

Západočeská univerzita v Plzni
Fakulta aplikovaných věd
Katedra matematiky

Diplomová práce

Použití Elo ratingu pro predikci výsledků utkání NBA

Prohlášení

Prohlašuji, že jsem diplomovou práci vypracoval samostatně a výhradně s použitím citovaných zdrojů informací.

V Plzni dne 2. června 2020

Bc. Jan Ondříček

Poděkování

Touto cestou bych rád poděkoval Ing. Patrice Markovi, PhD. za vedení této diplomové práce, odborný dohled, ochotu a čas, který této práci věnoval. Poděkování patří především také celé mojí rodině za neustálou podporu během mého studia.

Abstrakt

Hlavním cílem této práce je navrhnout změny a vylepšení Elo rating systému pro predikci výsledků utkání NBA a následně srovnat výsledky predikce s modelem vybrané sázkové kanceláře a s existujícím modelem od FiveThirtyEight založeným na Elo ratingu. V práci je prezentován historický vývoj a matematický popis principu Elo rating systému. Původní model Elo ratingu (navržený pro šach) je zde rozšiřován postupným přidáváním celkem 5 faktorů, za účelem zahrnout do modelu více informací o utkáních a pokusit se eliminovat některé jeho nedostatky. Jedná se o faktory zachování ratingu po skončení sezony, vlivu domácího prostředí, vlivu back-to-back utkání, sezonního poklesu maximální možné změny ratingu týmů po odehraném utkání a faktoru převzatého od FiveThirtyEight, který v maximální možné změně ratingu týmů zohledňuje také výsledné skóre utkání a naplnění jeho očekávání. Vybrané parametry všech modelů jsou optimalizovány a na základě kritérií kvality modelů je vybrán nejlepší z uvažovaných modelů. Predikční schopnost tohoto vybraného modelu je demonstrována na fiktivním použití proti sázkové kanceláři, a to jak z pohledu sázkaře, tak z pohledu sázkové kanceláře. Kromě toho je predikční schopnost vybraného modelu srovnána na základě kritérií kvality modelu s modelem sázkové kanceláře a také s existujícími modely od FiveThirtyEight.

Klíčová slova: Elo rating, NBA, basketbal, predikce, optimalizace modelu, predikční schopnost

Abstract

The main purpose of this thesis is to propose changes and improvements of the Elo rating system for a prediction of the NBA games results and compare the prediction outcomes with a model of selected bookmaker and with the FiveThirtyEight model based on the Elo rating. This thesis presents a historical development and a description of the mathematical background of the Elo rating system. The original Elo rating model (designed for chess) is extended by gradual addition of 5 factors in order to include more matches information and try to eliminate its shortcomings. These factors are about the maintaining the rating after the end of the season, the influence of the home advantage, the influence of back-to-back games, the seasonal decrease of the maximum possible change of the team rating after the game and the factor taken from FiveThirtyEight, which takes the final score of the game and its expectation fulfillment into account. The selected parameters of all models are optimized and based on the criterion of quality of the models the best of the considered models is selected. The predictive ability of this model is demonstrated on the fictitious use against the bookmaker, both from the bettor's point of view and from the bookmaker's point of view. In addition, the predictive ability of the selected model is compared with bookmaker's model and also with FiveThirtyEight models, based on the criterion of quality.

Keywords: Elo rating, NBA, basketball, prediction, model optimization, predictive ability

Obsah

1	Úvod	1
2	Popis získaných dat	3
2.1	Průběh soutěže a utkání	3
2.2	Data	4
3	Elo rating	7
3.1	Historický vývoj	7
3.2	Matematický popis	8
3.2.1	Normální rozdělení výkonnosti	9
3.2.2	Logistické rozdělení výkonnosti	10
3.2.3	Přepočet ratingu	13
3.2.4	Koeficient rozvoje K	14
4	Elo rating pro NBA	16
4.1	Model 1	16
4.2	Model 2	17
4.3	Model 3	18
4.4	Model 4	21
4.5	Model 5	23
4.6	Model Elo538	25
5	Optimalizace a výběr modelu	27
5.1	Kritéria kvality modelů	27
5.1.1	Přesnost	28
5.1.2	Logaritmická ztrátová funkce	28
5.1.3	Brierovo skóre	29
5.1.4	Kalibrační poměr	30
5.1.5	AUC-ROC	30
5.2	Optimalizace parametrů	33
5.2.1	Model 1	33
5.2.2	Model 2	34
5.2.3	Model 3	36
5.2.4	Model 4	38
5.2.5	Model 5	41
5.3	Výběr modelu	42

6	Predikční schopnost modelu	44
6.1	Sezona 2016/2017	45
6.1.1	Kritéria kvality modelů	45
6.1.2	Výběr utkání pro fiktivní sázení	46
6.1.3	Náhodný sázející	48
6.2	Sezona 2017/2018	51
6.2.1	Kritéria kvality modelů	51
6.2.2	Výběr utkání pro fiktivní sázení	51
6.2.3	Náhodný sázející	52
6.3	Sezona 2018/2019	55
6.3.1	Kritéria kvality modelů	55
6.3.2	Výběr utkání pro fiktivní sázení	55
6.3.3	Náhodný sázející	56
7	Závěr	59
	Literatura a další zdroje	62
	Příloha A	65

Seznam tabulek

2.1	Aktuální názvy všech současných týmů NBA	4
2.1a	Východní konference	4
2.1b	Západní konference	4
2.2	Získané zápasové statistiky a jejich vysvětlivky	5
2.3	Ukázka formátu vstupních dat (pět utkání, každé utkání na dva řádky) .	6
3.1	Výkonnostní třídy USCF	8
3.2	Pravděpodobnost vítězství v utkání v závislosti na rozdílu ratingů	12
3.2a	Normální rozdělení	12
3.2b	Logistické rozdělení	12
5.1	Kritéria kvality modelu 1, pro optimální parametry – optimalizace	34
5.2	Kritéria kvality modelu 2, pro optimální parametry – optimalizace	35
5.3	Kritéria kvality modelu 3, pro optimální parametry – optimalizace	37
5.4	Kritéria kvality modelu 4, pro optimální parametry – optimalizace	39
5.5	Kritéria kvality modelu 5, pro optimální parametry – optimalizace	42
5.6	Kritéria kvality všech uvažovaných modelů, včetně Elo538 – optimalizace	43
6.1	Kritéria kvality modelu 5, Elo538 a Bet365, pro sezonu 2016/2017 – predikce	46
6.2	Srovnání zisku po skončení sezony 2016/2017 SK modelu 5 a Bet365, vždy z 10 000 simulací pro různé hodnoty parametru alternativního roz- dělení p	49
6.3	Kritéria kvality modelu 5, Elo538 a Bet365, pro sezonu 2017/2018 – predikce	51
6.4	Srovnání zisku po skončení sezony 2017/2018 SK modelu 5 a Bet365, vždy z 10 000 simulací pro různé hodnoty parametru alternativního roz- dělení p	53
6.5	Kritéria kvality modelu 5, Elo538, RAPTOR a Bet365, pro sezonu 2018/2019 – predikce	55
6.6	Srovnání zisku po skončení sezony 2018/2019 SK modelu 5 a Bet365, vždy z 10 000 simulací pro různé hodnoty parametru alternativního roz- dělení p	57

Seznam obrázků

3.1	Pravděpodobnost vítězství v utkání v závislosti na rozdílu ratingů	12
4.1	Průměrný počet back-to-back situací na tým, za jednu sezonu (82 utkání)	21
4.2	Graf MOV multiplikátoru v závislosti na MOV a d_w	23
5.1	Logaritmická ztrátová funkce	29
5.2	Bierovo skóre	30
5.3	Ukázka různých příkladů AUC-ROC, zdroj [21]	31
5.4	Logaritmická ztrátová funkce v závislosti na K a p v okolí nalezeného optima pro model 1	34
5.5	Logaritmická ztrátová funkce v závislosti na K a p v okolí nalezeného optima pro model 2	36
5.6	Logaritmická ztrátová funkce v závislosti na K a p v okolí nalezeného optima pro model 3	37
5.7	Srovnání optimálního koeficientu rozvoje K pro model 4 a model 3 v závislosti na počtu utkání v sezoně	39
5.8	Logaritmická ztrátová funkce v závislosti na a a p v okolí nalezeného optima pro model 4	40
5.9	Logaritmická ztrátová funkce v závislosti na k a a v okolí nalezeného optima pro model 4	40
5.10	Vývoj ratingu vybraných týmů podle modelu 5 v celém sledovaném období	43
6.1	Počet vsazených utkání a průměrná návratnost v závislosti na L pro sezonu 2016/2017	47
6.2	Průměrný zisk po skončení sezony 2016/2017 pro SK modelu 5 a Bet365 pro různé hodnoty parametru alternativního rozdělení p	50
6.3	Počet simulací z celkových 10 000, ve kterých byla pro různé hodnoty parametru alternativního rozdělení p v sezoně 2016/2017 jedna sázková kancelář úspěšnější než druhá uvažovaná	50
6.4	Počet vsazených utkání a průměrná návratnost v závislosti na L pro sezonu 2017/2018	52
6.5	Průměrný zisk po skončení sezony 2017/2018 pro SK modelu 5 a Bet365 pro různé hodnoty parametru alternativního rozdělení p	54
6.6	Počet simulací z celkových 10 000, ve kterých byla pro různé hodnoty parametru alternativního rozdělení p v sezoně 2017/2018 jedna sázková kancelář úspěšnější než druhá uvažovaná	54

6.7	Počet vsazených utkání a průměrná návratnost v závislosti na L pro sezonu 2018/2019	56
6.8	Průměrný zisk po skončení sezony 2018/2019 pro SK modelu 5 a Bet365 pro různé hodnoty parametru alternativního rozdělení p	58
6.9	Počet simulací z celkových 10 000, ve kterých byla pro různé hodnoty parametru alternativního rozdělení p v sezoně 2018/2019 jedna sázková kancelář úspěšnější než druhá uvažovaná	58

Kapitola 1

Úvod

Predikce výsledků sportovních utkání není důležitou oblastí pouze pro sázkové kanceláře a jejich zákazníky, ale také pro sportovní fanoušky, majitele a manažery týmů, trenéry atd. V současnosti existuje velké množství predikčních modelů pro různá sportovní odvětví a zabývá se jimi mnoho vědeckých výzkumů. Tato práce se zabývá predikcí výsledků utkání nejvyšší profesionální basketbalové ligy v Severní Americe – NBA (Nation Basketball Asociation). Modely pro predikci výsledků utkání NBA jsou často založeny na strojovém učení (viz například [3]), ale mezi oblíbené metody patří také modely založené na tzv. Elo ratingu, ze kterého vychází také původní predikční model webu `FiveThirtyEight.com` [26]. Model pro predikci utkání NBA založený na Elo ratingu byl také vyhodnocen v bakalářské práci [19] jako nejúspěšnější z analyzovaných modelů, proto bude použit i v této práci. Hlavním cílem bude navrhnout změny a vylepšení původního Elo rating systému spolu s optimalizací vybraných parametrů.

Původní Elo rating systém pro hodnocení hráčů šachu, který ve své knize [6] představil Arpad Elo, předpokládal normální rozdělení náhodné veličiny popisující výkonnost hráčů, resp. týmů. Poznamenal však, že kromě normálního rozdělení může být použito také logistické, které postupně začalo být šachovými federacemi upřednostňováno a dnes je již základem běžně používaného modelu Elo rating systému. Matematický popis Elo rating systému bude, stejně jako jeho historický vývoj, obsahem kapitoly 3.

V roce 2015 začala webová stránka `FiveThirtyEight.com` používat pro predikci výsledků utkání NBA rozšířený Elo rating systém [26], ze kterého budeme v této práci vycházet. Rozšíření tohoto modelu spočívá především v zohlednění výsledného skóre utkání a naplnění jeho očekávání, při aktualizaci týmového ratingu (hodnocení) po odehraném utkání. V kapitole 4 bude uvedeno několik modelů s postupně přidávanými faktory, za účelem zahrnout do původního modelu Elo rating pro šach více informací o utkáních a pokusit se eliminovat některé nedostatky původního modelu, resp. více přizpůsobit model realitě. Některé faktory budou inspirovány právě již existujícím a používaným modelem od `FiveThirtyEight`. Uvažované modely a jejich rozšíření budou naprogramovány v softwaru MATLAB R2018a a aplikovány na reálná data.

Zajištěna budou data ze všech základních částí sezon 2012/2013 až 2018/2019, konkrétně výsledky všech utkání a zápasové statistiky, jejichž popis bude obsahem kapi-

toly 2. Tato data budou podle sezon rozdělena do tří datových sad: *učení*, *optimalizace* a *predikce*. V optimalizační sadě budou v kapitole 5 optimalizovány vybrané parametry všech uvažovaných modelů a následně bude vybrán nejlepší z uvažovaných optimalizovaných modelů. K tomu využijeme vybraná kritéria kvality predikčních modelů, která použil ve své práci [18] Rogier Noordman pro srovnání predikčních modelů výsledků fotbalových utkání a vybraná kritéria, která použila Stephanie Ann Kovalchik ve svém článku [13] pro srovnání predikčních modelů výsledků tenisových utkání.

Predikční schopnost našeho vybraného modelu bude v kapitole 6 hodnocena ze tří různých pohledů. V predikční sadě dat budou nejprve srovnána kritéria kvality (stejná jako pro výběr modelu) vybraného modelu s modelem vybrané sázkové kanceláře (zprostředkovaně pomocí získaných vypsání kurzů), s původním modelem od FiveThirtyEight, založeným na Elo ratingu a pro nejnovější data také s aktuálně používaným predikčním modelem od FiveThirtyEight, modelem RAPTOR. Dále bude po vzoru článku [15] analyzováno, zda je v jednotlivých sezonách z predikční sady dat možné vybrat utkání, na která vsadit u sázkové kanceláře tak, aby bylo možné díky výběru vhodných zápasů – na základě odhadů pravděpodobností vítězství podle našeho modelu – dosáhnout zisku na konci sezony. Nakonec budou „vypsány“ decimální kurzy podle odhadnutých pravděpodobností naším modelem a budou generováni „náhodní sázející“, díky čemuž bude možné srovnání našeho modelu z pohledu sázkové kanceláře s modelem vybrané sázkové kanceláře.

Kapitola 2

Popis získaných dat

2.1 Průběh soutěže a utkání

Tato práce je zaměřena pouze na National Basketball Association (dále jen NBA), což je nejvyšší mužská profesionální basketbalová liga v Severní Americe, která je považována za nejprestižnější basketbalovou soutěž na světě. NBA se účastní celkem 30 týmů (29 z USA a 1 z Kanady), které jsou rozděleny do 2 konferencí (východní a západní) po 15 týmech, viz tabulka 2.1. Každá konference je pak rozdělena do 3 divizí po 5 týmech. Během běžné základní části sezony odehrává každý tým 2 utkání s týmem z opačné konference, a to jedno domácí a jedno venkovní (celkem tedy 30 utkání). S týmy ze stejné divize odehrává každý tým 4 utkání, a to 2 domácí a 2 venkovní (celkem tedy 16 utkání). S týmy z ostatních dvou divizí odehrává každý tým 3 nebo 4 utkání. Celkově tak každý tým odehrává 82 utkání během základní části sezony, která probíhá obvykle od října do dubna. Následně se z každé konference účastní 8 nejlepších týmů po základní části play-off v rámci dané konference. Každé kolo play-off se hraje na 4 vítězná utkání. Vítězové play-off ze západní a východní konference se pak utkávají ve finále NBA, hraném opět na 4 vítězná utkání.

Základní hrací doba každého utkání trvá 48 minut, konkrétně 4 čtvrtiny po 12 minutách. Tým, který po uplynutí této hrací doby získá více bodů, je vítězem utkání. V případě nerozhodného stavu následuje prodloužení trvající 5 minut. Stejně dlouhá prodloužení se opakují, dokud jeden z týmů nezíská po daných 5 minutách více bodů a stane se tak vítězem utkání. Vítězství (resp. prohry) v případném prodloužení mají v soutěži stejnou váhu jako v základní hrací době, resp. nemají vliv na ohodnocení obou týmů, jako je tomu například v hokejových nebo fotbalových soutěžích, kde se obvykle v případě nerozhodného stavu po základní hrací době rozdělují body mezi oba soupeře. Na základě toho nebudou v této práci rozlišována vítězství v základní hrací době a v prodloužení.

Tabulka 2.1: Aktuální názvy všech současných týmů NBA

Zkratka	Název týmu	Zkratka	Název týmu
ATL	Atlanta Hawks	DAL	Dallas Mavericks
BKN	Brooklyn Nets	DEN	Denver Nuggets
BOS	Boston Celtics	GSW	Golden State Warriors
CLE	Cleveland Cavaliers	HOU	Houston Rockets
DET	Detroit Pistons	LAC	Los Angeles Clippers
CHA	Charlotte Hornets	LAL	Los Angeles Lakers
CHI	Chicago Bulls	MEM	Memphis Grizzlies
IND	Indiana Pacers	MIN	Minnesota Timberwolves
MIA	Miami Heat	NOP	New Orleans Pelicans
MIL	Milwaukee Bucks	OKC	Oklahoma City Thunder
NYK	New York Knicks	PHX	Phoenix Suns
ORL	Orlando Magic	POR	Portland Trail Blazers
PHI	Philadelphia 76ers	SAC	Sacramento Kings
TOR	Toronto Raptors	SAS	San Antonio Spurs
WAS	Washington Wizards	UTA	Utah Jazz

(a) Východní konference

(b) Západní konference

2.2 Data

Historická data výsledků utkání byla získána z oficiálních webových stránek NBA [34]. Tento zdroj je tak považován za nejlepší možný. Pro každé jednotlivé utkání byla pro oba soupeřící týmy získána data obsahující informace o výsledných zápasových statistikách, které jsou vysvětleny v tabulce 2.2. Data byla získána ze všech utkání základních částí sedmi sezon, a to od počátku sezony 2012/2013 do konce sezony 2018/2019. Tato doba byla považována za dostatečně reprezentativní vzorek a rozsah těchto dat vyhovuje metodám použitým v této práci. Jak již bylo uvedeno v odstavci 2.1, každý tým odehrává během základní části sezony 82 utkání, což znamená celkově 1 230 utkání za jednu sezonu. Výjimkou byla sezona 2012/2013, ve které se v závěru sezony neodehrálo utkání mezi Bostonem Celtics a Indianou Pacers v důsledku teroristického útoku na Bostonský maraton, což bylo jediné neodehrané utkání v období získaných dat. Celkový počet utkání, ze kterých byla získána data, je tak 8 609. Důležitou úpravou dat bylo sjednocení názvů všech týmů, jelikož se některé názvy během pozorovaného období změnily. Konkrétně se jedná o přejmenování týmu New Orleans Hornets na jeho aktuální název New Orleans Pelicans od sezony 2013/2014 a přejmenování týmu Charlotte Bobcats na jeho aktuální název Charlotte Hornets od sezony 2014/2015. V tabulce 2.3 je ukázka formátu získaných a upravených dat z pěti utkání. Kompletní data jsou v příloženém souboru `NBAstats.xlsx` v příloze A.1.

Tabulka 2.2: Získané zápasové statistiky a jejich vysvětlivky

Zkratka	Význam	Zkratka	Význam
TEAM	Zkratka názvu týmu	3P%	Úspěšnost 3bodových střel
OPP	Zkratka názvu týmu soupeře	FTM	Počet úspěšných trestných hodů
H/A	Domácí/venkovní prostředí	FTA	Počet trestných hodů celkem
DATE	Datum utkání	FT%	Úspěšnost trestných hodů
W/L	Výhra/prohra	OREB	Počet útočných doskoků
PTS	Celkový počet získaných bodů	DREB	Počet obranných doskoků
+/-	Rozdíl počtu získaných bodů	REB	Počet doskoků celkem
FGM	Počet proměněných střel z pole	AST	Počet asistencí
FGA	Počet pokusů střel z pole	STL	Počet zisků
FG%	Úspěšnost střel z pole	BLK	Počet bloků
3PM	Počet úspěšných 3bodových střel	TOV	Počet ztrát
3PA	Počet pokusů 3bodových střel	PF	Počet osobních faulů

Kromě zápasových statistik byla získána také data otevíracích decimálních kurzů na vítězství týmů v jednotlivých utkáních do rozhodnutí (ODDS). Získané decimální kurzy byly přiřazeny k jednotlivým utkáním, viz tabulka 2.3 (kompletní data v příloze A.1). Použity byly kurzy sázkové kanceláře Bet365, která je jednou z největších online sázkových kanceláří na světě. Kurzy byly získány z webové stránky `indatabet.com` [33], kde jsou shromažďovány historické kurzy více sázkových kanceláří. Chybějící kurzy sázkové kanceláře Bet365 byly dohledány ze serveru `oddsportal.com` [35], kde byla rovněž namátkově ověřena správnost získaných kurzů. Z pozorovaného období nebyly dostupné otevírací kurzy sázkové kanceláře Bet365 pouze ze tří utkání. K těmto utkáním byly dohledány kurzy sázkové kanceláře Pinnacle rovněž ze zdroje `oddsportal.com` [35] a jsou barevně vyznačeny v příloženém souboru A.1. Kompletní data ze všech sledovaných sezon byla importována do softwaru MATLAB R2018a a jsou k nalezení v příloženém souboru `NBAstats.mat`, viz příloha A.2.

Tabulka 2.3: Ukázka formátu vstupních dat (pět utkání, každé utkání na dva řádky)

TEAM	H/A	OPP	GAME DATE	W/L	PTS	FGM	FGA	FG%	3PM	3PA	3P%	FTM	FTA	FT%	OREB	DREB	REB	AST	STL	BLK	TOV	PF	+/-	ODDS
DAL	H	SAS	13.4.2016	L	91	31	77	40,3	12	39	30,8	17	21	81,0	11	30	41	20	8	2	12	23	-5	1,47
SAS	A	DAL	13.4.2016	W	96	35	74	47,3	6	23	26,1	20	26	76,9	8	34	42	24	8	3	15	19	5	2,75
DEN	A	POR	13.4.2016	L	99	40	94	42,6	7	25	28,0	12	18	66,7	20	33	53	24	10	6	20	19	-8	4,75
POR	H	DEN	13.4.2016	W	107	39	89	43,8	11	36	30,6	18	22	81,8	13	32	45	19	11	8	17	19	8	1,20
UTA	A	LAL	13.4.2016	L	96	39	83	47,0	9	30	30,0	9	15	60,0	8	30	38	27	10	1	14	16	-5	1,41
LAL	H	UTA	13.4.2016	W	101	41	85	48,2	6	25	24,0	13	15	86,7	8	39	47	19	6	3	14	17	5	3,00
PHX	H	LAC	13.4.2016	W	114	45	99	45,5	7	26	26,9	17	22	77,3	18	40	58	26	10	5	17	15	9	1,47
LAC	A	PHX	13.4.2016	L	105	43	96	44,8	9	31	29,0	10	12	83,3	6	34	40	23	9	6	15	20	-9	2,75
MEM	A	GSW	13.4.2016	L	104	42	98	42,9	8	24	33,3	12	16	75,0	13	27	40	25	12	1	13	18	-21	21,00
GSW	H	MEM	13.4.2016	W	125	46	87	52,9	20	47	42,6	13	16	81,3	12	39	51	35	7	7	17	14	21	1,00

Kapitola 3

Elo rating

3.1 Historický vývoj

Elo rating je systém hodnocení hráčů na základě porovnání jejich relativní síly¹ ve hrách s nulovým součtem. Autorem tohoto systému byl maďarský fyzik Arpad Elo (1903–1992). Jakožto hráč šachu a účastník United States Chess Federation (USCF) navrhl Arpad Elo tento systém původně právě pro hodnocení hráčů šachu.

Předchůdcem Elo ratingu byl tzv. Harkness systém navržený Kennethem Harknessem [11], který USCF používala od roku 1950 do roku 1960. Tento systém hodnotil hráče na základě průměrného ratingu jeho soupeřů v turnaji. Pokud hráč dosáhl v turnaji 50% skóre, jeho nový rating bude průměr ratingů jeho soupeřů. Pokud dosáhl vyššího než 50% skóre, k průměrnému ratingu jeho soupeřů se navíc přičte 10 bodů za každý jeden procentní bod nad 50%. Pokud naopak dosáhne nižšího skóre, za každý procentní bod pod 50% se mu od průměrného ratingu soupeřů 10 bodů odečte. S ohledem na to, že Harkness systém nebyl považován za příliš přesný [29], byl Arpad Elo požádán USCF, aby model vylepšil. Nový systém navrhl a upravil tak, aby se hodnocení hráčů příliš nelišilo od čísel, na která byli zvyklí, resp. tak, aby byly zachovány klasifikační třídy hráčů, viz tabulka 3.1, a popsal ho v knize [6]. Mimo jiné uvedl, že:

„Proces hodnocení hráčů lze přirovnat k měření pozice korku houpajícího se nahoru a dolu na hladině rozbourené vody pomocí metru uvázaného na šňůře, která se kymácí ve větru.“ – Arpad Elo, Chess Life, 1962

Tímto slavným citátem naznačil, jak obtížné a nepřesné je ratingové hodnocení hráčů, a to nejen jeho novým systémem.

V roce 1970 byl Elo systém přijat také Mezinárodní šachovou federací (FIDE), která jej používá dodnes. USCF v průběhu let upravovala a vylepšovala výpočet ratingu, nyní používá systém ovlivněný Glicko ratingem, jehož nejvýznamnější změnou je popis každého hráče třemi parametry: rating, odchylka ratingu a volatilita [9]. Postupně se Elo rating začal rozšiřovat a využívat také v ostatních hrách, kde se střetávají dva hráči, resp. týmy, a lze je porovnávat na základě ratingu, jako je fotbal, tenis, basketbal nebo vi-

¹Relativní síla znamená sílu určitého hráče vzhledem k síle jiného hráče.

deohry. Tato práce je zaměřena na použití Elo rating pro basketbal, proto od této chvíle bude používán termín *tým* ve smyslu basketbalového týmu, jako analogie k *hráči* šachu.

Tabulka 3.1: Výkonnostní třídy USCF

Kategorie	Rating	Kategorie	Rating
Senior master	2 400 a více	Class E	1 000–1 199
National master	2 200–2 399	Class F	800–999
Expert	2 000–2 199	Class G	600–799
Class A	1 800–1 999	Class H	400–599
Class B	1 600–1 799	Class I	200–399
Class C	1 400–1 599	Class J	100–199
Class D	1 200–1 399		

3.2 Matematický popis

Elo rating je v současnosti velmi populární v mnoha sportech především kvůli jeho poměrně jednoduché matematické stránce. Základní myšlenkou celého systému Elo rating je měření výkonu týmu nikoliv absolutně pro jednotlivé týmy, ale odvozeně od vítězství nebo prohry s ostatními týmy, tedy relativně vzhledem ratingu soupeřů. Rating je číslo reprezentující sílu, resp. kvalitu týmu. Tým s vyšší hodnotou ratingu je tak považován za „lepší“ a má tudíž vyšší pravděpodobnost ve vzájemném utkání s jiným týmem zvítězit. Elo rating se vyvíjí v čase a přizpůsobuje se výkonům daných týmů, přičemž je jeho nová hodnota přepočítávána po každém odehraném utkání, a to pouze pro dvojici soupeřících týmů. Rating se zvyšuje, pokud týmy vítězí a snižuje se, pokud týmy prohrávají. Výhodou je to, že hodnota ratingu se zvýší vítěznému týmu po odehraném utkání o stejnou hodnotu, o kterou se rating sníží týmu poraženému, přičemž tato hodnota závisí právě na hodnotě ratingu soupeře. Jinými slovy vítězstvím nad silnějším týmem se zvýší rating vítězného týmu více, než vítězstvím nad týmem slabším a naopak, viz příklad 3.3. Elo rating tedy automaticky zohledňuje tzv. *strength of schedule* (sílu harmonogramu), tedy tým, který porazí silnější soupeře, by měl skončit s vyšším ratingem. Nevýhodou ratingových metod je to, že jedno číslo (rating) reprezentuje všechny možné faktory, které daný tým a utkání ovlivňují. Slabší týmy občas porázejí silnější a naopak, tým může mít špatný den nebo naopak, někteří hráči týmu mohou být zranění, hru může také ovlivnit výkon trenérů nebo rozhodčích, unavenost, domácí prostředí nebo vzdálenost od něj atd. Z toho důvodu je potřeba na výkonnost týmů nahlížet jako na náhodnou veličinu. Určením rozdělení těchto náhodných veličin můžeme odhadnout pravděpodobnost vítězství týmu v utkání, přestože se spolu daní soupeři ještě nikdy předtím neutkali, což je další důležitou výhodou Elo ratingu.

3.2.1 Normální rozdělení výkonnosti

Arpad Elo ve své práci [6] předpokládal, že výkonnost jednotlivých hráčů šachu se řídí normálním rozdělením² a provedl rozsáhlé studie, aby toto tvrzení potvrdil. V následující části se budeme věnovat využití normálního rozdělení pro odhad pravděpodobnosti vítězství na základě Elo ratingu. Při popisu budeme vycházet ze zdroje [4].

Jako střední hodnota normálního rozdělení byla uvažována hodnota ratingu daného týmu a jako směrodatná odchylka velikost jedné třídy (viz tabulka 3.1), tedy 200. Označme X_i náhodnou veličinu popisující výkonnost týmu i a X_j náhodnou veličinu popisující výkonnost týmu j . Rating týmu i , resp. j , označme r_i , resp. r_j . Potom

$$X_i \sim N(\mu_i = r_i, \sigma_i = 200), \quad (3.1)$$

$$X_j \sim N(\mu_j = r_j, \sigma_j = 200). \quad (3.2)$$

Nyní lze odhadnout pravděpodobnost, že tým i porazí tým j , jako

$$P_{ij} = P(X_i > X_j) = P(X_i - X_j > 0) = 1 - F_{X_i - X_j}(0), \quad (3.3)$$

stejně jako pravděpodobnost, že tým j porazí tým i , tedy

$$P_{ji} = P(X_i < X_j) = P(X_i - X_j < 0) = F_{X_i - X_j}(0) = 1 - P_{ij}, \quad (3.4)$$

kde $F_{X_i - X_j} = F_X$ je distribuční funkce náhodně veličiny $X_i - X_j = X$. Nyní je potřeba určit rozdělení této veličiny X . Za předpokladu, že se náhodné veličiny X_i a X_j řídí normálním rozdělením a jsou vzájemně nezávislé, platí obecně, že jejich rozdíl se řídí normálním rozdělením

$$X \sim N(\mu = \mu_i - \mu_j, \sigma^2 = \sigma_i^2 + \sigma_j^2), \quad (3.5)$$

což je dokázáno např. v [14]. Střední hodnota tohoto rozdělení se tedy rovná hodnotě $\mu = \mu_i - \mu_j = r_i - r_j$ a směrodatná odchylka $\sigma = \sqrt{\sigma_i^2 + \sigma_j^2}$. Protože $\sigma_i = \sigma_j = 200$, dostáváme $\sigma = \sqrt{2\sigma_i^2} = \sqrt{2}\sigma_i = \sqrt{2} \cdot 200 \doteq 282,84$. Dostáváme tedy, že

$$X \sim N(\mu = r_i - r_j, \sigma = 282,84). \quad (3.6)$$

Pravděpodobnost, že tým i porazí v utkání tým j lze odhadnout (viz rovnost 3.3) jako $P_{ij} = 1 - F_X(0)$, což lze pro zjednodušení přepsat jako

$$P_{ij} = F_Y(r_i - r_j), \quad (3.7)$$

$$P_{ji} = 1 - P_{ij}, \quad (3.8)$$

kde F_Y je distribuční funkce náhodné veličiny

$$Y \sim N(\mu = 0, \sigma = 282,84). \quad (3.9)$$

²Více o normálním rozdělení např. v [12].

Pravděpodobnost výhry týmu i nad týmem j lze tedy určit jako distribuční funkci normálního rozdělení náhodné veličiny Y v bodě rozdílu jejich ratingů $r_i - r_j$. Na obrázku 3.1 je vykreslena distribuční funkce $F_Y(r_i - r_j)$ náhodné veličiny Y , tedy je zde graficky znázorněna pravděpodobnost vítězství týmu i nad týmem j v závislosti na rozdílu jejich ratingů. Pro srovnání jsou v tabulce 3.2 dopočítány pravděpodobnosti pro některé vybrané hodnoty rozdílů ratingů. Na základě znalosti ratingů dvou soupeřů lze tedy odhadnout pravděpodobnost vítězství v jakémkoliv utkání, přestože spolu dvojice týmů nikdy předtím nemusela sehrát žádné utkání. Toto lze rovněž použít v soutěžích pro jednotlivce nebo pro více než dva soupeře, a to s použitím průměrného ratingu soupeřů [4].

3.2.2 Logistické rozdělení výkonnosti

Elo také uvedl [6], že kromě normálního rozdělení pravděpodobnosti náhodné veličiny výkonnosti hráče, může být použito také logistické rozdělení³. Na tuto změnu přešla USCF i díky v té době zvýšené dostupnosti počítačů. Sledováním velkého počtu výsledků zjistila, že použitím logistického rozdělení dosahují nejpřesnějších výsledků [23]. FIDE rovněž používá ve svém hodnocení aproximaci logistického rozdělení [30]. Jedná se tedy o nahrazení distribuční funkce normálního rozdělení distribuční funkcí logistického rozdělení. Toto rozdělení se používá např. v logistické regresi a obecný tvar jeho distribuční funkce je

$$L(x; \mu, s) = \frac{1}{1 + e^{-\frac{(x-\mu)}{s}}}. \quad (3.10)$$

Střední hodnota je stejně jako v předchozím odstavci ve výrazu (3.9) uvažována $\mu = 0$. Logistická křivka se ale pro Elo rating obvykle vyskytuje se základem 10, namísto e . Tuto úpravu lze formálně zapsat pomocí volby parametru $s = \frac{400}{\ln 10}$, kde hodnota 400 je použita opět, podobně jako v odstavci 3.2.1, kvůli zachování velikosti jedné třídy [20], viz tabulka 3.1. Obecně jsou čísla 10 a 400 libovolná. Jejich význam je takový, že pro každých 400 ratingových bodů, o kterých má tým i více než tým j , se zvýší šance na vítězství týmu i 10krát – což je ukázáno v příkladu 3.2 – přičemž šanci na vítězství hráče i definujeme jako $o_i = \frac{P_{ij}}{1-P_{ij}}$. Po dosazení hodnot μ a s do (3.10) dostáváme logistickou funkci $L(x)$ ve tvaru

$$L(x) = \frac{1}{1 + 10^{-\frac{x}{400}}}. \quad (3.11)$$

Pravděpodobnost, že tým i porazí v utkání tým j , pak lze určit podobně jako v odstavci 3.2.1, tedy dosazením rozdílu $r_i - r_j$ do logistické funkce $L(x)$, tedy

$$P_{ij} = L(r_i - r_j) = \frac{1}{1 + 10^{-\frac{r_i - r_j}{400}}}, \quad (3.12)$$

a pravděpodobnost, že tým j porazí tým i jako

$$P_{ji} = 1 - P_{ij} = 1 - L(r_i - r_j) = 1 - \frac{1}{1 + 10^{-\frac{r_i - r_j}{400}}} = \frac{1}{1 + 10^{-\frac{r_j - r_i}{400}}}. \quad (3.13)$$

³Více o logistickém rozdělení např. v [2].

Získáváme tak běžně používaný vzorec pro výpočet pravděpodobnosti vítězství z Elo ratingu. Na obrázku 3.1 je vykreslena distribuční (logistická) funkce $L(r_i - r_j)$, tedy je zde graficky znázorněna pravděpodobnost vítězství týmu i nad týmem j v závislosti na rozdílu jejich ratingů. Zachycen je zde rozdíl v získané pravděpodobnosti s použitím normálního rozdělení s distribuční funkcí F_Y a logistického rozdělení s distribuční funkcí L . V tabulce 3.2 jsou pak pro vybrané rozdíly ratingů dopočítány pro srovnání pravděpodobnosti vítězství. Konkrétní postup výpočtu pravděpodobnosti je v příkladu 3.1.

Příklad 3.1. Uvažujme týmy i a j , které se mají utkat. Tým i má aktuální rating $r_i = 1450$. Tým j má aktuální rating $r_j = 1600$. Odhad pravděpodobnosti vítězství týmu i v utkání proti týmu j určíme jako

$$P_{ij} = L(r_i - r_j) = L(1450 - 1600) = \frac{1}{1 + 10^{-\frac{1450-1600}{400}}} = 29,66\%.$$

Pravděpodobnost vítězství týmu j proti týmu i je pak

$$P_{ji} = 1 - P_{ij} = \frac{1}{1 + 10^{-\frac{1600-1450}{400}}} = 70,34\%.$$

Pravděpodobnost vítězství týmu i je 29,66% a pravděpodobnost vítězství týmu j je 70,34%.

Příklad 3.2. Uvažujme týmy i a j , které mají shodný rating. Rozdíl jejich ratingů je tedy $r_i - r_j = 0$. Pravděpodobnost vítězství ve vzájemném utkání obou těchto týmů je

$$P_{ij} = P_{ji} = L(r_i - r_j) = L(0) = \frac{1}{1 + 10^{-\frac{0}{400}}} = 50\%.$$

Šanci o_i , že zvítězí tým i , lze vypočítat jako

$$o_i = \frac{P_{ij}}{1 - P_{ij}} = \frac{0,5}{1 - 0,5} = 1.$$

Šance, že zvítězí tým i je tedy 1:1.

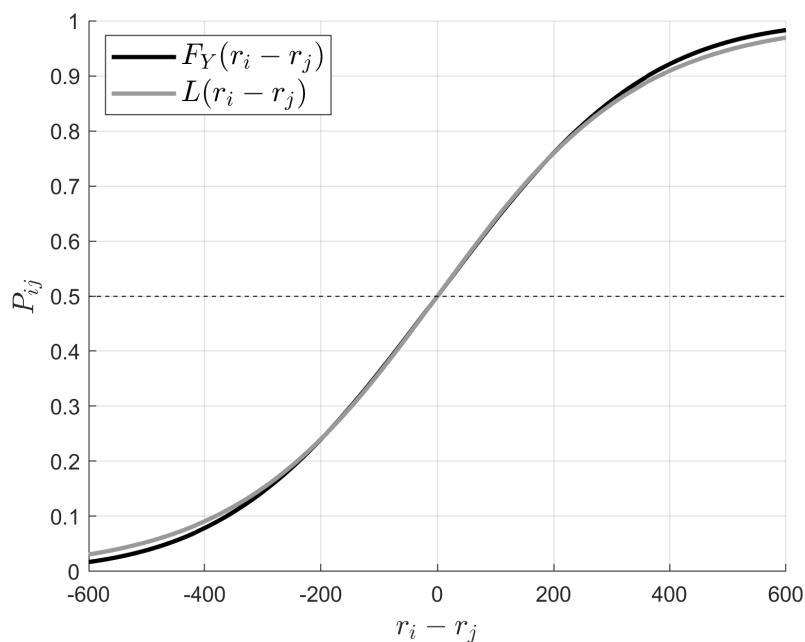
Nyní uvažujme, že rozdíl mezi ratingy týmů i a j je $r_i - r_j = 400$. Pravděpodobnost, že zvítězí tým i je

$$P_{ij} = L(r_i - r_j) = L(400) = \frac{1}{1 + 10^{-\frac{400}{400}}} = \frac{1}{1 + \frac{1}{10}} = 90,91\%,$$

Šanci o_i , že zvítězí tým i , vypočítáme opět jako

$$o_i = \frac{P_{ij}}{1 - P_{ij}} = \frac{0,9091}{1 - 0,9091} = 10.$$

Šance, že zvítězí tým i je tedy 10:1. Závěr je takový, že zvýší-li se rozdíl ratingů o 400 ve prospěch týmu i , šance na jeho vítězství v utkání se zvýší 10krát. Pokud se naopak rating o 400 sníží, šance na vítězství je 10krát menší. Toto platí pro jakýkoliv rozdíl mezi ratingy a jeho změnu o 400 bodů.



Obrázek 3.1: Pravděpodobnost vítězství v utkání v závislosti na rozdílu ratingů

Tabulka 3.2: Pravděpodobnost vítězství v utkání v závislosti na rozdílu ratingů

$r_i - r_j$	Pravděpodobnost	$r_i - r_j$	Pravděpodobnost
-600	1,69 %	-600	3,07 %
-400	7,86 %	-400	9,09 %
-300	14,44 %	-300	15,10 %
-200	23,98 %	-200	24,03 %
-150	29,79 %	-150	29,66 %
-100	36,18 %	-100	35,99 %
-75	39,54 %	-75	39,37 %
-50	42,98 %	-50	42,85 %
-25	46,48 %	-25	46,41 %
0	50,00 %	0	50,00 %
25	53,52 %	25	53,59 %
50	57,02 %	50	57,15 %
75	60,46 %	75	60,63 %
100	63,82 %	100	64,01 %
150	70,21 %	150	70,34 %
200	76,02 %	200	75,97 %
300	85,56 %	300	84,90 %
400	92,14 %	400	90,91 %
600	98,31 %	600	96,93 %

(a) Normální rozdělení

(b) Logistické rozdělení

3.2.3 Přepočet ratingu

V předchozím odstavci bylo ukázáno, jak na základě Elo ratingu vypočítat pravděpodobnost vítězství jednotlivých týmů v daném utkání. Bylo zde také uvedeno, že rating se zvyšuje pokud týmy vítězí a snižuje se, pokud týmy prohrávají. Otázkou ale zůstává jak, o kolik a jakou zde hrají roli získané odhady pravděpodobnosti vítězství. Této otázce se budeme věnovat v následujících částech, ve kterých budeme vycházet ze zdrojů [19] a [28].

Přístup Elo ratingu je velmi jednoduchý. Základní myšlenka je porovnat očekávané výsledky utkání se skutečnými výsledky. Pokud skutečný výsledek překročí očekávání, je to indikátor toho, že rating týmu je příliš nízký a je potřeba ho zvýšit. Naopak pokud výsledek týmu nenaplní očekávání, je potřeba rating snížit. Velikost úpravy ratingu je úměrná tomu, v jaké míře tým překročil, nebo nedosáhl očekávaného výsledku. Pokud tedy v utkání zvítězí favorit, zvýší se jeho rating, resp. sníží se rating poraženého týmu méně, než pokud v utkání zvítězí dle očekávání slabší tým, viz příklad 3.3. Nejprve je potřeba stanovit počáteční rating, který může být pro všechny týmy na začátku sledovaného období stejný. Pokud je však cílem sledovat historický vývoj ratingu od začátku působení daného týmu v soutěži, existuje několik metod, pomocí kterých lze odhadnout počáteční rating na základě úspěšnosti v několika prvních utkáních a ratingu soupeřů z těchto utkání, více v [20]. Vzorec tzv. průběžné metody pro přepočet ratingu týmu i je

$$r_{i_{new}} = r_{i_{old}} + K \cdot (S_{ij} - P_{ij}), \quad (3.14)$$

kde je

$r_{i_{new}}$	nový rating týmu i ,
$r_{i_{old}}$	předchozí rating týmu i ,
K	koeficient rozvoje (podrobněji v odstavci 3.2.4)
S_{ij}	skutečný výsledek týmu i v utkání proti týmu j ,
P_{ij}	očekávaný výsledek týmu i v utkání proti týmu j .

Tento přepočet ratingu (resp. jeho aktualizace) může být proveden po každém odehraném utkání, nebo po určitém počtu odehraných utkání, což se používá např. v šachu v případě vícekolového turnaje. Hodnota S_{ij} , tedy skutečný výsledek týmu i v utkání proti týmu j , může v našem případě (pro basketbal) nabývat pouze hodnot

$$S_{ij} = \begin{cases} 1, & \text{pokud tým } i \text{ zvítězil nad týmem } j, \\ 0, & \text{pokud tým } i \text{ prohrál s týmem } j. \end{cases} \quad (3.15)$$

V případě her, kde je uvažována také remíza, může S_{ij} nabývat také hodnoty 0,5. Pokud dochází k přepočtu ratingu po určitém počtu odehraných utkání, může tato hodnota reprezentovat součet výsledků ze všech proběhlých utkání. Hodnota P_{ij} , tedy očekávaný výsledek, představuje pravděpodobnost vítězství odvozenou v odstavci 3.2.2. Například pokud má tým pravděpodobnost výhry 70 %, očekává se, že v utkání získá

0,7 bodu. Pokud by v utkání skutečně zvítězil, překonal by tak očekávaný výsledek o $S_{ij} - P_{ij} = 1 - 0,7 = 0,3$ bodu. Pokud by prohrál, nedosáhne očekávání o $S_{ij} - P_{ij} = 0 - 0,7 = -0,7$ bodu. Jelikož $P_{ij} + P_{ji} = 1$ je zřejmé, že rating vítězného týmu se zvýší vždy o stejnou hodnotu, o kterou se rating poraženého týmu sníží, což je ilustrováno v příkladu 3.3. Opět, pokud k přepočtu ratingu dochází po více utkáních, může být namísto hodnoty P_{ij} použit součet všech těchto pravděpodobností, resp. očekávaných výsledků.

3.2.4 Koeficient rozvoje K

Koeficient rozvoje K , neboli tzv. K -faktor, je poslední nevysvětlenou proměnnou v uvedeném vzorci (3.14). K -faktor určuje maximální možnou hodnotu zvýšení, resp. snížení ratingu během jeho aktualizace. Tento faktor určuje to, jak moc výsledek každého nového utkání ovlivňuje změnu ratingu z poslední hodnoty $r_{i_{old}}$ na novou hodnotu $r_{i_{new}}$. Pro velké hodnoty K je rating citlivý na každé nové utkání a znamená to tedy velkou změnu ratingu po každé nové hře. Tedy pokud se týmu nevydaří jedno utkání tak, jak je jeho zvykem, může toto utkání příliš ovlivnit jeho rating. Naopak pokud K je příliš malé, model nebude schopný rychle reagovat na zlepšení, resp. zhoršení síly týmu a rating se tak nebude příliš měnit. Pro $K = 0$ by se rating neměnil nikdy.

Vhodná volba koeficientu K je velmi důležitá, ale také složitá. Tato hodnota se může lišit pro různé sporty, ale také existují různé způsoby její volby. Obecně tato hodnota nemusí být konstantní. Například v ČR pro národní Elo rating šachových hráčů závisí K na věku hráčů. Pro hráče starší než 20 let je $K = 15$, ale pro hráče ve věku 18–20 je používána hodnota $K = 20$, protože se předpokládá, že u mladších hráčů nastávají větší změny ve výkonnosti [36]. Ve FIDE nebo USCF má ze stejného důvodu vliv také hodnota aktuálního ratingu [20]. V basketbalu může být hodnota K například na začátku sezony vyšší, aby se ratingy ovlivněné předchozími sezonami rychleji přizpůsobily aktuální výkonnosti týmů, která může být výrazně ovlivněna změnami v týmech mezi sezonami. Naopak v průběhu sezony se může koeficient K postupně snižovat z důvodu stabilizace síly týmů. Hodnota K může být také závislá například na důležitosti utkání, nebo na výsledném skóre [26]. Obecně platí, že koeficient rozvoje nemusí být konstantní, ale může být funkcí několika proměnných, viz odstavce 4.4 a 4.5.

Příklad 3.3. Uvažujme týmy i a j , které se mají utkat. Tým i má aktuální rating $r_{i_{old}} = 1450$. Tým j má aktuální rating $r_{j_{old}} = 1600$. Pravděpodobnost vítězství obou týmů byla odhadnuta v příkladu 3.1. Pravděpodobnost, že v utkání zvítězí tým i je $P_{ij} = 29,66\%$ a pravděpodobnost, že zvítězí tým j je $P_{ji} = 70,34\%$. Uvažujme koeficient rozvoje $K = 20$. Provedme aktualizaci Elo ratingu na základě výsledku tohoto utkání pro následující dvě situace.

1. V utkání zvítězil tým i .

$$r_{i_{new}} = r_{i_{old}} + K \cdot (S_{ij} - P_{ij}) = 1\,450 + 20 \cdot (1 - 0,2966) = 1\,464,07$$

$$r_{j_{new}} = r_{j_{old}} + K \cdot (S_{ji} - P_{ji}) = 1\,600 + 20 \cdot (0 - 0,7034) = 1\,585,93$$

2. V utkání zvítězil tým j .

$$r_{i_{new}} = r_{i_{old}} + K \cdot (S_{ij} - P_{ij}) = 1\,450 + 20 \cdot (0 - 0,2966) = 1\,444,07$$

$$r_{j_{new}} = r_{j_{old}} + K \cdot (S_{ji} - P_{ji}) = 1\,600 + 20 \cdot (1 - 0,7034) = 1\,605,93$$

Je zřejmé, že v první variantě, kdy zvítězil „slabší“ tým i (s nižší odhadnutou pravděpodobností vítěství), se jeho rating po utkání zvýšil o 14,07 bodu, tedy o stejnou hodnotu, o kterou se snížil rating poraženého týmu j . Ve druhé variantě, kdy zvítězil dle očekávání favorit utkání tým j , se jeho rating zvýšil pouze o 5,93 bodu a opět pouze o stejnou hodnotu se snížil rating poraženého týmu i . Ilustrován je zde fakt, že (jak již bylo uvedeno v odstavci 3.2.3) pokud v utkání zvítězí dle odhadnuté pravděpodobnosti očekávaný favorit, zvýší se jeho rating, resp. sníží se rating poraženého týmu méně, než pokud v utkání zvítězí dle odhadnuté pravděpodobnosti slabší tým. Pokud by oba soupeři měli před utkáním stejný rating, tedy odhadnutá pravděpodobnost vítězství obou týmů by byla 50 %, hodnota, o kterou by se zvýšil rating vítězi, resp. snížil rating poraženého, by byla shodná nezávisle na tom, který z týmů by zvítězil. Z příkladu je také zřejmé, že zvýšení ratingu vítěze a snížení ratingu poraženého je vždy o stejnou hodnotu, tedy změna ratingu mezi vítězem a poraženým je symetrická. To znamená, že pokud bude počáteční rating pro všechny týmy v soutěži na začátku sledovaného období např. 1 500, průměrný rating v soutěži bude po celou dobu 1 500 (pokud zůstane počet týmů stejný), jelikož nové ratingové body nepřibývají ani neubývají, pouze se přesouvají od poražených týmů k vítězným.

Kapitola 4

Elo rating pro NBA

Jedním z hlavních cílů této práce bylo navrhnout změny, resp. vylepšení modelu používajícího Elo rating pro predikci výsledků utkání NBA. V této kapitole budou uvedeny všechny uvažované modely v této práci. Popsány zde budou modely od základní verze Elo ratingu (viz kapitola 3) po komplikovanější modely zahrnující více navržených, nebo převzatých parametrů. V kapitole 5 pak budou vybrané parametry optimalizovány a bude vybrán jeden vhodný model, který bude použit k predikci a fiktivnímu použití proti sázkové kanceláři v kapitole 6, čímž zprostředkovaně ověříme jeho predikční schopnost a to, jak by si vedl ve srovnání s modely, které běžně používají sázkové kanceláře. Parametry budou přidávány tak, aby byla pokud možno zachována jednoduchost modelu Elo rating, což je jedna z jeho největších výhod. Všechny popsané modely byly implementovány v softwaru MATLAB R2018a. K získání odhadů pravděpodobností vítězství v jednotlivých utkáních, pomocí uživatelsky zvoleného modelu a hodnot jeho parametrů, byla naprogramována funkce `elo.m`, viz příloha A.3.

4.1 Model 1

Prvním implementovaným modelem byla nejzákladnější varianta Elo ratingu. Na začátku sledovaného období byla stanovena hodnota $R_0 = 1500$ jako počáteční rating pro všechny týmy. Tato hodnota je libovolná, pouze určuje úroveň ratingu a průměrnou hodnotu ratingu po celé sledované období. Jelikož se týmy v soutěži během sledovaného období neměnily (žádný nový tým v soutěži nepřibyl, ani soutěž neopustil), průměrný rating po celou dobu je právě 1500, viz příklad 3.3. Změnou hodnoty R_0 by byly dosaženy totožné výsledky, jelikož při odhadu pravděpodobnosti vítězství v utkání a následném přepočtu ratingu nezáleží na úrovni ratingu dvou soupeřících týmu, ale na jejich rozdílu, viz kapitola 3. Při odhadu pravděpodobnosti vítězství v utkání bylo předpokládáno logistické rozdělení výkonnosti, viz výraz (3.12). K přepočtu ratingu docházelo po každém odehraném utkání podle vzorce (3.14). V tomto modelu byl uvažován **konstantní koeficient rozvoje K** po celé sledované období, který bude optimalizován v kapitole 5. Oproti základní verzi Elo ratingu např. pro šach, je zde navíc použit pouze jediný parametr, a to **parametr p** . Tento parametr, inspirovaný Elo modelem od FiveThirtyEight [26], viz odstavec 4.6, udává poměr zachování ratingu po skončení sezony. V této práci bude koncem sezony vždy myšlen konec základní části sezony – play-off zde nebude vůbec uvažováno, a to z důvodu odlišnosti systému soutěže a dalších vnějších vlivů. Parametr

p upravuje rating každého týmu na začátku sezony tak, že částečně vrací jeho hodnotu k počátečnímu ratingu, resp. k průměru, a to z důvodu očekávané změny síly týmů mezi sezonami. Úpravu ratingu každého týmu i před začátkem nové sezony s lze zapsat jako

$$r_i^{(s)} = p \cdot r_i^{(s-1)} + (1 - p) \cdot R_0, \quad (4.1)$$

kde je

$r_i^{(s)}$	nový rating týmu i před začátkem nové sezony s ,
$r_i^{(s-1)}$	rating týmu i po skončení předešlé sezony $s - 1$,
p	parametr zachování ratingu po skončení sezony,
R_0	počáteční (průměrný) rating, v našem případě $R_0 = 1\,500$.

Stejně jako koeficient rozvoje K , bude i parametr p **uvažován konstantní** a bude optimalizován v kapitole 5.

4.2 Model 2

Tento model je rozšířením modelu 1 z odstavce 4.1. Rozšíření tohoto modelu, inspirované Elo modelem od FiveThirtyEight [26], viz odstavec 4.6, spočívá v přidání parametru H . Jedná se o parametr výhody domácího prostředí.

Hlavní příčinou výhody domácího prostředí může být např. podpora domácích diváků, zaujatost rozhodčích, cestování venkovního týmu, znalost prostředí atd. Předpoklad výhody domácího prostředí potvrdili např. už v roce 1977 ve svém článku Schwartz a Barsky [24]. Výhodou domácího prostředí se dodnes zabývá mnoho výzkumů. V této práci bylo pouze ověřeno na dostupných datech, tedy celkem 8 609 utkáních, že domácí týmy skutečně vítězí častěji. K ověření byl použit předpoklad, že pokud domácí výhoda neexistuje a všechny týmy hrají během sezony stejný počet utkání doma i venku, měly by domácí týmy zvítězit pouze v 50 % ze všech 8 609 utkání. Jinými slovy, počet vítězství domácích týmů by měl pocházet ze základního souboru s binomickým rozdělením $\text{Bi}(n = 8\,609; \pi_0 = 0,5)$, což bylo testováno. Použita byla stejná metoda jako v článku [22], kde nebyl brán ohled na splnění přesných předpokladů pro použití binomického rozdělení a testování jeho parametrů. Například zde nebude řešeno to, že pravděpodobnosti vítězství nejsou pro všechny týmy stejné a jsou závislé na síle soupeře, že některé týmy proti sobě nehrají v sezoně stejný počet utkání doma a venku, nebo že jsou pozorování závislá. Je proto potřeba považovat toto testování pouze jako přibližné, což je pro účel v této práci dostačující. Nechť X je náhodná veličina popisující počet vítězství domácích týmů $X \sim \text{Bi}(n = 8\,609; \pi)$. Na hladině významnosti $\alpha = 5\%$ testujeme hypotézu

$$H_0 : \pi = 0,5, \quad \text{proti alternativě} \quad H_1 : \pi > 0,5.$$

Odhad parametru π je relativní četnost vítězství domácích, kterých bylo celkem $X = 5\,056$, vzhledem ke všem utkáním, tedy $\hat{\pi} = \frac{X}{n} = \frac{5\,056}{8\,609} \doteq 0,59$. Pro velká n (doporučuje se $n > \frac{9}{\hat{\pi}(1-\hat{\pi})}$ [12]), lze použít aproximaci založenou na normálním rozdělení.

Testové kritérium je pak, viz [1],

$$T = \frac{X - n\pi_0}{\sqrt{n\pi_0(1 - \pi_0)}} = \frac{5\,056 - 8\,609 \cdot 0,5}{\sqrt{8\,609 \cdot 0,5(1 - 0,5)}} \doteq 16,20. \quad (4.2)$$

Kritická hodnota je $u_{1-\alpha} = u_{0,95} \doteq 1,64$. Protože $T = 16,20 > 1,64$, zamítáme na hladině významnosti $\alpha = 5\%$ hypotézu H_0 a přikláníme se k alternativě, že $\pi > 0,5$, tedy relativní četnost vítězství domácích týmů je vyšší než 50 %, resp. domácí týmy vítězí častěji a na prostředí tak záleží.

Parametr $H > 0$ určuje počet bodů Elo ratingu, který bude vždy navíc přičten k ratingu domácího týmu. Díky tomu se změní rozdíl v ratingu soupeřících týmů, což má vliv na odhad pravděpodobnosti vítězství v utkání. Při odhadu pravděpodobnosti vítězství v utkání bylo předpokládáno logistické rozdělení výkonnosti, viz výraz (3.12). Tento výraz byl rozšířen právě o přidaný parametr H . Odhad pravděpodobnosti P_{ij} vítězství týmu i v utkání proti týmu j lze tedy formálně zapsat jako

$$P_{ij} = \frac{1}{1 + 10^{-\frac{(r_i + h_i \cdot H) - (r_j + h_j \cdot H)}{400}}}, \quad (4.3)$$

kde $H > 0$ je parametr výhody domácího prostředí a h_i je indikátor domácího prostředí týmu i , tedy

$$h_i = \begin{cases} 1, & \text{pokud je tým } i \text{ domácí,} \\ 0, & \text{pokud je tým } i \text{ hostující.} \end{cases} \quad (4.4)$$

Parametr H byl uvažován konstantní a shodný pro všechny týmy. K přepočtu ratingu docházelo opět po každém odehraném utkání podle vzorce (3.14). Jako počáteční rating všech týmů byla zvolena opět hodnota $R_0 = 1\,500$. V tomto modelu byl rovněž uvažován **konstantní koeficient rozvoje K** po celé sledované období. Opět byl použit – stejně jako v odstavci 4.1 – také **konstantní parametr p** . Všechny tři parametry K , p a H , budou optimalizovány v kapitole 5.

4.3 Model 3

Tento model je rozšířením modelu 2 z odstavce 4.2. Rozšíření tohoto modelu spočívá v přidaném **parametru B** . Jedná se o parametr penalizace *back-to-back* utkání. Tento pojem označuje den po sobě jdoucí utkání, tedy utkání, kdy tým neměl ani jeden den odpočinku (odehraje dvě utkání ve dvou dnech).

Několikrát během sezony jsou týmy NBA nuceny hrát dvě po sobě jdoucí utkání. Z hlediska cestování může mít při *back-to-back* zápasech významnější vliv vzdálenost, resp. délka letu (a změna časových pásem), kterou musí týmy absolvovat – hostující i domácí. Významnou roli hraje také zdraví hráčů. Například nejvíce vytížení hráči (z hlediska počtu odehraných minut v utkání) nemusí být schopni odehrát dvě po sobě jdoucí utkání

se stejnou očekávanou výkonností a stejným počtem minut. Příkladem může být hráč Kawhi Leonard, který je považován za jednoho z nejlepších hráčů současnosti (NBA Finals MVP 2019 a NBA All-Star Game MVP 2020), a který aktuálně již 3. sezonu v řadě vůbec nenastupuje ze zdravotních důvodů do druhého z back-to-back utkání. NBA se za účelem eliminace vlivu back-to-back utkání snaží při sestavování rozpisů utkání snižovat počet výskytů back-to-back situací, což lze pozorovat na obrázku 4.1. Podobný vliv nemusí mít pouze back-to-back utkání, ale samozřejmě také celkový zápasový rozvrh týmu i soupeře z hlediska počtu dní odpočinku mezi zápasy. Vliv počtu dní odpočinku v NBA dokázali ve své práci např. Entine a Small [7]. Zde se budeme zabývat pouze back-to-back zápasy. Jejich vliv na vítězství v utkání byl v této práci opět pouze ověřen na dostupných datech, tedy celkem 8 609 utkáních.

K ověření byl použit tzv. *test homogeneity dvou binomických rozdělení* [1]. Stejně jako při ověřování domácí výhody v odstavci 4.2 zde nebyl brán ohled na splnění přesných předpokladů pro toto testování, proto je opět potřeba považovat výsledky pouze jako přibližné, což je pro účel v této práci dostačující. Jelikož back-to-back situacím jsou častěji vystavovány hostující týmy, byly díky této nevyváženosti testovány dvě varianty.

(1) V první variantě byla nejprve vybrána pouze utkání, ve kterých ani jeden z týmů nehrál back-to-back – takových utkání bylo $m = 5\,580$. Z těchto utkání vyhráli domácí $X = 3\,225$ utkání, relativní četnost tak byla $\hat{\pi}_1 = \frac{X}{m} = \frac{3\,225}{5\,580} \doteq 0,5780$. Následně byla vybrána data, ve kterých hráli **back-to-back utkání pouze hosté**, takových bylo $n = 1\,887$. Z těchto utkání vyhráli domácí $Y = 1\,188$ utkání, relativní četnost tak byla $\hat{\pi}_2 = \frac{Y}{n} = \frac{1\,188}{1\,887} \doteq 0,6296$. Nyní předpokládáme, podobně jako v odstavci 4.2, že X a Y jsou náhodné veličiny popisující počet vítězství domácích týmů (v jednotlivých skupinách výběru utkání), a že $X \sim \text{Bi}(m = 5\,580; \pi_1)$ a $Y \sim \text{Bi}(n = 1\,887; \pi_2)$. Na hladině významnosti $\alpha = 5\%$ byla testována hypotéza

$$H_0 : \pi_1 = \pi_2, \quad \text{proti alternativě} \quad H_1 : \pi_1 < \pi_2.$$

Jinými slovy byla testována hypotéza, že domácí (nehrající back-to-back) vítězí stejně často proti týmům, které nehrají back-to-back, jako proti týmům které back-to-back hrají – proti alternativě, že domácí (nehrající back-to-back) proti týmům, které hrají back-to-back, vítězí častěji. Označme $z_1 = \frac{X+Y}{m+n} = \frac{3\,225+1\,188}{5\,580+1\,887} \doteq 0,5910$. Při velkých hodnotách m a n lze opět vycházet z aproximace normálním rozdělením a testová statistika je pak, viz [1],

$$T_1 = \frac{\hat{\pi}_1 - \hat{\pi}_2}{\sqrt{z_1(1-z_1)\left(\frac{1}{m} + \frac{1}{n}\right)}} = \frac{0,5780 - 0,6296}{\sqrt{0,5910(1-0,5910)\left(\frac{1}{5\,580} + \frac{1}{1\,887}\right)}} \doteq -3,94. \quad (4.5)$$

Kritická hodnota pro tuto jednostrannou alternativu je $-u_{1-\alpha} = -u_{0,95} \doteq -1,64$. Protože $T_1 = -3,94 < -1,64$, zamítáme na hladině významnosti $\alpha = 5\%$ hypotézu H_0 a přikláníme se k alternativě, že $\pi_1 < \pi_2$, tedy že domácí týmy (nehrající back-to-back) vítězí proti týmům, které hrají back-to-back častěji, než proti týmům nehrajícím back-to-back. Back-to-back faktor má tedy vliv na vítězství v utkání.

(2) Ve druhé variantě byla nejprve vybrána opět pouze utkání, ve kterých ani jeden z týmů nehrál back-to-back utkání, kterých bylo $m = 5\,580$. Z těchto utkání vyhráli domácí $X = 3\,225$ utkání, relativní četnost tak byla $\hat{\pi}_1 = \frac{X}{m} = \frac{3\,225}{5\,580} \doteq 0,5780$. Následně byla vybrána data, ve kterých hráli **back-to-back utkání pouze domácí**, takových bylo $o = 600$. Z těchto utkání vyhráli domácí $Z = 325$ utkání, relativní četnost tak byla $\hat{\pi}_3 = \frac{Z}{o} = \frac{325}{600} \doteq 0,5417$. Nyní předpokládáme, že X a Z jsou náhodné veličiny popisující počet vítězství domácích týmů (v jednotlivých skupinách výběru utkání), a že $X \sim \text{Bi}(m = 5\,580; \pi_1)$ a $Z \sim \text{Bi}(o = 600; \pi_3)$. Na hladině významnosti $\alpha = 5\%$ byla testována hypotéza

$$H_0 : \pi_1 = \pi_3, \quad \text{proti alternativě} \quad H_1 : \pi_1 > \pi_3.$$

Jinými slovy byla testována hypotéza, že domácí nehrající back-to-back vítězí stejně často proti ostatním týmům (nehrajícím back-to-back), jako když back-to-back hrají – proti alternativě, že pokud domácí nehrají back-to-back utkání (proti týmům nehrajícím back-to-back) vítězí častěji, než pokud hrají back-to-back. Označme nyní $z_2 = \frac{X+Z}{m+o} = \frac{3\,225+325}{5\,580+600} \doteq 0,5744$. Při velkých hodnotách m a o lze opět vycházet z aproximace normálním rozdělením a testová statistika je pak, viz [1],

$$T_2 = \frac{\hat{\pi}_1 - \hat{\pi}_3}{\sqrt{z_2(1-z_2)\left(\frac{1}{m} + \frac{1}{o}\right)}} = \frac{0,5780 - 0,5417}{\sqrt{0,5744(1-0,5744)\left(\frac{1}{5\,580} + \frac{1}{600}\right)}} \doteq 1,71. \quad (4.6)$$

Kritická hodnota pro tuto jednostrannou alternativu je $u_{1-\alpha} = u_{0,95} \doteq 1,64$. Protože $T_2 = 1,71 > 1,64$, zamítáme na hladině významnosti $\alpha = 5\%$ hypotézu H_0 a přikláníme se k alternativě, že $\pi_1 > \pi_3$, tedy že pokud domácí nehrají back-to-back utkání (proti týmům nehrajícím back-to-back), vítězí častěji, než pokud hrají back-to-back. Back-to-back faktor má tedy vliv na vítězství v utkání.

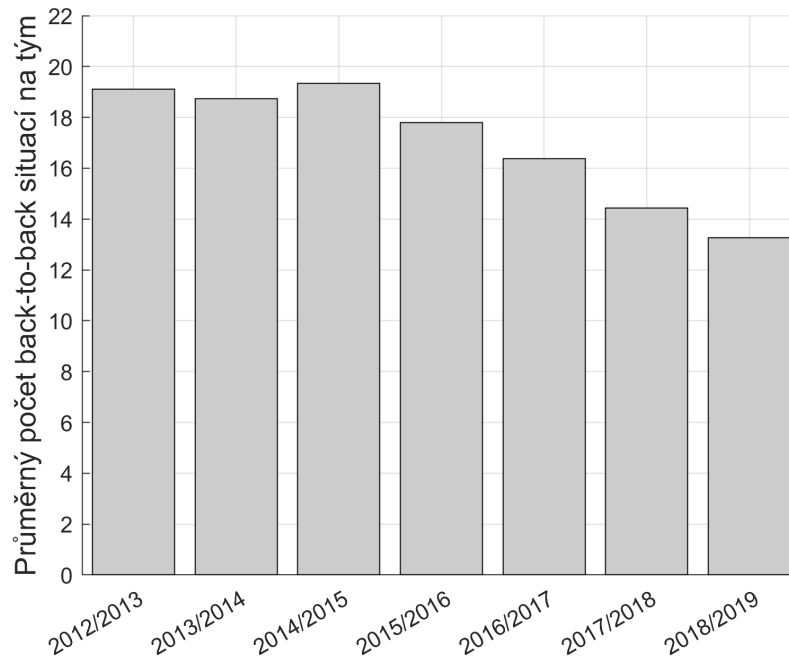
Parametr $B > 0$ určuje počet bodů Elo ratingu, který bude vždy odečten od ratingu týmu hrajícího utkání druhý den po sobě, tedy back-to-back utkání. Z logiky věci se jedná o penalizační parametr, jelikož jde v případě back-to-back utkání o nevýhodu (viz výše), proto je tento parametr odečítán. Opět se díky tomuto přídatnému parametru změní rozdíl v ratingu soupeřících týmů, což má vliv na odhad pravděpodobnosti vítězství. Předpokládáno bylo opět logistické rozdělení výkonnosti a vzorec (4.3) pro odhad pravděpodobnosti vítězství byl tak rozšířen jako

$$P_{ij} = \frac{1}{1 + 10^{-\frac{(r_i + h_i \cdot H - b_i \cdot B) - (r_j + h_j \cdot H - b_j \cdot B)}{400}}}, \quad (4.7)$$

kde $B > 0$ je parametr penalizace back-to-back utkání a b_i je indikátor back-to-back utkání týmu i , tedy

$$b_i = \begin{cases} 1, & \text{pokud tým } i \text{ hraje druhé utkání ve dvou dnech,} \\ 0, & \text{pokud tým } i \text{ nehraje druhé utkání ve dvou dnech.} \end{cases} \quad (4.8)$$

Parametr B byl uvažován konstantní a shodný pro všechny týmy. K přepočtu ratingu docházelo opět po každém odehraném utkání podle vzorce (3.14). Jako počáteční rating všech týmů byla zvolena opět hodnota $R_0 = 1\,500$. V tomto modelu byl rovněž uvažován **konstantní koeficient rozvoje K** po celé sledované období. Opět byly použity – stejně jako v odstavci 4.2 – **konstantní parametry p a H** . Všechny čtyři parametry K , p , H a B , budou optimalizovány v kapitole 5.



Obrázek 4.1: Průměrný počet back-to-back situací na tým, za jednu sezonu (82 utkání)

4.4 Model 4

Tento model je opět dalším rozšířením předchozích modelů, tedy rozšířením modelu 3 z odstavce 4.3. Na rozdíl od předchozích dvou se toto rozšíření netýká zvyšování/snižování ratingu na základě přidaných faktorů, ale týká se koeficientu rozvoje K . Jak již bylo uvedeno v odstavci 3.2.4, tento koeficient určuje maximální možnou hodnotu zvýšení, resp. snížení ratingu při jeho aktualizaci po utkání. Určuje tedy, jak moc výsledek utkání ovlivňuje změnu ratingu z jeho předchozí hodnoty. Do této chvíle byl uvažován koeficient K konstantní. Nyní bude tento **koeficient K klesající** v závislosti na počtu odehraných utkání v sezoně.

Parametr p popsaný v odstavci 4.1 slouží k úpravě ratingu každého týmu po skončení sezony, resp. na začátku sezony následující. Tato úprava spočívá v částečném navrácení ratingu k průměru, a to z důvodu očekávané změny síly týmů mezi sezonami. Přesto je nevýhodou Elo ratingu pro NBA to, že pokud nastanou v týmu velké změny (např. přestup jednoho z nejlepších hráčů historie, LeBrona Jamese), což se ve většině případů děje v přestávce mezi sezonami, může trvat delší dobu, než se tento fakt projeví v ratingu

daného týmu. Aby se po začátku nové sezony ratingy jednotlivých týmů rychleji přizpůsobovaly jejich skutečné síle, a naopak v průběhu sezony byly ratingy stabilnější a méně citlivé na případné odchylky od stabilní výkonnosti týmů, byl koeficient K uvažován klesající v průběhu každé jednotlivé sezony. Hodnota koeficientu K byla proto při každém utkání vynásobena *multiplikátorem poklesu* $m(x)$. Jedná se o funkci mocninného poklesu. Tento typ poklesu byl inspirován modelem pro predikci tenisového US Open 2016 od FiveThirtyEight [16], kde byla – podobně jako zde – použita varianta mocninného poklesu koeficientu K . V modelu 4 byla tedy použita hodnota koeficientu K ve tvaru

$$K = k \cdot m(x), \quad (4.9)$$

kde k je konstanta a

$$m(x) = \frac{1}{x^a}, \quad (4.10)$$

kde $a \in \langle 0, 1 \rangle$ je konstanta a x udává, kolikáté utkání v dané sezoně týmy hrají. V NBA se nehrají pravidelná soutěžní kola, kterých by se účastnily všechny týmy. Není tak výjimkou, že jednotlivé týmy mají v průběhu sezony rozdílný počet odehraných utkání. Jedná se však pouze o malé rozdíly. Proto je x uvažováno jako aritmetický průměr počtu, kolikáté utkání v dané sezoně týmy i a j hrají, tedy

$$x = \frac{n_i + n_j}{2}, \quad (4.11)$$

kde $n_i \in \{1, 2, \dots, 82\}$ udává, kolikáté utkání v dané sezoně hraje tým i a $n_j \in \{1, 2, \dots, 82\}$ udává, kolikáté utkání v dané sezoně hraje tým j . Po dosažení je tedy výsledná hodnota koeficientu K v tomto modelu rovna

$$K = k \cdot \left(\frac{n_i + n_j}{2} \right)^{-a}. \quad (4.12)$$

Parametry k a a jsou uvažovány konstantní a shodné pro všechny týmy. K přepočtu ratingu docházelo opět po každém odehraném utkání podle vzorce (3.14), přičemž pravděpodobnost vítězství byla odhadována stejně jako v modelu 3 podle vzorce (4.7). Jako počáteční rating všech týmů byla zvolena opět hodnota $R_0 = 1500$. Stejně jako v odstavci 4.3 byly použity **konstantní parametry p , H a B** . Všech pět parametrů k , p , H , B a a , bude optimalizováno v kapitole 5.

4.5 Model 5

Tento model je kombinací modelu založeného na Elo ratingu od FiveThirtyEight [26], viz odstavec 4.6, a předchozího modelu 4 z odstavce 4.4. Již dva parametry inspirované modelem od FiveThirtyEight, a to p a H , byly v našem postupně rozšiřovaném modelu použity. Nyní z tohoto modelu bude použita poslední část. Koeficient rozvoje K v Elo modelu od FiveThirtyEight nemá tu vlastnost, že je v průběhu sezony klesající, tak jako v modelu 4 z odstavce 4.4, přesto není konstantní. Tento model používá koeficient K , který je závislý na rozdílu výsledného skóre v utkání obou týmů a naplnění jeho očekávání. Koeficient rozvoje K je konstruován tak, aby se týmům po utkání rating zvyšoval, resp. snižoval více, pokud byl vyšší rozdíl ve výsledném skóre. Podle [26] funguje tento systém tak, že je v každé hře přiřazen ke koeficientu K multiplikátor, který je závislý na tzv. *margin of victory* (MOV), což je rozdíl výsledného skóre ve prospěch vítěze utkání, formálně zapsáno

$$MOV = |PTS_i - PTS_j|, \quad (4.13)$$

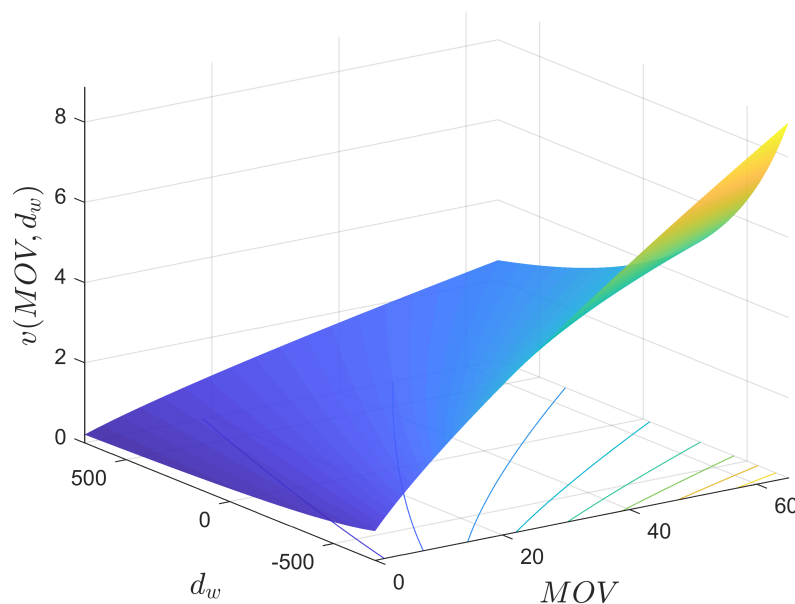
kde je

PTS_i celkový počet získaných bodů týmem i ,

PTS_j celkový počet získaných bodů týmem j ,

a na rozdíl ratingů týmů i a j z pohledu vítězného týmu, který označme jako d_w . Tento multiplikátor, resp. funkci, budeme nazývat jako *MOV multiplikátor* a značit jako v . Graf $v(MOV, d_w)$ je vykreslen na obrázku 4.2. Jeho exaktní předpis používaný FiveThirtyEight je

$$v(MOV, d_w) = \frac{(MOV + 3)^{0,8}}{7,5 + 0,006 \cdot d_w}. \quad (4.14)$$



Obrázek 4.2: Graf MOV multiplikátoru v závislosti na MOV a d_w

Odvození tohoto tvaru MOV multiplikátoru ani parametrů, které obsahuje, není zveřejněno. V [26] je pouze naznačeno, že se jedná o podíl MOV a „projekce MOV“ z rozdílu ratingů – tedy odhadu očekávaného MOV na základě rozdílu ratingů. Naplnění, resp. nenaplnění očekávaného MOV odhadnutého z rozdílu ratingů soupeřících týmů se tak odráží na výsledné velikosti koeficientu K . Dále autor modelu Nate Silver uvádí, že je takto MOV multiplikátor konstruován ve snaze minimalizovat autokorelaci, která v Elo modelech vzniká nejen díky tomu, že favorité mají tendenci vyhrávat častěji, ale z hlediska MOV mají tendenci vyhrávat s vyšším rozdílem skóre, než prohrávají [25]. Podobný tvar multiplikátoru navrhl Nate Silver také pro NFL (National Football League), rovněž bez zveřejnění jeho odvození. O jeho odvození se však nedávno pokusil ve svém pracovním textu Steven Morse [17].

Díky nejasnosti odvození MOV multiplikátoru nebudou optimalizovány jeho parametry, ale bude takto odvozený použit ve snaze rozšířit model 4 z odstavce 4.4 o další faktor, tedy rozdíl výsledného skóre MOV. V modelu 5 byla tedy použita hodnota koeficientu rozvoje K ve tvaru

$$K = k \cdot m(x) \cdot v(MOV, d_w). \quad (4.15)$$

V námi uvažovaném modelu, který zahrnuje rovněž parametry H a B , je pak – podobně jako ve vzorci (4.7) pro odhad pravděpodobnosti vítězství – rozdíl ratingů z pohledu vítězného týmu d_w , týmů i a j , o tyto parametry rozšířen. Formálně tento rozdíl ratingů z pohledu vítězného týmu zapíšeme jako

$$d_w = W_i \cdot [(r_i + h_i \cdot H - b_i \cdot B) - (r_j + h_j \cdot H - b_j \cdot B)], \quad (4.16)$$

kde

$$W_i = \begin{cases} 1, & \text{pokud byl tým } i \text{ vítězem utkání,} \\ -1, & \text{pokud byl tým } i \text{ poraženým utkání.} \end{cases} \quad (4.17)$$

Po dosazení do (4.15) je tedy výsledná hodnota koeficientu K v tomto modelu rovna

$$\begin{aligned} K &= k \cdot \frac{1}{x^a} \cdot \frac{(MOV + 3)^{0,8}}{7,5 + 0,006 \cdot d_w} = \\ &= k \cdot \left(\frac{n_i + n_j}{2}\right)^{-a} \cdot \frac{(|PTS_i - PTS_j| + 3)^{0,8}}{7,5 + 0,006 \cdot W_i [(r_i + h_i \cdot H - b_i \cdot B) - (r_j + h_j \cdot H - b_j \cdot B)]}. \end{aligned} \quad (4.18)$$

Opět jsou v koeficientu K uvažovány **konstantní parametry** k a a , které jsou shodné pro všechny týmy. K přepočtu ratingu docházelo rovněž po každém odehraném utkání podle vzorce (3.14), přičemž pravděpodobnost vítězství byla odhadována stejně jako v modelu 3, podle vzorce (4.7). Jako počáteční rating všech týmů byla zvolena opět hodnota $R_0 = 1500$. Stejně jako v předchozích modelech byly použity **konstantní parametry** p , H a B . Všech pět parametrů k , p , H , B a a , bude optimalizováno v kapitole 5.

4.6 Model Elo538

V tomto odstavci bude popsán Elo model navržený FiveThirtyEight, který budeme označovat jako model *Elo538*. Jedná se o nejpoužívanější model založený na Elo ratingu pro NBA, který byl do roku 2015 využíván k predikci výsledků utkání webovou stránkou FiveThirtyEight.com [26]. Postupně však tento model začali jeho tvůrci rozšiřovat o hodnocení jednotlivých hráčů a předpověď a projekci jejich výkonů do výkonu týmů. Nyní tento server používá predikční model RAPTOR, ve kterém bylo již úplně upuštěno od hodnocení využívajícího Elo rating, více v [27]. Přesto FiveThirtyEight stále aktualizují a zveřejňují rating a predikce také podle jejich původního Elo rating modelu [31].

Z původního modelu Elo538 byly v námi rozšiřovaných modelech postupně využívány některé parametry, k jejichž optimalizaci dojde v kapitole 5. Jedná se o parametr zachování ratingu po skončení sezony p , o parametr domácí výhody H a v neposlední řadě o koeficient rozvoje K , resp. hodnotu k rozšířenou o MOV multiplikátor $v(MOV, d_w)$. Podle modelu Elo538 jsou optimální hodnoty parametrů

- $k = 20$,
- $H = 100$,
- $p = 0,75$,

přičemž koeficient rozvoje K je rozšířený o MOV multiplikátor, viz model 5 z odstavce 4.5, tedy

$$K = k \cdot v(MOV, d_w) = 20 \cdot \frac{(MOV + 3)^{0,8}}{7,5 + 0,006 \cdot d_w}. \quad (4.19)$$

Na rozdíl od námi rozšířeného modelu 5 z odstavce 4.5 se zde nevyskytuje multiplikátor poklesu $m(x)$ a parametr penalizace back-to-back utkání B , proto je v tomto případě hodnota rozdílu ratingu z pohledu vítězného týmu pouze

$$d_w = W_i \cdot [(r_i + h_i \cdot 100) - (r_j + h_j \cdot 100)]. \quad (4.20)$$

Model Elo538 se na rozdíl od námi uvažovaných modelů učí už od začátku působení daného týmu v NBA, resp. dříve v ABA (American Basketball Association). Znamená to, že počáteční rating, který je v tomto modelu uvažován jako $R_0 = 1\,300$, byl každému týmu přidělen už na začátku jeho působení a od této chvíle byl rating počítán, což přináší možnost historického srovnání všech týmů. Díky historickým změnám v účastnících NBA a jejich počtu není zachován průměrný rating všech týmů 1 300, ale tento průměr je uvažován jako 1 505, viz [26]. Z tohoto důvodu je (na rozdíl od vzorce (4.1)) v tomto modelu úprava ratingu každého týmu i před začátkem nové sezony ve tvaru

$$r_i^{(s)} = p \cdot r_i^{(s-1)} + (1 - p) \cdot 1\,505 = 0,75 \cdot r_i^{(s-1)} + 0,25 \cdot 1\,505. \quad (4.21)$$

Dalším rozdílem modelu Elo538 je aktualizace ratingu i během play-off, zatímco ve všech 5 námi rozšířených modelech není play-off vůbec uvažováno.

Pro porovnání modelu Elo538 s námi navrženými 5 modely byla dodatečně získána zveřejněná data Elo ratingu, resp. odhadnutých pravděpodobností vítězství (ke kterému docházelo díky nepřítomnosti parametru B podle vzorce (4.3)), ze zdroje [31], viz příloha A.7. Z těchto dat byly vybrány pouze záznamy ze základních částí sledovaných sezon tak, aby bylo možné srovnat úspěšnost predikce s ostatními modely. Jak již bylo uvedeno výše, v tomto modelu byly uvažovány optimální parametry podle FiveThirtyEight, nebyly zde tedy žádné parametry optimalizovány.

Kapitola 5

Optimalizace a výběr modelu

V této kapitole budou optimalizovány parametry pro jednotlivé modely 1 až 5, které byly popsány v kapitole 4. Jak již bylo uvedeno v odstavci 2.2, byla získána data ze základních částí celkem sedmi sezon, konkrétně od sezony 2012/2013 do sezony 2018/2019. Tato data byla rozdělena do tří datových sad: *učení*, *optimalizace* a *predikce*. Na začátku sledovaného období byl každému týmu přidělen stejný rating, a to $R_0 = 1500$ pro všech pět modelů. Díky tomu není vhodné optimalizovat a vybírat model na základě predikčních výsledků modelů z počátku sledovaného období – tyto výsledky, resp. odhadnuté pravděpodobnosti vítězství, jsou zkreslené kvůli stejné „startovací pozici“ všech týmů. Je tedy potřeba nechat se modely určitý čas učit, resp. počítat rating bez toho, abychom sledovali jejich úspěšnost predikce. Tato doba učení byla zvolena jako jedna sezona, tedy sezona 2012/2013, což bylo považováno za dostatečně dlouhou dobu pro ustálení modelů. Následující tři sezony, tedy 2013/2014, 2014/2015 a 2015/2016, byly zvoleny jako dostatečně dlouhá doba pro optimalizaci modelů a zároveň dost dlouhá na to, aby se parametry modelů nepřizpůsobily např. pouze jedné sezoně, ve které mohl být náhodný výkyv. Pro tyto tři sezony budou na základě kritéria logaritmické ztrátové funkce, viz odstavec 5.1.2, optimalizovány parametry jednotlivých modelů. Pro získané optimální parametry budou vypočítána kritéria určující kvalitu modelů, viz odstavec 5.1, na základě kterých bude vybrán nejlepší model. Predikční schopnost vybraného modelu bude hodnocena v kapitole 6. Model bude použit k predikci a fiktivnímu použití proti sázkové kanceláři v predikční sadě dat, tedy ve zbývajících třech sezonách 2016/2017, 2017/2018 a 2018/2019, čímž zprostředkovaně ověříme, jak by si model vedl ve srovnání s modely, které běžně používají sázkové kanceláře. Kromě srovnání predikční schopnosti vybraného modelu s modelem sázkové kanceláře bude predikční schopnost srovnána také s modelem Elo538, viz odstavec 4.6.

5.1 Kritéria kvality modelů

Obecně je cílem každého predikčního modelu jeho úspěšnost mimo tréninkovou, resp. optimalizační sadu dat, tedy na datech, která byla skryta. Kvalita konkrétního modelu může být hodnocena pomocí různých metod a metrik, záleží především na cíli studie. V této práci bylo použito pět kritérií kvality modelu. Z hlediska měření úspěšnosti odhadnutých pravděpodobností vítězství v utkáních byla použita kritéria *logaritmická ztrátová funkce*, *Brierovo skóre* a *přesnost*. Z hlediska kalibrace byl použit *kalibrační poměr* a z hlediska

diskriminace kritérium *AUC-ROC*. Tato kritéria byla vypočítána pro všechny uvažované modely s optimálními parametry z kapitoly 4, a to vždy pro optimalizační sadu dat, přičemž kritérium *logaritmická ztrátová funkce* bylo použito pro optimalizaci parametrů modelů. Rovněž byla tato kritéria vypočítána v kapitole 6 pro predikční sadu dat vybraného modelu, pro zhodnocení jeho predikční schopnosti

5.1.1 Přesnost

Kritérium přesnosti je nejjednodušším kritériem kvality predikčního modelu. Jedná se o pouhé procentuální vyjádření, v kolika utkáních model správně určil vítěze utkání. Jako predikovaný vítěz byl vždy určen tým s vyšší odhadnutou pravděpodobností vítězství v daném utkání. Pro žádný z modelů nenastala ani v jednom utkání situace, že by byla pro oba soupeřící týmy odhadnuta stejná pravděpodobnost vítězství – v takovém případě by tato utkání byla z výpočtu kritéria vyřazena. Vzorec pro určení přesnosti *PR* je tedy

$$PR = \frac{\sum_{i=1}^N y_i}{N}, \quad (5.1)$$

kde N je celkový počet utkání, pro která bylo kritérium počítáno a y_i je indikátor vítězství týmu s vyšší odhadnutou pravděpodobností v utkání i , tedy

$$y_i = \begin{cases} 1, & \text{pokud zvítězil tým s vyšší odhadnutou pravděpodobností,} \\ 0, & \text{pokud nezměřil tým s vyšší odhadnutou pravděpodobností.} \end{cases} \quad (5.2)$$

Čím vyšší je tato hodnota, tím byl model úspěšnější. Toto kritérium však není vhodné pro určování přesnosti odhadovaných pravděpodobností pomocí modelu, a to z důvodu popsaném v následujícím odstavci.

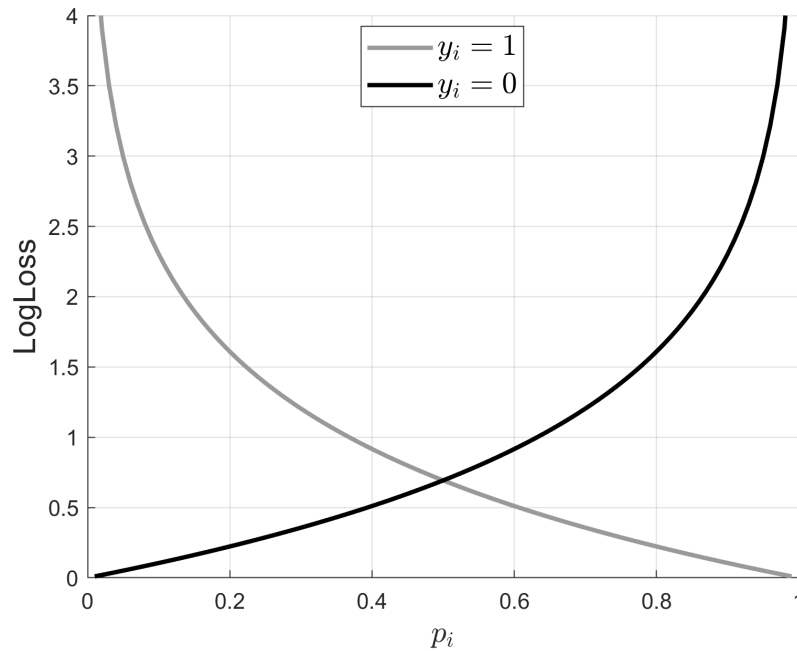
5.1.2 Logaritmická ztrátová funkce

Logaritmická ztrátová funkce (známá jako *LogLoss*) je křížová entropie, více v [10]. Oproti obyčejné přesnosti popsané v odstavci 5.1.1 má toto kritérium schopnost určit kvalitu modelu na základě přesnosti odhadovaných pravděpodobností. Uvažujme dva modely pro odhad pravděpodobnosti vítězství. První model odhadne, že pravděpodobnost vítězství týmu A je 51 %. Druhý model odhadne, že pravděpodobnost vítězství týmu A je 99 %. V případě výhry týmu A jsou podle kritéria obyčejné přesnosti oba tyto modely stejně přesné, protože oba předpovídaly, že tým A vyhraje. Oproti tomu logaritmická ztrátová funkce dokáže zachytit fakt, že druhý model byl v predikci přesnější. Vzorec pro výpočet kritéria *LogLoss* je ve tvaru

$$LL = -\frac{1}{N} \sum_{i=1}^N [y_i \cdot \ln(p_i) + (1 - y_i) \cdot \ln(1 - p_i)], \quad (5.3)$$

kde y_i je opět indikátor vítězství týmu s vyšší odhadnutou pravděpodobností v utkání i a p_i je odpovídající pravděpodobnost. N je celkový počet utkání, pro která bylo kritérium počítáno (pro toto kritérium je přípustná shodnost odhadnutých pravděpodobností).

Z uvedeného vzorce je zřejmé, že pokud se odhadnutá pravděpodobnost vítězství $p_i \rightarrow 1$, tedy $\ln(p_i) \rightarrow 0$, a zároveň tento tým skutečně zvítězí, tedy $y_i = 1$, pak se LogLoss nebude zvětšovat a naopak. Na obrázku 5.1 je logaritmická ztrátová funkce vykreslena. Je zde vidět, že čím nižší je logaritmická ztrátová funkce, tím přesnější je predikce, tedy tím lepší je model. Lze zde také pozorovat její vlastnost, že vysoce penalizuje velké chyby predikce [18].



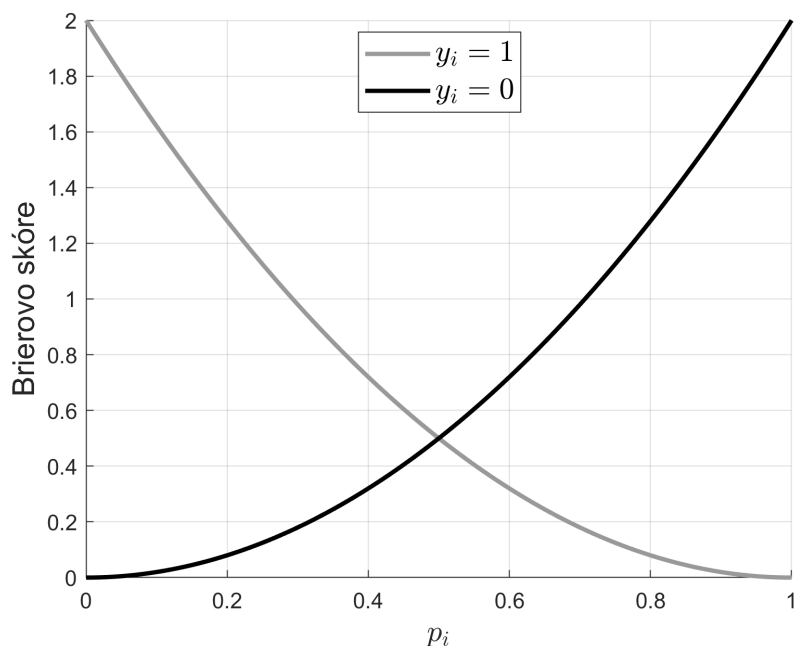
Obrázek 5.1: Logaritmická ztrátová funkce

5.1.3 Brierovo skóre

Brierovo skóre je alternativou k logaritmické ztrátové funkci. Jedná se o střední kvadratickou chybu mezi odhadovanou pravděpodobností vítězství a skutečným výsledkem utkání. Vzorec pro výpočet Brierova skóre BS je

$$BS = \frac{1}{N} \sum_{i=1}^N [(p_i - y_i)^2 + ((1 - p_i) - (1 - y_i))^2], \quad (5.4)$$

kde y_i je opět indikátor vítězství týmu s vyšší odhadnutou pravděpodobností v utkání i a p_i je odpovídající pravděpodobnost. N je celkový počet utkání, pro která bylo kritérium počítáno (pro toto kritérium je přípustná shodnost odhadnutých pravděpodobností). Z uvedeného vzorce je zřejmé, že čím větší budou odchylky odhadnutých pravděpodobností od skutečných výsledků, tím se bude BS zvyšovat. Na obrázku 5.2 je vidět, že čím nižší je hodnota Brierova skóre, tím přesnější je predikce, tedy tím lepší je model. Velké chyby predikce jsou opět hodně penalizovány, ale ne tolik jako u logaritmické ztrátové funkce [18].



Obrázek 5.2: Bierovo skóre

5.1.4 Kalibrační poměr

Kalibrace modelu je vlastnost týkající se shody mezi pozorovanými výsledky a predikcí. Model je dobře kalibrovaný, pokud při uvážení všech utkání, pro která byla odhadnuta pravděpodobnost vítězství jednoho z týmů jako hodnota p , bude poměr skutečných vítězství těchto týmů roven přibližně p . Tedy například pokud budeme uvažovat 100 utkání, ve kterých byla odhadnuta pravděpodobnost vítězství jednoho z týmů vždy 75 %, měl by počet vítězství těchto týmů být roven přibližně $100 \cdot 0,75 = 75$. K určení, zda byla tato podmínka splněna, byl použit tzv. kalibrační poměr. Vzorec pro výpočet kalibračního poměru KP je

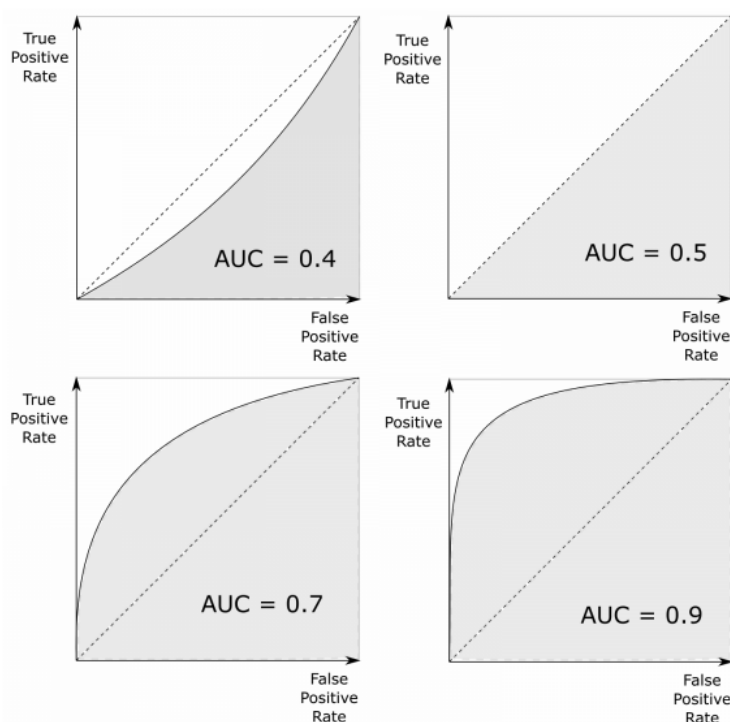
$$KP = \frac{\sum_{i=1}^N p_i}{\sum_{i=1}^N y_i}, \quad (5.5)$$

kde y_i je indikátor vítězství týmu s vyšší odhadnutou pravděpodobností v utkání i a p_i je odpovídající pravděpodobnost. N je celkový počet utkání, pro která bylo kritérium počítáno (podobně jako pro kritérium přesnosti, by zde byla vyřazena utkání se stejně odhadnutou pravděpodobností vítězství pro oba týmy). Pokud je model dobře kalibrovan, kalibrační poměr KP bude blízký hodnotě 1 [13].

5.1.5 AUC-ROC

Kritérium AUC-ROC slouží k hodnocení kvality modelu z hlediska toho, jak je model schopný rozlišovat mezi třídami (vítězství/prohra). Tato vlastnost se nazývá diskriminace. Jedná se o kombinaci ROC (Receiver Operating Characteristic) křivky a AUC (Area under the Curve). Křivka ROC zobrazuje relativní četnost skutečně pozitivních pozorování, tzv. *true positive rate* (TPR), v závislosti na relativní četnosti falešně pozitivních

pozorování, tzv. *false positive rate* (FPR) [18]. Pomocí této křivky lze tedy zachytit vztah mezi tzv. specificitou a senzitivitou modelu. AUC je pak obsah plochy pod touto křivkou [8]. Na obrázku 5.3 jsou znázorněny různé příklady AUC-ROC.



Obrázek 5.3: Ukázka různých příkladů AUC-ROC, zdroj [21]

Čím vyšší je hodnota AUC, tím vyšší je přesnost predikce, resp. tím lepší je diskriminace modelu. V případě, že $AUC = 1$, jedná se o perfektní přesnost. Jestliže $AUC = 0,5$, pak je úspěšnost modelu srovnatelná s úspěšností modelu „házení mince“, což bude ilustrováno v příkladu 5.1. K výpočtu tohoto kritéria byla použita funkce *perfcurve* v MATLABu.

Příklad 5.1. Uvažujme, že máme pouze 8 utkání dvou týmu. Označme

- TP správně predikovaná výhra,
- TN správně predikovaná prohra,
- FP chybně predikovaná výhra,
- FN chybně predikovaná prohra.

Výpočet poměru *true positive rate* je pak $TPR = \frac{TP}{TP+FN}$ a výpočet poměru *false positive rate* je $FPR = \frac{FP}{FP+TN}$. Pro různá rozhodovací pravidla o vítězi utkání by pro model s úspěšností srovnatelnou s „házením mince“ byla úspěšnost predikce v průměru 50%. Tabulky zachycují četnosti TP , TN , FP a FN , by pak pro různá rozhodovací pravidla mohly být např. v následujících tvarech:

	Výhra	Prohra
Výhra – predikovaná	2	2
Prohra – predikovaná	2	2

Pro tuto variantu je $TPR = \frac{TP}{TP+FN} = \frac{2}{2+2} = \frac{1}{2}$ a $FPR = \frac{FP}{FP+TN} = \frac{2}{2+2} = \frac{1}{2}$.

	Výhra	Prohra
Výhra – predikovaná	1	1
Prohra – predikovaná	3	3

Pro tuto variantu je $TPR = \frac{TP}{TP+FN} = \frac{1}{1+3} = \frac{1}{4}$ a $FPR = \frac{FP}{FP+TN} = \frac{1}{1+3} = \frac{1}{4}$.

	Výhra	Prohra
Výhra – predikovaná	3	3
Prohra – predikovaná	1	1

Pro tuto variantu je $TPR = \frac{TP}{TP+FN} = \frac{3}{3+1} = \frac{3}{4}$ a $FPR = \frac{FP}{FP+TN} = \frac{3}{3+1} = \frac{3}{4}$.

Je zřejmé, že ve všech příkladech pro model s úspěšností srovnatelnou s „házením mince“ by body $[FPR; TPR]$ křivky ROC ležely **na přímce** $FPR = TPR$. Obsah pod touto křivkou je $AUC = 0,5$. Pokud naopak budeme uvažovat model, který je úspěšnější než model srovnatelný s „házením mince“, tabulky zachycují četnosti TP , TN , FP a FN , by pro různá rozhodovací pravidla mohly být např. v následujících tvarech:

	Výhra	Prohra
Výhra – predikovaná	3	1
Prohra – predikovaná	1	3

Pro tuto variantu je $TPR = \frac{TP}{TP+FN} = \frac{3}{3+1} = \frac{3}{4}$ a $FPR = \frac{FP}{FP+TN} = \frac{1}{1+3} = \frac{1}{4}$.

	Výhra	Prohra
Výhra – predikovaná	4	1
Prohra – predikovaná	1	2

Pro tuto variantu je $TPR = \frac{TP}{TP+FN} = \frac{4}{4+1} = \frac{4}{5}$ a $FPR = \frac{FP}{FP+TN} = \frac{1}{1+2} = \frac{1}{3}$.

	Výhra	Prohra
Výhra – predikovaná	4	0
Prohra – predikovaná	0	4

Pro tuto variantu je $TPR = \frac{TP}{TP+FN} = \frac{4}{4+0} = 1$ a $FPR = \frac{FP}{FP+TN} = \frac{0}{0+4} = 0$.

Ve všech těchto případech by body $[FPR; TPR]$ křivky ROC ležely **nad přímkou** $FPR = TPR$, a tedy $AUC > 0,5$. V případě perfektní přesnosti je pak $AUC = 1$.

5.2 Optimalizace parametrů

Parametry jednotlivých modelů byly optimalizovány z hlediska minimalizace logaritmické ztrátové funkce (viz odstavec 5.1.2) v sezonách optimalizační sady dat, tedy 2013/2014, 2014/2015 a 2015/2016. Optimalizace byla provedena v softwaru MATLAB R2018a, ve kterém byla naprogramována funkce `optimize.m`, viz příloha A.5. Tato funkce po výpočtu vrátí nalezené optimální parametry konkrétního modelu a tabulku, která obsahuje vypočítané hodnoty kritérií kvality v optimalizační sadě dat pro model s nalezenými parametry. K optimalizaci parametrů byla v MATLABu použita funkce `fmincon`.

Funkce `fmincon` slouží k hledání lokálního minima nelineární funkce s omezeními. Při použití této funkce bylo pro všech 5 modelů použito její výchozí nastavení algoritmu – tedy algoritmus *metody vnitřního bodu*. Tato metoda je popsána např. v [5]. Jako zastavovací kritéria algoritmu byly použity výchozí hodnoty, viz [32]. Ve všech případech se pomocí funkce `fmincon` podařilo nalézt parametry modelů, ve kterých měla logaritmická ztrátová funkce minimum v optimalizační sadě dat vybraných 3 sezon, což bylo celkem 3 690 utkání.

5.2.1 Model 1

V tomto modelu byly optimalizovány parametry K a p , viz odstavec 4.1. S využitím funkce `fmincon` bylo metodou vnitřního bodu nalezeno lokální minimum logaritmické ztrátové funkce, a to po 21 iteracích. Jako počáteční byly nastaveny hodnoty $K_0 = 75$ a $p_0 = 0,5$, tedy středy intervalů zvolených omezení $K \in \langle 0; 150 \rangle$ a $p \in \langle 0; 1 \rangle$. Optimalizace byla úspěšně zastavena, protože funkce byla ve všech přípustných směrech neklesající vzhledem k výchozím hodnotám optimalizačních a omezujících tolerancí, viz [32]. Jako minimalizační, resp. optimalizační parametry, byly nalezeny hodnoty (zaokrouhleno na dvě desetinná místa)

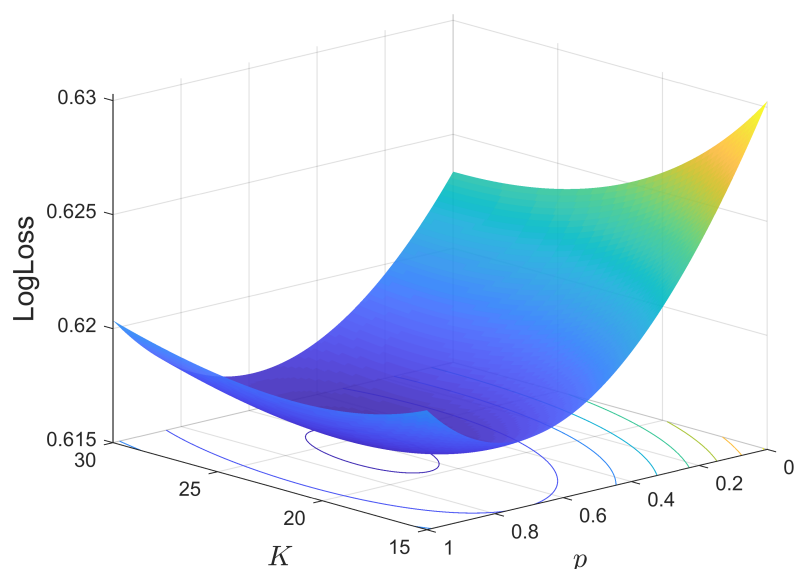
- $K = 23,39$,
- $p = 0,64$.

Ve snaze ověřit, že se nejedná pouze o lokální minimum, které je ovlivněné zvolenými počátečními podmínkami, byly testovány další různé kombinace počátečních hodnot parametrů. Testováno bylo dalších pět sad počátečních hodnot parametrů, přičemž počáteční hodnoty jednotlivých parametrů byly vybírány vždy náhodně z rovnoměrného rozdělení na intervalech zvolených omezení jednotlivých parametrů, viz výše. Pro všechny testované kombinace vstupních parametrů byly nalezeny – na dvě desetinná místa – stejné optimalizované parametry se stejnými hodnotami kritérií kvality modelu. V tabulce 5.1 jsou, kromě hodnoty logaritmické ztrátové funkce, uvedena všechna získaná kritéria kvality modelu pro nalezené optimální parametry (zaokrouhlené na dvě desetinná místa).

Tabulka 5.1: Kritéria kvality modelu 1, pro optimální parametry – optimalizace

	Hodnota kritéria
LogLoss	0,6158
Brierovo skóre	0,4273
AUC-ROC	0,7255
Kalibrační poměr	1,0006
Přesnost	0,6572

Na obrázku 5.4 je vykreslena LogLoss funkce v závislosti na parametrech K a p . Z obrázku je zřejmé, že LogLoss funkce je v okolí optima více citlivá na parametr p , než na parametr K . Ani v jednom případě však nelze mluvit o zásadní citlivosti, při které by malá změna v okolí optima způsobila výrazný nárůst LogLoss funkce, což je výhodou modelu.



Obrázek 5.4: Logaritmičká ztrátová funkce v závislosti na K a p v okolí nalezeného optima pro model 1

5.2.2 Model 2

Ve druhém modelu, který je rozšířený o parametr domácí výhody H , viz odstavec 4.2, byly optimalizovány všechny tři uvažované parametry: K, p, H . S využitím funkce `fmincon` bylo metodou vnitřního bodu nalezeno lokální minimum logaritmičké ztrátové funkce, a to po 30 iteracích. Jako počáteční byly nastaveny hodnoty $K_0 = 75$, $p_0 = 0,5$ a $H_0 = 75$, tedy středy intervalů zvolených omezení $K \in \langle 0; 150 \rangle$, $p \in \langle 0; 1 \rangle$ a $H \in \langle 0; 150 \rangle$. Optimalizace byla opět úspěšně zastavena, protože funkce byla ve všech přípustných směrech neklesající vzhledem k výchozím hodnotám optimalizačních a omezujících tolerancí, viz [32]. Jako minimalizační, resp. optimalizační parametry, byly nalezeny hodnoty (zaokrouhloeno na dvě desetinná místa)

- $K = 24,97$,

- $p = 0,62$,
- $H = 67,26$.

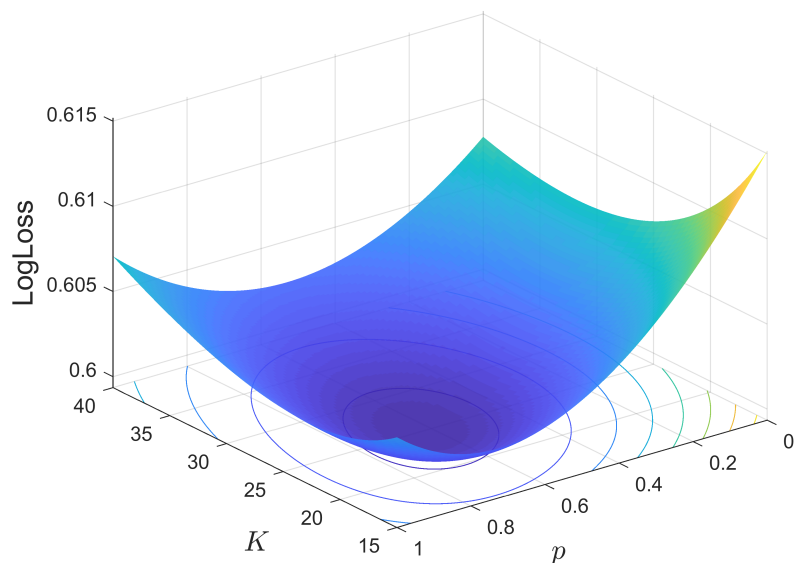
Ve snaze ověřit, že se nejedná pouze o lokální minimum, které je ovlivněno zvolenými počátečními podmínkami, byly testovány další různé kombinace počátečních hodnot parametrů. Testováno bylo dalších pět sad počátečních hodnot parametrů, přičemž počáteční hodnoty jednotlivých parametrů byly vybírány vždy náhodně z rovnoměrného rozdělení na intervalech zvolených omezení jednotlivých parametrů, viz výše. Pro všechny testované kombinace vstupních parametrů byly nalezeny – na dvě desetinná místa – stejné optimalizované parametry se stejnými hodnotami kritérií kvality modelu. Nalezené hodnoty K a p jsou blízké hodnotám z modelu 1. Nalezená optimální hodnota koeficientu K byla vyšší o 1,58 a hodnota parametru p nižší o 0,02. V tabulce 5.2 jsou, kromě hodnoty logaritmické ztrátové funkce, uvedena všechna získaná kritéria kvality modelu pro nalezené optimální parametry (zaokrouhlené na dvě desetinná místa).

Tabulka 5.2: Kritéria kvality modelu 2, pro optimální parametry – optimalizace

	Hodnota kritéria
LogLoss	0,5994
Brierovo skóre	0,4132
AUC-ROC	0,7263
Kalibrační poměr	0,9974
Přesnost	0,6753

Kritérium LogLoss je pro tento model s optimálními parametry nižší, než pro předchozí model s optimálními parametry – bez parametru H . Kromě kalibračního poměru dosahují lepších hodnot také všechna ostatní kritéria (AUC-ROC a přesnost jsou maximalizační kritéria). Kalibrační poměr je ale stále přibližně rovný hodnotě 1.

Logaritmická ztrátová funkce byla v okolí nalezeného optima vykreslena v závislosti na všech možných kombinacích dvou parametrů. Z grafů bylo zjištěno, že hodnota LogLoss funkce není v okolí optima zásadně citlivá na změny parametrů. Malá změna parametrů v okolí optima tedy nezpůsobí výrazný nárůst LogLoss funkce, což je výhodou tohoto modelu. Největší citlivost byla zjištěna na změnu parametru p , zejména v kombinaci se změnou hodnoty parametru K , viz obrázek 5.5. Naopak nejméně byla hodnota LogLoss funkce citlivá na změnu parametru H .



Obrázek 5.5: Logaritmická ztrátová funkce v závislosti na K a p v okolí nalezeného optima pro model 2

5.2.3 Model 3

Třetí model je rozšířením předchozího modelu o back-to-back faktor B , viz odstavec 4.3. V tomto modelu byly optimalizovány všechny čtyři uvažované parametry: K , p , H , B . S využitím funkce `fmincon` bylo metodou vnitřního bodu nalezeno lokální minimum logaritmické ztrátové funkce, a to po 54 iteracích. Jako počáteční byly nastaveny hodnoty $K_0 = 75$, $p_0 = 0,5$, $H_0 = 75$ a $B_0 = 75$, tedy středy intervalů zvolených omezení $K \in \langle 0; 150 \rangle$, $p \in \langle 0; 1 \rangle$, $H \in \langle 0; 150 \rangle$ a $B \in \langle 0; 150 \rangle$. Optimalizace byla opět úspěšně zastavena, protože funkce byla ve všech přípustných směrech neklesající vzhledem k výchozím hodnotám optimalizačních a omezujících tolerancí, viz [32]. Jako minimalizační, resp. optimalizační parametry, byly nalezeny hodnoty (zaokrouhlo na dvě desetinná místa)

- $K = 25,13$,
- $p = 0,62$,
- $H = 61,73$,
- $B = 32,81$.

Ve snaze ověřit, že se nejedná pouze o lokální minimum, které je ovlivněno zvolenými počátečními podmínkami, byly testovány další různé kombinace počátečních hodnot parametrů. Testováno bylo dalších pět sad počátečních hodnot parametrů, přičemž počáteční hodnoty jednotlivých parametrů byly vybírány vždy náhodně z rovnoměrného rozdělení na intervalech zvolených omezení jednotlivých parametrů, viz výše. Pro všechny testované kombinace vstupních parametrů byly nalezeny – na dvě desetinná místa – stejné optimalizované parametry se stejnými hodnotami kritérií kvality modelu. Nalezená hodnota parametru K je blízká hodnotě z modelu 2, resp. je vyšší pouze o 0,16 a parametr p je na 2 desetinná místa shodný s hodnotou z modelu 2. Na úkor přidání parametru B se snížila oproti předchozímu modelu hodnota H o 5,53. Závislost těchto dvou parametrů

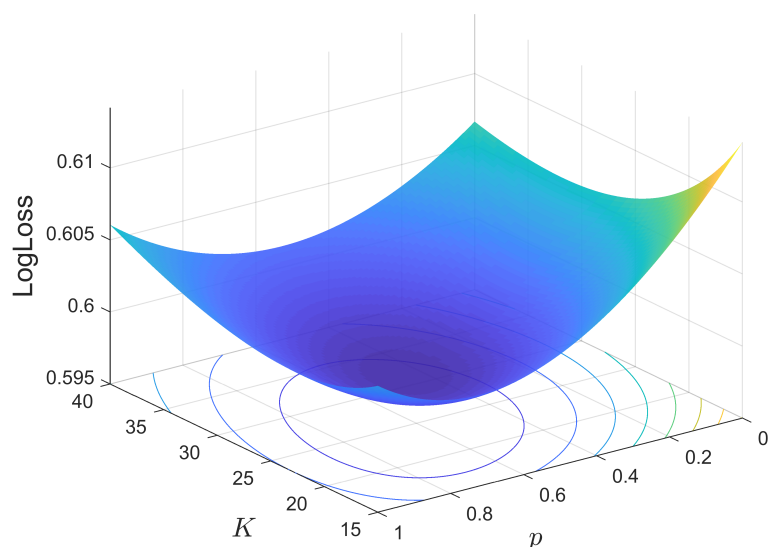
je způsobena tím, že domácí týmy nejsou tak často vystavovány back-to-back situacím jako týmy hostující, viz odstavec 4.3. V tabulce 5.3 jsou, kromě hodnoty logaritmické ztrátové funkce, uvedena všechna získaná kritéria kvality modelu pro nalezené optimální parametry (zaokrouhlené na dvě desetinná místa).

Tabulka 5.3: Kritéria kvality modelu 3, pro optimální parametry – optimalizace

	Hodnota kritéria
LogLoss	0,5983
Brierovo skóre	0,4124
AUC-ROC	0,7272
Kalibrační poměr	1,0014
Přesnost	0,6740

Pro tento model s optimálními parametry došlo oproti modelu 2 opět ke snížení minimalizačních kritérií LogLoss a Brierova skóre, zvýšení maximalizačního kritéria AUC-ROC a hodnota kalibračního poměru je v absolutní hodnotě blíže hodnotě 1. Pouze hodnota kritéria přesnosti se oproti modelu 2 zhoršila, jelikož je nižší (o 0,0013) – avšak toto kritérium není příliš vypovídající o přesnosti odhadnutých pravděpodobností, viz odstavec 5.1.2.

Logaritmická ztrátová funkce byla v okolí nalezeného optima vykreslena v závislosti na všech možných kombinacích dvou parametrů. Z grafů bylo zjištěno, že hodnota LogLoss funkce není v okolí optima zásadně citlivá na změny parametrů. Malá změna parametrů v okolí optima tedy nezpůsobí výrazný nárůst LogLoss funkce, což je výhodou tohoto modelu. Největší citlivost byla zjištěna opět na změnu parametru p , zejména v kombinaci se změnou hodnoty parametru K , viz obrázek 5.6. Na změnu rozšiřujícího parametru B v okolí optima byla zjištěna ještě nižší citlivost, než na změnu parametru H .



Obrázek 5.6: Logaritmická ztrátová funkce v závislosti na K a p v okolí nalezeného optima pro model 3

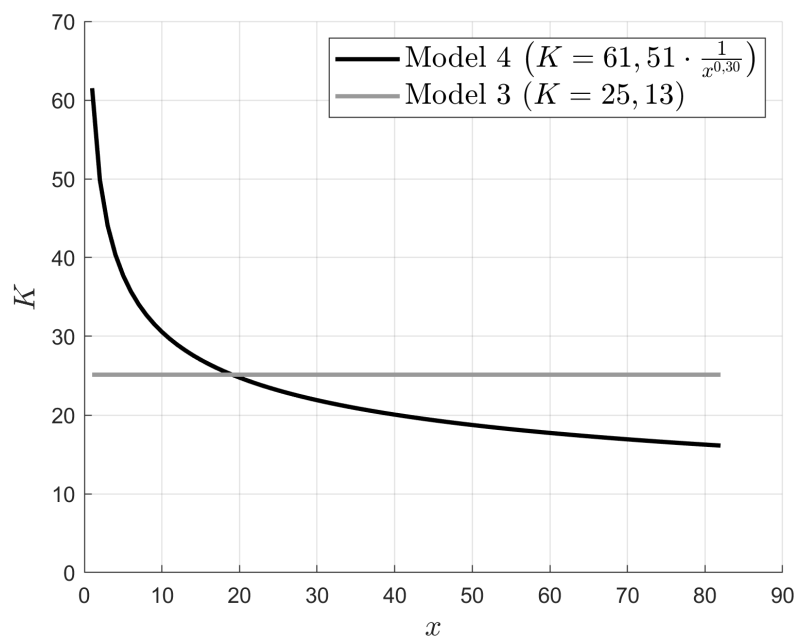
5.2.4 Model 4

Tento model je opět rozšířením předchozího modelu, tentokrát o multiplikátor mocninného poklesu parametru K , viz odstavec 4.4. Díky tomuto multiplikátoru je koeficient rozvoje K klesající v průběhu každé jednotlivé sezony. Koeficient K je ve tvaru $K = k \cdot m(x) = k \cdot \frac{1}{x^a}$, kde je x je průměrný počet odehraných utkání dvou soupeřících týmu v dané sezoně, viz výraz (4.11). Místo koeficientu rozvoje K se tak zde objevují dva parametry: k , p . V tomto modelu tedy bylo optimalizováno všech pět parametrů k , p , H , B , a .

S využitím funkce `fmincon` bylo metodou vnitřního bodu nalezeno lokální minimum logaritmické ztrátové funkce, a to po 62 iteracích. Jako počáteční byly nastaveny hodnoty $k_0 = 75$, $p_0 = 0,5$, $H_0 = 75$, $B_0 = 75$ a $a_0 = 0,5$, tedy středy intervalů zvolených omezení $k \in \langle 0; 150 \rangle$, $p \in \langle 0; 1 \rangle$, $H \in \langle 0; 150 \rangle$, $B \in \langle 0; 150 \rangle$ a $a \in \langle 0; 1 \rangle$. Optimalizace byla opět úspěšně zastavena, protože funkce byla ve všech přípustných směrech neklesající vzhledem k výchozím hodnotám optimalizačních a omezujících tolerancí, viz [32]. Jako minimalizační, resp. optimalizační parametry, byly nalezeny hodnoty (zaokrouhlo na dvě desetinná místa)

- $k = 61,51$,
- $p = 0,66$,
- $H = 62,08$,
- $B = 32,06$,
- $a = 0,30$.

Ve snaze ověřit, že se nejedná pouze o lokální minimum, které je ovlivněno zvolenými počátečními podmínkami, byly testovány další různé kombinace počátečních hodnot parametrů. Testováno bylo dalších pět sad počátečních hodnot parametrů, přičemž počáteční hodnoty jednotlivých parametrů byly vybírány vždy náhodně z rovnoměrného rozdělení na intervalech zvolených omezení jednotlivých parametrů, viz výše. Pro všechny testované kombinace vstupních parametrů byly nalezeny – na dvě desetinná místa – stejné optimalizované parametry se stejnými hodnotami kritérií kvality modelu. Nalezený parametr k je vzhledem k přidanému multiplikátoru neporovnatelný s hodnotou K v modelu 3. Parametry p , H a B zůstaly přibližně stejné, resp. parametr p vzrostl o 0,04, parametr H vzrostl o 0,35 a parametr B klesl o 0,75. Na obrázku 5.7 je zobrazena získaná optimální hodnota koeficientu rozvoje K v závislosti na průměrném počtu odehraných utkání dvou soupeřících týmu v dané sezoně. Z grafu jde vidět, že přibližně do 19. utkání je pro model 4 hodnota koeficientu K vyšší než pro model 3. Díky tomu se ratingy jednotlivých týmů v tomto modelu můžou rychleji přizpůsobovat jejich skutečné síle po začátku nové sezony a naopak později být stabilnější.



Obrázek 5.7: Srovnání optimálního koeficientu rozvoje K pro model 4 a model 3 v závislosti na počtu utkání v sezoně

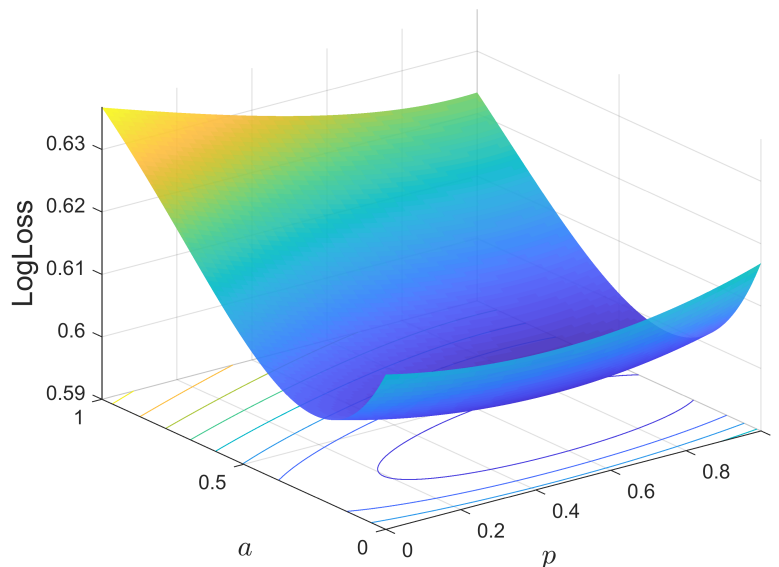
V tabulce 5.4 jsou, kromě hodnoty logaritmické ztrátové funkce, uvedena všechna získaná kritéria kvality modelu pro nalezené optimální parametry (zaokrouhlené na dvě desetinná místa).

Tabulka 5.4: Kritéria kvality modelu 4, pro optimální parametry – optimalizace

	Hodnota kritéria
LogLoss	0,5968
Brierovo skóre	0,4112
AUC-ROC	0,7290
Kalibrační poměr	1,0083
Přesnost	0,6726

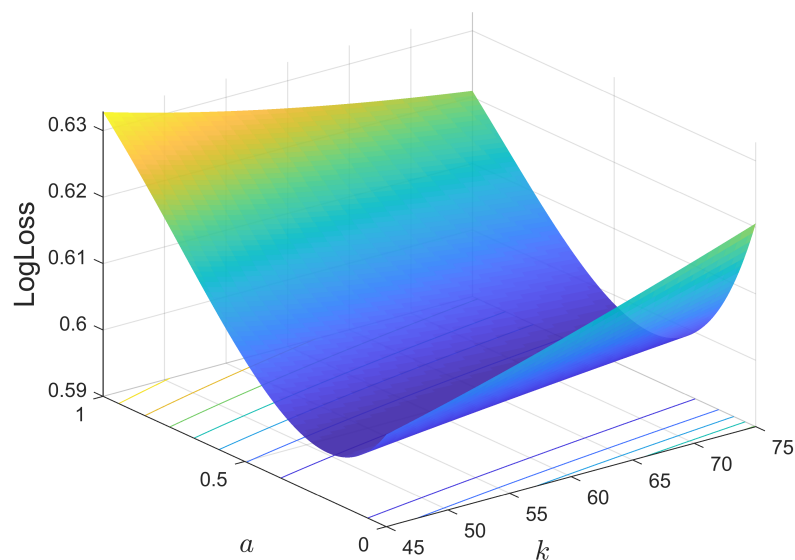
Oproti modelu 3, došlo opět ke snížení minimalizačních kritérií LogLoss a Brierova skóre, stejně jako ke zvýšení maximalizačního kritéria AUC-ROC. Přestože je hodnota kalibračního poměru v absolutní hodnotě vzdálenější od hodnoty 1, je tento poměr stále přibližně roven hodnotě 1. Dále došlo k mírnému poklesu kritéria přesnosti, konkrétně se jedná o pokles o 0,0014 – avšak toto kritérium není příliš vypovídající o přesnosti odhadnutých pravděpodobností, viz odstavec 5.1.2.

Logaritmická ztrátová funkce byla v okolí nalezeného optima vykreslena v závislosti na všech možných kombinacích dvou parametrů. Z grafů bylo zjištěno, že hodnota LogLoss funkce je více než na změnu parametru p v okolí optima (tak jak tomu bylo v předchozích modelech) citlivá na změnu parametru a , viz obrázek 5.8.



Obrázek 5.8: Logaritická ztrátová funkce v závislosti na a a p v okolí nalezeného optima pro model 4

Na obrázku 5.9 je pak zobrazena LogLoss funkce v okolí nalezeného optima v závislosti na parametru a v kombinaci s parametrem k . Na samostatné změny parametrů k , H a B v okolí nalezeného optima je LogLoss funkce opět citlivá jen minimálně. Přestože byla citlivost na změnu parametru a vyšší než na ostatní parametry, nejedná se v blízkém okolí optima o zásadní citlivost, díky které by při malé změně parametru a LogLoss funkce výrazně vzrostla.



Obrázek 5.9: Logaritická ztrátová funkce v závislosti na k a a v okolí nalezeného optima pro model 4

5.2.5 Model 5

Posledním uvažovaným modelem bylo rozšíření modelu 4 o tzv. MOV multiplikátor $v(MOV, d_w)$ koeficientu K navrženým FiveThirtyEight, viz odstavec 4.5. Díky tomuto multiplikátoru je koeficient rozvoje K v každém jednotlivém utkání závislý na MOV , což je tzv. *margin of victory* (rozdíl výsledného skóre zápasu, ve prospěch vítězného týmu) a navíc na rozdílu aktuálního ratingu dvojice soupeřících týmů d_w . Jak již bylo uvedeno v odstavci 4.5, díky nejasnosti odvození MOV multiplikátoru navrženého FiveThirtyEight nebyly optimalizovány parametry, které tento multiplikátor obsahuje, ale byl pouze použit ve svém základním navrženém tvaru, viz výraz (4.14), ve snaze rozšířit model o další faktory. Koeficient rozvoje je tak v tomto modelu ve tvaru

$$K = k \cdot m(x) \cdot v(MOV, d_w) = k \cdot \frac{1}{x^a} \cdot \frac{(MOV + 3)^{0,8}}{7,5 + 0,006 \cdot d_w}, \quad (5.6)$$

více v odstavci 4.5, výraz (4.18). Optimalizovány tak byly opět stejné parametry jako v předchozím modelu, tedy parametry: k , p , H , B , a .

S využitím funkce `fmincon` bylo metodou vnitřního bodu nalezeno lokální minimum logaritmické ztrátové funkce, a to po 60 iteracích. Jako počáteční byly nastaveny hodnoty $k_0 = 75$, $p_0 = 0,5$, $H_0 = 75$, $B_0 = 75$ a $a_0 = 0,5$, tedy středy intervalů zvolených omezení $k \in \langle 0; 150 \rangle$, $p \in \langle 0; 1 \rangle$, $H \in \langle 0; 150 \rangle$, $B \in \langle 0; 150 \rangle$ a $a \in \langle 0; 1 \rangle$. Optimalizace byla opět úspěšně zastavena, protože funkce byla ve všech přípustných směrech neklesající vzhledem k výchozím hodnotám optimalizačních a omezujících tolerancí, viz [32]. Jako minimalizační, resp. optimalizační parametry, byly nalezeny hodnoty (zaokrouhleno na dvě desetinná místa)

- $k = 51,61$,
- $p = 0,75$,
- $H = 62,53$,
- $B = 30,69$,
- $a = 0,29$.

Ve snaze ověřit, že se nejedná pouze o lokální minimum, které je ovlivněno zvolenými počátečními podmínkami, byly testovány další různé kombinace počátečních hodnot parametrů. Testováno bylo dalších pět sad počátečních hodnot parametrů, přičemž počáteční hodnoty jednotlivých parametrů byly vybírány vždy náhodně z rovnoměrného rozdělení na intervalech zvolených omezení jednotlivých parametrů, viz výše. Pro všechny testované kombinace vstupních parametrů byly nalezeny – na dvě desetinná místa – stejné optimalizované parametry se stejnými hodnotami kritérií kvality modelu. Nalezený parametr k se díky přidanému MOV multiplikátoru snížil o 9,90. Hodnoty ostatních parametrů se změnilo jen mírně, nejvýraznější změnu (vzhledem k významu) zaznamenal parametr p , a to nárůst o 0,09. Parametr H se zvýšil pouze o 0,45 a parametr B se snížil o 1,37. Hodnota parametru a zůstala rovněž téměř stejná, snížila se o 0,01. V tabulce 5.5 jsou,

kromě hodnoty logaritmické ztrátové funkce, uvedena všechna získaná kritéria kvality modelu pro nalezené optimální parametry (zaokrouhlené na dvě desetinná místa).

Tabulka 5.5: Kritéria kvality modelu 5, pro optimální parametry – optimalizace

	Hodnota kritéria
LogLoss	0,5926
Brierovo skóre	0,4075
AUC-ROC	0,7349
Kalibrační poměr	1,0067
Přesnost	0,6743

Oproti modelu 4 došlo opět ke snížení minimalizačních kritérií LogLoss a Brierova skóre, stejně jako ke zvýšení maximalizačního kritéria AUC-ROC. Oproti modelu 4 došlo tentokrát také ke zvýšení kritéria přesnosti a hodnota kalibračního poměru je v absolutní hodnotě bližší hodnotě 1.

Logaritmická ztrátová funkce byla v okolí nalezeného optima vykreslena v závislosti na všech možných kombinacích dvou parametrů. Stejně jako v případě modelu 4 bylo z grafů ověřeno, že hodnota LogLoss funkce je v okolí nalezeného optima více citlivá na změnu parametru a než na změnu všech ostatních parametrů. Míra citlivosti změny všech parametrů je obdobná jako v modelu 4. Tedy kromě změny parametru a je hodnota LogLoss funkce v okolí nalezeného optima jen mírně citlivá na změny parametru p . Na samostatné změny parametrů k , H a B v okolí nalezeného optima je LogLoss funkce opět citlivá jen minimálně. Opět platí, že přestože je citlivost na změnu parametru a vyšší než na ostatní parametry, nejedná se v blízkém okolí optima o zásadní citlivost, díky které by při malé změně parametru a LogLoss funkce výrazně vzrostla.

5.3 Výběr modelu

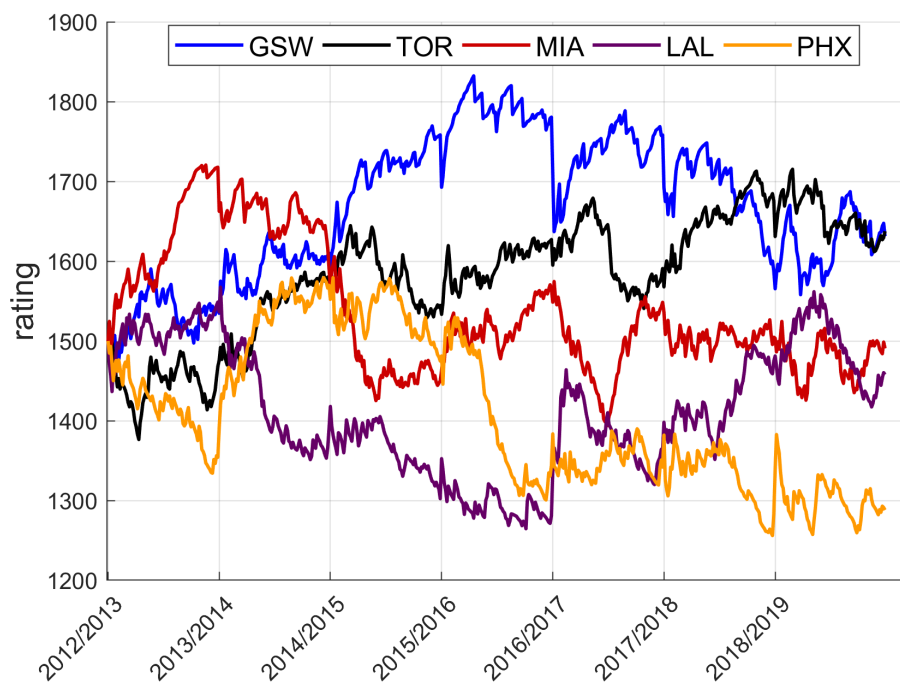
Pro všechny uvažované modely s nalezenými optimálními parametry byla dopočítána všechna uvažovaná kritéria kvality modelů, která jsou v jednotlivých tabulkách 5.1–5.5. Kritéria kvality modelů byla počítána v optimalizační sadě dat, tedy ve všech utkáních základních částí sezon 2013/2014, 2014/2015 a 2015/2016, což bylo celkem 3 690 utkání, přičemž předchozí sezona 2012/2013 sloužila k učení modelu – na začátku této sezony byl všem týmům přiřazen stejný počáteční rating $R_0 = 1\,500$. V tabulce 5.6 jsou shrnuty všechny výsledné hodnoty kritérií kvality jednotlivých modelů včetně modelu Elo538, což je přesný Elo model navržený FiveThirtyEight, tedy bez rozšíření a optimalizace parametrů, viz odstavec 4.6.

Tabulka 5.6: Kritéria kvality všech uvažovaných modelů, včetně Elo538 – optimalizace

	Model 1	Model 2	Model 3	Model 4	Model 5	Elo538
LogLoss	0,6158	0,5994	0,5983	0,5968	0,5926	0,5977
Brierovo skóre	0,4273	0,4132	0,4124	0,4112	0,4075	0,4118
AUC-ROC	0,7255	0,7263	0,7272	0,7290	0,7349	0,7332
Kalibrační poměr	1,0006	0,9974	1,0014	1,0083	1,0067	1,0294
Přesnost	0,6572	0,6753	0,6740	0,6726	0,6743	0,6732

Poznámka: Tučně zvýrazněné jsou nejlepší hodnoty pro každé kritérium kvality modelu – pro LogLoss funkci a Brierovo skóre jsou to hodnoty minimální, pro AUC-ROC a přesnost jsou to hodnoty maximální a pro kalibrační poměr je to hodnota nejbližší (z pohledu absolutní hodnoty) hodnotě 1.

Je zřejmé, že s postupným rozšiřováním modelu (od modelu 1 až do posledního rozšíření, tedy modelu 5) o přídatné parametry vždy klesala hodnota LogLoss funkce a Brierova skóre. Stejně tak se vždy zlepšovala, tedy zvyšovala hodnota AUC-ROC. Tyto hodnoty tak vyšly nejlepší pro model 5. Hodnota kritéria přesnosti byla pro model 5 až jako druhá nejvyšší, po modelu 2 – avšak toto kritérium není příliš vypovídající o přesnosti odhadnutých pravděpodobností, viz odstavec 5.1.2. Přestože hodnota kalibračního poměru pro model 5 není z hlediska absolutní hodnoty ze všech uvažovaných modelů nejbližší hodnotě 1 (je až 4. v pořadí), je stále přibližně rovna této hodnotě (1,0067). Na základě získaných hodnot kritérií kvality uvažovaných modelů v optimalizační sadě dat bylo rozhodnuto o **výběru modelu 5 jako nejlepšího** z uvažovaných. Na obrázku 5.10 je graf vývoje ratingu 5 vybraných reprezentativních týmů⁴ v celém sledovaném období, podle vybraného modelu 5.



Obrázek 5.10: Vývoj ratingu vybraných týmů podle modelu 5 v celém sledovaném období

⁴ Jedná se o týmy Golden State Warriors (GSW), Toronto Raptors (TOR), Miami Heat (MIA), Los Angeles Lakers (LAL) a Phoenix Suns (PHX).

Kapitola 6

Predikční schopnost modelu

V předchozí kapitole byl vybrán model 5 s nalezenými optimálními parametry $k = 51, 61$, $p = 0, 75$, $H = 62, 53$, $B = 30, 69$ a $a = 0, 29$, jako nejlepší model ze všech uvažovaných pro optimalizační sadu dat. V následujících odstavcích bude použitelnost tohoto modelu demonstrována na fiktivním použití proti sázkové kanceláři, čímž bude zprostředkovaně ověřena predikční schopnost tohoto modelu ve srovnání s modelem sázkové kanceláře. Pro sezony v predikční sadě dat, tedy 2016/2017, 2017/2018 a 2018/2019, ve kterých bude demonstrována predikční schopnost modelu, nebylo žádné utkání, pro které by nebyl dohledán vypsaný otevírací kurz sázkové kanceláře Bet365, proto budeme model sázkové kanceláře označovat jako *model Bet365*. Pro každou sezonu z predikční sady dat (2016/2017, 2017/2018 a 2018/2019) bude predikční schopnost optimalizovaného modelu 5 demonstrována jednotlivě. Pro každou sezonu budeme ve třech odstavcích zkoumat predikční schopnost modelu z různých pohledů. Nejprve budou pro srovnání – podobně jako pro optimalizační sadu dat v kapitole 5 – dopočítány hodnoty kritérií kvality modelu také pro jednotlivé sezony v predikční sadě dat, a to optimalizovaného modelu 5, modelu Bet365 (zprostředkovaně pomocí vypsaných kurzů), modelu Elo538 (původního Elo modelu navrženého od FiveThirtyEight, viz odstavec 4.6) a pro sezonu 2018/2019 také nejnovějšího predikčního modelu pro NBA od FiveThirtyEight, a to modelu RAPTOR⁵. K tomuto účelu byla v MATLABu naprogramována funkce `porovnaní_kriterii.m`, viz příloha A.9. Ve druhém odstavci bude analyzováno, zda je možné vybrat utkání, na která v sezoně fiktivně vsadit u sázkové kanceláře Bet365 tak, aby bylo možné dosáhnout díky výběru vhodných zápasů – na základě odhadnutých pravděpodobností pomocí optimalizovaného modelu 5 – zisku na konci dané sezony, resp. „porazit“ sázkovou kancelář. K tomuto účelu byla v MATLABu naprogramována funkce `vyber_zapasu.m`, viz příloha A.10. Třetí odstavec bude věnován predikční schopnosti z pohledu sázkové kanceláře. Na základě odhadnutých pravděpodobností pomocí optimalizovaného modelu 5 budou „vypsány“ decimální kurzy na všechna utkání, vždy s odpovídající marží sázkové kanceláře Bet365. Generování pak budou „náhodní sázející“, kteří budou fiktivně sázet vždy u sázkové kanceláře Bet365 a zároveň budou identicky sázet u virtuální kanceláře s námi vypsanými kurzy. Díky tomu bude možné srovnat

⁵Soubory obsahující data odhadů pravděpodobností vítěství všech týmů ve všech utkáních ze sledovaného období jsou přiloženy jako `data_model5.mat` (tento soubor obsahuje také vypsané kurzy sázkovou kanceláři Bet365, viz příloha A.6), `data_elo538.mat` (příloha A.7) a `data_raptor538.mat` (příloha A.8)

úspěšnost predikce optimalizovaného modelu 5 a modelu Bet365. K tomuto účelu byla v MATLABu naprogramována funkce `nahodny_sazejici.m`, viz příloha A.11.

6.1 Sezona 2016/2017

6.1.1 Kritéria kvality modelů

Aby bylo možné optimalizovaný model 5 a model Elo538 srovnat – z hlediska kritérií kvality modelu – s modelem Bet365, bylo potřeba dopočítat kritéria kvality také pro model Bet365, tedy z dostupných vypsáných kurzů určit odhadnuté pravděpodobnosti vítězství sázkovou kanceláří Bet365. Pokud budeme předpokládat konstantní marži sázkové kanceláře na obě varianty vítězství/prohra (což v realitě nemusí být pravda), odhad pravděpodobnosti vítězství q_i týmů i sázkovou kanceláří lze vyjádřit ze vzorce pro určení (pro nás známého) decimálního kurzu

$$o_i = \frac{1 - \zeta}{q_i}, \quad (6.1)$$

kde o_i je vypsáný decimální kurz sázkovou kanceláří na vítězství týmu i a ζ je hrubá zisková míra sázkové kanceláře, tzv. *marže*. V tomto odhadu nejsou zohledněny další faktory, které ovlivňují výši vypsáného kurzu sázkovou kanceláří, např. kurzy konkurence nebo objem sázek na dané varianty. Výpočet odhadnuté pravděpodobnosti sázkové kanceláře z vypsáných kurzů je ilustrován v příkladu 6.1.

Příklad 6.1. Uvažujeme utkání mezi týmy i a j , přičemž na vítězství týmu i byl sázkovou kanceláří vypsán kurz $o_i = 1,5$ a na vítězství týmu j byl sázkovou kanceláří vypsán kurz $o_j = 2,7$. Odhad konstantní marže sázkové kanceláře v daném utkání lze určit ze soustavy rovnic

$$o_i = \frac{1 - \zeta}{q_i}, \quad (6.2)$$

$$o_j = \frac{1 - \zeta}{q_j}, \quad (6.3)$$

$$q_i + q_j = 1, \quad (6.4)$$

podmínka $q_i + q_j = 1$ se zde vyskytuje díky tomu, že se jedná o pravděpodobnosti, a to pouze dvou možných výsledků (buď vyhraje tým i , nebo tým j). Po dosazení rovnic (6.2) a (6.3) do rovnice (6.4) a následné úpravě, získáváme vzorec pro odhad marže ve tvaru

$$\zeta = 1 - \frac{o_i \cdot o_j}{o_i + o_j}. \quad (6.5)$$

V našem případě je tedy po dosazení odhad konstantní marže sázkové kanceláře

$$\zeta = 1 - \frac{o_i \cdot o_j}{o_i + o_j} = 1 - \frac{1,5 \cdot 2,7}{1,5 + 2,7} \doteq 3,57\%.$$

Odhadnuté pravděpodobnosti sázkové kanceláře pak lze jednoduše vyjádřit z rovnic (6.2) a (6.3) jako

$$q_i = \frac{1 - \zeta}{o_i} = \frac{1 - 0,0357}{1,5} \doteq 64,29\%,$$

$$q_j = \frac{1 - \zeta}{o_j} = \frac{1 - 0,0357}{2,7} \doteq 35,71\%.$$

Tímto postupem určené odhady pravděpodobnosti vítězství sázkovou kancelář Bet365 byly použity k určení kritérií kvality modelu Bet365.

Hodnoty kritérií kvality modelu 5 s optimálními parametry nalezenými v odstavci 5.2.5 jsou v porovnání s hodnotami kritérií kvality modelů Bet365 a Elo538 pro sezonu 2016/2017 v tabulce 6.1. V této sezoně bylo 8 utkání, pro která byl odhad pravděpodobnosti vítězství Bet365 shodně 50% pro oba týmy – tato utkání byla vyřazena z výpočtu kritérií přesnost a kalibrační poměr pro model Bet365.

Tabulka 6.1: Kritéria kvality modelu 5, Elo538 a Bet365, pro sezonu 2016/2017 – predikce

	Model 5	Elo538	Bet365
LogLoss	0,6252	0,6299	0,6117
Brierovo skóre	0,4357	0,4396	0,4248
AUC-ROC	0,6893	0,6866	0,7060
Kalibrační poměr	1,0308	1,0574	1,0227
Přesnost	0,6504	0,6488	0,6604

Hodnoty všech kritérií kvality optimalizovaného modelu 5 byly lepší než pro původní model Elo538. Naopak byly hodnoty všech kritérií horší než pro model sázkové kanceláře Bet365. Predikční schopnost modelu 5 tak byla vyšší než modelu Elo538, ale nižší než sázkové kanceláře Bet365. To, že bude predikční schopnost našeho modelu nižší než Bet365 bylo možné předpokládat. Na rozdíl od sázkové kanceláře, která používá komplexnější predikční modely zahrnující mnohem víc faktorů a investování finančních prostředků, náš model používá pouze historické výsledky utkání, informaci o prostředí (domáci/hosté), back-to-back situaci a rozdílů skóre. Přesto se však podařilo z hlediska kritérií kvality modelu v porovnání s původním modelem Elo538 přiblížit v této sezoně odhadům sázkové kanceláře.

6.1.2 Výběr utkání pro fiktivní sázení

Druhým bodem bylo analyzování, zda je možné vybrat utkání, na která v sezoně vsadit u sázkové kanceláře Bet365, tak, aby bylo možné dosáhnout díky výběru vhodných zápasů – na základě odhadnutých pravděpodobností pomocí našeho modelu 5 – zisku na konci dané sezony, resp. „porazit“ sázkovou kancelář. Utkání, na která vsadit, byla vybírána na základě očekávané hodnoty návratnosti μ z daného utkání, viz [15]. Přírodním požadavkem je sázet pouze na utkání, pro která je očekávaná hodnota návratnosti

alespoň 1. Požadovaná očekávaná hodnota návratnosti L však může být i vyšší než 1. Pravidlo, podle kterého vybírat utkání, na která vsadit, lze zapsat vzhledem k požadované očekávané míře návratnosti $L \geq 1$ jako

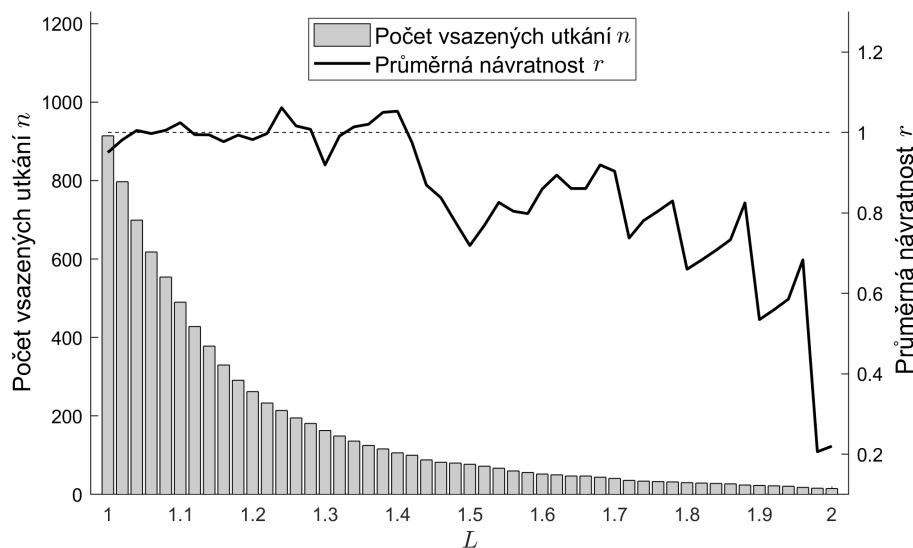
$$\mu = p_i \cdot o_i \geq L, \quad (6.6)$$

kde p_i je odhadnutá pravděpodobnost vítězství v utkání týmu i , v našem případě modelem 5 a o_i je vypsany decimální kurz na vítězství v utkání týmu i , v našem případě sázkovou kanceláří Bet365. V každém utkání je podmínka (6.6) splněna vždy nejvýše pro jeden tým. S rostoucí hodnotou L z logiky věci klesá počet utkání splňující tuto podmínku, resp. klesá počet vybraných utkání, na která vsadit. Na obrázku 6.1 je v závislosti na různých hodnotách L graf počtu vsazených utkání n v sezoně 2016/2017 (1230 utkání) a průměrná návratnost r z jednoho vsazeného utkání, při sázce vždy jedné jednotky na tým i (splňující podmínku (6.6)), při vypsání kurzu o_i na vítězství týmu i , tedy

$$r = \frac{1}{n} \cdot \sum_{k=1}^n r_k, \quad (6.7)$$

kde r_k je návratnost z k tého vsazeného utkání, tedy

$$r_k = \begin{cases} o_i, & \text{pokud tým } i \text{ ve vsazeném utkání } k \text{ zvítězil,} \\ 0, & \text{pokud tým } i \text{ ve vsazeném utkání } k \text{ prohrál.} \end{cases} \quad (6.8)$$



Obrázek 6.1: Počet vsazených utkání a průměrná návratnost v závislosti na L pro sezonu 2016/2017

Je zřejmé, že průměrná návratnost r je vyšší než 1 pouze pro ojedinělé hodnoty parametru L . Nejvyšší hodnota r byla dosažena pro $L = 1,24$, a to $r = 1,06$. Pro tuto hodnotu byl počet vsazených utkání $n = 214$, tedy 17,4% ze všech utkání v sezoně 2016/2017. Vzhledem k ojedinělosti hodnot parametru L , pro které je $r > 1$, lze konstatovat, že pomocí optimalizovaného modelu 5 nelze „porazit“ v sezoně 2016/2017

sázkovou kancelář Bet365, resp. model 5 je možné použít pro zisk proti Bet365 pouze pro ojedinělé hodnoty L .

6.1.3 Náhodný sázející

Třetím bodem bylo zhodnocení odhadnutých pravděpodobností optimalizovaným modelem 5 z pohledu sázkové kanceláře ve srovnání se sázkovou kancelář Bet365. Na základě odhadnutých pravděpodobností vítězství pomocí optimalizovaného modelu 5 byly „vypsány“ decimální kurzy na všechna utkání v sezoně 2016/2017, vždy s odpovídající marží sázkové kanceláře Bet365 (průměrná marže Bet365 v sezoně 2016/2017 byla 3,93%), která byla pro jednotlivá utkání odhadnuta podle vzorce (6.1). Vlastní „vypsané“ decimální kurzy u_i na vítězství týmu i pak byly pro jednotlivá utkání dopočteny podobně jako ve vzorci (6.1), tedy

$$u_i = \frac{1 - \zeta}{p_i}, \quad (6.9)$$

kde ζ je odpovídající marže Bet365 na dané utkání a p_i je pravděpodobnost vítězství v utkání týmu i , odhadnutá pomocí optimalizovaného modelu 5. Cílem bylo porovnat, zda by za předpokladu „náhodného sázejícího“ byla po skončení sezony více v zisku sázková kancelář Bet365 s jejich vypsány kurzy o_i (resp. s odhadem pravděpodobností podle modelu Bet365), nebo virtuální sázková kancelář s „vypsány“ kurzy u_i , které byly určeny na základě odhadů pravděpodobností podle optimalizovaného modelu 5.

V softwaru MATLAB bylo generováno 10 000 „náhodných sázejících“, resp. simulací. V každé simulaci bylo o každém utkání sezony 2016/2017 nejprve rozhodnuto, zda bude na utkání vsazeno, či nikoliv, a to na základě generování náhodné hodnoty z alternativního rozdělení s různými hodnotami parametrů p . Konkrétně bylo testován vždy 10 000 simulací pro každou hodnotu parametru alternativního rozdělení $p \in \left\{ \frac{1}{10}, \frac{2}{10}, \frac{3}{10}, \dots, 1 \right\}$, resp. pro $p = 1$ je vsazeno vždy na všechna utkání. Následně bylo pro každou simulaci v každém utkání, pro které bylo rozhodnuto o sázce, určeno, na který tým „náhodný sázející“ vsadí, přičemž o týmu bylo rozhodnuto na základě generování náhodné hodnoty z alternativního rozdělení s parametrem 0,5. Poslední rozhodnutí „náhodného sázejícího“ spočívalo v určení velikosti částky, kterou vsadí. Tato hodnota byla generována z rovnoměrného rozdělení na intervalu od 0 do 1 jednotky.

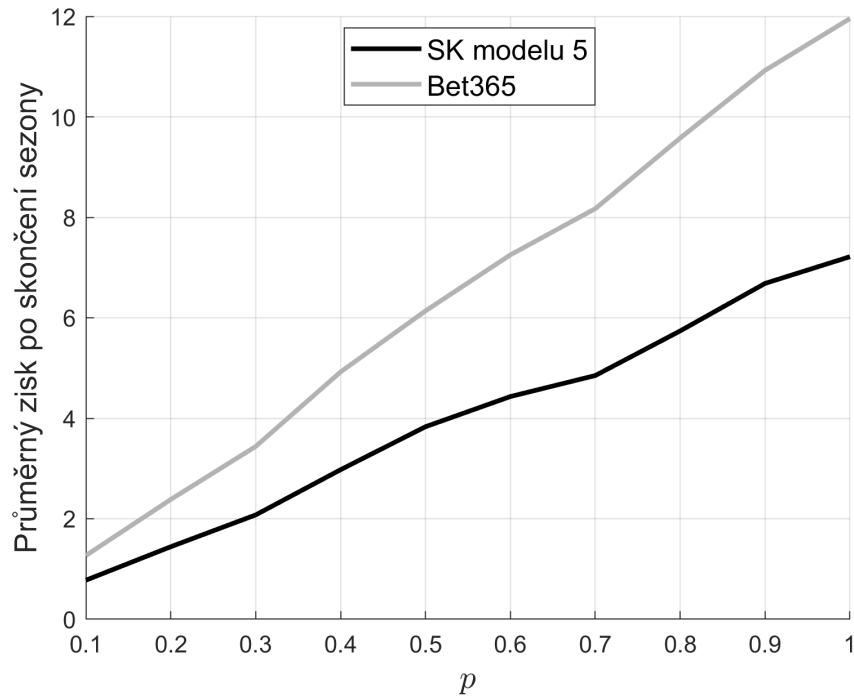
Na každé vybrané utkání a tým „náhodný sázející“ fiktivně vždy vsadil u sázkové kanceláře Bet365 s decimálním kurzem o_i – vypsáným touto sázkovou kancelář na vítězství týmu i – a zároveň stejnou částku na vítězství stejného týmu i u virtuální sázkové kanceláře s decimálním kurzem u_i , který byl určen na základě odhadnutých pravděpodobností podle optimalizovaného modelu 5. Tuto virtuální sázkovou kancelář označme jako *SK modelu 5*. V tabulce 6.2 jsou pro srovnání hodnoty průměrného zisku sázkové kanceláře Bet365 a SK modelu 5 po skončení sezony 2016/2017 na jednoho „náhodného sázejícího“ z celkově 10 000 simulací pro každou uvažovanou hodnotu parametru p (parametru

alternativního rozdělení pro rozhodování o vsazení na utkání). Vzhledem k náhodnosti výše sázek slouží tyto hodnoty průměrných zisků primárně pro porovnání sázkové kanceláře Bet365 a SK modelu 5, nelze je interpretovat přímo jako výše zisků. Zároveň jsou v tabulce uvedeny výběrové směrodatné odchylky z daných simulací. Ve sloupci *úspěšnější v porovnání* je pak zachycen procentuální podíl z daných 10 000 simulací, kdy byla na konci sezony 2016/2017 úspěšnější daná sázková kancelář v porovnání s druhou uvažovanou. Termín *úspěšnější* zde znamená vyšší zisk, resp. nižší ztrátu po skončení sezony.

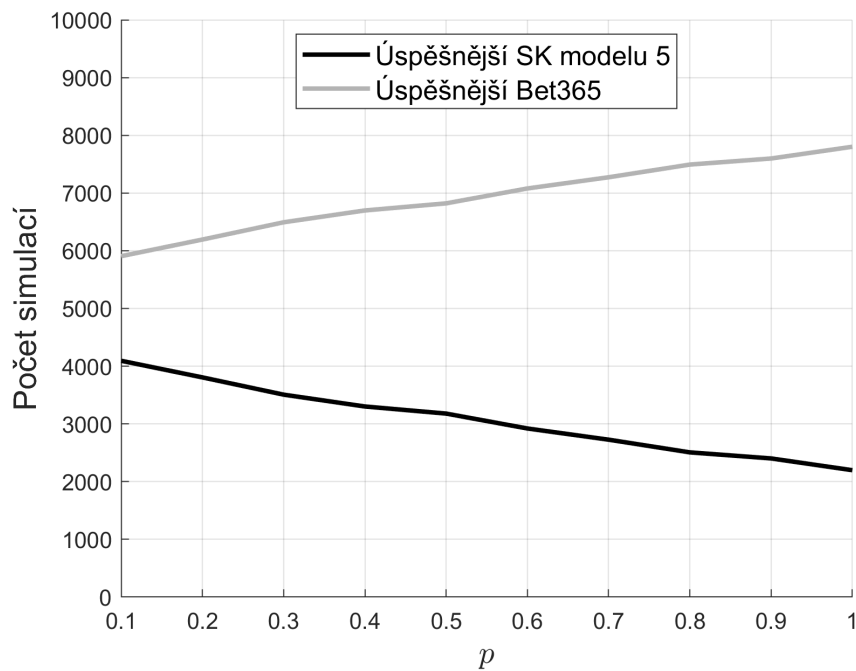
Tabulka 6.2: Srovnání zisku po skončení sezony 2016/2017 SK modelu 5 a Bet365, vždy z 10 000 simulací pro různé hodnoty parametru alternativního rozdělení p

p	Průměrný zisk		Výběrová směrod. odch.		Úspěšnější v porovnání	
	SK modelu 5	Bet365	SK modelu 5	Bet365	SK modelu 5	Bet365
0,1	0,77	1,27	7,98	8,04	40,88 %	59,12 %
0,2	1,43	2,38	11,06	11,14	38,07 %	61,93 %
0,3	2,06	3,44	13,50	13,69	35,00 %	65,00 %
0,4	2,95	4,92	15,57	15,68	33,87 %	67,13 %
0,5	3,80	6,14	17,34	17,42	31,62 %	68,38 %
0,6	4,40	7,25	18,87	18,95	29,06 %	70,94 %
0,7	4,81	8,17	20,00	20,06	27,09 %	72,91 %
0,8	5,69	9,58	21,47	21,52	24,85 %	75,15 %
0,9	6,63	10,93	22,45	22,51	23,78 %	76,22 %
1,0	7,16	11,96	23,41	23,44	21,83 %	78,17 %

Pozitivním zjištěním bylo, že pro všechny testované hodnoty parametru p byl průměrný zisk po skončení sezony kladný. Průměrný zisk byl však pro všechny hodnoty p vyšší pro sázkovou kancelář Bet365 než pro SK modelu 5. S rostoucí hodnotou parametru p – tedy s rostoucím počtem utkání, na která bylo rozhodnuto vsadit – se zvyšoval průměrný zisk pro obě sázkové kanceláře, avšak rozdíl v průměrných ziscích mezi Bet365 a SK modelu 5 se zvyšoval, což je ilustrováno také na obrázku 6.2. Pro všechny hodnoty parametru p byla vždy z 10 000 simulací úspěšnější Bet365 ve více simulacích, viz tabulka 6.2. S rostoucí hodnotou parametru p se podíl úspěšnějších simulací Bet365 oproti SK modelu 5 výrazně zvyšoval, což je ilustrováno na obrázku 6.3.



Obrázek 6.2: Průměrný zisk po skončení sezony 2016/2017 pro SK modelu 5 a Bet365 pro různé hodnoty parametru alternativního rozdělení p



Obrázek 6.3: Počet simulací z celkových 10 000, ve kterých byla pro různé hodnoty parametru alternativního rozdělení p v sezoně 2016/2017 jedna sázková kancelář úspěšnější než druhá uvažovaná

6.2 Sezona 2017/2018

6.2.1 Kritéria kvality modelů

Stejným postupem jako v odstavci 6.1.1 byla dopočítána kritéria kvality modelu Bet365 pro sezonu 2017/2018. Hodnoty kritérií kvality modelu 5 s optimálními parametry nalezenými v odstavci 5.2.5 jsou v porovnání s hodnotami kritérií kvality modelů Bet365 a Elo538 pro sezonu 2017/2018 v tabulce 6.3. V této sezoně bylo 13 utkání, pro která byl odhad pravděpodobnosti vítězství Bet365 shodně 50 % pro oba týmy – tato utkání byla vyřazena z výpočtu kritérií přesnost a kalibrační poměr pro model Bet365.

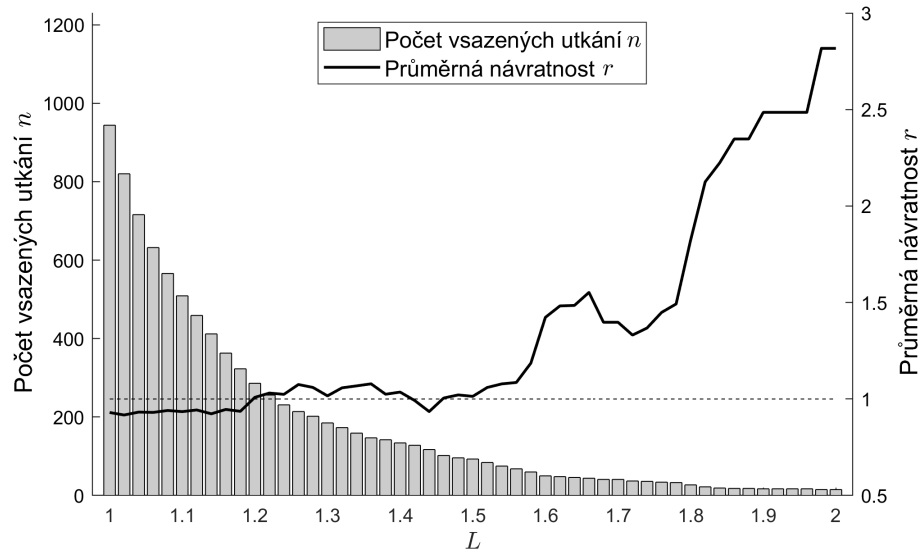
Tabulka 6.3: Kritéria kvality modelu 5, Elo538 a Bet365, pro sezonu 2017/2018 – predikce

	Model 5	Elo538	Bet365
LogLoss	0,6166	0,6203	0,5984
Brierovo skóre	0,4282	0,4305	0,4119
AUC-ROC	0,6994	0,7028	0,7269
Kalibrační poměr	1,0143	1,0289	0,9935
Přesnost	0,6537	0,6618	0,6804

Hodnoty kritérií LogLoss, Brierovo skóre a kalibrační poměr, byly pro optimalizovaný model 5 lepší, než pro původní model Elo538. Hodnota AUC-ROC byla mírně horší, ale rozdíl oproti modelu Elo538 byl pouze o 0,0034 – byla tedy srovnatelná. Podobně na tom byla hodnota přesnosti, která je nižší o 0,0081. Pro většinu kritérií byly tedy hodnoty optimalizovaného modelu 5 lepší, než původního modelu Elo538, a to včetně hodnoty LogLoss, která byla v této práci vybrána k optimalizaci parametrů. Hodnoty všech kritérií kvality pro model Bet365 jsou opět lepší než pro zbývající dva modely. Poznamenejme, že hodnoty všech kritérií pro všechny 3 modely byly lepší, než pro předchozí sezonou 2016/2017. Z pohledu kritérií kvality modelů tak byla predikční schopnost (všech modelů) pro sezonu 2017/2018 vyšší.

6.2.2 Výběr utkání pro fiktivní sázení

Stejně jako v odstavci 6.1.2 bylo analyzováno, zda je možné vybrat utkání, na která v sezoně 2017/2018 vsadit u sázkové kanceláře Bet365, tak, aby bylo možné dosáhnout díky výběru vhodných zápasů zisku na konci dané sezony, resp. „porazit“ sázkovou kancelář. Na obrázku 6.4 je v závislosti na různých hodnotách L graf počtu vsazených utkání n v sezoně 2017/2018 (1230 utkání) a průměrná návratnost r z jednoho vsazeného utkání při sázce vždy jedné jednotky na tým i , splňující podmínku (6.6).



Obrázek 6.4: Počet vsazených utkání a průměrná návratnost v závislosti na L pro sezonu 2017/2018

Průměrná návratnost r se zvyšuje s rostoucí hodnotou parametru L , což ale znamená snižující se počet vsazených utkání n . Pro většinu testovaných hodnot parametru L byla však hodnota $r > 1$, což platí – až na výjimku $L \in \langle 1,42; 1,44 \rangle$ – pro hodnoty $L \geq 1,20$. Vzhledem k tomu tedy můžeme konstatovat, že pro vyšší hodnoty L lze pomocí optimalizovaného modelu 5 „porazit“ v sezoně 2017/2018 sázkovou kancelář Bet365, resp. pro vyšší hodnoty L je možné model 5 použít pro zisk proti Bet365.

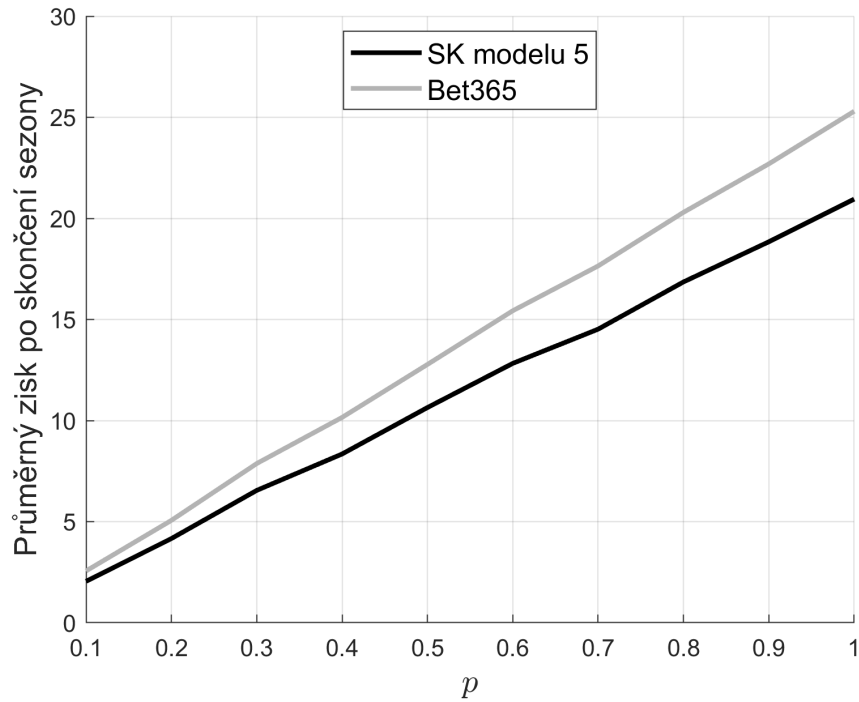
6.2.3 Náhodný sázející

Posledním bodem bylo stejně jako v odstavci 6.1.3 zhodnocení odhadnutých pravděpodobností optimalizovaným modelem 5 z pohledu sázkové kanceláře (SK modelu 5) ve srovnání se sázkovou kancelář Bet365, a to na základě „vypsání“ decimálních kurzů na všechna utkání s odpovídající marží sázkové kanceláře Bet365 (průměrná marže Bet365 v sezoně 2017/2018 byla 3,96%). Opět bylo generováno 10 000 „náhodných sázejících“ (simulací), kteří se rozhodovali podle stejných pravidel jako v odstavci 6.1.3. Každý „náhodný sázející“ opět fiktivně vždy vsadili stejnou částku na stejné utkání a stejný tým u sázkové kanceláře Bet365 a zároveň u SK modelu 5. Testováno bylo znovu několik parametrů p alternativního rozdělení pro rozhodování o vsazení na utkání. V tabulce 6.4 jsou pro srovnání hodnoty průměrných zisků, výběrových směrodatných odchylek a podílů úspěšností sázkových kancelář v porovnání s druhou uvažovanou. Termín *úspěšnější* zde opět znamená vyšší zisk, resp. nižší ztrátu po skončení sezony.

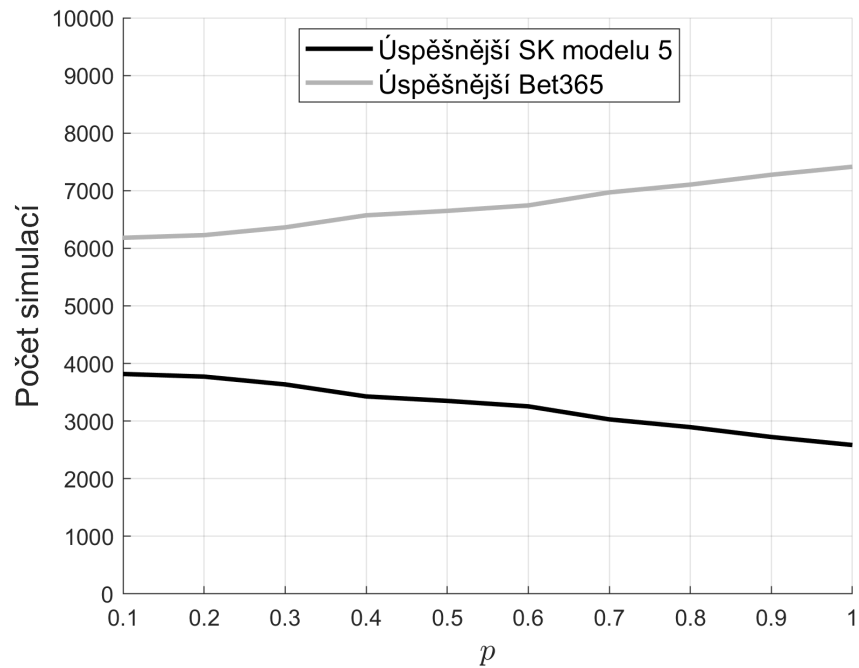
Tabulka 6.4: Srovnání zisku po skončení sezony 2017/2018 SK modelu 5 a Bet365, vždy z 10 000 simulací pro různé hodnoty parametru alternativního rozdělení p

p	Průměrný zisk		Výběrová směrod. odch.		Úspěšnější v porovnání	
	SK modelu 5	Bet365	SK modelu 5	Bet365	SK modelu 5	Bet365
0,1	2,05	2,56	7,67	7,86	38,21 %	61,79 %
0,2	4,17	5,06	10,64	10,89	37,76 %	62,24 %
0,3	6,55	7,88	12,95	13,17	36,39 %	63,61 %
0,4	8,35	10,16	14,82	15,11	34,33 %	65,67 %
0,5	10,65	12,78	16,40	16,80	33,56 %	66,44 %
0,6	12,84	15,43	17,91	18,28	32,69 %	67,31 %
0,7	14,53	17,65	19,40	19,74	30,35 %	69,65 %
0,8	16,87	20,30	20,47	20,84	29,05 %	70,95 %
0,9	18,86	22,70	21,53	21,91	27,27 %	72,73 %
1,0	20,97	25,30	22,49	22,85	25,94 %	74,06 %

Stejně jako v předchozí sezoně byl průměrný zisk po skončení sezony 2017/2018 pro všechny testované hodnoty parametru p kladný. Průměrný zisk byl opět pro všechny hodnoty p vyšší pro sázkovou kancelář Bet365 než pro SK modelu 5. Rozdílem však byla výše průměrného zisku, která byla pro všechny hodnoty p vyšší než v předchozí sezoně. I přes výrazně vyšší hodnoty průměrného zisku pro tuto sezonu, se s rostoucí hodnotou p zvyšoval rozdíl průměrného zisku mezi Bet365 a SK modelu 5 přibližně stejně jako v předchozí sezoně, což je ilustrováno na obrázku 6.5. Pro všechny hodnoty parametru p byla vždy z 10 000 simulací úspěšnější Bet365 ve více simulacích, viz tabulka 6.4. S rostoucí hodnotou p se podíl úspěšnějších simulací Bet365 oproti SK modelu 5 zvyšoval méně výrazně než v předchozí sezoně, což je ilustrováno na obrázku 6.6.



Obrázek 6.5: Průměrný zisk po skončení sezony 2017/2018 pro SK modelu 5 a Bet365 pro různé hodnoty parametru alternativního rozdělení p



Obrázek 6.6: Počet simulací z celkových 10 000, ve kterých byla pro různé hodnoty parametru alternativního rozdělení p v sezoně 2017/2018 jedna sázková kancelář úspěšnější než druhá uvažovaná

6.3 Sezona 2018/2019

6.3.1 Kritéria kvality modelů

Stejným postupem jako pro předchozí sezony byla dopočítána kritéria kvality modelu Bet365 pro sezonu 2018/2019. Hodnoty kritérií kvality modelu 5 s optimálními parametry nalezenými v odstavci 5.2.5 jsou v porovnání s hodnotami kritérií kvality modelů Bet365 a Elo538 pro sezonu 2018/2019 v tabulce 6.5. Pro tuto nejaktuálnější sezonu byla navíc ze zdroje [31] získána data (viz příloha A.8) nejnovějšího modelu RAPTOR od FiveThirtyEight, kterou používají pro predikci utkání od této sezony. Model je založený na projekci jednotlivých hráčů a jejich výkonů do výkonu týmů. V tomto modelu bylo již úplně upuštěno od hodnocení využívajícího Elo rating, více v [27]. Pro srovnání s námi optimalizovaným modelem 5 byla dopočítána kritéria kvality modelu také pro tento model (stejně jako pro ostatní modely pouze pro utkání ze základní části sezony). V této sezoně bylo 14 utkání, pro která byl odhad pravděpodobnosti vítězství Bet365 shodně 50 % pro oba týmy – tato utkání byla vyřazena z výpočtu kritérií přesnost a kalibrační poměr pro model Bet365.

Tabulka 6.5: Kritéria kvality modelu 5, Elo538, RAPTOR a Bet365, pro sezonu 2018/2019 – predikce

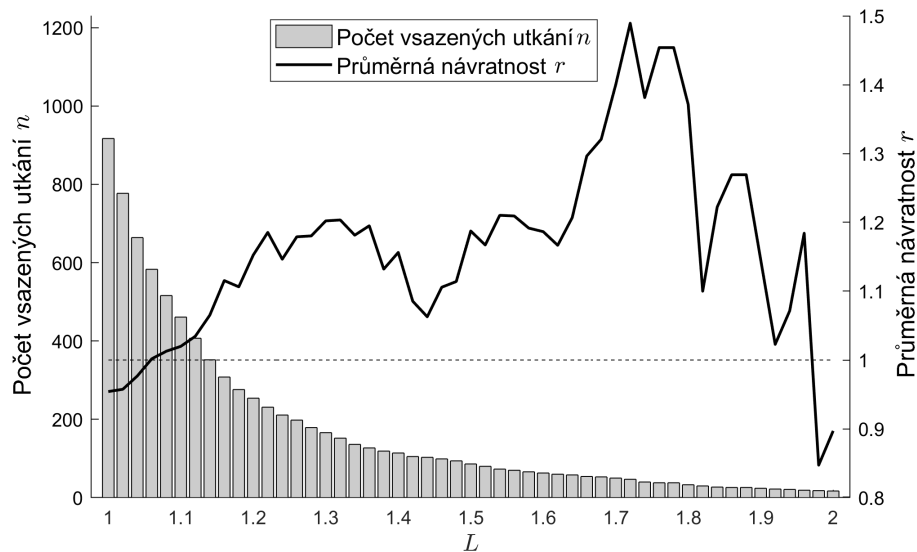
	Model 5	Elo538	RAPTOR	Bet365
LogLoss	0,6137	0,6171	0,6133	0,6002
Brierovo skóre	0,4253	0,4283	0,4245	0,4131
AUC-ROC	0,7023	0,6999	0,7099	0,7222
Kalibrační poměr	1,0178	1,0417	1,0814	1,0103
Přesnost	0,6577	0,6569	0,6561	0,6719

Hodnoty všech kritérií kvality optimalizovaného modelu 5 byly lepší než pro původní model Elo538. Ve srovnání s modelem RAPTOR byla hodnota LogLoss pouze o 0,0004 vyšší a hodnota Brierova skóre pouze o 0,0008 vyšší. Hodnoty kalibračního poměru a přesnosti byly pro náš optimalizovaný model dokonce lepší než pro model RAPTOR (což ale pro tuto sezonu platí i pro model Elo538). Optimalizovaný model 5 se tedy v první sezoně, ve které FiveThirtyEight používají model RAPTOR, tomuto modelu z hlediska kritérií kvality modelu velmi přiblížil a v některých kritériích ho dokonce předčil. Hodnoty všech kritérií pro model Bet365 jsou opět lepší, než pro zbývající 3 modely. Poznamenejme, že kromě hodnoty kalibračního poměru (který je však téměř shodný), byly všechny hodnoty kritérií kvality modelu 5 lepší než v předchozích sezonách.

6.3.2 Výběr utkání pro fiktivní sázení

Stejně jako v odstavcích 6.1.2 a 6.2.2 bylo analyzováno, zda je možné vybrat utkání, na která v sezoně 2018/2019 vsadit u sázkové kanceláře Bet365, tak, aby bylo možné dosáhnout díky výběru vhodných zápasů zisku na konci dané sezony, resp. „porazit“

sázkovou kancelář. Na obrázku 6.7 je v závislosti na různých hodnotách L graf počtu vsazených utkání n v sezoně 2018/2019 (1230 utkání) a průměrná návratnost r z jednoho vsazeného utkání při sázce vždy jedné jednotky na tým i , splňující podmínku (6.6).



Obrázek 6.7: Počet vsazených utkání a průměrná návratnost v závislosti na L pro sezonu 2018/2019

Je zřejmé, že průměrná návratnost r je vyšší než 1 pro většinu testovaných hodnot parametru L . Konkrétně $r > 1$ pro hodnoty $L \in \langle 1,06; 1,96 \rangle$. Vysoké hodnoty parametru L odpovídají velmi nízkému počtu sázek během sezony, proto nejsou příliš vypovídající. Nejvyšší hodnota r byla dosažena pro $L = 1,72$, a to $r = 1,49$. Pro tuto hodnotu byl počet vsazených utkání $n = 47$, tedy 3,82 % ze všech utkání v sezoně 2018/2019. Vzhledem k tomu, že $r > 1$ pro většinu hodnot L , lze konstatovat, že pomocí optimalizovaného modelu 5 je možné „porazit“ v sezoně 2018/2019 sázkovou kancelář Bet365, resp. je možné model 5 použít pro zisk proti Bet365.

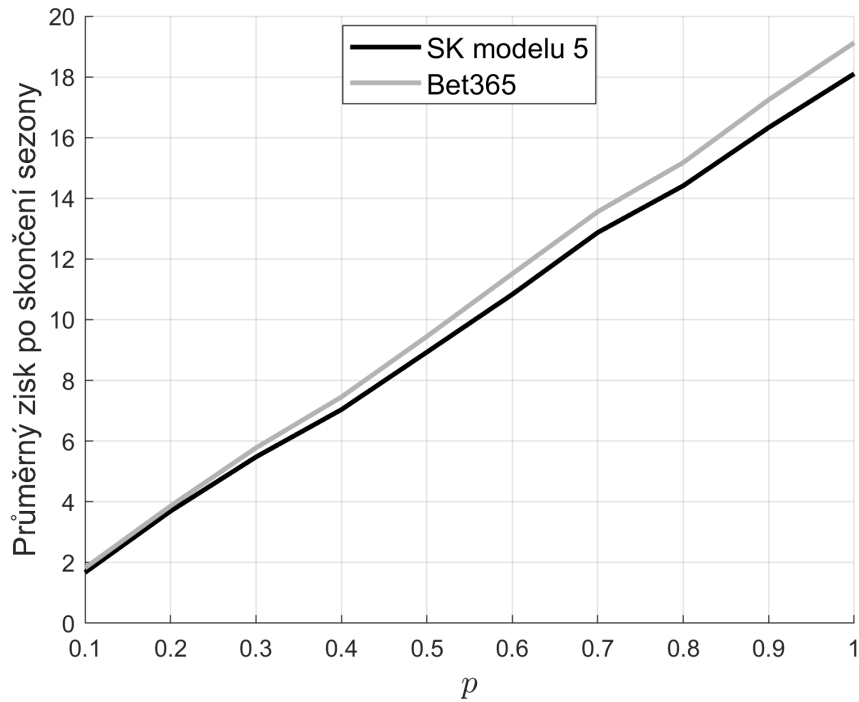
6.3.3 Náhodný sázející

Posledním bodem bylo stejně jako v odstavcích 6.1.3 a 6.2.3 zhodnocení odhadnutých pravděpodobností optimalizovaným modelem 5 z pohledu sázkové kanceláře (SK modelu 5) ve srovnání se sázkovou kancelář Bet365, a to na základě „vypsání“ decimálních kurzů na všechna utkání s odpovídající marží sázkové kanceláře Bet365 (průměrná marže Bet365 v sezoně 2018/2019 byla 4,01 %). Opět bylo stejným postupem jako v odstavci 6.1.3 generováno 10 000 „náhodných sázejících“ (simulací). V tabulce 6.6 jsou pro srovnání hodnoty průměrných zisků, výběrových směrodatných odchylek a podílů úspěšností sázkových kancelář v porovnání s druhou uvažovanou. Termín *úspěšnější* zde opět znamená vyšší zisk, resp. nižší ztrátu po skončení sezony.

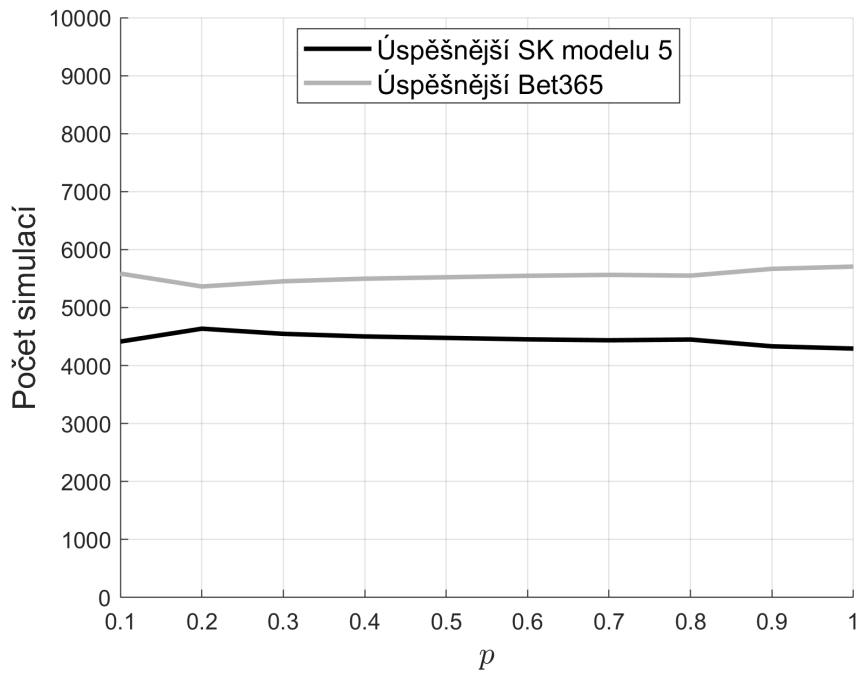
Tabulka 6.6: Srovnání zisku po skončení sezony 2018/2019 SK modelu 5 a Bet365, vždy z 10 000 simulací pro různé hodnoty parametru alternativního rozdělení p

p	Průměrný zisk		Výběrová směrod. odch.		Úspěšnější v porovnání	
	SK modelu 5	Bet365	SK modelu 5	Bet365	SK modelu 5	Bet365
0,1	1,67	1,82	7,72	8,03	44,14 %	55,86 %
0,2	3,69	3,86	10,85	11,31	46,36 %	53,64 %
0,3	5,48	5,78	13,20	13,75	45,47 %	54,53 %
0,4	7,04	7,45	15,14	15,84	45,02 %	54,98 %
0,5	8,93	9,45	17,06	17,74	44,77 %	55,23 %
0,6	10,84	11,52	18,49	19,18	44,52 %	55,48 %
0,7	12,88	13,56	19,54	20,29	44,36 %	55,64 %
0,8	14,41	15,18	20,90	21,64	44,49 %	55,51 %
0,9	16,33	17,25	21,94	22,64	43,32 %	56,68 %
1,0	18,10	19,13	22,96	23,71	42,93 %	57,07 %

Průměrný zisk byl po skončení sezony 2018/2019 opět pro všechny hodnoty parametru p kladný a vyšší pro sázkovou kancelář Bet365 než pro SK modelu 5. Průměrný zisk v sezoně 2018/2019 byl nižší než v sezoně předchozí. Na rozdíl od předchozích dvou sezon je zde výrazně nižší rozdíl v průměrném zisku SK modelu 5 a Bet365. Navíc se rozdíl v průměrném zisku s rostoucí hodnotou p , tedy s vyšším počtem vsazených utkání, zvyšuje jen mírně, což je ilustrováno také na obrázku 6.8. Například pro hodnotu $p = 1$, tedy v případě, že je vsazeno na všechna utkání a rozdíl v průměrných ziscích je nejvyšší, činí tento rozdíl pouze 1,03. Pro všechny hodnoty parametru p byla vždy z 10 000 simulací úspěšnější sázková kancelář Bet365 ve více simulacích, viz tabulka 6.6. S rostoucí hodnotou p se zvyšoval (v porovnání s předchozími sezonami) podíl úspěšnějších simulací Bet365 oproti SK modelu 5 jen minimálně, což je ilustrováno na obrázku 6.9. Lze tedy konstatovat, že pro sezonu 2018/2019 se podařilo s využitím optimalizovaného modelu 5 nejvíce přiblížit – z hlediska zisku z náhodně generovaných sázejících – sázkové kanceláři Bet365.



Obrázek 6.8: Průměrný zisk po skončení sezony 2018/2019 pro SK modelu 5 a Bet365 pro různé hodnoty parametru alternativního rozdělení p



Obrázek 6.9: Počet simulací z celkových 10 000, ve kterých byla pro různé hodnoty parametru alternativního rozdělení p v sezoně 2018/2019 jedna sázková kancelář úspěšnější než druhá uvažovaná

Kapitola 7

Závěr

Tato práce se zabývala predikcí výsledků utkání NBA s využitím systému Elo rating, který byl vyhodnocen jako nejúspěšnější v bakalářské práci [19]. Hlavním cílem této práce bylo navrhnout změny a vylepšení původního Elo rating systému (pro šach), optimalizovat vybrané parametry a srovnat výsledky predikce pomocí upraveného optimalizovaného modelu s modelem sázkové kanceláře Bet365 a již existujícím modelem Elo538 pro predikci výsledků utkání NBA od FiveThirtyEight, založeným na Elo ratingu.

Nejprve byl v kapitole 3 popsán historický vývoj Elo ratingu spolu s matematickým popisem tohoto systému. Díky tomu byl objasněn princip fungování Elo ratingu a postupy, na základě kterých bylo možné odhadovat pravděpodobnosti vítězství týmů v jednotlivých utkáních.

V kapitole 4 bylo navrženo 5 rozšiřujících modelů pro odhad pravděpodobnosti vítězství v utkáních NBA. Všechny modely byly naprogramovány v softwaru MATLAB R2018a. V jednotlivých modelech byly postupně přidávány faktory za účelem zahrnout do původního modelu více informací a pokusit se eliminovat některé nedostatky původního modelu. Do původního modelu pro šach byl v prvním modelu přidán pouze parametr, díky kterému se týmům po začátku nové sezony zachovává určitá část jejich hodnocení, tzv. ratingu, z předchozí sezony. V dalším modelu byl navíc přidán parametr zvyšující pravděpodobnost vítězství domácím týmům, přičemž výhoda domácího prostředí byla ověřena na dostupných datech. Ve třetím modelu byl přidán parametr, který penalizuje týmy vystavené back-to-back situaci, tedy hrající druhé utkání ve dvou dnech. Vliv tohoto faktoru na výkony týmů byl rovněž ověřen na dostupných datech. Čtvrtý uvažovaný model byl rozšířen o tzv. multiplikátor mocninného poklesu. Díky tomuto faktoru se pro všechny týmy v průběhu každé jednotlivé sezony postupně snižuje v modelu tzv. koeficient rozvoje K , který udává maximální možnou hodnotu zvýšení, resp. snížení ratingu týmu při jeho aktualizaci, tedy po každém odehraném utkání. Díky tomu se ratingy jednotlivých týmů v tomto modelu můžou rychleji přizpůsobovat jejich skutečné síle po začátku nové sezony a naopak později být stabilnější a méně citlivé na případné odchylky od stabilní výkonnosti týmů. Pátý a poslední model byl rozšířen o tzv. MOV multiplikátor, převzatý z modelu Elo538. Díky tomuto multiplikátoru závisí velikost koeficient rozvoje K v jednotlivých utkáních také na rozdílu výsledného skóre v utkání a naplnění jeho očekávání.

Vybrané parametry všech uvažovaných modelů byly s využitím MATLABu optimalizovány na reálných datech, a to na základě minimalizace logaritmické ztrátové funkce v kapitole 5. Pro všechny modely s nalezenými optimálními parametry byla dopočítána, kromě logaritmické ztrátové funkce, také další kritéria kvality modelu, konkrétně Brierovo skóre, AUC-ROC, kalibrační poměr a přesnost. Na základě získaných hodnot byl vybrán model 5 s optimálními parametry $k = 51, 61$, $p = 0, 75$, $H = 62, 53$, $B = 30, 69$ a $a = 0, 29$, viz odstavec 4.5, jako nejlepší z uvažovaných modelů.

Predikční schopnost optimalizovaného modelu 5 byla demonstrována v kapitole 6 ze tří různých pohledů, a to celkem ve třech sezonách 2016/2017, 2017/2018 a 2018/2019, pro každou sezonu jednotlivě. Nejprve byla pro každou sezonu srovnána kritéria kvality optimalizovaného modelu 5 s již existujícím a používaným modelem Elo538 a s modelem sázkové kanceláře Bet365 (zprostředkovaně pomocí vypsání kurzů). Kritéria kvality optimalizovaného modelu 5 dosahovala ve všech třech sezonách lepších hodnot, než u původního modelu Elo538, ale horších hodnot, než u sázkové kanceláře Bet365. Z hlediska kritérií kvality modelů se tedy v porovnání s původním modelem Elo538 podařilo přiblížit odhadům sázkové kanceláře. V sezoně 2018/2019 byla kritéria srovnána také s nejnovějším predikčním modelem od FiveThirtyEight, modelem RAPTOR, kterému se model 5 z hlediska kritérií kvality velmi přiblížil a v některých kritériích ho dokonce předčil. Druhým pohledem bylo analyzování, zda je možné vybrat utkání, na která v sezoně fiktivně vsadit u sázkové kanceláře Bet365, tak, aby bylo možné dosáhnout díky výběru vhodných zápasů – na základě odhadnutých pravděpodobností pomocí optimalizovaného modelu 5 – zisku na konci dané sezony, resp. „porazit“ sázkovou kancelář. Z výsledků lze konstatovat, že v sezonách 2017/2018 a 2018/2019 bylo možné sázkovou kancelář „porazit“. V sezoně 2016/2017 nikoliv, protože pouze pro ojedinělé hodnoty požadované očekávané návratnosti L by průměrná návratnost byla vyšší než 1, viz odstavec 6.1.2. Nakonec byla predikční schopnost modelu 5 demonstrována z pohledu sázkové kanceláře. Na základě odhadnutých pravděpodobností pomocí optimalizovaného modelu 5 byly „vypsány“ decimální kurzy na všechna utkání, vždy s odpovídající marží sázkové kanceláře Bet365. Generováno pak bylo 10 000 „náhodných sázejících“, kteří fiktivně sázeli vždy u sázkové kanceláře Bet365 a zároveň identicky sázeli u virtuální kanceláře s námi vypsány kurzy. Pozitivním zjištěním bylo, že průměrný zisk z jednoho sázejícího byl pro obě sázkové kanceláře vždy kladný. Pro sázkovou kancelář Bet365 byl však ve všech případech vyšší. S rostoucím počtem utkání, na která bylo rozhodnuto vsadit (resp. s rostoucí hodnotou parametru alternativního rozdělení pro rozhodování o vsazení na utkání), se ve všech sezonách průměrný zisk obou kancelářů zvyšoval. Rozdíl v průměrném zisku mezi sázkovými kancelářemi se však také zvyšoval (ve prospěch Bet365), stejně jako podíl úspěšnějších simulací (s vyšším ziskem, resp. nižší ztrátou) sázkové kanceláře Bet365, oproti virtuální sázkové kanceláři s námi „vypsány“ kurzy. V sezoně 2018/2019 byl však rozdíl v průměrných ziscích obou sázkových kancelářů jen minimální a zvyšoval se jen mírně. Pro tuto sezonu se tak podařilo s využitím modelu 5 nejvíce přiblížit sázkové kanceláři Bet365.

Vzhledem k tomu, že pouze v sezoně 2016/2017 se nepodařilo „porazit“ sázkovou kancelář, bylo by možné se v této sezoně pokusit identifikovat „problematická“ utkání, ve kterých byla nadhodnocována odhadnutá pravděpodobnost a pokusit se u těchto utkání odhalit, co bylo příčinou. V případě zjištění systematické příčiny by pak např. bylo možné přidat do modelu další faktor, který by tuto příčinu zohledňoval. Dalším námětem by mohlo být rozšíření modelu o další faktory, které nezohledňuje, např. použití dalších zápasových statistik, nebo rozšíření již používaných. Např. vliv domácího prostředí nemusí být uvažován pouze konstantní, ale může být jiný pro každý tým, nebo závislý na vzdálenosti měst, ve kterých týmy hrají. Stejně tak pro vliv back-to-back situací, který by mohl být rozšířen na vliv počtu dní odpočinku mezi zápasy obou soupeřících týmů, nebo zohlednění ratingu (síly) soupeře z posledního utkání, což může mít rovněž vliv na unavenost týmu. V neposlední řadě by námětem mohla být úprava multiplikátoru poklesu koeficientu K , který by bylo možné navrhnout tak, aby lépe odpovídal fázi sezony a zohledňoval například úmyslné prohrávání slabých týmů v pokročilé fázi sezony za účelem lepší pozice při draftové loterii.

Literatura a další zdroje

- [1] Anděl, J. *Statistické metody*. 4. vyd. Matfyzpress, 2007. ISBN 80-7378-003-8.
- [2] Balakrishnan, N. *Handbook of the logistic distribution*. Marcel Dekker, Inc., 1992. ISBN 0-8247-8587-8.
- [3] Beckler, M. – Wang, H. – Papamichael, M. NBA Oracle. *Zuletzt Besucht Am*. 2013, 17, 2008–2009. Dostupné z: https://www.mbeckler.org/coursework/2008-2009/10701_report.pdf.
- [4] Bester, D. W. – Maltitz, M. J. Introducing momentum to the elo rating system. *University of the Free State: Department of Mathematical Statistics and Actuarial Science*. 2013. Dostupné z: https://www.ufs.ac.za/docs/librariesprovider2/mathematical-statistics-and-actuarial-science-documents/technical-reports-documents/teg418-2069-eng.pdf?sfvrsn=243cf921_0.
- [5] Byrd, R. H. – Hribar, M. E. – Nocedal, J. An interior point algorithm for large-scale non-linear programming. *SIAM Journal on Optimization*. 1999, 9, 4, s. 877–900. ISSN 1052-6234. DOI: 10.1137/S1052623497325107.
- [6] Elo, A. E. *The rating of chessplayers, past and present*. Arco Pub., 1978. ISBN 0-668-04721-6.
- [7] Entine, O. A. – Small, D. S. The role of rest in the NBA home-court advantage. *Journal of Quantitative Analysis in Sports*. 2008, 4, 2. ISSN 1559-0410. DOI: 10.2202/1559-0410.1106.
- [8] Fawcett, T. An introduction to ROC analysis. *Pattern recognition letters*. 2006, 27, 8, s. 861–874. ISSN 0167-8655. DOI: 10.1016/j.patrec.2005.10.010.
- [9] Glickman, M. E. The glicko system. *Boston University*. 1998. Dostupné z: http://www.echecsonline.net/joueurs/doc/The_Glicko_system.pdf.
- [10] Godoy, D. *Understanding binary cross-entropy/log loss: a visual explanation* [online]. Towards Data Science, 2018. [cit. 3.3.2020]. Dostupné z: <https://towardsdatascience.com/understanding-binary-cross-entropy-log-loss-a-visual-explanation-a3ac6025181a>.
- [11] Harkness, K. *Official chess handbook*. D. McKay Co., 1956.
- [12] Hátle, J. – Likeš, J. *Základy počtu pravděpodobnosti a matematické statistiky*. SNTL Praha, 1974.
- [13] Kovalchik, S. A. Searching for the GOAT of tennis win prediction. *Journal of Quantitative Analysis in Sports*. 2016, 12, 3, s. 127–138. ISSN 1559-0410. DOI: 10.1515/jqas-2015-0059.

-
- [14] Lemons, D. S. – Langevin, P. *An introduction to stochastic processes in physics*. JHU Press, 2002. ISBN 0-8018-6866-1.
- [15] Marek, P. – Šedivá, B. – Toupal, T. Modeling and prediction of ice hockey match results. *Journal of quantitative analysis in sports*. 2014, 10, 3, s. 357–365. ISSN 1559-0410. DOI: 10.1515/jqas-2013-0129.
- [16] Morris, B. – Bialik, C. – Boice, J. *How We're Forecasting The 2016 U.S. Open* [online]. FiveThirtyEight, 2016. [cit. 10. 4. 2020]. Dostupné z: <https://fivethirtyeight.com/features/how-were-forecasting-the-2016-us-open/>.
- [17] Morse, S. *Autocorrelation in Elo ratings* [online]. GitHub, 2019. [cit. 11. 4. 2020]. Dostupné z: <https://stmorse.github.io/journal/Elo-2.html>.
- [18] Noordman, R. Improving the estimation of outcome probabilities of football matches using in-game information. *Amsterdam School of Economics, Faculty of Economics and Business*. 2019. Dostupné z: <https://www.scisports.com/wp-content/uploads/2019/10/Noordman-Rogier-12366315-MSc-ETRICS.pdf>.
- [19] Ondříček, J. Vybrané prediktivní modely pro výsledky zápasů NBA. Bakalářská práce, Západočeská univerzita v Plzni, Fakulta aplikovaných věd, Plzeň, 2018. Dostupné z: <https://dspace5.zcu.cz/handle/11025/32351>.
- [20] Špaček, J. Modelování a odhadování výsledků tenisových zápasů. Diplomová práce, Západočeská univerzita v Plzni, Fakulta aplikovaných věd, Plzeň, 2018. Dostupné z: <https://otik.uk.zcu.cz/handle/11025/32031>.
- [21] Parkes, D. *The ROC Curve* [online]. The blog of Duncan Parkes, 2018. [cit. 4. 4. 2020]. Dostupné z: <https://deparkes.co.uk/2018/02/16/the-roc-curve/>.
- [22] Pollard, R. Home advantage in soccer: A retrospective analysis. *Journal of sports sciences*. 1986, 4, 3, s. 237–248. ISSN 0264-0414. DOI: 10.1080/02640418608732122.
- [23] Ross, D. *Arpad Elo and the Elo Rating System* [online]. ChessBase Magazine, 2007. [cit. 15. 3. 2020]. Dostupné z: <https://web.archive.org/web/20140729203449/https://en.chessbase.com/post/arpad-elo-and-the-elo-rating-system>.
- [24] Schwartz, B. – Barsky, S. F. The home advantage. *Social forces*. 1977, 55, 3, s. 641–661. ISSN 0037-7732. DOI: 10.1093/sf/55.3.641.
- [25] Silver, N. *Introducing NFL Elo Ratings* [online]. FiveThirtyEight, 2014. [cit. 11. 4. 2020]. Dostupné z: <https://fivethirtyeight.com/features/introducing-nfl-elo-ratings/>.
- [26] Silver, N. – Fischer-Baum, R. *How We Calculate NBA Elo Ratings* [online]. FiveThirtyEight, 2015. [cit. 25. 3. 2020]. Dostupné z: <https://fivethirtyeight.com/features/how-we-calculate-nba-elo-ratings/>.
-

-
- [27] Silver, N. – Jay, B. – Paine, N. *How Our NBA Predictions Work* [online]. FiveThirtyEight, 2019. [cit. 12. 4. 2020]. Dostupné z: <https://fivethirtyeight.com/methodology/how-our-nba-predictions-work/>.
- [28] *Elo Ratings for NBA Teams* [online]. Practically Predictable, 2018. [cit. 22. 3. 2020]. Dostupné z: <http://practicallypredictable.com/2018/04/15/elo-ratings-for-nba-teams/#win-probabilities-and-home-court-advantage>.
- [29] *Elo rating system* [online]. Wikipedia, 2020. [cit. 20. 3. 2020]. Dostupné z: https://en.wikipedia.org/wiki/Elo_rating_system.
- [30] *FIDE Rating Regulations effective from 1 July 2017* [online]. International Chess Federation, 2017. [cit. 20. 3. 2020]. FIDE handbook. Dostupné z: <https://handbook.fide.com/chapter/B022017>.
- [31] *FiveThirtyEight data* [online]. GitHub, 2020. [cit. 12. 4. 2020]. NBA Elo forecast data. Dostupné z: <https://github.com/fivethirtyeight/data/tree/master/nba-forecasts>.
- [32] *fmincon* [online]. MathWorks. [cit. 15. 4. 2020]. Find minimum of constrained nonlinear multivariable function. Dostupné z: <https://www.mathworks.com/help/optim/ug/fmincon.html>.
- [33] *NBA 3in1 database* [online]. Indatabet, 2019. [cit. 11. 11. 2019]. Opening odds data. Dostupné z: <https://www.indatabet.com/bskb-free-3in1.html>.
- [34] *NBA Advanced Stats* [online]. NBA Stats, 2020. [cit. 3. 1. 2020]. Box Scores data. Dostupné z: <http://stats.nba.com/teams/boxscores/>.
- [35] *NBA Results and Historical Odds* [online]. Odds Portal, 2019. [cit. 11. 11. 2019]. Opening odds data. Dostupné z: <http://www.oddsportal.com/basketball/usa/nba/results/>.
- [36] *Šachový ELO systém* [online]. 2007. [cit. 25. 3. 2020]. Dostupné z: https://web.archive.org/web/20110303152542/http://www.juniorbb.iprofil.cz/ELO_CZ.pdf.

Příloha A

Zde je uveden seznam souborů, které se nacházejí na přiloženém CD (ve složce Příloha A) a jejich stručný popis. Jedná se datové soubory a funkce použité pro získání výsledků této práce. Všechny funkce byly napsány a použity v softwaru MATLAB R2018a.

A.1 NBAstats.xlsx

Soubor obsahující přehledně upravená vstupní data včetně decimálních kurzů. Data z každé jednotlivé sezony jsou na samostatném listu vždy v jedné tabulce. Soubor obsahuje rovněž list s tabulkou dat ze všech pozorovaných sezon dohromady. Záznam o každém utkání je na dva řádky. Soubor byl vytvořen v softwaru Excel 2016.

A.2 NBAstats.mat

Soubor obsahující tabulku se všemi získanými daty ve formátu požadovaném pro použití dat v naprogramovaných funkcích v softwaru MATLAB R2018a. Záznam o každém utkání je na dva řádky. Pro zjednodušení práce s daty byly v tomto souboru dat navíc prohozeny řádky tak, aby vždy první z dvojice řádků odpovídající jednomu utkání, odpovídal domácímu týmu. Data lze v MATLABu načíst příkazem `load('NBAstats.mat')`.

A.3 elo.m

Funkce určená k odhadu pravděpodobností pro libovolnou volbu sezon (s učením od počátku sledovaného období, tedy od začátku sezony 2012/2013) pomocí modelu 1 až modelu 5, a to pro libovolné parametry daných modelů. Pro zvolené sezony rovněž vypočítá kritéria kvality modelu. S využitím této funkce byly optimalizovány parametry jednotlivých modelů.

```
[LogLoss, kriteria, ppst] = elo(data, sezony, model, par)
```

- `data` je tabulka se všemi vstupními daty v požadovaném formátu, kterou lze načíst pomocí příkazu `load('NBAstats.mat')`, viz příloha A.2
- `sezony` jsou zvolené sezony, pro které budou odhadovány pravděpodobnosti vítězství v jednotlivých utkáních a počítány kritéria kvality modelu, a to s učením od počátku sledovaného období, tedy od začátku sezony 2012/2013 – např. volba ve tvaru `{'2018/2019'}` nebo pro více sezon `{'2013/2014', '2014/2015', '2015/2016'}`

- `model` je číslo zvoleného modelu, tedy hodnota 1, 2, 3, 4, nebo 5, viz kapitola 4
- `par` je vektor zvolených parametrů v pořadí k (resp. K), p , H , B , a (pro modely s méně parametry bude vektor obsahovat méně prvků, ale pořadí požadovaných parametrů zůstane zachováno) – např. pro optimální parametry modelu 5, ve tvaru $[51,62 \ 0,75 \ 62,53 \ 30,69 \ 0,29]$, viz odstavec 5.2.5
- `LogLoss` je hodnota logaritmické ztrátové funkce ve zvolených sezonách (slouží pro optimalizaci)
- `kriteria` je tabulka se všemi uvažovanými hodnotami kritérií kvality modelu ve zvolených sezonách
- `ppst` je tabulka s výslednými odhady pravděpodobností vítězství týmů a se skutečnými výsledky ve všech utkáních ve zvolených sezonách

A.4 `parametry.m`

Pomocná funkce pro přiřazení parametrů zvolenému modelu, která je spouštěna pouze z výše uvedené funkce `elo.m`.

A.5 `optimize.m`

Funkce určená k optimalizaci parametrů modelu 1 až modelu 5, pro libovolnou volbu sezon v optimalizační sadě dat (s učením od počátku sledovaného období, tedy od začátku sezony 2012/2013). Pro zvolené sezony rovněž vypočítá kritéria kvality modelu s nalezenými optimálními parametry.

```
[par_opt, LogLoss, exitflag, output, kriteria] = optimize(data, model, sezony, random)
```

- `data` je tabulka se všemi vstupními daty v požadovaném formátu, kterou lze načíst pomocí příkazu `load('NBAsstats.mat')`, viz příloha A.2
- `model` je číslo zvoleného modelu, tedy hodnota 1, 2, 3, 4, nebo 5.
- `sezony` jsou zvolené sezony (optimalizační sada dat), ve kterých budou optimalizovány parametry, a to s učením od počátku sledovaného období, tedy od začátku sezony 2012/2013 – např. pro námi uvažovanou optimalizační sadu dat ve tvaru $\{ '2013/2014' , '2014/2015' , '2015/2016' \}$
- `random` je indikátor náhodných počátečních hodnot z rovnoměrných rozdělení na intervalech jednotlivých parametrů $k \in \langle 0; 150 \rangle$, $p \in \langle 0; 1 \rangle$, $H \in \langle 0; 150 \rangle$, $B \in \langle 0; 150 \rangle$ a $a \in \langle 0; 1 \rangle$

- pro volbu `random = 1` jsou počáteční hodnoty generovány náhodně z rovnoměrných rozdělení na výše uvedených intervalech
 - pro volbu `random = 0` jsou počáteční hodnoty nastaveny jako středy výše uvedených intervalů
 - ve zdrojovém kódu funkce `optimize.m` lze v případě potřeby nastavit vlastní hodnoty počátečních hodnot i rozmezí omezujících intervalů jednotlivých parametrů
- `par_opt` je vektor nalezených optimálních parametrů v pořadí k (resp. K), p , H , B , a (pro modely s méně parametry bude vektor obsahovat méně prvků, ale pořadí parametrů zůstane zachováno), pro které má logaritmická ztrátová funkce lokální minimum ve zvolených sezonách
 - `LogLoss` je hodnota logaritmické ztrátové funkce ve zvolených sezonách pro nalezené optimální parametry
 - `exitflag` je indikátor zastavení funkce `fmincon`
 - `output` obsahuje informace o optimalizačním procesu funkce `fmincon`
 - `kriteria` je tabulka se všemi hodnotami kritérií kvality modelu ve zvolených sezonách pro nalezené optimální parametry

A.6 `data_model5.mat`

Soubor obsahující tabulku s výslednými odhady pravděpodobností vítězství týmů ve všech utkáních ze sledovaného období, podle optimalizovaného modelu 5. Tento soubor rovněž obsahuje vypsané kurzy sázkové kanceláře Bet365. Data lze v MATLABu načíst příkazem `load('data_model5.mat')`.

A.7 `data_elo538.mat`

Soubor obsahující tabulku s odhady pravděpodobností vítězství týmů ve všech utkáních ze sledovaného období, podle modelu Elo538 od FiveThirtyEight. Data lze v MATLABu načíst příkazem `load('data_elo538.mat')`.

A.8 `data_raptor538.mat`

Soubor obsahující tabulku s odhady pravděpodobností vítězství týmů ve všech utkáních ze sezony 2018/2019, podle modelu RAPTOR od FiveThirtyEight. Data lze v MATLABu načíst příkazem `load('data_raptor538.mat')`.

A.9 porovnaní_kriterii.m

Funkce určená k porovnání kritérií kvality optimalizovaného modelu 5, modelu Elo538, modelu Bet365, popřípadě modelu RAPTOR (v sezoně 2018/2019), ve zvolených sezonách. Pro spuštění této funkce je potřeba mít ve stejném adresáři soubory s daty: data_model5.mat, data_elo538.mat a data_raptor538.mat.

```
kriteria = porovnaní_kriterii(sezony)
```

- sezony jsou požadované sezony, pro které budou získána kritéria kvality výše uvedených modelů – např. ve tvaru { '2018/2019' } nebo pro více sezon { '2013/2014', '2014/2015', '2015/2016' }
- kriteria je tabulka se všemi hodnotami kritérií kvality výše uvedených modelů, ve zvolených sezonách

A.10 vyber_zapasu.m

Funkce určená k výběru utkání pro fiktivní sázení proti sázkové kanceláři Bet365 (podle pravidla (6.6)) ve zvolených sezonách a výpočtu průměrné návratnosti z jednotlivých vsazených utkání v daných sezonách pro různé hodnoty L . Pro spuštění této funkce je potřeba mít ve stejném adresáři soubor s daty data_model5.mat.

```
tabulka = filtr_zapasu(sezony)
```

- sezony jsou zvolené sezony – např. ve tvaru { '2018/2019' } nebo pro více sezon { '2016/2017', '2017/2018', '2018/2019' }
- tabulka je tabulka s výslednými hodnotami počtu vsazených utkání n a průměrné návratnosti r , pro různé hodnoty L ve zvolené sezoně
- funkce rovněž vykreslí graf zachycující výsledné hodnoty počtu vsazených utkání n a průměrné návratnosti r pro různé hodnoty L ve zvolené sezoně

A.11 nahodny_sazejici.m

Funkce určená k vypsání kurzů podle optimalizovaného modelu 5 – s marží odpovídající sázkové kanceláři Bet365 – a následnému generování 10 000 „náhodných sázejících“ (simulací) podle pravidel popsaných v odstavci 6.1.3, kteří budou fiktivně sázet vždy u sázkové kanceláře Bet365 a zároveň budou identicky sázet u virtuální sázkové kanceláře s námi vypsány kurzy. Funkce poskytne srovnání průměrných zisků ve zvolených sezonách obou sázkových kanceláří, výběrových směrodatných odchylek a podílů úspěšností sázkové kanceláře v porovnání s druhou uvažovanou, pro různé hodnoty parametru alternativního rozdělení pro rozhodování o vsazení na utkání. Pro spuštění této funkce

je potřeba mít ve stejném adresáři soubor s daty `data_model5.mat`.

```
tabulka_zisk = nahodny_sazejici(sezony, seed)
```

- sezony jsou zvolené sezony – např. ve tvaru { '2018/2019' } nebo pro více sezon { '2016/2017', '2017/2018', '2018/2019' }
- seed je volitelný parametr „výchozího číslo“, od kterého jsou generovány náhodné hodnoty – pro námi uvažované simulace byl použit `seed = 222`
- `tabulka_zisk` je tabulka s výsledným srovnáním průměrných zisků ve zvolené sezoně, výběrových směrodatných odchylek a podílů úspěšností sázkové kanceláře v porovnání s druhou uvažovanou, pro různé hodnoty parametru alternativního rozdělení pro rozhodování o vsazení na utkání
- funkce rovněž vykreslí graf průměrných zisků a graf podílů úspěšností sázkových kanceláří v závislosti na parametru alternativního rozdělení pro rozhodování o vsazení na utkání