

Reminder:

Multiple Linear Regression

Test-train data split to
avoid overfitting

Multiple Linear Regression Model

$$Y = \beta_0 + \beta_1 x_1 + \beta_2 x_2 + \beta_3 x_3 + \dots + \beta_k x_k + \varepsilon$$

One can also use powers and products of other variables or even non-linear functions like $\exp(x_i)$ or $\log(x_i)$

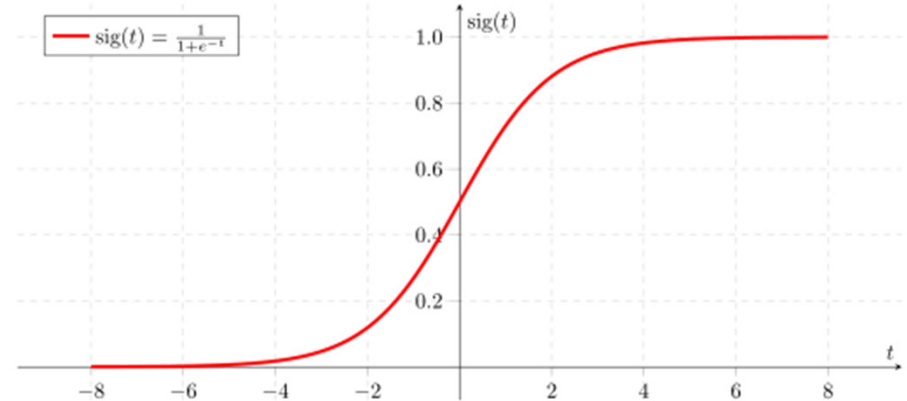
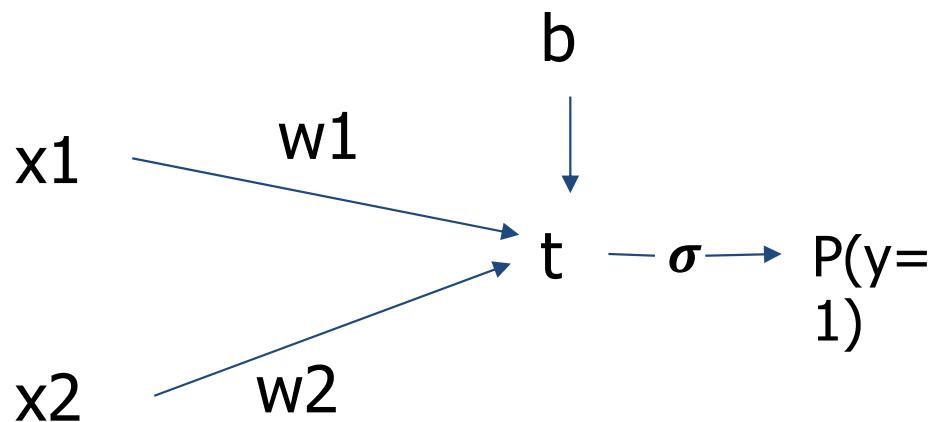
instead of x_3, \dots, x_k .

Example: the general two-variable quadratic regression has 6 constants:

$$Y = \beta_0 + \beta_1 x_1 + \beta_2 x_2 + \beta_3 (x_1)^2 + \beta_4 (x_2)^2 + \beta_5 (x_1 x_2) + \varepsilon$$

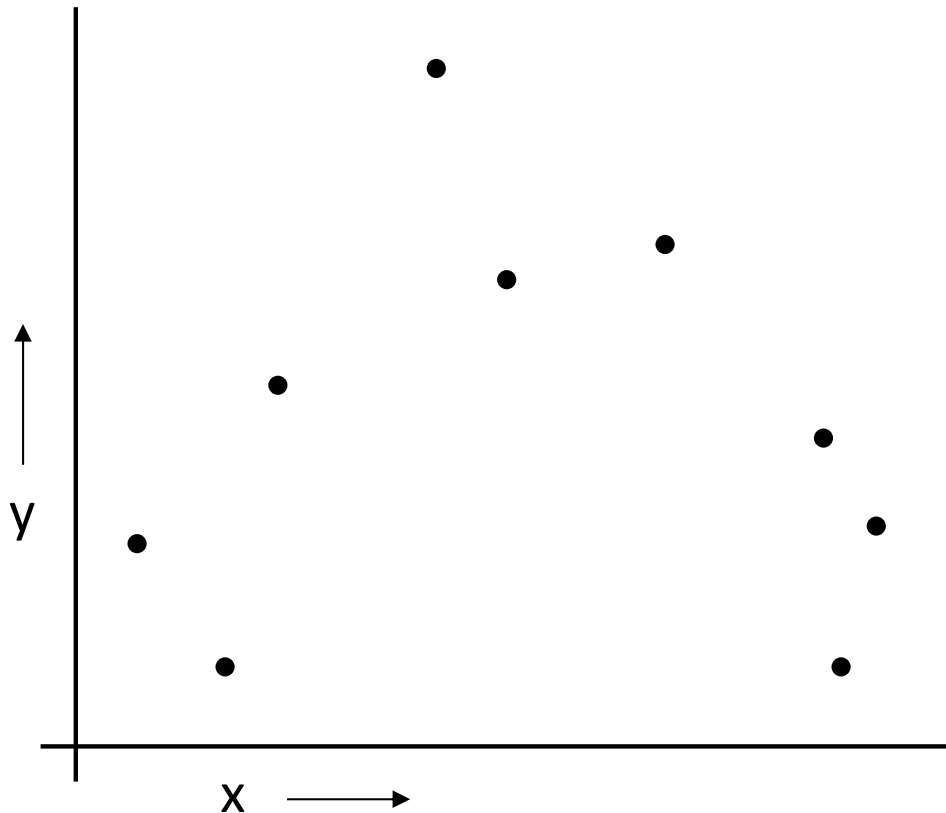
Logistic Regression

$$P(y=1) = \sigma(x_1 * w_1 + x_2 * w_2 + b)$$



How to know when to stop
adding new variables
or model parameters
in any data fitting algorithm such as
multiple linear regression?

A Regression Problem

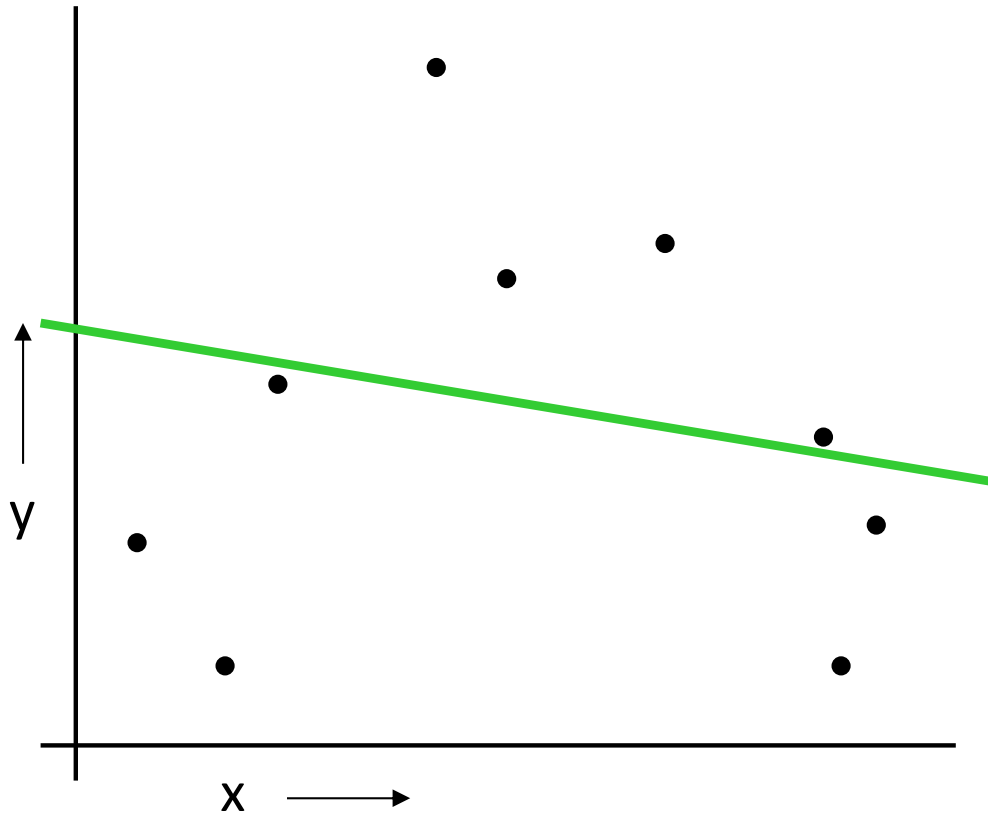


$$y = f(x) + \text{noise}$$

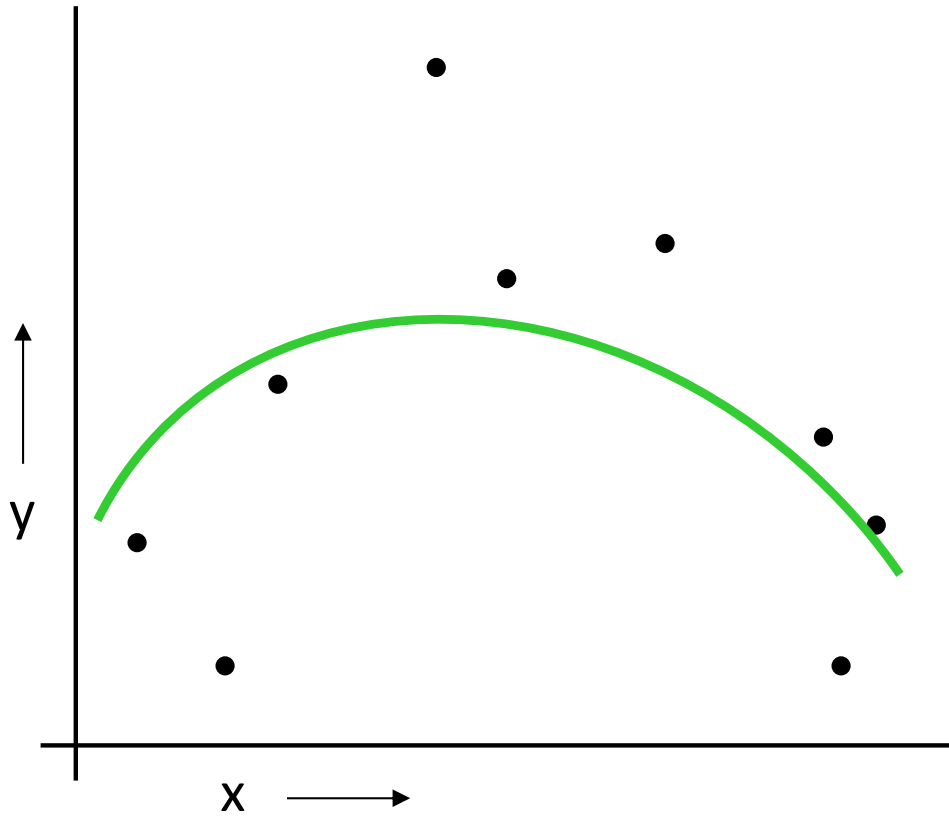
Can we learn f from this data?

Let's consider three methods...

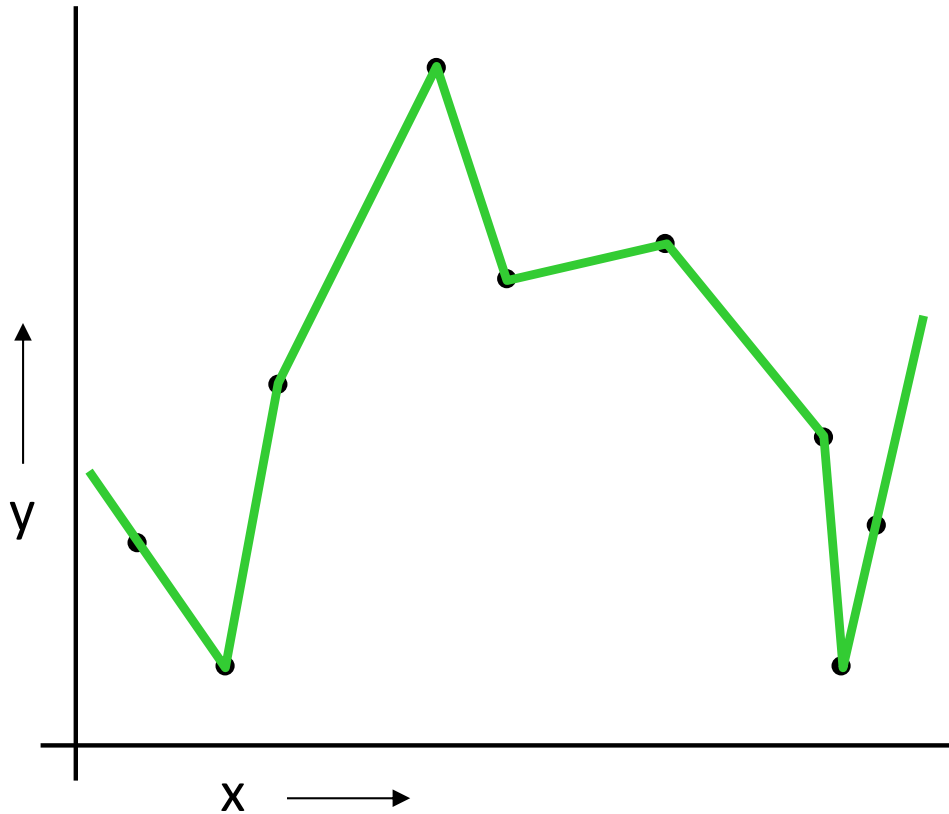
Single Variable Linear Regression



2-variable Linear Regression with x and x^2

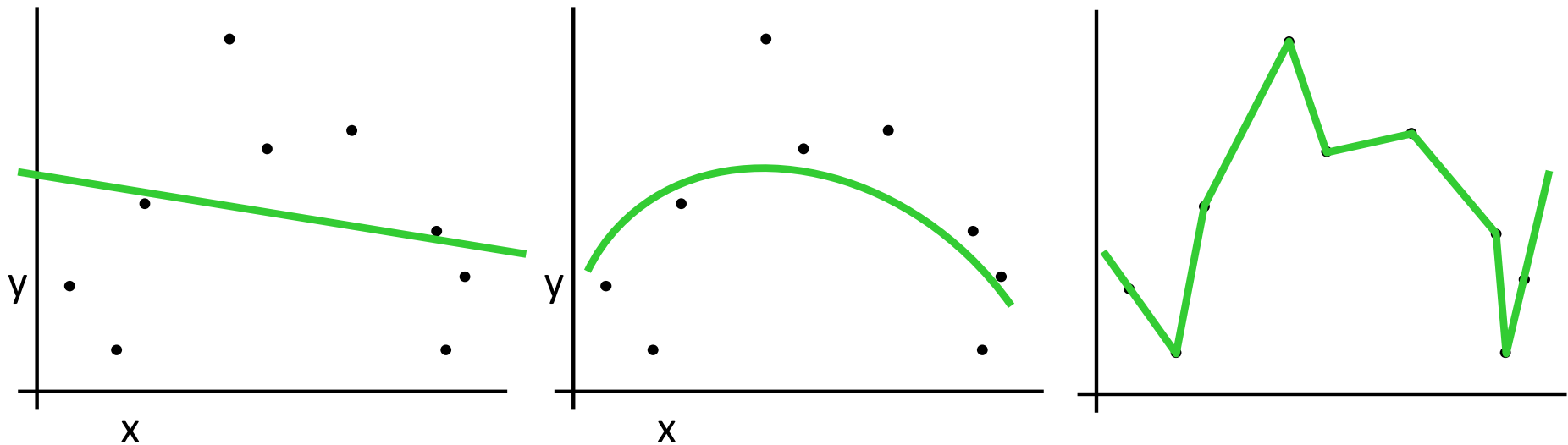


Join-the-dots



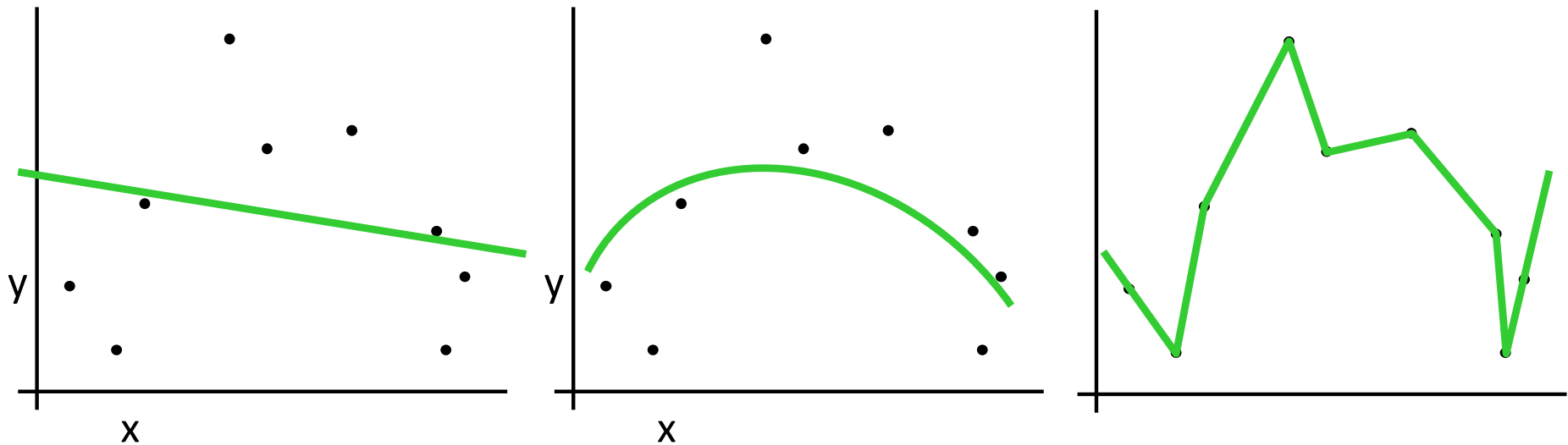
Also known as **piecewise linear nonparametric regression** if that makes you feel better

Which is best?



Why not choose the method with the best fit to the data?

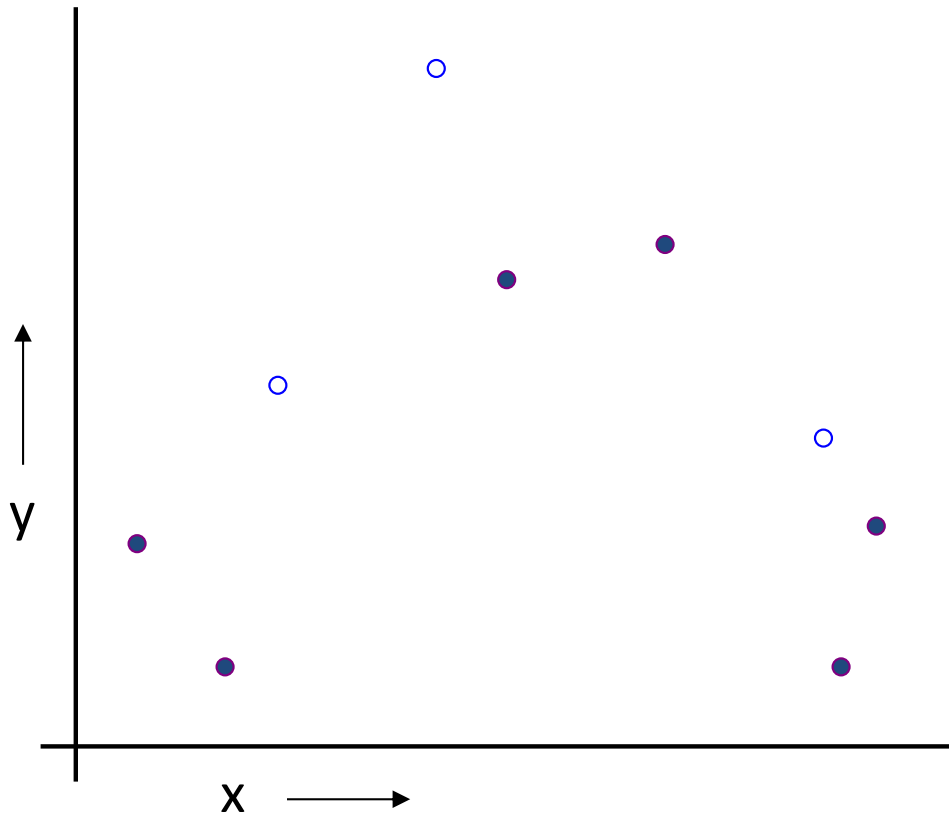
What do we really want?



Why not choose the method with the best fit to the data?

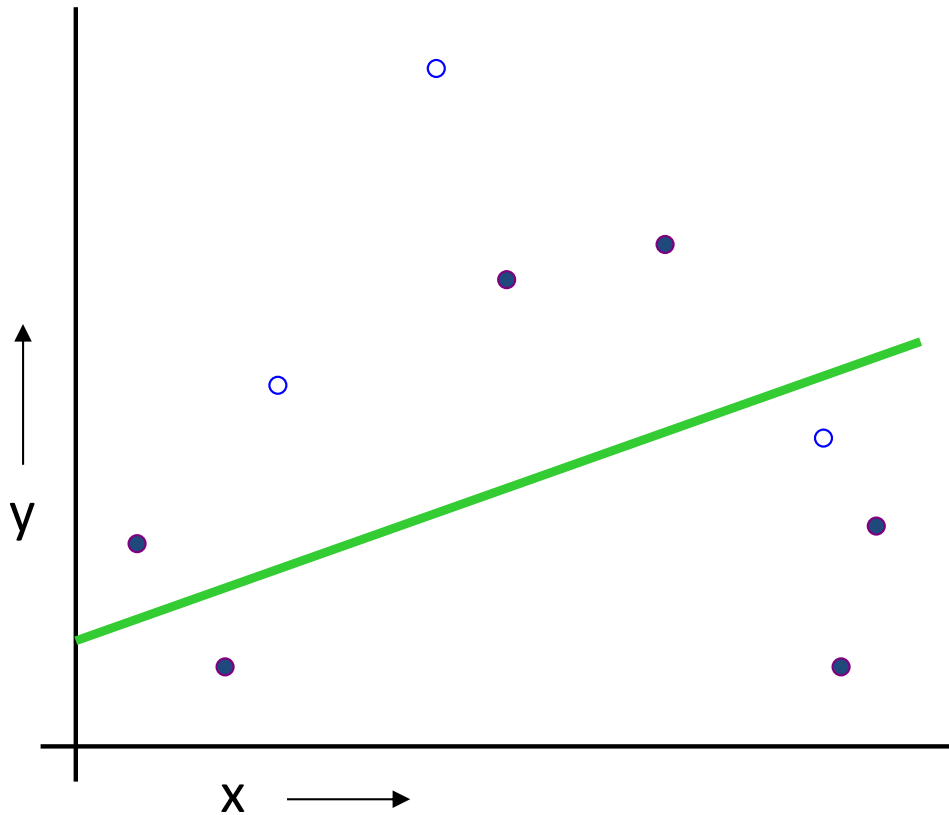
“How well are you going to predict future data drawn from the same distribution?”

The test set method



1. Randomly choose 30% of the data to be in a **test set**
2. The remainder is a **training set**

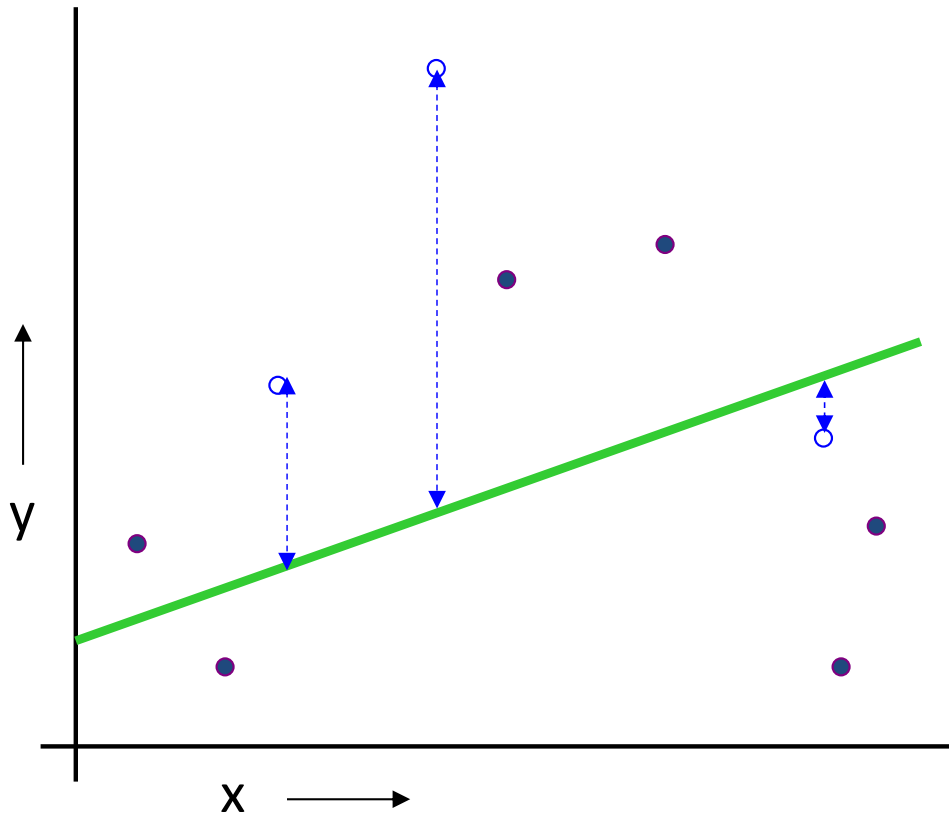
The test set method



(Linear regression example)

1. Randomly choose 30% of the data to be in a **test set**
2. The remainder is a **training set**
3. Perform your regression on the **training set**

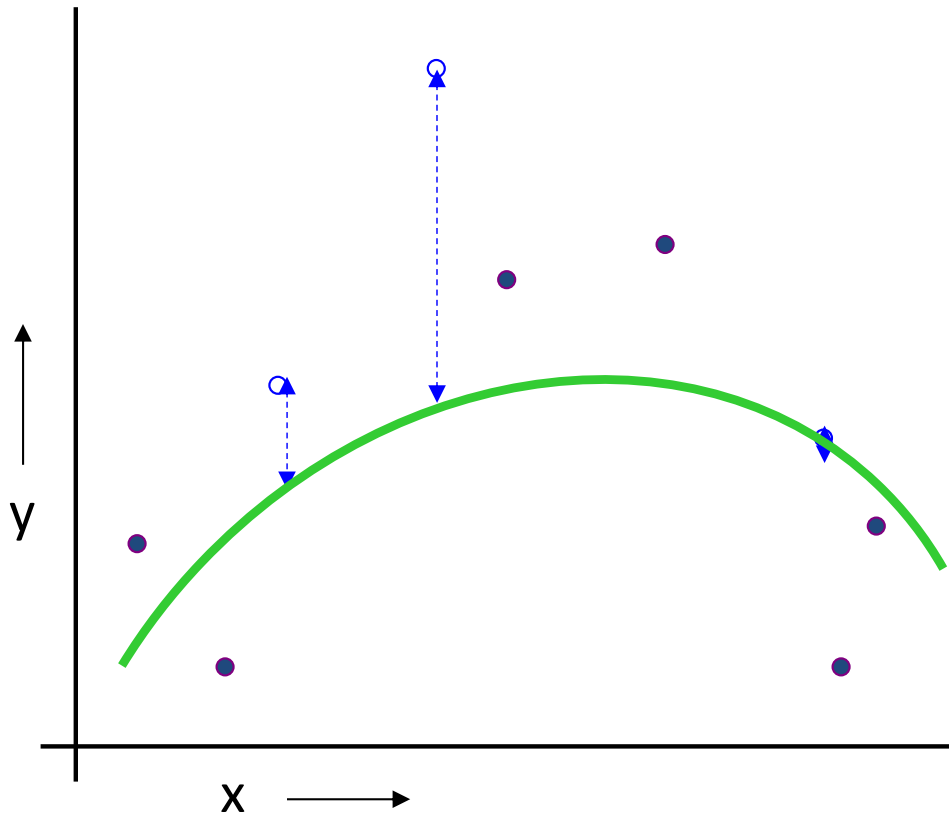
The test set method



(Linear regression example)
Mean Squared Error = 2.4

1. Randomly choose 30% of the data to be in a **test set**
2. The remainder is a **training set**
3. Perform your regression on the training set
4. Estimate your future performance with the test set

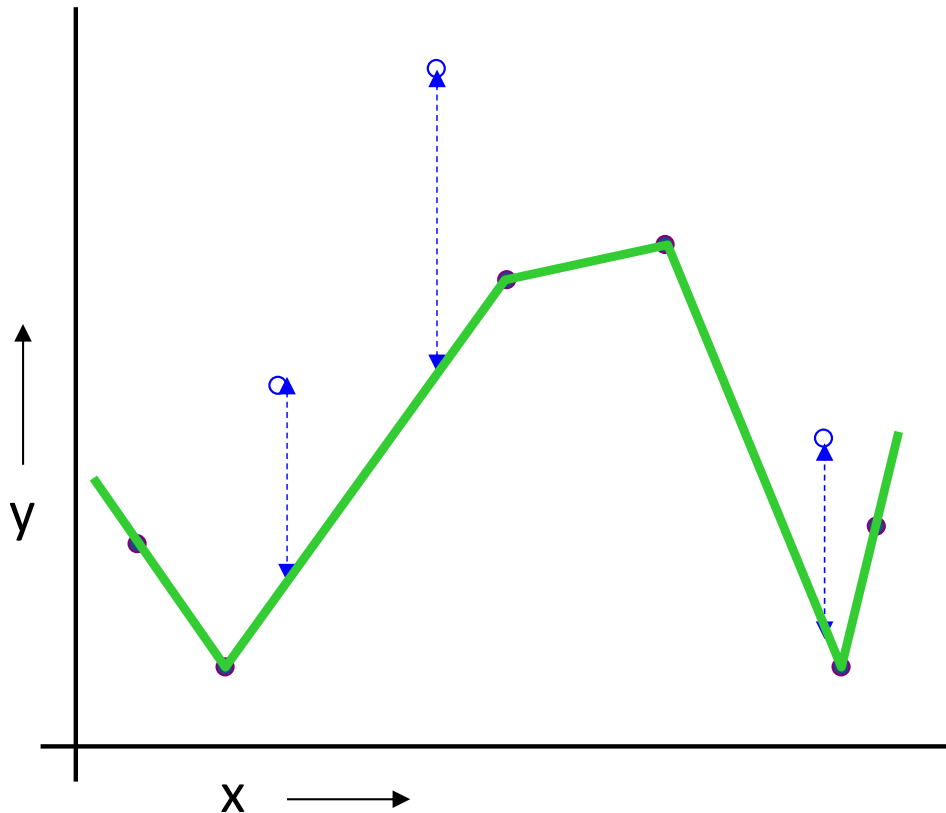
The test set method



(Quadratic regression example)
Mean Squared Error = 0.9

1. Randomly choose 30% of the data to be in a **test set**
2. The remainder is a **training set**
3. Perform your regression on the training set
4. Estimate your future performance with the test set

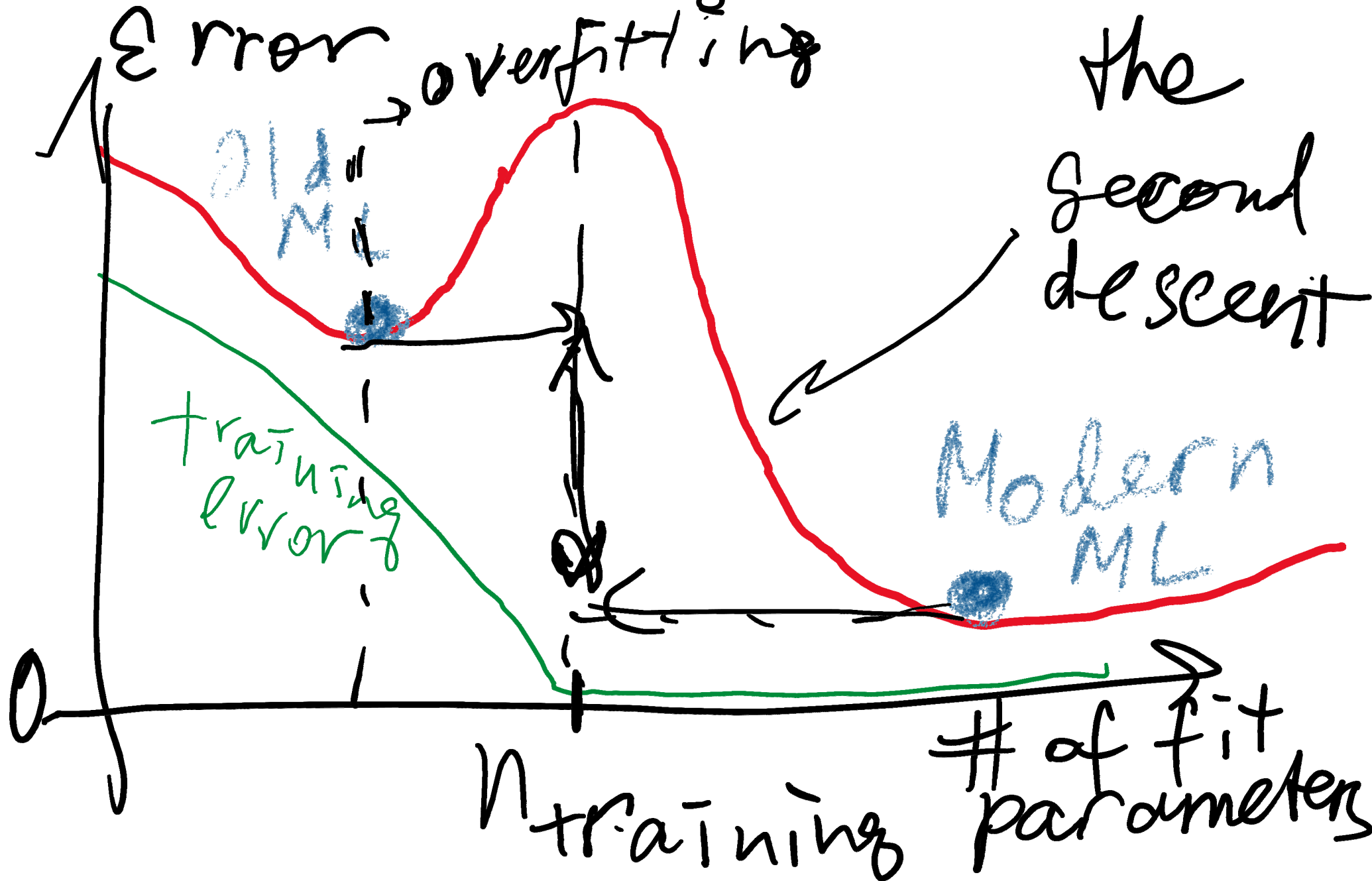
The test set method



(Join the dots example)
Mean Squared Error = 2.2

1. Randomly choose 30% of the data to be in a **test set**
2. The remainder is a **training set**
3. Perform your regression on the training set
4. Estimate your future performance with the test set

Double descend- the main reason modern
Machine Learning works so well



R^2 and Adjusted R^2

The **coefficient of multiple determination R^2**

$$R^2 = \frac{SS_R}{SS_T} = 1 - \frac{SS_E}{SS_T}$$

The **adjusted R^2** is

$$R^2_{\text{adj}} = 1 - \frac{SS_E/(n - p)}{SS_T/(n - 1)} \approx 1 - \frac{\sum \epsilon^2}{\sum y^2} \quad (12-23)$$

- The adjusted R^2 statistic penalizes **adding terms** to the MLR model.
- It can help guard against **overfitting** (including regressors that are not really useful)

How to know where to stop adding variables?

- Adding new variables x_i to MLR
watch the adjusted R^2
- Once the adjusted R^2
no longer increases = stop.
Now you did the best you can.

Credit: XKCD
comics

WHY ARE THERE SLAVES IN THE BIBLE

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

QUESTIONS

FOUND IN GOOGLE AUTOCOMPLETE

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

WHY AREN'T MY
ARMS GROWING



WHY ARE THERE WEEKS
WHY DO I FEEL DIZZY

WHY AREN'T ECONOMISTS RICH

WHY DO AMERICANS CALL IT SOCCER

WHY ARE MY EARS RINGING

WHY ARE THERE SO MANY AVENGERS

WHY ARE THE AVENGERS FIGHTING THE X MEN

WHY IS WOLVERINE NOT IN THE AVENGERS

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 ARE THERE MALE AND FEMALE BIKES

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 LYING GOOD

WHY ARE THERE
GHOSTS



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 AREN'T
THERE GUNS IN
HARRY POTTER



WHY ARE ULTRASOUNDS IMPORTANT

WHY ARE ULTRASOUND MACHINES EXPENSIVE

WHY IS STEALING WRONG

WHY ARE THERE FIREWORKS

WHY ARE THERE 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 BRIDESMAIDS

WHY DO DYING PEOPLE REACH UP

WHY AREN'T THERE VARICOSE ARTERIES

WHY ARE OLD KINGDOMS 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
SQUIRRELS



WHY IS SEX
SO IMPORTANT



Matlab exercise on #2 on MLR

- Every group works with
g0=2907; g1=1527; g2=2629; g3=2881;
g4=1144; g5=1066;
- Compute **Multiple Linear Regression (MLR)**:
where
y=exp_t (g0); x1= exp_t (g1); x2= 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

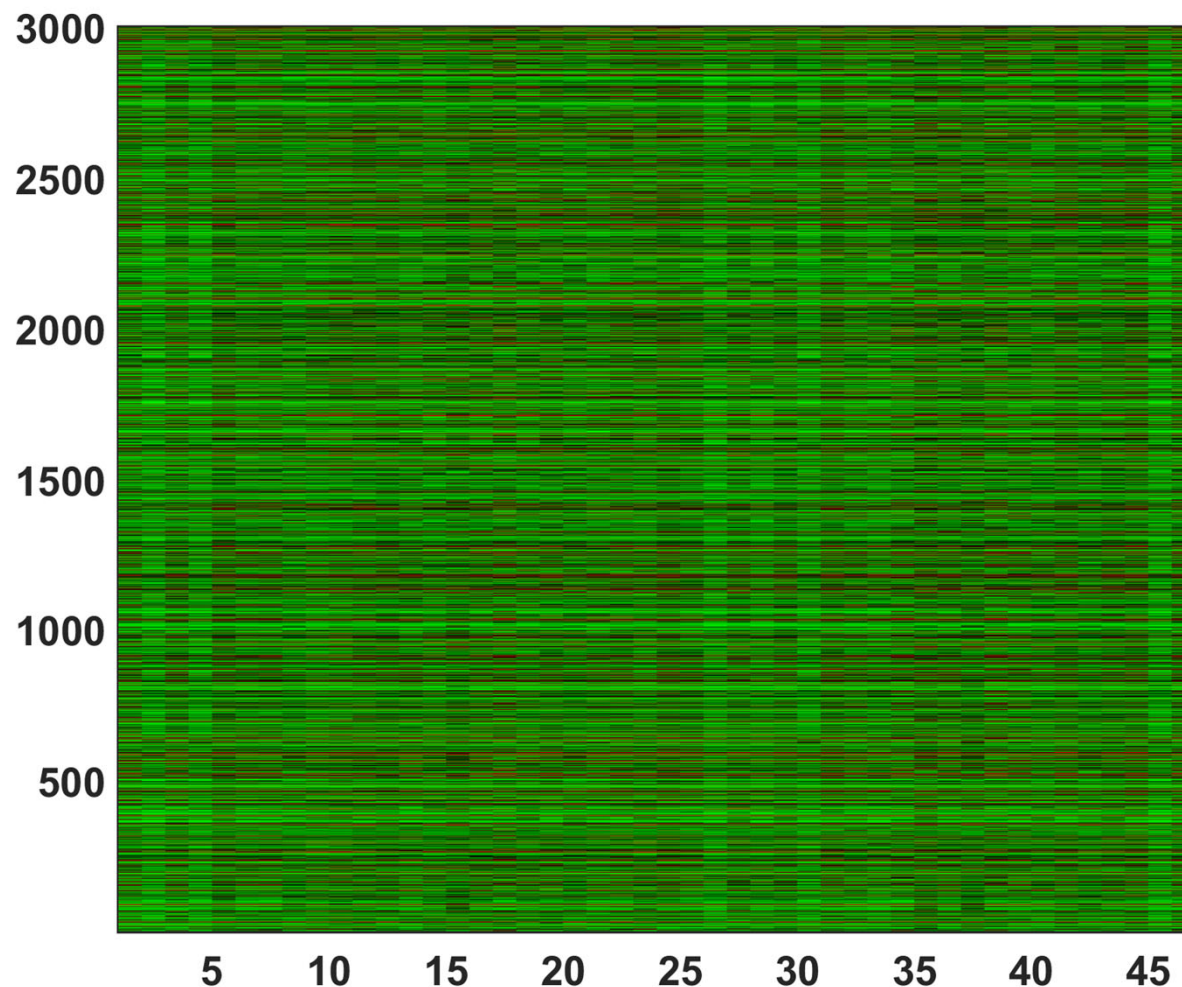
How I did it

- `g0=2907; g1=1527; g2=2629; g3=2881;g4=1144; g5=1066;`
- `y=exp_t(g0,:)' ;`
- `%% first use one x to predict y`
- `x=exp_t(g1,:)' ;`
- `figure; plot(x,y,'ko')`
- `lm=fitlm(x,y)`
- `y_fit=lm.Fitted;`
- `hold on;`
- `plot(x,lm.Fitted,'r-');`
- `%% now use 2 x's to predict y`
- `x=[exp_t(g1,:)', exp_t(g2,:)]';`
- `lm2=fitlm(x,y)`
- `y_fit=lm2.Fitted;`
- `hold on; plot(x(:,1),y_fit,'gd');`
- `%% now use m x's to predict y`
- `corr_matrix=corr(exp_t');`
- `g0=2907;`
- `[u v]=sort(corr_matrix(g0,:), 'descend');`
- `x=[exp_t(v(2:m+1),:)]';`
- `lm3=fitlm(x,y)`
- `y_fit=lm3.Fitted;`
- `plot(x(:,1),y_fit,'s');`

Clustering analysis of gene expression data

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

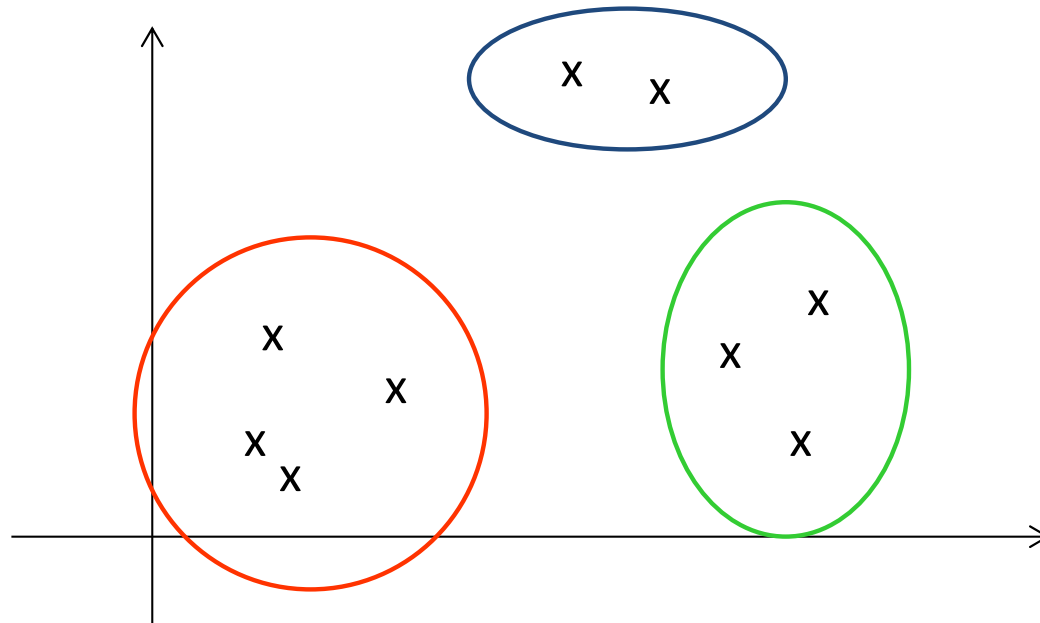
How to find the entire groups of mutually correlated genes if you have **many genes** and **many samples**?



Clustering to the rescue!

What is clustering?

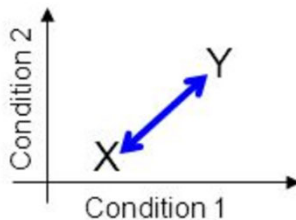
- The goal of **clustering** is to
 - group data points that are close (or **similar**) to each other
 - Usually, one needs to identify such groups (or clusters) in an **unsupervised** manner
 - Sometimes one takes into account **prior information** (Bayesian methods)
- Need to define some **distance 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?

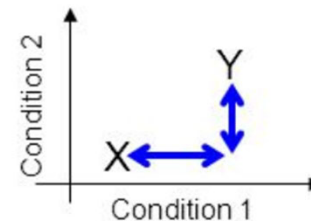
- Euclidean distance:
 - Most commonly used distance
 - Sphere shaped cluster
 - Corresponds to the geometric distance into the multidimensional space

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



- 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

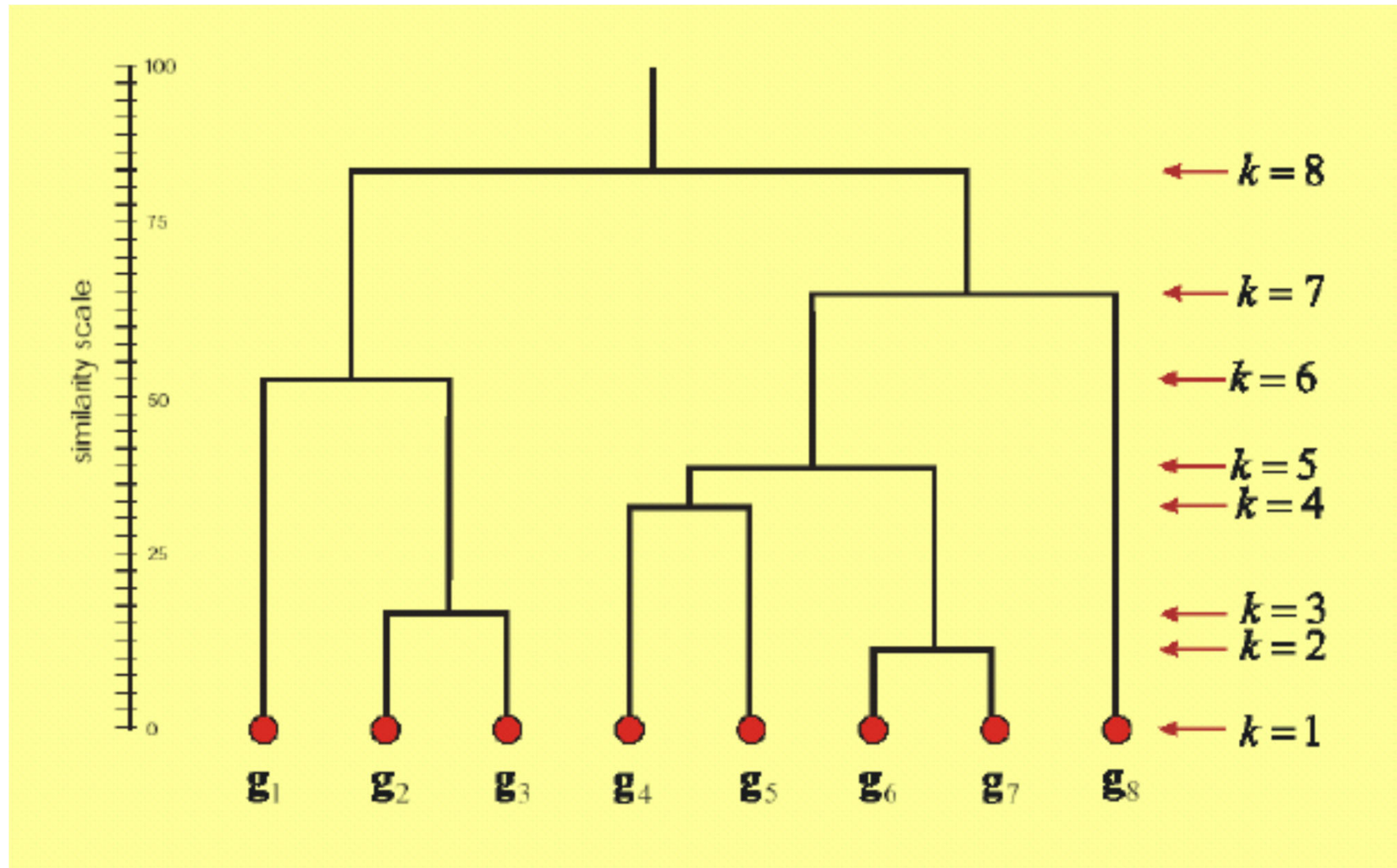
- Hierarchical if one doesn't know in advance the # of clusters
 - Agglomerative: start with N clusters and gradually merge them into 1 cluster
 - Divisive: start with 1 cluster and gradually break it up into N clusters
- Non-hierarchical algorithms
 - K-means clustering:
 - Iteratively apply the following 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

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



Clustering in Matlab

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	<p>A distance function specified using @:</p> <pre>D = pdist(X,@distfun)</pre> <p>A distance function must be of form</p> <pre>d2 = distfun(XI,XJ)</pre> <p>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,:).</p>

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

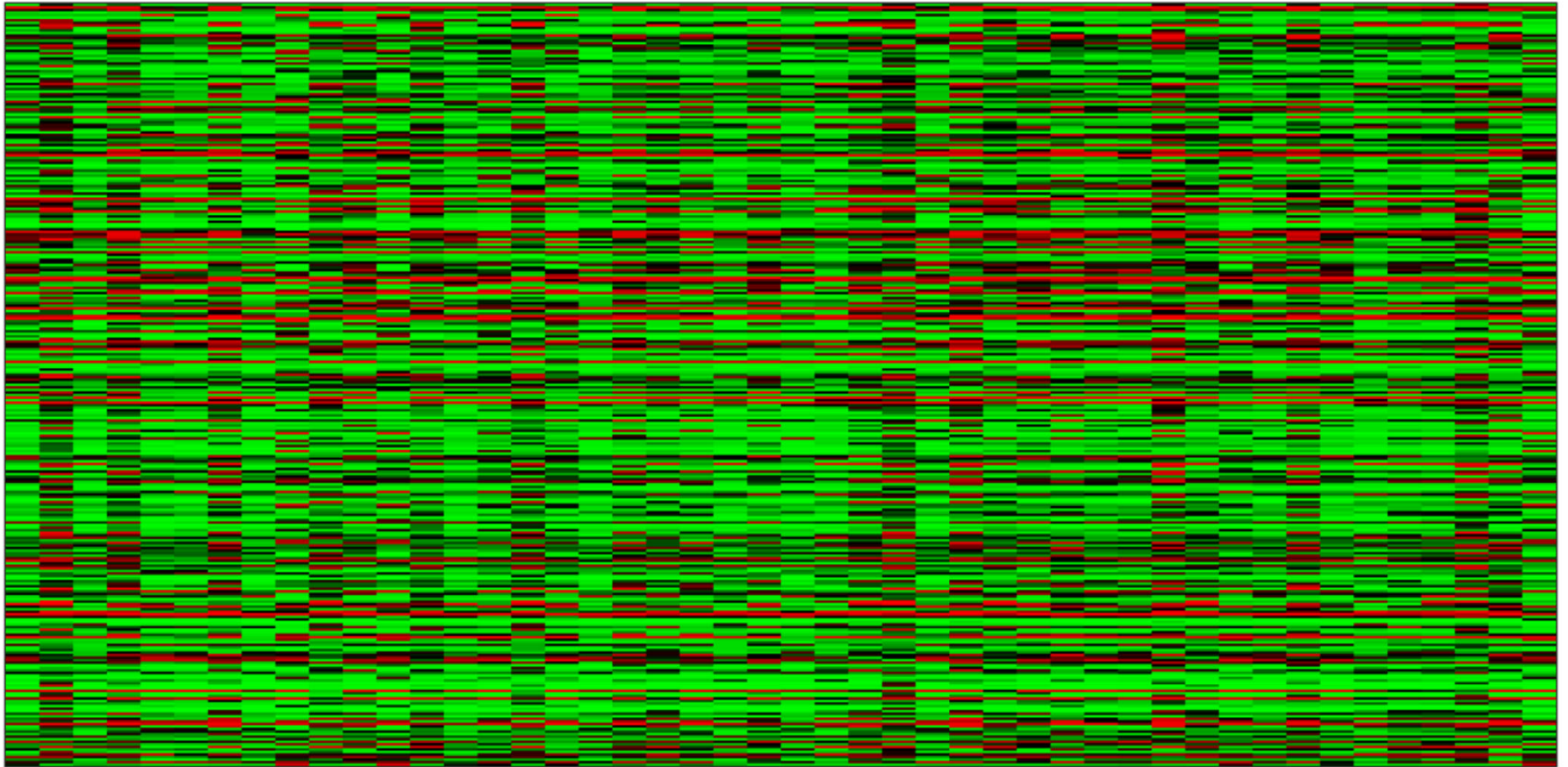
X	Matrix with two or more rows. The rows represent observations, the columns represent categories or dimensions.																
method	<div>Algorithm for computing distance between clusters.</div> <table><tr><th>Method</th><th>Description</th></tr><tr><td>'average'</td><td>Unweighted average distance (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></table> <div>Default: 'single'</div>	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

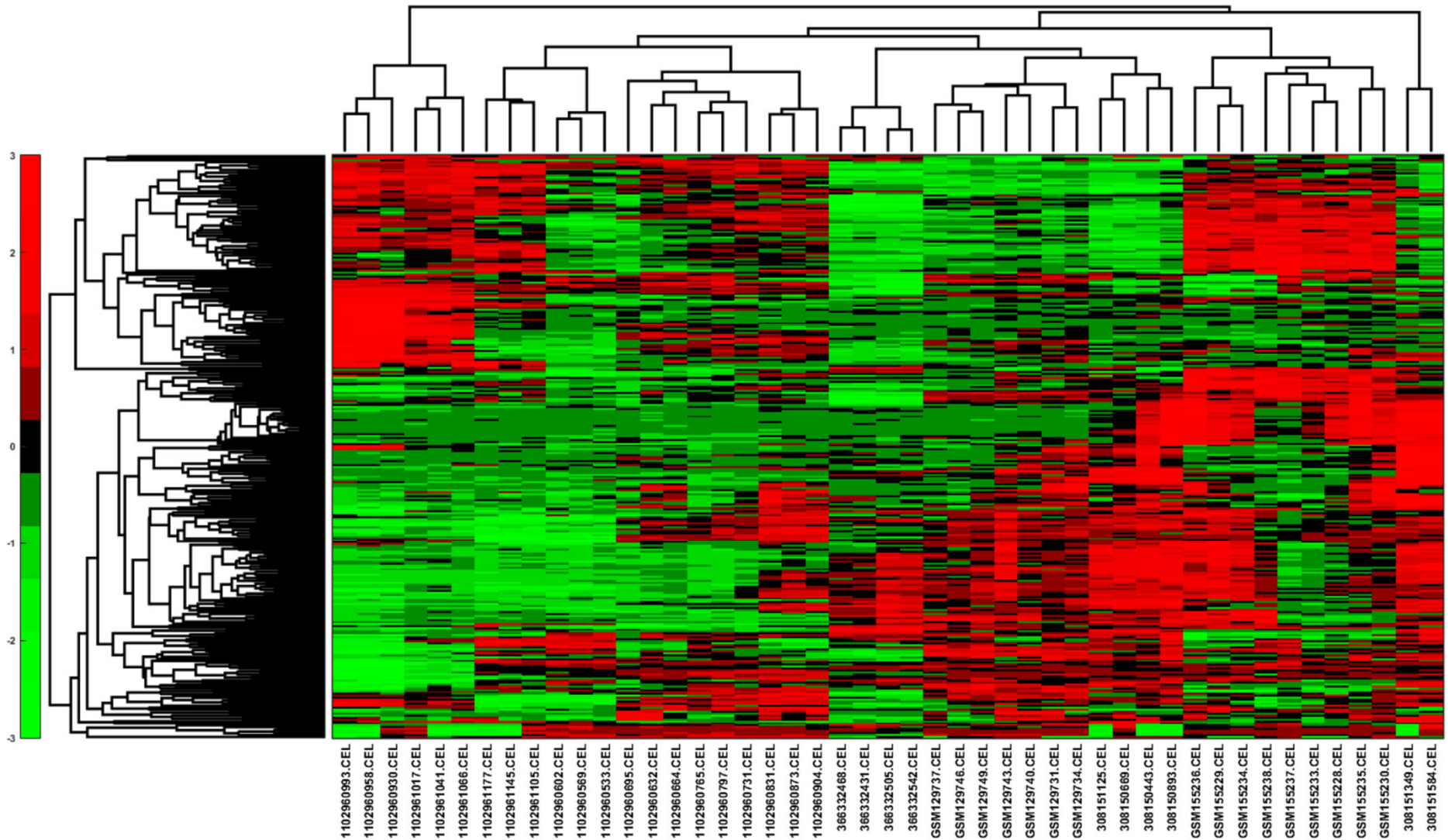
- Each group will analyze a **cluster of genes** identified in the T cell expression table
- Analyze the table of **top 100 genes by variance** in 47 samples
- Cluster them using:
 - Group 1: 'linkage', 'average', 'RowPDistValue', 'euclidean',
 - Group 2: 'linkage', 'single', 'RowPDistValue', 'cityblock',
 - Group 3: 'linkage', 'average', 'RowPDistValue', 'correlation',
 - Group 4: 'linkage', 'single', 'RowPDistValue', 'euclidean',
 - Group 5: '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)`

```
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');
set(CGobj1,'RowLabels',gene_names1,'ColumnLabels',array_names,'linkage',
'average','RowPDist','correlation');
```

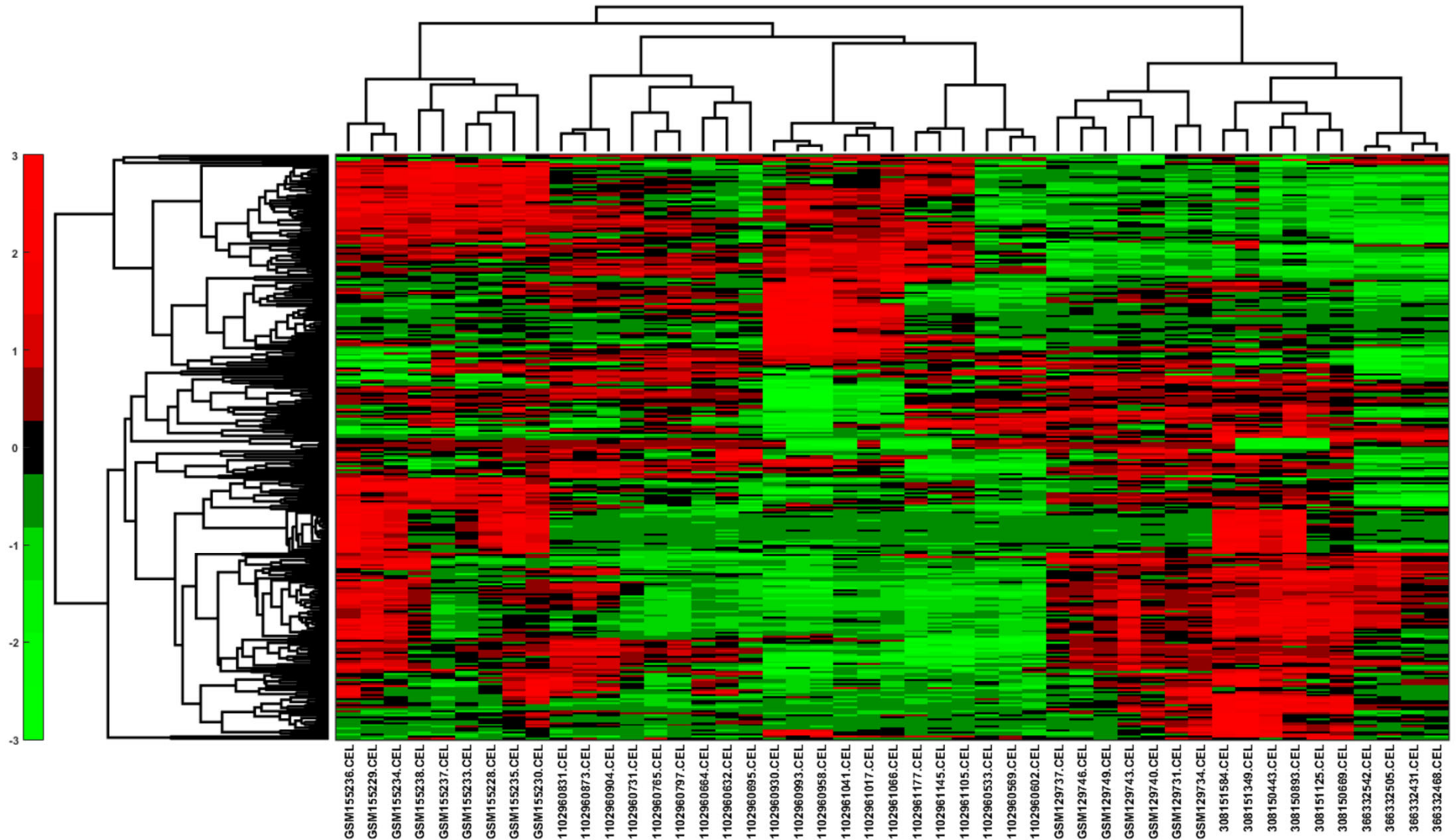
Before clustering



UPGMA hierarchical clustering, Euclidian distance



UPGMA hierarchical clustering, correlation distance



Search for shared biological functions

- copy the list of displayed genes
- go to "Start Analysis" on <https://david.ncifcrf.gov/tools.jsp>
- Paste genes from gene list displayed by Matlab into the box in the left panel of the website
- select ENSEMBL_GENE_ID and “gene list” radio button
- Click "Functional Annotation Clustering"
- Select groups in “Annotation Summary Results” which have many genes from your list. Definitely select “PUBMED_ID” and interaction databases like “Biogrid”
- First look at "Functional Annotation Chart" rectangular button below to display all overrepresented terms. Sort by “Benjamini” correction for multiple hypotheses testing
- Select "Functional Annotation Clustering" rectangular button below to display annotation results for gene list broken into multiple groups (clusters) each with related biological functions
- Write down the # of genes in the cluster and the top functions in two most interesting clusters

%%%

**%Which biological functions are
overrepresented in different clusters?**

%1) Pick a cluster:

%2) Select a node on the tree of rows,

%3) Right click

**%4) Choose “export group info” into
the workspace**

%5) Name it gene_list

**%Run the following two Matlab
commands to display genes**

g1=gene_list.RowNodeNames;

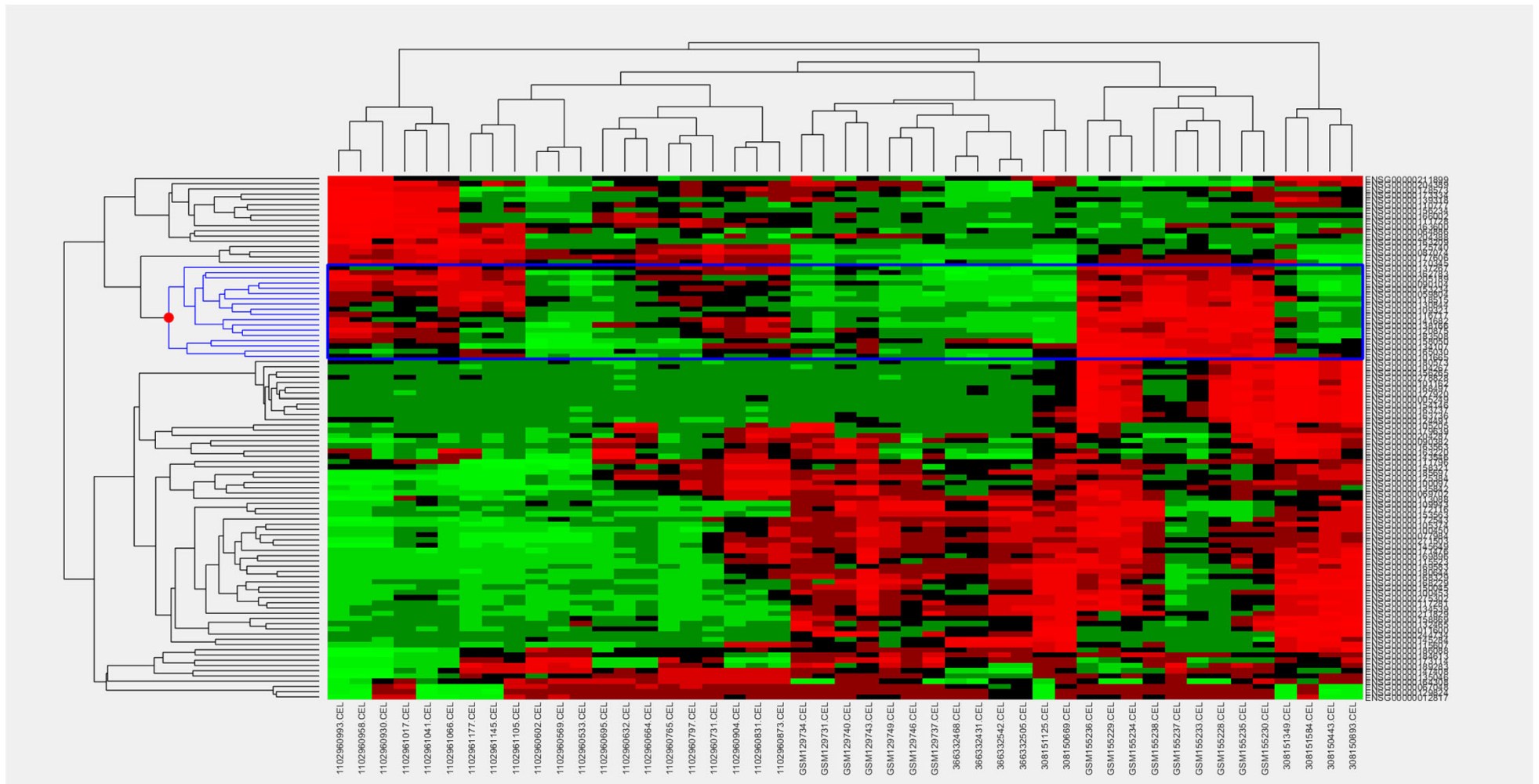
for m=1:length(g1);

disp(g1{m});

end;













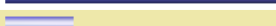











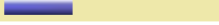

% select ENSEMBL_GENE_ID and “gene list” radio button
% Click "Functional Annotation Clustering"
% Select groups in “Annotation Summary Results”
% which have many genes from your list.
% Definitely select “PUBMED_ID” and
% interaction databases like “Biogrid”
% First look at "Functional Annotation Chart" rectangular button below
% to display all overrepresented terms.
% Sort by “Benjamini” correction for multiple hypotheses testing
% Select "Functional Annotation Clustering" rectangular button below
% to display annotation results for gene list broken into multiple groups
% (clusters) each with related biological functions
% Write down the # of genes in the cluster and the top functions
% in two most interesting clusters

Using options:
'linkage', 'average', 'RowPDistValue', 'euclidean',



54 chart records

[Download File](#)

Sublist	Category	Term	RT	Genes	Count	%	P-Value	Benjamini
<input type="checkbox"/>	GOTERM_CC_DIRECT	nucleus	RT		16	88.9	8.1E-7	3.7E-5
<input type="checkbox"/>	PIR_SUPERFAMILY	dual specificity protein phosphatase (MAP kinase phosphatase)	RT		3	16.7	4.0E-5	8.0E-5
<input type="checkbox"/>	GOTERM_MF_DIRECT	protein tyrosine/threonine phosphatase activity	RT		3	16.7	3.4E-5	1.3E-3
<input type="checkbox"/>	GOTERM_MF_DIRECT	MAP kinase tyrosine phosphatase activity	RT		3	16.7	3.4E-5	1.3E-3
<input type="checkbox"/>	GOTERM_MF_DIRECT	MAP kinase tyrosine/serine/threonine phosphatase activity	RT		3	16.7	5.9E-5	1.5E-3
<input type="checkbox"/>	INTERPRO	Mitogen-activated protein (MAP) kinase phosphatase	RT		3	16.7	3.3E-5	1.9E-3
<input type="checkbox"/>	SMART	RHOD	RT		3	16.7	2.5E-4	4.8E-3
<input type="checkbox"/>	INTERPRO	Rhodanese-like domain	RT		3	16.7	2.2E-4	6.2E-3
<input type="checkbox"/>	SMART	DSPc	RT		3	16.7	8.4E-4	8.0E-3
<input type="checkbox"/>	INTERPRO	Dual specificity phosphatase, catalytic domain	RT		3	16.7	6.0E-4	9.2E-3
<input type="checkbox"/>	INTERPRO	Dual specificity phosphatase, subgroup, catalytic domain	RT		3	16.7	6.6E-4	9.2E-3
<input type="checkbox"/>	GOTERM_BP_DIRECT	endoderm formation	RT		3	16.7	5.6E-5	1.1E-2
<input type="checkbox"/>	UP_KW_CELLULAR_COMPONENT	Nucleus	RT		13	72.2	1.5E-3	1.3E-2
<input type="checkbox"/>	SMART	PTPc motif	RT		3	16.7	2.3E-3	1.5E-2
<input type="checkbox"/>	GOTERM_MF_DIRECT	phosphoprotein phosphatase activity	RT		3	16.7	8.0E-4	1.5E-2
<input type="checkbox"/>	INTERPRO	Protein-tyrosine phosphatase, catalytic	RT		3	16.7	1.4E-3	1.6E-2
<input type="checkbox"/>	UP_KW_PTM	Ubl conjugation	RT		7	38.9	4.5E-3	1.9E-2
<input type="checkbox"/>	UP_KW_PTM	Isopeptide bond	RT		6	33.3	5.4E-3	1.9E-2
<input type="checkbox"/>	INTERPRO	Protein-tyrosine phosphatase, active site	RT		3	16.7	2.1E-3	2.0E-2
<input type="checkbox"/>	INTERPRO	Protein-tyrosine/Dual specificity phosphatase	RT		3	16.7	2.8E-3	2.3E-2
<input type="checkbox"/>	UP_SEQ_FEATURE	DOMAIN:Rhodanese	RT		3	16.7	1.9E-4	2.4E-2
<input type="checkbox"/>	KEGG_PATHWAY	MAPK signaling pathway	RT		5	27.8	5.9E-4	2.8E-2
<input type="checkbox"/>	GOTERM_MF_DIRECT	myosin phosphatase activity	RT		3	16.7	2.4E-3	3.6E-2
<input type="checkbox"/>	GOTERM_MF_DIRECT	protein tyrosine phosphatase activity	RT		3	16.7	4.2E-3	5.3E-2
<input type="checkbox"/>	GOTERM_CC_DIRECT	nucleoplasm	RT		10	55.6	2.3E-3	5.4E-2
<input type="checkbox"/>	GOTERM_BP_DIRECT	negative regulation of MAPK cascade	RT		3	16.7	7.0E-4	6.8E-2

Gene list being analyzed

Clustering options and stringency

score for the group based on the EASE scores of each term members. The higher, the more enriched.

Every term in the annotation cluster

Related Term Search

Genes involved in individual term

ALL genes involved in this annotation cluster

Functional Annotation Clustering

Current Gene List: demolist1
171 DAVID IDs

Options Classification Stringency: High

Rerun using options Create Sublist

[Download File](#)

A group of terms having similar biological meaning due to sharing similar gene members

Annotation Cluster	Enrichment Score	Term	RT	Count	EASE Score
Annotation Cluster 1 Enrichment Score: 3.69					
<input type="checkbox"/>	SP_PIR_KEYWORDS	chromoprotein	RT	7	1.1E-5
<input type="checkbox"/>	SP_PIR_KEYWORDS	metalloprotein	RT	8	4.7E-5
<input type="checkbox"/>	SP_PIR_KEYWORDS	iron	RT	9	2.1E-4
<input type="checkbox"/>	GOTERM_MF_ALL	iron ion binding	RT	10	2.5E-4
<input type="checkbox"/>	SP_PIR_KEYWORDS	heme	RT	7	3.5E-4
<input type="checkbox"/>	GOTERM_MF_ALL	tetrapyrrole binding	RT	6	1.3E-3
<input type="checkbox"/>	GOTERM_MF_ALL	heme binding	RT	6	1.3E-3
Annotation Cluster 2 Enrichment Score: 3.52					
<input type="checkbox"/>	SP_PIR_KEYWORDS	antibiotic	RT	5	2.2E-4
<input type="checkbox"/>	SP_PIR_KEYWORDS	antimicrobial	RT	5	2.4E-4
<input type="checkbox"/>	GOTERM_BP_ALL	defense response to bacteria	RT	6	5.4E-4
Annotation Cluster 3 Enrichment Score: 2.66					
<input type="checkbox"/>	UP_SEQ_FEATURE	domain:Ig-like C2-type 1	RT	8	5.4E-4
<input type="checkbox"/>	UP_SEQ_FEATURE	domain:Ig-like C2-type 2	RT	8	5.4E-4
<input type="checkbox"/>	INTERPRO_NAME	Immunoglobulin	RT	6	3.6E-2
Annotation Cluster 4 Enrichment Score: 2.63					

EASE Score, the modified Fisher Exact P-Value. They are identical to that in the Chart Report. The smaller, the more enriched.

Functional Annotation Clustering
























[Help and Manual](#)






















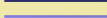



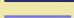



Current Gene List: List_3
Current Background: Homo sapiens
18 DAVID IDs






















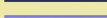



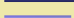



☒ Options Classification Stringency Medium ▾

25 Cluster(s)

 [Download File](#)

Annotation Cluster 1		Enrichment Score: 5.2			Count	P_Value	Benjamini
<input type="checkbox"/>	DISGENET	Juvenile arthritis	RT		7	1.5E-8	4.7E-7
<input type="checkbox"/>	DISGENET	Juvenile psoriatic arthritis	RT		7	1.5E-8	4.7E-7
<input type="checkbox"/>	DISGENET	Polyarthritis, Juvenile, Rheumatoid Factor Negative	RT		7	1.5E-8	4.7E-7
<input type="checkbox"/>	DISGENET	Polyarthritis, Juvenile, Rheumatoid Factor Positive	RT		7	1.5E-8	4.7E-7
<input type="checkbox"/>	DISGENET	Juvenile-Onset Still Disease	RT		7	1.8E-8	4.7E-7
<input type="checkbox"/>	KEGG_PATHWAY	MAPK signaling pathway	RT		5	5.9E-4	2.8E-2
<input type="checkbox"/>	BIOGRID_INTERACTION	mitogen-activated protein kinase 1(MAPK1)	RT		4	3.8E-3	1.0E0
<input type="checkbox"/>	WIKIPATHWAYS	MAPK signaling pathway	RT		3	5.8E-2	6.9E-1
<input type="checkbox"/>	GAD_DISEASE_CLASS	UNKNOWN	RT		5	1.5E-1	9.9E-1
Annotation Cluster 2		Enrichment Score: 2.83			Count	P_Value	Benjamini
<input type="checkbox"/>	INTERPRO	Mitogen-activated protein (MAP) kinase phosphatase	RT		3	3.3E-5	1.9E-3
<input type="checkbox"/>	GOTERM_MF_DIRECT	protein tyrosine/threonine phosphatase activity	RT		3	3.4E-5	1.3E-3
<input type="checkbox"/>	GOTERM_MF_DIRECT	MAP kinase tyrosine phosphatase activity	RT		3	3.4E-5	1.3E-3
<input type="checkbox"/>	PIR_SUPERFAMILY	dual specificity protein phosphatase (MAP kinase phosphatase)	RT		3	4.0E-5	8.0E-5
<input type="checkbox"/>	GOTERM_BP_DIRECT	endoderm formation	RT		3	5.6E-5	1.1E-2
<input type="checkbox"/>	GOTERM_MF_DIRECT	MAP kinase tyrosine/serine/threonine phosphatase activity	RT		3	5.9E-5	1.5E-3
<input type="checkbox"/>	PUBMED_ID	27880917	RT		4	1.7E-4	2.5E-2
<input type="checkbox"/>	UP_SEQ_FEATURE	DOMAIN:Rhodanese	RT		3	1.9E-4	2.4E-2
<input type="checkbox"/>	INTERPRO	Rhodanese-like domain	RT		3	2.2E-4	6.2E-3
<input type="checkbox"/>	SMART	RHOD	RT		3	2.5E-4	4.8E-3

Annotation Cluster 3		Enrichment Score: 2.43	G		Count	P_Value	Benjamini
<input type="checkbox"/>	DISGENET	Arsenic Poisoning, Inorganic	RT		3	3.5E-3	4.6E-2
<input type="checkbox"/>	DISGENET	Nervous System, Organic Arsenic Poisoning	RT		3	3.5E-3	4.6E-2
<input type="checkbox"/>	DISGENET	Arsenic Poisoning	RT		3	3.5E-3	4.6E-2
<input type="checkbox"/>	DISGENET	Arsenic Encephalopathy	RT		3	3.5E-3	4.6E-2
<input type="checkbox"/>	DISGENET	Arsenic Induced Polyneuropathy	RT		3	3.5E-3	4.6E-2
<input type="checkbox"/>	DISGENET	Dermatologic disorders	RT		3	5.1E-3	5.6E-2
Annotation Cluster 4		Enrichment Score: 2.26	G		Count	P_Value	Benjamini
<input type="checkbox"/>	PUBMED_ID	19322201	RT		7	1.3E-8	5.9E-6
<input type="checkbox"/>	BIOGRID_INTERACTION	ELAV like RNA binding protein 1(ELAVL1)	RT		7	4.4E-3	1.0E0
<input type="checkbox"/>	UCSC_TFBS	CEBPA	RT		7	1.8E-1	1.0E0
<input type="checkbox"/>	UCSC_TFBS	CDPCR3HD	RT		7	6.5E-1	1.0E0
<input type="checkbox"/>	UCSC_TFBS	FOX3D	RT		5	7.4E-1	1.0E0
Annotation Cluster 5		Enrichment Score: 2.14	G		Count	P_Value	Benjamini
<input type="checkbox"/>	GOTERM_BP_DIRECT	negative regulation of transcription from RNA polymerase II promoter	RT		6	1.4E-3	9.1E-2
<input type="checkbox"/>	BIOGRID_INTERACTION	retinoid X receptor alpha(RXRA)	RT		3	6.1E-3	1.0E0
<input type="checkbox"/>	GOTERM_MF_DIRECT	protein heterodimerization activity	RT		3	4.5E-2	3.7E-1
Annotation Cluster 6		Enrichment Score: 1.95	G		Count	P_Value	Benjamini
<input type="checkbox"/>	REACTOME_PATHWAY	Generic Transcription Pathway	RT		7	2.8E-3	1.7E-1
<input type="checkbox"/>	REACTOME_PATHWAY	RNA Polymerase II Transcription	RT		7	4.6E-3	1.7E-1
<input type="checkbox"/>	REACTOME_PATHWAY	Gene expression (Transcription)	RT		7	8.2E-3	2.0E-1
<input type="checkbox"/>	GAD_DISEASE_CLASS	UNKNOWN	RT		5	1.5E-1	9.9E-1
Annotation Cluster 7		Enrichment Score: 1.76	G		Count	P_Value	Benjamini
<input type="checkbox"/>	PUBMED_ID	18029348	RT		6	1.8E-5	3.4E-3
<input type="checkbox"/>	UP_KW_PTM	Isopeptide bond	RT		6	5.4E-3	1.9E-2
<input type="checkbox"/>	PUBMED_ID	15342556	RT		3	7.9E-3	4.8E-1
<input type="checkbox"/>	PUBMED_ID	26496610	RT		3	1.0E-1	1.0E0
<input type="checkbox"/>	GOTERM_MF_DIRECT	metal ion binding	RT		4	4.5E-1	1.0E0
<input type="checkbox"/>	UCSC_TFBS	TAL1ALPHA47	RT		3	7.9E-1	1.0E0

Annotation Cluster 3		Enrichment Score: 2.43	G		Count	P_Value	Benjamini
<input type="checkbox"/>	DISGENET	Arsenic Poisoning, Inorganic	RT		3	3.5E-3	4.6E-2
<input type="checkbox"/>	DISGENET	Nervous System, Organic Arsenic Poisoning	RT		3	3.5E-3	4.6E-2
<input type="checkbox"/>	DISGENET	Arsenic Poisoning	RT		3	3.5E-3	4.6E-2
<input type="checkbox"/>	DISGENET	Arsenic Encephalopathy	RT		3	3.5E-3	4.6E-2
<input type="checkbox"/>	DISGENET	Arsenic Induced Polyneuropathy	RT		3	3.5E-3	4.6E-2
<input type="checkbox"/>	DISGENET	Dermatologic disorders	RT		3	5.1E-3	5.6E-2
Annotation Cluster 4		Enrichment Score: 2.26	G		Count	P_Value	Benjamini
<input type="checkbox"/>	PUBMED_ID	19322201	RT		7	1.3E-8	5.9E-6
<input type="checkbox"/>	BIOGRID_INTERACTION	ELAV like RNA binding protein 1(ELAVL1)	RT		7	4.4E-3	1.0E0
<input type="checkbox"/>	UCSC_TFBS	CEBPA	RT		7	1.8E-1	1.0E0
<input type="checkbox"/>	UCSC_TFBS	CDPCR3HD	RT		7	6.5E-1	1.0E0
<input type="checkbox"/>	UCSC_TFBS	FOX D3	RT		5	7.4E-1	1.0E0
Annotation Cluster 5		Enrichment Score: 2.14	G		Count	P_Value	Benjamini
<input type="checkbox"/>	GOTERM_BP_DIRECT	negative regulation of transcription from RNA polymerase II promoter	RT		6	1.4E-3	9.1E-2
<input type="checkbox"/>	BIOGRID_INTERACTION	retinoid X receptor alpha(RXRA)	RT		3	6.1E-3	1.0E0
<input type="checkbox"/>	GOTERM_MF_DIRECT	protein heterodimerization activity	RT		3	4.5E-2	3.7E-1
Annotation Cluster 6		Enrichment Score: 1.95	G		Count	P_Value	Benjamini
<input type="checkbox"/>	REACTOME_PATHWAY	Generic Transcription Pathway	RT		7	2.8E-3	1.7E-1
<input type="checkbox"/>	REACTOME_PATHWAY	RNA Polymerase II Transcription	RT		7	4.6E-3	1.7E-1
<input type="checkbox"/>	REACTOME_PATHWAY	Gene expression (Transcription)	RT		7	8.2E-3	2.0E-1
<input type="checkbox"/>	GAD_DISEASE_CLASS	UNKNOWN	RT		5	1.5E-1	9.9E-1
Annotation Cluster 7		Enrichment Score: 1.76	G		Count	P_Value	Benjamini
<input type="checkbox"/>	PUBMED_ID	18029348	RT		6	1.8E-5	3.4E-3
<input type="checkbox"/>	UP_KW_PTM	Isopeptide bond	RT		6	5.4E-3	1.9E-2
<input type="checkbox"/>	PUBMED_ID	15342556	RT		3	7.9E-3	4.8E-1
<input type="checkbox"/>	PUBMED_ID	26496610	RT		3	1.0E-1	1.0E0
<input type="checkbox"/>	GOTERM_MF_DIRECT	metal ion binding	RT		4	4.5E-1	1.0E0
<input type="checkbox"/>	UCSC_TFBS	TAL1ALPHA E47	RT		3	7.9E-1	1.0E0

Credit: XKCD
comics

WHY ARE THERE SLAVES IN THE BIBLE

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

QUESTIONS

FOUND IN GOOGLE AUTOCOMPLETE

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

WHY AREN'T MY
ARMS GROWING



WHY ARE THERE WEEKS
WHY DO I FEEL DIZZY

WHY AREN'T ECONOMISTS RICH

WHY DO AMERICANS CALL IT SOCCER

WHY ARE MY EARS RINGING

WHY ARE THERE SO MANY AVENGERS

WHY ARE THE AVENGERS FIGHTING THE X MEN

WHY IS WOLVERINE NOT IN THE AVENGERS

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 ARE THERE MALE AND FEMALE BIKES

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 LYING GOOD

WHY ARE THERE
GHOSTS



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 AREN'T
THERE GUNS IN
HARRY POTTER



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

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 BRIDESMAIDS

WHY DO DYING PEOPLE REACH UP

WHY AREN'T THERE VARIKOSE ARTERIES

WHY ARE OLD KINGDOMS 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 IS SEX
SO IMPORTANT

