

## RNA SEQUENCING DATA ANALYSIS

### Differential Expression in RNA-seq:

**Defining Differential Expression:** Differential expression in the context of RNA-seq refers to the changes in gene expression levels between different experimental groups. This could be different treatments, conditions, time points, or any other experimental variable.

**Statistical Significance:** The determination of differential expression relies on statistical tests to ensure that the observed changes in expression are not due to random chance. This involves comparing the read counts of a particular gene across different conditions.

**Read Counts and Normalization:** RNA-seq data consists of read counts, which are the number of times a particular sequence has been read during sequencing. These counts need to be normalized to account for differences in sequencing depth and gene length across samples.

**Methods for Analysis:** The document likely describes various statistical methods and tools used to analyze differential expression in RNA-seq data. Common approaches include the use of models based on distributions like Poisson or Negative Binomial, and tools like DESeq, edgeR, or Limma.

**Challenges in Analysis:** Differential expression analysis in RNA-seq data can be challenging due to issues like low counts in some genes, overdispersion (variance larger than the mean), and the need for multiple testing corrections due to the large number of genes tested simultaneously.

**Interpretation of Results:** The results from differential expression analysis often include lists of genes with their respective p-values and fold-changes, indicating the level of expression change and its statistical significance. These results need careful interpretation, considering the biological context and the quality of the data.

### Statistical Tools for RNA-seq:

**Nature of RNA-seq Data:** RNA-seq data is fundamentally different from other types of gene expression data, like microarrays. It is composed of discrete read counts, representing the number of times a particular RNA molecule is sequenced. This discrete nature requires specific statistical treatment.

**Normalization Methods:** The document probably explains normalization techniques. Normalization is critical in RNA-seq analysis to account for various factors like differences in sequencing depth, RNA composition of the samples, and gene length. Techniques such as RPKM (Reads Per Kilobase of transcript, per Million mapped reads) or TPM (Transcripts Per Million) might be discussed.

**Statistical Distributions for Modeling RNA-seq Data:** Since RNA-seq data are count data, they are often modeled using specific statistical distributions. The Poisson distribution was initially used, but it often underestimates the variance in RNA-seq data. Therefore, the Negative Binomial distribution is more commonly employed, as it can handle overdispersion (where variance exceeds the mean).

**Bioinformatics Tools for Differential Expression Analysis:** There are several bioinformatics tools specifically designed for RNA-seq data analysis. Tools like DESeq, edgeR, and Cufflinks are mentioned, explaining how they use statistical models to identify differentially expressed genes between different conditions or treatments.

**Dealing with Technical Variability and Batch Effects:** RNA-seq experiments can be subject to technical variations and batch effects. The document might describe statistical methods to adjust for these, ensuring that the results are reflective of biological differences rather than technical artifacts.

**Multiple Testing Correction:** When testing thousands of genes for differential expression, the document would discuss the importance of correcting for multiple hypothesis testing, using methods like the False Discovery Rate (FDR) to control for type I errors (false positives).

**Statistical Significance and Biological Relevance:** Finally, the section might emphasize the distinction between statistical significance and biological relevance, noting that a statistically significant result does not always imply a biologically meaningful difference.

### RPKM (Reads Per Kilobase) and Length Bias:

**Understanding RPKM:** RPKM is a normalization method used in RNA-seq data analysis. It corrects for two major factors: the total number of reads in a sample (sequencing depth) and the length of the RNA transcript. The formula for RPKM is:

$$RPKM = \frac{\text{Number of reads mapped to a gene}}{(\text{Length of the gene in kilobases}) \times (\text{Total reads in millions})}$$

**Length Bias in RNA-seq:** Length bias refers to the tendency for longer transcripts to yield more reads than shorter ones, simply due to their size. Without correction, this bias can lead to misleading interpretations, where longer genes appear to be more highly expressed than they actually are.

**Importance of Normalization:** The section likely emphasizes the importance of normalization in RNA-seq data analysis. Normalization, such as that provided by RPKM, ensures that the read counts are comparable across genes and across samples, allowing for more accurate determination of gene expression levels.

**Limitations of RPKM:** Despite its usefulness, RPKM has limitations. It does not account for variations in sequencing depth or library size across different samples, which can lead to inaccuracies when comparing expression levels between samples. Methods like TPM (Transcripts Per Million) and normalization techniques used in tools like DESeq and edgeR are sometimes preferred for this reason.

**Comparing Across Samples:** The document might discuss the challenges and solutions in comparing expression levels across different samples using RPKM. It likely mentions the need for additional normalization steps or statistical models to make accurate cross-sample comparisons.

**Contextualizing RPKM within RNA-seq Analysis:** This section may place RPKM within the larger context of RNA-seq data analysis, highlighting its role as one of many tools and considerations necessary for accurate interpretation of RNA-seq data.

## Statistical Distributions

**Poisson Distribution:** Initially, the Poisson distribution was commonly used in RNA-seq data analysis. This distribution is suitable for modeling count data where the mean and variance are equal, which is often the assumption for RNA-seq read counts. The Poisson model is simple and mathematically convenient but has limitations.

**Limitations of the Poisson Model:** The Poisson distribution often falls short in RNA-seq data analysis due to the phenomenon of overdispersion – where the observed variance in read counts is greater than the mean. This overdispersion can lead to an increased rate of false positives if not properly accounted for.

**Negative Binomial Distribution:** To address overdispersion, the Negative Binomial distribution is frequently used in RNA-seq data analysis. Unlike the Poisson distribution, the Negative Binomial allows the variance to be greater than the mean, making it more flexible and suitable for modeling RNA-seq read counts.

**Modeling Overdispersion:** The Negative Binomial distribution includes an extra parameter to model overdispersion, providing a better fit for RNA-seq data compared to the Poisson model. This improved fit reduces the likelihood of false positives in differential expression analysis.

**Software Implementations:** Bioinformatics tools like DESeq, edgeR, and others use the Negative Binomial distribution for differential expression analysis. These tools not only implement statistical models but also provide methods for normalization and other crucial steps in RNA-seq data analysis.

**Other Distributions and Models:** The section might also mention other distributions or statistical models used in specific contexts within RNA-seq data analysis, although the Negative Binomial distribution is the most commonly used for differential expression analysis due to its flexibility in handling overdispersion.

## Likelihood Ratio Test:

**Basics of Likelihood Ratio Test (LRT):** The LRT is a statistical method used to compare the fit of two models to a given dataset. In the context of RNA-seq, it's used to compare a simpler model (null hypothesis) that assumes no difference in expression across conditions, against a more complex model (alternative hypothesis) that allows for differential expression.

**Application in RNA-seq:** In RNA-seq data analysis, the LRT is used to determine whether the inclusion of a specific variable (like treatment condition, time point, etc.) significantly improves the model fit, suggesting differential expression of genes.

**Statistical Models:** The LRT typically involves fitting models to RNA-seq count data, often using Negative Binomial distribution, which accounts for overdispersion in the data. The test evaluates if the additional complexity (more parameters) in the alternative model significantly improves the fit to the data compared to the null model.

**Calculating the Test Statistic:** The test statistic for the LRT is calculated as two times the difference in the log-likelihoods of the two models. This statistic follows a chi-square distribution, allowing for the determination of a p-value.

**Interpreting Results:** A low p-value from the LRT suggests that the more complex model (which may include terms for different experimental conditions or treatments) provides a significantly better fit, indicating differential gene expression. High p-values indicate insufficient evidence to reject the null hypothesis of no differential expression.

**Advantages and Limitations:** The document might discuss the advantages of LRT, such as its flexibility and robustness in model comparison. However, it might also cover limitations, such as sensitivity to sample size and assumptions about the underlying statistical model.

**Role in Differential Expression Analysis:** The LRT is a fundamental tool in differential expression analysis in RNA-seq, helping identify genes whose expression levels significantly vary across different conditions or treatments.

## Fisher's Exact Test:

**Basics of Fisher's Exact Test:** Fisher's Exact Test is a statistical significance test used for analyzing categorical data, specifically in 2x2 contingency tables. It's used to examine the significance of the association between two kinds of classifications.

**Application in RNA-seq Analysis:** In RNA-seq, Fisher's Exact Test can be used to assess whether the observed distribution of reads across different conditions is significantly different than expected by chance. This is particularly relevant for datasets with *small sample sizes or where replicates are not available*.

**Calculating the Test:** The test calculates the probability (p-value) of observing the data assuming that there is no association between the two variables (e.g., gene expression and treatment condition). The p-value is derived from the hypergeometric distribution.

**Advantages:** Fisher's Exact Test is exact, as it does not rely on large sample approximations like the chi-square test. This makes it suitable for small sample sizes where other tests may be less reliable.

**Limitations:** While Fisher's Exact Test is useful for small sample sizes, it has limitations in larger datasets or those with replicates. In such cases, tests that take into account the variability across replicates, such as those based on the Negative Binomial distribution, may be more appropriate.

**Contextual Use:** The document might discuss when Fisher's Exact Test is most appropriately used in RNA-seq data analysis, often as a complementary approach to other statistical tests, or in specific scenarios where other tests are not applicable.

**Interpreting Results:** Interpretation of the results from Fisher's Exact Test in RNA-seq involves understanding the p-values in the context of biological significance and the experimental setup, keeping in mind the limitations of the test.

• all tables with probabilities less than or equal to the observed table are used to compute p value

Fisher's exact test									
$p=0.02447552$ ( $k=6$ )					$p=0.004079254$ ( $k=1$ )				
	Cancer	Healthy				Cancer	Healthy		
Smokers	6	1	n=7		Smokers	1	6	n=7	
Non smokers	1	5	N-n		Non smokers	6	0	N-n	
	K	N-K	N = 13			K	N-K	N = 13	
$p=0.0005827506$ ( $k=7$ )					$p\text{-value} = 0.02447552 + 0.0005827506 + 0.004079254 = 0.029$				
	Cancer	Healthy				Cancer	Healthy		
Smokers	7	0	n=7						
Non smokers	0	6	N-n						
	K	N-K	N = 13						

## Challenges and Limitations

**Complexity of Biological Systems:** RNA-seq data is a reflection of complex biological systems. Variability in gene expression can be influenced by numerous factors, including genetic differences, environmental conditions, and experimental procedures.

**Technical Variability:** Technical variations can arise from differences in sample preparation, sequencing depths, and sequencing platforms. These variations can affect the accuracy and reproducibility of RNA-seq results.

**Biases in Sequencing Data:** RNA-seq data can have inherent biases, such as GC-content bias, sequence-specific bias, and length bias. Correcting these biases is crucial for accurate data interpretation.

**Handling of Low-Abundance (low) Transcripts:** Detecting and quantifying low-abundance transcripts is challenging. Low-expression genes can be difficult to distinguish from background noise, leading to potential false negatives.

**Overdispersion and Model Fitting:** RNA-seq data often exhibits overdispersion, where the observed variance in read counts is higher than expected. This necessitates the use of appropriate statistical models, such as the Negative Binomial distribution, for accurate analysis.

**Multiple Testing Problem:** Given the large number of genes analyzed in RNA-seq experiments, the issue of multiple testing is significant. Correction methods (e.g., False Discovery Rate) are required to control the rate of false positives.

**Interpretation of Differential Expression:** Determining biological relevance from statistically significant differential expression results is not straightforward. It requires careful consideration of the biological context and additional validation experiments.

**Computational Resources and Expertise:** RNA-seq data analysis requires substantial computational resources and bioinformatics expertise. The complexity of data processing and analysis can be a limiting factor, especially in resource-limited settings.

**Data Integration and Comparative Analysis:** Integrating RNA-seq data with other types of genomic or proteomic data for a more comprehensive understanding of biological processes is challenging but essential for holistic insights.

**Evolution of Analytical Tools and Methods:** The field of RNA-seq is rapidly evolving, with continual improvements in sequencing technologies and analytical methods. Staying updated with these advancements is crucial for researchers.

## Practical Examples and Calculations:

**Example Datasets:** The document may introduce example RNA-seq datasets derived from real experiments. These examples could include data from different experimental conditions, treatments, or time points.

**Data Preprocessing Steps:** Practical guidance on preprocessing steps such as quality control, alignment of reads to a reference genome, and quantification of gene expression levels might be provided, along with examples of commands or scripts used in popular bioinformatics tools.

**Normalization Calculations:** Demonstrations of how to perform normalization, such as calculating RPKM/FPKM (Fragments Per Kilobase of transcript per Million mapped reads) or TPM (Transcripts Per Million), with step-by-step calculations on the example data.

**Differential Expression Analysis:** Detailed examples of how to conduct differential expression analysis, possibly using statistical software or bioinformatics tools like DESeq2, edgeR, or Limma. This might include code snippets and interpretations of output results (like p-values, fold changes, etc.).

**Addressing Technical Variability:** Examples of how to identify and correct for batch effects and technical variability in RNA-seq data, perhaps including illustrations of before-and-after comparisons using principal component analysis (PCA) or other methods.

**Statistical Tests Application:** Practical application of statistical tests discussed in the document (like the Likelihood Ratio Test or Fisher's Exact Test) on RNA-seq data, including how to interpret the statistical significance and biological relevance of the results.

Example Scenario: RPKM

**Objective:** To identify differentially expressed genes in liver tissue between two groups of mice - those treated with a new drug (Treatment) and untreated (Control).

**Data:** RNA-seq data from liver tissues of 5 treated and 5 untreated mice.

**Step 1:** Data Preprocessing

**Quality Control:** Assess the quality of the raw sequencing data using a tool like FastQC. Check for issues like adapter content, GC content, and sequence quality.

**Read Alignment:** Align the reads to a mouse reference genome using an aligner like STAR.

**Quantification:** Quantify the number of reads mapping to each gene using a tool like featureCounts.

**Step 2:** Normalization

**Calculate RPKM:** For each gene in each sample, calculate RPKM to normalize read counts.

Example Calculation: RPKM Calculation for a Gene

- Read Count for Gene X in Sample 1: 800 reads
- Length of Gene X: 1.5 kb
- Total Reads in Sample 1: 30 million
- RPKM for Gene X in Sample 1:  $\frac{800}{1.5 \times 30} = 17.78$

- Formula:  $RPKM = \frac{\text{Read Count}}{\text{Gene Length (kb)} \times \text{Total Reads (millions)}}$
- Example: If a gene has 1,000 reads, a length of 2 kb, and the total reads are 20 million, its RPKM =  $\frac{1000}{2 \times 20} = 25$

Step 3: Differential Expression Analysis

**Using DESeq2:** Load the RPKM normalized data into DESeq2.

**Model Fit:** Fit a model considering the treatment as the variable of interest.

**Statistical Testing:** Perform a Wald test or Likelihood Ratio Test for each gene to determine differential expression.

**Step 4:** Results Interpretation

**Identify Significant Genes:** Select genes with adjusted p-values (False Discovery Rate) below a threshold (e.g., 0.05).

**Biological Interpretation:** Analyze the biological functions of significantly expressed genes using databases like Gene Ontology.

**Step 5:** Validation

**qPCR Validation:** Select a few differentially expressed genes for validation using quantitative PCR (qPCR). Compare qPCR results with RNA-seq findings for consistency.

**Functional Assays:** If a gene of interest is implicated in drug metabolism, perform functional assays to confirm its role.

**Step 6:** Further Analysis

**Pathway Analysis:** Use tools like KEGG to identify pathways significantly enriched with differentially expressed genes.

**Data Visualization:** Create visualizations such as heat maps or volcano plots to display the results.

Chi-square score:

DF:

Significance Level:

- ☐ 0.01
- ☒ 0.05
- ☐ 0.10

The P-Value is .050044. The result is *not* significant at p < .05.

Calculate