

Functional genomics + Data mining

BIO337 Systems Biology / Bioinformatics – Spring 2014

Edward Marcotte, Univ of Texas at Austin

Edward Marcotte/Univ. of Texas/BIO337/Spring 2014

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.

Edward Marcotte/Univ. of Texas/BIO337/Spring 2014

Adapted from Wikipedia

Functional genomics + Data mining

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

Edward Marcotte/Univ. of Texas/BIO337/Spring 2014

Adapted from Wikipedia

Functional genomics + Data mining



Edward Marcotte/Univ. of Texas/BIO337/Spring 2014

www.sparkpeople.com

Adapted from Wikipedia

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

Edward Marcotte/Univ. of Texas/BIO337/Spring 2014

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).

Edward Marcotte/Univ. of Texas/BIO337/Spring 2014

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

Edward Marcotte/Univ. of Texas/BIO337/Spring 2014

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

Edward Marcotte/Univ. of Texas/BIO337/Spring 2014

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.”

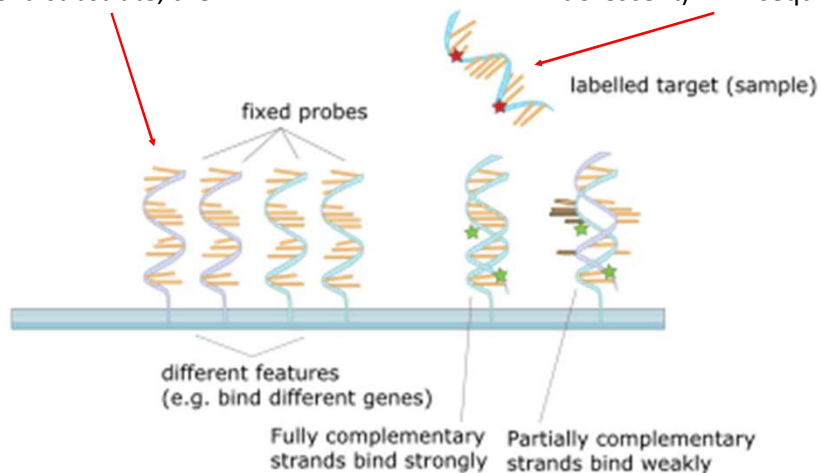
Edward Marcotte/Univ. of Texas/BIO337/Spring 2014

Nature 2000

Refresher: Profiling mRNA expression with DNA microarrays

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

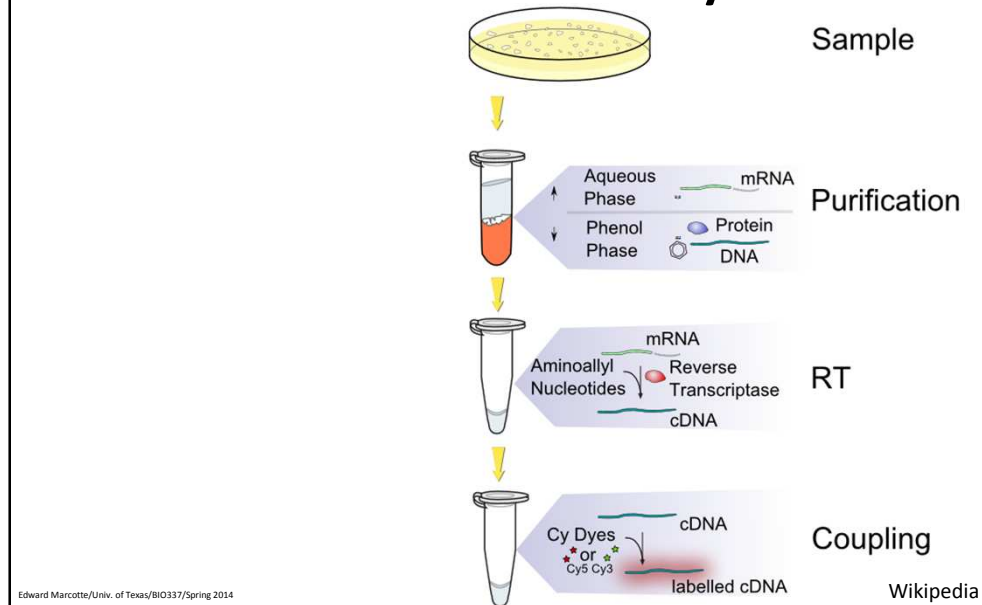
...probed with a labeled (usually fluorescent) DNA sequence



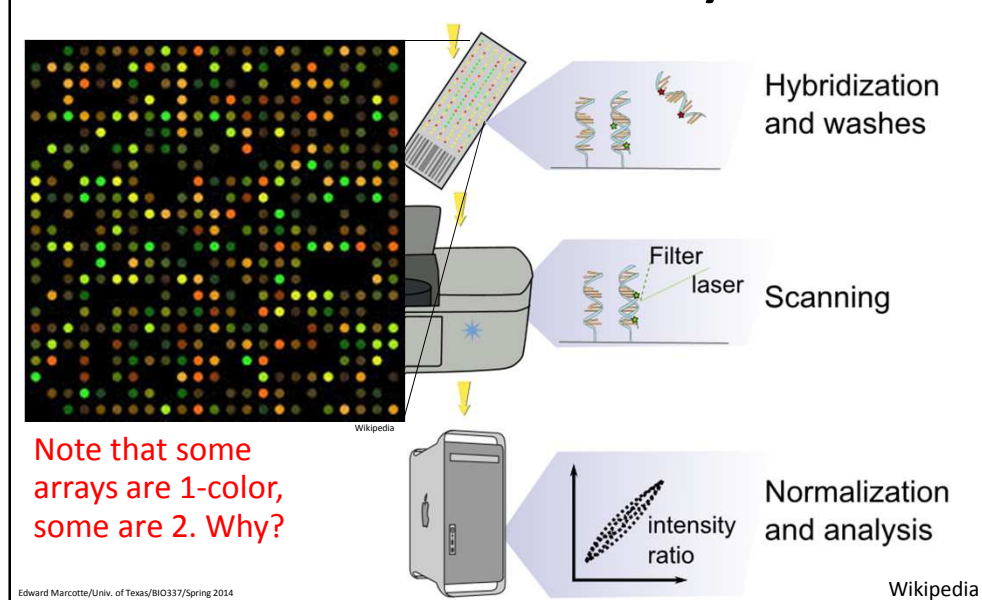
Edward Marcotte/Univ. of Texas/BIO 337/Spring 2014

Wikipedia

Refresher: Profiling mRNA expression with DNA microarrays



Refresher: Profiling mRNA expression with DNA microarrays



Back to diffuse large B-cell lymphoma...

96 patient biopsies
(normal and malignant lymphocyte samples)

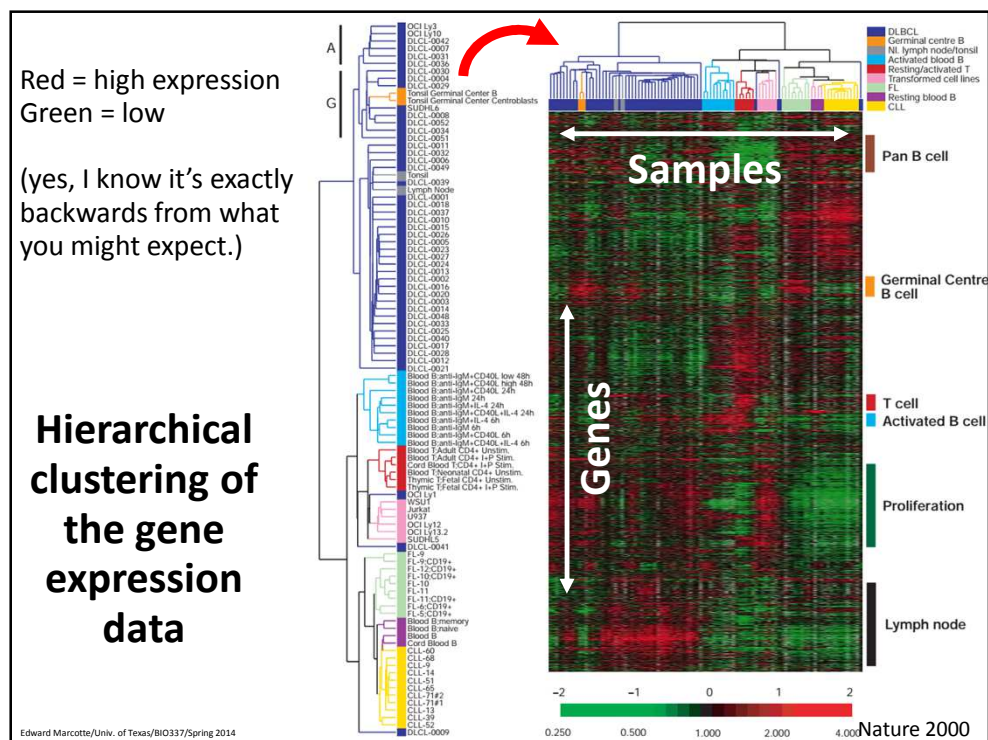
Extract mRNA from each sample

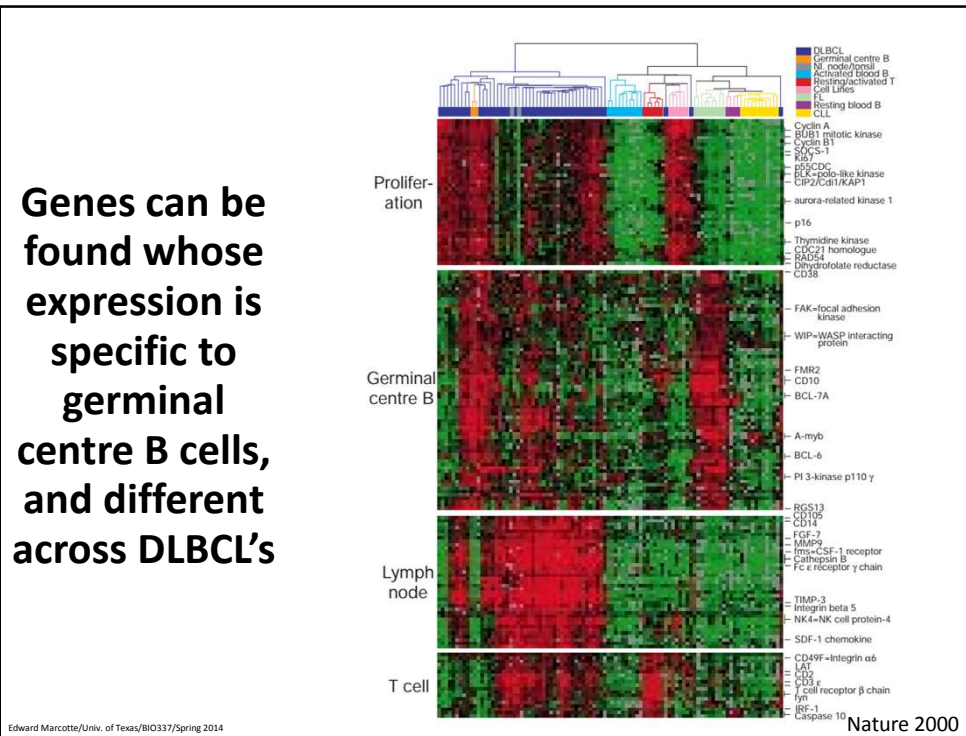
Perform DNA microarray experiment on each to measure mRNA abundances (~1.8 million total gene expression measurements)

Cluster samples by their expression patterns

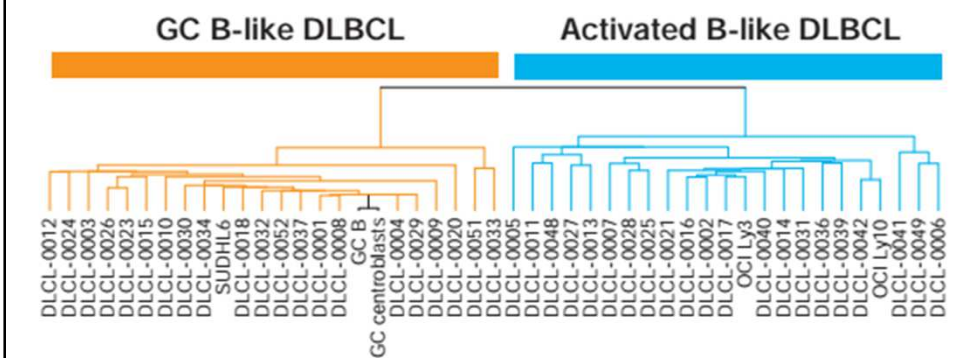
Edward Marcotte/Univ. of Texas/BIO337/Spring 2014

Nature 2000



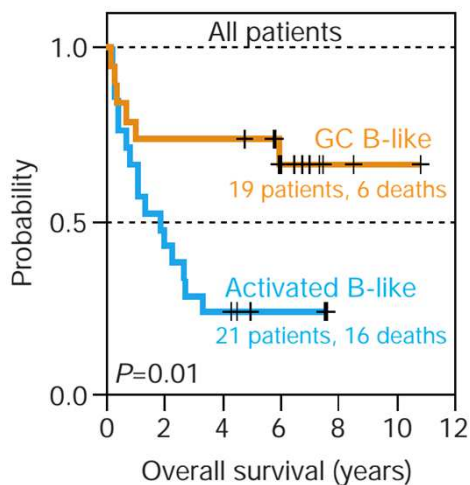


We can break up the DLBCL's according the germinal B-cell specific gene expression:



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

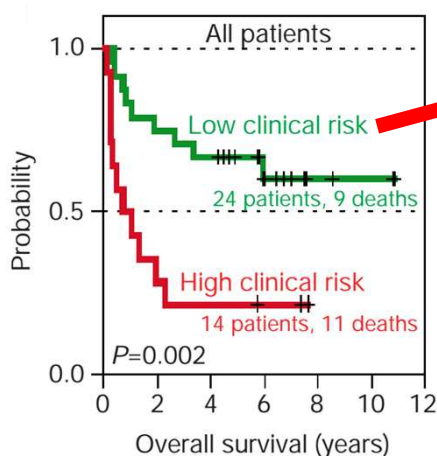
Kaplan-Meier plot
of patient survival



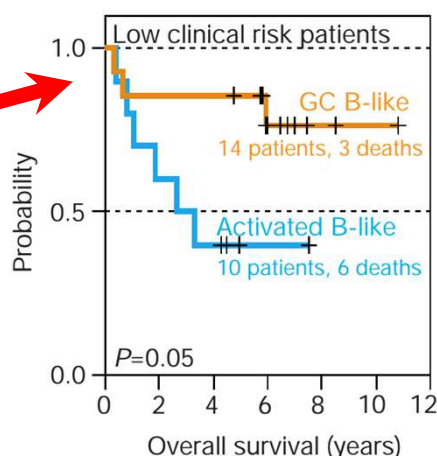
Edward Marcotte/Univ. of Texas/BIO337/Spring 2014

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

Edward Marcotte/Univ. of Texas/BIO337/Spring 2014

Nature 2000

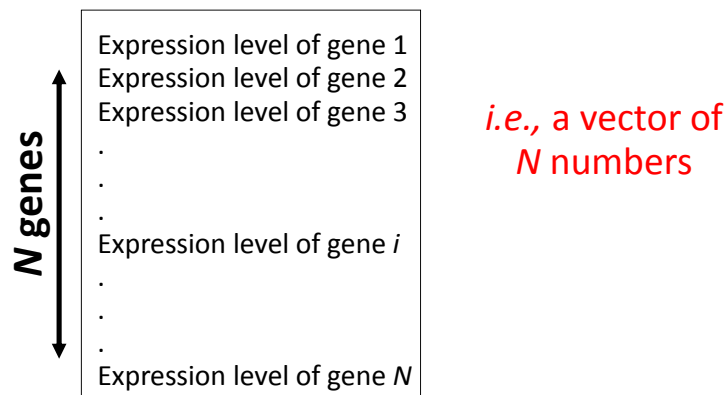
Gene expression, and other molecular measurements, provide far deeper phenotypes for cells, tissues, and organisms than traditional measurements

Now, tons of work using these approaches to diagnose specific forms of disease, as well as to discover functions of genes and many other applications

Edward Marcotte/Univ. of Texas/BIO337/Spring 2014

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:

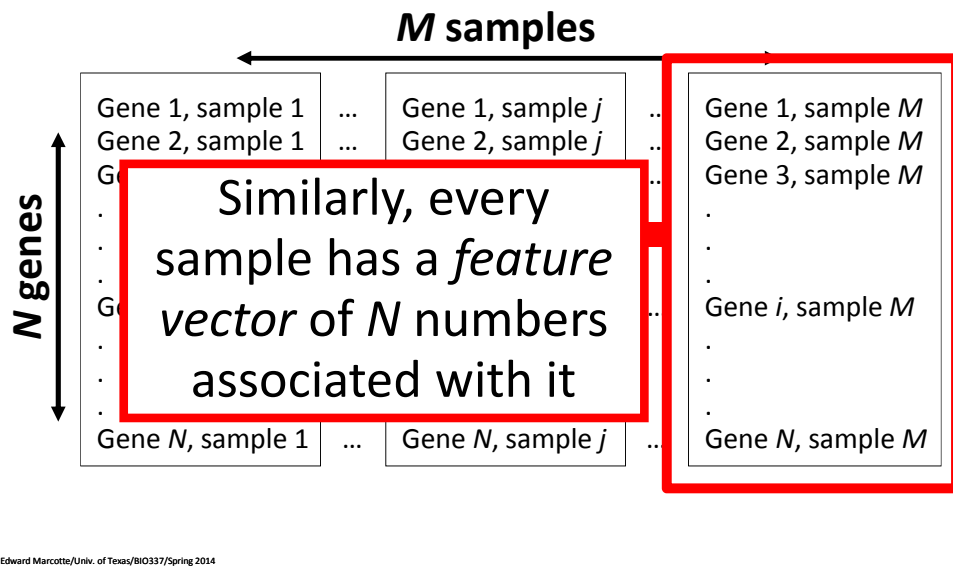


For yeast, $N \sim 6,000$

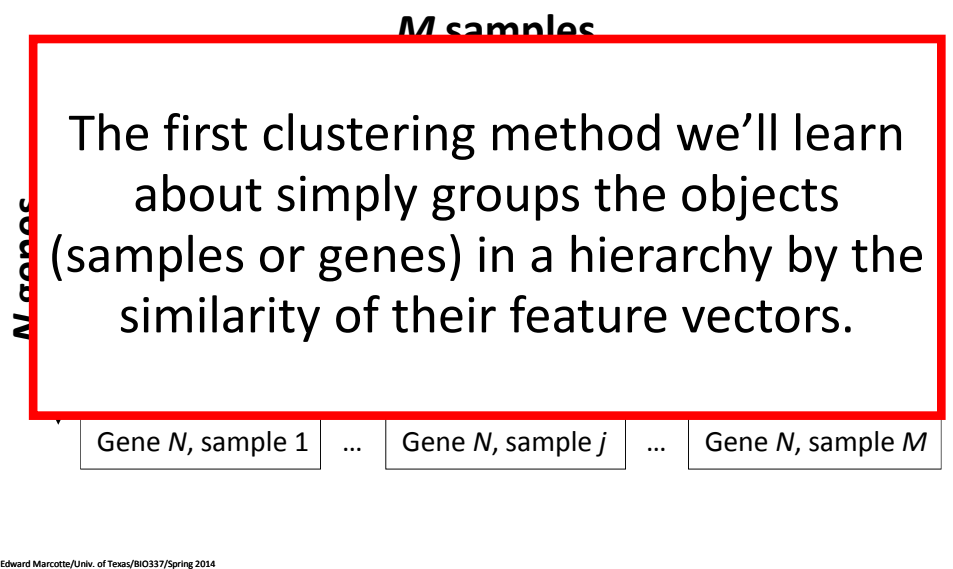
For human, $N \sim 22,000$

Edward Marcotte/Univ. of Texas/BIO337/Spring 2014

So, how does clustering work?



So, how does clustering work?



A hierarchical clustering algorithm

Start with each object in its own cluster

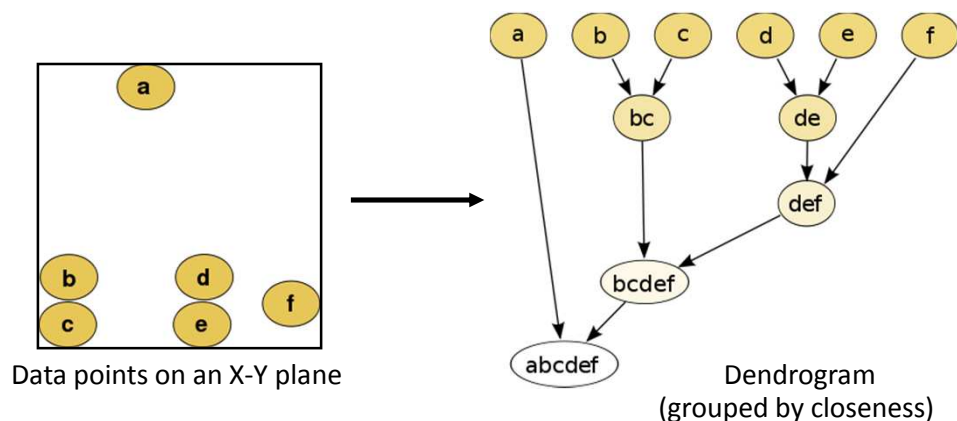
Until there is only one cluster left, repeat:
 Among the current clusters, find the two
 most similar clusters
 Merge those two clusters into one

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

Edward Marcotte/Univ. of Texas/BIO337/Spring 2014

Hierarchical clustering

Conceptually



Edward Marcotte/Univ. of Texas/BIO337/Spring 2014

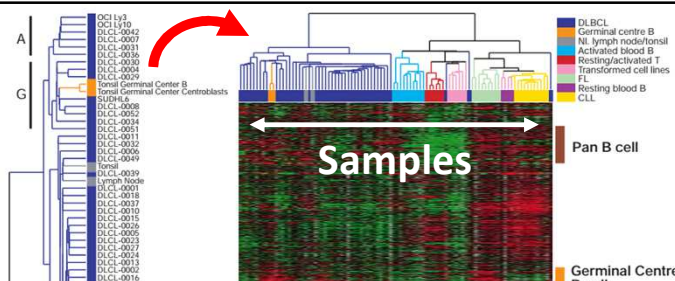
Wikipedia

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)

Edward Marcotte/Univ. of Texas/BIO337/Spring 2014

~2 -1 0 1 2
0.250 0.500 1.000 2.000 4.000 Nature 2000

K-means clustering is a common alternative clustering approach

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

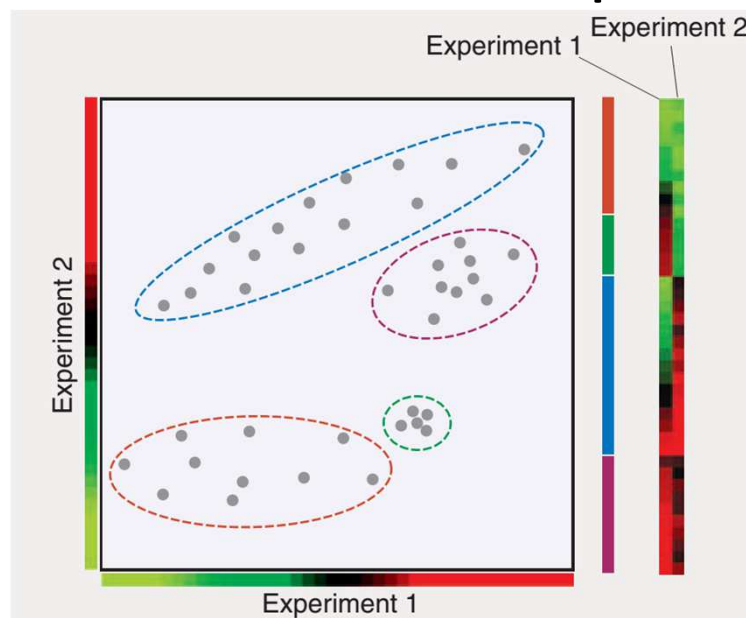
The basic algorithm:

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

Edward Marcotte/Univ. of Texas/BIO337/Spring 2014

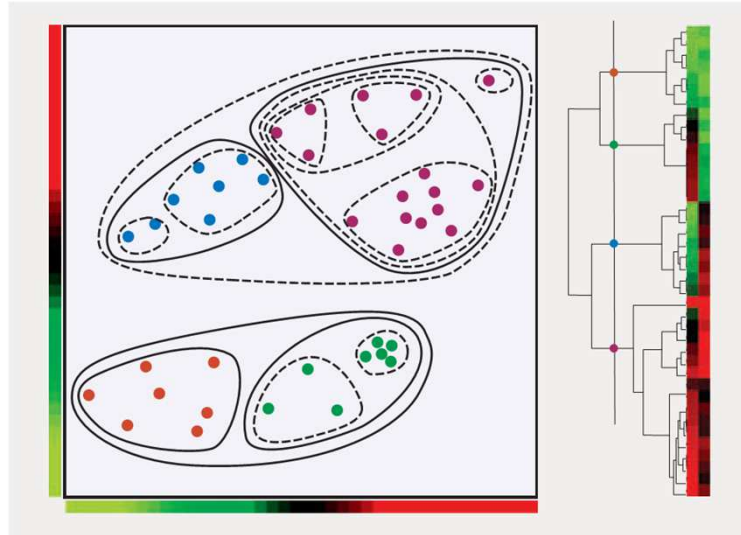
A 2-dimensional example



Edward Marcotte/Univ. of Texas/BIO337/Spring 2014

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

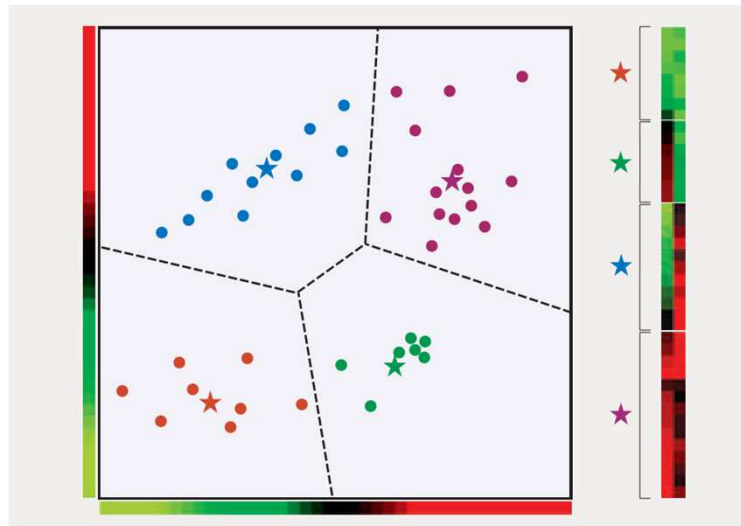
A 2-dimensional example: hierarchical



Edward Marcotte/Univ. of Texas/BIO337/Spring 2014

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

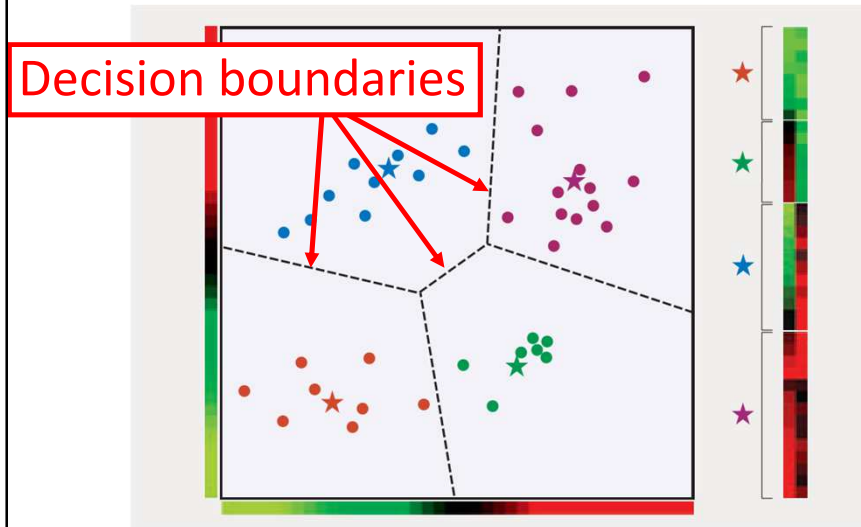
A 2-dimensional example: k -means



Edward Marcotte/Univ. of Texas/BIO337/Spring 2014

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

A 2-dimensional example: k -means



Edward Marcotte/Univ. of Texas/BIO337/Spring 2014

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'll meet again: 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)

Edward Marcotte/Univ. of Texas/BIO337/Spring 2014

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: 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?

Edward Marcotte/Univ. of Texas/BIO337/Spring 2014

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

Edward Marcotte/Univ. of Texas/BIO337/Spring 2014

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 .

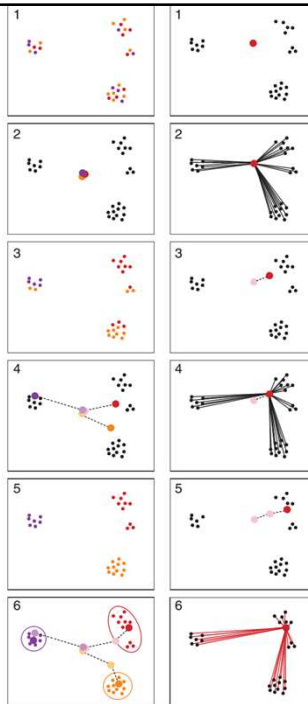
g_i = gene i

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

Edward Marcotte/Univ. of Texas/BIO337/Spring 2014

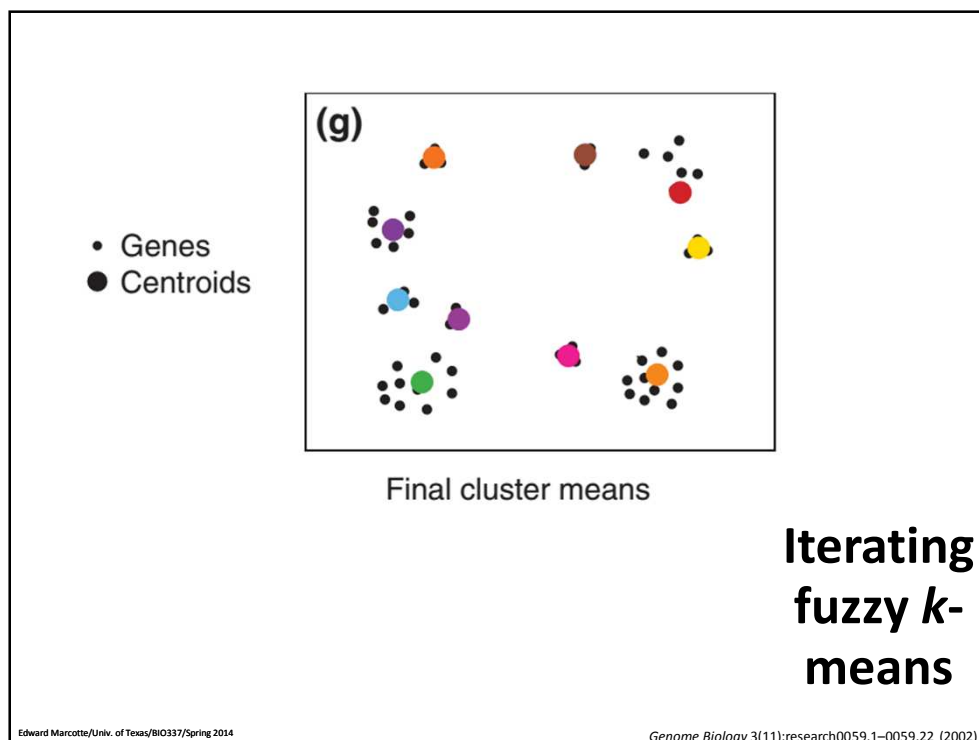
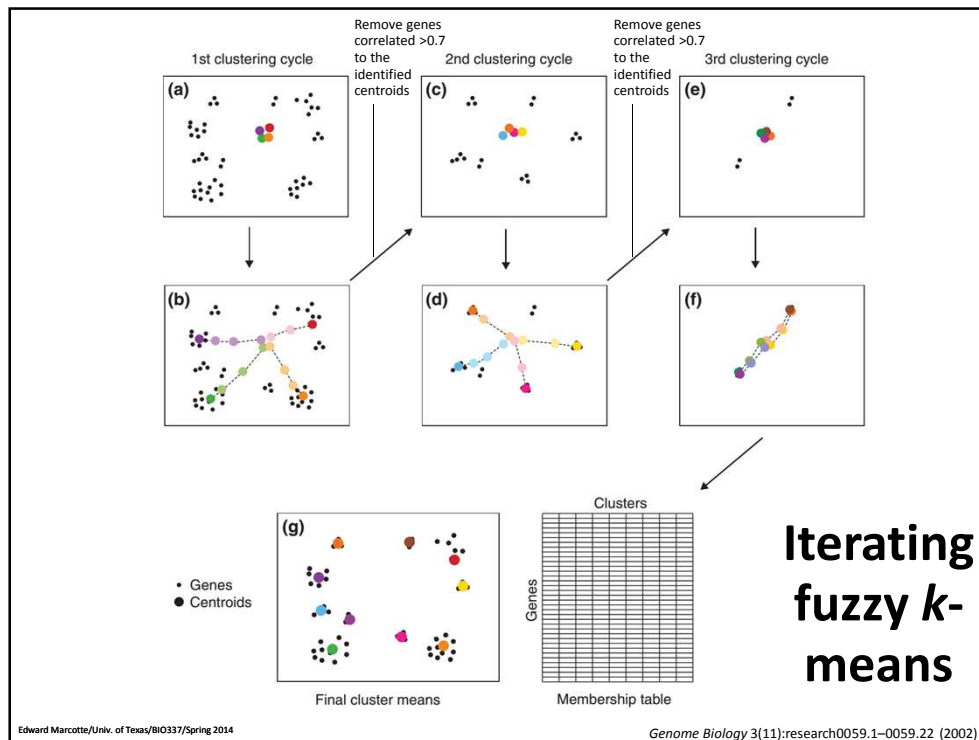
k -means



Fuzzy k -means

• Gene
• Centroids

Genome Biology 3(11):research0059.1–0059.22 (2002)



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*

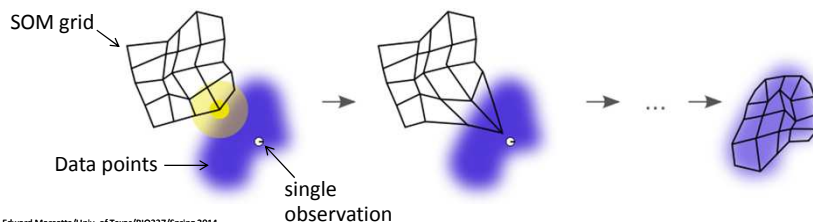
Edward Marcotte/Univ. of Texas/BIO337/Spring 2014

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

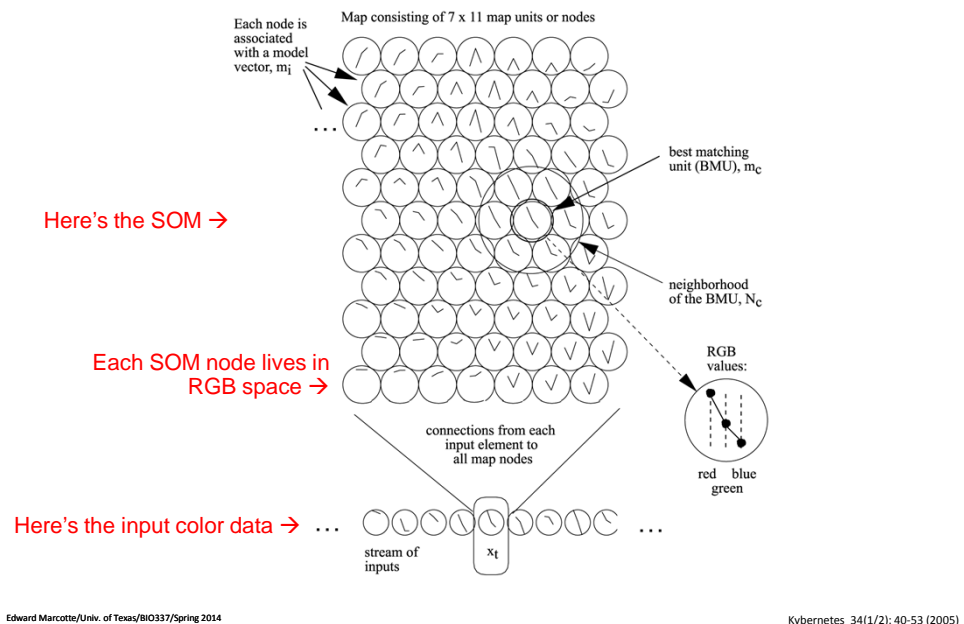
1. First, SOM nodes are arbitrarily positioned in data space. Then:
 2. Choose a training data point. Find the node closest to that point.
 3. Move its position closer to the training data point.
 4. Move its grid neighbors closer too, to a lesser extent.
- Repeat 2-4. After many iterations, the grid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



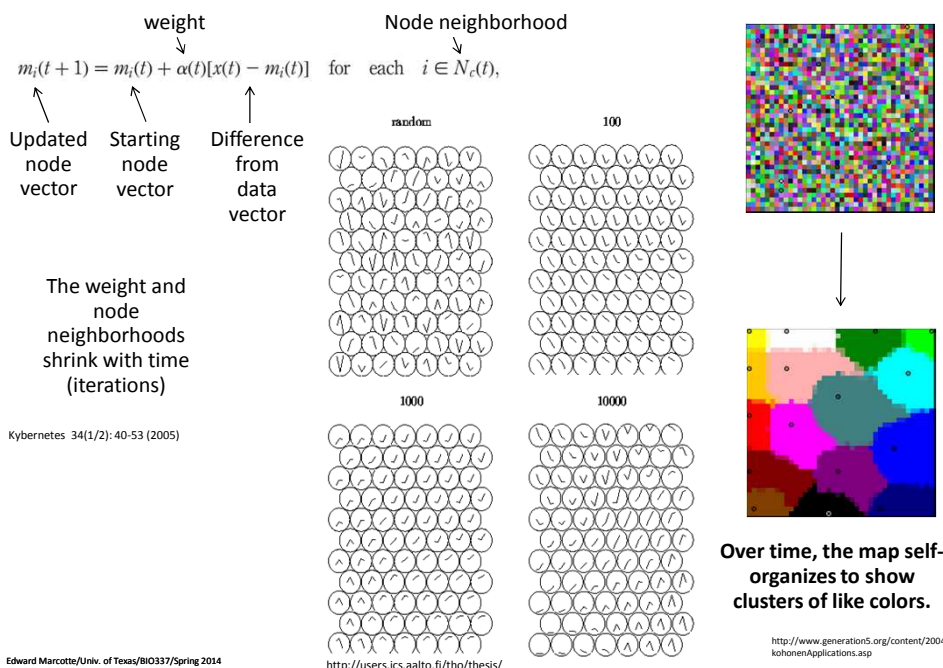
Edward Marcotte/Univ. of Texas/BIO337/Spring 2014

Wikipedia

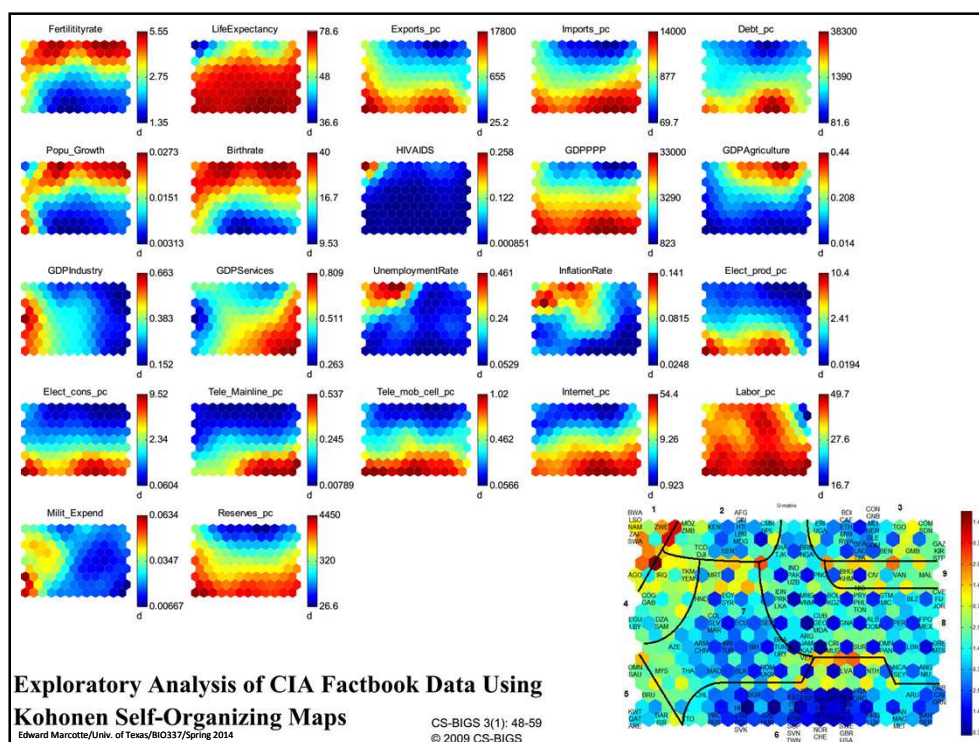
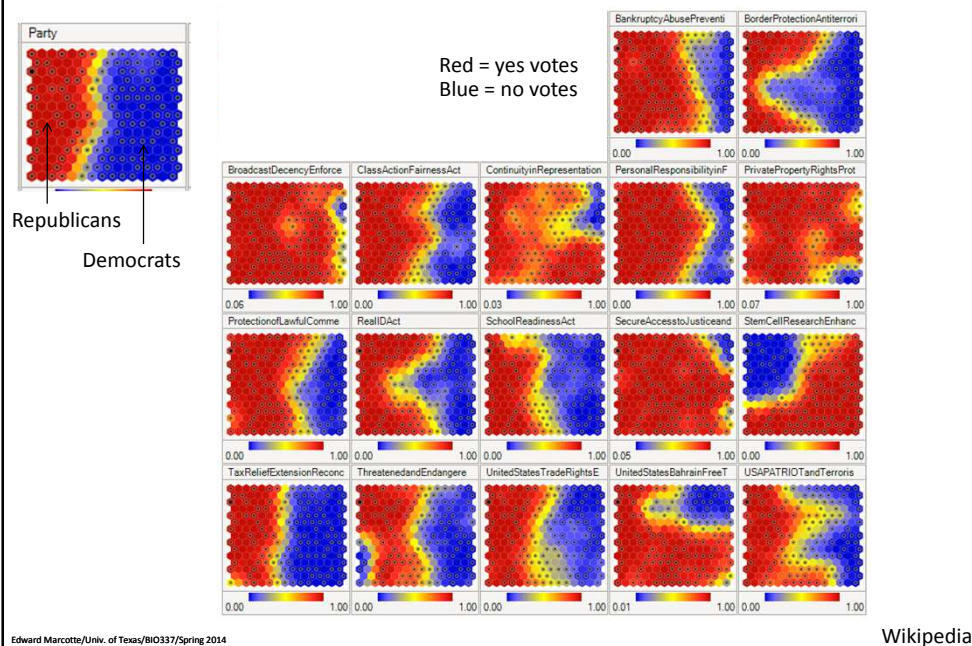
Here's an example using colors. Each color has an RGB vector. Take a bunch of random colors and organize them into a map of similar colors:



Iteratively test new colors, update the map using some rule



A SOM of U.S. Congress voting patterns



The diagram illustrates the interconnectedness of various fields of knowledge. The fields are represented as nodes, and the lines connecting them show the relationships between these fields. The fields are arranged in a circular pattern, with each field connected to multiple other fields, creating a dense network. The fields include:

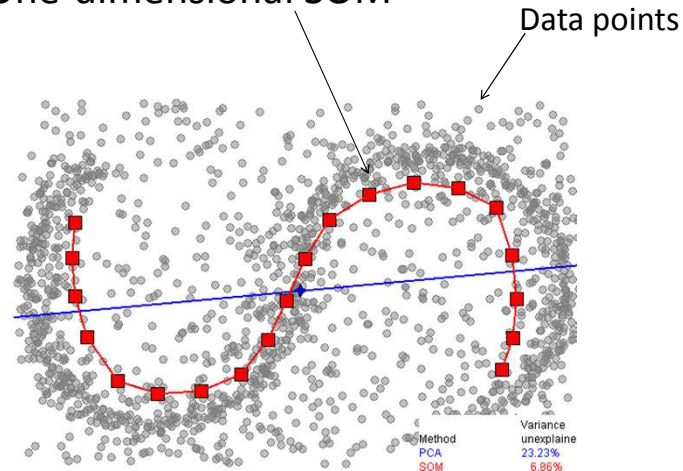
- History
- Geography
- Technology
- Culture
- Politics and government
- Education
- Society
- Law
- Nobility
- Philosophy
- Mathematics
- Religion
- Biology
- Science
- Physics
- Medicine
- Geology and M
- Sports
- Literature and Media
- Art and Architecture

The diagram is set against a background of a world map, with the fields of knowledge arranged in a circular pattern around the center. The lines connecting the fields are of varying thickness, indicating the strength or frequency of the relationships between them. The overall structure of the diagram suggests a highly interconnected and complex world of knowledge.

Wikipedia

SOMs can accommodate unusual data distributions

One-dimensional SOM



Edward Marcotte/Univ. of Texas/BIO337/Spring 2014

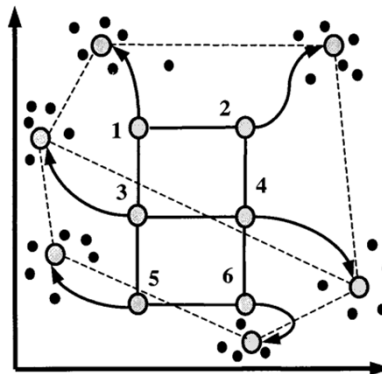
Wikipedia

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*†¶



Edward Marcotte/Univ. of Texas/BIO337/Spring 2014

Mitotic cell cycle

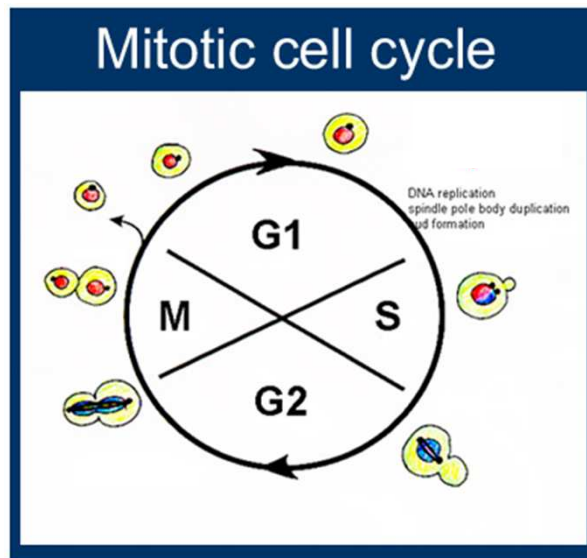


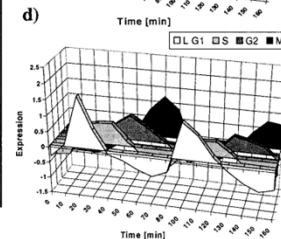
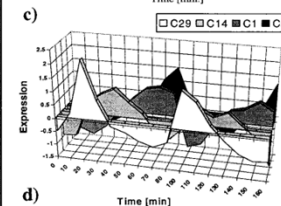
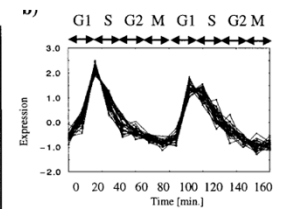
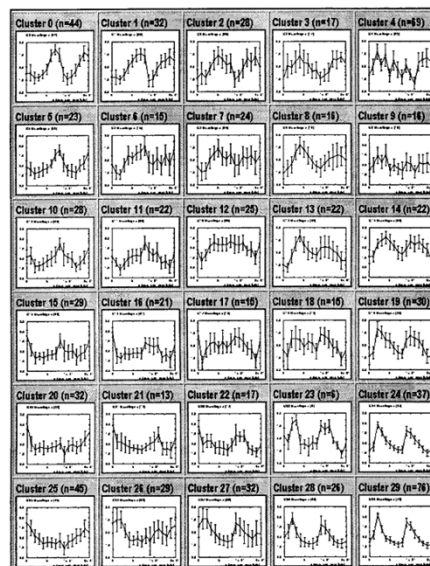
Image: <http://www.utoronto.ca/andrewslab/overview-Aux1.htm>

Edward Marcotte/Univ. of Texas/BIO337/Spring 2014

A biological example, analyzing mRNA expression

Yeast cell division cycle

Synchronized cells
↓
Collect mRNAs at
time points
↓
DNA microarrays



Edward Marcotte/Univ. of Texas/BIO337/Spring 2014