

UQ-Brisbane SISG

Module 9: Gene Expression and Epigenetic Profiling



Monday February 13, 2017
 "Transcriptomics Study Design"

Greg Gibson

Georgia Institute of Technology
 greg.gibson@biology.gatech.com

Content of the Course

Monday Morning	1a.	Experimental Design (GG)
	1b.	RNASeq (JP)
Monday Afternoon	2a.	Data Normalization (GG)
	2b.	Hypothesis testing (JP)
Tuesday Morning	3a.	Epigenomics (GG)
	3b.	eQTL analysis (JP)
Tuesday Afternoon	4a.	Colocalization and Immunogenomics (GG)
	4b.	Single Cell RNASeq (JP)

Content of the Lecture

Fundamentals of Gene Expression Profiling

Accessing and Depositing data in GEO or ArrayExpress

Fundamentals of Hypothesis Testing

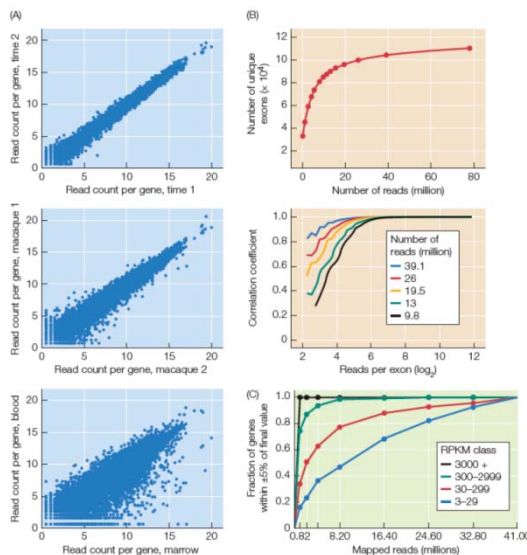
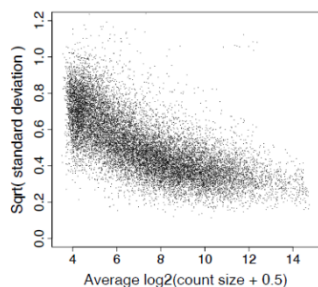
Volcano Plots and False Discovery Rates

Applications of Gene Expression Analysis

1. Atlases of gene expression for functional annotation
2. Identification of differentially expressed genes
3. Assembly of networks of co-regulated genes
4. Investigation of regulatory mechanisms
5. Evolutionary and ecological genomics
6. Clinical genomics
7. Quantitative basis of complex traits

3 Additional Analytical Concerns with RNASeq

1. Low abundance transcripts have high variability (the distribution approximates a negative binomial)
2. Abundance estimates may be a function of read depth
3. Use of Counts per Million (cpm) requires adjustment for over-representation of high abundance transcripts (eg using TMM in EdgeR)



Microarray vs RNASeq

Advantages of Microarrays

- Less expensive
- Better sensitivity for low abundance
- Computationally simpler
- Better-defined statistical properties
- Perfectly good for most applications

Disadvantages of Microarrays

- Only for humans, model organisms
- Different platforms give different results
- Large technical batch effects
- Sensitivity to polymorphism
- Low consistency of analytics among groups

Advantages of RNA-Seq

- Disruptive technology
- Unbiased by prior gene knowledge
- Alternative exon usage
- Allele specific expression (ASE)
- High repeatability

Disadvantages of RNA-Seq

- More opportunity to screw up the analysis
- Oversold resolution of exon level and ASE
- Short read alignment biases
- Sensitivity to polymorphism
- Low consistency of analytics among groups

Design Biases

Molecular biologists tend not to be used to thinking about things like:

- statistical power: how many replicates do I need to see an effect?
 - batch effects: there are lots of ways to obscure what you want to see
 - cost: each sample typically costs \$500 or more
- think about these things before you start the project!!*

At the design step, avoid confounding biological factors:

- don't contrast bloods from young males and old females
 - don't contrast hearts from normal mice and livers from obese ones
- as far as possible, balance all biological factors*

Be aware of the potential for technical confounding:

- date of RNA extraction or hybridization
- batch of arrays
- person who did the hybridization
- scanning software

Levels of Replication

Often you will have a fixed budget that constrains how many arrays can be processed. So your first task is to determine what levels of replication you can afford, and how they will impact statistical power.

Technical Replication:

- RNA preparation (eg. from adjacent biopsies)
- cDNA synthesis and labeling (pooling minimizes outlier effects)
- array hybridization (with commercial arrays, quality generally very high)
- duplicate probes for the same gene

Biological Replication:

- Fixed effects:
- gender
 - treatment (drug, growth regimen, tissue)
 - time of sampling (repeated measures in some cases)
 - genotype (IF specifically chosen and resampled)

- Random effects
- individual from a population
 - field plot

Statistical Power

Y —○— Power at Alpha=0.0001
 —+— Power at Alpha=0.001
 —◇— Power at Alpha=0.01

Power is the concept of how often you expect to detect an effect at a certain significance level, given a number of samples. It is a function of:

- the sample size
- the magnitude of the difference between classes
- the variance within the classes being compared

Since two of these parameters vary for each gene, Power in a microarray experiment is usually assessed in terms of the effect size (amount of variance explained), not as a magnitude of difference.
 But, biologically it is not clear what effect size is important for any given gene.

GEO and ArrayExpress

A GEO record

<https://www.ncbi.nlm.nih.gov/geo/browse/>

Description of the study

Citation

Address when uploaded

List of samples

Downloadable files

A GEO platform

Data table header descriptions

ID	Species	Source	Search_Key	Transcript	ILMN_Gene	Source_Reference_ID	RefSeq_ID	Unigene_ID	Entrez
ILMN_172081	Homo sapiens	RefSeq	ILMN_24919	ILMN_24919	LOC23317	HM_932624.1			23117
ILMN_195180	Homo sapiens	Unigene	ILMN_127219	ILMN_127219	HL555036		HL375026		HL375026
ILMN_1854174	Homo sapiens	RefSeq	ILMN_139282	ILMN_139282	FCOR2B	HM_938851.1			2213
ILMN_1790093	Homo sapiens	RefSeq	ILMN_3008	ILMN_3008	TRN144	NP_017583.3			54763

A GEO sample

Sample GSM426853 Query Datasets for GSM426853

Status Public on Dec 01, 2009
 Title A0942
 Sample type RNA
 Source Name Leukocytes, Agadr, Urban
 Organism Homo sapiens
 Characteristics geographic location: Agadr
 lifestyle: urban
 gender: male
 tissue: peripheral blood
 cell type: leukocytes
 total RNA
 Extracted molecule Total RNA was extracted from leukocyte samples using Ambion's
 Extraction protocol Leukolocks kit. Quality control was performed using Agilent's
 Bioanalyzer.
 Label Biotin
 Label protocol cDNA and cDNA labeling and amplification were all performed using a
 single kit: Ambion's Illumina TotalPrep RNA Amplification kit.
 Hybridization protocol Illumina's BeadArray protocol
 Scan protocol Illumina's BeadArray protocol and BeadStudio (Gene expression Module)
 Description A094
 Data processing Raw data was log2 transformed and median-centered using JMP
 Genomics (SAS, Cary NC).
 Submission date Jul 13, 2009
 Last update date Nov 12, 2009
 Contact name Yousef Idaghdour
 E-mail(s) idaghd@ncsu.edu
 Organization name NCSU
 Department Genetics
 Lab Gibson
 Street address South Gardner Hall
 City Raleigh
 State/province NC
 ZIP/postal code 27695
 Country USA
 Platform ID GPL6947
 Series (1) GSE17005 Geographical Genomics of Human Leukocytes Gene
 Expression Variation

Data table header descriptions

Data table header descriptions

ID_REF
VALUE log2 transformed and median-centered signal

Data table

ID_REF	VALUE
ILMN_1762337	-0.270341813
ILMN_2055271	0.212533766
ILMN_1736007	-0.065211006
ILMN_2393229	-0.294753746
ILMN_1800310	-0.179540575
ILMN_1779670	-0.156079383
ILMN_2321282	-0.070913334
ILMN_1671474	0.082277201
ILMN_1772582	0.227446423
ILMN_1735698	0.000532764
ILMN_1653355	0.197030952
ILMN_1717783	-0.284910419
ILMN_1705025	0.025293988
ILMN_1814316	-0.022732659
ILMN_2359168	0.04036411
ILMN_1731507	-0.480941688
ILMN_1787689	-0.046885413
ILMN_1745607	-0.099216214
ILMN_2130495	0.043699207
ILMN_1668111	-0.218402127

Total number of rows: 48803

Table truncated, full table size 1210 kbytes.