

## **Hong Kong Horse Racing Acquisition and Processing**

Jinru KUANG, Zhengyi LIU,

Xiaqing WU, Chenyan LYU

Department of Media and Communication

COM5507 Social Media Data Acquisition and Processing

Xiaofan LIU

December 13, 2024

**Author List and Contributions**

Members	Contributions
Jinru KUANG	25%
Zhengyi LIU	25%
Xiaqing, WU	25%
Chenyan LYU	25%
Total	100%

## Abstract

This report presents an initial exploration into the acquisition, processing, and preliminary analysis of Hong Kong horse racing data, setting the stage for more advanced predictive modeling in the future. Two primary data sources were selected to ensure both reliability and richness: the Hong Kong Jockey Club's official website, which offers authoritative race records, historical datasets, and detailed horse and jockey profiles, and MaJing, a local platform that provides supplementary insights and reference materials. Data extraction and preprocessing were carried out using Python, encompassing steps such as data cleaning, normalization, and restructuring, to ensure consistency and usability. Subsequently, fundamental visualization techniques—including histograms, scatter plots—were employed to uncover key patterns within the data. While the processed datasets and initial visualizations offer a useful starting point for newcomers to Hong Kong horse racing, it is important to note that this report represents only the initial stage of a more comprehensive research effort. Future work will focus on constructing and refining predictive models to facilitate more accurate forecasts and data-driven decision-making in the domain of Hong Kong horse racing.

## Introduction

### Background and Motivation

Horse racing, as a unique social event, embodies a vibrant aspect of Hong Kong's rich cultural tapestry. This exhilarating spectacle is not just about the speed of the horses or the strategy of betting, it's an occasion that brings together people from diverse backgrounds, creating a shared experience filled with excitement, anticipation, and community spirit. Both tourists and local residents flock to the races, drawn by the electrifying atmosphere that permeates the racecourse. The sights and sounds—the thundering hooves, the cheering crowds, and the elegant fashion—combine to create a memorable experience that captures the essence of Hong Kong's dynamic culture.

The motivation behind this project stems from the desire to help horse racing enthusiasts make more informed judgments regarding their betting strategies. By undertaking a comprehensive process that includes data scraping, selection, cleaning, analysis, and visualization, we aim to provide valuable insights into horse racing in Hong Kong. Our goal is to present clear and relevant information that aids bettors and fans in navigating the complexities of the sport. Through this meticulous approach, we aspire to create a resource that

empowers participants in this exhilarating arena, ultimately enhancing their overall experience and confidence in their betting choices.

## **Others Work and Our Angle**

Before delving deeper into our research, it is worth mentioning that there have been some excellent studies conducted in this field. Although most of the studies utilize machine learning algorithms to predict horse racing outcomes, we can still draw inspiration from the data types they employ to enhance the efficiency and accuracy of our data collection process. Gupta and Singh (2023) explored the impact of multiple factors on horse racing outcomes in their study, including information on horse characteristics and historical performance, to identify key variables that can improve prediction accuracy. Chavan et al. (2023) analyzed horse race performance and other related information, focusing on identifying the main factors that affect performance and providing relevant data analysis results. Both studies utilized horse racing data, with a focus on analyzing the impact of historical race records on predicting horse racing outcomes.

In contrast to these two articles from India, our angle is to leverage the insights from various studies to conduct research specific to horse racing in Hong Kong, where existing literature is

relatively scarce. Therefore, we are drawing on the data types identified in previous research and concentrating on comprehensive scraping of detailed information, historical race records, and relevant statistical data concerning current horses in Hong Kong.

## Methods

### Data Source

Regarding data source selection, we chose the Hong Kong Jockey Club and Majing website as the authoritative data source for our research. The Hong Kong Jockey Club provides up-to-date, official information, including basic information about the horses, race records, and relevant background information. At the same time, the rich historical data on *Majing* also provides us with a deep perspective, enabling us to analyze recent and long-term competition trends.

The URLs we scraped are as follows:

<https://racing.hkjc.com/racing/information/chinese/Horse>SelectHorse.aspx>

<https://racing.on.cc/cgi-bin/srh/search/search.cgi>

### W1-Hong Kong Jockey Club

On the Hong Kong Jockey Club website, we scraped the basic information and recent 3 season records of all active horses in Hong Kong based on the word count of their names.

Targeted web pages are as follows:

**賽事資料(本地) - 馬匹資料 - 馬匹資料**

[返回馬匹搜尋頁](#)

馬名字數			
• 未命名馬匹	• 二字馬	• 三字馬	• 四字馬

二字馬 (依評分由大至小橫向排列)							
• 安遇 (104)	• 安騁 (104)	• 氣勢 (98)	• 增有 (89)	• 球星 (83)			
• 瑪瑙 (82)	• 耀寶 (78)	• 晴王 (77)	• 玩笑 (76)	• 米奇 (76)			
• 優才 (73)	• 十力 (72)	• 伊臣 (71)	• 馬力 (71)	• 祝願 (71)			
• 豐盛 (71)	• 快路 (70)	• 將義 (70)	• 飲杯 (69)	• 增勁 (69)			
• 俊才 (68)	• 銀進 (67)	• 奔放 (67)	• 堅闖 (66)	• 增強 (66)			
• 佳登 (66)	• 福星 (66)	• 赤壁 (66)	• 翠紅 (65)	• 星價 (65)			
• 三強 (65)	• 快搏 (65)	• 安都 (64)	• 凡喜 (64)	• 旗鷹 (63)			
• 論文 (62)	• 初戀 (61)	• 爵登 (61)	• 大才 (61)	• 攻頂 (60)			
• 麋峯 (60)	• 勝意 (58)	• 爆笑 (58)	• 銀騰 (57)	• 安泰 (56)			
• 帥炸 (56)	• 同喜 (56)	• 爆熱 (56)	• 定數 (55)	• 真感 (54)			
• 神馳 (53)	• 上校 (53)	• 摩界 (53)	• 飛雲 (53)	• 黑白 (53)			

Selecting horses

安遇 (G236)		出生地 / 馬齡	練馬師	賀賢
		毛色 / 性別	棕 / 閹	馬主
• 往績紀錄		進口類別	自購馬	現時評分
• 馬匹評分/體重/名次		今季獎金*	\$187,250	季初評分
• 所跑途程賽績紀錄		總獎金*	\$14,181,025	父系
• 晨操紀錄		冠/亞/季-總出賽次數*	7-2-2-28	母系
• 傷患紀錄		最近十個賽馬日	: 2	外祖父
• 撤還紀錄		出賽場數		同父系馬
• 海外賽績紀錄		現在位置	香港	美麗之光
• 血統簡評		(到達日期)	(24/09/2024)	瀏覽
• 其他馬匹		進口日期	16/10/2021	

\*包括本地及海外賽績及獎金

Horse detailed information

馬匹近三季往績紀錄 - 安遇														近三季往績		三季前往績		所有往績	
場次	名次	日期	馬場/跑道/ 賽道	途程	場地 狀況	賽事 班次	檔位	評分	練馬師	騎師	頭馬 距離	獨贏 賠率	實際 賠率	沿途 走位	完成 時間	排位	配備	賽事 重播	
														▼	▲	▼	▲	▼	▲
24/25 馬季																			
190	05	17/11/24	沙田草地"B+2"	2000	好/快	G2	6	104	賀賢	布文	5	34	123	10 10 10 9 5	2.00.51	1174	B	▲	
153	08	03/11/24	沙田草地"C+3"	1800	好/快	G3	8	105	賀賢	霍兆聲	3-1/2	42	121	12 12 12 12 8	1.46.97	1161	B	▲	
117	12	20/10/24	沙田草地"B+2"	1200	好/快	G2	6	107	賀賢	班德禮	7-1/2	183	116	12 13 12	1.08.79	1168	B2/TT-	▲	
23/24 馬季																			
719	07	02/06/24	沙田草地"B"	1600	好/快	G3	4	111	文家良	田泰安	13-3/4	26	127	4 3 6 7	1.36.15	1115	TT	▲	
566	13	07/04/24	沙田草地"B+2"	1600	好	G2	4	113	文家良	田泰安	17-1/4	10	123	7 6 10 13	1.36.79	1125	TT	▲	
452	07	25/02/24	沙田草地"A+3"	2000	好	G1	2	114	文家良	田泰安	5-1/2	16	126	5 6 7 7 7	2.01.20	1146	TT	▲	
395	08	04/02/24	沙田草地"B+2"	1800	好	G3	12	114	文家良	梁家俊	3	15	132	2 3 5 5 8	1.47.89	1155	TT	▲	
330	08	10/01/24	跑馬地草地"B"	1800	好	G3	12	114	文家良	梁家俊	4-1/4	11	135	7 8 8 7 8	1.48.99	1164	TT	▲	
238	08	10/12/23	沙田草地"A"	1600	好	G1	1	114	文家良	梁家俊	4-3/4	24	126	6 6 7 8	1.34.88	1164	TT	▲	
144	01	05/11/23	沙田草地"C+3"	1800	好	G3	6	110	文家良	梁家俊	1-1/4	12	128	3 4 4 4 1	1.47.40	1164	TT	▲	
088	02	15/10/23	沙田草地"A+3"	1600	好	G2	4	105	文家良	梁家俊	1/2	14	117	3 3 4 2	1.34.28	1156	TT	▲	
036	06	24/09/23	沙田草地"C"	1400	好	G3	6	105	文家良	梁家俊	10-1/2	4.4	123	5 6 4 6	1.23.59	1153	TT	▲	
22/23 馬季																			
834	01	16/07/23	沙田草地"A"	1600	好	1	7	100	文家良	梁家俊	▼位	4.9	128	9 9 5 1	1.34.93	1165	TT	▲	
749	01	14/06/23	跑馬地草地"B"	1650	好	2	10	89	文家良	梁家俊	4-3/4	3.2	124	5 4 1 1	1.38.87	1153	TT	▲	
576	01	12/04/23	跑馬地草地"B"	1650	好	2	6	83	文家良	梁家俊	1	5.9	120	4 4 3 1	1.39.08	1138	TT	▲	
509	08	19/03/23	沙田草地"A"	2000	好	4YO	8	83	文家良	蘇兆輝	3-3/4	124	126	2 2 2 2 9	2.03.37	1152	B/H/T/T	▲	
453	12	26/02/23	沙田草地"A+3"	1800	好/快	4YO	9	83	文家良	蘇兆輝	8-1/2	28	126	7 6 11 12	1.47.62	1152	B/H/T/T	▲	
406	02	08/02/23	跑馬地草地"B"	1650	好	3	5	81	文家良	蘇兆輝	▼位	2.2	132	2 2 2 2	1.39.62	1153	B/T/T	▲	
328	01	11/01/23	跑馬地草地"B"	1650	好	3	12	74	文家良	蘇兆輝	1	4.2	131	2 2 1 1	1.39.53	1137	B/T/T	▲	
250	01	14/12/22	跑馬地草地"B"	1650	好	3	3	67	文家良	蘇兆輝	2-1/4	3.5	126	2 3 4 1	1.39.13	1129	B/T/T	▲	
166	03	12/11/22	沙田全天候	1650	好	3	6	68	文家良	潘明輝	9	11	121	2 2 2 3	1.38.22	1133	TT	▲	
116	05	23/10/22	沙田草地"B+2"	1600	好/快	3	11	69	文家良	潘明輝	2	5.8	124	2 2 2 5	1.33.80	1136	TT	▲	
044	06	25/09/22	沙田草地"C"	1600	好/快	3	2	69	文家良	潘明輝	2-3/4	4.2	124	3 5 5 6	1.35.23	1143	TT	▲	
010	03	11/09/22	沙田草地"A"	1400	好	3	12	69	文家良	潘明輝	1-1/4	8.4	122	4 2 2 3	1.22.58	1140	TT	▲	

Recent 3 seasons racing records

The inspecting the element windows are as follows:

```

<table width="760" border="0" class="bigborder border="0" cellspacing="1" cellpadding="0">
  <tr>
    <td class="subheader">馬名字數</td>
  </tr>
  <tr>
    <td>
      <table width="760" border="0" border=0 cellspacing=1 cellpadding=0>
        <tr>
          <td bgcolor="#FFFFFF" width="190" class="table_text_two">
            <ul><li class="table_text_two"><a href="/racing/information/chinese/Horse>SelectHorsebyChar.aspx?ordertype=0" class="table_text_two">未命名馬匹</a></li></ul>
          </td>
          <td bgcolor="#FFFFFF" width="190" class="table_text_two">
            <ul><li class="table_text_two"><a href="/racing/information/chinese/Horse>SelectHorsebyChar.aspx?ordertype=2" class="table_text_two">二字馬</a></li></ul>
          </td>
          <td bgcolor="#FFFFFF" width="190" class="table_text_two">
            <ul><li class="table_text_two"><a href="/racing/information/chinese/Horse>SelectHorsebyChar.aspx?ordertype=3" class="table_text_two">三字馬</a></li></ul>
          </td>
          <td bgcolor="#FFFFFF" width="190" class="table_text_two">
            <ul><li class="table_text_two"><a href="/racing/information/chinese/Horse>SelectHorsebyChar.aspx?ordertype=1" class="table_text_two">四字馬</a></li></ul>
          </td>
        </tr>
      </table>
    </td>
  </tr>
  <br />
<table width="760" border="0" class="bigborder border="0" cellspacing="1" cellpadding="0">

```

The character of the horse's name

```
<tr>
    <td class="subheader">二字馬 (依詳分由大至小橫向排列)</td>
</tr>
<tr>
    <td>
        <table width="260" border="0" cellspacing="1" cellpadding="0">
            <tr>
                <td width=20% bgcolor="#FFFFFF" class="table_eng_text">
                    <table width="100%" cellspacing="0" cellpadding="0" border="0">
                        <tr>
                            <td class="table_text_two" hcolor="#ffffff" align="left" style="padding:4px 2px; width:100px;">
                                <ul><li class="table_text_two" style="padding:4px 2px;"><a href="/racing/information/chinese/Horse/Horse.aspx?HorseId=HK_2021_6238" class="ltable_text" style="color:#000000;">安遇安明氣勢
```

Each horse

## Specific information

馬西逐季往績記錄 - 安通	
<tr>	
<td &gt;<="" border="0" style="width:1000" td=""></td>	
<tr>	
<td align="left" class="table_text_two" style="width:639">&lt;b&gt;馬西逐季往績記錄 - 安通&lt;/b&gt;</td>	<b>馬西逐季往績記錄 - 安通</b>
</tr>	
<table cellpadding="0" cellspacing="0" border="0" width="329" align="right">	
<tr value="1">	
<td style="text-align:center">&lt;img src="/racing/kkic.com/racing/content/Tables/StaticFile/Chinese/hd_3r.htm_a.jpg" alt="第三季目標" style="width:111px; height: 24px;" border="0" /&gt;&lt;br/&gt;</td>	 
<td style="text-align:center">&lt;a href="/racing/information/Chinese/Boxcar_boxcar.html?boxcarID=192_2021_52296&amp;viewOrder=2"&gt;&lt;img src="/racing/kkic.com/racing/content/Tables/StaticFile/Chinese/hd_b2r.htm.jpg" alt="第三季往績" style="width:107px; height: 24px;" border="0" id="hd_b2r_boxc2r_1" /&gt;&lt;br/&gt;</td>	<a href="/racing/information/Chinese/Boxcar_boxcar.html?boxcarID=192_2021_52296&viewOrder=2"> 
<td style="text-align:center">&lt;a href="/racing/information/Chinese/Boxcar_boxcar.html?boxcarID=192_2021_52296&amp;viewOrder=1"&gt;&lt;img src="/racing/kkic.com/racing/content/Tables/StaticFile/Chinese/hd_all.htm.jpg" alt="所有往績" style="width:111px; height: 24px;" border="0" id="hd_all_boxc1r_1" /&gt;&lt;br/&gt;</td>	<a href="/racing/information/Chinese/Boxcar_boxcar.html?boxcarID=192_2021_52296&viewOrder=1"> 
</tr>	
</table>	
</tr>	
<table border="0" cellpadding="1" cellspacing="1" class="borderheader" width="1000">	
<tr>	
<td class="tableheader" style="width:50">馬西</td>	馬西
<td class="tableheader" style="width:50">各季</td>	各季
<td class="tableheader" style="width:50">各次</td>	各次
<td class="tableheader" style="width:100">馬西/各季</td>	馬西/各季
<td class="tableheader" style="width:100">馬西/各次</td>	馬西/各次
<td class="tableheader" style="width:50">各季次</td>	各季次
<td class="tableheader" style="width:50">各次</td>	各次
<td class="tableheader" style="width:50">部分</td>	部分
<td class="tableheader" style="width:100">總馬西</td>	總馬西
<td class="tableheader" style="width:100">各季</td>	各季
<td class="tableheader" style="width:100">各次</td>	各次
<td class="tableheader" style="width:50">總次</td>	總次
<td class="tableheader" style="width:50">各季次</td>	各季次
<td class="tableheader" style="width:100">各季</td>	各季
<td class="tableheader" style="width:100">各次</td>	各次
<td class="tableheader" style="width:50">部分</td>	部分
<td class="tableheader" style="width:100">總馬西</td>	總馬西
<td class="tableheader" style="width:100">各季</td>	各季
<td class="tableheader" style="width:100">各次</td>	各次
<td class="tableheader" style="width:50">總次</td>	總次
<td class="tableheader" style="width:50">各季次</td>	各季次
<td class="tableheader" style="width:100">各季</td>	各季
<td class="tableheader" style="width:100">各次</td>	各次
<td class="tableheader" style="width:50">部分</td>	部分
<td class="tableheader" style="width:100">總馬西</td>	總馬西
<td class="tableheader" style="width:100">各季</td>	各季
<td class="tableheader" style="width:100">各次</td>	各次
<td class="tableheader" style="width:50">總次</td>	總次
<td class="tableheader" style="width:50">各季次</td>	各季次
<td class="tableheader" style="width:100">各季</td>	各季
<td class="tableheader" style="width:100">各次</td>	各次
<td class="tableheader" style="width:50">部分</td>	部分
<td class="tableheader" style="width:100">總馬西</td>	總馬西
<td class="tableheader" style="width:100">各季</td>	各季
<td class="tableheader" style="width:100">各次</td>	各次
<td class="tableheader" style="width:50">總次</td>	總次
<td class="tableheader" style="width:50">各季次</td>	各季次
<td class="tableheader" style="width:100">各季</td>	各季
<td class="tableheader" style="width:100">各次</td>	各次
<td class="tableheader" style="width:50">部分</td>	部分
<td class="tableheader" style="width:100">總馬西</td>	總馬西
<td class="tableheader" style="width:100">各季</td>	各季
<td class="tableheader" style="width:100">各次</td>	各次
<td class="tableheader" style="width:50">總次</td>	總次
<td class="tableheader" style="width:50">各季次</td>	各季次
<td class="tableheader" style="width:100">各季</td>	各季
<td class="tableheader" style="width:100">各次</td>	各次
<td class="tableheader" style="width:50">部分</td>	部分
<td class="tableheader" style="width:100">總馬西</td>	總馬西
<td class="tableheader" style="width:100">各季</td>	各季
<td class="tableheader" style="width:100">各次</td>	各次
<td class="tableheader" style="width:50">總次</td>	總次
<td class="tableheader" style="width:50">各季次</td>	各季次
<td class="tableheader" style="width:100">各季</td>	各季
<td class="tableheader" style="width:100">各次</td>	各次
<td class="tableheader" style="width:50">部分</td>	部分
<td class="tableheader" style="width:100">總馬西</td>	總馬西
<td class="tableheader" style="width:100">各季</td>	各季
<td class="tableheader" style="width:100">各次</td>	各次
<td class="tableheader" style="width:50">總次</td>	總次
<td class="tableheader" style="width:50">各季次</td>	各季次
<td class="tableheader" style="width:100">各季</td>	各季
<td class="tableheader" style="width:100">各次</td>	各次
<td class="tableheader" style="width:50">部分</td>	部分
<td class="tableheader" style="width:100">總馬西</td>	總馬西
<td class="tableheader" style="width:100">各季</td>	各季
<td class="tableheader" style="width:100">各次</td>	各次
<td class="tableheader" style="width:50">總次</td>	總次
<td class="tableheader" style="width:50">各季次</td>	各季次
<td class="tableheader" style="width:100">各季</td>	各季
<td class="tableheader" style="width:100">各次</td>	各次
<td class="tableheader" style="width:50">部分</td>	部分
<td class="tableheader" style="width:100">總馬西</td>	總馬西
<td class="tableheader" style="width:100">各季</td>	各季
<td class="tableheader" style="width:100">各次</td>	各次
<td class="tableheader" style="width:50">總次</td>	總次
<td class="tableheader" style="width:50">各季次</td>	各季次
<td class="tableheader" style="width:100">各季</td>	各季
<td class="tableheader" style="width:100">各次</td>	各次
<td class="tableheader" style="width:50">部分</td>	部分
<td class="tableheader" style="width:100">總馬西</td>	總馬西
<td class="tableheader" style="width:100">各季</td>	各季
<td class="tableheader" style="width:100">各次</td>	各次
<td class="tableheader" style="width:50">總次</td>	總次
<td class="tableheader" style="width:50">各季次</td>	各季次
<td class="tableheader" style="width:100">各季</td>	各季
<td class="tableheader" style="width:100">各次</td>	各次
<td class="tableheader" style="width:50">部分</td>	部分
<td class="tableheader" style="width:100">總馬西</td>	總馬西
<td class="tableheader" style="width:100">各季</td>	各季
<td class="tableheader" style="width:100">各次</td>	各次
<td class="tableheader" style="width:50">總次</td>	總次
<td class="tableheader" style="width:50">各季次</td>	各季次
<td class="tableheader" style="width:100">各季</td>	各季
<td class="tableheader" style="width:100">各次</td>	各次
<td class="tableheader" style="width:50">部分</td>	部分
<td class="tableheader" style="width:100">總馬西</td>	總馬西
<td class="tableheader" style="width:100">各季</td>	各季
<td class="tableheader" style="width:100">各次</td>	各次
<td class="tableheader" style="width:50">總次</td>	總次
<td class="tableheader" style="width:50">各季次</td>	各季次
<td class="tableheader" style="width:100">各季</td>	各季
<td class="tableheader" style="width:100">各次</td>	各次
<td class="tableheader" style="width:50">部分</td>	部分
<td class="tableheader" style="width:100">總馬西</td>	總馬西
<td class="tableheader" style="width:100">各季</td>	各季
<td class="tableheader" style="width:100">各次</td>	各次
<td class="tableheader" style="width:50">總次</td>	總次
<td class="tableheader" style="width:50">各季次</td>	各季次
<td class="tableheader" style="width:100">各季</td>	各季
<td class="tableheader" style="width:100">各次</td>	各次
<td class="tableheader" style="width:50">部分</td>	部分
<td class="tableheader" style="width:100">總馬西</td>	總馬西
<td class="tableheader" style="width:100">各季</td>	各季
<td class="tableheader" style="width:100">各次</td>	各次
<td class="tableheader" style="width:50">總次</td>	總次
<td class="tableheader" style="width:50">各季次</td>	各季次
<td class="tableheader" style="width:100">各季</td>	各季
<td class="tableheader" style="width:100">各次</td>	各次
<td class="tableheader" style="width:50">部分</td>	部分
<td class="tableheader" style="width:100">總馬西</td>	總馬西
<td class="tableheader" style="width:100">各季</td>	各季
<td class="tableheader" style="width:100">各次</td>	各次
<td class="tableheader" style="width:50">總次</td>	總次
<td class="tableheader" style="width:50">各季次</td>	各季次
<td class="tableheader" style="width:100">各季</td>	各季
<td class="tableheader" style="width:100">各次</td>	各次
<td class="tableheader" style="width:50">部分</td>	部分
<td class="tableheader" style="width:100">總馬西</td>	總馬西
<td class="tableheader" style="width:100">各季</td>	各季
<td class="tableheader" style="width:100">各次</td>	各次
<td class="tableheader" style="width:50">總次</td>	總次
<td class="tableheader" style="width:50">各季次</td>	各季次
<td class="tableheader" style="width:100">各季</td>	各季
<td class="tableheader" style="width:100">各次</td>	各次
<td class="tableheader" style="width:50">部分</td>	部分
<td class="tableheader" style="width:100">總馬西</td>	總馬西
<td class="tableheader" style="width:100">各季</td>	各季
<td class="tableheader" style="width:100">各次</td>	各次
<td class="tableheader" style="width:50">總次</td>	總次
<td class="tableheader" style="width:50">各季次</td>	各季次
<td class="tableheader" style="width:100">各季</td>	各季
<td class="tableheader" style="width:100">各次</td>	各次
<td class="tableheader" style="width:50">部分</td>	部分
<td class="tableheader" style="width:100">總馬西</td>	總馬西
<td class="tableheader" style="width:100">各季</td>	各季
<td class="tableheader" style="width:100">各次</td>	各次
<td class="tableheader" style="width:50">總次</td>	總次
<td class="tableheader" style="width:50">各季次</td>	各季次
<td class="tableheader" style="width:100">各季</td>	各季
<td class="tableheader" style="width:100">各次</td>	各次
<td class="tableheader" style="width:50">部分</td>	部分
<td class="tableheader" style="width:100">總馬西</td>	總馬西
<td class="tableheader" style="width:100">各季</td>	各季
<td class="tableheader" style="width:100">各次</td>	各次
<td class="tableheader" style="width:50">總次</td>	總次
<td class="tableheader" style="width:50">各季次</td>	各季次
<td class="tableheader" style="width:100">各季</td>	各季
<td class="tableheader" style="width:100">各次</td>	各次
<td class="tableheader" style="width:50">部分</td>	部分
<td class="tableheader" style="width:100">總馬西</td>	總馬西
<td class="tableheader" style="width:100">各季</td>	各季
<td class="tableheader" style="width:100">各次</td>	各次
<td class="tableheader" style="width:50">總次</td>	總次
<td class="tableheader" style="width:50">各季次</td>	各季次
<td class="tableheader" style="width:100">各季</td>	各季
<td class="tableheader" style="width:100">各次</td>	各次
<td class="tableheader" style="width:50">部分</td>	部分
<td class="tableheader" style="width:100">總馬西</td>	總馬西
<td class="tableheader" style="width:100">各季</td>	各季
<td class="tableheader" style="width:100">各次</td>	各次
<td class="tableheader" style="width:50">總次</td>	總次
<td class="tableheader" style="width:50">各季次</td>	各季次
<td class="tableheader" style="width:100">各季</td>	各季
<td class="tableheader" style="width:100">各次</td>	各次
<td class="tableheader" style="width:50">部分</td>	部分
<td class="tableheader" style="width:100">總馬西</td>	總馬西
<td class="tableheader" style="width:100">各季</td>	各季
<td class="tableheader" style="width:100">各次</td>	各次
<td class="tableheader" style="width:50">總次</td>	總次
<td class="tableheader" style="width:50">各季次</td>	各季次
<td class="tableheader" style="width:100">各季</td>	各季
<td class="tableheader" style="width:100">各次</td>	各次
<td class="tableheader" style="width:50">部分</td>	部分
<td class="tableheader" style="width:100">總馬西</td>	總馬西
<td class="tableheader" style="width:100">各季</td>	各季
<td class="tableheader" style="width:100">各次</td>	各次
<td class="tableheader" style="width:50">總次</td>	總次
<td class="tableheader" style="width:50">各季次</td>	各季次
<td class="tableheader" style="width:100">各季</td>	各季
<td class="tableheader" style="width:100">各次</td>	各次
<td class="tableheader" style="width:50">部分</td>	部分
<td class="tableheader" style="width:100">總馬西</td>	總馬西
<td class="tableheader" style="width:100">各季</td>	各季
<td class="tableheader" style="width:100">各次</td>	各次
<td class="tableheader" style="width:50">總次</td>	總次
<td class="tableheader" style="width:50">各季次</td>	各季次
<td class="tableheader" style="width:100">各季</td>	各季
<td class="tableheader" style="width:100">各次</td>	各次
<td class="tableheader" style="width:50">部分</td>	部分
<td class="tableheader" style="width:100">總馬西</td>	總馬西
<td class="tableheader" style="width:100">各季</td>	各季
<td class="tableheader" style="width:100">各次</td>	各次
<td class="tableheader" style="width:50">總次</td>	總次
<td class="tableheader" style="width:50">各季次</td>	各季次
<td class="tableheader" style="width:100">各季</td>	各季
<td class="tableheader" style="width:100">各次</td>	各次
<td class="tableheader" style="width:50">部分</td>	部分
<td class="tableheader" style="width:100">總馬西</td>	總馬西
<td class="tableheader" style="width:100">各季</td>	各季
<td class="tableheader" style="width:100">各次</td>	各次
<td class="tableheader" style="width:50">總次</td>	總次
<td class="tableheader" style="width:50">各季次</td>	各季次
<td class="tableheader" style="width:100">各季</td>	各季
<td class="tableheader" style="width:100">各次</td>	各次
<td class="tableheader" style="width:50">部分</td>	部分
<td class="tableheader" style="width:100">總馬西</td>	總馬西
<td class="tableheader" style="width:100">各季</td>	各季
<td class="tableheader" style="width:100">各次</td>	各次
<td class="tableheader" style="width:50">總次</td>	總次
<td class="tableheader" style="width:50">各季次</td>	各季次
<td class="tableheader" style="width:100">各季</td>	各季
<td class="tableheader" style="width:100">各次</td>	各次
<td class="tableheader" style="width:50">部分</td>	部分
<td class="tableheader" style="width:100">總馬西</td>	總馬西
<td class="tableheader" style="width:100">各季</td>	各季
<td class="tableheader" style="width:100">各次</td>	各次
<td class="tableheader" style="width:50">總次</td>	總次
<td class="tableheader" style="width:50">各季次</td>	各季次
<td class="tableheader" style="width:100">各季</td>	各季
<td class="tableheader" style="width:100">各次</td>	各次
<td class="tableheader" style="width:50">部分</td>	部分
<td class="tableheader" style="width:100">總馬西</td>	總馬西
<td class="tableheader" style="width:100">各季</td>	各季
<td class="tableheader" style="width:100">各次</td>	各次
<td class="tableheader" style="width:50">總次</td>	總次
<td class="tableheader" style="width:50">各季次</td>	各季次
<td class="tableheader" style="width:100">各季</td>	各季
<td class="tableheader" style="width:100">各次</td>	各次
<td class="tableheader" style="width:50">部分</td>	部分
<td class="tableheader" style="width:100">總馬西</td>	總馬西
<td class="tableheader" style="width:100">各季</td>	各季
<td class="tableheader" style="width:100">各次</td>	各次
<td class="tableheader" style="width:50">總次</td>	總次
<td class="tableheader" style="width:50">各季次</td>	各季次
<td class="tableheader" style="width:100">各季</td>	各季
<td class="tableheader" style="width:100">各次</td>	各次
<td class="tableheader" style="width:50">部分</td>	部分
<td class="tableheader" style="width:100">總馬西</td>	總馬西
<td class="tableheader" style="width:100">各季</td>	各季
<td class="tableheader" style="width:100">各次</td>	各次
<td class="tableheader" style="width:50">總次</td>	總次
<td class="tableheader" style="width:50">各季次</td>	各季次
<td class="tableheader" style="width:100">各季</td>	各季
<td class="tableheader" style="width:100">各次</td>	各次
<td class="tableheader" style="width:50">部分</td>	部分
<td class="tableheader" style="width:100">總馬西</td>	總馬西
<td class="tableheader" style="width:100">各季</td>	各季
<td class="tableheader" style="width:100">各次</td>	各次
<td class="tableheader" style="width:50">總次</td>	總次
<td class="tableheader" style="width:50">各季次</td>	各季次
<td class="tableheader" style="width:100">各季</td>	各季
<td class="tableheader" style="width:100">各次</td>	各次
<td class="tableheader" style="width:50">部分</td>	部分
<td class="tableheader" style="width:100">總馬西</td>	總馬西
<td class="tableheader" style="width:100">各季</td>	各季
<td class="tableheader" style="width:100">各次</td>	各次
<td class="tableheader" style="width:50">總次</td>	總次
<td class="tableheader" style="width:50">各季次</td>	各季次
<td class="tableheader" style="width:100">各季</td>	各季
<td class="tableheader" style="width:100">各次</td>	各次
<td class="tableheader" style="width:50">部分</td>	部分
<td class="tableheader" style="width:100">總馬西</td>	總馬西
<td class="tableheader" style="width:100">各季</td>	各季
<td class="tableheader" style="width:100">各次</td>	各次
<td class="tableheader" style="width:50">總次</td>	總次
<td class="tableheader" style="width:50">各季次</td>	各季次
<td class="tableheader" style="width:100">各季</td>	各季
<td class="tableheader" style="width:100">各次</td>	各次
<td class="tableheader" style="width:50">部分</td>	部分
<td class="tableheader" style="width:100">總馬西</td>	總馬西
<td class="tableheader" style="width:100">各季</td>	各季
<td class="tableheader" style="width:100">各次</td>	各次
<td class="tableheader" style="width:50">總次</td>	總次
<td class="tableheader" style="width:50">各季次</td>	各季次
<td class="tableheader" style="width:100">各季</td>	各季
<td class="tableheader" style="width:100">各次</td>	各次
<td class="tableheader" style="width:50">部分</td>	部分
<td class="tableheader" style="width:100">總馬西</td>	總馬西
<td class="tableheader" style="width:100">各季</td>	各季
<td class="tableheader" style="width:100">各次</td>	各次
<td class="tableheader" style="width:50">總次</td>	總次
<td class="tableheader" style="width:50">各季次</td>	各季次
<td class="tableheader" style="width:100">各季</td>	各季
<td class="tableheader" style="width:100">各次</td>	各次
<td class="tableheader" style="width:50">部分</td>	部分
<td class="tableheader" style="width:100">總馬西</td>	總馬西
<td class="tableheader" style="width:100">各季</td>	各季
<td class="tableheader" style="width:100">各次</td>	各次
<td class="tableheader" style="width:50">總次</td>	總次
<td class="tableheader" style="width:50">各季次</td>	各季次
<td class="tableheader" style="width:100">各季</td>	各季
<td class="tableheader" style="width:100">各次</td>	各次
<td class="tableheader" style="width:50">部分</td>	部分
<td class="tableheader" style="width:100">總馬西</td>	總馬西
<td class="tableheader" style="width:100">各季</td>	各季
<td class="tableheader" style="width:100">各次</td>	各次
<td class="tableheader" style="width:50">總次</td>	總次
<td class="tableheader" style="width:50">各季次</td>	各季次
<td class="tableheader" style="width:100">各季</td>	各季
<td class="tableheader" style="width:100">各次</td>	各次
<td class="tableheader" style="width:50">部分</td>	部分
<td class="tableheader" style="width:100">總馬西</td>	總馬西
<td class="tableheader" style="width:100">各季</td>	各季
<td class="tableheader" style="width:100">各次</td>	各次
<td class="tableheader" style="width:50">總次</td>	總次
<td class="tableheader" style="width:50">各季次</td>	各季次
<td class="tableheader" style="width:100">各季</td>	各季
<td class="tableheader" style="width:100">各次</td>	各次
<td class="tableheader" style="width:50">部分</td>	部分
<td class="tableheader" style="width:100">總馬西</td>	總馬西
<td class="tableheader" style="width:100">各季</td>	各季
<td class="tableheader" style="width:100">各次</td>	各次
<td class="tableheader" style="width:50">總次</td>	總次
<td class="tableheader" style="width:50">各季次</td>	各季次
<td class="tableheader" style="width:100">各季</td>	各季
<td class="tableheader" style="width:100">各次</td>	各次
<td class="tableheader" style="width:50">部分</td>	部分
<td class="tableheader" style="width:100">總馬西</td>	總馬西
<td class="tableheader" style="width:100">各季</td>	各季
<td class="tableheader" style="width:100">各次</td>	各次
<td class="tableheader" style="width:50">總次</td>	總次
<td class="tableheader" style="width:50">各季次</td>	各季次
<td class="tableheader" style="width:100">各季</td>	各季
<td class="tableheader" style="width:100">各次</td>	各次
<td class="tableheader" style="width:50">部分</td>	部分
<td class="tableheader" style="width:100">總馬西</td>	總馬西
<td class="tableheader" style="width:100">各季</td>	各季
<td class="tableheader" style="width:100">各次</td>	各次
<td class="tableheader" style="width:50">總次</td>	總次
<td class="tableheader" style="width:50">各季次</td>	各季次
<td class="tableheader" style="width:100">各季</td>	各季
<td class="tableheader" style="width:100">各次</td>	各次
<td class="tableheader" style="width:50">部分</td>	部分
<td class="tableheader" style="width:100">總馬西</td>	總馬西
<td class="tableheader" style="width:100">各季</td>	各季
<td class="tableheader" style="width:100">各次</td>	各次
<td class="tableheader" style="width:50">總次</td>	總次
<td class="tableheader" style="width:50">各季次</td>	各季次
<td class="tableheader" style="width:100">各季</td>	各季
<td class="tableheader" style="width:100">各次</td>	各次
<td class="tableheader" style="width:50">部分</td>	部分
<td class="tableheader" style="width:100">總馬西</td>	總馬西
<td class="tableheader" style="width:100">各季</td>	各季
<td class="tableheader" style="width:100">各次</td>	各次
<td class="tableheader" style="width:50">總次</td>	總次
<td class="tableheader" style="width:50">各季次</td>	各季次
<td class="tableheader" style="width:100">各季</td>	各季
<td class="tableheader" style="width:100">各次</td>	各次
<td class="tableheader" style="width:50">部分</td>	部分
<td class="tableheader" style="width:100">總馬西</td>	總馬西
<td class="tableheader" style="width:100">各季</td>	各季
<td class="tableheader" style="width:100">各次</td>	各次
<td class="tableheader" style="width:50">總次</td>	總次
<td class="tableheader" style="width:50">各季次</td>	各季次
<td class="tableheader" style="width:100">各季</td>	各季
<td class="tableheader" style="width:100">各次</td>	各次
<td class="tableheader" style="width:50">部分</td>	部分
<td class="tableheader" style="width:100">總馬西</td>	總馬西
<td class="tableheader" style="width:100">各季</td>	各季
<td class="tableheader" style="width:100">各次</td>	各次
<td class="tableheader" style="width:50">總次</td>	總次
<td class="tableheader" style="width:50">各季次</td>	各季次
<td class="tableheader" style="width:100">各季</td>	各季
<td class="tableheader" style="width:100">各次</td>	各次
<td class="tableheader" style="width:50">部分</td>	部分
<td class="tableheader" style="width:100">總馬西</td>	總馬西
<td class="tableheader" style="width:100">各季</td>	各季
<td class="tableheader" style="width:100">各次</td>	各次
<td class="tableheader" style="width:50">總次</td>	總次
<td class="tableheader" style="width:50">各季次</td>	各季次
<td class="tableheader" style="width:100">各季</td>	各季
<td class="tableheader" style="width:100">各次</td>	各次
<td class="tableheader" style="width:50">部分</td>	部分
<td class="tableheader" style="width:100">總馬西</td>	總馬西
<td class="tableheader" style="width:100">各季</td>	各季
<td class="tableheader" style="width:100">各次</td>	各次
<td class="tableheader" style="width:50">總次</td>	總次
<td class="tableheader" style="width:50">各季次</td>	各季次
<td class="tableheader" style="width:100">各季</td>	各季
<td class="tableheader" style="width:100">各次</td>	各次
<td class="tableheader" style="width:50">部分</td>	部分
<td class="tableheader" style="width:100">總馬西</td>	總馬西
<td class="tableheader" style="width:100">各季</td>	各季
<td class="tableheader" style="width:100">各次</td>	各次
<td class="tableheader" style="width:50">總次</td>	總次
<td class="tableheader" style="width:50">各季次</td>	各季次
<td class="tableheader" style="width:100">各季</td>	各季
<td class="tableheader" style="width:100">各次</td>	各次
<td class="tableheader" style="width:50">部分</td>	部分
<td class="tableheader" style="width:100">總馬西</td>	總馬西
<td class="tableheader" style="width:100">各季</td>	各季
<td class="tableheader" style="width:100">各次</td>	各次
<td class="tableheader" style="width:50">總次</td>	總次
<td class="tableheader" style="width:50">各季次</td>	各季次
<td class="tableheader" style="width:100">各季</td>	各季
<td class="tableheader" style="width:100">各次</td>	各次
<td class="tableheader" style="width:50">部分</td>	部分
<td class="tableheader" style="width:100">總馬西</td>	總馬西
<td class="tableheader" style="width:100">各季</td>	各季
<td class="tableheader" style="width:100">各次</td>	各次
<td class="tableheader" style="width:50">總次</td>	總次
<td class="tableheader" style="width:50">各季次</td>	各季次
<td class="tableheader" style="width:100">各季</td>	各季
<td class="tableheader" style="width:100">各次</td>	各次
<td class="tableheader" style="width:50">部分</td>	部分
<td class="tableheader" style="width:100">總馬西</td>	總馬西
<td class="tableheader" style="width:100">各季</td>	各季
<td class="tableheader" style="width:100">各次</td>	各次
<td class="tableheader" style="width:50">總次</td>	總次
<td class="tableheader" style="width:50">各季次</td>	各季次
<td class="tableheader" style="width:100">各季</td>	各季
<td class="tableheader" style="width:100">各次</td>	各次
<td class="tableheader" style="width:50">部分</td>	部分
<td class="tableheader" style="width:100">總馬西</td>	總馬西
<td class="tableheader" style="width:100">各季</td>	各季
<td class="tableheader" style="width:100">各次</td>	各次
<td class="tableheader" style="width:50">總次</td>	總次
<td class="tableheader" style="width:50">各季次</td>	各季次
<td class="tableheader" style="width:100">各季</td>	各季

## Racing records

## Data scraping process.

## **(1) Webpage structure and its techniques**

The Hong Kong Jockey Club categorized horse information by the number of characters in the horse names (such as "二字馬" and "三字馬"), providing a main URL for horse selection according to character count. This structure includes HTML tables that display the horse names and links to detailed pages for each horse, where relevant information is organized within specific HTML tables.

## (2) Data Scraping Methods

During the data scraping process, we used the following Python libraries:

```
import requests
from bs4 import BeautifulSoup
import pandas as pd
import time
```

**Requests:** Used to send HTTP requests to retrieve the HTML content of a webpage.

**BeautifulSoup**: used to parse HTML and extract the required data.

**Pandas**: Used to organize the scraped data into a DataFrame and save it as a CSV file.

**Time**: To manage delays between requests and avoid overwhelming the server.

During the scraping process, we first define the categories of horses, build URLs for each type of horse, and send a GET request to visit the corresponding page. Here are the main steps in the scraping process:

**Send request**: Use `requests.get()` to send HTTP requests and obtain responses, and use headers to simulate browser requests to avoid being recognized by websites as crawlers.

```
# 发送 HTTP GET 请求以获取马匹列表
headers = {
    'User-Agent': 'Mozilla/5.0 (Windows NT 10.0; Win64; x64) AppleWebKit/537.36 (KHTML, like Gecko) Chrome/89.0.4389.90 Safari/537.36'
}
response = requests.get(trainer_url, headers=headers)
```

**HTML parsing**: Use BeautifulSoup to parse HTML and use the `.find()` and `.find_all()` methods to locate specific HTML elements.

```
# 解析 HTML 内容
soup = BeautifulSoup(response.text, 'html.parser')

# 提取所有马匹的信息
horse_links = []

# 查找当前马匹类型部分
horse_table = soup.find('td', class_='subheader', string=f'{horse_name} (依评分由大至小横向排列)')
if horse_table:
    # 找到包含具体马匹信息的表格
    inner_table = horse_table.find_next('table')
    rows = inner_table.find_all('tr')
```

**Iteratively capture horse information**: Iterate through each horse's links to get its detailed information.

```

# 遍历每匹马并提取详细信息
unique_horse_names = set()
for horse_name, horse_url in horse_links:
    if horse_name in unique_horse_names:
        continue # 如果马匹名称已存在，跳过
    unique_horse_names.add(horse_name) # 添加到集合中，以防重复

# 解析马匹基本信息
horse_info = {'马匹名称': horse_name}

details = horse_soup.find_all('table', class_='table_top_right table_eng_text')
for detail in details:
    rows = detail.find_all('tr')
    for row in rows:
        cols = row.find_all('td')
        if len(cols) >= 3:
            key = cols[0].get_text(strip=True)
            value = cols[2].get_text(strip=True)
            horse_info[key] = value

# 提取近三季往绩记录
performance_table = horse_soup.find('table', class_='bigborder')
if performance_table:
    records = performance_table.find_all('tr')[2:] # 跳过前两行表头

```

**Data storage:** Store the scraped information in a list, and finally, use Pandas to save the

data as a CSV file.

### (3) Scraping results.

In the end, we captured a total of 1235 horses in Hong Kong, along with their recent 3

seasons records:

A	B	C	D	E	F	G	H	I	J	K	L	M	N	O	P	Q	R					
	马匹名称	生地	馬色	性别	進口類別	今季獎金*	總獎金*	季	總出賽	勝率	塞馬日	位置	到達	進口日期	練馬师	馬主	現時評分	季初評分	父系	母系	外祖父	同父系馬
1	安通	愛爾蘭 / 條 / 關	自購馬	0	\$13,993	\$77,2-2-27	2	香港(24/01/06/10/2022智晉賢	王晉賢	104	107	Churchill	Enrol	Pivotal	美麗之光女							
2	安鵠	美國 / 4 級 / 闊	自購馬	\$3,975	\$0	\$11,664	4-4-1-1-9	2	香港(21/01/06/10/2022蔡約翰	王賢駿	104	85	Lemon	DraAlluvial	Danehill	沒有						
3	安鶴	澳洲 / 4 級 / 闊	自購新馬	\$4,228	\$0	\$8,012	60-6-2-0-8	2	香港(19/08/05/2022蔡約翰	李國能	98	80	Rich	Emu	Could Be General	N玲瓏寶美						
4	氣勢	澳洲 / 4 級 / 闊	自購新馬	\$1,041	\$0	\$4,819	20-4-2-5-14	1	香港(07/11/13/09/2022智晉	曾水達	91	95	Street	BoKellys	OcDansili	K17級愛丸						
5	增有	澳洲 / 4 級 / 闊	自購新馬	\$322	\$0	\$9,468	80-3-3-1-26	1	香港(22/10/08/05/2022大衛希斯	香港足	83	76	Rubick	Emirates	Dubawi	星河勇士K						
6	球星	澳洲 / 4 級 / 闊	自購新馬	\$1,041	\$0	\$4,819	20-4-2-5-14	1	香港(22/10/08/05/2022大衛希斯	香港足	83	76	Ali	Tee	Hilbaran	Frosted						
7	瑪瑞	澳洲 / 6 級 / 闊	自購馬	\$322	\$0	\$9,468	80-3-3-1-26	1	從化(13/12/09/12/2022羅富全	Cheng	92	85	Emu	Fastnet	R	新鵬快寶						
8	晴王	紐西蘭 / 4 級 / 闊	自購新馬	\$140,400	\$6,294	15-5-1-3-21	1	香港(07/02/23/03/2022鄭俊偉	吳廷傑	78	78	Frosted	Darci	DanDarci	Bru							
9	羅寶	澳洲 / 4 級 / 闊	自購新馬	\$9,583	\$0	\$3,563	20-2-0-0-3	2	香港(04/01/18/03/2022羅富全	許智亨	77	64	Toronado	Liberty	W Statue Of	雄伟乐开						
10	玩笑	英國 / 3 級 / 闊	自購馬	\$0	\$0	0	0-0-0-0	0	香港(04/04/04/08/2022蔡約翰	蕭劍平	76	76	Siyouni	Illaunmor	Shamardal	縱橫大進削						
11	米奇	愛爾蘭 / 4 級 / 闊	自購馬	\$0	\$0	0	0-0-0-0	0	香港(24/01/07/2022伍鵬志	梁麟炳	73	73	Soldier's	Parle	Moi Giant's C	沒有						
12	優才	澳洲 / 4 級 / 闊	自購馬	\$0	\$0	0	0-0-0-1	1	香港(24/01/24/06/2022葉楚航	蕭劍平	72	70	Real	Stee	Northern	幸運震撼						
13	將義	紐西蘭 / 4 級 / 闊	自購馬	\$626	\$0	\$998	100	0-1-1-4	1	香港(02/02/21/03/2022鄭俊偉	七俠五聖	72	70	U S Navy	Pinstripe	Pins	時時如意勝					
14	十四	愛爾蘭 / 4 級 / 闊	自購馬	\$0	\$0	0	0-0-0-0	0	香港(02/02/21/03/2022伍鵬志	誠信園	72	72	Justify	Hence	Galliano	風雨起時天						
15	馬力	澳洲 / 4 級 / 闊	自購馬	\$2,541	\$0	\$2,541	60-1-0-0-1	1	香港(18/01/18/03/2022羅富全	利子厚	71	65	Harry	Ang	Romneya	Red Ransom						
16	祝福	澳洲 / 4 級 / 闊	自購新馬	\$11,910	\$0	\$2,197	35-2-1-1-5	1	香港(19/11/19/12/2022羅富全	羅富全	71	65	Flying	Ar Set The T	Reset	美麗毛搗炮						
17	豐盈	澳洲 / 4 級 / 闊	自購馬	\$653	\$0	\$653	700	0-1-0-5	2	香港(15/01/01/03/2022羅康祐	多多少少	71	69	National	Choice	W Choisis	星運燭					
18	快路	愛爾蘭 / 4 級 / 闊	自購馬	\$374	\$0	\$3,718	32-2-1-0-13	2	從化(11/11/05/11/2022告東尼	雷永光	70	69	Galileo	G Providence	Acclamati	沒有						
19	飲杯	澳洲 / 4 級 / 闊	自購新馬	\$1,324	\$0	\$225	40-2-1-1-8	1	香港(26/10/17/11/2022大衛希斯	雷志祥	69	59	Fastnet	R My	Sabel Savabeel	星價領創						
20	伊臣	澳洲 / 6 級 / 闊	自購新馬	\$213	\$0	\$7,505	75-5-3-1-28	1	香港(26/14/07/2022鮑本輝	李伊潤	69	69	Written	T Prete	A P Choisir	皆益善住						
21	增勁	英國 / 4 級 / 闊	自購馬	\$0	\$0	0	0-0-0-1	1	香港(29/01/14/10/2022智晉賢	曾水達	69	69	Lope	de V	Suffused Champs	El 悅悠乾坤						
22	俊才	愛爾蘭 / 4 級 / 闊	自購馬	\$0	\$0	0	0-0-0-0	0	香港(06/06/06/07/2022鮑本輝	梁麟炳	68	68	Phoenix	O Alfee	Kentucky	祥麟將軍						
23	銀進	澳洲 / 6 級 / 闊	自購新馬	\$1,305	\$0	\$5,527	80-4-4-4-31	1	香港(09/10/9/11/2022方喬柏	陳永光	67	60	Hinchinbr	Morinda	Desert	Ki 沒有						
24	奔放	愛爾蘭 / 4 級 / 闊	自購馬	\$0	\$0	0	0-0-0-0	0	從化(21/11/07/09/2022告東尼	蕭劍平	67	67	Muhhaar	Bahjee	Pivotal	沒有						
25	堅闢	紐西蘭 / 4 級 / 闊	自購新馬	\$1,454	\$0	\$2,554	20-2-2-2-13	1	香港(11/01/29/06/2022羅富全	張頤信	66	57	Var Decre	Udditaga D' Cash	穿甲戰	堅闢						
26	增強	澳洲 / 4 級 / 闊	自購新馬	\$1,444	\$0	\$2,205	45-2-2-3-7	1	從化(08/12/08/08/2022智晉賢	曾水達	66	52	Capitalis	Sebrina	Sebrina	吉祐灰福						
27	快搏	澳洲 / 7 級 / 闊	自購馬	\$65,100	\$0	\$4,865	17-3-2-4-37	1	從化(11/10/10/2022鮑本輝	達建書	66	70	Written	T Best Yet	Economist	皆益善伊巨						
28	佳登	愛爾蘭 / 4 級 / 闊	自購馬	\$0	\$0	0	0-0-0-0	0	從化(18/10/07/09/2022鮑本輝	晴朗的	66	66	Kodi	Bear	Malilla	Red Club						
29	福星	澳洲 / 5 級 / 闊	自購新馬	\$0	\$0	\$1,769	10-2-0-2-9	0	從化(11/10/7/11/2022鮑本輝	鄭永安	66	66	Star Turn	Madame	Be Publishin	美麗緣分						

A	B	C	D	E	F	G	H	I	J	K	L	M	N	O	P	Q	R	S	
1	马匹名称	场次	名次	日期	赛道	途程	场地状况	赛事班次	档位	评分	练马师	骑师	头马距离	独赢赔率	实际负磅	沿途走位	完成时间	排位体重	配备
2	安逸	153	08	03/11/24	沙田草地	1200	好/快	G3	6	105	智臂	霍宏毅	5-1/2	183	116	12.12.1.46.97	1161	B	
3	安逸	117	12	20/10/24	沙田草地	1200	好/快	G2	6	107	智臂	班德祖	5-1/2	183	116	12.13.12.1.08.79	1168	B2/TT-	
4	安逸	719	07	02/06/24	沙田草地	1600	好/快	G3	4	111	文家良	田泰安	13-3/4	26	127	4.5 6.7	1.38.15	1115	TT
5	安逸	566	13	07/04/24	沙田草地	1600	好/快	G2	4	113	文家良	田泰安	17-1/4	10	123	7.6 10.11.38.79	1125	TT	
6	安逸	452	07	25/02/24	沙田草地	2000	好	G1	14	114	文家良	田泰安	5-1/2	16	126	5.5 6.7	1.20.21.146	1146	TT
7	安逸	395	08	04/02/24	沙田草地	1800	好	G3	12	114	文家良	梁家俊	5-1/2	15	132	2.5 5.5	1.47.89	1155	TT
8	安逸	330	08	10/01/24	跑馬地草	1800	好	G3	12	114	文家良	梁家俊	4-1/4	11	135	7.8 8.7	1.48.99	1164	TT
9	安逸	238	08	10/12/23	跑馬地草	1600	好	G1	14	114	文家良	梁家俊	4-3/4	24	126	6.5 7.8	1.34.88	1164	TT
10	安逸	144	01	05/11/23	沙田草地	1800	好	G3	6	110	文家良	梁家俊	1-1/4	12	128	3.4 4.4	1.47.40	1164	TT
11	安逸	688	02	15/10/23	沙田草地	1600	好	G2	4	105	文家良	梁家俊	1/2	14	117	3.5 4.2	1.34.28	1156	TT
12	安逸	039	06	24/09/23	沙田草地	1400	好	G3	6	105	文家良	梁家俊	10-1/2	4.4	123	5.5 4.6	1.28.59	1153	TT
13	安逸	834	01	16/07/23	沙田草地	1600	好	G1	1	100	文家良	梁家俊	4-3/4	9	128	9.5 5.1	1.34.93	1165	TT
14	安逸	749	01	14/06/23	跑馬地草	1650	好	G2	10	89	文家良	梁家俊	4-3/4	3.2	124	5.4 1.1	1.38.87	1153	TT
15	安逸	576	01	12/04/23	跑馬地草	1650	好	G2	6	83	文家良	梁家俊	5.9	120	4.4 5.1	1.39.08	1138	TT	
16	安逸	509	09	19/03/23	沙田草地	2000	好	G7	6	83	文家良	薛兆輝	5-3/4	24	126	2.2 2.2	2.03.37	1152	B-H-/TT
17	安逸	453	12	26/02/23	沙田草地	1800	好/快	G7	6	83	文家良	薛兆輝	5-1/2	28	126	7.6 6.11.11.47.62	1152	B/H1-/TT	
18	安逸	406	02	08/02/23	跑馬地草	1650	好	G3	5	81	文家良	薛兆輝	2.2	132	2.2 2.2	1.39.62	1153	B/TT	
19	安逸	328	01	11/01/23	跑馬地草	1650	好	G3	12	74	文家良	薛兆輝	4.2	131	2.2 2.1	1.39.53	1137	B/TT	
20	安逸	250	01	14/12/22	跑馬地草	1650	好	G3	6	67	文家良	薛兆輝	2-1/4	3.5	126	2.5 4.1	1.38.13	1129	B1/TT
21	安逸	166	03	12/11/22	沙田全天	1650	好	G3	6	68	文家良	潘明辉	1	211	2.2 2.3	1.38.22	1133	TT	
22	安逸	116	05	23/10/22	沙田草地	1600	好/快	G3	11	69	文家良	潘明辉	2	5.8	124	2.2 2.5	1.33.80	1136	TT
23	安逸	044	06	25/09/22	沙田草地	1600	好/快	G3	2	69	文家良	潘明辉	4-3/4	4.2	124	3.5 5.6	1.35.23	1143	TT
24	安逸	010	03	11/09/22	沙田草地	1400	好	G3	12	69	文家良	潘明辉	1-1/4	8.4	122	4.2 2.3	1.22.58	1140	TT
25	安逸	153	01	03/11/24	沙田草地	1800	好/快	G3	13	86	莱约翰	田泰安	3/4	11	115	3.4 5.5	1.46.41	1081	--
26	安逸	095	04	13/10/24	沙田草地	2000	好	G3	5	85	莱约翰	潘领	1/2	2.3	135	3.4 4.3	2.01.52	1072	--
27	安逸	528	12	24/03/24	沙田草地	2000	好/快	G7	12	85	莱约翰	莫雅	1-6/2	6.8	126	5.5 5.1	11.2.02.49	1053	--
28	安逸	471	03	03/03/24	沙田草地	1800	好	G7	12	85	莱约翰	莫雅	3-3/4	7.9	126	4.3 2.4	1.48.17	1069	--
29	安逸	358	02	21/01/24	沙田草地	2000	好	G3	4	83	莱约翰	田泰安	2.9	131	3.4 3.3	2.02.60	1062	--	
30	安逸	275	01	23/12/23	沙田草地	2000	好	G3	2	76	莱约翰	田泰安	3/4	2.7	115	5.4 4.3	2.01.77	1067	--
31	安逸	237	08	10/12/23	沙田草地	1800	好	G3	2	76	莱约翰	麦道明	3-1/4	2.7	133	5.5 6.9	1.49.40	1047	--
32	安逸	181	01	19/11/23	沙田草地	2000	好/快	G3	12	70	莱约翰	艾道臣	1-1/4	8.8	126	6.5 5.6	2.00.93	1065	--
33	安逸	087	01	15/10/23	沙田草地	1800	好/快	G3	11	64	莱约翰	潘领	1/2	7.8	123	6.5 5.2	1.48.03	1063	--
34	策势	711	02	09/11/24	沙田草地	1800	好	G3	5	95	莱约翰	潘领	2.2	132	2.5 4.3	0.56.39	1211	--	
35	策势	102	03	13/10/24	沙田草地	2000	好/快	G7	12	88	莱约翰	潘领	1.9	124	2.2	1.08.29	1218	--	
36	策势	045	01	22/09/24	沙田草地	1200	黏	G3	5	80	莱约翰	潘领	1/2	5.5	134	1.1	1.09.65	1202	--
37	策势	353	01	21/01/24	沙田草地	1200	好	G3	2	73	莱约翰	麦道明	2.9	129	2.2	1.08.93	1194	--	
38	策势	305	01	01/01/24	沙田草地	1200	好	G3	2	66	莱约翰	麦道明	5.1	122	3.3	1.10.12	1201	--	
39	策势	260	02	17/12/23	沙田草地	1200	好	G3	9	65	莱约翰	麦道明	3-1/4	5.9	123	2.1	1.09.60	1195	--
40	策势	196	01	26/11/23	沙田草地	1800	好	G3	4	58	莱约翰	麦道明	3/4	6.6	133	2.1	0.57.39	1199	--
41	策势	140	01	05/11/23	沙田草地	1800	好	G3	4	52	莱约翰	田泰安	1-1/4	10	123	4.2	0.56.29	1195	--
42	增有	722	09	09/11/24	沙田草地	1400	好/快	G3	2	63	智臂	潘领	3-3/4	24	123	12.11.10.20.93	1070	B/H/XB	
43	增有	102	09	13/10/24	沙田草地	1200	好/快	G3	2	65	智臂	布文	5-3/4	33	131	12.12.9	1.09.21	1068	B/H/XB
44	增有	628	04	14/07/24	沙田草地	1600	好	G1	6	66	智臂	田泰安	3-3/4	8.8	126	10.10.9	1.34.55	1068	B/H/XB
45	增有	787	05	01/07/24	沙田草地	1400	好	G3	4	67	智臂	布文	3-1/4	8.7	135	7.6 5.5	1.21.23	1064	B/H/XB
46	增有	719	02	02/06/24	沙田草地	1600	好/快	G3	2	67	智臂	巴度	3	8.2	115	3.4 4.5	1.34.43	1062	B/H/XB
47	增有	677	02	19/05/24	沙田草地	1400	好	G3	1	67	智臂	潘领	2.9	122	6.5 6.2	1.21.59	1061	B/H/XB	
48	增有	625	04	28/04/24	沙田草地	1400	黏	G3	10	97	智臂	布文	1-1/4	7.6	132	13.12.11.1.22.55	1058	B/H/XB	

## W2-Majing

More comprehensive race data on horse racing is available on the Majing website, which is the reason for choosing to crawl the data from the second website.

馬經																			
主頁		新聞專欄		賽事資料		統計資料		綜合貼士		賠率賽果		晨操分析		其他		資料搜尋		轉播賽	
赛事列表	日期	场	次	赛	程	前	期	地	公	配	球	日	期	假	休	档	期	间	三甲席
赛事列表	12/11/24	258	3	80-60	夜	BC 1650	好	好地	67	HT	124	1109	渣打	7	222	1	1	14048	▲ 头生仔 麥兜
赛事列表	12/11/24	257	4	60-64	夜	BC 1200	好	好地	57	BT	132	108	又	4	-4.4	1	1	11011	▲ 斯特波
赛事列表	12/11/24	256	3	80-60	夜	BC 1000	好	好地	74	XT	134	1106	丽	6	-6.6	1	1/2	0.5711	▲ 年少好
赛事列表	12/11/24	255	4	60-64	夜	BC 1650	好	好地	45	B2T	119	1077	丽	7	-6.6	1	1/2	0.5711	▲ 英雄豪傑
赛事列表	12/11/24	254	3	80-60	夜	BC 1200	好	好地	57	BT	124	1052	又	12	-4.4	1	1	12.77	▲ 天生奇才
赛事列表	12/11/24	253	4	60-64	夜	BC 1200	好	好地	45	SW	120	1111	清	2	-9.0	1	1	12.078	▲ 普通人
赛事列表	12/11/24	252	5	60-64	夜	BC 1200	好	好地	52	T	127	1083	清	2	-3.4	1	1	11.007	▲ 莫近身
赛事列表	12/11/24	251	5	40-60	夜	BC 1000	好	好地	35	E	132	1109	田	11	-10.10	1	1	0.5765	▲ 有你我
赛事列表	12/08/24	249	2	105-80	日	田 A 1400	好	好地	86	118	1112	丽	10	7.88	1	1	12.129	▲ 律敦特	
赛事列表	12/08/24	248	3	80-60	日	田 A 1200	好	好地	133	T	126	1103	香	1	-4.3	1	1/2	2.0051	▲ 白田健
赛事列表	12/08/24	247	4	60-64	日	田 A 1000	好	好地	127	SW	126	1233	丽	5	-2.2</				

**URL:** <https://racing.on.cc/cgi-bin/srh/search/search.cgi> / (First website)

<https://racing.on.cc/cgi-bin/srh/search/onlinesrh.cgi> (Final website)

Web Page Structure of Majing website: HTML

The inspecting the element windows are as follows:

```
<html> <scroll>
  > <head> </head>
  > <body bgcolor="#FFFFFF" leftmargin="0" topmargin="0" marginwidth="0" marginheight="0" onload="init_mainNav(); style>
    > <div id="menu"> </div>
    > <div class="wrapHeader" id="wrapHeader" style> </div>
    > <table id="maintable" width="970" border="0" cellpadding="0" cellspacing="0" class="scriptloaded">
      > <tbody>
        > <tr> </tr>
        > <tr style="display: none;"> </tr>
        > <tr> </tr>
        > <tr> </tr>
        > <tr> </tr>
        > <tr>
          > <td colspan="2">
            > <table width="970" border="0" cellpadding="0" cellspacing="0">
              > <tbody>
                > <tr>
                  > <td> </td>
                  > <td valign="top">
                    <!-- start of left content -->
                    > <table width="650" border="0" cellpadding="0" cellspacing="0" class="ctable">
                      > <tbody>
                        > <tr> </tr>
                        > <tr> </tr>
                        > <tr> </tr>
```

## Data Scraping Process.

### (1) Page Loading and Element Location

Since it is impossible to directly obtain the content of the target web page, we use the

Selenium to crawl the web page data.

```
1  from selenium import webdriver
2  from selenium.webdriver.common.by import By
3  from selenium.webdriver.chrome.service import Service
4  from selenium.webdriver.chrome.options import Options
5  from selenium.webdriver.support.ui import WebDriverWait
6  from selenium.webdriver.support import expected_conditions as EC
```

### (2) Form Submission

Through the Elements, find the name of the start year, the end year, and the submit button

on the entry website. The script selects start and end years from dropdowns and submits the form by locating and clicking the submit button.

Accessed pages

### (3) Navigation to Result Page

After submitting, it will be redirected to a second category page as follows:

冠	258 次
亞	260 次
季	256 次
負	2312 次

After selecting and clicking on the specific links, it will be redirected to the target page.

## (4) Data Extraction

The script locates a target table on the target page with Xpath and extracts rows and

cells. Then, store the data in a list, and convert it into a Pandas DataFrame.

```
# 获取表格内容
try:
    # 使用提供的XPath定位表格
    table = driver.find_element(By.XPATH, '//*[@id="maintable"]/tbody/tr[5]/td/table/tbody/tr[2]/table/tbody/tr[4]/td[2]/table/tbody/tr/td/table/tbody/tr')
    rows = table.find_elements(By.TAG_NAME, "tr") # 获取所有行
```

## (5) Data Storage

The DataFrame is saved as a CSV file, ensuring proper encoding for international

characters.

```
# 检查是否获取到数据
if table_data:
    # 将数据转换为DataFrame
    df = pd.DataFrame(table_data)

    # 保存为CSV文件，使用 utf-8-sig 编码，确保中文字符在Excel中正常显示
    df.to_csv("racing_data2.csv", index=False, encoding="utf-8-sig")

    print("数据已保存为racing_data.csv")
```

## (6) Scraping results

Through the above process, we have collected nearly 30,000 pieces of race data from 4

seasons (2021-2024) of Hong Kong horse racing.

	A	B	C	D	E	F	G	H	I	J	K	L	M	N	O	P	Q	R	S	T	U	V	21
1	馬匹	0	1	2	3	4	5	6	7	8	9	10	11	12	13	14	15	時間	三甲馬	騎師	賠		
2	日期	場	班次	路程	地	分	配備	磅	排位	騎	練	槽	走位	名	距離	距離	時間	時間	三甲馬	騎師	賠		
3	競駿先鋒	11/17/24	192	3 80-60	田日 B 1200	好地	70 T	122	1157 鐘	東	4-11	1	1月2日 10:08:81	▲	樂勝天下	健康快車	10.881	17					
4	銀狼奔騰	11/17/24	191	3 80-60	田日 B 1400	好地	64 T1	122	1121 鐘	游	14-34	1	1月2日 12:16:66	▲	同慶美麗	玩笑	12.166	70					
5	浪漫勇士	11/17/24	190	G2	田日 B 2000	好地	133 T	128	1174 莎	沈	16-52	1	4/14/15:59:70	▲	嘉惠傳承	勸斗雲	15.970	11					
6	波濤氣勢	11/17/24	189	G2	田日 B 1600	好地	125 SW	128	1239 步	娘	9-21	1	3月4日 13:22:82	▲	禦勝輝煌	銀之霸河	13.282	3.5					
7	嘉應高飛	11/17/24	188	G2	田日 B 1200	好地	127	123	1120 潘	希	10-33	1	3/14/1 07:43	▲	熾愛心	驕陽朝曉	10.743	11					
8	愛心神鵠	11/17/24	187	3 80-60	田日 B 2000	好地	70	123	1090 明	9-10 107	1 頸	20146	▲	共享富裕	金陽飛羽	20.146	5.3						
9	穩步生風	11/17/24	186	3 80-60	田日 B 1200	好地	69 BT	124	1173 潘	葉	13-84	1	4/12/0 15:55:83	▲	英雄豪邁	金陽飛羽	0.083	18					
10	伯求開心	11/17/24	185	4 80-60	田日 B 1400	好地	55 T	123	1128 安	葉	25-65	1	2/13/0 16:53:63	▲	開心樂開	余達有盈	2.163	2					
11	順利得勝	11/17/24	184	4 80-60	田日 B 1200	好地	52 T1	128	1101 豪	葉	4-42	1	2/10/0 14:04:94	▲	順利得勝	余達有盈	1.094	40					
12	豐意龍	11/17/24	183	4 80-60	田日 B 1400	好地	50 T	134	1203 潘	賢	12-29 8	1 頸	1-21.84	▲	超風速	勝達贏	21.84	16					
13	珠得金	11/17/24	182	4 80-60	田日 B 1200	好地	43 T	112	1061 挑	伍	2-11	1	3月4日 10:09:03	▲	自己話事	勒德倚儂	10.093	7.4					
14	路進飛	11/13/24	181	4 80-60	谷夜 BC 1200	好地	52	127	1149 潘	1	1-35	1 頸	10.06	▲	醒目先鋒	競勝天下	11.006	2.4					
15	炳炳有希望	11/13/24	180	4 80-60	谷夜 BC 1200	好地	52 B-T	128	1145 豪	廖	1-76	1 頸	11.008	▲	真心傳奇	苗霸霸王	11.008	4.6					
16	知道長勝	11/13/24	179	4 80-60	谷夜 BC 1000	好地	47 T	123	1134 禧	沈	5-98	1 頸	0.5723	▲	精益求精	光裕大師	0.5723	3.7					
17	天天得勝	11/13/24	178	4 80-60	谷夜 BC 1650	好地	51 T	126	1085 澤	葉	67-75	1	1/3/1 13:39:90	▲	歡喜福星	超勁寶寶	13.990	4.1					
18	勝意	11/13/24	177	4 80-60	谷夜 BC 1650	好地	51 BT	128	1115 潘	黎	72-54	1	1/14/1 13:39:73	▲	風起雲湧	爆熱	13.973	2.8					
19	波誠駒	11/13/24	176	5 40-60	谷夜 BC 2200	好地	36 T	130	1176 紳	東	31-11	1 頸	21.678	▲	神威金剛	大登基	21.678	3.8					
20	追遊武將	2011/9/24	175	3 80-60	田日 A 1600	好快	76 SW	131	1192 柒	東	42-22	1	3月4日 13:33:89	▲	雙盛	久久為尊	13.389	2.8					
21	羅威	2011/9/24	174	3 80-60	田日 A 1600	好快	70 XT	125	1188 楊	羅	5-44 11	1 頸	108553	▲	河奸運	風雲武士	1.055	2.5					
22	羅城造將	2011/9/24	173	4 80-60	田日 A 1400	好快	45 B	124	1100 蔡	羅	5-44 11	1 頸	1-14.43	▲	保和傳奇	大登基	12.143	10					
23	堅威	2011/9/24	172	3 105-80	田日 A 1000	好快	94	123	1199 明	蔡	25-75	1	1/14/1 13:30:34	▲	堅威大將	飛騰轉翔	12.034	6.8					
24	東街控制	2011/9/24	171	2 100-80	田日 A 1000	好快	90 T	127	1083 紹	平	5-31	1 短頸	0.5636	▲	策勢	我為您	0.5636	4					
25	旺財	2011/9/24	170	4 60-40	泥日 A 1650	好泥	50 HT	127	1108 澤	伍	45-51	1	2/13/0 13:38:37	▲	速達欢呼	天威	1.3837	5.6					
26	富喜來	2011/9/24	169	4 60-40	田日 A 2000	好泥	47 BT	124	1091 潘	沈	87-76	1 短頸	2.0186	▲	慶喜家	盈益善	2.0186	2.3					
27	手銳之星	2011/9/24	168	4 60-40	田日 A 1200	好地	48 H	124	1172 田	莫	10-33	1	1月2日 10:09:54	▲	馬鳳凰	祥龍駒	1.0954	2.9					
28	喜應喬岳	2011/9/24	167	4 60-40	田日 A 1200	好地	52	128	1079 明	算	7-34	1 短頸	1.0971	▲	至黃高超	添輝焯	1.0971	12					

The raw dataset

## Difficulties and Solutions

Although the Majing website has no obvious anti-crawl mechanism, we still encountered a few minor problems in crawling the data.

### (1) Simulator jumping.

**Difficulty:** During the crawling process, we can't reach the target page directly because copying the URL of the target page directly will only result in a blank page without data. Therefore, in order to ensure that the data is visible, we can only access the target page by searching for the jumps.

The screenshot shows the Majing website's search interface. At the top, there is a navigation bar with links like '主頁', '新聞專欄', '賽事資料', '統計資料', '綜合貼士', '賠率賽果', '晨報分析', '其他', '資料搜尋', and '轉播賽'. Below the navigation bar is a search bar with the placeholder '互動搜尋' and a '列印' button. The main content area contains several filter categories: '新聞專欄', '賽事資料', '統計資料', '綜合貼士', '賠率賽果', '晨報分析', and '資料搜尋'. Each category has a list of items, such as '即時快訊', '馬匹賽道成績', '賽事備忘', '張基', '卡洛斯', '匡公', '金鉅', '騎練成績表', '騎師/練馬師', '騎練行程成績', '五程績', '積分榜 / 詳細版 /', '強化版', '騎練半Q', '騎練半串Q', '名家存三T貼士', '名家存兩場貼士', '馬房綜合貼士', '馬房次級料', '三T神算', '騎馬該串Q', '陪率逐勝', '陪率逐贏 / 位置 /', '組合獨贏 / 單贏 / 位置Q /', '二重彩 / 單T / 百連理 /', '孖寶 / 駕牢比較', '即時追彩', '詳報逐彩', and '大跑頭', '內膽', '每日競操', '快跑 / 跑步 / 溶水', '出賽馬群探', '出賽馬分數晨報', and '騎練結束'. There is also a '馬匹資訊' section with links like '飛起邊隻先', '互動搜尋', '馬匹資料', '騎師資料', and '練馬師資料'.

Pages accessed by copying URL

**Solution:** Use selenium to simulate a human search from the search portal and reach our target page through two webpage jumps to get our target data.

```

# Step 1: 打开搜索页面
driver.get("https://racing.on.cc/cgi-bin/srh/search/search.cgi")
time.sleep(2) # 等待页面加载

# Step 2: 找到“赛事年度（开始年份）”选择框并选择“2024”
start_year_select = WebDriverWait(driver, 10).until(
    EC.presence_of_element_located((By.NAME, "fryear")) # 根据实际的字段名称来获取
)
start_year_select = Select(start_year_select)
start_year_select.select_by_visible_text("2024") # 选择2024

# Step 3: 找到“赛事年度（结束年份）”选择框并选择“2025”
end_year_select = WebDriverWait(driver, 10).until(
    EC.presence_of_element_located((By.NAME, "toyear")) # 根据实际的字段名称来获取
)
end_year_select = Select(end_year_select)
end_year_select.select_by_visible_text("2025") # 选择2025

# Step 4: 找到提交按钮，先使用滚动和显式等待确保按钮可点击
submit_button = WebDriverWait(driver, 10).until(
    EC.element_to_be_clickable((By.NAME, "submitButton"))
)

# Step 5: 等待页面跳转到结果页面
time.sleep(3) # 等待页面跳转

# Step 6: 在结果页面中找到链接文本为“192”的元素并点击
try:
    # 根据链接文本来找到并点击
    link = WebDriverWait(driver, 10).until(
        EC.presence_of_element_located((By.LINK_TEXT, "192")) # 查找链接文本为“192”的元素
    )
    link.click() # 点击链接

    # 等待页面跳转
    time.sleep(3) # 等待页面加载到目标页面

```

## (2) Ad blocking.

**Difficulty:** In the search page, affected by pop-up ads, unable to successfully click to the submit button.



**Solution:** Use code to make sure the button is visible.

```

# Step 4: 找到提交按钮，先使用滚动和显式等待确保按钮可点击
submit_button = WebDriverWait(driver, 10).until(
    EC.element_to_be_clickable((By.NAME, "submitButton"))
)

# 确保按钮可见
driver.execute_script("arguments[0].scrollIntoView();", submit_button)

```

## Data Processing

### **MaJing**

(1) Integrate four datasets from the website.

Because we have obtained the data of "冠", "亚", "季" and "亚" from MaJing, we need to

integrate the data.

```
import pandas as pd

# File paths
file_paths = [
    '/Users/Keanu/Desktop/racing_data/csv/冠.csv',
    '/Users/Keanu/Desktop/racing_data/csv/亚.csv',
    '/Users/Keanu/Desktop/racing_data/csv/季.csv',
    '/Users/Keanu/Desktop/racing_data/csv/负.csv'
]

# Load each CSV file into a DataFrame and concatenate them
combined_all_df = pd.concat([pd.read_csv(file) for file in file_paths], axis=0, ignore_index=True).drop_duplicates()

# Set output file path and export the combined DataFrame to CSV
output_file_path = '/Users/Keanu/Desktop/all.csv'
combined_all_df.to_csv(output_file_path, index=False, encoding='utf-8')

print(f"The combined table has been saved to: '{output_file_path}'")
```

(2) Sort by Date(日期) & Event(场) In Descending Order.

The original dates were mixed up and needed to be reorganized.

```
import pandas as pd

# Load the CSV file
file_path = '/Users/Keanu/Desktop/all.csv'
data = pd.read_csv(file_path)

# Convert '日期' column to datetime
data['日期'] = pd.to_datetime(data['日期'], format='%m/%d/%y', errors='coerce')

# Drop rows with invalid dates
data = data.dropna(subset=['日期'])

# Sort by '日期' in descending order
df_sorted = data.sort_values(by='日期', ascending=False)

# Sort by '場' in descending order while preserving the order of '日期'
df_sorted = df_sorted.sort_values(by=['日期', '場'], ascending=[False, False])

# Reset the index
df_sorted.reset_index(drop=True, inplace=True)

# Save the sorted data to a new CSV file
df_sorted.to_csv('/Users/Keanu/Desktop/all_sorted.csv', index=False)

print("Data has been sorted and saved as 'all_sorted.csv'")
```

(3) Check Missing Values.

```

import pandas as pd

# Load the CSV file
file_path = '/Users/Keanu/Desktop/all_sorted.csv'
data = pd.read_csv(file_path)

# Check for missing values in each column
missing_values = data.isnull().sum()
print("Missing values in each column:\n", missing_values)

# Extract rows with missing values
rows_with_missing = data[data.isnull().any(axis=1)]

# Save the rows with missing values to a new CSV file
rows_with_missing.to_csv('/Users/Keanu/Desktop/rows_with_missing_values.csv', index=False)

print("Rows with missing values have been saved as 'rows_with_missing_values.csv'")

```

#### (4) Follow the Result.

```

Missing values in each column:
馬匹      0
日期      0
場        0
班次      0
路程      0
場地      10
分        0
配備      5061
磅        0
排位重    0
騎        0
練        0
檔        0
走位      0
名        0
距離      21
時間      25
冠        0
亞        0
季        0
頭馬      12
賠        0
dtype: int64
Rows with missing values have been saved as 'rows_with_missing_values.csv'

```

#### (5) Replace missing values in the “配備” column with None. According to the material

and actual competitions, it is normal to have no equipment, so the "equipment" column is processed as "None", and then drop rows with other missing values.

```

import pandas as pd

# Load the CSV file
file_path = '/Users/Keanu/Desktop/all_sorted.csv'
data = pd.read_csv(file_path)

# Replace missing values in '配備' column with 'None'
data['配備'].fillna('None', inplace=True)

# Drop rows with any other missing values
data = data.dropna()

# Remove rows where '名' column is equal to 90
data = data[data['名'] != 90]

# Check for missing values in each column
missing_values = data.isnull().sum()
print("Missing values in each column:\n", missing_values)

# Save the cleaned data to a new CSV file
data.to_csv('/Users/Keanu/Desktop/all_sorted_cleaned.csv', index=False)

print("Cleaned data has been saved as 'all_sorted_cleaned.csv'")

```

(6) Remove all other rows with missing values.

(7) Output the cleaned data.

cleaned_all_sorted																					
馬匹	日期	場	班次	路程	場地	分	配備	磅	排位重	騎	練	檔	走位	名	距離	時間	冠	亞	季	頭馬	賠
得道獵王	2024-11-20	201	3 80-60	谷夜 C3 1200	好黏	69	T	126	1196	鍾 東	10	-11	1	頭	1.10.17	▲	伊臣	勇威神駒	1.10.17	4.4	
寶賢得得	2024-11-20	201	3 80-60	谷夜 C3 1200	好黏	75	T	128	1129	弘 柏	11	-12 12	11	4 1/2	1.10.88	得道獵王	伊臣	勇威神駒	1.10.17	136.0	
金發盛世	2024-11-20	201	3 80-60	谷夜 C3 1200	好黏	66	B	126	1067	湯 呂	8	-6 8	10	3 1/2	1.10.72	得道獵王	伊臣	勇威神駒	1.10.17	36.0	
多利神駒	2024-11-20	201	3 80-60	谷夜 C3 1200	好黏	62	T	122	1034	禮 賢	7	-11 11	4	1 1/4	1.10.38	得道獵王	伊臣	勇威神駒	1.10.17	35.0	
金蓮來	2024-11-20	201	3 80-60	谷夜 C3 1200	好黏	60	XT	120	1102	森 文	12	-2 2	6	2	1.10.47	得道獵王	伊臣	勇威神駒	1.10.17	28.0	
無限美麗	2024-11-20	201	3 80-60	谷夜 C3 1200	好黏	65	None	125	1153	布 蔡	4	-9 9	5	1 3/4	1.10.43	得道獵王	伊臣	勇威神駒	1.10.17	8.6	
多利彩駒	2024-11-20	201	3 80-60	谷夜 C3 1200	好黏	69	XT	129	1085	班 廉	6	-4 5	9	3 1/2	1.10.72	得道獵王	伊臣	勇威神駒	1.10.17	6.3	
金牌活力	2024-11-20	201	3 80-60	谷夜 C3 1200	好黏	63	B	123	1147	湯 鄭	5	-10 10	7	2 3/4	1.10.61	得道獵王	伊臣	勇威神駒	1.10.17	63.0	
友愛心得	2024-11-20	201	3 80-60	谷夜 C3 1200	好黏	66	SW	126	1185	潘 希	9	-3 3	12	7 1/4	1.11.33	得道獵王	伊臣	勇威神駒	1.10.17	7.7	
駿駿皇者	2024-11-20	201	3 80-60	谷夜 C3 1200	好黏	74	T	134	1075	聲 游	3	-5 4	8	3 1/4	1.10.70	得道獵王	伊臣	勇威神駒	1.10.17	5.5	
伊臣	2024-11-20	201	3 80-60	谷夜 C3 1200	好黏	69	None	127	1145	明 姚	2	-7 6	2	頭	1.10.20	得道獵王	▲	勇威神駒	1.10.17	9.4	
勇威神駒	2024-11-20	201	3 80-60	谷夜 C3 1200	好黏	66	T	126	1137	麥 沈	1	-8 7	3	1	1.10.34	得道獵王	伊臣	▲	1.10.17	5.5	
美麗奔馳	2024-11-20	200	2 100-80	谷夜 C3 1200	好黏	90	T	126	1096	艾 東	7	-10 9	1	頭	1.09.43	▲	傑出漢子	精算驕雪	1.09.43	26.0	
傑出漢子	2024-11-20	200	2 100-80	谷夜 C3 1200	好黏	83	B	119	1097	潘 呂	1	-6 3	2	頭	1.09.47	美麗奔馳	▲	精算驕雪	1.09.43	20.0	
我為您	2024-11-20	200	2 100-80	谷夜 C3 1200	好黏	80	None	114	1166	明 伍	6	-1 1	7	3 3/4	1.10.04	美麗奔馳	傑出漢子	精算驕雪	1.09.43	29.0	
輝煌精英	2024-11-20	200	2 100-80	谷夜 C3 1200	好黏	83	H	119	1220	梁 丁	10	-2 2	8	4	1.10.06	美麗奔馳	傑出漢子	精算驕雪	1.09.43	9.3	
占士德	2024-11-20	200	2 100-80	谷夜 C3 1200	好黏	84	T1SW	120	1111	周 容	12	-4 4	11	7	1.10.56	美麗奔馳	傑出漢子	精算驕雪	1.09.43	64.0	
華麗再勝	2024-11-20	200	2 100-80	谷夜 C3 1200	好黏	88	T	124	1160	麥 沈	8	-7 7	5	2 1/4	1.09.80	美麗奔馳	傑出漢子	精算驕雪	1.09.43	9.4	
昇淵駒	2024-11-20	200	2 100-80	谷夜 C3 1200	好黏	82	None	118	1188	班 廉	3	-8 6	4	2 1/4	1.09.78	美麗奔馳	傑出漢子	精算驕雪	1.09.43	7.4	
精算驕雪	2024-11-20	200	2 100-80	谷夜 C3 1200	好黏	88	P	124	1091	潘 蔡	4	-11 10	3	1 1/4	1.09.64	美麗奔馳	傑出漢子	▲	1.09.43	2.2	
幸運遇見	2024-11-20	200	2 100-80	谷夜 C3 1200	好黏	91	None	127	1122	禮 希	2	-5 8	12	7 3/4	1.10.66	美麗奔馳	傑出漢子	精算驕雪	1.09.43	5.9	
小霸王	2024-11-20	200	2 100-80	谷夜 C3 1200	好黏	99	BT	132	1188	鍾 東	9	-12 12	6	3 3/4	1.10.04	美麗奔馳	傑出漢子	精算驕雪	1.09.43	15.0	
魅力實駒	2024-11-20	200	2 100-80	谷夜 C3 1200	好黏	92	None	128	1197	度 蔡	5	-9 11	9	4 1/2	1.10.16	美麗奔馳	傑出漢子	精算驕雪	1.09.43	28.0	
人和家興	2024-11-20	200	2 100-80	谷夜 C3 1200	好黏	94	BT	130	1202	董 希	11	-3 5	10	5 1/2	1.10.31	美麗奔馳	傑出漢子	精算驕雪	1.09.43	45.0	
勁好運	2024-11-20	199	3 80-60	谷夜 C3 1650	好黏	62	T	120	1094	潘 雅	6	6 5 4	2	短頭	1.39.60	開心勇駒	▲	浪漫老撾	1.39.62	2.8	
浪漫老撾	2024-11-20	199	3 80-60	谷夜 C3 1650	好黏	72	SW2T	129	1152	梁 巫	2	8 8 8	3	1 1/4	1.39.80	開心勇駒	勁好運	▲	1.39.62	16.0	
好好心得	2024-11-20	199	3 80-60	谷夜 C3 1650	好黏	74	T-B	131	1080	禮 游	4	9 9 11	8	4 3/4	1.40.38	開心勇駒	勁好運	浪漫老撾	1.39.62	46.0	
飛輪步	2024-11-20	199	3 80-60	谷夜 C3 1650	好黏	74	T	131	1159	布 容	10	4 4 5	11	12 1/4	1.41.55	開心勇駒	勁好運	浪漫老撾	1.39.62	47.0	
獨步天下	2024-11-20	199	3 80-60	谷夜 C3 1650	好黏	63	None	120	1137	潘 呂	8	7 6 6	10	7 1/2	1.40.79	開心勇駒	勁好運	浪漫老撾	1.39.62	49.0	

## Hong Kong Jockey Club

(1) Data format conversion. Separate the mixed up information about total appearances

and ranking information.

```

import pandas as pd

# 读取csv文件
file_path = '/Users/Keanu/Desktop/racing_data/csv/HM_Horses.csv'
df = pd.read_csv(file_path)

# 定义一个函数来处理“场-季-场-练-比赛次数”列，并创建新的列
def process_race_counts(race_counts_str):
    if pd.isna(race_counts_str):
        if df.isna(race_counts_str): # 处理缺失值
            return 0, 0, 0, 0 # 返回默认值
        else:
            return 0, 0, 0, 0 # 处理以Nan

    try:
        parts = race_counts_str.split('-')
        if len(parts) == 4:
            wins = int(parts[0])
            seconds = int(parts[1])
            thirds = int(parts[2])
            total_races = int(parts[3].replace("N", "")) # Remove 'N' and convert to int
            return wins, seconds, thirds, total_races
        else:
            return 0, 0, 0, 0 # handle cases with incorrect format
    except ValueError: # handle cases with non-numeric data
        return 0, 0, 0, 0

# 用函数处理“场-季-场-练-比赛次数”列，并创建新的列
df[['wins', 'seconds', 'thirds', 'total_races']] = df['场-季-场-练-比赛次数'].apply(lambda x: pd.Series(process_race_counts(x)))

# 删掉原始的“场-季-场-练-比赛次数”列
df = df.drop(['场-季-场-练-比赛次数'], axis=1)

# 保存新的DataFrame到CSV文件
new_file_path = '/Users/Keanu/Desktop/processed_data.csv'
df.to_csv(new_file_path, index=False)

print(f"Processed data saved to {new_file_path}")

```

(2) Check missing values.

```

import pandas as pd

# Load the csv file
file_path = '/Users/Keanu/Desktop/processed_data.csv'
data = pd.read_csv(file_path)

# Check for missing values in each column
missing_values = data.isnull().sum()
print("Missing values in each column:\n", missing_values)

# Extract rows with missing values
rows_with_missing = data[data.isnull().any(axis=1)]

# Save the rows with missing values to a new CSV file
rows_with_missing.to_csv('/Users/Keanu/Desktop/rows_with_missing_values.csv', index=False)

print("Rows with missing values have been saved as 'rows_with_missing_values.csv'")

```

### (3) Follow the result.

```

Missing values in each column:
馬匹名稱          0
出生地 / 馬齡      0
毛色 / 性別        0
進口類別          0
今季獎金*          0
總獎金*          0
最近十個賽馬日出賽場數    0
現在位置(到達日期)  0
進口日期          0
練馬師            0
馬主              0
現時評分          15
季初評分          15
父系              0
母系              0
外祖父            0
同父系馬          0
自購馬來港前賽事片段  1235
wins              0
seconds            0
thirds             0
total_races        0
dtype: int64
Rows with missing values have been saved as 'rows_with_missing_values.csv'

```

### (4) Remove all rows with missing values.

Because the column "自購馬來港前賽事片段" appears as a video link in the original

webpage, it needs to be deleted first. Then drop rows with any other missing values.

```

import pandas as pd

# Load the csv file
file_path = '/Users/Keanu/Desktop/processed_data.csv'
data = pd.read_csv(file_path)

# Remove the specified column
data = data.drop(columns=['自購馬來港前賽事片段'])

# Drop rows with any other missing values
data = data.dropna()

# Save the cleaned data to a new CSV file
output_file_path = '/Users/Keanu/Desktop/cleaned_data.csv'
data.to_csv(output_file_path, index=False)

# Check for missing values in each column
missing_values = data.isnull().sum()
print("Missing values in each column:\n", missing_values)

```

### (5) Output the cleaned data.

## **Analysis and Results**

From the data we have acquired, we propose four questions and make analysis.

**Q1:** What is the probability of the highest rating horse running first in past seasons (2021-2024)?

**Q2:** What range of "Win Odds" is suitable for betting in past seasons (2021-2024)?

**Q3:** What is the average rating (from current rating) of each sire and dam's sire of the

horses currently in Hong Kong Jockey Club?

**Q4:** Which horses are now in the top ten in the Hong Kong Jockey Club's winning rate

(Champion & Top3) ?

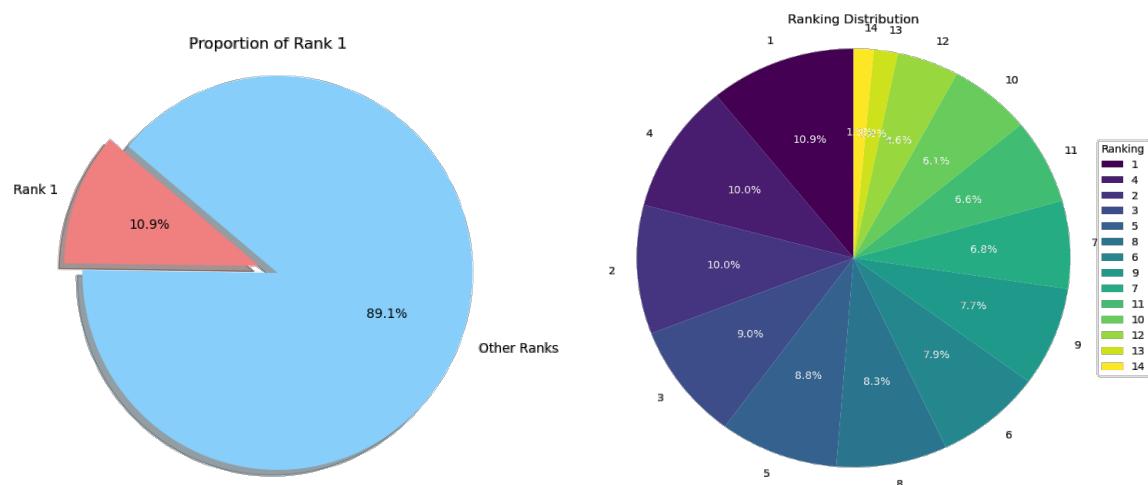
## Result 1

We want to know the performance of the highest rating horses in each race. Is it likely to

always win?

The results are as follows: The probability of champion is not very ideal, but it is normal.

Although it seems that the highest rating horse has a low chance of winning the first prize, we can still see that the horses with high ratings are generally ranked relatively high. The ratios for first to fourth place are all around 10%, and the payouts are generally given to the top four, which means that horses with higher ratings can generally meet the public's expectations of them.



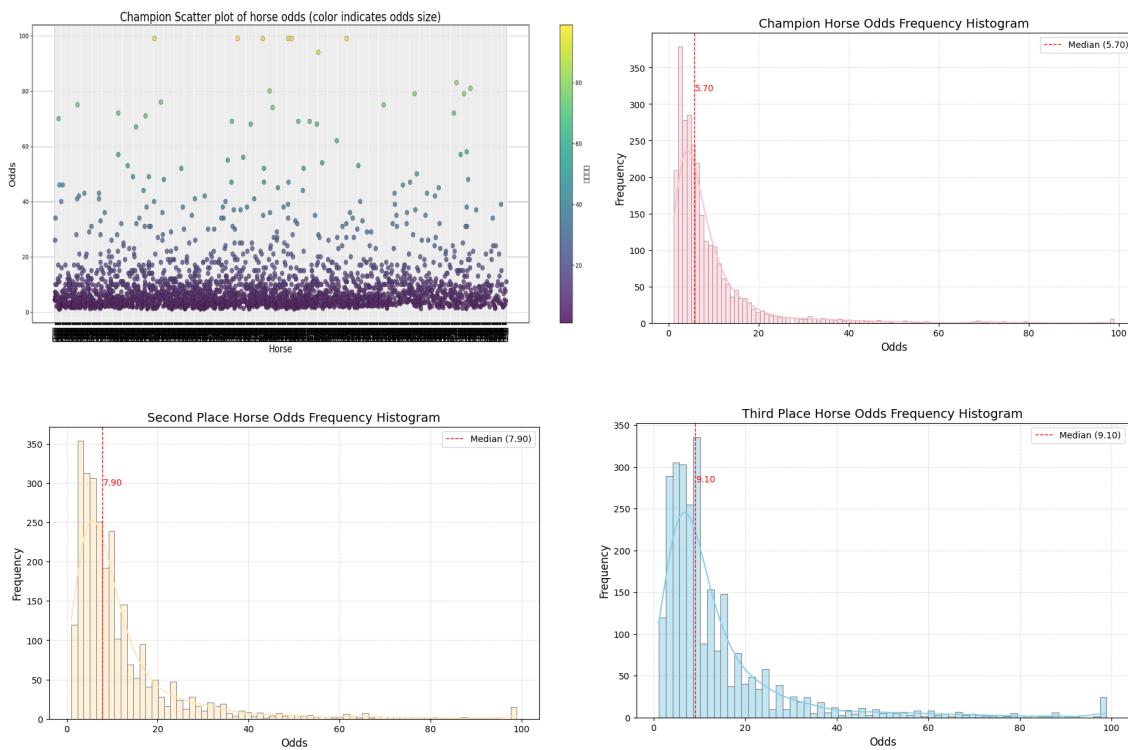
## Result 2

We first made a scatter plot of the championship odds and found that the odds were too concentrated, so we made a frequency histogram later.

Since the data distribution is skewed, we marked the median. The champion odds are concentrated at 5.70 or even lower, the second odds are concentrated at 7.90, and the third

place odds are concentrated at 9.10.

From the visualization we can see that the winning horses generally have lower odds. For example, when we reviewed the competition on December 8, we found that the "浪漫勇士", which had won six games in a row, had its win and position odds even dropped to 1.0-1.1, which was almost a non-profitable state.

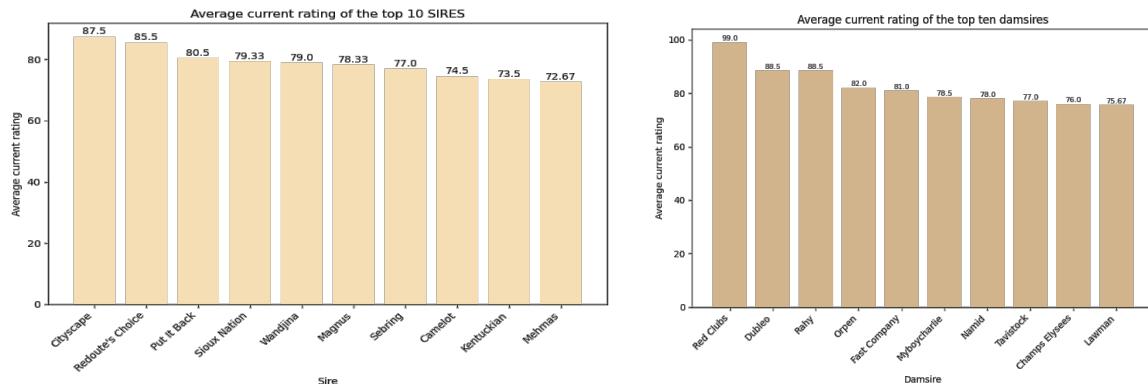


### Result 3

From the official information, we know that the sire's bloodline has a great influence on the horse's ability, so we try to compare and rank their sire's current ratings and dam's sire's current ratings. We have selected 20 pedigree horses with the highest average scores for the

time being. We can pay attention to the results of horses with the same origin and consider

combining this information when betting in the future.



## Result 4

The winning rate of a small number of times is unstable, only horses that have competed at least 10 times were selected for analysis. We initially tried to choose more than just three races, but the conditions were still not representative enough, so after comparing 3, 10, and 15 races' data, we finally chose “10” as the dividing point.

```
import pandas as pd

# 读取处理后的Excel文件
file_path = '/content/cleaned_data.xlsx'
df = pd.read_excel(file_path)

# 清理掉"total_races"大于9的数据
df_cleaned = df[df['total_races'] > 9]

# 保存清理后的数据到新的Excel文件
new_file_path = '/content/cleaned_races.xlsx'
df_cleaned.to_excel(new_file_path, index=False)
print(f"Cleaned data saved to {new_file_path}")

# 计算两种胜率, 即 冠/总次数 的比率, 以及 (冠+亚+季)/总次数 的比率
# 读取处理后的csv文件
file_path = '/content/cleaned_races.csv'
df = pd.read_csv(file_path)

# 计算两种胜率
df['win_rate'] = df['wins'] / df['total_races']
df['top3_rate'] = (df['wins'] + df['seconds'] + df['thirds']) / df['total_races']

# 打印结果或保存到新的csv文件
print(df[['wins', 'seconds', 'thirds', 'total_races', 'win_rate', 'top3_rate']])

# 保存结果到新的csv文件(可选)
new_file_path = '/content/win_rates.csv'
df.to_csv(new_file_path, index=False)
print(f"Win rates saved to {new_file_path}")
```

Because there are two betting ways, "win" and "place", we conducted a winning rate analysis on the current participating horses. That is, we listed the probability of winning the

championship and the probability of getting the top three. The results are shown in the chart.

Some horses really stand out.

Top 10 Horses by Win Rate			Top 10 Horses by Top 3 Rate		
	马匹名称	win_rate		马匹名称	top3_rate
170	浪漫勇士	0.750000	215	手機錶霸	1.000000
171	金鑽貴人	0.666667	277	營造組裝	0.909091
184	驕陽明駒	0.600000	170	浪漫勇士	0.900000
174	錶之銀河	0.538462	184	驕陽明駒	0.900000
193	合夥奔馳	0.538462	171	金鑽貴人	0.875000
176	永遠美麗	0.500000	172	加州星球	0.857143
172	加州星球	0.464286	190	美麗第一	0.846154
39	堅又威	0.416667	176	永遠美麗	0.833333
47	占士德	0.400000	2	球星	0.785714
175	維港智能	0.388889	174	錶之銀河	0.769231

## Conclusions

This study has meticulously acquired and processed Hong Kong horse racing data from reliable sources, namely the Hong Kong Jockey Club's official website and MaJing. Through a combination of data acquisition methods, including Python libraries like Requests, BeautifulSoup, Pandas, and Selenium, we successfully extracted and cleaned a significant dataset comprising the details and records of many horses. A total of nearly 30,000 pieces of data were obtained. This comprehensive effort allowed us to conduct exploratory analyses, explore horse performance, and provide insights into horse racing trends. For instance, performance of Highest Rating Horses underscores the importance of ratings as a factor in predicting race outcomes but also highlights the variability inherent in horse racing.

However, key predictive factors such as jockey experience, track conditions, and training regimes were not included, which may affect the actual reference significance. Moreover, this research relies on a lot of historical data and may not fully explain the dynamics and real-time changes in horse racing. The static nature of the past results limits their applicability to real-time strategy.

In general, our data may be suitable for novices just as a reference, it is still very limited, and there is no direct guarantee that a bet on a certain horse will win, and we need to build models to further prediction. Maybe machine learning techniques can be applied to enhance prediction accuracy.

## Codebook

### Hong Kong Jockey Club

Source: <https://racing.hkjc.com/racing/information/chinese/Horse>SelectHorse.aspx>

Data format : csv

Data collection time : 9/11/2024

**Table 1:** The basic information on the 1,235 horses in Hong Kong.

Variable Name	Description	Data Type	Example
馬匹名称	Name of the horse.	String	安遇
出生地/馬齡	Birthplace of the horse and its age in years.	String	愛爾蘭/5
毛色/性別	The coat color of the horse and its sex.	String	棕/閹
進口類別	Category under which the horse was imported.	String	自購新馬
今季獎金	Prize money earned in the current season.	Numeric	\$3,975,000(HKD)
總獎金	Total prize money earned throughout the horse's racing career.	Numeric	\$13,993,775(HKD)
冠-亞-季-總出賽次數	Number of wins, places (second), shows (third), and total number of races started.	Numeric	7-2-2-27
最近十個賽馬日出賽場數	The number of races started in the last ten race days.	Numeric	2
現在位置 (到達日期)	Current location of the horse and the date of arrival.	String	Hong Kong (24/09/2024)
進口日期	Date the horse was imported to HK.	Date	16/10/2021
練馬師	Name of the horse's trainer.	String	賀賢

馬主	Name of the horse's owner.	String	王賢訊
現時評分	Current rating of the horse.	Numeric	104
季初評分	Rating of the horse at the beginning of the season.	Numeric	107
父系	Name of the horse's sire (father).	String	Churchill
母系	Name of the horse's dam (mother).	String	Enrol
外祖父	Name of the horse's grandsire (mother's father).	String	Pivotal
同父系馬	List of other horses sired by the same sire.	String	祥勝將軍戰騎飛

**Table 2:** Recent 3 seasons race records for the above 1,235 horses.

Variable Name	Description	Data Type	Example
马匹名称	The name of the horse participating in the race	String	安遇
场次	The race meeting number	Numeric	243
名次	The finishing position of the horse in the race	Numeric	44
日期	The date of the race	Date	2024/3/24
赛道	The type of track where the race was held	String	沙田草地"A"
途程	The distance of the race (in meters)	Numeric	1200
場地狀況	The condition of the track on the race day	String	好/快
賽事班次	The type or class of the race	String	G1
档位	The starting position of the horse in the race	Numeric	2
評分	The rating score of the horse after the race	Numeric	85
练马师	The name of the trainer responsible for the horse	String	蔡約翰
騎師	The name of the jockey riding the horse	String	麥道朗
頭馬距離	The distance to the winning horse	String, Numeric	“頸位” or “3-3/4”
獨贏賠率	The odds for the horse to win the race	Numeric	2.9
實際負磅	The actual weight carried by the horse during the race (in kilograms)	Numeric	126
沿途走位	The riding position of the horse during the race	Numeric	4 6 6 7 7 4

完成時間	The time taken for the horse to finish the race	String	2.28.11
排位體重	The horse's weight before the race (in kilograms)	Numeric	1072
配備	Equipment used for the horse during the race	String	B2/TT-

Note:

It wasn't used for data analysis because the retired horse's records were not included, and

MaJing's race data was more comprehensive.

## MaJing

Data sources: Majing Website (<https://racing.on.cc/index.html>)

Sample size: 32377(Total); 2705(Championship); 2708(Runner-up); 2700(Quarter);

24264(Loss)

Data format : csv

Data collection time : 20/11/2024

Describe: This table records the horse racing data in the last four seasons in Hong Kong.

(Incomplete data for 24/25 season) In order to better analyze and organize the data, this table is categorized into 4 sub-tables for Championship(冠), Runner-up(亞), Quarter(季), and Loss (负).

**Table 3:** Four seasons' horse racing records in Hong Kong.

Variable Name	Description	Data Type	Example
馬匹	The name of the horse participating in the race	String	浪漫勇士
日期	The date of the race	Date	11/17/24
場	The race meeting number	Numeric	190
班次	The type or class of the race	String	G2
路程	The type of track where the race was held and The distance of the race (in meters)	String	田日 B 2000
場地	The condition of the track on the race day	String	好地
分	The rating score of the horse after the race	Numeric	133
配備	Equipment used for the horse during the race	String	T
磅	The actual weight carried by the horse during the race (in kilograms)	Numeric	128
排位重	The horse's weight before the race (in kilograms)	Numeric	1174
騎	The name of the jockey riding the horse	String	麥
練	The name of the trainer responsible for the horse	String	沈
檔	The starting position of the horse in the race	Numeric	1
走位	The riding position of the horse during the race	Numeric	6 5 2
名	The finishing position of the horse in the race	Numeric	1
距離	The distance to the winning horse	String	4 ¼ or 頸
時間	The time taken for the horse to finish the race	Numeric	1.59.70
冠	The name of the champion horse	String	浪漫勇士 or ▲
亞	The name of the runner-up horse	String	嘉應傳承 or ▲
季	The name of the third-place horse	String	勦斗雲 or ▲
頭馬	The time taken for the champion horse to finish the race	Numeric	1.59.70
賠	The win odds of the horse for the game	Numeric	1.1

Note:

- (1) Not all horses are equipped, so the “配備” column may be empty.
- (2) “▲” in “冠”, “亞” or “季” columns means that the horse in this line get first, second or third place.
- (3) Due to data source, trainer and jockey names are incomplete (only family name).

## References

- Gupta, M. ., & Singh, L. . (2023). Predicting Outcomes of Horse Racing using Machine Learning. International Journal on Recent and Innovation Trends in Computing and Communication, 11(9), 38–47. <https://doi.org/10.17762/ijritcc.v11i9.8119>
- Gupta, M. ., & Singh , L. . (2023). Horse Race Results Prediction Using Machine Learning Algorithms With Feature Selection. International Journal of Intelligent Systems and Applications in Engineering, 12(2s), 132–139. Retrieved from <https://ijisae.org/index.php/IJISAE/article/view/3565>