

# **Statistical Analysis of RNA-Seq Data**

**Jeff Glaubitz, Qi Sun**

**Bioinformatics Facility  
Cornell University**

# RNA-Seq Data

If you have a reference genome



**STAR**



**HTSeq**

(built into STAR)



Without reference genome

**Trinity**

Build reference



**Bowtie2**

Map reads



**RSEM**

Count reads



**DESeq2 or EdgeR**

# RNA-Seq Statistics:

- Normalization between samples;
- **D**ifferentially **E**xpressed Genes (**DE**);

Genes	Control	Treated
Gene A	10	30
Gene B	30	90
Gene C	5	15
Gene D	1	3
Gene N	80	240

**126**

**378**

# Normalization

Genes	Control	Treated
Gene A	10	30
Gene B	30	90
Gene C	5	15
Gene D	1	3
Gene N	80	240

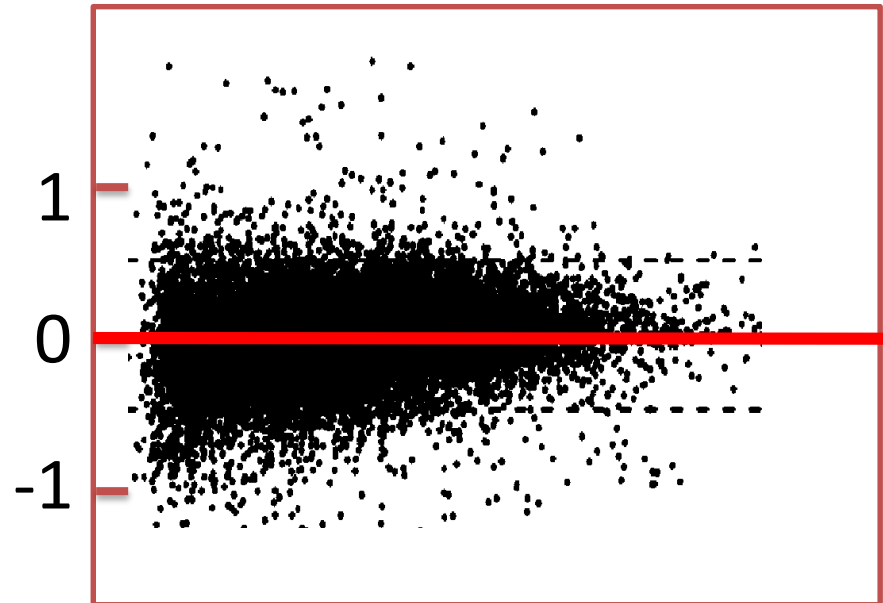
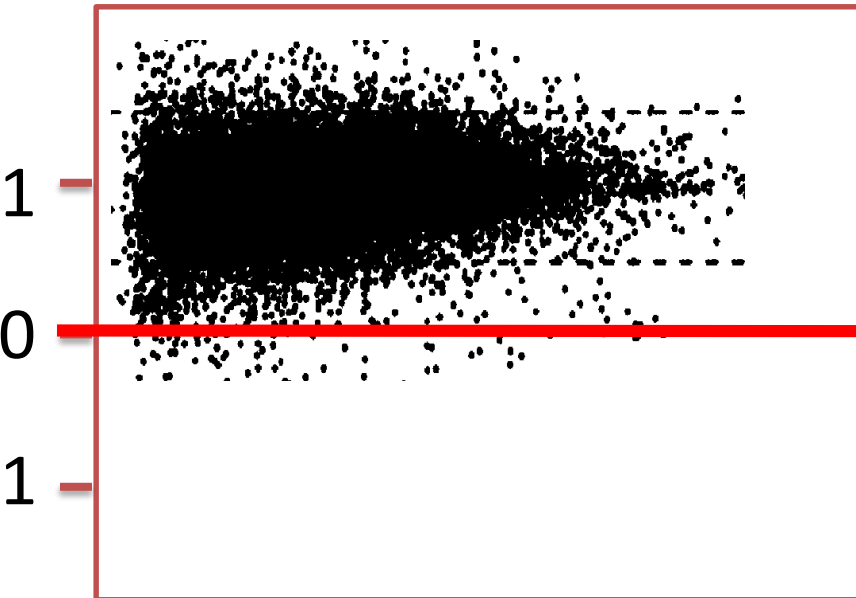
126

378

## MA Plots between samples

Before normalization

After normalization



- Y axis: log ratio of expression level between two conditions;
- With the assumption that most genes are expressed equally, the log ratio should mostly be close to 0

# Simple normalization

**CPM** (Count Per Million Reads)

Normalized by:

- Total fragment count;

**FPKM** (Fragments Per Kilobase Of Exon Per Million Fragments)

Normalized by:

- Total fragment count;
- Gene length (kb);

CPM : Not normalized by gene length. Longer genes tend to have higher CPM values than shorter genes. But that is ok, as in RNA-Seq experiments, we do not compare between genes, only compare the same gene between different samples.

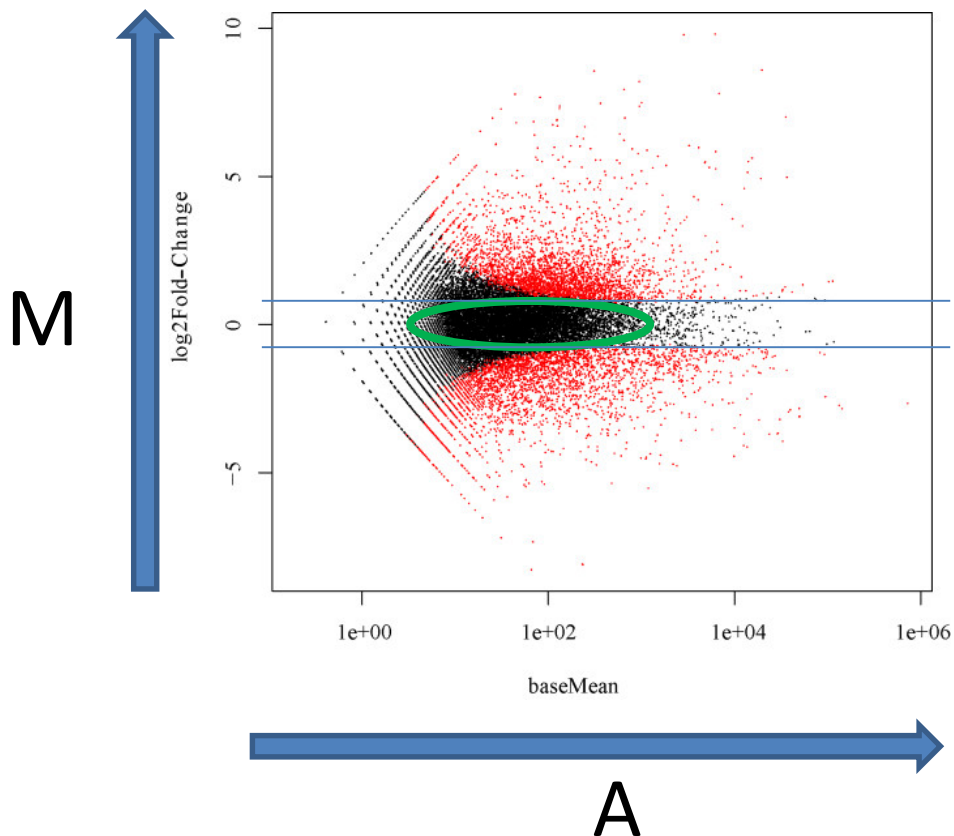
# Simple normalization could fail

Genes	Control	Treated
Gene A	10	30
Gene B	30	90
Gene C	5	15
Gene D	1	3
Gene N	1000	240

**1046**                      **378**

# TMM normalization

(Trimmed mean of M-values )



$$M = \log_2(\text{Test}/\text{Test\_total}) - \log_2(\text{Ref}/\text{Ref\_total})$$

$$A = 0.5 * \log_2(\text{Test}/\text{Test\_total} * \text{Ref}/\text{Ref\_total})$$

Effective library size

# DESeq2 normalization

## 1. For each gene, calculate geometric mean

	Sample 1	Sample 2	Sample 3	Sample 4	Sample 5	Sample 6	Geomean
Gene1	34	56	23	12	10	30	<b>23</b>
Gene 2	10	6	7	11	12	8	<b>9</b>
.....							
Gene n	65	78	67	34	56	23	<b>50</b>

## 2. For each gene, calculate ratio to geometric mean

	Sample 1	Sample 2	Sample 3	Sample 4	Sample 5	Sample 6
Gene1	1.5	2.4	1.0	0.5	0.4	1.3
Gene 2	1.1	0.7	0.8	1.3	1.4	0.9
.....						
Gene n	1.3	1.6	1.4	0.7	1.1	0.5

## 3. Take median of these ratio as sample normalization factor

<b>1.3</b>	<b>1.6</b>	<b>1</b>	<b>0.7</b>	<b>1.1</b>	<b>0.9</b>
------------	------------	----------	------------	------------	------------



# Differentially Expressed Genes

## Expression level of gene 1

**Control:**

*Repeat 1* **24**

*Repeat 2* **25**

*Repeat 3* **27**

**Treated:**

*Repeat 1* **23**

*Repeat 2* **26**

*Repeat 3* **102**

**Question : is this a DE gene?**

You might get different answers depending of which software you run.

# Available RNA-seq analysis packages for DE

**TABLE 2.** A summary of the recommendations of this paper

	Agreement with other tools <sup>a</sup>	WT vs. WT FPR <sup>b</sup>	Fold-change threshold (T) <sup>c</sup>	Tool recommended for: (# good replicates per condition) <sup>d</sup>		
				≤3	≤12	>12
<i>DESeq</i>	Consistent	Pass	0	-	-	Yes
			0.5	-	Yes	Yes
			2.0	Yes	Yes	Yes
<i>DESeq2</i>	Consistent	Pass	0	-	-	Yes
			0.5	Yes	Yes	Yes
			2.0	Yes	Yes	Yes
<i>EBSeq</i>	Consistent	Pass	0	-	-	Yes
			0.5	-	Yes	Yes
			2.0	Yes	Yes	Yes
<i>edgeR (exact)</i>	Consistent	Pass	0	-	-	Yes
			0.5	Yes	Yes	Yes
			2.0	Yes	Yes	Yes
<i>Limma</i>	Consistent	Pass	0	-	-	Yes
			0.5	-	Yes	Yes
			2.0	Yes	Yes	Yes
<i>cuffdiff</i>	Consistent	Fail				
<i>BaySeq</i>	Inconsistent	Pass				
<i>edgeR (GLM)</i>	Inconsistent	Pass				
<i>DEGSeq</i>	Inconsistent	Fail				
<i>NOISeq</i>	Inconsistent	Fail				
<i>PoissonSeq</i>	Inconsistent	Fail				
<i>SAMSeq</i>	Inconsistent	Fail				

From: Schurch *et al.* 2016. *RNA* 22:839-851

# Why DESeq2?

1. Top method recommended by Schurch *et al.* (2016), along with *EdgeR (exact)*
2. Cutting-edge tool widely used and accepted: 11,934 citations (Google Scholar on Oct 25, 2019)
3. Documentation (and papers) very thorough and well-written
4. The first author (Mike Love) provides amazing support! Most questions that you Google (e.g., [support.bioconductor.org](http://support.bioconductor.org)) are clearly and definitively answered by the author himself.
5. See <https://mikelove.wordpress.com/2016/09/28/deseq2-or-edger/>
6. R functions in *DESeq2* package are intuitive to R users (and modifiable). Defining the experimental design is easy and intuitive, even for complex, multifactor designs:

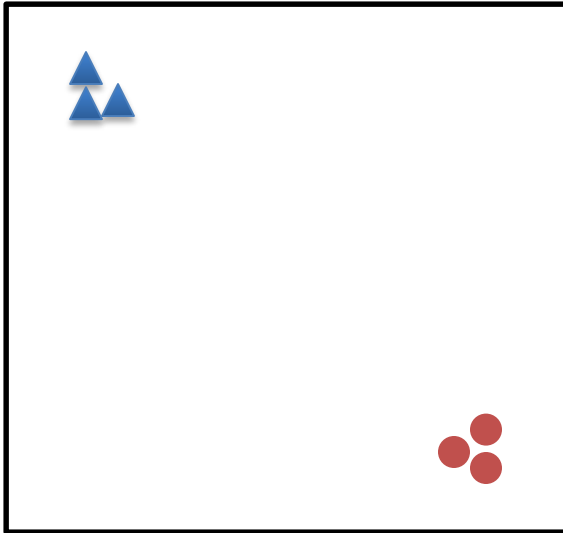
```
design= ~ batch + weight + genotype + treatment + genotype:treatment
```

# Experimental design

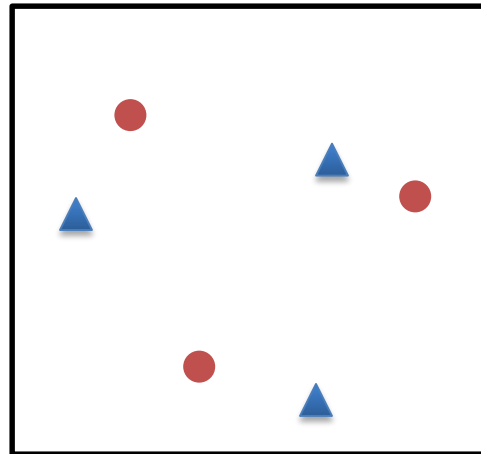
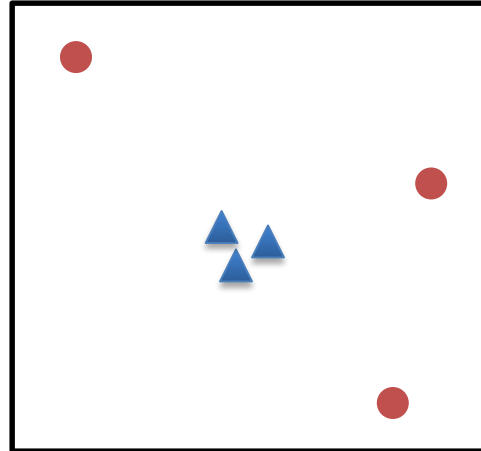
▲: Control samples

●: Treated samples

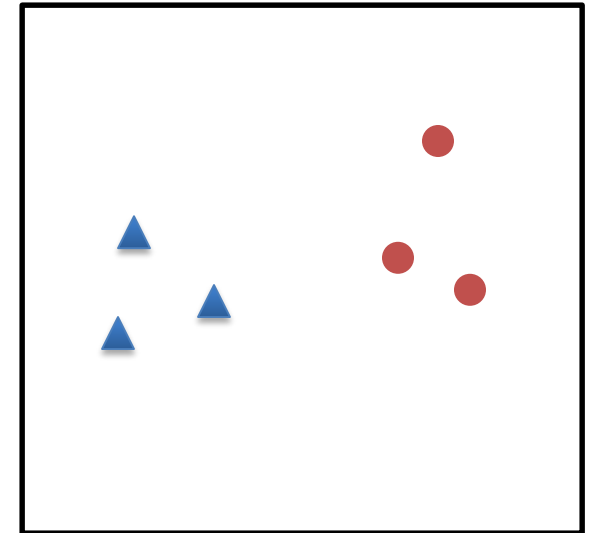
## PCA Plots



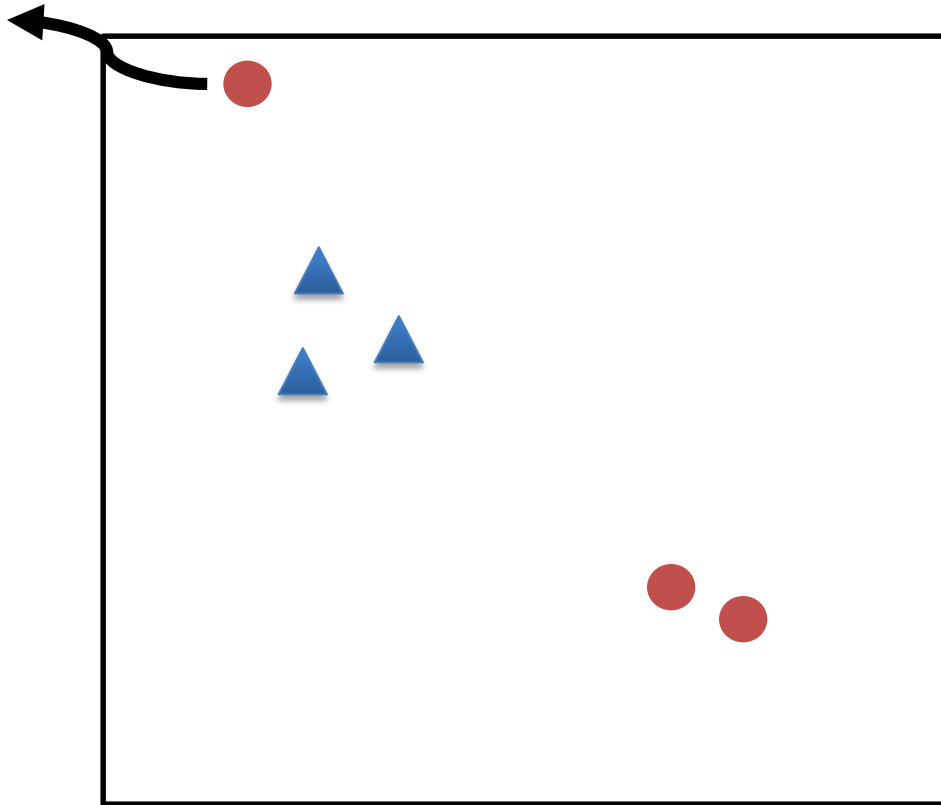
Too many DE genes



Too few DE genes



Remove outlier samples

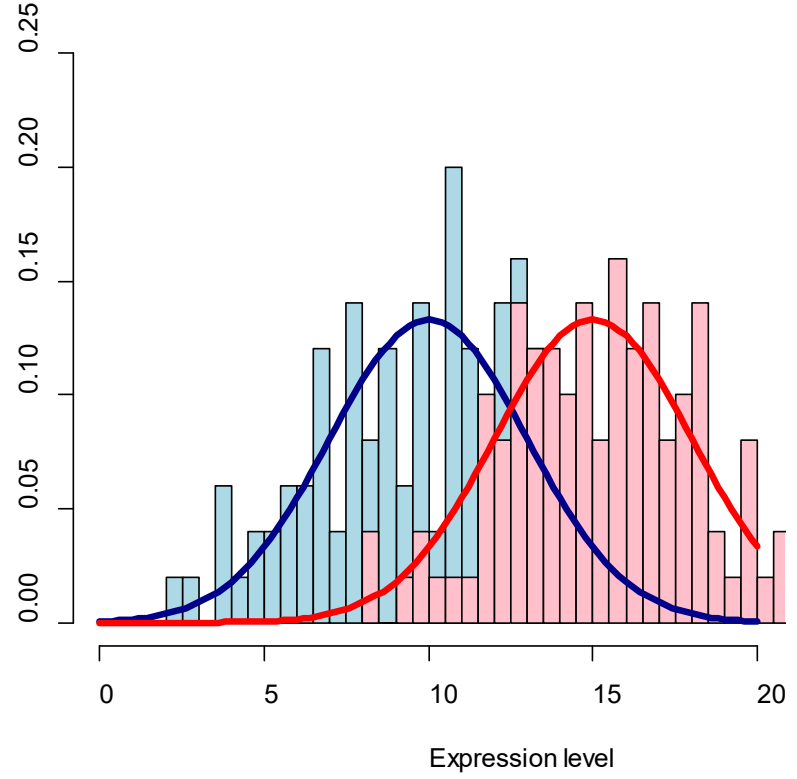
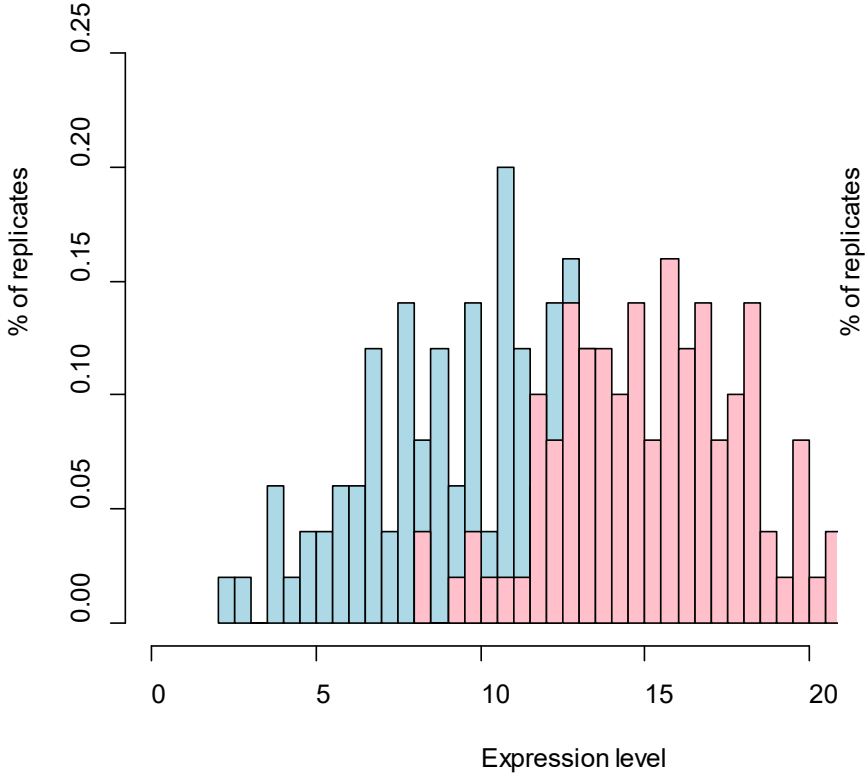


# Biological vs. technical replicates



Scenario	Replicate Type
Split tissue sample evenly into 2 RNA preps	Technical
Split RNA sample into two library preps	Technical
Split library across two sequencing flow cells	Technical
RNA prep from different leaves on same plant	Technical/Biological
Different clones of the same genotype in same treatment condition	Biological
Different genotypes in same treatment condition	Biological

# Differentially expressed genes

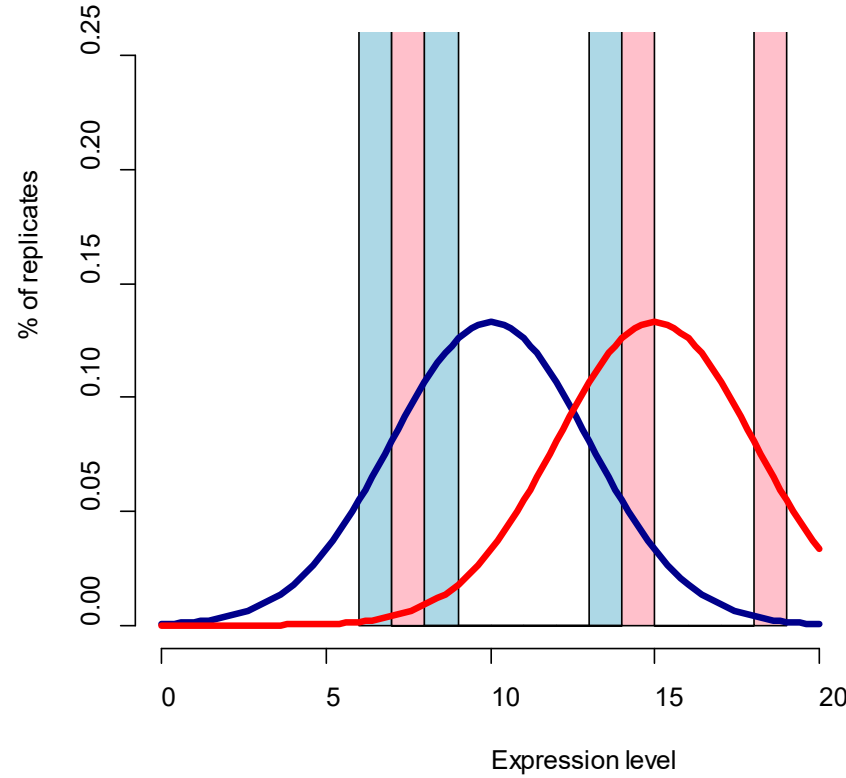
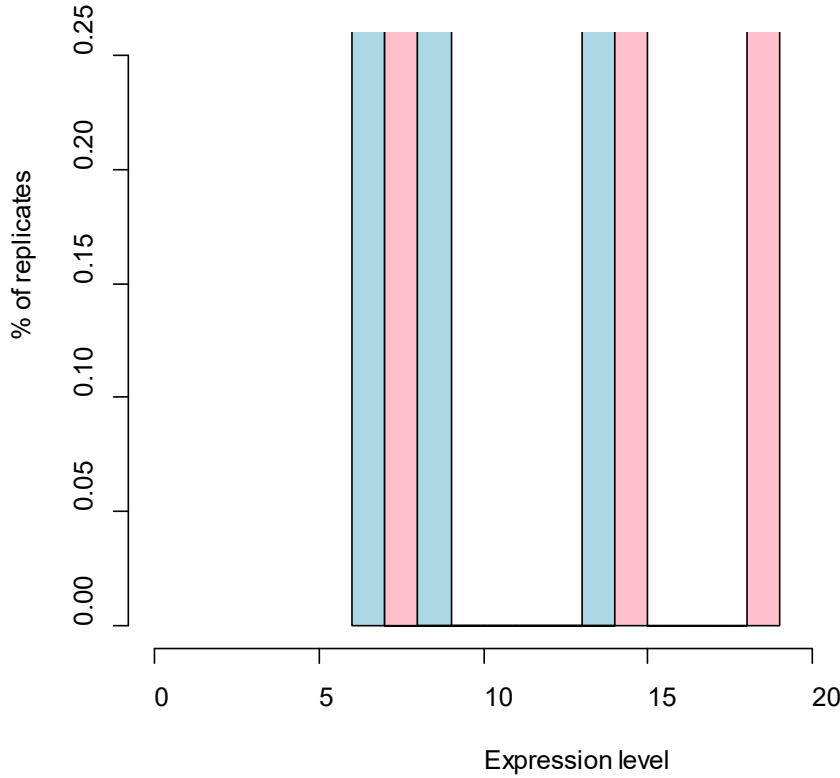
If we could do 100 biological replicates,





**Distribution of Expression Level of A Gene**

-  Control samples
-  Treated samples

The reality is, often we can only afford 3 replicates,



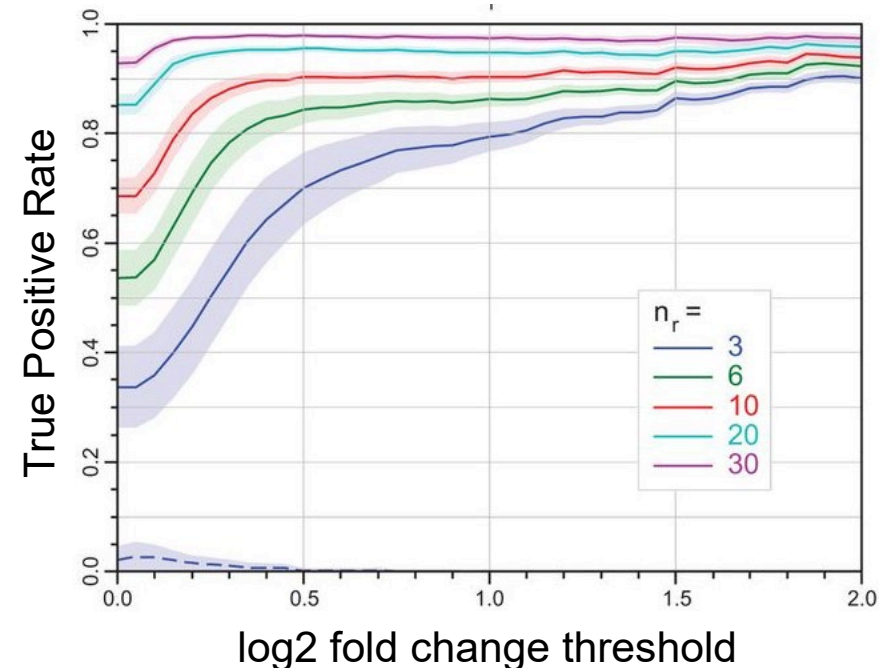
### Distribution of Expression Level of A Gene

-  Control samples
-  Treated samples

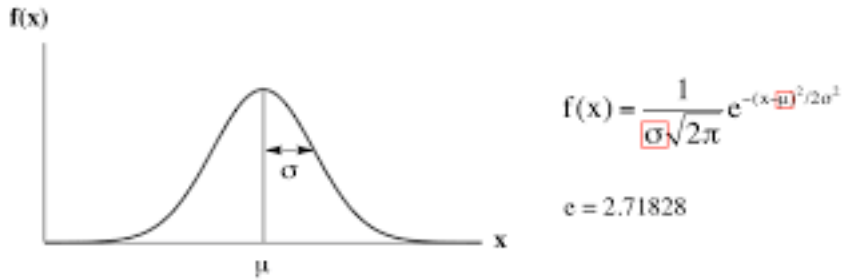


# How many biological replicates?

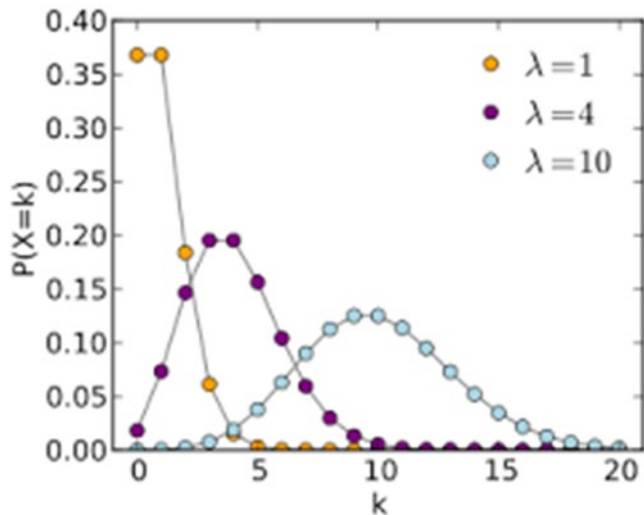
- 3 replicates are the *bare minimum* for publication
- Schurch *et al.* (2016) recommend at least 6 replicates for adequate statistical power to detect DE
- Depends on biology and study objectives
- Trade off with sequencing depth
- Some replicates might have to be removed from the analysis because poor quality (outliers)



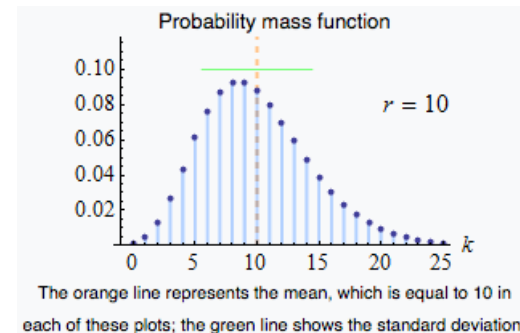
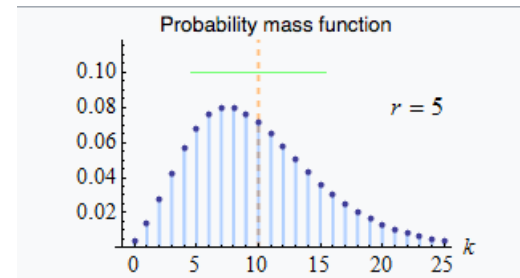
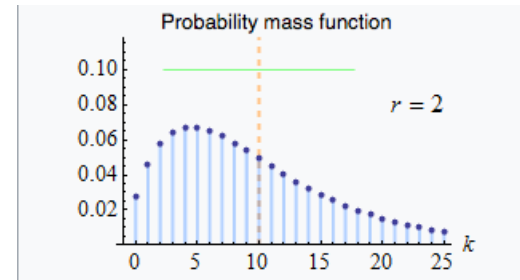
# Hypothesis tests require accurate statistical model



Gaussian (Normal)

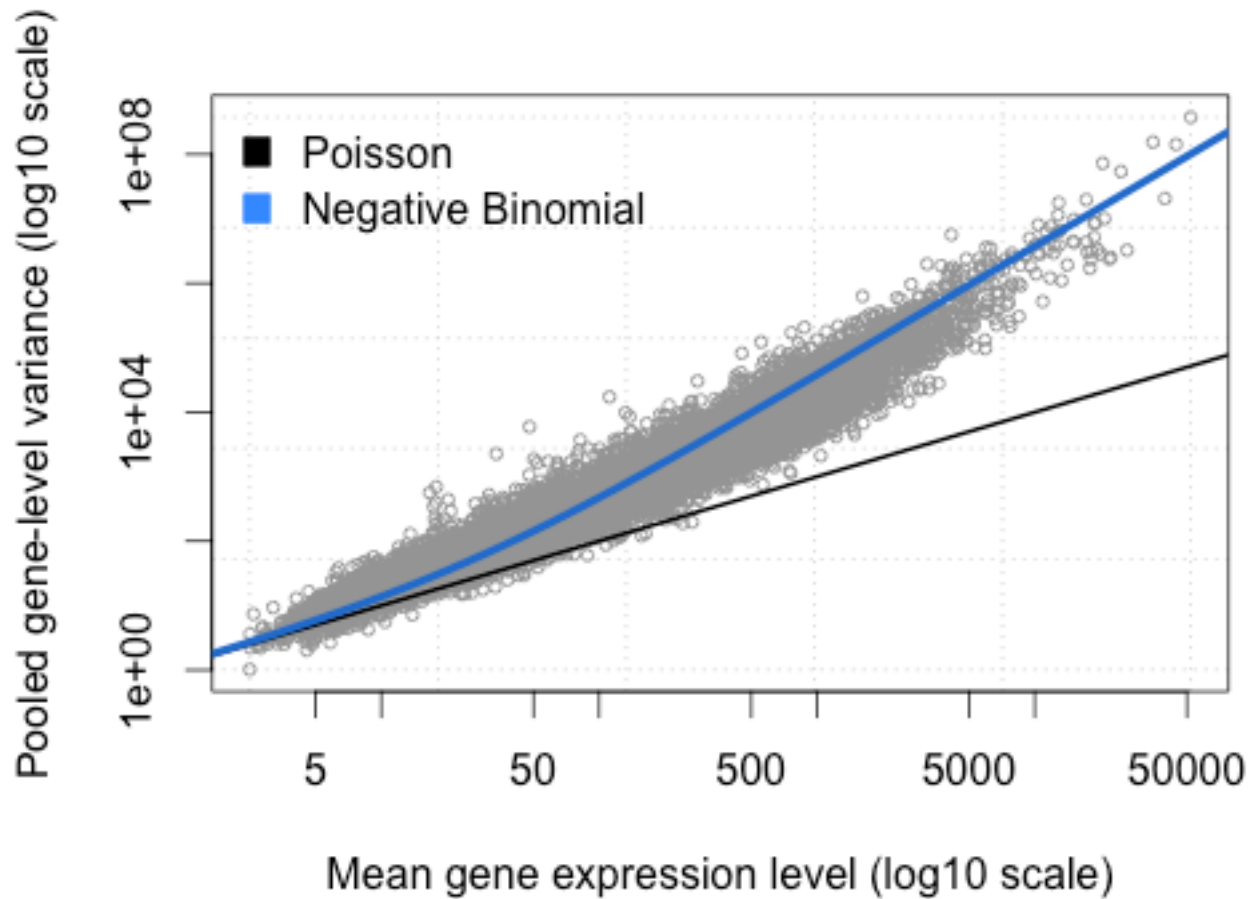


Poisson (variance=mean)



Negative binomial  
(variance > mean)

# Negative binomial best fit for RNA-Seq data



# DESeq2 fits a negative binomial model

Raw count for gene  $i$  in sample  $j$

Controls the variance

$$K_{ij} \sim \text{NB}(\text{mean} = \mu_{ij}, \text{dispersion} = \alpha_i)$$

$$\mu_{ij} = s_{ij} q_{ij}$$

Normalization ("size") factor

Normalized count

Design matrix

- Control or Treatment?
- Batch (e.g., flow cell or plate)
- Other co-factors (e.g., sex)

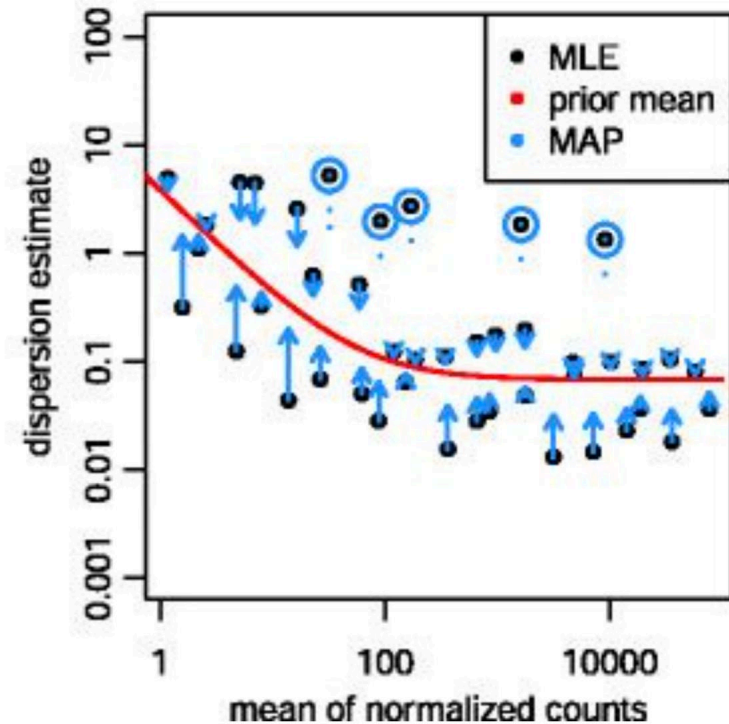
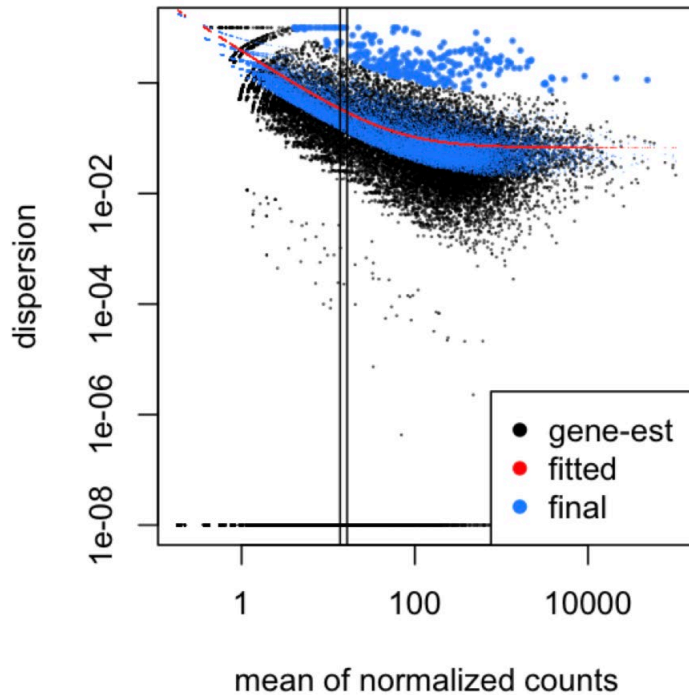
$$\log q_{ij} = \sum_r x_{jr} \beta_{ir}$$

Coefficient  $r$

GLM coefficients

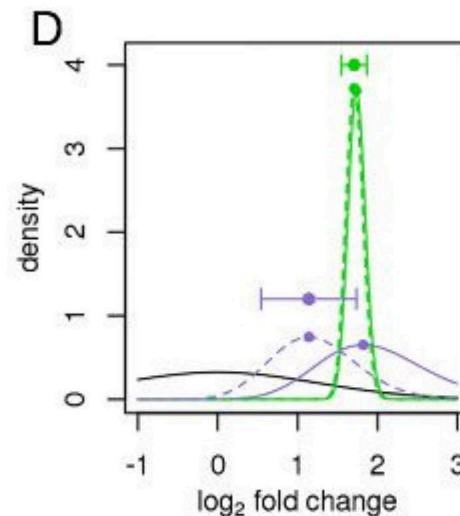
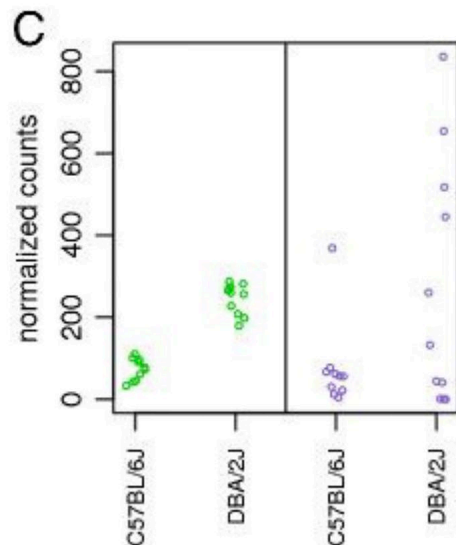
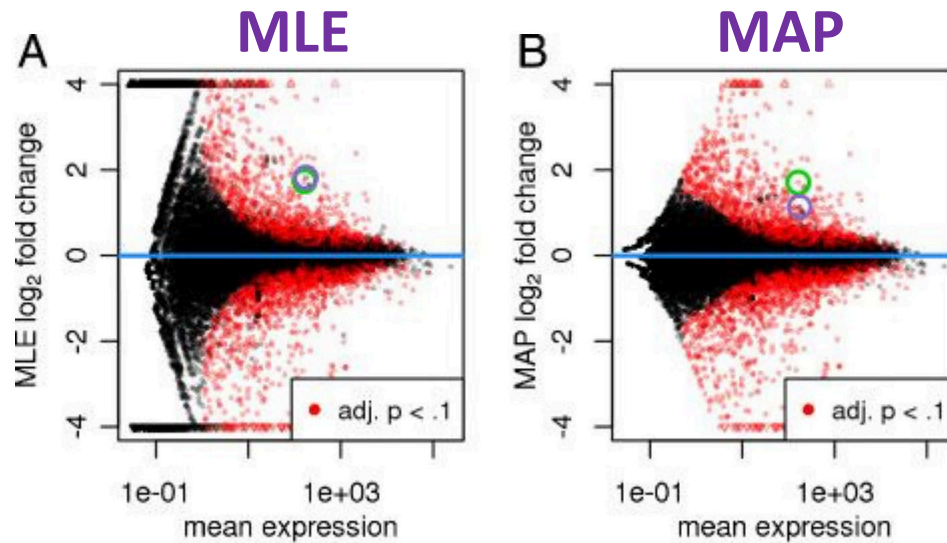
- One for each Design matrix element = strength of effect
- Overall expression strength of gene
- **log2 fold change**

# DESeq2: Empirical Bayes shrinkage of dispersion



- Not enough replicates to estimate variance ("dispersion") for individual genes
- Borrow information from genes of similar expression strength among the replicates
- Genes with very high dispersion left as is (violate model assumptions?)

# DESeq2: Empirical Bayes shrinkage of fold change



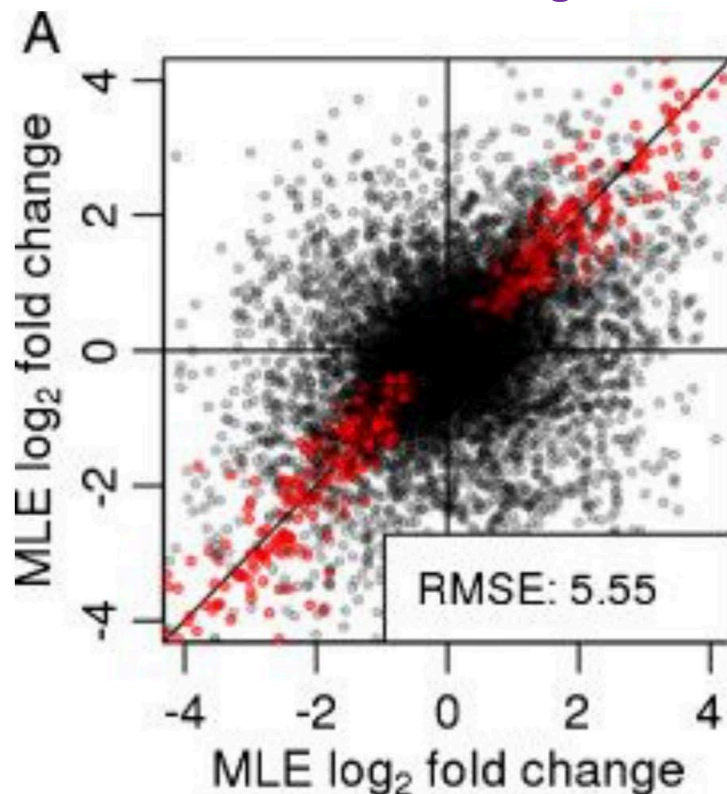
Normalized Counts

Likelihood & Posterior Densities

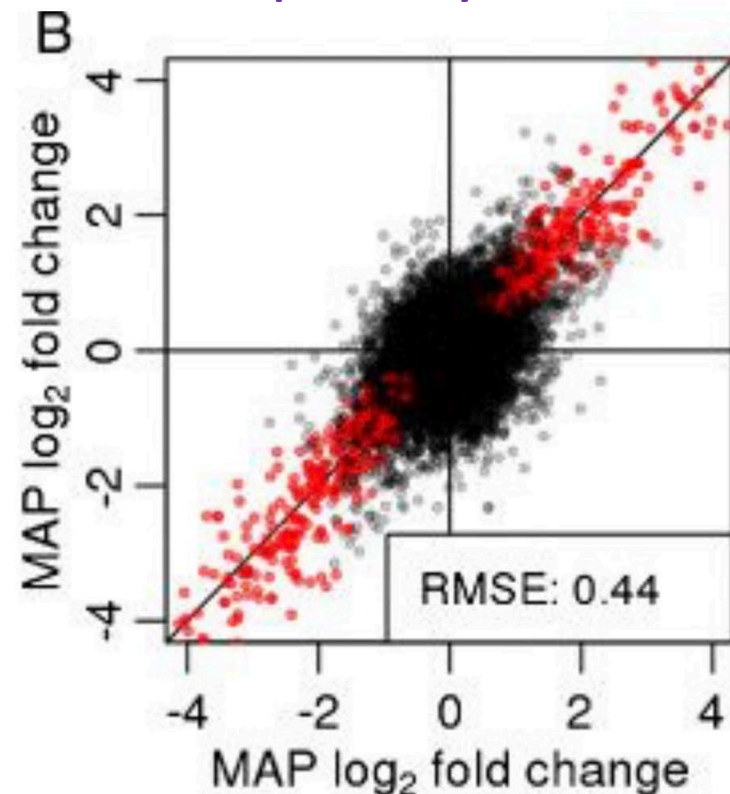
# DESeq2: Empirical Bayes shrinkage of log fold change improves reproducibility

- Large data-set split in half  $\rightarrow$  compare log<sub>2</sub> fold change estimates for each gene

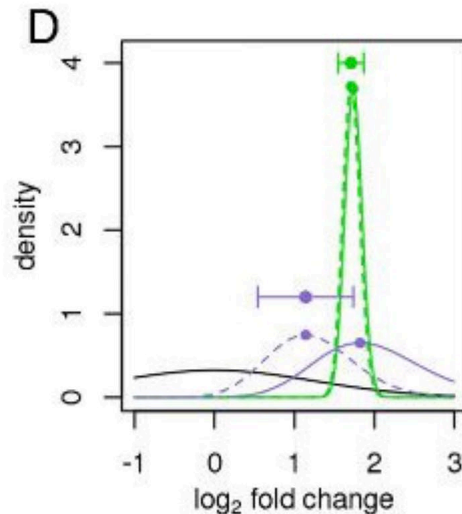
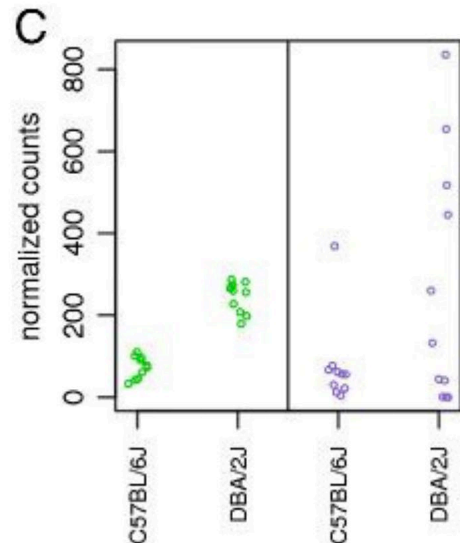
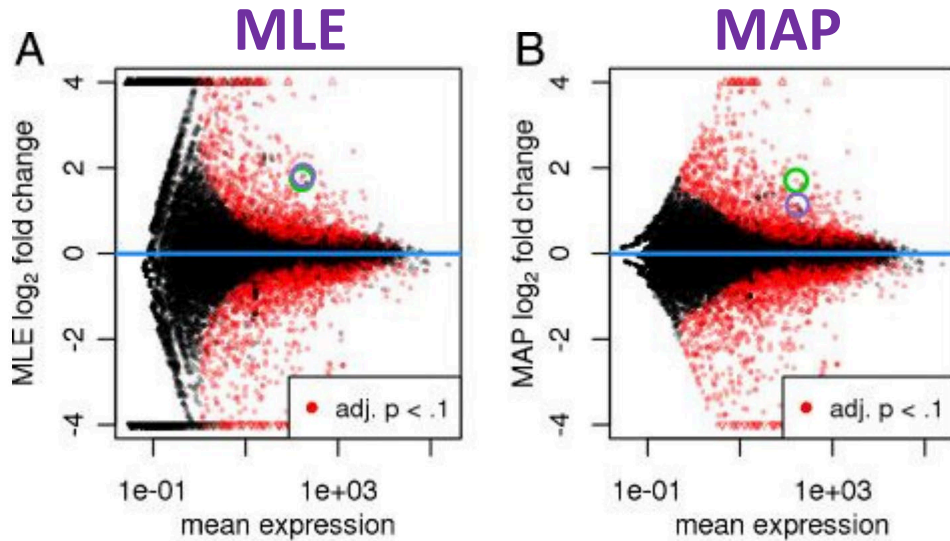
Before shrinkage:



After empirical Bayesian shrinkage:



# DESeq2: Statistical test for DE



## Test for DE:

- $(shrunkenLFC) / (stdErr) = Z \text{ stat}$
- Z stat follows std. normal dist.
- $p$  value for Z stat (LFC) obtained from standard normal distribution
- $p$  values adjusted for multiple testing using Benjamini and Hochberg (1995) procedure
  - Controls *false discovery rate* (FDR)

Normalized Counts

Likelihood & Posterior Densities



# False Discovery Rate

		Truth		Total
		Different	Same	
Experiment	Different	TP	FP	R
	Same	FN	TN	m - R
	Total	P	N	m

- **m**: total number of tests (e.g., genes)
- **N**: number of true null hypotheses
- **P**: number of true alternate hypotheses
- **R**: number of rejected null hypotheses (“discoveries”)
- **TP**: number of true positives (“true discoveries”)
- **TN**: number of true negatives
- **FP**: number of false positives (“false discoveries”) (Type I error)
- **FN**: number of false negatives (Type II error)
- **FDR = “false discoveries” / “discoveries” =  $FP / (FP + TP)$**

# DESeq2: Automated independent filtering of genes

- DESeq2 automatically omits weakly expressed genes from the multiple testing procedure
  - Fewer tests increase statistical power → more discoveries
- LFC estimates for weakly expressed genes very noisy
  - Very little chance that these will be detected as DE
- Threshold overall counts (filter statistic) optimized for target FDR (default FDR = 0.1)

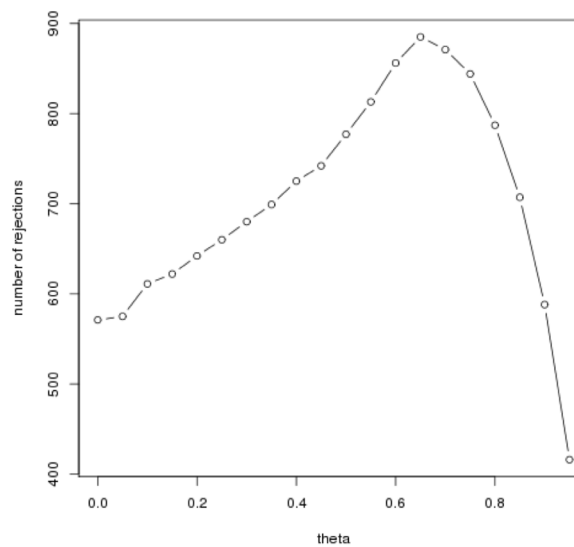
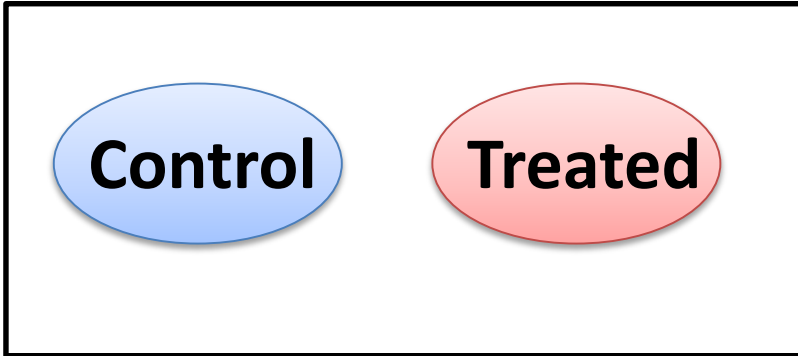


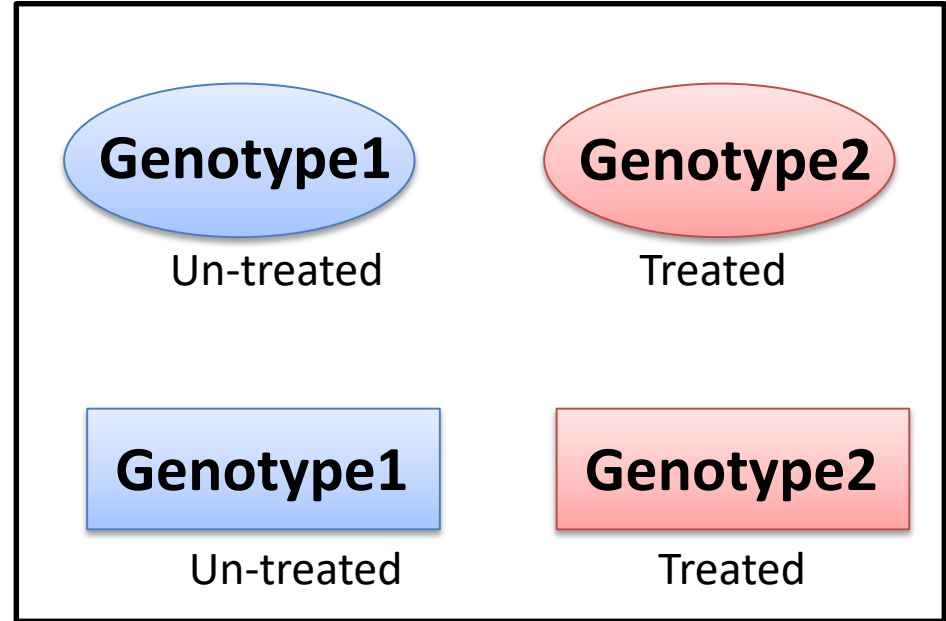
Figure 1: The number of rejected tests for FDR less than 0.1 plotted over theta, the quantiles of the filter statistic.

# Type of analyses

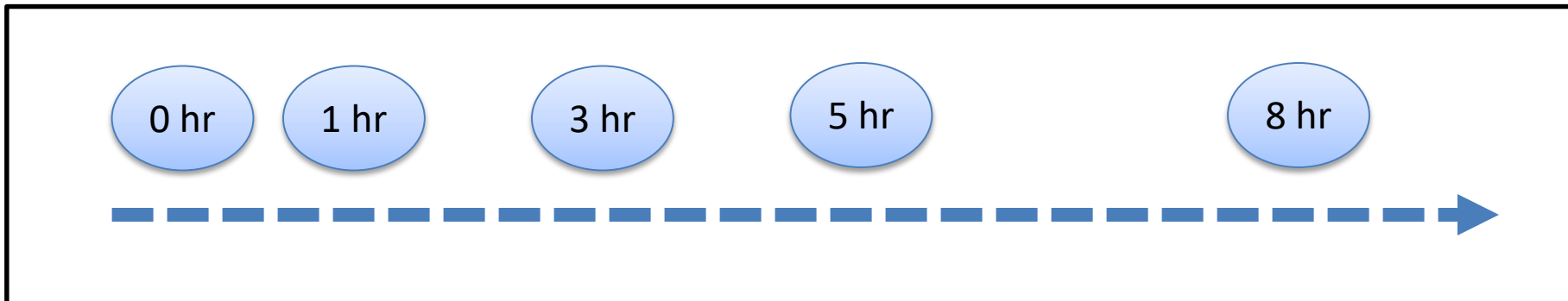
## One factor



## Two factors



## Time series



## DESeq2: Design specifications

```
dds <- DESeqDataSetFromMatrix(countData = cts, colData = coldata,  
design= ~ treatment)
```

```
dds <- DESeqDataSetFromMatrix(countData = cts, colData = coldata,  
design= ~ batch + treatment)
```

*# Model genotype by treatment interaction:*

```
dds <- DESeqDataSetFromMatrix(countData = cts, colData = coldata,  
design= ~ batch + genotype + treatment + genotype:treatment)
```

*# Likelihood ratio test for genotype by treatment interaction:*

```
ddsLRT <- DESeq(dds, test="LRT", reduced= ~ batch + genotype + treatment )
```

```
resLRT <- results(ddsLRT)
```

# DESeq2: Output of DE analysis

1	gene	baseMean	log2FoldChange	lfcSE	stat	pvalue	padj
2	gene27816	416.9148177	-2.511490321	0.428727267	-5.858013979	4.68E-09	5.11E-05
3	gene27620	101.8253191	-2.778767979	0.553389846	-5.021357008	5.13E-07	0.001440413
4	gene31204	365.5088989	-1.143004071	0.227873895	-5.015950038	5.28E-07	0.001440413
5	gene4446	125.745322	-3.205715488	0.637910711	-5.025335725	5.03E-07	0.001440413
6	gene1					20E-06	0.002553055
7	gene1					40E-06	0.002553055
8	gene3					77E-06	0.002756968
9	gene8					23E-06	0.003038386
10	gene2					89E-06	0.005937898
11	gene2					99E-06	0.007929629
12	gene9777	207.8848249	0.90630494	0.202189437	4.482454445	7.38E-06	0.007929629
13	gene21278	357.2070995	-1.375680007	0.309791352	-4.440666269	8.97E-06	0.008159457
14	gene34591	77.02015308	-2.724251177	0.632181028	-4.309289679	1.64E-05	0.013754978

**Get list of interesting genes by filtering on:**

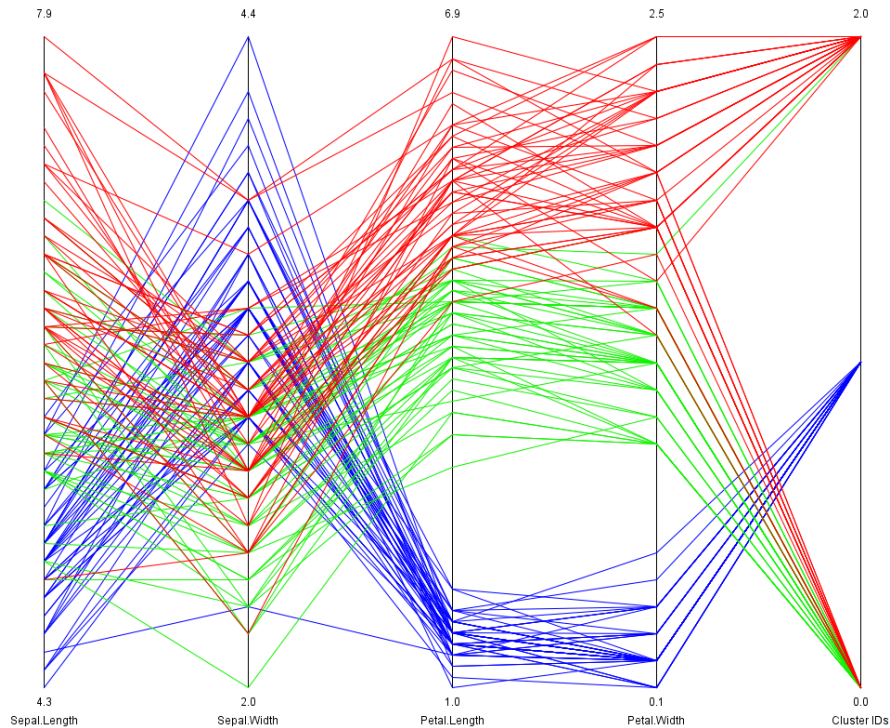
- 1.  $p_{adj}$  (FDR) < 0.05, and/or**
- 2.  $\log_2\text{FoldChange}$  < -1 or >1, and/or**
- 3. baseMean (optional)**

...bottom of file = genes excluded from multiple testing:

22902	gene599	0.182683434	1.286409862	4.409114495	0.291761501	0.770468983	NA
22903	gene11602	0.175175275	1.286401869	4.409114509	0.291759687	0.77047037	NA
22904	gene30325	0.168776837	1.606989135	4.409114563	0.364469807	0.715507216	NA
22905	gene35203	0.159702673	1.846921927	4.40661107	0.419125241	0.675124605	NA
22906	gene25371	0.153270142	1.286411877	4.409114495	0.291761958	0.770468634	NA
22907	gene7678	0.141308727	0.93343057	4.407959907	0.211760222	0.832294103	NA
22908	gene13239	0.132221267	1.492522653	2.686412371	0.55558211	0.578496564	NA
22909	gene1935	0.116143364	0.740578083	4.395665987	0.168479153	0.866206343	NA
22910	gene26270	0.107670322	2.315104965	4.402746292	0.525832017	0.599004927	NA
22911	gene30327	0.060455387	0.580738721	4.403240387	0.131888943	0.895072134	NA
22912	gene26805	0.013434773	1.286431679	4.40911449	0.291766449	0.770465199	NA

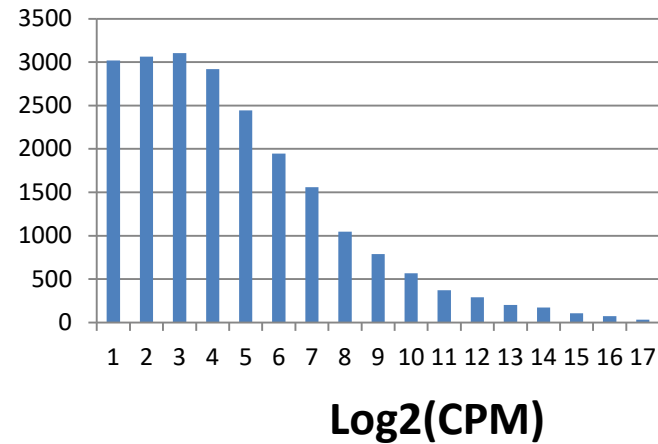
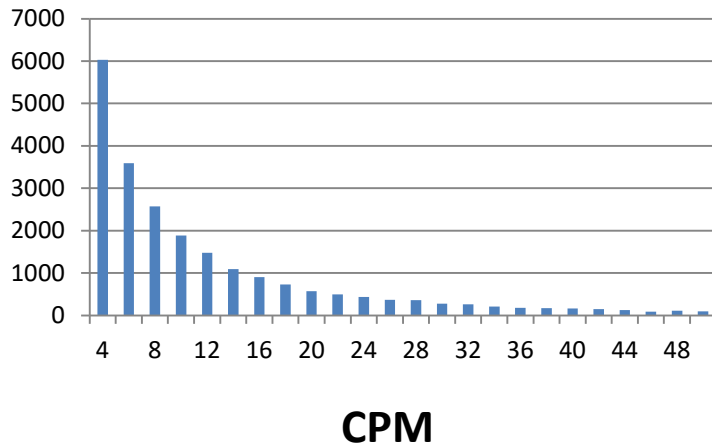
# Clustering analysis

1. Hierarchical
2. K-means
3. Co-expression network



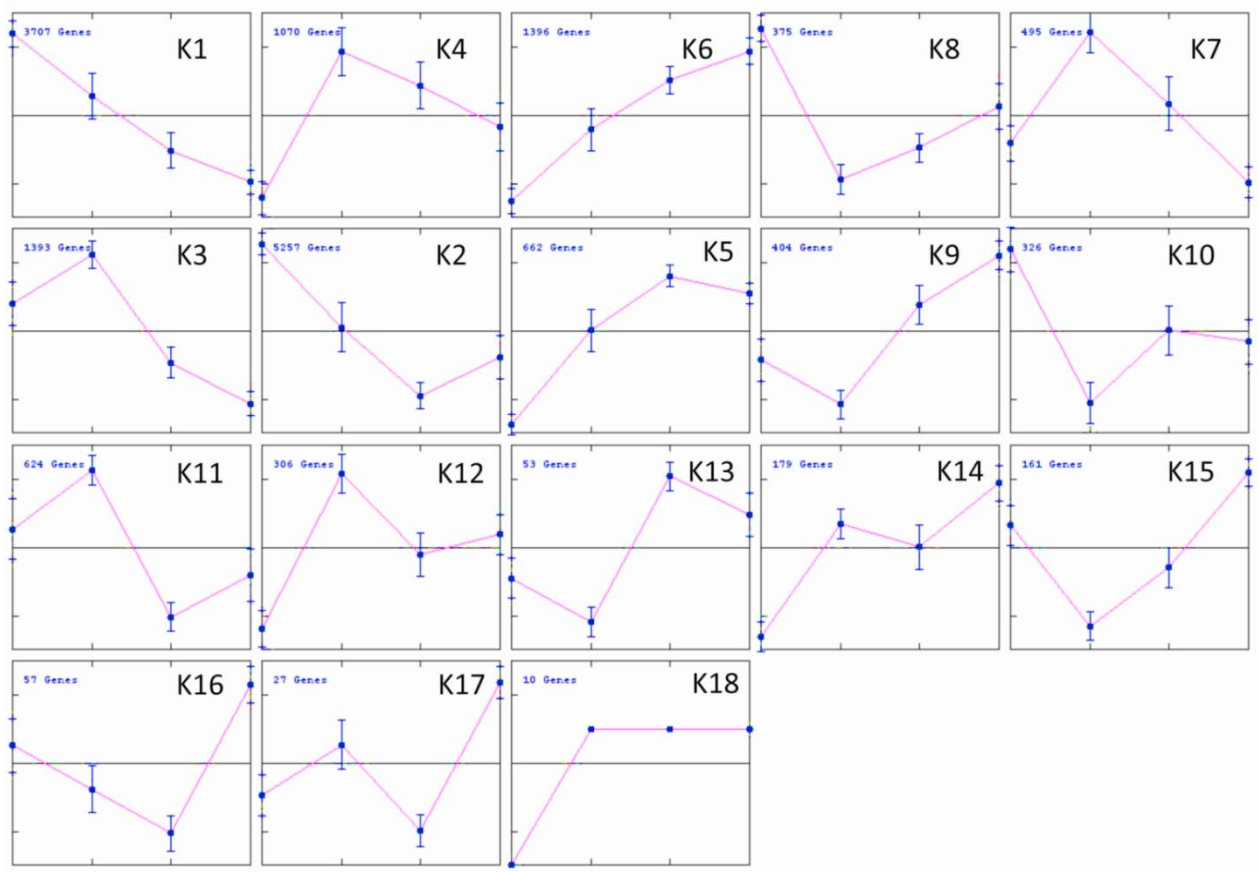
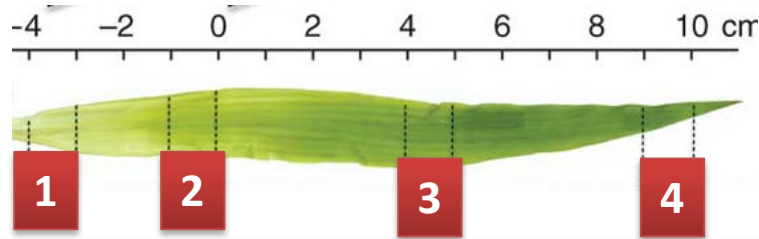
## Prepare data for clustering

**Step 1. LOG transformation of CPM value to improve the distribution**



**Step 2. Remove genes with no variation across samples**

# Clustering analysis on multiple conditions of RNA-seq data



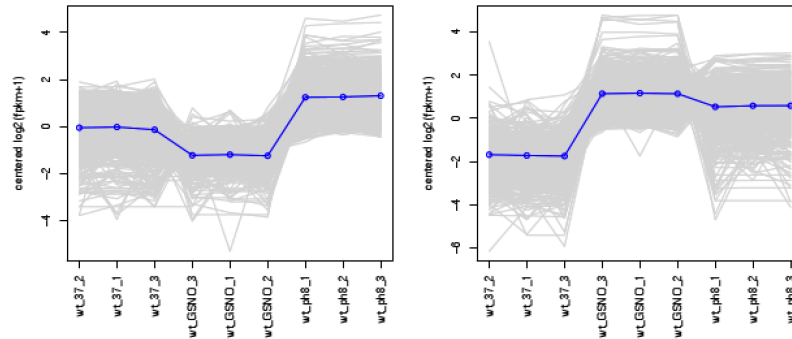


# Hierarchical clustering

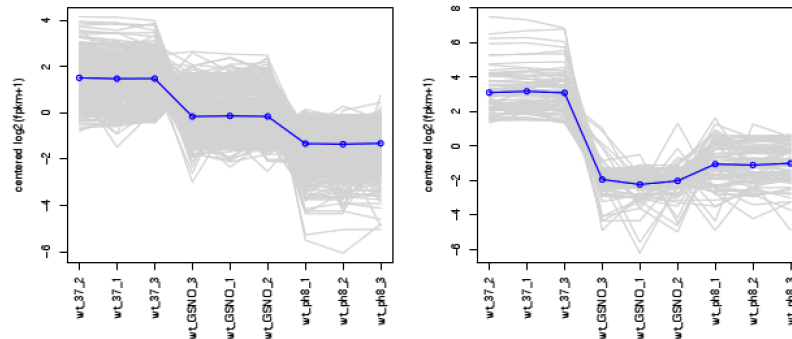
# K-means clustering

```
$TRINITY_HOME/Analysis/DifferentialExpression/  
define_clusters_by_cutting_tree.pl -R  
diffExpr.P0.001_C2.matrix.RData -K 18
```

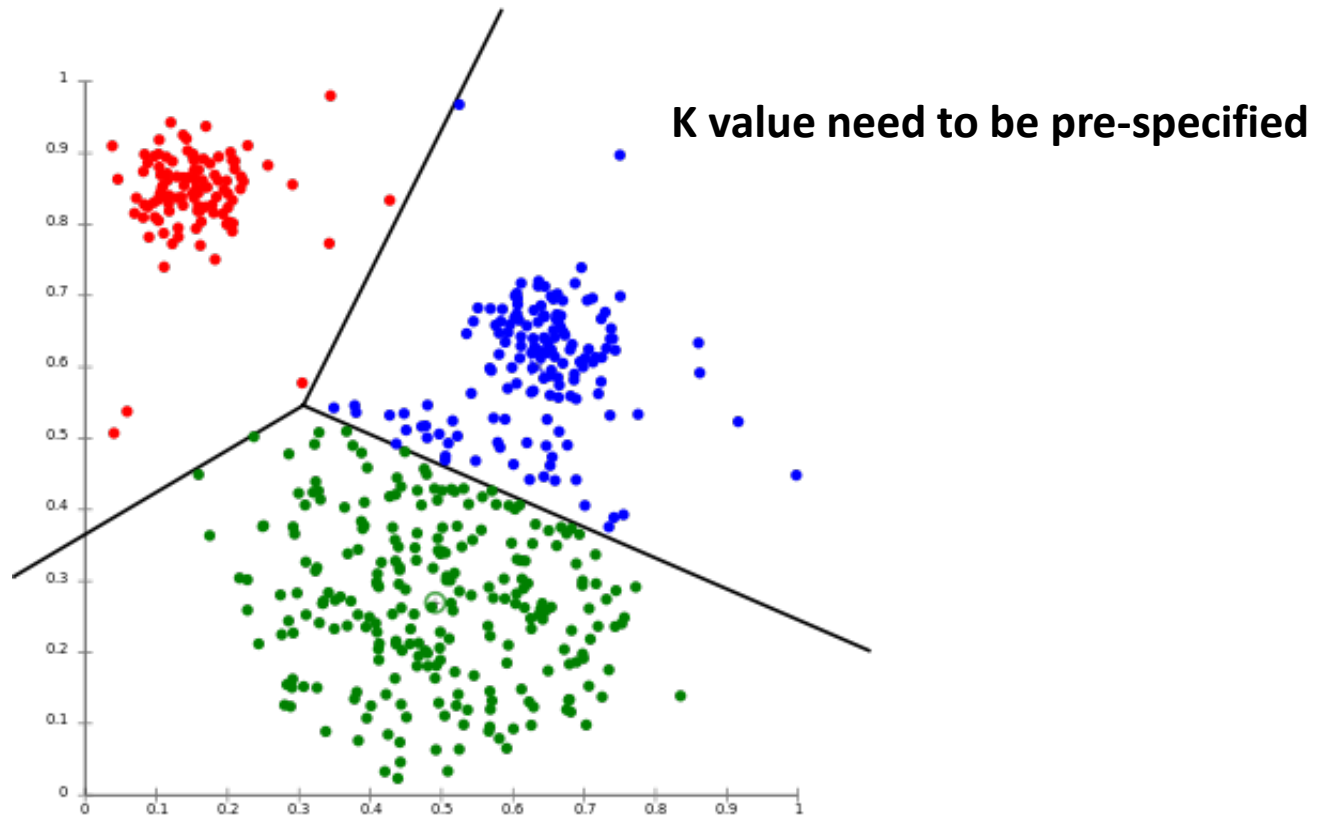
subcluster\_1\_log2\_medianCentered\_fpkm\_matrix, 428 tra | subcluster\_2\_log2\_medianCentered\_fpkm\_matrix, 794 tra



subcluster\_3\_log2\_medianCentered\_fpkm\_matrix, 528 tra | subcluster\_4\_log2\_medianCentered\_fpkm\_matrix, 78 tra



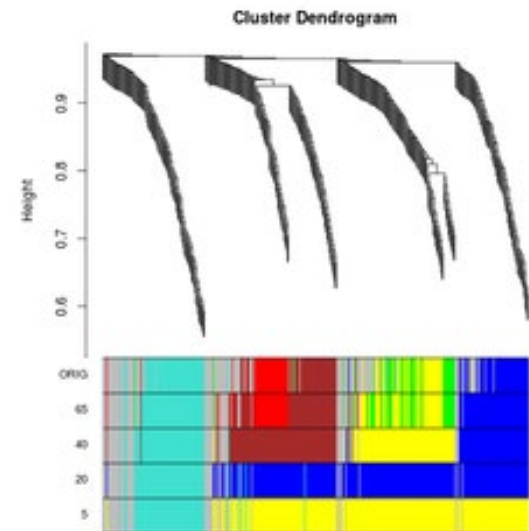
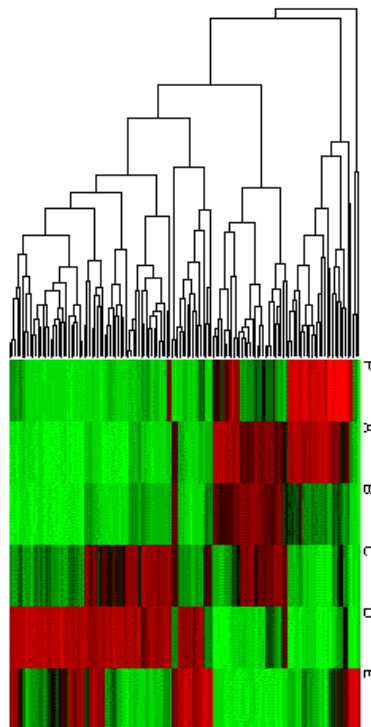
# K-means clustering



# Co-expression network modules

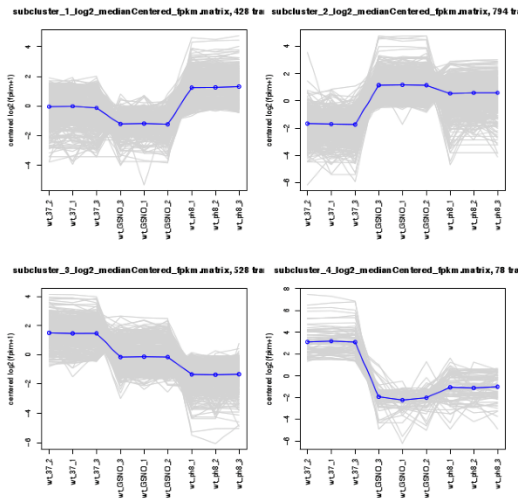
## WGCNA (weighted correlation network analysis)

- transform the initial distance matrix into Topological Overlap Matrix

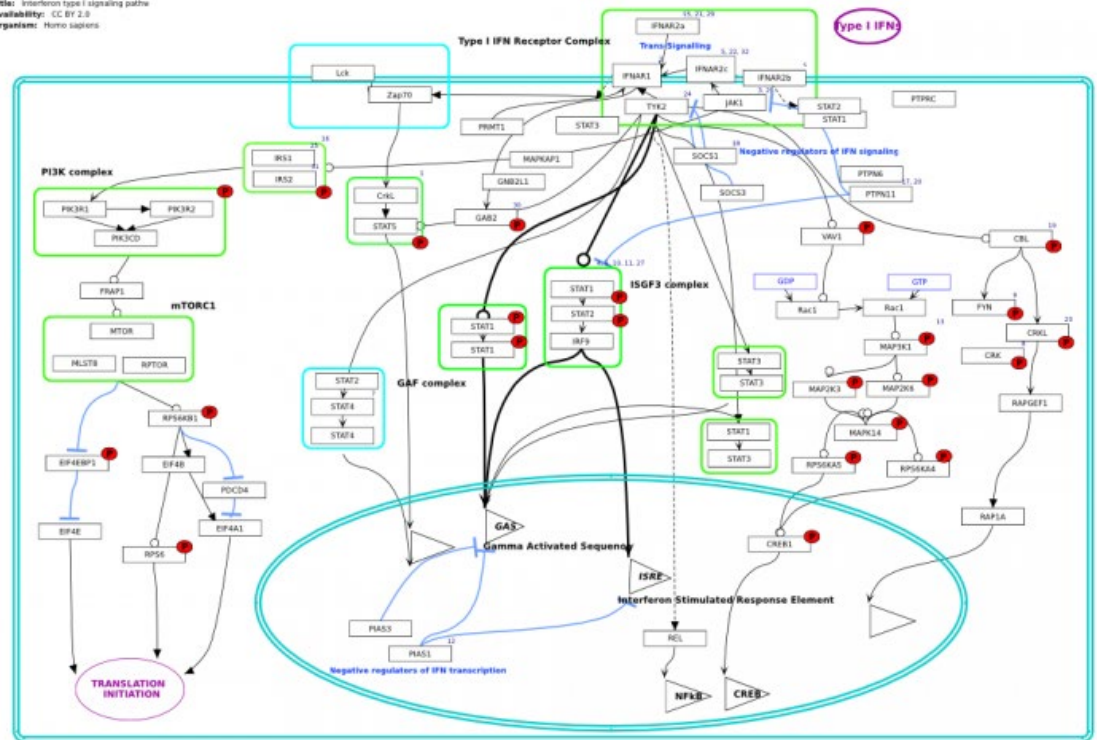


# Gene Set Enrichment Analysis

Will be covered in this workshop:  
**Genome Annotation And Sequence Based Gene  
Function Prediction ([December 12 and 19 2018](#))**



Title: Interferon type I signaling pathway  
Availability: CC BY 2.0  
Organism: Homo sapiens



<https://www.ebi.ac.uk/training/online/course/functional-genomics-ii-common-technologies-and-data-analysis-methods/gene-set-enrichment>