

# Unsupervised clustering

# Selecting a distance

- Unsupervised method need a criteria to assess profile similarity (and dissimilarity) : a distance

Table 1 Gene expression similarity measures

Manhattan distance

(city-block distance, L1 norm)

$$d_{fg} = \sum_c |e_{fc} - e_{gc}|$$

Euclidean distance

(L2 norm)

$$d_{fg} = \sqrt{\sum_c (e_{fc} - e_{gc})^2}$$

Mahalanobis distance

$$d_{fg} = (\mathbf{e}_f - \mathbf{e}_g)' \boldsymbol{\Sigma}^{-1} (\mathbf{e}_f - \mathbf{e}_g), \text{ where } \boldsymbol{\Sigma} \text{ is the (full or within-cluster) covariance matrix of the data}$$

Pearson correlation

(centered correlation)

$$d_{fg} = 1 - r_{fg}, \text{ with } r_{fg} = \frac{\sum_c (e_{fc} - \bar{e}_f)(e_{gc} - \bar{e}_g)}{\sqrt{\sum_c (e_{fc} - \bar{e}_f)^2 \sum_c (e_{gc} - \bar{e}_g)^2}}$$

Uncentered correlation

(angular separation, cosine angle)

$$d_{fg} = 1 - r_{fg}, \text{ with } r_{fg} = \frac{\sum_c e_{fc} e_{gc}}{\sqrt{\sum_c e_{fc}^2 \sum_c e_{gc}^2}}$$

Spellman rank correlation

As Pearson correlation, but replace  $e_{gc}$  with the rank of  $e_{gc}$  within the expression values of gene  $g$  across all conditions  $c = 1 \dots C$

Absolute or squared correlation

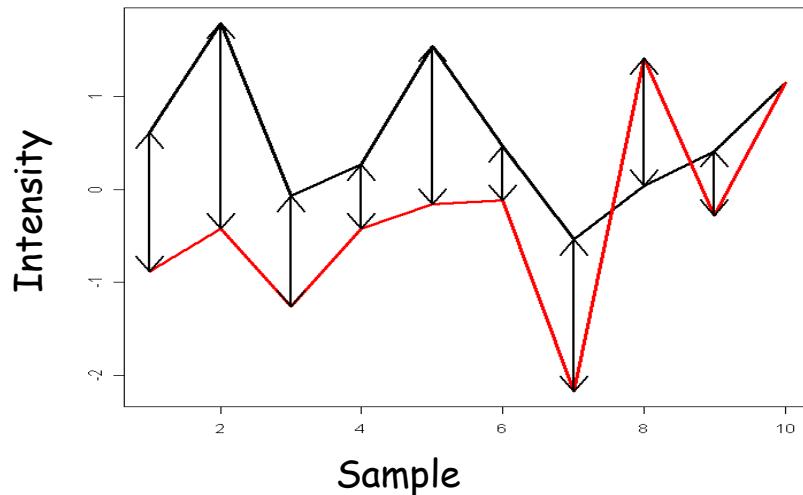
$$d_{fg} = 1 - |r_{fg}| \text{ or } d_{fg} = 1 - r_{fg}^2$$

$d_{fg}$ : distance between expression patterns for genes  $f$  and  $g$ .  $e_{gc}$ : expression level of gene  $g$  under condition  $c$ .

How does gene expression clustering work ? P. D'Haeseleer. Nat Biotech. 23(12). 2005.

# Euclidean distance

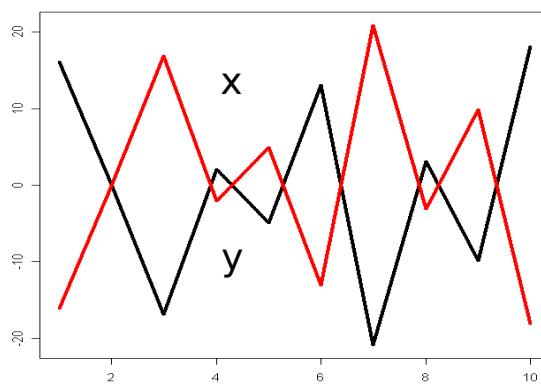
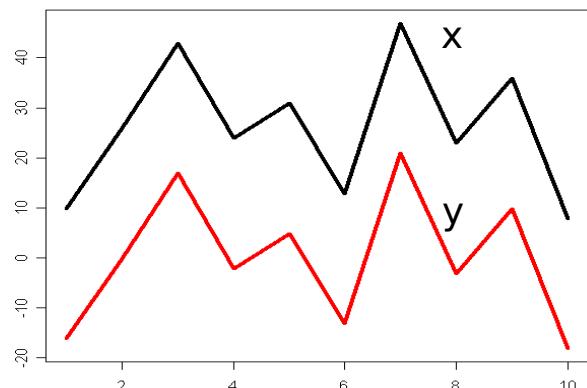
$$d(x, y) = \sqrt{\sum_{i=1}^n (x_i - y_i)^2}$$



$d(x, y) = 82.2$

$d(x, y) = 80.4$

● Limitations



- Highly sensitive to global expression level
- Does not measure anti-correlation
- No upper bound

# Pearson correlation coefficient ( $r$ )

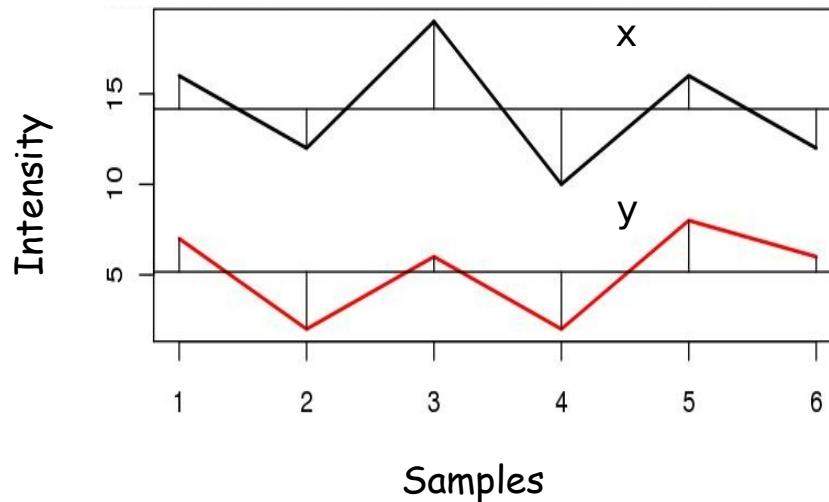
- Obtained by dividing the covariance of the genes by the product of their standard deviations
  - Covariance is a measure of how much two variables change together

$$cov(x, y) = \frac{1}{n-1} \sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y})$$

$$r = \frac{cov(x, y)}{\sqrt{var(x) var(y)}}$$

- Limits : -1 and 1
- Distance is obtained using:

$$d = (1 - r) / 2$$



# Spearman correlation coefficient (*rho*)

- Used to measure the degree of correspondence between two rankings

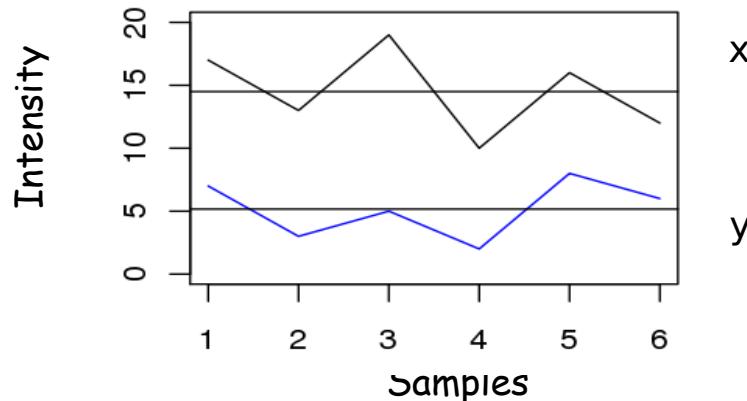
$$\rho = 1 - \frac{6 \sum d_i^2}{n(n^2 - 1)}$$

Intensities

	1	2	3	4	5	6
x	17	13	19	10	16	12
y	7	3	5	2	8	6

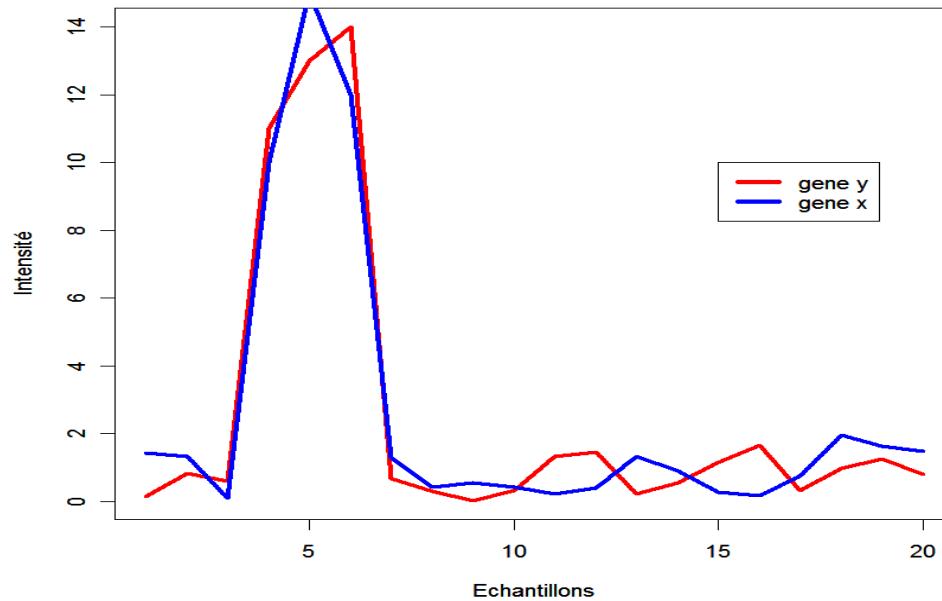
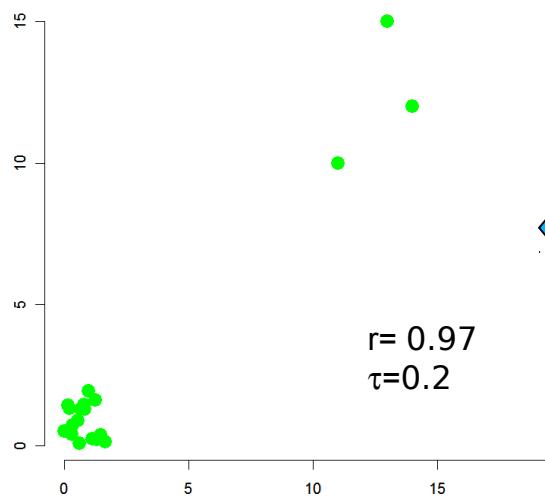
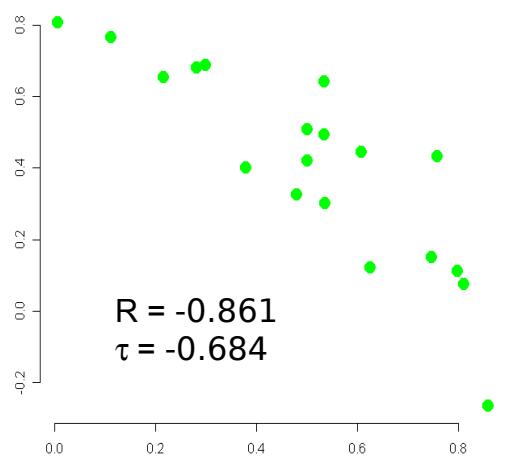
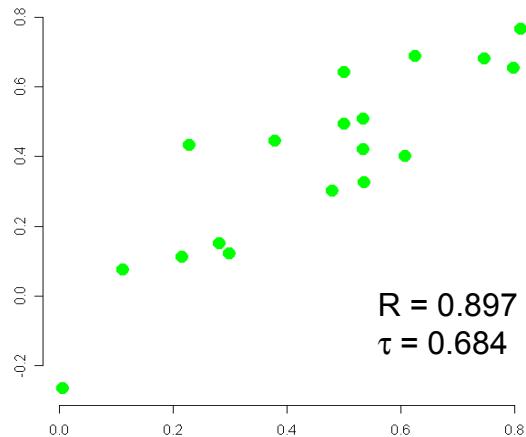
Ranks

	1	2	3	4	5	6
x	5	3	6	1	4	2
y	5	2	3	1	6	4
d	0	1	3	0	-2	-2



$$\rho = 1 - \frac{108}{210} = 0.485$$

# Pearson and spearman



---

# Hierarchical clustering

# Distance matrix

## Expression matrix

	E1	E2	E3	E4	E5	E6	E7	E8	E9	E10
g1	-0.733	-1.098	-0.741	1.778	0.119	-0.510	0.446	0.225	0.244	0.006
g2	-0.660	-0.772	0.586	1.201	-1.785	-0.108	-0.191	0.160	1.585	1.840
g3	0.753	-1.088	-1.672	-0.659	-0.709	-0.496	0.158	-0.024	0.355	1.418
g4	0.498	1.006	0.331	-0.960	0.925	-1.174	-0.790	-0.334	0.338	-0.105
g5	-1.020	-0.450	1.332	0.153	-0.611	0.515	-0.323	-1.014	0.831	0.656
g6	0.616	-0.366	0.038	2.375	-0.526	1.089	-1.356	0.591	-2.298	0.204
g7	0.670	0.691	-0.534	1.611	1.376	-0.334	-0.420	1.892	0.234	0.788
g8	-0.978	-0.279	-0.654	-0.337	-0.722	2.244	0.376	-0.145	-1.190	0.634
g9	2.669	0.790	-0.445	0.712	-1.073	1.010	-1.598	-1.398	0.944	-0.483
g10	0.318	0.414	1.462	0.251	1.033	-1.032	0.841	0.662	-0.209	-0.485
g11	-0.737	0.139	0.439	1.524	-0.337	0.835	-0.001	-1.082	-0.795	0.943
g12	-0.771	-0.577	0.221	-1.217	0.302	0.804	-0.424	-0.465	-0.826	0.258
g13	1.355	2.217	0.566	-1.470	-0.051	1.060	1.871	0.552	-0.555	-0.999
g14	0.640	2.381	-0.928	0.577	2.518	-0.164	-1.731	0.302	1.725	-2.903
g15	0.896	0.132	-0.368	1.224	-0.306	-1.046	-0.869	-0.969	1.453	-1.264
g16	0.093	-0.749	0.923	-0.157	1.180	-0.520	-0.098	-0.429	-0.011	-0.450
g17	2.415	-2.154	-0.900	-0.969	-0.970	-2.180	0.619	0.057	0.170	-0.979
g18	-0.580	-2.411	0.959	1.523	0.047	-1.019	-0.727	0.741	-0.294	1.205
g19	-0.085	1.080	-0.350	-0.332	-0.493	-0.547	-0.478	-0.304	0.766	-0.271
g20	-1.002	-0.954	-0.674	-0.737	-0.427	0.919	1.140	-0.468	0.041	-2.620

## Distance matrix (samples)

	E1	E2	E3	E4	E5	E6	E7	E8	E9
E2	<b>6.370</b>								
E3	<b>6.867</b>	<b>6.548</b>							
E4	<b>7.277</b>	<b>8.038</b>	<b>6.120</b>						
E5	<b>6.531</b>	<b>5.244</b>	<b>5.461</b>	<b>6.797</b>					
E6	<b>7.635</b>	<b>6.183</b>	<b>5.893</b>	<b>6.964</b>	<b>7.054</b>				
E7	<b>6.869</b>	<b>7.203</b>	<b>5.177</b>	<b>7.990</b>	<b>6.463</b>	<b>5.958</b>			
E8	<b>5.839</b>	<b>6.323</b>	<b>5.027</b>	<b>5.631</b>	<b>4.373</b>	<b>6.270</b>	<b>4.755</b>		
E9	<b>5.804</b>	<b>6.248</b>	<b>6.080</b>	<b>6.676</b>	<b>5.942</b>	<b>7.574</b>	<b>6.659</b>	<b>6.038</b>	
E10	<b>8.250</b>	<b>9.004</b>	<b>5.719</b>	<b>6.475</b>	<b>8.211</b>	<b>6.866</b>	<b>6.784</b>	<b>5.935</b>	<b>7.655</b>

	g1	g2	g3	g4	g5	g6	g7	g8	g9	g10	g11	g12	g13	g14	g15	g16	g17	g18	g19
g2	<b>3.401</b>																		
g3	<b>3.443</b>	<b>3.716</b>																	
g4	<b>4.188</b>	<b>4.864</b>	<b>3.873</b>																
g5	<b>3.452</b>	<b>2.650</b>	<b>4.035</b>	<b>3.658</b>															
g6	<b>4.022</b>	<b>5.078</b>	<b>5.138</b>	<b>5.343</b>	<b>4.824</b>														
g7	<b>3.315</b>	<b>4.622</b>	<b>4.291</b>	<b>3.781</b>	<b>4.823</b>	<b>4.122</b>													
g8	<b>4.012</b>	<b>4.525</b>	<b>3.919</b>	<b>4.885</b>	<b>3.542</b>	<b>4.088</b>	<b>5.035</b>												
g9	<b>5.257</b>	<b>5.176</b>	<b>4.681</b>	<b>4.383</b>	<b>4.707</b>	<b>4.884</b>	<b>5.192</b>	<b>5.339</b>											
g10	<b>3.532</b>	<b>4.772</b>	<b>4.615</b>	<b>2.694</b>	<b>3.769</b>	<b>4.882</b>	<b>3.400</b>	<b>4.886</b>	<b>5.445</b>										
g11	<b>2.979</b>	<b>3.464</b>	<b>4.190</b>	<b>4.201</b>	<b>2.473</b>	<b>3.240</b>	<b>4.200</b>	<b>2.881</b>	<b>4.662</b>	<b>3.919</b>									
g12	<b>3.791</b>	<b>4.470</b>	<b>3.530</b>	<b>3.183</b>	<b>2.706</b>	<b>4.452</b>	<b>4.654</b>	<b>2.410</b>	<b>5.073</b>	<b>3.684</b>	<b>3.084</b>								
g13	<b>5.822</b>	<b>6.447</b>	<b>5.454</b>	<b>4.234</b>	<b>5.341</b>	<b>6.115</b>	<b>5.324</b>	<b>4.699</b>	<b>5.381</b>	<b>3.886</b>	<b>5.275</b>	<b>4.544</b>							
g14	<b>6.061</b>	<b>7.599</b>	<b>6.991</b>	<b>4.523</b>	<b>6.595</b>	<b>6.972</b>	<b>5.051</b>	<b>7.249</b>	<b>5.527</b>	<b>5.393</b>	<b>6.699</b>	<b>6.318</b>	<b>6.078</b>						
g15	<b>3.363</b>	<b>4.313</b>	<b>4.191</b>	<b>3.272</b>	<b>3.867</b>	<b>5.032</b>	<b>4.238</b>	<b>5.478</b>	<b>3.201</b>	<b>3.931</b>	<b>4.199</b>	<b>4.620</b>	<b>5.579</b>	<b>4.442</b>					
g16	<b>3.072</b>	<b>4.431</b>	<b>3.868</b>	<b>2.337</b>	<b>2.863</b>	<b>4.625</b>	<b>3.849</b>	<b>4.244</b>	<b>4.811</b>	<b>2.063</b>	<b>3.416</b>	<b>2.488</b>	<b>4.658</b>	<b>5.311</b>	<b>3.342</b>				
g17	<b>4.858</b>	<b>5.819</b>	<b>3.673</b>	<b>4.765</b>	<b>5.786</b>	<b>6.407</b>	<b>5.866</b>	<b>6.299</b>	<b>5.462</b>	<b>4.919</b>	<b>6.246</b>	<b>5.325</b>	<b>6.026</b>	<b>7.351</b>	<b>4.398</b>	<b>4.449</b>			
g18	<b>2.887</b>	<b>3.418</b>	<b>4.230</b>	<b>4.844</b>	<b>3.695</b>	<b>4.171</b>	<b>4.187</b>	<b>5.013</b>	<b>6.206</b>	<b>4.118</b>	<b>3.814</b>	<b>4.191</b>	<b>7.094</b>	<b>7.609</b>	<b>4.766</b>	<b>3.486</b>	<b>5.324</b>		
g19	<b>3.410</b>	<b>3.778</b>	<b>3.288</b>	<b>1.979</b>	<b>2.969</b>	<b>4.891</b>	<b>3.781</b>	<b>3.997</b>	<b>3.749</b>	<b>3.235</b>	<b>3.472</b>	<b>3.103</b>	<b>4.101</b>	<b>4.730</b>	<b>2.573</b>	<b>2.937</b>	<b>4.744</b>	<b>4.749</b>	
g20	<b>4.082</b>	<b>5.730</b>	<b>4.873</b>	<b>4.849</b>	<b>4.366</b>	<b>5.830</b>	<b>5.959</b>	<b>3.905</b>	<b>5.731</b>	<b>4.591</b>	<b>4.766</b>	<b>3.640</b>	<b>4.734</b>	<b>6.100</b>	<b>4.552</b>	<b>3.876</b>	<b>5.136</b>	<b>5.809</b>	<b>4.011</b>

## Distance matrix (genes)

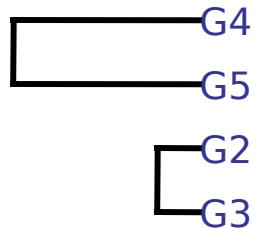
# Hierarchical clustering

	G1	G2	G3	G4	G5	G6
G1	0					
G2	0,1	0				
G3	0,16	0,03	0			
G4	0,9	0,8	0,5	0		
G5	0,25	0,4	0,77	0,12	0	
G6	1,3	1,5	0,95	1,0	1,1	0

	G1	G2-G <sub>3</sub>	G4	G5	G6
G1	0				
G2-G <sub>3</sub>	0,13	0			
G4	0,9	0,65	0		
G5	0,25	0,58	0,12	0	
G6	1,3	1,22	1,0	1,1	0

	G1	G2-G <sub>3</sub>	G4-G <sub>5</sub>	G6
G1	0			
G2-G <sub>3</sub>	0,13	0		
G4-G <sub>5</sub>	0,57	0,61	0	
G6	1,3	1,2	1,0	0

G2  
G3



# Agglomeration methods

	G1	G2	G3	G4	G5	G6
G1	0					
G2	0,1	0				
G3	0,16	0,03	0			
G4	0,9	0,8	0,5	0		
G5	0,25	0,4	0,77	0,12	0	
G6	1,3	1,5	0,95	1,0	1,1	0

Single

	G1	G2-G3	G4	G5	G6
G1	0				
G2-G3	0,1	0			
G4	0,9	0,5	0		
G5	0,25	0,4	0,12	0	
G6	1,3	0,95	1,0	1,1	0

Complete

	G1	G2-G3	G4	G5	G6
G1	0				
G2-G3	0,16	0			
G4	0,9	0,8	0		
G5	0,25	0,77	0,12	0	
G6	1,3	1,5	1,0	1,1	0

Average

	G1	G2-G3	G4	G5	G6
G1	0				
G2-G3	0,13	0			
G4	0,9	0,65	0		
G5	0,25	0,58	0,12	0	
G6	1,3	1,22	1,0	1,1	0

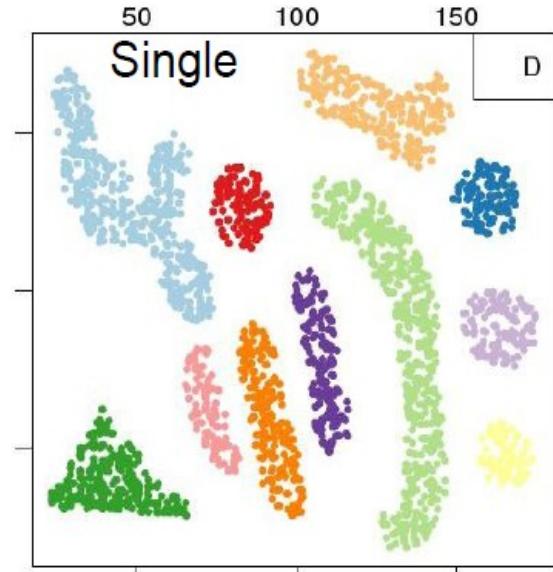
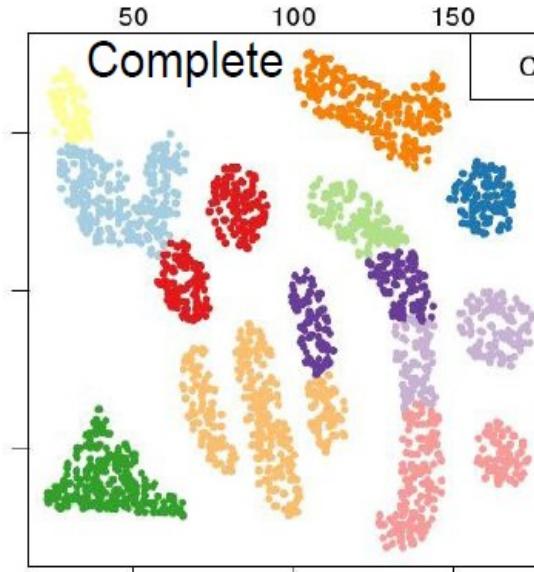
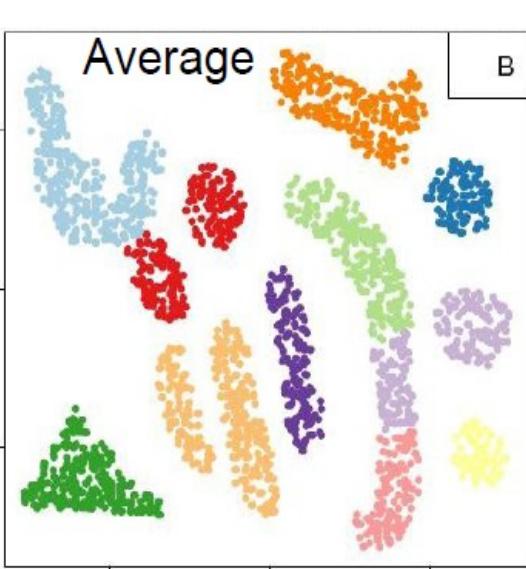
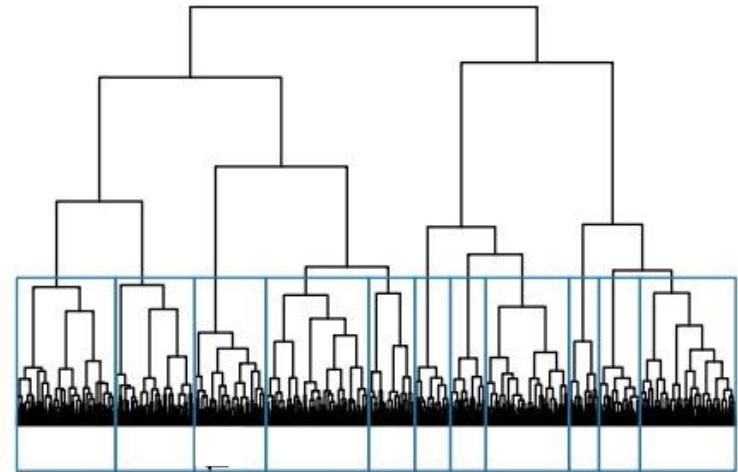
# Which agglomeration method ?

- An artificial 2D dataset

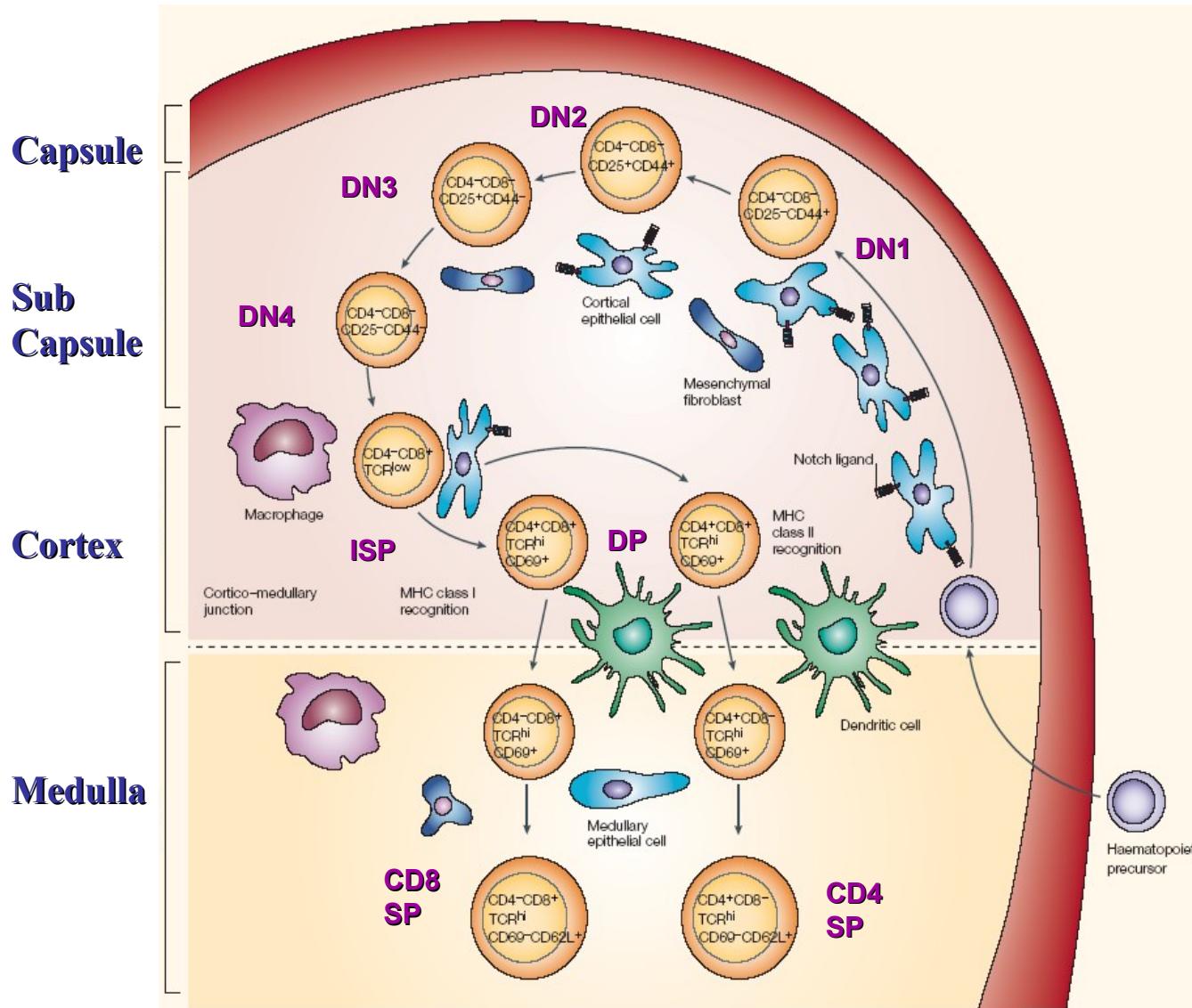
- Here single criteria performs best

- Microarrays

- Most generally average



# The Thymus project



(From Zúñiga-Pflücker JC, *Nature Reviews Immunology* 4, 67-72; 2004)



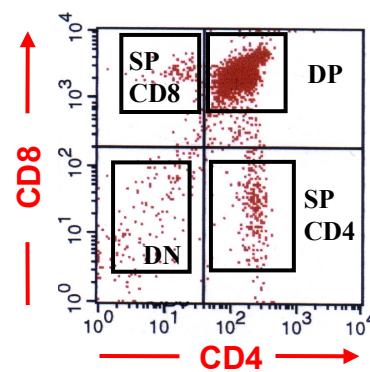
# Genetically engineered mice

## Thymi from WT and KO mice

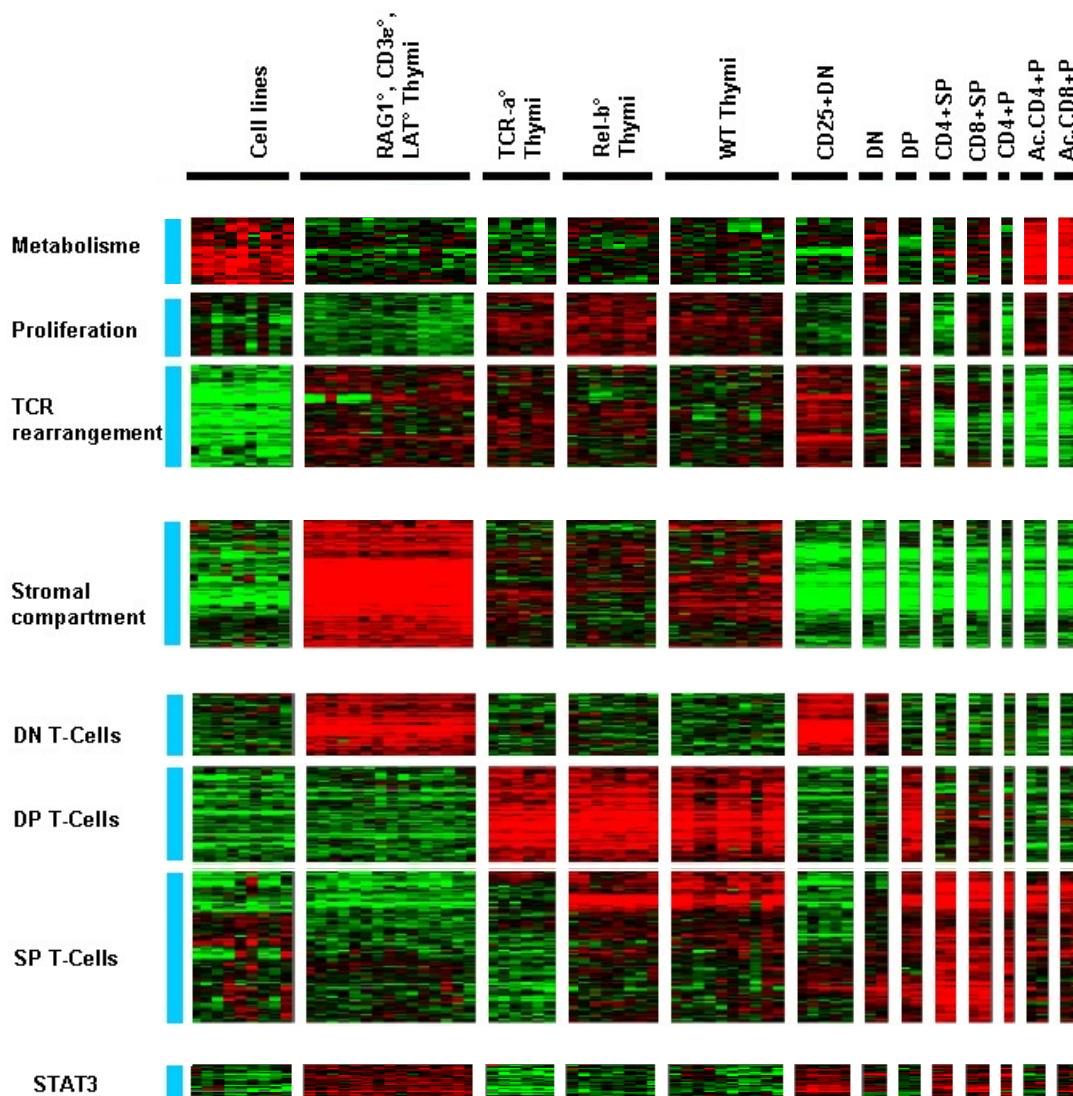
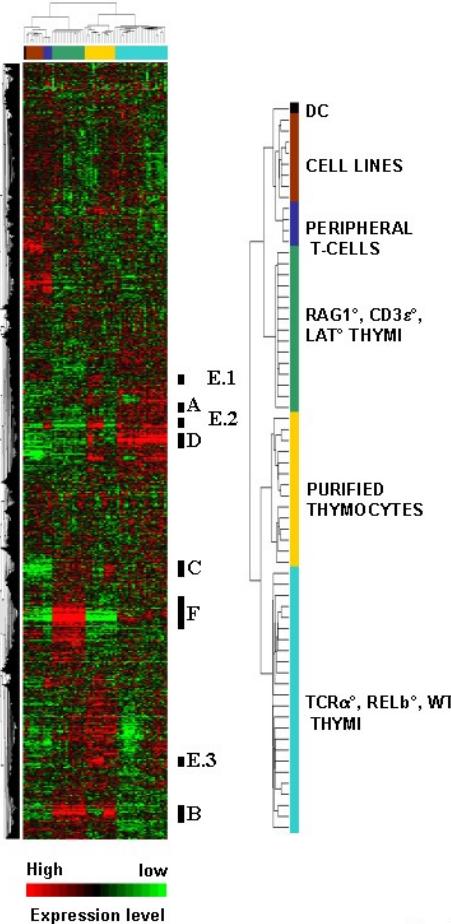
- ✓ Thymi from KO mice are enriched in particular cell types



## Purified thymocytes

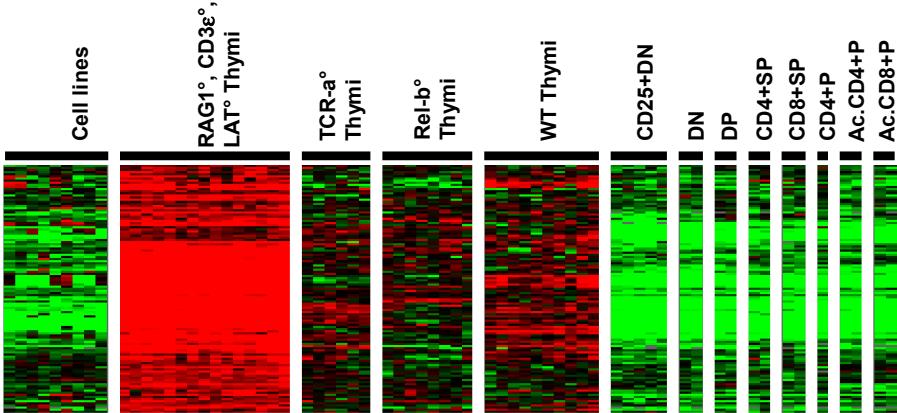


# Virtual microdissection



Puthier et al. A General Survey of Thymocyte Differentiation by Transcriptional Analysis of Knockout Mouse Models. *J. Immunol.*, 2004, 173(10):6109

# Virtual microdissection



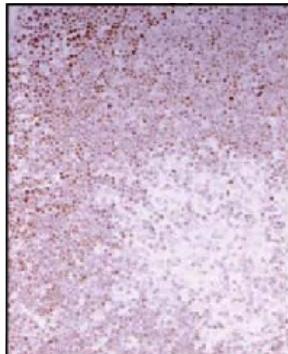
## SPATIAL

(Flommerfelt FA, Genes & Immunity, 2000, 1:391)



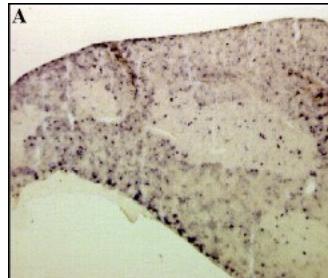
## CD83

(Fujimoto Y, Cell, 2002, 108: 755)



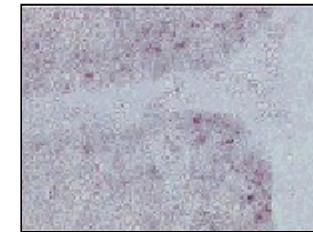
## TECK

(Wurbel MA, Eur.J. Immunol, 2000, 30:262)



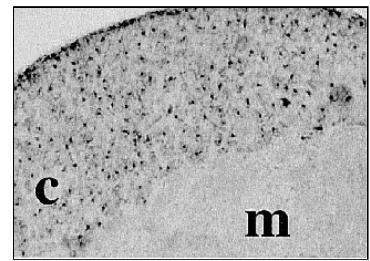
## SDF1

(Conrad CB, Eur.J. Immunol, 2000, 30:3371)



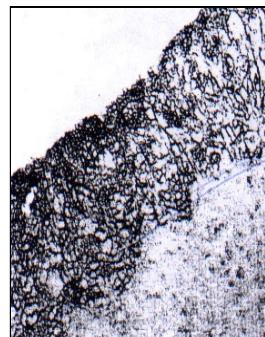
## TSSP

(Carrier A et al, Immunogenetics, 1999, 50 :255-270)



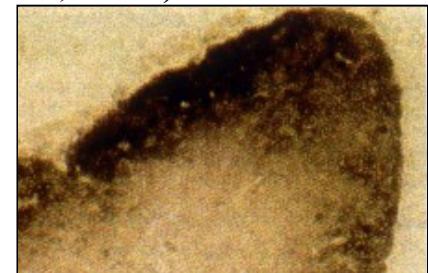
## Ly75/CD205

Kraal G, J. Exp. Med, 1986, 163:981)



## TSCOT

(Kim GM, J. Immunol, 2000, 164:3185)



# K- means algorithm

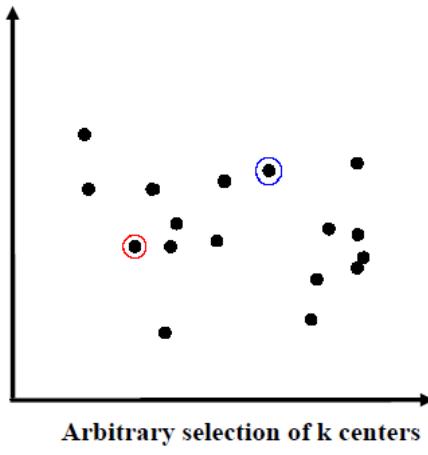
---

- K-means

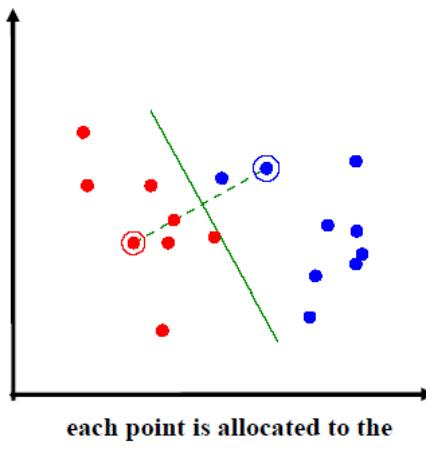
- Méthode de partitionnement
  - Définir k-classes de patients ou de gènes pour les analyser individuellement
- Définir a priori le nombre de classes (k)
  - ✚ Cet argument peut être assez difficile à déterminer lorsqu'on s'intéresse aux gènes (voire aux échantillons)
- L'initialisation se fait au hasard (plusieurs solutions possibles)

# K-means algorithm

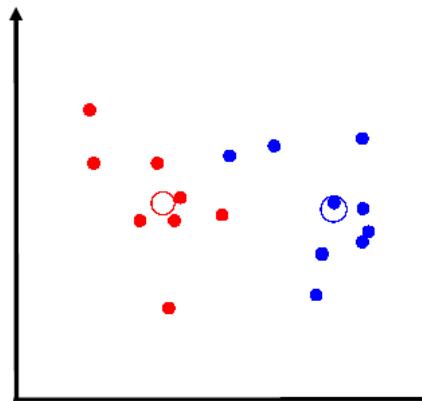
- The number of classes is given as input to the algorithm (here  $k=2$ )



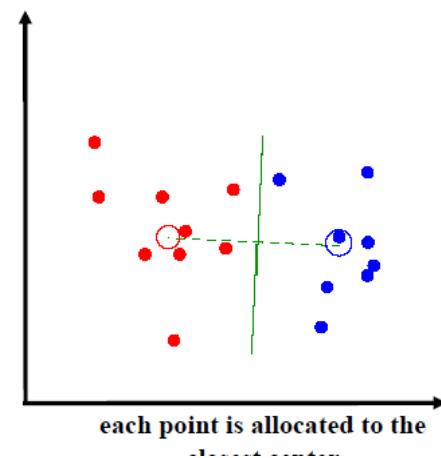
Arbitrary selection of  $k$  centers  
 $k=2$



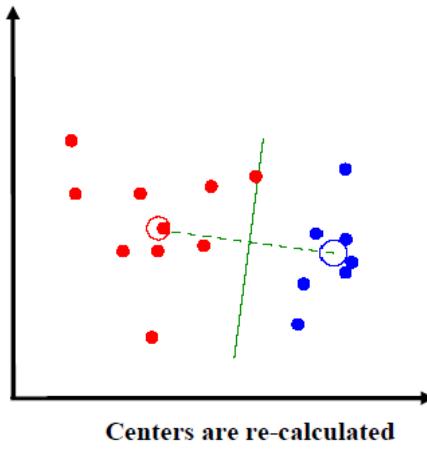
each point is allocated to the  
closest center



Centers are calculated



each point is allocated to the  
closest center



Centers are re-calculated

• • •