

# Multiple group comparisons for RNA-Seq and stable effect size estimates

Michael Love

Irizarry Lab

Department of Biostatistics

Dana Farber Cancer Institute &  
Harvard School of Public Health

DESeq2

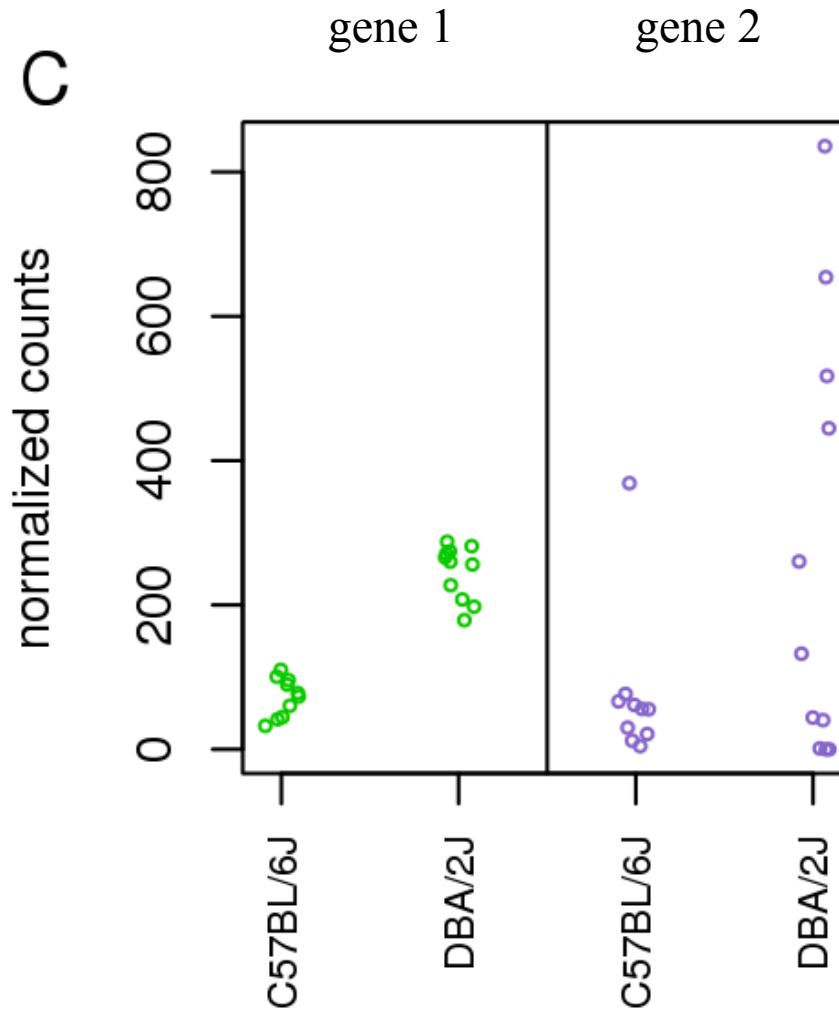
4/2013 Bioc package

2/2014 preprint bioRxiv

Simon Anders, EMBL

Wolfgang Huber, EMBL

# Estimating effect sizes for counts



*effect* of:

- treatment
- species
- tissue, etc.

# Modeling differences in counts

For read count  $K$  for gene  $i$ , sample  $j$ ...

mean dispersion

$$K_{ij} \sim NB(\mu_{ij}, \alpha_i)$$

e.g.

size factor

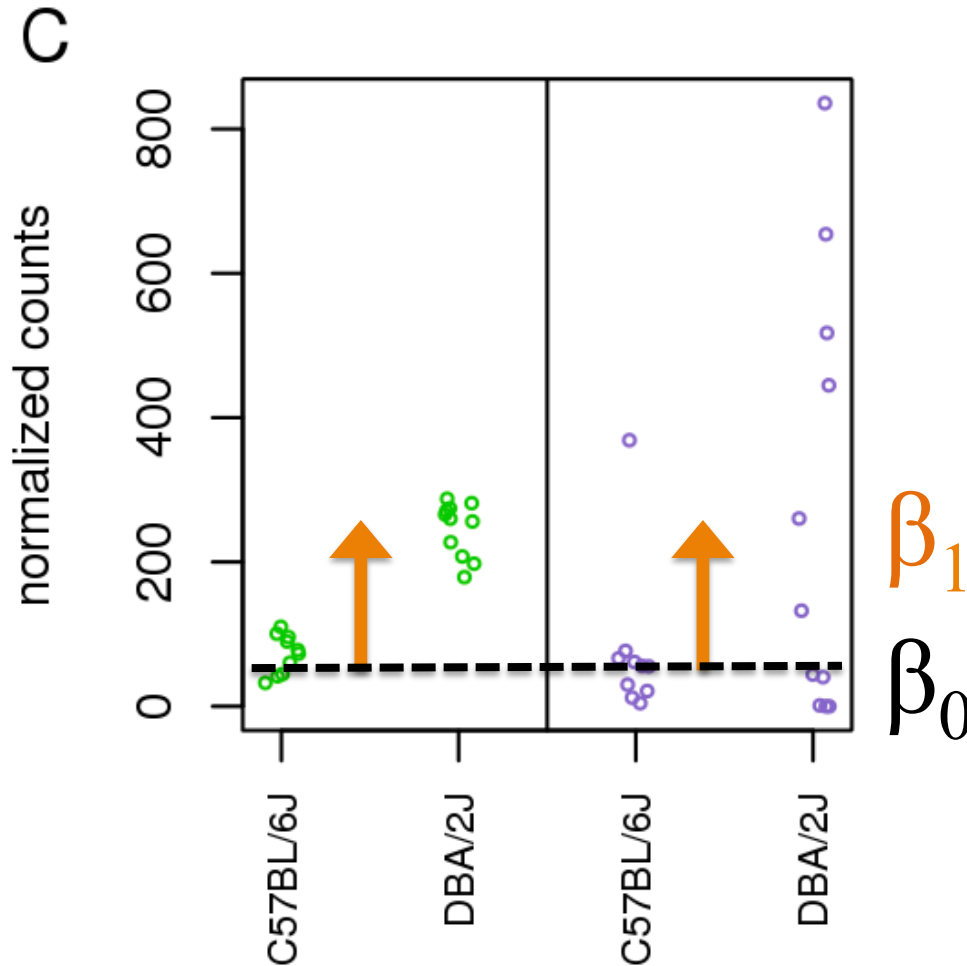
$$\mu_{ij} = s_{ij} q_{ij}$$

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

predictors

$\log q_1$	1	0	
$\log q_2$	1	0	$\beta_0$
$\log q_3$	1	1	$\beta_1$
$\log q_4$	1	1	

# Two genes with equal effect size

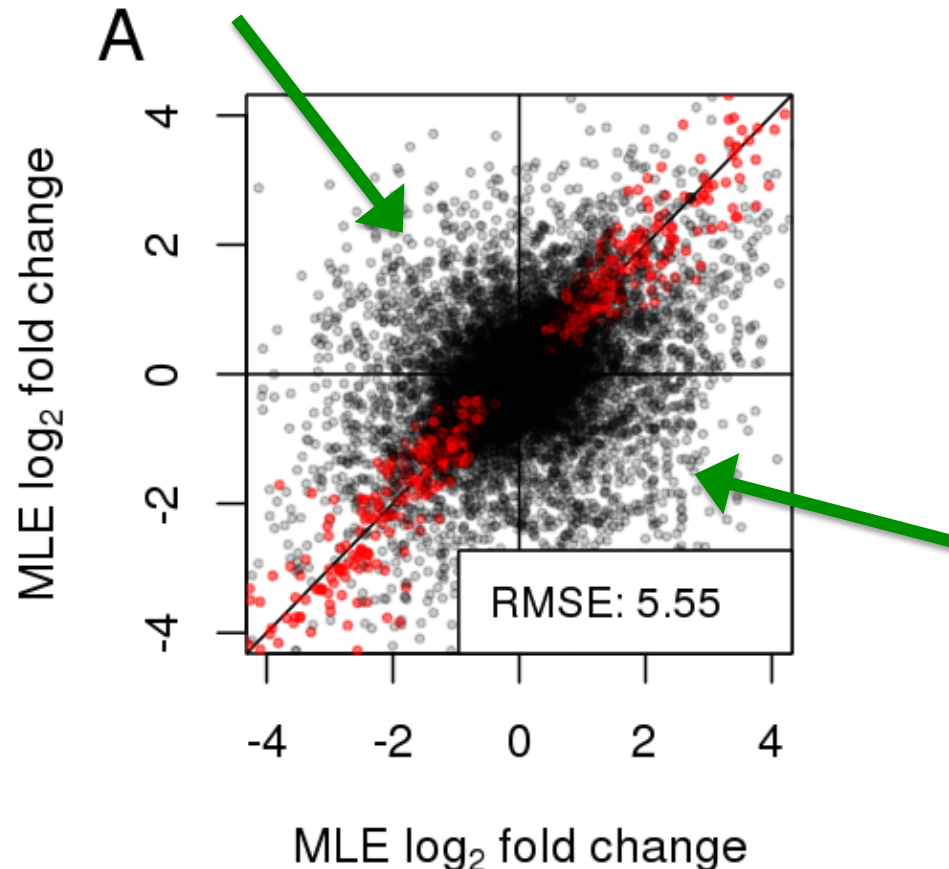


# Maximum likelihood estimates can have high variance

Split 10 vs 11  
samples:

- 5 vs 5
- 5 vs 6

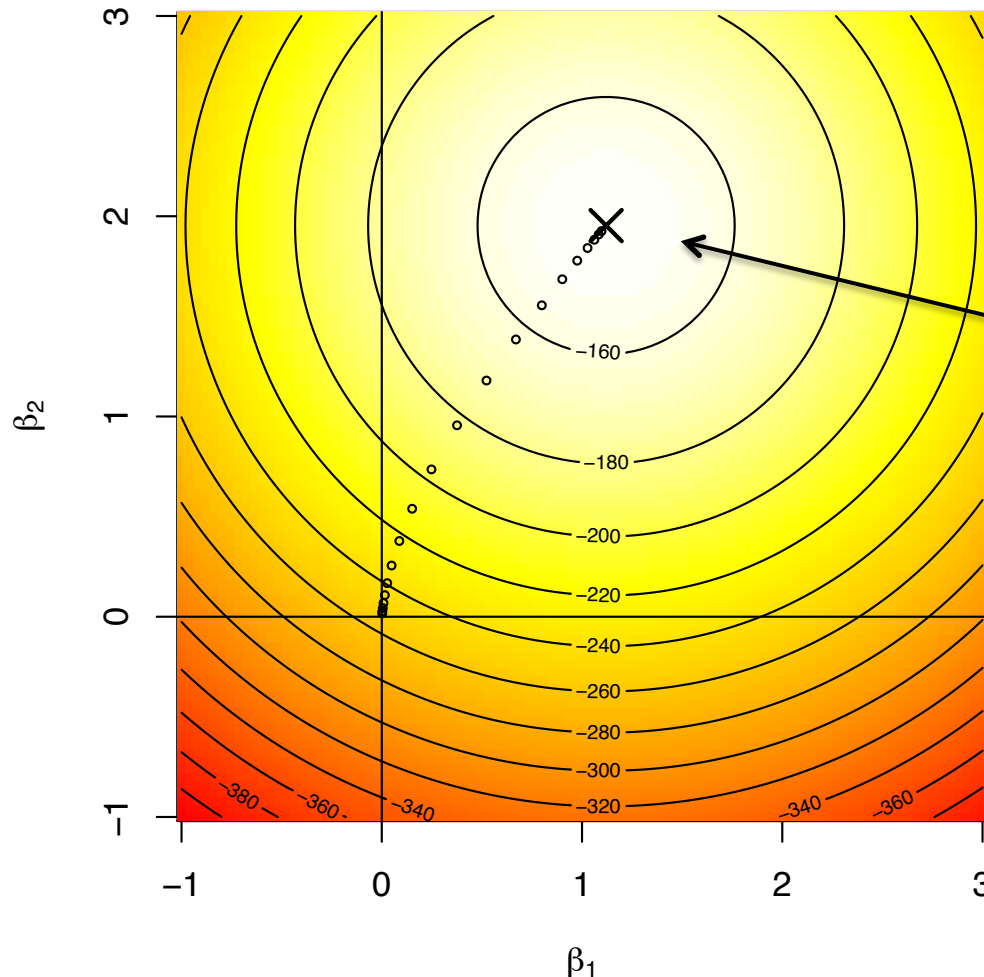
For every gene,  
we get a  $\beta$  from  
both comparisons.



...especially for genes with low counts,  
but also for genes with high variance

# Moderation stabilizes effect sizes

Moderation through a prior:  $\beta_{ir} \sim N(0, \sigma_r^2)$  ...for  $r \neq 0$

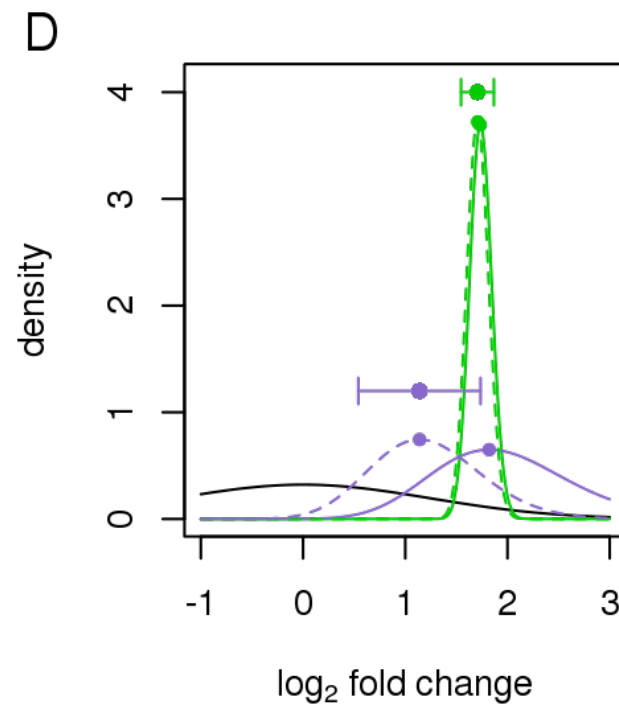
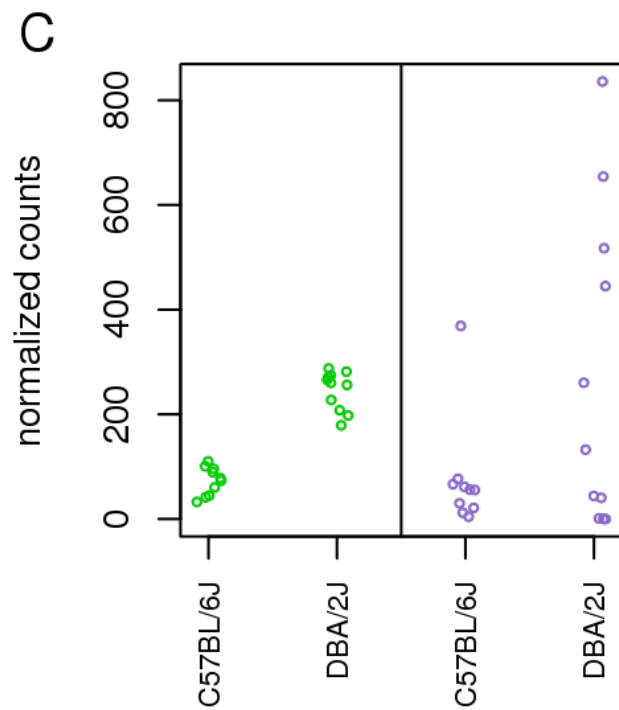
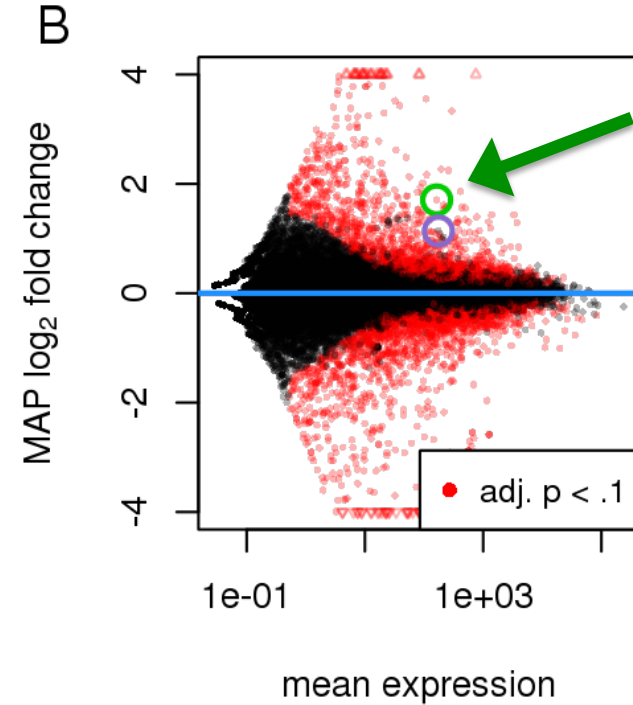
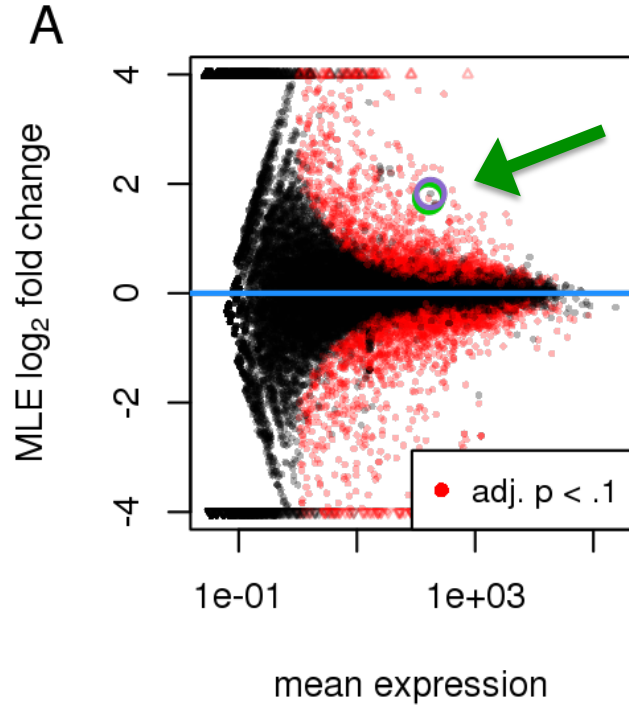


MLE

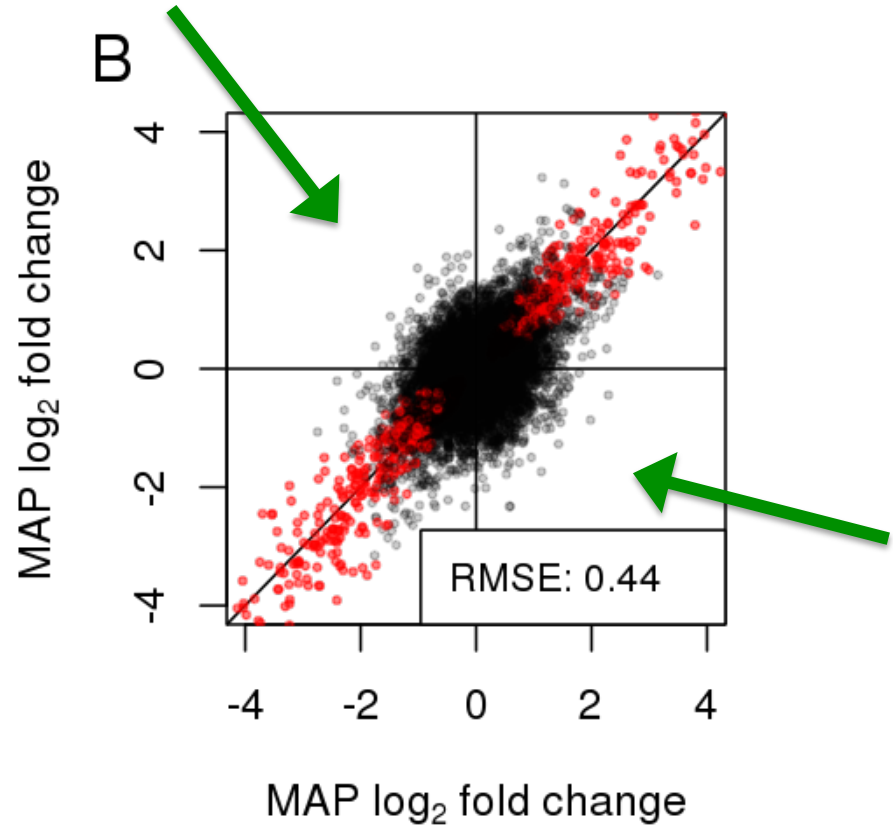
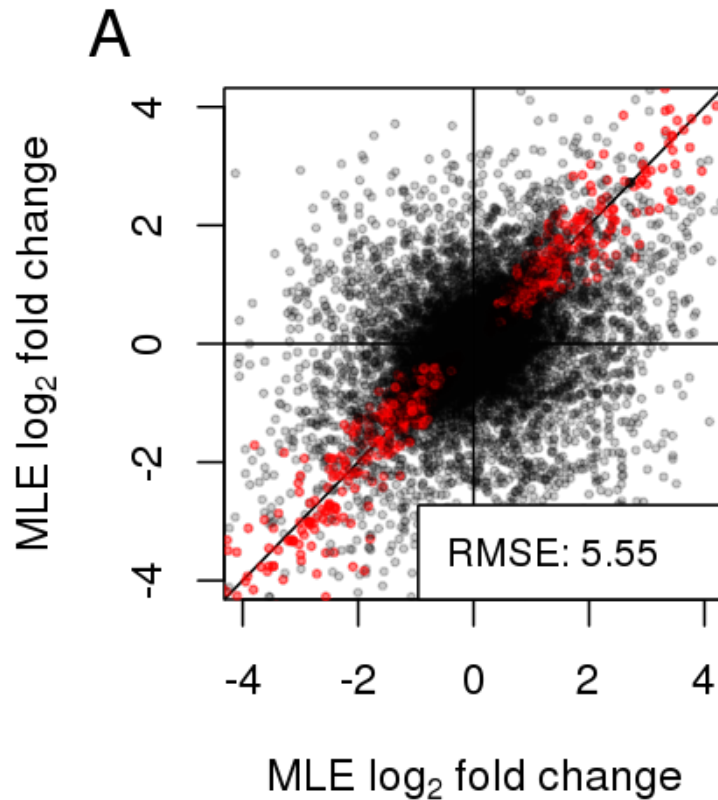
estimated  
from MLE  
effect sizes

the max posterior is  
pulled in towards 0

(here with  $\sigma_1 = \sigma_2$ )







...introduces bias but reduces mean squared error.  
for an estimator:  $MSE = \text{bias}^2 + \text{variance}$

...but moderation with  $\geq 3$  groups  
is not symmetric

		$\triangle$	$\triangle$
A	1	0	0
	1	0	0
B	1	1	0
	1	1	0
C	1	0	1
	1	0	1

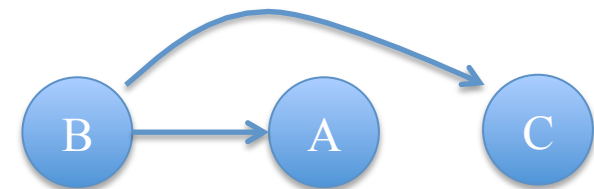
$\beta_0$   
 $\beta_{B-A}$   
 $\beta_{C-A}$



$\neq$

		$\triangle$	$\triangle$
A	1	1	0
	1	1	0
B	1	0	0
	1	0	0
C	1	0	1
	1	0	1

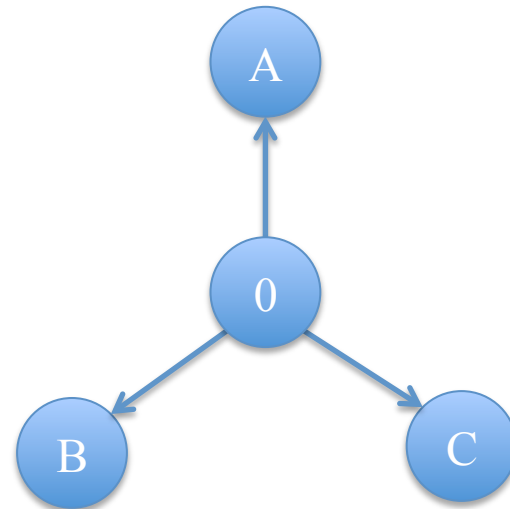
$\beta_0$   
 $\beta_{A-B}$   
 $\beta_{C-B}$



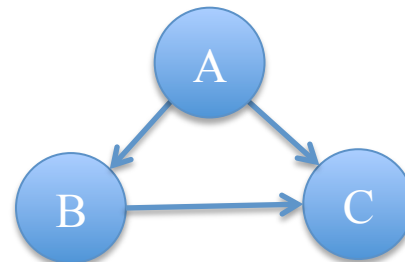
# A solution

Expand by adding a column,  
shrinking all groups towards an intercept.

$$\begin{array}{l} \text{A} \left\{ \begin{array}{cccc} \triangle & \triangle & \triangle & \triangle \\ 1 & 1 & 0 & 0 \\ 1 & 1 & 0 & 0 \end{array} \right. \quad \begin{array}{l} \beta_0 \\ \beta_A \end{array} \\ \text{B} \left\{ \begin{array}{cccc} 1 & 0 & 1 & 0 \\ 1 & 0 & 1 & 0 \end{array} \right. \quad \begin{array}{l} \beta_B \end{array} \\ \text{C} \left\{ \begin{array}{cccc} 1 & 0 & 0 & 1 \\ 1 & 0 & 0 & 1 \end{array} \right. \quad \begin{array}{l} \beta_C \end{array} \end{array}$$



Estimate prior using all MLE contrasts:



# Statistical inference

Test a Wald statistic for each coefficient:

$$W = \frac{\hat{\beta}_{ir}}{\widehat{\text{SE}}(\hat{\beta}_{ir})}$$

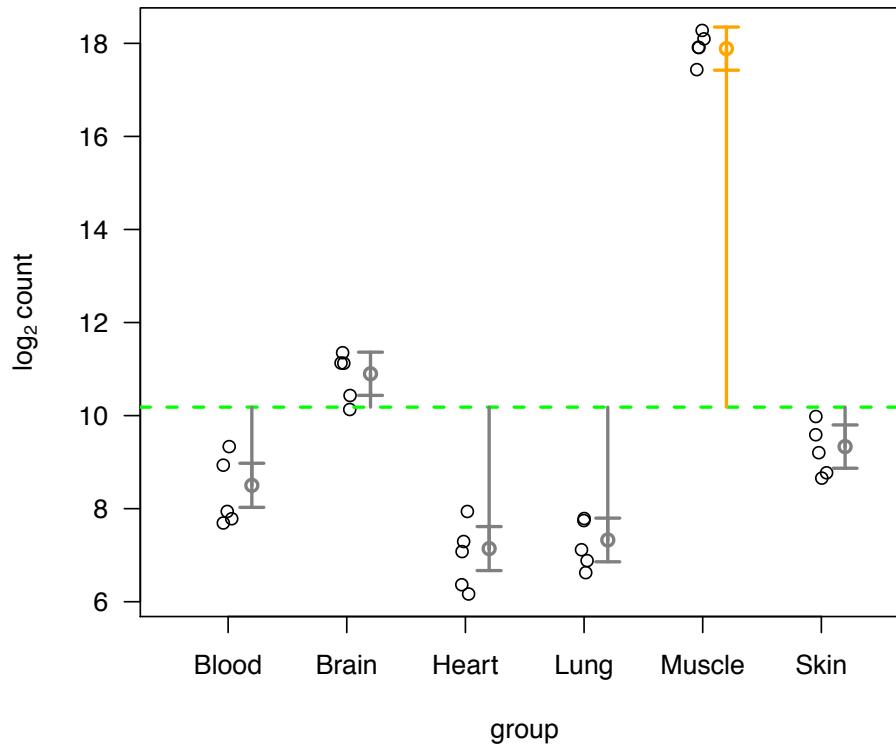
Contrasts can be used to compare multiple levels of a factor,  
e.g.,  $\mathbf{c}^t = [0, 0, -1, 1]$

$$\beta_i^c = \vec{c}^t \vec{\beta}_i$$

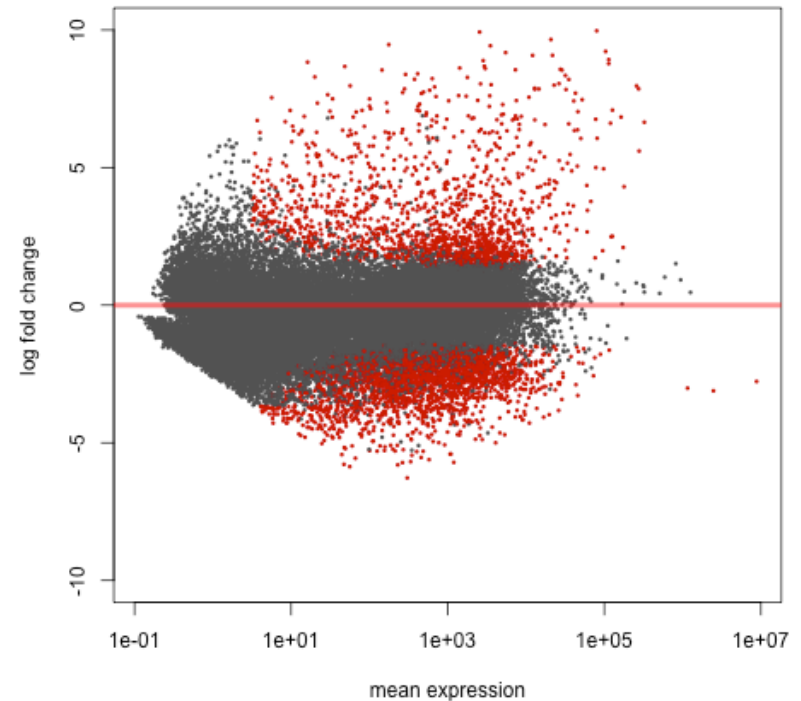
$$\text{SE}(\beta_i^c) = \sqrt{\vec{c}^t \Sigma_i \vec{c}}$$

# GTEx: tissue specific

one gene, highest Wald statistic



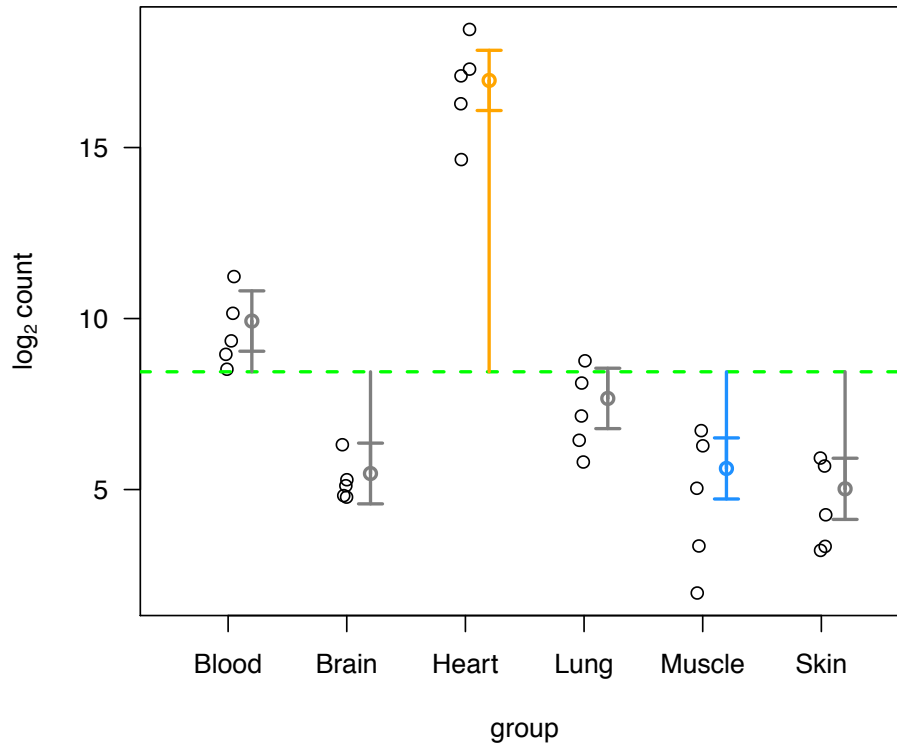
MA-plot all genes



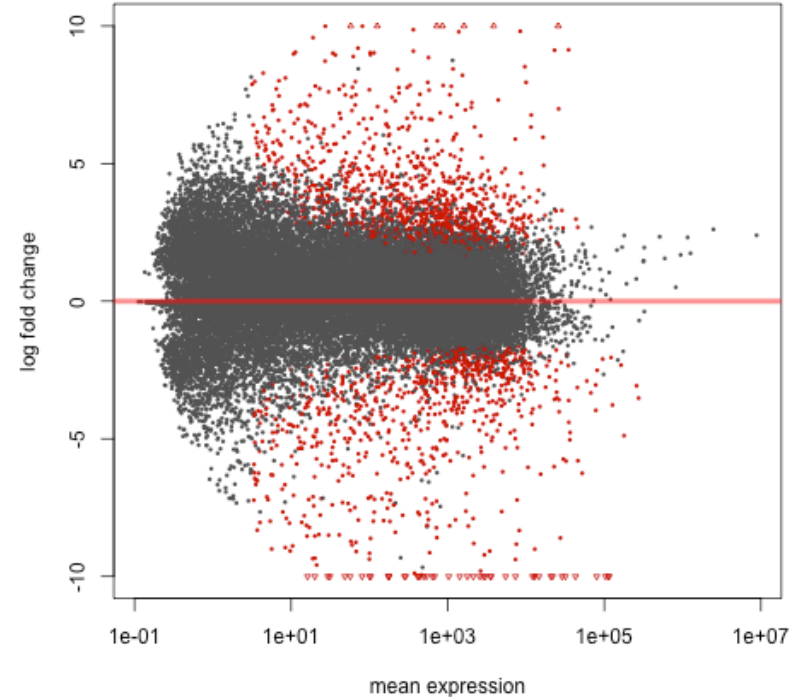
red points:  
more than doubling  
adjusted  $p$ -value  $< 0.1$

# GTEX: contrast two tissues

one gene, highest Wald statistic



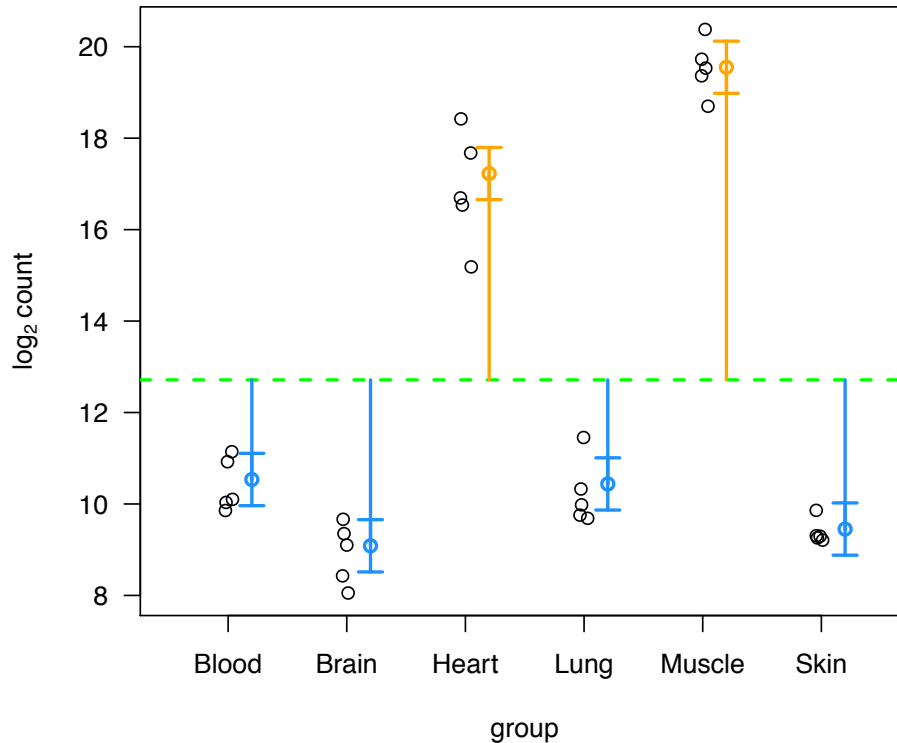
MA-plot all genes



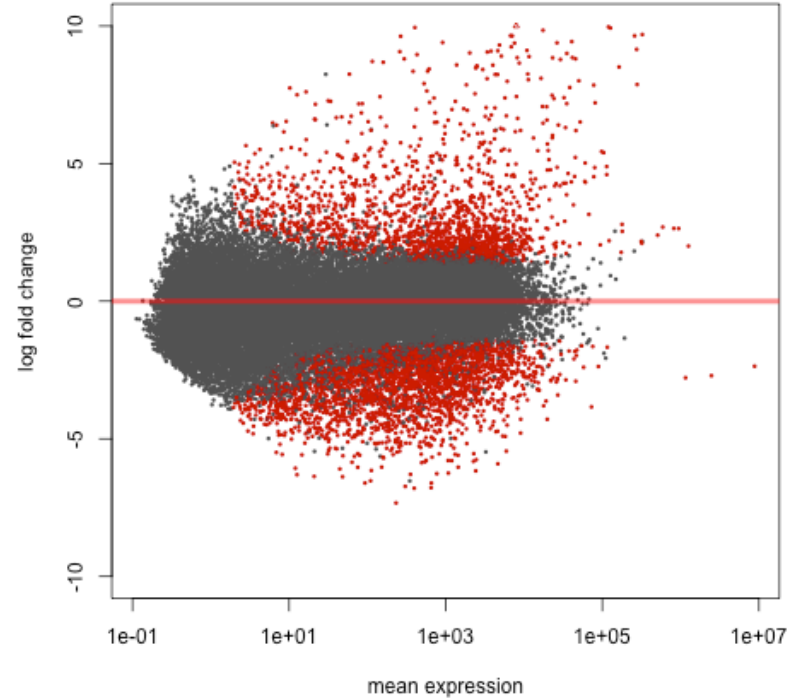
red points:  
more than doubling  
adjusted  $p$ -value  $< 0.1$

# GTEx: contrast 2 vs 4 tissues

one gene, highest Wald statistic



MA-plot all genes



red points:  
more than doubling  
adjusted  $p$ -value < 0.1

# Acknowledgments

- Simon Anders, EMBL
- Wolfgang Huber, EMBL
- Rafael Irizarry, DFCI/HSPH
- Martin Vingron & Knut Reinert, IMPRS, MPIMG, FU

## Related work:

- For RNA-Seq: GFOLD, BitSeq, ShrinkBayes, NPEBSeq
- Genereally: selection bias, "Winner's curse", adaptive shrinkage, Brad Efron, Noah Simon, Matthew Stephens et al.
- Background on bias-variance trade-off: "Elements of Stat. Learning"

## More:

- Effect size shrinkage in DESeq2



- Preprint available on **bioRxiv**  
(search "biorxiv DESeq2")
- Code for GTEx example:  
`github.com/mikelove/multigroup`



# GLM with ridge regularization

$$\vec{\beta}_i = \arg \max_{\vec{\beta}} \left( \overset{\text{log likelihood}}{\sum_j \log f_{\text{NB}} \left( K_{ij}; \mu_j(\vec{\beta}), \alpha_i \right)} + \overset{\text{log prior}}{\Lambda(\vec{\beta})} \right)$$

maximum *a posteriori*

$$\beta_{ir} \sim N(0, \sigma_r^2) \longrightarrow \Lambda(\vec{\beta}) = \sum_r \frac{-\beta_r^2}{2\sigma_r^2}$$

the "penalty term" for large  $\beta$

# simulated multifactor

