QBIO 490: Directed Research - Multi-Omic Analysis

Fall 2024 Review Project

Due: Tuesday, November 19th (11:59 pm). Submit your GitHub link to Brightspace, with all your code and code outputs in a folder called r_review_name within your qbio_490_name repo. Please email extension requests (include the reason for your extension and a proposed new due date) to Mahija and Wade by **Thursday, November 21st 11:59 pm**. This is a hard deadline, and no requests will be accepted after this date, except for reasons of emergency or illness.

Purpose:

This review project is meant to recap the analyses we've performed so far in R. It's also intended to rehash various parts of scientific writing and communication. For this project, please do your own work and submit your own written report, but you are more than encouraged to discuss ideas and debug code in groups! Note there are *three parts* to this assignment.

Overview:

In the first part, you will be answering short questions about R and TCGA. In the second part, you will choose one of two analyses of SKCM clinical, transcriptomic, and epigenomic data to explore a predetermined question about SKCM. In the third and final part, you will briefly write up your interpretations.

Part 1: Review Questions

General Concepts

- 1. What is TCGA and why is it important?
- 2. What are some strengths and weaknesses of TCGA?

Coding Skills

- 1. What commands are used to save a file to your GitHub repository?
- 2. What command(s) must be run in order to use a package in R?
- 3. What command(s) must be run in order to use a *Bioconductor* package in R?
- 4. What is boolean indexing? What are some applications of it?
- 5. Draw a mock up (just a few rows and columns) of a sample dataframe. Show an example of the following and explain what each line of code does.
 - a. an ifelse() statement
 - b. boolean indexing

Part 2: SKCM Analysis

Before starting your analysis, you may find it helpful to read the following review article on SKCM to get a broad understanding of the cancer pathogenesis and possible treatment options. This may be especially helpful with understanding why each clinical variable was collected and what they mean. https://www.ncbi.nlm.nih.gov/pmc/articles/PMC3004577/

In this project, you will conduct multi-omic analyses to explore the following research question:

What are the differences between metastatic and non-metastatic SKCM across the epigenome and do these have any effect on the transcriptome?

Exploration of Methylation Patterns and Effect on Transcription

To do this, you must include at least the following analyses (at least 6 plots):

- 1. Difference in survival between metastatic and non-metastatic patients (KM plot)
- 2. Differential expression between non-metastatic and metastatic patients controlling for treatment effects, race, gender, and vital status (DESeq2 + Volcano plot)
 - a. Treatments must include radiation, chemotherapy, immunotherapy, molecular therapy, vaccine
 - b. If you run this on CARC, it may take up to 1-2 hours
- 3. Naive differential methylation between non-metastatic and metastatic patients (Volcano plot)
- 4. Direct comparison of methylation status to transcriptional activity across non-metastatic vs metastatic patients
- 5. Visualization of CpG sites and protein domains for 3 genes for a few genes (use UCSC genome browser)

All of your code can be in a R Notebook or R script, which you will push to GitHub and provide a repo link to Brightspace. As a part of the grading, we will check that your code runs with no errors starting from a clean environment. However, you can assume that any of the csv's we save in class are present (brca_clinical_data, brca_rna_clinical, brca_rna_genes, brca_rna_counts, brca_methylation_clinical, brca_methylation_betas, and brca_cpg_sites). Remember to comment your code so other people can follow along.

Technical Tips:

- The accession code for SKCM is TCGA-SKCM
- The following commands can be used to access the drug and radiation dataframes once SKCM clinical data has been downloaded from TCGA:

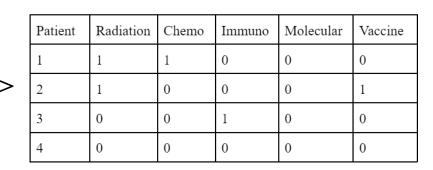
```
rad <- clinical.BCRtab.all$clinical_radiation_skcm[-c(1,2),]
drug <- clinical.BCRtab.all$clinical_drug_skcm[-c(1,2),]</pre>
```

- Metastasis status should be based on the rna_se@colData\$definition column.
 - Only consider "Metastatic" or "Primary solid Tumor" samples
- Be careful about what "barcode" columns you use! The patient id, sample id, and sample barcode columns are all named slightly differently across the different dataframes.

 Double check that the columns you are using to match index values are correct!
- For DESeq2 data preprocessing:
 - Use the rna se clinical data (rna se@colData).
 - Filter out genes with a total expression across all patients of < 20
 - Threshold padj values at 0.05 and log2FoldChange at |1|

- Since there are 5 different treatments and each individual may have multiple treatments, you must use a technique called **one-hot encoding** where you create a column for each treatment and give a 1/0 value for whether each patient underwent that treatment.
 - o For example:

Patient	Treatment
1	Radiation, chemo
2	Radiation, Vaccine
3	Immunotherapy
4	None



Part 3: Results and Interpretations

For each analysis, include an image of the relevant plot you created in Part 2 and a 3-4 sentence description answering the following question:

• Analyze the plot. What conclusions can you and can you not draw about differences

between metastatic and non-metastatic TCGA SKCM patients? Why?
1) Difference in survival between metastatic and non-metastatic patients
2) Expression differences between metastatic and non-metastatic patients
3) Methylation differences between metastatic and non-metastatic patients
4) Direct comparison of transcriptional activity to methylation status for 10 genes
5) Visualization of CpG sites and protein domains for 3 genes (use UCSC genome browser) for a few genes. Describe at least one academic article (research or review) that either supports or doesn't support your final conclusion for one of the genes. If previously published work doesn't support your analysis, explain why this might be the case.
At the end of your report, include a References page of all the articles you used. Any citation format works, as long as you are consistent (all MLA, APA, etc.). Reminder: we are permitting the use of properly attributed AI work on the coding portion of this assignment (ie part 2), but not on any written portions (parts 1 and 3).