BGGN 239 – Week 2

Bulk and single-cell gene expression analysis

Ferhat Ay ferhatay@lji.org

Associate Professor of Computational Biology, LJI Department of Pediatrics & BISB PhD Program, UCSD

RNA-seq data flow



Slide adapted from Priya Pantham, PhD

Reads and mapping them





Read mapping – popular tools

• STAR, Bowtie, HISAT2, TopHat

Slide adapted from Priya Pantham, PhD

Visualizing mapped reads



Issues:

- 1. Gene length (does not matter much for across sample comparisons)
- 2. Read depth
- 3. GC content, mappability

- Reads per million or Counts per million
- Does not account for transcript length
- OK to use for sequencing protocols where reads are generated irrespective of gene length

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

- RPKM: Reads Per Kilobase of transcript per Million mapped reads
- FPKM*: Fragments Per Kilobase of transcript per Million mapped reads
- FPKM (or RPKM) attempts to normalize for gene size and library depth

*Fragments can mean either individual reads (SE) or paired-end reads that map together (PE)

 $\label{eq:RPKM} RPKM = \frac{\text{Number of reads mapped to gene} \times 10^3 \times 10^6}{\text{Total number of mapped reads} \times \text{gene length in bp}}$

- TPM: Transcripts per million (Transcripts Per Kilobase Million)
- Another form of normalization for gene length and sequencing depth, but in a slightly different order

$$\mathrm{TPM} = A imes rac{1}{\sum(A)} imes 10^6$$

Where
$$A = \frac{\text{total reads mapped to gene} \times 10^3}{\text{gene length in bp}}$$

$$\mathrm{TPM} = rac{RPKM}{\sum(RPKM)} imes 10^6$$

Which distribution better captures count data?



Dispersion matters!

- --: estimate using edgeR
- -----: fit to real RNA- Seq data
- ——: Poisson variance for each mean



- α : "dispersion" $\alpha = (\mu v) / \mu^2$ (squared coefficient of variation of extra-Poisson variability)
- Dispersion is a measure of the spread or variability in the data
- Biological Data is often 'overdispersed'. With increasing mean the variance grows disproportionally
- Negative binomial model can account for this overdispersion

SLE paper



SLE paper



IFN-20

Homework #1

SLE mini project

- Do the same exercise (SLE.qmd) for another cell type and answer the questions below.
 - How many genes were differentially expressed at adjusted p-value cutoff of 5%? how may up and down-regulated?
 - How many genes remain when you filter with log2FC greater than 1 versus absolute log2FC greater than 1?
 - Write the resulting short list of genes (p.adj <0.05 and log2FC >1) in a csv file.
 - From that short list, select one gene and write 3-4 sentences about how this gene in this specific cell type may be relevant to SLE. Ask Google and ChatGPT for help if you like.

Questions?

Functional analysis of gene sets

GO Term Enrichment Analysis



Homework #2

GO Term enrichment

- In the enrichGO function, try setting universe = names(sig_genes) instead of universe = names(all_genes_list). What happened? How many terms are statistically significant now?
- In the enrichGO function, set ont = "CC" rather than ont = "BP". What did this do? Do you believe BP or CC will be more relevant for most use-cases?
- Go to the NYU link and select one other visualization you like to use. Add a code chunk that generates this visualization. <u>https://learn.gencore.bio.nyu.edu/rna-seq-analysis/over-representation-analysis/</u>

Thank You!

- Paramita Dutta LJI
- Priya Pantham UCSD
- Barry Grant UCSD

Resources

- <u>https://learn.gencore.bio.nyu.edu/rna-seq-analysis/over-representation-analysis/</u>
- <u>https://allisonhorst.com/r-packages-functions</u>
- <u>http://bioconductor.org/packages/devel/bioc/vignettes/DESeq2/inst/doc/DE</u>
 <u>Seq2.html</u>
- <u>http://yulab-smu.top/biomedical-knowledge-mining-book/enrichment-overview.html</u>

single-cell!

Evolution of single cell techniques (till 2018)

<u>Valentine Svensson</u> [⊡], <u>Roser Vento-Tormo</u> & <u>Sarah A Teichmann</u> [⊡]

CITE-seq workflow - wet

CITE-seq workflow - dry

Things you can do with ADTs

Shared nearest neighbor graph (SNN)

Louvain algorithm

- Is a greedy algorithm
- Weighted graphs
- This algorithm has been widely utilized in many application domains because of:
 - Its rapid convergence properties
 - High modularity
 - Hierarchal partitioning

Louvain algorithm

Louvain algorithm

Algorithm [edit]

The value to be optimized is modularity, defined as a value in the range [-1/2, 1] that measures the density of links inside communities compared to links between communities.^[1] For a weighted graph, modularity is defined as:

$$Q = rac{1}{2m}\sum_{ij}igg[A_{ij}-rac{k_ik_j}{2m}igg]\delta(c_i,c_j),$$

where

- A_{ij} represents the edge weight between nodes i and j;
- k_i and k_j are the sum of the weights of the edges attached to nodes i and j, respectively;
- *m* is the sum of all of the edge weights in the graph;
- c_i and c_j are the communities of the nodes; and
- δ is Kronecker delta function ($\delta(x,y)=1$ if x=y,0 otherwise).

Based on the above equation, the modularity of a community c can be calculated as:

 $Q_c = rac{\Sigma_{in}}{2m} - (rac{\Sigma_{tot}}{2m})^2,$

where

• Σ_{in} is the sum of edge weights between nodes within the community c (each edge is considered twice); and

• Σ_{tot} is the sum of all edge weights for nodes within the community (including edges which link to other communities).

Resources

- Review: <u>https://www.nature.com/articles/nprot.2017.149</u>
- 10x datasets: <u>https://www.10xgenomics.com/resources/datasets</u>
- CITE-seq: <u>https://www.nature.com/articles/nmeth.4380</u>
- Ab-seq: <u>https://www.nature.com/articles/srep44447</u>
- Louvain algorithm: <u>https://towardsdatascience.com/louvain-algorithm-93fde589f58c</u>