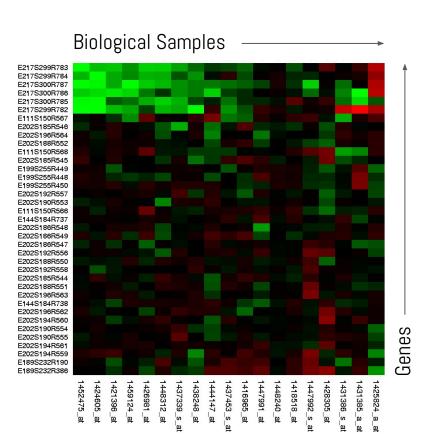
Gene-level Qs:

- What's expressed (& by how much) in a given context/condition?
- 2. What's differentially expressed between two (or more) contexts/conditions?

Group-level Qs:

- 1. Are there groups of genes that respond similarly to changing contexts (across samples)?
- 2. Are there groups of samples that have very similar gene expression profiles?

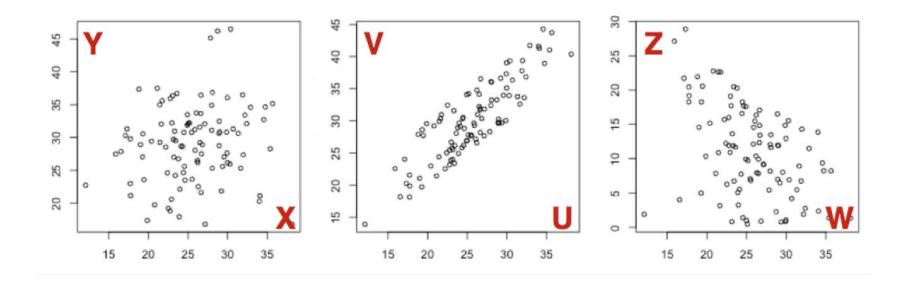
Calculating "distance" between genes or samples



Variab ♦		Attributes / Features									
x	10	8	13	9	11	14	6	4	12	7	5
у	8.04	6.95	7.58	8.81	8.33	9.96	7.24	4.26	10.84	4.82	5.68

Calculating "distance"/"similarity" between genes or samples

Variat		Attributes / Features									
x	I.	I.a.		l				I	12	7	5
y	8.04	6.95	7.58	8.81	8.33	9.96	7.24	4.26	10.84	4.82	5.68



Similarity measures

Pearson Correlation Coefficient

 Measures 'linear' relationship between variables.

$$r = rac{\sum_{i=1}^{n}(x_i - ar{x})(y_i - ar{y})}{\sqrt{\sum_{i=1}^{n}(x_i - ar{x})^2}\sqrt{\sum_{i=1}^{n}(y_i - ar{y})^2}}$$

where:

- n is the sample size
- x_i, y_i are the single samples indexed with i
- $ullet ar x = rac{1}{n} \sum_{i=1}^n x_i$ (the sample mean); and analogously for ar y

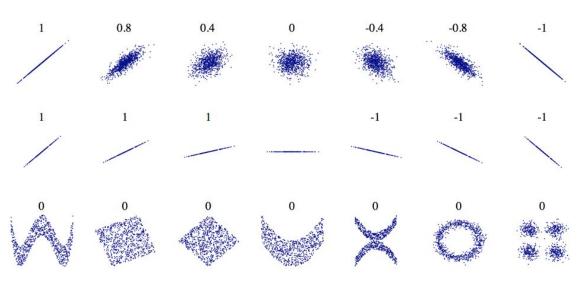
$$m{r} = rac{1}{n-1} \sum_{i=1}^n \left(rac{x_i - ar{x}}{s_x}
ight) \left(rac{y_i - ar{y}}{s_y}
ight).$$

Distance measures

Pearson Correlation Coefficient

 Measures 'linear' relationship between variables.

$$m{r} = rac{1}{n-1} \sum_{i=1}^n \left(rac{x_i - ar{x}}{s_x}
ight) \left(rac{y_i - ar{y}}{s_y}
ight)$$



$$-1 \le r \le +1$$

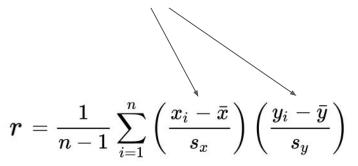
-1 is total -ve correlation | 0 is no correlation | +1 is total +ve correlation

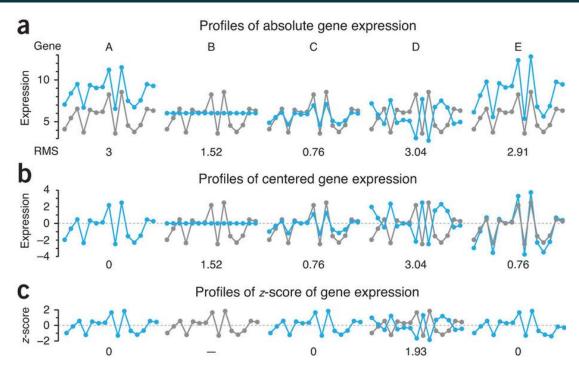
Distance measures

Pearson Correlation Coefficient

Captures the relationship between
 2 vectors after centering each
 vector by its mean and scaling by
 its standard deviation.

 The final quantities for each vector are called z-scores.



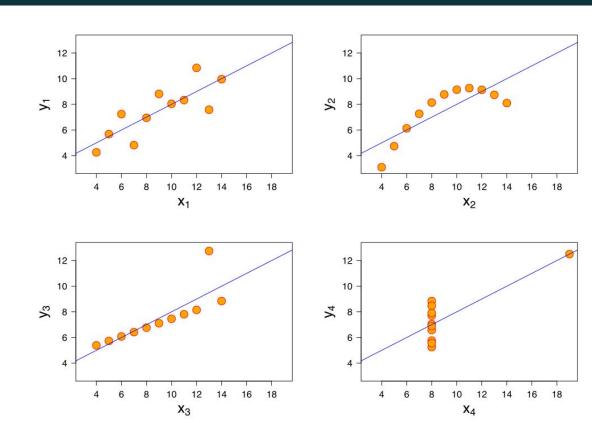


Anscombe's quartet: "calculation are exact; graphs are rough!"

11 datapoints

- Mean (x) = 9
- Var(x) = 11
- Mean (y) = 7.50
- Var (y) ~ 4.12
- Cor(x, y) = 0.816
- Linear regression line:

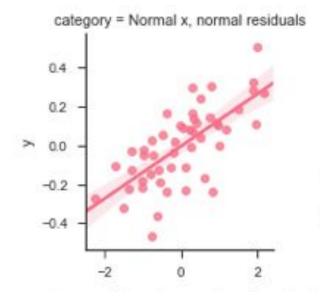
$$\circ$$
 y = 3.00 + 0.500x

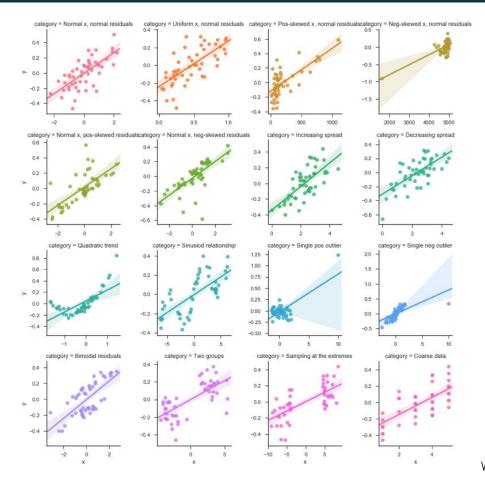


Anscombe, F. J. (1973). "Graphs in Statistical Analysis". American Statistician 27 (1): 17–21.

What does a correlation coefficient tell you about the data?

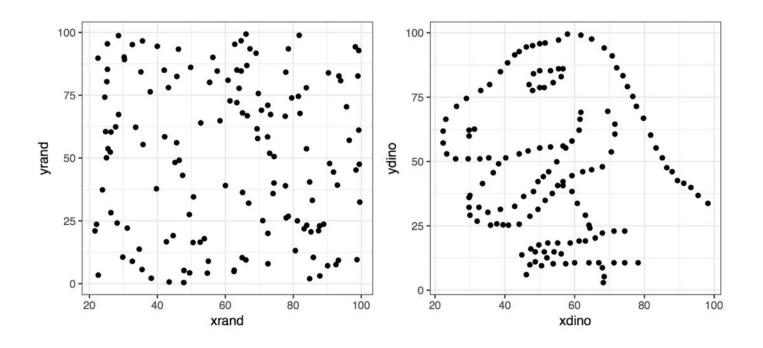
Correlation = 0.7





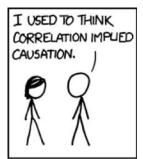
What does a correlation coefficient tell you about the data?

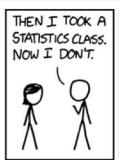
Correlation = -0.06

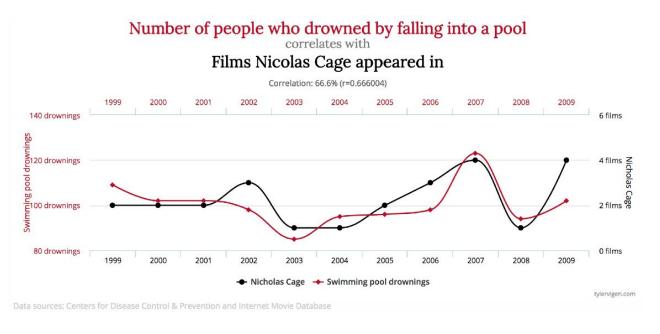


Spurious correlations

What does Nicholas Cage have to do with people drowning in swimming pools?









Many distance measures

Pearson Correlation Coefficient

Spearman Rank Correlation

Euclidean Distance

Mutual Information

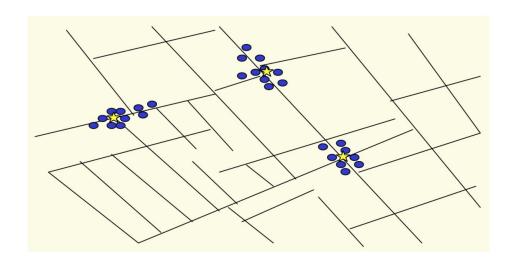
. . .

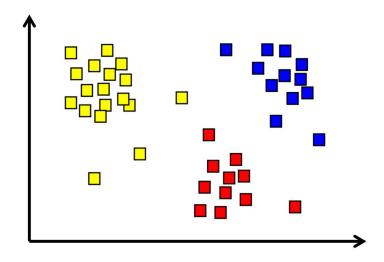
$$d = \sqrt{\sum_{i=1}^{n} (x_i - y_i)^2}$$

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

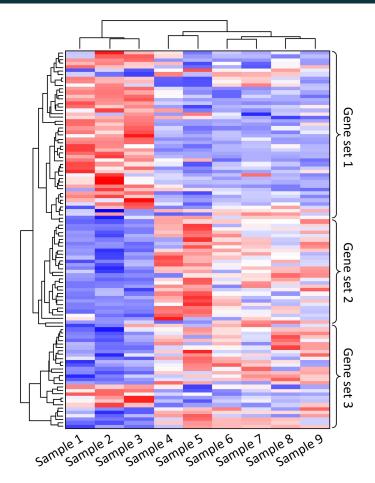
$$\rho = 1 - \frac{6\sum_{i=1}^{n} [rank(x_i) - rank(y_i)]}{n(n^2 - 1)}$$

Clustering





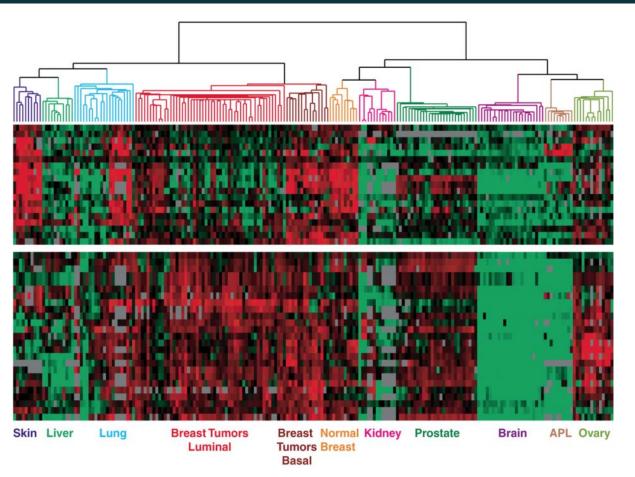
Clustering gene-expression profiles



Group-level Qs:

- 1. Are there groups of genes that respond similarly to changing contexts (across samples)?
- 2. Are there groups of samples that have very similar gene expression profiles?

Clustering gene-expression profiles



Group-level Qs:

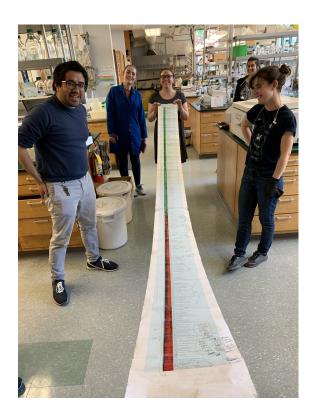
- 1. Are there groups of genes that respond similarly to changing contexts (across samples)?
- 2. Are there groups of samples that have very similar gene expression profiles?

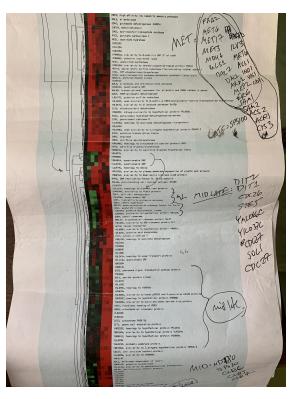
Clustering gene-expression profiles



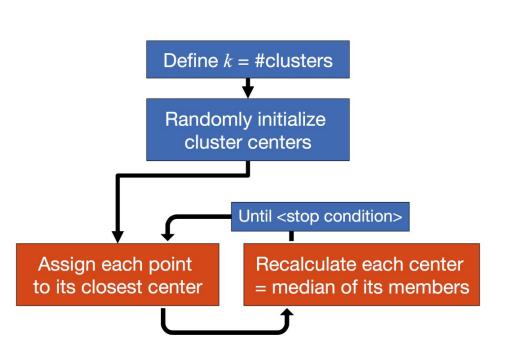
Inspired by @UCSDCooperLab's question about origins of the red/green color scheme in microarray clustering, I present THE FIRST dna microarray cluster analysis made by me in 1997 for

ncbi.nlm.nih.gov/m/pubmed/97841... w/handwritten notes from Pat Brown and the late Ira Herskowitz.





K-means clustering



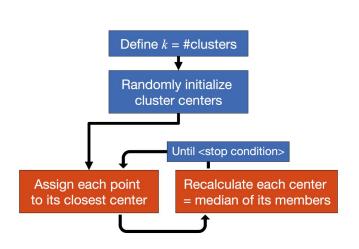
Conceptually similar to Expectation-Maximization, alternating between 2 two steps:

- E step: Creates a function for the expectation of the log-likelihood evaluated using the current estimate for the parameters.
- M step: Computes parameters maximizing the expected log- likelihood found on the E step.

These parameter-estimates are then used to determine the distribution of the latent variables in the next E step.

Checkout http://www.naftaliharris.com/blog/visualizing-k-means-clustering/

K-means clustering



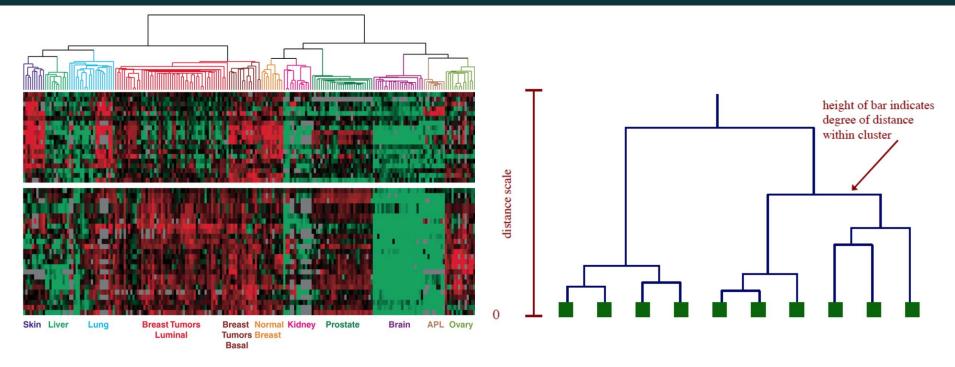
Stopping condition

- Until the change in centers is less than <constant>.
- Until all genes get assigned to the same partition.
- Until some minimal number of genes (e.g. 90%) get assigned to the same partition twice in a row.

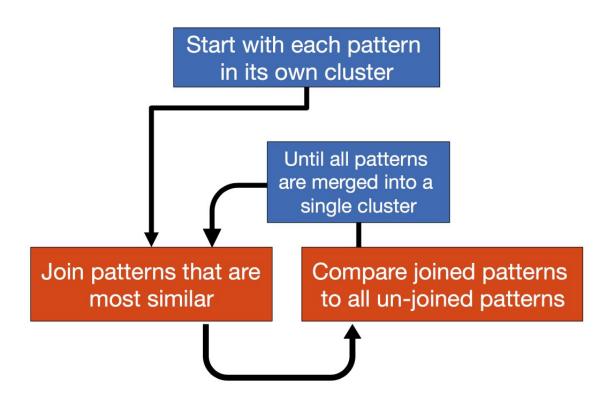
Some issues

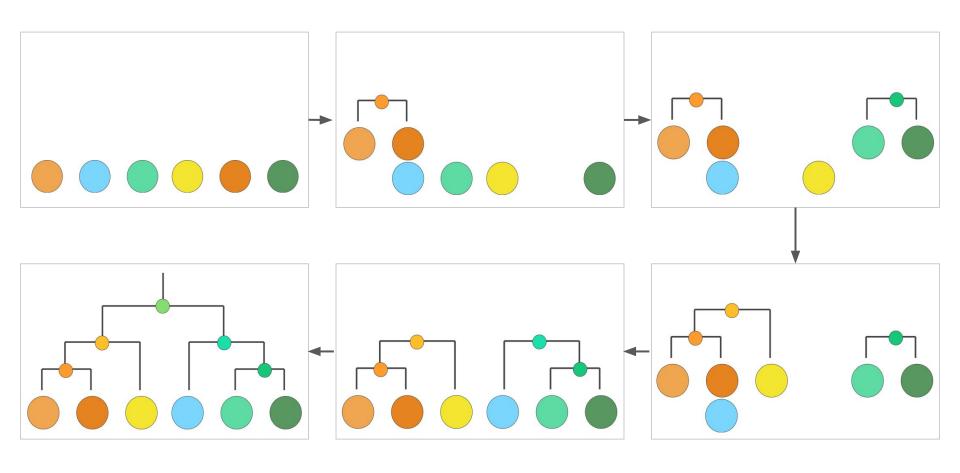
- Have to set k ahead of time.
- Works well if clusters of approx. similar sizes.
- Each gene only belongs to 1 cluster.

Genes assigned to clusters on the basis of all experiments;
 Experiments assigned to clusters based on all genes.



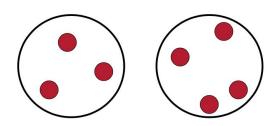
- Imposes hierarchical structure on all of the data.
- Easy visualization of similarities and differences between genes (experiments) and clusters of genes (experiments).



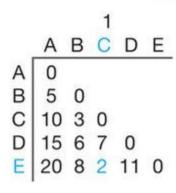


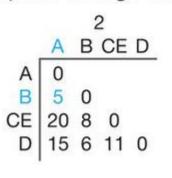
Linkage criteria:

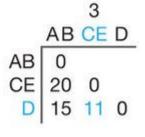
- Single/Minimum linkage (nearest neighbors)
- Complete/Maximum linkage (farthest neighbors)
- Average linkage (average of all pairs)

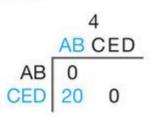


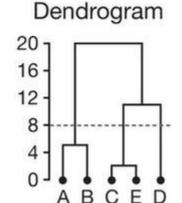
Complete linkage clustering of 5 objects





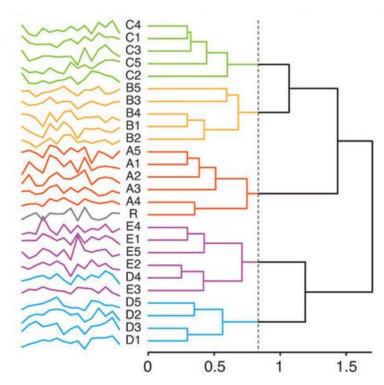






Complete linkage clustering

Tends to create balanced dendrograms by first clustering objects into small nodes and then clustering the nodes



Single linkage clustering

Tends to create stringy dendrograms by first creating a few nodes and then adding objects to them one at a time

