

Course	Section	CRN	Date	Day	Start Time	End Time	Room	Exam Type
BIOE 505	CB	54270	12/12/2023	T	8:00 AM	11:00 AM	1302 Everitt Laboratory	Non Combined



	p1 (10)	p2 (10)	P3 (15)	P4 (10)	P5 (10)	P6 (10)
mean score	13.68	9.95	14.58	9.21	9.95	8.47
std score	2.69	0.23	0.90	1.40	0.23	1.58
points subtracted	1.32	0.05	0.42	0.79	0.05	1.53
	15	10	15	10	10	10

1. **(15 points)** If the letters of ILLINI are randomly ordered, all orderings being equally likely, what is the probability that not a single position has the same letter as in the original order? Hint: 3 Is (and 2 Ls) are identical.

**Answer:** Three letters I must go in places of L, L, and N. Once I pick where N goes, the rest is determined. There are 3 places to put N. There are 3 solutions. The total number of ways to order these 6 letters is  $6!/(3!*2!*1!)=6*5*4/2=60$ . Hence the probability is  $3/60=1/20=0.05$

4. **(10 points)** The common logarithm (base 10) of the expression level (mRNA copies/cell) of a cancer driver gene in a randomly selected cell is normally distributed with mean  $\mu = 4$ , and standard dev.  $\sigma = 1$ .

**(a) (5 points)** What is the probability that the expression level measured in a given cell is between 1000 and 1,000,000 mRNA copies/cell?

**Answer:**  $P(-1 < Z < 2) = P(Z < 2) - P(Z < -1) = 0.97725 - 0.15866 = 0.81859$

**(b) (5 points)** Expression level of this gene was measured in 6 individual cells. What is the probability that in exactly 3 cells gene's expression level within these bounds?

**Answer:**  $(6! / (3!3!)) * (0.81859^3) * (1 - 0.81859)^3 = 0.0655$

6. **(10 points)** In a data communication system, several messages that arrive at a node are bundled into a packet before they are transmitted over the network. Assume the messages arrive according to a Poisson process with the mean rate equal to two messages per five minutes. Six messages are required to form a packet and the packet is formed immediately after the last message has arrived.

**(a) (5 points)** What is the probability that a time interval between two consecutive messages is longer than 4 minutes?

**Answer:**  $\lambda = 2 \text{ message}/5 \text{ minutes} = 0.4 \text{ messages/minute}$ .

Exponential distribution  $P(X > 4) = \exp(-0.4 * 4) = \exp(-1.6) = 0.2019$

**(b) (5 points)** What is the mean time until a packet is formed, that is, until exactly six messages have arrived at the node?

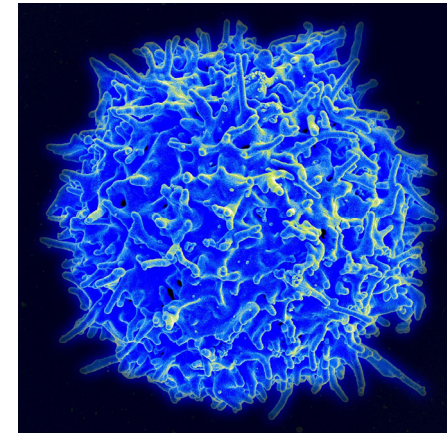
**Answer:** Using Erlang distribution with  $r=6$ ,  $\lambda=0.4$  one gets  $(6/0.4) \text{ minutes} = 15 \text{ minutes}$

# Clustering analysis of gene expression data

Chapter 11 in  
Jonathan Pevsner,  
Bioinformatics and Functional Genomics,  
3<sup>rd</sup> edition  
(Chapter 9 in 2<sup>nd</sup> edition)

# Human T cell expression data

- The matrix contains **47 expression samples** from Lukk et al, Nature Biotechnology 2010
- All samples are **from T cells in different individuals**
- Only the **top 3000 genes** with the largest variability **were used**
- The value is **log2 of gene's expression level** in a given sample as measured by the microarray technology



## A global map of human gene expression

Margus Lukk, Misha Kapushesky, Janne Nikkilä, Helen Parkinson, Angela Goncalves, Wolfgang Huber, Esko Ukkonen & Alvis Brazma

Affiliations | Corresponding author

*Nature Biotechnology* **28**, 322–324 (2010) | doi:10.1038/nbt0410-322

Although there is only one human genome sequence, different genes are expressed in many different cell types and tissues, as well as in different developmental stages or diseases. The structure of this 'expression space' is still largely unknown, as most transcriptomics experiments focus on sampling small regions. We have constructed a global gene expression map by integrating microarray data from 5,372 human samples representing 369 different cell and tissue types, disease states and cell lines. These have been compiled in an online resource (<http://www.ebi.ac.uk/gxa/array/U133A>) that allows the user to search for a gene of interest and





# Correlated pairs plausible biological connection based on short description

**g1=1994; g2=188; group 1**

**g1=2872; g2=1269; group 2**

**g1=1321; g2=10; group 3**

**g1= 886; g2=819; group 4**

**g1=2138; g2=1364; group 5**

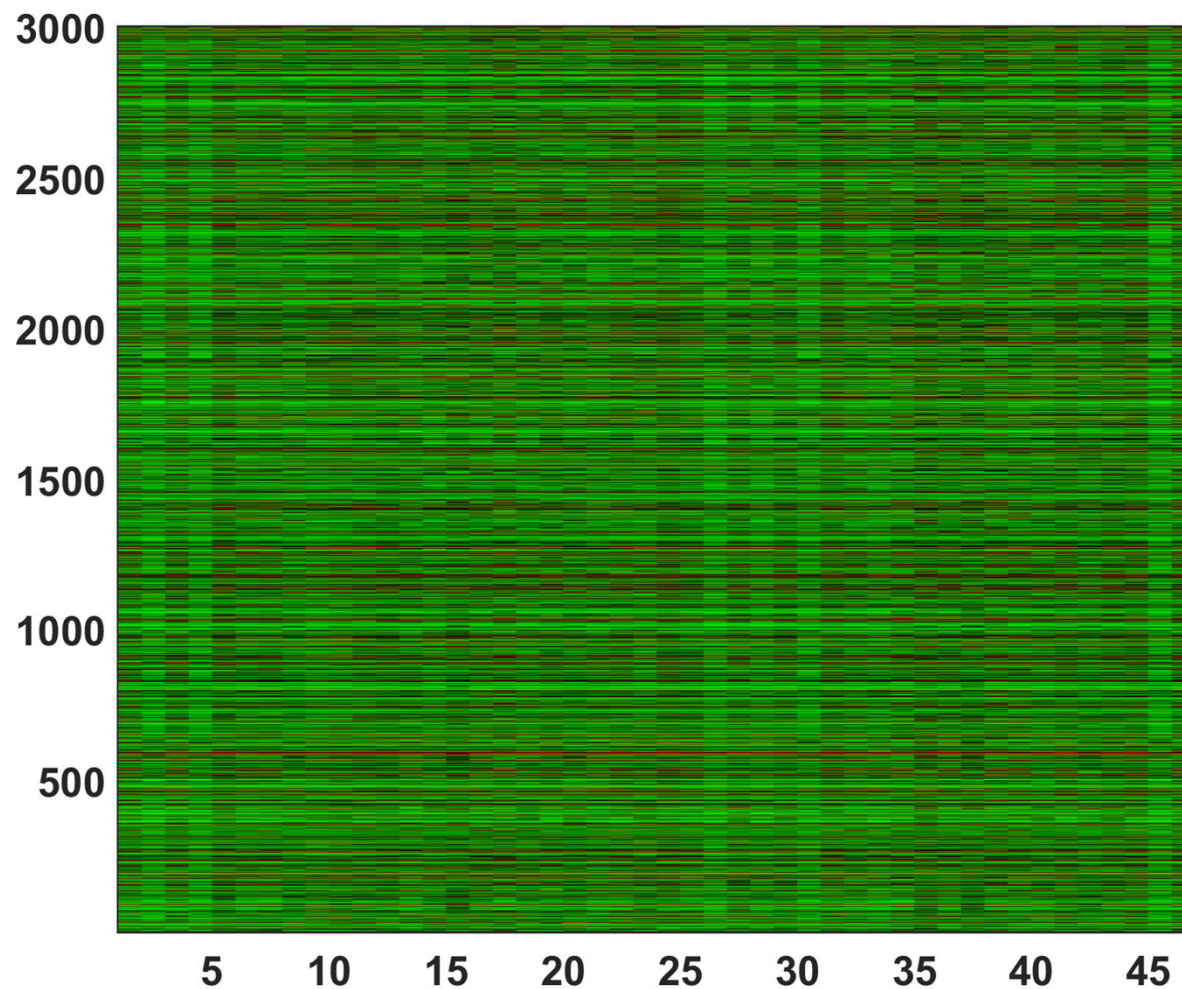
**no obvious biological common function**

```
g1=1+floor(rand.*3000); g2=1+floor(rand.*3000);  
disp([g1, g2])
```

# Matlab exercise

- Every group works with  
g0=2907; g1=1527; g2=2629; g3=2881;  
g4=1144; g5=1066;
- Compute **Multiple Linear Regression (MLR)**,  
where  $y = \text{exp\_t}(g0)$ ;  
 $x1 = \text{exp\_t}(g1)$ ;  $x2 = \text{exp\_t}(g2)$ ;
- **How much better** the MLR did compared to the  
Single Linear Regression (SLR)?
- **Continue increasing** the number of genes in x  
until **R\_adj** starts to decrease

How to find the entire groups of mutually correlated genes if you have **many genes** and **many samples**?



Clustering to the rescue!

# Clustering is a part of Machine Learning

- ***Supervised Learning:***

A machine learning technique whereby a system uses a set of human-labelled training examples to learn how to correctly perform a task

**Example:** a sample of cancer expression profiles each annotated with cancer type

**Goal:** predict cancer type based on expression pattern

- ***Unsupervised Learning (including clustering):***

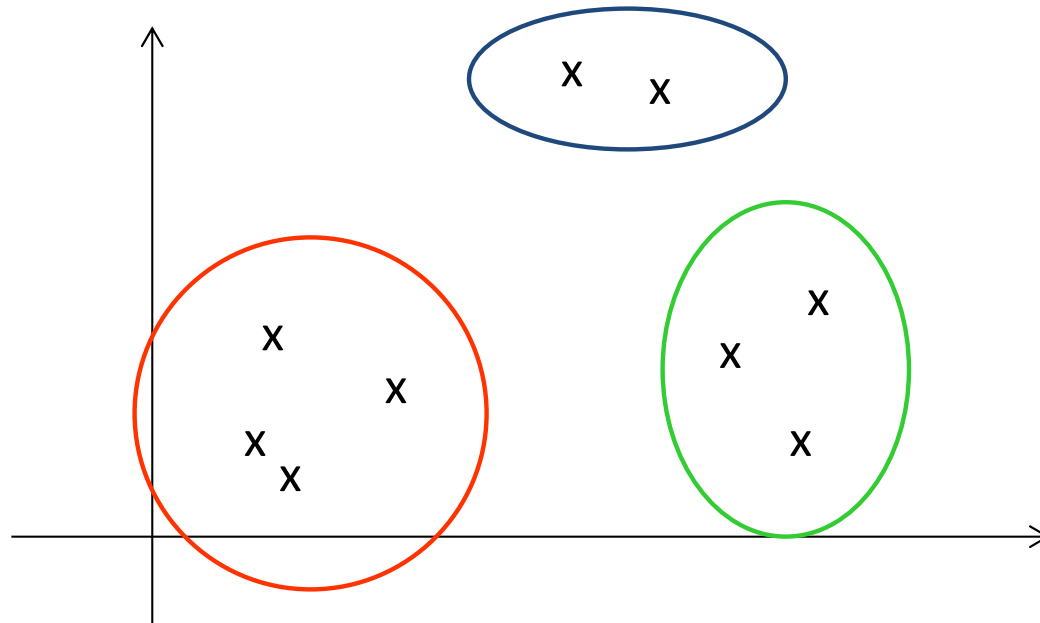
In machine learning, unsupervised learning is a class of problems in which one seeks to determine how the data are organized. One only has unlabeled examples.

**Example:** a sample of breast cancer expression profiles.

**Goal:** Identify several different (yet unknown) subtypes with potentially different treatments

# What is clustering?

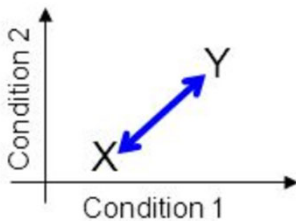
- The goal of **clustering** is to
  - group data points that are close (or **similar**) to each other
  - Usually, one needs to identify such groups (or clusters) in an **unsupervised** manner
  - Sometimes one takes into account **prior information** (Bayesian methods)
- Need to define some **distance  $d_{ij}$**  between **objects  $i$  and  $j$**
- Clustering is easy in **2 dimensions** but **hard in 3000 dimensions** -> need to somehow **reduce dimensionality**



# How to define the distance?

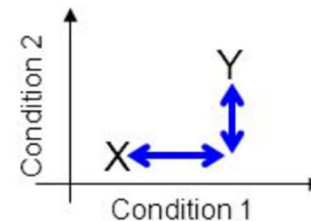
- Euclidean distance:
  - Most commonly used distance
  - Sphere shaped cluster
  - Corresponds to the geometric distance into the multidimensional space

$$d(X, Y) = \sqrt{\sum_i (x_i - y_i)^2}$$



- City Block (Manhattan) distance:
  - Sum of differences across dimensions
  - Less sensitive to outliers
  - Diamond shaped clusters

$$d(X, Y) = \sum_i |x_i - y_i|$$



The Canberra distance metric is calculated in R by

$$\sum \left( \frac{|x_i - y_i|}{|x_i + y_i|} \right).$$

Correlation coefficient distance

$$d(X, Y) = 1 - \rho(X, Y) = 1 - \frac{\text{Cov}(X, Y)}{\sqrt{\text{Var}(X) \cdot \text{Var}(Y)}}$$

# Common types of clustering algorithms

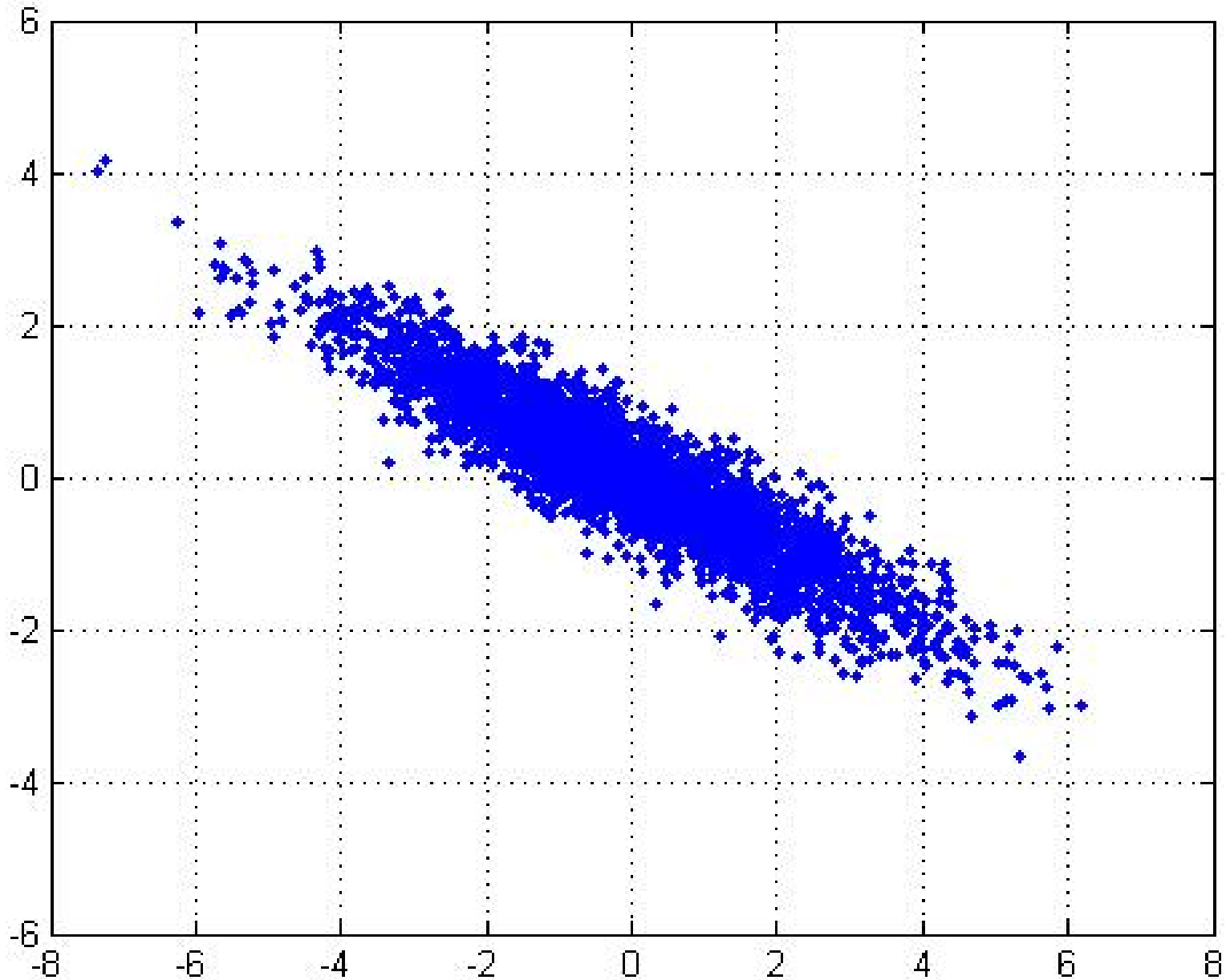
- Hierarchical if one doesn't know in advance the # of clusters
  - Agglomerative: start with  $N$  clusters and gradually merge them into 1 cluster
  - Divisive: start with 1 cluster and gradually break it up into  $N$  clusters
- Non-hierarchical algorithms
  - K-means clustering:
    - Iteratively apply the following two steps:
    - Calculate the centroid (center of mass) of each cluster
    - Assign each to the cluster to the nearest centroid
  - Principal Component Analysis (PCA)
    - plot pairs of top eigenvectors of the covariance matrix  $\text{Cov}(X_i, X_j)$  and uses visual information to group



# The Principal Components

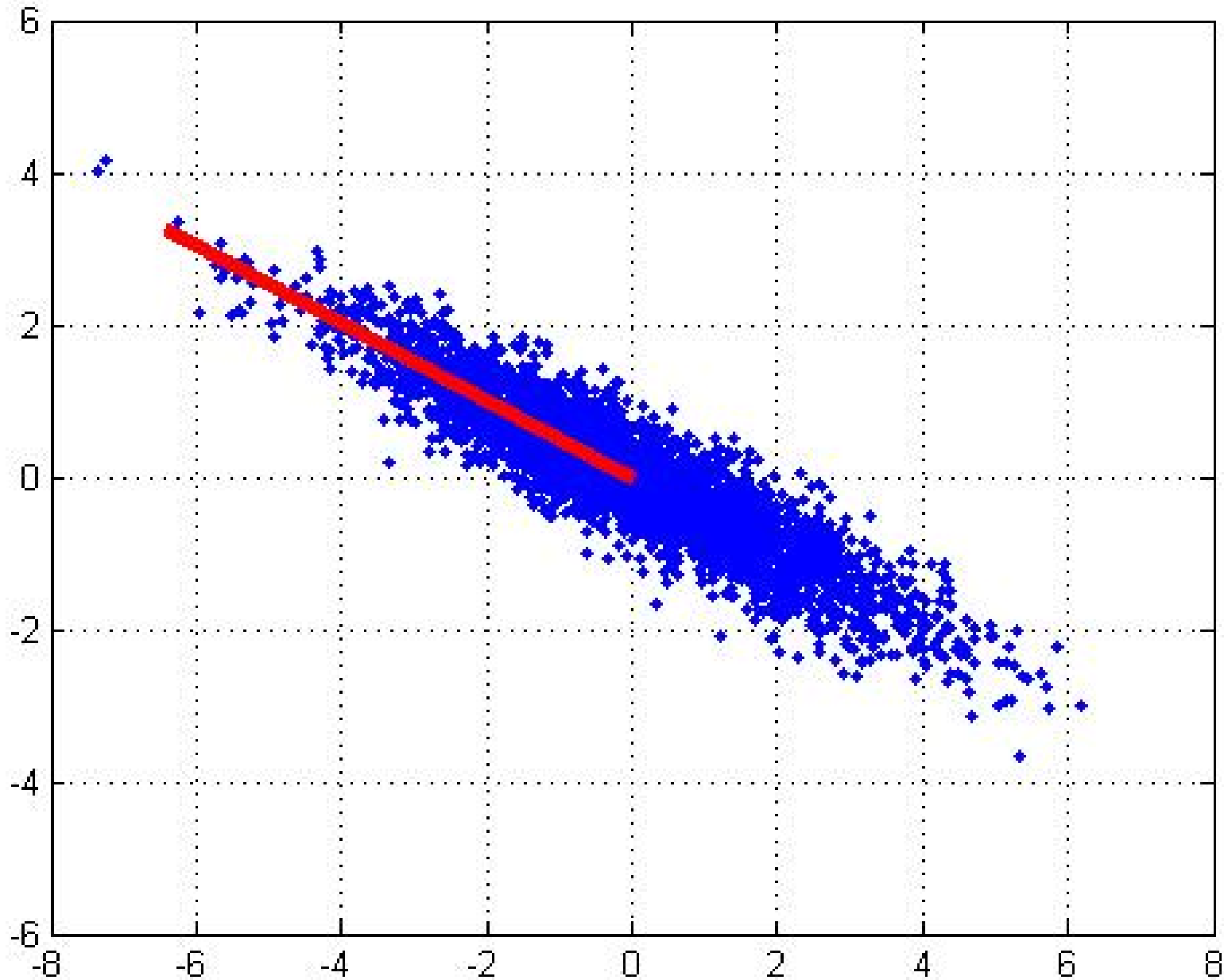
- **Vectors** originating from the center of mass
- Principal component #1 points in the direction of the **largest variance**.
- Each subsequent principal component...
  - is **orthogonal** to the previous ones, and
  - points in the directions of the **largest variance of the residual subspace**

# 2D Gaussian dataset



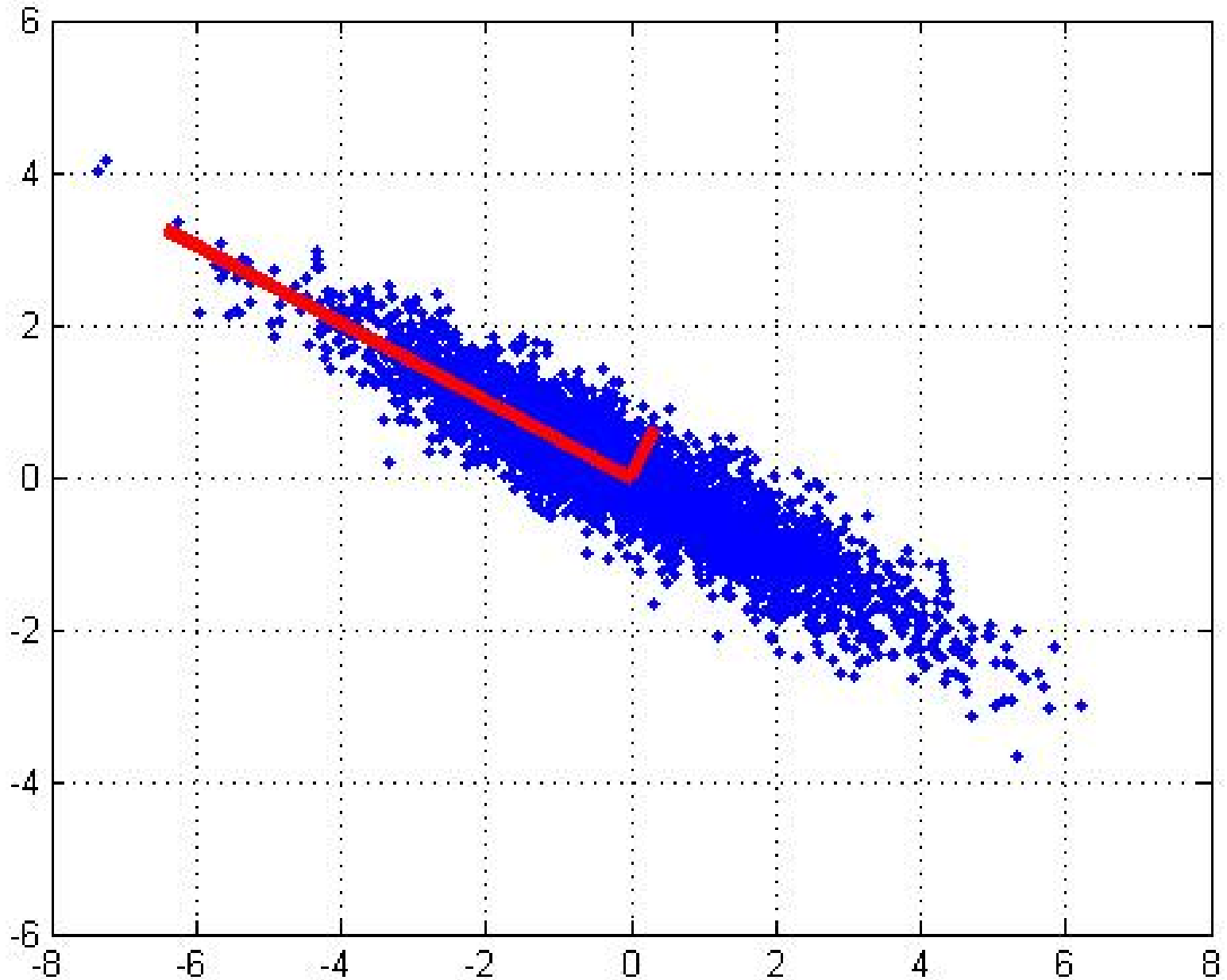
Adapted from lectures Prof. Pat Virtue at CMU based on original slide from Barnabas Poczos

# 1<sup>st</sup> PCA axis



Adapted from lectures Prof. Pat Virtue at CMU based on original slide from Barnabas Poczos

# 2<sup>nd</sup> PCA axis



Adapted from lectures Prof. Pat Virtue at CMU based on original slide from Barnabas Poczos

# Data for PCA

$$\mathcal{D} = \{\mathbf{x}^{(i)}\}_{i=1}^N \quad \mathbf{X} = \begin{bmatrix} (\mathbf{x}^{(1)})^T \\ (\mathbf{x}^{(2)})^T \\ \vdots \\ (\mathbf{x}^{(N)})^T \end{bmatrix}$$

We assume the data is **centered**

$$\mu = \frac{1}{N} \sum_{i=1}^N \mathbf{x}^{(i)} = \mathbf{0}$$

**Q:** What if your data is **not** centered?

**A:** Subtract off the sample mean

# Sample Covariance Matrix

The sample covariance matrix is given by:

$$\Sigma_{jk} = \frac{1}{N} \sum_{i=1}^N (x_j^{(i)} - \mu_j)(x_k^{(i)} - \mu_k)$$

Since the data matrix is centered, we rewrite as:

$$\Sigma = \frac{1}{N} \mathbf{X}^T \mathbf{X}$$

$$\mathbf{X} = \begin{bmatrix} (\mathbf{x}^{(1)})^T \\ (\mathbf{x}^{(2)})^T \\ \vdots \\ (\mathbf{x}^{(N)})^T \end{bmatrix}$$

# PCA algorithm

PCA algorithm( $\mathbf{X}$ ,  $k$ ): top  $k$   
eigenvalues/eigenvectors

- $\{ \lambda_i, \mathbf{u}_i \}_{i=1:m}$  = eigenvectors/eigenvalues of  $\Sigma$   
...  $\lambda_1 \geq \lambda_2 \geq \dots \geq \lambda_m$
- **PCA** basis vectors = the eigenvectors of  $\Sigma$
- Larger eigenvalue  $\Rightarrow$  more important  
eigenvectors

# PCA and units

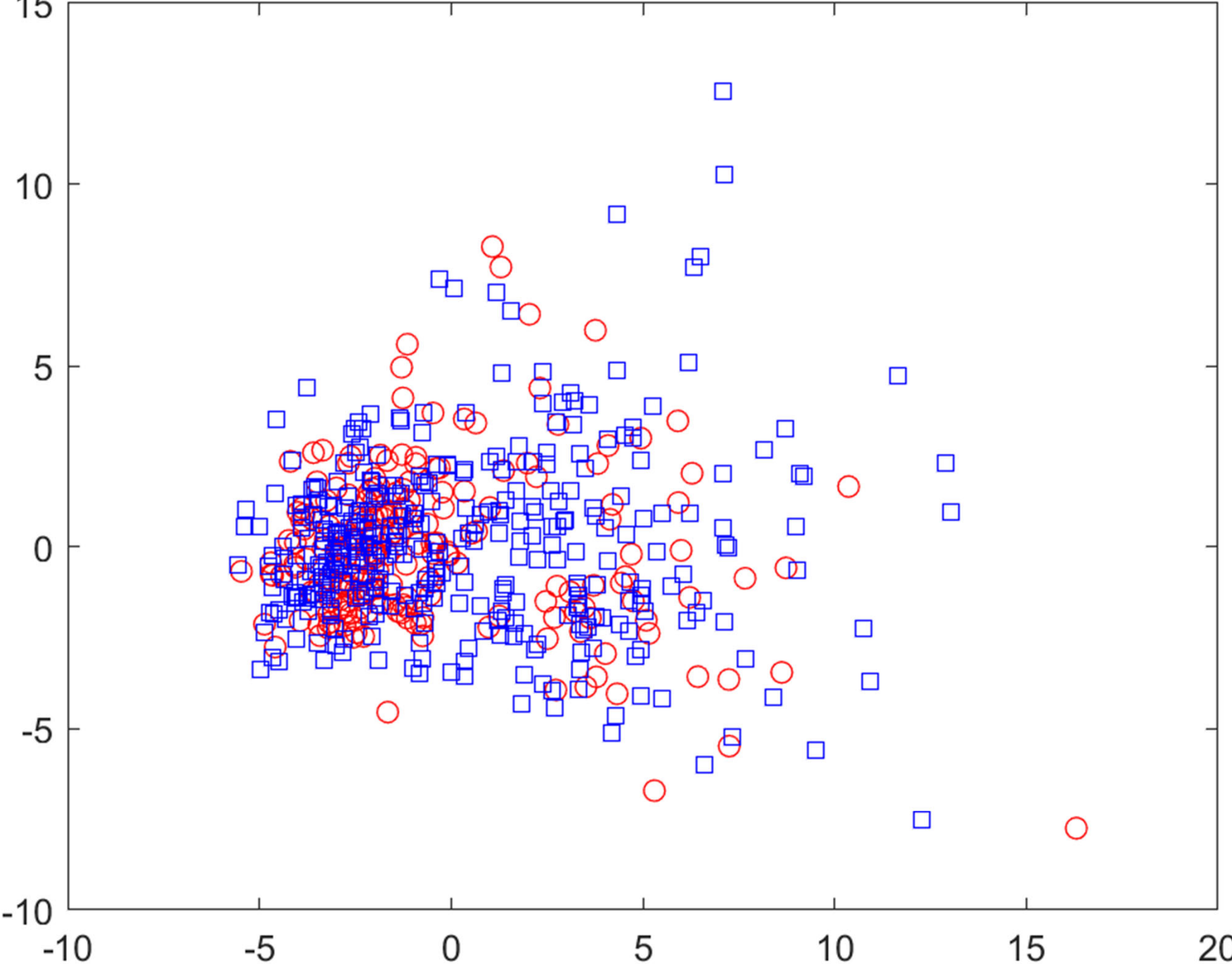
- When different variables have different units (like temperature and mass), the meaning of principal components is a somewhat arbitrary
- One way of making the PCA less arbitrary is to use variables scaled so as to have unit variance, by standardizing the data
- Before making PCA of  $X$  transform it using  $Z = \text{zscore}(X)$ ;



# Group project 4

- load cancer\_wdbc.mat
- `Z=zscore(cancerwdbc);`
- `[coeff_z, score_z, latent_z] = pca(Z);`
- `ic=find(cancer_yn==1); whos ic;`  
`inc=find(cancer_yn==0); whos inc;`
- `figure; plot(score_z(ic,1), score_z(ic,2),'ro'); hold on;`  
`plot(score_z(inc,1), score_z(inc,2),'bs');`  
`title('PC2 vs PC1');`
- Plot pairs of `score_z` components
  - 1<sup>st</sup> principal component vs 2<sup>nd</sup> principal component.
  - 1<sup>st</sup> principal component vs 3<sup>rd</sup> principal component
  - 3<sup>rd</sup> principal component vs 2<sup>nd</sup> principal component

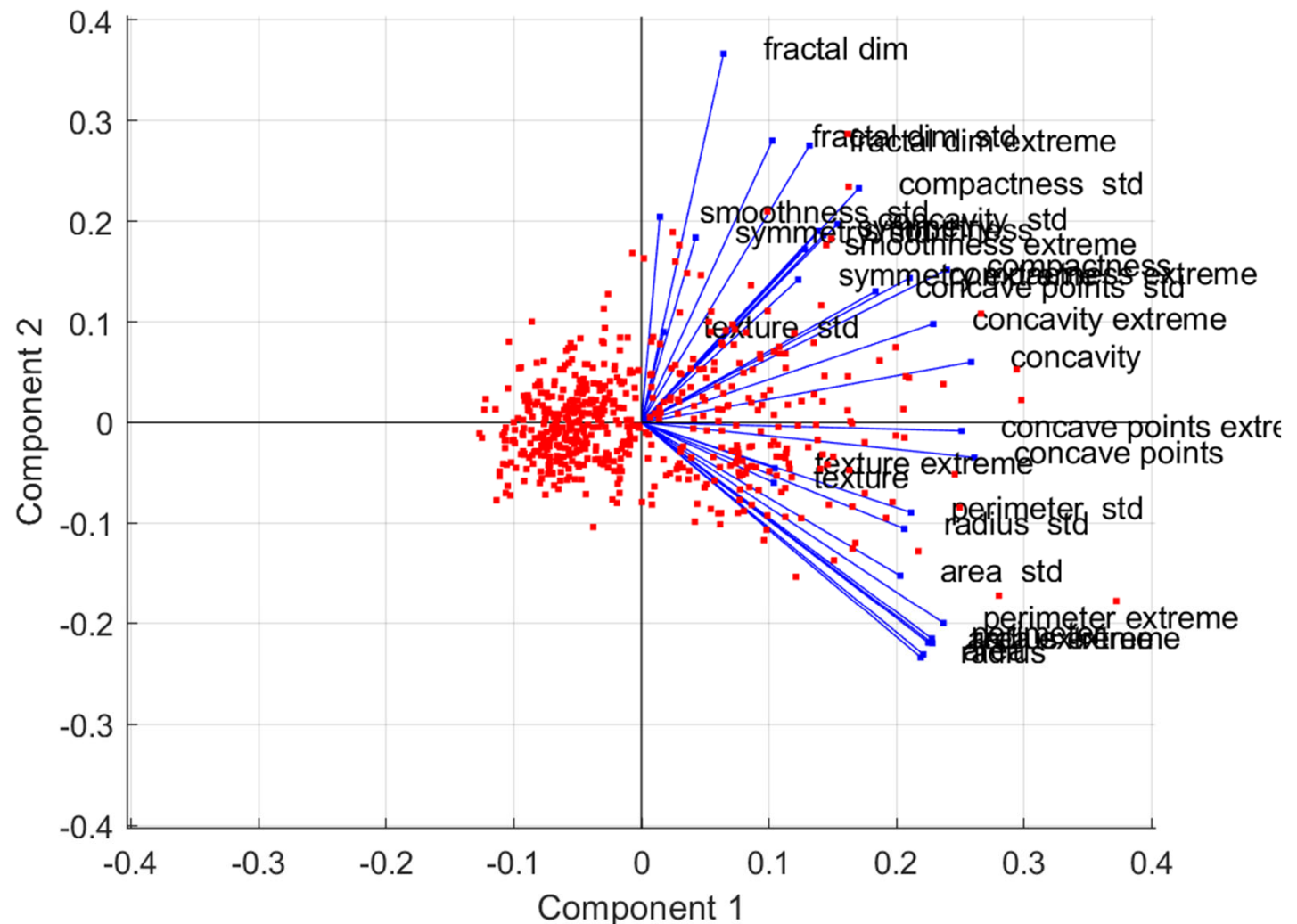
**PC2 vs PC1**



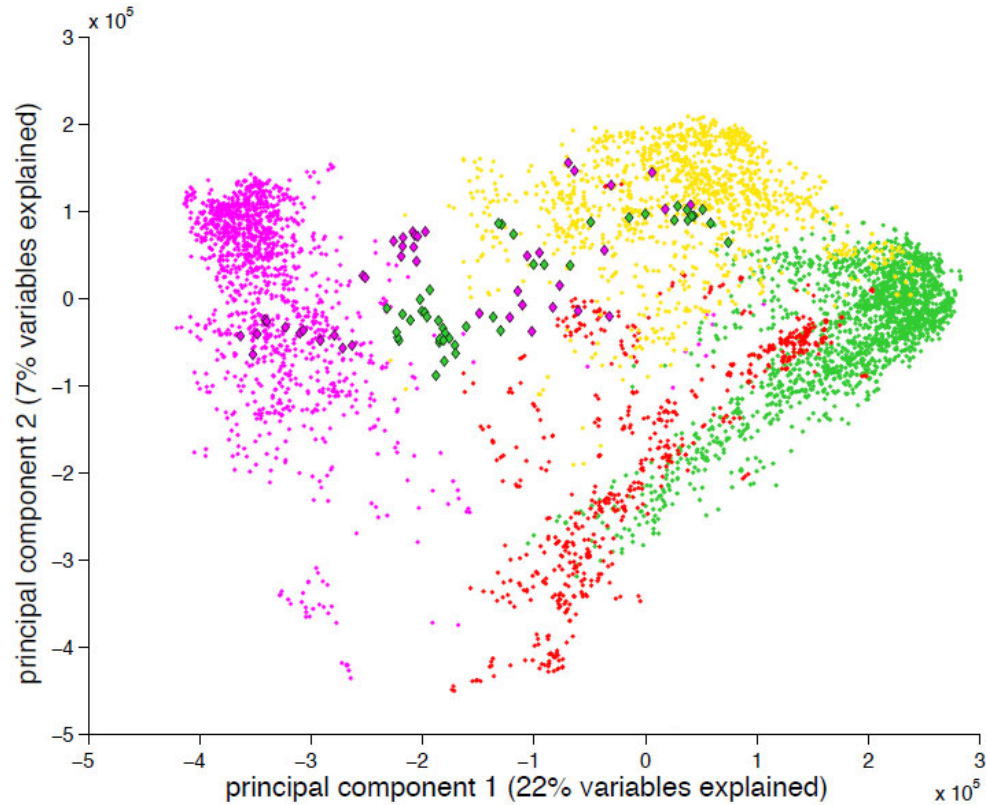
# Which variables contribute to which PC?

## Add loadings (coeff eigenvectors)

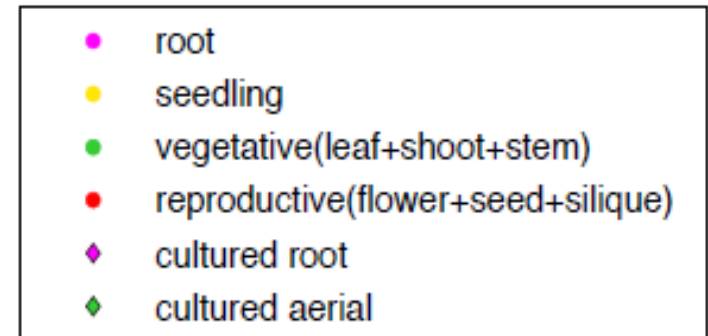
- `figure; biplot(coeff_z(:,1:2),'scores',score_z(:,1:2), 'VarLabels' feature names);`



# Example of Principal Component Analysis (PCA) clustering



7000 gene expression  
samples of model plant  
*Arabidopsis thaliana*



[Plant J.](#) 2016 Mar 25. doi: 10.1111/tpj.13175. [Epub ahead of print]

**Large-scale atlas of microarray data reveals the distinct expression landscape of different tissues in Arabidopsis.**

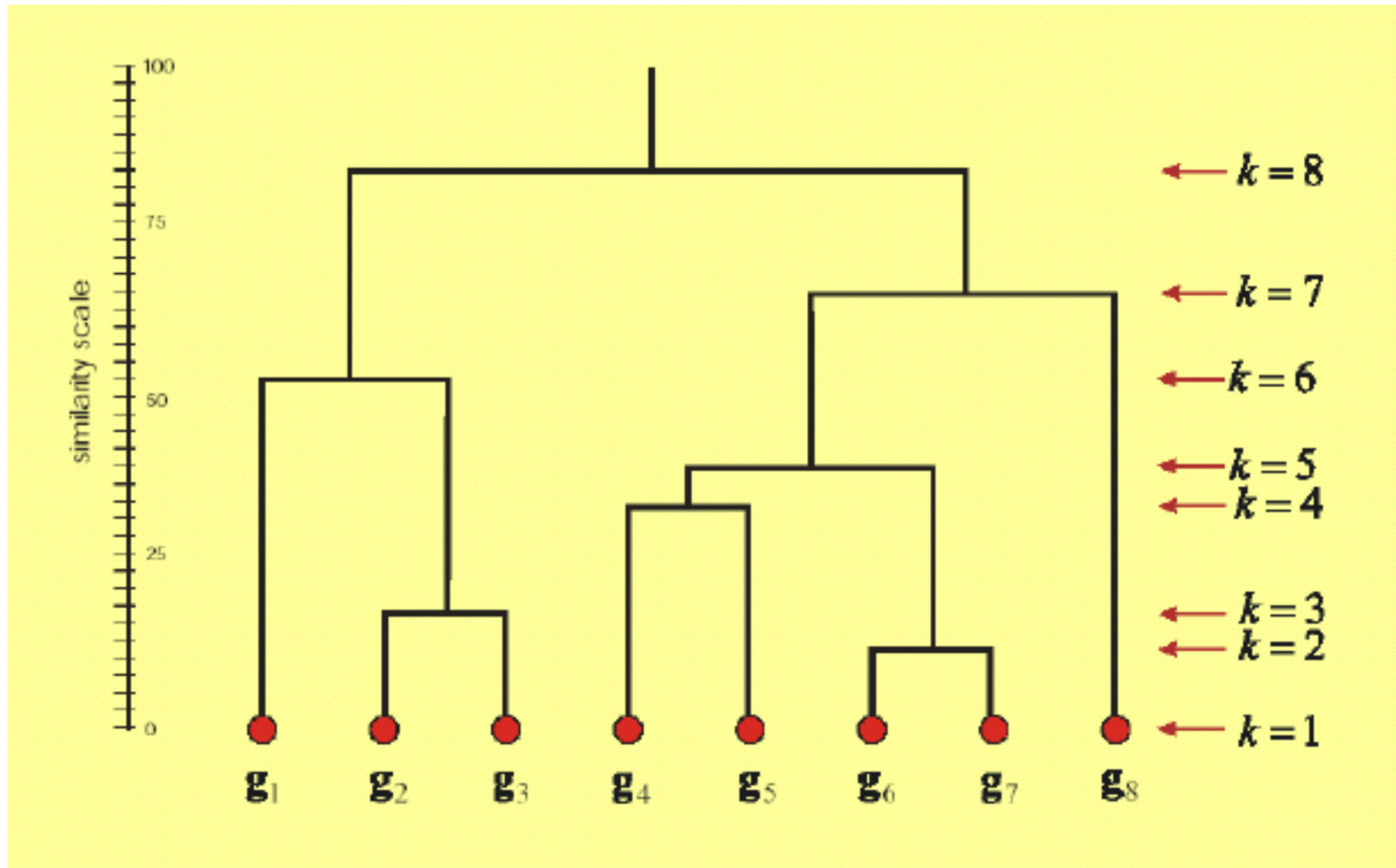
[He F](#)<sup>1</sup>, [Yoo S](#)<sup>2,3</sup>, [Wang D](#)<sup>4</sup>, [Kumari S](#)<sup>5</sup>, [Gerstein M](#)<sup>4</sup>, [Ware D](#)<sup>5,6</sup>, [Maslov S](#)<sup>1,7</sup>.

# Hierarchical clustering

# UPGMA algorithm

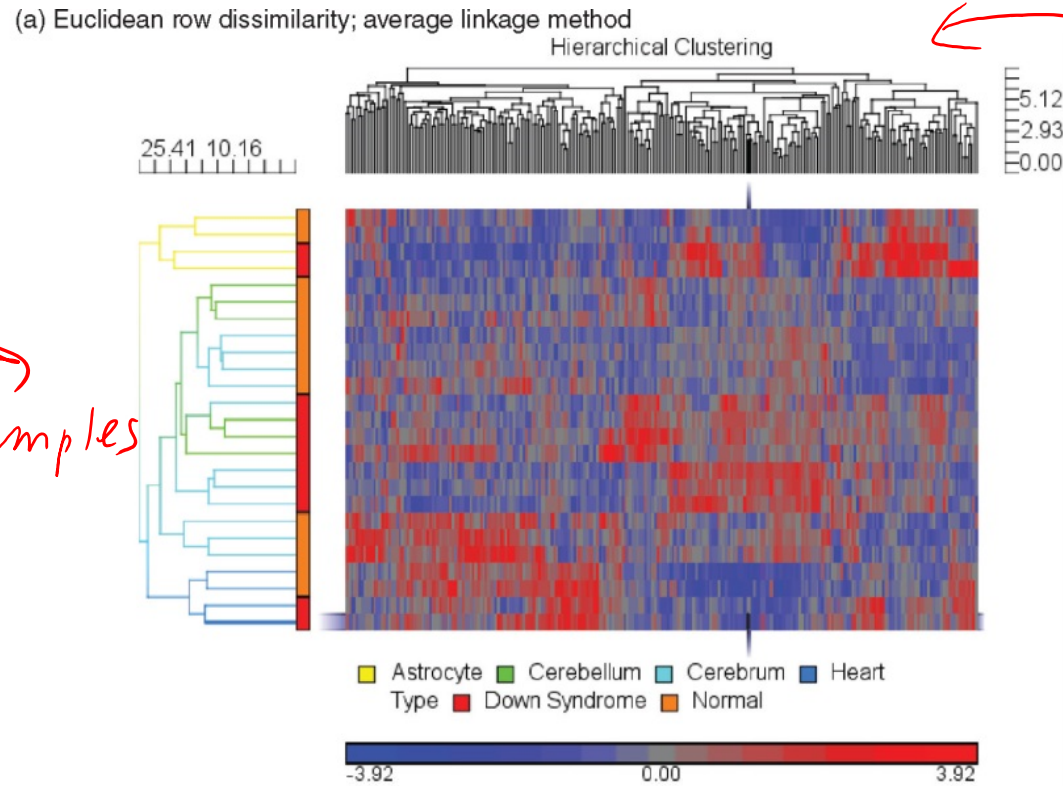
- Hierarchical agglomerative clustering algorithm
- **UPGMA** = **U**nweighted **P**air **G**roup **M**ethod with **A**rithmetic mean
- **Iterative** algorithm:
- Start with a **pair with the smallest  $d(X,Y)$**
- **Cluster these two together** and replace it with their arithmetic mean  $(X+Y)/2$
- **Recalculate all distances to this new “cluster node”**
- **Repeat** until all nodes are merged

# Output of UPGMA algorithm



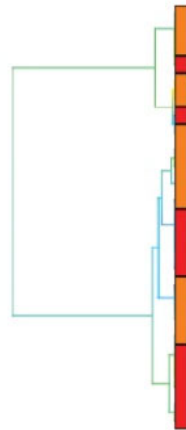
UPGMA  
algorithm

25 samples

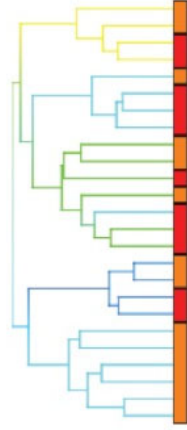


250 genes  
on  
chromosome  
21

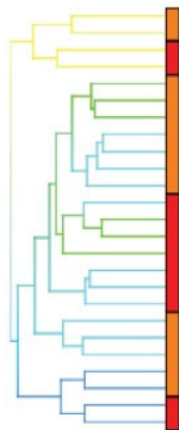
(b) Canberra  
dissimilarity



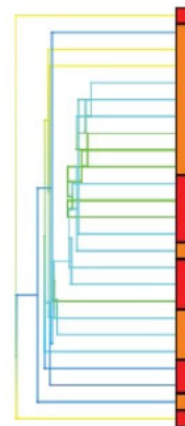
(c) Pearson's  
Dissimilarity



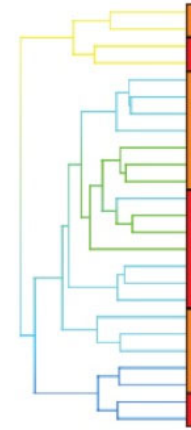
(d) City Block



(e) Euclidean,  
centroid linkage



(f) Euclidean,  
complete-linkage



**FIGURE 11.16** Hierarchical clustering of 250 chromosome 21 transcripts in 25 samples using Partek software. (a) Hierarchical clustering of microarray data using the default settings of Euclidean dissimilarity for rows (samples) and columns (transcripts). Colors correspond to expression intensity values.



Credit: XKCD  
comics

WHY ARE THERE SLAVES IN THE BIBLE

WHY DO TWINS HAVE DIFFERENT FINGERPRINTS  
WHY ARE AMERICANS AFRAID OF DRAGONS

WHY IS HTTPS CROSSED OUT IN RED  
WHY IS THERE A LINE THROUGH HTTPS  
WHY IS THERE A RED LINE THROUGH HTTPS ON FACEBOOK  
WHY IS HTTPS IMPORTANT

# QUESTIONS

FOUND IN GOOGLE AUTOCOMPLETE



WHY ARE THERE WEEKS  
WHY DO I FEEL DIZZY

WHY AREN'T ECONOMISTS RICH

WHY ARE THERE SO MANY CROWS IN ROCHESTER, MN  
WHY IS THERE PHLEGM

WHY DO AMERICANS CALL IT SOCCER

WHY IS PSYCHIC WEAK TO BUG

WHY ARE MY EARS RINGING

WHY DO CHILDREN GET CANCER

WHY ARE THERE SO MANY AVENGERS

WHY IS POSEIDON ANGRY WITH ODYSSEUS

WHY ARE THE AVENGERS FIGHTING THE X MEN  
WHY IS WOLVERINE NOT IN THE AVENGERS

WHY IS THERE ICE IN SPACE

## WHY ARE THERE ANTS IN MY LAPTOP

WHY IS EARTH TILTED  
WHY IS SPACE BLACK  
WHY IS OUTER SPACE SO COLD  
WHY ARE THERE PYRAMIDS ON THE MOON  
WHY IS NASA SHUTTING DOWN



WHY IS THERE AN OWL IN MY BACKYARD  
WHY IS THERE AN OWL OUTSIDE MY WINDOW  
WHY IS THERE AN OWL ON THE DOLLAR BILL  
WHY DO OWLS ATTACK PEOPLE

WHY ARE AK 47s SO EXPENSIVE

WHY ARE THERE HELICOPTERS CIRCLING MY HOUSE

WHY ARE THERE GODS

WHY ARE THERE TWO SPOCKS

WHY IS MT VESUVIUS THERE

WHY DO THEY SAY T MINUS

WHY ARE THERE OBELISKS

WHY ARE WRESTLERS ALWAYS WET

WHY ARE OCEANS BECOMING MORE ACIDIC

WHY IS ARWEN DYING

WHY AREN'T MY QUAIL LAYING EGGS  
WHY AREN'T MY QUAIL EGGS HATCHING  
WHY AREN'T THERE ANY FOREIGN MILITARY BASES IN AMERICA

WHY ARE CIGARETTES LEGAL  
WHY ARE THERE DUCKS IN MY POOL  
WHY IS JESUS WHITE  
WHY IS THERE LIQUID IN MY EAR  
WHY DO Q TIPS FEEL GOOD  
WHY DO GOOD PEOPLE DIE



WHY ARE ULTRASOUNDS IMPORTANT  
WHY ARE ULTRASOUND MACHINES EXPENSIVE  
WHY IS STEALING WRONG

WHY ARE DOGS AFRAID OF FIREWORKS  
WHY IS THERE NO KING IN ENGLAND

WHY DO WHALES JUMP  
WHY ARE WITCHES GREEN  
WHY ARE THERE MIRRORS ABOVE BEDS

WHY DO I SAY UH  
WHY IS SEA SALT BETTER

WHY ARE THERE TREES IN THE MIDDLE OF FIELDS  
WHY IS THERE NOT A POKEMON MMO  
WHY IS THERE LAUGHING IN TV SHOWS  
WHY ARE THERE DOORS ON THE FREEWAY  
WHY ARE THERE SO MANY SVCHOST.EXE RUNNING  
WHY AREN'T THERE ANY COUNTRIES IN ANTARCTICA  
WHY ARE THERE SCARY SOUNDS IN MINECRAFT  
WHY IS THERE KICKING IN MY STOMACH  
WHY ARE THERE TWO SLASHES AFTER HTTP

WHY ARE THERE CELEBRITIES  
WHY DO SNAKES EXIST

WHY DO OYSTERS HAVE PEARLS  
WHY ARE DUCKS CALLED DUCKS  
WHY DO THEY CALL IT THE CLAP  
WHY ARE KYLE AND CARTMAN FRIENDS  
WHY IS THERE AN ARROW ON AANG'S HEAD  
WHY ARE TEXT MESSAGES BLUE

WHY ARE THERE MUSTACHES ON CLOTHES  
WHY ARE THERE MUSTACHES ON CARS  
WHY ARE THERE MUSTACHES EVERYWHERE  
WHY ARE THERE SO MANY BIRDS IN OHIO  
WHY IS THERE SO MUCH RAIN IN OHIO  
WHY IS OHIO WEATHER SO WEIRD

## WHY ARE THERE MALE AND FEMALE BIKES

WHY ARE THERE BRIDESMAIDS  
WHY DO DYING PEOPLE REACH UP  
WHY AREN'T THERE VARICOSE ARTERIES  
WHY ARE OLD KUNGONS DIFFERENT

WHY ARE THERE SQUIRRELS



WHY IS PROGRAMMING SO HARD  
WHY IS THERE A 0 OHM RESISTOR  
WHY DO AMERICANS HATE SOCCER  
WHY DO RHYMES SOUND GOOD  
WHY DO TREES DIE  
WHY IS THERE NO SOUND ON CNN  
WHY AREN'T POKEMON REAL  
WHY AREN'T BULLETS SHARP  
WHY DO DREAMS SEEM SO REAL

WHY ARE THERE TINY SPIDERS IN MY HOUSE  
WHY DO SPIDERS COME INSIDE  
WHY ARE THERE HUGE SPIDERS IN MY HOUSE  
WHY ARE THERE LOTS OF SPIDERS IN MY HOUSE  
WHY ARE THERE SPIDERS IN MY ROOM  
WHY ARE THERE SO MANY SPIDERS IN MY ROOM

WHY DO SPIDER BITES ITCH

WHY IS DYING SO SCARY

WHY IS THERE NO GPS IN LAPTOPS

WHY DO KNEES CLICK

WHY AREN'T THERE E GRADES

WHY IS ISOLATION BAD

WHY DO BOYS LIKE ME  
WHY DON'T BOYS LIKE ME  
WHY IS THERE ALWAYS A JAVA UPDATE  
WHY ARE THERE RED DOTS ON MY THIGHS

WHY IS SEX SO IMPORTANT



WHY AREN'T THERE DINOSAUR GHOSTS

WHY ARE THERE FEMALE MR NIMES

WHY IS THERE HELL IF GOD FORGIVES

WHY IS GPS FREE

# Clustering

## Matlab demo

# Choices of distance metrics in `clustergram(... 'RowPDistValue' ...,` `'ColumnPDistValue' ...)`

Metric	Description
'euclidean'	Euclidean distance (default).
'seuclidean'	Standardized Euclidean distance. Each coordinate difference between rows in X is scaled by dividing by the corresponding element of the standard deviation <code>S=nansd(X)</code> . To specify another value for S, use <code>D=pdist(X, 'seuclidean', S)</code> .
'cityblock'	City block metric.
'minkowski'	Minkowski distance. The default exponent is 2. To specify a different exponent, use <code>D = pdist(X, 'minkowski', P)</code> , where P is a scalar positive value of the exponent.
'chebychev'	Chebychev distance (maximum coordinate difference).
'mahalanobis'	Mahalanobis distance, using the sample covariance of X as computed by <code>nancov</code> . To compute the distance with a different covariance, use <code>D = pdist(X, 'mahalanobis', C)</code> , where the matrix C is symmetric and positive definite.
'cosine'	One minus the cosine of the included angle between points (treated as vectors).
'correlation'	One minus the sample correlation between points (treated as sequences of values).
'spearman'	One minus the sample Spearman's rank correlation between observations (treated as sequences of values).
'hamming'	Hamming distance, which is the percentage of coordinates that differ.
'jaccard'	One minus the Jaccard coefficient, which is the percentage of nonzero coordinates that differ.
custom distance function	A distance function specified using @: <code>D = pdist(X, @distfun)</code>  A distance function must be of form  <code>d2 = distfun(XI, XJ)</code>  taking as arguments a 1-by-n vector XI, corresponding to a single row of X, and an m2-by-n matrix XJ, corresponding to multiple rows of X. <code>distfun</code> must accept a matrix XJ with an arbitrary number of rows. <code>distfun</code> must return an m2-by-1 vector of distances d2, whose kth element is the distance between XI and XJ(k, :).

# Choices of hierarchical clustering algorithm in `clustergram( ...'linkage',...)`

X	Matrix with two or more rows. The rows represent observations, the columns represent categories or dimensions.																
method	<p>Algorithm for computing distance between clusters.</p> <table border="1"><thead><tr><th>Method</th><th>Description</th></tr></thead><tbody><tr><td>'average'</td><td>Unweighted average distance (UPGMA)</td></tr><tr><td>'centroid'</td><td>Centroid distance (UPGMC), appropriate for Euclidean distances only</td></tr><tr><td>'complete'</td><td>Furthest distance</td></tr><tr><td>'median'</td><td>Weighted center of mass distance (WPGMC), appropriate for Euclidean distances only</td></tr><tr><td>'single'</td><td>Shortest distance</td></tr><tr><td>'ward'</td><td>Inner squared distance (minimum variance algorithm), appropriate for Euclidean distances only</td></tr><tr><td>'weighted'</td><td>Weighted average distance (WPGMA)</td></tr></tbody></table> <p><b>Default:</b> 'single'</p>	Method	Description	'average'	Unweighted average distance (UPGMA)	'centroid'	Centroid distance (UPGMC), appropriate for Euclidean distances only	'complete'	Furthest distance	'median'	Weighted center of mass distance (WPGMC), appropriate for Euclidean distances only	'single'	Shortest distance	'ward'	Inner squared distance (minimum variance algorithm), appropriate for Euclidean distances only	'weighted'	Weighted average distance (WPGMA)
Method	Description																
'average'	Unweighted average distance (UPGMA)																
'centroid'	Centroid distance (UPGMC), appropriate for Euclidean distances only																
'complete'	Furthest distance																
'median'	Weighted center of mass distance (WPGMC), appropriate for Euclidean distances only																
'single'	Shortest distance																
'ward'	Inner squared distance (minimum variance algorithm), appropriate for Euclidean distances only																
'weighted'	Weighted average distance (WPGMA)																

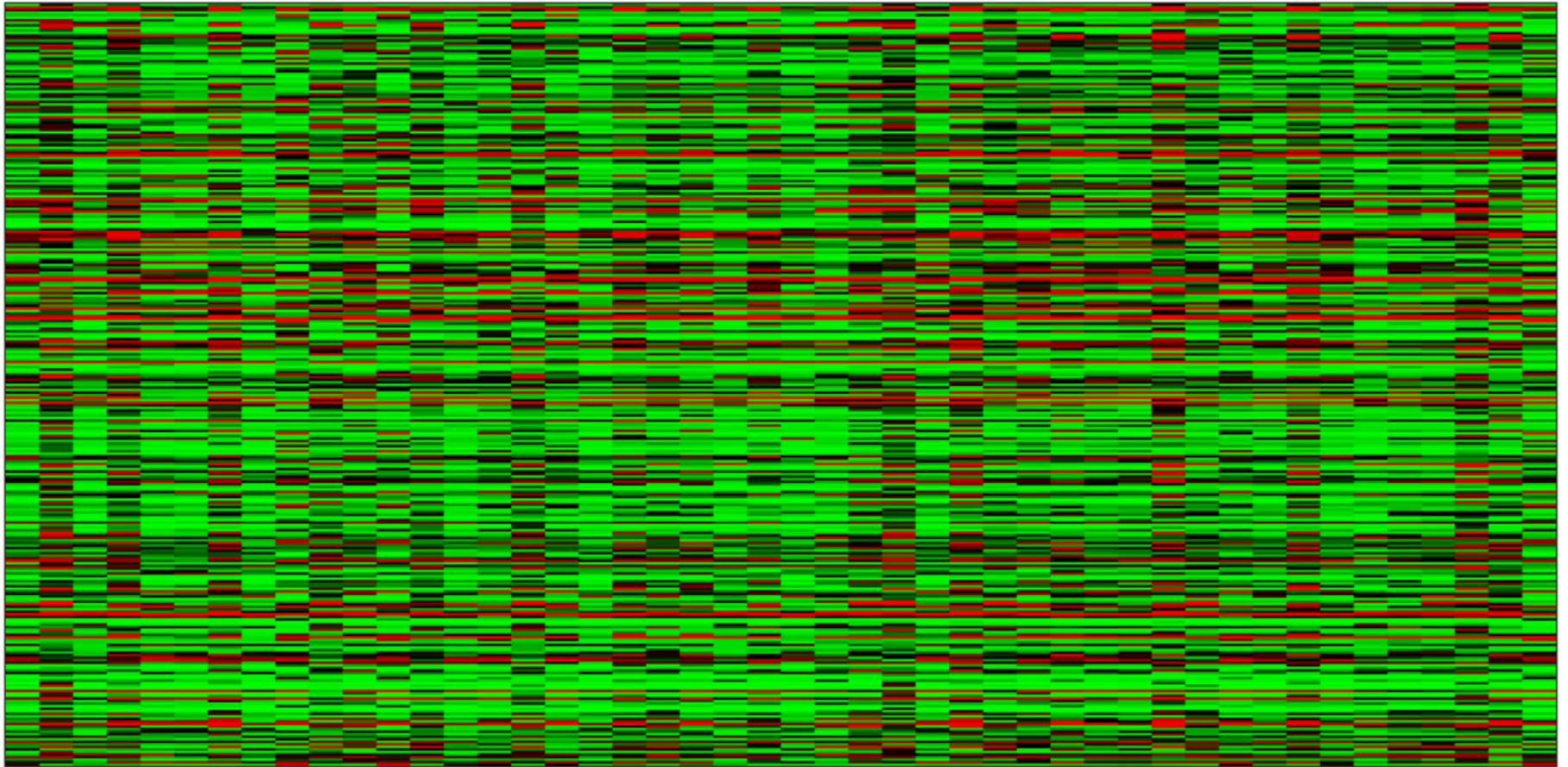
# Clustering group exercise

- Each group will analyze a **cluster of genes** identified in the T cell expression table
- Analyze the table of **top 100 genes by variance** in 47 samples
- Cluster them using:
  - Group 1: UPGMA = 'linkage', 'average', 'RowPDistValue', 'euclidean',
  - Group 2: 'linkage', 'single', 'RowPDistValue', 'cityblock',
  - Group 3: 'linkage', 'average', 'RowPDistValue', 'correlation',
  - Group 4: UPGMA = 'linkage', 'single', 'RowPDistValue', 'euclidean',
  - Group 5: UPGMA = 'linkage', 'weighted', 'RowPDistValue', 'correlation',
- Use clustergram(..., 'Standardize','Row',  
'linkage', as specified for your group,  
'RowPDistValue' as specified for your group,  
'RowLabels',gene\_names1,'ColumnLabels', array\_names)

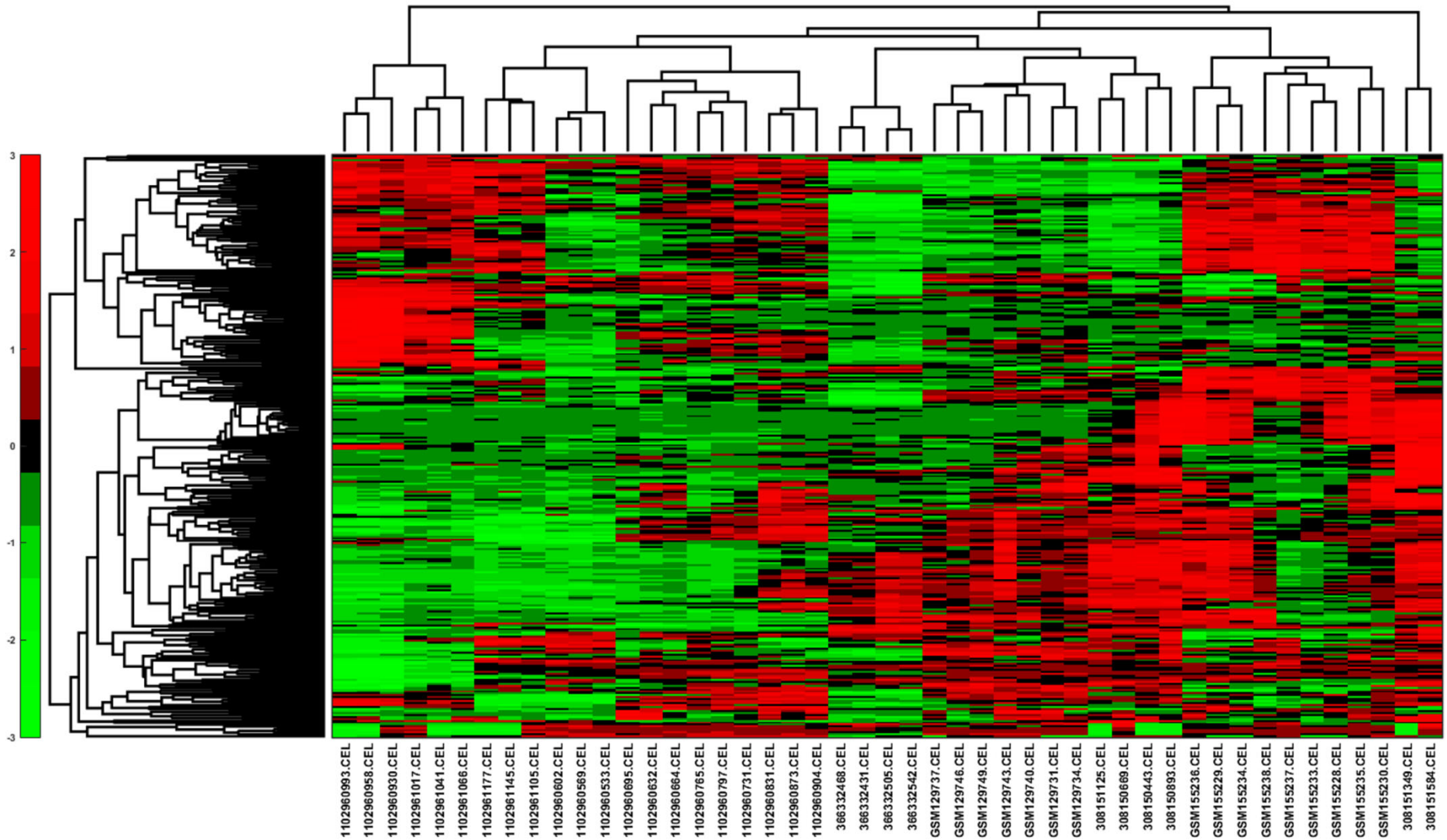
# Matlab code

- `load expression_table.mat`
- `gene_variation=std(exp_t)';`
- `[a,b]=sort(gene_variation,'descend');`
- `ngenes=100;`
- `exp_t1=exp_t(b(1:ngenes),:);`
- `gene_names1=gene_names(b(1:ngenes));`
- **%%% for group 1**
- `CGobj1 = clustergram(exp_t1, 'Standardize','Row',  
'RowLabels',  
gene_names1,'ColumnLabels',array_names)`
- `set(CGobj1,'RowLabels',gene_names1,'ColumnLabels',  
array_names,'linkage',  
'average','RowPDist','euclidean');`

Before clustering



# UPGMA hierarchical clustering, Euclidian distance





# UPGMA hierarchical clustering, correlation distance

