

# 6.874/... Recitation 2

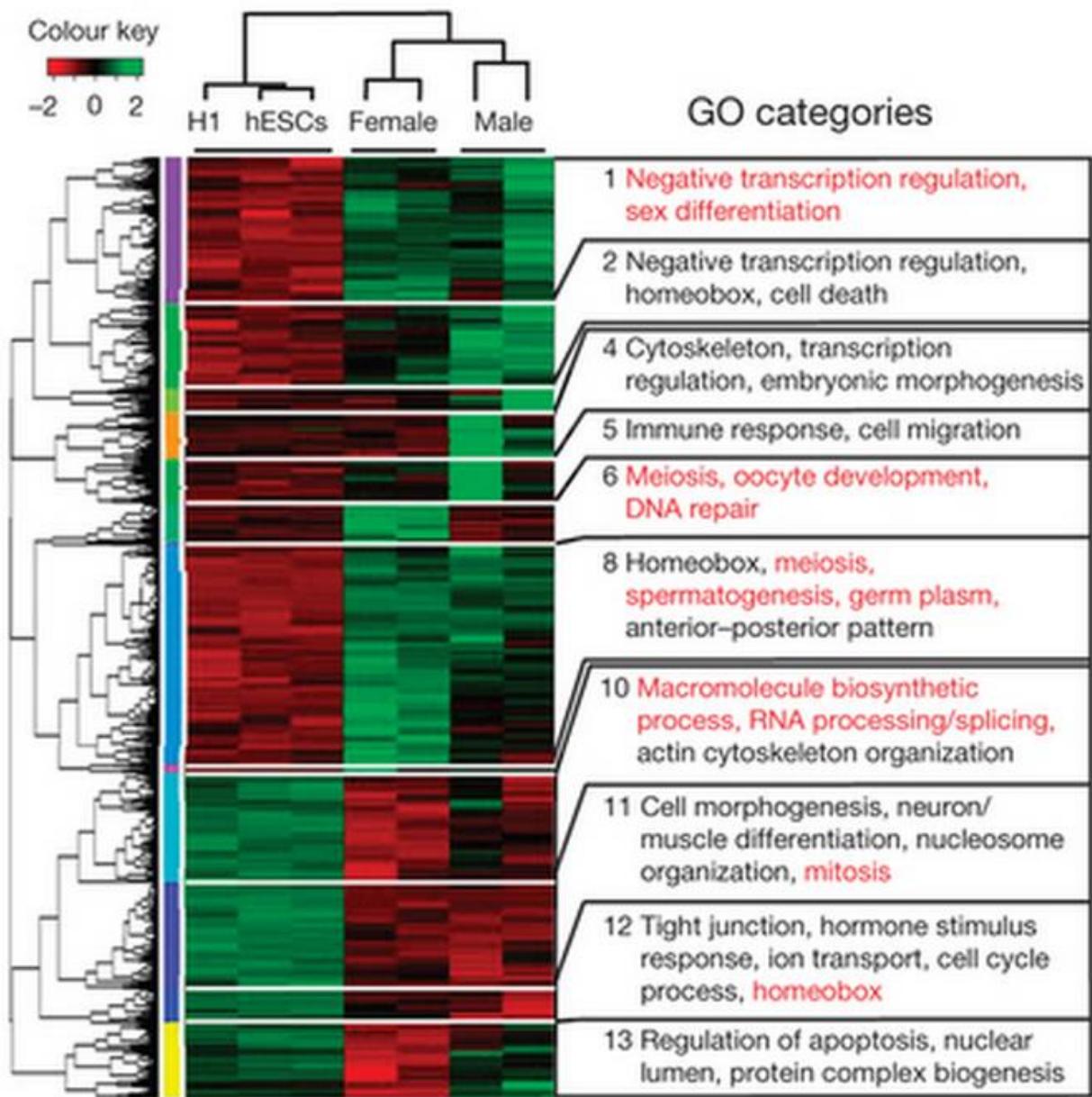
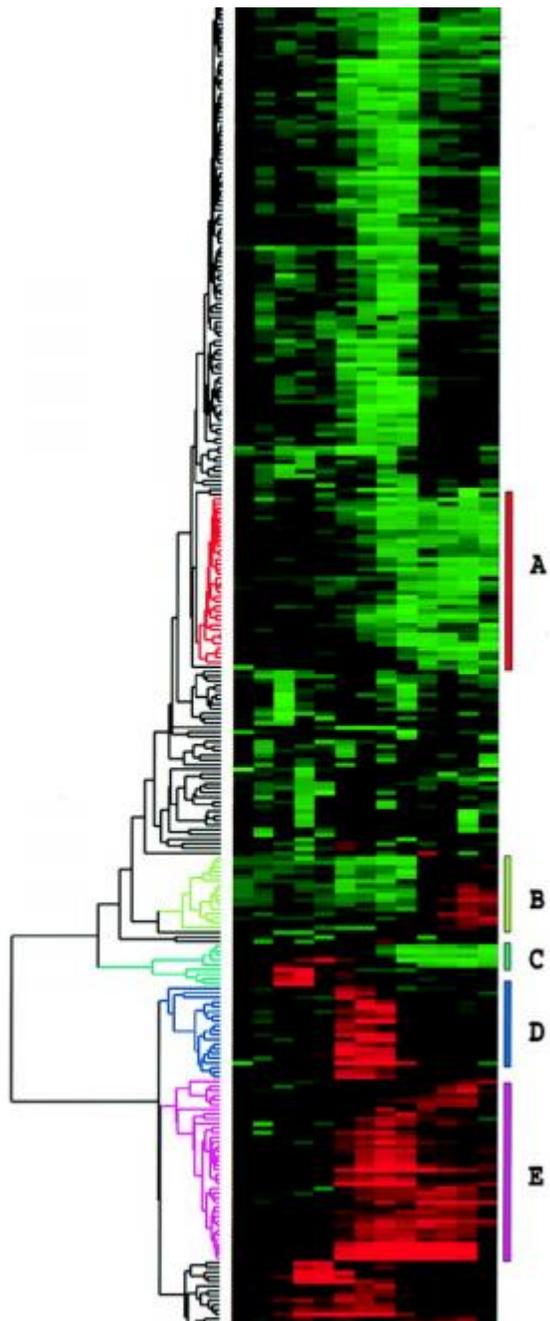
Courtesy of an MIT Teaching Assistant.

# Reminders

- Pset 1 posted – due Feb 20<sup>th</sup> (no extra problem)
- Pset 2 posted – due Mar 13<sup>th</sup>
- Project teams due – Feb 25<sup>th</sup>
  - Interests and background directory has been posted
- Lecture videos will be posted on MITx soon – next week?

# Today

- Clustering (6.874 topic)
- Biology review
- Alignment



Courtesy of Macmillan Publishers Limited. Used with permission.  
 Source: Gkoutela, Sofia, Ziwei Li, et al. "The Ontogeny of cKIT+ Human Primordial Germ Cells Proves to be a Resource for Human Germ Line Reprogramming, Imprint Erasure and in Vitro Differentiation." *Nature Cell Biology* 15, no. 1 (2013): 113-22.

# Clustering – K-means

- Group points together based on how ‘close’ they are to each other
- Dataset of unlabelled points:  $X = \{x_1, x_2, \dots, x_N\}, x_n \in R^d$
- Assume K clusters – each is defined by a centroid  $\mu_k$
- $r_{nk} = 1$  if  $x_n$  belongs to cluster k
- Find unknowns  $\mu_k$  and  $r_{nk}$

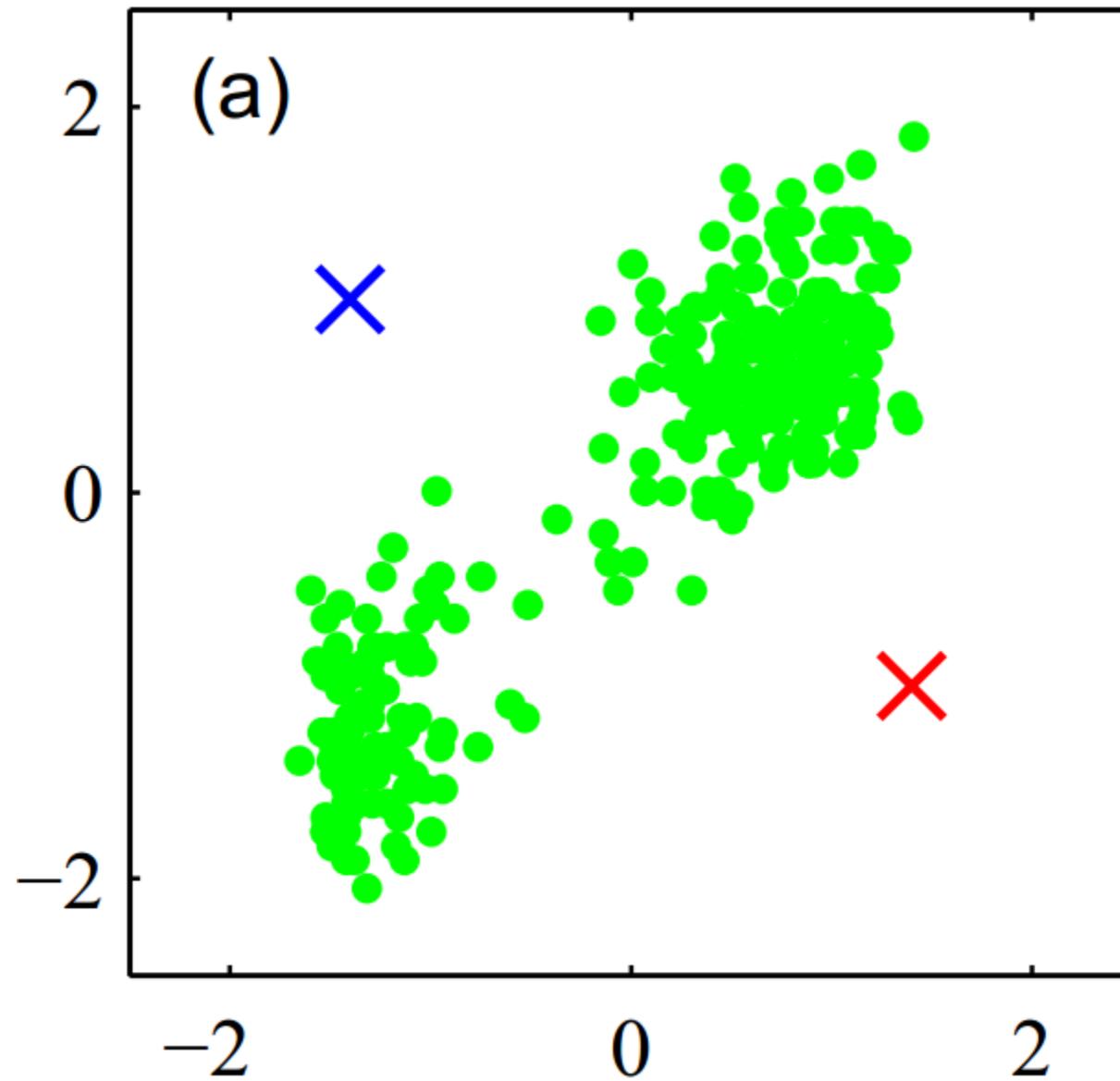
$$J = \sum_{n=1}^N \sum_{k=1}^K r_{nk} \| \mathbf{x}_n - \mu_k \|^2$$

---

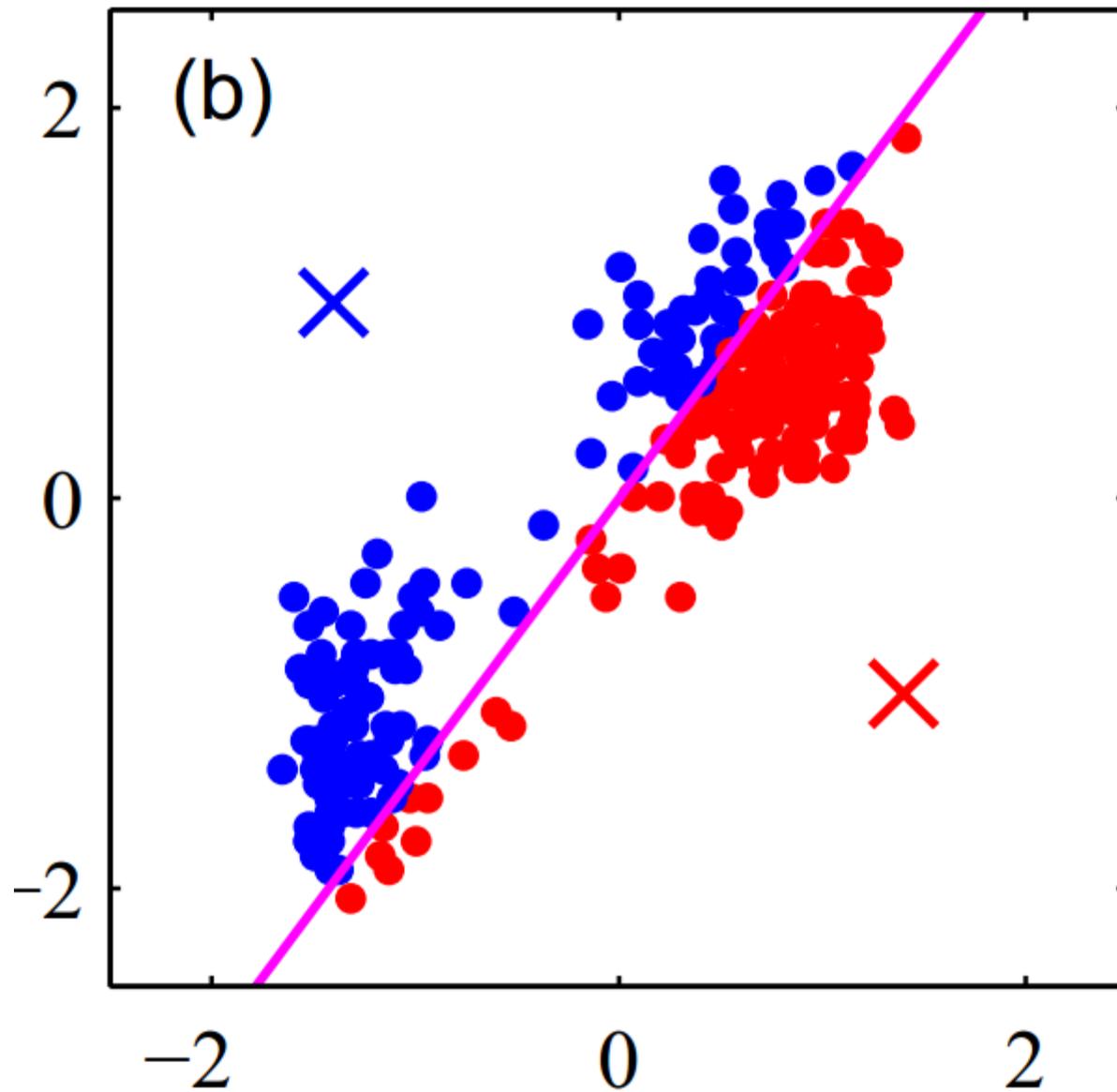
## Algorithm 10.1 *K*-Means Clustering

---

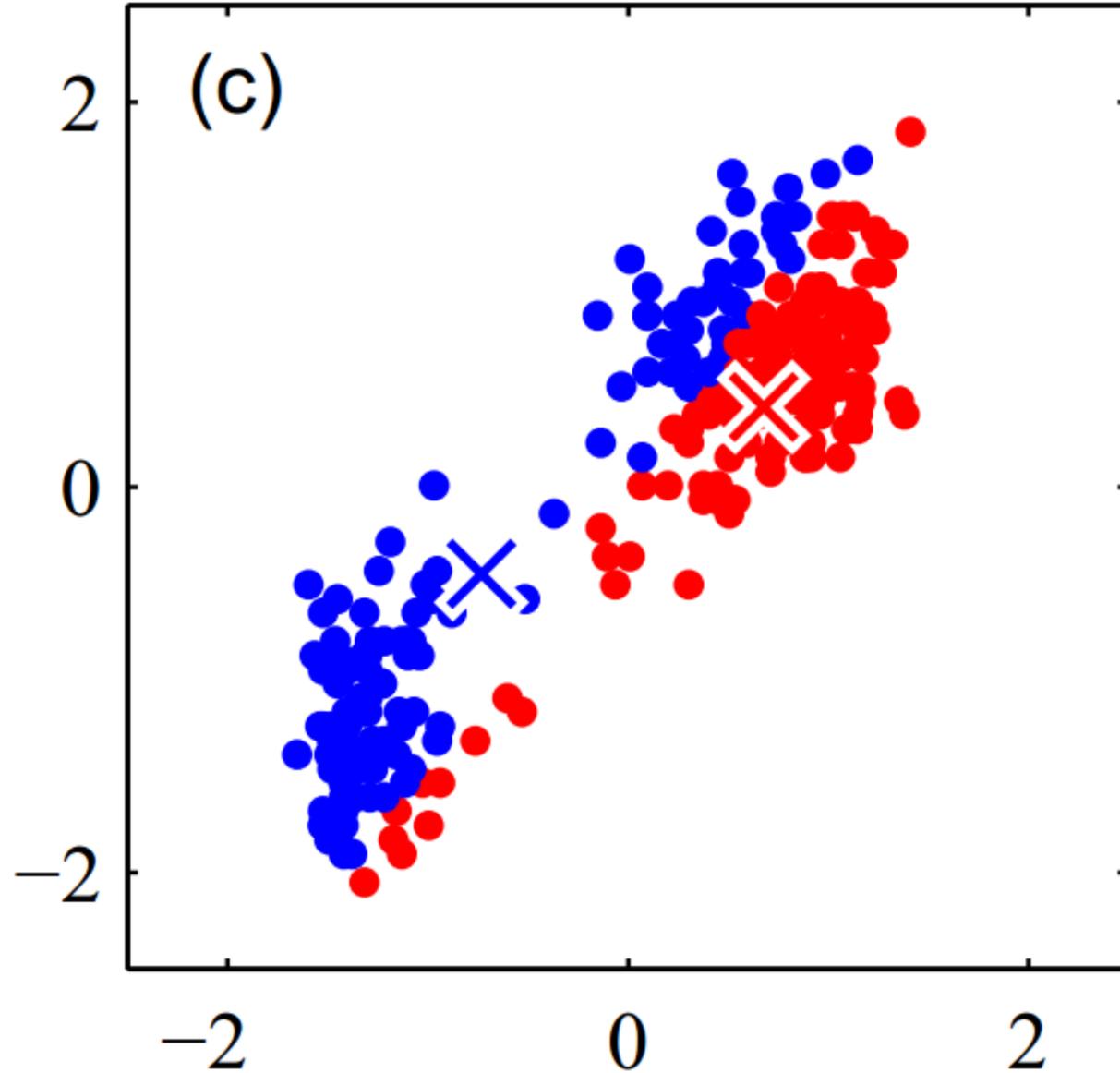
1. Randomly assign a number, from 1 to  $K$ , to each of the observations. These serve as initial cluster assignments for the observations.
2. Iterate until the cluster assignments stop changing:
  - (a) For each of the  $K$  clusters, compute the cluster *centroid*. The  $k$ th cluster centroid is the vector of the  $p$  feature means for the observations in the  $k$ th cluster.
  - (b) Assign each observation to the cluster whose centroid is closest (where *closest* is defined using Euclidean distance).



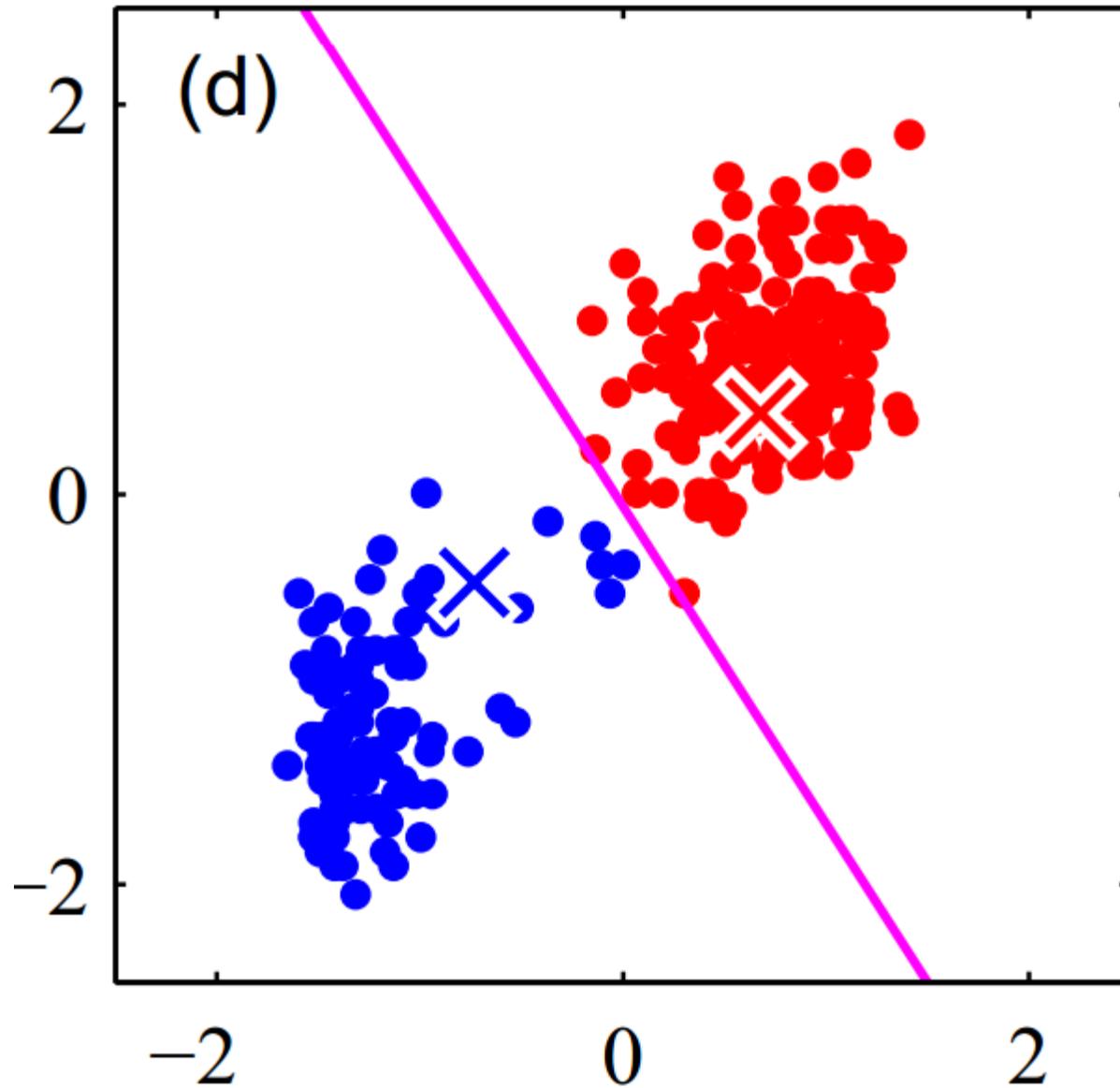
© unknown source. All rights reserved. This content is excluded from our Creative Commons license. For more information, see <http://ocw.mit.edu/help/faq-fair-use/>.



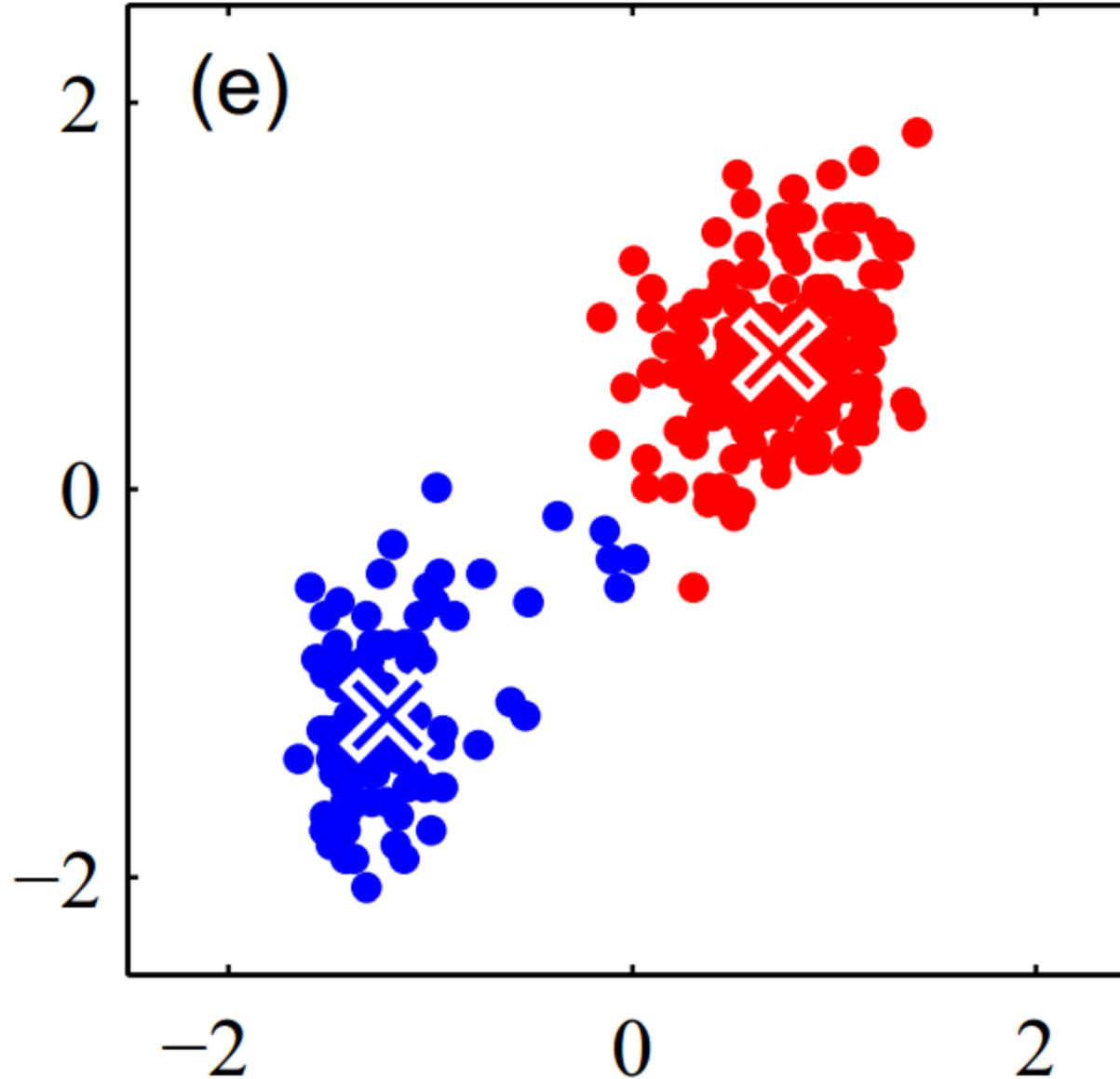
© unknown source. All rights reserved. This content is excluded from our Creative Commons license. For more information, see <http://ocw.mit.edu/help/faq-fair-use/>.



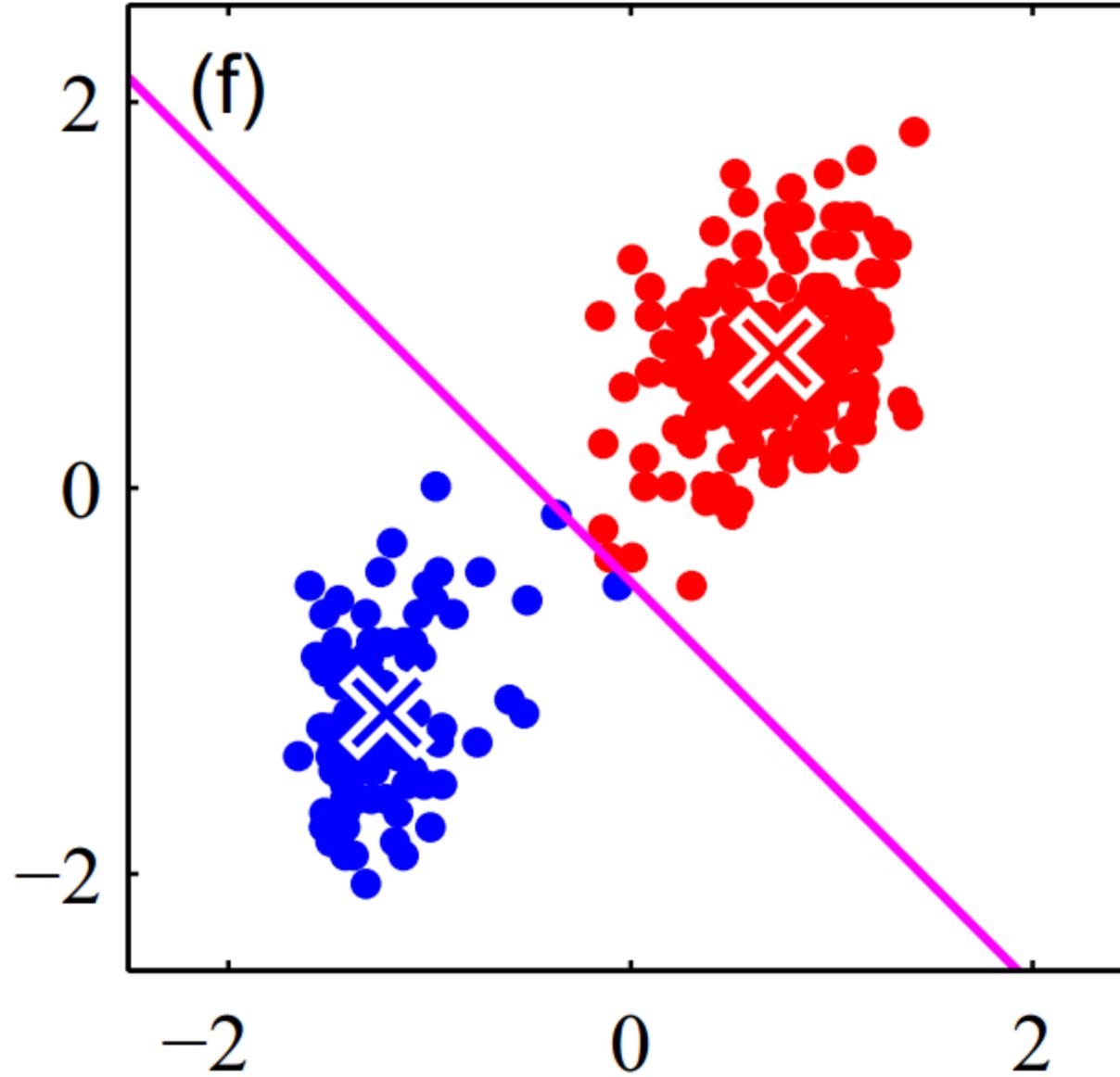
© unknown source. All rights reserved. This content is excluded from our Creative Commons license. For more information, see <http://ocw.mit.edu/help/faq-fair-use/>.



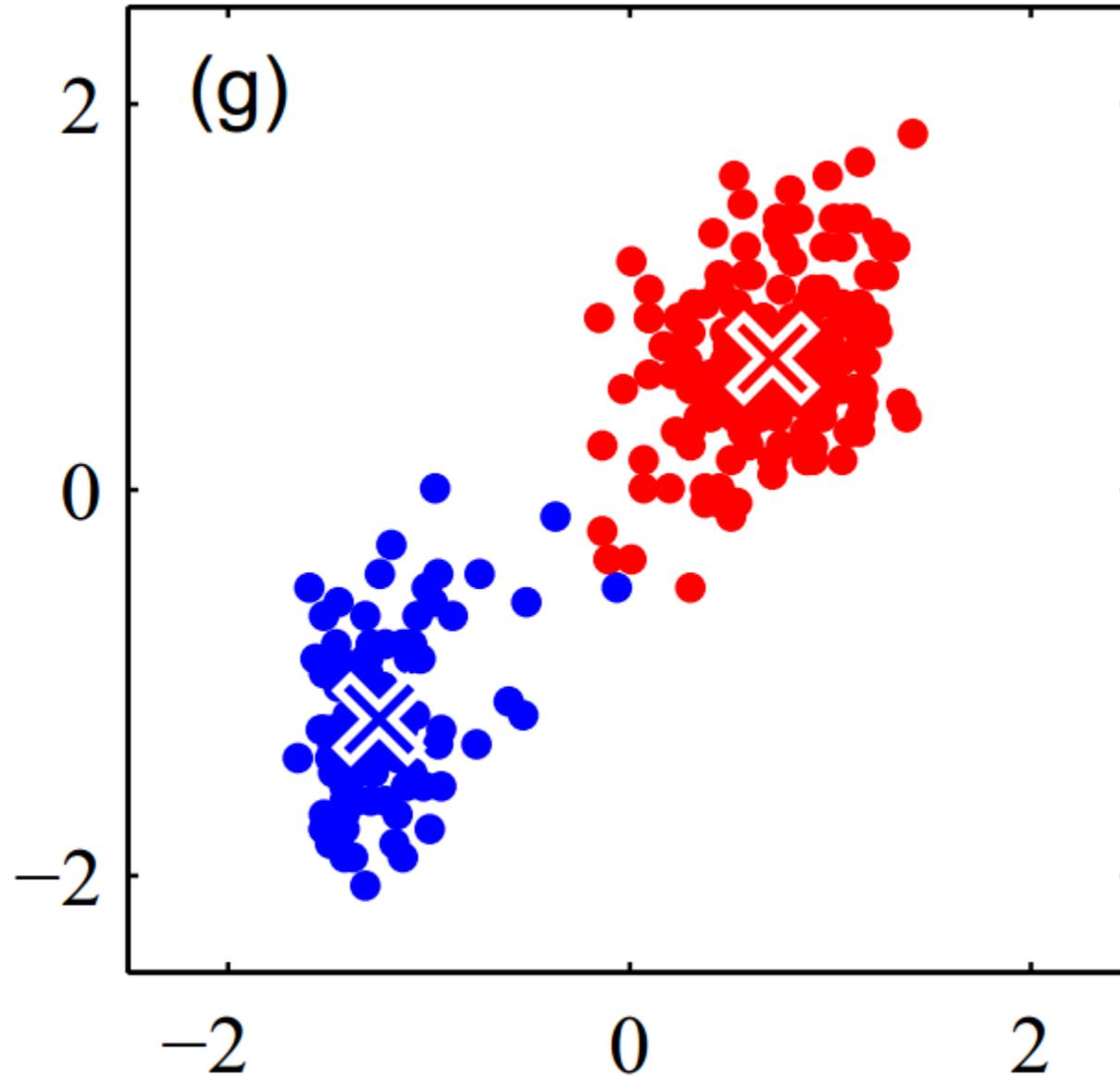
© unknown source. All rights reserved. This content is excluded from our Creative Commons license. For more information, see <http://ocw.mit.edu/help/faq-fair-use/>.



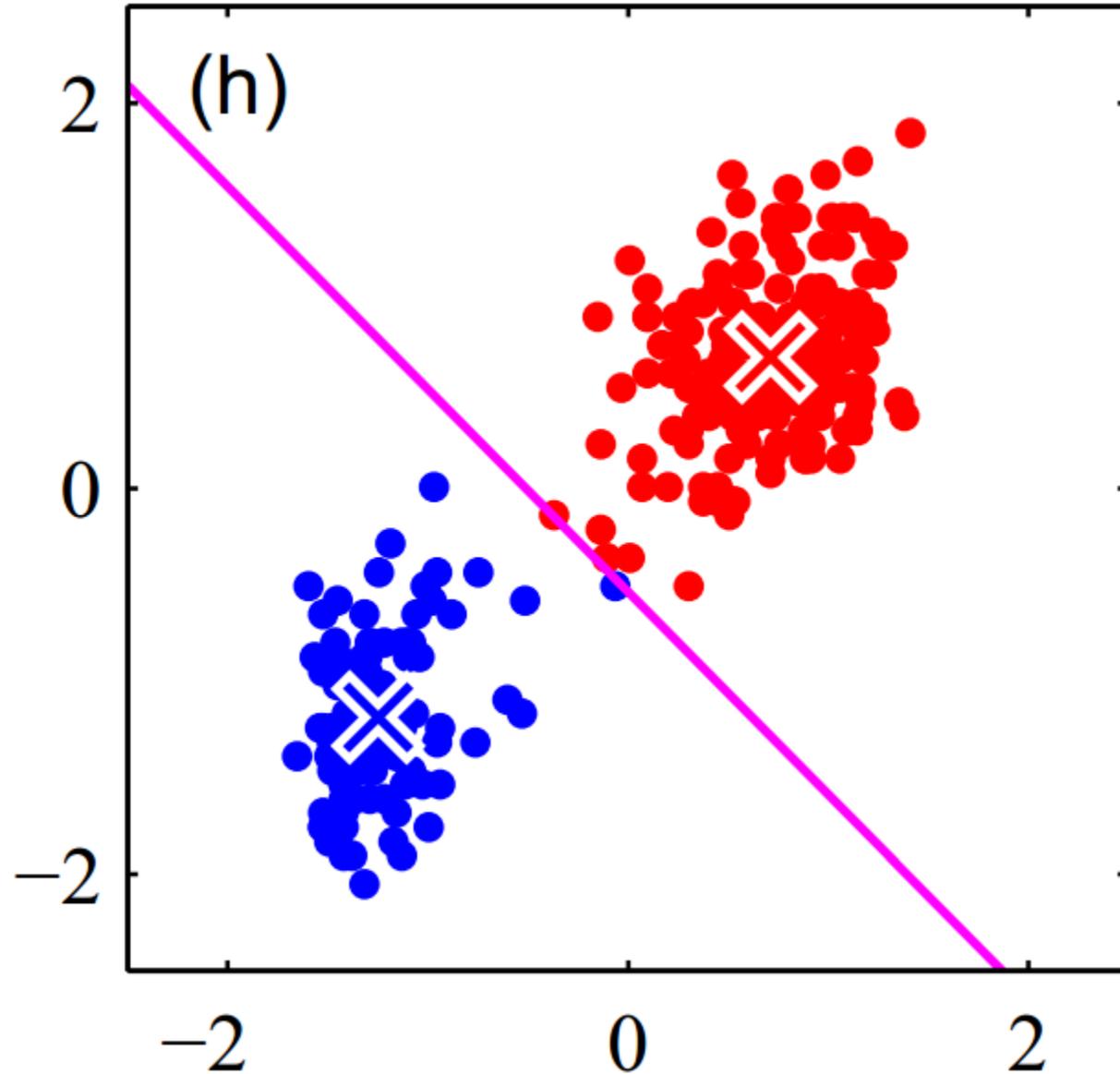
© unknown source. All rights reserved. This content is excluded from our Creative Commons license. For more information, see <http://ocw.mit.edu/help/faq-fair-use/>.



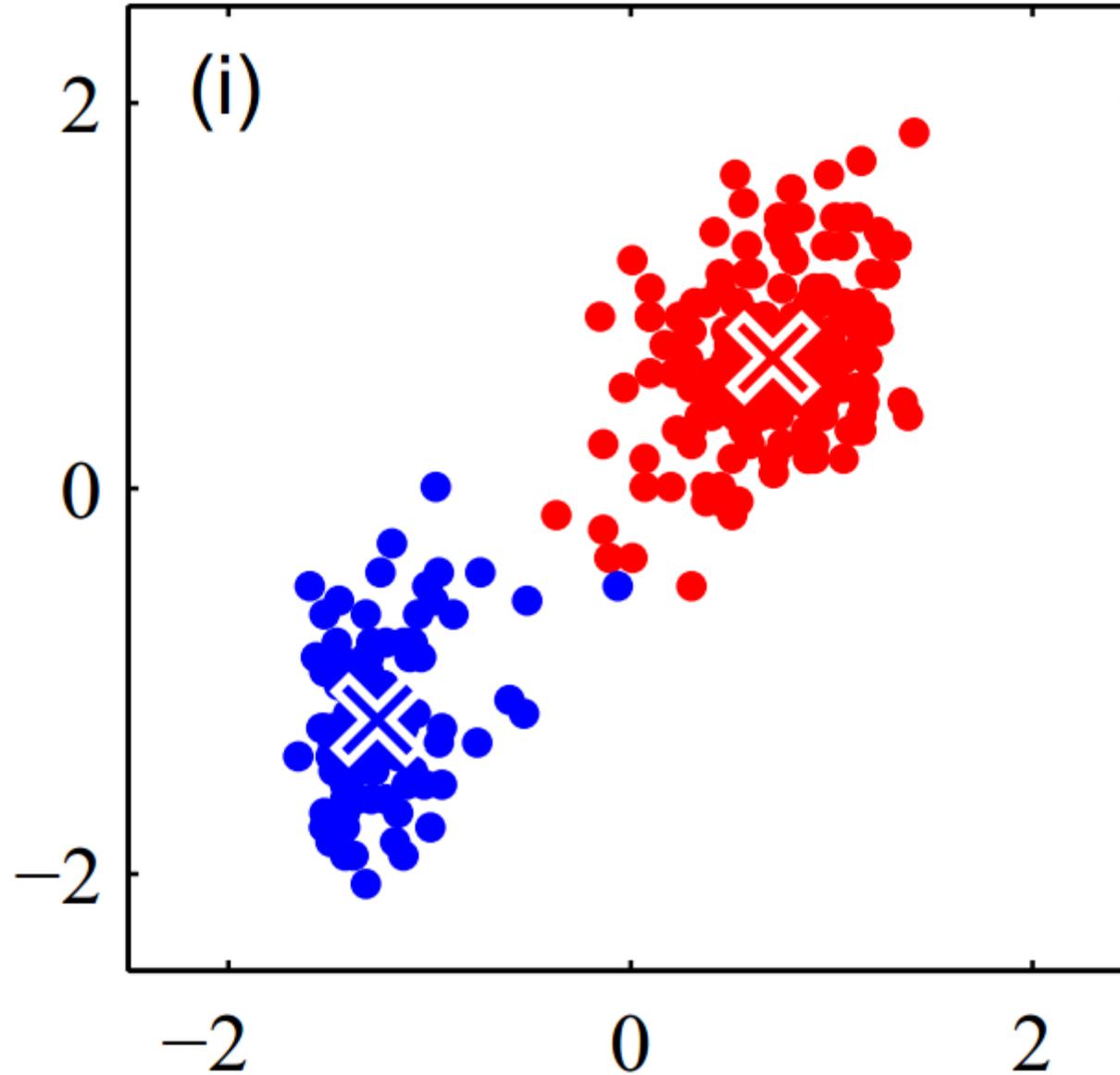
© unknown source. All rights reserved. This content is excluded from our Creative Commons license. For more information, see <http://ocw.mit.edu/help/faq-fair-use/>.



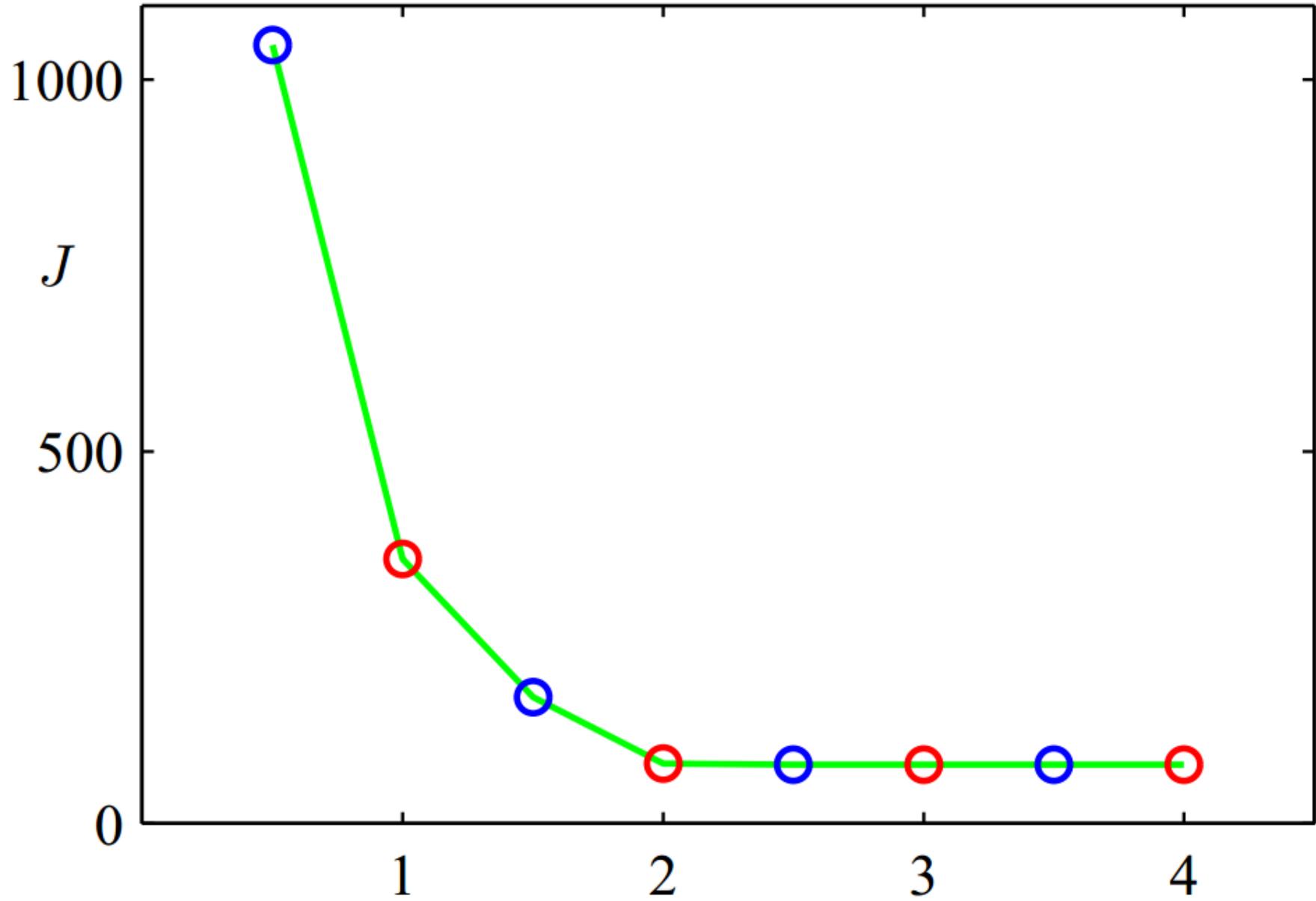
© unknown source. All rights reserved. This content is excluded from our Creative Commons license. For more information, see <http://ocw.mit.edu/help/faq-fair-use/>.



© unknown source. All rights reserved. This content is excluded from our Creative Commons license. For more information, see <http://ocw.mit.edu/help/faq-fair-use/>.



© unknown source. All rights reserved. This content is excluded from our Creative Commons license. For more information, see <http://ocw.mit.edu/help/faq-fair-use/>.



© unknown source. All rights reserved. This content is excluded from our Creative Commons license. For more information, see <http://ocw.mit.edu/help/faq-fair-use/>.

**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	$d_{fg} = (e_f - e_g)' \Sigma^{-1} (e_f - e_g)$ , where $\Sigma$ is the (full or within-cluster) covariance matrix of the data
Pearson correlation (centered correlation)	$d_{fg} = 1 - r_{fg}$ , 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)	$d_{fg} = 1 - r_{fg}$ , 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} $ or $d_{fg} = 1 - r_{fg}^2$

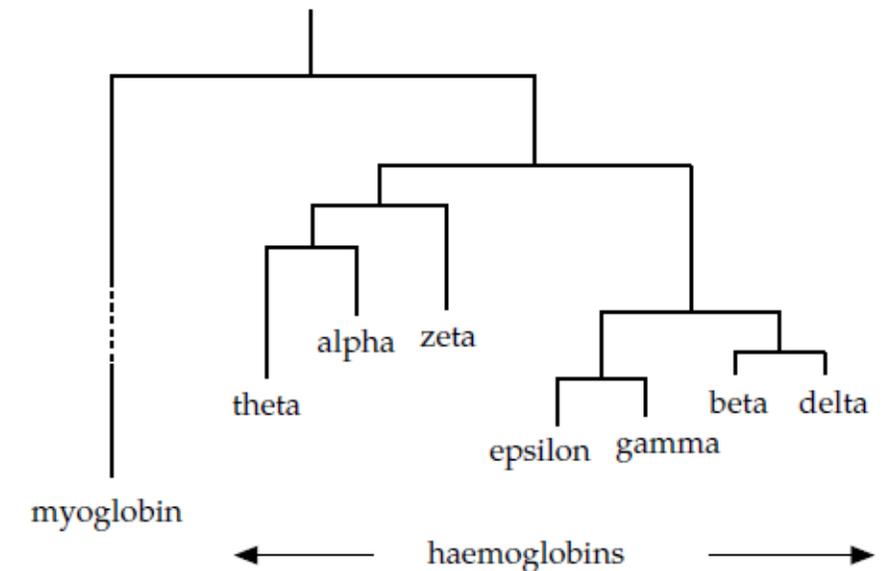
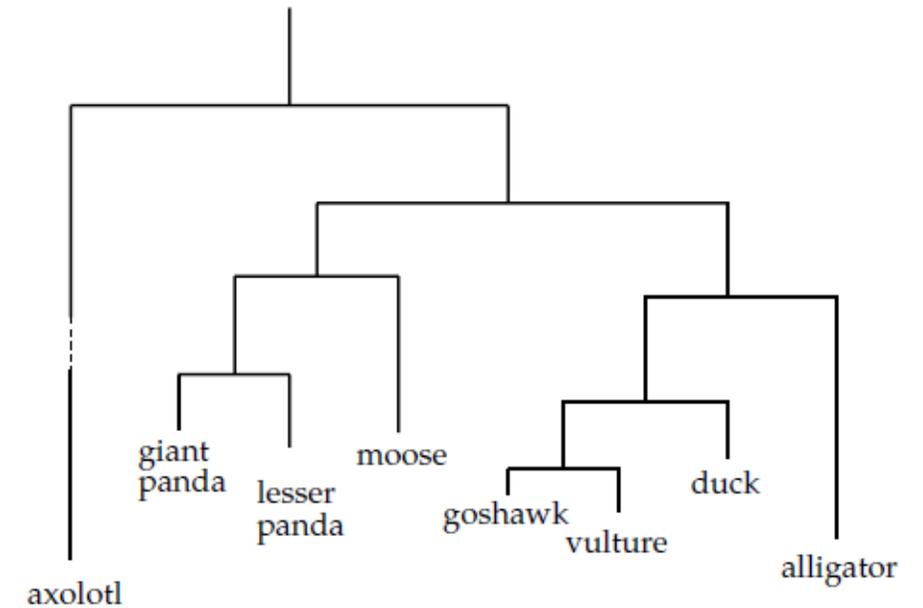
$d_{fg}$ , distance between expression patterns for genes  $f$  and  $g$ .  $e_{gc}$ , expression level of gene  $g$  under condition  $c$ .

Courtesy of Macmillan Publishers Limited. Used with permission.  
Source: D'haeseleer, Patrik. "How Does Gene Expression Clustering Work?." *Nature Biotechnology* 23, no. 12 (2005): 1499-1502.

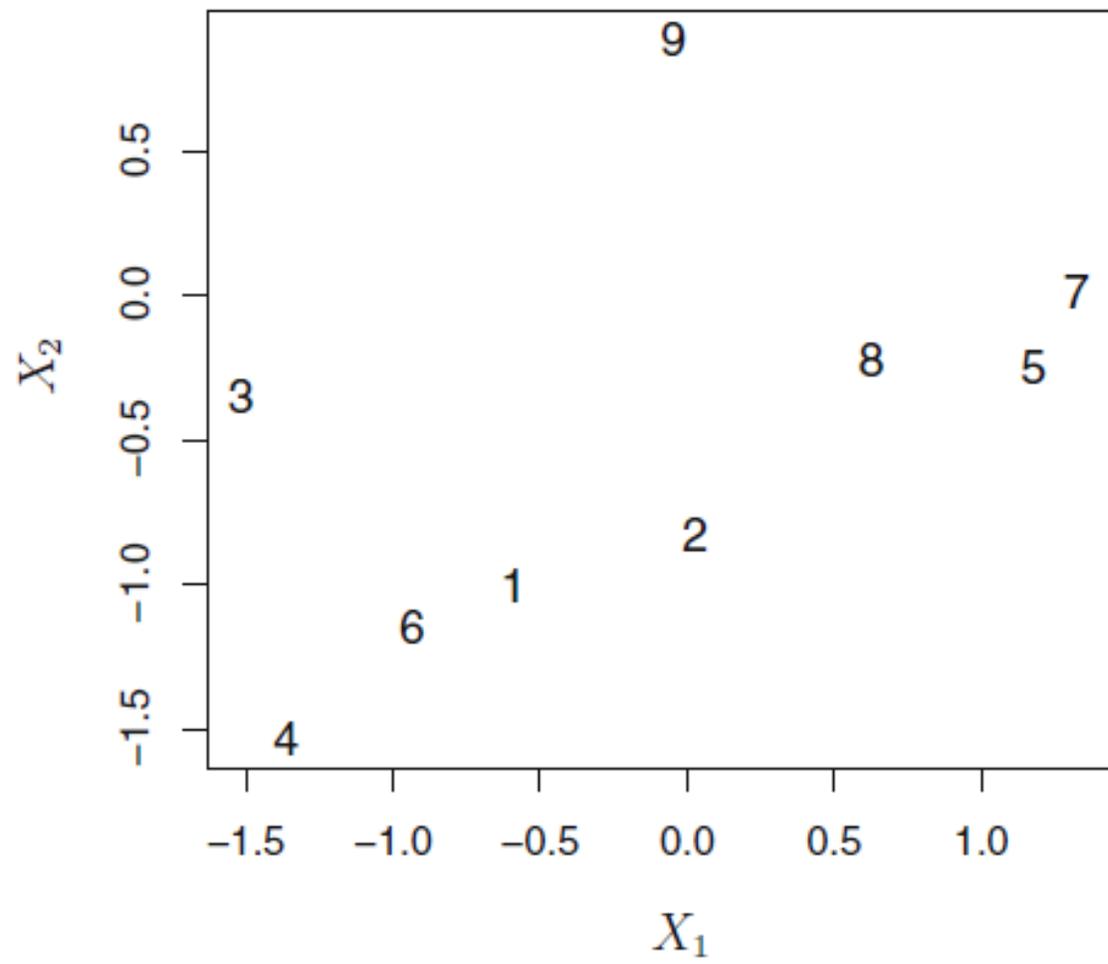
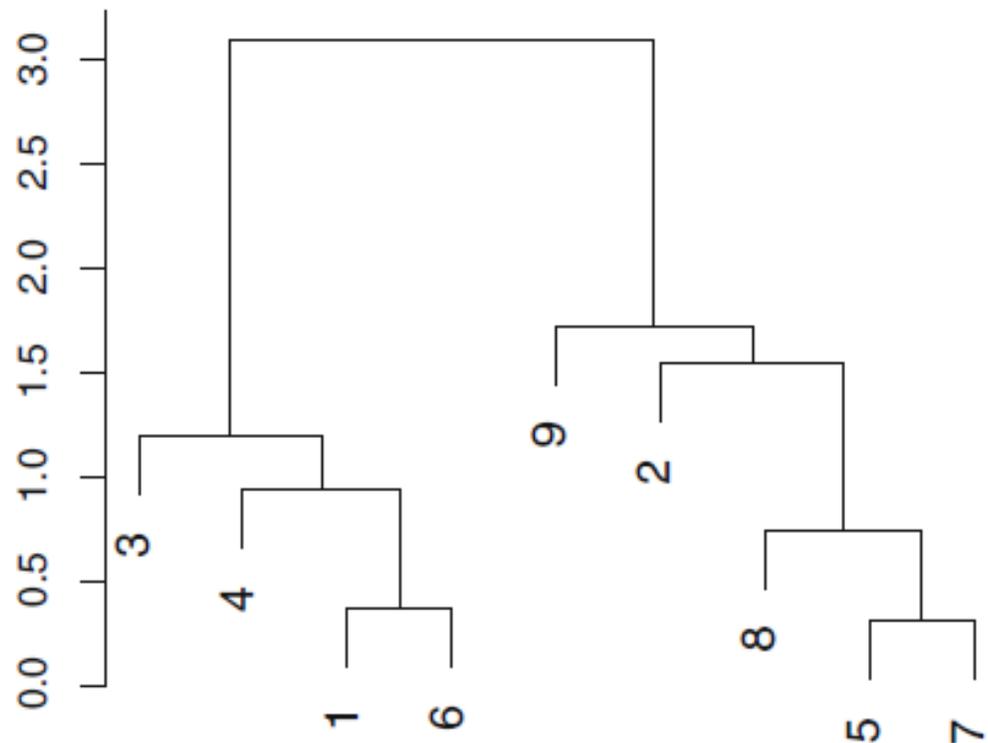
**How does gene expression clustering work?**  
Patrik D'haeseleer  
*Nature Biotechnology* 23, 1499 - 1501 (2005)  
doi:10.1038/nbt1205-1499

# Hierarchical clustering

- Organize data in a tree
  - Leaves are individual genes/species
  - Path lengths between leaves are distances
  - Similar points should lie in same lower subtrees
- Used to reveal evolutionary history of sequences



© unknown source. All rights reserved. This content is excluded from our Creative Commons license. For more information, see <http://ocw.mit.edu/help/faq-fair-use/>.



# Hierarchical clustering

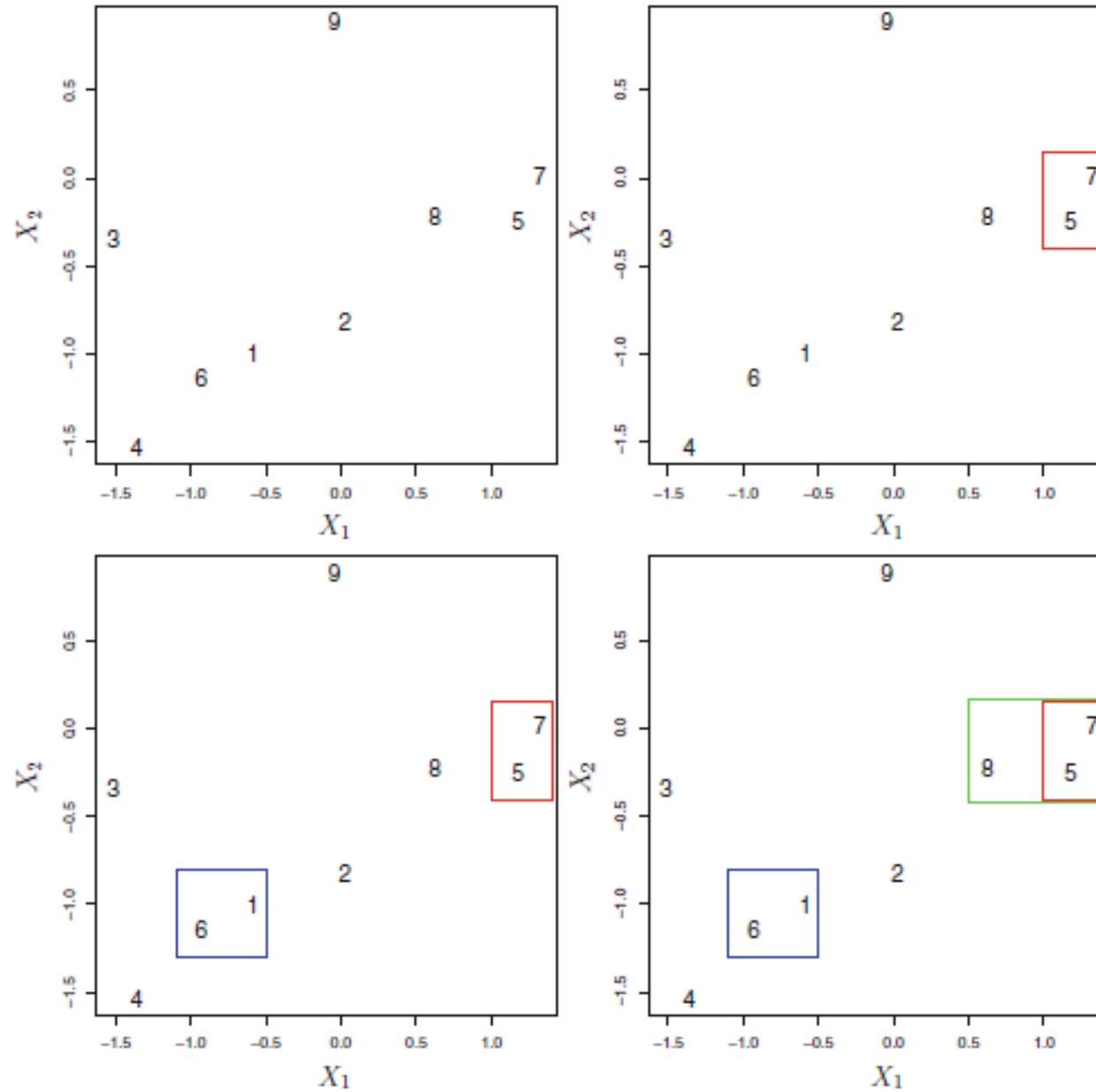
- What dissimilarity measure should be used?
  - To compare individual points
- What type of linkage should be used?
  - To compare clusters with each other
- Where do we cut dendrogram to obtain clusters?

---

## Algorithm 10.2 *Hierarchical Clustering*

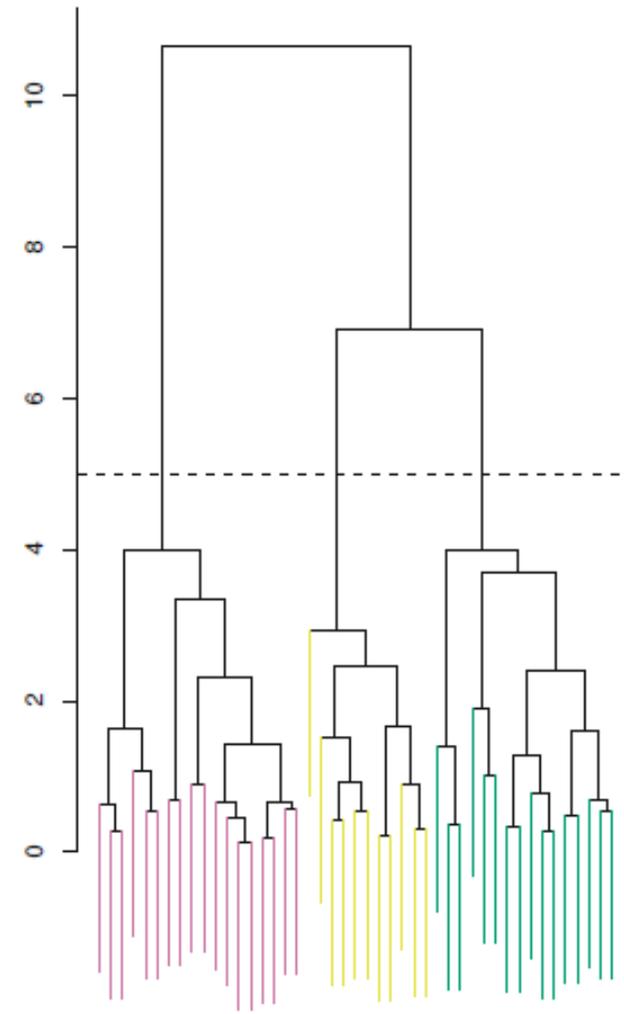
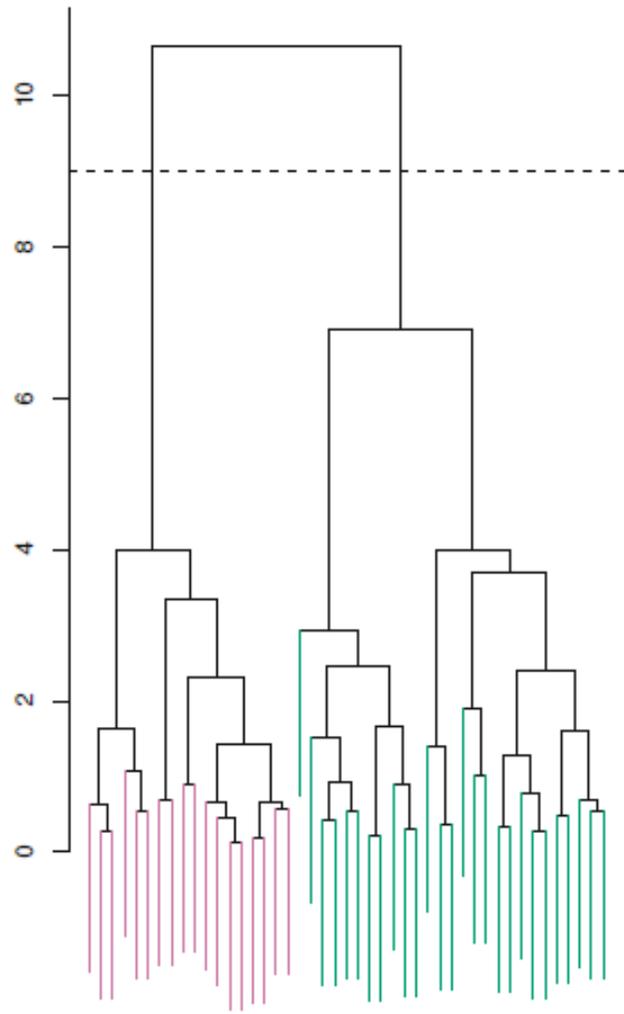
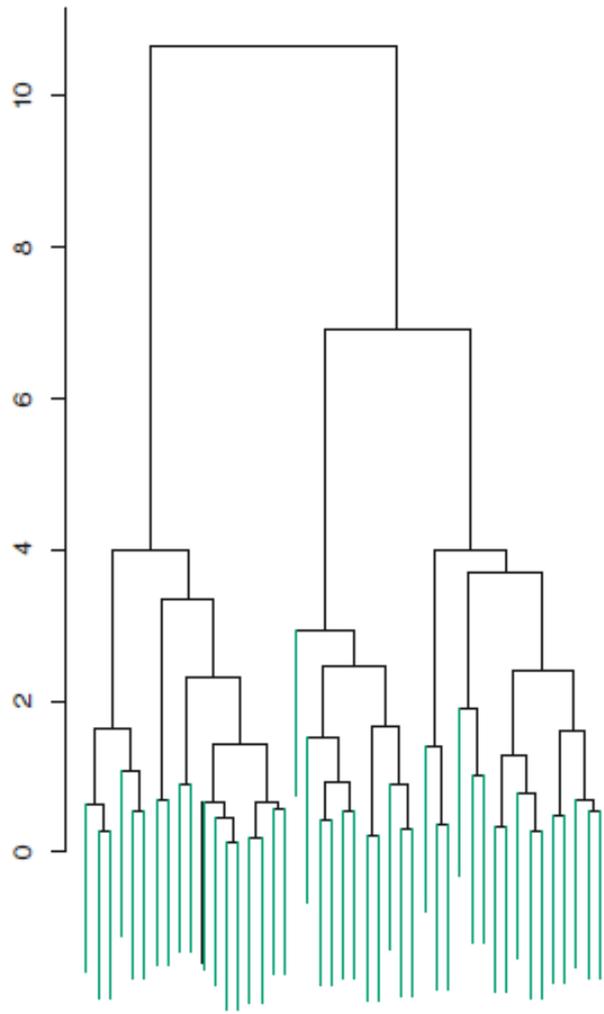
---

1. Begin with  $n$  observations and a measure (such as Euclidean distance) of all the  $\binom{n}{2} = n(n-1)/2$  pairwise dissimilarities. Treat each observation as its own cluster.
  
2. For  $i = n, n-1, \dots, 2$ :
  - (a) Examine all pairwise inter-cluster dissimilarities among the  $i$  clusters and identify the pair of clusters that are least dissimilar (that is, most similar). Fuse these two clusters. The dissimilarity between these two clusters indicates the height in the dendrogram at which the fusion should be placed.
  - (b) Compute the new pairwise inter-cluster dissimilarities among the  $i-1$  remaining clusters.



© unknown source. All rights reserved. This content is excluded from our Creative Commons license. For more information, see <http://ocw.mit.edu/help/faq-fair-use/>.

Complete	$d(A, B) = \max_{x \in A, y \in B} d(x, y)$	Maximal intercluster dissimilarity
Single	$d(A, B) = \min_{x \in A, y \in B} d(x, y)$	Minimal intercluster dissimilarity
Average	$d(A, B) = \sum_{x \in A, y \in B} \frac{d(x, y)}{ A  B }$	Average intercluster dissimilarity (UPGMA)
Centroid	$d(A, B) = d\left(\frac{\sum_{x \in A} \frac{x}{ A }, \sum_{y \in B} \frac{y}{ B }}{\quad}\right)$	Dissimilarity between centroids of each cluster



© unknown source. All rights reserved. This content is excluded from our Creative Commons license. For more information, see <http://ocw.mit.edu/help/faq-fair-use/>.

# Model selection

- Hierarchical clustering – cutoff tree at certain point
- K-means – how to choose K?
  - Balance number of clusters (# of parameters –  $K=n$  is uninformative) and the variance of the clusters
  - BIC Score – general model selection criterion
  - $BIC = -2 \times \text{loglikelihood} + d \times \log(N)$
  - Can use to decide whether to split a cluster
  - Computer BIC score of cluster and two potential child clusters – if BIC score is lower after split, do not accept split

# Biclustering

- Simultaneous clustering of rows and columns of a matrix
- Biclusters – subset of rows which exhibit similar behavior across a subset of columns, or vice versa

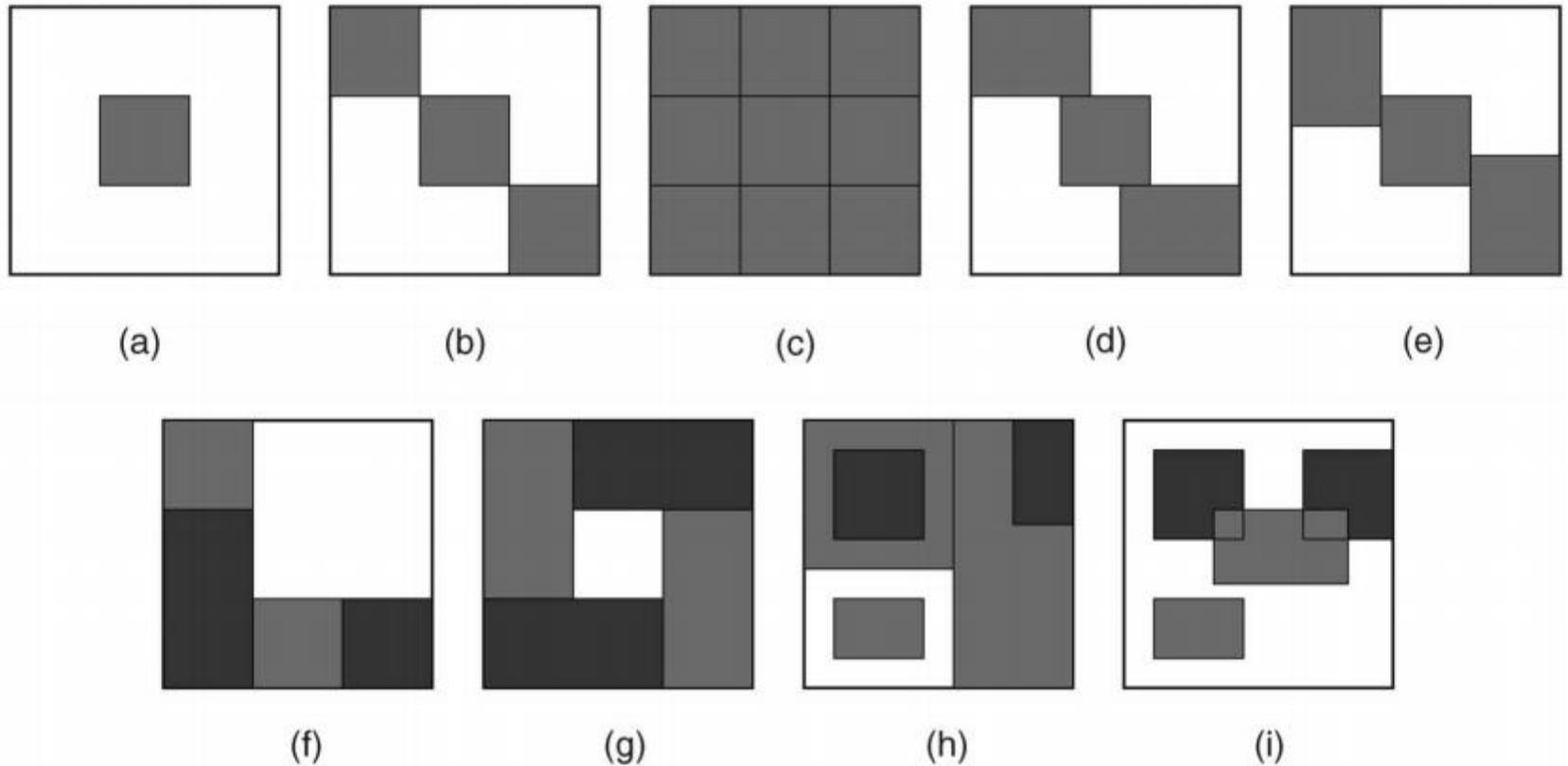
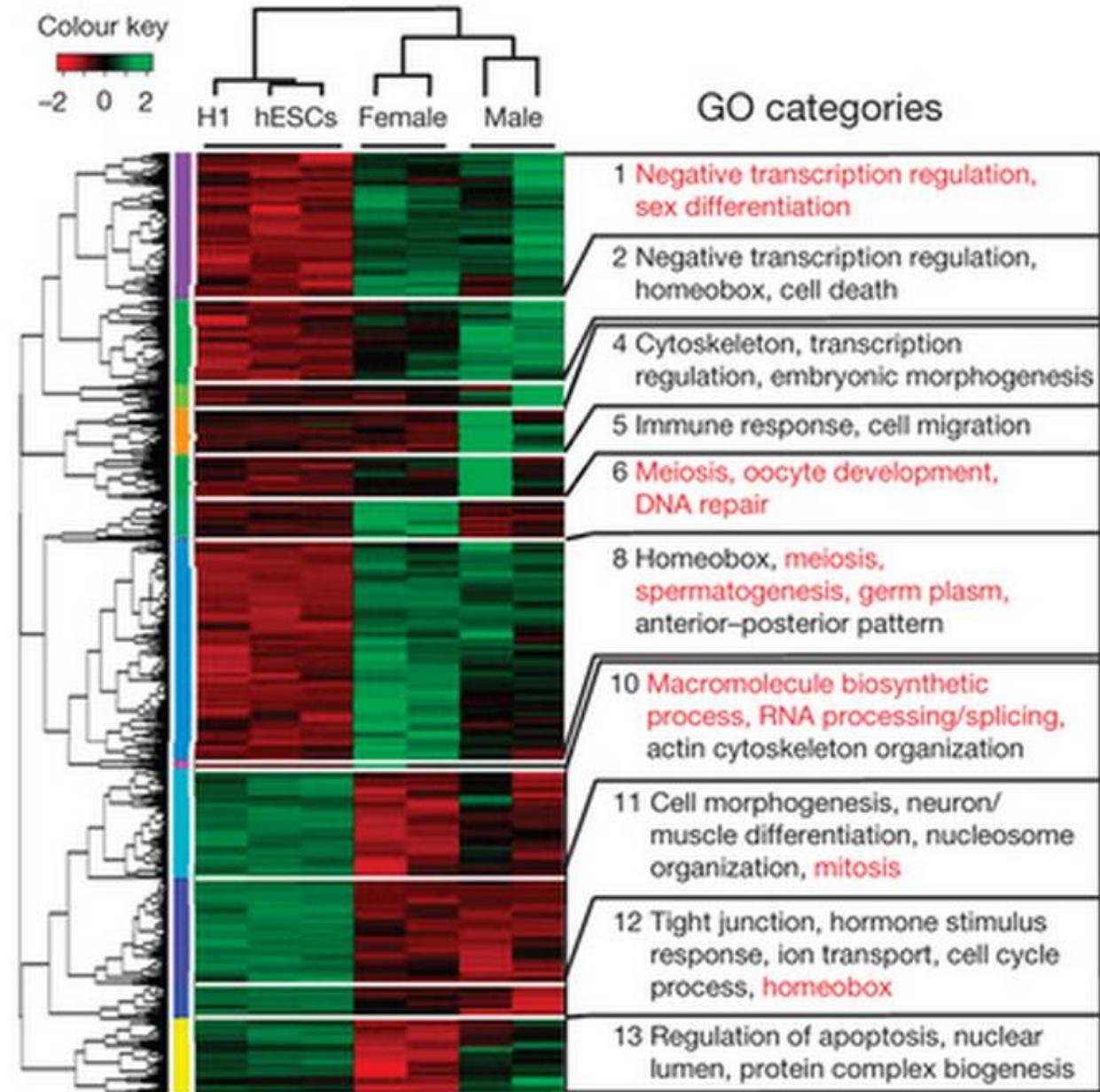


Fig. 4. Bicluster structure. (a) Single bicluster, (b) exclusive row and column biclusters, (c) checkerboard structure, (d) exclusive rows biclusters, (e) exclusive columns biclusters, (f) nonoverlapping biclusters with tree structure, (g) nonoverlapping nonexclusive biclusters, (h) overlapping biclusters with hierarchical structure, and (i) arbitrarily positioned overlapping biclusters.

© IEEE. All rights reserved. This content is excluded from our Creative Commons license. For more information, see <http://ocw.mit.edu/help/faq-fair-use/>.  
 Source: Madeira, Sara C., and Arlindo L. Oliveira. "Biclustering Algorithms for Biological Data Analysis: A Survey." *Computational Biology and Bioinformatics*, IEEE/ACM Transactions on 1, no. 1 (2004): 24-45.



Courtesy of Macmillan Publishers Limited. Used with permission.

Source: Gkoutela, Sofia, Ziwei Li, et al. "The Ontogeny of cKIT+ Human Primordial Germ Cells Proves to be a Resource for Human Germ Line Reprogramming, Imprint Erasure and in Vitro Differentiation." *Nature Cell Biology* 15, no. 1 (2013): 113-22.

# Biology Review

# Selection

- Negative selection (purifying/natural selection) – removal of deleterious traits
- Positive selection – increases prevalence of adaptive traits
- Thinking about selection happening at different levels
  - *Protein level: Sequence -> Structure -> Function*
  - RNA level: splicing, degradation/processing (NMD)
  - DNA level: DNA-protein binding sites

# Synonymous/Non-synonymous mutations

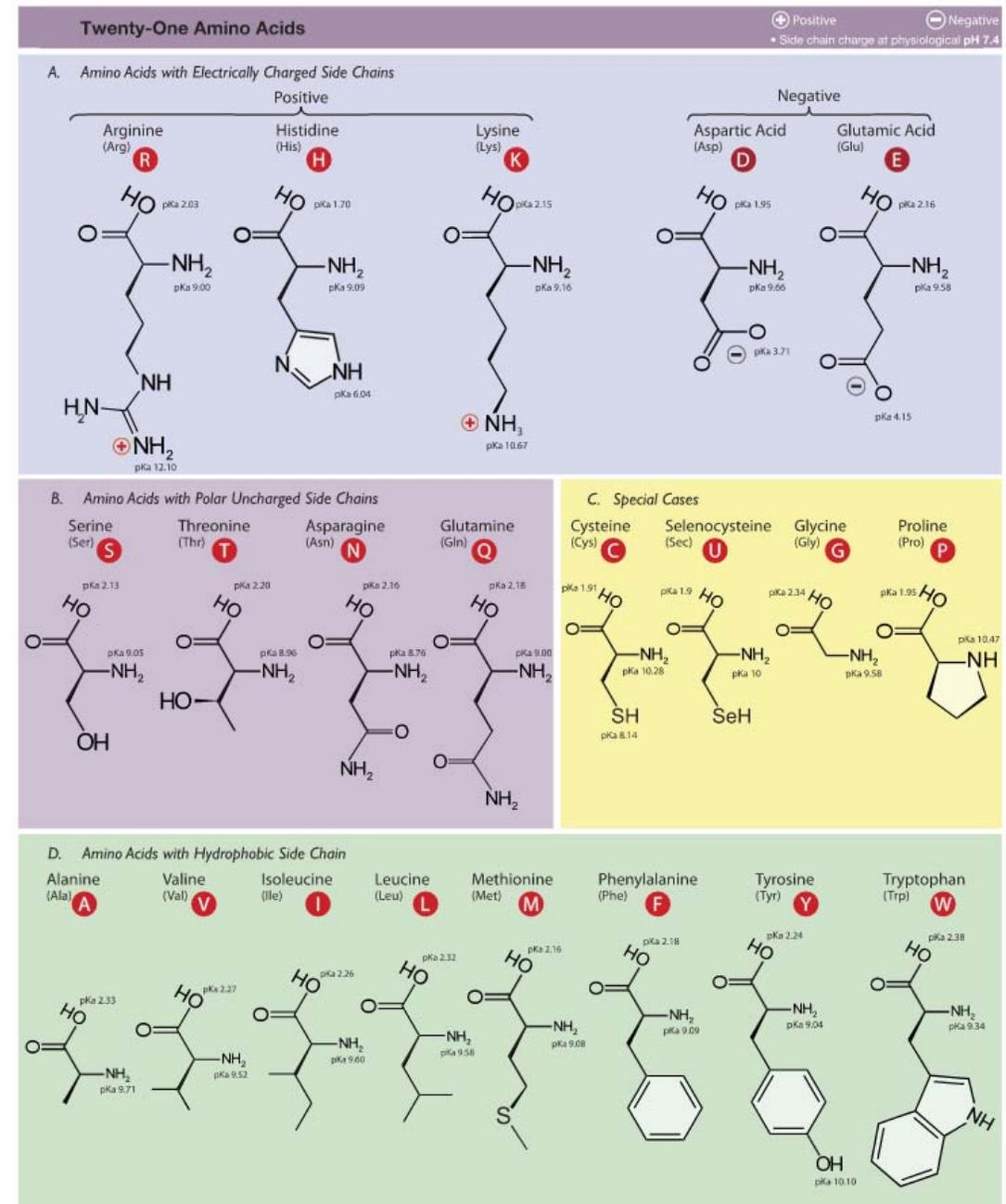
- Redundancy built into the genetic code
- Synonymous – one base changes for another in an exon, but the resulting amino acid sequence is unchanged
- Non-synonymous – new AA
- Can affect splicing, mRNA processing - so may not be silent

		Second Position									
		U		C		A		G			
First Position (5' end)	U	UUU UUC UUA UUG	Phe  Leu	UCU UCC UCA UCG	Ser	UAU UAC UAA UAG	Tyr  Stop Stop	UGU UGC UGA UGG	Cys  Stop Trp	U C A G	
	C	CUU CUC CUA CUG	Leu	CCU CCC CCA CCG	Pro	CAU CAC CAA CAG	His  Gln	CGU CGC CGA CGG	Arg	U C A G	
	A	AUU AUC AUA AUG	Ile  Met	ACU ACC ACA ACG	Thr	AAU AAC AAA AAG	Asn  Lys	AGU AGC AGA AGG	Ser  Arg	U C A G	
	G	GUU GUC GUA GUG	Val	GCU GCC GCA GCG	Ala	GAU GAC GAA GAG	Asp  Glu	GGU GGC GGA GGG	Gly	U C A G	
										Third Position (3' end)	

Image by MIT OpenCourseWare.

# Side-chain biochemistry

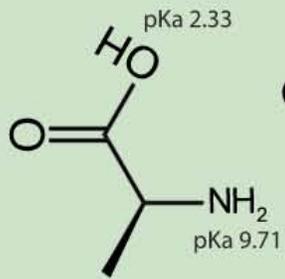
- Amino acids classified by properties of side chains
  - Grouped by general properties
- Substitutions of amino acid with another of similar chemical properties may conserve protein function



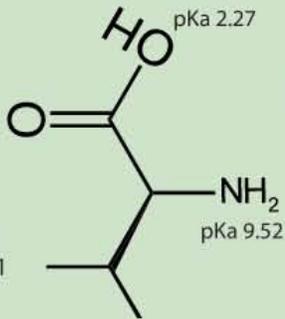
# Side chain size (Trp – W)

## D. Amino Acids with Hydrophobic Side Chain

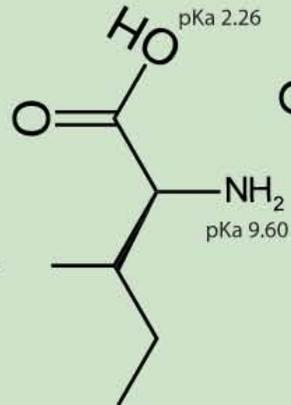
Alanine  
(Ala) **A**



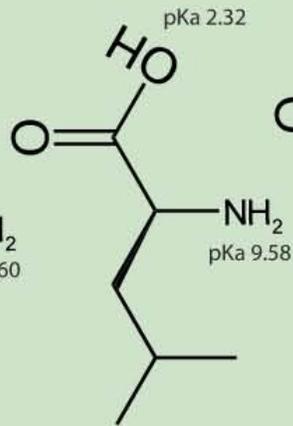
Valine  
(Val) **V**



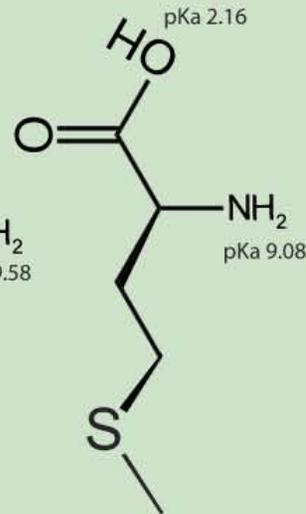
Isoleucine  
(Ile) **I**



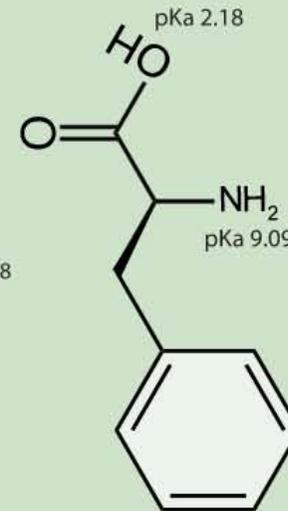
Leucine  
(Leu) **L**



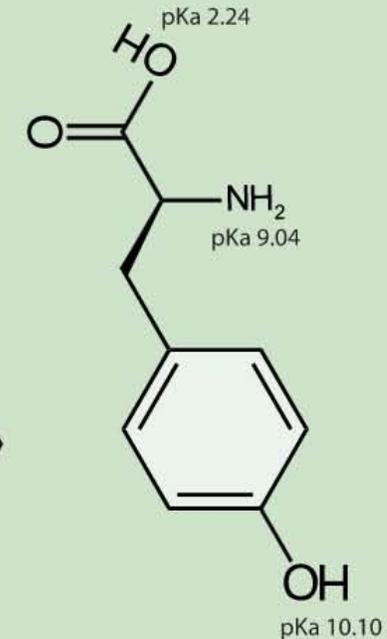
Methionine  
(Met) **M**



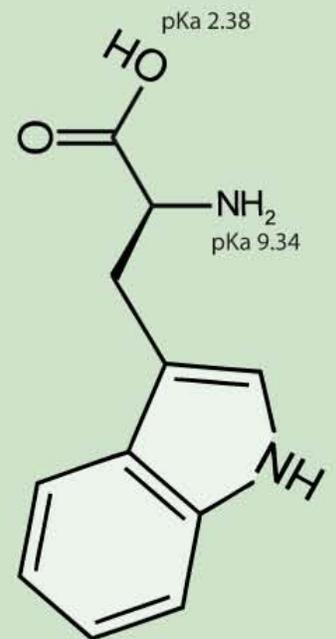
Phenylalanine  
(Phe) **F**



Tyrosine  
(Tyr) **Y**

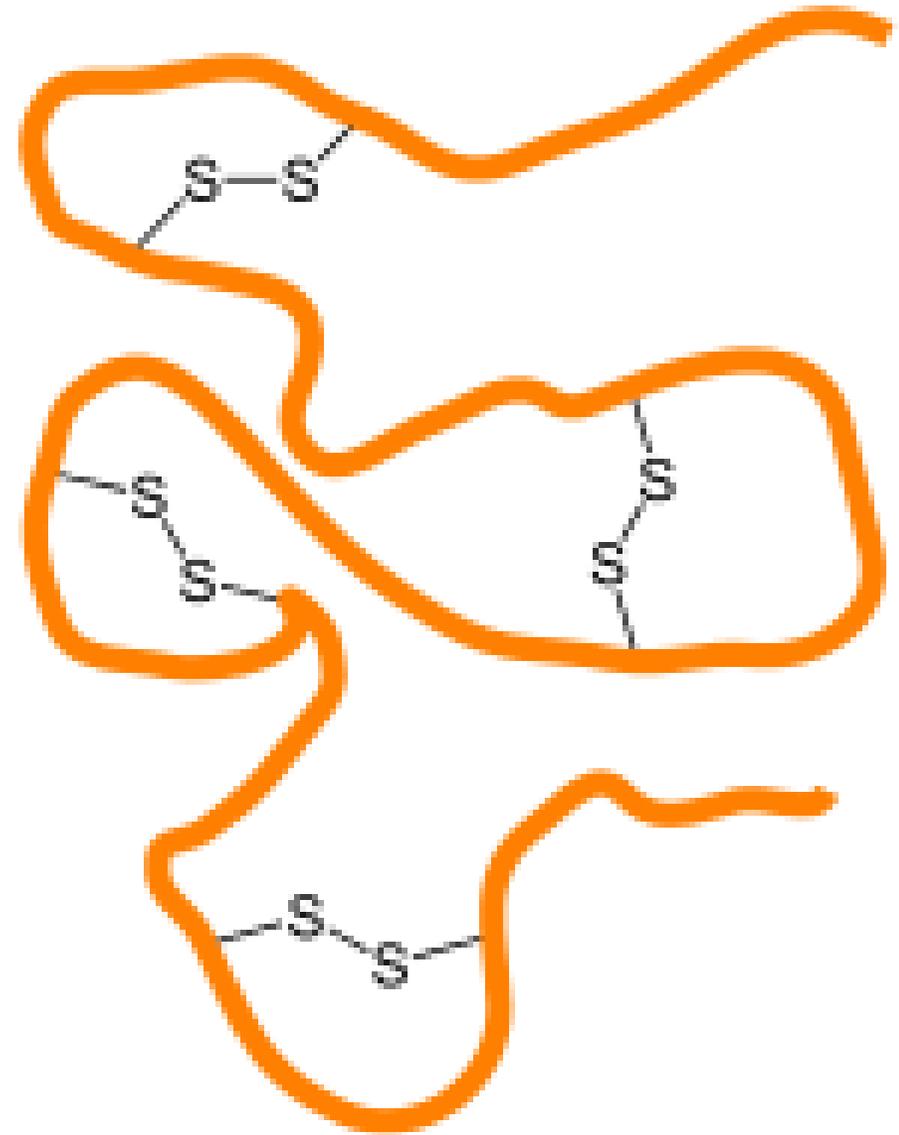
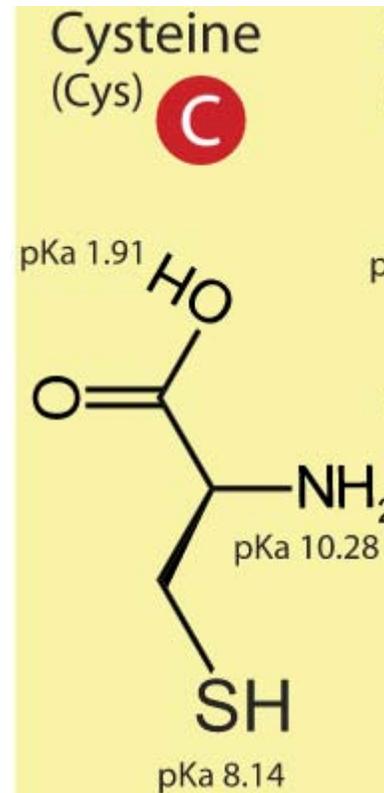


Tryptophan  
(Trp) **W**



# Disulfide bond

- Important in protein folding – holds two distant portions of protein together
- Occurs between Cys residues



# Next-generation sequencing

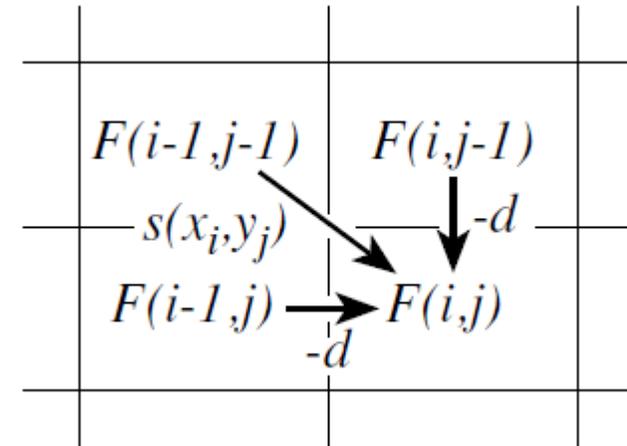
- Sequencing is always of DNA
  - Need to convert RNA to DNA by *reverse transcription* (RT)
- Illumina is current leader in the field
  - 8 lanes on a flow cell
  - Each lane can sequence 200 million 100bp reads – 20 Gbps!
  - Can sequence multiple samples per lane by barcoding
  - Requires (heterogeneous) population of cells to get enough DNA for sample
- Single cell sequencing applications are becoming more common (RNAseq)
- Single molecule technologies are still being developed – PacBio

# Alignment

# Alignment

$$F(i, j) = \max \begin{cases} F(i-1, j-1) + s(x_i, y_j), \\ F(i-1, j) - d, \\ F(i, j-1) - d. \end{cases}$$

$$F(i, j) = \max \begin{cases} 0, \\ F(i-1, j-1) + s(x_i, y_j), \\ F(i-1, j) - d, \\ F(i, j-1) - d. \end{cases}$$



# Local alignment example

Do a local alignment between these using PAM250 and gap penalty -2:

**AWEK**

**FWEF**

	C	S	T	F	A	G	N	D	E	Q	H	R	K	M	I	L	V	F	Y	W	
C	12																				C
S	0	2																			S
T	-2	1	3																		T
F	-3	1	0	6																	F
A	-2	1	1	1	2																A
G	-3	1	0	-1	1	5															G
N	-4	1	0	-1	0	0	2														N
D	-5	0	0	-1	0	1	2	4													D
E	-5	0	0	-1	0	0	1	3	4												E
Q	-5	-1	-1	0	0	-1	1	2	2	4											Q
H	-3	-1	-1	0	-1	-2	2	1	1	3	6										H
R	-4	0	-1	0	-2	-3	0	-1	-1	1	2	6									R
K	-5	0	0	-1	-1	-2	1	0	0	1	0	3	5								K
M	-5	-2	-1	-2	-1	-3	-2	-3	-2	-1	-2	0	0	6							M
I	-2	-1	0	-2	-1	-3	-2	-2	-2	-2	-2	-2	-2	2	5						I
L	-6	-3	-2	-3	-2	-4	-3	-4	-3	-2	-2	-3	-3	4	2	6					L
V	-2	-1	0	-1	0	-1	-2	-2	-2	-2	-2	-2	-2	2	4	2	4				V
F	-4	-3	-3	-5	-4	-5	-4	-6	-5	-5	-2	-4	-5	0	1	2	-1	9			F
Y	0	-3	-3	-5	-3	-5	-2	-4	-4	-4	0	-4	-4	-2	-1	-1	-2	7	10		Y
W	-8	-2	-5	-6	-6	-7	-4	-7	-7	-5	-3	2	-3	-4	-5	-2	-6	0	0	17	W
	C	S	T	F	A	G	N	D	E	Q	H	R	K	M	I	L	V	F	Y	W	

# Local alignment solution

	Gap	A	W	E	K
Gap	0	0	0	0	-8
F	0	0	0	0	0
W	0	0	17	15	13
E	0	0	15	21	19
F	0	0	13	19	17

alignment:

W E  
W E

# Global alignment solution

	Gap	A	W	E	K
Gap	0	-2	-4	-6	-8
F	-2	-4	-2	-4	-6
W	-4	-6	13	11	9
E	-6	-4	-6	17	11
F	-8	-6	-4	15	12

alignment:            A    W    E    K  
                      F    W    E    F

	<b>Global</b>	<b>Semiglobal</b>	<b>Local (gapped)</b>
<b>Penalties at edges?</b>	Yes	No	No
<b>Reset to 0 instead of including negative entries?</b>	No	No	Yes
<b>End of alignment</b>	Bottom right entry	Highest score entry in bottom row or rightmost column	Highest score entry in matrix

# Reminders

- Pset 1 posted – due Feb 20<sup>th</sup> (no extra problem)
- Pset 2 posted – Due Mar 13<sup>th</sup>
- Project teams due – Feb 25<sup>th</sup>
  - Interests and background directory has been posted
- Lecture videos will be posted on MITx soon – next week?

MIT OpenCourseWare  
<http://ocw.mit.edu>

7.91J / 20.490J / 20.390J / 7.36J / 6.802 / 6.874 / HST.506 Foundations of Computational and Systems Biology  
Spring 2014

For information about citing these materials or our Terms of Use, visit: <http://ocw.mit.edu/terms>.