Intermediate RNA-seq data analysis using R

Michela Traglia, Min-Gyoung Shin Bioinformatics Core, GIDB February 1st, 2024

GLADSTONE INSTITUTES

Requirements for the demo

Please install the following library in Rstudio

install.packages("magrittr")

install.packages("statmod")

install.packages("tidyverse")

install.packages("ggplot2")

install.packages("BiocManager")

BiocManager::install("edgeR")

BiocManager::install("org.Mm.eg.db")

Verify the installation:

library(magrittr)

library(statmod)

library(tidyverse)

library(ggplot2)

library(edgeR)

library(library(edgeR))

Materials for this workshop

Please download the compressed file 2023Sept_intermediatRNAseq.zip

Double click on the file, the unzipped folder includes:

- This presentation with concepts
- Hands-on session files:
 - handson.R
 - targets.txt
 - DE resutls.txt
 - GSE60450_Lactation-GenewiseCounts.txt.gz

Introductions

Min-Gyoung Shin

Bioinformatician II

Michela Traglia

Senior Statistician

Workshop outline

- Intro to a real experiment Approach for Differentially Expressed Gene analysis: edgeR
- Filtering genes
- Normalization
 - Demo I
- Exploratory visualization: MDS PCA
- Fit the model for DEG
- Compare groups and visualize the DEG
 - Demo II

Assumed background

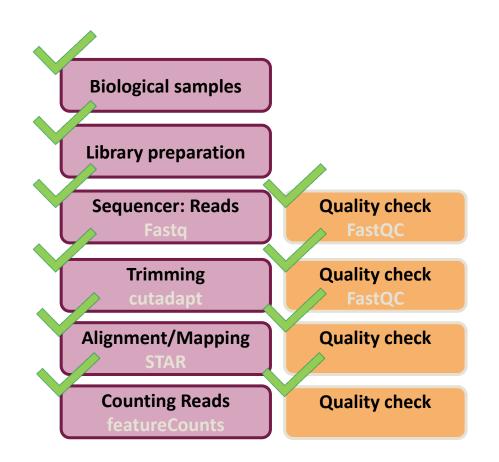
- o Familiarity with R and RStudio
- Familiarity with RNA-seq protocol
- o Familiarity with basic concepts of statistics and hypothesis testing

Poll 1

How familiar are you with the RNAseq pipeline?

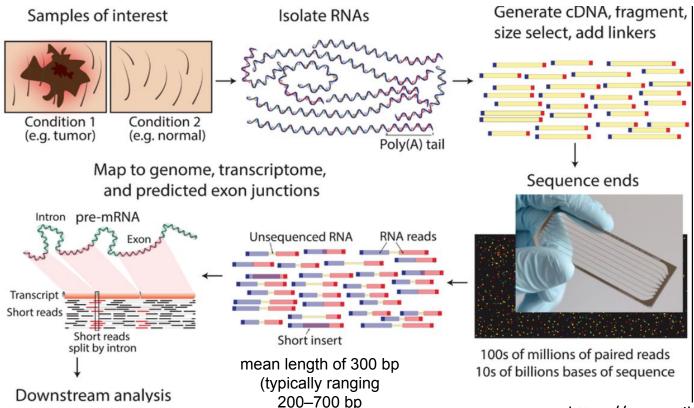
- 1) I attended the Introduction to RNAseq workshop last week but I never run the pipeline
- I attended other courses/ read papers and pipeline documentation only
- 1) I have experience running at least one RNAseq pipeline

RNA-seq - analysis workflow



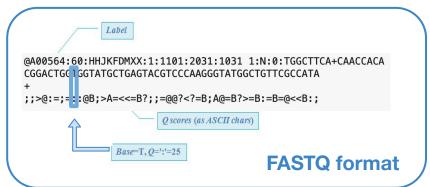


Typical RNA-seq protocol

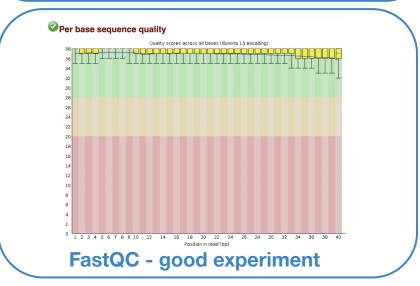


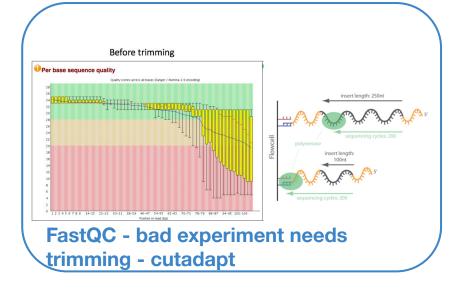
https://www.wikiwand.com/en/RNA-Seq

Bioinformatic pipeline - summary

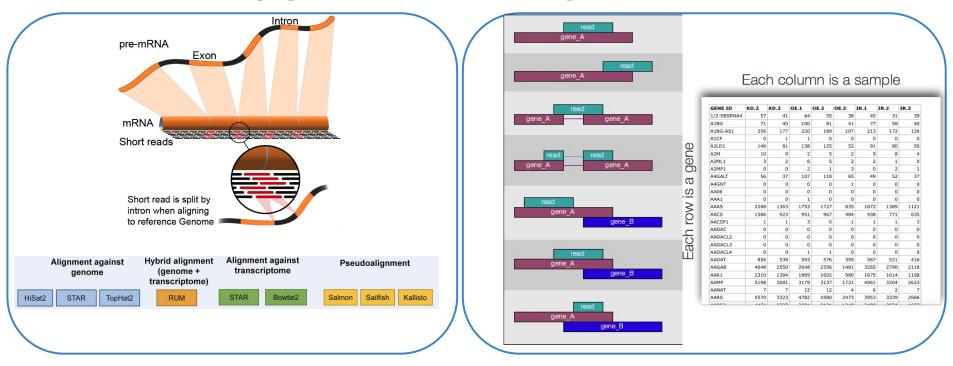








Bioinformatic pipeline - summary



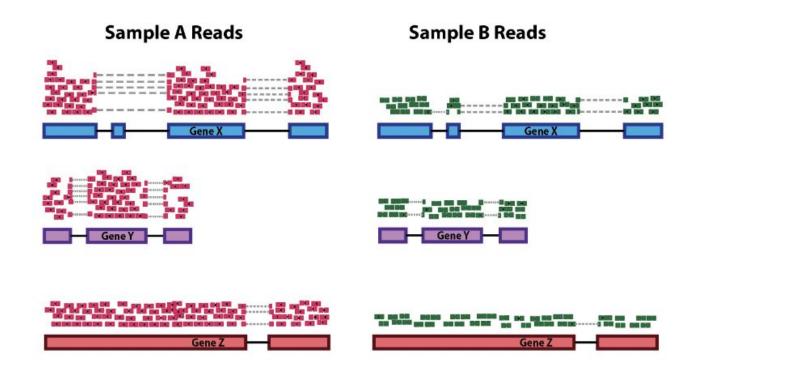
Alignment and pseudo-alignment tools

Feature counts

Please find more details on our wiki page

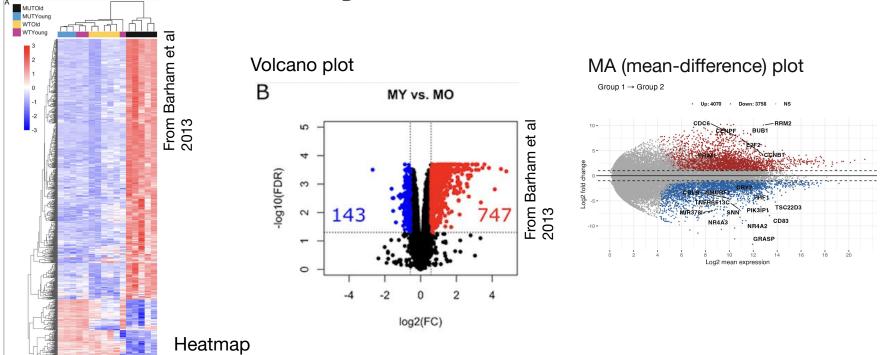
Introduction-to-RNA-Seq-Analysis

Many steps to calculate the Differentially Expressed Genes (DEG) between sample A and B



We need to make the counts comparable across samples

Goal of DEG analysis



Identify genes (and molecular pathways) that are differentially expressed (DE) between two or more biological conditions

Reference for the workshop



F1000Research 2016, 5:1438 Last updated: 06 DEC 2018



SOFTWARE TOOL ARTICLE

From reads to genes to pathways: differential expression analysis of RNA-Seq experiments using Rsubread and the edgeR quasi-likelihood pipeline [version 2; referees: 5 approved]

Yunshun Chen^{1,2}, Aaron T. L. Lun ¹⁰³, Gordon K. Smyth ^{101,4}

¹The Walter and Eliza Hall Institute of Medical Research, Parkville, Victoria, 3052, Australia

²Department of Medical Biology, The University of Melbourne, Victoria, 3010, Australia

³Cancer Research UK Cambridge Institute, University of Cambridge, Li Ka Shing Centre, Cambridge, UK

⁴Department of Mathematics and Statistics, The University of Melbourne, Victoria, 3010, Australia

Dataset

Transcriptome analysis of luminal and basal cell subpopulations in the lactating versus pregnant mammary gland

- o GEO (gene expression omnibus) accession: GSE60450
- o Tissue of origin: Mammary glands of mouse
- Cell types: Basal stem-cell enriched cells (B) and committed luminal cells (L)
- Biological conditions: Virgin, Lactating (2 day) and Pregnant (18.5 day)
- \circ # of groups: 2 cell types (B/L) x 3 conditions (V/L/P) = 6 groups
- o # of replicates: 2 of each group
- Illumina Hiseq sequencer about 30 million 100bp single-end reads for each sample.

 https://www.nct

https://www.ncbi.nlm.nih.gov/geo

Files for the hands-on session

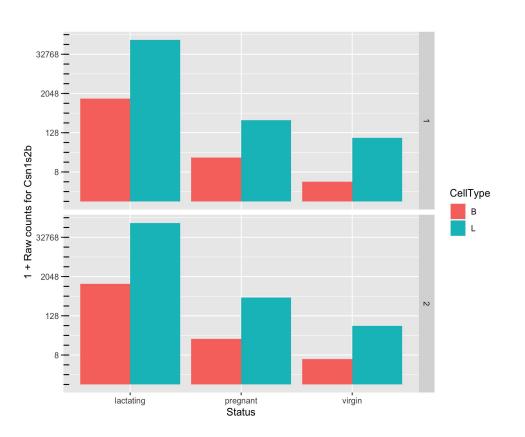
	Status	CellType	SRA	GEO	
	virgin	В	SRR1552450	GSM1480297	MCL1.DG
targets.txt	virgin	В	SRR1552451	GSM1480298	MCL1.DH
3	pregnant	В	SRR1552452	GSM1480299	MCL1.DI
Phenofile	pregnant	В	SRR1552453	GSM1480300	MCL1.DJ
1 Herieme	lactating	В	SRR1552454	GSM1480301	MCL1.DK
	lactating	В	SRR1552455	GSM1480302	MCL1.DL
	virgin	L	SRR1552444	GSM1480291	MCL1.LA

GSE60450_Lactation-GenewiseCounts.txt.gz

Counts for each sample for each gene (Entrez Gene Identifiers)

	Length	MCL1.DG	MCL1.DH	MCL1.DI	MCL1.DJ	MCL1.DK	MCL1.DL	MCL1.LA	MCL1.LB
497097	3634	438	300	65	237	354	287	0	0
100503874	3259	1	0	1	1	0	4	0	0
100038431	1634	0	0	0	0	0	0	0	0
19888	9747	1	1	0	0	0	0	10	3
20671	3130	106	182	82	105	43	82	16	25
27395	4203	309	234	337	300	290	270	560	464

Potential biological questions



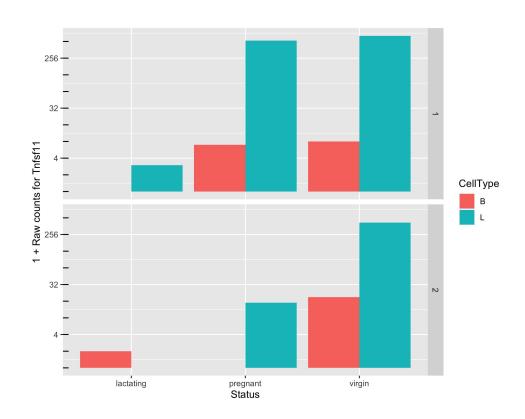
Which comparisons are we interested in?

Example:

- 1. B vs L,
- 2. B.lactating vs L.pregnant,
- 3. ...
- 4. All of them

Or not interested in comparisons but in gene expression of one sample

Potential issues

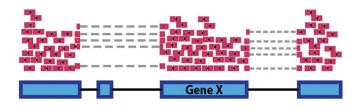


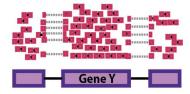
Can we make reliable inferences for genes with very low counts? What should we consider "very low"?

Possible analyses

Within a sample

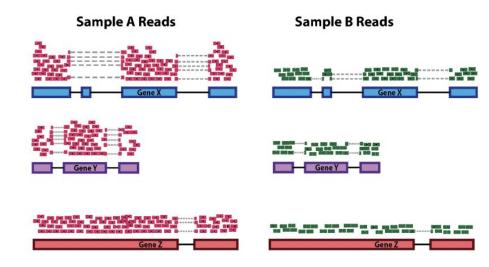
Sample A Reads





Estimate and compare gene expressions across genes (features)

Between samples



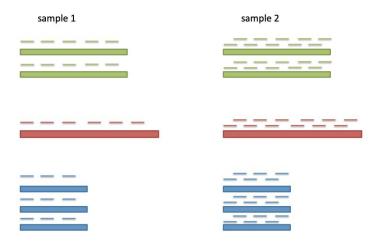
DEG: Compare gene expressions across samples

Bias: Gene-length and sequence depth

At the same expression level, a long gene will have more reads than a shorter gene



Higher sequencing depth, higher counts



Different sequence depths cause variation in counts

Variation in sequencing depths => Need to normalize counts

Group	Total counts
B.virgin	23085177
B.virgin	21628857
B.pregnant	23919152
B.pregnant	22490570
B.lactating	21382233
B.lactating	19884434

Group	Total counts
L.virgin	20213223
L.virgin	21509988
L.pregnant	22073815
L.pregnant	21837341
L.lactating	24638939
L.lactating	24581591

Library size about 20M

Source of technical variability

Identify and correct technical biases removing the least possible biological signal.

- Gene length
- Library size or sequence depth (number of mapped reads)
- RNA sample composition
- Batch effects
- 0 ...

Poll 2

Within same sample gene expression quantification might be biased by:

- 1. Gene length
- 2. Sequencing depth
- 3. Both

Poll 3

Between samples differentially gene expression (DEG) might be biased by:

- 1. Gene length
- 2. Sequencing depth
- 3. Both

Several normalization methods to remove technical bias

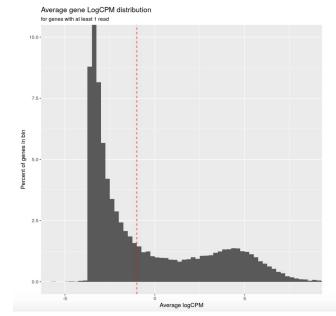
- Normalized expression are necessary to remove technical biases
 - Depth of sequencing correction
 - Depth of sequencing and gene length correction
 - RNA sample composition
- To remove between samples batch effects
- Various normalized gene expression units such as RPM (or CPM), RPKM, FPKM, TPM, TMM (edgeR), DESeq
- Measure of the abundance of gene or transcripts.

Sequence depth/Library size - Count Per Million mapped reads

$$RPM \ or \ CPM = \frac{Number \ of \ reads \ mapped \ to \ gene \times 10^6}{Total \ number \ of \ mapped \ reads}$$

Sequenced one library with 5 million(M) reads. Total 4 M matched to the genome sequence 5000 reads matched to a given gene $\frac{5000 \times 10^6}{4 \times 10^9} = 1250$

Filter on count-per-million (CPM) values to avoid favoring genes that are expressed in larger libraries over those expressed in smaller libraries



A gene at least 10–15 counts in at least some libraries before it is considered to be expressed -> Identifying the CPM that corresponds to 10-15 counts

For comparison within replicates or same group, NOT for within sample, NOT for DEG

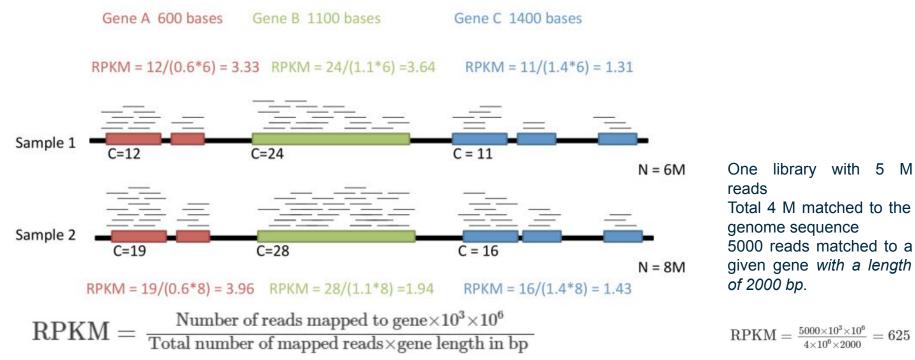
Common methods to normalize the counts considering gene length

Commonly used normalization method that includes <u>sequence</u> <u>depth and gene length</u> correction:

- RPKM /FPKM (Reads/Fragments Per Kilobase per Million)
- TPM (Transcripts Per kilobase Million)

RPKM: seq depth and gene length bias

Reads/Fragments Per Kilobase per Million



RPKM and FPKM -> the sum of the normalized reads in each sample may be different, and this makes it harder to compare samples directly.

TPM: seq depth and gene length bias

TPM: Normalize for gene length first to get the Reads Per Kilobase, sum up all the RPK in a sample (across all genes), and then normalize for sequencing depth

$$TPM = 10^6 * \frac{\text{reads mapped to transcript/transcript length}}{\text{Sum(reads mapped to transcript/transcript length)}}$$

Represent the relative abundance of a transcript among a population of sequenced transcripts

The sum of all TPMs in all samples are the same -> to compare the proportion of reads that mapped to a gene in each sample.

TPM is equal to 10⁶ (1 million) divided by the number of annotated transcripts in a given annotation so it is constant and proportional to the relative RNA molar concentration

Don't use these methods for between condition/groups comparison!

When total RNA composition is similar across samples, these methods could be potentially used.

But you can't assume it!

Cells don't necessarily produce similar levels of RNA/cell between cell types, disease states or developmental stages is not always valid



Methods for between samples comparisons / DEG

Observed counts depend on total reads sequenced and RNA sample composition

Sample-to-sample variability in total RNA concentration

- Need to normalize for <u>difference in total reads between samples</u>.
 - o Might be enough if total nucleic acid is the same in both samples.
 - Example: technical replicates

Need to account for difference in RNA sample composition

RNA composition -> scaling factors

Normalize for RNA composition differences by scaling factor

A few highly differentially expressed genes may have a strong influence on read counts -> minimizing effect of such genes

Assumption: A majority of transcripts is not differentially expressed

Adjust counts such that for most genes, counts are not differential.

Approach to identify DE genes: edgeR

- Bioconductor package edgeR utilizes a theoretical model that captures some of the known processes leading to noise in counts data. (null model)
- Assume that the data is generated according to this model.
- Given any observed level of difference in mean expression levels of a gene, compute the probability that the observation will result from the null model (p value)
- If the probability is very low (e.g., p < 0.05), infer that something may be happening that we did not account for in the null model. (e.g., biological processes in L cells for milk production)

Need to correct for multiple testing

P value represents the chance that we may be wrong in calling something significantly differential. Example:

- P = 0.01 means 1% chance that we may be wrong.
- P = 0.50 means 50% chance that we may be wrong.

More than 20k genes under consideration

- => if a certain difference in expression levels has only 1% chance of happening given the null model, it might be observed for 200 genes even if the null model were true for all the genes.
- => 200+ false positives

Hence, there is a need to adjust the p-values.

- The more genes we test, the more we must adjust.
- Reduce the number of tests by filtering out "uninteresting" genes.

"Uninteresting" genes Filtering

- Biological point of view: minimal expression level of a gene -> translation into a protein -> biologically relevant
- Statistical point of view: low counts -> not enough statistical evidence.

Genes with consistently low counts are very unlikely be assessed as significantly DE

	SRR1039508	SDD1030500	SRR1039512	SDD1030513	SDD1030516	count outlie		
ENSG00000000003	67	44	87	40	1138	7		
ENSG00000000005	0	0	0	0	0			
ENSG00000000419	467	515	621	365	587			
ENSG00000000457	260	211	263	164	245	Genes with		
ENSG00000000460	2	5	1	0	1	zero counts		
		Genes with low mean normalized counts						

Between sample normalization

Normalize for RNA composition by a set of scaling factors that minimize the log-fold changes between the samples for most genes

- o Reference sample: have the closest average expressions to the mean of all samples
- oTest samples: other samples

<u>Scaling factor</u>: weighted mean of log ratios between the test and reference, from a gene set removing most/lowest expressed genes (avg read counts) and genes with highest/lowest log ratios (differences in expression)

Normalization with edgeR: Trimmed Mean of M-values

- 1. Choose a reference sample.
- 2. Compute the M and A values for all genes.

M=log FC between ref and test; A=avg count gene between ref and test

- 3. Filter genes that fall in the tails of M and A distributions.
- 4. Estimate variance of M values.
- 5. Estimate TMM --- the weighted average of trimmed M-values.
- 6. Size factor is 2^{TMM} .
- 7. Adjust such that these multiply to 1.

23218026	1.2368993
21768136	1.2139485
24091588	1.1255640
22656713	1.0698261
21522033	1.0359212
20008326	1.0872153
20384562	1.3684449
21698793	1.3653200
22235847	1.0047431

21982745 24719697

24652963

0.9232822

0.5291015

0.5354877

lib.size norm.factors

TMM is implemented in edgeR and performs better for between-samples comparisons, when comparing the samples from different tissues or genotypes.

Other approaches to normalization

- Relative Log Expression (RLE) approach by Anders and Huber (2010)
 - o Reference: geometric mean of all samples
 - Normalization factor: median ratio of each sample to the reference
 - o RLE and TMM give similar results with real and simulated data
 - R package <u>DESeq</u>
- OUpper quartile normalization by Bullard et al (2010)
 - Normalization factor: 75% quantile of the counts for each sample
 - Not recommended in general
 - Control genes (housekeeping genes, spike-in) to estimate technical noise
 (RUVSeq 2014 Remove Unwanted Variation)

Best practice to choose a normalization

An effective normalization should result in a <u>stabilization</u> of read counts across samples (eliminate composition biases between libraries)

- oTC, RPKM, UQ Adjustment of distributions, implies a similarity between RNA repertoires expressed
- oDESeq, TMM More robust ratio of counts using several samples, suppose that the majority of the genes are not DE
- RUVSeq Powerful when a large set of control genes can be identified

Main normalization approaches summary

Normalization method	Description	Accounted factors	Recommendations for use
CPM (counts per million)	counts scaled by total number of reads	sequencing depth	gene count comparisons between replicates of the same samplegroup; NOT for within sample comparisons or DE analysis
TPM (transcripts per kilobase million)	counts per length of transcript (kb) per million reads mapped	sequencing depth and gene length	gene count comparisons within a sample or between samples of the same sample group; NOT for DE analysis
RPKM/FPKM (reads/fragments per kilobase of exon per million reads/fragments mapped)	similar to TPM	sequencing depth and gene length	gene count comparisons between genes within a sample; NOT for between sample comparisons or DE analysis
DESeq2's median of ratios [1]	counts divided by sample-specific size factors determined by median ratio of gene counts relative to geometric mean per gene	sequencing depth and RNA composition	gene count comparisons between samples and for DE analysis; NOT for within sample comparisons
EdgeR's trimmed mean of M values (TMM) [2]	uses a weighted trimmed mean of the log expression ratios between samples	sequencing depth, RNA composition	gene count comparisons between samples and for DE analysis; NOT for within sample comparisons



Break (10 min) Demo I

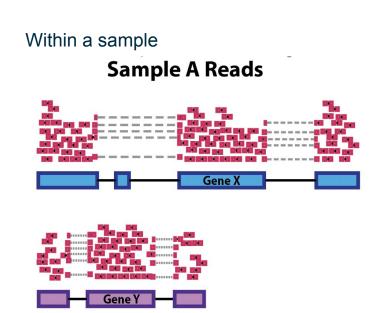
Poll 4

Which normalization is suitable for within-sample gene

expression quantification:

1. TPM

1. TMM (edgeR)



Poll 5

Which normalization is suitable for comparing gene expression between disease samples vs control samples

- 1. RPKM
- 2. TMM (edgeR)





Demo I

Hands-on session

- Load and reformat the data
- oExploratory visualization : MA plot
- Create DGElist object and retrieve gene symbols
- oFilter genes with inadequate information
- Exploratory visualization : MDS and PCA plots
- Define and fit a model
- oHypothesis testing (four example hypotheses)
- Save results as a table and explore in Excel

Dataset

Transcriptome analysis of luminal and basal cell subpopulations in the lactating versus pregnant mammary gland

- o GEO (gene expression omnibus) accession: GSE60450
- o Tissue of origin: Mammary glands of mouse
- Cell types: Basal stem-cell enriched cells (B) and committed luminal cells (L)
- Biological conditions: Virgin, Lactating (2 day) and Pregnant (18.5 day)
- \circ # of groups: 2 cell types (B/L) x 3 conditions (V/L/P) = 6 groups
- o # of replicates: 2 of each group
- Illumina Hiseq sequencer about 30 million 100bp single-end reads for each sample.

 https://www.nct

https://www.ncbi.nlm.nih.gov/geo



MDS and PCA plots

Assessing overall similarity across samples

- O Which samples are similar to each other, which are different?
- Does this fit to the expectation from the experiment's design?
- O What are the major sources of variation in the dataset?

Expression level of Casein varies in a way that is strongly indicative of the effect of CellType and Status.

Why are the B.lactating samples not close to B.virgin and B.pregnant samples?

Could it be due to batch effects?

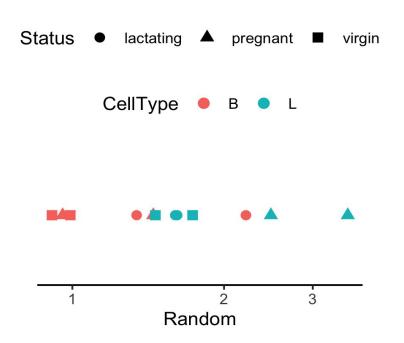
The distance between each pair of samples can be interpreted as the leading log-fold change between the samples for the genes that best distinguish that pair of samples.

Casein

10000

Expression appears to vary across samples but...

In general, the way expression appears to vary across samples could be dominated by noise, batch effects, real signal, etc.



Identify the source fof technical variability!



Fitting the model

How to model the normalized RNA-seq red counts

Total number of reads for a sample ~ millions

Counts per gene ~ tens /hundreds /thousands.

The chance of a given read to be mapped to any specific gene is rather small.

Discrete events sampled out of a large pool with low probability are usually modeled with Poisson distribution

Overdispersion of read counts between samples

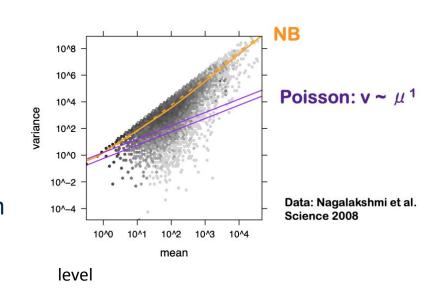
The number of reads that are mapped into a gene was first modeled using a Poisson distribution

Assumption: assumes that mean and variance are the same

The variance grows faster than the mean in RNAseq data.

Overdispersion in RNA-seq data

- -> counts from <u>biological replicates</u> vary so tend to have variance exceeding the mean (highly expressed genes)
- -> underestimation of the biological variance increased the probability to falsely declare a gene DE when it is non-DE (type I error rate)

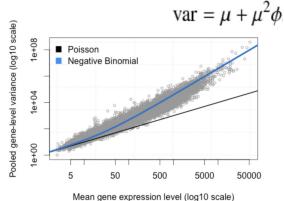


Three dispersion estimates

The negative binomial (NB) distribution proposed as an alternative to model the read counts for each gene <u>in each sample</u>

The variance is always larger than the mean for the negative binomial ⇒ suitable for RNA-seg data

Many genes, <u>few biological samples</u> - difficult to estimate ϕ on a gene-by-gene basis Using information *across all genes* for stable estimates of ϕ .



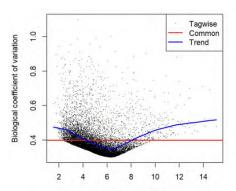
Empirical Bayes estimates of dispersion parameters: Learning from the experience of others

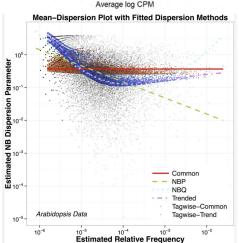
Dispersion accounts for variability between biological replicates

- Common dispersion: a global dispersion estimate averaged across the genes – not enough
- Trended dispersion: dispersion of a gene is predicted from its abundance – similar abundant genes
- <u>Tagwise dispersion</u>: measure of the degree of consistent inter-library variation for that tag

Empirical Bayes estimates need to be controlled for the possibility of outlier genes with exceptionally large or small individual dispersions (robust=TRUE)

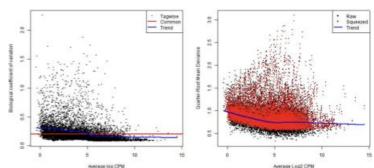
plotBCV (biological coeff variation)





Fitting the model : glmQLFit to manage the variance of gene counts

- NB dispersions higher for genes with very low counts and decrease smoothly with abundance and asymptotically to a constant value for genes with larger counts.
- Extended NB model to account for gene-specific variability from both biological and technical sources (quasi-likelihood)
- 1) NB dispersion trend is used to describe the overall biological variability across all genes (fit GLM)
- For each gene-specific variability above and below the overall level (deviance) is picked up by the QL dispersion



edgeR fitting models

- o Classic (pairwise comparisons between two or more groups), glm and glmQL
- QL for bulk RNA-seq:
 - + stricter error rate control (more rigorous dispersion and uncertainty)
 - + speed improvement compared to other quasi-methods
 - + appropriate for multiple treatment factors and with small # of biological replicates
 - + relative changes in expression levels between conditions (not absolute)

Limma package for large scale datasets – high overlap across methods

Get the DE genes - glmQLFTest

- Identifies differential expression based on statistical significance regardless of how small the difference might be -> 5000 DE genes between condition and control groups
- Interested only in genes with large expression changes -> subset of genes more biologically meaningful.
- Modify the statistical test to evaluate variability as well as the magnitude of change of expression values -> expression changes greater than a specified threshold
- Not equivalent to a simple fold change cutoff: "the fold-change below which we are definitely not interested in the gene"
- The total number of DE genes identified at an FDR of 5% can be shown with decideTestsDGE() – set cutoff

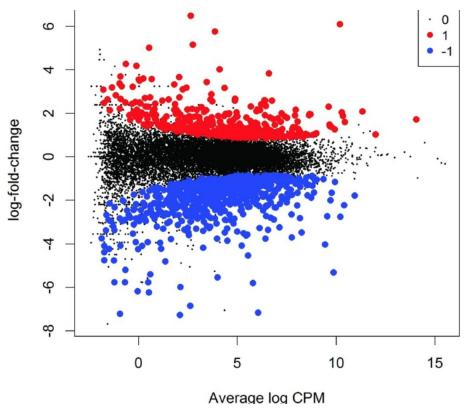
Get the DE genes - glmQLFTest

	genes	logFC	logCPM	F	PV alue	FDR
PCDHA10	PCDHA10	-3.602	5.676	499.9	8.164e-11	1.354e-06
CHGA	CHGA	2.923	5.976	185.4	1.972e-08	0.0001635
ARRB1	ARRB1	-3.914	5.015	158.3	4.627e-08	0.0002019
TSSC2	TSSC2	3.175	3.301	156.8	4.869e-08	0.0002019

Complicated contrasts - makeContrasts()

- between lactating and pregnant mice is the same for basal cells as it is for luminal cells
- the interaction effect between mouse status and cell type

MD plot: Over and under expressed genes



Library size-adjusted log-fold change between two libraries (the difference) vs the average log-expression across those libraries (the mean).

Log-fold change and average abundance of each gene

In summary

- Raw counts are not comparable across samples/genes within a sample
- Several normalization methods: some more suitable for DEG
- Estimate the dispersion, visualize the technical variability
- o Fit the model: counts variance exceeding the mean
- Make all the comparisons that you wish using complex contrast
- Visualize the DEG and get a list for pathway analysis

Your feedback is important to us!

At the end of the hands-on session:

Please take the survey ~3 min: https://www.surveymonkey.com/r/F75J6VZ

Real data might need additional analyses choices that need experience.

Consult with the <u>Gladstone Bioinformatics core</u> for such scenarios and data.

Winter-Spring workshops schedule coming soon - Data Science Training Program



Demo II

Hands-on session

oLoad and reformat the data

- Exploratory visualization : MA plot
- oCreate DGElist object and retrieve gene symbols
- oFilter genes with inadequate information
- Define and fit a model
- Hypothesis testing (four example hypotheses)
- Save results as a table and explore in Excel

Exploratory visualization : MDS and PCA plots

