

Man v Machine: Greyhound Racing Predictions

MSc Research Project
Data Analytics

Alva Lyons
x15014274

School of Computing
National College of Ireland

Supervisor: Mr. Oisín Creaner

National College of Ireland
Project Submission Sheet – 2015/2017
School of Computing



Student Name:	Alva Lyons
Student ID:	x15014274
Programme:	Data Analytics
Year:	2016
Module:	MSc Research Project
Lecturer:	Mr. Oisín Creaner
Submission Due Date:	21/12/2016
Project Title:	Man v Machine: Greyhound Racing Predictions
Word Count:	5964

I hereby certify that the information contained in this (my submission) is information pertaining to research I conducted for this project. All information other than my own contribution will be fully referenced and listed in the relevant bibliography section at the rear of the project.

ALL internet material must be referenced in the bibliography section. Students are encouraged to use the Harvard Referencing Standard supplied by the Library. To use other author's written or electronic work is illegal (plagiarism) and may result in disciplinary action. Students may be required to undergo a viva (oral examination) if there is suspicion about the validity of their submitted work.

Signature:	
Date:	20th September 2018

PLEASE READ THE FOLLOWING INSTRUCTIONS:

1. Please attach a completed copy of this sheet to each project (including multiple copies).
2. **You must ensure that you retain a HARD COPY of ALL projects**, both for your own reference and in case a project is lost or mislaid. It is not sufficient to keep a copy on computer. Please do not bind projects or place in covers unless specifically requested.
3. Assignments that are submitted to the Programme Coordinator office must be placed into the assignment box located outside the office.

Office Use Only	
Signature:	
Date:	
Penalty Applied (if applicable):	

Man v Machine: Greyhound Racing Predictions

Alva Lyons

x15014274

MSc Research Project in Data Analytics

21st January 2017

Research Questions

Can a system be built using machine learning techniques to predict greyhound racing results? When bench-marked against a human expert how effective is this system?

Abstract

The purpose of this research is to ascertain if machine learning techniques can prove advantageous in predicting the outcome of greyhound races. The main focus of this research is in bridging the gap between existing sports prediction models which use manual feature selection to creating a model built from machine chosen subsets by algorithmically sub-setting the feature space. Feature selection is the process of sub-setting the feature space by analysing the relevance of features both to each other and to the predicted variable so that only the most relevant features are used within the modelling framework. The reason for introducing the greyhound racing expert is to test whether the system can outperform the average social gambler who tends to make their betting selection based on tips given to them by domain experts.

1 Introduction

The greyhound racing industry in Ireland is controlled by the Irish Greyhound Board (IGB). It is estimated that 720,000 visitors annually attend IGB controlled greyhound stadia. Shelbourne Park is the premier greyhound stadium in Ireland and hosts one of the world's richest greyhound races, The Irish Derby, every September. Racing takes place in Shelbourne Park every Wednesday, Thursday and Saturday.

This research will attempt to predict the finishing position of a greyhound in a given race. The data used in this research comprises of 64,908 observations of 10,986 races ran in Shelbourne Park between January 2009 and August 2016. The prediction rate of this model is bench-marked against that of the stadium's resident greyhound expert who is employed by IGB to predict the winning greyhound, the top 2 and the top 3 finishing greyhounds for the top of the race card for each race on a given race night.

The use of machine learning techniques in sports prediction is not a new phenomenon but rather it has gained many more practitioners since the spread of online gambling markets. While the use of machine learning techniques is prevalent in predicting horse racing (Butler et al. (1998), Silverman and Suchard (2013), Davoodi and Khanteymoori (2010), Williams and Li (2008) etc.) results there have only three documented cases of utilising machine learning in predicting the outcome of greyhound races. While the two sports are often synonymous there are distinct differences between the two which ensures that modelling concepts need to be amended. A greyhound race result is the outcome of 6 greyhounds chasing a mechanical hare in their attempts to catch it; while a horse race result is the outcome of the interactions between a jockey and its mount as they traverse the race course. While this might seem trivial, the key difference is apparent when considering a model's attempts at predicting the finishing positions of competitors in a race. A greyhound is bred to chase the hare and will continue its mission even if the race has already been won. On the other hand, a jockey which surmises it has no chance of finishing in the first x positions, may choose to pull back so that the horse's handicap rating is not affected for its next race. This nuance is one of the factors that led this researcher to choose greyhound racing as the sport of choice for this research.

The seminal paper in the field of greyhound racing predictions, Chen et al. (1994), dates back to 1994 and uses the knowledge of a greyhound expert for feature selection in choosing which performance variables to use when running machine learning techniques on their dataset. Both of the follow up studies utilise a similar feature selection in their models (Schumaker and Johnson (2008), Johansson and Sönströd (2003)). This research paper uses feature selection algorithms to limit the problem space of the domain in order to avoid the subjectivity of human interactions within the modelling process.

Additionally, this paper combines various data mining techniques from sentiment analysis through to deep learning ensemble methods in its attempts to test if a machine learnt model can out-perform an expert in the area of greyhound racing predictions. The rest of this document is laid out as follows:

- **Section 2** discusses the related work in the field of sports predictions and highlights the role this research plays within this field.
- **Section 3** discusses the methodology framework used in completing this research.
- **Section 4** reviews and justifies the implementation steps carried out in this research.
- **Section 5** evaluates the results of the prediction algorithms.
- **Section 6** concludes the research and discusses potential future work to be carried out.

2 Related Work

2.1 Sports Predictions

The literature and academic work produced in the field of sports predictions is far reaching. Using historical results data to predict the outcome of sporting events has gained

exposure due to the growth of on-line betting markets and the large volumes of historical data which are easily accessible.

2.2 Horse Racing’s Favourite Long-Shot Bias

The sport of horse racing is often synonymous with greyhound racing and includes a vast amount of literature and analysis on sports prediction. Many of the earlier models focus on predicting the efficiency of betting markets in the hopes to apply profit to bets. Central to these predictions is the perceived notion of the favourite long-shot bias which is based on the phenomenon that continual betting on an outside selection, or long-shot, will return less than betting on a favourite despite the higher odds attached to the long shot, Leighton Vaughan Williams (1997). Kanto et al. (1992) discovered that most gamblers are willing to risk losses by over valuing long shots due to the enticement of added gain when winning on a higher priced selection.

Cain et al. (2003) carried out analysis to test whether the favourite long shot bias that is prevalent in horse racing betting markets is found in other sports. The results of their analysis shows that betting repeatedly on favourites results in smaller losses than betting on long shots in most sports; with the exception of soccer and greyhound racing, Lyons (2016).

2.3 Problem Space and Feature Set

An important step in the data mining process is choosing which features to include in your model. Feature selection can either be done manually through the use of domain knowledge or algorithmically with the use of machine learning methods. The race card available on tracks includes 50 variables which could potentially affect the outcome of a race. Adding all of these variables into a model would increase it’s complexity and be algorithmically inefficient, Lyons (2016).

Many of the works performed on predicting results of horse and greyhound races focus on the model used for prediction and it’s tuning parameters rather than the selection of the feature subset (Pudaruth et al. (2013), Davoodi and Khanteymoori (2010), Williams and Li (2008)). Their feature subset are listed but the motivation behind choosing which features to include in their model is not elaborated on. This lack of formal explanation for the feature subset doesn’t allay the assumption of human subjectivity in the choosing of what performance variables affect the outcome of these sporting events.

McCabe and Trevathan (2008)’s paper focuses more on the feature set than the model used in sports prediction. This paper provides an interesting discussion on why variables were included in the model however they are very vague on the potential ”subjective” variables not added. Similar to the papers listed above feature selection in the research by McCabe and Trevathan is a manual process and does not use machine learning to choose the optimal subset of features to include in the modelling.

2.3.1 Historical Feature Selection Techniques in Greyhound Racing Predictions

Chen et al. (1994) in their prediction of greyhound racing results chose their feature set following discussions with domain experts who informed them which 10 performance variables they believed were most important in predicting winners. They admit that while this is not optimal it is a consequence of their chosen algorithms being unable to handle noisy data. Remarkably neither Johansson and Sönströd (2003) nor Schumaker and Johnson (2008), in their follow up studies, chose to research further attempts at feature selection. Rather they used a similar feature subset to those used in the study by Chen et al.. (Lyons (2016))

2.3.2 Bridging The Research Gap

This research attempts to apply various feature selection algorithms to the transformed dataset in order to ascertain which features have a greater impact on influencing a greyhounds finishing position within a race. The features which are extracted as relevant to this domain problem are then chosen as the final dataset to be used in the model. The choice of using a neural network in the modelling phase of this research is to test if the use of machine based feature selection can outperform those as used by Chen et al. (1994) and Johansson and Sönströd (2003). As this research is focusing on classification rather than regression modelling the model choice of Schumaker and Johnson (2008) (Support Vector Regression) is discounted from the outset.

3 Methodology

The methodology used in this research is Knowledge Discovery in Databases (KDD). The KDD methodology allows for an iterative approach to the processes involved in extracting knowledge from raw data. Initial plans were to utilise the SEMMA notation, as developed by SAS, but the sequential nature of this methodology couldn't rival the flexibility and interactivity of KDD, Azevedo and Santos (2008). KDD focuses on the entire process from data selection through pre-processing, extraction, data mining to interpretation, Fayyad et al. (1996). An illustration of the KDD methodology as it pertains to this research is shown in Figure 1

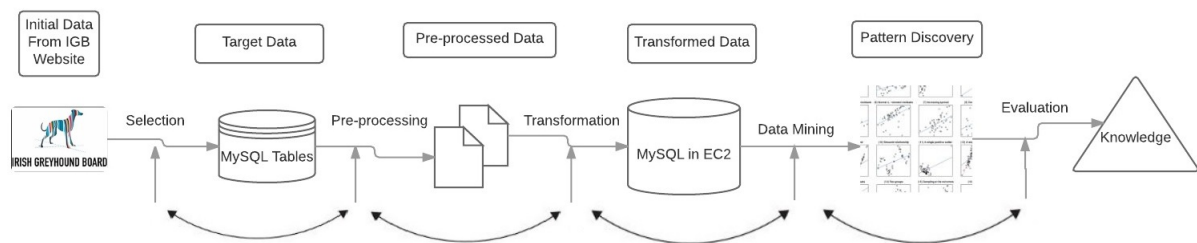


Figure 1: KDD Methodology

3.1 Selection

The raw data used in this research is extracted from the Irish Greyhound Board’s (IGB) website ¹ using a Python script. This data consists of 383,746 observations of 28,271 races ran throughout Ireland between 1st April 2007 and 3rd September 2016. This data is collected from multiple embedded pages and loaded into 6 tables in a MySQL database. The flow of this script is illustrated in the diagram in Appendix A.

3.2 Pre-Processing

The benefits of a thorough pre-processing phase ensure a strong knowledge of the dataset is gained before transformations commence. The data from IGB’s website contains numerous inaccuracies and missing data points ensuring the pre-processing phase of the KDD methodology plays an integral role in this research. Errors in the data are discovered when the data is examined using visualization and descriptive statistics.

3.2.1 Dealing with Missing Values

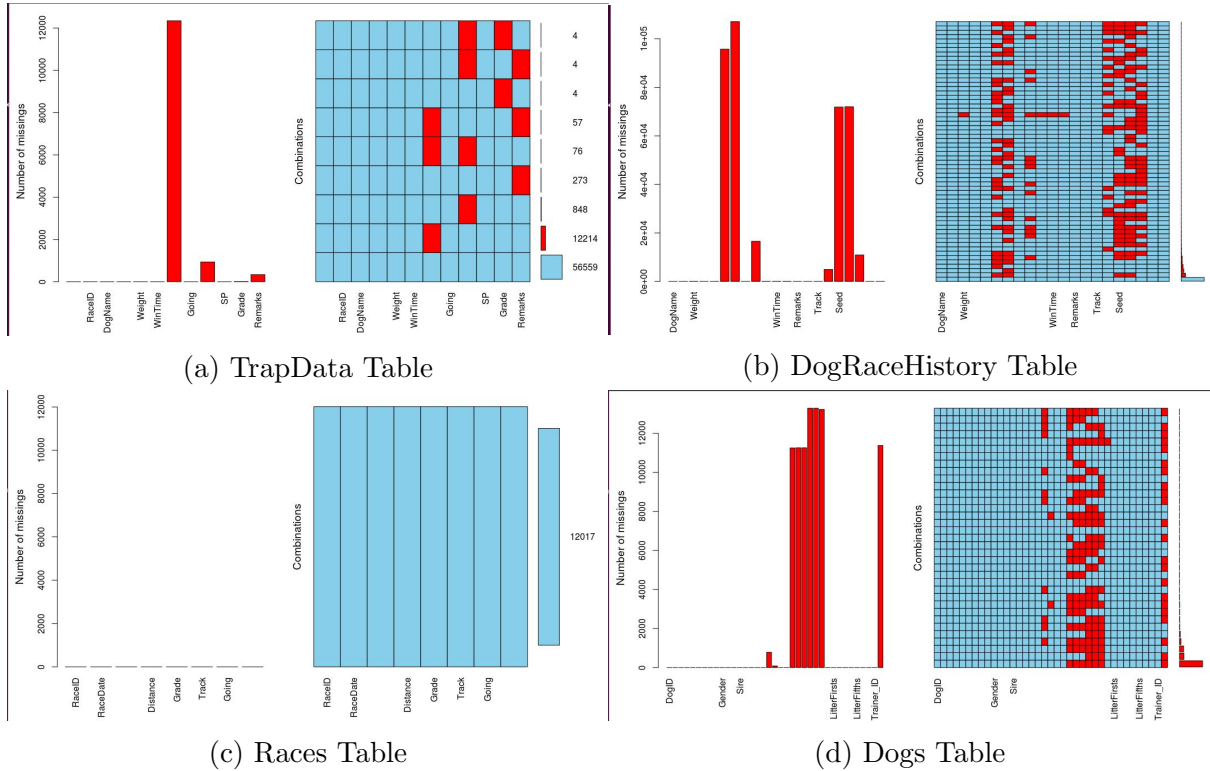


Figure 2: Missing Values in Raw Data Tables

In dealing with missing values it is necessary to ascertain why a value is missing or incorrect and to decide if action is to be taken, Witten et al. (2011). The proportion of missing values to actual data can be seen in Figure 2. In these diagrams the red squares represent missing values while the blue squares represent the presence of data. It is evident that the DogRaceHistory table, in particular, has a large proportion of missing

¹www.igb.ie/results

values. As this table is integral to determining the performance history of a greyhound it is important to ascertain the best method of handling this missing data.

Domain knowledge plays an important part in deciding what steps to take in handling missing values. For instance, a NULL value in the Seed column does not represent a true missing value. A greyhound's seed indicates their preferred running style; Inside (I) seeded greyhounds tend to run toward the bend; Middle(M) seeded greyhounds tend to run in the centre of the track; Wide(W) seeded greyhounds run toward the rails. The lack of one of these characters in the seed column is likely to indicate that a greyhound has no preferred running style. For this reason the missing data points in this column are replaced with "A" (Any).

Further exploratory analysis of the DogRaceHistory table shows up explanations for some of the high volume of missing data points within this dataset. Table 1 depicts the percentage of missing data in this table that can be attributed to the inclusion of time trials.

Table 1: % of Missing Values in DogRaceHistory Table that are attributed to Time Trials

DogRaceHistory Table	
Field	Responsible for % of Missing Data
Weight	100
NumberOfDogs	100
WinTime	100
Going	100
PlacedDistance	99.9
RunnerGrade	86.90
RaceGrade	73.1
SP	70.5
Remarks	61.2
SectionalPosition	46.9
SectionalTime	44.5
EstimatedTime	8.4

Time trials take place on race nights before racing commences whereby 1 or more greyhounds run around the track in a non competitive setting to see how fast they can chase the mechanical hare. Due to the non competitive nature of these events it is deemed appropriate to remove these from the dataset. Removal of time trials significantly reduce the number of missing data points in the DogRaceHistory table.

3.2.2 Dealing With Incorrect Values

Outlier detection is performed on the dataset to allow for purification of the data before modelling commences, Hodge and Austin (2004). For instance, in the *Weight* column of

the DogRaceHistory table the minimum weight of a greyhound is listed as 0lbs while the maximum is listed as 677lb. The average weight of a greyhound is 65-75lbs. The outlier values are imputed by taking the mean value of a greyhounds weight in its preceding and succeeding 2 races and inputting this as the new value in the weight column. In the instances where a greyhound does not have any other races against it the mean value of the weights of its competitors is imputed as its new value. Similar imputations are performed on other outliers within the dataset.

3.2.3 Pruning The Dataset

The track ratings of greyhound stadia in Ireland differ depending on the ground conditions; as such the *RunnerGrade* variable has varying significance depending on where a race took place. As this research attempts to predict results in Shelbourne Park it is deemed appropriate to limit the race history of greyhounds to their performance at this race track. Additionally, only A grade (middle distance) races are used in this research. The reason for omitting sprint races lies in the short distance between the first bend and the finish line. In sprint races if a runner gets knocked at the first bend their chances of recovery are limited Lyons (2016). The remaining dataset consists of 64,908 observations of 10,986 races ran between January 2009 and August 2016.

3.2.4 Limitations Of The Dataset Discovered During Pre-Processing

During the pre-processing phase limitations of the dataset are exposed. These limitations lie in the scraped data being from a view of the IGB's database at the time of scraping. While this doesn't affect historical race content it does ensure that owner and trainer data is inadmissible for modelling as there is no way of ascertaining whether the greyhound was attached to its current owner or trainer at the time of each historical race. Including the data in these two tables in the modelling phase of this research could potentially lead to incorrect predictions based on corrupt data.

Similiarly, litter distribution data in the Dogs table, which depicts the number of starts and race placings of a greyhound's siblings, is also inadmissible as it is that of a view of the to-date total of the litter at the time of scraping rather than the time of racing.

While omitting this data does have its benefits in that it cuts down the processing time of feature selection in the data mining phase; it restricts the model, in that it does not have access to the same data available in real time to the greyhound expert the model will be bench-marked against.

3.3 Transformation

3.3.1 Text Analysis

The remarks column in the DogRaceHistory table provides a shorthand of comments on how a greyhound ran in a given race eg. FAw (Fast Away), BBkd (Badly Baulked), TRec (Track Record) etc. The full list of 320 possible remarks can be found on the IGB's website ².

²<https://www.igb.ie/upload/pdf/Abbreviations.pdf>

In order to analyse these remarks across the dataset the basic premise of sentiment analysis is performed such that the text is classified as expressing a positive or negative tenet, Liu (2010). While sentiment analysis deals with "the computational treatment of opinion, sentiment, and subjectivity in text", Pang and Lee (2008), and is considered to be more suitable for text mining of unstructured datasets; the simplicity and power of this method is deemed appropriate to analyse this variable.

The first step in utilising the premise of sentiment analysis is to create domain specific lexicons of the shorthands used in the remarks column. Once the lexicons are created confirmation is received from 2 experts working within the greyhound racing industry to ensure subjectivity is minimized during this phase. These dictionaries are created by assigning each shorthand comment into one of 5 categories; Very Positive, Positive, Neutral, Negative, Very Negative.

The motivation for utilising this variable and performing text analysis lies with the possibility that the greyhound's ability is not properly reflected by it's finishing position. For instance, a greyhound may only finish 5th in a race despite being quick out of the traps due to being impeded by another greyhound. By the same respect the 1st placed finisher in this race might have missed the early fighting and received a clear run despite being slow away.

In order to run sentiment analysis on this column the 5 dictionaries are loaded into MySQL tables and the remarks column is scored using an SQL statement which scans each of these tables and matches the word in remarks to those in the dictionaries. A scoring is given to the words depending on which category they fall into:

- Very Positive = +2 points
- Positive = +1 point
- Neutral = 0 points
- Negative = -1 point
- Very Negative = -2points

The scoring for each remark is totalled and set as the *RemarkScore* for each greyhound in a given race.

3.3.2 Feature Engineering

Feature engineering is the task of extracting and transforming raw data into features that better represent the problem domain in order to improve a model's prediction when new data is introduced, Brownlee (2014). The importance of feature engineering lies in finding the best representation of variables in a dataset in order to bridge the gap between the features in the initial problem domain to the structure of features needed for the solution architecture, Dash and Liu (2003). There are many elements to feature engineering from framing the problem domain to data cleansing and formatting. The element discussed in this research pertains to the manual construction of new features from the raw data.

The raw data is transformed to create meaningful information from each dog's race history. While Chen et al. (1994) and Schumaker and Johnson (2008) average their variables over 7 races, this research looks to emulate the on track race card; which provides historical data on each greyhound's last 5 races; by averaging the greyhound history statistics over 5 races. This ensures that the model has access to the same data as the ordinary social gambler. Where a greyhound has run less than 5 races the data is averaged over the number of runs of that greyhound up to a maximum of 5 races. The below formula shows how the rolling averages are calculated; this example calculates a greyhounds average position at the first bend in each of it's last n races.

$$BreakAvg5 = \frac{FirstBend_r1 + FirstBend_r2 + \dots + FirstBend_rn}{n}$$

Similar formulae are used to transform other variables in the raw data. Figure 3 depicts a table of the transformations that took place in this phase of the research methodology.

Figure 3: Variables Created From Raw Data

Transformations	
Field	Description
DogsAge	Age of dog at time of the race. Subtracting RaceDate from WhelpDate.
1st/2nd/3rd/4th Bend	The SectionalPosition column is a string of 4 digits which when split represent the greyhound's position in a race at each of the first 4 bends
BreakAvg5	The dog's average position at the first bend over its last 5 races.
Avg2ndBend	Average position at 2nd bend in last 5 races.
WinPercent5	Average win percentage in last 5 races.
PlacedPercent5	Percentage of 1st/2nd place finishes in last 5 races.
ShowPercent5	Percentage of top 3 finishes in last 5 races.
EstTimeAvg5	Average finishing time in last 5 races.
AvgRemarks5	Average Remarks score over 5 races. *see Section 3.3.1
FinishingPositionAvg5	Average Finishing Position in last 5 races.
RankedGradeAvg5	A scoring on RunnerGrade v RaceGrade - how a dog ran in each of its last 5 races - if runner grade is better than race grade additional points are given.
PrizeMoneyWonAvg5	Average Prize Money won over the last 5 races.
SecTimeAvg5	Average time taken to reach the start line for the first time in last 5 races.
PlacedPercent	Overall percentage of top 2 finishes.
ShowPercent	Overall percentage of Top 3 finishes.
OverallAvgTime	Average time across all races.
DaysSinceLastRace	The number of days between races.

3.4 Data Mining

Data mining is the process of analysing datasets to find unobserved and often unsuspected relationships within the data by combining statistics, artificial intelligence and machine learning features, (Hand et al. (2001)). Fayyad et al. 1996 address the importance of understanding the data mining activity before including it in the KDD process. Similar to the KDD methodology, the choosing of an algorithm to use in tackling a prediction problem can involve many interactions and iterations before knowledge is gleaned. An important first step is to decide which data mining process of predictive analysis is required in ascertaining the value of the predictor variable. **Classification analysis** deals with predicting which category or class an object falls into. The required output is a discrete variable. **Regression analysis** is used to predict missing or unavailable numerical data values; the output variable is a continuous variable. Han (2005)

The application of predicting the outcome of a competitive event does not strictly fall into either a classification or a regression problem and as a result both regression and classification techniques are possible within the realms of this research domain. As a classification problem the output variable can be a matter of predicting the binary output of win or lose. As a regression problem it is possible to look at the finishing order of a race with the view to regressing on the FinishingPosition variable.

This research approaches the problem of predicting greyhound racing results as a classification problem. However, rather than choosing binary classification of "win" or "lose" it attempts are made to classify a greyhound's Finishing Position. The reasoning for not choosing binary classification is partly due to class imbalance; for each race 6 greyhounds are entered and 5 greyhounds cannot win as such the number of observations in the lose class in the training set is larger than that of winners and random predictions could result in a higher rate of prediction due to chance alone. Additionally, by choosing to classify the problem using the Finishing Position as the predictor variable this allows for testing how wrong a predicted class is. For instance, incorrectly predicting a 1st place finisher will finish in 2nd place is "less wrong" than predicting the same greyhound will finish in 6th place. The choice of algorithms and justification for their uses is discussed in the implementation section of this paper.

4 Implementation

4.1 Tools Used

The tools used in implementing this research are:

- Python (Version 2.7.12)
- MySQL
- R (Version 3.3.1)
- R Studio 64bit
- Amazon EC2

Python is used in the selection phase of this research to scrape the raw data due to the power of its BeautifulSoup library; which provides an easy to use framework for parsing HTML into a tree representation. The data mining phase utilises R; a statistical programming language which is widely used for data analysis, Lantz (2013).

4.2 Examining The Dataset

Once the feature engineering phase is complete the next step is to combine and explore the dataset. A flattened correlation matrix of the processed dataset is produced in R using the *corrplot* library (See Figure 4).

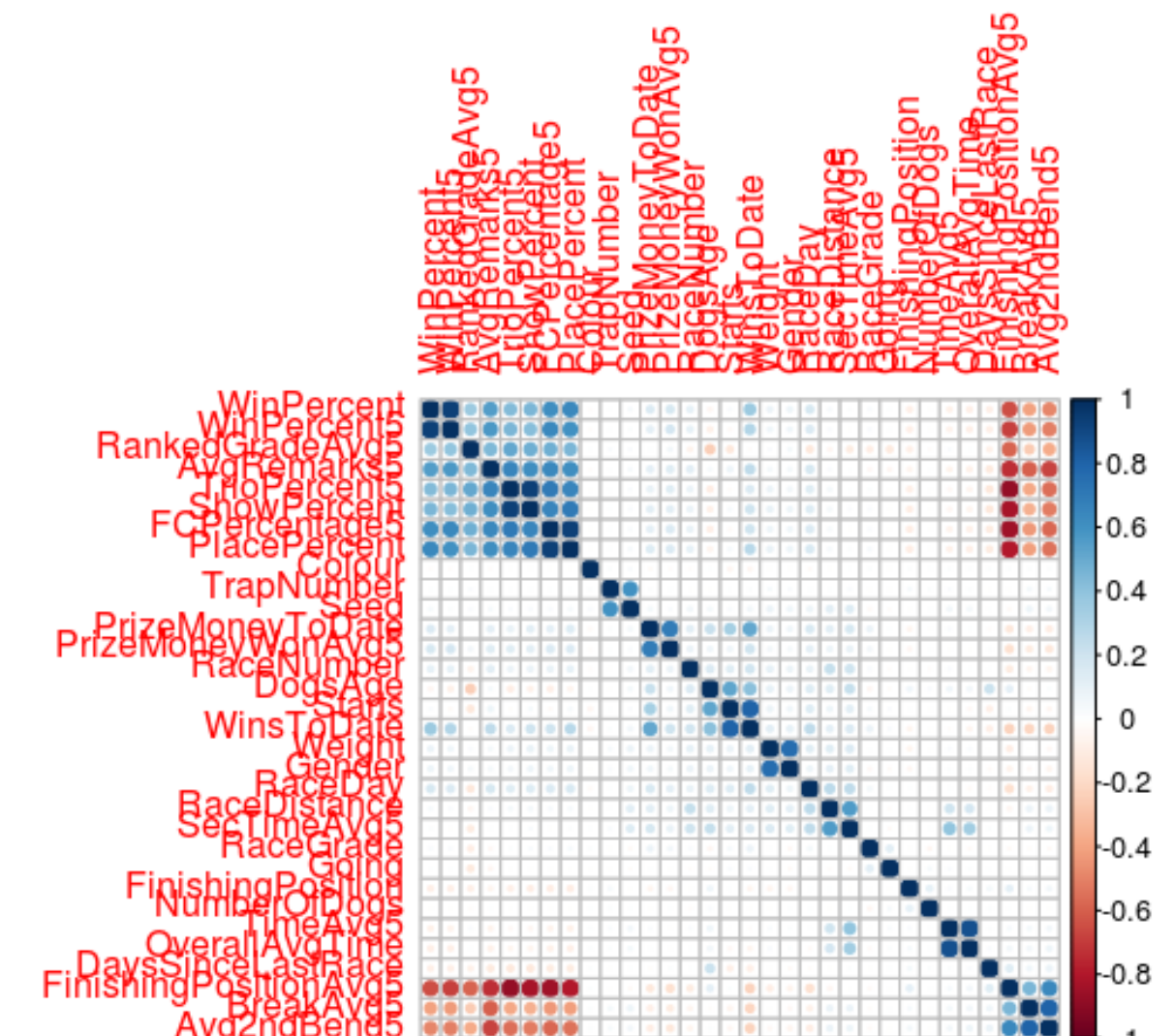


Figure 4: Flattened Correlation Matrix

As is evident in this visual representation of the correlation between features there are a number of strong correlations amongst variables. While some correlations are expected; such as the positive correlation between the number of starts of a greyhound and the number of wins; others are unexpected, such as the negative correlation between a greyhound's win percentage (the percentage of wins over all races) and its finishing position in its last 5 races. A negative correlation depicts an inverse proportionality

between variables in that as one increases the other will decrease. This would suggest that over time a greyhound's recent form has an obvious affect on its overall performance in that rather than remaining consistent it will either improve or deteriorate over time.

4.3 Feature Selection

Feature selection is the process of binning variables into subsets of relevant and irrelevant features such that only the most relevant features are used within the modelling framework, Dash and Liu (1997). A feature is deemed to be relevant if it affects the target problem in any way. The benefits of feature selection lie in reducing the complexity and run-time of the machine learning algorithm. The reducing of complexity allows for better understanding of the patterns that arise in the data mining process. Additionally, feature selection when performed correctly, can improve model performance.

4.3.1 Methods of Feature Selection

There are three categories of feature selection methods; Wrapper, Filter and Embedded.

- **Filter Methods** - are concerned with exploring only the inherent features of a dataset. They are based on statistical tests and are independent of the variable to be predicted.
- **Wrapper Methods** - Unlike filter methods wrapper methods are used to find features subsets which interact with the variable to be predicted. In this way the choosing of a wrapper method is closely linked to the choosing of a modelling algorithm as the feature subset space is wrapped around the classifying model, Saeys et al. (2007).
- **Embedded** - Embedded methods are an extension of the Wrapper Method framework and attempt to combine the best properties of the preceding two methods. The feature selection is 'embedded' in the modelling algorithm which runs feature selection and prediction concurrently.

This research focuses on wrapper and embedded methods as they interact with the variable to be predicted, are less likely to get stuck in a local optima and model feature dependencies. The limitations of these methods, however, lie in the increased risk of over-fitting, Saeys et al. (2007).

4.3.2 Wrapper with Boruta Package

Several packages are available in order to ascertain the importance of independent variables in predicting the dependent variable. The *Boruta* package in R comprises of a wrapper algorithm which utilises random forests in order to extract relevant features from a data set. This is achieved by comparing a variable's importance against importance that is achievable at random, Kurasa and Rudnicki (2016). The application of this package on the dataset did not dramatically reduce the feature space; removing only one variable (*Colour*) from the original 30 inputted into the algorithm.

4.3.3 Wrapper With caret & randomForest Packages

The *caret*, an acronym for **C**lassification and **R**egression **T**raining, package contains functions to organise a model's training approach, Kuhn (2016), and utilises a number of other packages in R. In this example, the *caret* package wraps around the *randomForest* package in order to rank the variable importance of the features in the dataset. Variable importance ratings are assigned to each feature and they are then ranked according to how important they are to the dependent variable, *FinishingPosition*. The ranked feature results are shown in Appendix B. Interestingly *Colour*, which is the only variable to be dismissed when running the Boruta package with random forests, ranks above 9 other features when the the *caret* package is wrapped around a random forest.

4.3.4 Embedded - Recursive Feature Elimination with caret

Recursive feature elimination (RFE) is an embedded method of feature selection. It attempts to find the optimal subset of features by iterating through all features, assigning weights to each feature depending on its value to the dependent variable. Features are eliminated based on their ranked weighting, in that those with a smaller weighting are eliminated first. Once a feature is pruned the remaining features are reassigned weights and the process is iterated until a stopping criterion is reached whereby the optimal number of features is selected, Guyon et al. (2002).

The *caret* packages provides a set of predefined functions to embed RFE with algorithmic functions; such as Naïve Bayes; Random Forests; and Bagged Trees. These 3 functions are modelled on the dataset in order to attempt to find the optimal subset of features to use in prediction.

1. **Naïve Bayes** is a classification algorithm which is based on Bayes Theorem and assumes independence amongst the feature space. The running of a recursive feature elimination using Naïve Bayes limits the feature subspace to 3 features; *FinishingPositionAvg5*, *OverallAvgTime*, *Avg2ndBend*. Examining the flattened correlation matrix in figure 4 tells us that the basic assumptions of Naïve Bayes are violated in the our dataset, in that the features are not independent. Therefore, the results of this analysis are dismissed for the remainder of this research.
2. **Random Forests** are the result of combining decision trees such that each tree depends on the values of a randomly sampled independent vector whereby the entire forest is distributed homogeneously Breiman (2001). The output of the random forest RFE is shown in Appendix C. This depicts the top 10 features selected to be *NumberOfDogs*, *TimeAvg5*, *OverallAvgTime*, *PrizeMoneyWonAvg5*, *DogsAge*, *FinishingPositionAvg5*, *SecTimeAvg5*, *PlacedPercent*, *RankedGradeAvg5* and *Avg2ndBend5*.
3. **Tree Bagging** is an ensemble method which uses decision trees to generate multiple versions of a predictor and aggregates the result, Breiman (1996). The output of the Tree Bagging RFE method is shown in Appendix D. The top 10 features returned embedding RFE with Tree Bagging are *DogsAge*, *Weight*, *TimeAvg5*, *SecTimeAvg5*, *OverallAvgTime*, *RaceNumber*, *RankedGradeAvg5*, *PrizeMoneyToDate*, *BreakAvg5* and *TrapNumber*.

4.3.5 Combined Feature Selection Results

The results of the caret & randomForest wrapper, embedded treebagging and embedded randomForest are combined to create an optimal subset of the data for use in the final stage of the modelling process. The top 10 features of each of these methods are combined in order to ascertain if there are any features which are prevalent across all feature selection methods. The ranked table is shown in Table 2.

Table 2: Top Ten Features Selected

Features Selected		
caret & randomForest	RFE with RandomForest	RFE with TreeBagging
DogsAge	NumberOfDogs	DogsAge
OverallAvgTime	TimeAvg5	Weight
SecTimeAvg5	OverallAvgTime	TimeAvg5
TimeAvg5	PrizeMoneyWonAvg5	SecTimeAvg5
Weight	DogsAge	OverallAvgTime
RankedGradeAvg5	FinishingPositionAvg5	RaceNumber
PrizeMoneyWonAvg5	SecTimeAvg5	PrizeMoneyToDate
BreakAvg5	PlacedPercent	BreakAvg5
RaceNumber	RankedGradeAvg5	RankedGradeAvg5
DaysSinceLastRace	Avg2ndBend5	TrapNumber

The features highlighted in red indicate those that are selected as a top 10 ranking feature across all 3 feature selection methods. These 5 variables are selected for the final feature subset. Five more features for this subset are selected by choosing the highest ranked features amongst the 3 methods deployed. The transformed variable, *AvgRemarks5*, resulting from the Text Analysis phase of this research (section 3.3.1), fails to make the cut despite finishing amongst the top 13 features across all feature selection methods.

Resulting Optimal Feature Subset: The final feature subset selected for the modelling phase consists of *DogsAge*, *OverallAvgTime*, *SecTimeAvg5*, *TimeAvg5*, *RankedGradeAvg5*, *Weight*, *RaceNumber*, *PrizeMoneyWonAvg5*, *BreakAvg5* and *Avg2ndBend5*.

5 Evaluation

5.1 Model Performance

For the purpose of modelling the dataset is split in a 60/20/20 ratio for training, validation and testing and all variables are scaled to reduce variance across the dataset. An important design criterion for model performance is choosing the correct parameters. It is important to ensure that while these parameters are tuned the test set is not utilised so as to avoid the model learning from iterations over the test set.

5.1.1 Neural Network

In order to emulate the research in the field of greyhound racing by Chen et al. (1994) and Johansson and Sönströd (2003); who used shallow neural networks in their predictions; the optimal feature subset chosen following feature selection is inputted into a Deep Learning Neural Network using the *H2O* package in R. Deep Learning reduces the complexity of an algorithm but is better suited to larger datasets.

This research initially utilises deep neural networks in its attempts at classifying the finishing position of a greyhound. A deep neural network is ran several times across different subsets of the data and the average prediction performance is found to be 18.92%. The use of shallow neural networks improves the performance of the feature subset; using just one hidden layer on subsets of the dataset improves the average prediction performance to 19.89%. The performance of the shallow neural network approaches the prediction rate of Chen et al. (1994), whose model correctly predicted 20% of winners. Similar comparisons cannot be made on Johansson and Sönströd’s paper as their model is assessed against monetary gain on bets rather than win percentages.

5.1.2 Random Forest

A random forest algorithm is ran on the feature subset to see if the use of an ensemble method can boost the prediction accuracy of the feature set. The random forest fails to improve upon the performance of the neural network; only predicting 19% of finishing positions correctly when tested against the validation set and 18.28% when run on the test set.

5.2 Model Evaluation

The breakdown of the model’s success rates as discussed in Section 5.1 are shown in table 3.

Table 3: Model Success Rates

Model	Success Rate
Random Forest	18.28%
Deep Neural Network	18.92%
Shallow Neural Network	19.89%

5.2.1 Greyhound Expert Prediction

The greyhound expert’s predictions are scraped from embedded pop-ups on the IGB’s website and the percentage of first place finishers correctly predicted is derived. A limitation of this research lies therein. This research, similar to the greyhound expert, attempts to predict the finishing order of greyhounds in a race; however the expert makes 3 predictions per race; what greyhound will win the race; what two greyhounds will finish in the top 2 in any order; and what 3 greyhounds will finish in the top 3 in any order. It is necessary to only choose the percentage of first place finishes correctly predicted as the other 2 predictions turn the problem from a 6 class multi-class problem to a binary

classification problem. The greyhound expert correctly predicted 23.7% of first placed finishers in the time frame used in this research (Jan. 2009 - Sept. 2016).

5.2.2 Random Chance

A typical greyhound race consists of 6 greyhounds chasing a mechanical hare. However, not all races will have 6 runners. The mean number of greyhounds per race in the final dataset is 5.864. The relative probability of correctly predicting winning greyhounds in a race should approach $\frac{1}{5.864}$ or 17.05% with increased data.

5.2.3 Comparative Performance in Greyhound Racing Prediction

Table 4: Comparative Performances

Model	Greyhound Expert	Random Chance	Chen et al.
(Deep) Neural Network	-4.78%	+1.87%	-1.08%
(Shallow) Neural Network	-3.81%	+2.84%	-0.11%

Neural Network: A neural network was chosen for the modelling phase of this project to analyse whether algorithmically limiting the feature space is more, less or equally as effective as using domain expert knowledge in choosing an optimal subset of features. As discussed in section 2.3.1, Johansson and Sönströd and Chen et al.; who limited their feature space by seeking advice on what variables to add into their model by the same experts they hoped their model would outperform; both used neural networks for their modelling phase. This research shows that algorithmic and manual feature selection in the domain of greyhound racing when combined with the neural network model have comparable results, with the manual feature selection performing 0.11% better; as shown in table 4.

6 Conclusion and Future Work

6.1 Conclusion

The aim of this research was threefold. The first was to bridge the gap of previous research in the domain of greyhound racing prediction modelling by algorithmically subsetting the feature space. The outcome of this research shows that the use of a machine learnt feature subset can be effectively used in place of a domain expert knowledge feature set without significant loss to prediction accuracy.

The second aim was to build a system to predict greyhound racing results. As discussed in section 5.2.2, with an increased number of races the random probability of correctly predicting the winner in a greyhound race should reach 17.05%. The models are tested across 20% of the dataset, 2197 races, which should bring the probability close to 17%. Table 3 shows that all 3 machine learning models preform better than random chance; thus highlighting that the use of machine learning techniques is advantageous in predicting greyhound racing results.

Finally, the system built was to be bench marked against a human expert in greyhound predictions. The system which is created in this research needs further tuning to rival the hit ratio of the human expert. While the model's performance is improved during tuning the results in table 4 depict that the use of shallow neural networks correctly predict 3.81% less winners than the human expert.

It can be concluded that machine learnt feature selection must at all times be accompanied by domain knowledge; it is in combining the two that an optimal feature set can be obtained. Although feature selection plays an important role in this research it is only 1 step within the iterative framework. It alone, cannot adequately account for a model's success or failure; rather the amalgamation of domain knowledge, feature engineering, feature selection and model selection when combined optimally ensure success.

6.2 Future Works

Future work for this research should concentrate on the modelling phase. Suggestions on possible work to bolster the success of this research are listed below:

1. Modelling in this paper is performed on each individual greyhound separately in order to ascertain their individual probability of winning a race given their historical performance data. The features subset developed in this research could be combined with conditional statistics to allow that the sum of probabilities for all greyhounds in a given race equal 1 so that within race competition can be accounted for.
2. Deep learning algorithms improve with added data. A future focus to improve model performance could be to generalise the feature selection across all 28,271 races scraped from tracks throughout Ireland.
3. Utilize the optimal feature subspace in regression modelling.

Acknowledgements

This research would not have been possible without the help of my supervisor **Mr. Oisín Creaner**, alongside **Mr. Michael Bradford** and **Dr. Simon Caton**, whose help and guidance throughout the completion of this Masters programme proved invaluable. I would like to take this opportunity to thank them sincerely for their generosity of time and spirit.

I would also like to thank the **Irish Greyhound Board** who generously allowed me to use their data in the completion of this Masters programme.

Without the continued support of my family and friends I would not have had the courage nor strength to be able to complete my Masters. Particularly, I would like to dedicate this thesis to my mam, **Eileen Lyons**, who sadly passed away before its completion. *Ar dheis Dé go raibh a h-anam.*

References

- Azevedo, A. and Santos, M. F. (2008). Kdd, semma and crisp-dm: a parallel overview., in A. Abraham (ed.), *IADIS European Conf. Data Mining*, IADIS, pp. 182–185.
URL: <http://dblp.uni-trier.de/db/conf/iadis/dm2008.htmlAzevedoS08>
- Breiman, L. (1996). Bagging predictors, *Machine Learning* **24**(2): 123–140.
URL: <http://dx.doi.org/10.1007/BF00058655>
- Breiman, L. (2001). Random forests, *Machine Learning* **45**(1): 5–32.
URL: <http://dx.doi.org/10.1023/A:1010933404324>
- Brownlee, J. (2014). Machine learning mastery.
URL: <http://machinelearningmastery.com/discover-feature-engineering-how-to-engineer-features-and-how-to-get-good-at-it/>
- Butler, J., Tsang, E. P. K. and Sq, C. C. (1998). EDDIE beats the bookies.
- Cain, M., Law, D. and Peel, D. (2003). The favourite-longshot bias, bookmaker margins and insider trading in a variety of betting markets, *Bulletin of Economic Research* **55**(3): 263–273.
URL: <http://dx.doi.org/10.1111/1467-8586.00174>
- Chen, H., Rinde, P. B., She, L., Sutjahjo, S., Sommer, C. and Neely, D. (1994). Expert prediction, symbolic learning, and neural networks: An experiment on greyhound racing., *IEEE Expert* **9**(6): 21–27.
URL: <http://dblp.uni-trier.de/db/journals/expert/expert9.htmlChenRSSSN94>
- Dash, M. and Liu, H. (1997). Feature selection for classification, *Intelligent Data Analysis* **1**: 131–156.
- Dash, M. and Liu, H. (2003). Consistency-based search in feature selection, *Artificial Intelligence* **151**(1–2): 155 – 176.
URL: <http://www.sciencedirect.com/science/article/pii/S0004370203000791>
- Davoodi, E. and Khanteymoori, A. R. (2010). Horse racing prediction using artificial neural networks, *Proceedings of the 11th WSEAS International Conference on Neural Networks and 11th WSEAS International Conference on Evolutionary Computing and 11th WSEAS International Conference on Fuzzy Systems*, NN’10/EC’10/FS’10, World Scientific and Engineering Academy and Society (WSEAS), Stevens Point, Wisconsin, USA, pp. 155–160.
URL: <http://dl.acm.org/citation.cfm?id=1863431.1863457>
- Fayyad, U., Piatetsky-Shapiro, G. and Smyth, P. (1996). The kdd process for extracting useful knowledge from volumes of data, *Commun. ACM* **39**(11): 27–34.
URL: <http://doi.acm.org/10.1145/240455.240464>
- Guyon, I., Weston, J., Barnhill, S. and Vapnik, V. (2002). Gene selection for cancer classification using support vector machines, *Machine Learning* **46**(1): 389–422.
URL: <http://dx.doi.org/10.1023/A:1012487302797>
- Han, J. (2005). *Data Mining: Concepts and Techniques*, Morgan Kaufmann Publishers Inc., San Francisco, CA, USA.

- Hand, D., Mannila, H. and Smyth, P. (2001). *Principles of Data Mining*, A Bradford book, MIT Press.
URL: <https://books.google.ie/books?id=SdZ-bhVhZGYC>
- Hodge, V. J. and Austin, J. (2004). A survey of outlier detection methodologies, *Artificial Intelligence Review* **22**(2): 85–126.
URL: <http://dx.doi.org/10.1007/s10462-004-4304-y>
- Johansson, U. and Sönströd, C. (2003). Neural networks mine for gold at the greyhound racetrack, *Proceedings of the International Joint Conference on Neural Networks* pp. 1798 – 1801 vol.3.
- Kanto, A. J., Rosenqvist, G. and Suvas, A. (1992). On utility function estimation of racetrack bettors, *Journal of Economic Psychology* **13**(3): 491 – 498.
URL: <http://www.sciencedirect.com/science/article/pii/016748709290006S>
- Kuhn, M. (2016). A short introduction to the caret package.
URL: <https://cran.r-project.org/web/packages/caret/vignettes/caret.pdf>
- Kursa, M. B. and Rudnicki, W. R. (2016). Boruta package in r - documentation.
URL: <https://cran.r-project.org/web/packages/Boruta/Boruta.pdf>
- Lantz, B. (2013). *Machine Learning with R*, Packt Publishing.
- Leighton Vaughan Williams, D. P. (1997). Why is there a favourite-longshot bias in british racetrack betting markets, *The Economic Journal* **107**(440): 150–158.
URL: <http://www.jstor.org/stable/2235276>
- Liu, B. (2010). Sentiment analysis and subjectivity, *Handbook of Natural Language Processing, Second Edition*. Taylor and Francis Group, Boca.
- Lyons, A. (2016). RIC research proposal for greyhound racing predictions modelling.
- McCabe, A. and Trevathan, J. (2008). Artificial intelligence in sports prediction, *Fifth International Conference on Information Technology: New Generations* .
- Pang, B. and Lee, L. (2008). Opinion mining and sentiment analysis, *Found. Trends Inf. Retr.* **2**(1-2): 1–135.
URL: <http://dx.doi.org/10.1561/15000000011>
- Pudaruth, S., Medard, N. and Dookhun, Z. B. (2013). Article: Horse racing prediction at the champ de mars using a weighted probabilistic approach, *International Journal of Computer Applications* **72**(5): 37–42. Full text available.
- Saeyns, Y., Inza, I. n. and Larrañaga, P. (2007). A review of feature selection techniques in bioinformatics, *Bioinformatics* **23**(19): 2507–2517.
URL: <http://dx.doi.org/10.1093/bioinformatics/btm344>
- Schumaker, R. P. and Johnson, J. W. (2008). An investigation of svm regression to predict longshot greyhound races, *Communications of the IIMA: Vol. 8: Iss. 2, Article 7* pp. 67–82.
URL: <http://scholarworks.lib.csusb.edu/ciima/vol8/iss2/7>

Silverman, N. and Suchard, M. (2013). Predicting horse race winners through a regularized conditional logistic regression with frailty, *Journal of Prediction Markets* **7**(1): 43–52.

URL: <http://EconPapers.repec.org/RePEc:buc:jpredm:v:7:y:2013:i:1:p:43-52>

Williams, J. and Li, Y. (2008). A case study using neural networks algorithms: Horse racing predictions in jamaica., *in* H. R. Arabnia and Y. Mun (eds), *IC-AI*, CSREA Press, pp. 16–22.

URL: <http://dblp.uni-trier.de/db/conf/icaai/icaai2008.htmlWilliamsL08>

Witten, I. H., Frank, E. and Hall, M. A. (2011). *Data Mining: Practical Machine Learning Tools and Techniques*, 3rd edn, Morgan Kaufmann Publishers Inc., San Francisco, CA, USA.

A Python Script Flow Chart

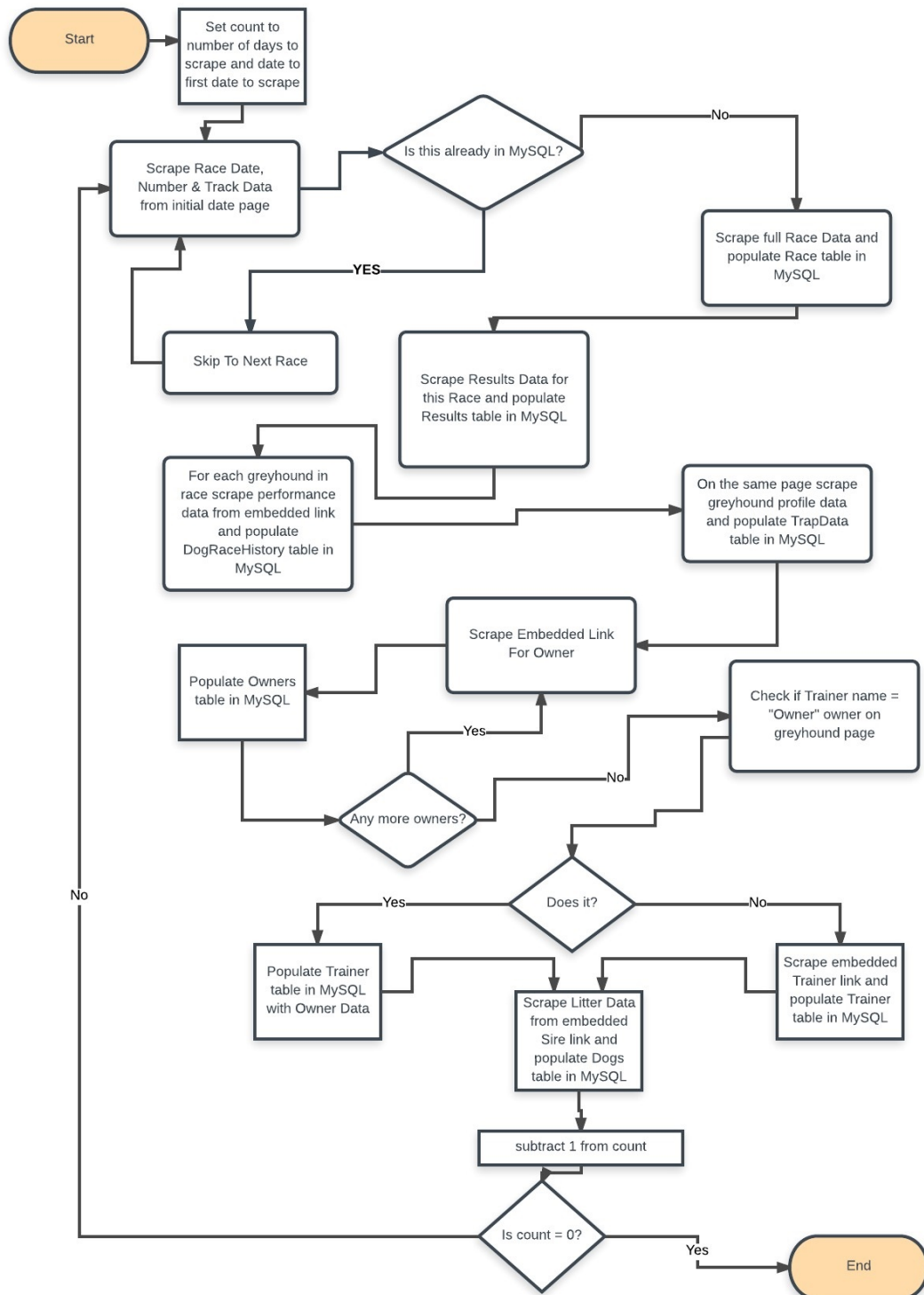


Figure 5: Python Script Flow Chart

B Caret and RandomForest Variable Importance Rankings

Table 5: Variable Importance Results with caret & randomForest packages

Variable Importance in order	
Rank	Variable
1	DogsAge
2	OverallAvgTime
3	SecTimeAvg5
4	TimeAvg5
5	Weight
6	RankedGradeAvg5
7	PrizeMoneyWonAvg5
8	BreakAvg5
9	RaceNumber
10	PrizeMoneyToDate
11	DaysSinceLastRace
12	Avg2ndBend5
13	AvgRemarks5
14	FinishingPositionAvg5
15	Starts
16	TrapNumber
17	ShowPercent
18	RaceGrade
19	PlacedPercent
20	WinPercent
21	Colour
22	RaceDay
23	ShowPercent5
24	PlacedPercent5
25	WinsToDate
26	WinPercent5
27	Going
28	RaceDistance
29	Seed
30	Gender

C RFE with Random Forest Results

Table 6: Optimal Features as Selected by RFE & Random Forest

Feature Importance Rankings in order	
Rank	Variable
1	NumberOfDogs
2	TimeAvg5
3	OverallAvgTime
4	PrizeMoneyWonAvg5
5	DogsAge
6	FinishingPositionAvg5
7	SecTimeAvg5
8	PlacedPercent
9	RankedGradeAvg5
10	Avg2ndBend5
11	ShowPercent
12	RaceGrade
13	AvgRemarks5
14	BreakAvg5
15	WinPercent
16	ShowPercent5
17	PlacedPercent5
18	TrapNumber
19	Weight
20	RaceDay
21	RaceDistance
22	Seed
23	RaceNumber
24	Going

D RFE with Tree Bagging Results

Table 7: Optimal Features as Selected by RFE & Tree Bagging

Feature Importance Rankings in order	
Rank	Variable
1	DogsAge
2	Weight
3	TimeAvg5
4	SecTimeAvg5
5	OverallAvgTime
6	RaceNumber
7	PrizeMoneyToDate
8	BreakAvg5
9	RankedGradeAvg5
10	TrapNumber
11	PrizeMoneyWonAvg5
12	AvgRemarks5
13	RaceGrade
14	Avg2ndBend5
15	DaysSinceLastRace
16	FinishingPositionAvg5
17	WinPercent
18	PlacedPercent
19	ShowPercent
20	Starts
21	Colour
22	RaceDay
23	ShowPercent5
24	PlacedPercent5
25	RaceDistance
26	Going
27	Seed
28	WinPercent5
29	WinsToDate
30	Gender
31	NumberOfDogs