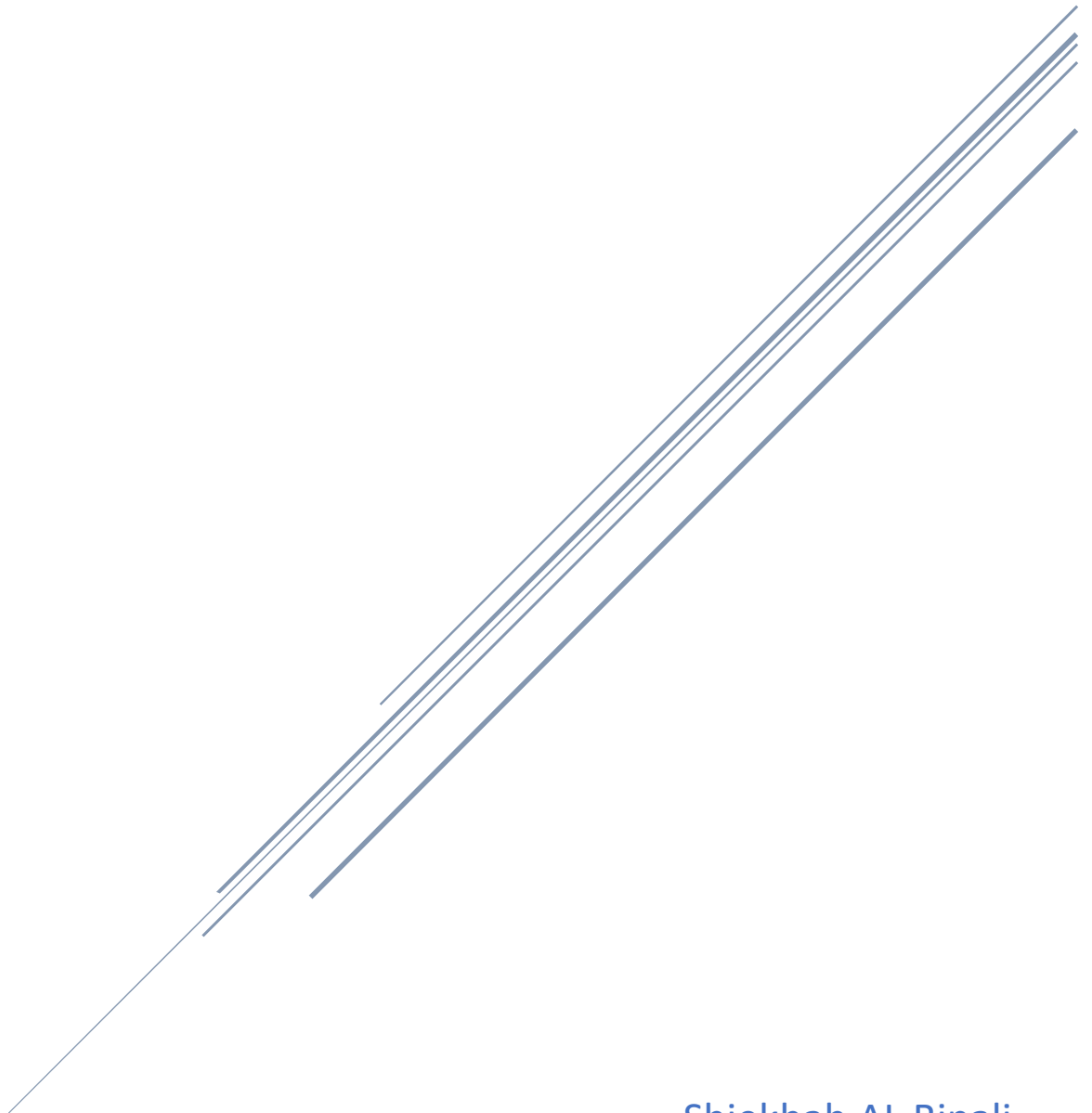# FINDING COMMON MUTATIONS FOR CANCER TREATMENT

KAIMRC training Problem Solution

Shiekhah AL Binali
Shiekhah.albinali@gmail.com

# Contents

# Figure List

# Introduction

Cancer is defined as an uncontrollable cellular proliferation; after scientists found the origin of the phenotypic expression is the DNA, research goals were focused to understand the involvement of gene expression and mutation in cancer development [1]. This disease has many contributed factors and untold clinical manifestations; thus, Genetic tools are used to evaluate cancer arise systematically; researchers have identified specific genes modification and determined how the gene changes can cause tumor growth in normal tissues. Cancer gene theory arranges a framework to interpret how hereditary and environmental factors contribute to cancers [1].

Germline mutation developed in the germ cells that construct sperm and oocytes and pass them vertically from generation to another[1]. If a change happened, a new variant is created. Small mutations can be detected using DNA sequencing while larger mutants need microscopy to be visualized[1]. Mutation can be categorized based on the amount of DNA sequence change, extensive mutations cause either insertion of a new DNA sequence or DNA sequence loss[1]. Mutation can change the gene function and affect protein structure; the change can happen directly by changing condons and change the amino acid[1]. Another way is to alter the RNA spliced way, where the splicing depends on the existence of short splice donor, splice acceptor and branch point[1].

The Concept of personalized medication has been used before to determine the suitable options for patients in case like determining the choice of antibiotics for bacterial pneumonia.[2] The goal for using personalized medicine on cancer patients is to increase efficiency of the treatment and decrease the toxicity by providing sufficient information to decide the appropriate time and type of the treatment.[2] Several techniques has been developed to define the genetic variability among patients and decode genetic signature to improve therapy.[2]

Having the data of patient and whether they respond or not to the treatment would contribute in enhancing the treatment, finding out the factors that are common in the responded group will give possible reasons about the reasons why the treatment did not benefit the non-respond group.

This report contains my attempt to solve a problem statement given by KAIMRC. The report starts with an introduction, followed by methodology section which describe the methods used to solve the problem statement, a discussion about the results is presented is also present. And finally, the conclusion section which has a brief description of the solution and the suggested future work.

# Solution

To be able to solve the problems statement, which is examining genomic data for 50 treated patients with only half of the patient who respond to the treatment. A brief description of the data is required. The data contain 50 files for 50 patients and sample information file. The computer used for this problem statement has windows 10 operating system with an i7 core processor, 16 GB Ram, 1.5 TB storage and RStudio 4.0.3 is used.

At the begging, downloading, and reading all patients' files in RStudio to understand the data. Since all data is in MAF version, there is no need to convert the data or prepare it before the analysis. All files have main parts such as Data, Variants.per.sample which helps in handling the data. Figure 1 presents the content of one of the patients' files.



Figure 1 Content of patient 0 file

After exploring the patients' files, to focus on the "Nonsynonymous Mutation" type, subsetting the "Silent" type from each of the original patients' files. Which results in having 50 altered files.



Figure 2 The data after subsetting the silent variants

After that, to find the 15 most frequent mutation, two possible solutions were applied. The first one is examining each patient file to find the top 15 mutation. And The other is determining the 15 frequent mutations for patients. For both ways. the frequent mutations were plotted and the plot for all patients is more accurate, thus the approach for each patient is discarded. Figure 3 shows the frequent mutation for a patient's data, and Figure 4 Shows the plot for frequent mutations for all patients.
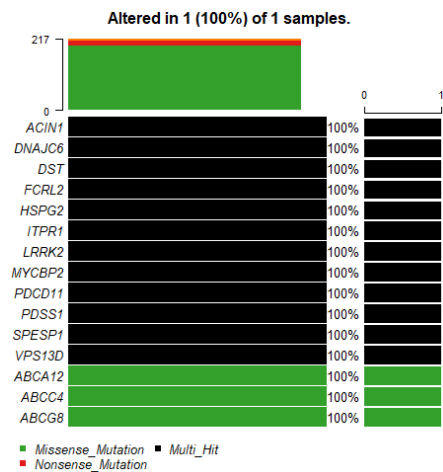


Figure 3 Frequent mutations for one patient



Figure 4 Frequent mutations for all patients

Next, the sample information is filtered based on response which leads to having two equally separate groups. Then based on the data for both groups, two MAF file are created for respond and non-respond groups. After that, mutation comparison between these two groups and the results is shown in figure 5.

Figure 5 The mutant comparison results for Respond/ Non-respond groups

After storing the comparison in a new file, presenting the comparison in a plot that shows the comparison results and the percent of cases. Figure 6 shows the comparison visualization.
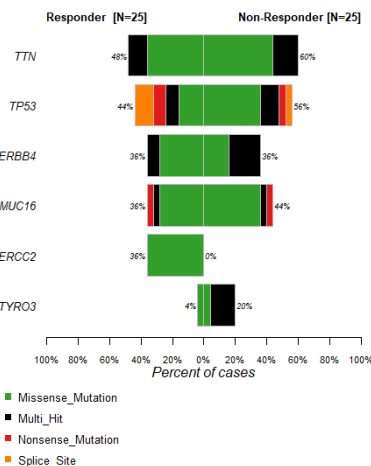


Figure 6 Mutation results in respond group compared to non-respond group

To differentiate between wild type and mutant samples, subsetting the data when the variant classification is "Missense Mutation" as mutant type and otherwise as wild type. After subsetting the samples, A plot is required to compare the results between the sample groups.

# Conclusion

This report presents an attempt to solve a problem statement to find mutations associated with treatment response for 50 patients. Genetic information is helpful in developing precision medication which handles the starting period and the appropriate type of treatment most suitable for each patient. This field has improved due to contributions of previous research and having patients' data. After solving the problem statement, The Responded group has more mutation samples than non-responded group which leads to considering this discovery as one of the potential reasons for the non-respond group. For possible future work, analyzing the non-responded group and comparing the non-respond mutant with the responded group and observe the differences, find a way to decrease the number of required mutation samples to detect the cancer subtype since it is a possible cause for non-responding to the treatment. Another possible future work is focusing on a specific variant classification and observe the results for both responding groups.

# References

[1]     F. Bunz, *Principles of cancer genetics*. 2008. doi: 10.1007/978-1-4020-6784-6.

[2]     *Bioinformatics in Cancer and Cancer Therapy*. 2009. doi: 10.1007/978-1-59745-576-3.