SUID: 387781563

# **Evolutionary Machine Learning – HW5**

**HW5:** Use Genetic Programming for learning a tree model for the same problem. Compare with the previous results (HW1-4).

#### Dataset:

Link:

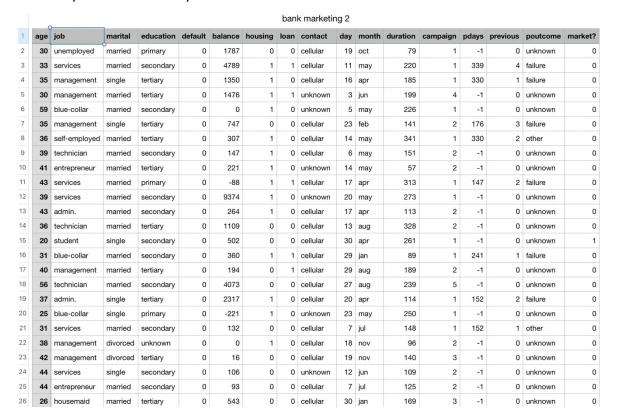
https://www.kaggle.com/chaithanya96/bankmarketing

The bank marketing dataset is the csv file with users' details used to predict the marketing decision. The decision yes or no is represented with binary numbers 1 and 0. The sample columns included are:

- 1. 'age',
- 2. 'job',
- 3. 'marital',
- 4. 'education',
- 5. 'default',
- 6. 'balance',
- 7. 'housing',
- 8. 'loan',
- 9. 'contact',
- 10. 'day',
- 11. 'month',
- 12. 'duration',
- 13. 'campaign',
- 14. 'pdays',
- 15. 'previous',
- 16. 'poutcome',
- 17. 'market?'

The screenshot of the sample dataset is attached below:

#### Chaithanya Chikkannaswamy



## **Code execution steps:**

## **Genetic Programming**

I have implemented the Genetic Programming to train a shallow feedforward neural network for a 2-class classification task using the deap library available in python. The accuracy is then compared with the accuracy calculated using Genetic Algorithm, CMA Evolutionary Strategy, particle swarm optimization and LCS implemented in previous HWs.

- 1. Deap python library is installed to implement the Genetic Programming approach to train the data.
- 2. Pandas data frame available in python is used to read the data set.
- 3. The data is split into training and testing data for the further implementation by using train\_test\_split() and specifying the test\_size as 0.2.
- 4. Once the training and testind data are ready, Genetic Programming is implemented on the dataset by specifying the parameters.
- 5. PrimitiveSetTyped() is used to add the different operators for float and Boolean data.
- 6. FitnessMAx and Individuals are determined using the base. Fitness and Tree representation of data i.e. gp. Primitive Tree.
- 7. The prediction accuracy is calculated for the model based.

#### Chaithanya Chikkannaswamy

The snapshot of data set after converting it to discreet data:

]:																		
		age	job	marital	education	default	balance	housing	loan	contact	day	month	duration	campaign	pdays	previous	poutcome	market?
	0	0	unemployed	married	primary	0	0	0	0	cellular	1	oct	0	1	-1	0	unknown	C
	1	0	services	married	secondary	0	0	1	1	cellular	0	may	0	1	339	4	failure	(
	2	0	management	single	tertiary	0	0	1	0	cellular	1	apr	0	1	330	1	failure	(
	3	0	management	married	tertiary	0	0	1	1	unknown	0	jun	0	4	-1	0	unknown	(
	4	1	blue-collar	married	secondary	0	0	1	0	unknown	0	may	0	1	-1	0	unknown	C
	4516	0	services	married	secondary	0	0	1	0	cellular	2	jul	0	5	-1	0	unknown	(
	4517	1	self-employed	married	tertiary	1	0	1	1	unknown	0	may	0	1	-1	0	unknown	(
	4518	1	technician	married	secondary	0	0	0	0	cellular	1	aug	0	11	-1	0	unknown	(
	4519	0	blue-collar	married	secondary	0	0	0	0	cellular	0	feb	0	4	211	3	other	(
	4520	0	entrepreneur	single	tertiary	0	0	1	1	cellular	0	apr	0	2	249	7	other	(
4	521 r	ows ×	17 columns															

The snapshot of data set after dropping the non-usable data and converting the required column for prediction to bool data type:

```
print(df.columns)
df[10] = df[10].astype('int').astype('bool')
print(df.head(10))
```

```
Int64Index([0, 1, 2, 3, 4, 5, 6, 7, 8, 9, 10], dtype='int64')
   0
            2
               3
                  4
                      5
                           6
                             7
                                  8
                                     9
                                            10
                                         False
1
   30 0
          1787
                0
                  0
                      19
                           79 1
                                  -1 0
   33 0
          4789
                                         False
2
                1
                   1
                      11
                          220 1
                                 339
                                      4
3
   35 0
          1350 1
                   0
                      16
                          185
                              1
                                 330 1
                                         False
4
   30 0
          1476
               1
                   1
                       3
                          199
                              4
                                  -1
                                         False
                       5
5
   59 0
             0
               1
                   0
                          226
                              1
                                  -1
                                      0
                                         False
   35 0
           747
                   0
                          141
                                         False
6
                0
                      23
                              2
                                 176
7
   36 0
           307
                   0
                              1
                                      2
                                         False
               1
                      14
                          341
                                 330
8
   39
       0
           147 1
                   0
                          151
                              2
                                         False
                      6
                                  -1
9
   41 0
           221
                   0
                              2
                                  -1 0
                                         False
                      14
                           57
                                      2
10
   43
           -88
                         313 1
                                         False
       0
                   1
                      17
                                 147
```

+ Code

+ Markdown

The snapshot of splitting the data set into train and test data:

```
train, test = train_test_split(df, test_size=0.2)
print(train.shape)
print(train.head(10))
train_data = train.values.tolist()
print(test.head(10))
```

```
(3616, 11)
      0 1
                2
                   3
                      4
                           5
                                6
                                   7
                                        8
                                           9
                                                  10
4256
          0
                       0
                          23
                               698
                                   2
      42
              1331
                    0
                                        -1
                                           0
                                               False
                                               False
4489
              315
                               130
      31
          0
                    0
                       0
                          30
                                        2
                                            1
                                    1
3895
      38
          0
                 0
                    1
                       0
                          17
                               213
                                    2
                                            0
                                               False
                                        -1
2834
      38
          0
                62
                    0
                       0
                          19
                               212
                                    1
                                        -1
                                            0
                                               False
                                    5
860
      51
          0
                 4
                    1
                       1
                           20
                                74
                                        -1
                                            0
                                               False
1890
                          8
                               376
                                    2
                                               True
      42
          0
              1205
                    0
                       0
                                        -1
                                            0
3909
               167
                                    1
      36
          0
                    1
                       0
                           28
                                57
                                        -1
                                            0
                                               False
841
      28
          1
              -298
                       0
                            3
                               559
                                    7
                                        -1
                                            0
                                               False
                    1
3664
      37
          0
               196
                    1
                       0
                            5
                                66
                                    3
                                       -1
                                            0
                                               False
1404
      29
               912
                               785
                                    1
                                               False
          0
                    1
                       0
                          13
                                        -1
                                            0
                                6
                                             9
                2
                           5
                                    7
                                                    10
      0 1
                   3
                      4
                                          8
509
      41
          0
              428
                    1
                       0
                          12
                                92
                                     1
                                          -1
                                              0
                                                 False
              -278
1689
      38
          0
                    1
                       0
                          28
                               143
                                     2
                                              0
                                                 False
                                          -1
3856
      47
          0
                -9
                    0
                       0
                               457
                                     2
                                                 False
                          14
                                          -1
                                              0
1636
      33
          0
                    1
                       0
                           14
                                72
                                     2
                                             1
                                                 False
               687
                                        370
628
      57
          0
              374
                    1
                       1
                           16
                               236
                                     1
                                          -1
                                              0
                                                 False
4468
                                                 False
      43
          0
              1577
                    1
                       0
                           19
                                87
                                     1
                                          -1
                                              0
84
      52
          0
              657
                           7
                               398
                                     2
                                         460
                                              2
                                                 True
                    0
                       0
3533
      42
          0
              3620
                    1
                       0
                          27
                                22
                                    16
                                          -1
                                              0
                                                 False
2731
      52
          0
                            6
                                                 False
                 9
                    0
                       1
                                44
                                     1
                                          -1
                                              0
3337
      37
          0
               215 1
                       0
                            6
                                61
                                              0
                                     1
                                          -1
                                                 False
```

## Chaithanya Chikkannaswamy

The snapshot of testing and training results of the data:

```
train_res = train[10]
test_res = test[10]
print(train_res)
```

```
4256
       False
4489
       False
3895
      False
2834
       False
860
       False
2284
       False
3116
      False
2272 False
2221 True
716
       False
Name: 10, Length: 3616, dtype: bool
```

The snapshot of the accuracy rate:

```
[72]: accuracy_score(test_res, predict(pop[0]))
```

Out[72] 0.7933701657458564

## **Conclusion:**

- 1. The **Genetic Programming** resulted in **79%** accuracy.
- 2. The **Learning Classifier System** resulted in the accuracy of **99%**, **88%**, **99%** and **99%** based on the columns.
- 3. The **Particle Swarm Optimization** resulted in the accuracy of 0.88 i.e. **88**%.
- 4. The **CMA-ES** resulted in 0.715 test accuracy.
- 5. The **Genetic Neural Network** resulted in a test accuracy of 0.89 and the **Sequential Neural Network** resulted in a test accuracy of 0.88.

## References

- 1. https://github.com/trevorstephens/gplearn/blob/master/gplearn/fitness.py
- 2. <a href="https://github.com/sighmin/gpstocks">https://github.com/sighmin/gpstocks</a>
- 3. <a href="https://docs.google.com/file/d/0B4JHGiC-rWKmbm9MYWx2b1VNV2c/edit">https://docs.google.com/file/d/0B4JHGiC-rWKmbm9MYWx2b1VNV2c/edit</a>
- 4. https://docs.google.com/file/d/0B4JHGiC-rWKmWThhbWtwVngxbnM/edit
- 5. Discussed with classmates.