# FOOTBALL (SOCCER) MATCH RESULT PROJECTIONS

Nikolai Stambler-Tennant

Brown University – DSI

October 19th, 2022

https://github.com/Niktennant/Premier-League-Result-Projections.git

| The problem? | Importance? | Type? | Data Origin? |
|---|---|---|---|
| • Predicting football matches outcomes - specifically in the Premier League. | • Billion $ industry, sports betting, personal interest | • Multi-class Classification (H : 1.0, D : 0.0, A : -1.0) | • Multiple sources<br>• Fantasy Premier League API<br>• FIFA 22 & FIFA 23<br>• PL Match Records |

# INTRODUCTION
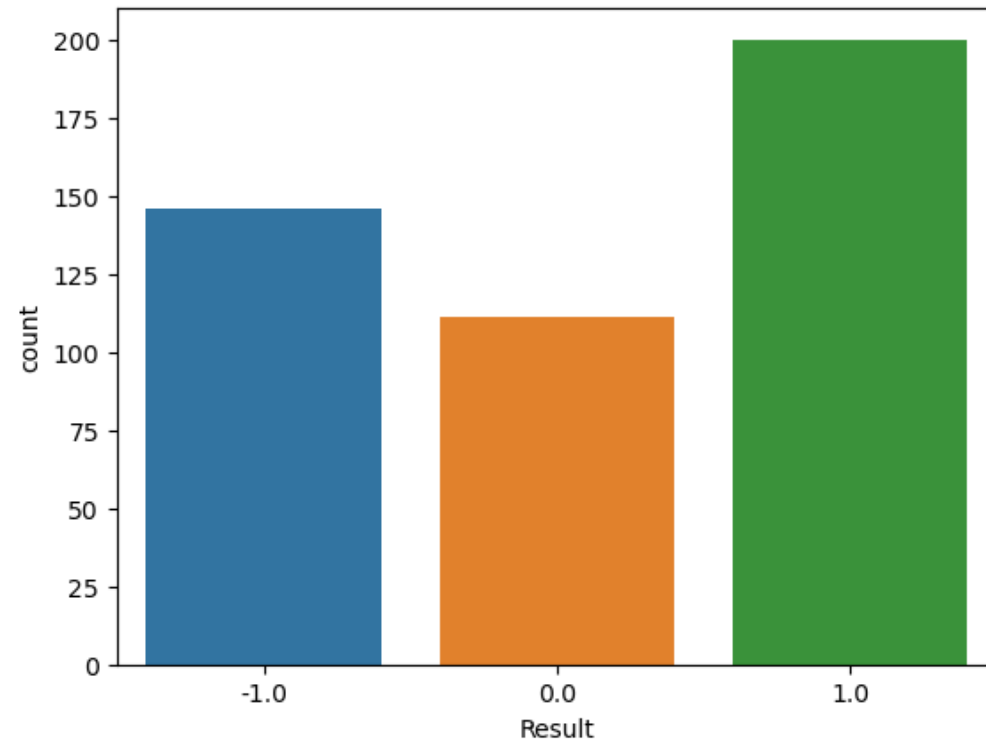
MY DATA

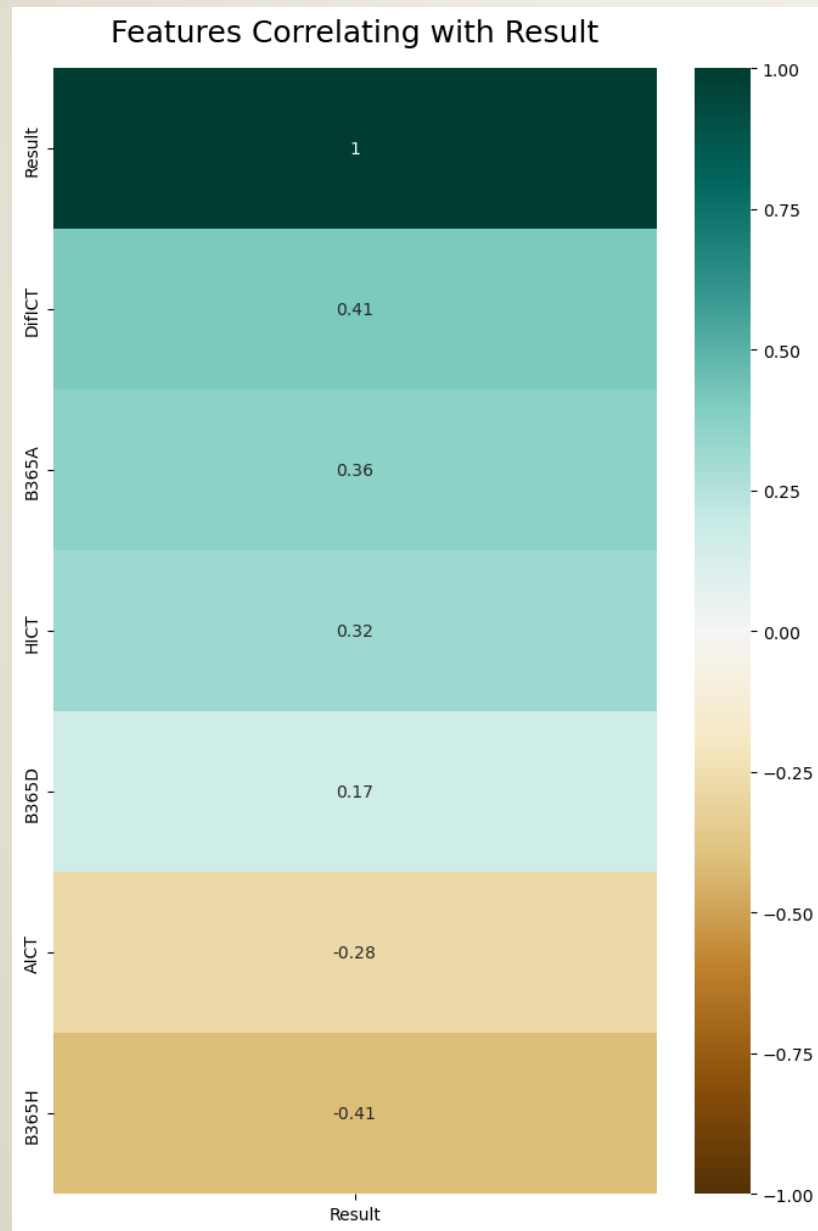| | Referee | B365H | B365D | B365A | HICT | AICT | DifICT | Result |
|---|---|---|---|---|---|---|---|---|
| 0 | M Oliver | 4.00 | 3.40 | 1.95 | 73000.0 | 79000.0 | -6000.0 | 1.0 |
| 1 | A Madley | 1.90 | 3.50 | 4.00 | 79000.0 | 76000.0 | 3000.0 | 1.0 |
| 2 | D Coote | 3.10 | 3.10 | 2.45 | 76000.0 | 76000.0 | 0.0 | -1.0 |
| 3 | J Moss | 1.25 | 5.75 | 13.00 | 83000.0 | 76000.0 | 7000.0 | 1.0 |
| 4 | M Dean | 3.10 | 3.20 | 2.37 | 75000.0 | 78000.0 | -3000.0 | 1.0 |
| ... | ... | ... | ... | ... | ... | ... | ... | ... |
| 452 | M Oliver | 2.50 | 3.60 | 2.62 | 80323.8 | 85198.8 | -4875.0 | 1.0 |
| 453 | P Tierney | 2.05 | 3.50 | 3.60 | 76195.1 | 76165.2 | 29.9 | 1.0 |
| 454 | C Kavanagh | 1.72 | 3.80 | 4.75 | 80176.2 | 76193.7 | 3982.5 | 1.0 |
| 455 | D Coote | 3.60 | 3.50 | 2.05 | 78201.2 | 82202.0 | -4000.8 | -1.0 |
| 456 | A Taylor | 3.30 | 3.40 | 2.20 | 76196.1 | 79181.8 | -2985.7 | 0.0 |

457 rows × 8 columns

# EDA - TARGET VARIABLE (RESULT)

- Home Wins – 200 (43.8%)

- Away Wins – 146 (31.9&)

- Draws – 111 (24.3%)



| | count | mean | std | min | 25% | 50% | 75% | max |
|---|---|---|---|---|---|---|---|---|
| **Result** | 457.0 | 0.118162 | 0.863006 | -1.0 | -1.0 | 0.0 | 1.0 | 1.0 |

Features Correlating with Result

EDA- PT. 2

# SPLITTING/PREPROCESSING

| SPLIT | PREPROCESSES | (FEATURES, DATA POINTS) | MISSING VALUES? |
|---|---|---|---|
| • Time Series (no group structure)<br><br>• Lagged data (5 matches)<br><br>• 60/20/20 split with df sorted by date (ascending order)<br>  • 271/91/90 | • StandardScaler: continuous features w/ normal/tailed distribution<br><br>• One-hot: 'Referee' - unordered categorical data<br><br>• Min/MaxScaler: HICT/AICT (73000- 86309) | • Before preprocess (457, 8)<br><br>• After lag (452, 38)<br><br>• After preprocessing (452, 58) | • Nan |