



VYSOKÉ UČENÍ TECHNICKÉ V BRNĚ

BRNO UNIVERSITY OF TECHNOLOGY

FAKULTA STROJNÍHO INŽENÝRSTVÍ

FACULTY OF MECHANICAL ENGINEERING

ÚSTAV MATEMATIKY

INSTITUTE OF MATHEMATICS

MODELOVÁNÍ PRAVDĚPODOBNOSTI SKÓROVÁNÍ VE SPORTU

MODELING OF SCORING PROBABILITY IN SPORT

BAKALÁŘSKÁ PRÁCE

BACHELOR'S THESIS

AUTOR PRÁCE

AUTHOR

Ondřej Hilscher

VEDOUCÍ PRÁCE

SUPERVISOR

Ing. Pavel Hrabec, Ph.D.

BRNO 2022

Zadání bakalářské práce

Ústav: Ústav matematiky
Student: Ondřej Hilscher
Studijní program: Matematické inženýrství
Studijní obor: bez specializace
Vedoucí práce: Ing. Pavel Hrabec, Ph.D.
Akademický rok: 2021/22

Ředitel ústavu Vám v souladu se zákonem 111/1998 o vysokých školách a se Studijním a zkušebním řádem VUT v Brně určuje následující téma bakalářské práce:

Modelování pravděpodobnosti skórování ve sportu

Stručná charakteristika problematiky úkolu:

Výsledné skóre v n kterých sportech není vždy objektivním měřítkem předvedeného výkonu. Zejména ve sportech, ve kterých je skórování "vzácnou" událostí (např. fotbal) mnohdy nereflktuje události na hřišti (velký vliv náhody). Pro objektivnější popis skutečného dění na hřišti lze použít některé statistické metody. Mezi příklady takových modelů lze zařadit např. tzv. "expected goals" ve fotbale i hokeji. Jedná se tedy v podstatě o převedení problému na modelování pravděpodobnosti skórování.

Cíle bakalářské práce:

Referenčně používaných modelů skórování ve zvoleném sportovním odvětví.
Seznámení s vhodnými nástroji matematické statistiky.
Analýza dat prostřednictvím vhodného softwarového nástroje (R, Python, ...).

Seznam doporučené literatury:

ANDRÁŠ, Jiří. Základy matematické statistiky. Vyd. 3. Praha: Matfyzpress, 2011. ISBN 978-80-7378-162-0.

AGRESTI, Alan. Categorical Data Analysis. 2nd ed. Hoboken: Wiley, 2002. Wiley series in probability and statistics. ISBN 04-713-6093-7.

Soccer analytics handbook (<https://github.com/devinpleuler/analytics-handbook>)

Termín odevzdání bakalářské práce je stanoven časovým plánem akademického roku 2021/22

V Brně, dne

L. S.

prof. RNDr. Josef Šlapal, CSc.
ředitel ústavu

doc. Ing. Jaroslav Katolický, Ph.D.
děkan fakulty

Abstrakt

Práce je zaměřena na modelování pravděpodobnosti skórování ve fotbale. V práci je popsán nutný matematický aparát potřebný k sestavení modelu logistické regrese a základní testy statistických hypotéz. Popsaný matematický aparát je poté aplikován na volně přístupná data z profesionálních fotbalových utkání. Výsledný model používá vysvětlující proměnné jako způsob zakončení, polohu na hřišti a zjednodušeně popsanou herní situaci.

Summary

This thesis aims for modelling of scoring probability in football. It describes necessary mathematical methods used in logistic regression model building and in basic statistical hypothesis tests. Afterwards the mathematical methods are used on available data from professional football matches. Resulting model uses shooting method, pitch location and simplified match situation as predictors.

Klíčová slova

logistická regrese, GLM, metoda maximální věrohodnosti, Expected goals

Keywords

logistic regression, GLM, maximum likelihood method, Expected goals

ONDŘEJ HILSCHER Modelování pravděpodobnosti skórování ve sportu. Brno: Vysoké učení technické v Brně, Fakulta strojního inženýrství, 2022. 23 s. Vedoucí diplomové práce Ing. Pavel Hrabec, Ph.D.

Prohlašuji, že jsem bakalářskou práci Modelování pravidel podobnosti skórování ve sportu napsal samostatně pod vedením Ing. Pavla Hrabce Ph.D. a Ing. Martina Roseckého, všechny použité materiály jsou uvedeny v seznamu literatury.

Ondřej Hilscher

Dìkuji v¹em, kteøí mne b¹hem vypracování podporovali a obohacovali mne svými zku¹enými a praktickými radami, zvlá¹tì Ing. Pavlu Hrabcovi Ph.D. a Ing. Martinu Ro-seckému, kteøí pomohli vytvoøit skvìlou atmosféru pro psaní práce, velké dìk patøí také mé rodinì, její¼ podpora mì provází celé studium.

Ondøej Hilscher

Obsah

1	Sport a statistika	2
2	Početné úvahy	3
2.1	Statistické modely v jednotlivých sportech	3
2.1.1	Expected goals method	4
3	Potřebný matematický aparát	6
3.1	Zobecněný lineární model	6
3.1.1	Binomické Logit Modely	6
3.2	Zobecněný lineární model pro binární data	7
3.2.1	Lineární Pravděpodobnostní Model	7
3.3	Metoda maximální věrohodnosti	7
3.3.1	Pomocná tvrzení	7
3.3.2	Věrohodnostní funkce a maximální věrohodné odhady	8
3.3.3	Věrohodnostní funkce veličiny s binomickým rozdělením	9
3.4	Logistická regrese	10
3.4.1	Interpretace parametrů	10
3.4.2	Fitování modelů logistické regrese	10
3.4.3	Testování podmodelu	11
4	Modelování pravděpodobnosti skórování	13
4.1	Triviální modely	13
4.2	Pokročilé modely	14
4.2.1	Základní modely logistické regrese	15
4.2.2	Pokročilé modely logistické regrese	17
5	Závěr	22

1. Sport a statistika

Pro mnoho lidí byly matematika a sport dlouhou dobu dva odlišné světy, a pro značnou část světové populace tomu tak stále je, ale aniž bychom si to uvědomovali, matematika, konkrétní statistika a statistické modely mohou mít a nikde jinde nemají zásadní vliv. Principy měření délky hodu, času, seřazení bodů i určení pořadí jsou tisíce let staré. Samozřejmě byly využívány jiné metody i jednotky, nicméně je vhodné zde uvést tvrzení, že matematika sport provází již od jeho raných kořenů. Dnešní pohled na matematiku, statistiku a její modely ve sportu se začal utvářet v druhé polovině minulého století, a už to byli nadšenci z řad fanoušků, trenéři, či samotní sportovci, začaly se uchovávat informace o turnajích, utkáních nebo závodech, bez kterých si dnes nedokážeme představit žádný televizní přenos. U každé sportovní události je uvedena grafika s mnoha statistickými údaji - nejdelší hod, počet střel, čas, úspěšných odpalů, držení míče, úspěšnost střelby, zákroků gólmána, největší vyvinutá rychlost během utkání, čas strávený na soupeřovi polovině, třetina a takto by šlo vyjmenovat spoustu dalších. Jistě si takovou tabulku každý, kdo viděl nějakou sportovní akci v televizi, vybaví.

Tato práce čtenáře uvede do prostřední datové analytiky sportu, ukáže, jak může práce se sportovními daty (zde využívány data [rmy Stats Bomb¹](#)) vypadat v praxi a představí jejich využití v oblasti sportu, a už z hlediska událostí uplynulých nebo těch budoucích. K tomu jsou využívány matematické resp. statistické nástroje, které jsou posléze aplikovány při vytváření predikčních modelů pravděpodobnosti skórování ve sportu. Modely jsou následně testovány, tak aby bylo zřejmé, který model by bylo vhodné použít, aby výsledky byly co nej přesnější a zda jsou jeho proměnné relevantní.

2. Poèáteèní úvahy

Data zmiòovaná výše (poèet støel, dr¾ení míèe...) uvádíme jako základní - popisná, kvantitativní, øeknou nám co a kolikrát se stalo, problém je, ¾e vzhledem k budoucnosti, pøítm zápasùm, je jejich vyu¾ití omezené, èasto zkreslují situaci na høíti, který tým byl lepš, kdo si zaslou¾il vyhrát. Kolikrát fanouci po zápasech diskutují ve stylù kdyby to tre l, vyhráli bychom nebo „sice nás pøehrávali, ale k ¾ádným vítím 1ancím jsme je nepustili.“ Tato práce se vinuje "pokroèilým" analytickým metodám, kromì kvantity vyu¾ívající i kvalitu, například kvalitu støelby, kvalitu samotných 1ancí ke skórování nebo bodového zisku. Další výhodou je možnost zpøesnìní jejich výsledkù pomocí parametrù vyskytujících se ve høe. Tyto pokroèilé, kvalitativní statistiky jsou èasto schovány za urèitou variací modelu "Oèekávané hodnoty"- Expected value model, který je modifikován speci ky pro každé sportovní odvíví. Zde je pozornost vinována pøevážnì úvazetody Oèekávaných gólù { The Expected goals method, na základì které se nehodnotí èistí výsledky utkání a kvantitativní statistiky uèinkujících, ale výkony mu¾stev, jednotlivcù, díky nim¾ lze predikovat jejich pùsobenì bhem dalších sportovních klání.

Expected value mù¾e být pro skauty, trenéry, hráèe a vlastnì i pro fanouky velmi u¾iteèná, její popularizace bhem posledních let a hlavní popularizace z ní odvozených modelù pomohla leckterému laikovi nahlédnout pod poklièku jednotlivým sportùm a ukázala spojitost mezi øíí èisel a pohybu. Zároveò týmy, které tyto modely dlouhodobì vyu¾ívají prokazatelnì dosahují, pro ni døív stí¾í dosa¾itelných, posunù v podobì vylepšeného skautingu hráèù, který mnohdy pøináší vyší ekonomické zisky, výkonostních posunù a s nimi souvisejícími pøípadnými postupy do vyších soutí¾í.

2.1. Statistické modely v jednotlivých sportech

Každý sport je speci ky svou krásou, pravidly, nicménì mezi jejich statistickými modely lze najít èasto nijaké podobnosti, samozøejmì, èím si jsou sporty podobnìjší, tím podobnìjší si budou i metody modelování pravdìpodobnosti skórování, na dalších pár øádcích je krátce pøedstaveno, jak jsou predikovány události průbìhu utkání.

V Baseballu jsou èasto zmiòovány "Oèekávané dobihy"- Run expectancy (RE). Výsledky RE ukazují pravdìpodobnost poètu dobihù (získaných bodù týmu) do konce smìny v závislosti zejména na poètu útoèících hráèù na metách a poètu outù, pøípadnì na historické úspìnosti nadhazovaèe nebo pálkáøe.

Baseball jako takový je považován za jeden z průkopnických sportù, co se analytických modelù týèe, první velké úspìchy jsou nyní již 20 let staré, mù¾eme je vidìt zdokumentované například ve snímku "Moneyball"(re¾ie Bennett Miller, 2011).

Mezi nejznámìjší modely vyu¾ívaných v americkém fotbale patøí "Expected points added"(EPA), vychází z aktuální pozice míèe na høíti, ukazují kolik bodù obvykle tým z dané akce získá pøi uskuteènìní urèitého jevu, jako je například zisk území díky pøebìhnutí urèité vzdálenosti s míèem nebo zatlaèení soupeøe pressingem zpìt ve høíti.

Hokejoví analytici pou¾ívají pøevážnì nám již známé "Expected Goals", zjišují pravdìpodobnost, ¾e støela èi 1ance skonèí v síti soupeøe, zohlednùjí pøitom historickou bilanci støel z dané pozice, mno¾ství hráèù v dráze støely, herní situaci, jedná-li se o nájezd, oslabení, pøesilovku 5 na 4, 5 na 3, hru v plném poètu hráèù a další parametry. Problém je, kdy¾ se

2.1. STATISTICKÉ MODELY V JEDNOTLIVÝCH SPORTECH

na základě této statistiky snaží srovnávat jednotlivé hráče, jelikož jejich časy strávené na ledě jsou různé, proto se zavádí "Expected goals per 60" - Očekávané góly na 60 minut, kde se nesrovnalosti s časem na ledě vyrovnají a je poté snáze porovnat výkony jednotlivců.

2.1.1. Expected goals method

„Metoda Očekávaných gólů {The Expected goals method (dále xG) je prosazována i ve fotbalové bublině, je to také dáno tím, že branek v hokeji nebo fotbale nepadá tolik jako třeba v házené, samotná událost vstřelení branky je v obou těchto sportech naprosto raritní věc, která má zcela zásadní vliv na celou hru. V jedné z nejpůvodnějších fotbalových lig světa, anglické Premier League v sezóně 2020/2021 bylo průměrně vstřeleno 2.69 gólů za 90minutové fotbalové utkání. Divák tak průměrně mezi dvěma góly čekal déle než půlhodinu. Právě fotbalu a metodě "Expected goals" je v textu věnována největší pozornost.

Základní princip této metody je přitom jednoduchý, stělem, kvantitativní informaci, je přizvána hodnota xG, kvalitativní informace. Vezmeme-li všechny střely jako počet pokusů a góly jako úspěšné pokusy tak díky informaci, že za sezóny 11/12 a 15/16 bylo ve skotské lize a čtyřech nejvyšších anglických fotbalových soutěžích v průměru potřeba 8,4427 střel na gól, bude xG každé střely rovno $\frac{1}{8,4427} = 0,1184$.

Nyní vezmeme pevně danou pozici na hrací ploše, ze které hráč vystřelil, v databázích statistik sportovních utkání lze nalézt tisíce střel přesně ze stejného místa na hřišti a hodnota xG (pravděpodobnost že padne branka) bude rovna:

$$xG = \frac{\text{celkový počet gólů z dané pozice}}{\text{celkový počet střel z dané pozice}} \quad (2.1)$$

Takto můžeme hodnotit každou fázi v zápase a zjistit, jaké xG oba týmy nasbírali { kolik gólů by průměrně padlo ze fází vytvořených během utkání.

Samotné modelování tak jednoduché není, jelikož je hledána pravděpodobnost uskutečnění určitého jevu, padne-li, či nepadne branka. Je možné využít machine learning modely, díky jejich přesnosti častěji využívané v běžné praxi, nebo také klasické modely jako je model logistické regrese, která oproti přesnějším modelům poskytuje lepší interpretovatelnost a proto je zde převážně věnována pozornost jí. Chceme využít co nejvíce informací, které máme k dispozici a ověřit, zda jejich zahrnutí do modelu je relevantní pro získané výsledky, aby hodnoty koeficientů regresních parametrů byly co nejpresnější. Mezi parametry, se kterými pracujeme, patří kromě pozice na hrací ploše také pod jakým úhlem střelec vidí branku, zakonejšuje-li slabší, silnější nohou, kolik hráč stojí v potenciální dráze střely, jedná-li se o volej, hlavičku a další. Z principu plyne, že pravděpodobnost vstřelení gólu do odkryté brány je daleko vyšší, než když hráč vystřelí z 25 metrů přes chumel protihráče. Co xG neuvádí, je informace, kdo danou střelu vykonal, xG říká kolikrát ze 100 pokusů by se za daných podmínek trel běžný profesionální fotbalista.

Ukážeme si, jak může xG vypadat v praxi. Vezmeme zápas Sparta { Slavia s výsledkem 1:1, kde Sparta přestřelí Slavii 9:5 a bude mít výrazně vyšší procento držení míče na svých kopačkách. Na první pohled by se zdálo, že letenčí byly lepší tým a remíza je

2. POĚÁTEĚNÍ ÚVAHY

pro ni smolný výsledek. Støely spar»anských hráèù ov'em byly z vit'ích vzdáleností od branky, a ne v¾dy z optimálních pozic, průmìrnì mìla spar»anská støela $xG = 0,1$ (1 z 10 takovýchto støel rozvlní sí» soupeøe), Slavie si z protiútokù vytvoøila 5 støel s průmìrným $xG=0,3$, tedy o něco kvalitnìjí 1ance. Sparta tedy za zápas získala $9 \cdot 0,1 = 0,9$ xG, zatímco se1ívaní $5 \cdot 0,3 = 1,5$ xG, vidíme, že Slavie byla v utkání nebezpeènìjí a mìla k výhøe blíže než její soupeø.

Pomocí xG lze sledovat jednotlivé hráèe, mùžeme mít hráèe kteøí dávají hromadu branek, ov'em jejich xG je daleko nižší, z toho usuzujeme, že hráè má skvilou formu nebo prostì jen 1tístí, ze statistického hlediska mu takové výkony, kdy pøekraèuje své xG, dlouhodobì nevydrží, jelikož bìhem delšího èasového horizontu své xG pøekonává pouze jeden hráè na svìtì { Lionel Messi^[4]. Naopak mùžeme najít útoèníky, kteøí dávají výraznì ménì branek, než je jejich xG, to pak mùžeme kluby odradit od jejich angažování, pøípadnì je donutit zapracovat na zakonèování hráèù. Dále bývá xG využíváno pro skautování skrytých talentù, èi pro pøedpovìi dalších zápasù, což je jeho veliká výhoda oproti èistì popisným statistikám.

3. Potøebný matematický aparát

V této kapitole jsou uvedeny matematické nástroje využívané posléze pro vytvoření modelu. Vychází ze zdrojů [1] [2] [3].

3.1. Zobecněný lineární model

Základní lineární regresní modely pracují se závislými proměnnými s normálním rozdělením, aby bylo možné modelovat závislé proměnné, které mají rozdělení jiné než normální, zavádíme Zobecněné lineární modely (angl. Generalized linear models - GLMs). Zobecněné lineární modely obsahují tři části, náhodnou, systematickou a link function .

Náhodná část Zobecněného lineárního modelu obsahuje závislou proměnnou y_i a nezávislé pozorováními (y_1, \dots, y_N) , s hustotou pravděpodobnosti:

$$f(y_i; \eta_i) = a(\eta_i) b(y_i) \exp[y_i Q(\eta_i)] \tag{3.1}$$

Hodnota parametru η_i je funkcí vysvětlujících veličin pro každé $i = 1, \dots, N$. Parametr Q_i označujeme jako přirozený parametr

V systematické části je vysvětlující proměnné pomocí lineárního modelu přirozen vektor $(\eta_1; \dots; \eta_N)$. Nechť x_{ij} označuje hodnotu vysvětlující proměnné j ($j = 0; 1; 2; \dots$) pro subjekt i . Pak

$$\eta_i = \sum_j x_{ij}; \quad i = 1; \dots; N: \tag{3.2}$$

Tato lineární kombinace vysvětlujících proměnných se nazývá lineární prediktor. Pro koeficient konstantního členu lineárního modelu η_0 mimo výjimky pro všechna i platí $x_{i0} = 1$. Aby byl model kompletní je nutné první dvě části propojit. Nechť $\eta_i = E(Y_i); i = 1; \dots; N$. Zavádíme monotónní diferencovatelnou funkci g , platí $\eta_i = g(\eta_i)$, kde g nazýváme link function. Díky g existuje vztah mezi η_i a vysvětlujícími proměnnými

$$g(\eta_i) = \sum_j x_{ij}; \quad i = 1; \dots; N: \tag{3.3}$$

3.1.1. Binomické Logit Modely

Uvažujme, že chceme modelovat, zda určitá událost nastane, takovéto případy představují binární závislé proměnné. Úspěch, uskutečnění jevu a neúspěch, jeho neuskutečnění reprezentujeme jako 1 a 0. Při binomickém rozdělení pro závislou proměnnou Y platí: $P(Y = 1) = p$, $P(Y = 0) = 1 - p$ a $E(Y) = p$. Jedná se o speciální případ binomického rozdělení s $n=1$ s hustotu pravděpodobnosti

$$f(y; p) = p^y (1 - p)^{1 - y} = (1 - p) \frac{p^y}{1 - p + p^y} = (1 - p) \exp\left\{y \log \frac{p}{1 - p}\right\} \tag{3.4}$$

pro $y = 0$ nebo 1. Je zřejmé, že se jedná o Zobecněný lineární model s hustou pravděpodobnosti (3.1). Parametru η přísluší $a(\eta) = 1 - p$, $b(y) = p^y$ a $Q(\eta) = \log\left[\frac{p}{1 - p}\right]$,

3. POTØEBNÝ MATEMATICKÝ APARÁT

pøirozený parametr $\log[\pi/(1-\pi)]$ vyjadøuje log odds ("zlogaritmovaná šance") šance, že výstup modelu bude 1, dále v textu se setkáváme s označením logit funkce. Modely využívající logit jako link function jsou popsány v kapitole 3.4, vlnované logistické regresi a jejím modelům, nebolilogit modelům.

3.2. Zobecnìný lineární model pro binární data

Nechť Y je binární závislá promìnná, udávající "uskuteènìní" nebo "neuskuteènìní" události, každé pozorování ukáže jeden z těchto výsledkù, oznaèené 1 a 0. Støední hodnota Y je $E(Y) = P(Y = 1)$. Aby byla zøejmá závislost Y na hodnotách $\mathbf{x} = (x_1; \dots; x_p)$ oznaèíme $P(Y = 1)$ jako $\pi(\mathbf{x})$. Rozptyl Y je

$$\text{var}(Y) = \pi(\mathbf{x})[1 - \pi(\mathbf{x})] \quad (3.5)$$

Jedná se o rozptyl binomického rozdílení $p=1$:

3.2.1. Lineární Pravdìpodobnostní Model

Regresní model binární závislé promìnné

$$\pi(\mathbf{x}) = \frac{1}{1 + \exp(-\beta_0 - \beta_1 x_1 - \dots - \beta_p x_p)} \quad (3.6)$$

je nazýván lineární pravdìpodobnostní model. Lineární funkce pokrývají svými hodnotami celou reálnou osu, ovšem pravdìpodobnost nabývá hodnot v mezích pouze od 0 do 1. V modelu (3.6) mùže pro některé hodnoty \mathbf{x} nastat $\pi(\mathbf{x}) < 0$ nebo $\pi(\mathbf{x}) > 1$, ovšem model mùže platit na omezené oblasti, pokud taková situace opravdu nastane je využívána snadná interpretace parametrù β_j , které vyjadøují zminu $\pi(\mathbf{x})$ při jednotkové zminì x_j .

3.3. Metoda maximální vìrohodnosti

3.3.1. Pomocná tvrzení

Definice 3.1 Nechť f je reálná funkce definovaná na otevøeném intervalu $I = (a,b)$, kde $1 < a < b < \infty$. Funkce f se nazývá konvexní, platí-li

$$f[\lambda x + (1-\lambda)y] \leq \lambda f(x) + (1-\lambda)f(y) \quad (3.7)$$

pro všechna x, y , taková, že $a < x < y < b$ a $0 < \lambda < 1$.

Funkce f se nazývá striktnì konvexní, platí-li pro všechna uvedená x, y ; v (3.7) ostrá nerovnost.

Funkce f se nazývá (striktnì) konkávní, je-li $-f$ (striktnì) konvexní.

Vìta 3.2 Nechť f je definována na I a nechť je tam spojitá. Nechť existuje na I spojitá f'' a nechť $f'' \geq 0$ existuje a je koneèná na I . Pak funkce je konvexní právi tehdy, platí-li $f'' \geq 0$ pro všechna $x \in I$:

Dùkaz. Viz Fichtengolc I (1958), odst. 143, str. 299, víta 2.

3.3. METODA MAXIMÁLNÍ VÍROHODNOSTI

Věta 3.3 (Jensenova nerovnost) Nechť f je konvexní funkce definovaná na I . Nechť X je náhodná veličina s konečnou střední hodnotou taková, že $P(X \in I) = 1$. Pak platí

$$f(E X) \leq E f(X); \quad (3.8)$$

Je-li f striktně konvexní, pak nerovnost (3.8) je ostrá s výjimkou případu, kdy veličina X je rovna konstantě s pravděpodobností 1. Nerovnost (3.8) se nazývá Jensenova.

Důkaz. Viz Lehmann (1991), str. 50, věta 6.3, nebo Rao (1978), kap. 1e.5.

3.3.2. Vírohodnostní funkce a maximálně vírohodné odhady

Uvažujme náhodný vektor $X = (X_1, \dots, X_n)^0$ se sdruženou hustotou $u(x; \theta)$, kde $\theta \in \Omega$, kde Ω je parametrický prostor parametrů: Nechť x má pevně danou hodnotu, pak se funkce $f(x; \theta)$ jakožto funkce nazývá vírohodnostní funkce. Hodnota $\hat{\theta}$ parametru θ , která maximalizujeme pro dané $X = x$ vírohodnostní funkci $p(x; \theta)$, se nazývá maximálně vírohodný odhad parametru θ :

Nechť X je náhodný vektor se sdruženou hustotou $u(x; \theta)$, kde $\theta \in \Omega \subset \mathbb{R}^m$. Uvažujme funkci $u: \Omega \rightarrow \mathbb{R}_+$, která zobrazuje na \mathbb{R}_k . Předpisem $\theta = u(\theta)$ každému $\theta \in \Omega$ přiřadíme $\theta \in \mathbb{R}^k$. Nechť $G(\theta) = \{x: u(x; \theta) > 0\}$; $u(\theta) = \int_{G(\theta)} u(x; \theta) dx$. Označme

$$M(x; \theta) = \sup_{\theta \in G(x)} f(x; \theta)$$

M jakožto funkci nazýváme vírohodnostní funkcí indukovanou parametrickou funkcí u . Hodnotu $\hat{\theta}$, která maximalizuje $M(X; \theta)$ označujeme jakožto maximálně vírohodný odhad parametrické funkce u .

Nechť θ je jednorozměrný parametr a platí následující předpoklady.

P_1 : Nechť Ω je parametrický prostor, který obsahuje takový neprázdný otevřený interval I , kde skutečná hodnota parametru θ patří do I .

P_2 : Nechť $X = (X_1, \dots, X_2)^0$, kde X_i jsou stejně rozdělené, nezávislé veličiny s hustotou $f(x; \theta)$ vzhledem k nějaké -konečné míře

P_3 : Nechť $M = x: f(x; \theta) > 0$ nezávisí na θ .

P_4 : Nechť $\theta_1, \theta_2 \in \Omega$. Pak $f(x; \theta_1) = f(x; \theta_2)$ skoro všude (vzhledem k míře) právě tehdy, je-li $\theta_1 = \theta_2$.

Tedy vzhledem k míře μ je sdružená hustota náhodného vektoru $f(x; \theta) = f(x_1; \theta) \dots f(x_n; \theta)$.

Věta 3.4 Jestliže $\theta \in \Omega$, pak pro každé takové pevné $\theta_0 \in \Omega$, platí

$$P_{\theta_0} f(X; \theta_0) > f(X; \theta) \text{! } 1 \quad (3.9)$$

3. POTØEBNÝ MATEMATICKÝ APARÁT

Dùkaz. Viz [3] kpt. 7.6, str. 149.

Pokraèujeme v úvaze, θ je jednorozmìrný parametr. Uvažujeme funkci promìnnéf $(x; \theta)$ pro pevné x . Funkce

$$L(x; \theta) = \ln f(x; \theta) \quad (3.10)$$

se nazývá logaritmická vùrohodnostní funkce. Èasto $L(x; \theta)$ znaèíme jen jako $L(\theta)$.

Vìta 3.5 Nech θ_0 jsou splnìny pøedpoklady P_4 . Nech θ_n na intervalu I existuje $f'(x; \theta) = \frac{\partial f(x; \theta)}{\partial \theta}$ pro skoro všechna x . Pak pro každé $\epsilon > 0$ pøi $n \rightarrow \infty$ platí, že s pravdìpodobností konvergující k jedné má vùrohodnostní rovnice

$$\frac{\partial L(X; \theta_n)}{\partial \theta} = 0 \quad (3.11)$$

takový koøen $\hat{\theta}_n = \hat{\theta}_n(X)$, $\theta_0 - \epsilon < \hat{\theta}_n < \theta_0 + \epsilon$, kde θ_0 je skuteèná hodnota parametru, v níž $f(x; \theta)$ nabývá svého maxima.

Dùkaz. Nech $\epsilon > 0$ je tak malé, $\theta_0 - \epsilon \in I$; $\theta_0 + \epsilon \in I$. De nujme

$$S_n = \{x : L(x; \theta_0 - \epsilon) > L(x; \theta_0) \text{ a } L(x; \theta_0) > L(x; \theta_0 + \epsilon)\} \quad (3.12)$$

Dle (3.9) platí $P_0(X \in S_n) \rightarrow 0$. Pro každé $x \notin S_n$ tedy s pravdìpodobností blížící se jedné existuje $\hat{\theta}_n$ takové, $\theta_0 - \epsilon < \hat{\theta}_n < \theta_0 + \epsilon$ a θ funkce $L(X; \theta)$ má lokální maximum v bodì $\theta = \hat{\theta}_n$. Pak $L'(X; \hat{\theta}_n) = 0$:

3.3.3. Vùrohodnostní funkce velièiny s binomickým rozdílením

Uvažujeme náhodnou velièinu $Bi(n; p)$; $0 < p < 1$, pro kterou platí

$$f(x) = P_p(X = x) = \binom{n}{x} p^x (1-p)^{n-x} \quad (3.13)$$

Hledejme maximální vùrohodný odhad \hat{p} , tedy pro $X=x$ jde o maximalizaci funkce $g(p) = \binom{n}{x} p^x (1-p)^{n-x}$. Nech $x \in 0; x \in n$. Pak

$$g'(p) = [x(1-p)^{-(x-1)} - (n-x)p^{x-1}(1-p)^{-x}] p^x (1-p)^{n-x-1} \quad (3.14)$$

nulovým bodem funkce $g'(p)$ je $\hat{p} = \frac{x}{n}$. Nyní si ukážeme, že v tomto bodì skuteènì funkce $g(p)$ nabývá svého maxima. De nujme

$$Z_i = \begin{cases} \frac{n}{x}; & i = 1; \dots; x; \\ \frac{n(1-p)}{n-x}; & i = x+1; \dots; n; \end{cases} \quad (3.15)$$

Aritmetický prùmìr tìchto èísel je

$$z = \frac{1}{n} \left[x \frac{n}{x} + (n-x) \frac{n(1-p)}{n-x} \right] = 1; \quad (3.16)$$

3.4. LOGISTICKÁ REGRESE

geometrický průměr je

$$z_G = \left(\frac{n}{x} \cdot \frac{n(1-x)}{n-x} \right)^{\frac{1}{n}} \quad (3.17)$$

z nerovnosti $z_G \geq z$ zjistíme, že

$$\left(\frac{n}{x} \cdot \frac{n(1-x)}{n-x} \right)^{\frac{1}{n}} \geq \frac{x}{n} \cdot \frac{n-x}{n} \quad (3.18)$$

zjevně při $x = \frac{1}{2}$ nastává rovnost a tedy z_G maximalizuje funkci $g(x)$:

Nastane-li jeden z případů $x = 0$ nebo $x = n$, pak maximálně věrohodný odhad parametru na intervalu $(0; 1)$ neexistuje, nastal by spor s předpokladem P_1 .

3.4. Logistická regrese

Pro práci s binárními závislými proměnnými, je vhodný nástroj logistická regrese, která je široce využívána napříč obory, ať už se jedná o medicínu, ekonomii nebo právě sport.

3.4.1. Interpretace parametrů

Předpokládejme binární závislou proměnnou Y a proměnnou X , na které je Y závislá, nechť $\pi(x) = P(Y = 1 | X = x) = \frac{1}{1 + \exp(-\eta(x))}$. Model Logistické regrese je

$$\pi(x) = \frac{\exp(\eta(x))}{1 + \exp(\eta(x))} \quad (3.19)$$

Je tedy zjevné, že logit je lineární

$$\text{logit}[\pi(x)] = \log \frac{\pi(x)}{1 - \pi(x)} = \eta(x) \quad (3.20)$$

Interpretace parametru

Nyní se zaměříme na $\eta(x)$ z rovnice (3.20), znaménko tohoto parametru určuje, zda $\pi(x)$ bude rostoucí či klesající s rostoucí x . Sklon logistické křivky roste s nárůstem j , pro $\eta(x) \rightarrow 0$ se logistická křivka přibližuje k horizontální rovné přímce. V případě nezávislosti Y na X platí $\eta(x) = 0$.

Aplikací exponenciály na obě strany rovnice (3.20) ukazujeme, že "odds" (odds) jsou exponenciální funkcí, odds se zvětší e -krát s každým jednotkovým navýšením proměnné x , e je poměr odds při $X = x + 1$ a odds při $X = x$. Vzhledem k nelineární závislosti $\pi(x)$ na x je změna $\pi(x)$ pro rozdílná x v logistické regresi "nerovnoměrná".

3.4.2. Fitování modelu logistické regrese

Uvažujme binárních nezávislých pozorování, nechť $\mathbf{x}_i = (x_{i1}; \dots; x_{ip})$ udává nastavení hodnot p vyvítlujících proměnných, $i=1, \dots, N$. Nabývají-li všechna pozorování různých hodnot, pak $N=n$. Model logistické regrese (3.19), uvažující regresní parametry jako konstantní, je

$$\pi(\mathbf{x}_i) = \frac{\exp\left(\sum_{j=1}^p \beta_j x_{ij}\right)}{1 + \exp\left(\sum_{j=1}^p \beta_j x_{ij}\right)} \quad (3.21)$$

3. POTŘEBNÝ MATEMATICKÝ APARÁT

Vírohodnostní rovnice

Když více než jedno pozorování nastane při pevné hodnotě, je vhodné zaznamenat počet pozorování y_i a počet úspěšných pokusů. Poté uvažujeme jako to počet úspěšných pokusů, namísto odezvy jednotlivých binárních proměnných. Pak Y_1, \dots, Y_N jsou nezávislé binomické proměnné se střední hodnotou $E(Y_i) = n_i(x_i)$, kde $n_1 + \dots + n_n = n$. Jejich marginální hustota pravděpodobnosti je úmírná součinu N binomických funkcí

$$\begin{aligned} \Psi &= \prod_{i=1}^N (x_i)^{y_i} [1 - (x_i)]^{n_i - y_i} = \\ &= \prod_{i=1}^N \left(\frac{x_i}{1 - (x_i)} \right)^{y_i} (1 - (x_i))^{n_i} \\ &= \exp \left(\sum_{i=1}^N y_i \log \frac{x_i}{1 - (x_i)} - \sum_{i=1}^N n_i \log [1 - (x_i)] \right) \end{aligned} \quad (3.22)$$

Pro model (3.21), je-li β_j logit, tak je exponenciální člen v posledním výrazu je roven $\exp(\beta_j x_{ij}) = \exp(\beta_j x_{ij}) / [1 + \exp(\beta_j x_{ij})]$, tak vírohodnostní funkce je rovna

$$L(\beta) = \prod_{j=1}^p \prod_{i=1}^N \frac{\exp(\beta_j x_{ij})^{y_{ij}}}{1 + \exp(\beta_j x_{ij})} n_i \log [1 + \exp(\beta_j x_{ij})]^{-n_i} \quad (3.23)$$

Vírohodnostní rovnice získáme, když položíme $\frac{\partial L(\beta)}{\partial \beta_j} = 0$: Platí

$$\frac{\partial L(\beta)}{\partial \beta_j} = \sum_{i=1}^N y_{ij} x_{ij} - \sum_{i=1}^N n_i x_{ij} \frac{\exp(\beta_j x_{ij})}{1 + \exp(\beta_j x_{ij})}; \quad (3.24)$$

vírohodnostní rovnice jsou

$$\sum_{i=1}^N y_{ij} x_{ij} - \sum_{i=1}^N n_i \hat{x}_{ij} = 0; \quad j = 1, \dots, p; \quad (3.25)$$

kde $\hat{x}_i = \exp(\beta_k \hat{x}_{ik}) / [1 + \exp(\beta_k \hat{x}_{ik})]$ je maximální vírohodný odhad (x_i) : Tyto rovnice jsou nelineární a vyžadují numerické řešení. Při modelování je v programovacím jazyce Python během vytváření modelu využívána metoda **IRLS** - Iteratively reweighted least squares^[10].

3.4.3. Testování podmodelu

K testování podmodelu je využíván test poměrem vírohodnosti (Likelyhood ratio test) který využívá logaritmickou vírohodnostní funkci, která pro binomické závislé proměnné vypadá následovně:

$$\begin{aligned} L(\beta) &= \log \prod_{i=1}^N \frac{Y_i (1 - \hat{x}_i)^{1 - Y_i}}{\hat{x}_i^{Y_i} (1 - \hat{x}_i)^{1 - Y_i}} \\ &= \sum_{i=1}^N (Y_i \log \hat{x}_i + (1 - Y_i) \log(1 - \hat{x}_i)) \\ &= \sum_{i=1}^N Y_i \log \frac{\hat{x}_i}{1 - \hat{x}_i} + \sum_{i=1}^N \log(1 - \hat{x}_i) \end{aligned} \quad (3.26)$$

3.4. LOGISTICKÁ REGRESE

Pozorované náhodné veličiny se v logistické vïrohodnostní funkci vyskytují v souèinech s výrazy $\log\left(\frac{y_i}{1-y_i}\right)$.

Uvažujme nejbohatší možný model, model s vïtší hodnotou vïrohodnostní funkce nelze vytvořit, takový model se nazývá saturaovaný. Saturaovaný model má právě tolik parametrů, kolik je různých hodnot vektorů x_i : Maximální hodnotu vïrohodnostní funkce v saturaovaném modelu označíme L_{\max} . Každý další model je podmodelem saturaovaného modelu. Pomocí deviance posoudíme pøíležitost bïžného modelu

$$D(b) = 2(L_{\max} - L(b)). \quad (3.27)$$

Èím je model pøíležitivší, hodnota deviance D klesá. Dále pøedpokládáme, že všechny vektory x_i jsou různé, pak má saturaovaný model parametrů $\beta_1; \dots; \beta_k$. Odhadem střední hodnoty y_i je \hat{y}_i , dle (3.26) platí

$$L_{\max} = \sum_{i=1}^n (y_i \log y_i + (1 - y_i) \log(1 - y_i)) = 0 \quad (3.28)$$

Oznaème odhad pravdïpodobnosti jedničky $\hat{y}_i = \hat{\pi}_i(x_i)$, dosazením do (3.28) vyjádøíme devianci v modelu logistické regrese jako

$$D(b) = -2L(b) = -2 \sum_{i=1}^n (y_i \log \hat{\pi}_i + (1 - y_i) \log(1 - \hat{\pi}_i)) \quad (3.29)$$

Test pomìrem vïrohodnosti

Uvažujme model M_1 s odhadem parametrů \hat{b} a jeho podmodel M_2 , který vznikl odebráním èásti regresorů model M_1 , s odhadem parametrů b . Při testu pomìrem vïrohodnosti testujeme, zda všechny parametry obsažené v modelu a zároveň vynechané v podmodelu M_2 jsou rovny nule, porovnáváme hodnoty logistické vïrohodnostní funkce pro \hat{b} a b pomocí statistiky

$$LR = 2(L(\hat{b}) - L(b)) \quad (3.30)$$

LR statistiku lze také vyjádřit pomocí deviance modelu a podmodelu

$$2(L(\hat{b}) - L(b)) = 2(L_{\max} - L(b)) - 2(L_{\max} - L(\hat{b})) = D(b) - D(\hat{b}) \quad (3.31)$$

Platí-li testovaný podmodel a jsou-li zároveň splněny podmínky regularity^[2], pak má statistika LR asymptoticky rozdělení χ^2_q , kde q je rozdíl počtu nezávislých parametrů v modelech, které porovnáváme.

4. Modelování pravděpodobnosti skórování

V této kapitole jsou uvedeny jednotlivé modely pro modelování pravděpodobnosti skórování ve fotbale, modelování Expected goals - xG. Data byla získána od rmy Statsbomb, která je volně distribuuje, jejich následné zpracování proběhlo v programovacím jazyce Python, s využitím knihoven jak pandas, numpy, statsmodels nebo FCPython. Zpracováno bylo 568 zápasů, ve kterých dohromady padlo 14775 strel, branka padla v 1753 případech.

Postupně budou uváděny modely od nejjednoduššího po modely komplexnější, které pak lze otestovat proti jednodušším modelům, a zjistit tak relevanci jejich parametrů.

4.1. Triviální modely

Chceme-li znát pravděpodobnost, že ze strel padne branka, tak se samozřejmě jako nejjednodušší úvaha nabízí poměr

$$xG = \frac{\text{celkový počet gólů}}{\text{celkový počet strel}}$$

Výsledkem tohoto triviálního modelu, je každá strela $xG = 0;1184$, jelikož z principu víme, že pravděpodobnost padnutí branky ze 2 metrů je větší, než, že hráč rozvlí sí z poloviny hřiště, je zjevné, že daný model lze sice využít, spoléhat na jeho jakoukoliv přesnost je ovšem na pováženou.

Proto dále uvažujeme opět poměr $\frac{\text{celkový počet gólů}}{\text{celkový počet strel}}$, nicméně již jej počítáme pro každou "čtverec" metr na metr na hrací ploše. Celou hrací plochu, kterou uvažujeme dělíme na čtverce 100 a 70 metrů, jsme rozdělili na síť těchto metrových čtverců, přičemž jejich hranice jsou celočíselné. Takže přesněji počítáme $\frac{\text{celkový počet gólů z jednoho čtverce na hřišti}}{\text{celkový počet strel z toho samého čtverce na hřišti}}$. Tento model je vyobrazen na 4.3. Pro snazší interpretaci byly na obrázcích 4.1 respektive 4.2, vyobrazeny všechny strel a góly z použité databáze.

Obrázek 4.1:
lokace jednotlivých strel

4.2. POKROÈILÉ MODELY

Obrázek 4.2:
lokace v'ech gólù

Na obrázku 4.3 ji¼ vidíme vykreslené jednoduché xG, tmavá oblast okolo brány potvrzuje, že nejvyšší pravdìpodobnost vstøelení branky je, kdy¼ jsme k ní co nejbli¼e, tmavší body u stran høi'ti jsou zpùsobené malým poètem støel z tìchto pozic v databázi. Vzhledem k nerovnomìrnému poètu støel z jednotlivých èástí høi'ti, je ovšem tento model znaènì nepøesný.

Obrázek 4.3: jednoduché xG

4.2. Pokroèilé modely

Nyní ji¼ zaèneme vyu¼ívat matematické nástroje popsané v kapitole 3. Binární data na výstupu modelu, padne/nepadne branka, modelujeme pomocí logistické regrese, kdy budeme postupnì vytváøet modely s rùznými promìnnými, jejich¼ poèet budeme postupnì navyšovat a vytváøet tak modely komplexnijší.

4. MODELOVÁNÍ PRAVDĚPODOBNOСТИ SKÓROVÁNÍ

4.2.1. Základní modely logistické regrese

xG v závislosti na vzdálenosti

Začneme modelem jedné proměnné X , bereme za ni vzdálenost od branky, nikoliv ovšem od jejího středu, ale kolmou vzdálenost od brankové čáry, takto uvažovanou vzdálenost značíme X i v dalších modelech. Program z dat získal následující hodnoty koeficientů a jejich intervalový odhad při hladině významnosti $\alpha = 0,05$.

Tabulka 4.1: xG v závislosti na X

Proměnná	koeficient	[0,025	0,975]
konstanta	0,7193	0,615	0,824
X	0,0924	0,085	0,100

Deviance modelu je 10057

Takže funkce xG jako taková vypadala následovně

$$xG = \frac{1}{1 + e^{0,7193 + 0,0924x}}; \quad (4.1)$$

můžeme ji vidět vykreslenou na následujícím obrázku.

Obrázek 4.4: xG v závislosti na vzdálenosti

Z obrázku 4.4 je zřetelně pozorovatelný nelineární vztah pravděpodobnosti padnutí branky a střelcoví vzdálenosti od brankové čáry, což ovšem model neuvazuje je, kde přesně se hráč nachází, může být ve stejné vzdálenosti od brankové čáry, nicméně může stát nikde u autové čáry nebo nikde uprostřed, oběma těmito případy model stejné xG ,

4.2. POKROÈILÉ MODELY

proto kromì kolmé vzdálenosti X , pøidáme do modelu úhel, pod kterým hráè brankovou konstrukci vidí, viz 4.5.

Obrázek 4.5: Úhel, pod kterým støelec vidí branku

xG v závislosti na vzdálenosti a úhlu

Nyní ji¾ v modelu krom X uva¾ujeme i úhel popsáný výše. Problematika stejných hodnot xG pro situace, kdy¾ je hráè v rohu høíti a kdy¾ je tisni pøed brankou, tímto zaniká (jejich vzdálenost X od brankové èáry je v takovém pøípadi stejná). Pøi hladinì hladinì významnosti $\alpha = 0:05$ byly koeficienty promìnných a jejich intervalové odhady následovni.

Tabulka 4.2: xG v závislosti na X a θ

Promìnná	koeficient	[0,025	0,975]
konstanta	2,0586	1,857	2,260
X	0,0503	0,041	0,059
	-1,7412	-1,957	-1,525

Deviance modelu je 9778

Znaménko koeficientu jednoduše interpretuje, zdali se s rostoucí promìnnou hodnota xG zvyšuje, resp. sni¾uje. Je-li koeficient promìnné kladný, tak s rostoucí promìnnou xG klesá, lze vidít u vzdálenosti X , èím dál je hráè od brankové èáry, tím bude menší pravdìpodobnost, ¾e skóruje. Z logiky vici plyne, bude-li znaménko koeficientu promìnné mínus, bude s rostoucí promìnnou rùst zároveň hodnota xG . Tento stav pozorujeme u úhlu θ , pravdìpodobnost støelení branky roste s velikostí úhlu, pod kterým hráè vidí brankovou konstrukci. Nikteøí trenèøi mláde¾e tento fakt vyu¾ívají pøi výchovì mladých fotbalových talentù, sna¾í se jim pøedat informaci, pokud mají dobrý výhled na branku

4. MODELOVÁNÍ PRAVDĚPODOBNOСТИ SKÓROVÁNÍ

stojí za to akci zákonit, jak si tuto informaci přebírají hráči je už na nich.

Obrázky 4.6 a 4.7 ukazují rozložení pravděpodobnosti skórování na hrací ploše, zlepšení oproti předchozím modelům je ihned pozorovatelné.

Obrázek 4.6: xG v závislosti na X a

Následující obrázek převádí obr. 4.6 do prostoru, jedná se o graf logistické funkce proměnných z tabulky.

$$xG = \frac{1}{1 + e^{2.0586 + 0.0503X - 1.7412}} \quad (4.2)$$

Obrázek 4.7: xG v závislosti na X a

4.2.2. Pokročilé modely logistické regrese

V předchozích modelech byly uvažovány pouze spojité proměnné, a už se jedná o vzdálenost či úhel. V této části budeme model s proměnnými, X postupně rozlišovat o

4.2. POKROÈILÉ MODELY

promìnné diskretní, mùže se jednat například o způsob zakonèení a dal'í. Zároveò budou modely testovány, jestli nám jejich rozl'íøení pøidá na jejich "kvalitì"nebo jestli pøidaná promìnná není pro modelování relevantní.

xG v závislosti na vzdálenosti, úhlu a "tlaku"

Zahrnování diskretních promìnných zaène pøidáním promìnné, která speci kuje, zda støílející hráè byl bìhem zakonèení pod tlakem èi nikoliv. Mùžeme ji speci kovat, jakožto okolnosti, které mohly ovlivnit hráèe bìhem støelby. Podstupoval zároveò souboj s protihráèem? Byl protihráèem napadán, že musel se støelou spìchat jinak by o míè pøíel nebo by jiže ke støele nemil prostor? Tyto a podobné otázky jsou zahrnuté v naší promìnné pojmenované "tlak", jelikož se jedná o diskretní promìnnou nabývá hodnot 1 - hráè byl pøi støele pod tlakem a 0 - hráè nebyl pøi støele pod tlakem, mìl na ni dostatek èasu a prostoru.

Tabulka 4.3: xG v závislosti na X , a tlaku

Promìnná	koefficient	[0,025	0,975]
konstanta	1,9392	1,736	2,142
X	0,0528	0,044	0,062
	-1,7850	-2,003	-1,567
tlak	0,6291	0,479	0,779

Deviance modelu je 9702,5

Vidíme, že pokud je hráè pod tlakem, 1ance na skórování se snižuje.

Nyní pomocí testu pomìrem vìrohodnosti zjistíme, zda je pøidaná promìnná, tedy tlak, pro model významná. Porovnáme pøedchozí model₁: $xG = f(X,)$ - model "níže" a poslední "rozsáhlejší" model_{v1}, $xG = f(X, , tlak)$ - model "vyšší".

Pro oba modely pomocí programu získáme logaritmicovou vìrohodnostní funkci (), bude urèena hodnota ² testu pomìrem vìrohodnosti, kterým je porovnáváme a výsledná p-hodnota bude porovnána s hladinou vìrohodnosti = 0.05.

$$L(v_1) = 4851,257934716189$$

$$L(n_1) = 4888,9960593471619$$

$$\text{výsledek } ^2 \text{ testu: } 75,47624959132008$$

$$p\text{-hodnota} = 4;07887732552535210^{-17}$$

Vidíme, p-hodnota < , odebráním promìnné tlaku, se sniží kvalita modelu, tato promìnná je pro model významná a proto budeme dále pracovat s modelem $xG = f(X, , tlak)$.

xG v závislosti na vzdálenosti, úhlu, tlaku a způsobu zakonèení

Nyní zahrneme způsob zakonèení, jakožto promìnné modelu. Uvažujeme zakonèení levou (lf) a pravou nohou (rf), tøetí způsob zakonèení, hlavou je zahrnut v konstantním èlenu.

4. MODELOVÁNÍ PRAVDĚPODOBNOSTI SKÓROVÁNÍ

Tabulka 4.4: xG v závislosti na X, , tlaku a způsobu zakonění

Proměnná	koefficient	[0,025	0,975]
konstanta	3,1537	2,877	3,431
X	0,0600	0,051	0,069
	-2,3270	-2,569	-2,085
tlak	0,2884	0,133	0,444
levá noha (lf)	-1,2302	-1,418	-1,042
pravá noha (rf)	-1,2242	-1,403	-1,046

Deviance modelu je 9477,2

Koefficienty u zakonění levou i pravou jsou záporné a podobné. Jelikož v modelu neuvažujeme kdo střílel, mají pro nás pravá i levá noha stejnou váhu, záporné znaménko souvisí s poslední možností zakonění a to hlavou, pravděpodobnost gólu vstřeleným nohou je vyšší oproti zakonění hlavou, proto po dosažení "jedničky" za jednu z proměnných lf, rf, což se v realitě rovná zakonění spodní končetinou, zvýšíme hodnotu xG střely. Test poměrem věrohodnosti vyšší následovní, za "vyšší" model χ^2 považujeme $xG = f(X, , tlak, rf, lf)$, za "nižší" model χ^2 bereme $xG = f(X, , tlak)$.

$$L(\chi^2) = 4738,609961574685$$

$$L(\chi^2) = 4851,257934716189$$

$$\text{výsledek } \chi^2 \text{ testu: } 225,295946$$

$$p\text{-hodnota} = 1:1956577 \cdot 10^{-49}$$

Opět uvažujeme hladinu významnosti $\alpha = 0,05$, p-hodnota $< \alpha$. Odebrání proměnných lf, rf značně ovlivnilo model, jsou pro něj významné. Tudíž dále budeme pracovat s modelem $xG = f(X, , tlak, rf, lf)$.

xG v závislosti na vzdálenosti, úhlu, tlaku, způsobu zakonění a straně hřiště, ze které bylo zakoněno

Nyní budeme zkoumat, zdali ovlivníme model, budeme-li jakožto proměnnou uvažovat stranu, ze které hráč zakoněl- pravá, levá strana hřiště. Jakožto proměnnou bereme například pravou stranu (rs), kterou stranu vezmeme model neovlivní, druhá strana se promítne do konstantního členu. Program vypočítal koefficienty a jejich intervalové odhady na hladině významnosti $\alpha = 0,05$ následovně:

Tabulka 4.5: xG v závislosti na X, , tlaku, způsobu zakonění a straně hřiště, ze které se střílelo

Proměnná	koefficient	[0,025	0,975]
konstanta	3,1980	2,915	3,481
X	0,0601	0,051	0,069
	-2,3232	-2,565	-2,081
tlak	0,2899	0,135	0,445
levá noha (lf)	-1,2236	-1,412	-1,036
pravá noha (rf)	-1,2295	-1,408	-1,051
pravá strana (rs)	-0,0866	-0,194	0,021

Deviance modelu je 9474,7

4.2. POKROÈILÉ MODELY

Testem pomìrem vùrohodnosti nyní bude zkontrolován vliv strany høi'ti na pravdìpodobnost skàování. "Vy11í model" $v_3 - xG=f(X, \text{tlak}, rf, lf, rs)$, porovnáme s modelem "ni¾4íím" $n_3 - xG=f(X, \text{tlak}, rf, lf)$.

$$L(v_3) = 4737,364539455062$$

$$L(n_3) = 4738,609961574685$$

výsledek χ^2 testu: 2,490844239246144

$$p\text{-hodnota} = 0,28781938827292586$$

Nyní ji¾ p-hodnota pøekroèila hladinu významnosti = 0,05, vliv strany, ze které bylo zakonèeno je pøi této hladinì zanedbatelný, dostateènì pøesné výsledky poskytuje "ni¾4í model" $n_3 - xG=f(X, \text{tlak}, rf, lf)$. Pravdìpodobnost skàování, xG bude tedy poèítáno v závislosti na vzdálenosti, úhlu, tlaku a zpùsobu zakonèení.

Tabulka 4.6: xG v závislosti na X, tlaku a zpùsobu zakonèení

Promìnná	koe ciny	[0.025	0.975]
konstanta	3,1537	2,877	3,431
X	0,0600	0,051	0,069
	-2,3270	-2,569	-2,085
tlak	0,2884	0,133	0,444
levá noha (lf)	-1,2302	-1,418	-1,042
pravá noha (rf)	-1,2242	-1,403	-1,046

Deviance modelu je 9477,2

Samotná funkce xG s danými koeficienty z tabulky vypadá následovně:

$$xG = \frac{1}{1 + e^{3,1537 + 0,06X - 2,3270 + 0,2884 \text{tlak} - 1,2302 lf - 1,2242 rf}} \quad (4.3)$$

Po testování jednotlivých modelů, zůstaly tyto proměnné jako¾to pro model relevantní, tohle tvrzení lze podpořit i shlédnutím jejich intervalových odhadů, kde žádný neobsahuje nulu.

Následující obrázky 4.8 ukazují některé případy výsledného modelu. Nejpatrnější je změna při zakonèování hlavičkou, kde je 1ance skórování ni¾4í, což je na grafu vidit jeho "zpløtí-ním" oproti dvom grafům nad ním. Zároveò je patrný vliv "tlaku" na hráèe, kdy graf potvrzuje, že hráè který není pod tlakem má z dané pozice vyšší 1anci na vstøelení branky.

4. MODELOVÁNÍ PRAVDĚPODOBNOTI SKÓROVÁNÍ

Obrázek 4.8: xG pro různé varianty výsledného modelu.

Obrázek 4.9: Porovnání dvou situací z obrázku 4.8

Na obrázku 4.9 jsou vyobrazeny situace, kdy model nabývá nejvyšších hodnot, pravděpodobnost padnutí gólu při střelbě levou nohou, když hráč je v klidu - modrá plocha, a situace, kdy nabývá hodnot nejnižších, pravděpodobnost padnutí branky při zakončení hlavou pod tlakem - zelená plocha. Rozdíl, jakožto výše zmíněné zploštiní grafu je u zelené plochy oproti ploše modré na první pohled zjevný. Graf dokazuje, že zakončení nohou v klidu má vyšší pravděpodobnost, že skončí brankou, oproti zakončení hlavou pod tlakem, uvažujeme-li zakončení ze stejné pozice na hřišti. Zelená plocha nabývá v každém bodě hrací plochy (souřadnice X,Y) nižší hodnoty než plocha modrá.

5. Závěr

Práce popisuje problematiku pravděpodobnosti skórování ve sportu, představuje jednotlivé metody používaných modelů ve sportech, zvláště podrobně pak ve fotbale. Zavede nejznámější metodu "očekávaných gólů", představí její výhody i nevýhody. Problematika je následně popsána ze statistického hlediska, zavedením zobecněných lineárních modelů, přesněji pak modelů logistické regrese, které jsou k interpretaci binomických závislých proměnných, jako je výskyt nebo nevýskyt určitého jevu, zde padnutí branky, vhodné. Představena je metoda maximální věrohodnosti, která je využita při výpočtu koeficientů jednotlivých proměnných modelu. Jednotlivé modely jsou poté charakterizovány pomocí jejich deviance a porovnány na základě věrohodnostních funkcí.

Představený matematický aparát je následně aplikován v programovacím jazyce Python, kde z open source dat [rmy Statsbomb](#) byly získány koeficienty jednotlivých spojitých i diskrétních vysvětlujících proměnných, pomocí kterých byly následně sestaveny modely pravděpodobnosti skórování v profesionálním fotbalu. Vypočtené modely s různými počty vysvětlujících proměnných byly následně pomocí testu věrohodnostním poměrem porovnávány. Bylo zjištěno, které vysvětlující proměnné mají statisticky významný vliv na pravděpodobnost vstřelení branky. Po testování několika modelů dochází práce k závěru, že z dostupných dat pravděpodobnost skórování ve fotbale nejlépe predikuje model využívající proměnné: způsob zakončení, herní situace a poloha na hrací ploše. Test věrohodnostním poměrem ukázal, že v profesionálním fotbale nehraje roli, jestli hráč zakončuje z levé nebo pravé strany hrací plochy.

Literatura

- [1] AGRESTI, Alan. Categorical data analysis. 2nd ed. Hoboken: Wiley, 2002, xv, 710 s. : grafy, tab. ISBN 0-471-36093-7.
- [2] ZVÁRA, Karel. Regrese. Praha: Matfyzpress, 2008, 253 s. : il. ; 24 cm. ISBN 978-80-7378-041-8.
- [3] ANDĚL, Jiří. Základy matematické statistiky. Vyd. 3. Praha: Matfyzpress, 2011, 358 s. : grafy, tab. ISBN 978-80-7378-162-0.
- [4] TIPPETT, James. The Expected Goals Philosophy. 1. Velká Británie: vydáno nezávisle, 2019. ISBN 978-1-08988-318-0.
- [5] Friends of Tracking, <https://www.youtube.com/channel/UCUBFJYcag8j2rm9HkrrA7w>
- [6] Premier League official website, <https://www.premierleague.com/>
- [7] Stats Bomb soccer data, <https://statsbomb.com/what-we-do/soccer-data/>
- [8] <http://allanderek.github.io/football-analysis/posts/shots-per-goal/>
- [9] Soccer analytics handbook <https://github.com/devinpleuler/analytics-handbook>
- [10] <https://www.statsmodels.org/>