Lecture 11. RNAseq

Michael Schatz

March 6, 2017 JHU 601.749: Applied Comparative Genomics



Assignment 4: Due March I

Assignment 4: Read mapping and variant calling

Assignment Date: Thursday, Feb. 22, 2018 Due Date: Thursday, Mar. 1, 2018 @ 11:59pm

Assignment Overview

In this assignment, you will align reads to a reference genome to call SNPs and short indels. Then, you will perform an experiment to empirically determine the "mappability" of a genomic region. Finally, you will investigate some empirical behavior of the binomial test for heterozygous variant calling.

As a reminder, any questions about the assignment should be posted to Plazza. Don't forget to read the Resources section at the bottom of the page!

Question 1. Small Variant Analysis [10 pts]

Download chromosome 22 from build 38 of the human genome from here: http://hgdownload.cse.ucsc.edu/goldenPath/hg38/chromosomes/chr22.fa.gz

Download the read set from here: http://schatzlab.cshl.edu/data/teaching/sample.tgz

For this question, you may find this tutorial helpful: http://clavius.bc.edu/~erik/CSHL-advanced-sequencing/freebayes-tutorial.html

 1a. How many reads align to the reference? How many reads did not align? How many aligned reads had a mate that did not align (AKA singletons)? Count each read in a pair separately. [Hint: Build the index using bowtie2-build, align reads using bowtie2, analyze with samtools flagstat.]

Assignment 5: Due March 8

Assignment 5: Genome Arithmetic

Assignment Date: Thursday, March 1, 2018 Due Date: Thursday, March 8, 2018 @ 11:59pm

Assignment Overview

In this assignment, you will call structural variants and analyze the properties of variants in the human genome. Make sure to show your work in your writeup! As before, any questions about the assignment should be posted to Piazza.

Question 1. Gene Annotation Preliminaries [10 pts]

Download the annotation of build 38 of the human genome from here: ftp://ftp.ensembl.org/pub/release-87/gtf/homo_sapiens/Homo_sapiens.GRCh38.87.gtf.gz

- Question 1a. How many many GTF data lines are in this file? [Hint: The first few lines in the file beginning with "#" are so-called "header" lines describing thing like the creation date, the genome version (more on that later in the course), etc. Header lines should not be counted as data lines.]
- Question 1b. How many annotated protein coding genes are on each autosome of the human genome? [Hint: Protein coding genes will have "gene" in the 3rd column, and contain the following text: gene_biotype "protein_coding"]
- Question 1c. What is the maximum, minimum, mean, and standard deviation of the span of protein coding genes? [Hint: use the genes identified in 1b]
- Question 1d. What is the maximum, minimum, mean, and standard deviation in the number of exons for protein coding genes? [Hint: you should separately consider each isoform for each protein coding gene]

Campylobacter jejuni RM1221 30.3%GC

					— — .
en mys. dr. 25	4×-75%	15. 47.55	. 47	<u></u>	40-50, 40-70, 40-50, 92-70

Mycobacterium smegmatis MC2 67.4%GC

-		_								
				-						
_		_								
			_							
a a state of the s	10 m	2.57	12 2 7 02	17 200	12 22 200	ার <u>্</u> র্যা ৫(১০	· · · · · · · · · · · · · · · · · · ·	v r7(x, − e	್ಲೇ ಇನ್ ರಿ ಗ್ರಿ	 10.77 20.02

P(heads) = 61/64 (95.4%) P(tails) = 3/64 (4.6%)

How many flips until my first tail?

\$./coinflip.pl 0.046875 1000

- 0: ННННННННННН 15
- 1: HHHHHHT 7
- 2: HHHHHHHHHH 12
- 3: НННННННННННННННННН 24
- 4: HT 2
- 5: НННННННННН 14
- 6: НННННННН 10
- 7: ННННННННННН 14
- 8: HHHHHT 6
- 9: HHHHHHHHH 11
- 11: НННННННННННННННННННННННННННННННННН
- 13: HHHT 4
- 14: ННННННННННН 15
- 15: ННННННННННННННННННННННННННННННННН
- 16: HHHHHT 6
- 17: НННННННННННННННННННННННННННННННН
- 18: ННННННННННННННННННННН 26
- 19: ННННННННН 12

P(heads) = 61/64 (95.4%) P(tails) = 3/64 (4.6%)

How many flips until my first tail?



P(heads) = 61/64 (95.4%) P(tails) = 3/64 (4.6%)

How many flips until my first tail?

Geometric Distribution: $P(X=x) = p_{heads}^{x-1}p_{tails}$



P(heads) = 61/64 (95.4%) P(tails) = 3/64 (4.6%)

How many flips until my first tail?

Geometric Distribution: $P(X=x) = p_{heads}^{x-1}p_{tails}$



Flips until heads

P(heads) = 61/64 (95.4%) P(tails) = 3/64 (4.6%)

How many flips until my first tail?

Geometric Distribution: $P(X=x) = p_{heads}^{x-1}p_{tails}$



Flips until heads

Stop Codon Frequencies



If the sequence is mostly A+T, then likely to form stop codons by chance!

In High A+T (Low G+C):

Frequent stop codons; Short Random ORFs; long ORFs likely to be true genes

In High G+C (Low A+T):

Rare stop codons; Long Random ORFs; harder to identify true genes

A relationship between GC content and coding-sequence length.

Oliver & Marín (1996) J Mol Evol. 43(3):216-23.



*-seq in 4 short vignettes



RNA-seq



Gene expression patterns of breast carcinomas distinguish tumor subclasses with clinical implications. Sørlie et al (2001) *PNAS*. 98(19):10869-74.

RNA-seq Overview



RNA-seq Overview



RNA-seq Overview



RNA-seq Challenges



Challenge I: Eukaryotic genes are spliced

RNA-Seq Approaches



Fig. 2 Read mapping and transcript identification strategies. Three basic strategies for regular RNA-seq analysis. **a** An annotated genome is available and reads are mapped to the genome with a gapped mapper. Next (novei) transcript discovery and quantification can proceed with or without an annotation file. Novel transcripts are then functionally annotated. **b** If no novel transcript discovery is needed, reads can be mapped to the reference transcriptome using an ungapped aligner. Transcript identification and quantification can occur simultaneously. **c** When no genome is available, reads need to be assembled first into contigs or transcripts. For quantification, reads are mapped back to the novel reference transcriptome and further analysis proceeds as in (**b**) followed by the functional annotation of the novel transcripts as in (**a**). Representative software that can be used at each analysis step are indicated in *bold text*. Abbreviations: *GFF* General Feature Format, *GTF* gene transfer format, *RSEM* RNA-Seq by Expectation Maximization

A survey of best practices for RNA-seq data analysis Conesa et al (2016) Genome Biology. doi 10.1186/s13059-016-0881-8

RNA-Seq Approaches



ated in bold text. Abbreviations: GFF General Feature Format, GTF gene transfer format,

RSEM RNA-Seg by Expectation Maximization

A survey of best practices for RNA-seq data analysis Conesa et al (2016) Genome Biology. doi 10.1186/s13059-016-0881-8

RNA-seq Challenges



Challenge I: Eukaryotic genes are spliced

Solution: Use a spliced aligner, and assemble isoforms

TopHat: discovering spliced junctions with RNA-Seq.

Trapnell et al (2009) Bioinformatics. 25:0 1105-1111



Challenge 2: Read Count != Transcript abundance



Counting Reads that align to a gene DOESN'T work!

- Overall Coverage: 1M reads in experiment 1 vs 10M reads in experiment 2
- Gene Length: gene 3 is 10kbp, gene 4 is 100kbp

1. RPKM: Reads Per Kilobase of Exon Per Million Reads Mapped (Mortazavi et al, 2008)



Counting Reads that align to a gene DOESN'T work!

- Overall Coverage: 1M reads in experiment 1 vs 10M reads in experiment 2
- Gene Length: gene 3 is 10kbp, gene 4 is 100kbp

1. RPKM: Reads Per Kilobase of Exon Per Million Reads Mapped (Mortazavi et al, 2008)

(Count reads aligned to gene) / (length of gene in kilobases) / (# millions of read mapped)

=> Wait a second, reads in a pair arent independent!



Counting Reads that align to a gene DOESN'T work!

- Overall Coverage: 1M reads in experiment 1 vs 10M reads in experiment 2
- Gene Length: gene 3 is 10kbp, gene 4 is 100kbp

1. RPKM: Reads Per Kilobase of Exon Per Million Reads Mapped (Mortazavi et al, 2008) => Wait a second, reads in a pair arent independent!

2. FPKM: Fragments Per Kilobase of Exon Per Million Reads Mapped (Trapnell et al, 2010) ⇒ Does a much better job with short exons & short genes by boosting coverage

 \Rightarrow Wait a second, FPKM depends on the average transcript length!



Counting Reads that align to a gene DOESN'T work!

- Overall Coverage: 1M reads in experiment 1 vs 10M reads in experiment 2
- Gene Length: gene 3 is 10kbp, gene 4 is 100kbp

1. RPKM: Reads Per Kilobase of Exon Per Million Reads Mapped (Mortazavi et al, 2008)

=> Wait a second, reads in a pair arent independent!

2. FPKM: Fragments Per Kilobase of Exon Per Million Reads Mapped (Trapnell et al, 2010)

=> Wait a second, FPKM depends on the average transcript length!

3. TPM: Transcripts Per Million (Li et al, 2011)

- ⇒ If you were to sequence one million full length transcripts, TPM is the number of transcripts you would have seen of type i, given the abundances of the other transcripts in your sample
- => Recommend you use TPM for all analysis, easy to compute given FPKM

$$\mathrm{TPM}_i = \left(\frac{\mathrm{FPKM}_i}{\sum_j \mathrm{FPKM}_j}\right) \cdot 10^6$$









Key point : The length of the actual molecule from which the fragments derive is crucially important to obtaining accurate abundance estimates.



The gene has three isoforms (red, green, blue) of the same length. Our initial expectation is all 3 isoforms are equally expressed

There are five reads (a,b,c,d,e) mapping to the gene.

- Read a maps to all three isoforms
- Read d only to red
- Reads b,c,e map to each of the three pairs of isoforms.

What is the most likely expression level of each isoform?

Models for transcript quantification from RNA-seq

Pachter, L (2011) arXiv. 1104.3889 [q-bio.GN]



The gene has three isoforms (red, green, blue) of the same length. Initially every isoform is assigned the same abundance (red=1/3, green=1/3, blue=1/3)

There are five reads (a,b,c,d,e) mapping to the gene. Read a maps to all three isoforms, read d only to red, and the other three (reads b,c,e) to each of the three pairs of isoforms.

During the expectation (E) step reads are proportionately assigned to transcripts according to the (current) isoform abundances (RGB): a=(.33,.33,.33), b=(0,.5,.5), c=(.5,.5), d=(1,0,0), e=(.5,.5,0)



The gene has three isoforms (red, green, blue) of the same length. Initially every isoform is assigned the same abundance (red=1/3, green=1/3, blue=1/3)

There are five reads (a,b,c,d,e) mapping to the gene. Read a maps to all three isoforms, read d only to red, and the other three (reads b,c,e) to each of the three pairs of isoforms.

During the expectation (E) step reads are proportionately assigned to transcripts according to the (current) isoform abundances (RGB): a=(.33,.33,.33), b=(0,.5,.5), c=(.5,.5), d=(1,0,0), e=(.5,.5,0)

Next, during the maximization (M) step isoform abundances are recalculated from the proportionately assigned read counts:

red: 0.47 = (0.33 + 0.5 + 1 + 0.5)/(2.33 + 1.33 + 1.33)blue: 0.27 = (0.33 + 0.5 + 0.5)/(2.33 + 1.33 + 1.33)green: 0.27 = (0.33 + 0.5 + 0.5)/(2.33 + 1.33 + 1.33)



The gene has three isoforms (red, green, blue) of the same length. Initially every isoform is assigned the same abundance (red=1/3, green=1/3, blue=1/3)

There are five reads (a,b,c,d,e) mapping to the gene. Read a maps to all three isoforms, read d only to red, and the other three (reads b,c,e) to each of the three pairs of isoforms.

During the expectation (E) step reads are proportionately assigned to transcripts according to the (current) isoform abundances (RGB): a=(.33,.33,.33), b=(0,.5,.5), c=(.5,.5), d=(1,0,0), e=(.5,.5,0)

Next, during the maximization (M) step isoform abundances are recalculated from the proportionately assigned read counts:

red: 0.47 = (0.33 + 0.5 + 1 + 0.5)/(2.33 + 1.33 + 1.33)blue: 0.27 = (0.33 + 0.5 + 0.5)/(2.33 + 1.33 + 1.33)green: 0.27 = (0.33 + 0.5 + 0.5)/(2.33 + 1.33 + 1.33)

Repeat until convergence!



The gene has three isoforms (red, green, blue) of the same length. Initially every isoform is assigned the same abundance (red=1/3, green=1/3, blue=1/3)

There are five reads (a,b,c,d,e) mapping to the gene. Read a maps to all three isoforms, read d only to red, and the other three (reads b,c,e) to each of the three pairs of isoforms.

During the expectation (E) step reads are proportionately assigned to transcripts according to the (current) isoform abundances (RGB): a=(.33,.33,.33), b=(0,.5,.5), c=(.5,.5), d=(1,0,0), e=(.5,.5,0)

Next, during the maximization (M) step isoform abundances are recalculated from the proportionately assigned read counts:

red: 0.47 = (0.33 + 0.5 + 1 + 0.5)/(2.33 + 1.33 + 1.33)blue: 0.27 = (0.33 + 0.5 + 0.5)/(2.33 + 1.33 + 1.33)green: 0.27 = (0.33 + 0.5 + 0.5)/(2.33 + 1.33 + 1.33)

Repeat until convergence!

Sailfish: Fast & Accurate RNA-seq Quantification



Sailfish enables alignment-free isoform quantification from RNA-seq reads using lightweight algorithms Patro et al (2014) Nature Biotechnology 32, 462–464 doi:10.1038/nbt.2862