Name: __Talha Rafique__

Date: __9/11/22__

# TCGA Website Scavenger Hunt

## TCGA (Home Page):

The Cancer Genome Atlas (TCGA), founded in December of 2005, is a cancer genomics program hosted by the __National Cancer Institute (NCI)__ and the National Human Genome Research Institute. The publicly available data from this project includes ____genomic____, epigenomic, __transcriptomic__, and proteomic data. This data was collected from 20,000 different samples that span 33 different cancer types, including breast cancer, which we will be focusing on this semester.

## Program History:

Describe one outcome or impact of TCGA: __One of the outcomes (which we in this class beneft off of as well) is the fact TCGA has established a rich genomic data resource for the broad research community (thats us!).__

Briefly skim the "Timeline & Milestones" page. When did TCGA publish their paper on breast cancer? __October__

Because TCGA is a public dataset, and one of the first of its kind, they faced some initial concerns regarding the ethics of releasing health data to the public. Choose one of the papers in the "Ethics & Policies" section to skim. What is one way that your paper addresses these privacy concerns? __The Data Use certification agreement addresses privacy concerns by bounding access to those who will use the data. Including Research use, Requester and Approved User Responsibilities, and non-Transferability.__

## TCGA Cancers Selected for Study:

List three criteria used to select which cancers to study: Poor Prognosis, Availability of samples meeting standards for patient consent, Overall public health impact

Open the breast ductal carcinoma page and read TCGA's provided background. List one interesting fact you found: __I found it interesting how the majority of breast cancers are ductal carcinoma__

## Publications by TCGA:

TCGA published (at least) one paper on each of their studied cancer types. These papers, called marker papers, include an early analysis of the data, including any molecular characterizations that were performed. Read the abstract of the 2012 breast ductal carcinoma cancer paper. List any genes you come across (these may be good starting points for your future analyses of this cancer): __TP53, PIK3CA, GATA3, MAP3K1__

## Using TCGA:

Go to the Genomic Data Commons (GDC) Data Portal via the link on TCGA home. This portal lets you view TCGA's data in a visual way. Let's explore this website. According to the Data Portal Summary, there are __72__ projects in the GDC data portal. Now click on the "Projects" tab. Notice that not all projects in this data portal are TCGA-affiliated, though TCGA does make up __33__ of the projects included.

## Using TCGA (Continued)

Under the "Program" tab, select just TCGA studies. According to the graph at the top of the page, __TP53__ is the most mutated gene in TCGA projects, affecting approximately __35__ % of cases.

Return to the GDC Portal home page. Now click the breast image in the diagram to the right of the page. This directs you to the "Exploration" tab and automatically selects all primary sites associated with breast cancers. Now select TCGA as the program, and TCGA-BRCA as the as the project. This is the data we will be focusing on this semester.

The table on this page shows each patient along with their data. Feel free to explore the data files by clicking on any of the links provided.

Now explore the Cases, Genes, Mutations, and OncoGrid tabs above the pie charts. What is one takeaway from the plots provided here: One takeaway from the plots provided is that the survival rate linearely decreases as time elapses for those who have breast cancer

As you can see, the GDC portal provides an overwhelming amount of information. Feel free to continue to explore it on your own time!

## Discussion:

Think through the following questions, and record your answers below:

1. What is the goal of TCGA?
   By making public the information that TCGA has they have provided more accurate and descriptive datasets that can help imrprove cancer treatment as well as help identify and prevent cancers such as breast cancer.
   _____

2. What are some ways that we use TCGA's data for our own cancer research? (Think about the types of data available and brainstorm some research questions that can be proposed given that data.)
   We can use TCGA to analyze for which genes breast cancer is more prevalent and utilizing programming techniques we can determine what patters in genomic sequences result in higher rates of or higher likleyhood of breast cancer.

3. What are the benefits and drawbacks of TCGA or other large publicly available datasets?
   The primary benefit is such a large dataset that is reflecting real patients, however the data could be spotty On top of that since its being shared amongst many researches findings and results could be biased since all the information is being pulled from this one data set.