

ZASTOSOWANIE DRZEW KLASYFIKACYJNYCH DO ANALIZY POKERA ONLINE

Marek Zaslona

Kolegium Zarządzania i Finansów
Szkola Główna Handlowa w Warszawie
e-mail: mz75781@doktorant.sgh.waw.pl

Tomasz Ząbkowski (ORCID: 0000-0003-1722-1179)

Wydział Zastosowań Informatyki i Matematyki
Szkola Główna Gospodarstwa Wiejskiego w Warszawie
e-mail: tomasz_zabkowski@sggw.pl

Streszczenie: Niniejsza publikacja stanowi próbę scharakteryzowania deterministycznych czynników wpływających na wygraną w pokera. Przeprowadzono analizę w oparciu o jedną z metod eksploracji danych – drzewa klasyfikacyjne. Wybór tej techniki podyktowany był wykorzystaniem danych jakościowych jako zmiennych objaśniających rozgrywkę pokerową oraz prostotą prezentacji otrzymanych wyników, nawet przy bardzo rozbudowanych drzewach. W badaniu odkryto kilka czynników, które w istotny sposób mają wpływ na przebieg gry.

Słowa kluczowe: eksploracja danych, drzewa klasyfikacyjne, poker

JEL classification: C15, C38, C44, C57

WPROWADZENIE I CEL BADANIA

Nie ulega wątpliwości, że w każdej grze karcianej, również w pokerze, występuje element losowy. W tym miejscu należy jednak zadać pytanie, czy oprócz elementu losowego występuje także element związany z umiejętnościami, uwzględniający doświadczenie gracza. Jeśli tak to, który z tych elementów ma większy wpływ na rezultat rozgrywki? Gdyby bowiem okazało się, że umiejętności gracza są bardziej istotne niż szczęśliwy traf, należałoby wtedy stwierdzić, że poker nie jest grą hazardową (losową), lecz grą umiejętności z elementem losowym.

Celem pracy jest zastosowanie drzew klasyfikacyjnych jako metody eksploracji danych w badaniu rozgrywki pokerowej w celu identyfikacji czynników determinujących wynik gry, co może być istotnym elementem strategii gracza. W pracy zostanie zweryfikowana hipoteza, że gra w pokera jest grą umiejętności z elementem losowym, a nie tylko grą losową i uzależniona jest w znacznym stopniu od właściwego wykorzystania przez gracza danych w postaci kolorów i figur otrzymanych kart.

ZAŁOŻENIA DO DANYCH

Wszystkie dane, które zostały poddane analizie, zostały uzyskane poprzez rozegranie gier według jednej odmiany pokera, *Texas Hold'em* w formacie turnieju jednostolikowego (*Sit'n'Go*) dla sześciu graczy na tym samym poziomie wpisowego. 30 tysięcy rozdań (rozegranych przez jednego gracza) było wybrane w sposób losowy, przedział czasowy z posiadanej przez gracza całej jego historii.

Tabela 1. Startowe karty gracza

Grupa	Karty gracza
1	AA, AKs, KK, QQ, JJ
2	AK, AQs, AJs, KQs, TT
3	AQ, ATs, KJs, QJs, JTs, 99
4	AJ, KQ, KTs, QTs, J9s, T9s, 98s, 88
5	A9s - A2s, KJ, QJ, JT, Q9s, T8s, 97s, 87s, 77, 76s, 66
6	AT, KT, QT, J8s, 86s, 75s, 65s, 55, 54s
7	K9s - K2s, J9, T9, 98, 64s, 53s, 44, 43s, 33, 22
8	A9, K9, Q9, J8, J7s, T8, 96s, 87, 85s, 76, 74s, 65, 54, 42s, 32s
9	Pozostałe karty, których nie ma w poprzednich grupach

Oznaczenia kart: A – As, K – Król, Q – Dama, J – Walet, T - Dziesiątka.

Źródło: opracowanie własne

Sklansky i Malmuth [Sklansky, Malmuth 1999] na bazie doświadczeń, statystyki i własnych spostrzeżeń stworzyli podział wszystkich możliwych kombinacji kart w celu łatwego określenia siły swojej ręki. Na tej podstawie opierają się wszystkie strategie gry w pokera [Harrington, Robertie 2006]. W tabeli 1 został przedstawiony podział kart gracza, z którymi może zacząć rozgrywkę pokerową. Karty te zostały podzielone na dziewięć grup ze względu na ich siłę tzn. karty w pierwszej grupie są najsilniejsze, a w dziewiątej najsłabsze. Litera „s” przy kartach oznacza, że obie karty są w tym samym kolorze (AKs może oznaczać zarówno Asa pik i Króla pik jak i Asa trefl i Króla trefl – analogicznie karo lub kier).

CHARAKTERYSTYKA DANYCH

Z trzydziestu tysięcy rozegranych rozdań do dalszej analizy zostało zakwalifikowanych 23481. Zakwalifikowane rozdania odbyły się podczas 346 jednostolikowych turniejów. Podczas zakwalifikowanych rozdań gracz zajął 215 razy miejsca płatne, tzn. takie, którego uzyskanie pozwala na otrzymanie nagrody.

Spośród rozdań, które zostały poddane analizie wynika, że co czwarte z nich (5873 rozdania z 23481) zostało wygrane przez gracza, którego gra została przeanalizowana.

Powody usunięcia 6506 rozdań:

- Zbędne rekordy: rozdania, których historia rozgrywki turniejowej nie uwzględniała całego przebiegu turnieju: brakowało początkowych lub końcowych rozdań.
- Rekordy z brakującymi wartościami: rozdania, które z niewyjaśnionych powodów nie miały ciągłej struktury. Problemem mogło być zerwanie połączenia z Internetem lub wyłączenie aplikacji podczas rozgrywki.
- Punkty oddalone: rozdania, które zostały rozegrane w turniejach, gdzie była inna liczba graczy, inne wpisowe lub inny rodzaj rozgrywki niż przyjęte w założeniach.

Jednym z problemów towarzyszących budowie drzew klasyfikacyjnych jest wybór cech (zmiennych), na których będzie dokonywany podział zbioru. Dlatego też nie wszystkie informacje, które były dostępne zostały wykorzystane w tej pracy. W analizie pominięte zostały informacje dotyczące m.in.:

1. Opisu rozgrywki: wysokość wpisowego, data i godzina rozpoczęcia danego rozdania, unikalny numer przypisany do każdego rozdania, numer wirtualnego stolika, na którym zostało rozegrane rozdanie.
2. Przebiegu rozgrywki: poziom gry, liczbę wyświetlanych kart wspólnych (*Flop, Turn, River*), na której rundzie rozdania wystąpiło zakończenie rozgrywki, kto wygrał dane rozdanie, jak wyglądał podział puli nagród (jeśli wygrał więcej niż jeden gracz).
3. Graczy: nazwa gracza, ilość żetonów posiadanych przez danego gracza, miejsce, na którym siedzi dany gracz, wielkość zakładu wniesionego przez gracza w poszczególnych turach licytacji, rodzaj reakcji gracza, układ kart posiadanych przez gracza na koniec rozgrywki (o ile dochodzi do pokazania kart).
4. Inne: rozmowy graczy poprzez wbudowany moduł służący do komunikacji między graczami podczas rozgrywki, informacje reklamowe automatycznie wysyłane do graczy przez aplikację, z której korzystają.

W celu analizy danych, pozycja danego gracza została najpierw przekształcona na pozycję danego gracza względem rozdającego, a później została wystandaryzowana ze względu na ilość graczy biorących udział w danym rozdaniu.

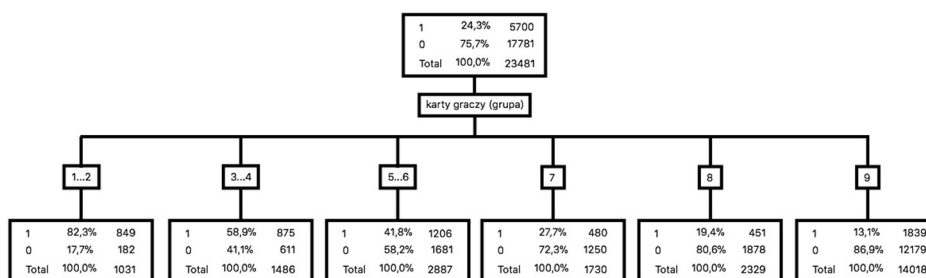
ANALIZA DANYCH Z WYKORZYSTANIEM DRZEW KLASYFIKACYJNYCH

Do analizy zostały wykorzystane drzewa klasyfikacyjne. Wybór taki podyktowany był prostotą prezentacji uzyskanych reguł klasyfikacyjnych, nawet przy rozbudowanych drzewach [Larose 2006]. Do analizy zgromadzonych danych zostały zastosowane dwa różne modele drzewiaste, gdzie skorzystano z indeksu Giniego lub testu Chi² jako metod wyboru zmiennych do podziału drzewa. Kryterium Giniego oparte jest na indeksie Giniego jako mierze koncentracji zmiennej losowej. Nadrzędnym celem w tym przypadku jest dokonanie podziału na możliwie jednorodny przypadki w węzłach potomnych. Z kolei w drugim przypadku ocena wykonywana jest przez obliczenie testu Chi-kwadrat (Pearsona) i wybierany jest predyktor o najniższej wartości poziomu p, a więc ten, który daje najbardziej istotny podział populacji.

Wszystkie przedstawione w tej pracy drzewa są drzewami niebinarnymi tzn. z każdego węzła mogły wychodzić więcej niż dwie gałęzie.

Punktem wyjściowym jest 23481 rozdań z czego 5700 (co stanowi 24,3% wszystkich rozdań) uzyskało sukces (wygrane rozdanie), natomiast reszta rozdań tzn. 17781 (co stanowi 75,7% wszystkich rozdań) reprezentuje porażkę (przegrane rozdanie). Podział puli (remis) zaliczany jest do sukcesu. Jako pierwszą zmienną występującą w rozdaniu, która została poddana analizie, wybrano karty gracza. Zaznaczyć należy, iż nie są istotne poszczególne karty, tylko ich połączenie. Każda para kart została zaklasyfikowana do jednej z 9 grup (zgodnie z tabelą 1). Po zastosowaniu drzewa otrzymano graf składający się z korzenia i sześciu liści (rysunek 1), jako „1” oznaczono sukces, a „0” porażkę.

Rysunek 1. Podział drzewa ze względu na otrzymane przez gracza karty



Źródło: opracowanie własne

Z analizy drzewa na rysunku 1 jednoznacznie wynika, że karty z pierwszych dwóch grup zdecydowanie zwiększają prawdopodobieństwo wygranej, i to ponad trzykrotnie względem ogółu (82,3% na sukces w pierwszym liściu w porównaniu do 24,3% ogółu). Ze względu na to, że wszystkie rozdania można przypisać do jednej z dwóch grup (zwycięstwo lub porażka), przy czym remis został zaliczony

do zwycięstwa, można zauważyć, że szansa na porażkę gracza, którego karty znajdują się w jednej z pierwszych dwóch grup zmniejsza się ponad trzykrotnie względem całej populacji (17,7% na porażkę w pierwszym liściu w porównaniu do 75,7% ogółu). Zdecydowanie inaczej wygląda sytuacja w ostatnim liściu, w którym znalazło się 14018 rozdań. Tam prawdopodobieństwo sukcesu wynosi 13,1%. Jest to o 11,2 punktów procentowych mniej porównując ze wszystkimi rozdaniem.

W dalszej części artykułu zostały zaprezentowane oraz ocenione wyniki klasyfikacji dla modeli korzystających atrybutów innych niż karty startowe gracza.

Do oceny modeli wykorzystano szereg miar takich jak trafność klasyfikacji, czułość i specyficzność. Punktem wyjścia było zbudowanie macierzy klasyfikacji zgodnie z tabelą 2.

Tabela 2. Macierz klasyfikacji

Wartości rzeczywiste	Wartości zakładane	
	Pozytywna (1)	Negatywna (0)
Pozytywna (1)	TP	FN
Negatywna (0)	FP	TN

Źródło: opracowanie własne

Oznaczenia wykorzystane w tabeli:

- TP (ang. *True positive*) – poprawna klasyfikacja do klasy pozytywnej
- FN (ang. *False negative*) – błędna klasyfikacja do klasy pozytywnej
- TN (ang. *True negative*) – poprawna klasyfikacja do klasy negatywnej
- FP (ang. *False positive*) – błędna klasyfikacja do klasy negatywnej

Następnie, na tej podstawie macierzy klasyfikacji wyznaczono następujące miary:

- Trafność/Dokładność (ang. *Accuracy*)

$$\text{trafność} = \frac{TP + TN}{TP + FN + FP + TN}$$

- Czułość/Wrażliwość (ang. *Sensitivity*)

$$\text{czułość} = \frac{TP}{TP + FN}$$

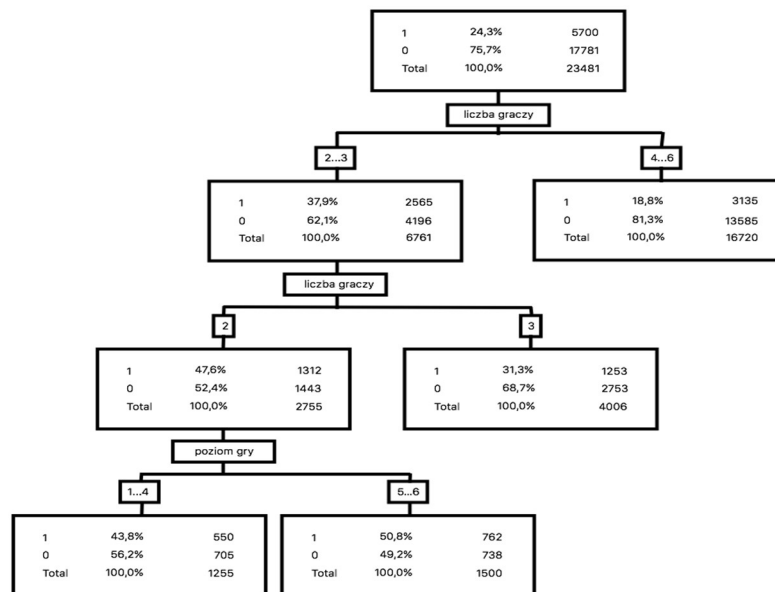
- Specyficzność/Swoistość (ang. *Specificity*)

$$\text{specyficzność} = \frac{TN}{TN + FP}$$

Dodatkowo można zauważyć następujące zależności: $FP = 1 - \text{specyficzność}$ oraz $FN = 1 - \text{czułość}$.

W wyniku analizy zbudowano szereg drzew, z różną liczbą zmiennych oraz o różnej głębokości. Jedno z przykładowych drzew zaprezentowane zostało na rysunku 2. W strukturze tej, „liczba graczy” oraz „poziom gry” są cechami, w oparciu o które dokonał się podział (biorąc pod uwagę kryterium Giniego). Dla przykładu, liczba graczy nie większa niż 2 oraz wysoki poziom gry warunkują ponad dwukrotnie większą szansę na wygraną.

Rysunek 2. Drzewo uzyskane w oparciu o dwie zmienne (poziom gry i liczba graczy)



Źródło: opracowanie własne

Zbiorcze zestawienie wyników dla wszystkich rozważanych drzew (łącznie 14), z różną liczbą zmiennych oraz o różnej głębokości przedstawiają tabela 3 oraz tabela 4. Przy czym tabela 3 zawiera wyniki dla zbioru uczącego, a tabela 4 dla zbioru walidacyjnego, gdzie przypisania dokonano w sposób losowy w proporcjach 60% i 40%, odpowiednio dla zbioru uczącego i walidacyjnego.

Z analizy tabel wynika, że tam gdzie występuje zmienna „Karty gracza”, trafność oraz czułość jest o kilka punktów procentowych większa niż w innych modelach. Porównując obie tabelę można zauważyć, że nie ma znacznej różnicy w miarach jakości modeli (poza drzewem nr 8), co wskazywuje na stabilność uzyskanych wyników.

Tabela 3. Zbiorcze porównanie wybranych drzew decyzyjnych – zbiór uczący

Nr drzewa	Zmienne	Trafność	Czułość	Specyficzność
1	Karty gracza	79,7%	30,2%	95,5%
2	Poziom gry	75,7%	0,0%	100 %
3	Liczba graczy	75,7%	0,0%	100 %
4	Pozycja gracza	75,7%	0,0%	100 %
5	Kolor kart	75,7%	0,0%	100 %
6	Karty gracza Liczba graczy	80,8%	36,4%	95,1%
7	Kolor kart Liczba graczy	75,9%	5,9%	98,3%
8	Poziom gry Liczba graczy	75,9%	2,2%	99,5%
9	Karty gracza Kolor kart	80,0%	28,1%	96,7%
10	Pozycja gracza Liczba graczy	76,2%	13,1%	96,4%
11	Pozycja gracza Poziom gry	75,9%	1,3%	99,9%
12	Poziom gry Pozycja gracza Liczba graczy	76,3%	12,5%	99,9%
13	Poziom gry Pozycja gracza Kolor kart Liczba graczy	76,3%	14,1%	96,3%
14	Poziom gry Pozycja gracza Kolor kart Karty gracza Liczba graczy	81,4%	38,6%	95,1%
Średnia		77,2%	13,0%	98,0%
Odchylenie		0,02	0,14	0,02

Źródło: opracowanie własne

Tabela 4. Zbiorcze porównanie wybranych drzew testowych – zbiór walidacyjny

Nr drzewa	Zmienne	Trafność	Czułość	Specyficzność
1'	Karty gracza	79,2%	29,3%	95,3%
2'	Poziom gry	75,6%	0%	100%
3'	Liczba graczy	75,6%	0%	100%
4'	Pozycja gracz	75,6%	0%	100%
5'	Kolor kart	75,6%	0%	100%
6'	Karty gracza Liczba graczy	80,4%	37,8%	94,2%
7'	Kolor kart Liczba graczy	75,7%	5,9%	98,2%
8'	Poziom gry Liczba graczy	75,9%	13,9%	95,8%

Nr drzewa	Zmienne	Trafność	Czułość	Specyficzność
9'	Karty gracza Kolor kart	79,6%	27,2%	96,5%
10'	Pozycja gracza Liczba graczy	76,0%	13,0%	96,3%
11'	Pozycja gracza Poziom gry	75,8%	1,4%	99,8%
12'	Poziom gry Pozycja gracza Liczba graczy	76,2%	12,5%	96,8%
13'	Poziom gry Pozycja gracza Kolor kart Liczba graczy	76,2%	12,5%	96,8%
14'	Poziom gry Pozycja gracza Kolor kart Karty gracza Liczba graczy	81,0%	36,3%	95,4%
Średnia		77,0%	13,6%	97,1%
Odchylenie		0,02	0,14	0,02

Źródło: opracowanie własne

WNIOSKI

W niniejszej pracy wykazano, że istnieją czynniki, które mają istotny wpływ na wygraną w rozgrywce pokerowej. Stopień tego wpływu jest zróżnicowany, przy czym duży wpływ na wynik rozgrywki mają startowe karty gracza i poziom gry. Umiarkowany wpływ jest związany z takimi cechami jak liczba graczy, pozycja gracza oraz kolor kart gracza. Dodatkowo, znaczenie tych czynników podkreśla fakt, iż wszystkie z nich występują w każdej odmianie pokera i można je odpowiednio stosować. Jednocześnie w niniejszej pracy potwierdzono użyteczność i zaletę wykorzystania drzew decyzyjnych jako sposobu przejrzystej prezentacji wyników, co sprzyja łatwości ich zrozumienia oraz interpretacji reguł klasyfikacyjnych.

Zaprezentowane wyniki dowodzą, że kluczowe jest, aby z posiadanych danych wywnioskować jak najwięcej, co pozwoli graczowi w rozgrywce podjąć lepszą decyzję względem innych graczy. Wykorzystywanie tak pojętej dodatkowej wiedzy do osiągnięcia celu można nazwać umiejętnością. Specyfika pokera sprawia, że podczas gry nie mamy pełnych informacji dotyczących wszystkich przeciwników, gdyż związane jest to z losowym przydziałem graczy do określonych stołów pokerowych. Dlatego istotne jest zwiększenie umiejętności i poprawienie gry nie koncentrując się na poszczególnych graczach, ale na ogóle zachowań czy charakterystyk dla większej liczby graczy.

Wyniki zaprezentowanych badań mają potencjał wykorzystania w kilku dziedzinach. Pierwszą z nich jest prawo. Zmiana definicji pokera i zaklasyfikowanie go jako gry umiejętności z elementem losowym wymuszałoby na prawodawcy zmianę m.in. ustawy o grach losowych, co w dalszej konsekwencji doprowadziłoby do organizacji wielu legalnych imprez (turniejów) pokerowych, w tym nawet rozgrywek na skalę krajową (ligi, mistrzostw Polski). Drugim obszarem zastosowania pracy jest rozwój strategii pokerowej. Zapoznanie się z wynikami pracy może posłużyć graczom pokerowym, zarówno początkującym jak i tym bardziej zaawansowanym, do osiągania lepszych wyników poprzez właściwą analizę czynników. Wreszcie lektura tej pracy mogłaby zostać uznana za przydatną w sferze szeroko pojętej kultury i obyczajowości. Wbrew pozorom, tego typu rozrywki, mają negatywny wizerunek, ponieważ postrzeganie pokera i negatywne stereotypy z nim związane są mocno zakorzenione w ogólnym przekazie społecznym. Jednak zaakcentowanie czynnika umiejętności jako rozgraniczającego grę amatorską, kojarzoną z hazardem, od profesjonalnej, wymagającej wiedzy matematycznej i dużego poziomu umiejętności, pomogłoby ukształtować bardziej pozytywny stosunek społeczeństwa do tej formy rozrywki.

Planowana jest kontynuacja badań związanych z analizą pokera on-line, mająca na celu próbę identyfikacji nowych czynników wpływających na rozgrywkę, a także analiza zagadnień takich jak ryzyko, czy zarządzanie budżetem gracza.

BIBLIOGRAFIA

- Berthet V. (2010) Best Hand Wins: How Poker Is Governed by Chance, *Chance* 23(3), 34-38.
- Billings D., Davidson A., Schaeffer J., Szafron S. (2000) The Challenge of Poker. *Artificial Intelligence*, 134(1-2), 201-240.
- Czajkowski M. (2015) Poker – gra szczęścia czy umiejętności. *Ekonomia*, 40, 33-56.
- Demski T. (2004) Drzewa klasyfikacyjne w przewidywaniu migracji klientów (churn). *SYSTEMY IT*, 53-57.
- Hand D., Mannila H., Smyth P. (2005) *Eksploracja danych*. Wydawnictwo Naukowo – Techniczne, Warszawa.
- Harrington D., Robertie B. (2006) *Harrington on Hold'em Expert Strategy for No-Limit Tournaments; Volume I: Strategic Play*. Creel Printing, Inc. Las Vegas, Nevada.
- Kantardzic M. (2003) *DATA MINING. Concepts, Models, Methods and Algorithms*. IEEE Computer Society, Sponser, University of Louisville.
- Larose D. T. (2006) *Odkrywanie wiedzy z danych. Wprowadzenie do eksploracji danych*. Wydawnictwo Naukowe PWN, Warszawa.
- Potter van Loon R. J. D., Van den Assem M. J., Van Dolder D. (2014) Beyond Chance? The Persistence of Performance in Online Poker, *SSRN*, 10(3), 1-35.
- Sklansky D., Malmuth M. (1999) *Hold'em poker for advanced players*. Two Plus Two Publishing LLC.

**APPLICATION OF CLASSIFICATION TREES
TO ANALYSE POKER GAME OUTCOME**

Abstract: The paper aims to characterize key factors determining poker game outcome. The analysis was based on classification trees and this was due to the qualitative data used as the explanatory variables. The method enables clear presentation of the results even in case of very complex tree structures. The study describes also a few other factors that significantly influence the game outcome.

Keywords: data mining, classification trees, poker