# Deployment Solution Architecture for Classification of Cancer Type with Gene Expression RNA-Seq Data

**Zhanyang Zhu [Zhanyang.Zhu@Gmail.com](mailto:Zhanyang.Zhu@Gmail.com) 12/06/2022 for UCSD Machine Learning Bootcamp Capstone Project**

## Design Requirements:

The tool will be web-based. There will be three steps:

1. Load the trained classification model (saved as pickle file). Follow the following link to see how the model is trained:

   https://github.com/ZhanyangZhuSD/UCSDMLCapstone/blob/main/GeneExpressionCancerRNA-Seq_FullData_XGBoost.ipynb

   - After loading the model, display the model accuracy and date when model is last update. New mode can be updated.
   - The feature (gene) list used in the model will be shown

2. Classification:
   In this step, user can upload their data file.
   - After loading the file, number of samples will be outputted.
   - Data will be transformed into proper format.
   - Classification result will be shown in table format.
   - User will have the ability to download the result into a csv file.

3. Distribution plot of predicted cancer types:
   This step will show the count distribution of the predicted cancer types.

## Implementation Design:

*StreamLit* (https://streamlit.io/) will be used in the implementation.

*Streamlit* is an open-source Python library that makes it easy to create and share beautiful, custom web apps for machine learning and data science.

Multiple pages will be created to match each steps described above. The distribution plot will be interactive by using *streamlit* plot function.

## Deployment Result:

The deployment architect design is described here <link>. Here are a few screen-captures of a classification run:
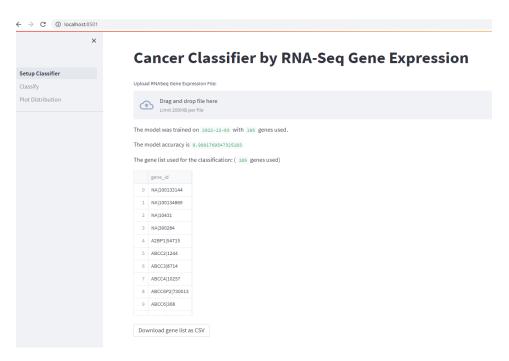
Step 0 - Start the *streamlit* server:

```
E:\UserData\Zhanyang\ML\DataSet4Projects\GeneExpressionCancerRNA-Seq\streamlit>streamlit run "00_Setup Classifier.py"

  You can now view your Streamlit app in your browser.

  Local URL: http://localhost:8501
  Network URL: http://192.168.1.113:8501
```
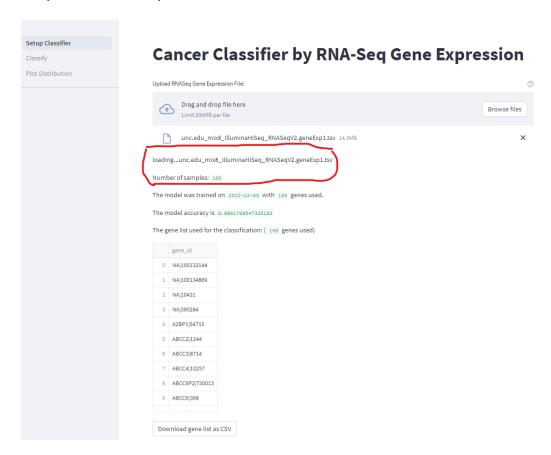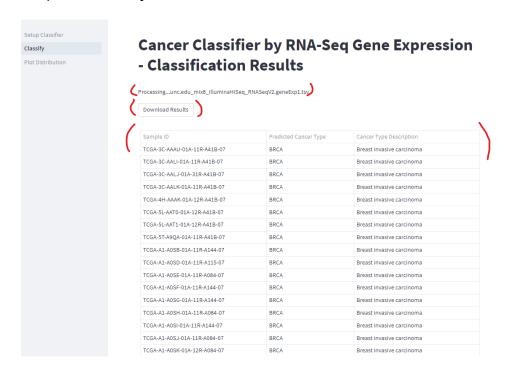
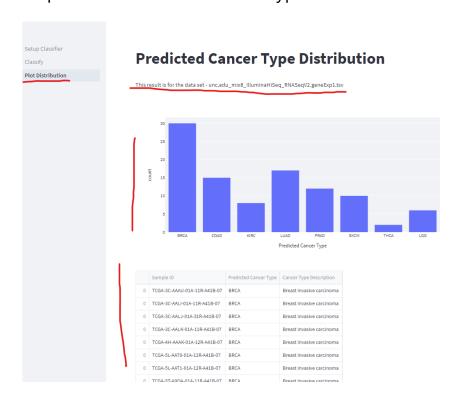Step 1 - Setup Pre-trained Classifier

## Step 2 – Select input file and load:



## Step 3 – Classify – the result can be download to a csv file

## Step 4 – Plot Predicted Cancer Type Distribution:



The graph is interactive.