

Department of AI, Faculty of ICT, University of Malta

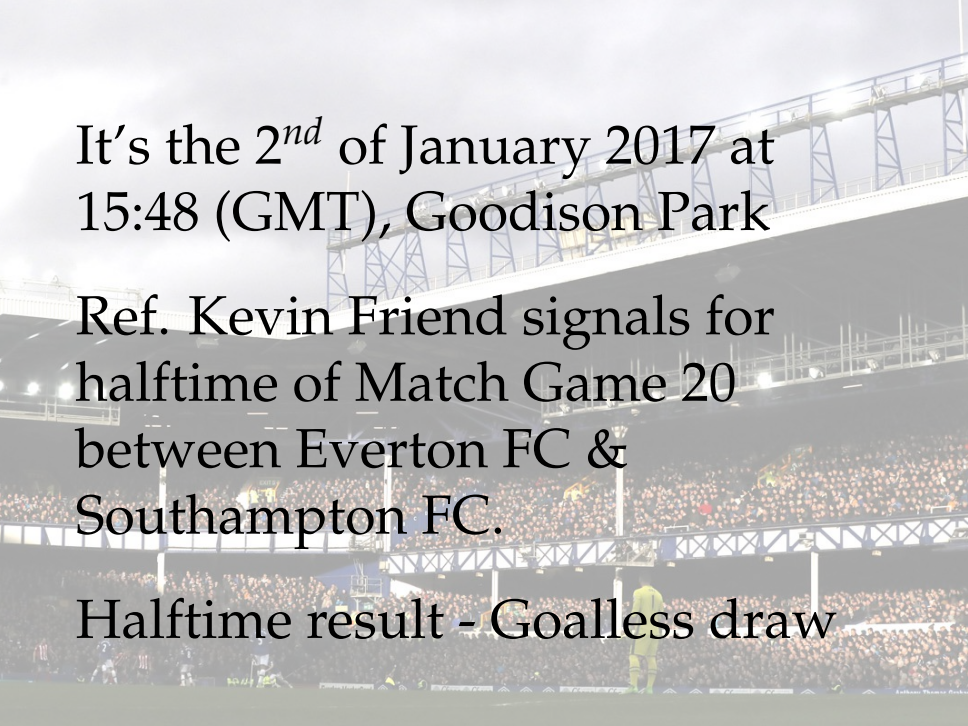
---

# Predictive Analysis of Football Matches using In-play Data

Matthew J. Zammit

matthew.zammit.09@um.edu.mt

September 11, 2018



It's the 2<sup>nd</sup> of January 2017 at  
15:48 (GMT), Goodison Park

Ref. Kevin Friend signals for  
halftime of Match Game 20  
between Everton FC &  
Southampton FC.

Halftime result - Goalless draw



# Match day - Halftime Analysis (1)

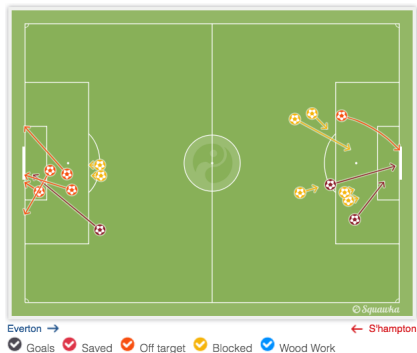


Figure: Shots per team

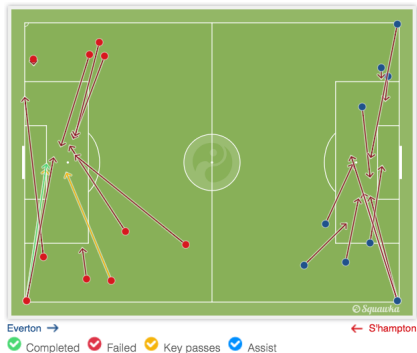


Figure: Crosses per team



# Match day - Halftime Analysis (2)

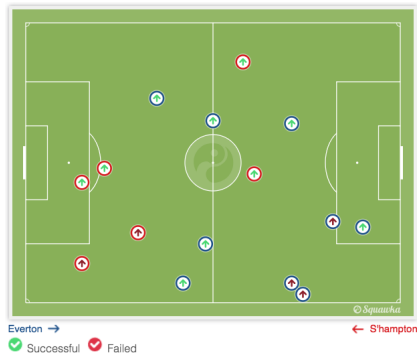


Figure: Dribbles per team

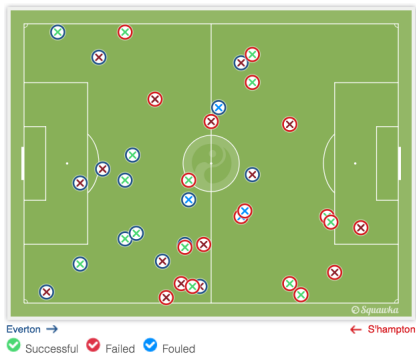


Figure: Tackles per team



# Match day - Halftime Analysis (3)

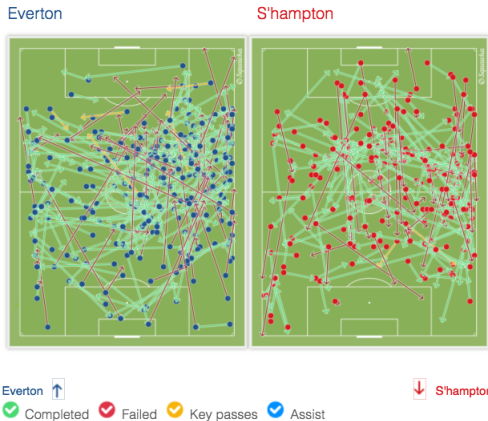


Figure: Passes per team



How do you think is the match going to end?

Match result market at halftime;  
Everton to win, odds at 2.86  
(0.35). Southampton to win, odds at  
4.00 (0.25). Ends in a draw, odds  
at 2.50 (0.40).



- Rise in popularity of betting exchanges through the Internet
- Prediction Markets, have been found to be accurate
- Sports data is recently being captured at precise and granular levels than ever before



# Why Machine Learning?

- Multitude of complex variables associated with a football match
- Difficult for humans to think in terms of probability and to react to market changes
- Emotions might hinder the performance of humans to make rational decisions





# Aims and Objectives

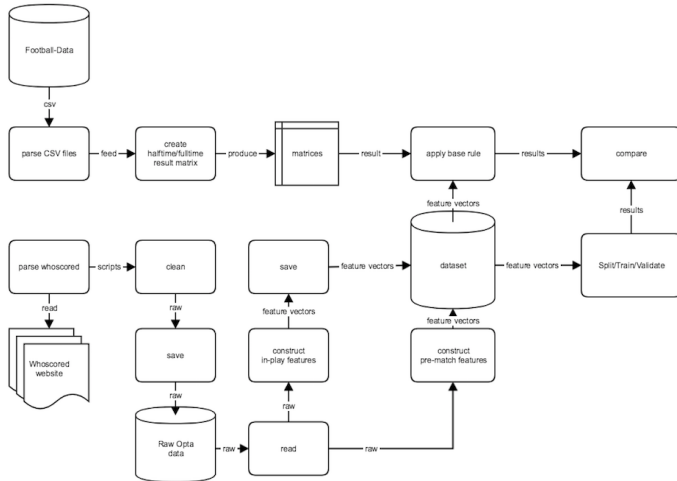
- Predict the fulltime result (H/D/A) of matches drawn at halftime using in-play match data
- Investigate whether using Feature Selection (FS) by a Genetic Algorithm (GA) process would remove certain predictors and increase classification accuracy
- Test if the addition of pre-match data to the in-play game statistics would improve accuracy rate
- Compare the probabilistic classification of the classifier with that of the implied probability from the betting exchange market



- No publicly accessible datasets available
- No previously conducted studies using same data to compare with
- Most similar study was carried out using a Case Based Reasoning approach on the over/under 2.5 goals market



# Dissertation Schematic





# Halftime/Fulltime Result Base Rule

- Data was retrieved from *Fooball-Data*
- We investigated the transition probability of the halftime result to that of the fulltime result of each match
- Consists of 77,553 match instances spanning several major and minor leagues across Europe over multiple seasons
- We found that for a high percentage of matches the fulltime result remained the same as that of the halftime

## Base rule

$$BR_{ftr}(R_{htr}) = R_{htr} \quad (1)$$



- *Football-Data* data was not granular enough for the study
- Data was parsed from *Whoscored* website and had to be engineered using conditional rules
- Main benefits why this site was chosen:
  - 1 Data recorded at a play-by-play rate
  - 2 Actions are labeled with a type, the x and y coordinates of the ball
  - 3 Opta as the source
  - 4 Continually being updated with data of major European competitions. Most importantly, English Premier League, Italian Serie A, Spanish La Liga, French Ligue and the German Bundesliga from 2009/10 to present



Dataset URL: <https://bit.ly/2QdlCs6>

- SHOT\_TOTAL
- SHOT\_ON\_GOAL
- ASSIST\_SHOT
- ASSIST\_INTENT'L
- ASSIST\_INTENT'L\_GOAL
- PASS\_TOTAL/SUCCESS
- PASS\_LONG
- PASS\_FORW/BACK
- PASS\_TRGT\_FINAL\_TRD
- PASS\_TRGT\_MID\_TRD
- PASS\_TRGT\_DEF\_TRD
- CORNER\_FAVOUR
- FOUL\_RECEIVED
- CROSS\_FAV\_TOTAL
- CROSS\_FAV\_SUCCESS
- OFFSIDE\_COMMITTED
- POSSESSION\_TOTAL
- POSSESSION\_ATT
- POSSESSION\_DEF
- INTERCEPTION
- CARD\_YELLOW/RED
- TACKLE\_TOT/SUCCESS
- DRIBBLE\_TOT/SUCCESS



# Feature and Target Vector

- The feature vector constitutes of the difference in the halftime statistics between the home and away team
- A positive value for a particular feature means that the home team had accumulated more of that statistic till the halftime than the away team
- The target vector consists of only one element for each feature vector. The value could be from the set  $\{0,1,2\}$ , where the elements represents home, draw and away win respectively



# Feature and Target Vector Example

match	target (FTR)	shotTotalDiff	shotOnGoalDiff	passTotalDiff	passLongDiff	passSuccessDiff	passBackwardDiff	passTargetFinalThirdDiff	...	tackleTotalDiff	tackleSuccessDiff	dribbleTotalDiff
Arsenal v Aston Villa	2	2	2	115	-10	123	44	22 ...		-9	-3	8
Arsenal v Cardiff	0	5	1	132	-4	138	49	64 ...		-9	-5	8
Arsenal v Chelsea	1	-4	-3	89	4	86	46	12 ...		-2	-5	1
Arsenal v Crystal Palace	0	4	2	297	-2	302	120	139 ...		0	-3	-6
Arsenal v Everton	1	-2	1	-123	-20	-131	-69	-2 ...		8	5	-4
Arsenal v Fulham	0	7	3	42	-19	32	12	95 ...		-5	-1	3
Arsenal v Hull	0	11	5	232	-20	243	86	123 ...		-7	-4	8
Arsenal v Liverpool	0	2	0	118	1	102	58	76 ...		12	10	-7
Arsenal v Man City	1	-6	-3	-80	0	-75	-37	-25 ...		-3	-1	7
Arsenal v Man Utd	1	6	4	33	1	27	1	52 ...		-6	-4	8
Arsenal v Newcastle	0	8	5	130	-1	133	48	50 ...		-1	-3	-1
Arsenal v Norwich	0	3	0	22	-8	27	16	23 ...		5	4	-3
Arsenal v Southampton	0	-1	-2	7	-3	8	-4	19 ...		-3	-4	6
Arsenal v Stoke	0	7	4	179	-9	184	94	45 ...		-7	-2	13
Arsenal v Sunderland	0	9	5	280	0	279	131	131 ...		-8	-4	9
Arsenal v Swansea	1	7	4	105	-4	110	52	143 ...		-8	-7	7
Arsenal v Tottenham	0	3	-2	-27	-14	-40	-18	-1 ...		12	5	-8
Arsenal v West Brom	0	4	4	182	15	176	87	55 ...		0	-1	-3
Arsenal v West Ham	0	1	2	147	2	146	55	50 ...		-6	-4	





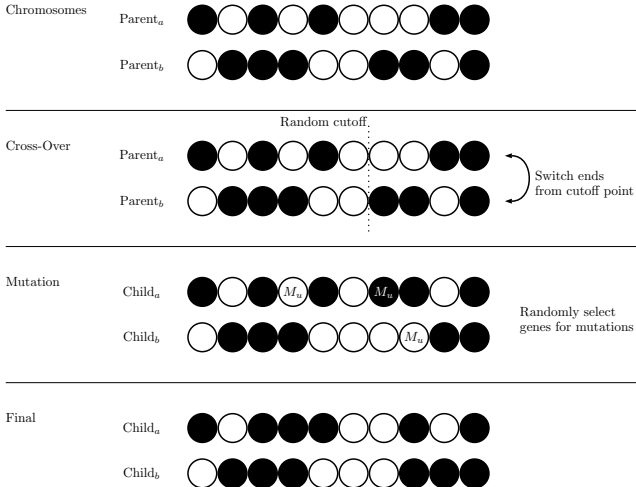
- Instances from the English 2015/16 season were used as a sample for initial experimentation
- Features were iteratively being added to the feature space depending on the accuracy and based on our football intuition
- Machine Learning Algorithms used; Neural Nets (NN), Naive Bayes (NB), Decision Tree (DT), Random Forest (RF)
- Dataset was normalized for algorithms which trained faster and perform better with scaled data
- Random Forest was found to be consistently accurate across all the tests



- Using only the Random Forest algorithm with custom GA separately for each league
- Investigate classification performance and predictors chosen with default parameter settings and with model tuning
  - Nested Cross Validation (CV)
  - Grid search used for parameter tuning
  - Fitness function promotes fewer predictors
  - Growth function for mutation rate
  - Stopping criteria: score of latest generation subtracted by the mean of the scores from the previous ten generations greater than a threshold (0.1)



# Genetic Algorithm - Feature selection

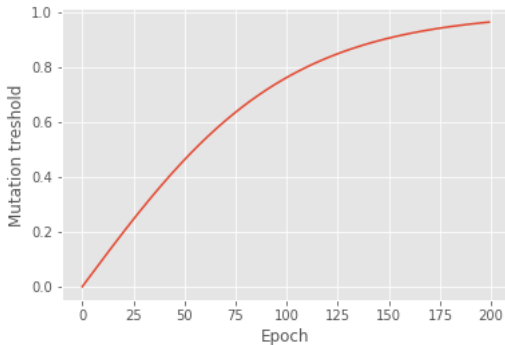




## Mutation Threshold

$$mt = \tanh\left(\frac{2i}{n}\right) \quad (2)$$

Where  $i$  is the current epoch and  $n$  is the maximum number of epochs.







## Accuracy

$$Accuracy = \frac{TP + TN}{TP + TN + FP + FN} \quad (3)$$

- TP - True Positive
- TN - True Negative
- FP - False Positive
- FN - False Negative



- Prematch data added to contextualise the instances
- Prematch feature set includes simple attributes:
  - Goals Scored
  - Goals Conceded
  - Points
- and computed ones
  - Team Form based on the teams' latest performances
  - Attacking Strength
  - Defensive Strength
- Inner partitioning loop customised to train/test on a seasonal basis because of temporal data
- Same as with in-play data, the feature vector consisted of the difference between the home and away team pre-match statistics



where,  $t$  denotes the team for which the strength is being calculated,  $n$  represents the match game and  $m$  describes the total number of teams.

$S$  and  $C$  represent the goals scored and goals conceded matrices, respectively.

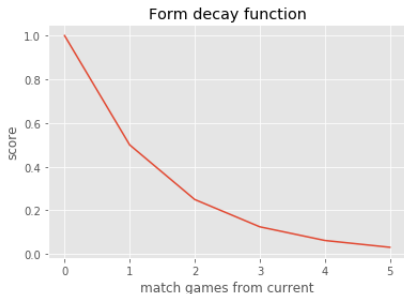
## Attack Strength

$$AttackStr(t, n, m) = \frac{\frac{1}{n} \sum_{i=1}^n S_{it}}{\frac{1}{mn} \sum_{j=1}^m \sum_{i=1}^n S_{ij}} \quad (4)$$

## Defence Strength

$$DefenceStr(t, n, m) = \frac{\frac{1}{n} \sum_{i=1}^n C_{it}}{\frac{1}{mn} \sum_{j=1}^m \sum_{i=1}^n C_{ij}} \quad (5)$$





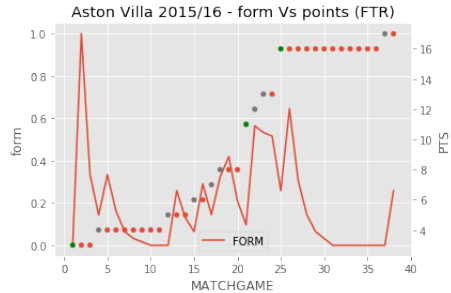
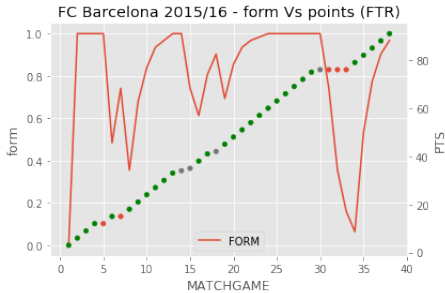
## Form

$$Form(t, j) = \frac{\sum_{i=0}^5 (\frac{1}{2})^i (y_{j-i,t})}{\sum_{i=0}^5 (\frac{1}{2})^i} \quad (6)$$

Importance is given to the result of the previous games depending on how recent they have been played by assigning them different weights.



# Pre-match Form - Examples





# Custom Inner Loop for Prematch data

Data Set	Split				
	1	2	3	4	
2011/12	$t_1$	$t_2$	$t_3$	$t_4$	
2012/13	$v_1$				
2013/14		$v_2$			
2014/15			$v_3$		
2015/16					$v_4$
Score	$S_1$	$S_2$	$S_3$	$S_4$	$\bar{S}$



## Brier Score Function

$$\frac{1}{n} \sum_{i=1}^n \sum_{j=1}^r (p_{ij} - o_{ij})^2 \quad (7)$$

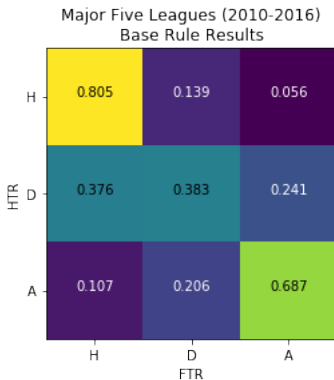
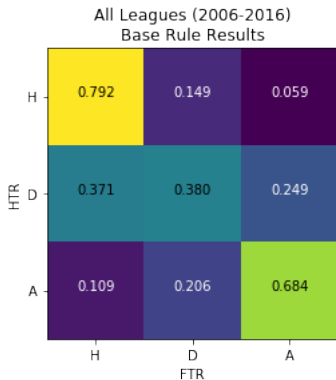
Where,  $n$  is the total number of instances and  $r$  is the number of possible outcomes (three is our case).

$p_{i,j}$  is the probability of the  $j^{th}$  outcome for the  $i^{th}$  instance from the model. For example, when  $i = 1$  the probability vector is  $[0.7, 0.2, 0.1]$ .

$o_{i,j}$  is the actual probability of the  $j^{th}$  outcome for the  $i^{th}$  instance after its occurrence. For example, for the same instance  $i = 1$ , the actual result was  $[1.0, 0.0, 0.0]$ .



# Base Rule Matrices Results

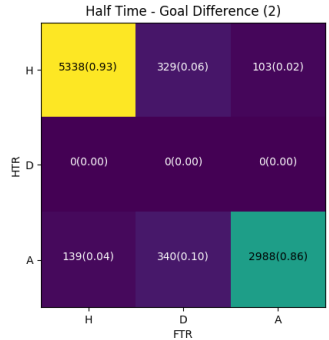
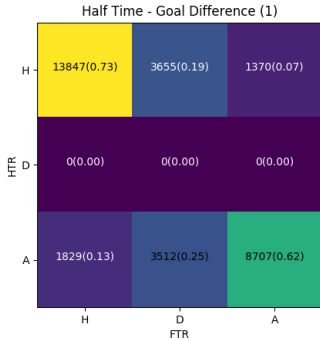


HTR - Halftime Result, FTR - Fulltime Result.

A darker color represents a lower value.

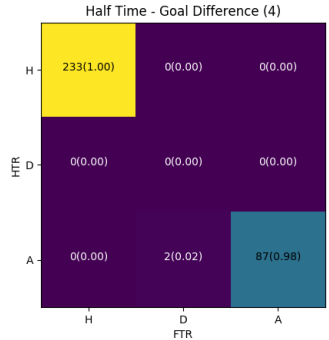
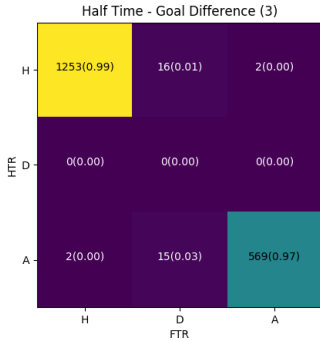


# BR Matrices Results by Goal Difference (1)



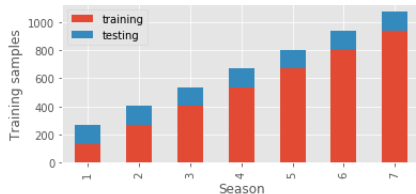
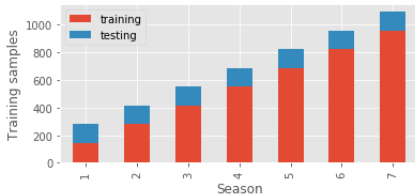
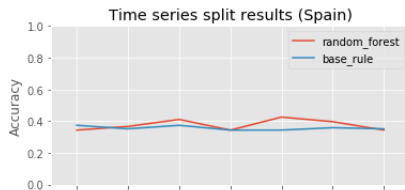
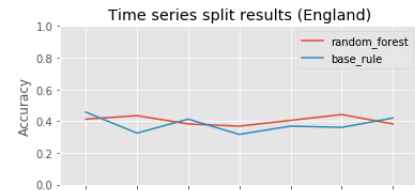


# BR Matrices Results by Goal Difference (2)





# Time Series Results

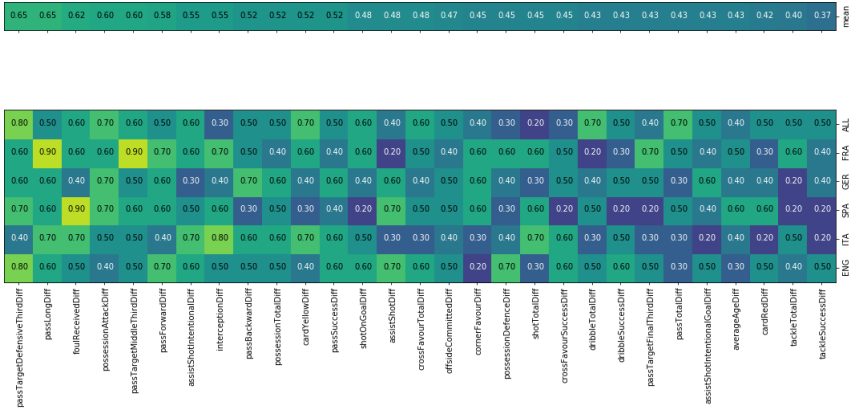


One tailed paired t-test showed that the accuracy of the time series random forest was not significantly different from that of the base-rule with t-statistic of 1.33 and p-value of 0.19.



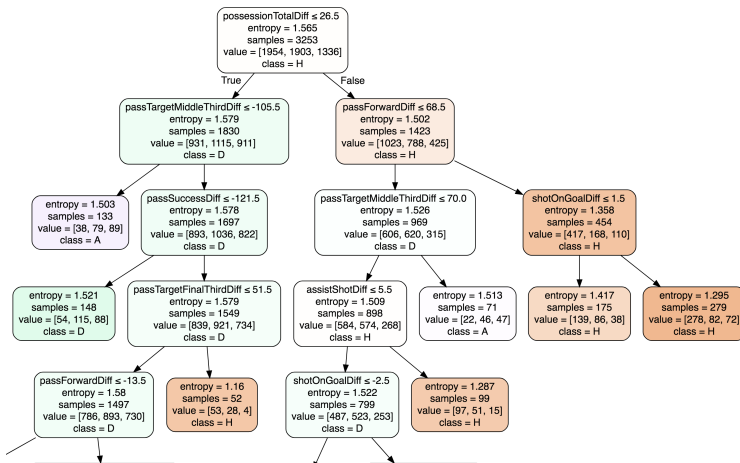


# Predictors Chosen from GA and Tuning





# Decision Tree Example





# Results Summary

Leagues	$RF_{HT+PM}$	$RF_{GA(T)}$	$RF_{GA}$	$RF_{TS}$	$RF_{10CV}$	BR
English Premier League	<b>0.482</b>	0.471 ( $\pm 0.041$ )	0.440 ( $\pm 0.034$ )	0.403 ( $\pm 0.028$ )	0.390	0.379 ( $\pm 0.052$ )
Italian Serie A	<b>0.485</b>	0.442 ( $\pm 0.051$ )	0.426 ( $\pm 0.039$ )	0.404 ( $\pm 0.035$ )	0.372	0.392 ( $\pm 0.018$ )
Spanish La Liga	0.438	<b>0.462</b> ( $\pm 0.038$ )	0.455 ( $\pm 0.031$ )	0.375 ( $\pm 0.035$ )	0.418	0.356 ( $\pm 0.013$ )
German Bundesliga	<b>0.449</b>	0.415 ( $\pm 0.037$ )	0.433 ( $\pm 0.047$ )	0.357 ( $\pm 0.043$ )	0.372	0.346 ( $\pm 0.053$ )
French Ligue 1	<b>0.458</b>	0.438 ( $\pm 0.036$ )	0.435 ( $\pm 0.040$ )	0.384 ( $\pm 0.046$ )	0.392	0.388 ( $\pm 0.037$ )
Mean	<b>0.461</b> ( $\pm 0.020$ )	0.450 ( $\pm 0.016$ )	0.438 ( $\pm 0.011$ )	0.384 ( $\pm 0.040$ )	0.389	0.371 ( $\pm 0.041$ )
All leagues	-	<b>0.434</b> ( $\pm 0.027$ )	0.407 ( $\pm 0.015$ )	-	-	-

- BR - Base Rule
- $RF_{10CV}$  - Random Forest 10-fold Cross Validation
- $RF_{TS}$  - Random Forest Time Series
- $RF_{GA}$  - Default Random Forest with Genetic Algorithm
- $RF_{GA(T)}$  - Random Forest with Genetic Algorithm & Tuned
- $RF_{HT+PM}$  - Random Forest & Genetic Algorithm with In-play & Pre-match Data



# Comparison with Betting Exchange (Italy)

Brier Score for the PM was 0.544 and for the RF was 0.623.

Match	Prediction	Home	Draw	Away
AC Milan vs Cagliari	Actual	<b>1</b>	0	0
	Random Forest	<b>0.45</b>	0.38	0.17
	BetFair	<b>0.61</b>	0.29	0.09
Crotone vs Empoli	Actual	<b>1</b>	0	0
	Random Forest	0.39	<b>0.42</b>	0.19
	BetFair	0.34	<b>0.41</b>	0.24
Empoli vs Udinese	Actual	<b>1</b>	0	0
	Random Forest	0.26	<b>0.43</b>	0.30
	BetFair	0.28	<b>0.43</b>	0.29
Lazio vs Chievo	Actual	0	0	<b>1</b>
	Random Forest	<b>0.54</b>	0.30	0.16
	BetFair	<b>0.62</b>	0.30	0.08
Lazio vs Crotone	Actual	<b>1</b>	0	0
	Random Forest	<b>0.60</b>	0.24	0.16
	BetFair	<b>0.68</b>	0.25	0.06
Napoli vs Pescara	Actual	<b>1</b>	0	0
	Random Forest	<b>0.56</b>	0.28	0.16
	BetFair	<b>0.75</b>	0.21	0.04

Match	Prediction	Home	Draw	Away
Napoli vs Pescara	Actual	<b>1</b>	0	0
	Random Forest	<b>0.56</b>	0.28	0.16
	BetFair	<b>0.75</b>	0.21	0.04
Palermo vs Inter	Actual	0	0	<b>1</b>
	Random Forest	0.19	<b>0.41</b>	0.40
	BetFair	0.12	0.33	<b>0.56</b>
Roma vs Cagliari	Actual	<b>1</b>	0	0
	Random Forest	<b>0.48</b>	0.36	0.16
	BetFair	<b>0.72</b>	0.21	0.06
Sampdoria vs Empoli	Actual	0	<b>1</b>	0
	Random Forest	0.38	<b>0.40</b>	0.22
	BetFair	<b>0.46</b>	0.36	0.18
Sassuolo vs Torino	Actual	0	<b>1</b>	0
	Random Forest	0.26	<b>0.39</b>	0.34
	BetFair	0.23	0.36	<b>0.40</b>
Udinese vs AC Milan	Actual	<b>1</b>	0	0
	Random Forest	0.32	<b>0.40</b>	0.27
	BetFair	0.27	<b>0.40</b>	0.33



# Comparison with Betting Exchange (Eng)

Brier Score for the PM was 0.622 and for the RF, 0.655.

Match	Prediction	Home	Draw	Away
Arsenal vs Burnley	Actual	<b>1</b>	0	0
	Random Forest	<b>0.70</b>	0.17	0.13
	BetFair	<b>0.70</b>	0.22	0.07
Burnley vs Southampton	Actual	<b>1</b>	0	0
	Random Forest	0.14	0.38	<b>0.48</b>
	BetFair	0.16	0.41	<b>0.44</b>
Hull vs Bournemouth	Actual	<b>1</b>	0	0
	Random Forest	0.35	0.26	<b>0.39</b>
	BetFair	0.29	<b>0.38</b>	0.32
Liverpool vs Swansea	Actual	0	0	<b>1</b>
	Random Forest	<b>0.71</b>	0.14	0.15
	BetFair	<b>0.66</b>	0.28	0.07
Man City vs Burnley	Actual	<b>1</b>	0	0
	Random Forest	<b>0.55</b>	0.31	0.14
	BetFair	<b>0.53</b>	0.33	0.14
Man City vs Tottenham	Actual	0	<b>1</b>	0
	Random Forest	<b>0.52</b>	0.22	0.26
	BetFair	<b>0.46</b>	0.34	0.20
Watford vs Middlesbrough	Actual	0	<b>1</b>	0
	Random Forest	<b>0.40</b>	0.32	0.27
	BetFair	0.31	<b>0.46</b>	0.23
West Ham vs Man Utd	Actual	0	0	<b>1</b>
	Random Forest	0.29	<b>0.36</b>	0.35
	BetFair	0.06	0.23	<b>0.71</b>

A photograph of a football player in a blue Everton kit, identified as Lukaku by the name on his back, celebrating a goal. He is on the grass, kneeling on one knee with his arms raised in a 'V' shape. The background shows a large stadium filled with spectators.

Everton won the match by three goals to nil.

Everton vs Southampton	Actual	1	0	0
	Random Forest	<b>0.59</b>	0.28	0.13
	BetFair	0.35	<b>0.40</b>	0.25



- We have derived a base rule for predicting fulltime results at the halftime interval of football matches
- Parsed in-play data from an Opta source and developed a dataset consisting of the differences between team statistics till the halftime for both pre-match and match day data
- Shown that random forest using both types of available data produced the best results
- Similar accuracy as the betting market when considering probabilities for predictions, in some cases out-performing the market and thus giving an edge to the user



- Addition of other predictors such as inclusion of key players in team, player individual form and their scoring and defensive abilities
- Rate or number of entrances into opposition penalty box and dangerous areas
- Split match data into several minute time-frames and investigate predictions along the time of play
- Investigate predictions on other markets such as over/under goals and next team to score
- Use predictions as part of a betting strategy



**Thank you for your  
attention**