

Functional genomics + Data mining

BCH364C/394P Systems Biology / Bioinformatics
Edward Marcotte, Univ of Texas at Austin

Functional genomics

= field that attempts to use the vast data produced by genomic projects (e.g. genome sequencing projects) to describe gene (and protein) functions and interactions.

Focuses on dynamic aspects, e.g. transcription, translation, and protein–protein interactions, as opposed to static aspects of the genome such as DNA sequence or structures.

Adapted from Wikipedia

Functional genomics + Data mining

= field that attempts to computationally discover patterns in large data sets

Adapted from Wikipedia

Functional genomics + Data mining



www.sparkpeople.com

Adapted from Wikipedia

**We're going to first learn
about clustering algorithms
& classifiers**

**We're going to first learn
about clustering algorithms
& classifiers**

Clustering = task of grouping a set of objects in such a way that objects in the same group (a **cluster**) are more similar (in some sense) to each other than to those in other groups (clusters).

Adapted from Wikipedia

We're going to first learn about clustering algorithms & classifiers

Classification = task of categorizing a new observation,
on the basis of a training set of data with observations
(or instances) whose categories are known

Adapted from Wikipedia

Let's motivate this with an example:

Distinct types of diffuse large B-cell lymphoma identified by gene expression profiling

Ash A. Alizadeh^{1,2}, Michael B. Eisen^{2,3,4}, R. Eric Davis⁵, Chi Ma⁵, Izidore S. Lossos⁶, Andreas Rosenwald⁵, Jennifer C. Boldrick¹, Hajeer Sabet⁵, Truc Tran⁵, Xin Yu⁵, John I. Powell⁷, Liming Yang⁷, Gerald E. Marti⁸, Troy Moore⁹, James Hudson Jr⁹, Lisheng Lu¹⁰, David B. Lewis¹⁰, Robert Tibshirani¹¹, Gavin Sherlock⁴, Wing C. Chan¹², Timothy C. Greiner¹², Dennis D. Weisenburger¹², James O. Armitage¹³, Roger Warnke¹⁴, Ronald Levy², Wyndham Wilson¹⁵, Michael R. Grever¹⁶, John C. Byrd¹⁷, David Botstein⁴, Patrick O. Brown^{1,18} & Louis M. Staudt⁵

Nature 2000

“Diffuse large B-cell lymphoma (DLBCL), the most common subtype of non-Hodgkin's lymphoma ... is one disease in which attempts to define subgroups on the basis of morphology have largely failed...”

“DLBCL ... is clinically heterogeneous:
40% of patients respond well to current therapy and have prolonged survival, whereas the remainder succumb to the disease.

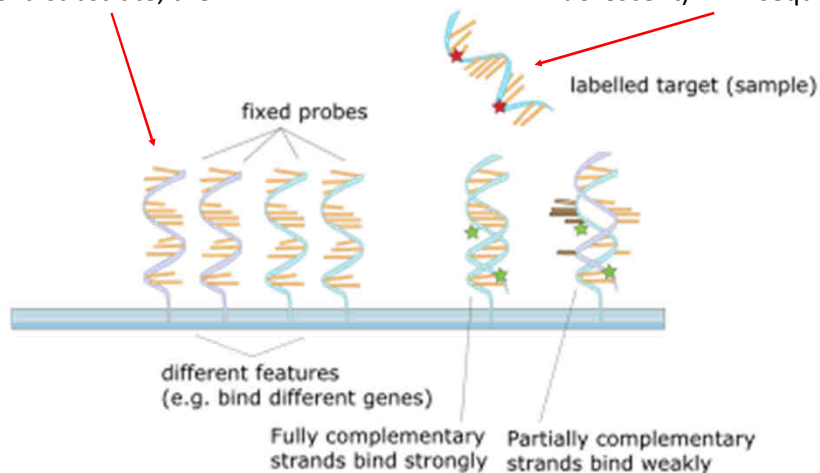
We proposed that this variability in natural history reflects unrecognized molecular heterogeneity in the tumours.”

Nature 2000

Blast from the past: Profiling mRNA expression with DNA microarrays

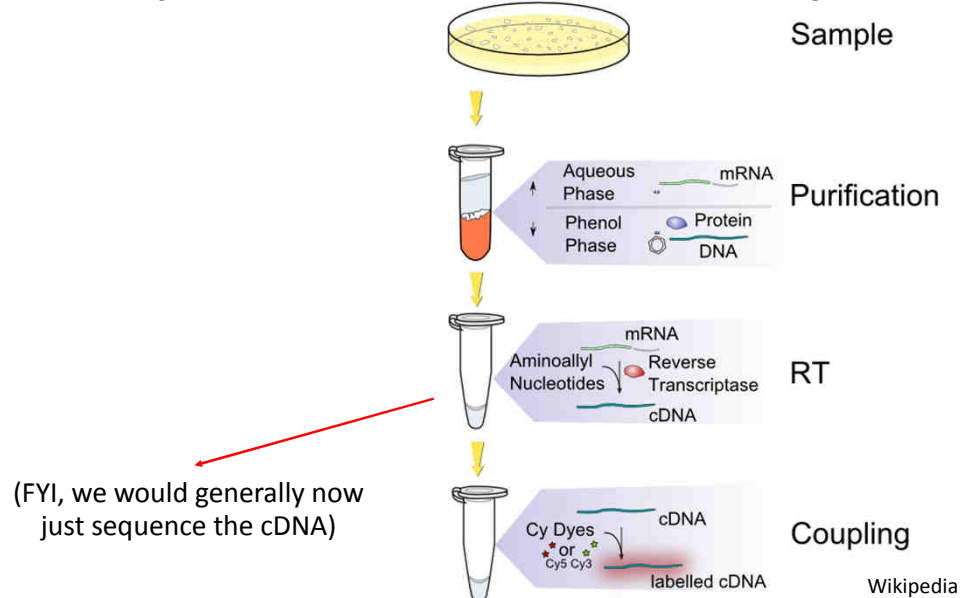
DNA molecules are attached to a solid substrate, then...

...probed with a labeled (usually fluorescent) DNA sequence

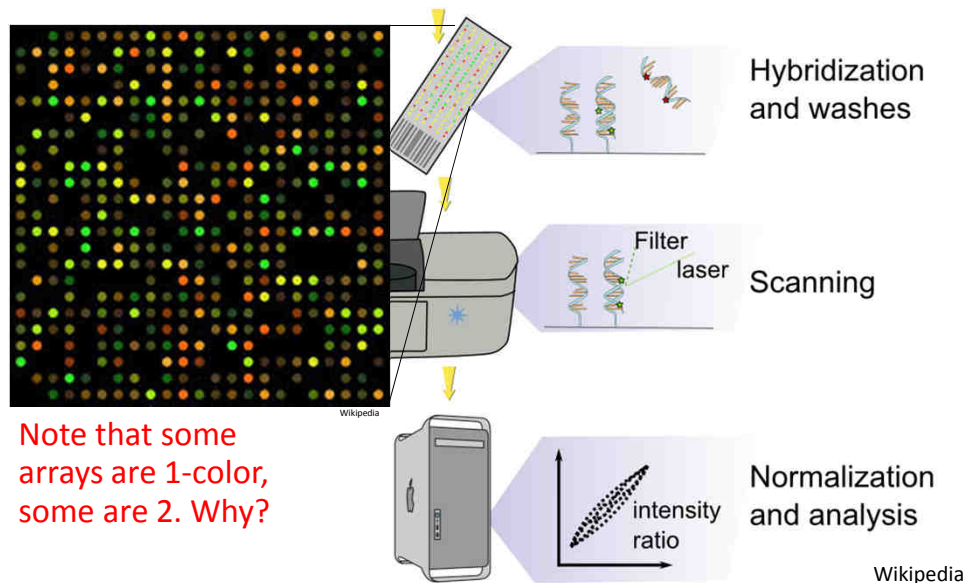


Wikipedia

Blast from the past: Profiling mRNA expression with DNA microarrays



Blast from the past: Profiling mRNA expression with DNA microarrays



Back to diffuse large B-cell lymphoma...

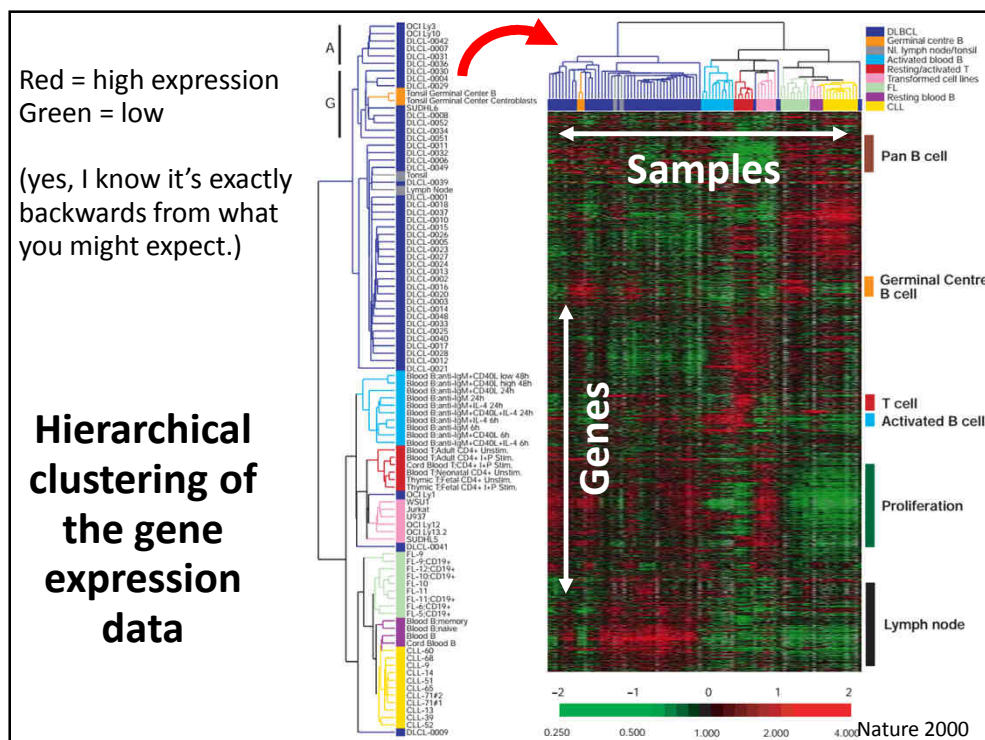
96 patient biopsies
(normal and malignant lymphocyte samples)

↓
Extract mRNA from each sample

↓
**Perform DNA microarray experiment on each to
measure mRNA abundances (~1.8 million total gene
expression measurements)**

↓
Cluster samples by their expression patterns

Nature 2000



Cell Type Legend:

- CLBLCL
- Germinal centre B
- 4 node/naïve
- Activated blood B
- Resting/activated T
- Cell Lites
- Resting blood B
- TLL

Proliferation

- Cyclin A
- BUB1 mitotic kinase
- CDC20 B1
- CDK5
- SKY
- SKP2
- CKK-polo-like kinase
- CDP2/CDK1/KAP1
- aurora-related kinase 1
- p16
- Thymidine kinase
- CDC22
- CDK21 homologue
- ROCK2
- Dihydrofolate reductase
- CD38

Germinal centre B

- FAK=focal adhesion kinase
- WIP=WASP interacting protein
- FMR2
- CD10
- BCL-7A
- A-myb
- BCL-6
- PI 3-kinase p110 γ

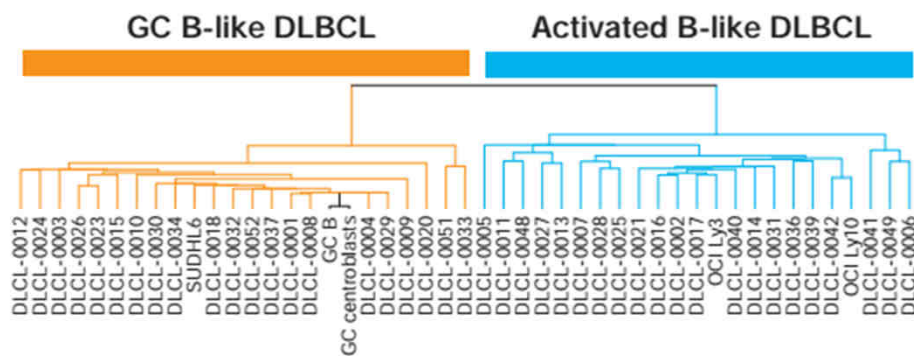
Lymph node

- RGS13
- CD105
- CD11
- FGF-7
- MMIP9
- Interleukin-5F-1 receptor
- Cathepsin B
- Fc ϵ receptor γ chain

T cell

- TIMP-3
- Integrin beta 5
- NK4=NK cell protein-4
- SDF-1 chemokine
- CD49F=Integrin $\alpha 6$
- CD3
- CD28
- cell receptor β chain
- ITP
- JRF-1
- Chemokine 10

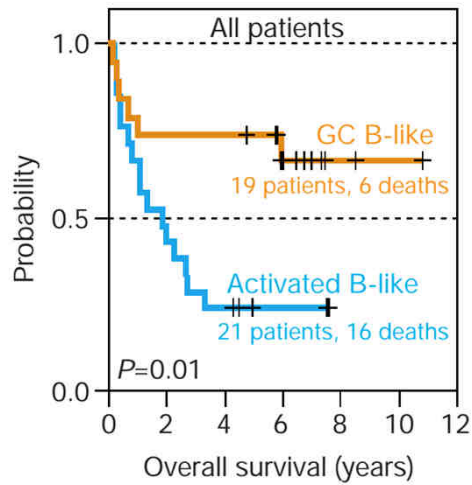
We can break up the DLBCL's according the germinal B-cell specific gene expression:



8

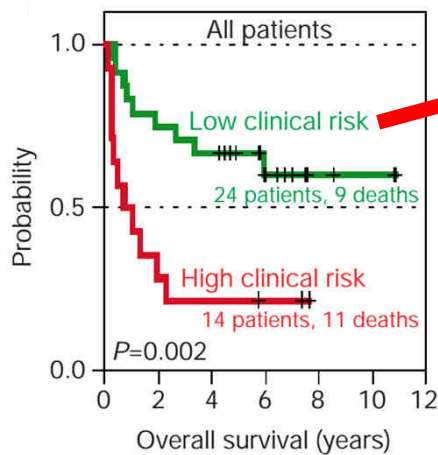
What good is this? These molecular phenotypes predict clinical survival.

Kaplan-Meier plot
of patient survival

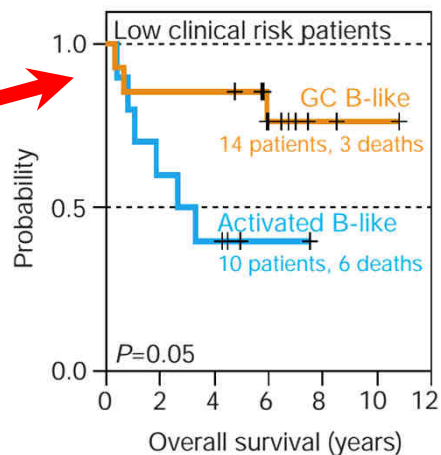


Nature 2000

What good is this? These molecular phenotypes predict clinical survival.



Grouping patients by clinical prognostic index



Regrouping low risk patients by gene expression

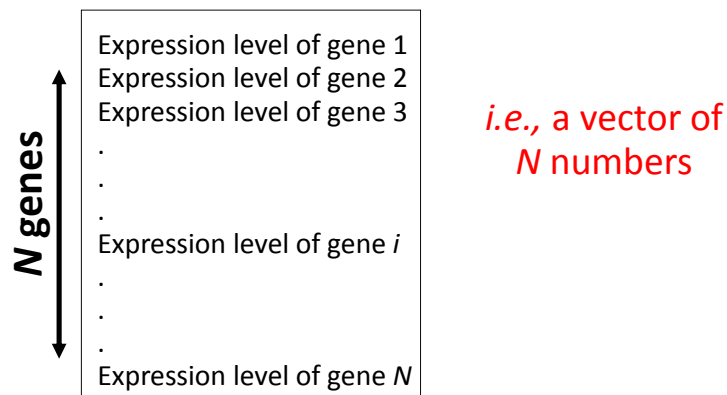
Nature 2000

Gene expression, and other molecular measurements, provide far deeper phenotypes for cells, tissues, and organisms than traditional measurements

These sorts of observations have now motivated tons of work using these approaches to diagnose specific forms of disease, as well as to discover functions of genes and many other applications

So, how does clustering work?

First, let's think about the data, e.g. as for gene expression.
From one sample, using DNA microarrays or RNA-seq, we get:



For yeast, $N \sim 6,000$
For human, $N \sim 22,000$

So, how does clustering work?

Every additional sample adds another column, giving us a matrix of data:

or data:

	M samples				
N genes	Gene 1, sample 1	...	Gene 1, sample j	...	Gene 1, sample M
	Gene 2, sample 1	...	Gene 2, sample j	...	Gene 2, sample M
	Gene 3, sample 1	...	Gene 3, sample j	...	Gene 3, sample M
	.		.		.
	.		.		.
	.		.		.
	Gene i , sample 1	...	Gene i , sample j	...	Gene i , sample M
	.		.		.
.		.		.	
.		.		.	
Gene N , sample 1	...	Gene N , sample j	...	Gene N , sample M	

For yeast, $N \sim 6,000$
For human, $N \sim 22,000$

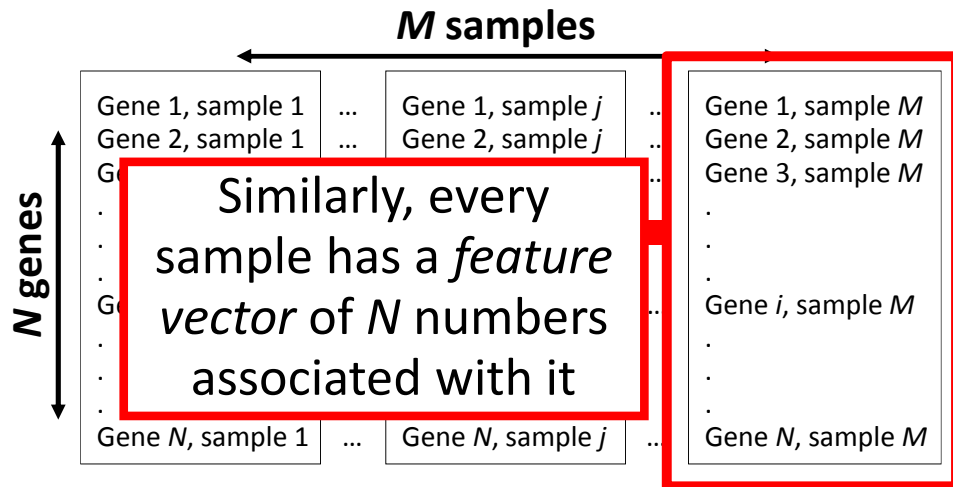
i.e., a matrix of N
 $\times M$ numbers

So, how does clustering work?

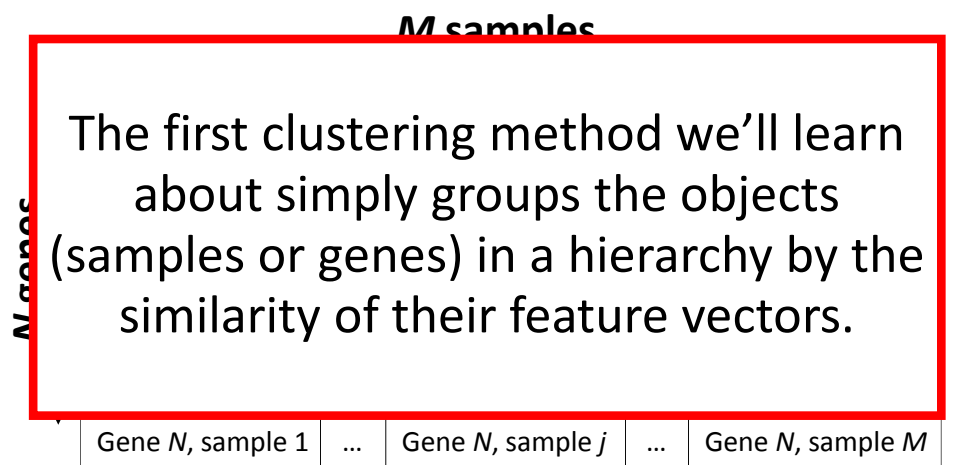
	M samples				
N genes	Gene 1, sample 1	...	Gene 1, sample j	...	Gene 1, sample M
	Gene 2, sample 1	...	Gene 2, sample j	...	Gene 2, sample M
	Gene 3, sample 1	...	Gene 3, sample j	...	Gene 3, sample M
	.		.		.
	.		.		.
	Gene i , sample 1	...	Gene i , sample j	...	Gene i , sample M
	.		.		.
	Gene N , sample 1	...	Gene N , sample j	...	Gene N , sample M

Every gene has a *feature vector*
of M numbers associated with it

So, how does clustering work?



So, how does clustering work?



A hierarchical clustering algorithm

Start with each object in its own cluster

Until there is only one cluster left, repeat:

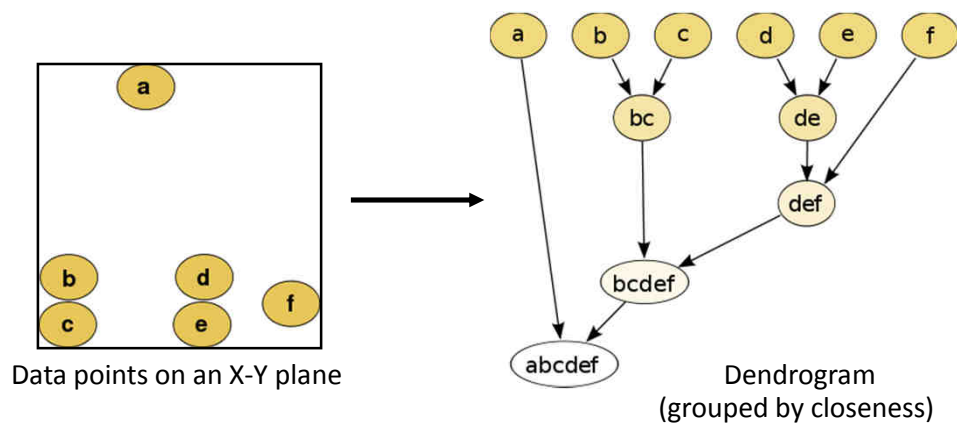
Among the current clusters, find the two most similar clusters

Merge those two clusters into one

We can choose our measure of similarity and how we merge the clusters

Hierarchical clustering

Conceptually



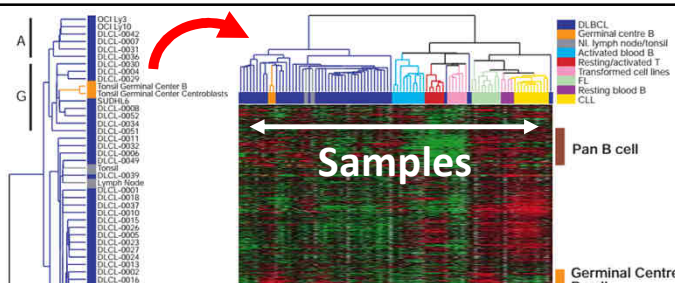
Wikipedia

We'll need to measure the similarity between feature vectors. Here are a few (of many) common distance measures used in clustering.

Names	Formula
Euclidean distance	$\ a - b\ _2 = \sqrt{\sum_i (a_i - b_i)^2}$
Manhattan distance	$\ a - b\ _1 = \sum_i a_i - b_i $
cosine similarity	$\frac{a \cdot b}{\ a\ \ b\ }$

Wikipedia

Back to the B cell lymphoma example



Hierarchical clustering

Similarity measure = Pearson correlation coefficient between gene expression vectors

**Similarity between clusters = average similarity
between individual elements of each cluster
(also called average linkage clustering)**



K-means clustering is a common alternative clustering approach

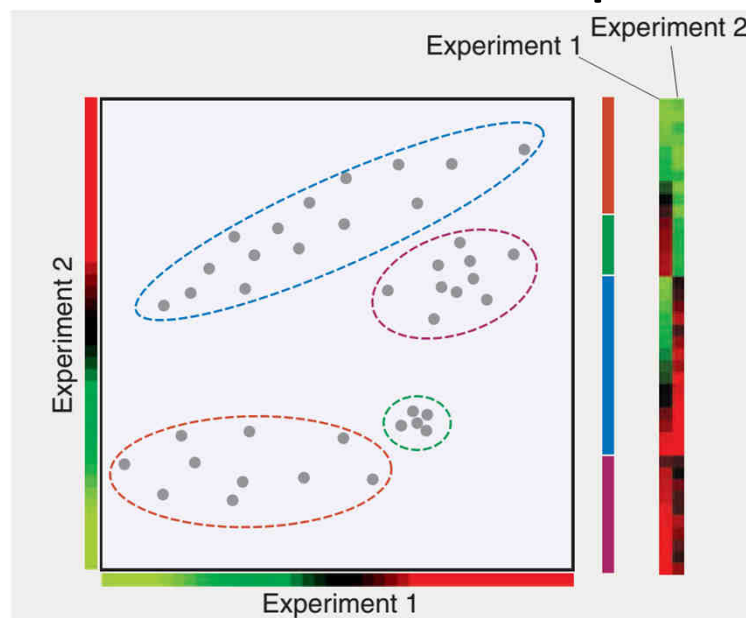
mainly because it's easy and can be quite fast!

The basic algorithm:

1. Pick a number (k) of cluster centers
2. Assign each gene to its nearest cluster center
3. Move each cluster center to the mean of its assigned genes
4. Repeat steps 2 & 3 until convergence

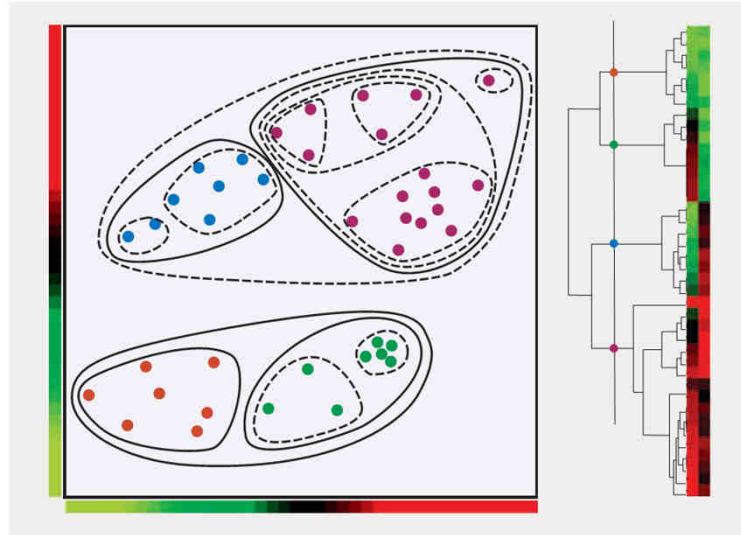
See the K-means example posted on the web site

A 2-dimensional example



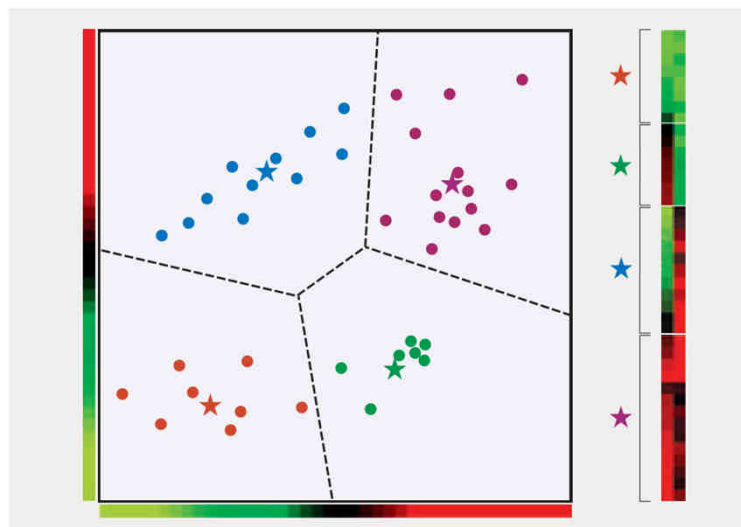
Nature Biotech 23(12):1499-1501 (2005)

A 2-dimensional example: hierarchical



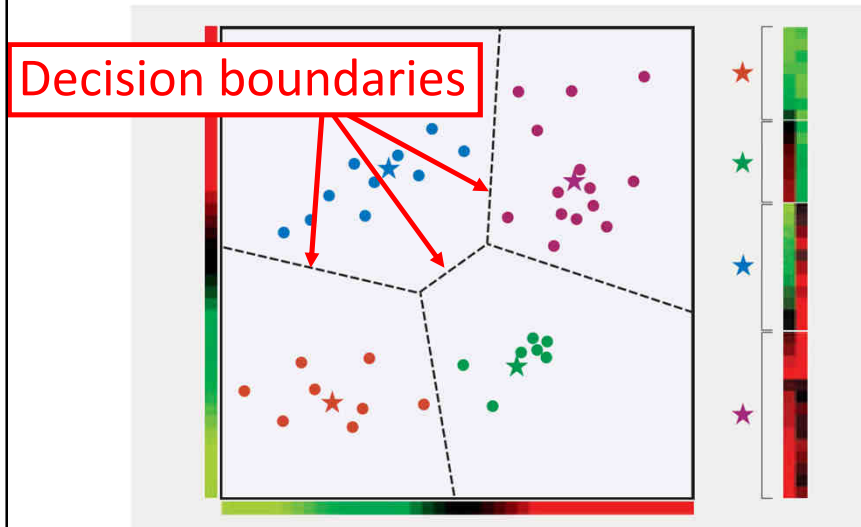
Nature Biotech 23(12):1499-1501 (2005)

A 2-dimensional example: k -means



Nature Biotech 23(12):1499-1501 (2005)

A 2-dimensional example: k -means



Nature Biotech 23(12):1499-1501 (2005)

Some features of K-means clustering

- Depending on how you seed the clusters, it may be stochastic. You may not get the same answer every time you run it.
- Every data point ends up in exactly 1 cluster (so-called *hard* clustering)
- Not necessarily obvious how to choose k
- Great example of something we've seen already: Expectation-Maximization (E-M) algorithms

EM algorithms alternate between assigning data to models (here, assigning points to clusters) and updating the models (calculating new centroids)

Some features of K-means clustering

- Depending on how you seed the clusters, it may be stochastic. You may not get the same answer every time you run it.
- Every data point ends up in exactly 1 cluster (so-called *hard* clustering)
- Not necessarily obvious how to choose k
- EM algorithm: alternating between assigning data points to clusters and updating the models (calculating new centroids)

Let's think about this aspect for a minute.

Why is this good or bad?

How could we change it?

k -means

The basic algorithm:

1. Pick a number (k) of cluster centers
2. Assign each gene to its nearest cluster center
3. Move each cluster center to the mean of its assigned genes
4. Repeat steps 2 & 3 until convergence

Fuzzy k -means

The basic algorithm:

1. Choose k . Randomly assign cluster centers.
2. Fractionally assign each gene to each cluster:

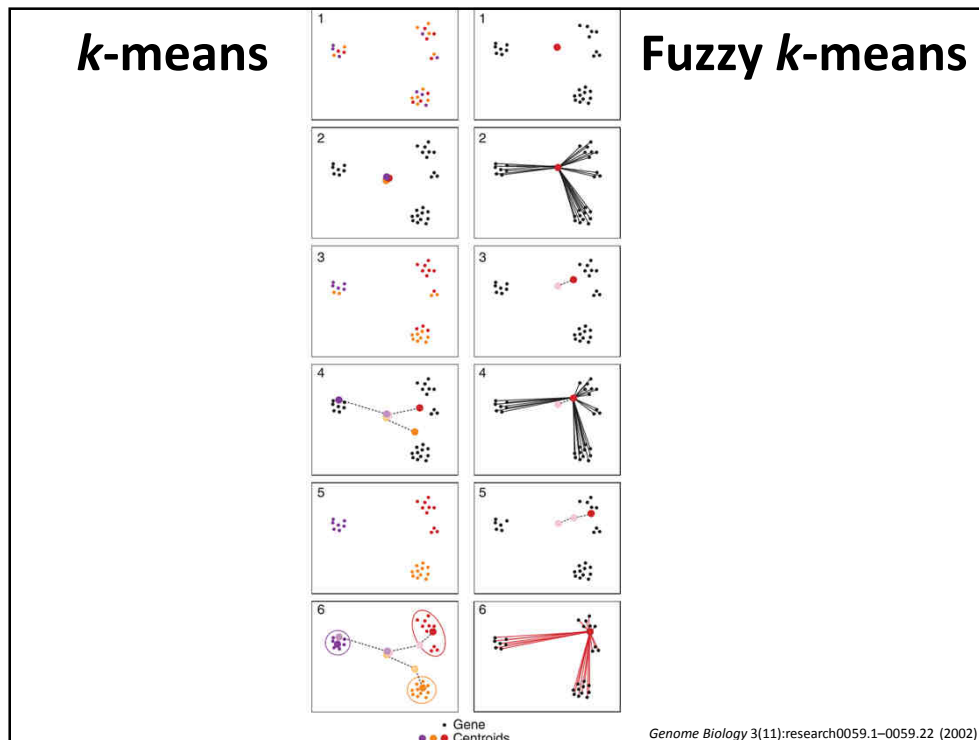
e.g. occupancy $(g_i, m_j) = \frac{e^{-\|g_i - m_j\|^2}}{\sum_j e^{-\|g_i - m_j\|^2}}$

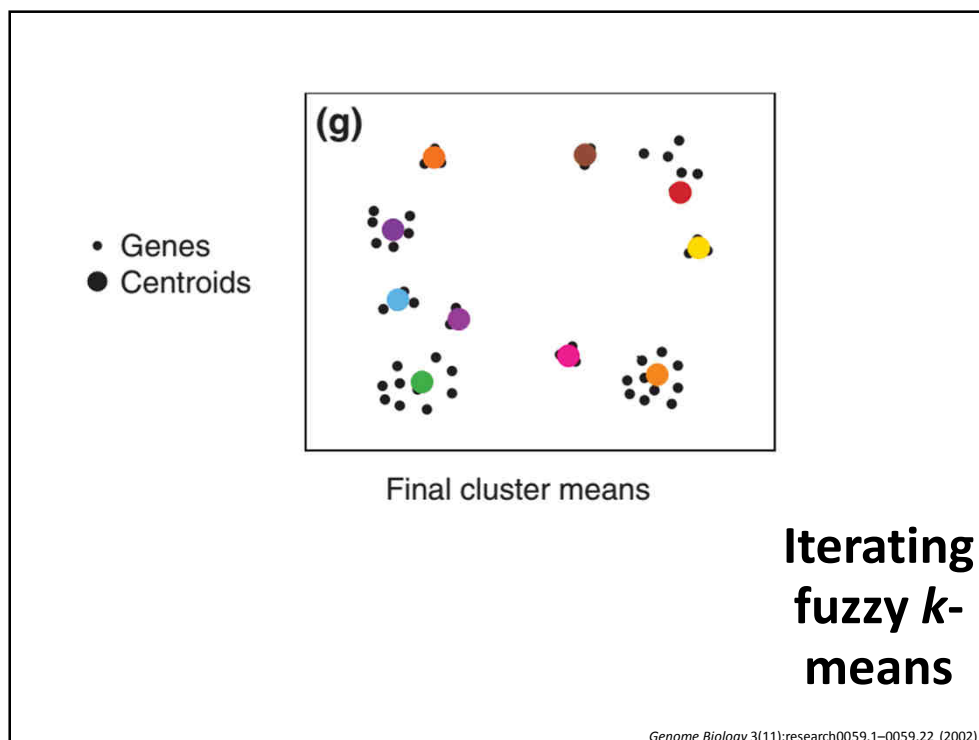
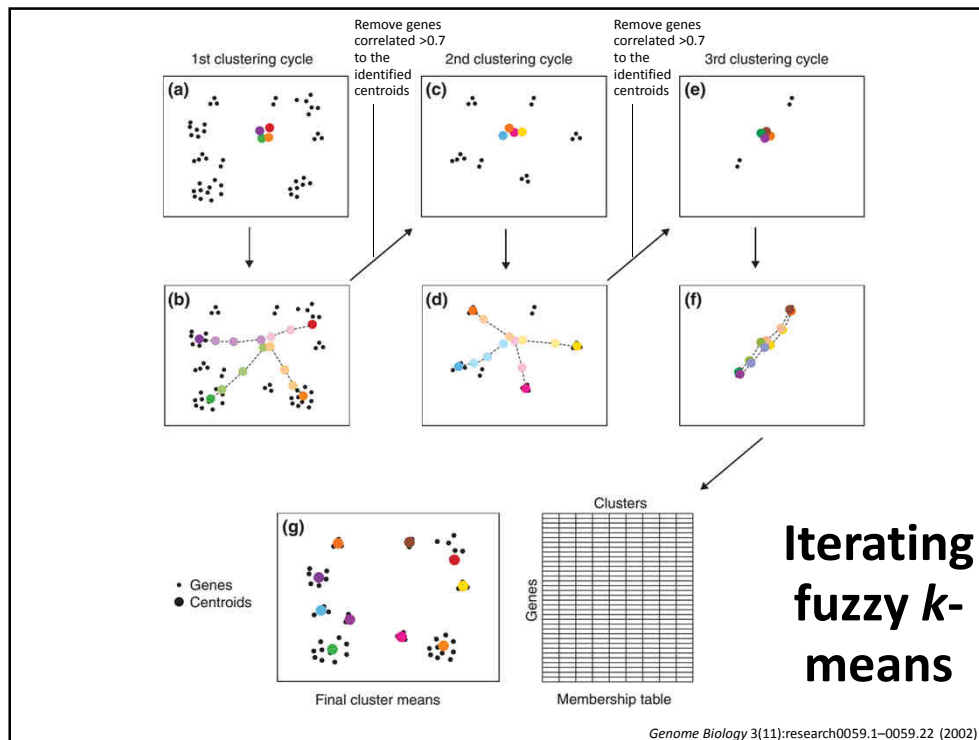
Note: $\|x\|$ is just shorthand for the length of the vector x .

g_i = gene i

m_j = centroid of cluster j

3. For each cluster, calculate weighted mean of genes to update cluster centroid
4. Repeat steps 2 & 3 until convergence





A fun clustering strategy that builds on these ideas: Self-organizing maps (SOMs)

- Combination of clustering & visualization
- Invented by Teuvo Kohonen, also called Kohonen maps



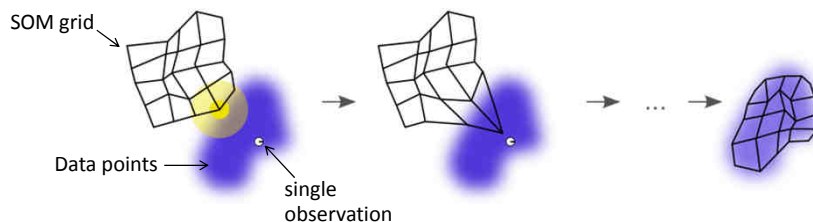
*Dr. Eng., Emeritus
Professor of the
Academy of Finland;
Academician*

A fun clustering strategy that builds on these ideas: Self-organizing maps (SOMs)

SOMs have:

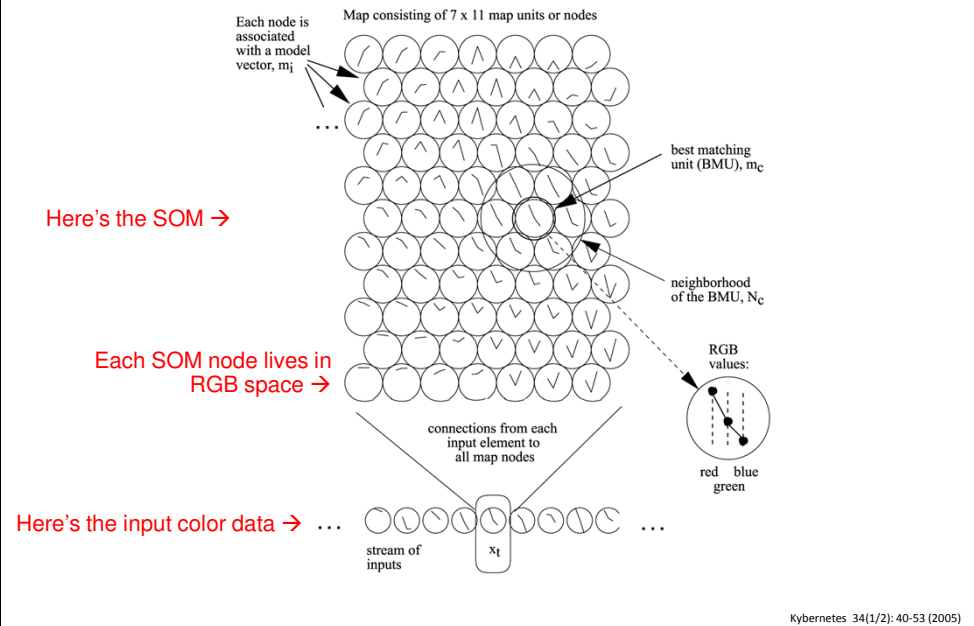
- your data (points in some high-dimensional space)
- a grid of nodes, each node also linked to a point someplace in data space

1. First, SOM nodes are arbitrarily positioned in data space. Then:
 2. Choose a training data point. Find the node closest to that point.
 3. Move its position closer to the training data point.
 4. Move its grid neighbors closer too, to a lesser extent.
- Repeat 2-4. After many iterations, the grid approximates the data distribution.

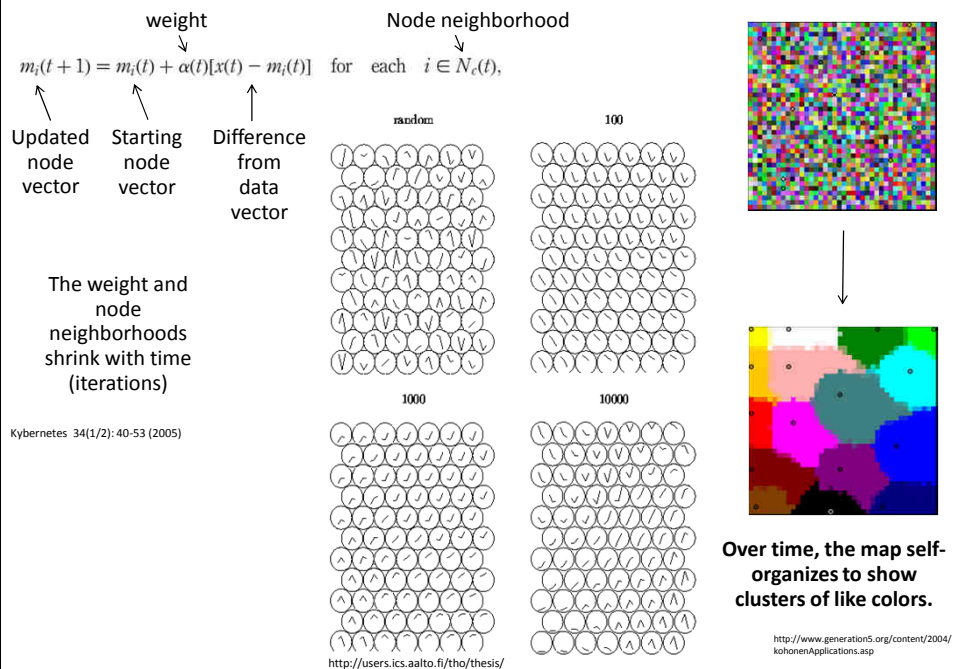


Wikipedia

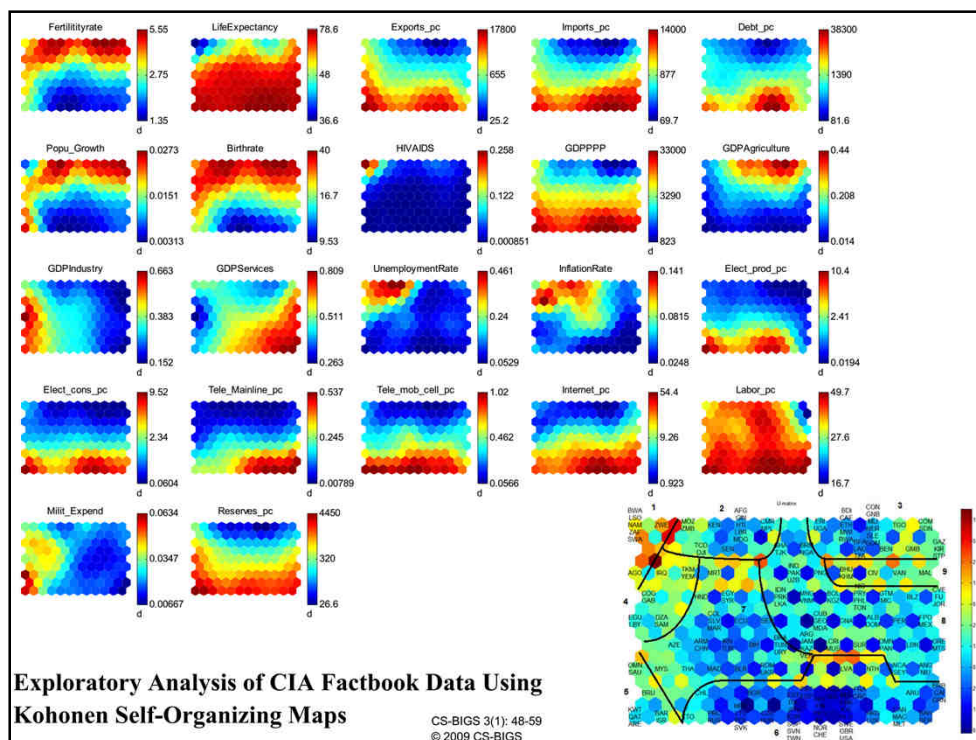
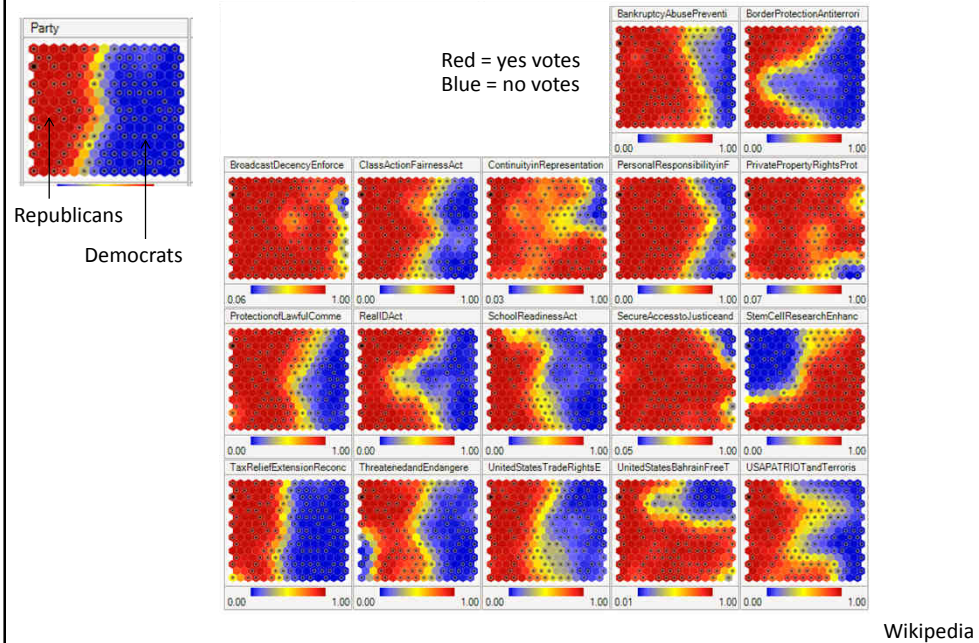
Here's an example using colors. Each color has an RGB vector. Take a bunch of random colors and organize them into a map of similar colors:

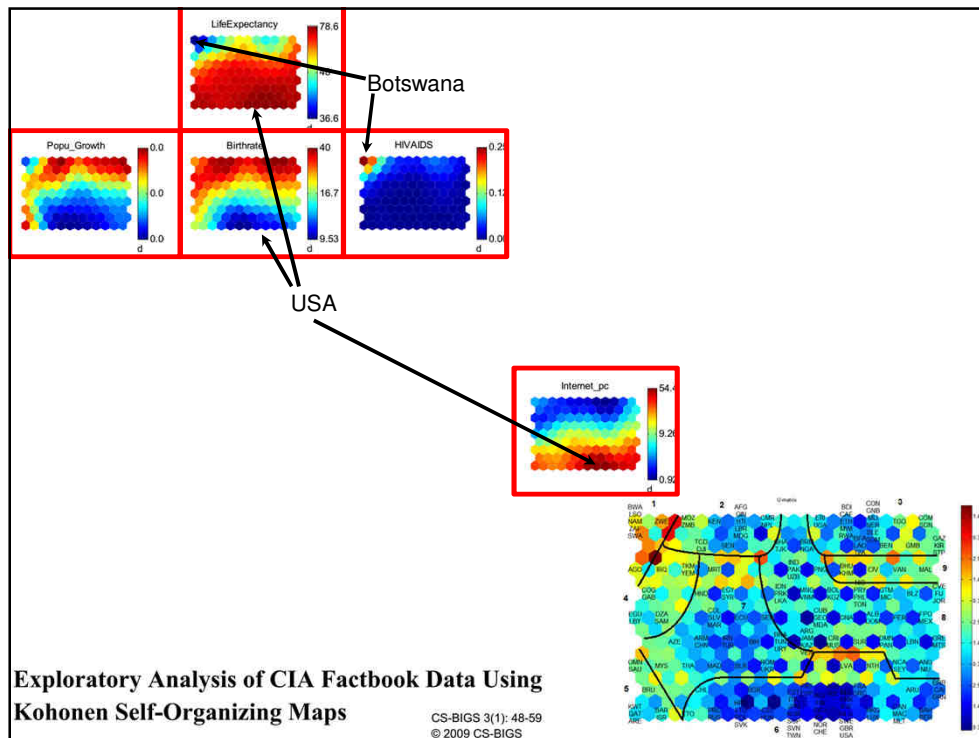


Iteratively test new colors, update the map using some rule

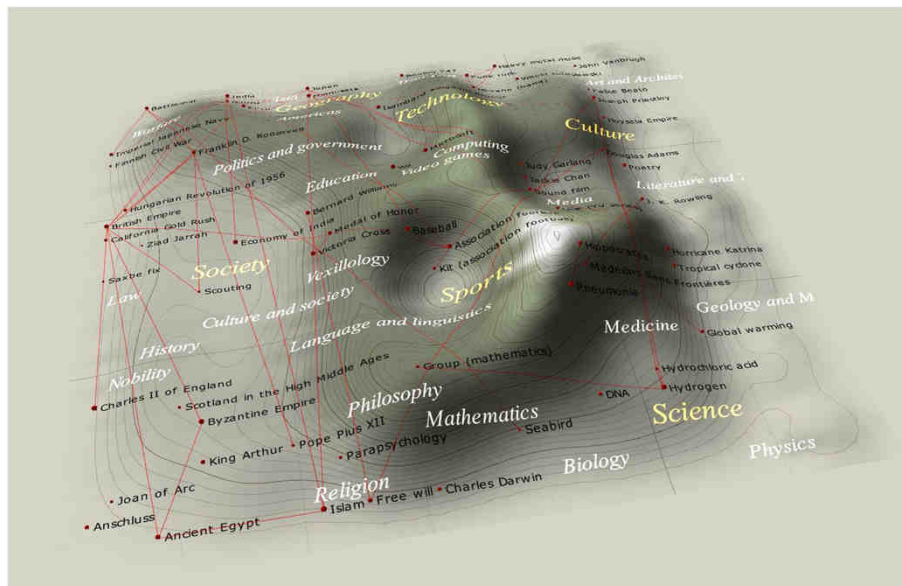


A SOM of U.S. Congress voting patterns





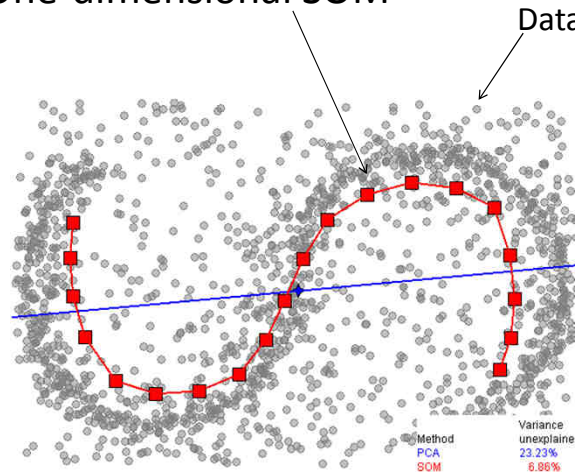
SOM of Wikipedia (from Wikipedia, naturally) (data = wiki article word frequency vectors)



Wikipedia

SOMs can accommodate unusual data distributions

One-dimensional SOM



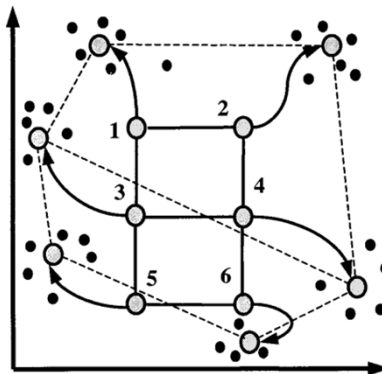
Wikipedia

A biological example, analyzing mRNA expression

Proc. Natl. Acad. Sci. USA
Vol. 96, pp. 2907-2912, March 1999
Genetics

Interpreting patterns of gene expression with self-organizing maps: Methods and application to hematopoietic differentiation

PABLO TAMAYO*, DONNA SLONIM*, JILL MESIROV*, QING ZHU†, SUTISAK KITAREEWAN‡, ETHAN DMITROVSKY‡,
ERIC S. LANDER*§, AND TODD R. GOLUB*†¶



Mitotic cell cycle

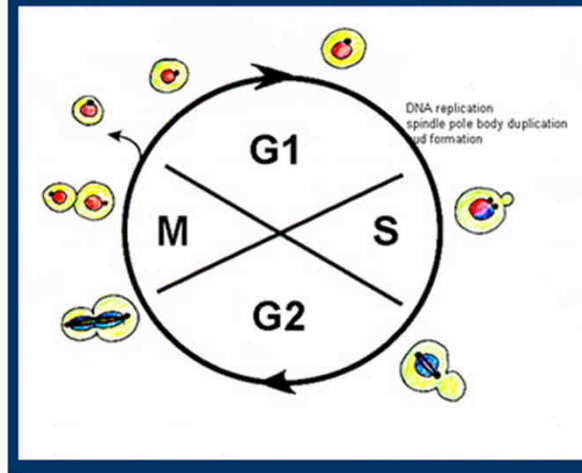
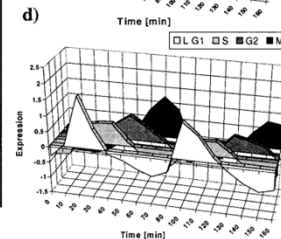
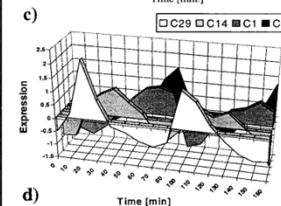
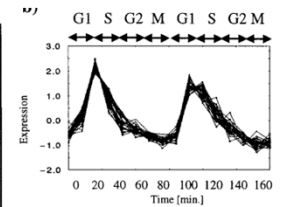
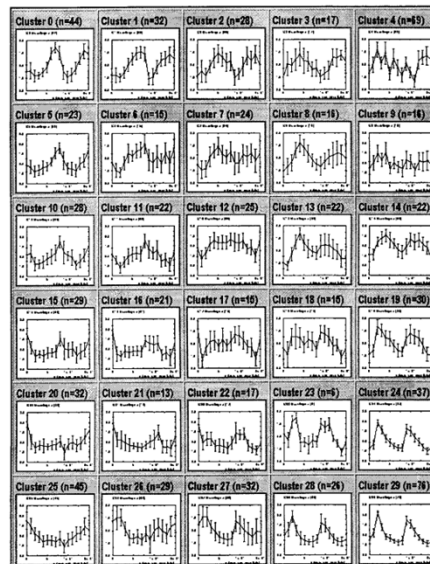


Image: <http://www.utoronto.ca/andrewslab/overview-Aux1.htm>

A biological example, analyzing mRNA expression

Yeast cell division cycle

Synchronized cells
↓
Collect mRNAs at
time points
↓
DNA microarrays



Finally, **t-SNE** is a nice way to visualize data in 2 or 3D
= *t-distributed stochastic neighbor embedding*

t-SNE tries to reproduce high-D *data neighborhoods* in a 2D or 3D picture by:

1. Defining a probability distribution over pairs of high-D objects such that “similar” objects have a high probability of being picked, whilst “dissimilar” objects have an extremely small probability of being picked
2. Defining a similar probability distribution over the points in the low-D map
3. Minimizing the Kullback–Leibler divergence between the two distributions by varying the locations of the points in the low-D map, i.e.

minimize this:
$$\sum_{i \neq j} p_{ij} \log \frac{p_{ij}}{q_{ij}}$$

Sum over all pairs of points

p_{ij} ← probability i and j are close in high-D space
 q_{ij} ← probability i and j are close in low-D space

van der Maaten & Hinton, Visualizing High-Dimensional Data Using t-SNE.
Journal of Machine Learning Research 9: 2579–2605 (Nov 2008)

**You can compute your own t-SNE embeddings
using the online tools at:**
<http://projector.tensorflow.org/>

There are also some great examples at:
<http://distill.pub/2016/misread-tsne/>

There are only a couple of parameters you can tweak, mainly perplexity,
which effectively captures the number of neighbors (often 5 to 50)