

# **Economic Analysis of Horseracing Betting Markets**

**Chi Zhang**

Doctor of Philosophy

University of York

Economics

January 2017

## **Abstract**

This thesis presents both empirical and theoretical studies on horseracing betting markets. The first two chapters mainly deal with the insider trading problem in the betting markets based on the Shin model (1993) and its extension (Jullien and Salanié (1994)) by employing a novel data set from Yorkshire racecourses during the 2013-2014 racing season. Apart from measuring the incidence of insider trading, we empirically test market efficiency. Our result demonstrates that the degree of insider trading based on the original Shin measure is slightly lower than the calculation based on its extension. We also find no evidence to confirm that the market is strongly efficient. The next chapter studies price-determining factors that affect the starting prices in the racing markets by utilising a unique cross-sectional and time series data set. We find strong evidence to suggest that the winning potential, the age of the horse, the weight the horse carries and the distance of the race are very significant factors in explaining the starting prices. Our findings also confirm that the condition of the turf, the size of the racecourse and the classification of the race have influences on the price. In the last chapter, we propose a theoretical model of how betting odds are adjusted by bookmakers in betting markets. We introduce the optimal stopping techniques into the betting literature for the first time through a two-horse simple benchmark model with both informed and uninformed noise punters. Our main finding shows that increased fraction of informed traders will initially lift the loss per trade to the bookmaker, but after reaching a certain point the loss declines. We also find out that as the fraction of noise traders goes up, the loss is incurred to learn, but the learning process is less informative and the costs are the same, so the decision of changing the prices for each horse is taken sooner.

# Contents

<b>Abstract</b>	<b>ii</b>
<b>Contents</b>	<b>iii</b>
<b>List of Figures</b>	<b>vi</b>
<b>List of Tables</b>	<b>viii</b>
<b>Acknowledgements</b>	<b>x</b>
<b>Declaration</b>	<b>xi</b>
<b>1 Introduction</b>	<b>1</b>
1.1 History of British Horse Racing . . . . .	1
1.2 Motivation . . . . .	3
1.3 Outline of the Thesis . . . . .	5
1.4 Appendix: Some Mathematical Preliminaries . . . . .	7
<b>2 An Examination of Insider Trading in Yorkshire Horseracing Betting Markets</b>	<b>14</b>
2.1 Introduction . . . . .	14
2.2 Literature Review . . . . .	16

2.3	The Betting System in Britain . . . . .	18
2.4	Empirical Analysis . . . . .	19
2.4.1	Data Sources . . . . .	19
2.4.2	Empirical Results . . . . .	23
2.5	Extension of the Shin Measure . . . . .	26
2.6	Conclusions . . . . .	34
2.7	Appendices . . . . .	35
<b>3</b>	<b>A Further Examination of Yorkshire Horseracing Betting Markets</b>	<b>50</b>
3.1	Introduction . . . . .	50
3.2	Empirical Analysis . . . . .	52
3.2.1	Data Analysis . . . . .	52
3.2.2	Test . . . . .	55
3.3	Conclusions . . . . .	68
<b>4</b>	<b>On the Determinants of Starting Prices in Horseracing Betting Markets</b>	<b>69</b>
4.1	Introduction . . . . .	69
4.2	Empirical Model . . . . .	72
4.2.1	Basic Fixed-Effects Model . . . . .	72
4.2.2	Driscoll and Kraay Standard Errors for Pooled OLS Estimation	74
4.2.3	Fixed-Effects Model with Driscoll and Kraay Standard Errors	76
4.3	Data Analysis and Empirical Results . . . . .	76
4.3.1	Data Description . . . . .	76
4.3.2	Variable Classification . . . . .	82
4.3.3	Statistic Tests . . . . .	89

4.3.4	Estimation Results . . . . .	91
4.4	Conclusions . . . . .	102
<b>5</b>	<b>Dynamic Pricing in Horseracing Betting Markets</b>	<b>103</b>
5.1	Introduction . . . . .	103
5.2	Literature Review . . . . .	107
5.3	General Model . . . . .	109
5.4	The Optimal Stopping Problem and Main Result . . . . .	117
5.5	Comparative Statics . . . . .	128
5.6	Conclusions . . . . .	135
5.7	Appendices . . . . .	136
<b>6</b>	<b>Concluding Remarks</b>	<b>144</b>
6.1	Conclusions . . . . .	144
6.2	Future Research . . . . .	147
	<b>Betting Glossary</b>	<b>149</b>
	<b>Bibliography</b>	<b>152</b>

# List of Figures

1.1	Sample Paths of an Arithmetic Brownian Motion . . . . .	12
2.1	Map of Yorkshire with Racecourses . . . . .	21
2.2	Distribution of Sum of Prices . . . . .	22
2.3	Distribution of Number of Runners . . . . .	22
2.4	Distribution of ZSP (2013) . . . . .	29
2.5	Distribution of ZOP (2013) . . . . .	29
2.6	ZSP Overlap ZOP (2013) . . . . .	30
2.7	Comparison between ZSP and ZOP (2013) . . . . .	31
2.8	The Density of the Difference between ZSP and ZOP (2013) . . . .	33
2.9	Example of the Original Data . . . . .	40
3.1	Distribution of ZSP (2014) . . . . .	57
3.2	Distribution of ZOP (2014) . . . . .	57
3.3	ZSP Overlap ZOP (2014) . . . . .	58
3.4	The Density of the Difference between ZSP and ZOP (2014) . . . .	58
3.5	Winning Starting Prices versus Favourite Starting Prices . . . . .	62
4.1	Scatter Plot of Horses Winning First Place . . . . .	79
4.2	Scatter Plot of Horses Winning Second Place . . . . .	80

4.3	Scatter Plot of Horses Winning Third Place . . . . .	80
4.4	£1 Level Stake . . . . .	81
4.5	Winning Percentage . . . . .	81
4.6	Winning Percentage vs. Level Stake . . . . .	82
4.7	Starting Prices for Horses 1 to 25 . . . . .	83
4.8	Starting Prices for Horses 26 to 50 . . . . .	84
4.9	Starting Prices for Horses 51 to 75 . . . . .	84
4.10	Starting Prices for Horses 76 to 100 . . . . .	85
5.1	Sample Path of an Arithmetic Brownian Motion . . . . .	115
5.2	The Loss Function . . . . .	129
5.3	Comparative Statics for $\mu$ . . . . .	132
5.4	Comparative Statics for $\sigma$ . . . . .	133
5.5	Comparative Statics for $\lambda$ . . . . .	134
5.6	The Loss Function for Different Values of $p_L$ . . . . .	139
5.7	The Loss Function for Different Values of $p_H$ . . . . .	140

# List of Tables

2.1	Distribution of Prices . . . . .	21
2.2	OLS Estimates of $z$ . . . . .	25
2.3	Robust Test of $z$ . . . . .	25
2.4	$z$ for Each Racecourse at Starting Prices (2013) . . . . .	28
2.5	$z$ for Each Racecourse at Opening Prices (2013) . . . . .	28
2.6	Summary Statistics for SP and OP . . . . .	33
2.7	Hypothesis Test Result . . . . .	34
3.1	Movers and Rates of Return (MOP) . . . . .	54
3.2	Movers and Rates of Return (Crafts' Ratio) . . . . .	54
3.3	$z$ for Each Racecourse at Starting Prices (2014) . . . . .	56
3.4	$z$ for Each Racecourse at Opening Prices (2014) . . . . .	56
3.5	Hypothesis Test Result . . . . .	59
3.6	Statistic Summary: Winning Starting Prices versus Favourite Starting Prices . . . . .	61
3.7	Regression Results - Winning Starting Prices versus Favourite Starting Prices . . . . .	62
3.8	Regression Results - Average Rate of Return at Starting Prices . . . . .	65
3.9	Regression Results - Average Rate of Return at Opening Prices . . . . .	66



3.10	Regression Results Including <i>MOP</i> . . . . .	67
4.1	Summary of Racecourses . . . . .	77
4.2	Summary Statistics on 100 Horses . . . . .	79
4.3	Variables in the Empirical Model . . . . .	88
4.4	Standard Hausman Test . . . . .	90
4.5	Fixed-effects Regression Results . . . . .	95
4.6	Fixed-effects Regression Results (Including Going Conditions) . . . . .	97
4.7	Fixed-effects Regression Results (Including British Racecourses) . . . . .	98
4.8	Fixed-effects Regression Results (Including Irish Racecourses) . . . . .	99
4.9	Fixed-effects Regression Results (Including Classification) . . . . .	100
4.10	Fixed-effects Regression Results (All) . . . . .	101
5.1	The Loss of Two Types of Errors . . . . .	121
5.2	Parameters for a Base-case Numerical Example . . . . .	129

## Acknowledgements

I am truly grateful for both of my supervisors - Professor Zaifu Yang and Professor Jacco Thijssen - for their patient guidance, practical suggestions and stimulating discussions during the past four years. They inspired me in so many ways and opened a door to the academic area. Not only did they provide professional advice, we also had a lot of fun during the lunch break. Apart from Zaifu and Jacco, my thesis advisory panel member Dr. Michael Thornton has been very helpful during my study. I thank him for his valuable comments and suggestions on my work.

I also owe a debt of gratitude to all my colleagues and friends. Without their generous help, I could not finish this big project. I would like to thank in no particular order: Xueqi Dong, Jiawen Li, Yu Wang, Jiayi Shi. Thanks for their company and kindness.

Last but not least, I would like to dedicate this thesis to my parents who love me selflessly. I thank them for supporting my decisions to pursuing a PhD degree unconditionally, for being patient with my sometimes irrational nonsense, and for always keeping me on the right path. To them, I cannot express my gratitude enough.

## **Declaration**

I, Chi Zhang, declare that this thesis titled, “Economic Analysis of Horseracing Betting Markets”, is a presentation of original work and I am the sole author. This work has not been presented for an award at this, or any other, University. All sources are acknowledged as References.

*This thesis is dedicated to my loving parents.*

# Chapter 1

## Introduction

*“It were not best that we should all think alike;  
it is difference of opinion that makes horse-races.”*

Mark Twain,

*Pudd’nhead Wilson’s Calendar, 1894.*

### 1.1 History of British Horse Racing

Sports betting is an ancient human activity prevailed in the United Kingdom for hundreds of years, which now contributes multi-billion pounds to the economy and has become a large and thriving industry. Horse racing is the second biggest spectator sport after football with a long history dating back to the Roman era. Racing has been part of Britain’s national heritage for ages and still serves as an important everyday pastime.

It was the Roman soldiers who brought the first race to take place in Yorkshire, Britain around 200 AD. Before that racing was only popular in Egypt, Syria and Ancient Greece. By the 9th/10th century horse racing had become quite popular in the UK. During the reign of Henry VIII in the 16th century, he passed breeding laws as to the breeding of horses and imported a large number of stallions and mares. That is when the formal race gathering began to be instigated and a trophy started to be presented to the winner of a race. Kiplingcotes Derby is widely known as the

world's oldest horse race in England since 1519. Newmarket, discovered by James I in 1605, has been generally considered the birthplace of horse race since then. Race meetings began to spring up elsewhere in the country and jockey weights were rigorously enforced in the 17th century. Gatherley, Yorkshire, Croydon and Theobalds on Enfield Chase were the major places for races to run for silver bells. Hundred years later, in the early 18th century Queen Anne set eyes on Ascot where the opening race nowadays at Royal Ascot is still called the Queen Anne Stakes to commemorate her. Between these two periods, Oliver Cromwell in 1654 banned all horse racing, along with other gambling activity that the public enjoyed. However, Charles II restored racing as soon as he claimed the throne. He also introduced the Newmarket Town Plate and wrote the rules himself<sup>1</sup>.

*“Articles ordered by His Majestie to be observed by all persons that put in horses to ride for the Plate, the new round heat at Newmarket set out on the first day of October, 1664, in the 16th year of our Sovereign Lord King Charles II, which Plate is to be rid for yearly, the second Thursday in October for ever.”*

King Charles II,  
*Rules of the Newmarket Town Plate.*

In 1740, the Parliament tried to restrain and prevent the excessive increase in horse racing by introducing an act, but it did not go well. So in 1750, the Jockey Club was founded to implement the rules of racing. In 1993, the British Horseracing Board (BHB) became responsible for race planning, training, financing and marketing, which removed the governance role of the Jockey Club, but it can still regulate the sport. In 2006, the Horseracing Regulatory Authority (HRA) was formed to carry out the regulatory process. In July 2007, after the merger of the HRA and the BHB, the British Horseracing Authority was formed to “provide the most compelling and attractive racing in the world; be seen as the world leader in race-day regulation; ensure the highest standards for the sport and participants, on and away from the

---

<sup>1</sup>The history of Great British racing. Accessible via: <http://www.greatbritishracing.com/about-great-british-racing/the-history-of-great-british-racing>.

racecourse; promote the best for the racehorse; and represent and promote the sport and the industry”<sup>2</sup>.

Horse racing is still considered as the Sport of Kings. The current Queen of England Elizabeth II has bred and owned horses to compete at Royal Ascot and some classic races. However, as the internet can be easily accessed, it is more accessible to the public nowadays. It has become a pastime activity to the masses - from rich to poor, lord to civilian, and professional to laymen - through bookmakers in the ring, betting shops in the street, newspapers, televisions and betting websites.

## **1.2 Motivation**

Economic impact measures the benefit of an event to the economy. In other words, consumers are willing to spend money now rather than keep it, which based on the theory that a pound flows into the market. The impact of sports economics not only incurs large sums of cash to the gambling market but also generates a positive economic outcome for the local economy in terms of facilities, tax incomes, jobs and infrastructure. In the literature, the seminal work of Simon Rottenberg (1956), which studied the restrictions on baseball players in the labour market, laid the foundation of sports economics. Over the last 30 years, the literature has thrived and flourished, especially in the 1990s. The earlier research has a strong preference on the major National American leagues. Most researchers focus on the data, which are collected from the National Football League (NFL), the National Basketball Association (NBA) and the Major League Baseball (MLB) and so on. These markets provide a natural laboratory for economists<sup>3</sup>. Horse racing plays an important role in the British history. Even to this day, it is still the second largest spectator sports after football in the United Kingdom. The revenues generated from the racing are a significant component of the economy.

---

<sup>2</sup>British Horseracing Authority. Accessible via: [https://en.wikipedia.org/wiki/British\\_Horseracing\\_Authority](https://en.wikipedia.org/wiki/British_Horseracing_Authority).

<sup>3</sup>Cooke. A. Sports Economics. Accessible via: <http://www.studyingeconomics.ac.uk/module-options/sports-economics/>.

According to the latest report on the economic impact of British racing by De-Loitte commissioned by the British Horseracing Authority (BHA) in 2013<sup>4</sup>, horseracing was worth about £1.1 billion to the UK economy - £3.45 billion if all direct and indirect expenditure of racing are included and contributed more than £275 million to the government in tax and over £1.4 billion in total in the last five years. There were 5.6 million attendance at 1369 fixtures over 60 racecourses and 17500 individual runners in 2012. As “[r]acing’s economic impact is estimated based on the direct expenditure of its participants and the associated expenditure of racing consumers”, it is obvious that horseracing markets alone play an important role in the British economy.

A majority literature on the horse racing betting markets focuses on market efficiency (see Vaughan Williams (2005) for a survey. Parimutuel: Ali (1977), Snyder (1978), Figlewski (1979), Hausch *et al.* (1981), Asch *et al.* (1984, 1986); Fixed odds: Dowie (1976), Crafts (1985), Cain *et al.* (1990), Gabriel and Marsden (1990)), famous favourite - longshot bias (FL bias; see Griffith (1949), McGlothlin (1956), Fabricand (1965), Ali (1977), Hausch, Ziemba, and Rubinstein (1981), Asch, Malkiel, and Quandt (1982), Henery (1985), Ziemba and Hausch (1984), Brown, D’Amato and Gertner (1994), and Shing and Koch (2008)) and insider trading (Shin (1991, 1992, 1993), Jullien and Salanié (1994), Fingleton and Waldron (1999)). As a betting market resembles a financial market in many ways, we utilise this market to study several topics related to insider trading problems, market efficiency theory, pricing factors and odds-setting procedure.

As mentioned before, a betting market is a simple example of a financial market, for instance, a large number of investors (bettors), extensive market information and ease of entry. It also offers a chance to study decision-making under uncertainty and risk. This market also raises the informational issues. The major difference is that there is a well-defined termination point at which the value of the asset is fixed at the end of the betting period. The added advantage is that we are no longer

---

<sup>4</sup>Economic Impact of British Racing 2013. Accessible via: <http://www.britishhorseracing.com/wp-content/uploads/2014/03/EconomicImpactStudy2013.pdf>.



bounded by future cash flow and net present value. Based on the above narrative, we pay our attention to this particular market.

### **1.3 Outline of the Thesis**

This thesis contributes to the horse racing betting literature by utilising both empirical and theoretical methods to analyse the following problems: insider trading, market efficiency, determinants of the betting odds and odds-setting mechanism.

In Chapter 2 we examine the insider trading problem in the Yorkshire horse racing betting markets. We employ a large new data set, collecting from Yorkshire racecourses during the 2013 racing season, to estimate the incidence of insider trading. Two methods are included. Based upon the Shin (1993) measure, our results suggest that the incidence of insider trading in Yorkshire in 2013 is around 1.7 percent and there is a strong positive correlation between the overall margin implicit in bookmakers' odds and the number of runners in a race. Based on the extended Shin measure, the weighted average degree of insider activity at starting prices is around 2.103 percent, which is slightly higher than that measured by the original method, and at opening prices, it is around 1.64 percent. In Chapter 3, we recollect the data from 9 Yorkshire racecourses during the 2014 racing season. We perform some tests to consolidate our results in the previous chapter. First of all, we apply the extension of the Shin measure to obtain the average degree of insider trading at starting and opening prices respectively. The degree is around 2.07 percent at starting prices and 1.65 percent at opening prices. They are pretty much the same when we use the 2013 data. Secondly, we find supportive evidence to confirm that insiders tend to place a bet at the early stage, which is consistent with the Crafts' (1985) hypothesis. Thirdly, the average rates of return at both starting and opening prices are significantly affected by the incidence of insider trading. Last but not least, we show that a gambler is doomed to lose even if he constantly bets on the favourite horse, and the average loss is around 0.3314 per stake.

Chapter 4 analyses the determinants that have influences on the starting prices of horses by utilising a unique cross-sectional and time series data set. A fixed-effects (within) regression model is estimated on a database of 100 horses and the corresponding results have partially identified the factors that can be used to explain the price setting mechanism. There is a plethora of empirical evidence to suggest that the winning potential, the age of the horse, the weight the horse carries and the distance of the race are very significant factors in determining the price of the horse. Our findings confirm that the condition of the turf, the size of the racecourse and the classification of the race affect the price as well. All these factors are quite persistent based on the estimation results.

In Chapter 5, we present a theoretical model of how betting odds are adjusted by bookmakers in betting markets. We introduce the optimal stopping techniques into the betting literature for the first time through a tractable two-horse setting with both informed and uninformed noise punters. A risk-neutral bookmaker who selects a stopping time decides when to adjust the betting odds for each horse in a race when the cumulative sales of tickets on horse A, net of the cumulative sales of tickets on horse B are modelled by an arithmetic Brownian motion. A costly learning process discloses what information the informed trader possesses. With sequential hypothesis testing, the bookmaker can declare that one of the two hypotheses is true with reasonable certainty and therefore changes the odds correspondingly. Our main finding shows that the increased fraction of informed traders will initially lift the loss per trade to the bookmaker, but after reaching a certain point the loss declines. One explanation is that the more informed traders, there is more information the bookmaker can get per time period, which leads to a decision being taken sooner. So the decreased loss should be expected, which is consistent with our assumption as well as our intuition about this market. We also find out that as the fraction of noise traders goes up, the loss is incurred to learn, but the learning process is less informative and costs are the same, so the decision of changing the prices for each horse is taken sooner.

## 1.4 Appendix: Some Mathematical Preliminaries

### Weierstrass Approximation Theorem

**Theorem 1.4.1.** (See Estep (2002)) Assume that  $f$  is continuous on a closed bounded interval  $I$ . Given any  $\epsilon > 0$ , there is a polynomial  $P_n$  with sufficiently high degree  $n$  such that

$$|f(x) - P_n(x)| < \epsilon \text{ for } a \leq x \leq b.$$

The following section is mainly for Chapter 5.

### Probability Spaces and Stochastic Processes

In this subsection we briefly review some basic concepts from general probability theory and the definition of stochastic processes.

**Definition 1.4.2.** (See Øksendal (2005)) If  $\Omega$  is a given set, then a  $\sigma$ -algebra  $\mathcal{F}$  on  $\Omega$  is a family  $\mathcal{F}$  of subsets of  $\Omega$  with the following properties:

(i)  $\emptyset \in \mathcal{F}$

(ii)  $F \in \mathcal{F} \Rightarrow F^C \in \mathcal{F}$ , where  $F^C = \Omega \setminus F$  is the complement of  $F$  in  $\Omega$

(iii)  $A_1, A_2, \dots \in \mathcal{F} \Rightarrow A := \cup_{i=1}^{\infty} A_i \in \mathcal{F}$

The pair  $(\Omega, \mathcal{F})$  is called a measurable space. A probability measure  $P$  on a measurable space  $(\Omega, \mathcal{F})$  is a function  $P: \mathcal{F} \rightarrow [0, 1]$  such that

(a)  $P(\emptyset) = 0, P(\Omega) = 1$

(b) if  $A_1, A_2, \dots \in \mathcal{F}$  and  $\{A_i\}_{i=1}^{\infty}$  is disjoint (i.e.  $A_i \cap A_j = \emptyset$  if  $i \neq j$ ) then

$$P(\cup_{i=1}^{\infty} A_i) = \sum_{i=1}^{\infty} P(A_i).$$

The triple  $(\Omega, \mathcal{F}, P)$  is called a probability space. It is called a complete probability space if  $\mathcal{F}$  contains all subsets  $G$  of  $\Omega$  with  $P$ -outer measure zero, i.e. with

$$P^*(G) := \inf \{P(F); F \in \mathcal{F}, G \subset F\} = 0.$$

The subsets  $F$  of  $\Omega$  which belong to  $\mathcal{F}$  are called  $\mathcal{F}$ -measurable sets. In a probability context these sets are called *events* and we use the interpretation

$$P(F) = \text{“the probability that the event } F \text{ occurs”}.$$

In particular, if  $P(F) = 1$  we say that “ $F$  occurs with probability 1”, or “almost surely (a.s.)”.

Given any family  $\mathcal{U}$  of subsets of  $\Omega$  there is a smallest  $\sigma$ -algebra  $\mathcal{H}_{\mathcal{U}}$  containing  $\mathcal{U}$ , namely

$$\mathcal{H}_{\mathcal{U}} = \bigcap \{ \mathcal{H}; \mathcal{H} \text{ } \sigma\text{-algebra of } \Omega, \mathcal{U} \subset \mathcal{H} \}$$

Thus, a collection  $\mathcal{U}$  of all open subsets of a topological space  $\Omega$ , then  $\mathcal{B} = \mathcal{H}_{\mathcal{U}}$  is called the *Borel  $\sigma$ -algebra* on  $\Omega$  and the elements  $B \in \mathcal{B}$  are called *Borel sets*.  $\mathcal{B}$  contains all open sets, all closed sets, all countable unions of closed sets, all countable intersections of such countable unions etc.

**Definition 1.4.3** (Borel function). (See Capinski and Kopp (2004)) For any interval  $I \in \mathbb{R}$ , if all the sets

$$f^{-1}(I) \in \mathcal{B}$$

we say that  $f$  is a Borel function.

Suppose that  $P$  and  $Q$  are two probability measures on  $(\Omega, \mathcal{F})$ . Then we define the following theorem.

**Theorem 1.4.4** (Radon-Nikodym Theorem). Suppose  $Q$  is absolutely continuous with respect to  $P$  (i.e.  $P \sim Q \iff [\forall A \in \mathcal{F}, P(A) = 0 \iff Q(A) = 0]$ ). Then there exists a random variable  $f$  such that

$$Q(F) = \int_F f dP, \forall F \in \mathcal{F}$$

The function  $f$  is called the Radon-Nikodym derivative of  $Q$  with respect to  $P$ . This can be written as

$$f(\omega) = \frac{dQ}{dP}(\omega)$$

The Radon-Nikodym Theorem shows how to change from one probability measure to another.

**Definition 1.4.5** (Conditional Expectations). (See Capinski and Kopp (2004)) For an integrable random variable  $\xi$  on a probability space  $(\Omega, \mathcal{F}, P)$  and an event  $B \in \mathcal{B}$  such that  $P(B) \neq 0$  the conditional expectations of  $\xi$  given  $B$  is defined by

$$\mathbb{E}(\xi | B) = \frac{1}{P(B)} \int_B \xi dP.$$

**Definition 1.4.6.** (See Øksendal (2005)) A stochastic process is a parametrized collection of random variables

$$\{X_t\}_{t \in T}$$

defined on a probability space  $(\Omega, \mathcal{F}, P)$  and assuming values in  $\mathbf{R}^n$ .

The following subsections are based on Thijssen (2013).

## Poisson Process

Poisson process is a simple example of stochastic process. Let  $(T_0, T_1, \dots)$  be a strictly increasing sequence of random variables with  $T_0 = 0$ . Define the indicator function

$$I_{t \geq T_n} = \begin{cases} 1 & \text{if } t \geq T_n(\omega) \\ 0 & \text{if } t < T_n(\omega) \end{cases}$$

The  $T_n$ -s describe the times at which the events happen.

**Definition 1.4.7.** The counting process associated with the sequence  $(T_0, T_1, \dots)$  is the process  $(N_t)_{t \geq 0}$  defined by  $N_t = \sum_{n \geq 0} I_{t \geq T_n}$ .

Note that  $(N_t)_{t \geq 0}$  takes values in  $\{0, 1, 2, \dots\}$ . Let  $T := \sup_n T_n$ . Thus  $(N_t)_{t \geq 0}$  is a counting process *without explosions* if  $T = \infty$ ,  $P$ -a.s. In such cases we never encounter sample paths with infinitely many defaults.

**Definition 1.4.8.** *A counting process without explosions is a Poisson process if*

(i)  $N_t - N_s$  is independent of  $N_s$ , for all  $s < t$ ;

(ii) for any  $s < t$  and  $u < v$ , with  $t - s = v - u$  it holds that  $N_t - N_s$  have the same distribution.

The following can now be shown.

**Theorem 1.4.9.** *Let  $(N_t)_{t \geq 0}$  be a Poisson process. Then there exists a  $\lambda \geq 0$  such that  $N_t \sim \text{Poiss}(\lambda t)$  for all  $t \geq 0$ .*

This theorem explains why we often say the  $(N_t)_{t \geq 0}$  is a Poisson process “with parameter  $\lambda$ ”. In differential form we can write

$$dN_t = \begin{cases} 1 & \text{w.p. } \lambda dt \\ 0 & \text{w.p. } 1 - \lambda dt \end{cases}$$

It can also be shown that the inter-arrival times between jumps is exponentially distributed with parameter  $1/\lambda$ .

## Binomial Tree

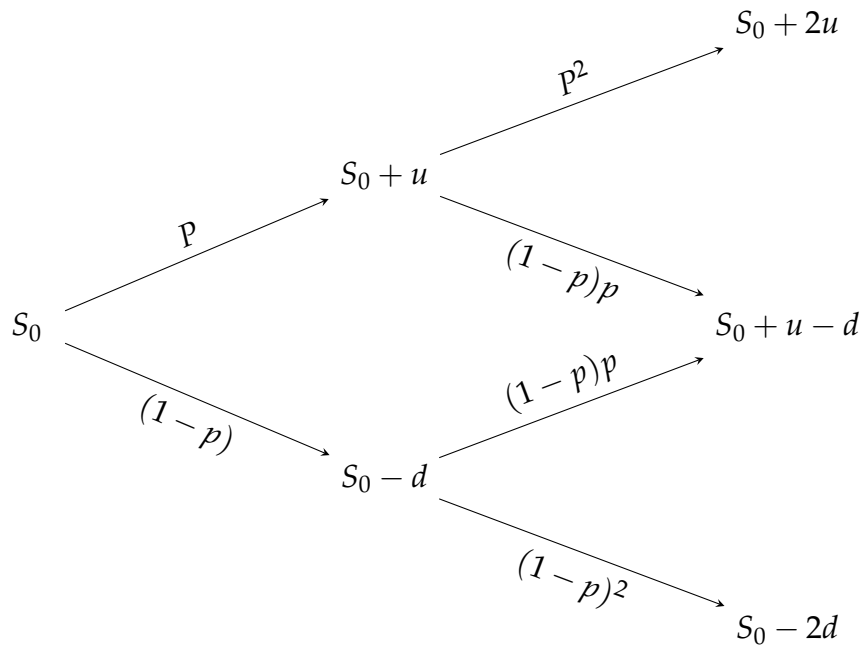
Suppose a coin is flipped infinitely many times, and if Heads (H) shows up we win  $u$ , and losing  $d$  if Tails (T) comes up. Assume that all coin flips are independent and that the probability of heads is  $0 < p < 1$ . Let  $n_H$  denote the number of Heads in  $n$  flips, so the probability of observing  $n_H$  Heads in  $n$  coin flips is

$$p(n_H) = p^{n_H} (1 - p)^{n - n_H}$$

Let  $X_1, X_2, \dots$  be a sequence of random variables, which indicate your winning for each flip

$$X_i = \begin{cases} u & \text{if Heads} \\ -d & \text{if Tails} \end{cases}$$

Total gains after the  $n$ -th coin flip are equal to  $S_n := \sum_{i=1}^n X_i$ . The sequence  $S_0, S_1, S_2, \dots$  describe a stochastic process, where  $S_0$  is the initial wealth. The evolution of  $S_n$  can be depicted in a tree diagram as below.



**Theorem 1.4.10** (Central Limit Theorem). Let  $X_1, X_2, \dots$  be a sequence of independent and identically distributed (i.i.d.) random variables with mean  $\mu$  and variance  $\sigma^2$ . Let  $\bar{X}_n := \frac{1}{n} \sum_{i=1}^n X_i$ . Then

$$\sqrt{n} \frac{\bar{X}_n - \mu}{\sigma} \xrightarrow{d} \mathcal{N}(0, 1), \text{ as } n \rightarrow \infty.$$

## Arithmetic Brownian Motion (ABM)

**Definition 1.4.11** (Brownian motion). A process  $(B_t)_{t \geq 0}$  taking values in  $\mathbf{R}$  is a Brownian motion if

(i)  $B_t - B_s$  is independent of  $B_s$ , for all  $0 \leq s < t$ ;

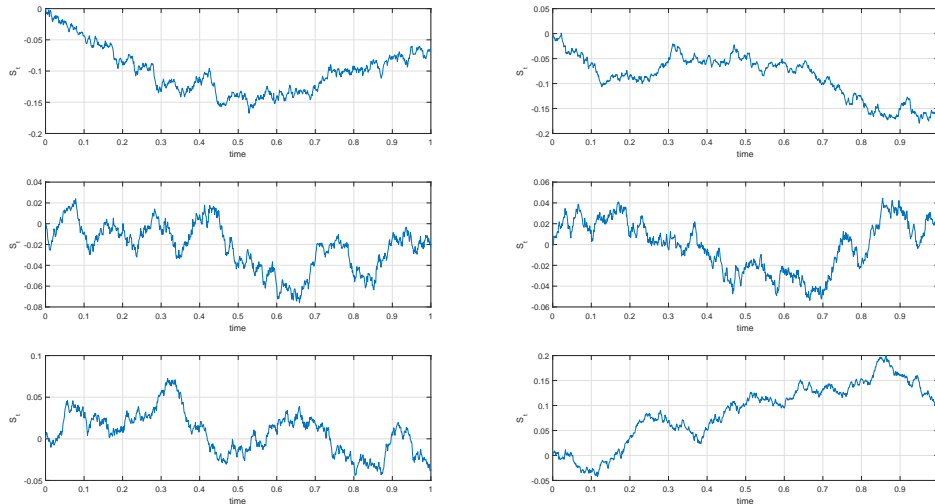
(ii)  $B_t - B_s \sim \mathcal{N}(0, t - s)$ , for all  $0 \leq s < t$ .

An arithmetic Brownian motion can be derived as the continuous time limit of a binomial tree with  $S_0 = y$ ,  $p = 1/2$ ,  $u = \mu dt + \sigma\sqrt{dt}$  and  $d = \mu dt - \sigma\sqrt{dt}$ , thus we define the ABM as

$$dS_t = \mu dt + \sigma dB_t$$

A few examples of the ABM are depicted in Figure 1.1.

Figure 1.1 Sample Paths of an Arithmetic Brownian Motion



## Stopping Time

A *stopping time* is a random variable  $\tau$  to which the question “has the event  $\{\tau \leq t\}$  occurred?” can be answered at any time  $t \geq 0$ . So for each  $y \in E$  and each  $t \geq 0$ , it holds that  $P_y(\tau \leq t) \in \{0, 1\}$ . Most stopping times depend on the underlying stochastic process  $(S_t)_{t \geq 0}$ . For example, if the firm decides to invest as soon as some pre-determined threshold (or “trigger”)  $S^*$  is reached, then the (random) time



at which investment takes place is the *first hitting time*

$$\tau_y(S^*) = \inf \{t \geq 0 \mid S_t \geq S^*, S_0 = y\}.$$

## Ito's Lemma

In this subsection we first define the following stochastic process,

$$Y_t = y + \int_0^t \mu(s, Y_s) ds + \int_0^t \sigma(s, Y_s) dB_s \quad (1.1)$$

where  $(B_t)_{t \geq 0}$  is a Brownian motion and  $Y_0 = y$ . Processes as in (1.1) are called *Ito Diffusions*. In differential notation we obtain

$$dY_t = \mu(t, Y_t) dt + \sigma(t, Y_t) dB_t \quad (1.2)$$

Equation (1.2) is called a *stochastic differential equation*(SDE), the function  $\mu(\cdot)$  is called the *trend* and  $\sigma(\cdot)$  is called the *volatility*. If the trend and volatility do not depend on  $t$ , we call the diffusion *time homogeneous*. In this thesis we just work with time homogeneous diffusion only.

**Theorem 1.4.12** (Ito's Lemma). *Let  $(Y_t)_{t \geq 0}$  follows an Ito diffusion (1.1) and let  $g(t, Y_t)$  be a twice continuously differentiable function. Then*

$$\begin{aligned} dX &= \frac{\partial g(\cdot)}{\partial t} dt + \frac{\partial g(\cdot)}{\partial Y} dY + \frac{1}{2} \frac{\partial^2 g(\cdot)}{\partial Y^2} dY^2 \\ &= \left[ \frac{\partial g(\cdot)}{\partial t} + \frac{\partial g(\cdot)}{\partial Y} \mu(\cdot) + \frac{1}{2} \frac{\partial^2 g(\cdot)}{\partial Y^2} \sigma(\cdot)^2 \right] dt + \frac{\partial g(\cdot)}{\partial Y} \sigma(\cdot) dB_t \end{aligned}$$

## **Chapter 2**

# **An Examination of Insider Trading in Yorkshire Horseracing Betting Markets**

### **2.1 Introduction**

The pioneering work of Bagehot (1971) shows his great interest in the field of the bid-ask spread in financial markets. Copeland and Galai (1983) then formalise the problem and shed lights on how a market maker tries to optimise the bid-ask spread by maximising the difference between the gain from liquidity traders and the loss to traders with superior information. Subsequent research followed the path to analyse the optimal bid-ask spread. In a seminal paper, Glosten and Milgrom (1985) have explored the idea that the way the specialists setting the bid-ask spread can be an informational phenomenon by assuming that they are risk-neutral and make zero profits. They also find that transaction prices can reflect insider information.

Most recently, researchers have studied both US and UK betting markets regarding the insider trading problem. Furthermore, this market can also be used to test market efficiency theory and herding behaviour and so on. Applying the theory of Glosten and Milgrom, Shin (1991, 1992, 1993) examines insider activity in betting

markets in the United Kingdom. He shows that the spread is increasing with the incidence of insider trading and this insider problem might partially explain market distortion. In the Shin's paper, he proposes an index to measure the incidence of insider trading in the British horse racing betting markets. His setting fits into the British horse racing system well, as in the UK betting odds are determined by bookmakers, in contrast to the system in North America where odds are derived from the parimutuel method.

The problem of insider trading may disturb financial markets in many ways. Bookmakers at racecourses in horse-racing betting markets face the same problem. There is a bunch of punters, some of whom may have inside information on which horse will win in a race, and on the contrary, some outsiders whose preferences are indifferent across all horses. The betting market shares many similarities with the financial market. In particular, in both types of markets, a large number of investors (punters) can easily access to readily cheaply public information before they buy state-contingent claims (betting tickets). In addition, the return on an investment (a bet) is uncertain because the winning horse is unknown at the moment. Therefore anyone who has inside information will earn extra profits. A bookmaker, on the other hand, is an intermediary just like the market maker in financial markets, who is responsible for setting the odds or the prices of tickets in the United Kingdom.

Due to the above characteristics, we aim to examine, in samples of 997 horse races, the degree of insider trading based on the original Shin measure (1993) and its extension by Jullien and Salanié (1994). The other purpose of the paper is to test the existence of strong form efficiency. The data set that we utilise is collected from nine racecourses in Yorkshire during the 2013 racing season. Our findings consolidate the result of the Shin model, and we conclude that the incidence of insider trading is around 1.7%. Based on the extension of the Shin measure, the degree of insider trading at Starting Prices is around 2.103%, while at Opening Prices 1.64%. We also demonstrate that the market is not strongly efficient.

The outline of this chapter is as follows. The related literature is summarised in Section 2.2. Section 2.3 outlines the betting system in the United Kingdom. In

Section 2.4, we analyse data on starting prices collected during the 2013 racing season in Yorkshire, United Kingdom and present the empirical results based on the Shin measure. Section 2.5 contains the empirical results based on the extension of the Shin measure. Section 2.6 delivers the concluding remarks.

## 2.2 Literature Review

The literature on the horse racing betting markets can be divided into two parts. The first part studies the incidence of insider trading and the second part focuses on testing market efficiency.

The incidence of insider trading in the betting market is measured by Shin (1991, 1992, 1993) who portrays a situation where bookmakers purposely raise risk premium to insure themselves against the presence of a certain percentage of informed traders. His model exhibits *favourite-longshot bias* in general, whereby the prices against favourites are relatively high than those against longshots. The other basic feature of the Shin model is that there is a strong positive correlation between the overall margin implicit in bookmakers' odds and the number of runners in a race (Shin (1993) and Fingleton and Waldron (1999)). Jullien and Salanié (1994) firstly revise the Shin's iterative procedure on linearised versions of the equations. They show that a standard nonlinear estimation procedure can be obtained by employing the quadratic formula. Fingleton and Waldron (1999) relax the assumptions in the Shin measure and develop a more general model of how the odds are determined by bookmakers given that the levels of insider trading vary from race to race on the part of punters. This model incorporates the bookmakers' attitudes towards risk and the possibility of anti-competition among them as well. It turns out that bookmakers are extremely risk-averse in the Irish market during the 1993 race season and a combination of monopoly rents and operating costs are taken up 4 percent of turnover and between 3.1 and 3.7 percent of all bets is placed by bettors with inside information. Vaughan Williams and Paton (1997) examines the favourite-longshot bias by using a large new data set, which is collected from

British racetracks. Their empirical tests identify this bias and characterise it as a rational response to bookmakers who face punters with the superior information. This is an adverse selection problem in the betting markets<sup>1</sup>. Law and Peel (2002) attempt to capture herding activity in the horse-race betting market and distinguish the shortening odds over the betting auction period caused by insider trading. Their conclusions suggest that a large plunge from opening to starting prices accompanied by a fall in the Shin measure would be a strong signal of herd behaviour.

Research on market efficiency theory begins with Fama in the early 1960s. A market is considered to be efficient if two conditions are satisfied: (1) prices fully reflect historical prices and public announcements and (2) no special groups or individuals have monopolistic access to some relevant information (inside information) and can achieve higher than normal profits. However, market efficiency *per se* is not testable and must proceed on the models of market equilibrium which can be described in terms of expected returns in the theoretical and empirical work. It has become conventional wisdom that the efficient market can be tested by utilising three subsets of information - weak form, semi-strong form and strong form. Dowie (1976) uses the term “equity” in place of “strong form efficiency” to test the betting market on horse races in Britain. His results include that there is no evidence of the existence of superior inside information holding by a small subset of investors. He is confident that the hypothesis of “strongly inefficient” betting markets is in doubt but supporting the belief that the market is “weakly efficient”, which is more reasonable and acceptable. A final point in the Dowie’s paper shows that even if the “outsider” obtains as good information as the “insider”, it does not mean he can exploit it effectively<sup>2</sup>. Crafts (1985) suggests that the way Dowie tests the “strong inefficient” market is not appropriate. To argue that, Crafts redesigns a test by using the same data source as in Dowie’s and gives evidence to prove that there exist profitable opportunities for insider trading<sup>3</sup>. Snyder (1978) demonstrates how

---

<sup>1</sup>Vaughan Williams, L. and Paton, D. (1997). Why is There a Favourite-Longshot Bias in Britain Racetrack Betting Markets. *The Economic Journal*, Vol. 107, No. 440, 150-158.

<sup>2</sup>Dowie, J. (1976). On the Efficiency and Equity of Betting Markets. *Economica*, Vol. 43, No. 170, 139-150.

<sup>3</sup>Crafts, N. F. R. (1985). Some Evidence of Insider Knowledge in Horse Race Betting in Britain. *Economica*, Vol. 52, No. 207, 295-304.

these three forms can be applied to horse racing, especially when the odds are determined by the parimutuel system. He shows through several weak and strong tests that horse race betting exhibits stable biases of the subjective and empirical probabilities of winning for both general public bettors and experts and we cannot expect an above average rate of return. Moreover, a great bias is reflected among the experts rather than the betting public. Semi-strong tests also shed light on the absence of an efficient market<sup>4</sup>. Gabriel and Marsden (1990, 1991) compare the returns to the bookie's starting prices and to the parimutuel tote bets, based on the data from the 1978 horse race season in Britain. Their analysis indicates that on average the Tote returns to the winning bets are persistently higher than the odds set by bookmakers given that the risk and the payoffs are widely available under both systems. This implies that the British horse-race betting market fails to satisfy the conditions of semi-strong and strong efficiency<sup>5</sup>.

## 2.3 The Betting System in Britain

In order to proceed the following analysis, a brief guide to the horse-racing betting system in Britain is presented here. A typical British punter can wager the money through two different betting mediums: the totalisator (or the "tote") and fixed odds betting. The former one is similar to the parimutuel method in North America. The payoff to a tote bet is based on the weight of money on each horse relative to the total bet in the pool. If a horse wins, "the winning pool is proportioned to those wagering on the winning horse after the takeout (taxes, track cut, owners' and trainers' cut) is removed" (Gabriel and Marsden (1990)). Another way is to place bets with bookmakers at the fixed odds. The odds-setting procedure is quite similar to pricing risky assets. Bets with a bookie can be accepted at specific but changing odds in a race. That is, it can be placed either at offered at the time of the bet or at the starting prices (SP). Bets that made subsequently will not affect

---

<sup>4</sup>Snyder, W. (1978). Horse Racing: Testing the Efficient Markets Model. *The Journal of Finance*, Vol. 33, No. 4, 1109-1118.

<sup>5</sup>Gabriel, P. E. and Marsden, J. R. (1990). An Examination of Market Efficiency in British Racetrack Betting. *Journal of Political Economy*, Vol. 98, No. 4, 874-885.

the return to any individual bet, but the odds *per se* will plunge as the number of bettors is increasing. The starting prices are defined as the odds at which a “sizable” bet could have been made *on the course* just before the race starts (Dowie (1976)). In the on-course market, as described by Gabriel and Marsden (1990), the starting prices are measured as “the average of a set of the largest bookmakers at the racetrack just before the off”. The on-course prices are reported to the off-course market, and in turn, the off-course prices are reflected as information is relayed to the on-course market. Therefore the actual starting price is determined by both off- and on-course behaviour. Forecast odds is conducted by leading sports news agencies (for example, *The Sporting Life*) and utilised by the off-course market. In this chapter, the fractional forms of odds will be converted to the implicit winning probability. A simple example of SP of two to one (2/1) corresponds to a “percentage” of  $1/(1+2)=1/3$ , which gives unit stakes required to yield a payout of 1 if the horse wins in the race. The shorter the odds, the higher the winning probability.

## 2.4 Empirical Analysis

### 2.4.1 Data Sources

The data set in this paper contains 997 races in Yorkshire in the 2013 horse race season, of which 11 races would be ruled out due to the fact that the sum of prices is actually less than one (see Table 2.1). Apparently, this is an anomaly in the pricing system because any price less than one would give a punter a chance to obtain risk-free profits. Note also that 97.693 percent of races falls in the range of 1.0 to 1.3 and the distribution is heavily skewed to the left as described in Figure 2.2. Data on the on-course starting prices (SP) for each race are collected from the website *The Sporting Life*<sup>6</sup> and include all horses (except non-runners) from nine racecourses in Yorkshire on standard race days. These nine racecourses are located in Beverley, Catterick, Doncaster, Pontefract, Redcar, Ripon, Thirsk, Wetherby and York respectively.

---

<sup>6</sup> Accessible via: <http://www.sportinglife.com/racing>.

Racing in Yorkshire has a long history. Yorkshire racing is a micro- miniature of British horseracing not only because it boasts more top racetracks than any other region of the UK, but also because this area contains a number of successful stables and the Northern Racing College. Nine racecourses host 180 days of racing throughout the year, from flat to jump races, from national championship races to most relaxed informality races for family and friends. In the 18th and 19th centuries, Beverley was a main horse training centre. It held flat races between 25th April and 24th September in 2013. Catterick Racecourse is one of the busiest racecourses in North Yorkshire and one of the homes of the Northern racing scene. It holds meetings all year round. Doncaster racecourse is one of the busiest racecourses in the UK holding the National Hunt and Flat races almost every month. Racecourse in Pontefract is also steeped in history with the first meeting in 1801. Our data contain races from 21st May to 21st October 2013. Redcar and Ripon racecourses are not just famous for flat racing, it is also the perfect venue for parties and friends and family gatherings. Thirsk Racecourse is benefited from its location and is renowned for being one of the prettiest racecourses in the country. Wetherby is one of the Country's leading jumping tracks, and it holds flat racing either. York Racecourse holds flat races from May until October each year, offering a world-class horseracing experience. Ebor Festival is the most famous race at York and the richest flat handicap in Europe, which has also been ranked as the best race in the world by the International Federation of Horseracing Authorities in 2014. All in all, Yorkshire is synonymous with the very best in British horseracing<sup>7</sup>. As we see from Figure 2.1, most racecourses adjacent to the A1/M1, which provide convenience for racegoers.

As shown by Table 2.1, there are 11 anomalies that need to be deleted as the sum of starting prices is less than 1. So there are 986 races left and in total 9678 runners. The actual race runs between 2 and 23 horses as depicted in Figure 2.3. The frequency distribution is summarised in terms of the number of runners ( $n$ ) in Figure 2.3. Our regression equation demonstrates a correlation between the sum of

---

<sup>7</sup> Accessible via: <http://goracing.co.uk/>.



Figure 2.1 Map of Yorkshire with Racecourses



Source: Google Map

prices and the number of runners, along with the distribution of the sum of prices, we are confident that this correlation is positive and can be identified.

Table 2.1 Distribution of Prices

<i>Sum of prices</i>	<i>Frequency</i>
< 1.0	<b>11</b>
1.0 ~ 1.1	<b>94</b>
1.1 ~ 1.2	<b>676</b>
1.2 ~ 1.3	<b>204</b>
1.3 ~ 1.4	<b>11</b>
1.4 ~ 1.5	<b>0</b>
1.5 ~ 1.6	<b>0</b>
> 1.6	<b>1</b>

Note: the first row shows that the sum of the starting prices is less than 1 in 11 races, which offers the punters an arbitrage opportunity. This is clearly an anomaly in the pricing system. We delete these anomalies from the sample set. Normally the sum of prices is larger than 1 and the excess part is the bookmaker's margin.

Figure 2.2 Distribution of Sum of Prices

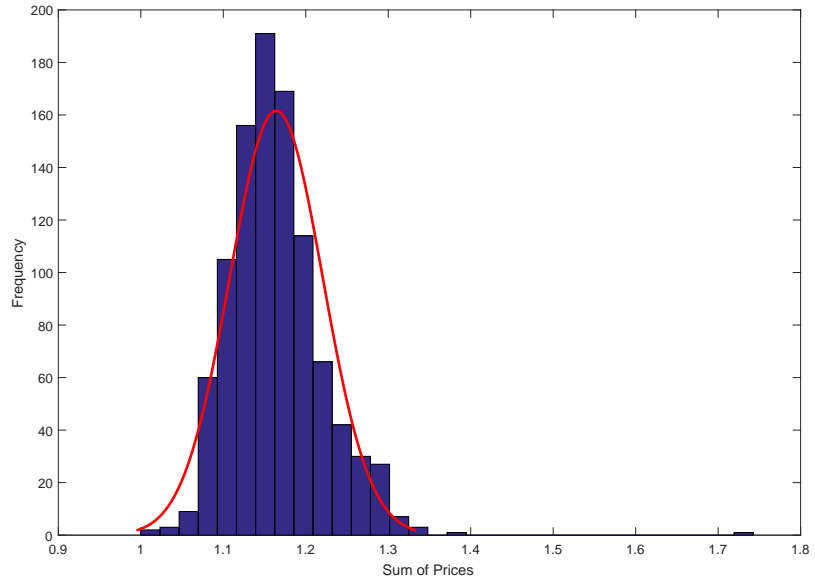
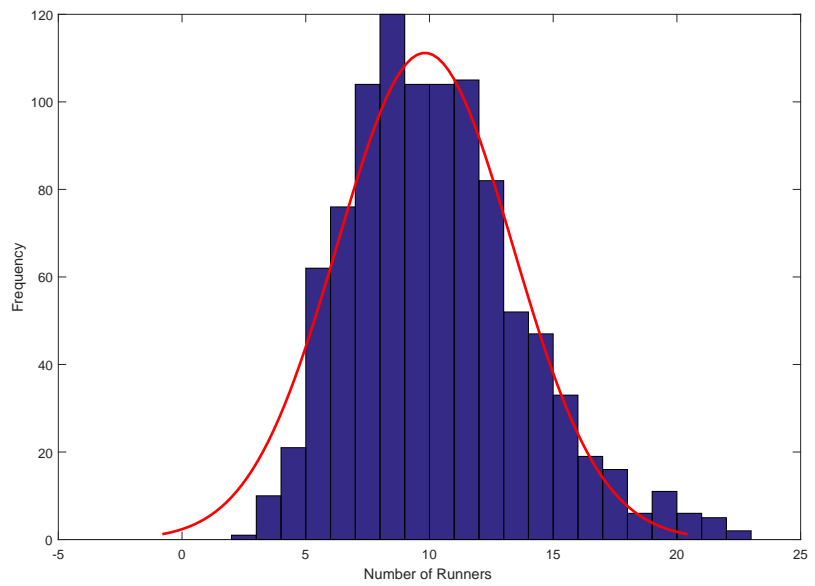


Figure 2.3 Distribution of Number of Runners



## 2.4.2 Empirical Results

By adopting the Shin model (1993) as the starting point, a bookmaker faces a bunch of informed punters and some outsiders. With probability  $z$ , an informed trader is drawn to place a bet. We assume that the insider is allowed to observe what the winning horse is and buy tickets up to £ 1 at the fixed prices. Outsiders are unknown to the identity of the winning horse. If there are  $n$  horses in a race, the  $i$ th horse wins with probability  $p_i$  ( $i = 1, 2, \dots, n$ ). With probability  $(1 - z)p_i$ , the outsider who believes that the horse  $i$  wins is chosen to be the punter. Let  $\mathbf{p}$  denote the vector of winning probabilities  $(p_1, p_2, \dots, p_n)$  and  $\boldsymbol{\pi}$  be the vector of starting prices  $(\pi_1, \pi_2, \dots, \pi_n)$ . With some twist of terminology, let  $Var(\mathbf{p})$  denote the ‘variance’ of  $\mathbf{p}$ , where  $Var(\mathbf{p}) = \frac{1}{n} \sum_i \left(p_i - \frac{1}{n}\right)^2$  and  $Var(\mathbf{p})$  is a scalar.  $Var(\mathbf{p})$ , introduced in the paper, is to help adjust the estimation of the degree of insider trading. How we define  $Var(\mathbf{p})$  will be presented in Appendix A. Suppose that the bookmaker’s winning bid is given by  $\beta$ . Since we assume the bookmaker’s revenue is one pound, the optimisation problem becomes to maximise the expected profits.

$$\begin{aligned} & \text{maximise} \quad 1 - \sum_{i=1}^n \frac{z p_i + (1 - z) p_i^2}{\pi_i} \\ & \text{s.t.} \quad \sum_{i=1}^n \pi_i \leq \beta \text{ and } 0 \leq \pi_i \leq 1 \text{ for all } i = 1, 2, \dots, n \end{aligned}$$

By solving the above problem, we get

$$\pi_i = \frac{\beta \cdot \sqrt{z p_i + (1 - z) p_i^2}}{\sum_{s=1}^n \sqrt{z p_s + (1 - z) p_s^2}} \quad (2.1)$$

and

$$\beta = \sum \pi_s = \left( \sum_s \sqrt{z p_s + (1 - z) p_s^2} \right)^2 \quad (2.2)$$

Substituting (2.2) into (2.1), we then get

$$\pi_i = \sqrt{z p_i + (1 - z) p_i^2} \left( \sum_s \sqrt{z p_s + (1 - z) p_s^2} \right) \quad (2.3)$$

To get our benchmark polynomial regression model, Appendix A provides us with the detailed information. With some twists, we obtain the following model.

$$D = z(n - 1) + \sum_{k=0}^K a_k n^k \text{Var}(\mathbf{p}) + \sum_{k=0}^K b_k n^k [\text{Var}(\mathbf{p})]^2 \quad (2.4)$$

By applying Weierstrass Approximation Theorem, we get

$$\begin{aligned} D = & z(n - 1) + a_0 \text{Var}(\mathbf{p}) + a_1 n \text{Var}(\mathbf{p}) + a_2 n^2 \text{Var}(\mathbf{p}) \\ & + b_0 [\text{Var}(\mathbf{p})]^2 + b_1 n [\text{Var}(\mathbf{p})]^2 + b_2 n^2 [\text{Var}(\mathbf{p})]^2 + \epsilon \end{aligned} \quad (2.5)$$

where  $D$  is the bookmaker's over-round, which equals the sum of starting prices minus one (e.g.  $\sum_s \pi_s - 1$ ).  $a_0, a_1, a_2, b_0, b_1, b_2$  are the slope parameters. Note also that there is no intercept term in this equation, which means if there is no runner in the race, there is no gain for the bookmakers.

The next step is to write the algorithm in MatLab<sup>8</sup> and run the regression. The results are shown on the next page. As we shall see, the estimate of  $z$  is stable after the third adjustment. The initial estimate of  $z$  is from the ordinary least squares regression of (2.5) in which  $\text{Var}(\hat{\boldsymbol{\pi}})$  is used as a proxy  $\text{Var}(\mathbf{p})$ .  $\hat{\pi}_i$  here is the normalised price of horse  $i$ , which equals  $\frac{\pi_i}{\sum_s \pi_s}$ .  $\hat{\boldsymbol{\pi}}$  is the vector of normalised prices. It appears in the first column of Table 2.2. This regression yields an initial estimate of  $z$  of around 1.77%. The estimates of the adjustments of regressors in Equation (2.5) appear in the successive columns of Table 2.2. It is worth noting that all estimates are significant at the 1% level, except for the coefficients of  $\text{Var}(\mathbf{p})$  which are insignificant at any level. The revised estimates of  $z$  converge to 1.7% after two iterations.

---

<sup>8</sup>I am grateful to Xueqi Dong for her help on the MatLab coding process.

Table 2.2 OLS Estimates of z

	<b>No adjust.</b>	<b>1st adjust.</b>	<b>2nd adjust.</b>	<b>3rd adjust.</b>
n-1	0.0177*** (71.494)	0.017*** (71.2831)	0.017*** (70.4653)	0.017*** (70.4705)
$Var(\mathbf{p})$	-0.7283 (-0.9411)	-0.3601 (-0.5715)	-0.3473 (-0.5532)	-0.3473 (-0.5533)
$nVar(\mathbf{p})$	1.0391*** (5.6326)	0.7298*** (4.842)	0.7257*** (4.8384)	0.7257*** (4.8385)
$n^2Var(\mathbf{p})$	-0.0972*** (-8.0785)	-0.057*** (-5.8466)	-0.0569*** (-5.8561)	-0.0569*** (-5.8560)
$[Var(\mathbf{p})]^2$	49.3209*** (6.5459)	24.1781*** (4.7248)	24.1441*** (4.7307)	24.1444*** (4.7306)
$n[Var(\mathbf{p})]^2$	-34.8936*** (-6.6316)	-17.7374*** (-5.3407)	-17.7452*** (-5.3521)	-17.7451*** (-5.3520)
$n^2[Var(\mathbf{p})]^2$	3.9318*** (5.5752)	1.8138*** (4.2184)	1.8197*** (4.234)	1.8197*** (4.2339)
	$R^2 = 0.6702$	$R^2 = 0.6785$	$R^2 = 0.6787$	$R^2 = 0.6787$

1. *t* values are in parentheses.

2. \*\*\* indicates significance at the 1% level.

Note: in this quadratic case, we adjust the regression results three times because the true winning probabilities are unknown. So the normalised price  $\hat{\pi}_s = \pi_s / (\sum_i \pi_i)$  can be our empirical proxy for the true probability. The iterative procedure is as follows. In the first step, we apply the original normalised starting prices to the Equation (2.5) and get the estimation of  $z$ . In the second step, we substitute  $z$  into Equation (2.18) in Appendix A. Then we use revised  $Var(\mathbf{p})$  to re-estimate  $z$ . In the last step, we repeat step 2 until the revised values of  $z$  converge.

Table 2.3 Robust Test of z

	<b>No adjust.</b>	<b>1st adjust.</b>	<b>2nd adjust.</b>	<b>3rd adjust.</b>
<b>Quadratic</b>	0.0177	0.017	0.017	0.017
<b>Cubic</b>	0.0179	0.0169	0.0169	0.0169
<b>Quartic</b>	0.0178	0.0169	0.0169	0.0169
<b>Quintic</b>	0.018	0.017	0.0171	0.0171
<b>Sextic</b>	0.0179	0.0171	0.0171	0.0171

Note: the OLS regression equation is  $D = z(n-1) + \sum_{k=0}^K a_k n^k Var(\mathbf{p}) + \sum_{k=0}^K b_k n^k [Var(\mathbf{p})]^2$ . In the quadratic case,  $K = 2$ . Here we test whether the quadratic case is enough for our research purpose. So we let  $K = 3, 4, 5, 6$ . It turns out that even the sextic case does not change the estimated value of  $z$  too much. In light of this evidence, we conclude that the quadratic specification can be a benchmark case.

Now the robustness of the estimate of  $z$  is presented as we change  $k$ , i.e., the degree of polynomial approximation. For each degree of the polynomial and for each round in the iteration, the estimate of  $z$  is recorded. These estimates of  $z$  are presented in Table 2.3. The columns indicate the number of iterations and the rows indicate the degree of the polynomial in the estimated equation. Notice that after the first iteration, the value of  $z$  converges to the same value and these estimates are virtually identical. In light of this evidence, we conclude that the quadratic case can be a benchmark case and it would appear that we have good reason to believe that the estimate of  $z$  is acceptable.

## 2.5 Extension of the Shin Measure

The formula proposed in this part can be used to calculate  $z$ , which is defined as the incidence of insider trading, for each race at opening and starting prices, respectively. Starting from the following two equations, we can link the probabilities of winning ( $\mathbf{p}$ ), the price of a ticket on each runner  $i$  ( $\pi_i$ ) and  $z$ .

$$\beta = \sum \pi_s = \left( \sum \sqrt{z p_s + (1-z) p_s^2} \right)^2$$

$$\pi_i = \sqrt{z p_i + (1-z) p_i^2} \left( \sum \sqrt{z p_s + (1-z) p_s^2} \right)$$

Rewriting the above two equations as

$$\frac{\pi_i}{\sqrt{\beta}} = \sqrt{z p_i + (1-z) p_i^2} \quad (2.6)$$

This is easily inverted to yield

$$p_i = p_i \left( \frac{\pi_i}{\sqrt{\beta}}, z \right) = \frac{\sqrt{z^2 + 4 \frac{\pi_i^2}{\beta} (1-z)}}{2(1-z)} \quad (2.7)$$

$$\sum_i p_i \left( \frac{\pi_i}{\sqrt{\beta}} z \right) = 1 \quad (2.8)$$

Since the sum of the probabilities must equal one, the estimation of  $z$  is done by solving Equation (2.8). A programme had been written to numerically solve the Shin model for the implied degree of insider trading in each race (see Appendix C). As we see from Table 2.4, for the sample of 986 races in 2013 Yorkshire horse-race season, the weighted average degree of insider trading at starting prices is 2.103%, which is slightly higher than that estimated under the Shin measure (1.7%) with a minimum value of 0.0004863 occurred at Thirsk racecourse and a maximum of 0.0488694 at Doncaster. Notice that we rule out 11 anomalies as the sum of prices is less than one in those races which lead to the negative values. Another advantage of this measure is that the weighted average degree of insider trading at opening odds can easily be calculated following the same way. As can be seen from Table 2.5, for the sample of 982 races, the average value is 1.64%, with a minimum value of 0.000182 at Beverley and a maximum of 0.053218 at Doncaster. At opening prices, 15 out of 997 races have the negative values that need to be removed from our sample set. For comparison, we only keep the races that have the positive degrees, so in total there are 981 races remained. The Shin measure of insider trading increases from opening to starting prices in 905 races, and decreases in the rest 76 races. As depicted by Figures 2.4 and 2.5, the distribution of  $zsp$  is obviously skewed to the left, so is the distribution of  $zop$ . Figures 2.6 and 2.7 describe a fact that the Shin measure of insider trading at opening odds is lower than at starting prices.

Table 2.4 z for Each Racecourse at Starting Prices (2013)

	<b>Mean</b>	<b>S.D</b>	<b>Min</b>	<b>Max</b>	<b>Obs</b>
<b>Beverley</b>	0.0232695	0.0056476	0.0015887	0.0462109	135
<b>Catterick</b>	0.0202013	0.005112	0.011867	0.0415089	125
<b>Doncaster</b>	0.021258	0.0056014	0.0118623	0.0488694	163
<b>Pontefract</b>	0.0206307	0.0059553	0.0047323	0.0455253	86
<b>Redcar</b>	0.0193934	0.0047464	0.012139	0.0313522	91
<b>Ripon</b>	0.0229091	0.0064364	0.0117921	0.0458977	107
<b>Thirsk</b>	0.0194775	0.0057961	0.0004863	0.0400381	97
<b>Wetherby</b>	0.021324	0.0070254	0.0131799	0.0455633	69
<b>York</b>	0.0194856	0.0037759	0.0142865	0.0346617	113

Note: this table summarises the average degree of insider trading for each racecourse. The degree of insider trading for each race is estimated based on the method proposed by Jullien and Salanié (1994). At starting prices, we rule out 11 races that have negative values of  $z$ , which leaves us 986 races for investigation. For comparison, we delete all the races that have a negative degree of insider trading.

Table 2.5 z for Each Racecourse at Opening Prices (2013)

	<b>Mean</b>	<b>S.D</b>	<b>Min</b>	<b>Max</b>	<b>Obs</b>
<b>Beverley</b>	0.0186074	0.0050934	0.000182	0.0371681	135
<b>Catterick</b>	0.0154645	0.0041342	0.0060564	0.0307309	125
<b>Doncaster</b>	0.0167703	0.0056476	0.0084087	0.0532178	163
<b>Pontefract</b>	0.0161782	0.0041872	0.0103552	0.0292508	85
<b>Redcar</b>	0.0150025	0.0032587	0.0088909	0.0235444	91
<b>Ripon</b>	0.0176198	0.0058318	0.0058933	0.0451168	107
<b>Thirsk</b>	0.0149567	0.0039767	0.010168	0.0298376	94
<b>Wetherby</b>	0.0164562	0.0050481	0.0107231	0.0374676	69
<b>York</b>	0.0156787	0.0039582	0.0015884	0.0384567	113

Note: this table summarises the average degree of insider trading for each racecourse. The degree of insider trading for each race is estimated based on the method proposed by Jullien and Salanié (1994). At opening prices, we rule out 15 races that have negative values of  $z$ , which leaves us 982 races for investigation. For comparison, we delete all the races that have a negative degree of insider trading.



Figure 2.4 Distribution of ZSP (2013)

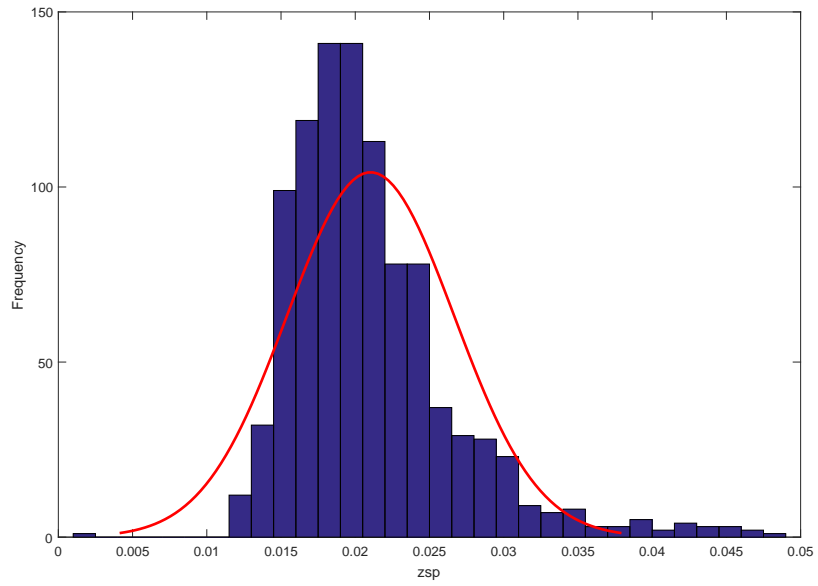


Figure 2.5 Distribution of ZOP (2013)

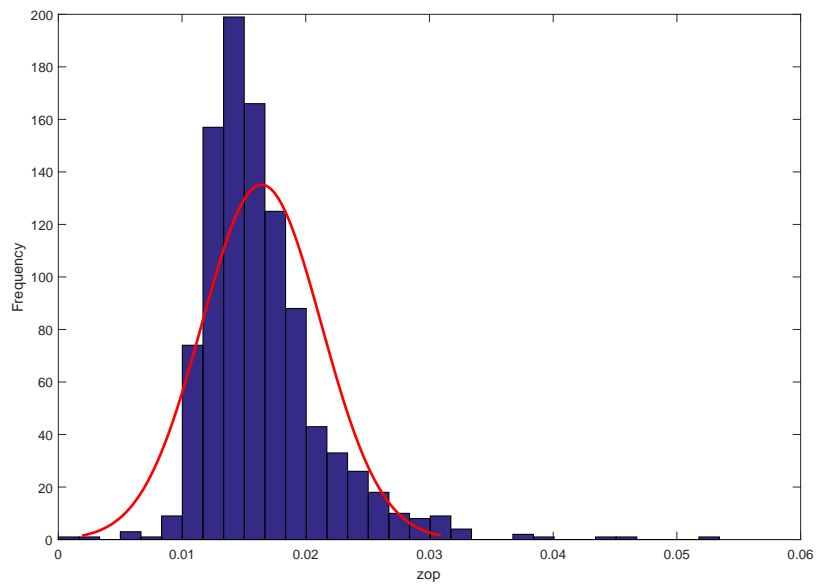
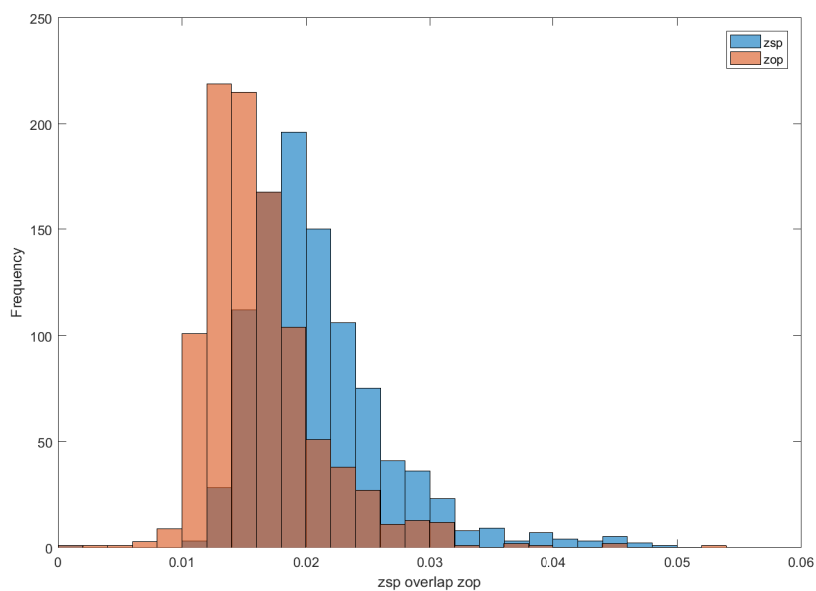
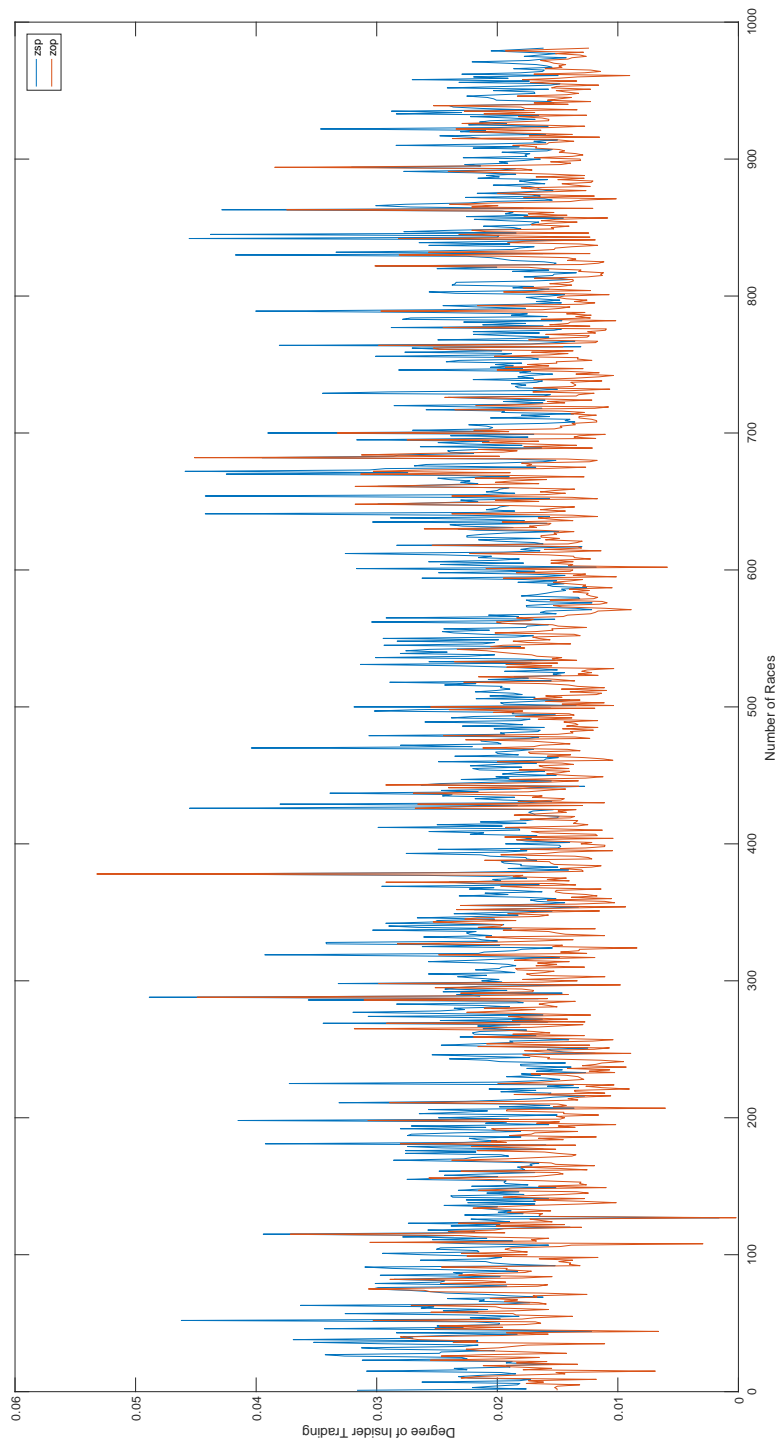


Figure 2.6 ZSP Overlap ZOP (2013)



Note: For comparison, we display two histograms in a single figure. It is obvious that the degree of insider trading at SP is, in general, larger than the degree at OP.

Figure 2.7 Comparison between ZSP and ZOP (2013)



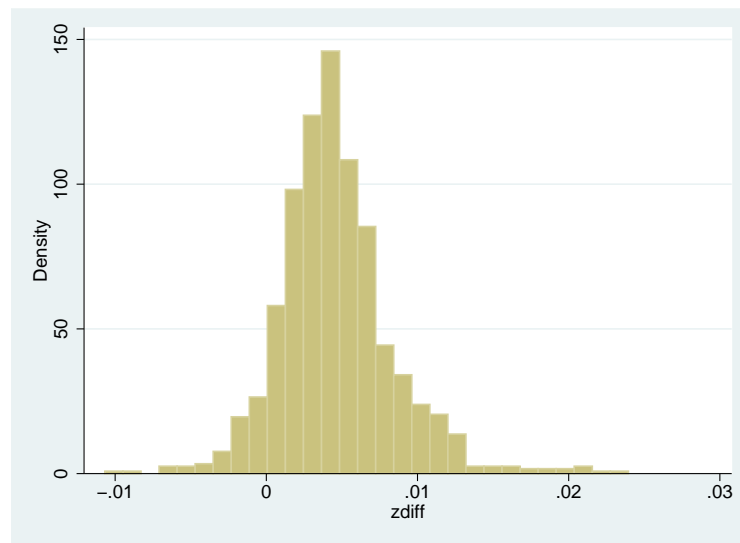
We conventionally test market efficiency theory in terms of information, that is, a market is efficient if the prices fully reflect all available information as described by Dowie (1976). The available information contains historical prices, public and inside or private information. In what follows, we address the issue of information efficiency in betting markets, which is related to opportunities to earn abnormal returns. In financial markets, weak form efficiency implies that there is no profitable strategy that can be used to yield above-average returns by predicting future prices from historical prices. Similarly, semi-strong form efficiency represents a situation where it is not possible to earn abnormal profits based on the information that is publicly available, while strong form efficiency denies the chance to get above-average or abnormal returns on the basis of all information including private information. Since the main purpose of this paper is to detect the incidence of insider trading in 2013 Yorkshire horse race betting markets, we test efficiency relating to the inside information, which is the existence of strong form efficiency. In the context of the horse race betting market, the concept of “efficiency” needs to restate. The added advantage of the betting markets is that there is a termination point for each bet. If the market is efficient, the starting prices should reflect past information on prices. More importantly, it should incorporate any inside or private information. So bettors with superior information cannot earn abnormal profits if they place a bet at the later stage, therefore the degree of insider trading at SP should accordingly decline. It is obvious that the incidence of insider trading at SP is higher than at OP from Figure 2.7. The reason why  $z_{sp}$  is higher than  $z_{op}$  is that the sum of starting prices is bigger than the sum of opening prices. As we see in Table 2.6, the adjusted final prices are larger than the opening prices in 3538 race runners. For the sample of 2920 runners, the starting prices equal the opening prices. The opening prices are higher in 3314 race horses. The following is a test for market efficiency.

If betting markets are efficient, the starting prices posted by on-course bookmakers should reflect the equilibrium prices that summarise all the available information. Punters with superior information, therefore, are not willing to place bets at the final set of prices. Next, we perform the paired t-test, also known as the paired-

Table 2.6 Summary Statistics for SP and OP

	<i>OP &gt; SP</i>		<i>OP = SP</i>		<i>OP &lt; SP</i>	
	<b>OP</b>	<b>SP</b>	<b>OP</b>	<b>SP</b>	<b>OP</b>	<b>SP</b>
<b>Obs</b>	3314	3314	2920	2920	3538	3538
<b>Mean</b>	0.1137	0.0969	0.9149	0.9149	0.1357	0.1609
<b>S.D.</b>	0.0921	0.0834	0.0916	0.0916	0.107	0.1192
<b>Min</b>	0.0066	0.0049	0.0049	0.0049	0.0039	0.0049
<b>Max</b>	0.6923	0.6522	0.7778	0.7778	0.9412	0.9524

Figure 2.8 The Density of the Difference between ZSP and ZOP (2013)



samples t-test or dependent t-test, which is used to determine whether there is a difference in the degree of insider trading in two related groups - OP and SP. To get a valid result, we have to check whether there are significant outliers in both groups and whether the distribution of the differences in the dependent variable is normally distributed. From Figure 2.8, we find that there is no significant outlier in our sample and the dependent variable is approximately normally distributed. The result of the paired t-test in Table 2.7 significantly reject the null hypothesis and accept the alternative that the mean of the degree of insider trading is not the same in two groups. And  $\mu_{zsp} - \mu_{zop} > 0$  is statistically significant. Thus we cannot conclude that this market is strongly efficient.

Table 2.7 Hypothesis Test Result

Paired t test						
Variable	Obs	Mean	Std. Err.	Std. Dev.	95% Conf. Interval	
<b>zsp13</b>	981	.0210303	.0001799	.0056337	.0206774	.0213833
<b>zop13</b>	981	.0164279	.0001544	.0048361	.0161249	.0167309
<b>diff</b>	981	.0046024	.0001212	.0037946	.0043646	.0048401
mean(diff) = mean(zsp13 - zop13)					t=37.9883	
$H_0$ : mean(diff)=0					degrees of freedom=980	
$H_a$ : mean(diff)<0		$H_a$ : mean(diff)≠0		$H_a$ : mean(diff)>0		
$Pr(T < t) = 1.0000$		$Pr( T  <  t ) = 0.0000$		$Pr(T > t) = 1.0000$		

## 2.6 Conclusions

Based on the Shin's (1993) theory of insider trading in a horse-race betting market, our study presents two methods to estimate the incidence of insider trading for nine Yorkshire on-course markets in 2013. The estimation under the original Shin measure is around 1.7%. It also demonstrates that there is a strong positive correlation between the sum of prices and the number of runners. Inside information on the part of punters might explain why total percent exceeds one in most races. The betting odds on the day diverged across bookies might be a possible explanation for the anomalies, but it is not always the case. The second method attempts to compute this degree of insider trading at opening and starting prices for each racetrack. From this perspective, we find that the weighted average value at SP is around 2.103% and at OP 1.64%. Lastly, there is no evidence to confirm that this market is strongly efficient.

## 2.7 Appendices

### Appendix A. Methodology: the Shin Model

Following the Shin's theoretical model, the first assumption is that bookmakers are perfectly competitive and risk neutral. The next assumption is that operating costs are negligible and any profits are bid away. Under this circumstance, if there are no taxes or levies and no inside information among punters, the sum of prices on all horses should have paid one pound for sure at the end of the race. The divergence of SP total percent from one shows how much the betting was over-round (i.e., the bookmakers' margin), which can be explained by the existence of inside information on the part of bettors. Thus the probability of the Insider chosen to place bets with the bookie is denoted by  $z$ . In what follows we assume there are  $n$  types of tickets ( $n$  runners) in a race. Denote  $\pi_i$  the price of the ticket and  $\pi_i$  corresponds to the raw odds of  $k$  to  $l$  ( $k/l$ ),  $0 \leq \pi_i \leq 1$  for all  $i$ . Now assume that any punter can only bet precisely £1 on a particular horse. If horse  $i$  is chosen to bet by a punter, this means  $1/\pi_i$  units of the  $i$ th ticket are sold in the market. Let  $\boldsymbol{\pi}$  be the vector of prices  $(\pi_1, \dots, \pi_n)$  and the implicit total percent  $\sum \pi_i \leq \beta$ , where  $\beta$  is an upper limit on the sum of prices for all horses in the race. Denote by  $p_i$  ( $i = 1, \dots, n$ ) the probability of the  $i$ th horse winning the race and by  $\mathbf{p}$  the vector of winning probabilities  $(p_1, \dots, p_n)$ , where  $0 \leq p_i \leq 1$  for all  $i$  and  $\sum p_i = 1$ . From the bookmaker's point of view, the Insider is chosen to be the punter with probability  $z$  and the Outsider to be the punter with probability  $(1 - z)p_i$ . Conditional on horse  $i$  winning, the bookie expects to pay  $z/\pi_i$  to the Insider and  $(1 - z)p_i/\pi_i$  to the Outsider.

The bookmaker's total unconditional expected liabilities are

$$\sum p_i \frac{(1 - z)p_i + z}{\pi_i} \quad (2.9)$$

Since the bookmaker's revenue is one pound by assumption, the expected profit is

$$1 - \sum \frac{(1 - z)p_i^2 + zp_i}{\pi_i} \quad (2.10)$$

which is maximised subject to  $\sum \pi_i \leq \beta$  and  $0 \leq \pi_i \leq 1$  for all  $i$ .

Notice it is assumed that  $\beta \geq 1$ , and negative total percent is ruled out in the model. The solution of  $\pi_i$  in terms of  $\beta$ ,  $z$  and  $\mathbf{p}$  is obtained by classical optimization.

$$\pi_i = \frac{\beta \sqrt{z p_i + (1-z) p_i^2}}{\sum \sqrt{z p_s + (1-z) p_s^2}} \quad (2.11)$$

Substituting (2.11) into (2.10), and since the expected profit must be zero in any equilibrium, we get the expression of  $\beta$ .

$$\beta = \sum \pi_s = \left( \sum \sqrt{z p_s + (1-z) p_s^2} \right)^2 \quad (2.12)$$

Substituting (2.12) into (2.11), we get the expression of  $\pi_i$  in terms of  $\mathbf{p}$  and the parameter  $z$ .

$$\pi_i = \sqrt{z p_i + (1-z) p_i^2} \left( \sum \sqrt{z p_s + (1-z) p_s^2} \right)^{-2} \quad (2.13)$$

The next step is to estimate the parameter  $z$  in which the Shin's empirical model employs the function  $F(p_i) = \sqrt{z p_i + (1-z) p_i^2}$  and its second order Taylor expansion at the point  $1/n$  which is  $F(1/n) + F'(1/n)(p_i - 1/n) + \frac{1}{2} F''(1/n)(p_i - 1/n)^2$ . Summing over  $i$ ,

$$\sum F(p_i) = n F(1/n) + F'(1/n) \sum (p_i - 1/n) + \frac{1}{2} F''(1/n) \sum (p_i - 1/n)^2 \quad (2.14)$$

Since  $\sum (p_i - 1/n) = 0$ , the second term disappears. With some twist of terminology, we define the 'variance' of  $\mathbf{p}$  by  $Var(\mathbf{p})$ , where  $Var(\mathbf{p}) = (1/n) \sum (p_i - 1/n)^2 = (\sum p_i^2 / n) - (1/n^2)$ . Rewriting Equation (2.14) as

$$\sum F(p_i) = \sqrt{1 + z(n-1)} + \frac{1}{2} n F''(1/n) Var(\mathbf{p}) \quad (2.15)$$



Now that the square of  $\sum F(p_i)$  is equal to the sum of prices  $\sum \pi_i$  according to the Equation (2.12), we firstly take square for both sides of Equation (2.15) and subtract one subsequently, then we get an expression for the over-round which is denoted by the deviation  $D$ . Thus,

$$D = z(n-1) + n\sqrt{1 + z(n-1)F''(1/n)Var(\mathbf{p})} + \frac{1}{4}n^2[F''(1/n)Var(\mathbf{p})]^2 \quad (2.16)$$

By applying the Weierstrass Approximation Theorem, we get

$$A = n\sqrt{1 + z(n-1)F''(1/n)} = \sum_{k=0}^K a_k n^k$$

$$B = \frac{1}{4}n^2[F''(1/n)]^2 = \sum_{k=0}^K b_k n^k$$

Substituting these into (2.16), we get a linear equation in terms of the variables  $(n-1)$ ,  $(n^k Var(\mathbf{p}))$  and  $(n^k [Var(\mathbf{p})]^2)$ , where  $k = 0, 1, \dots, K$ .

$$D = z(n-1) + \sum_{k=0}^K a_k n^k Var(\mathbf{p}) + \sum_{k=0}^K b_k n^k [Var(\mathbf{p})]^2 \quad (2.17)$$

The above equation is the Shin's empirical measurement of the incidence of insider trading, which will be copied in the present paper.<sup>9</sup> The next problem is to find out the values of the adjustment terms  $Var(\mathbf{p})$  and  $[Var(\mathbf{p})]^2$ .

In practice, we do not know the winning probability vector  $\mathbf{p}$  for sure, but the normalised prices  $\hat{\pi}_s = \pi_s / (\sum_i \pi_i)$  can be a good proxy for the true probability. Let  $\hat{\boldsymbol{\pi}}$  be the vector of normalised prices. Based on Equation (2.13), the sum of squares of the normalised prices  $\hat{\boldsymbol{\pi}}_s$  is as follows.

$$\sum \hat{\pi}_s^2 = \sum \left( \frac{\pi_s}{\sum_i \pi_i} \right)^2 = \sum \frac{\pi_s^2}{\beta^2} = \sum \frac{(z p_i + (1-z) p_i^2) \beta}{\beta^2} = \frac{z + (1-z) \sum p_i^2}{\beta}$$

---

<sup>9</sup>Shin, H. S. (1993). Measuring the Incidence of Insider Trading in a Market for State-Contingent Claims. *The Economic Journal*, Vol. 103, No. 420, 1141-1153.

and

$$\begin{aligned}
Var(\hat{\boldsymbol{\pi}}) &= \frac{\sum(\hat{\pi}_s - 1/n)^2}{n} \\
&= \frac{\sum[(\hat{\pi}_s)^2 + 1/n^2 - 2\hat{\pi}_s/n]}{n} \\
&= \frac{\sum(\hat{\pi}_s)^2}{n} + \frac{(1/n^2)n}{n} - \frac{(2/n)\sum\hat{\pi}_s}{n} \\
&= \frac{\sum(\hat{\pi}_s)^2}{n} - \frac{1}{n^2}
\end{aligned}$$

In a similar way,  $Var(\mathbf{p})$  becomes  $\frac{\sum p_i^2}{n} - \frac{1}{n^2}$ . Substituting out  $\sum \hat{\pi}_s^2$  and rearranging, we obtain an expression for  $Var(\mathbf{p})$  in terms of  $Var(\hat{\boldsymbol{\pi}})$ ,

$$Var(\mathbf{p}) = \frac{\beta}{1-z} Var(\hat{\boldsymbol{\pi}}) + \frac{\beta - 1 - z(n-1)}{(1-z)n^2} \quad (2.18)$$

Thus, the iterative procedure for the estimation of  $z$  is as follows. In *Step 1*, utilising the original vector  $Var(\hat{\boldsymbol{\pi}})$  as a proxy for  $Var(\mathbf{p})$  in Equation (2.17), by running an ordinary least squares regression we get the initial value of  $z$ . In *Step 2*, substituting the initial  $z$  into Equation (2.18), we can calculate  $Var(\mathbf{p})$ . Then, applying the revised values of  $Var(\mathbf{p})$  to Equation (2.17) and re-estimate, we derive a revised estimate of  $z$ . In *Step 3*, repeat *Step 2* until the revised values of  $z$  converge.

## Appendix B. Algorithm in MatLab

The dataset is collected from a website called *the sporting life* (available from the author on request). The significant advantage is that it is cheap and publicly available. The disadvantage is the data are short lived, which is the main reason that we only include Yorkshire racecourses in this study. The way to clean data is as follows.

In the first step, we copy and paste from the website to an excel sheet. The problem we encounter in this stage is that odds are in fractional form and the excel sheet is in text form, which we cannot use directly to calculate  $z$ . So the following example shows how we pick out columns from excel and transfer them into math mode to do the calculation.

Using York Racecourse as an example, the rest just follow the same procedure. The example of the form of the original data is shown by Figure 2.9, which is the race taking place on 15th May 2013. If we only calculate the incidence of insider trading at Starting Prices, all we need is the first column: number of horses in a race; and the last column: SP. If we are also interested in the incidence of insider trading at Opening Prices, the last lines in the third column under each horse are very important. The MatLab codes show this picking process and how to transfer odds format into  $(1/(SP+1))$ .

In the second stage, we count the number of runners in each race and sum each horse's prices. Then we put all data together to form a new file "data". The last step is to calculate variance and run the regression.

Figure 2.9 Example of the Original Data

York Racecourse 15/5/2013					
#	Horse	On-course prices	Tote win	Tote pl.	SP
1	First Mohican	3/1 10/3 7/2 4/1 9/2 5/1	3.60	1.50	3/1
2	Lahaag	9/2 5/1 9/2 5/1	5.20	1.90	9/2
3	Clayton	6/1 11/2 5/1 9/2 5/1 9/2 5/1 11/2	6.40	2.20	6/1
4	Prompter	16/1 14/1 12/1 14/1 12/1 11/1	21.00	6.10	16/1
5	Fluidity	50/1 40/1 33/1	50.60	11.50	50/1
6	Ruscello	7/1 15/2 8/1 17/2 9/1 17/2 8/1	7.70	2.70	7/1
7	Alfred Hutchinson	33/1	53.00	10.80	33/1
8	Itlaaq	20/1 18/1 16/1 20/1 18/1 16/1	22.70	5.30	20/1
10	Bridle Belle	8/1 17/2 9/1 10/1 11/1	10.20	2.90	8/1
11	Silvery Moon	20/1	21.50	4.80	20/1

### First Step.

```
% MatLab Code
clear all;
clc

%read excel table
[num1,horse,row1]=xlsread('York 2013','A:A');
[num2,op,row2]=xlsread('York 2013','C:C');
[num3,sp,row3]=xlsread('York 2013','F:F');

save('horse','horse');
save('op','op');
save('sp','sp');
load horse
load op
load sp

N = size(sp,1);
N1 = size(horse,1);
N2 = size(op,1);
SP = zeros(N1,1);
Horse = zeros(N1,1);
OP = zeros(N1,1);

%pick out SP
for i=1:N
a=sp(i);
b=cell2mat(a);
c=str2num(b);
if isempty(c)
SP(i+1)=0;
else
SP(i+1)=c;
end
end

%pick out OP
for i=1:N2
a2=op(i);
b2=cell2mat(a2);
```

```

c2=str2num(b2);
if isempty (c2)
OP(i+1)=0;
else
OP(i+1)=c2;
end
end

%pick out horse
for i=1:N1
a1=horse(i);
b1=cell2mat(a1);
c1=str2num(b1);
if isempty(c1)
Horse(i)=0;
else
if c1<1 || c1>100
Horse(i)=0;
else
Horse(i)=c1;
end
end
end

SP=real(SP);
OP=real(OP);
table1=[Horse,OP,SP];
save('table1',table1)

load table1
table2=table1(3:3481,:);
M=[0,0,0];
L=3479; %change number here
for i=1:L
if i>=L
break
end
if all(table2(i,:)==M,2) && all(table2(i+1,:)==M,2) && ...
all(table2(i+2,:)==M,2)

```

```

table2(i:i+2,:)=[];
L=L-3;
end
end
A=table1(:,1);
B=table1(:,2);
C=table1(:,3);
I=find(A);
for i=1:1389
if A(I(i)+1)==0
B1(i,:)=B(I(i+1)-1)
else
B1(i,:)=B(I(i))
end
end

B1(1390)=20;
A1=A(I);
C1=C(I);
table3=[A1,B1,C1];
save('table3',table3)

load table3
horse=table2(:,1);
op=table2(:,2);
op=1./(op+1);
sp=table2(:,3);
sp=1./(sp+1);
table4=[horse,op,sp];
save('table4',table4)

```

### **Second Step.**

```

%calculate op total price;
load table4
table5=[horse,op];
save('table5',table5)
result1=zeros(1000,40);
table5=[table5;[0,0]];

```

```

n=size(table5,1);
j=0;
h=1;
for i=1:n-1
a=table5(i,1)
b=table5(i+1,1)
if b>a
h=h+1;
else
j=j+1;
result1(j,1:h)=table5(i+1-h:i,2);
h=1;
end
end
save('result1',result1)
N=nnz(result1(:,1));
final1=zeros(N,40);
for i=1:N
final1(i,1)=nnz(result1(i,:));
final1(i,2)=sum(result1(i,:));
end
final1(:,3:40)=result1(1:N,1:38)
save('final1',final1)

%calculate sp total price;
table6=[horse,sp];
save('table6',table6)
load table6
result2=zeros(1000,40);
table6=[table6,[0,0]];
n1=size(table6,1);
j1=0;
h1=1;
for i=1:n-1
a1=table6(i,1)
b1=table6(i+1,1)
if b1>a1
h1=h1+1;

```



```

else
j1=j1+1;
result2(j1,1:h1)=table6(i+1-h1:i,2);
h1=1
end
end
save('result2',result2')
N1=nnz(result2(:,1));
final2=zeros(N1,40);
for i=1:N1
final2(i,1)=nnz(result2(i,:));
final2(i,2)=sum(result2(i,:));
end
final2(:,3:40)=result2(1:N1,1:38)
save('final2',final2')
A=final1;
A(:,23:40)=[];
B=final2;
B(:,23:40)=[];
save('A',A)
save('B',B)

```

### Third Step.

```

function [NV,N2V,N3V,N4V,N5V,N6V,V2,NV2,N2V2,N3V2,N4V2,N5V2,N6V2]
= calV(V,cleandata)
n=cleandata(:,1);
N=size(V,1);
for i=1:N
NV(i,1) = V(i,1) * n(i);
N2V(i,1) = n(i)^2 * V(i,1);
N3V(i,1) = (n(i)^3) * V(i,1);
N4V(i,1) = (n(i)^4) * V(i,1);
N5V(i,1) = (n(i)^5) * V(i,1);
N6V(i,1) = (n(i)^6) * V(i,1);
V2(i,1) = V(i,1)^2;
NV2(i,1) = V2(i,1) * n(i);
N2V2(i,1) = V2(i) * (n(i)^2);

```

```

N3V2(i,1) = V2(i) * (n(i)3);
N4V2(i,1) = V2(i) * (n(i)4);
N5V2(i,1) = V2(i) * (n(i)5);
N6V2(i,1) = V2(i) * (n(i)6);
end
end

% load the rearranged data
load data
cleandata=data;
cleandata(:,1)=[];
TF1 = cleandata(:,2)-1<0 ;
cleandata(TF1,:) = [] ;
save ('cleandata','cleandata');
beta=cleandata(:,2);
n=cleandata(:,1);
N=size(cleandata,1);
Runner=cleandata(:,1)-1;
D=cleandata(:,2)-1;
for i=1:N
V(i,1) = sum((cleandata(i,3 : (n(i) + 2))/cleandata(i,2)).2)/n(i) - 1/(n(i)2);
end

[NV,N2V,N3V,N4V,N5V,N6V,V2,NV2,N2V2,N3V2,N4V2,N5V2,N6V2 ]
= calV(V,cleandata);
save('var1','D','Runner','V','NV','N2V','N3V','N4V','N5V','N6V',
'V2','NV2','N2V2','N3V2','N4V2','N5V2','N6V2');
X=[Runner,V,NV,N2V,N3V,N4V,N5V,N6V,V2,NV2,N2V2,N3V2,N4V2,N5V2,N6V2];
b=regstats(D,X,eye(15));
Bb(1)=b;
BNV(:,1)=NV;
BN2V(:,1)=N2V;
BN3V(:,1)=N3V;
BN4V(:,1)=N4V;
BN5V(:,1)=N5V;
BN6V(:,1)=N6V;
BV(:,1)=V2;
BNV2(:,1)=NV2;

```

```

BN2V2(:,1)=N2V2;
BN3V2(:,1)=N3V2;
BN4V2(:,1)=N4V2;
BN5V2(:,1)=N5V2;
BN6V2(:,1)=N6V2;

for i=1:3
Z(i)=b.beta(1);
AV(:,i) = beta / (1 - Z(i)) . * V + (beta - 1 - Z(i) * (n - 1)) ./ (n.^2 * (1 - Z(i)));
[NV,N2V,N3V,N4V,N5V,N6V,V2,NV2,N2V2,N3V2,N4V2,N5V2,N6V2 ]
= calV(AV(:,i),cleandata);
save(strcat('var',num2str(i+1)), 'D','Runner','V','NV','N2V','N3V','N4V','N5V','N6V',
'V2','NV2','N2V2','N3V2','N4V2','N5V2','N6V2');
X=[Runner,AV(:,i),NV,N2V,N3V,N4V,N5V,N6V,V2,NV2,N2V2,N3V2,N4V2,N5V2,N6V2];
b=regstats(D,X,eye(15));
Bb(i+1)=b;
BNV(:,i+1)=NV;
BN2V(:,i+1)=N2V;
BN3V(:,i+1)=N3V;
BN4V(:,i+1)=N4V;
BN5V(:,i+1)=N5V;
BN6V(:,i+1)=N6V;
BV2(:,i+1)=V2;
BNV2(:,i+1)=NV2;
BN2V2(:,i+1)=N2V2;
BN3V2(:,i+1)=N3V2;
BN4V2(:,i+1)=N4V2;
BN5V2(:,i+1)=N5V2;
BN6V2(:,i+1)=N6V2;
end
Z(4)=b.beta(1);
save ('Bb','Bb')

```

## Appendix C. The Extension of the Shin Measure

We follow the same step to deal with the original data. The extension of the Shin measure follows:

```
function [ p ] =indp(z, beta, pai)
% here pai is scalar

$$p = (z^2 / (4 * (1 - z)^2) + pai.^2 / (beta * (1 - z))).(1/2) - z / (2 * (1 - z));$$

end

function [ rp ] =p( z, beta, pai)
% rp the is the sum of probablities minus 1
% beta is scarlar
% pai is a vector
p=indp(z, beta, pai);
rp=sum(p)-1;
end

load A
NN1=size(A, 1);
Z1=zeros(NN1, 1);
for i=1:NN1
beta=A(i, 2);
N=nnz(A(i, 3:25))+2;
pai=A(i, 3:N);
f=@(z) p( z, beta, pai);
Z1(i)=fzero(f, 0.2);
end
save Z1
xlswrite('zop13.xlsx', Z1)

load B
NN=size(B, 1);
Z=zeros(NN, 1);
for i=1:NN
beta=B(i, 2);
N=nnz(B(i,3:25))+2;
pai=B(i, 3:N);
f=@(z) p( z, beta, pai);
```

```
Z(i)=fzero(f, 0.2);  
end  
save Z  
xlswrite('zsp13.xlsx', Z)
```

## Chapter 3

# A Further Examination of Yorkshire Horseracing Betting Markets

*A gambler is doomed to lose even if he  
constantly herds with favourable odds.*

### 3.1 Introduction

The majority of economic papers on betting markets over the decades have been concentrated on the discovery and explanation of anomalies that approve of being supportive of market inefficiency. Another part of the literature has been focusing on market microstructure theory, which investigates the theoretical conditions under which the problem of insider trading deviates the equilibrium prices. As betting markets are short-lived, they take advantage of not considering its fundamentals or future dividend when pricing state-contingent claims (the price of bets) because a commonly acknowledged outcome is yielded at the end of the betting period. Betting markets also share other features with financial markets, for example, a publicly available cheap source of information (Racing Post, the Sporting Life etc.), with hundreds of thousands of investors. In this paper, we continue utilising

Yorkshire on-course horse-race betting markets to consolidate the results in the previous paper and test some other results that widely spread in this field.

On-course horseracing betting markets are usually bounded in time and have cash constraint. The betting period lasts for 20 to 30 minutes, and punters cannot withdraw their decisions once the purchase is completed. Further, transactions can only be made while the market is open, and thus any deviations of the final set of prices from predicted equilibrium levels can be considered as the outcome of some inefficiency (Schnytzer and Snir (2007)). Also if the market is efficient, starting prices should reflect equilibrium prices that summarise all the available information. Bettors with superior information are supposed not to take bets in the later stage, thus the incidence of insider trading at SP should not higher than at OP. Our findings demonstrate that the degree of the insider at SP is not only higher at the average level but also larger at each racetrack, which means the Yorkshire on-course betting markets are inefficient based on this data set. The paired t-test can confirm this result.

Another interesting finding is that our data set can provide sufficient evidence of the Crafts' (1985) result that insiders tend to bet early with bookmakers. His hypothesis also includes that "horses that exhibit large decreases in the odds against winning in the betting period, known as plungers, have been hypothesised to be indicative of insider activity" (Law and Peel (2002)). Our empirical results do not conform to the hypothesis that plunger could be an indicator of insider activity in terms of Law and Peel's (2002) definition of a big mover. Compared with the Shin measure of insider trading, movement factor plays a less important role in explaining the rates of return.

Last but not least, we prove that the degree of insider activity has significantly negative effects on the average rates of return at both opening and starting odds from the punter's perspective, which verifies a well-known stylised fact that bookmakers respond to insider trading by raising their overall margin as well as biasing the odds they set. A gambler is doomed to lose even if he constantly herds with favourite

odds. The average loss to one unit stake on the favourite horse is estimated as 0.3314.

The outline of this paper is as follows. In section 3.2, we describe the dataset and the empirical results. Section 3.3 summarises the empirical findings of the paper.

## 3.2 Empirical Analysis

### 3.2.1 Data Analysis

The dataset in this paper includes 1265 races and 12205 runners at Yorkshire race-tracks from 1st January 2014 to 28th December 2014, as reported by the *Sporting Life*. For each horse we collect its opening odds and the starting odds. Opening odds are the first set of odds posted by the on-course bookmakers. The last set of odds posted by bookmakers before the race commences is known as starting odds. We also collect the odds of favourite horses in each race, denoting as favourite odds in this paper.

In order to proceed with our analysis, the following variables need to be defined:

$OO_i^j$  - The opening odds for horse  $i$  in race  $j$ .

$SO_i^j$  - The starting odds for horse  $i$  in race  $j$ .

$FO_i^j$  - The favourite odds for horse  $i$  in race  $j$ .

The relationship between prices and odds is as follows.

$$price = \frac{1}{1 + odds}$$

Thus bets on the winner earn  $1 - price > 0$  and bets on all other horses lose the stake.

$OP_i^j = \frac{1}{OO_{i+1}^j}$  - The opening price of horse  $i$  winning race  $j$ .



$SP_i^j = \frac{1}{SO_{i+1}^j}$  - The starting price of horse  $i$  winning race  $j$ .

$sdummy_i^j$  generates 1 if horse  $i$  wins race  $j$  and  $-1$  otherwise.

$fdummy_i^j$  generates 1 if horse  $i$  is the favourite horse in race  $j$  and 0 otherwise.

Indicating by  $sdummy$  and  $fdummy$ , we can identify the winning prices and the favourite prices from each race.

$RateO$  and  $RateS$  are the rates of returns at opening and starting prices respectively.

To capture the magnitude of the change between opening and starting prices, the measure of movement is defined as follows.

$$mop = \ln \left( \frac{1 - op}{1 - sp} \right)$$

Table 3.1 reports the summary statistics when we introduce the measure of movement ( $mop$ ). It is positive if the starting prices are greater than the opening prices, 0 if  $sp = op$  and negative otherwise. In other words, odds are shortened during the betting period. We then define a big mover as one for which the measure of movement ( $mop$ ) exceeds 0.05, which is similar to that reported by Law and Peel (2002). With this definition we have 835 runners (6.84% of the sample of horses) regarded as big movers from opening to starting prices. Table 3.1 also summarises the average rates of return to one unit stake at OP and SP based on the  $mop$ . The prominent feature of this table is that there are no positive returns at starting prices for all ranges of movers. Only two positive returns occur at opening prices fell in the range  $mop > 0$ , which is consistent with the Crafts' (1985) empirical findings.

In order to confirm the above results, we also include the Crafts' ratio and define a big mover as the one in which the ratio of probabilities derived from starting to opening prices exceeds 1.5. Table 3.2 is shown to describe the relationship between average rates of returns and movers. Combining both ratios in Table 3.2, we find that the rates of return at opening prices are positive and recorded at starting prices are negative for all ranges of movers, indicating that there is supportive evidence

Table 3.1 Movers and Rates of Return (MOP)

<b>Range: mop</b>	<i>N</i>	<i>RateS</i>	<i>RateO</i>
> 0.05	835	-0.0535982	0.0084638
0 ~ 0.05	3363	-0.0096033	0.0074582
0	3823	-0.0184617	-0.0184617
-0.05 ~ 0	3879	-0.0112051	-0.0250121
< -0.05	305	-0.0116164	-0.0654312

Note that  $mop = \ln(1 - op) - \ln(1 - sp)$  is the measure of movement from opening to starting prices. This ratio is positive if starting odds are larger than opening odds, which means horses exhibit large decreases in the odds against winning, also known as plungers or steamers. As pointed out by Law and Peel (2004), plungers are the ones for which  $mop > 0.05$ . They also suggest that big movers from opening to starting odds could be an indicator of insider activity. Positive rates of return are shown at opening odds when the ratio exceeds 0.

Table 3.2 Movers and Rates of Return (Crafts' Ratio)

<b>Crafts Ratio</b>	<i>N</i>	<i>RateS</i>	<i>RateO</i>	<i>mop</i>	<b>Range:mop</b>
< 0.5	3	-0.0297563	-0.070915	-0.043787	<0
$\geq 0.5 < 1$	4181	-0.0112218	-0.0279277	-0.0203571	
1	3823	-0.0184617	-0.0184617	0	0
> 1 < 1.1	842	-0.0055099	0.0088754	0.0232533	0 ~ 0.5
$\geq 1.1 < 1.2$	1430	-0.0126534	0.0070402	0.0265483	
$\geq 1.2 < 1.3$	875	-0.0235477	0.0050862	0.0372001	
$\geq 1.3 < 1.4$	497	-0.0366034	0.0003032	0.0459985	
$\geq 1.4 < 1.5$	164	-0.0100774	0.0321869	0.052975	
$\geq 1.5 < 2$	354	-0.0340944	0.0107195	0.0522392	>0.05
$\geq 2$	36	-0.0499587	0.0259487	0.0868293	
<i>ALL</i>	12205	-0.0159473	-0.0127332	0.0047535	

Note that the Crafts' ratio equals  $SP/OP$ . The last column of the above table shows the *mop*. Combining both ratios, we find that the rates of returns recorded at opening odds are positive and recorded at starting prices are negative for all ranges of movers. The empirical results are in accordance with the findings reported by Crafts (1985).

of insiders tending to place bets at the early stage, which is consistent with the Crafts' (1985) results.

There are three differences from the original Crafts' ratio. The first difference is that forecast odds are excluded as the data that we use are not collected day by day. The data set is all the historical prices including opening odds and starting odds. Since opening odds reflect the forecast odds, informed information and public

trends in the press (McCririck (1992)), and they are also the first odds posted by the on-course bookmakers, it is more precise than forecast odds when we test whether the markets are efficient. The second difference is that we exclude the extent of movement from FP to SP and its reverse ratio in the present study. The third difference is that we do not exclude the odds that are larger than 10/1 (0.090909). Instead we bring in the measurement of movement (*mop*).

### **3.2.2 Test**

#### **ZSP and ZOP**

Following the extension of the Shin measure from the previous chapter, we can get the degree of insider trading at SP and OP respectively using the 2014 Yorkshire racing data. For our sample of 1265 races, the estimated degrees of insider trading at both starting and opening prices are given by the following tables 3.3 and 3.4, respectively. The weighted average degree at the level of SP is estimated as 2.07%, with a minimum value of 0.0014691 and a maximum of 0.095421, of which 11 negative values are deleted due to the fact that the sum of prices does not exceed one. At opening prices the weighted average value is 1.65%, with a minimum value of 0.0006524 and a maximum of 0.0959448. The negative degree of insider trading is yielded in 21 out of the 1265 races. The degree of insider trading decreases in 108 races from opening to starting prices and increases in the remaining 1134 races. Figures 3.1, 3.2 and 3.3 describe the same features for *zsp* and *zop*.

#### **Market Efficiency Test**

As we mentioned in the previous chapter, if betting markets are efficient, the starting prices posted by on-course bookmakers should reflect equilibrium prices that summarise all the available information. Punters with superior information, therefore, are not willing to place bets at the final set of prices. We also run the paired t-test to determine whether the mean of the difference of *zsp* and *zop* is the same. Figure 3.4 demonstrates that the dependent variable is approximately

Table 3.3 z for Each Racecourse at Starting Prices (2014)

	<b>Mean</b>	<b>S.D</b>	<b>Min</b>	<b>Max</b>	<b>Obs</b>
<b>Beverley</b>	0.0239868	0.0084675	0.0132433	0.0952211	131
<b>Catterick</b>	0.0202186	0.0092139	0.0025023	0.0917412	179
<b>Doncaster</b>	0.0207472	0.0081152	0.0014691	0.0759641	256
<b>Pontefract</b>	0.0202059	0.0049563	0.0132023	0.0440676	102
<b>Redcar</b>	0.0177758	0.004952	0.0032617	0.0403167	134
<b>Ripon</b>	0.0217448	0.0053365	0.0143872	0.049095	105
<b>Thirsk</b>	0.0213681	0.0105228	0.0080064	0.095421	107
<b>Wetherby</b>	0.0219475	0.0064094	0.0116605	0.0469832	125
<b>York</b>	0.018299	0.0032622	0.0082536	0.0325894	115

Note: this table summarises the average degree of insider trading for each racecourse. The degree of insider trading for each race is estimated based on the method proposed by Jullien and Salanié (1994). At starting prices, we rule out 11 races that have negative values of z, which leaves us 1254 races for investigation. For comparison, we delete all the races that have a negative degree of insider trading.

Table 3.4 z for Each Racecourse at Opening Prices (2014)

	<b>Mean</b>	<b>S.D</b>	<b>Min</b>	<b>Max</b>	<b>Obs</b>
<b>Beverley</b>	0.0190567	0.0064088	0.0035896	0.0668446	131
<b>Catterick</b>	0.0160973	0.0094212	0.0022536	0.0959448	177
<b>Doncaster</b>	0.0163892	0.0075561	0.0006524	0.0719377	254
<b>Pontefract</b>	0.0163647	0.0045722	0.0102278	0.033492	101
<b>Redcar</b>	0.0137984	0.0040969	0.0008263	0.0314539	131
<b>Ripon</b>	0.0165539	0.0044739	0.0065818	0.0366767	104
<b>Thirsk</b>	0.017703	0.0106591	0.0018628	0.0911478	106
<b>Wetherby</b>	0.0178254	0.0067011	0.0093708	0.042789	125
<b>York</b>	0.0154542	0.0028539	0.006057	0.0261539	115

Note: this table summarises the average degree of insider trading for each racecourse. The degree of insider trading for each race is estimated based on the method proposed by Jullien and Salanié (1994). At opening prices, we rule out 21 races that have negative values of z, which leaves us 1244 races for investigation. For comparison, we delete all the races that have a negative degree of insider trading.

Figure 3.1 Distribution of ZSP (2014)

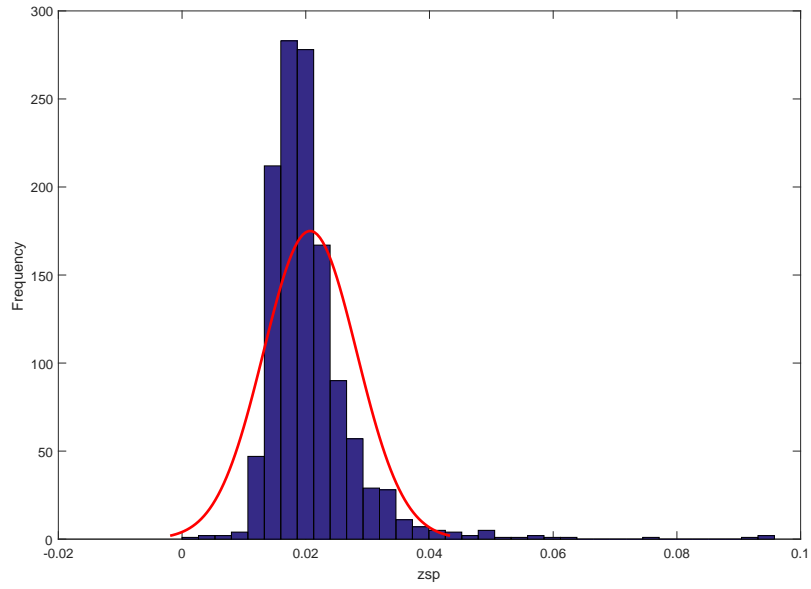


Figure 3.2 Distribution of ZOP (2014)

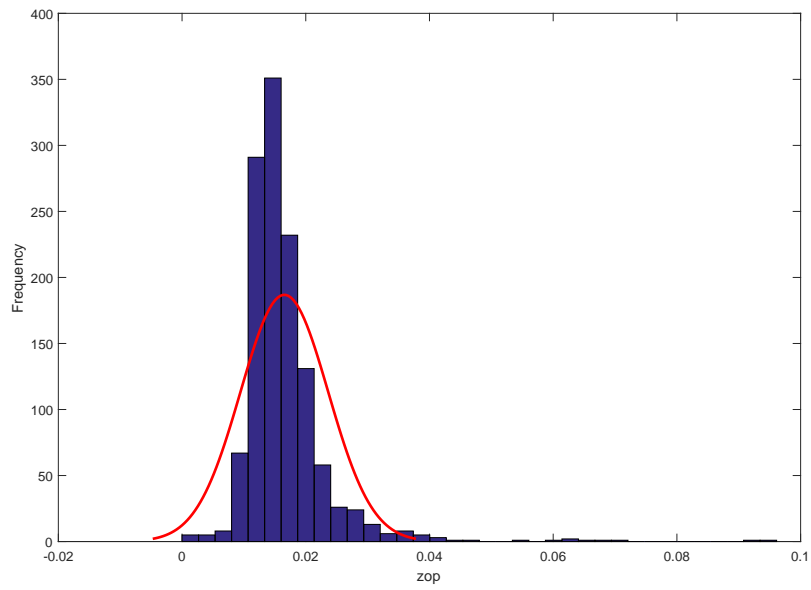
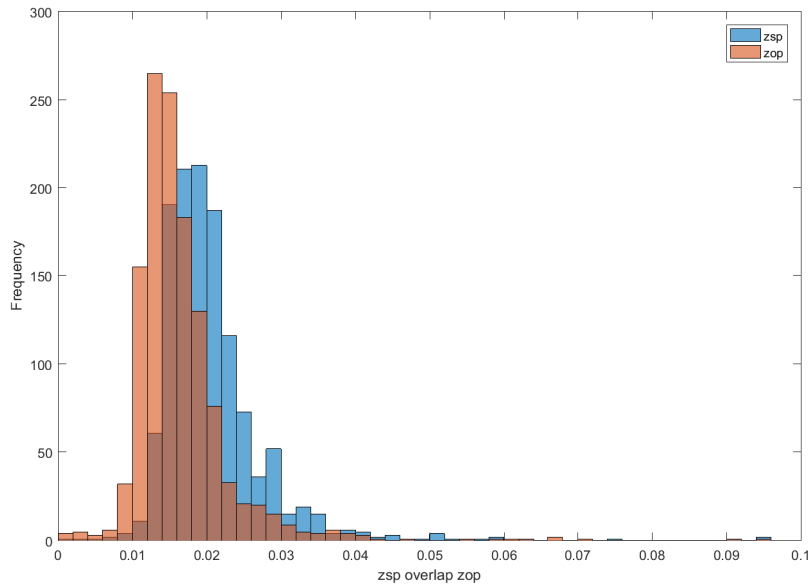


Figure 3.3 ZSP Overlap ZOP (2014)



Note: For comparison, we display two histograms in a single figure. It is obvious that the degree of insider trading at SP is, in general, larger than the degree at OP.

normally distributed and there is no significant outlier in our sample set. Thus the result that we get from Table 3.5 is valid. We can conclude that there is a difference in two related groups and the market is not strongly efficient.

Figure 3.4 The Density of the Difference between ZSP and ZOP (2014)

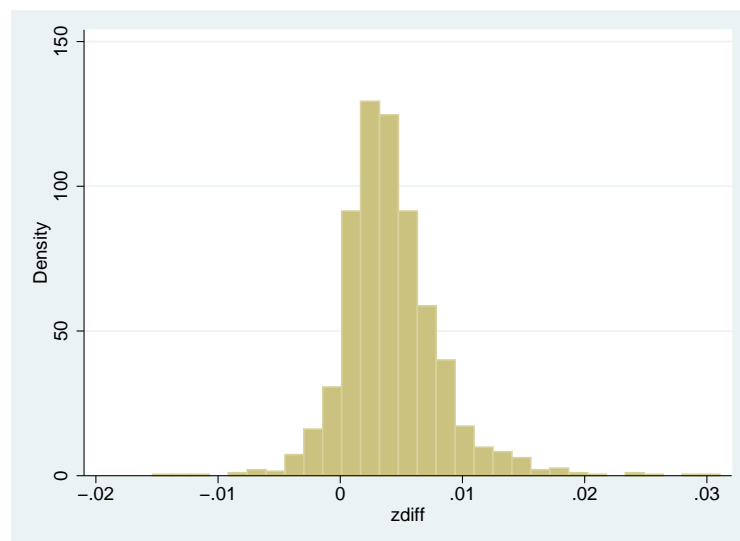


Table 3.5 Hypothesis Test Result

Paired t test						
Variable	Obs	Mean	Std. Err.	Std. Dev.	95% Conf. Interval	
<b>zsp14</b>	1244	.0207016	.0002138	.0075416	.0202821	.0211211
<b>zop14</b>	1244	.0165373	.0002011	.0070915	.0161429	.0169318
<b>diff</b>	1244	.0041642	.000115	.0040565	.0039386	.0043899
mean(diff) = mean(zsp14 - zop14)					t=36.2070	
$H_0$ : mean(diff)=0					degrees of freedom=1243	
$H_a$ : mean(diff)<0		$H_a$ : mean(diff)≠0		$H_a$ : mean(diff)>0		
$Pr(T < t) = 1.0000$		$Pr(  T   <   t  ) = 0.0000$		$Pr(T > t) = 1.0000$		

### Relationship between Winning Prices and Favourite Prices

Winning prices are the set of prices that represent a contingent claim to 1 if horse  $i$  wins race  $j$  at the level of starting prices. Since the winning prices indicate the actual results in each race, it can be considered as the unadjusted probabilities of winning. Favourite prices are the set of prices that represent a contingent claim to 1 if horse  $i$  is chosen to be the favourite horse in race  $j$ . Due to the longshot-favourite bias, bookmakers tend to underprice the favourites in most cases, which means the favourite prices are always the highest prices in each race (High price, high winning probability, shorter odds, lower return, vice versa). Our sample contains 1265 races, of which 95 observations are deleted due to the fact that there is more than one or no favourite horse in these races. The average value of the winning prices is estimated as 0.2385, with a minimum value of 0.0099 and a maximum of 0.90909. The average value of the favourite prices is estimated as 0.356, with a minimum value of 0.11765 and a maximum of 0.90909. As we see from Table 3.6, among the whole sample, 390 out of 1170 races have the same prices. This means that the favourite horses are actually the winners in these 390 races. For the majority, the winning prices are less than the favourite prices. Figure 3.5 also shows that the favourite prices are larger than the winning prices. The winning prices are greater than the favourite prices in only 4 races. Normally, the favourite horse is the one with the shortest odds in the race. In reality, there do exist anomalies, but as can be seen from Table 3.6, the difference between the favourite prices and the winning

prices is very small. It approximately equals to 0. These results confirm that in order to offset the negative effects of insider trading, bookmakers tend artificially to bias the odds to distract punters and guarantee their profits or balance their book.

In order to obtain significance levels, we run ordinary least squares (OLS) regression equations in which the winning prices are the dependent variable and independent variables included the favourite prices and the Shin measure of insider trading at starting odds.

Equation (1) in Table 3.7 is the winning starting prices regressed on a constant over the whole sample. The estimated value is 0.2385, as reported above. Equation (2) tests the relationship between the winning starting prices and the favourite starting prices. The coefficient on the favourite prices is significantly positive but less than 1, indicating the longshot-favourite bias and that market's favourite horses might not be the best choice for outsiders. Taking the degree of insider trading at starting odds into account, the adjusted *R – squared* is improving and *zsp* plays an important role in explaining the winning probability.

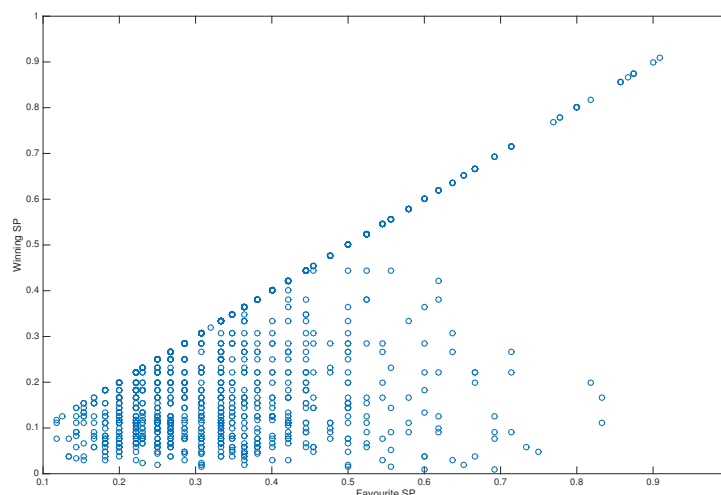


Table 3.6 Statistic Summary: Winning Starting Prices versus Favourite Starting Prices

	$WinSP - FavSP > 0$		$WinSP - FavSP = 0$		$WinSP - FavSP < 0$	
	Win SP	Favourite SP	Win-Fav	Win SP	Favourite SP	Win-Fav
<b>Obs</b>	4	4	4	390	390	776
<b>Mean</b>	0.45455	0.45454	9.98e-06	0.4053	0.4053	0.1535
<b>S.D.</b>	0	0	0	0.1617	0.1617	0.09196
<b>Min</b>	0.45455	0.45454	9.98e-06	0.11765	0.11765	0.0099
<b>Max</b>	0.45455	0.45454	9.98e-06	0.909	0.909	0.54545

Note that 95 observations are deleted due to the fact that there is more than one or no favourite horse in these races, which leaves us 1170 observations for investigation. For the sample of 390 races, the favourite horse is actually the winning horse. For the rest, betting on the favourite horse is not the best choice for punters and the average loss is  $(0.45454 \times 4 + 0.3308 \times 776) / (4 + 776) = 0.3314$ .

Figure 3.5 Winning Starting Prices versus Favourite Starting Prices



Note: this scatterplot shows that the winning prices are always less than the favourite prices. Recall that a high price implies a high winning probability, thereby giving the shorter odds in the fraction form. So conventionally the favourite horse in a race is always the one that has the shortest odds. The winning horse is not always the favourite horse.

Table 3.7 Regression Results - Winning Starting Prices versus Favourite Starting Prices

<i>Eqs.</i>	<i>Dependent Variable</i>	<i>Constant</i>	<i>Favourite Prices</i>	<i>zsp</i>	$\bar{R}^2$	<i>N</i>
(1)	<i>Winprice</i>	0.2385*** (48.28)				1170
(2)	<i>Winprice</i>	-0.0316** (-3.05)	0.7584*** (28.09)		0.4026	1170
(3)	<i>Winprice</i>	-0.0633*** (-5.62)	0.6333*** (19.33)	3.6919*** (6.52)	0.4231	1170

*t* scores are in parentheses.

\*\*\* indicates significance at the 0.1% level.

\*\* indicates significance at the 1% level.

### Test of the Shin Measure of Insider Trading on Average Returns

As can be seen from Table 3.8, Equation (1) is the average rate of return at starting prices regressed on a constant for the whole sample. The mean value is -0.0167 which indicates the average loss to one unit stake on every runner in a race. Equation (2) adds the winning prices to test the relationship between the average rate of return and the implicit winning probability on each horse for the whole sample. Notice that the winning prices can be regarded approximately as the winning probability based on starting odds. The coefficient on the winning prices is significantly negative, indicating horses with a higher probability of winning lower the average rate of return. The negative coefficient on the favourite prices in Equation (3) has the similar effects as that in Equation (2). Putting the Shin measure of insider trading at starting odds in Equation (4) significantly changes the coefficients on the winning and favourite prices, which is not surprising at all because insider activity is expected to depress the average rate of return on a massive scale. Bookmakers ensure their profits by artificially setting biased odds so that the winning probability and the favourite prices no longer matter in this case. These two coefficients are relatively small compared to the significantly negative one on the  $zsp$ , which is consistent with the Shin model. Taking into consideration of *dummy*, which takes the value of 1 if the Shin measure of insider trading rises between opening and closing odds and is 0 otherwise, and  $zop$  (the Shin measure of insider trading at opening odds) plays no role in turning the tide, but the adjusted  $R - squared$  has been indeed improved.

Table 3.9 reports 8 regression equations where the average rate of return at opening prices is the dependent variable and the explanatory variables are the same as that in Table 3.8. Comparing with these two tables, we find out that the results have not been changed too much except  $zop$  has more influence on the average rate of return at opening odds.

Thus we conclude that the extended Shin measure of insider trading at both opening and closing odds do significantly affect the average rate of return at both levels.

## Movement Factor

As mentioned above, the movement factor ( $mop$ ) is defined as the natural logarithm of returns at OP over returns at SP, which gives more weight to the higher implied win probability following the change in odds. A big mover is defined as when  $mop \geq 0.5$ . As shown in Table 3.10, Equation (1) regresses the rates of return at starting prices on a constant, the measure of movement,  $fdummy$  (which takes the value of 1 if a horse is chosen to be the favourite by the market and 0 otherwise) and  $wdummy$  (1 if a horse wins a race and 0 otherwise). Equation (2) has the same explanatory variables but a different dependent variable which is the rates of return at opening odds. The results do not demonstrate a big difference except that the coefficient on  $mop$  in the first equation is significantly negative while in the second equation it is significantly positive. Dummy variables on the winning horse are of vital importance in explaining the returns. Downsizing data set by filtering if  $mop$  greater than 0.5 does not change the results too much. The coefficient on  $mop$  becomes a negative value and thus has the negative effects on the rates of return at opening prices. Lowering the returns at the early stage does not help to attract punters with inside information take bets. Therefore based on our methodology and data set, we cannot conclude that the plunger is regarded as an indicator of insider activity.

Table 3.8 Regression Results - Average Rate of Return at Starting Prices

<i>Eqs</i>	<i>Dependent Variable</i>	<i>Constant</i>	<i>Winning Prices</i>	<i>Favourite Prices</i>	<i>zsp</i>	<i>zop</i>	<i>dummy</i>	$\bar{R}^2$	<i>N</i>
(1)	AR	-0.0167*** (-112.87)							1170
(2)	AR	-0.014*** (-59.12)	-0.0109*** (-13.31)					0.131	1170
(3)	AR	-0.0115*** (-31.74)	-0.0048*** (-4.71)	-0.0114*** (-9.37)				0.1911	1170
(4)	AR	-0.0055*** (-47.77)	0.0016*** (5.27)	0.0064*** (16.91)	-0.6687*** (-114.52)			0.9339	1170
(5)	AR	-0.0056*** (-48.79)	0.0015*** (5.23)	0.0062*** (16.53)	-0.6661*** (-115.53)		0.0008*** (6.09)	0.9359	1170
(6)	AR	-0.0055*** (-48.31)	0.0015*** (5.29)	0.0062*** (16.57)	-0.6886*** (-71.07)	0.0273** (4.89)	0.0006*** (4)	0.9363	1170

1. *t* scores are in parentheses.

2. \*\*\* indicates significance at the 0.1% level.

\*\* indicates significance at the 0.5% level.

Note: AR is the average rate of return at starting prices for each race.

Winning prices are the set of prices which represent a contingent claim to 1 if horse *i* won race *j*.

Favourite prices are the set of prices which represent a contingent claim to 1 if horse *i* is chosen to be the favourite horse by the market in race *j*.

*zsp* is the Shin measure of insider trading at starting odds.

*zop* is the Shin measure of insider trading at opening odds.

*dummy* takes the value 1 when the change in the Shin measure of insider trading from opening to starting odds is positive, that is,  $zsp < zop$ .

*N* is the number of observations.

Table 3.9 Regression Results - Average Rate of Return at Opening Prices

<i>Eqs</i>	<i>Dependent Variable</i>	<i>Constant</i>	<i>Winning Prices</i>	<i>Favourite Prices</i>	<i>zsp</i>	<i>zop</i>	<i>dummy</i>	$\bar{R}^2$	<i>N</i>
(1)	ARO	-0.0133*** (-86.41)							1170
(2)	ARO	-0.0108*** (-43.42)	-0.0109*** (-12.29)					0.1137	1170
(3)	ARO	-0.0085*** (-21.93)	-0.0051*** (-4.5)	-0.0104*** (-7.92)				0.1582	1170
(4)	ARO	-0.0042*** (-40.98)	0.0016*** (5.69)	0.0046*** (13.39)		-0.6873*** (-133.19)		0.948	1170
(5)	ARO	-0.0041*** (-40.92)	0.0016*** (5.77)	0.0047*** (13.52)		-0.6913*** (-130.09)	0.00038** (3.03)	0.9484	1170
(6)	ARO	-0.0037*** (-12.21)	0.00025 (0.31)	0.0035** (3.47)	-0.5288*** (-34.42)			0.5822	1170
(7)	ARO	-0.0035*** (-12.31)	0.0007 (0.94)	0.0046*** (4.95)	-0.5448*** (-38.49)		-0.0047*** (-14.72)	0.6473	1170
(8)	ARO	-0.0044*** (-41.28)	0.0015*** (5.27)	0.0041*** (11.64)	0.0585*** (6.55)	-0.7374*** (-84.17)	0.00076*** (5.56)	0.9502	1170

1. *t* scores are in parentheses.

2. \*\*\* indicates significance at the 0.1% level.

\*\* indicates significance at the 0.5% level.

Note: ARO is the average rate of return at opening prices for each race.

Winning prices are the set of prices which represent a contingent claim to 1 if horse *i* won race *j*. Here we report winning at opening prices.

Favourite prices are the set of prices which represent a contingent claim to 1 if horse *i* is chosen to be the favourite horse by the market in race *j*.

*zsp* is the Shin measure of insider trading at starting odds.

*zop* is the Shin measure of insider trading at opening odds.

*dummy* takes the value 1 when the change in the Shin measure of insider trading from opening to starting odds is positive, that is,  $zsp < zop$ .

*N* is the number of observations.

Table 3.10 Regression Results Including MOP

<i>Eqs</i>	<i>Dependent Variable</i>	<i>Constant</i>	<i>mop</i>	<i>fdummy</i>	<i>wdummy</i>	$\bar{R}^2$	<i>N</i>
(1)	<i>RateS</i>	-0.08699*** (-115.81)	-0.4787*** (-21.96)	-0.2098*** (-88.17)	0.9339*** (400.89)	0.9295	12205
(2)	<i>RateO</i>	-0.08688*** (-114.18)	0.2053*** (9.3)	-0.2106*** (-87.37)	0.9334*** (395.52)	0.9278	9306
<i>Movement: mop</i> $\geq$ 0.05							
(3)	<i>RateS</i>	-0.109*** (-13.62)	-1.3122*** (-16.3)	-0.1107*** (-15.11)	0.9386*** (117.66)	0.9433	835
(4)	<i>RateO</i>	-0.07978*** (-8.94)	-0.8809*** (-9.82)	-0.1182*** (-14.48)	0.9318*** (104.84)	0.9296	835

1. *t* scores are in parentheses.

2. \*\*\* indicates significance at the 0.1% level.

3. *N* is the number of observations.

4. *fdummy* takes 1 if a horse is chosen to be the favourite by the market, and 0 otherwise.

5. *wdummy* takes 1 if a horse wins a race, and 0 otherwise.

### 3.3 Conclusions

Shin (1993) has constructed a well-known theoretical model, which allows estimation of the incidence of insider activity in the betting, based on the system in the United Kingdom, from the profile of starting prices. Since this paper is an extension of the previous one, we consolidate our results by applying a new data set collected from Yorkshire racecourses in 2014. Given that this is the up-to-date information, the time period in this set starts on the 1st of January 2014 and ends on the 28th of December 2014. Besides starting odds, our new data set consists of opening odds, winning prices and favourite prices. We conclude that the weighted average degree of insider trading at starting odds is around 2.07% while at opening odds 1.65%, which is very similar to the results that reported in the second chapter. There is also no obvious evidence of market efficiency, but there is supportive evidence of the Crafts' hypothesis that insiders tend to bet early with on-course bookmakers.

Significance tests above suggest that of particular interest from the punter's perspective the average rates of return at both starting and opening prices are spectacularly affected by the incidence of insider trading, which is consistent with the results reported by Shin (1991, 1992, 1993) in a different way. Betting on the favourite horse might not always be the best choice for outsiders because there is an average loss which equals 0.3314 per stake.

Finally, due to the constraints on the on-course betting markets that transactions are made in cash and decisions are irreversible, some improvements to the model and results are very likely if we explore exotic betting markets, that is, online markets in the future study.



# Chapter 4

## On the Determinants of Starting Prices in Horseracing Betting Markets

### 4.1 Introduction

As already acknowledged, panel data analysis has become increasingly popular in the field of social science. A panel contains observations from different objects over multiple periods of time. The use of panel data analytic applications by researchers on economics, political science, sociology and psychology is extremely common. In financial markets, panel data analysis mainly focuses on corporate finance and asset pricing. But in sports betting (or gambling) markets, in particular in horseracing betting markets, there is little literature employing cross-sectional and time-series data sets to tackle the related issues. Although there is an extensive literature of empirical studies on betting markets testing theories in financial markets, for instance, market efficiency theory (see Vaughan Williams (1999) for a survey. Parimutuel: Ali (1977), Snyder (1978), Figlewski (1979), Hausch et al. (1981), Asch et al. (1984, 1986); Fixed odds: Dowie (1976), Crafts (1985), Gabriel and Marsden (1990)), insider trading problem (Shin (1991, 1992, 1993), Jullien and Salanié

(1994), Fingleton and Waldron (1999), Cain et al. (2001)) and so on, it has provided little guidance, to the best of our knowledge, to researchers as to determinants that bookmakers will be taken into account when pricing horses. In this paper, we examine the factors that affect the starting prices of horses through a unique publicly available panel data set. The uniqueness of this data set will be described in Subsection 3.1, but the reason why we are interested in this market is as follows.

Sports betting is a large and prominent industry around the globe. In the United Kingdom, the gross gambling yield (excluding the National Lottery) is 7.1 billion pounds during the period October 2013 - September 2014<sup>1</sup>. Horse racing in the UK and Ireland is exceptionally popular, in fact, it is the second biggest spectator sport after football in Great Britain in terms of attendances. According to the 2013 study on economic impact of British racing conducted by Deloitte, there are around 5.6 million attendances at 1,369 fixtures<sup>2</sup>. Our attentions paid to the horseracing betting markets not only because of the size of this industry, but the fact that this market shares a lot of similarities with financial markets, such as a large number of participants, the availability of public information and the uncertain future value of the asset and so on and so forth. Unlike asset pricing, which is the mainstream in the literature of financial markets, horse-pricing mechanism is barely reached in economics. The main purpose of this paper is to provide empirical evidence on the determinants of starting prices in horse-racing betting markets. The second purpose of this paper is to show that these factors are quite persistent.

As mentioned before, financial markets and betting markets happen to share some similar characteristics as to the structure of the market. In contrast to the betting system in North America in which the betting odds are determined by the pari-mutuel method, the on-course horse-racing betting market in Britain and Ireland is well suited for our investigation (Shin (1993)) due to the fact that bookmakers set odds. This allows us to use methods that have been widely applied

---

<sup>1</sup>Gambling Commission - Industry Statistics April 2010 to September 2014.

<sup>2</sup>Economic Impact of British Racing 2013. Accessible via: <http://www.britishhorseracing.com/wp-content/uploads/2014/03/EconomicImpactStudy2013.pdf>.

in finance to analyse the problem of pricing a state-contingent claim, which in this paper is the price of the horse.

In what follows other features of the on-course horse-racing betting markets and some terminologies should be brought in. On-course betting markets are usually bounded in time and have cash constraint, the betting period lasts for 20 to 30 minutes, and punters cannot withdraw their decisions once the purchase is completed (Schnytzer and Snir (2008)). Furthermore, transactions can only be made while the market is open and in the end there is a definite and commonly acknowledged outcome, which rules out the possibility of taking its future dividends or fundamentals into account. Opening prices are the set of prices which are first available to punters to bet on the course, while starting prices are defined as “*the odds at which a ‘sizeable’ bet could have been made just before the race commences*” (Dowie (1976)). Odds of  $a/b$  correspond to the price  $\pounds \frac{b}{a+b}$  in this paper, which means the punter pays  $\pounds \frac{b}{a+b}$  for a ticket in return for  $\pounds 1$  if that particular horse wins.

The contribution of this paper is that we introduce panel data analysis to find factors that have impacts on the starting prices, which has never been done before in the betting literature. Due to lack of resources for some horses, our panel data set is collected over unequally spaced time intervals. The panel is unbalanced as each horse races at different time periods. The results we find in this paper are quite fascinating. We provide empirical evidence to support the following factors that are very significant in explaining the price of the horse, including the winning potential (defined as  $1 - \frac{\text{rank}}{\#\text{runners}}$ ), the age of the horse, the weight the horse carries in each race and the distance of the race. Our findings confirm that the condition of the turf, the size of racetracks and the classification system affect the price as well. Based on our estimation results, we can conclude that all these factors are persistent.

The rest of the paper is structured as follows. Section 4.2 gives the outline of our empirical model. Data description, variable classification and estimation results are presented in Section 4.3 with relevant tables and figures. Section 4.4 concludes.

## 4.2 Empirical Model

### 4.2.1 Basic Fixed-Effects Model

In this chapter, we are interested in explaining starting prices for each horse. By adopting the fixed-effects regression model in Frees (2004), we consider each horse to be a *subject* and differentiate among horses with the index  $i$ , where  $i$  ranges from 1 to  $N$ .  $N$  is the total number of horses. Each horse is observed  $T_i$  times and  $t$  indexes the historical order of the races that horse  $i$  participated to. With these indices, let  $y_{it}$  denote the response of the  $i$ th horse at the  $t$ th race. For each response  $y_{it}$ , there is a set of explanatory variables. If we assume there are  $K$  explanatory variables  $x_{it,1}, x_{it,2}, \dots, x_{it,K}$  that may vary by horse  $i$  and race  $t$ , the  $K$  explanatory variables can be expressed as a  $K \times 1$  column vector

$$\mathbf{x}_{it} = \begin{pmatrix} x_{it,1} \\ x_{it,2} \\ \vdots \\ x_{it,K} \end{pmatrix}$$

Alternatively, we use the expression  $\mathbf{x}_{it} = (x_{it,1}, x_{it,2}, \dots, x_{it,K})'$ , where the prime means transpose. Thus, the data for the  $i$ th horse contains

$$\begin{aligned} &\{x_{i1,1}, \dots, x_{i1,K}, y_{i1}\} \\ &\quad \vdots \\ &\{x_{iT_i,1}, \dots, x_{iT_i,K}, y_{iT_i}\} \end{aligned}$$

We allow the panel to be unbalanced and the number of responses varies by horse, indicated by  $T_i$ . If  $T_i = T$  for each horse, the panel is balanced.

Consider a fixed-effects model, the relationships between the dependent and the independent variables follow the *regression equation*

$$y_{it} = \alpha_i + \beta_1 x_{it,1} + \beta_2 x_{it,2} + \dots + \beta_K x_{it,K} + \epsilon_{it} \quad (4.1)$$

where the intercept terms  $\{\alpha_i\}$  are allowed to vary by horse, also known as the unobserved time-invariant individual effect, and  $\beta_1, \beta_2, \dots, \beta_K$  are the slope parameters associated with the  $K$  explanatory variables. The error terms  $\epsilon_{it}$  are the regression disturbance. More compactly, the parameters can be expressed as a  $K \times 1$  column vector

$$\boldsymbol{\beta} = \begin{pmatrix} \beta_1 \\ \beta_2 \\ \vdots \\ \beta_K \end{pmatrix}$$

With this notation, we rewrite (4.1) as

$$y_{it} = \alpha_i + \mathbf{x}'_{it}\boldsymbol{\beta} + \epsilon_{it} \quad (4.2)$$

In this study, to investigate the relationship between starting prices and pricing factors, we fit the following regression model:

$$\begin{aligned} SP_{it} = & \alpha_i + \beta_{LAGWIN1}LAGWIN_{i,t-1} + \beta_{LAGWIN2}LAGWIN_{i,t-2} \\ & + \beta_{AGE}AGE_{it} + \beta_{AGE2}AGE2_{it} \\ & + \beta_{WEIGHT}WEIGHT_{it} + \beta_{DIST}DIST_{it} + \epsilon_{it} \end{aligned} \quad (4.3)$$

where  $i = 1, \dots, 100$  denotes the horses and  $t$  is the index of historical order of the races that horse  $i$  participated to.

Equation (4.3) is our *basic fixed-effects model*. We also extend the explanatory variables to include GOING dummies, RACETRACK dummies and CLASS dummies. The full list of variables can be found in Table 4.3.

The fixed-effects model allows us to analyse the impact of variables that vary over time and it explores the relationship between the response and the explanatory variables within a horse. Each horse has its own individual characteristics that may affect starting prices, such as the sex of the horse that remains unchangeable.

Thus, employing the fixed-effects would remove the effect of the time-invariant characteristics. We can assess the net effect of the explanatory variables on the outcome variable. There are two important assumptions. First of all, we assume that there is no correlation between the scalar disturbance term and the explanatory variables. Secondly, we allow that  $\alpha_i$  and  $\mathbf{x}_{it}$  can be correlated. Since each horse is different, the horse's error term and the constant should not be correlated with the others. To be more precise, we will present the results of Hausman's test in Subsection 4.3.3 to show that the fixed-effects regression model is appropriate.

## 4.2.2 Driscoll and Kraay Standard Errors for Pooled OLS Estimation

Returning to our regression model in Equation (4.2), we now re-organise the model as follows.

$$y_{it} = \alpha + \mathbf{x}'_{it}\boldsymbol{\beta} + (\alpha_i - \alpha + \epsilon_{it}) = \alpha + \mathbf{x}'_{it}\boldsymbol{\beta} + u_{it} \quad (4.4)$$

Equation (4.4) therefore can be expressed as

$$\mathbf{y} = \mathbf{W}\boldsymbol{\theta} + \mathbf{u} \quad (4.5)$$

where  $\mathbf{y} = (y_{1t_{11}} \cdots y_{1T_1} \quad y_{2t_{21}} \cdots y_{NT_N})'$ , and  $\mathbf{u} = (u_{1t_{11}} \cdots u_{1T_1} \quad u_{2t_{21}} \cdots u_{NT_N})'$ .  $\mathbf{W} = (\mathbf{w}_{1t_{11}} \cdots \mathbf{w}_{1T_1} \quad \mathbf{w}_{2t_{21}} \cdots \mathbf{w}_{NT_N})'$ , where  $\mathbf{w}_{it}$  is a  $(K + 1) \times 1$  vector of explanatory variables whose first element is 1 and the rest is  $K$  independent variables for horse  $i$  at race  $t$ .  $\boldsymbol{\theta} = (\alpha \quad \boldsymbol{\beta}')$  is a  $(K + 1) \times 1$  vector of unknown coefficients. In this case, we are assuming that  $E(\mathbf{u} | \mathbf{W}) = 0$ . Thus the error terms  $\mathbf{u}$  are strictly exogenous.  $\boldsymbol{\theta}$  can consistently be estimated by OLS regression, which yields

$$\hat{\boldsymbol{\theta}} = (\mathbf{W}'\mathbf{W})^{-1}\mathbf{W}'\mathbf{y} \quad (4.6)$$

Since we allow  $\alpha_i$  and  $\mathbf{x}_{it}$  to be correlated,  $u_{it}$  and  $\mathbf{x}_{it}$  can also be correlated. We therefore cannot use OLS regression because the estimation is inconsistent. The

disturbances  $u_{it}$  themselves are allowed to be autocorrelated, heteroskedastic, and cross-sectionally dependent.

Therefore, this study reports Driscoll and Kraay standard errors (Driscoll and Kraay (1998)) to mitigate the presence of cross-sectional dependence which has been confirmed by Pesaran's (2004) CD test.

Driscoll and Kraay standard errors for the coefficient estimates are obtained as follows. It is the square roots of the diagonal elements of the asymptotic (robust) covariance matrix<sup>3</sup>

$$Cov(\hat{\theta}) = (\mathbf{W}'\mathbf{W})^{-1} \hat{\mathbf{S}}_T (\mathbf{W}'\mathbf{W})^{-1} \quad (4.7)$$

where  $\hat{\mathbf{S}}_T$  is defined as in Newey and West (1987):

$$\hat{\mathbf{S}}_T = \hat{\mathbf{\Omega}}_0 + \sum_{l=1}^{L(T)} (\hat{\mathbf{\Omega}}_l + \hat{\mathbf{\Omega}}_l') \omega(l) \quad (4.8)$$

and

$$\omega(l, L(T)) = 1 - \frac{l}{1 + L(T)} \quad (4.9)$$

In the above expressions, let  $L(T)$  denote the lag length up to which the residuals may be autocorrelated and  $\omega(l, L(T))$  be the modified Bartlett weights that ensure positive semi-definiteness of  $\hat{\mathbf{S}}_T$ . It also "smooths the sample autocovariance function such that higher order lags get less weight" (Hoechle (2007)).

The  $(K + 1) \times (K + 1)$  matrix  $\hat{\mathbf{\Omega}}_l$  is defined as

$$\hat{\mathbf{\Omega}}_l = \sum_{t=l+1}^T \mathbf{h}_t(\hat{\theta}) \mathbf{h}_{t-l}(\hat{\theta})' \quad (4.10)$$

with

$$\mathbf{h}_t(\hat{\theta}) = \sum_{i=1}^{N(t)} \mathbf{h}_{it}(\hat{\theta})$$

---

<sup>3</sup>See Hoechle (2007).

Note that the individual orthogonality conditions  $\mathbf{h}_{it}(\hat{\boldsymbol{\theta}})$  in (4.10) runs from 1 to  $N(t)$ , which admits the following representation for pooled OLS estimation:

$$\mathbf{h}_{it}(\hat{\boldsymbol{\theta}}) = \mathbf{w}_{it}\hat{\epsilon}_{it} = \mathbf{w}_{it}(y_{it} - \mathbf{w}'_{it}\hat{\boldsymbol{\theta}})$$

$N$  varies with  $t$ , which allows the Driscoll and Kraay's original estimator to be ready for use with the unbalanced panel.

### 4.2.3 Fixed-Effects Model with Driscoll and Kraay Standard Errors

For the fixed-effects (within) regression with Driscoll and Kraay standard errors, the estimator is achieved in two steps. In step 1, we eliminate  $\alpha_i$  by transforming all model variables as follows.

$$y_{it} - \bar{y}_i = (\mathbf{x}_{it} - \bar{\mathbf{x}}_i)' \boldsymbol{\beta} + (\epsilon_{it} - \bar{\epsilon}_i) \quad (4.11)$$

where  $\bar{y}_i = T_i^{-1} \sum_{t=1}^{T_i} y_{it}$ ,  $\bar{\mathbf{x}}_i = T_i^{-1} \sum_{t=1}^{T_i} \mathbf{x}_{it}$  and  $\bar{\epsilon}_i = T_i^{-1} \sum_{t=1}^{T_i} \epsilon_{it}$ . Since  $\alpha_i$  is constant,  $\bar{\alpha}_i = \alpha_i$ , and therefore the effect is eliminated.

In step 2, we estimate the *within* transformed regression model in (4.11) by pooled OLS estimation with Driscoll-Kraay standard errors as in the above description.

## 4.3 Data Analysis and Empirical Results

### 4.3.1 Data Description

To identify the determinants that can be utilised to explain starting prices in horse-racing betting markets, we pick 100 horses with its lifetime results to form a longitudinal panel from a website *sportinglife.com*, which is publicly available to every individual. The data are comprised of combined races running at 101



Table 4.1 Summary of Racecourses

<b>Racecourse</b>	<b>Course Type</b>	<b>Races</b>
<i>United Kingdom</i>		
Aintree	National Hunt	99
Ascot	Flat/National Hunt	101
Cheltenham	National Hunt	226
Doncaster	Flat/National Hunt	56
Haydock Park	Flat/National Hunt	78
Kempton Park	Flat/National Hunt/All Weather	110
Lingfield	Flat/National Hunt/All Weather	68
Newbury	Flat/National Hunt	62
Sandown Park	Flat/National Hunt	68
York	Flat	64
<i>Ireland</i>		
Curragh	Flat	60
Fairyhouse	Flat/National Hunt	57
Leopardstown	Flat/National Hunt	121
Navan	Flat/National Hunt	51
Punchestown	National Hunt	129

racecourses in both the United Kingdom and Ireland during the period from 2006 to 2015, of which we pick up 10 racecourses from the UK including Aintree, Ascot, Cheltenham, Doncaster, Haydock Park, Kempton Park, Lingfield Park, Newbury, Sandown Park and York, and 5 racetracks (Curragh, Fairyhouse, Leopardstown, Navan and Punchestown) from Ireland to detect whether larger and competitive racecourses could affect the price of the horse. The reason why these 15 racecourses are described as larger and competitive is because they all contain more than 50 races and in total it accounts for 52.39% out of our 2577 samples. As a matter of fact, those 15 racecourses are indeed relatively large in terms of capacity, history and big events and all course types are included. Table 4.1 gives the specific details.

The uniqueness of this paper lies in the way we arrange the data in the form of a panel. Firstly, we label horses numerically so that 100 horses would be 100 objects in a cross-sectional dataset. For each horse we have all the entries and the day the horse entered the race has the value of 1 and then labels sequentially up to the latest race. To be precise, it is the separate cross-sectional dataset, which

is similar to the panel data. Therefore the method in the study is the Panel Data analysis, which has been widely used in the empirical research in recent years.

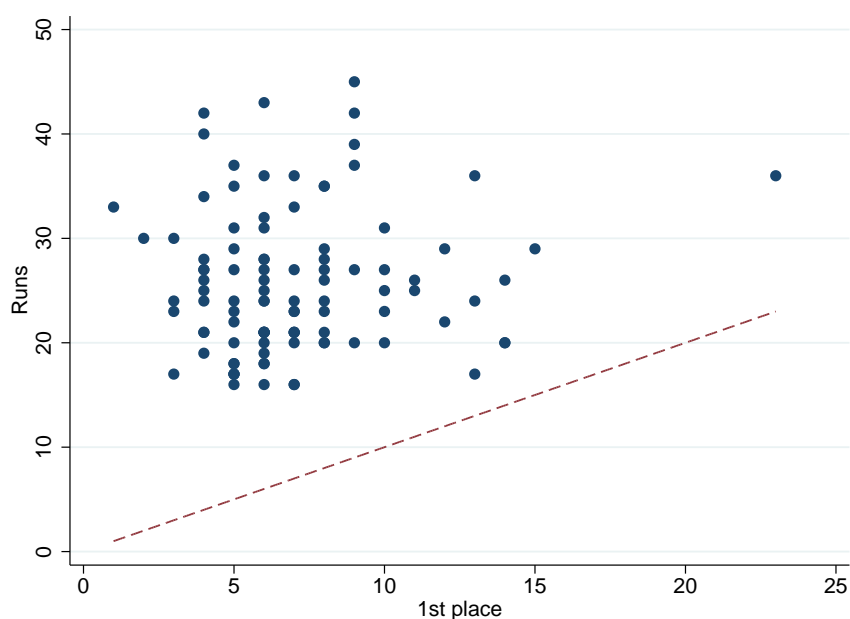
For each horse the following variables are collected: Starting Prices (SP), age, position and a total number of horses in a race, the weight the horse carries, the names of racecourses, the distance, the condition of the surface (i.e., going), and the classification. In our sample, a larger number of horses start its career at the age of four, and some of them start early (only at the age of two). The youngest horse in the sample is only four-year-old and the oldest one is thirteen years old in 2015. Table 4.2 and Figures 4.1, 4.2, 4.3 and 4.4 summarise the data, listing its lifetime results and some features of the horses. As can be seen from Table 4.2, the data include 100 observations with an average of 25.77 races ran by each horse. More specific, although this dataset is still unbalanced, we choose horses with at least 16 runs but no more than 45 races just in order to, to some extent, mitigate the effects caused by the unbalanced data. The place the horse was obtained can also be found in Table 4.2. Horses win the first place 7.05 times on average and all of them win at least once, in other words, there are no maiden horses in the sample. Horses get the second place with an average of 4.14 times and the third 3.31 times. The win percentage, as the name implies, represents the percentage of winning for each horse and equals to **First/Runs** in this study. It can also be used to measure the quality of the horses and averages 28.86% of winning. Red dash lines in Figures 4.1, 4.2 and 4.3 are 45-degree reference lines. As we see, all the dots are above the line which makes perfect sense as no horse can be guaranteed to win or get the second place or the third place every time they run a race. The outlier in Figure 4.1 represents the highest winning percentage in our sample set.

The principle of the level stake is that the total amount of betting funds is divided by a number of bets that you want to stake and keeps staking on each bet for evermore. It is the simplest among all staking plans and is somewhat beneficiary of its simplicity. Once the amount has been worked out to be staked on one bet, for the next thousand bets punter places the same amount, as long as we keep a

Table 4.2 Summary Statistics on 100 Horses

	Runs	First	Second	Third	Win %	£1 Level Stake
<b>Obs</b>	100	100	100	100	100	100
<b>Mean</b>	25.77	7.05	4.14	3.31	28.86	12.83
<b>S.D.</b>	6.8547	3.2579	2.2474	2.1115	0.1383	17.6644
<b>Min</b>	16	1	0	0	3.03	-17.61
<b>Max</b>	45	23	11	10	76.47	74.88

Figure 4.1 Scatter Plot of Horses Winning First Place



sensible number of bets under control, the chance of going bust is small<sup>4</sup>. As can be seen from Figure 4.4, red dash line is a benchmark, which represents the average amount of money staked on each horse in our data set. On average £12.83 staked on each horse with a minimal value of minus £17.61 and a maximum £74.88. The higher the level stake, the horse is more favourable in the market, so the winning percentage is high accordingly as Figure 4.5 describes. Notice that the winning percentage and the level stake are positively correlated, except that the level stake can go negative.

<sup>4</sup> Accessible via: <http://www.ukhorseracing.co.uk/faq/stakingplans.asp>.

Figure 4.2 Scatter Plot of Horses Winning Second Place

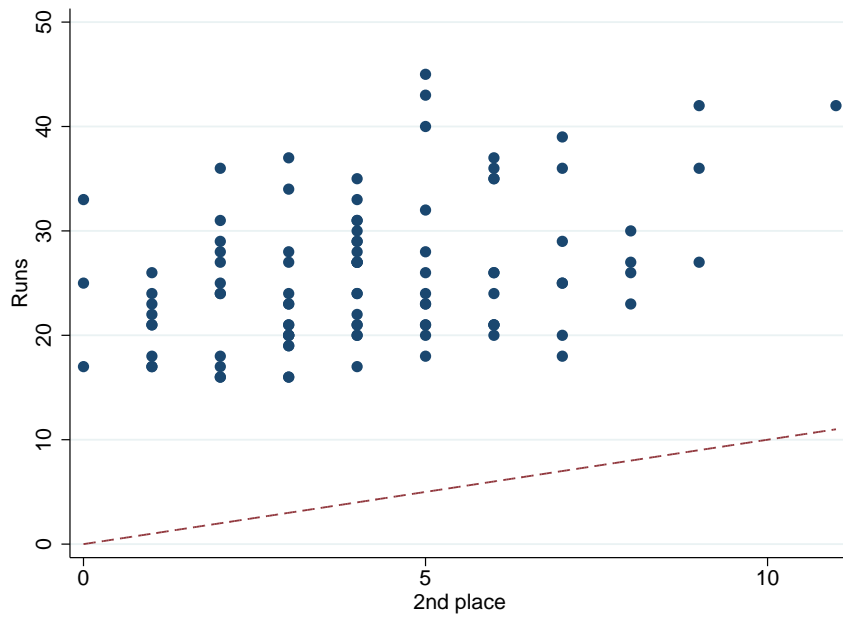


Figure 4.3 Scatter Plot of Horses Winning Third Place

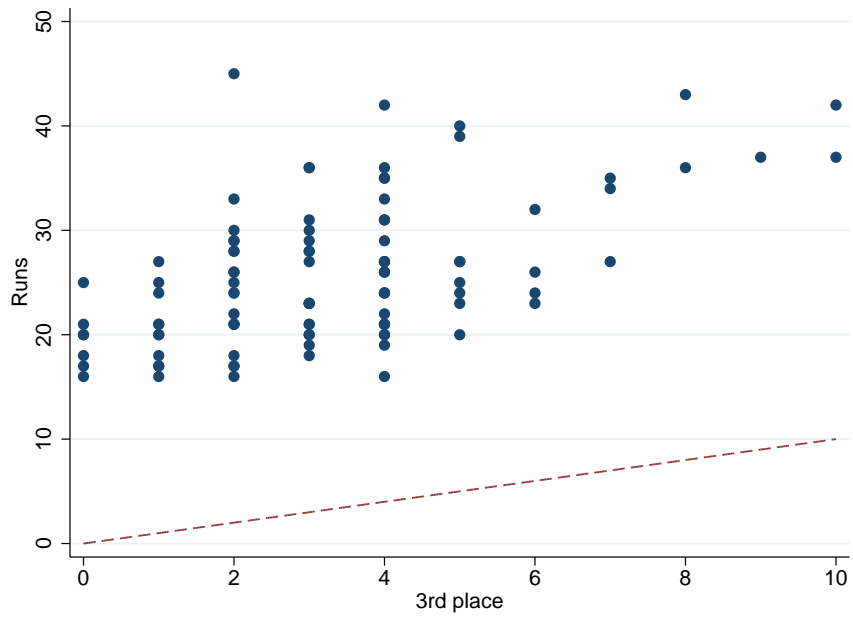


Figure 4.4 £1 Level Stake

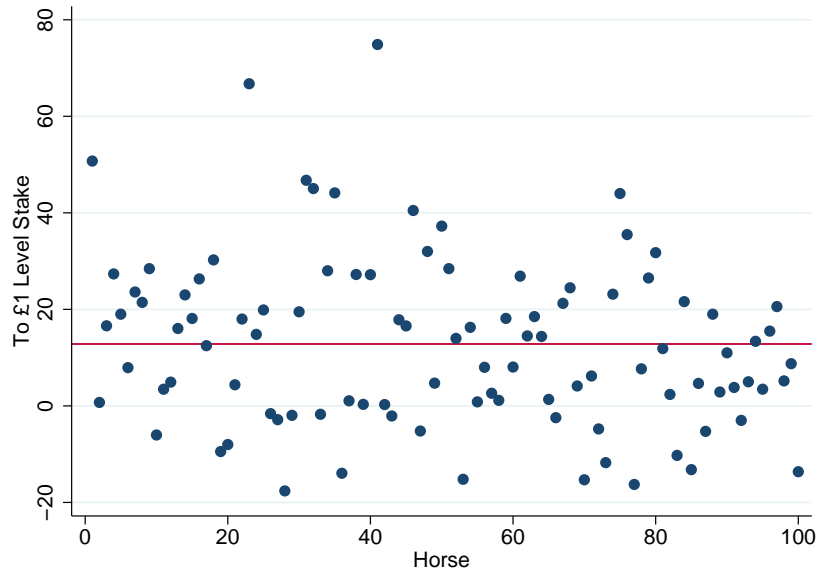


Figure 4.5 Winning Percentage

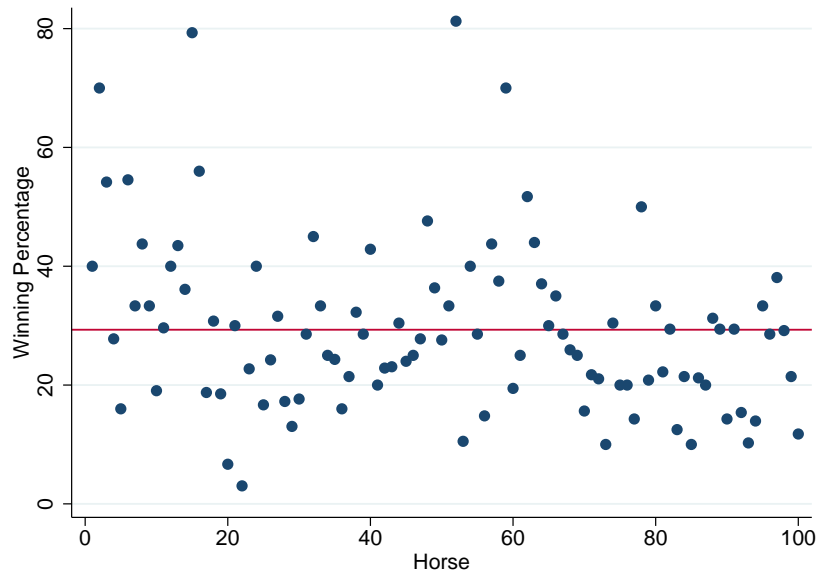
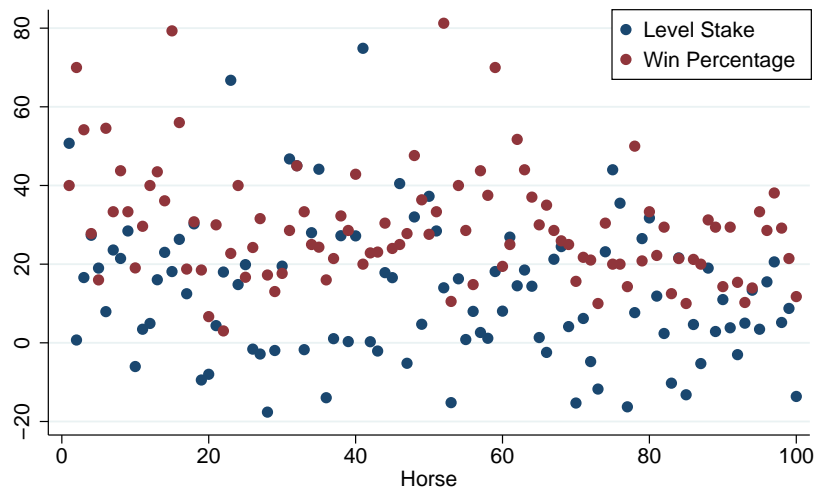


Figure 4.6 Winning Percentage vs. Level Stake

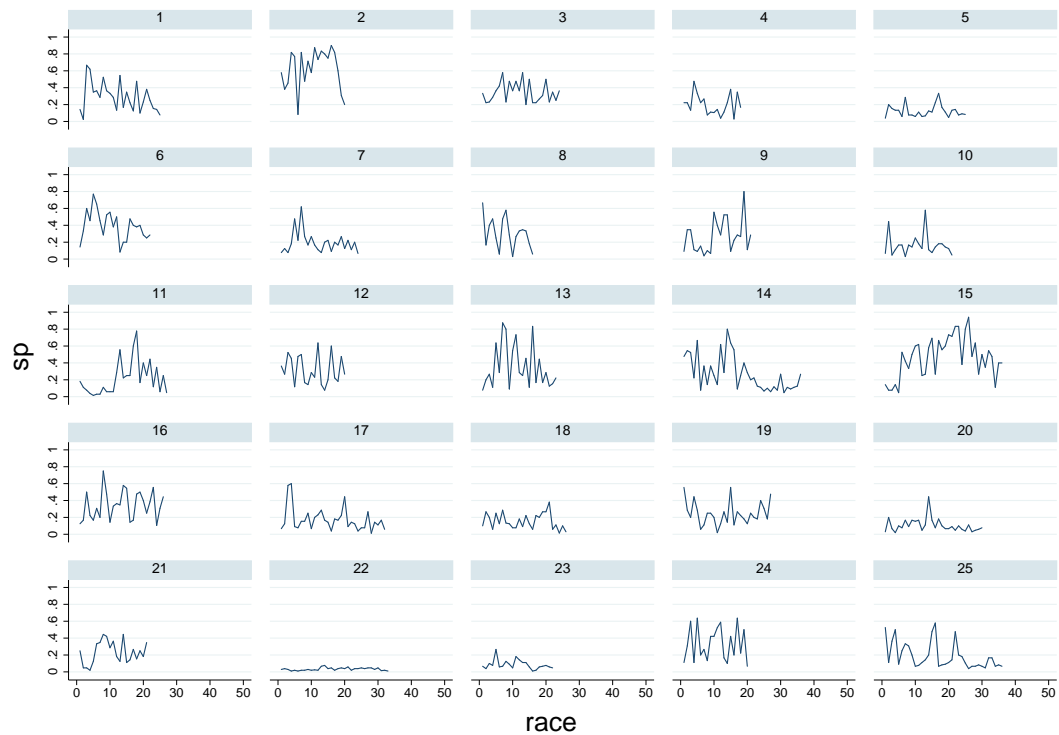


### 4.3.2 Variable Classification

#### Dependent Variable

As we attempt to identify the relationship between the starting price of the horse and factors that influence this price, starting odds are the best choice of the outcome variable in the model with a slight change of its format. As we have discussed in Section 4.1, for example, odds of 2/1 correspond to the price £0.3333 ( $1/(1+2)$ ), which means you pay £0.3333 for a ticket in return for £1 from bookmaker if that horse wins. Otherwise, you just lose £0.3333. This price can also be roughly considered as the subjective probability of winning for each horse in a race. The higher the price, the bigger the subjective winning probability is. So the horse is favoured in the market. Figures 4.7, 4.8, 4.9 and 4.10 below depict the changing starting prices during the entire races for 100 horses. Although most subfigures in these four figures fluctuate significantly, there are some representative subfigures, for example, horse 22 in Figure 4.7, horse 75 in Figure 4.9 and horse 76 in Figure 4.10, which show that the prices keep at a low level and do not change too much. Those are the ones who have a lower probability to win, also known as the underdog. The other feature is that almost all the horses start at a lower price because it is the

Figure 4.7 Starting Prices for Horses 1 to 25



Graphs by horse

first time they enter a race. A low price is anticipated due to the fact that there is no available information beforehand for both bookmakers and punters. In other words, there is no previous record, i.e. no rank, which serves as a very important factor in our study. For dependent variable, we should also notice that there are 9 missing data out of 2577 samples.

### Explanatory Variables

Variables that affect the performance of the horse can be divided into three parts: the first part is related to the horse *per se*, which includes its own age at the different period of time, the weight the horse carries and the winning potential. Winning potential, which represents the percentage of winning, can be used to reflect the long-term quality of a horse. The second part specifies the details of racetracks, which contains the names of racecourses, the distance the horse runs, post position, weather, going conditions and the classification system. The final part concerns jockey's characteristics. In what follows we discuss each of these components.

Figure 4.8 Starting Prices for Horses 26 to 50

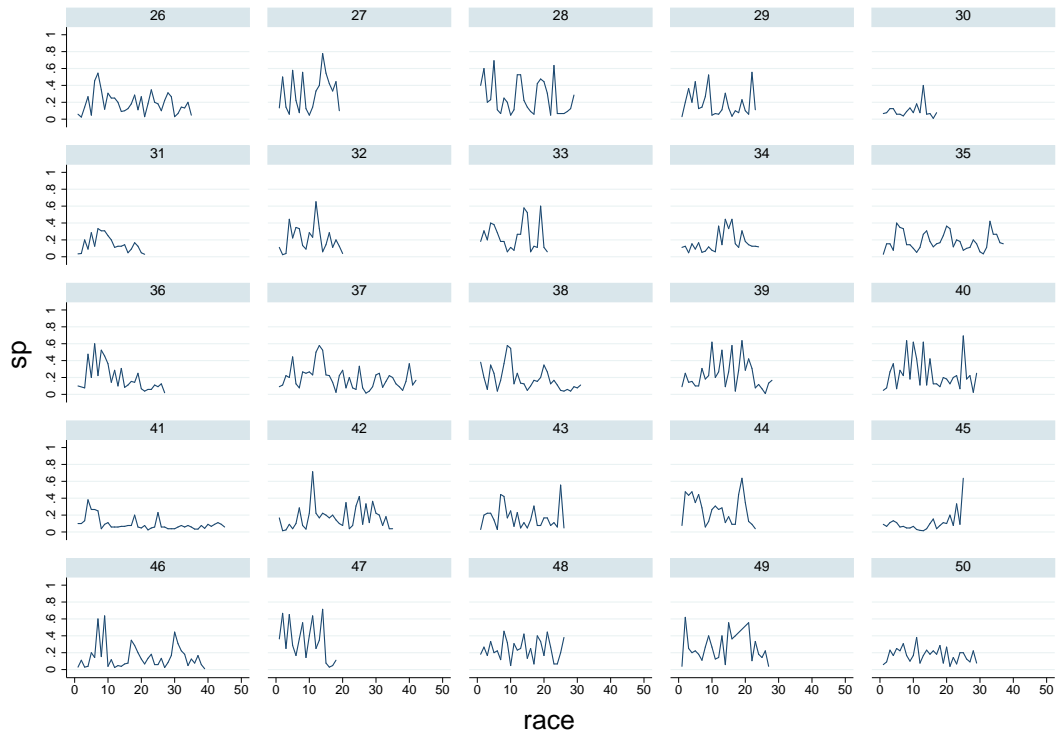


Figure 4.9 Starting Prices for Horses 51 to 75

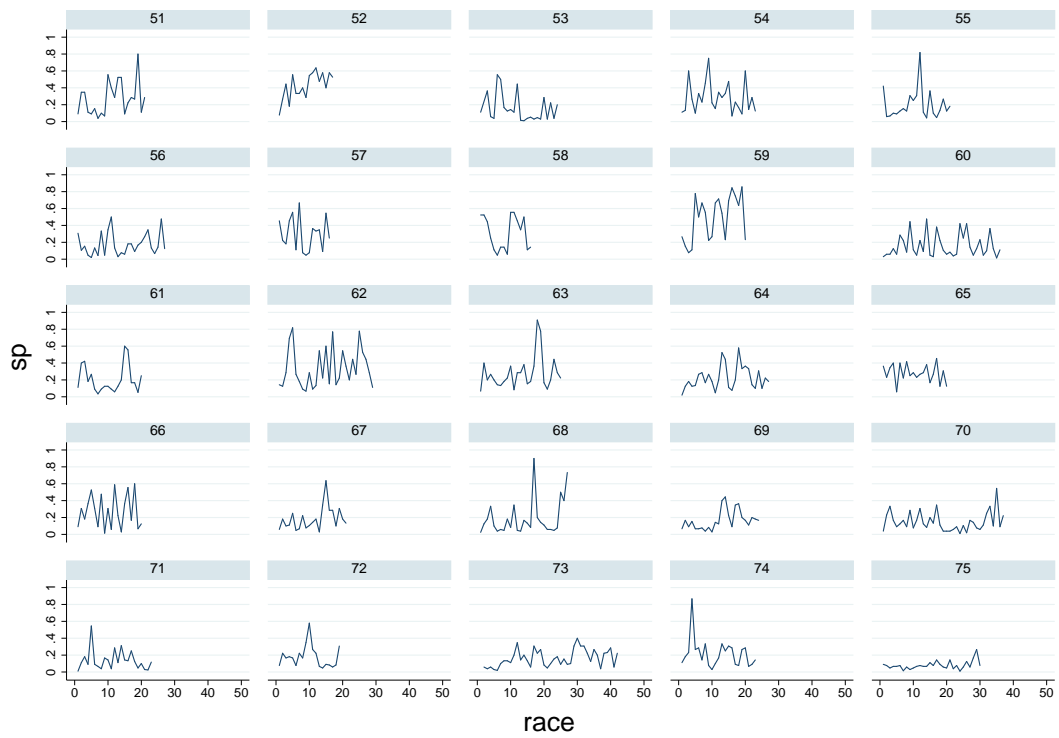
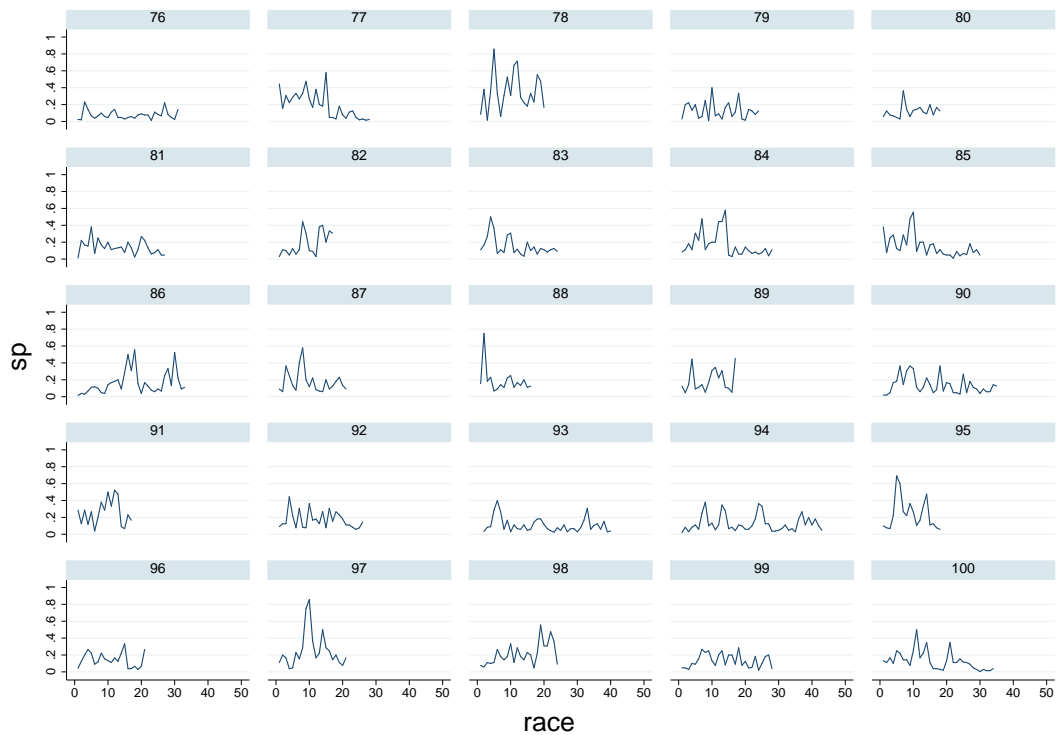




Figure 4.10 Starting Prices for Horses 76 to 100



Graphs by horse

As measured by its winning potential, high-quality horses always attract more attention from both punters and bookmakers. Measurements of winning potential include races won and earnings. The more the level stake, the higher quality the horse has. In this paper **winning potential** is equal to one minus its final position divided by total runners in the race. If we hold total runners constant, the higher value implies the higher rankings (1 is the highest ranking in this case). Given the fact that the total number of runners varies from race to race in reality, 75 percent could still imply the first place if the total runners equal to four. But still we can conclude that, without loss of generality, the higher the value of winning potential, the better the horse performs.

Age is assumed, to some extent, to be a factor that has an influence on the performance of the horse. As races are restricted to particular groups, such as maiden, juvenile or state-bred, by simply including the age factor to the model we can control for the imprecision caused by the different level of the age of the horse, especially in the case of juvenile races when horses have fewer past performance

records. The effect of age on winning probability may be positive or negative. The older the horse is, the more experience it can get, thus the higher the winning probability. Up to a certain point, horses are too old to surpass other runners. Thus we put age squared into the regression model due to the fact that the relationship between the effect of the age and the price is non-linear.

The effect of weight is to make the field more competitive, so the racing secretary assigns more weight to better horses. Consequently, a higher weight would lead to a decline in the winning probability, *ceteris paribus*. However, higher weight levels always follow the higher quality horses, so there may have a positive relationship between the weight and the horse. To conclude, weight to a horse is likely to be a double-edged sword.

When it comes to the objective conditions of a horse race, taking racecourses, distance, post position, surface and weather into consideration seems to be a workable method. Unfortunately, post position (which is the horse's stall in the starting gate and theoretically an inside post position would improve the probability of a horse winning because a slightly shorter race distance is involved) and weather are currently unavailable in this study.

A horse running at a longer distance may not perform very well as it introduces more uncertainty. With longer-distance races, horses are needed to be equipped with the different requirements of pace, stamina and speed and also there is more time for random events to occur. Luck, jockey's quality or trainer's strategy may become more important. Therefore, the length of the race will have a negative impact on the performance.

Racecourses *per se* implicitly reflect the wealthiness of the neighbourhood and in turn, prosperous people tend to place more bets than poor folks, which would affect the bookmaker's pricing mechanism. Larger and more competitive racecourses could result in a bigger margin, so we pick out 15 racecourses from both the UK and Ireland to track its influence by setting dummy variables. It takes the value of one if horses run at these racecourses and zero otherwise.

As can be imagined that the going conditions can dramatically affect the outcome of a horse race, so our model consists of the following ten goings: Good, Good to Firm, Good to Soft, Good to Yielding, Heavy, Soft, Soft to Heavy, Standard, Yielding and Yielding to soft. As an indicator variable, it equals one if a horse had run on any of aforementioned condition, and otherwise takes the value of zero. Note that there are 19 missing data in this group.

Horse racing classification system in the UK and Ireland is also very important. Understanding the class system can be of benefit to the game players. “Each horse has its own optimum level, whether it improves from a lower grade, it will eventually reach a class of race where it simply cannot match the performance of other contestants in the race”<sup>5</sup>. In other words, if a horse enters in a grade that is too high, it will be outclassed by speedier runners. There are six levels of the class included in this paper. Similarly, we assign value one to horses that belong to any level of the group and zero otherwise. One issue on the classification might be that there is a great deal of missing data (639 missing observations) in the sample.

The last component is jockey’s characteristics, which may be of secondary importance in determining a horse’s overall performance. Provided that the data on jockeys are not available, so this factor will not be considered in our model for now.

All variables are listed in Table 4.3.

---

<sup>5</sup> Accessible via: <http://www.horseracing-and-bettingsystems.com/horse-racing-class.html>.

Table 4.3 Variables in the Empirical Model

Name	Brief Explanation	Name	Brief Explanation
Starting Prices (SP) The Age of the Horse (AGE) AGE <sup>2</sup>	Avg £0.2076; equals 1/(1+odds)	Kempton Park (KEMPTON) Lingfield Park (LINGFIELD) Newbury (NEWBURY) Sandown Park (SANDOWN) York (YORK)	4.27% of all races 2.64% of all races 2.41% of all races 2.64% of all races 2.48% of all races
<b>Winning Potential:</b> LAGWIN1 LAGWIN2 Weight (WEIGHT) Distance (DIST)	Equals 1-(position/total runners) Lag behind one period Lag behind two periods Changed from imperial units to metric units and took logarithm Took logarithm and noted that 1 furlong = 1/8 mile	Curragh (CURRAGH) Fairyhouse (FAIRYHOUSE) Leopardstown (LEOPARDSTOWN) Navan (NAVAN) Punchestown (PUNCHESTOWN)	2.33% of all races 2.21% of all races 4.7% of all races 1.98% of all races 5% of all races
<b>Going Dummies (2558 obs):</b> Good (GD) Good to Firm (GDFM) Good to Soft (GDSFT) Good to Yielding (GDYLD) Heavy (HVVY) Soft (SFT) Soft to Heavy (SFTHVY) Standard (STD) Yielding (YLD) Yielding to Soft (YLDSFT)	25.76% of all races 11.61% of all races 15.21% of all races 2.19% of all races 9.3% of all races 18.96% of all races 3.36% of all races 7.7% of all races 2.74% of all races 1.92% of all races	<b>Classification Dummies (1938 obs):</b> C1 C2 C3 C4 C5 C6	38.7% of all races 17.7% of all races 12.23% of all races 17.34% of all races 10.63% of all races 34.06% of all races
<b>Racetrack Dummies (2577 obs):</b> Aintree (AINTREE) Ascot (ASCOT) Cheltenham (CHELTENHAM) Doncaster (DONCASTER) Haydock Park (HAYDOCK)	3.84% of all races 3.92% of all races 8.77% of all races 2.17% of all races 3.03% of all races		

### 4.3.3 Statistic Tests

Test for fixed versus random effects, tests for cross-sectional dependence and for time fixed effects are firstly presented just in order to fully appreciate the estimation results.

#### Hausman Test

In order to test for the presence of horse-specific fixed effects, it is common to run a Hausman test and its null hypothesis states that the random effects model is valid. The Hausman test statistic for fixed versus random effects specification is 51.41 (with P-value = 0.0000) as depicted in Table 4.4, which means that the null hypothesis of no fixed-effects is rejected at the 1% level of significance. Therefore, we can conclude that the pooled OLS estimation gives inconsistent coefficient estimates by the standard Hausman test. As a result, the fixed effects regression (4.2) should be applied.

Obviously, the above standard Hausman test that we perform is inconsistent if cross-sectional dependence is present. Furthermore, we have proved the existence of cross-sectional dependence in the following subsection. Thus, we perform a Hausman test with Driscoll and Kraay standard errors. The test statistic for fixed versus random effects specification is 11.56 (with  $p$ -value=0). The null hypothesis of no FE can be rejected at the 1% level of significance.

Combining all the Hausman test results, we reject the null hypothesis at any level of significance. We therefore conclude that statistical inference from the FE regression tends to be valid even if cross-sectional dependence is present.

#### Cross-sectional Dependence Test

In order to see if the residuals from a fixed effects estimation of the regression model (4.2) are correlated across horses, Pesaran's Cross-sectional Dependence test (Pesaran's CD test) is performed. The null hypothesis of the CD test states that the residuals are cross-sectionally uncorrelated. Correspondingly, the test's alternative

Table 4.4 Standard Hausman Test

	Coefficients			
	(b) FE	(B) RE	(b-B) Difference	sqrt(diag(V_b-V_B)) S.E.
LAGWIN1	.082758	.0937364	-.0109784	.0018605
LAGWIN2	.0283163	.0399569	-.0116406	.0018816
AGE	-.0019486	-.0047611	.0028125	.0034257
AGE <sup>2</sup>	-.0004269	-.000215	-.0002119	.0002698
WEIGHT	.4823472	.5627896	-.0804425	.0290235
DIST	-.102901	-.0872259	-.0156751	.0094921

**b** = consistent under  $H_0$  and  $H_a$ ; obtained from *xtreg*

**B** = inconsistent under  $H_a$ , efficient under  $H_0$ ; obtained from *xtreg*

Test:  $H_0$ : difference in coefficients not systematic

$$chi2(6) = (\mathbf{b} - \mathbf{B})'[(\mathbf{V}_b - \mathbf{V}_B)^{-1}](\mathbf{b} - \mathbf{B}) = 51.41$$

$$Prob > chi2 = 0.0000$$

hypothesis is that cross section dependence is present. In this study, we perform the standard Cross-sectional Dependence Test in Eviews 9, Pesaran's statistic of cross-sectional independence equals 2.19 (with  $p$ -value = 0.0285). Therefore, it comes as no surprise that Pesaran's CD test rejects the null hypothesis of cross-sectionally independence at the 5% level. Compared with other standard errors' calibration, Driscoll and Kraay standard errors are much more appropriate since they are very robust to very general forms of cross-sectional and temporal dependence.

### Race-fixed Effects

In order to test whether or not race-fixed effects are needed, we use the command **testparm** in STATA. It is a joint test to see if the dummies for all races are equal to zero (i.e., no race fixed effects are needed). There is evidence that the F statistic for race fixed effects is 1.26 (with  $p$ -value = 0.1241), therefore, the null hypothesis of no race fixed effects is failed to reject. To conclude, the race fixed effects should not be included in the model.

#### 4.3.4 Estimation Results

The model (4.2) is estimated by utilising 100 horses' lifetime results in the study dataset. The associated empirical results are displayed in the following six tables. We report five sections of fixed-effects (within) regression results as well as Driscoll and Kraay standard errors.

The first section runs the fixed-effects (within) regression equation without any dummy variables. The first column in Table 4.5 describes the dependent variable in this model. Adjusted R-squared and the number of observations are reported in the last two columns respectively. The rest are the explanatory variables. In examining the results displayed below, it is noted that two standard errors - FE standard errors and DK standard errors - are showed in the two parentheses respectively. 14 different equations are regressed to test whether or not adding or deleting variables would have an influence on the outcome variable (i.e., changing the sign of the coefficient or its significance). The F-test of the null hypothesis that all the coefficients in the model are equal to zero is rejected at the 0.1% level of significance for most of the equations, except for equations (1) and (10). The first equation is the starting prices regressed on a constant over the entire sample. The estimated value is 0.2076, which also is the average value of SP. Note that the number of observations is 2568 and there are 9 missing data. The reason why we include the constant in the regression model is because whether or not the bookmaker considers price factors, there is always a starting price for each horse. It is crystal clear from equations (2) to (7) that the signs of the coefficients on *LAGWIN1* are all positive at the 1% level of significance under both FE and DK standard errors, which is consistent with our prior expectation that there is a positive relation between the winning potential and the price. In other words, higher values of the winning potential imply higher rankings in the last race, therefore a higher price should be anticipated when bookmakers set up the price of the horse in the present race. Similar analysis is applied to *LAGWIN2* where the winning potential is lagged behind two races. The signs of the coefficients are statistically significant and positive, but the values are lower than the coefficients on *LAGWIN1* with

approximately 0.05 decrease. Based on this finding, we can conclude that when pricing a horse in a race, the information from the last race is more needed by bookmakers than the one before. The effect of the age on the price is insignificant at any level of significance, so does the coefficient of age squared, under both methods of standard errors, which are inconsistent with *a priori* expectation that there is a non-linear relationship between the age and the price. As mentioned before, the weight the horse is assigned to carry may have a positive or negative impact on the horses' performance, that is, the horses' prices. Results in Table 4.5 present a positive effect on the price. As the weight has been changed from imperial units to metric units and then has taken the logarithm of kilogrammes, the coefficient, for example in Equation (9), on *WEIGHT* means that Starting Prices increase, on average, by 0.004188% when one more kilogramme is added to the horses. Even when we add more predictors to the model, the coefficients do not change too much and they all are statistically significant at the 1% level. For the distance factor, we prove that the statistical inference on the *DIST* is consistent with a *priori* theoretical expectation is that the longer distance leads to uncertainty which would undermine the horse's performance. Therefore, it has a negative effect on the price and prices should reduce correspondingly. For example, suppose that all other factors hold constant, one extra furlong leads to a decline in price by 0.001029% in Equation (7). Last but not least, the number of observations is changing as some of the variables have missing data, thus STATA automatically deletes it. The adjusted R-squared reported in Table 4.5 is obtained from the horse-demeaned OLS regression, which suggests that absorbing horse dummies helps speed up exploratory work and provide quick feedback about whether a dummy variable approach is worthwhile. It is noted that the estimated coefficients and standard errors will not be altered in this method. Furthermore, a low value of the adjusted  $R^2$  should be expected in the context of horse races since there are other unpredicted factors existed, such as weather and a variety of constraints placed on the competing horses. In the meantime, it is worth noticing that the adjusted  $R^2$  is improving from equation to equation and the highest value 0.3211 is achieved in Equation (7) when all explanatory variables are contained in the model.



The second section runs the fixed-effects (within) regression model with GOING dummy variables. GOING represents the condition of the surface at racetracks. Results in Table 4.6 explicitly shed light on how ten different going conditions affect the main variables. Signs on *LAGWIN1*, *LAGWIN2*, *AGE*, *AGE*<sup>2</sup>, *WEIGHT* and *DIST* have not changed at all, so do the significance. The adjusted *R*<sup>2</sup> (equals 0.3404) is improved in this case. Our results also weakly support the assumption that going conditions can dramatically affect the outcome of a horse race as 7 going dummies are statistically significant at the 1%, 5% and 10% level respectively.

The third section can be divided into two parts: the regression equation with ten British racecourses' dummies is displayed in Table 4.7 and with five Irish racecourses' dummies presented in Table 4.8. As can be seen from Table 4.7, when adding British racecourse dummies, the significance of *AGE* and *AGE*<sup>2</sup> is not improving. The last feature in this table is that almost all racecourses reveal a negative impact on the price except for Lingfield Park racecourse, of which DONCASTER dummy variable is not significant at any level under both methods. Therefore, competing horses at the larger and more competitive racecourses in Britain tend to be set lower prices to attract more punters. For racecourses in Ireland, the signs and the significance of all the coefficients remain the same. 3 out of 5 dummies show no significance. *Curragh* is statistically significant at the 5% level under the FE standard error, while under the DK standard error it is statistically significant at the 1% level. *Punchestown* is statistically significant at the 5% level only under the DK standard error. To sum up, the factor of racecourses in Britain is of great importance than in Ireland when it comes to setting the price of the horse.

The fourth section runs the regression equation with the classification dummies and the corresponding results are reported in Table 4.9. Class 6 is automatically deleted by STATA because of the collinearity problem. Class 1, Class 2, Class 3 and Class 4 are statistically significant, not only with FE standard errors but also with DK standard errors. Including class dummies makes the coefficients on *AGE* and *AGE*<sup>2</sup> become significant and it also demonstrates a simple positive relationship

between the age and the price and a negative correlation between  $AGE^2$  and the price.

The fifth section adds all 37 available predictors to the model and also takes horse dummies into consideration. The results can be found in Table 4.10. The first thing that should be noticed in this regression is the significantly improved adjusted R-squared (0.472), which is the highest value among all the equations. Thus, it may be concluded that the model is explaining a statistically significant amount of the variation in horse racing prices. The coefficient on  $AGE$  is significant under the calibration of both standard errors at the 1% level, so does  $AGE^2$ . In comparison with the classification, going conditions are not important anymore. Racecourse dummies still play an important role in this case.

Table 4.5 Fixed-effects Regression Results

Eq.	DV	Constant	LAGWIN1	LAGWIN2	AGE	AGE <sup>2</sup>	WEIGHT	DIST	AdjR <sup>2</sup>	Obs
(1)	SP	0.2076 (0.003)*** (0.0071)***							0.2484	2568
(2)	SP	0.15015 (0.0083)*** (0.0106)***	0.0984 (0.0124)*** (0.0121)***						0.2786	2333
(3)	SP	0.1357 (0.0113)*** (0.0105)***	0.0893 (0.0128)*** (0.0106)***	0.0349 (0.013)*** (0.0102)***					0.2934	2130
(4)	SP	0.1686 (0.0352)*** (0.0328)***	0.0832 (0.0128)*** (0.0112)***	0.0284 (0.0131)** (0.01)**	0.0011 (0.0111) (0.0085)	-0.00087 (0.00088) (0.00066)			0.2977	2130
(5)	SP	-1.4685 (0.245)*** (0.2247)***	0.0796 (0.0127)*** (0.0114)***	0.0254 (0.01296)* (0.0095)**	-0.0154 (0.0113) (0.0095)	0.0002 (0.00088) (0.00075)	0.4045 (0.0599)*** (0.0547)***		0.3129	2129
(6)	SP	-1.549 (0.2441)*** (0.22)***	0.0873 (0.0126)*** (0.0109)***	0.0332 (0.0128)*** (0.0091)***			0.4804 (0.0609)*** (0.0533)***	-0.1205 (0.0187)*** (0.0279)***	0.3195	2129
(7)	SP	-1.5725 (0.2444)*** (0.2161)***	0.0827 (0.0127)*** (0.0115)***	0.0283 (0.0129)** (0.0094)**	-0.0019 (0.0115) (0.0102)	-0.0004 (0.00088) (0.00076)	0.4823 (0.0615)*** (0.0547)***	-0.1029 (0.0203)*** (0.0257)***	0.3211	2129

Notes: 1. \*\*\* indicates significance at the 1% level.

\*\* indicates significance at the 5% level.

\* indicates significance at the 10% level.

2. Standard errors of parameter estimates are in the first parenthesis below each coefficient.

Driscoll-Kraay standard errors are reported in the second parenthesis.

Table 4.4: Fixed-effects Regression Results (continued)

Eq.	DV	Constant	LAGWIN1	LAGWIN2	AGE	AGE <sup>2</sup>	WEIGHT	DIST	AdjR <sup>2</sup>	Obs
(8)	SP	0.1325 (0.027)*** (0.0553)**			0.0327 (0.009)*** (0.0145)**	-0.003 (0.0007)*** (0.0008)***			0.2546	2568
(9)	SP	-1.5489 (0.2251)*** (0.2344)***					0.4188 (0.0537)*** (0.0552)***		0.2662	2567
(10)	SP	0.2317 (0.0425)*** (0.124)*						-0.0089 (0.0157) (0.0445)	0.2482	2568
(11)	SP	-1.6 (0.227)*** (0.1987)***			0.0166 (0.0092)* (0.015)	-0.002 (0.0007)*** (0.0009)**	0.4259 (0.0554)*** (0.047)***		0.2718	2567
(12)	SP	0.1621 (0.0447)*** (0.1621)			0.0359 (0.0098)*** (0.0132)***	-0.0032 (0.0007)*** (0.0008)***		-0.0152 (0.0183) (0.0307)	0.2545	2568
(13)	SP	-1.6746 (0.2276)*** (0.2063)***					0.485 (0.0569)*** (0.0455)***	-0.0564 (0.0164)*** (0.044)	0.2694	2567
(14)	SP	-1.665 (0.2278)*** (0.1897)***			0.0264 (0.0098)*** (0.015)*	-0.0025 (0.0007)*** (0.0009)***	0.4673 (0.0572)*** (0.0469)***	-0.0538 (0.0187)*** (0.03)*	0.274	2567

Notes: 1. \*\*\* indicates significance at the 1% level.

\*\* indicates significance at the 5% level.

\* indicates significance at the 10% level.

2. Standard errors of parameter estimates are in the first parenthesis below each coefficient.  
Driscoll-Kraay standard errors are reported in the second parenthesis.

Table 4.6 Fixed-effects Regression Results (Including Going Conditions)

<i>SP</i>	<i>Coefficient</i>	<i>Standard Errors</i>	<i>Driscoll-Kraay Standard Errors</i>
<i>Constant</i>	-1.6233	0.2459***	0.2178***
<i>LAGWIN1</i>	0.0853	0.0126***	0.0108***
<i>LAGWIN2</i>	0.0312	0.0128**	0.0089***
<i>AGE</i>	0.0013	0.0115	0.00997
<i>AGE<sup>2</sup></i>	-0.0006	0.00088	0.0007
<i>WEIGHT</i>	0.4708	0.0615***	0.0559***
<i>DIST</i>	-0.0902	0.0202***	0.0254***
<i>GD</i>	0.0214	0.0284	0.0254
<i>GDFM</i>	0.0631	0.0293**	0.0253**
<i>GDSFT</i>	0.0297	0.0288	0.0223
<i>GDYLD</i>	0.0092	0.0366	0.0269
<i>HVY</i>	0.0735	0.0298**	0.0276**
<i>SFT</i>	0.0628	0.0287**	0.0298**
<i>SFTHVY</i>	0.0741	0.0342**	0.0308**
<i>STD</i>	0.1037	0.0306***	0.0333***
<i>YLD</i>	0.0626	0.0339*	0.0305**
<i>YLDSFT</i>	0.0764	0.0379**	0.0248***

- Notes: 1. \*\*\* indicates significance at the 1% level.  
 \*\* indicates significance at the 5% level.  
 \* indicates significance at the 10% level.  
 2. There are 2117 observations in this model.  
 3. Adjusted R-squared equals 0.3404.

Table 4.7 Fixed-effects Regression Results (Including British Racecourses)

<i>SP</i>	<i>Coefficient</i>	<i>Standard Errors</i>	<i>Driscoll-Kraay Standard Errors</i>
<i>Constant</i>	-1.6162	0.236***	0.1911***
<i>LAGWIN1</i>	0.0866	0.0122***	0.012***
<i>LAGWIN2</i>	0.0396	0.0124***	0.0097***
<i>AGE</i>	-0.0064	0.0112	0.0104
<i>AGE</i> <sup>2</sup>	0.00016	0.0009	0.0007
<i>WEIGHT</i>	0.4792	0.0595***	0.0462***
<i>DIST</i>	-0.0738	0.0198***	0.0238***
<i>AINTREE</i>	-0.1158	0.0167***	0.0158***
<i>ASCOT</i>	-0.042	0.016***	0.0105***
<i>CHELTENHAM</i>	-0.1484	0.0124***	0.0125***
<i>DONCASTER</i>	-0.0175	0.0205	0.0226
<i>HAYDOCK</i>	-0.0477	0.0183***	0.0114***
<i>KEMPTON</i>	-0.0444	0.0162***	0.0135***
<i>LINGFIELD</i>	0.0563	0.0205***	0.0149***
<i>NEWBURY</i>	-0.0599	0.0208***	0.0235**
<i>SANDOWN</i>	-0.0659	0.0197***	0.0158***
<i>YORK</i>	-0.0804	0.0201***	0.0219***

Notes: 1. \*\*\* indicates significance at the 1% level.

    \*\* indicates significance at the 5% level.

    \* indicates significance at the 10% level.

2. There are 2129 observations in this model.

3. Adjusted R-squared equals 0.3763.

Table 4.8 Fixed-effects Regression Results (Including Irish Racecourses)

<i>SP</i>	<i>Coefficient</i>	<i>Standard Errors</i>	<i>Driscoll-Kraay Standard Errors</i>
<i>Constant</i>	-1.5663	0.2455***	0.2207***
<i>LAGWIN1</i>	0.0839	0.0127***	0.0113***
<i>LAGWIN2</i>	0.0279	0.0129**	0.0097***
<i>AGE</i>	-0.00085	0.0115	0.0099
<i>AGE<sup>2</sup></i>	-0.00049	0.00088	0.00076
<i>WEIGHT</i>	0.4825	0.0617***	0.0569***
<i>DIST</i>	-0.1059	0.0204***	0.0272***
<i>CURRAGH</i>	-0.0548	0.0223**	0.0191***
<i>FAIRYHOUSE</i>	-0.0313	0.0229	0.0303
<i>LEOPARDSTOWN</i>	-0.00199	0.0165	0.0101
<i>NAVAN</i>	0.0132	0.0242	0.0185
<i>PUNCHESTOWN</i>	-0.0273	0.0167	0.0128**

Notes: 1. \*\*\* indicates significance at the 1% level.

    \*\* indicates significance at the 5% level.

    \* indicates significance at the 10% level.

2. There are 2129 observations in this model.

3. Adjusted R-squared equals 0.3231.

Table 4.9 Fixed-effects Regression Results (Including Classification)

<i>SP</i>	<i>Coefficient</i>	<i>Standard Errors</i>	<i>Driscoll-Kraay Standard Errors</i>
<i>Constant</i>	-0.863	0.2923***	0.2963***
<i>LAGWIN1</i>	0.0921	0.0137***	0.0075***
<i>LAGWIN2</i>	0.0617	0.0137***	0.0125***
<i>AGE</i>	0.0538	0.0129***	0.0112***
<i>AGE<sup>2</sup></i>	-0.00358	0.00099***	0.0009***
<i>WEIGHT</i>	0.2816	0.0732***	0.057***
<i>DIST</i>	-0.07228	0.0229***	0.0225**
<i>C1</i>	-0.2595	0.0267***	0.0251***
<i>C2</i>	-0.1964	0.0259***	0.023***
<i>C3</i>	-0.1358	0.0257***	0.0206***
<i>C4</i>	-0.0818	0.0246***	0.0199***
<i>C5</i>	-0.0397	0.0249	0.0206*

- Notes: 1. \*\*\* indicates significance at the 1% level.  
 \*\* indicates significance at the 5% level.  
 \* indicates significance at the 10% level.
2. There are 1674 observations in this model.  
 3. Adjusted R-squared equals 0.4243.  
 4. C6 is omitted due to collinearity.



Table 4.10 Fixed-effects Regression Results (All)

<i>SP</i>	<i>Coefficient</i>	<i>Standard Errors</i>	<i>Driscoll-Kraay Standard Errors</i>
<i>Constant</i>	-1.0846	0.2919***	0.3183***
<i>LAGWIN1</i>	0.0889	0.0132***	0.0083***
<i>LAGWIN2</i>	0.0689	0.0133***	0.0111***
<i>AGE</i>	0.0401	0.0127***	0.012***
<i>AGE<sup>2</sup></i>	-0.0025	0.00097***	0.00089***
<i>WEIGHT</i>	0.3193	0.0729***	0.0653***
<i>DIST</i>	-0.0448	0.0225**	0.0196**
<i>GD</i>	0.0076	0.0281	0.0246
<i>GDFM</i>	0.0336	0.02908	0.0288
<i>GDSFT</i>	0.0202	0.0282	0.0222
<i>GDYLD</i>	-0.0698	0.0432	0.0323**
<i>HVY</i>	0.0379	0.0307	0.0275
<i>SFT</i>	0.0244	0.0286	0.0295
<i>SFTHVY</i>	0.0277	0.0446	0.0436
<i>STD</i>	0.08	0.0315**	0.0344**
<i>YLD</i>	0.0407	0.0412	0.0458
<i>YLDSFT</i>	0.0031	0.0629	0.0291
<i>AINTREE</i>	-0.0853	0.0173***	0.018***
<i>ASCOT</i>	-0.032	0.0155**	0.011***
<i>CHELTENHAM</i>	-0.1233	0.0138***	0.021***
<i>DONCASTER</i>	-0.0106	0.019	0.0228
<i>HAYDOCK</i>	-0.0389	0.0172**	0.01***
<i>KEMPTON</i>	-0.051	0.0163***	0.0141***
<i>LINGFIELD</i>	0.0159	0.0222	0.0184
<i>NEWBURY</i>	-0.0554	0.0197***	0.0222**
<i>SANDOWN</i>	-0.0474	0.0187**	0.0159***
<i>YORK</i>	-0.0629	0.0189***	0.0162***
<i>CURRAGH</i>	-0.0387	0.0328	0.0189**
<i>FAIRYHOUSE</i>	-0.0535	0.0418	0.0457
<i>LEOPARDSTOWN</i>	0.0039	0.0264	0.0288
<i>NAVAN</i>	0.061	0.0409	0.0339*
<i>PUNCHESTOWN</i>	-0.0229	0.0281	0.0292
<i>C1</i>	-0.2094	0.0264***	0.0172***
<i>C2</i>	-0.1625	0.0253***	0.02***
<i>C3</i>	-0.1137	0.025***	0.0175***
<i>C4</i>	-0.0657	0.0239***	0.0203
<i>C5</i>	-0.0301	0.0241	0.0214

Notes: 1. \*\*\* indicates significance at the 1% level.  
 \*\* indicates significance at the 5% level.  
 \* indicates significance at the 10% level.  
 2. There are 1665 observations in this model.  
 3. Adjusted R-squared equals 0.472.

## 4.4 Conclusions

This paper aims at providing empirical evidence on the determinants of starting prices in horse-racing betting markets by utilising a unique cross-sectional and time series data set. Those determinants are discovered through an empirical analysis of 100 horses with its lifetime results from the year 2006 to 2015. We contribute to the literature by answering the following two questions: (i) what are the factors that can be used to explain starting prices in horse-racing betting markets? and (ii) are these factors persistent?

Before we answer the above questions, three tests are performed because of the uniqueness of our dataset. The result of Hausman test suggests that we should employ the fixed-effects regression model. Pesaran's cross-sectional dependence test confirms the presence of cross-sectional dependence so that Driscoll and Kraay standard errors are reported as they are robust to very general forms of cross-sectional and temporal dependence.

There is a plethora of empirical evidence to suggest that the winning potential (defined as 1-position/total runners in this paper) is quite significant in explaining the price of the horse. Other factors, which include the age of the horse, the weight the horse carries and the distance of the racetrack, are also important. Our findings show that the condition of the turf, the size of the racetrack and the classification of the race could influence the price as well. All the estimation results from Table 4.5 to Table 4.10 confirm that these factors are persistent.

# Chapter 5

## Dynamic Pricing in Horseracing Betting Markets

### 5.1 Introduction

As horse racing betting markets provide economists with a ready-made laboratory and brand new data sets, there are numerous empirical papers, which have examined these markets from different perspectives by testing market efficiency theories, explaining the *favourite-longshot bias*, or detecting the incidence of insider trading and decision making under uncertainty and so on. The literature documenting these problems covers both pari-mutuel markets, which prevail in North America, and bookmakers markets in the United Kingdom and Ireland. This bookmaking system can be traced back to the late 18th century and on average the bookmakers are the winners in the *punter-versus-bookmaker* battle ever since the Parliament passed the Gaming Act in 1845<sup>1</sup>, which aroused our curiosity about the bookmakers' behaviour. In this paper, we will analyse the bookmaking market where the bookmaker serves as a market intermediary who sets betting odds or prices for each horse in any race. We aim to understand how odds are adjusted.

---

<sup>1</sup>See "All bets are off: why bookmakers aren't playing fair" from *The Guardian* (<https://www.theguardian.com/global/2015/aug/02/betting-horses-gambling-bookmakers-accounts-closed>).

Despite the popularity of sports betting, there is not much theoretical work on how bookmakers optimally set odds which guarantee non-zero profits. A series of papers by Shin (1991, 1992 and 1993) proposes a theoretical model that determines the prices of state contingent claims in the mimic market - the British horseracing betting market - with the presence of bettors who have superior information. Fingleton and Waldron (1999) generalise the Shin model by assuming that bookmakers are risk-averse instead and face transaction costs in a non-competitive market. Recently, Sandford and Shea (2013) address the problem of how bookmakers set lines<sup>2</sup>. They mainly deal with sports like football or basketball which has point spreads<sup>3</sup>. The focus of our paper is horserace betting, which is quite different from football betting especially when it comes to adjusting odds during the betting period<sup>4</sup>.

The objective of this paper is to

- Set up the general equilibrium framework, and
- Try and explain how the bookmaker adjusts the betting odds in a two-horse benchmark model.

It is worth noting that in reality a rational bookmaker initially does not change the odds constantly, especially not because of one punter's particular behaviour unless the bettor shows the appearance of bias in favour of one certain type of bet. This is based on our observation of the market in the past few years for research purpose. One possible explanation is that the informed punters do not show up until just before the end of the betting period, which is not the case. Another explanation is that the bookmaker is willing to take initial losses to collect information about the

---

<sup>2</sup>Line betting is a form of betting, which is used for events with two possible outcomes. In order to equalise the money on both sides, the bookmaker handicaps a team by setting a margin, which is referred to as the line. Accessible via: <http://bet-types.com.au/line-betting/>

<sup>3</sup>Point spread is defined as a measurement to determine the likely gap between the two teams in the final result. It attempts to find which team is more likely to win, and by how much. Accessible via: <https://www.sportingcharts.com/dictionary/nfl/point-spread.aspx>

<sup>4</sup>Horserace betting and football betting differ in two ways: (i) the outcome for horse race is "win or lose", whilst for football betting there is a point spread, which attempts to determine the likely gap between two teams and to find which team is more likely to win, and by how much; (ii) horses' in-play performance has no effect on the odds-adjusting procedure, but football players' performance influences the odds tremendously.

informed traders as they are the main source of the bookmaker's loss. We call this the *learning process*. Our benchmark model, which focuses on the statistical losses, offers a simple rule to find an optimal time to change prices.

The contribution of this paper is twofold. First, a general equilibrium model is presented by making use of the basic economic assumption that profit should at least equal cost in order to keep the market functioning. More specifically, the basic principle for bookmakers is to balance the total gains and the payoffs for the winning horse. Our intention is to discover the prices at each stage of the betting process dynamically. To the best of our knowledge, the model developed in this paper is the first of this kind for this market. In order to make the general equilibrium model tractable and obtain rich results, we start with a two-horse benchmark model by assuming that two types of punters exist in the market, i.e., informed and uninformed noise traders. Second, we make a methodological contribution by introducing optimal stopping techniques into the betting literature.

Now we explain how the process of optimally adjusting odds can be formulated as an optimal stopping problem. We present a continuous-time, two-state model for a two-horse benchmark framework in which the bookmaker changes the odds of horse A if the "high" state reveals or the odds of horse B should be adjusted if the "low" state prevails. No matter which horse is chosen to adjust, the other changes automatically. We assume that the cumulative sales of tickets on horse A, net of the cumulative sales of tickets on horse B can be modelled by an arithmetic Brownian motion with either a positive or negative trend. These are the two states of the world: *high state* implies that informed traders know that horse A is going to win; *low state* means horse B is more favourable among insiders. The true state remains unknown *ex ante*, so the bookmaker has no information about which horse the informed traders are backing. These two states decide the direction in which the odds are changed. This information is gradually disclosed through sequential observation of trades, acts as if a positive trend is observed at a certain point, the cumulative sales on horse A exceed the tickets sold on horse B, the bookmaker, therefore, determines that the informed punters know that horse A is the winning

horse and adjusts the price accordingly. If a negative trend is revealed, the odds of horse B should be changed.

Our model is analysed using Peskir and Shiryaev's (2006) study of the Bayesian problem of minimising the non-discounted cost of incorrect decisions in sequential hypothesis testing in a setting in which the true state is slowly revealed through a Brownian motion. Instead of analysing the non-discounted cost, we model our Bayesian problem with Poisson jumps. By adopting Thijssen and Bregantini's (2016) model, we obtain analytical results, and we show that there is a unique solution to the problem that gives both upper and lower boundaries in terms of the posterior belief in informed punters knowing that horse A wins. The upper boundary represents the trigger beyond which the bookmaker changes the odds of horse A, while the lower boundary means below which the odds of horse B should be changed. The "continuation region" is between these two boundaries, which implies the bookmaker should keep observing.

The main analytical result of this study is that the decision bounds get wider as the fraction of informed traders becomes larger, but the loss is non-monotonic even though the wider bounds provide more information per time period. This happens because of two opposing effects. It is quite straightforward that more insiders in the market cause tremendous losses to the bookmaker, so, on the one hand, the loss is increasing. On the other hand, a large proportion of informed traders means the learning process is very informative, so the decision will be reached sooner, which gives a declined loss function. In contrast with the standard real options literature, which shows a higher uncertainty implies wider decision bounds and a higher value of the project, our model finds that the posterior bounds get narrower in the volatility as uncertainty enters as noise traders in the observations. More uninformed traders provide less useful information, thereby explaining the narrow bounds. The loss is increasing in the volatility as we expect the decision should be taken at the early stage. Otherwise, a higher volatility makes the learning process less informative, which leads to an increased loss.

The rest of the paper is structured as follows. We review related literature in Section 5.2. In Sections 5.3 and 5.4, we present our model with main results. The comparative statics analysis is effected in Section 5.5. Conclusions are given in Section 5.6.

## 5.2 Literature Review

This paper analyses an insider trading problem using optimal stopping theory. Before we proceed to the analysis, it is helpful to briefly review the closely related lines of research.

*Insider Trading.* The idea of insider trading is first introduced by Bagehot in 1971. He specifically mentions that “every time one investor benefits from a trade, ..., another loses” and the market maker always loses to traders with special information and gains with liquidity-motivated transactors. The gains from noise traders must exceed the losses to the informed traders. In our study, the bookmaker is also losing to the insiders and breaks even on the uninformed one. Ever since then, the presence of traders with superior information gets wide attention in the area of the finance literature. Kyle (1985) develops a dynamic model of insider trading where prices follow a Brownian motion and the market maker can only see the order imbalance. With sequential auctions, the informed trader makes positive profits while the noise traders provide camouflage, but in the end all private information is incorporated into prices. In the Glosten-Milgrom (1985) model, traders with special information and liquidity traders arrive as a Poisson process. The bid-ask spread can be explained by adverse selection as well as the exogenous arrival pattern of traders. The presence of insiders leads to a positive spread. Since the transaction prices are informative, the spread declines with the trade. Back and Baruch (2004) show that the Kyle and Glosten-Milgrom models were essentially modelling the same phenomenon and they prove that the equilibria of the Glosten-Milgrom model can converge to the equilibrium of the Kyle model. Shin (1991, 1992 and 1993) extends the Glosten-Milgrom model into the horse-racing context and provides

the empirical support for measuring the degree of insider trading in the British horse-racing betting markets. Fingleton and Waldron (1999) generalise the Shin model. They show that informed trading, the bookmaker's attitude towards risk and the bookmaker margin are three factors that determine optimal odds. Our model only includes insider trading because we assume that the margin is competed away. Schnytzer and Shilony (1995) provide the empirical evidence for the presence of inside information on the Melbourne horse betting markets. Based on the two segregated markets, they show that bettors with inside information, even someone who is exposed to 'second hand' information, can change their behaviour and get a rise in payoffs. Positive information can be of significance to predict a race's outcome. Gabriel and Marsden (1990) discover that Tote (Totalisator) returns on winning bets exceed starting price returns which are paid out by bookmakers. Cain, Law and Peel (2001) further prove that bookmakers pay more generously on favourites than the Tote, but less generously on high odds bets. The difference between those two papers is that the latter considers the incidence of insider trading in the paper<sup>5</sup>.

*Optimal Stopping Problem.* There is a vast of literature that analyses optimal stopping problems in real options models (Dixit and Pindyck (1994)), but there is no such analysis in the context of horse-racing betting markets. We introduce the optimal stopping model into the field to figure out the price adjusting mechanism - when and how the bookmaker changes the odds to minimise the losses. Our paper is an extension of the simple sequential hypothesis testing with which the non-discounted expected costs are minimised. It was pioneered by Wald (1945) and further developed by Shirayev (1967) who laid the foundation of this subject. The optimal stopping problem in continuous time is analysed by Shirayev (1978), Øksendal (2003) and Peskir and Shirayev (2006). Since the solution of the sequential hypothesis testing problem can be reduced to the solution of the *Stefan problem* with two boundaries and the solution method includes solving a partial differential equation with smooth fit conditions, Shirayev (1967) proves the uniqueness of the

---

<sup>5</sup>Cain, M., Law, D. and Peel, D. A. (2001). The incidence of insider trading in betting markets and the Gabriel and Marsden anomaly. *The Manchester School*. Vol. 69, No. 2, p. 197.



solution and provides analytical results. This is the starting point of our model. The distinguishing feature of our setting, however, is that we are detecting a positive or negative trend. Mostly, our paper is based on Thijssen and Bregantini (2016) who study a two-sided investment/abandonment problem through a costly sequential experiment that can provide the information about the true states of nature by maximising the discounted expected payoffs, but in this study we minimise the discounted cost of errors.

### 5.3 General Model

#### *The Key Features of the Market.*

The horse-racing betting market in the United Kingdom has some striking features because of the odds-setting procedure in which bookmakers determine odds rather than the parimutuel method in North America in which odds are proportional to money wagered. A great number of traders are in the market, among which some of them are endowed with superior information, and the rest are just noise traders. Odds are changing all the time during the betting period until the race commences. These provide us with a replicated market, which is a particularly simple example of a contingent claims market. Consider an  $n$ -horse race which corresponds to  $n$  states of the world in a market for contingent claims, the  $i$ th state corresponds to the outcome in which the  $i$ th horse wins the race<sup>6</sup>. The main difference is that in this market there is no need to take its future dividends or fundamentals into account because of the fact that by the end of the race there is a commonly acknowledged outcome, which gives certain payoffs. Another difference is that the on-course horse racing market usually convenes for around half an hour. So in general we present a finite time model.

#### *Dates and Assets (Horses)*

---

<sup>6</sup>Shin. H. S. (1993). Measuring the Incidence of Insider Trading in a Market for State-Contingent Claims. *The Economic Journal*, Vol. 103, No. 420, p. 1142.

In the on-course horse racing betting market, there are two significant dates - at time  $t = 1$  the market opens for trading; and at time  $t = T$  the market closes right before the race commences. During the betting period there are finite dates at which the only one risk-neutral bookmaker that we assume in the paper changes betting odds based on the observations of the cumulative sales for each horse. Let  $t = \{1, 2, \dots, T\}$  denote a set of times and at each time there is a new set of prices for each horse in which it could be the change for one horse or adjusting many horses simultaneously.

Horses are indexed by  $i = \{1, 2, \dots, N\}$ , which implies there are  $N$  horses in a race. Let  $P_i^t$  denote the price of the ticket for horse  $i$  at time  $t$ <sup>7</sup>. Note that there are two terminologies; one is the opening odds which are the set of prices for  $N$  horses at time  $t = 1$  and the other is called the starting prices which are the prices at time  $t = T$ . Denote  $Q_i^t$  be the number of tickets sold for horse  $i$  at time  $t$ .

#### *Bookmaker's Problem*

*No-Arbitrage Condition.* “Arbitrage” in finance theory is defined as a trading strategy that generates a risk-less profit. In the context of horse racing betting market, the arbitrage opportunity implies that if, at any time, the sum of all the odds is no larger than £1, there is a chance that punters can buy all horses and get £1 for sure in the end no matter which horse wins. In order to avoid that, at each point in time, the bookmaker should set the prices which satisfy the following conditions:

$$P_1^1 + P_2^1 + \dots + P_N^1 \geq 1$$

$$P_1^2 + P_2^2 + \dots + P_N^2 \geq 1$$

⋮

$$P_1^T + P_2^T + \dots + P_N^T \geq 1$$

In general, the goal of the bookmaker is to maximise the profits or at least make ends meet. For convenience, we assume zero expected profits condition in our

---

<sup>7</sup>For simplicity, we rearrange  $P_i^t = \frac{1}{\text{odds}+1}$ , so  $0 < P_i^t \leq 1$ . In this setup, the bookmaker pays out £1 if horse  $i$  wins, otherwise (s)he gains  $P_i^t$ .

model and the conditions above ensure that there is no arbitrage opportunity either.

So the total revenue is given by

$$\begin{aligned}
& P_1^1 * Q_1^1 + P_2^1 * Q_2^1 + \dots + P_N^1 Q_N^1 + \\
& P_1^2 * Q_1^2 + P_2^2 * Q_2^2 + \dots + P_N^2 Q_N^2 + \\
& \vdots \\
& P_1^T * Q_1^T + P_2^T * Q_2^T + \dots + P_N^T Q_N^T \\
& = \sum_{i=1}^N \sum_{t=1}^T P_i^t Q_i^t
\end{aligned}$$

Since we assume the winning horse gets everything, the cost for the bookmaker is unknown *ex ante*, but it will be realised at the end of the race. Therefore, the *ex post* total payout is  $\sum_{t=1}^T Q_i^t$  if horse  $i$  wins. In both theory and reality, the total revenue should be no less than the total payout for the bookmaker, that is,  $\sum_{i=1}^N \sum_{t=1}^T P_i^t Q_i^t \geq \sum_{t=1}^T Q_i^t$  for  $i = 1, 2, \dots, N$ . Otherwise there is no incentive to keep the market functioning. However, the winning horse is not directly observable *ex ante*. How does the bookmaker adjust the odds over the betting period under this circumstance? It is when the zero expected profits condition kicks in. In general this condition should only be satisfied at final time  $T$ , but in our model we assume that the bookmaker has to make sure the expected profits are zero at each point in time when (s)he changes the prices, that is,  $ER^t = EC^t$ , where  $ER^t$  denotes the expected value of revenue at time  $t$  and  $EC^t$  is the expected value of cost at time  $t$ . For any  $i$  and  $t$  ( $\forall i, \forall t$ ), the following condition meets

$$\begin{aligned}
E \left[ \sum_{s=1}^T P_i^s Q_i^s \mid Q_i^1, \dots, Q_i^{t-1} \right] &= \sum_{s=1}^{t-1} P_i^s Q_i^s + E \left[ \sum_{s=t-1}^T P_i^s Q_i^s \mid Q_i^1, \dots, Q_i^{t-1} \right] \\
&= E \left[ \sum_{s=1}^T Q_i^s \mid Q_i^1, \dots, Q_i^{t-1} \right]
\end{aligned}$$

### Two-horse Framework

In what follows we describe the market in which there is a two-horse race surrounding. Consider a bookmaker who sets betting odds for two horses in the race. Unlike betting activity in other sports games where bookmakers are allowed to

change odds during the game based mainly on players' performance (football, tennis etc.), we pay particular attention to horse racing because of the short racing period during which it is not suitable to adjust the price. If horses' in-game performance has little effect on the price-adjusting process, bookmakers can only rely on the demanding for each horse to update the betting odds. At any point in time, the bookmaker faces a *two-sided* decision: either changing the price of horse A or the price of horse B.

Since the purpose of the economic activity is all about the profitability of the project, there is no doubt that the bookmaker always tries to guarantee his(her) profits by setting prices in which the summation of all the odds should be larger or equal to one no matter which horse wins, otherwise there exists an arbitrage opportunity. The over-roundness which is a means of expressing to what extent the odds are in favour of the bookmaker is the bookmaker's profit margins, also known as profits in an economic sense. For convenience, we simply assume there is no net profit for the bookmaker and if there are two horses in a race,  $P_A + P_B = 1$  where  $P_A$  denotes the price for horse A, and  $P_B$  the price for horse B. Therefore, adjusting one price would cause a change in the other automatically. It is worth noting that the cost per ticket is £1, which means the bookie pays out £1 for the winning horse per ticket sold. If all punters in the market are noise traders who are indifferent between two horses, the only criterion for them to pick a horse is to observe the given odds. In this case, the fair price is one half each and there is also no need to change the price because the uninformed putters go with each horse with probability  $1/2$ . However, this is not the truth - bookmakers change the odds all the time, so there must be some traders who are privileged. Those people are endowed with information of which horse is going to win, also known as *Informed traders*. Because of the existence of insiders, bookmakers face uncertainty when setting up the betting odds.

#### *Informed Traders and Uninformed 'Noise' Traders*

Informed traders have empirically been proved their existence in this market. Suppose that the fraction of informed traders is denoted by  $\mu$ , thus the proportion

of uninformed punters is  $\sigma \leq 1 - \mu$ . Insiders are assumed to know which horse is going to win, while the uninformed ones, whose valuations are uniformly distributed on the interval  $[0, 1]$ , have their own subjective winning probabilities for each horse respectively. It is no surprise that the bookmaker is definitely losing to informed bettors, but (s)he can try to break even on the uninformed.

### *Discrete-time versus Continuous-time*

The model we described above is obviously a finite discrete-time model with  $t = 1, 2, \dots, T$ , but the model we are going to present is an infinite continuous-time model. The following argument shows how we link them. The general model provides a framework that satisfies the equilibrium condition in the market, but in this paper we are more interested in the one period, for example, say from the opening odds to the first time the bookmaker changes the prices.

In what follows we show how a discrete time model can converge to a continuous time arithmetic Brownian motion. Suppose  $t$  is the first time that the bookmaker stops the learning process and makes the decision. Now consider a time interval  $[0, t]$ , which we partition into  $n$  parts of equal length  $dt = t/n$ . We will be interested in the limit as  $n \rightarrow \infty$ , which implies that  $dt \rightarrow 0$ . Now consider a binomial tree with cumulative sales going up or down. Over a small time interval  $dt$ , we expect the net cumulative sales to go up by  $u = \theta\mu dt + \sigma\sqrt{dt}$  and to go down by  $d = \theta\mu dt - \sigma\sqrt{dt}$ , where  $\mu$  is the fraction of informed traders and  $\sigma \leq 1 - \mu$  is the fraction of uninformed noise traders as defined in the previous subsection.  $\theta$  is the hypothesis that we are going to test. The net cumulative sales depend on what the informed traders know. If they know that horse A is going to win, all the net cumulative sales will go up. Otherwise they will go down. So  $\theta$  is going to be either 1 or  $-1$ . The  $\sqrt{dt}$  term here is to keep the variance of the cumulative relative sales finite when we take the limit  $dt \rightarrow 0$ . Each of these movements is assumed to occur with probability  $p = \frac{1}{2}$ . Let  $X_1, X_2, \dots, X_n$  be a sequence of independent

and identically distributed (*i.i.d*) random variables with

$$X_i = \begin{cases} \theta\mu dt + \sigma\sqrt{dt} & \text{w.p. } 1/2 \\ \theta\mu dt - \sigma\sqrt{dt} & \text{w.p. } 1/2 \end{cases}$$

Note that  $E[X_i] = \theta\mu dt$  and  $Var[X_i] = \sigma^2 dt$ .<sup>8</sup> Applying the Central Limit Theorem, we get<sup>9</sup>

$$\sqrt{n} \frac{\bar{X}_n - \theta\mu dt}{\sigma\sqrt{dt}} \xrightarrow{d} N(0,1)$$

where  $\bar{X}_n = \frac{1}{n} \sum_{i=1}^n X_i$ .

Thus, the distribution of  $\bar{X}_n$  becomes approximately normal if  $n$  gets large enough, which can be written as<sup>10</sup>

$$\bar{X}_n \overset{A}{\sim} N(\theta\mu dt, \sigma^2 dt/n)$$

Let the sequence  $S_1, S_2, \dots$  describe a stochastic process with  $S_n(0) = 0$ . The net cumulative sales after  $n$ -th trade are equal to

$$S_n(t) := \sum_{i=1}^n X_i = n\bar{X}_n \overset{A}{\sim} N(\theta\mu ndt, n^2\sigma^2 dt/n)$$

That is,

$$S(t) = \lim_{n \rightarrow \infty} S_n(t) \overset{A}{\sim} N(\theta\mu t, \sigma^2 t)$$

We construct Brownian motion as a limit of binomial trees in discrete time. As  $n \rightarrow \infty$ , for every  $t$ ,  $S(t)$  has a normal distribution with mean  $\theta\mu t$  and variance  $\sigma^2 t$ , which makes it an arithmetic Brownian motion ( $dS_t = \theta\mu dt + \sigma dB_t$ , where  $B_t$  is a standard Brownian motion.). Sample path of an arithmetic Brownian motion with

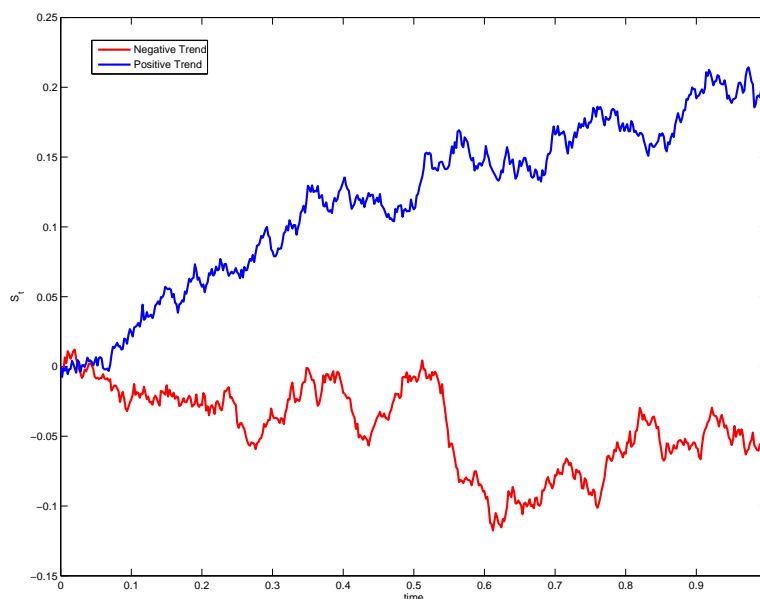
---

<sup>8</sup>  $E[X_i] = 1/2(\theta\mu dt + \sigma\sqrt{dt}) + 1/2(\theta\mu dt - \sigma\sqrt{dt}) = \theta\mu dt$ ;  
 $E[X_i^2] = 1/2(\theta\mu dt + \sigma\sqrt{dt})^2 + 1/2(\theta\mu dt - \sigma\sqrt{dt})^2 = \sigma^2 dt$ ;  
 $Var[X_i] = E[X_i^2] - (E[X_i])^2 = \sigma^2 dt$ .

<sup>9</sup>Here  $\xrightarrow{d}$  stands for *convergence in distribution*.

<sup>10</sup>Here  $\overset{A}{\sim}$  stands for "is asymptotically distributed as".

Figure 5.1 Sample Path of an Arithmetic Brownian Motion



a positive trend ( $\theta = 1$ ) is drawn in blue in Figure 5.1, while the red line depicts a negative trend ( $\theta = -1$ ). Following the sample path, the bookmaker can decide whether (s)he faces a positive trend or a negative one at the quickest possible point in time.

### *Jump Process*

For analytical convenience, we analyse a Poisson process with intensity  $\lambda$  that jumps from 0 to 1. There are two reasons for introducing the jump process. First, the infinite horizon model that we have already set up in the previous subsection is not realistic as the betting period lasts for about 30 minutes in the on-course horserace gambling markets. We know for sure that betting takes place over a fixed time interval, the race will start at a certain point of time and the market will be closed at the starting prices. However, the number of bets placed during that period is unknown *ex ante*. To model this, we use an infinite horizon model, which gets killed by the first jump of a Poisson process with intensity  $\lambda$ . If a finite time deadline is imposed as we did in our general equilibrium framework, the problem becomes much more difficult and can no longer be analytically solved. A jump

process is just a mechanism to introduce a finite time in the sense that we can still model our problem in infinite time horizon with all the analytical tools.

Second, note that “real time” and “model time” are not the same in our model. Time here is defined by numbers of trades, which are uniformly spread out, because we assume that the same amount of trade happens in every interval of time in the infinite time horizon model. But of course in the real world, punters come and go with no particular order. The trade does not happen at regular time intervals so that the bookmaker does not know exactly how many trades are going to be made and how fast the trade will come. So we redefine the time, not the real time, but the time between trades which are to some extent random. We assume those times are regularly interspaced. The intensity  $\lambda$ , mathematically, plays the same role as a discount factor. Uncertainty over the number of trades that a bookmaker will make during the fixed and known period of time that bets can be made, has the same effect as having a discount factor in an infinite horizon model. Therefore our model which has a random cut-off point is discounted at  $\lambda$ . In order to get analytical results, we assume that there is a random process that determines when the decision should be made. There is no discounting *per se*, but since trades do not appear regularly and at some stage there are no more trades to be done, that is when the bookmaker determines to end the learning period and changes the price.

The following shows how we mathematically transfer  $\lambda$  to a “discount rate”. The way we deal with this problem is by assuming there is a Poisson process  $(q_t)_{t \geq 0}$ , independent of  $(S_t)_{t \geq 0}$ , with intensity  $\lambda$  and  $q_0 = 0$ . Define

$$\bar{\tau} = \inf \{t \geq 0 \mid q_t = 1\}$$

Assume that the expected value of the project equals

$$F(s) = E_s \left[ \int_0^{\bar{\tau}} S_t dt \right]$$



Since  $\bar{\tau}$  follows an exponential distribution with parameter  $\lambda$ , we can write

$$\begin{aligned} F(s) &= E_s \left[ \int_0^\infty S_t P(\bar{\tau} \geq t) dt \right] \\ &= E_s \left[ \int_0^\infty e^{-\lambda t} S_t dt \right] \end{aligned}$$

Our model can be used to explain why initially the price does not change very often. One reason could be that the punters do not show up until just before the race. Another explanation is that in the first period the bookmaker is trying to collect the information about the informed traders. It can also be used to explain why the bookmaker always sets the odds which add up to something more than one. It is because in the learning period the bookmaker is going to make a loss, and in order to compensate the loss, the bookmaker should set the odds that give (s)he a profit margin to cover the learning loss. In theory, this profit margin would be completely competed away.

## 5.4 The Optimal Stopping Problem and Main Result

### *The Stochastic Process*

Recall that the bookmaker knows that there is a fraction  $\mu$  of insiders, but that (s)he does not know which horse the insiders are backing. If the initial price is such that noise traders are indifferent, then their trades should be white noise. Let  $S_t$  denote the cumulative sales of tickets on horse A, net of cumulative sales of tickets on horse B. If we assume that punters arrive sequentially in continuous time, then we observe

$$dS_t = \theta \mu dt + \sigma dB_t \tag{5.1}$$

where  $\mu > 0$  and  $\sigma > 0$  are constant, and  $(B_t)_{t \geq 0}$  is a standard Brownian motion.

The hypotheses that the bookmaker is testing are  $H_0 : \theta = -1$  and  $H_1 : \theta = 1$ . The intuition behind this is that if insiders know that horse A(B) will win we should have a positive(negative) trend. The decision depends on the state of the world,  $\theta$ : conditional on the event  $\{\theta = 1\}$  the bookmaker makes a decision to increase the price of horse A as all informed traders are backing horse A, thereby decreasing the price for horse B accordingly. However if the state of the world turns out to be  $\{\theta = -1\}$ , the price of horse B should be changed. The state is not known, but the bookmaker knows  $\pi$ , the prior probability that the state is  $\{\theta = 1\}$ .

### *Statistical Errors*

As mentioned before, the bookmaker changes the odds based on which state of the world is revealing. If the decision is correct, zero profits should be expected. There is also a chance that the decision is wrong. What the bookmaker really cares about is the loss of statistical errors.

Let *Type I error* denote the situation where all the informed traders know that horse B is going to win, but the bookmaker thinks they favour horse A conditional on the event  $\{\theta = 1\}$ . A *Type II error* represents the opposite situation. In this case the insiders are backing horse A but the bookmaker thinks horse B is their choice. If these two types of errors occur, for example, the true state is  $\{\theta = -1\}$ , but the bookmaker decides it is  $\{\theta = 1\}$ , *vice versa*, a huge loss is on the way.

### *The Loss Function*

In order to derive a proper loss function, it seems appropriate to exclude two extreme cases in the first step. The first case is that a rational bookmaker will never set a price equal to £1 (£0) for either horse A or horse B because this price immediately reveals the winning horse to the market. On the one hand, from punters' point of view, given that the odds are properly determined, if they purchase the winning horse, there is no profit to earn as the cost of the ticket now equals the return. If punters buy the other one, this is definitely going to be a loss. So no punter is willing to be involved in such a market except for the informed traders

who happen to have different information about the winning horse. From the bookmaker's perspective, on the other hand, (s)he takes a risk of setting the wrong prices, which implies tremendous losses to the insiders.

In reality, a rational bookmaker would never set a price equal to one to break down the market. (S)he is more willing to take losses to collect information about the true state in which they make zero profits theoretically, and normally the price will not be changed based on one person's particular behaviour. Recall that this is the bookmaker's *learning process*.

The second case is that if we assume all punters are informed, the bookmaker will learn in the first trade what the information the traders have and there will be the loss if the price of the winning horse is other than £1. This is not the case in our model either.

In what follows we derive the optimal price for horse A at time  $t$  if  $\{\theta = 1\}$  and the loss of the *Type I error*. In the current study we just consider if and only if the fraction  $\mu$  is informed, but the bookmaker does not know what they are - there are two possibilities, either a fraction of  $\mu$  of punters knows that horse A is the winning horse or they know that is horse B. Let  $P_A^* \in (0, 1)$  denote the adjusted price for horse A at time  $t$  conditional on  $\{\theta = 1\}$  and we assume that the valuation of the uninformed trader, denoted by  $V$ , is uniformly distributed on the interval  $[0, 1]$ , where 1 indicates that the uninformed punter is sure that horse A is going to win and 0 indicates the other direction, that is, they are sure that horse B is going to win. For informed traders, the loss the bookmaker makes is  $\mu(P_A^* - 1)$ , whilst for uninformed traders, the expected value is  $1/2$  and if the opening odds are such that the noise traders are indifferent, which is  $1/2$ , with probability  $1/2$  the bookmaker encounters a noise trader who buys horse A, and with probability  $1/2$  (s)he buys horse B. Once the bookmaker made a decision to set another price  $P_A^*$  based on sequential observations, the uninformed traders who are willing to buy the ticket of horse A at that price are the ones who have a valuation that is higher than  $P_A^*$ . Punters who have a valuation that is less than  $P_A^*$  are the ones who will buy horse B. So with probability  $1 - P_A^*$ , an uninformed punter buys a horse-A ticket and the

loss is  $P_A^* - 1$  provided that the bookmaker decides that horse A is going to win, and with probability  $P_A^*$ , the bookmaker gets  $1 - P_A^*$  from horse B and pays 0. In total, the expected profit of a trade is

$$E(\text{profit}) = \mu(P_A^* - 1) + (1 - \mu)[(1 - P_A^*)(P_A^* - 1) + P_A^*(1 - P_A^* - 0)] \quad (5.2)$$

Setting (5.2) equal to 0, we obtain

$$P_A^* = \frac{1}{2(1 - \mu)} \quad (5.3)$$

where  $0 \leq \mu < \frac{1}{2}$ . As can be seen from (5.3), if  $\mu$  goes to 1, which implies all punters in the market are informed, the price the bookmaker should set goes to infinity. In this case no one is willing to participate in the market. If  $\mu = \frac{1}{2}$ ,  $P_A^* = 1$  and the market also breaks down. The market is only functioning when the fraction of  $\mu$  of informed traders is less than half and the bookmaker gets a chance to break even on the uninformed ones.

To sum up, the expected profit equals 0 if the bookmaker sets the price  $P_A^*$  conditional on the right decision (s)he makes and the true state is indeed  $\{\theta = 1\}$  - all insiders are backing horse A. However if the *Type I error* occurs, with probability  $\mu$ , an informed punter appears and buys horse B instead of horse A at price  $1 - P_A^*$  and the loss to them is  $(1 - P_A^* - 1)$ , and with probability  $1 - \mu$ , an uninformed punter shows up and buys horse A with probability  $1 - P_A^*$  and horse B with probability  $P_A^*$ , the expected loss from the uninformed traders equals  $(1 - \mu)[(1 - P_A^*)(P_A^* - 0) + P_A^*(1 - P_A^* - 1)]$ . So the loss to the *Type I error* is as follows.

$$\begin{aligned} E_A(\text{loss}) &= \mu(1 - P_A^* - 1) + (1 - \mu)[(1 - P_A^*)(P_A^* - 0) + P_A^*(1 - P_A^* - 1)] \\ &= -\mu P_A^* + (1 - \mu)P_A^*(1 - 2P_A^*) \\ &= -\frac{\mu}{2(1 - \mu)} + \frac{1 - \mu}{2(1 - \mu)} \left[ 1 - \frac{2}{2(1 - \mu)} \right] \\ &= -\frac{\mu}{1 - \mu} \end{aligned} \quad (5.4)$$

In expectation if  $\mu = 0$ , the price  $P_A^* = \frac{1}{2}$  and the punters are uniformly distributed, the bookmaker makes zero loss. As  $\mu$  goes to  $\frac{1}{2}$ , the price goes to 1 and the loss per ticket goes to 1 as well. So the loss is bounded between 0 and 1, which is consistent with our assumption.

The loss of the *Type II error* is also  $\mu / (1 - \mu)^{11}$ . Table 5.1 summarises the loss of two types of errors under different states of the world.

Table 5.1 The Loss of Two Types of Errors

		State of the world	
		$\theta = 1$	$\theta = -1$
Decision	Adjust $H_A$	0	$\frac{\mu}{1-\mu}$
	Adjust $H_B$	$\frac{\mu}{1-\mu}$	0

If the posterior belief in the event  $\{\theta = 1\}$  at time  $t$  is  $p_t \in (0, 1)$ , the expected loss on the *Type I error*, denoted by  $L_1$ , is

$$\begin{aligned}
 L_1(p_t) &= p_t \times 0 + (1 - p_t) \left( \frac{\mu}{1 - \mu} \right) \\
 &= (1 - p_t) \left( \frac{\mu}{1 - \mu} \right)
 \end{aligned} \tag{5.5}$$

Similarly the the expected loss on the *Type II error*, denoted by  $L_0$ , is

$$\begin{aligned}
 L_0(p_t) &= p_t \left( \frac{\mu}{1 - \mu} \right) + (1 - p_t) \times 0 \\
 &= p_t \left( \frac{\mu}{1 - \mu} \right)
 \end{aligned} \tag{5.6}$$

Note that at time  $t$  the bookmaker will choose to change the price of horse A if, and only if,

$$p_t > \bar{p} \equiv \frac{1}{2} \tag{5.7}$$

The stochastic process also indicates what the opening price should be at time 0. The way we are modelling the net cumulative sales is a Brownian motion and

<sup>11</sup>See Appendix A for more details.

initially what the bookmaker is trying to make sure is that the expected profit the (s)he makes on the uninformed traders is zero. So if the opening odds for both horses are  $\frac{1}{2}$  and  $\frac{1}{2}$  respectively given that there is no informed trader in the market, the uninformed trader buys horse A with probability  $\frac{1}{2}$  and horse B with probability  $\frac{1}{2}$ . This gives us a binomial tree - the net cumulative sales could either go up or down. In this case the loss per trade is always going to be  $\frac{1}{2}$  no matter what information they possess or which horse they are going to bet because the bookmaker is always getting  $\frac{1}{2}$  by selling both types of the tickets. For informed traders, the bookmaker should have charged either 1 or 0 on a horse-A ticket, but if (s)he charged  $\frac{1}{2}$  for horse A instead, the loss is also always going to be  $\frac{1}{2}$  per trade on the informed traders. So the expected loss, as long as the bookmaker does not make the decision, is  $\frac{\mu}{2} + (1 - \mu) \left[ \frac{1}{2} \left( \frac{1}{2} - 1 \right) + \frac{1}{2} \left( \frac{1}{2} - 0 \right) \right] = \frac{\mu}{2}$ .

Therefore until a decision is reached the bookmaker is assumed to take a loss for this learning process, from time 0 to  $\tau$  the loss function is

$$\begin{aligned} L^*(p) &= \inf_{\tau \in \mathcal{T}} E_p \left[ \int_0^\tau e^{-\lambda t} \frac{\mu}{2} dt + e^{-\lambda \tau} \min \{L_1(p), L_0(p)\} \right] \\ &= E_p \left[ \int_0^{\tau^*} e^{-\lambda t} \frac{\mu}{2} dt + e^{-\lambda \tau^*} \min \{L_1(p_{\tau^*}), L_0(p_{\tau^*})\} \right] \end{aligned} \quad (5.8)$$

where  $\mathcal{T}$  is the set of all stopping times. This is our optimal stopping problem.

Since  $(p_t)_{t \geq 0}$  is Markovian, the optimal stopping problem (5.8) can be written as

$$\begin{aligned} L^*(p) &= \frac{1}{\lambda} \cdot \frac{\mu}{2} + \inf_{\tau} E_p \left[ e^{-\lambda \tau} \min \left\{ L_0(p_\tau) - \frac{\mu}{2\lambda}, L_1(p_\tau) - \frac{\mu}{2\lambda} \right\} \right] \\ &= \frac{1}{\lambda} \cdot \frac{\mu}{2} + \inf_{\tau} E_p \left[ e^{-\lambda \tau} \min \left\{ p_\tau \left( \frac{\mu}{1-\mu} \right) - \frac{\mu}{2\lambda}, (1-p_\tau) \left( \frac{\mu}{1-\mu} \right) - \frac{\mu}{2\lambda} \right\} \right] \\ &= \frac{1}{\lambda} \cdot \frac{\mu}{2} + \inf_{\tau} E_p \left[ e^{-\lambda \tau} \min \{G_0(p_\tau), G_1(p_\tau)\} \right] \end{aligned} \quad (5.9)$$

Note that since we assume  $\mu \in [0, \frac{1}{2})$ ,  $G_0(\cdot)$  and  $G_1(\cdot)$  are increasing and decreasing, respectively. That condition also ensures that  $G_0(0) < G_1(0)$  and  $G_0(1) > G_1(1)$ . Note that  $\bar{p} = 1/2$  is the unique point where  $G_0(\bar{p}) = G_1(\bar{p})$ .

This learning process reveals information about the true state of nature. We model the optimal stopping problem (5.8) as one of Bayesian sequential testing of two simple hypothesis in continuous time. In the Bayesian formulation of the problem it is assumed that we observe a trajectory of the Brownian motion  $S = (S_t)_{t \geq 0}$  with drift  $\theta\mu$  where the random variable  $\theta$  may be 1 or  $-1$  with probability  $p$  or  $1 - p$ , respectively.

Following Peskir and Shiyayev (2006), uncertainty is modelled on a *probability-statistical* space  $(\Omega; \mathcal{F}; P_\pi, \pi \in (0, 1))$  where for fixed  $\pi \in (0, 1)$  the probability measure  $P_\pi$  is obtained as follows.

$$P_\pi = \pi P_1 + (1 - \pi) P_{-1}$$

where  $P_1$  and  $P_0$  denote the distributions of the observed process under  $H_1$  and  $H_0$ , respectively, with  $P_1(\theta = 1) = P_{-1}(\theta = -1) = 1$ <sup>12</sup>. Recall that  $\theta$  is the hypothesis that we are testing, which takes two values 1 and  $-1$  with probabilities  $P_p(\theta = 1) = p$  and  $P_p(\theta = -1) = 1 - p$ .

Since our task is to test sequentially the hypotheses  $H_1$  and  $H_0$  with a minimal loss, we observe the continuous process  $(S_t)_{t \geq 0}$ . This process generates the natural filtration  $\mathcal{F}_t^S = \sigma(S_s : 0 \leq s \leq t)$ , which is augmented with the  $P_\pi$  null set.

---

<sup>12</sup>Suppose informed traders are not perfectly informed, then this  $P_1(\theta = 1) = P_{-1}(\theta = -1) = 1$  is no longer what we are sequentially testing. Suppose informed traders are correct with probability  $\eta > 1/2$ , then  $P(\text{punter buys horse A} \mid \text{informed}) = \eta P(\text{horse A wins}) + (1 - \eta) P(\text{horse B wins})$ . We know that

$$P(\text{horse A wins} \mid \text{punter buys horse A}) = \frac{P(\text{punter buys horse A} \mid \text{horse A wins}) P(\text{horse A wins})}{P(\text{punter buys horse A})}$$

By applying Bayes' rule, the probability that horse A wins given that punter buys horse A equals the probability that punter buys horse A conditional on horse A wins times the probability of horse A wins divided by the probability that punter buys horse A. In this chapter, we assume informed traders know which horse wins. So the probability that informed punter buys horse A conditional on the probability that horse A wins is 1. If the trader is not fully informed, the probability is not 1, but  $\eta$ . If horse A wins, the probability that the informed traders will buy horse A is  $\eta$  and they will buy horse B with probability  $1 - \eta$ .  $\frac{2\mu}{\sigma}$  that shows up in the volatility of the geometric Brownian motion is going to change. The distance between  $\mu$  and  $-\mu$  becomes narrower because  $\eta \leq 1$ . Most analysis can go through, but we are learning less from the observations.

Let  $p_t = P_\pi(\theta = 1 \mid \mathcal{F}_t^S)$  with  $t \geq 0$ , thus

$$\begin{aligned} p_t &= \frac{\pi \cdot \exp\left\{-\frac{(S_t - \mu t)^2}{2\sigma^2 t}\right\}}{\pi \cdot \exp\left\{-\frac{(S_t - \mu t)^2}{2\sigma^2 t}\right\} + (1 - \pi) \cdot \exp\left\{-\frac{(S_t + \mu t)^2}{2\sigma^2 t}\right\}} \\ &= \left[1 + \frac{1 - \pi}{\pi} \exp\left\{-\frac{2\mu}{\sigma^2} S_t\right\}\right]^{-1} \end{aligned} \quad (5.10)$$

and let the Radon-Nikodym derivative

$$\varphi_t = \frac{d(P_1 \mid \mathcal{F}_t^S)}{d(P_{-1} \mid \mathcal{F}_t^S)} \quad (5.11)$$

defines the *likelihood ratio process* between the two hypotheses  $(\varphi_t)_{t \geq 0}$ . It is a well-known fact that<sup>13</sup>

$$\varphi_t = \exp\left\{\frac{2\mu}{\sigma^2} S_t\right\} \quad (5.12)$$

and since  $(\pi < 1)$

$$p_t = \frac{\pi d(P_1 \mid \mathcal{F}_t^S)}{\pi d(P_1 \mid \mathcal{F}_t^S) + (1 - \pi) d(P_{-1} \mid \mathcal{F}_t^S)} \quad (5.13)$$

it follows that

$$p_t = \frac{\frac{\pi}{1 - \pi} \varphi_t}{1 + \frac{\pi}{1 - \pi} \varphi_t} \quad (5.14)$$

Taking Ito differentials of the right-hand side in (5.12), we find that

$$d\varphi_t = \frac{2\mu}{\sigma^2} \varphi_t (dS_t + \mu dt) \quad (5.15)$$

which follows the geometric Brownian motion  $\frac{d\varphi}{\varphi} = \left(\theta \frac{2\mu^2}{\sigma^2} + \mu\right) dt + \frac{2\mu}{\sigma} dB$ .

From Ito's lemma the right-hand side of (5.14) follows that

$$dp_t = \frac{2\mu^2}{\sigma^2} p_t(1 - p_t)(1 - 2p_t)dt + \frac{2\mu}{\sigma^2} p_t(1 - p_t)dS_t \quad (5.16)$$

---

<sup>13</sup>See Shiriyayev, A. N. (1978).



Note that in our study  $p_t = P_p(\theta = 1 | \mathcal{F}_t^S) = E_p[\theta | \mathcal{F}_t^S]$ , then we consider the process

$$\bar{B}_t = \sigma^{-1} \left( X_t + \mu t - 2\mu \int_0^t p_s ds \right) \quad (5.17)$$

Combined (5.16) and (5.17), we obtain that  $(p_t)_{t \geq 0}$  follows the stochastic differential equation

$$dp_t = \frac{2\mu}{\sigma} p_t(1 - p_t) d\bar{B}_t, \quad \text{with } p_0 = \pi \quad (5.18)$$

where  $(\bar{B}_t)_{t \geq 0}$  is also a standard Brownian motion, called the *innovation process* (See Poor and Hadjiliadis, 2009). Using (5.12) and (5.13) it can be verified that the process  $(p_t)_{t \geq 0}$  is time-homogeneous and strongly Markovian under  $P_p$  with respect to the natural filtration. Note that if  $\theta = 1$ ,  $p_t \xrightarrow{a.s.} 1$  and if  $\theta = -1$ ,  $p_t \xrightarrow{a.s.} 0$  as  $t \rightarrow \infty$ . So, as  $t \rightarrow \infty$ ,  $\text{Var}(dp_t) \rightarrow 0$  holds in either cases since  $\text{Var}(dp_t) = (2\mu/\sigma)^2 p_t^2(1 - p_t)^2 dt$ .

In what follows we adapt Thijssen and Bregantini's model (2016) as the starting point. With all the setup in mind we see that the closer  $(p_t)_{t \geq 0}$  gets to either 0 or 1 the less likely that the loss will decrease upon continuation. This suggests that there exist points  $p_B \in (0, \frac{1}{2})$  and  $p_A \in (\frac{1}{2}, 1)$ , thus the state space  $(0, 1)$  can be divided into three regions. The first one is called the *continuation region*, which is a region around  $\bar{p}$  where keeping the price unchanged is optimal. It is denoted by

$$\mathcal{C} = \{p \in (0, 1) \mid L^*(p) < \min(L_0(p), L_1(p))\} = (p_B, p_A)$$

where  $p_B$  and  $p_A$ , with  $0 < p_B < \bar{p} < p_A < 1$ , are the boundaries for changing the price of horse B and changing the price of horse A, respectively. When  $p$  gets lower enough, we enter the *adjusting horse B region*, where the cumulative sales of horse B are larger than of horse A so that we change the price of horse B. This region is denoted by

$$\mathcal{D}_B = \{p \in (0, 1) \mid L^*(p) = L_0(p)\} = (0, p_B]$$

Conversely, when  $p$  gets large enough we enter the *adjusting horse A region*, where we should increase the price of horse A, and this region is denoted by

$$\mathcal{D}_A = \{p \in (0,1) \mid L^*(p) = L_1(p)\} = [p_A, 1)$$

It follows that the stopping time

$$\tau^* = \inf\{t \geq 0 : p_t \notin (p_B, p_A)\} \quad (5.19)$$

is optimal in (5.8). The next step is to find a function  $L^* \in \mathcal{C}^2$  that solves the following *free-boundary problem*

$$\begin{cases} \mathcal{L}L^* - \lambda L^* = 0 & \text{for } p \in (p_B, p_A) \\ L^*(p_B; p_A) = G_0(p_B) \\ L^*(p_A; p_B) = G_1(p_A) \\ L^{*'}(p_B; p_A) = G_0'(p_B) & \text{(smooth fit)} \\ L^{*'}(p_A; p_B) = G_1'(p_A) & \text{(smooth fit)} \end{cases} \quad (5.20)$$

Here  $\mathcal{L}$  defines the characteristic operator (see Øksendal, 2005) of  $(p_t)_{t \geq 0}$ , i.e. for any  $\psi \in \mathcal{C}^2$ ,

$$\mathcal{L}\psi(p) = \frac{1}{2} \left( \frac{2\mu}{\sigma} \right)^2 p^2 (1-p)^2 \psi''(p) \quad (5.21)$$

In order to derive the loss function of (5.8), we introduce the parameter

$$\gamma := \frac{1}{2} \sqrt{1 + 2\lambda \left( \frac{\sigma}{\mu} \right)^2} > \frac{1}{2}$$

There are two *fundamental solutions* to the differential equation  $\mathcal{L}\psi - \lambda\psi = 0$ :

$$\hat{\psi}(p) = \sqrt{p(1-p)} \left( \frac{p}{1-p} \right)^\gamma \quad (5.22)$$

$$\check{\psi}(p) = \sqrt{p(1-p)} \left( \frac{1-p}{p} \right)^\gamma \quad (5.23)$$

Note that  $\hat{\psi}(\cdot)$  is increasing and  $\check{\psi}(\cdot)$  is decreasing and the general solution to  $\mathcal{L}\psi - \lambda\psi = 0$  is of the form

$$\psi(p) = \hat{A}\hat{\psi}(p) + \check{A}\check{\psi}(p) \quad (5.24)$$

where  $\hat{A}$  and  $\check{A}$  are arbitrary constants. Furthermore, it is easily obtained that

$$\hat{\psi}'(p) = \hat{\psi}(p) \frac{1/2 + \gamma - p}{p(1-p)} > 0, \check{\psi}'(p) = \check{\psi}(p) \frac{1/2 - \gamma - p}{p(1-p)} < 0$$

$$\hat{\psi}''(p) = \hat{\psi}(p) \frac{\gamma^2 - 1/4}{p^2(1-p)^2} > 0$$

and

$$\check{\psi}''(p) = \check{\psi}(p) \frac{\gamma^2 - 1/4}{p^2(1-p)^2} > 0$$

In Appendix C we show the proof of the following proposition.

**Proposition 5.4.1.** *Suppose that*

1.  $0 \leq \mu < \frac{1}{2}$
2.  $\frac{1-\mu}{2\lambda} > \frac{1}{2} + \frac{1}{4\gamma}$

*In the problem of testing two simple hypotheses  $H_0 : \theta = -1$  and  $H_1 : \theta = 1$  on the observations of the process given by  $dS_t = \theta\mu dt + \sigma dB_t$ , the optimal decision rules  $\delta_p^* = (\tau_p^*, d_p^*)$  exists and is*

$$\tau_p^* = \inf \{ t \geq 0 \mid p_t^p \notin (p_B^*, p_A^*) \},$$

$$d_p^* = \begin{cases} 1, & p_{\tau^*}^p \geq p_A^* \\ -1, & p_{\tau^*}^p \leq p_B^* \end{cases} \quad (5.25)$$

and the loss function  $L$  is explicitly given by

$$L^*(p) = \begin{cases} p \cdot \left(\frac{\mu}{1-\mu}\right) & \text{if } p \in (0, p_B^*] \\ \frac{\mu}{2\lambda} + \hat{v}_{p_B^*, p_A^*}(p)G_1(p_A^*) + \check{v}_{p_B^*, p_A^*}(p)G_0(p_B^*) & \text{if } p \in (p_B^*, p_A^*) \\ (1-p) \cdot \left(\frac{\mu}{1-\mu}\right) & \text{if } p \in [p_A^*, 1) \end{cases} \quad (5.26)$$

where

$$\hat{v}_{p_B^*, p_A^*}(p) := \sqrt{\frac{p(1-p)}{p_A^*(1-p_A^*)} \frac{\left(\frac{1-p_B^*}{p_B^*} \frac{p}{1-p}\right)^\gamma - \left(\frac{p_B^*}{1-p_B^*} \frac{1-p}{p}\right)^\gamma}{\left(\frac{1-p_B^*}{p_B^*} \frac{p_A^*}{1-p_A^*}\right)^\gamma - \left(\frac{p_B^*}{1-p_B^*} \frac{1-p_A^*}{p_A^*}\right)^\gamma}} \quad (5.27)$$

<sup>14</sup>and

$$\check{v}_{p_B^*, p_A^*}(p) := \sqrt{\frac{p(1-p)}{p_B^*(1-p_B^*)} \frac{\left(\frac{1-p}{p} \frac{p_A^*}{1-p_A^*}\right)^\gamma - \left(\frac{p}{1-p} \frac{1-p_A^*}{p_A^*}\right)^\gamma}{\left(\frac{1-p_B^*}{p_B^*} \frac{p_A^*}{1-p_A^*}\right)^\gamma - \left(\frac{p_B^*}{1-p_B^*} \frac{1-p_A^*}{p_A^*}\right)^\gamma}} \quad (5.28)$$

are the expected discount factors of first reaching  $p_B^*$  and  $p_A^*$ , respectively, given the current posterior probability  $p$ .

## 5.5 Comparative Statics

In this section we obtain analytical results on the posterior bounds  $p_A^*$  and  $p_B^*$ , as well as the loss function given different parameters. In order to assess the quantitative effects on the bounds, we consider the following base-case scenario. Table 5.2 provides the parameters for a base-case scenario as well as the exact values of the posterior bounds of  $p_A^*$  and  $p_B^*$ , respectively. For this particular case we obtain  $p_B^* = 0.2527$  and  $p_A^* = 0.7473$ .<sup>15</sup> For different values of the posterior belief in the event  $\{\theta = 1\}$  the value of the loss is given by Figure 5.2. Note that the curve line between the thresholds  $p_B^*$  and  $p_A^*$  shows the loss of waiting for beliefs.

<sup>14</sup>See Appendix D.

<sup>15</sup>All calculations are done by MatLab.

Table 5.2 Parameters for a Base-case Numerical Example

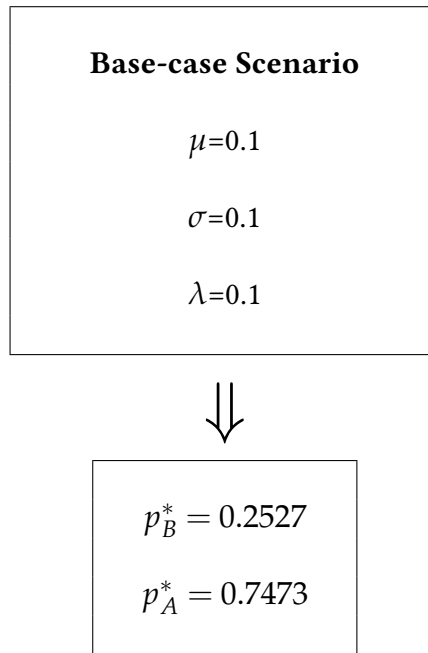
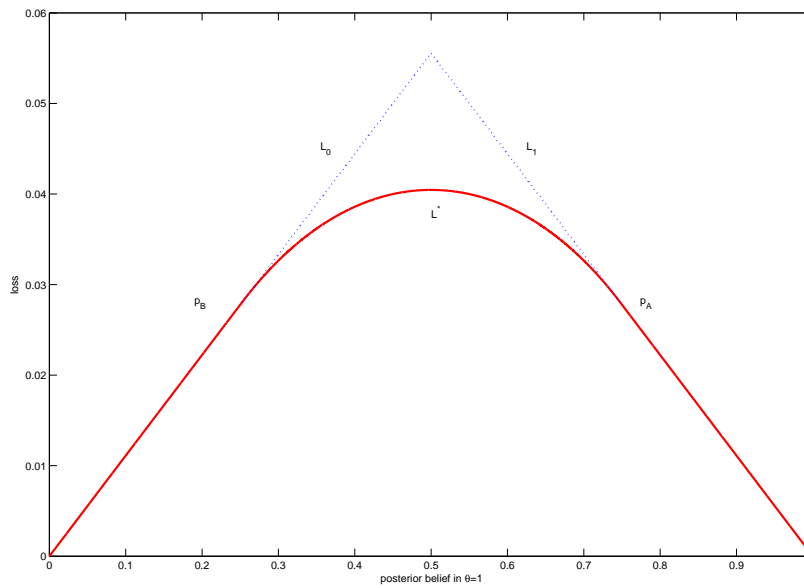


Figure 5.2 The Loss Function



The loss per trade given the different values of the posterior belief in the event  $\{\theta = 1\}$  is provided in Figure 5.2. The continuation region is between the threshold  $p_B$  and  $p_A$  as the loss is higher if the bookmaker decides to stop learning and takes action of adjusting the price. In order to guarantee the minimum loss for the

bookmaker, instead of making a decision at the moment it is optimal to continue learning. The maximum loss happens at the point  $\frac{1}{2}$  because under this circumstance the bookmaker is unable to distinguish which horse is treated as the winning horse from the informed traders, the opening odds keep at  $\frac{1}{2}$  so that the loss to the winning horse is one half no matter what information they possess. Note that this figure contains the loss from both sides - whether the bookmaker determines to change the price of horse A or horse B. It comes as no surprise that the figure is symmetric because we assume for the bookmaker to achieve zero profits, the sum of the prices of the two horses is equal to 1.

*$p_B$  is decreasing and  $p_A$  is increasing in the proportion of insiders (conditional on  $\theta = 1$ ),  $\mu$ , provided that  $p_B < 1/2 < p_A$ .* The comparative statics for the parameter  $\mu$  are described in Figure 5.3. Recall that the parameter  $\mu$  denotes the fraction of the informed traders in the learning process over a unit of time. It is easily seen that  $p_B$  is decreasing and  $p_A$  is increasing in  $\mu$ , the posterior bounds become wider because the higher value of  $\mu$  implies the learning process is more informative. The non-monotonicity in the loss function can be explained from two perspectives. On the one hand, a large proportion of insiders in the market would lead to a huge loss to the bookmaker. On the other hand, since the bookmaker can get more information per time period from the insiders, decisions will be taken sooner, which means the loss declines. The base-case scenario explicitly shows that the maximum loss occurs when the fraction  $\mu = 16\%$ .

*$p_B$  is increasing and  $p_A$  is decreasing in the volatility,  $\sigma$ , provided that  $p_B < 1/2 < p_A$ .* The comparative statics for the parameter  $\sigma$  are depicted in Figure 5.4. As  $p_B$  is increasing and  $p_A$  is decreasing, the posterior bounds narrow to the point 0.5 which implies a decision is reached earlier. This is partially different from the standard literature on real options where more uncertainty widens the decision bounds. In this paper, more uncertainty lowers the bounds, but increases the loss. Since this volatility part also represents the fraction of noise traders in the market, the larger the value of  $\sigma$ , the more the noise traders and the less the informed traders. As more uninformed traders make the learning process less informative

and the signal is more noising, the loss is incurred in order to learn, but the costs are the same, so we expect the decision is taken sooner, thereby lowering the loss accordingly.

$p_B$  is decreasing and  $p_A$  is increasing in the jump process,  $\lambda$ . The comparative statics for the parameter  $\lambda$  are described in Figure 5.5. As  $p_B$  is decreasing and  $p_A$  is increasing, the posterior bounds widen which implies the learning process becomes more informative. A higher  $\lambda$  reduces the expected cost of keeping the current odds and makes waiting longer more attractive, but the higher  $\lambda$  which acts as a discount factor, the less we care about the future, the faster a decision is reached. These are two opposing effects, one of which apparently dominates. In our setting, the presence of a running loss is important. So the loss is decreasing and the expected time to decision is increasing as in expectation it will take longer to reach if the bounds widen.

We can also derive the expected time to decision. The procedure is standard, which follows Poor and Hadjiliadis (2009). Combining THEOREMS 4.13, 4.15 and 4.17, we obtain<sup>16</sup>

$$\begin{aligned} E_p[\tau^*] &= (1-p)E_0[\tau^*] + pE_1[\tau^*] \\ &= \frac{2\sigma^2}{(2\mu)^2} \cdot \log \left[ \left( \frac{p}{1-p} \right)^{1-2p} \left( \frac{1-p_B}{p_B} \right)^{1-2p_B} \right] \\ &\quad + \frac{2\sigma^2}{(2\mu)^2} \cdot \frac{p-p_B}{p_A-p_B} \cdot \log \left[ \left( \frac{p_B}{1-p_B} \right)^{1-2p_B} \left( \frac{1-p_A}{p_A} \right)^{1-2p_A} \right] \end{aligned}$$

---

<sup>16</sup>See Appendix E for proof.

Figure 5.3 Comparative Statics for  $\mu$

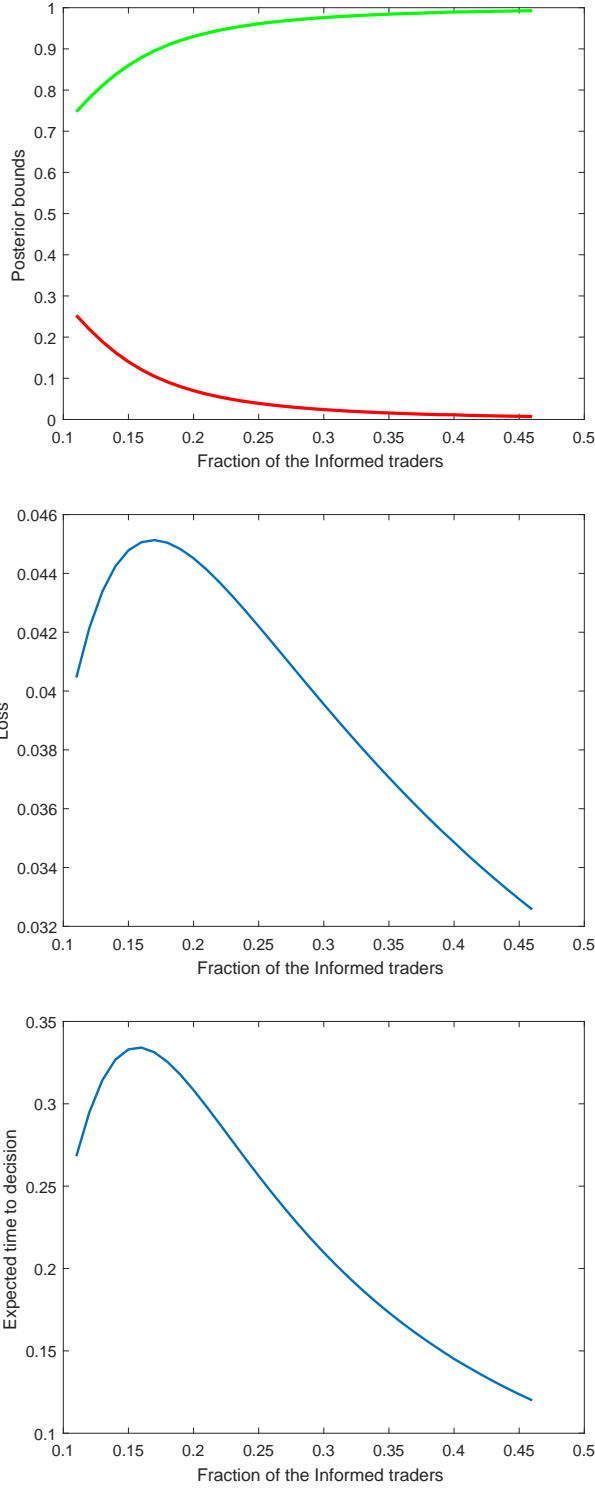




Figure 5.4 Comparative Statics for  $\sigma$

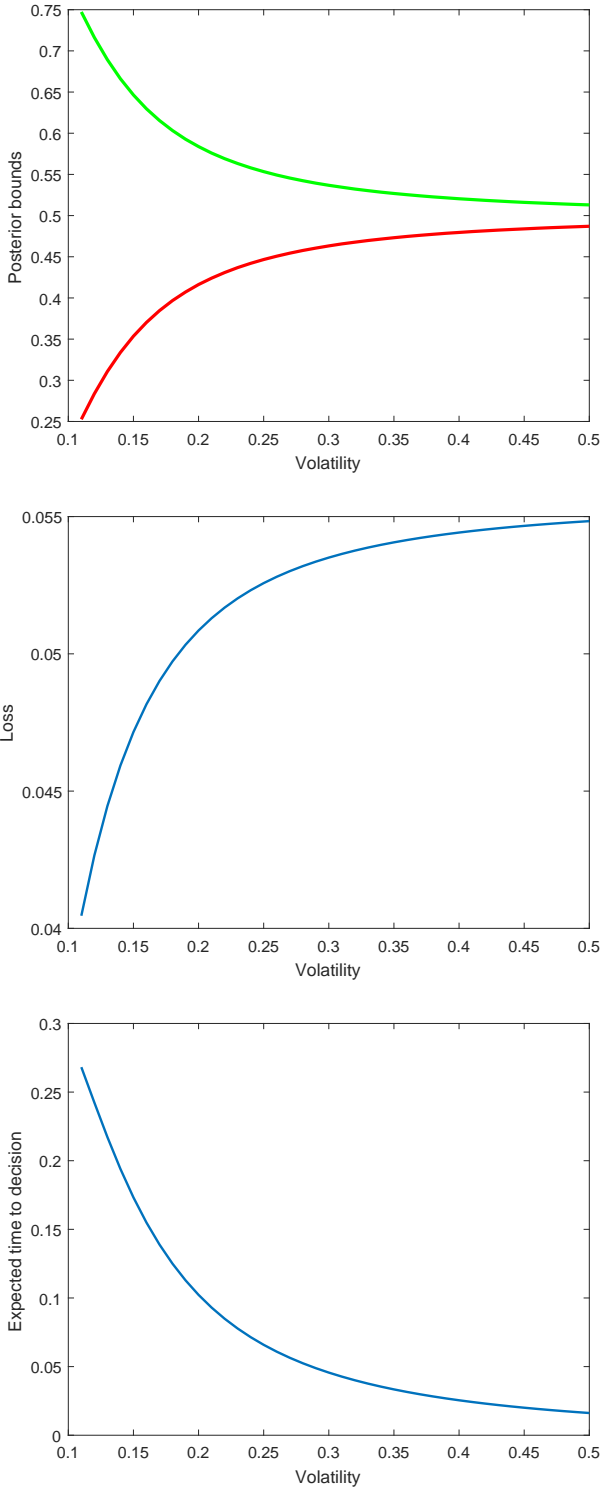
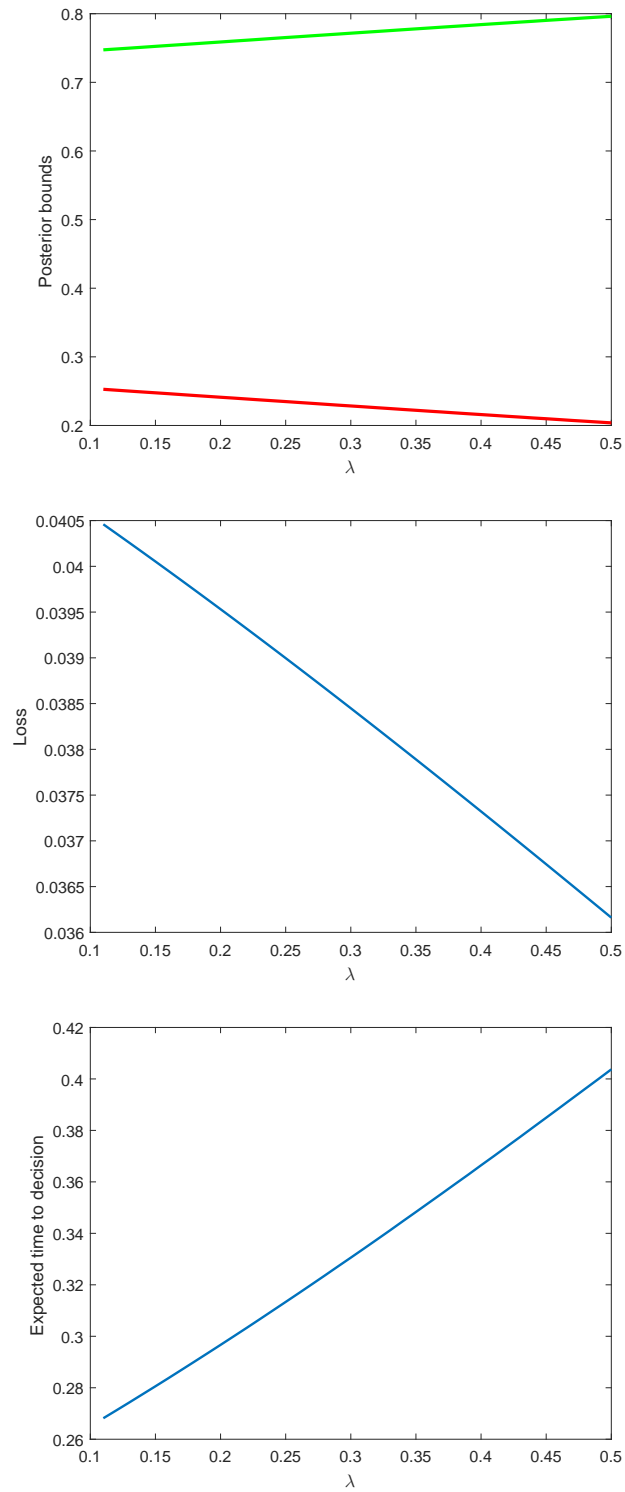


Figure 5.5 Comparative Statics for  $\lambda$



## 5.6 Conclusions

The Bayesian model presented in this paper focuses on the optimal stopping problem for the bookmaker who makes a decision on when to change the odds for each horse in a race, given that there are informed traders in the market and the information on which horse the informed traders prefer is unknown *ex ante*. The bookmaker goes through this *so-called* learning process to assess whether (s)he has gathered enough statistical evidence to determine the favourite horse among insiders at each point in time by applying sequential hypothesis test.

In this 2-horse benchmark model, we have provided analytical results that can give posterior bounds for adjusting the price of horse A or the price of horse B when the learning stops. Once a decision is taken, no matter what it is, it is guaranteed that the loss to the bookmaker is minimised.

We have also investigated the sensitivity of the model's parameters to the proposed solution to (5.8), which is consistent with our intuition. We find that (i) as the fraction of the informed traders becomes larger, the posterior bounds get wider and the bookmaker gets more information per time period along with the non-monotonicity loss. One explanation is that the more informed traders, there is more information the bookmaker can get per time period, which leads to a decision being taken sooner. So the decreased loss should be expected, which is consistent with our assumption as well as our intuition about this market. (ii) the decision bounds narrow in the volatility term, as in noise traders in this model. A higher volatility means less information, which leads to a decision is taken sooner; (iii) the posterior bounds widen in the jump process, the bigger the value of  $\lambda$  the wider the decision bounds, the lower the loss.

## 5.7 Appendices

### Appendix A. Loss of the *Type II error*

Similarly, let  $P_B^*$  denote the adjusted price for horse B at time  $t$  given the event  $\{\theta = -1\}$ . The valuation of the uninformed trader,  $V$ , is also uniformly distributed on the interval  $[0, 1]$ . For the informed trader, the bookmaker loses  $\mu(P_B^* - 1)$ . For the uninformed trader, he buys a horse-2 ticket if and only if his own valuation is higher than  $P_B^*$ . With probability  $1 - P_B^*$ , an uninformed trader buys horse B and the bookmaker loses  $P_B^* - 1$  on this kind. With probability  $P_B^*$ , a horse-A ticket is bought and the loss of this kind is 0. So the expected profit of a trade is

$$E(\text{profit}) = \mu(P_B^* - 1) + (1 - \mu)[(1 - P_B^*)(P_B^* - 1) + P_B^*(1 - P_B^* - 0)]$$

Setting the above equation equal to 0, we can get

$$P_B^* = \frac{1}{2(1 - \mu)}$$

So the loss to the *Type II error* follows the same analysis.

$$\begin{aligned} E_B(\text{loss}) &= \mu(1 - P_B^* - 1) + (1 - \mu)[(1 - P_B^*)(P_B^* - 0) + P_B^*(1 - P_B^* - 1)] \\ &= -\mu P_B^* + (1 - \mu)P_B^*(1 - 2P_B^*) \\ &= -\frac{\mu}{2(1 - \mu)} + \frac{1 - \mu}{2(1 - \mu)} \left[1 - \frac{2}{2(1 - \mu)}\right] \\ &= -\frac{\mu}{(1 - \mu)} \end{aligned}$$

where  $0 \leq \mu < \frac{1}{2}$ .

## Appendix B. Proof of Eqs. (5.15) and (5.16)

Applying Ito's lemma to equation (5.12) yields

$$\begin{aligned}
 d\varphi_t &= \frac{2\mu^2}{\sigma^2} \varphi_t dS_t + \frac{1}{2} \left( \frac{2\mu}{\sigma^2} \right)^2 \varphi_t (dS_t)^2 \\
 &= \frac{2\mu^2}{\sigma^2} \varphi_t dS_t + \frac{1}{2} \left( \frac{2\mu}{\sigma^2} \right)^2 \varphi_t \sigma^2 dt \\
 &= \frac{2\mu}{\sigma^2} \varphi_t (dS_t + \mu dt)
 \end{aligned}$$

Similarly applying Ito's lemma to equation (5.14) yields

$$\begin{aligned}
 dp_t &= -\frac{\pi(\pi-1)}{(\pi\varphi_t - \pi + 1)^2} d\varphi_t + \frac{1}{2} \frac{2\pi^2(\pi-1)}{(\pi\varphi_t - \pi + 1)^3} d\varphi_t^2 \\
 &= -\frac{\pi(\pi-1)}{(\pi\varphi_t - \pi + 1)^2} \frac{2\mu}{\sigma^2} \varphi_t (dS_t + \mu dt) \\
 &\quad + \frac{1}{2} \frac{2\pi^2(\pi-1)}{(\pi\varphi_t - \pi + 1)^3} \left( \frac{2\mu}{\sigma^2} \varphi_t (dS_t + \mu dt) \right)^2 \\
 &= \frac{2\mu}{\sigma^2} \left( \frac{\pi\varphi_t}{\pi\varphi_t - \pi + 1} \right) \left( \frac{1-\pi}{\pi\varphi_t - \pi + 1} \right) (dS_t + \mu dt) \\
 &\quad - \left( \frac{2\mu}{\sigma^2} \right)^2 \sigma^2 \left( \frac{\pi\varphi_t}{\pi\varphi_t - \pi + 1} \right)^2 \left( \frac{1-\pi}{\pi\varphi_t - \pi + 1} \right) dt \\
 &= \frac{2\mu}{\sigma^2} p_t (1-p_t) dS_t + \frac{2\mu^2}{\sigma^2} p_t (1-p_t) (1-2p_t) dt
 \end{aligned}$$

## Appendix C. Proof of Proposition 1

For  $p_L \leq \bar{p}$ , we define a mapping in the  $p \mapsto \check{L}(p; p_L)$  by

$$\check{L}(p; p_L) = \hat{A}(p_L)\hat{\psi}(p) + \check{A}(p_L)\check{\psi}(p) \quad (5.29)$$

where the constants  $\hat{A}(p_L)$  and  $\check{A}(p_L)$  are given by

$$\hat{A}(p_L) = \frac{\check{\psi}(p_L)}{2\gamma} \left[ G'_0(p_L) - \frac{1/2 - \gamma - p_L}{p_L(1 - p_L)} G_0(p_L) \right] \quad (5.30)$$

and

$$\check{A}(p_L) = \frac{\hat{\psi}(p_L)}{2\gamma} \left[ \frac{1/2 + \gamma - p_L}{p_L(1 - p_L)} G_0(p_L) - G'_0(p_L) \right]. \quad (5.31)$$

Notice that  $\mathcal{L}\check{L}(p; p_L) - \lambda\check{L}(p; p_L) = 0$  for all  $p \in (0, 1)$ . Furthermore, we show that  $\check{L}(p_L; p_L) = G_0(p_L)$  and  $\check{L}'(p_L; p_L) = G'_0(p_L) > 0$ . With condition 2, it is ensured that  $G_0(\bar{p}) < 0$  and since  $G'_0(p) = \frac{\mu}{1-\mu} > 0$ , it implies  $G_0(\cdot)$  is monotonically increasing. With  $p_L \in (0, \bar{p})$ , it follows  $G_0(p_L) < 0$  and

$$\frac{\partial \hat{A}(p_L)}{\partial p_L} = -\frac{\check{\psi}(p_L)}{2\gamma} \cdot \frac{\gamma^2 - \frac{1}{4}}{p_L^2(1 - p_L)^2} \cdot G_0(p_L) > 0$$

and

$$\frac{\partial \check{A}(p_L)}{\partial p_L} = \frac{\hat{\psi}(p_L)}{2\gamma} \cdot \frac{\gamma^2 - \frac{1}{4}}{p_L^2(1 - p_L)^2} \cdot G_0(p_L) < 0$$

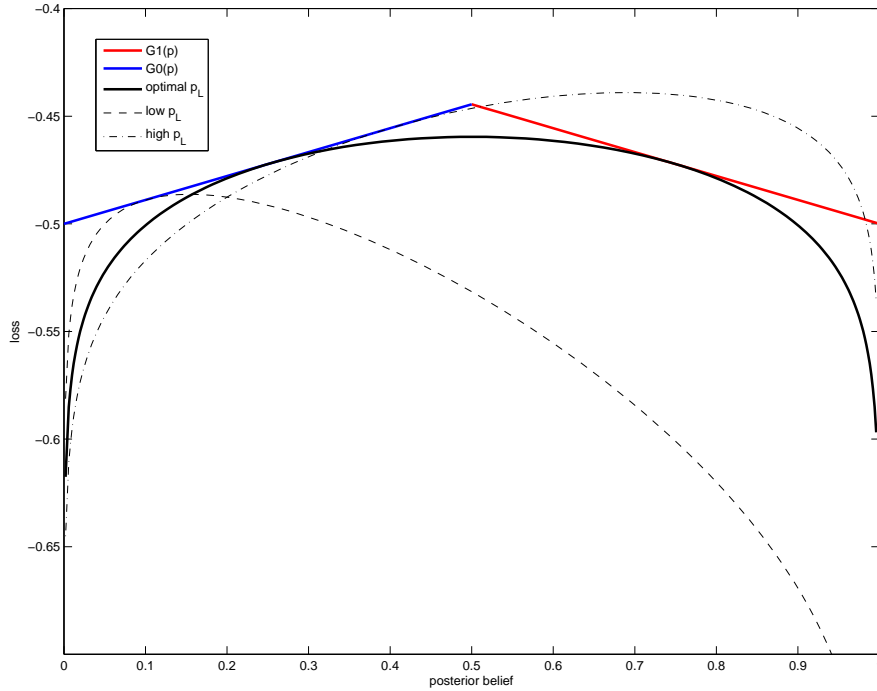
Therefore, it is easily verified that  $\check{A}(p_L) < 0$  for all  $p_L$ . Since  $\check{A}(p_L)$  is decreasing in  $p_L$ , and  $\check{A}(\cdot)$  is strictly monotone between 0 and  $\bar{p}$ , it is seen that  $\check{A}(\bar{p}) < 0$ . Condition 2 also ensures that  $\hat{A}(\bar{p}) < 0$ . Since  $\hat{A}(\cdot)$  is monotonically increasing on  $(0, \bar{p})$ ,  $\hat{A}(p_L) < 0$ . Thus, the function  $p \mapsto \check{L}(p; \bar{p})$  is concave on  $(0, 1)$ <sup>17</sup> and it satisfies  $\check{L}(0+; \bar{p}) = \check{L}(1-; \bar{p}) = -\infty$ .

<sup>17</sup> $\check{L}''(p; \bar{p}) = \hat{A}(\bar{p})\hat{\psi}''(p) + \check{A}(\bar{p})\check{\psi}''(p) < 0$ .

Since  $\hat{A}(p_L)$  increases and  $\check{A}(p_L)$  decreases in  $p_L$ , it holds that  $\frac{\partial \check{L}(p; p_L)}{\partial p_L} = \frac{\partial \hat{A}(p_L)}{\partial p_L} \hat{\psi}(p) + \frac{\partial \check{A}(p_L)}{\partial p_L} \check{\psi}(p) > 0$  for all  $p > p_L$ <sup>18</sup>.

So there is a unique point  $p_H \in (\bar{p}, 1)$  such that  $\check{L}'(p_H; p_L) = G_1'(p_H)$ , ensuring that  $\check{L}(p_H; p_L)$  increases in  $p_L$ , but decreases in  $p_H$ . Then we get the existence of a unique point  $p_B \in (0, \bar{p})$  for which there is  $p_A \in (\bar{p}, 1)$  such that  $\check{L}(p_A; p_B) = G_1(p_A)$  and  $\check{L}'(p_A; p_B) = G_1'(p_A)$ .

Figure 5.6 The Loss Function for Different Values of  $p_L$



For  $p_H \geq \bar{p}$ , we can also define the mapping in the  $p \rightarrow \hat{L}(p; p_H)$  by

$$\hat{L}(p; p_H) = \hat{B}(p_H) \hat{\psi}(p) + \check{B}(p_H) \check{\psi}(p) \quad (5.32)$$

<sup>18</sup>For every  $p > p_L$ ,  $\hat{\psi}(p) > \hat{\psi}(p_L)$  and  $\check{\psi}(p) < \check{\psi}(p_L)$ , then  $\frac{\partial \hat{A}(p_L)}{\partial p_L} \hat{\psi}(p) > \frac{\partial \hat{A}(p_L)}{\partial p_L} \hat{\psi}(p_L)$  and  $\frac{\partial \check{A}(p_L)}{\partial p_L} \check{\psi}(p) > \frac{\partial \check{A}(p_L)}{\partial p_L} \check{\psi}(p_L)$ , therefore  $\frac{\partial \check{L}(p; p_L)}{\partial p_L} = \frac{\partial \hat{A}(p_L)}{\partial p_L} \hat{\psi}(p) + \frac{\partial \check{A}(p_L)}{\partial p_L} \check{\psi}(p) > \frac{\partial \hat{A}(p_L)}{\partial p_L} \hat{\psi}(p_L) + \frac{\partial \check{A}(p_L)}{\partial p_L} \check{\psi}(p_L) = 0$ .

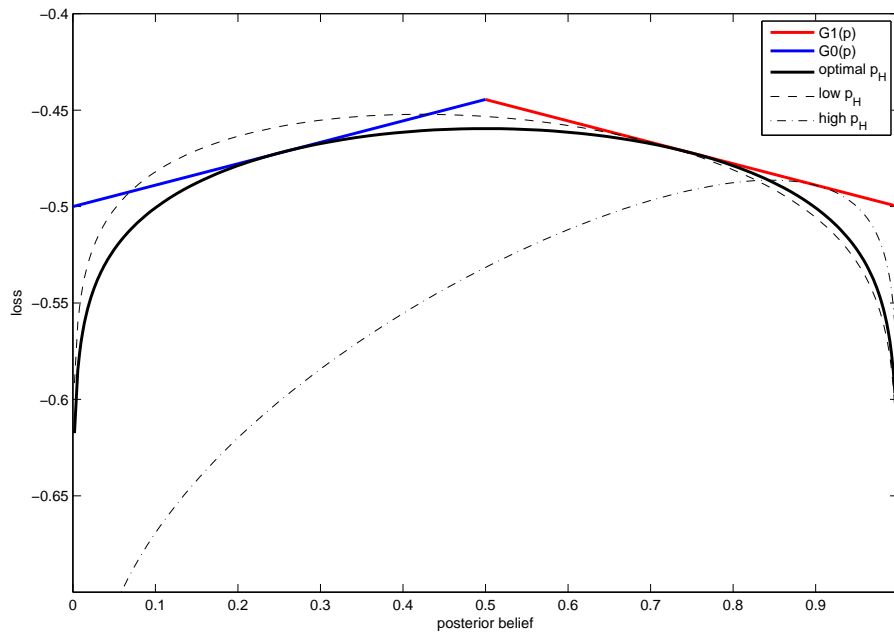
where the constants  $\hat{B}(p_H)$  and  $\check{B}(p_H)$  are given by

$$\hat{B}(p_H) = \frac{\check{\psi}(p_H)}{2\gamma} \left[ G_1'(p_H) - \frac{1/2 - \gamma - p_H}{p_H(1-p_H)} G_1(p_H) \right] \quad (5.33)$$

and

$$\check{B}(p_H) = \frac{\hat{\psi}(p_H)}{2\gamma} \left[ \frac{1/2 + \gamma - p_H}{p_H(1-p_H)} G_1(p_H) - G_1'(p_H) \right]. \quad (5.34)$$

Figure 5.7 The Loss Function for Different Values of  $p_H$



It is easily verified that  $\hat{L}(p_H; p_H) = G_1(p_H)$  and  $\hat{L}'(p_H; p_H) = G_1'(p_H)$ . Using a similar method it holds that  $\hat{B} < 0$ ,  $\check{B} < 0$ , and  $\hat{L}$  is concave on  $(0, 1)$ .

So, in order to find the unique solution  $(p_B, p_A)$  to the optimal stopping problem (5.9), the following conditions need to be satisfied

$$\begin{cases} \hat{L}(p_B; p_A) = G_0(p_B) \\ \check{L}(p_A; p_B) = G_1(p_A) \\ \hat{L}'(p_B; p_A) = G_0'(p_B) \\ \check{L}'(p_A; p_B) = G_1'(p_A) \end{cases}$$



## Appendix D. Expected Discount Factors

Recall from equation (5.24) that the general solution to the second order stochastic differential equation  $\mathcal{L}\psi - \lambda\psi = 0$  is

$$\psi(p) = \hat{A}\hat{\psi}(p) + \check{A}\check{\psi}(p) \quad (5.35)$$

We need that  $\psi(p_B) = 0$  and  $\psi(p_A) = 0$ , so

$$\hat{A}\hat{\psi}(p_B) + \check{A}\check{\psi}(p_B) = 0 \quad (5.36)$$

which means

$$\hat{A} = -\frac{\check{\psi}(p_B)}{\hat{\psi}(p_B)}\check{A} \quad (5.37)$$

Substituting (5.37) into (5.36), we obtain

$$\psi(p) = -\frac{\check{\psi}(p_B)}{\hat{\psi}(p_B)}\check{A}\hat{\psi}(p) + \check{A}\check{\psi}(p) \quad (5.38)$$

Therefore the discount factor is

$$\begin{aligned} \hat{v}_{p_B, p_A}(p) &= \frac{\psi(p)}{\check{\psi}(p_A)} \\ &= \frac{-\frac{\check{\psi}(p_B)}{\hat{\psi}(p_B)}\check{A}\hat{\psi}(p) + \check{A}\check{\psi}(p)}{-\frac{\check{\psi}(p_B)}{\hat{\psi}(p_B)}\check{A}\hat{\psi}(p_A) + \check{A}\check{\psi}(p_A)} \\ &= \frac{-\check{\psi}(p_B)\hat{\psi}(p) + \hat{\psi}(p_B)\check{\psi}(p)}{-\check{\psi}(p_B)\hat{\psi}(p_A) + \hat{\psi}(p_B)\check{\psi}(p_A)} \\ &= \sqrt{\frac{p(1-p)}{p_A(1-p_A)} \frac{(\frac{1-p_B}{p_B} \frac{p}{1-p})^\gamma - (\frac{p_B}{1-p_B} \frac{1-p}{p})^\gamma}{(\frac{1-p_B}{p_B} \frac{p_A}{1-p_A})^\gamma - (\frac{p_B}{1-p_B} \frac{1-p_A}{p_A})^\gamma}} \end{aligned} \quad (5.39)$$

## Appendix E. Expected Time to Decision

**Theorem 5.7.1.** (See Poor and Hadjiliadis (2009)) Consider the hypothesis pair

$$H_0 : Z_t = \sigma W_t + \mu_0 t, t \geq 0$$

versus

$$H_1 : Z_t = \sigma W_t + \mu_1 t, t \geq 0 \quad (5.40)$$

and suppose  $(T, \delta)$  is a sequential decision rule with error probabilities

$$P_0(\delta_T = 1) = \alpha \quad \text{and} \quad P_1(\delta_T = 0) = \gamma, \quad (5.41)$$

where  $\alpha, \gamma \in (0, 1)$ . Then

$$E_0\{T\} \geq -\frac{2\sigma^2}{(\mu_1 - \mu_0)^2} \left[ \alpha \log \left( \frac{1 - \gamma}{\alpha} \right) + (1 - \alpha) \log \left( \frac{\gamma}{1 - \alpha} \right) \right] \quad (5.42)$$

and

$$E_1\{T\} \geq -\frac{2\sigma^2}{(\mu_1 - \mu_0)^2} \left[ (1 - \gamma) \log \left( \frac{1 - \gamma}{\alpha} \right) + \gamma \log \left( \frac{\gamma}{1 - \alpha} \right) \right] \quad (5.43)$$

with equality if  $(T, \delta)$  is the sequential probability ratio test (SPRT).

**Theorem 5.7.2.** (See Poor and Hadjiliadis (2009)) Consider the hypotheses (5.40), and denote by  $(T, \delta)$  the SPRT( $A, B$ ) with  $0 < A \leq 1 \leq B < \infty$ , and  $A < B$ . Then

$$P_0(\delta_T = 1) = \frac{1 - A}{B - A} \quad (5.44)$$

and

$$P_1(\delta_T = 0) = A \frac{B - 1}{B - A} \quad (5.45)$$

where

$$A = \frac{1 - \pi}{\pi} \frac{\pi_L}{1 - \pi_L} \quad (5.46)$$

and

$$B = \frac{1 - \pi}{\pi} \frac{\pi_H}{1 - \pi_H} \quad (5.47)$$

Combining 5.7.1 and 5.7.2, we obtain

$$\begin{aligned} \alpha &= P_0(\delta_T = 1) = \frac{1 - A}{B - A} \\ &= \frac{(\pi - \pi_L)(1 - \pi_H)}{(1 - \pi)(\pi_H - \pi_L)} \end{aligned} \quad (5.48)$$

and

$$\begin{aligned} \gamma &= P_1(\delta_T = 0) = A \frac{B - 1}{B - A} \\ &= \frac{\pi_L(\pi_H - \pi)}{\pi(\pi_H - \pi_L)} \end{aligned} \quad (5.49)$$

So

$$\begin{aligned} E_\pi(T) &= \pi \cdot E_1(T) + (1 - \pi) \cdot E_0(T) \\ &= \frac{2\sigma^2}{(\mu_1 - \mu_0)^2} \cdot \log \left[ \left( \frac{\pi}{1 - \pi} \right)^{1-2\pi} \left( \frac{1 - \pi_L}{\pi_L} \right)^{1-2\pi_L} \right] \\ &\quad + \frac{2\sigma^2}{(\mu_1 - \mu_0)^2} \cdot \frac{\pi - \pi_L}{\pi_H - \pi_L} \cdot \log \left[ \left( \frac{\pi_L}{1 - \pi_L} \right)^{1-2\pi_L} \left( \frac{1 - \pi_H}{\pi_H} \right)^{1-2\pi_H} \right] \end{aligned}$$

# Chapter 6

## Concluding Remarks

### 6.1 Conclusions

In order to summarise this thesis, we present the main research questions of each chapter in turn and show what the limitations are and how these issues can be extended for future research.

This thesis is a collection of empirical and theoretical studies on the horserace betting markets, which contributes to the literature by answering the following questions: (i) Do insider trading problems exist in the horseracing betting markets? If so, can we find out the degree of the insider trading? (ii) Are the racing betting markets efficient? (iii) What kind of determinants will affect the starting prices? Are these factors persistent? (iv) How do bookmakers adjust the betting odds for each horse in a race?

To provide answers to the first two questions that we aim to investigate, we present the empirical analysis in Chapters 2 and 3 by applying large novel data sets, which collected respectively from the Yorkshire on-course racing markets in 2013 and 2014, to the Shin measure (1993) as well as the extension of the Shin (Jullien and Salanié (1994)). As is highlighted by the theoretical results of Shin (1993), the way to infer the incidence of insider trading in the British betting markets is similar to the bid-ask spread (hinted at by Bagehot (1971), developed by Copeland and Galai

(1983), and further analysed by Glosten and Milgrom (1985) and Kyle (1985)) quoted by market makers in financial markets. His empirical investigation concludes that insider trading problem do exist in betting markets and the degree of insider trading is around 2% when the weekly data are utilised. He also shows a strong positive relationship between the sum of prices and the number of competitors in the race.

Our empirical results in Chapter 2 demonstrate that the incidence of insider trading is around 1.7%, based on the nine Yorkshire racecourses data in the 2013 race season. The correlation between the total prices and the number of runners is also positive, which can be verified in the study. Besides that, we attempt to compute the degree of insider trading at opening and starting prices respectively for each racetrack based on the extension of the Shin model that allows doing the calculation at the individual level. We find that the weighted average degree at the starting prices is 2.103%, which is slightly higher than the value that we compute under the original Shin. We also report that the weighted average degree of insider trading at the opening prices is 1.64%. We find no evidence to confirm that the market is strongly efficient.

Chapter 3 comes as an extension of Chapter 2, which is used to consolidate our results by employing a different data set. The data set is also collected nine racecourses in Yorkshire, but the racing season that we focus on in this chapter was in 2014. The estimation results demonstrate that the weighted average degree of insider trading at the starting prices is around 2.068%, while at the opening prices 1.65%. The results in both chapters are pretty much the same. We also extend our analysis to other issues in this chapter. In particular, we provide empirical evidence that a gambler is doomed to lose even if he constantly herds with the favourite horse, and in our sample set, the winning starting prices equal the favourite starting prices in 390 races, which implies in these races the winning horse is actually the favourite horse. For the rest 780 races, the winning horse is not the favourite horse. On average the loss to the favourite horse is approximately 0.3314. Our finding shows that there is supportive evidence that informed punters tend to bet early with on-course bookmakers, which is consistent with the Crafts' empirical findings

(1985). Regarding the market efficiency test, we demonstrate that Yorkshire on-course betting markets are not strongly efficient. We also show that the average rate of return at both starting and opening prices is significantly negatively affected by the incidence of insider trading.

We aim to answer question (ii) in Chapter 4. In our setup, a unique cross-sectional and time series data set is employed to find out the determinants that might affect starting prices in horserace betting markets. In order to fully appreciate our empirical results, we perform three tests, that is, Hausman test, Pesaran's cross-sectional dependence test and time-fixed effects tests. Winning potential is a very important variable in this chapter, which is defined as one minus the horse's position over total runners in the race. Our findings show that the winning potential is significant in explaining the price of the horse. Besides this factor, the age of the horse, the weight the horse carries and the distance of the racecourse are quite important as well. Other factors like the condition of the turf, the size of the racetrack and the classification of the race also have influences on the starting prices. All the estimation results confirm that these determinants are persistent.

The last question is presented in Chapter 5. To do so, we firstly develop a general equilibrium framework, then in order to make the general model tractable and obtain rich analytical results, we analyse a two-horse benchmark model by introducing optimal stopping theory into the betting literature. We are interested in how bookmakers adjust the betting odds. By taking informed and uninformed punters into consideration, a risk neutral bookmaker who decides an optimal time to change the price for each horse in a race given that the cumulative sales of tickets on horse A, net of the cumulative sales of tickets on horse B are modelled by an arithmetic Brownian motion. Analytical solutions provide a complete understanding of the model in this study. Our findings highlight the following properties. As the fraction of informed traders gets bigger, the posterior bounds become wider and the bookmaker could receive more information per time period. This gives us the non-monotonicity loss happens because of two opposing effects. On the one hand, a large proportion of insiders causes a huge loss to the bookmaker, which

explains why the loss is increasing at the initial stage. On the other hand, more informed traders imply that there is more information the bookmaker can get per time period, which leads to a decision being taken sooner. This explains why the loss is declining after reaching a maximum point  $\mu = 16\%$ . As  $\sigma$  goes up, there are more noise traders in the market. So the signal that the bookmaker observed is more noising. The loss is incurred to learn, but the process is not informative and the costs are the same. We, therefore, expect a decision is reached earlier. The decision bounds get widening as the intensity  $\lambda$  becomes larger. The loss is declining because of the same reason that we have identified for  $\mu$ .

## 6.2 Future Research

Besides the questions we first set out to investigate, this thesis has raised additional problems for future research. In Chapters 2 and 3 we applied the new data sets to the Shin model, we contributed to the literature by confirming the results presented by Shin and other researchers who also followed the Shin measure. Another contribution is how we deal with data in MatLab as the data we utilised is publicly available and short lived. Although we agree that Shin did capture the insider trading problem in British horseracing betting markets, his model is restricted by lots of assumptions. It would be worthwhile to re-modify and see whether a more general model can be established. Furthermore, due to the constraints on the on-course betting markets that transactions are made in cash, we are looking forward to some improvements if online betting markets are included in the future study. Is it difficult or easy to detect the degree of insider trading when punters prefer to bet online with a large amount of money? What do betting companies react to such behaviour? Does prize money have anything to do with the insider trading?

In Chapter 4 we found out the factors that influence the starting prices. As we mentioned earlier, jockey's characteristics had not been included in the empirical analysis due to the lack of data. Other factors like weather or a bombshell could

also affect the odds. The empirical method applies ordinary least squares equations to panel data. Do other models suit our research purpose? These will improve the accuracy of our empirical model which provide an agenda for future research.

All of these questions lead to our final problem. That is, how do bookmakers adjust the odds of the horses during the betting period? Analytical results only provide a basic understanding of a two-horse benchmark model in this thesis. Can we extend this model? Can we obtain rich results from our general equilibrium framework? Both questions are equivalently important and are challenging us in the future.



# Betting Glossary

A glossary of used betting terms is as follows<sup>1</sup>.

## **Back**

To bet or wager a horse is to back it.

## **Betting ring**

The main area at a racecourse where the bookmakers operate.

## **Bookmaker**

The person who takes your bet and pays you if you win.

## **Decimal odds**

See odds.

## **Dutch book**

To bet on a number of horses, at varying odds, such that whichever bet wins, a set profit is guaranteed.

## **Evens**

The fractional odds 1/1.

## **Favourite**

The most popular horse in a race. It will have the shortest odds. There may be more than one horses in a race. See Joint-favourites.

---

<sup>1</sup>Accessible via: <http://www.racinguk.com/about-us/horseracing-betting-terms>;  
<http://www.tophorseracinglinks.com/html/glossary.htm>.  
<http://onlinebookmaker.com/betting-terms-glossary/list/2/d>.

**Fixed-odds betting**

The dividend is fixed at the odds when you placed you bet.

**Fractional odds**

See odds.

**Joint-favourites**

When a bookmaker cannot separate two/three horses for favouritism, they are made joint favourites.

**Long odds**

For example 10/1 is longer than 2/1 - long odds are applied to competitors that the bookmaker thinks are less likely to win.

**Long-shot**

A long-shot horse in a race is the one that has long odds and is therefore deemed to have little chance to win the race.

**Mutuel pool**

Short for 'parimutuel pool'. Sum of the bets on a race or event, such as the win pool, daily double pool, exacta pool, etc.

**Odds**

The bookmakers' view of the chance of a competitor winning (adjusted to include a profit). The figure or fraction by which a bookmaker or totalisator offers to multiply a bettor's stake, which the bettor is entitled to receive (plus his or her own stake) if the horse that (s)he selects wins. Odds can be expressed as a fraction, for example, 5/1 or 5/2. Odds can also be expressed as a decimal - so 5/1 would be 5.0, and 5/2 would be 2.5.

**Oddsbroker**

Similar to bookmaker. A person who sets the betting odds.

**Off-course bookmaker**

A bookmaker who is not present at the racecourse or other event, for example, in a high street betting shop or online.

**On-course bookmaker**

On the racecourse or at the event.

**Opening odds**

The first available odds at the racecourses.

**Over-broke**

A bookmaker makes a mistake by calculating the odds at an event that add up to less than 100%.

**Over-round**

The opposite of the above. The bookmaker's profit, which is determined by how much over 100% the total odds add up to.

**Payout**

What you get back from the bookmaker - your winnings and returned stake.

**Punt**

Another term for bet or wager.

**Punter**

A person who places a bet.

**Short odds**

For example 2/1 is shorter than 4/1 - short odds are applied to competitors that the bookmaker thinks are more likely to win.

**Starting Price or SP**

The price on a horse at the start of the race when the book closes.

**Stake**

The amount of money you bet.

**Tote**

The organisation appointed to receive bets and supply dividends in proportion to the amount of money wagered.

**Wager**

Bet, lay or gamble.

# Bibliography

- Ali, M. M. (1977). Probability and utility estimates for racetrack bettors. *Journal of Political Economy*, 85(4):803–815.
- Asch, P., Malkiel, B. G., and Quandt, R. E. (1984). Market efficiency in racetrack betting. *Journal of Business*, 57(2):165–175.
- Asch, P., Malkiel, B. G., and Quandt, R. E. (1986). Market efficiency in racetrack betting: Further evidence and a correction. *Journal of Business*, 59(1):157–160.
- Back, K. and Baruch, S. (2004). Information in securities markets: Kyle meets glosten and milgrom. *Econometrica*, 72(2):433–465.
- Bagehot, W. (1971). The only game in town. *Financial Analysts Journal*, 27(2):12–14, 22.
- Baltagi, B. (2008). *Econometric Analysis of Panel Data*. John Wiley & Sons.
- Bolton, R. N. and Chapman, R. G. (1986). Searching for positive returns at the track: A multinomial logit model for handicapping horse races. *Management Science*, 32(8):1040–1060.
- Bregantini, D. (2014). *Applications of continuous time stochastic processes in sequential clinical research design and econometrics*. PhD thesis, University of York, Department of Economics and Related Studies.
- Brown, L. D., D’Amato, R., and Gertner, R. (1994). Racetrack betting: Do bettors understand the odds? *CHANCE*, 7(3):17–23.

- Cain, M., Law, D., and Peel, D. A. (2001). The relationship between two indicators of insider trading in British racetrack betting. *Economica*, 68(269):97–104.
- Capinski, M. and Kopp, E. (2004). *Measure, Integral and Probability*. Springer, London.
- Copeland, T. E. and Galai, D. (1983). Information effects on the bid-ask spread. *Journal of Finance*, 38(5):1457–1469.
- Crafts, N. F. (1985). Some evidence of insider knowledge in horse race betting in Britain. *Economica*, 52(207):295–304.
- Décamps, J.-P., Mariotti, T., and Villeneuve, S. (2005). Investment timing under incomplete information. *Mathematics of Operations Research*, 30(2):472–500.
- Dixit, A. K. and Pindyck, R. S. (1994). *Investment under Uncertainty*. Princeton university press.
- Dowie, J. (1976). On the efficiency and equity of betting markets. *Economica*, 43(170):139–150.
- Driscoll, J. C. and Kraay, A. C. (1998). Consistent covariance matrix estimation with spatially dependent panel data. *Review of Economics and Statistics*, 80(4):549–560.
- Estep, D. (2002). *Practical Analysis in One Variable*. Springer New York.
- Fama, E. F. (1970). Efficient capital markets: A review of theory and empirical work. *Journal of Finance*, 25(2):383–417.
- Figlewski, S. (1979). Subjective information and market efficiency in a betting market. *Journal of Political Economy*, 87(1):75–88.
- Fingleton, J. and Waldron, P. (1999). Optimal determination of bookmakers' betting odds: Theory and tests. Trinity Economic Paper Series, Technical Paper No. 96/9, Trinity College, Dublin.
- Frees, E. W. (2004). *Longitudinal and Panel Data: Analysis and Applications in the Social Sciences*. Cambridge University Press.

- Gabriel, P. and Marsden, J. R. (1991). An examination of efficiency in British racetrack betting: Errata and corrections. *Journal of Political Economy*, 99(3):657–659.
- Gabriel, P. E. and Marsden, J. R. (1990). An examination of market efficiency in British racetrack betting. *Journal of Political Economy*, 98(4):874–885.
- Glosten, L. R. and Milgrom, P. R. (1985). Bid, ask and transaction prices in a specialist market with heterogeneously informed traders. *Journal of Financial Economics*, 14(1):71–100.
- Gramm, M., McKinney, C. N., Owens, D. H., and Ryan, M. E. (2007). What do bettors want? Determinants of pari-mutuel betting preference. *American Journal of Economics and Sociology*, 66(3):465–491.
- Griffith, R. M. (1949). Odds adjustments by American horse-race bettors. *The American Journal of Psychology*, 62(2):290–294.
- Hausch, D. B., Ziemba, W. T., and Rubinstein, M. (1981). Efficiency of the market for racetrack betting. *Management science*, 27(12):1435–1452.
- Henery, R. J. (1985). On the average probability of losing bets on horses with given starting price odds. *Journal of the Royal Statistical Society*, 148(4):342–349.
- Hill, T. P. (2009). Knowing when to stop: How to gamble if you must - the mathematics of optimal stopping. *American Scientist*, 97(2):126–133.
- Hoechle, D. (2007). Robust standard errors for panel regressions with cross-sectional dependence. *The Stata Journal*, 7(3):281–312.
- Jullien, B. and Salanié, B. (1994). Measuring the incidence of insider trading: A comment on shin. *The Economic Journal*, 104(427):1418–1419.
- Kano, S. and Ohta, M. (2005). Estimating a matching function and regional matching efficiencies: Japanese panel data for 1973-1999. *Japan and the World Economy*, 17:25–41.

- Kuypers, T. (2000). Information and efficiency: An empirical study of a fixed odds betting market. *Applied Economics*, 32(11):1353–1363.
- Kwon, H. D. and Lippman, S. A. (2011). Acquisition of project-specific assets with bayesian updating. *Operations Research*, 59(5):1119–1130.
- Kyle, A. S. (1985). Continuous auctions and insider trading. *Econometrica*, 53(6):1315–1335.
- Law, D. and Peel, D. A. (2002). Insider trading, herding behaviour and market plungers in the British horse–race betting market. *Economica*, 69(274):327–338.
- Levitt, S. D. (2004). Why are gambling markets organised so differently from financial markets? *The Economic Journal*, 114(495):223–246.
- Lipster, R. and Shiriyayev, A. (1977). *Statistics of Random Processes I*. New York : Springer-Verlag.
- McCricrick, J. (1992). *World of Betting*. London: Stanley Paul.
- Newey, W. K. and West, K. D. (1987). A simple, positive semi-definite, heteroskedasticity and autocorrelation consistent covariance matrix. *Econometrica*, 55(3):703–708.
- Øksendal, B. (2005). *Stochastic Differential Equations*. Berlin, Heidelberg, New York: Springer.
- Peskir, G. and Shiriyayev, A. (2006). *Optimal Stopping and Free Boundary Problems*. Basel: Birkhäuser Verlag.
- Poor, H. V. and Hadjiliadis, O. (2009). *Quickest Detection*. Cambridge University Press Cambridge.
- Rottenberg, S. (1956). The baseball players' labour market. *Journal of Political Economy*, 64(3):242–258.
- Ryan, R. and Lippman, S. A. (2003). Optimal exit from a project with noisy returns. *Probability in the Engineering and Informational Sciences*, 17(04):435–458.

- Sandford, J. and Shea, P. (2013). Optimal setting of point spreads. *Economica*, 80(317):149–170.
- Schnytzer, A. and Snir, A. (2008). Herding in imperfect betting markets with inside traders. *Journal of Gambling Business and Economics*, 2(2):1–15.
- Shin, H. S. (1991). Optimal betting odds against insider traders. *The Economic Journal*, 101(408):1179–1185.
- Shin, H. S. (1992). Prices of state contingent claims with insider traders, and the favourite-longshot bias. *The Economic Journal*, 102(411):426–435.
- Shin, H. S. (1993). Measuring the incidence of insider trading in a market for state-contingent claims. *The Economic Journal*, 103(420):1141–1153.
- Shing, H.-F. and Koch, A. (2008). Bookmaker and pari-mutuel betting: Is a (reverse) favourite-longshot bias built-in? *The Journal of Prediction Markets*, 2(2):29–50.
- Shiryaev, A. N. (1967). Two problems of sequential analysis. *Kibernetika*, 3(2):63–69.
- Shiryaev, A. N. (1978). *Optimal Stopping Rules*. New York: Springer-Verlag.
- Snyder, W. W. (1978). Horse racing: Testing the efficient markets model. *Journal of Finance*, 33(4):1109–1118.
- Thaler, R. H. and Ziemba, W. T. (1988). Anomalies: Parimutuel betting markets: Racetracks and lotteries. *The Journal of Economic Perspectives*, 2(2):161–174.
- Thijssen, J. (2003). *Investment under uncertainty, market evolution and coalition and spillovers in a game theoretic perspective*. PhD thesis, Tilburg University, Center for Economic Research.
- Thijssen, J. (2013). Principles of investment under uncertainty. Lecture Notes, University of York, Department of Economics and Related Studies.
- Thijssen, J. and Bregantini, D. (2016). Costly sequential experimentation and project valuation with an application to health technology assessment. Working paper, University of York.



- Torres-Reyna, O. (2007). Panel data analysis: Fixed and random effects using stata. Lecture Notes, Princeton University.
- Wald, A. (1973). *Sequential Analysis*. London: John Wiley & Sons.
- Williams, L. V. (1999). Information efficiency in betting markets: A survey. *Bulletin of Economic Research*, 51(1):1–39.
- Williams, L. V. and Paton, D. (1997). Why is there a favourite-longshot bias in British racetrack betting markets? *The Economic Journal*, 107(440):150–158.
- Wooldridge, J. M. (2010). *Econometric Analysis of Cross Section and Panel Data*. MIT press.