

## **MLBA ASSIGNMENT-2**

### **(Group-9)**

#### **★ Readme:**

We have to give command like:

```
!python3 /content/group_9_priyankaboral_mt19127.py "kaggle_train.csv"
"kaggle_test.csv" "output.txt"
```

**A. To run the code named Priyanka\_Boral\_MT19127.py file with best MCC**

(i) first we have to go at Scripts folder in our desktop where python is installed along with two test and train folder kaggle\_train.csv, kaggle\_test.csv and Priyanka\_Boral\_MT19127\_first.py

(ii) Path where python is present:

C:\Users\akanksha\AppData\Local\Programs\Python\Python37\Scripts

(iii)      >ipython3      Priyanka\_Boral\_MT19127.py      kaggle\_train.csv  
kaggle\_test.csv latest.csv metric\_values.txt

(iv) You get the output file **latest1.csv** where all prediction are present

(v) metrics were written in **metric\_values.txt** file

**Output file submitted using this code: latest1.csv**

**B. To run the code named Priyanka\_Boral\_MT19127.py file with 2nd best MCC**

(i) first we have to go at Scripts folder in our desktop where python is installed along with two test and train folder kaggle\_train.csv, kaggle\_test.csv and Priyanka\_Boral\_MT19127\_second.py

(ii) Path where python is present:

C:\Users\akanksha\AppData\Local\Programs\Python\Python37\Scripts

(iii) **>ipython3 Priyanka\_Boral\_MT19127.py kaggle\_train.csv kaggle\_test.csv latest2.csv metric\_values.txt**

(iv) You get the output file **latest2.csv** where all prediction are present

(v) metrics were written in **metric\_values.txt** file

**Output file submitted using this code: latest2.csv**

### **C. To run the code named Priyanka\_Boral\_MT19127.py file with 3rd best MCC**

(i) first we have to go at Scripts folder in our desktop where python is installed along with two test and train folder kaggle\_train.csv, kaggle\_test.csv and Priyanka\_Boral\_MT19127\_first.py

(ii) Path where python is present:

C:\Users\akanksha\AppData\Local\Programs\Python\Python37\Scripts

(iii) **>ipython3 Priyanka\_Boral\_MT19127\_first.py kaggle\_train.csv kaggle\_test.csv latest3.csv metric\_values.txt**

(iv) You get the output file **latest3.csv** where all prediction are present

(v) metrics were written in **metric\_values.txt** file

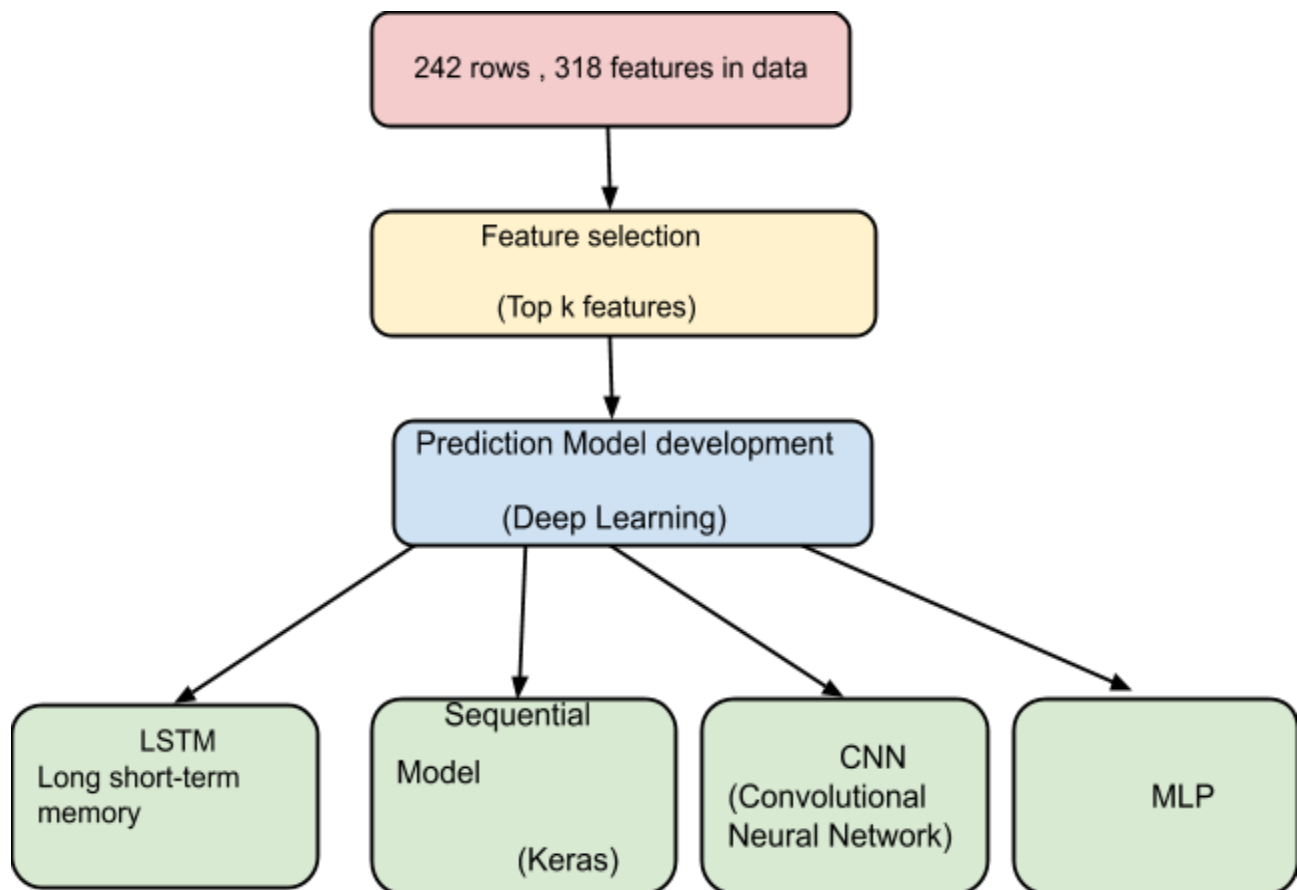
**Output file submitted using this code: latest3.csv**

Report contains information regarding all the techniques used for this assignment and analysis.

★ **Aim:** To develop deep learning models to classify high and low risk cancer patients.

★ **Preprocessing and Methodology:**

Complete description through flow chart



## Preprocessing and Feature Selection:

### 1. Feature Selection:

(i) SelectKBest RFE(Recursive feature elimination) for 100 feature extraction

SelectKBest RFE(Recursive feature elimination) for 82 feature extraction

SelectKBest RFE(Recursive feature elimination) for 64 feature extraction

(ii) `SelectKBest, f_classif sel_f = SelectKBest(score_func=f_classif, k=100)`

`SelectKBest, f_classif sel_f = SelectKBest(score_func=f_classif, k=82)`

`SelectKBest, f_classif sel_f = SelectKBest(score_func=f_classif, k=64)`

2. Various models are developed for prediction of high and low risk cancer patients.

3. Finally after obtaining predictions from the developed model, their id's and labels are stored in .csv file.

### **Models developed using deep learning techniques:**

- 1. LSTM (Long short-term memory)**
- 2. Sequential model (Keras)**
- 3. CNN (Convolutional neural network)**
- 4. MLP (Multi-layered perceptron)**

All the above models are giving different results on using various parameters and on various runs. Best validation scores which are obtained on the parameters are mentioned in the below table.

**For validation MCC (Mathews correlation coefficient) is used.**

<b>Classifiers</b>	<b>Validation Accuracy</b>	<b>Validation MCC</b>	<b>Leaderboard</b>
--------------------	----------------------------	-----------------------	--------------------

<b>1. LSTM</b> ( <code>model.add(LSTM((20 0),input_shape=(1, 318),activation='r elu',return_sequen ces=True))</code>	55.2	22	63.025
<b>2. Sequential Model</b> (without feature selection)  = <code>keras.Sequential([ keras.layers.Dense (120,activation='r elu',input_shape=( 318,))</code>	66	23	63.985
<b>3. Sequential Model</b> with feature selection= <code>keras.Sequential([ keras.layers.Dense (120,activation='r elu',input_shape=( 318,))</code>	54.6	27.32	58.644
<b>4. CNN (Convolutional Neural Network using conv 1D)</b>  Without feature selection  <code>model.add(Conv1D(6 4, kernel_size=5, activation='relu', input_shape=(318,1 )))</code>	53.06	17.6	50.120

<b>5. CNN (Convolutional Neural Network using conv 1D)</b>  With feature selection  <pre>model.add(Conv1D(6 4, kernel_size=5, activation='relu', input_shape=(318,1 )))</pre>	52.3	16.4	54.355
<b>MLPClassifier(solver='adam', alpha=1e-5,  hidden_layer_sizes=(1,318), random_state=0)</b>	51	2.4	-----

- For validation accuracy, cross validation K-fold (5-fold) has been used.
- In the above table result results of each model are shown but are checked through different parameters like by changing various layers, dropouts, batches, epochs etc.

## CONCLUSION:

- From above observation, it is concluded that in our case, **Sequential model** (`keras.Sequential([keras.layers.Dense(120,activation='relu',input_shape=(318,))` with validation MCC=35.2.
- 2nd best result is obtained, using **Sequential model** (`keras.Sequential([keras.layers.Dense(60,activation='relu',input_shape=(318,))` with validation MCC= 27.32.
- 3rd best result is obtained using CNN `model.add(Conv1D(64, kernel_size=5, activation='relu', input_shape=(318,1)))` with validation MCC=27.

Note: The results may vary on different parameters and on various runs.



**-Priyanka Boral (MT19127)-Reecha Kumari Giri(MT19134)- Akanksha Dewangan(MT19049)**