

Model 3. LR-RS-BT-0.25-1

Trains model 3 based on:

- Regular Linear Regression
- RankStandardised Data
- Both Teams data
- FS_Val 0.25
- FS_Rule 1

Author: Lang (Ron) Chen 2021.12-2022.1

0. Import Libraries

```
In [1]: import pandas as pd
import os
from sklearn.model_selection import train_test_split
from sklearn import linear_model
from sklearn.metrics import mean_squared_error, r2_score
import numpy as np
import pickle

from BrownlowPredictorTools.predict import predict
from BrownlowPredictorTools.test import test
from BrownlowPredictorTools.return_tp import return_tp
from BrownlowPredictorTools.wholeseason import wholeseason
from BrownlowPredictorTools.feature_selection2 import feature_selection2
```

```
In [2]: choice = 'RankStandardisedData'
```

```
In [3]: filelist = os.listdir(f'./Data/{choice}')
filelist.sort()
filelist = filelist[1:]
# Remove the first file (an ipynb checkpoint file)
```

1. Feature Selection

```
In [4]: # As we need to perform tests to evaluate this final model, still need to use do this

# Gets list of emperical test games (full 2021 season)
final_test_games = [file for file in filelist if '2021' in file]

# Gathers full games list (except 2021) and performs a single Train-Test Split (note differ
test_train_games = [file for file in filelist if '2021' not in file]
train_games, test_games = train_test_split(test_train_games, train_size = 0.8, test_size =
```

```
In [5]: # Read in pre-prepared sample data of trained data only
# (the same rows as if we used concatenated all the data from the train_games list)
train_data = pd.read_csv('./Models/TrainingData/M3_Data.csv')
```

```
In [6]: # Select Columns of Both Teams Stats only
```

```
cols = [col for col in train_data.columns if ('BTRS' in col or 'Winloss' in col)]

# Select Columns with correlation higher than 0.25 only
corr = dict()
for col in cols:
    corr[col] = train_data[[col, 'Brownlow Votes']].corr(method = 'pearson').loc[col]['Brownlow Votes']

corr = list(corr.items())

selected_features = [col[0] for col in corr if col[1] > 0.25]
selected_features
```

Out[6]:

```
['Kicks BTRS',
 'Disposals BTRS',
 'Contested Possessions BTRS',
 'Effective Disposals BTRS',
 'Score Involvements BTRS',
 'Metres Gained BTRS',
 'Behind Assists BTRS']
```

2.Trains Models

No need to run feature_selection2 because this utilises rule 1

In [7]:

```
# Trains LR model
traindata_x_1 = train_data[selected_features]
traindata_x_1.index = range(0,len(traindata_x_1))
traindata_y_1 = train_data['Brownlow Votes']
traindata_y_1.index = range(0,len(traindata_y_1))

lm_1 = linear_model.LinearRegression()
traindata_x_1 = traindata_x_1.replace((np.inf, -np.inf, np.nan), 0).reset_index(drop=True)
model_1 = lm_1.fit(traindata_x_1, traindata_y_1)
```

In [8]:

```
# Get predictions and observations
predictions_1, testdata_y_1 = predict(test_games, lm_1, selected_features, choice)
```

In [9]:

```
# Get True Positive/True Negative results
result1_1, result2_1 = test(predictions_1, testdata_y_1, 4)
```

In [10]:

```
# TP/TN based on what was predicted
result1_1
```

Out[10]:

```
[0.9617637661092768,
 0.015230588987112579,
 0.013526467142400682,
 0.009479177761209927],
[0.6899563318777293,
 0.13100436681222707,
 0.10480349344978165,
 0.07423580786026202],
[0.5545851528384279,
 0.12663755458515283,
 0.1222707423580786,
 0.1965065502183406],
[0.3231441048034934,
 0.11790393013100436,
 0.2183406113537118,
 0.3406113537117904]]
```

```
In [11]: # TP/TN based on what was observed
result2_1
```

```
Out[11]: [[0.9617637661092768,
0.01682820321652998,
0.013526467142400682,
0.007881563531792523],
[0.6244541484716157,
0.13100436681222707,
0.12663755458515283,
0.11790393013100436],
[0.5545851528384279,
0.10480349344978165,
0.1222707423580786,
0.2183406113537118],
[0.388646288209607,
0.07423580786026202,
0.1965065502183406,
0.3406113537117904]]
```

```
In [12]: # Only the True Positive Values
return_tp(result1_1)
```

```
Out[12]: (0.9617637661092768,
0.13100436681222707,
0.1222707423580786,
0.3406113537117904)
```

3. Summary Observations

1. Emperical Experiment

```
In [13]: # Runs the season 2021 data onto predictor and gets top players
leaderboard1 = wholeseason(final_test_games, lm_1, selected_features, choice)
```

```
In [14]: leaderboard1[0:15]
```

```
Out[14]: [('Clayton Oliver', 33),
('Jack Steele', 33),
('Oliver Wines', 32),
('Darcy Parish', 32),
('Jackson Macrae', 31),
('Jarryd Lyons', 30),
('Sam Walsh', 28),
('Christian Petracca', 27),
('Touk Miller', 27),
('Travis Boak', 23),
('Rory Laird', 23),
('Marcus Bontempelli', 23),
('Cameron Guthrie', 21),
('Dayne Zorko', 21),
('Zachary Merrett', 19)]
```

1. Predictor's r scores

```
In [15]: print(lm_1.score(traindata_x_1, traindata_y_1))

0.1367427545098272
```

4. Picklising

In [16]:

```
with open('./Models/M3.pickle', 'wb') as f:  
    pickle.dump([lm_1, selected_features, choice], f)
```