

Exploring Clinical Computational Cancer Genomics Data (Tumor / Normal whole genome sequencing)

Overall question: Can you discover any mutations that are associated with treatment response?

- 1. Read and Merge the Mutation Annotation Format (MAF) files.** Each of these 50 files contains the genomic mutations observed in a different patient's tumor, obtained by biopsy and sequenced with whole-exome sequencing. Each row in a MAF file corresponds to a different mutation.
- 2. Subset for mutations that are not of the Variant Classification "Silent".** For the purposes of this analysis, we will restrict ourselves to substitutions which result in changes to the produced protein ("nonsynonymous mutations").
- 3. Find the 15 most common mutations.** Gene names are included in the column Hugo_Symbol and protein changes are stored in the column Protein_Change.
- 4. Perform a statistical test to explore if any mutated genes are enriched in patients who either responded or not.** Response labels for individual patients are found in the file data/sample-information.tsv.
- 5. Create a scatter plot of genes with the number of mutated patients on the x-axis and your results from question 4 on the y-axis.** Can the figure in any way improve readability? If so, recreate the plot using your suggestion(s).
- 6. How many samples are wild-type versus mutant with respect to the most significantly enriched gene from Question 4? Plot the number of nonsynonymous mutations per megabase in the mutant vs. wild-type samples. Is there a significant difference in the number of mutations between the two groups?** Information on the number of nonsynonymous mutations per megabase for each patient can be found in the file data/sample-information.tsv.
- 7. Write any conclusions that you have made based on your analysis. How might this analysis be improved or expanded upon? Please include all requested figures in your report.**