

データマイニング手法を用いた競馬研究

A Study of Horse Racing Using Data Mining Methods

金枝 隆之介[†] 華山 宣胤^{††}

Ryunosuke kanaeda,[†] Nobutane Hanayama,^{††}

^{†,††}尚美学園大学 芸術情報学部 情報表現学科

^{†,††}Department of Information Technology, Shobi University

概要: 競馬については、競走馬に関するものやレースに関するものなど入手可能なデータは膨大であるが、どのデータがレース結果の予想に有用であるかははっきりしていない。そこで、本研究では、データマイニングの手法を用いて競馬データを分析し、勝馬を予測するシステムを構築する。また、構築したシステムによる予想と実際の結果の比較検討を行う。

1. はじめに

競馬データは、競走馬に関するデータ (馬データ) とレースに関するデータ (レースデータ) に大別される。馬データは出走馬の血統、体重の変化、前レースまでの成績などの競走馬に関するものであり、レースデータは馬場の開催地、距離、出走条件、馬場状態などのレースに関するものである。そして、これらのデータに基づいて競馬専門紙やスポーツ新聞などではレース結果の予想が行われている。しかし、それらデータの項目数は膨大であり、どのデータがレース結果の予想に有用であるかははっきりしていない。実際、競馬予想の専門家によって、参考としている項目が異なり、また解釈の方法もまちまちである。つまり、“何から何がわかるか分からないデータベース”といえる。そこで本研究では、膨大な競馬データにデータマイニング手法を適用することにより、統計的な分析結果に基づく勝馬予測システムを構築し、その有用性を検証する。

2. 分析手順

提案する分析手順は下記のとおりである。

- (1) 競馬データの取得
- (2) データの標準化
- (3) 標準化されたデータのクラスター分析
- (4) クロス集計表の作成

(5) 多群判別分析による新規レースの予想

2.1. 競馬データの取得

使用するデータは、Yahoo!スポーツにより公開されている 2007 年度の GI データを使用する。その後必要のないデータの除外を行う。

図表 1. 馬データ (Excel)

馬名	性別	年齢	性別 (種別)	種別	斤量	出走回数	出走回数 (出走回数)	出走回数 (出走回数)	出走回数 (出走回数)	出走回数 (出走回数)
10月1日	A	2	1	3	55	482	-4	3	153	
10月2日	A	2	1	3	57	488	-5	14	445	
10月3日	A	2	1	3	57	441	-2	4	187	
10月4日	A	2	1	3	57	478	10	8	23	
10月5日	A	2	1	3	57	472	2	8	228	
10月6日	A	2	1	3	57	482	-2	5	23	
10月7日	A	2	1	3	57	288	2	1	18	
10月8日	A	2	1	3	57	458	-4	18	293	
10月9日	A	2	1	3	57	472	-4	7	87	
10月10日	A	2	1	3	57	488	-4	13	78	
10月11日	A	2	1	3	57	480	4	7	237	
10月12日	A	2	1	3	57	480	2	18	293	
10月13日	A	2	1	3	57	480	8	13	291	
10月14日	A	2	1	3	57	482	2	15	87	
10月15日	A	2	1	3	57	510	2	11	532	
10月16日	A	2	1	3	57	482	4	5	211	
10月17日	A	2	1	3	57	478	-4	18	1538	
10月18日	A	2	1	3	57	380	-2	13	137	

図表 2. レースデータ (Excel)

レース名	中央	中京	京都	東京	名産 (0)	種別	出走回数 (0)	出走回数 (0)	出走回数 (0)	出走回数 (0)	出走回数 (0)	出走回数 (0)
高松宮記念	1	0	0	0	0	1200	0	2	1	0	1	15
桜花賞	0	1	0	0	0	1600	2	1	0	1	0	15
皐月賞	0	0	1	0	0	2000	0	1	0	1	1	15
天皇賞 (春)	0	0	0	1	0	3200	0	2	1	1	1	15
NHKマイルカップ	0	0	0	1	0	1800	0	1	0	0	0	15
フィリッパマイル	0	0	0	1	0	1600	2	2	1	1	1	15
新緑賞 (オークス)	0	0	0	1	0	2400	2	1	0	1	1	15
東京賞 (牡)	0	0	0	0	1	2400	0	1	0	1	1	15
宝塚記念 (牡)	0	0	0	1	0	1600	0	2	1	1	1	15
宝塚記念	0	1	0	0	0	2200	0	2	1	0	0	15
スプリングステークス	0	0	1	0	0	1200	2	2	1	0	0	15
秋賞	0	0	0	1	0	2000	2	2	0	1	1	15
菊花賞	0	0	0	1	0	3000	0	1	0	1	1	15
天皇賞 (秋)	0	0	0	0	1	2500	0	2	1	1	0	15
E1杯 (スズマ杯)	0	0	0	0	1	2200	2	2	1	1	1	15
マイルチャンピオンシップ	0	0	0	1	0	1600	0	2	1	0	1	15
マイルチャンピオンシップ	0	0	0	1	0	2400	0	2	1	1	1	15
有馬記念	0	0	1	0	0	2500	0	2	1	1	0	15

2.2. データの標準化

競馬データには様々なタイプや単位の項目が含まれるため、下記の計算式によりデータの標準化を行う。

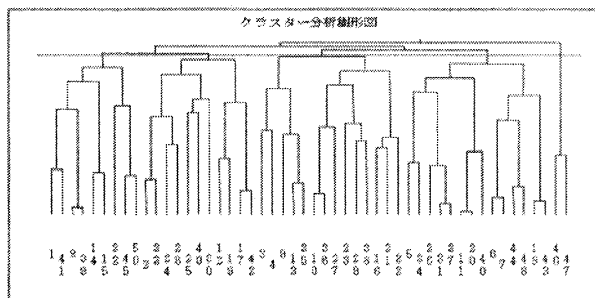
標準化データ

$$= (\text{個々のデータ} - \text{平均値}) / \text{標準偏差}$$

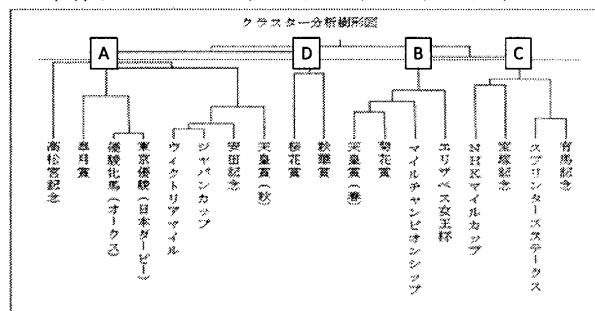
2.3. 標準化されたデータのクラスター分析

標準化された馬データ、レースデータのクラスター分析を行う。

図表 3. 馬データのデンドログラム



図表 4. レースデータのデンドログラム



2.4. クロス集計表の作成

クラスター分析の結果に基づいて「レースクラスター」×「馬クラスター」別の勝率算出する (図表 5)。図表 5 は、各馬クラスターに属する馬が、各レースクラスターに出走した場合の勝率を表している。

図表 5. 2007 年データに基づいた馬クラスター別の勝利率

馬分類 \ レース分類	1	2	3	4	5	総計
A	17%	6%	0%	22%	0%	44%
B	6%	0%	0%	17%	0%	22%
C	0%	6%	11%	6%	0%	22%
D	0%	0%	0%	11%	0%	11%
総計	22%	11%	11%	56%	0%	100%

2.5. 多群判別分析による新規レースの予想

多群判別分析を用いて、新規レース (2008 年ジャパンカップ) に出走した馬を図表 5 の A~D の馬クラスターに分類 (判別) する (図表 6)。

図表 6. 2008 年のジャパンカップに出走した各場の判別結果と実際の順位 (結果)

着順	馬名 (G1 出走馬)	属するクラスター
1	スクリーンヒーロー	1
2	ディープスカイ	5
3	ウオッカ	5
4	マツリダゴッホ	1
5	オウケンブルー	1
6	メイショウサムソン	4
7	ネヴァクション	2
8	アサクサキングス	1
9	パープルムーン	1
10	トーホウアラン	1
11	オースミグラスワン	5
12	アドマイヤモナーク	5
13	シックスティーズアイコン	1
14	ペイパルブル	1
15	トーセンキャプテン	1
16	ダイワワイルドボア	1
17	コスモバルク	5

3. 分析結果の検討

図表 5 と 6 から、2008 年のジャパンカップの勝ち馬は、「馬クラスター 4」に属するメイショウサムソンと予想されたが、実際は「馬クラスター 1」に属するスクリーンヒーローであった。しかし、図表 5 から「馬クラスター 1」に属する馬も「馬クラスター 4」に匹敵する勝率であることは注目すべきであろう。また、勝率が「馬クラスター 5」に属する 2 頭が上位に位置しているが、ディープスカイは「馬クラスター 1 および 4」に近く、ウオッカは「馬クラスター 4」に近かったことも注目すべきであろう。

4. まとめ

本研究では、勝ち馬予想手順を提案したが、結果は実際の勝ち馬予想に有用な段階に達しなかった。今後は用いる項目の選択方法を検討するなどの改善を行う予定である。

参考文献

- [1] 月本 洋, 実践データマイニング—金融・競馬予測の科学, オーム, 1999
- [2] 小岩井 弥, 競馬データを 120% 活用するクォンツ分析理論, 毎日コミュニケーションズ, 1996
- [3] 竹内 光悦, 酒折 文武, Excel で学ぶ理論と技術 多変量解析入門, ソフトバンククリエイティブ, 2006