

České vysoké učení technické v Praze
Fakulta jaderná a fyzikálně inženýrská

Katedra softwarového inženýrství

Obor: Aplikace informatiky v přírodních vědách



Analýza příčin vzniku shrinku produktů společnosti na základě logistických dat

Root Cause Analysis of Shrinkage Based on Logistics Data

DIPLOMOVÁ PRÁCE

Vypracoval: Bc. Anna Gruberová

Vedoucí práce: Ing. Martin Plajner, Ph.D.

Rok: 2023

České vysoké učení technické v Praze

Fakulta jaderná a fyzikálně inženýrská

Katedra softwarového inženýrství

Akademický rok 2022/2023

ZADÁNÍ DIPLOMOVÉ PRÁCE

Student:	Bc. Anna Gruberová
Studijní program:	Aplikace informatiky v přírodních vědách
Název práce:	Analýza příčin vzniku shrinku produktů společnosti na základě logistických dat
Název práce anglicky:	Root Cause Analysis of Shrinkage Based on Logistics Data

Pokyny pro vypracování:

1. Seznamte se s procesy v dodavatelském řetězci. Pozornost věnujte štihlé logistice a plýtvání v logistických procesech.
2. Provedte přípravu dat – definujte shrink, rozpoznajte a označte jej v obdržených datech.
3. Na základě předchozího kroku najděte a analyzujte příčiny, které vedou ke vzniku shrinku produktů dané společnosti. Zjištěné výsledky vhodně vizualizujte.
4. Navrhněte automatizaci tohoto postupu, kterou by bylo možné aplikovat i na logistická data dalších společností.

Doporučená literatura:

- [1] JIRSÁK, P., MERVART M., VINŠ M. *Logistika pro ekonomy - vstupní logistika*. 1. vydání. Praha: Wolters Kluwer ČR, 2012. ISBN 978-80-7357-958-6.
- [2] SUTHERLAND J., BENNET B. *The Seven Deadly Wastes of Logistics: Applying Toyota Production System Principles to Create Logistics Value*. Bethlehem, PA: Lehigh University, 2007.
- [3] HASTIE, T., TIBSHIRANI R., FRIEDMAN J. H. *The elements of statistical learning: data mining, inference, and prediction*. 2nd ed. New York: Springer, c2009. Springer series in statistics. ISBN 978-0-387-84857-0.

Jméno a pracoviště vedoucího práce:

Ing. Martin Plajner, Ph.D.

Oddělení matematické teorie rozhodování, Ústav teorie informace a automatizace AV ČR, v.v.i.

.....
vedoucí práce

Datum zadání diplomové práce: 12. 10. 2022

Termín odevzdání diplomové práce: 3. 5. 2023

Doba platnosti zadání je dva roky od data zadání.

.....
garant oboru

.....
vedoucí katedry

.....
děkan

V Praze dne 12. 10. 2022

Prohlášení

Prohlašuji, že jsem svou bakalářskou práci vypracovala samostatně a použila jsem pouze podklady (literaturu, projekty, SW atd.) uvedené v příloženém seznamu.

V Praze dne

.....

Bc. Anna Gruberová

Poděkování

Chtěla bych poděkovat za vedení mé bakalářské práce, za cenné rady a připomínky k tvorbě této práce a za čas strávený touto pomocí.

Bc. Anna Gruberová

Název práce:

Analýza příčin vzniku shrinku produktů společnosti na základě logistických dat

Autor: Bc. Anna Gruberová

Studijní program: Aplikace přírodních věd

Obor: Aplikace informatiky v přírodních vědách

Druh práce: Diplomová práce

Vedoucí práce: Ing. Martin Plajner, Ph.D.
Oddělení matematické teorie rozhodování, Ústav teorie informace a automatizace AV ČR, v.v.i.

Konzultant: –

Abstrakt:

Klíčová slova: Datová analýza, Logistika

Title:

Root Cause Analysis of Shrinkage Based on Logistics Data

Author: Bc. Anna Gruberová

Abstract:

Key words: Data Analysis, Logistics

Obsah

Úvod	11
1 Logistika	13
1.1 Definice Logistiky	13
1.2 Štíhlá logistika	13
1.2.1 Toyota Production System	14
1.2.2 Plýtvání v logistických procesech	18
2 Použité metody	19
2.1	19
2.2 Redukce dimenzionality	19
2.2.1 Analýza hlavních komponent	19
2.2.2 Korepondenční analýza	21
2.2.3 Metoda GUHA	22
2.2.4 Genetický algoritmus	23
3 Shrink	25
3.1 Definice	25
3.2 Typy shrinků	25
4 Zpracování dat	29
4.1 Popis obdržených dat	29
4.2 Období jednoho měsíce - damages shrinky	32
4.2.1 Předzpracování dat	32
4.2.2 Klasifikace dat	42
5 Korelační analýza	45
5.1 Postup	45
5.2 Implementace	47
5.2.1 Vstupy a výstupy	47
5.2.2 Popis funkcí a struktura kódu	48
5.2.3 Testování	51
5.3 Výsledky	51
Závěr	55
Literatura	56

Přílohy	60
A Obsah přiloženého CD	61

Úvod

Tato diplomová práce se zabývá

První kapitola se věnuje

V následující kapitole

Ve třetí kapitole se nachází

Poslední kapitola

- Data - jak jsou data uložena v DB u zákazníka - provázané - SQL příkazy - výběr proměnných - target hodnoty - cost, množství - produktová Hierarchie - Vymazání outlierů - outlier metody - businessově

- Co vysvětluje target - Miner - PCA

- Korelační analýza mezi produkty v rámci kategorie - korelace - Rozčlenění produktů

- Vizualizace dat - Jak funguje PBI - Seznam metri

Kapitola 1

Logistika

1.1 Definice Logistiky

Logistika zahrnuje všechny operace, které se týkají doručení zboží nebo služeb od výrobce k zákazníkovi, s výjimkou samotné výroby zboží nebo provádění služby. Výrobou je naopak rozuměno vše, co mění podobu materiálu. Během výroby se však logistika uplatňuje, například jako přesun materiálu nebo polotovarů mezi jednotlivými výrobními zařízeními. Operace lze rozdělit do tří hlavních toků: materiálový, informační a finanční tok. Materiálový obsahuje všechny pohyby týkající se fyzického materiálu, tedy jeho získávání, přesuny a skladování, a to jak mezi zákazníky, dodavateli či výrobními areály a sklady, tak i vnitřní pohyby mezi produkčními linkami nebo skladovými pozicemi. Informační tok popisuje procesy vznikající během materiálového toku, dále se do něj řadí analýzy již proběhlých toků a plánování a předpovědi budoucích toků. Poslední kategorie, finanční tok mapuje náklady způsobené předešlými dvěma zmíněnými toky.[1]

Pojem logistika je úzce propojen s pojmem Supply Chain Management (SCM)¹. Zatímco logistika se zabývá toky zboží, služeb či lidí, Supply Chain Management zahrnuje operace logistiky, navíc ale sleduje vztahy mezi procesory, které koordinuje a optimalizuje za účelem naplnění určitých cílů. Tímto cílem bývá často snížení nákladů v rámci částí procesu nebo zvýšení konkurenceschopnosti podniku [3]. Supply Chain Management se tedy prolíná s pojmem logistika a bývají často zaměňovány. Důvodem může být i to, že se jedná o nový pojem, který byl poprvé použitý v roce 1982.[2]

1.2 Štíhlá logistika

Štíhlost neboli „lean“ je koncept neustálého vylepšování procesu vytváření produktu nebo služby pomocí odstranění jakéhokoli plýtvání. Plýtváním rozumíme jakoukoli

¹Do češtiny lze Supply Chain Management přeložit jako řízení či správa dodavatelského řetězce. V českém prostředí se používá jak anglická tak česká podoba.

činnost, která v očích zákazníka nezvyšuje hodnotu produktu a tedy není ochotný za tuto činnost zaplatit ve formě vyšší prodejní ceny. Z této definice plýtvání je patrné, že pohled zákazníka hraje důležitou roli při vytváření hodnoty produktu ve štihlých systémech.[14, ?]

Svůj původ nachází štíhlá logistika na začátku 20. století, kdy Henry Ford zavedl pohyblivou montážní linku při výrobě automobilu Ford modelu T. Tato linka měla za následek několikanásobné snížení výrobního času a odstartovala sériovou výrobu aut. Díky čemuž se snížila prodejní cena, a automobily tak byly dostupné nejen nejbohatší vrstvě společnosti. Po druhé světové válce navázala automobilová společnost Toyota Motor Company na Fordovu efektivní montážní linku a vytvořila systém nazvaný Toyota Production System (TPS), který je přímým předchůdcem štíhlé logistiky.[13]

1.2.1 Toyota Production System

Toyota Production System je založen na pěti základních principech. Nejdůležitějším krokem je odstranit plýtvání. Je třeba se soustředit na jednotlivé procesy a na vazby mezi nimi. Pomocí metody genchi genbutsu² se nasbírají data a informace o procesech přímo na místě, kde procesy probíhají, aby případné problémy a zdroje plýtvání mohly být přesně určeny. Po této analýze se aplikuje přístup řešení problémů zvaný kaizen³, jehož cílem je kontinuální zlepšování procesů. Posledním z principů je dodržování vzájemného respektu mezi všemi oddělení společnosti, jak vedoucími pracovníky, tak zaměstnanci u výrobních linek. [13]

V TPS je plýtvání rozděleno do tří kategorií - Muda (plýtvání), Mura (nevyváženost) a Muri (přetěžování) [16]. V následující části jsou podrobněji popsány jednotlivé typy.

Muda

Japonské označení Muda v překladu znamená plýtvání, neúčinnost či marnost. Muda zahrnuje všechny činnosti, které nepřispívají ke zvyšování hodnoty produktu. Mudu lze rozdělit na dva podtypy – 1. typ zahrnuje aktivity, které jsou nezbytné pro koncového zákazníka, např. testování, zda je produkt nebo služba bezpečná. Druhý typ obsahuje ty procesy, které již zákazník nepotřebuje, či dokonce nechce, neboť mohou mít vliv na rychlost výroby produktu (výkonu služby) nebo přímo na jeho kvalitu.

Taiichi Ohno, manažer ve společnosti Toyota, identifikoval sedm typů plýtvání, někdy nazývané *seven deadly wastes*. Klasifikace a popis včetně příkladů je uveden níže [13]:

1. **Nadprodukce** – Pokud je vyrobeno více produktů, než je možné expedovat k zákazníkovi, nebo více materiálu, než kolik je požadováno k další výrobě či okamžité spotřebě.

²Genchi v překladu znamená skutečná lokace a genbutsu skutečná věc.

³Kaizen je japonský překlad slova zlepšení.

2. **Zpoždění/čekání** – Jakákoli prodleva mezi dvěma na sebe navazujícími procesy, např. čekání jedné montážní linky na meziprodukty z jiné linky vlivem rozdílných výrobních časů nebo vlivem nedostatečné výrobní kapacity jednoho ze strojů, dále sem patří také čekání zaměstnanců z důvodu kontroly odvedené práce, pomalého načítání počítačového programu nebo čekání na konkrétní instrukce k výkonu práce [14].
3. **Transport** – Zbytečný přesun produktů, materiálů nebo informací. Tento transport navíc může vést k poškození produktu. Příkladem tohoto typu plýtvání může být situace, kdy materiál, který je nejvíce potřebný pro výrobu produktů je umístěn v největší vzdálenosti, nebo pokud přístup k jedné položce ve skladu je blokován jinými položkami.
4. **Pohyb** – Zbytečný pohyb lidí, vzniklý špatným rozmístěním objektů v prostoru, např. nepřiměřeně dlouhotrvající chůze, natahování se pro předměty, vyhýbání se lidem či předmětům.
5. **Skladování** – Pokud je naskladněno více surovin, rozpracovaných výrobků a hotových produktů, než kolik je požadováno, např. předčasná dovážka položek do skladu, chyba v dodávce, naskladnění položek do zásoby tzv. pro jistotu nebo z důvodu množstevní slevy.
6. **Nadbytečné zpracování** – Při výrobě dochází k použití více energie nebo prostředků než nutné, nebo je vytvořen koncový produkt, který má vyšší hodnotu, než jaký je dohodnutý a požadovaný standard.
7. **Defekty** – Produkty či meziprodukty, které je nutné přepracovat nebo odstranit z výroby z důvodu vady.

Tyto podoby plýtvání aplikované v TPS byly inspirací pro identifikaci sedmi typů plýtvání v logistice [13, 4]:

1. **Nadprodukce** – V případě logistiky je nadprodukce chápána jako doručení produktů dříve nebo ve větším množství než bylo požadováno.
2. **Zpoždění/čekání** – Jakákoli prodleva mezi dvěma na sebe navazujícími procesy, např. čekání na převoz meziproduktů mezi dvěma výrobními linkami, příjezd kamionu mimo časové okno, doba mezi příjezdem kamionu a jeho naložením nebo čas mezi přijetím objednávky a zahájením její realizace.
3. **Transport** – Zbytečný přesun produktů, materiálů nebo informací, např. materiál, který je nejvíce potřebný pro výrobu produktů je umístěn v největší vzdálenosti, nebo pokud přístup k jedné položce ve skladu je blokován jinými položkami.
4. **Pohyb** – Zbytečný pohyb lidí, např. vzniklý špatnou organizací předmětů ve skladu, kdy položky, ke kterým se nejčastěji přistupuje, jsou v méně přístupných pozicích skladu, nebo dokonce sklad není strukturovaný vůbec, nebo nutnost změnit trasu při převozu položek ve skladu kvůli nedostatečně širokým uličkám.

5. **Skladování** – Pokud je naskladněno více surovin, rozpracovaných výrobků a hotových produktů, než kolik je požadováno, např. předčasná dovážka položek do skladu, chyba v dodávce, naskladnění položek do zásoby tzv. pro jistotu.
6. **Prostor** – Neoptimální využití dostupného místa, např. nedostatečná výška regálů ve skladech, nevyužitá kapacita regálů, neoptimální naložení kamionu, přetížení dostupných kapacit.
7. **Defekty** – Činnosti, které způsobí nutnost opakovat určitý proces, znehodnocení produktu nebo zvýší náklady, např. špatné zavezení produktu, špatné nebo chybějící označení produktu, chyby v evidenci.

V devadesátých letech, kdy se metody TPS začaly aplikovat ve společnostech, byl mezi sedm typů plýtvání Muda začleněn osmý typ - Dovednosti. V tomto případě dochází k neefektivitě kvůli nevyužití lidského potenciálu a talentů jednotlivých zaměstnanců. K tomu může docházet například striktním rozdělením na manažery a zaměstnance, kde role zaměstnanců je poslouchat nařízení shora a vykonávat práci tak, jak byla navržena vedoucími pracovníky. Avšak právě zaměstnanci pracující přímo v terénu lépe identifikují případné problémy a snadněji naleznou řešení díky svým zkušenostem.[14]

Mura

Mura lze přeložit jako nestejnomyšlnost, nevyrovnanost a nepravidelnost. Jedná se o plýtvání vznikající špatnou provázaností jednotlivých procesů a to jak interních, tak externích. Následkem nevyváženosti je pak vznik plýtvání Muda. [16, 4]

Plýtvání v podobě Mura se rozlišuje jak v procesech informačního, tak hmotného toku. V případě informačního toku je nejvýznamnějším zdrojem plýtvání situace, kdy je chybně predikována poptávka mezi jednotlivými články logistického řetězce. Ignorování vztahů mezi procesy může vést k chybovosti i v řádu desítek procent. Informace, jejichž opomíjení způsobuje chybovost předpovídání poptávky, mohou být např. v jaké fázi životního cyklu se výrobek nachází, plánování promoakcí nebo výrobní a logistická omezení dodavatelů. Další zdroj Mura v informačním toku je nedostatečná znalost stavu zásob mezi dodavatelem a odběratelem. Následkem čehož dochází k méně častým závazkům avšak s větším objemem, což vede k vyšším pojistným zásobám ve skladech. Většinu zmíněných situací lze eliminovat aplikováním konceptu „Just-in-Time“ do jednotlivých procesů. Plýtvání také vzniká při administrativě, pokud nejsou vhodně standardizované dokumenty používané v logistickém řetězci. Příkladem může být špatná evidence pohybů ve skladu či tvorba objednávek. Nesjednocenost v administrativních procesech vede ke zpomalení navažujících činností nebo dokonce k chybám, které způsobí nemožnost dokončení celého procesu. Pak je nutné vybrané procesy provést znovu a napravit chyby.[4]

Plýtvání v hmotném toku je přímým důsledkem chyb vznikajících v informačním toku. Lze identifikovat i takové zdroje plýtvání, které nesouvisejí přímo s informačním tokem, a to například dodržování různých standardů přepravních prostředků na

straně dodavatele a odběratele. To má pak za následek nadbytečné překládání materiálu do podoby, kterou druhá strana používá a se kterou je schopna následně efektivněji manipulovat.[4]

Muri

Pojem Muri označuje přetěžování. Muri často vzniká při snaze zvýšit produktivitu a odstranit tak předešlé typy plýtvání, v konečném důsledku může ale vést k výrazně větší chybovosti i celkovému selhání. Přetížení mohou být zaměstnanci, ale i stroje. V obou případech vytížení na více než 100 % se může projevit na snížení kvality výstupu. Lidé mohou být méně pozorní a může docházet k nehodám, které mohou v menší či větší míře negativně ovlivnit i větší část logistického řetězce. Stroje mohou produkovat zmetkové výrobky, nebo může dojít k jejich poškození až zničení.[4, 16]

Příklad plýtvání Muda, Mura a Mudi

Všechny tři zmíněné typy plýtvání Muda, Mura a Muri jsou navzájem propojené. Tuto skutečnost je třeba brát v potaz při řešení zefektivňování procesů a eliminaci plýtvání. Pro představu je uvedena následující situace. Společnost potřebuje zákazníkovi přivést šest tun materiálu, uloženého ve stejných jednotunových kontejnerech. Možné způsoby řešení této úlohy jsou znázorněny na obr. 1.1.[16]

Nejjednodušší možností je naložit na jeden kamion veškerý požadovaný materiál. V takovém případě společnost ušetří na počtu vozidel a eliminuje tak plýtvání přepravními prostředky, ušetří čas při nakládce a vykládce, protože není nutné obsluhovat více vozidel, zároveň . Na druhou stranu ale hrozí přetížení kamionu. Následkem přetížení se může zvýšit riziko nehody vozidla, firma může být pokutována nebo vozidlu nemusí být umožněn vjezd na určitá místa.

Opačným extrémem je použít tři kamiony, každý se dvěma tunami materiálu. Potom ale není efektivně využita dostupná kapacita a je patrné, že dochází k mnoha druhům plýtvání typu Muda.

Třetí možností je využití dvou kamionů, kdy první je naložen čtyřmi a druhý dvěma tunami. Toto rozložení nepodléhá žádným pravidlům a patrně proces nakládky není dostatečně spjatý s ostatními procesy nebo neprobíhá správný přenos informací o požadavcích mezi jednotlivými procesy. Nakládka a vykládka prvního velmi naloženého kamionu vyžaduje více času než druhého kamionu. Z toho plyne, že buď není možné v dostupném čase stihnout obsloužit první kamion a dochází k přetížení, anebo v případě druhého kamionu je zbude velké množství času a zaměstnanci zbytečně čekají. Z této volby plyne, že plýtvání typu Mura může způsobit Mudu i Muru.[16]

Optimální řešení je naložit dva kamiony po třech tunách, což je jejich ideální kapacita. V takovém případě společnost minimalizuje za daných podmínek všechny tři typy plýtvání. V reálném světě jsou situace mnohonásobně komplexnější a ne vždy existuje jednoznačné optimální řešení, které je navíc snadno dosažitelné. Důležité ale



Muri = overburdened



Mura = unevenness, fluctuation, variation



Muda = waste



No Muri, Mura, or Muda

Obrázek 1.1: Příklady plýtvání Muda, Mura a Muri při transportu šesti tun materiálu.[16]

je soustředit se na všechny tři typy současně, protože optimalizace pouze jednoho kritéria může způsobit jiný druh plýtvání nebo kolaps části systému.

V roce 2011 bylo realizováno dotazníkové šetření Vysokou školou ekonomickou v Praze, které mapovalo, kolik procent logistických expertů se zabývá odstraněním zmiňovaných tří typů plýtvání. Plýtvání Muda se snaží odstranit z logistických procesů 72 % respondentů, Murou se zabývá 39 % a plýtvání Muri řeší 30 % dotazovaných.[4]

1.2.2 Plýtvání v logistických procesech

Tato sekce se zabývá třinácti vybranými logistickými procesy z hlediska plýtvání, jak jsou uvedeny v knize *Logistika pro ekonomy – Vstupní logistika*. Analýza vychází z již zmíněného dotazníkového šetření z roku 2011.

Kapitola 2

Použité metody

2.1

2.2 Redukce dimenzionality

2.2.1 Analýza hlavních komponent

Analýza hlavních komponent (anglicky *Principal component analysis*, dále jako PCA) je statistická metoda využívaná pro extrakci proměnných, redukcí vícedimenzionálních dat nebo vizualizaci dat. Lze ji aplikovat pouze na kvantitativní data s numerickými, spojitými hodnotami, neboť metoda využívá lineární algebraické techniky, jako je například kovarianční matice, pro jejíž výpočet se předpokládají spojitě hodnoty.

Jednotlivá pozorování obsažená v datech bývají popsána několika různými proměnnými. Tyto proměnné jsou často vzájemně korelované a obsahují šum. Metoda PCA umožňuje extrahovat pouze důležité informace z proměnných a snížit šum. k tomu je třeba vypočítat nové ortogonální proměnné, nazývané hlavní komponenty, které se získají jako lineární kombinace původních proměnných [6]. Hlavní komponenty reprezentují směry největšího rozptylu původních dat a jsou řazeny podle své významnosti. Jinými slovy, první hlavní komponenta zachycuje co nejvíce variability v datech, druhá hlavní komponenta zachycuje co nejvíce variability, která nebyla zachycena první hlavní komponentou, pro zbylé komponenty analogicky.

Princip

Předpokládáme množinu dat $\mathbf{X} = (\mathbf{x}_1, \dots, \mathbf{x}_N)$, kde N je počet pozorování a každý vektor \mathbf{x}_i přísluší jednomu pozorování popsanému M proměnnými. \mathbf{X} je potom matice rozměru $N \times M$ vstupních dat. Dále je definovaný výběrový průměr $\bar{\mathbf{x}}$ jako

$$\bar{\mathbf{x}} = \frac{1}{N} \sum_{i=1}^N \mathbf{x}_i, \quad (2.1)$$

a výběrová kovarianční matice \mathbf{C}

$$\mathbf{C} = \frac{1}{N} \sum_{i=1}^N (\mathbf{x}_i - \bar{\mathbf{x}}). \quad (2.2)$$

První hlavní komponentu, která popisuje největší rozptyl dat označíme y_{1i} a vypočteme následovně jako lineární kombinaci původních proměnných

$$y_{1i} = \mathbf{a}_1^\top (\mathbf{x}_i - \bar{\mathbf{x}}), \quad \text{pro } i = 1, \dots, N, \quad (2.3)$$

kde $\mathbf{a}_1 = (a_{11}, \dots, a_{M1})^\top$ je vektor vah.

Optimální vektor \mathbf{a}_1 je takový vektor, který maximalizuje výběrový rozptyl nové proměnné y_{1i} za podmínky $\mathbf{a}_1^\top \mathbf{a}_1 = 1$. Pakliže je výběrový rozptyl y_{1i} definován jako

$$D(y_{11}, \dots, y_{1N}) = \mathbf{a}_1^\top \mathbf{C} \mathbf{a}_1 \quad (2.4)$$

můžeme maximalizační úlohu vyřešit pomocí metody Lagrangeových multiplikátorů. Lagrangeova funkce s parametrem λ_1 má následující tvar

$$\mathcal{L}(\mathbf{a}_1, \lambda_1) = \mathbf{a}_1^\top \mathbf{C} \mathbf{a}_1 - \lambda_1 (\mathbf{a}_1^\top \mathbf{a}_1 - 1). \quad (2.5)$$

Derivaci funkce položíme rovnou nule

$$\begin{aligned} \frac{\partial \mathcal{L}}{\partial \mathbf{a}_1} &= 2\mathbf{C}\mathbf{a}_1 - 2\lambda_1 \mathbf{a}_1 \stackrel{!}{=} 0 \\ (\mathbf{C} - \lambda_1 \mathbf{I})\mathbf{a}_1 &\stackrel{!}{=} 0, \end{aligned}$$

kde \mathbf{I} je jednotková matice.

Řešíme soustavu lineárních rovnic pro neznámý parametr \mathbf{a}_1 , která má řešení právě tehdy, když je matice $\mathbf{C} - \lambda_1 \mathbf{I}$ singulární, tedy platí, že její determinant je roven nule. λ_1 je pak největší vlastní číslo matice \mathbf{C} a \mathbf{a}_1 vlastní vektor příslušný tomuto vlastnímu číslu. Toto tvrzení se matematicky zapíše následovně

$$\mathbf{C}\mathbf{a}_1 = \lambda_1 \mathbf{a}_1. \quad (2.6)$$

Po vynásobení vektorem \mathbf{a}_1 zleva získáme řešení pro maximální rozptyl proměnné y_{1i}

$$D(y_{11}, \dots, y_{1N}) = \mathbf{a}_1^\top \mathbf{C} \mathbf{a}_1 = \lambda_1. \quad (2.7)$$

Druhá hlavní komponenta

$$y_{2i} = \mathbf{a}_2^\top (\mathbf{x}_i - \bar{\mathbf{x}}), \quad \text{pro } i = 1, \dots, N, \quad (2.8)$$

se vypočte obdobným způsobem s přidáním podmínkou ortogonalitě vzhledem k první hlavní komponentě – druhá hlavní komponenta nesmí být korelovaná s předchozí, první hlavní komponentou. Potom popisuje druhý největší možný rozptyl v datech. Znázornění dvou hlavních komponent ve dvoudimenzionálním prostoru je vyobrazeno na obrázku 2.1. Vektor \mathbf{a}_2 se opět získá jako jednotkový vlastní vektor kovarianční matice \mathbf{C} příslušící druhému největšímu vlastnímu číslu λ_2 . [6]

Obrázek 2.1: Znázornění dvou hlavních komponent na dvoudimenzionálních datech.

Získání předpisů pro další hlavní komponenty je analogické. Obecně lze zapsat metodu PCA a převod původních proměnných následujícím maticovým zápisem

$$\mathbf{Y} = \mathbf{X}\mathbf{A}, \quad (2.9)$$

kde \mathbf{X} je matice vstupních dat, \mathbf{A} je matice vlastních vektorů kovarianční matice \mathbf{C} . Pro matici \mathbf{A} zároveň platí $\mathbf{C} = \mathbf{A}\mathbf{A}^\top$, kde \mathbf{A} je diagonální matice vlastních čísel \mathbf{C} . [7]

2.2.2 Korespondenční analýza

Vícenásobná korespondenční analýza (anglicky *Multiple correspondence analysis*, dále jako MCA) je metoda, která umožňuje popsat vztahy mezi daty, které jsou popsány kategorickými proměnnými, vytvořením kontingenční tabulky. V případě, že se popisuje vzájemná relace pouze dvou proměnných, se použije základní korespondenční analýza¹. MCA je alternativou k PCA, pokud jsou analyzovanými daty kategorická data. [8]

Značení

Nechť \mathbf{N} je matice dat s rozměry $I \times J$, kde I odpovídá počtu pozorování a J je počet kategorií. Matice \mathbf{N} je převedena na korespondenční matici \mathbf{P} vydělením matice \mathbf{N} jejím celkovým součtem $n = \sum_{i=1}^I \sum_{j=1}^J n_{ij} = \mathbf{1}_I^\top \mathbf{N} \mathbf{1}_J$. To zaručuje, že součet prvků matice \mathbf{P} je roven jedné. Tyto kroky lze shrnout následujícím matematickým zápisem

$$\mathbf{P} = \frac{1}{n} \mathbf{N}, \quad \mathbf{P} = \{p_{ij}\}, \quad \sum_{i=1}^I \sum_{j=1}^J p_{ij} = 1. \quad (2.10)$$

Součet i tého řádku, resp. součet j tého sloupce je značen následovně

$$r_i = \sum_{j=1}^J p_{ij} \quad \text{pro } i = 1, \dots, I,$$

$$c_j = \sum_{i=1}^I p_{ij} \quad \text{pro } j = 1, \dots, J.$$

Vektor $\mathbf{r} = \mathbf{P} \mathbf{1}_J$ obsahuje všechny řádkové součty matice \mathbf{P} , analogicky vektor $\mathbf{c} = \mathbf{P}^\top \mathbf{1}_I$ obsahuje všechny sloupcové součty téže matice.

Pro další výpočty zavedeme značení pro diagonální matice, které mají na diagonále řádkový, resp. sloupcový součet

$$\mathbf{D}_r = \text{diag}(\mathbf{r}), \quad \text{resp.} \quad \mathbf{D}_c = \text{diag}(\mathbf{c}). \quad (2.11)$$

¹anglicky *correspondence analysis* (CA)

Výpočetní algoritmus základní korespondenční analýzy [10, 9]

Označme $\mathbf{S} = \{s_{ij}\}$ následující matici

$$\mathbf{S} := \mathbf{D}_r^{-\frac{1}{2}}(\mathbf{P} - \mathbf{r}\mathbf{c}^\top)\mathbf{D}_c^{-\frac{1}{2}}. \quad (2.12)$$

Po té provedeme singulární rozklad této matice

$$\mathbf{S} = \mathbf{U}\mathbf{\Delta}\mathbf{V}^\top, \quad (2.13)$$

kde $\mathbf{\Delta} = \mathbf{\Delta}^2$ je matice vlastních čísel λ_k pro $k = 1, \dots, K$, kde $K = \min\{I-1, J-1\}$. Potom rozměry matice \mathbf{U} , resp. \mathbf{V} jsou $I \times k$, resp. $J \times k$. Dále platí $\mathbf{U}^\top \mathbf{U} = \mathbf{V}^\top \mathbf{V} = \mathbf{I}$.

Korespondenční analýza měří míru váženého rozptylu, tzv. inercií pomocí vlastních čísel λ_k matice \mathbf{S} , λ_k se pak nazývají hlavní inercie. Celková inercie je rovna

$$I = \sum_{k=1}^K \lambda_k = \sum_{i=1}^I \sum_{j=1}^J s_{ij}^2. \quad (2.14)$$

Hlavní komponenta řádků \mathbf{F} je rovna

$$\mathbf{F} = \mathbf{D}_r^{-\frac{1}{2}} \mathbf{U} \mathbf{\Delta}. \quad (2.15)$$

Hlavní komponenta sloupců \mathbf{G} je rovna

$$\mathbf{G} = \mathbf{D}_c^{-\frac{1}{2}} \mathbf{V} \mathbf{\Delta} \quad (2.16)$$

Výpočetní algoritmus MCA

Předpokládejme, že původní matice kategorických dat má tvar $N \times Q$, tj. N pozorování a Q proměnných. Matici dat převedeme na indikátorovou matici. Indikátorová matice \mathbf{Z} je vytvořena tak, že kategorická data jsou rozepsána do pomocných proměnných. Pokud q tá proměnná je má J_q typů kategorií, tak příslušná indikátorová matice bude mít $J = \sum_{q=1}^Q J_q$ sloupců a N . Tzn. počet proměnných byl tímto rozepsáním rozšířen z počtu původních Q proměnných na J proměnných. První způsob MCA aplikuje základní algoritmus korespondenční analýzy na matici \mathbf{Z} , takto se získají souřadnice pro N pozorování a J kategorií.

2.2.3 Metoda GUHA

Uvažujeme potenciálně nekonečnou množinu objektů. Metoda předpokládá výsledky existenci pozorování této množiny reprezentované v matici dat. Cílem je získat všechny zajímavé vztahy, které jsou pravdivé pro celou množinu objektů

Tvrzení pravdivé o celé množině objektů verifikovaným teoretickým tvrzením. Tvrzení na datech z pozorování nazýváme observačním tvrzením. Důležité je, že mezi teoretickými a observačními tvrzeními obvykle platí vztah 1:1. Tento vztah je zprostředkován statistickými testy hypotéz.

GUHA procedura generuje a verifikuje napozorovaná data s množinou možných relevantních tvrzení, výstupem jsou potom všechna prostá observační tvrzení. *Relevantní tvrzení je prosté, je-li pravdivé v daných datech a zároveň neplyne zjevným způsobem z jiného, kratšího a ve výstupu již uvedeného tvrzení.*

Metoda pracuje s následujícími pojmy

- Matice dat a atributy - Řádky matice jsou jednotlivá pozorování. Atributem se rozumí sledovaná vlastnost, jedná se o sloupec matice.
- Základní booleovský atribut - Jedná se o výraz $A(a)$, kde A je atribut a a je vlastní podmnožina A . a může obsahovat více prvků než jeden.
- Booleovský atribut - Každý zákl. b. atribut je booleovský atribut. Boolovské atributy jsou negace, konjunkce a disjunkce základních boolovských atributů.
- Literál - Základní boolovský atribut nebo jeho negace
- Dílčí cedent - jedná se o konjunkci nebo disjunkci literálů.
- Cedent - Jedná se o konjunkci dílčích cedentů. Příkladem cedentu je booleovský atribut.

Procedury

4ft-Miner pracující s GUHA asociačními pravidly a s podmíněnými GUHA asociačními pravidly $\phi \approx \psi/\chi$

Asociační pravidlo - Asociačním pravidlem se tedy obvykle rozumí výraz $A \rightarrow S$ kde A i S jsou konjunkce dvojic atribut-hodnota.

Konfidence asociačního pravidla $A \rightarrow S$ je podíl (počet řádků splňujících A i S)/(počet řádků splňujících S)

2.2.4 Genetický algoritmus

Genetický algoritmus byl inspirován evoluční teorií a přirozeným výběrem. Je založen na teorii, že jedinci, kteří jsou lepší než ostatní mají větší šanci na to, předat svou genetickou informaci dál. Jejich geny tak budou základem nové generace. Podobně jako ve zmíněné teorii, algoritmus používá následující informace: genetickou reprezentaci pomocí bitového řetězce, funkci pro vyhodnocení tzv. fitness funkci, kombinování genů a mutaci.

Nejprve se vytvoří náhodná populace. Od této inicializační populace se postupně se iteruje dokud již další změny nevedou k lepšímu řešení, nebo dokud neskončí počet iterací. Jedna iterace je analogií k jedné evoluční generaci.

Bitový řetězec populace se vyhodnotí pomocí cílové funkce (objective function). Hodnota cílové funkce pak určuje hodnotu fitness daného řešení. Fitness je možné uvažovat jako minimalizační nebo maximalizační kritérium.

Dále se vyberou dva rodiče podle své fitness hodnoty, kteří se mezi sebou zkříží a jejich potomek vstoupí do další iterace. Jeden ze způsobů výběru je vybrat k náhodných populací a z nich vybrat populaci s nejvyšším fitness. Křížení rodičů je pravděpodobnostní, takže v některých případech ,uže vzniknout potomek stejný jako jeho rodič. Obvykle se hyperparametr pro křížení se nastavuje na vyšší hodnoty např. 80 nebo 90 %. Dalším parametrem je mutace, ta určuje zda přenesený bit zmutuje nebo ne. Jeho hodnota se nastavuje jako $1/L$, kde L je délka řetězce.

Kapitola 3

Shrink

Cílem této práce je analyzovat shrinky produktů, které byly zaznamenány v datech dané společnosti, a zjistit příčiny jejich vzniku. V následující části je vysvětlen pojem shrink a popsány kategorie, které vybraná společnost rozeznává ve svých datech.

3.1 Definice

Slovem shrink se označuje ztráta zisku z neuskutečněného prodeje hotového produktu. Tento produkt je vyroben, či naskladněn, ale z nějakého důvodu nemohl být prodán zákazníkovi. Tímto důvodem může být například poničení produktu, jeho ztráta nebo prošlá doba spotřeby. Za shrink produktu lze označovat i stav, kdy cena produktu je neplánovaně snížena v důsledku zmíněných důvodů. Shrinkem je potom rozdíl plánované prodejní ceny a ceny, za kterou byl produkt skutečně prodán.[17]

3.2 Typy shrinků

Vybraná společnost rozlišuje ve svých datech tři kategorie shrinku – inventory, damages, price downs. Dále se budu věnovat popisu jednotlivých typů v rámci kategorií. Každý typ má přiřazeno jednoznačné identifikační číslo, podle kterého je zaznamenáván v databázi. Z důvodů anonymizace dat v práci nejsou uvedené přesné hodnoty těchto ID, namísto toho jsou uvedeny názvy, které definují shrinky.

Damages

Do kategorie damages (česky škody) jsou řazeny zbylé důvody k odstranění produktu z prodeje z důvodu degradace produktu. V následující tabulce 3.1 jsou vypsané všechny typy, které moho být evidovány.

Tabulka 3.1: Přehled jednotlivých typů shrinků z kategorie damages.

Název	Popis
Poškození	Odpis zboží, které bylo poškozené. Např. nedopečené, spálené, špatně vyrobené nebo poškozené zaměstnancem nebo zákazníkem (kdy nelze uplatnit reklamaci na zákazníka.)
Prošlé a zkažené zboží	Odpis zboží, kterému prošla doba spotřeby (v případě výrobků, kde je datum uvedené), zkažené či shnilé zboží (ovoce, zelenina) nebo ztvrdlé pečivo.
Potravinová banka	Odpis potravinářského zboží, které bylo darováno potravinovým bankám. Jedná se o produkty, které nebylo možné zařadit znovu do oběhu.
Zvířecí útulky	Odpis potravinářského zboží, které bylo darováno do útulků zvířat. Jedná se o produkty, které nebylo možné zařadit znovu do oběhu.
Uznané zákaznické reklamace	Odpis zboží, které zákazník reklamoval a reklamace byla uznána, ale zároveň nelze toto zboží reklamovat u dodavatele.
Neupl. reklamace distribučního centra	Odpis zboží, které fyzicky nedorazilo z distribučního centra a nebylo možné ho reklamovat z důvodu nesplnění limitu pro vytvoření reklamace na distribučním centru. Také obsahuje odpisy neprodaných EXIT položek po ukončení výprodeje.
Kompostéry	Odpis zboží, které je prošlé nebo poškozené a které prodejna zlikviduje v kompostéru.
Poškozeno vnějšími vlivy	Odpis zboží, které bylo poškozeno nebo zničeno vlivem třetí strany (výbuch, vytopení, poškození majetku vloupáním) nebo přírodními živly. Zboží se tedy na prodejně nenachází a nemůže proto být zlikvidováno.
Zničeno rozmražením	

Inventory

Pojem inventory, který lze do češtiny přeložit jako inventář, sdružuje všechny shrinky týkající se změn v inventáři, tj. stavu zásob či inventory. V tabulce 3.2 se nachází přehled všech evidovaných typů.

Tabulka 3.2: Přehled jednotlivých typů shrinků z kategorie inventory.

Název	Popis
Inventura+	Kladné připsání zboží během inventory.
Inventura−	Záporné odepsání zboží během inventory.
Inventura	Velká inventura skladu.
Oprava inventory	Dodatečné opravy, které bylo třeba provést po dokončení velké inventory.
Částečná inventura	Odpis, nebo naskladnění zboží při inventuře položek.
Neuznané reklamace distribučního centra	Odpis zboží, které bylo fyzicky dodané distribučním centrem na prodejnu, ale prodejna jej vrátila, ale distribuční centrum vratku neuznalo.
Inventura	Starší verze ID používaného pro inventuru.
Partes merchandise	Odpis prokazatelně ukradeného zboží nebo i ztraceného zboží.

Price down

Kapitola 4

Zpracování dat

Úvod ke kapitole....

4.1 Popis obdržených dat

Všechna data poskytnutá společností jsou uložena v databázi, ke které byl zhotoven omezený přístup pro účely získání dat pro analýzy shrinku produktů společnosti. Zároveň s možností přístupu jsem obdržela i tabulku, která stručně komentuje všechny tabulky v databázi a sloupce v jednotlivých tabulkách. Celkem se v databázi nachází přes čtyři sta tabulek, z nichž bylo potřeba vybrat pouze ty, které obsahují relevantní data pro úlohu shrinků.

Z důvodu ochrany dat nelze uvádět přesné názvy tabulek, nicméně pro lepší orientaci v textu, každé použité tabulce přiřadím název, který odpovídá obsaženým datům v tabulce.

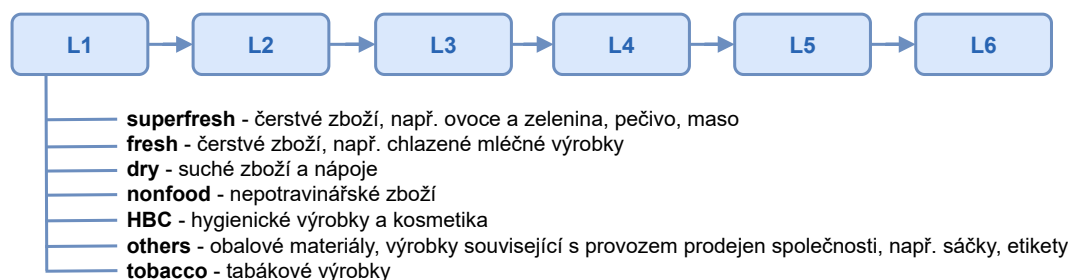
Produkty

Základní číselník s údaji o produktech, se nachází v tabulce **produkt** se 27 sloupci. Pro analýzu vzniku shrinků jsem z této tabulky vybrala jako možné významné údaje následující sloupce:

Tabulka 4.1: Vybrané sloupce v tabulce **produkt**.

Název	Popis
ID produktu	
ID prodejní varianty	Určuje o jaký typ balení daného produktu se jedná
Expirace	Expirace produktu ve dnech (hodnoty 0, 999 a NULL označují neomezenou dobu trvanlivosti)
ID kategorie	Kategorie produktu v číselné struktuře (pro lepší interpretaci, o jakou kategorie se jedná)

Tabulka `produkt_kategorie` obsahuje převod z číselné struktury, která popisuje kategorii produktu na 4BOX strukturu. 4BOX struktura se skládá ze šesti úrovní. Hierarchie těchto úrovní je znázorněná na obrázku 4.1 spolu se stručným popisem složek první úrovně.



Obrázek 4.1: Hierarchie struktury 4BOX používané pro kategorizaci produktů.

V tabulce 4.2 jsou uvedeny počty podkategorií pro každou z kategorií L1 a také procentuální zastoupení kategorie L1 v rámci produktového portfolia vybrané společnosti. Zastoupení je odvozeno podle počtu produktů v kategorii.

Tabulka 4.2: Počet podkategorií v L1 kategorii a zastoupení L1 kategorie v rámci produktového portfolia.

L1	Počty kategorií					Zastoupení kategorie
	L2	L3	L4	L5	L6	
nonfood	1	7	27	76	179	76,12%
dry	3	13	33	147	494	7,28%
HBC	1	4	21	59	193	7,07%
fresh	5	11	27	111	469	4,27%
superfresh	6	10	31	92	271	4,04%
others	4	4	4	5	5	1,06%
tobacco	1	1	1	3	8	0,17%

Kategorie L6 je přímo napojená na hodnotu číselné struktury, která je uvedena v číselníku produktů (v tabulce `produkt`). Pro získání všech úrovní kategorizace 4BOX k danému produktu je třeba vyhledat v tabulce `produkt` číselné ID kategorie daného produktu a napojit ho na kategorii L6 v tabulce `produkt_kategorie`. V této tabulce je pak hodnota nejbližší nadřazené kategorie L5. Poté je potřeba opět vyhledat v tabulce tuto hodnotu a zjistit její nadřazenou kategorii L4. Takto se postupuje dokud není dosaženo úrovně L1.

Další tabulka, se kterou jsem pracovala obsahuje informace o velikosti a hmotnosti produktů. Tato tabulka je důležitá z toho důvodu, že některé položky jsou vážené a pokud se udává jejich množství udává se v gramech, zatímco nevážené položky jsou uváděny v kusech. Aby bylo možné porovnávat oba číselné údaje, ke každému váženému produktu existuje přepočítaný počet kusů (ozn. SKU). K tomu jsou

využity údaje o počtu kusů na vychystávací jednotku (dále označeno jako SKU_{VJ}) a hmotnost jedné vychystávací jednotky daného produktu (ozn. $hmotnost_{VJ}$). Vychystávací jednotka je jednotka množství používaná pro vychystávání produktů – jeho balení a transport. Postup pro přepočty je následovný:

$$SKU_{vazena} = \frac{hmotnost_{vazena}}{hmotnost_{VJ}} \cdot SKU_{VJ}.$$

Tabulky transakcí

V tabulce **transakce** se nachází údaje o všech provedených transakcích, a to jak skladové transakce, tak prodeje a další pohyby na prodejnách. V případě prodejů prodejen jsou údaje agregované podle prodejny, konkrétního produktu a dne transakce, tzn. v této tabulce nelze rozlišit konkrétní prodeje na jednotlivých pokladnách, ale pouze souhrn za jeden den. Tabulka obsahuje údaje za posledních dvanáct měsíců.

Tabulka transakcí obsahuje 21 sloupců, jako možné podstatné sloupce pro analýzu jsem vybrala následující sloupce:

- ID transakce, jedinečné pro každou transakci.
- ID produktu, kterého se transakce týká. Každá transakce obsahuje údaje pouze o jediném produktu.
- ID prodejny (skladů i prodejen). Z důvodů ochrany dat jsem původní ID převedla na hodnoty od 0 do p , kde p je počet prodejen.
- Datum transakce (tzv. business datum, pokud samotná transakce proběhne až po půlnoci uvedeného dne, tak se posílá s datem z předchozího dne, neboť do toho dne businessově náleží.)
- ID promoce, příznak zda a v jaké promoční akci se produkt nacházel v čase uvedeném v datu transakce.
- ID shrinku, neboli *motive type*, z hlediska hledání shrinku produktu se jedná o klíčový sloupec, neboť obsahuje označení jednotlivých typů shrinků. Celkem je identifikováno 17 typů shrinků. V databázi parametr *motive type* označuje i jiná data než data týkající se shrinků, z toho plyne, že bylo třeba vyfiltrovat pouze ta data, která obsahují oněch 17 identifikačních čísel patřící shrinkům. Z důvodu anonymity dat jsem původní ID přechíslovala na celočíselné hodnoty od 0 až do počtu shrinků (pro obě hlavní kategorie).
- Množství produktu uvedené v transakci. U kusových produktů se jedná o celočíselný údaj u vážených to je desetinné číslo.
- Hodnota transakce v nákladové ceně v českých korunách (desetinné číslo).
- Hodnota transakce v prodejní ceně včetně DPH – v případě prodejů se jedná o skutečnou cenu, u zbylých transakcích je uvedena odpovídající cena podle ceníku.

Tabulku, která obsahuje údaje o jednotlivých prodeích na prodejnách společnosti, jsem pro účely této práce nazvala `transakce_prodeje`. Celkem obsahuje třináct sloupců, z nichž jsem vybrala následující:

4.2 Období jednoho měsíce - damages shrinky

Následující část text bude věnována rozboru dat pro shrinky typu damages v období jednoho kalendářního měsíce.

4.2.1 Předzpracování dat

Vzhledem k vysokému počtu dat pro jeden kalendářní rok, v roce 2022 bylo v databázi evidováno přes 32 milionů záznamů o týkající se shrinků, jsem se rozhodla provést analýzu na měsíčním výběru dat z tohoto období. Jako zkoumaný měsíc jsem vybrala měsíc říjen, neboť v porovnání s letními měsíci a Vánocemi se v říjnu nevyskytují významné sezónní výkyvy.

Zkoumaná říjnová data obsahují 2 712 669 řádků a patnáct sloupců. Každý řádek odpovídá jednomu záznamu v databázi shrinku daného produktu. Sledované údaje ve sloupcích jsou:

- ID prodejny, kategorická proměnná,
- ID produktu, kategorická proměnná,
- datum transakce, kategorická proměnná,
- typ shrinku, kategorická proměnná,
- L1, kategorická proměnná,
- L2, kategorická proměnná,
- L4, kategorická proměnná,
- L5, kategorická proměnná,
- L6, kategorická proměnná,
- expirace, kategorická proměnná,
- množství, spojitá proměnná,
- ztracená nákladová cena, spojitá proměnná,
- den v týdnu, kategorická proměnná,
- číslo den, kategorická proměnná,
- období v měsíci (rozdělení měsíce na pět částí), kategorická proměnná.

Původní sloupec datum jsem rozdělila na tři jiné proměnné, a to den v týdnu, číslo dne a období v měsíci a sloupec datum jsem vynechala. Z důvodu vysokého počtu záznamů a odlišné povahy dvou typů shrinků jsem data dále rozdělila na shrinky typu damages a shrinky typu inventory.

Výběr dat

Nejprve jsem graficky analyzovala zastoupení shrinků v závislosti na vybraných proměnných pomocí nástroje Power BI, viz obr. 4.2. V návaznosti na zjištěné zastoupení shrinků v datech jsem se rozhodla vybrat pouze ty typy shrinků, které tvoří více jak jedno procento z celkových nákladů (tj. náklady činily alespoň jeden milion korun). Vynechala jsem tedy shrinky s označením 5 až 9 a naopak shrinky 0 až 4 byly ponechány. Obdobně jsem přistupovala k záznamům i z hlediska kategorie produktu úrovně L1, jelikož z grafu je patrné, že majoritní zastoupení mají pouze dvě kategorie, a to kategorie superfresh a fresh produktů. Všechny záznamy se zbylými kategoriemi (HBC, others, nonfood, dry food a tobacco) jsem z datasetu odstranila. Těmito kroky jsem zredukovala původní počet řádků datasetu na 1 393 223 řádků.



Obrázek 4.2: Zastoupení shrinků typu damage a zastoupení kategorie L1 v datech z října roku 2022.

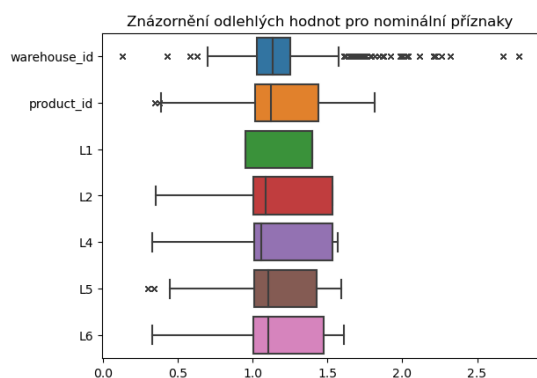
Jako cílové sloupce (*target* sloupce) jsem určila sloupec s typem shrinku, množstvím produktu a nákladovou cenou. Zbylých jedenáct sloupců slouží jako vysvětlující proměnné, dále budou označovány jako příznaky pro cílový sloupec. Všechny vybrané příznaky jsou kategorické proměnné, které lze dále rozdělit na nominální a ordinální. Nominální proměnné jsou ID prodejny, ID produktu, kategorie L1, L2, L4, L5 a L6. Ordinální proměnné jsou expirace, den v týdnu, číslo dne a období měsíce. Ordinální příznaky jsem přeznačila tak, aby každá obsahovala pouze hodnoty od nuly do n_p , kde n_p je počet kategorií v p -tém příznaku.

Pro další výpočty bylo vhodné přesunout se z nominálních kategorických hodnot na číselné hodnoty. Pro tyto účely jsem zvolila metodu *target encoding*. Neboť toto kódování na numerické hodnoty zachovává velikost datového souboru, to je klíčové vzhledem k tomu, že nominální proměnné ve zkoumaných datech obsahují velký počet kategorií. Např. počet unikátních produktů v datech je 19 026, což odpovídá stejnému počtu kategorií pro tuto proměnnou. Pokud bych použila one-hot kódování¹ mohlo by dojít k zásadnímu zvýšení počtu sloupců v datech, v tomto

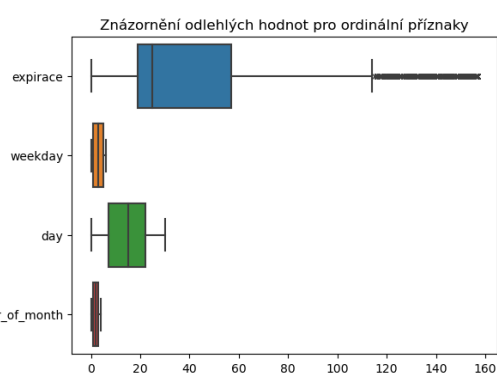
¹One-hot kódování převádí kategorické hodnoty na numerické takovým způsobem že pro každou kategorii vytvoří samostatný sloupec s binárními hodnotami, kde 1 odpovídá dané kategorii a 0

případě až o desítky tisíc. *Target kódování* je podobné převodu, který jsem použila pro ordinální proměnné. Avšak na rozdíl od něj, hodnota, která je kategorii přiřazena, souvisí se zastoupením této skupiny v cílovém sloupci a nesouvisí s uspořádáním hodnot uvnitř příznaku. Nevýhodou je, že takto upravená data mohou být náchylná na overfitting, proto je potřeba při predikování použít křížovou validaci.[19]

Dále jsem se zabývala identifikací odlehlých hodnot. Nejprve jsem vizualizovala hodnoty pomocí grafu, obrázky 4.3 a 4.4. Z grafu je patrné, že problémová je proměnná `warehouse_id`, která označuje ID prodejny. Prodejny, které tvoří outliersy mohou být malé prodejny, které kvůli menšímu počtu celkových produktů nevidují větší počet shrinků.



Obrázek 4.3: Znázornění odlehlých hodnot pro nominální příznaky.



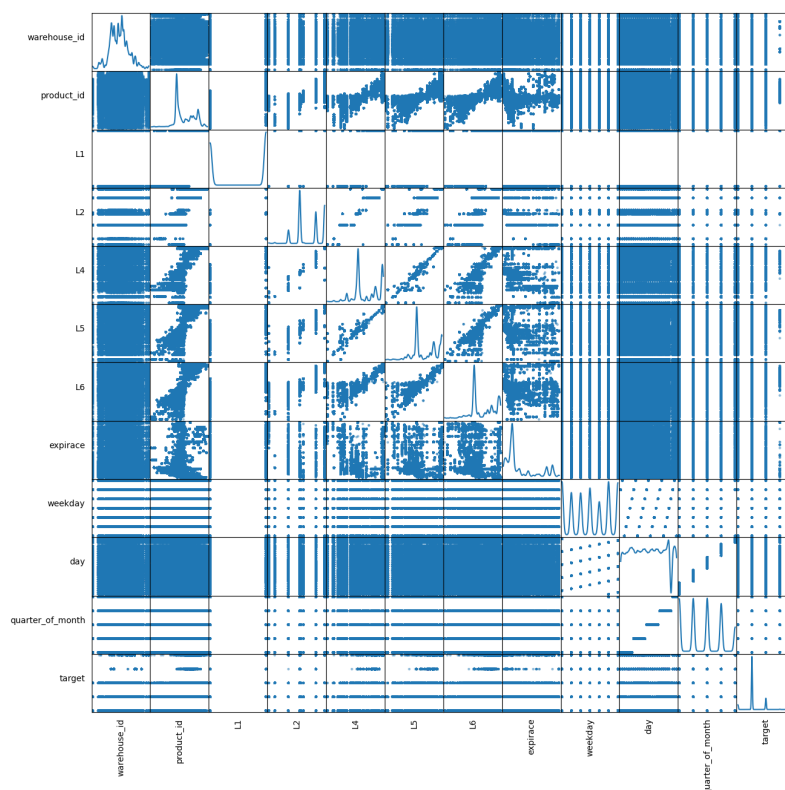
Obrázek 4.4: Znázornění odlehlých hodnot pro ordinální příznaky.

Pomocí Tukeyho testu jsem identifikovala přes 150 000 outlierů pro příznak ID prodejny (`warehouse_id`), čímž se dataset zredukoval na 1 218 453 řádků. S tímto krokem klesl i počet ostatních outlierů.

V dalším kroku jsem se zaměřila na míru korelace mezi proměnnými. Vizualizovala jsem data pomocí scatter matice pro všechny proměnné, matice je možné vidět na obr. č. 4.5. Z této matice můžeme na první pohled vidět, že příznaky odpovídající 4BOX kategorizaci a ID produktu vykazují závislost, což plyne z definice uspořádání této hierarchické kategorizace. V následujících krocích je cílem vybrat tu kategorii, která nejlépe popisuje data ve vztahu k shrinkům.

Jako první metodu jsem zvolila χ^2 test. Vzhledem k vysokému počtu dat je matice příliš řídká, a proto nejsou výsledné hodnoty vypovídající a test je tedy pro tuto úlohu nespolehlivý. Jiným měřítkem pro korelaci mezi proměnnými je Pearsonův korelační koeficient. Výslednou matici popisující korelační vztahy mezi příznaky jsem vizualizovala teplotní mapou, která je zobrazena na obrázku 4.6. Z výsledků je opět patrné, že mezi jednotlivými kategoriemi produktů a produkty je silná korelace. Toto zjištění je zcela logické, neboť se jedná o stromovou strukturu kategorií. Zároveň existuje korelace mezi produktovými kategoriemi a expirací produktu. p-hodnota odpovídající jednotlivým koeficientům byla vždy nulová, kromě pro koeficient týkající se dvojice proměnných `expirace` a `ID prodejny` a `expirace` a pořadí

zbylým kategoriím.



Obrázek 4.5: Scatter matice příznaků.

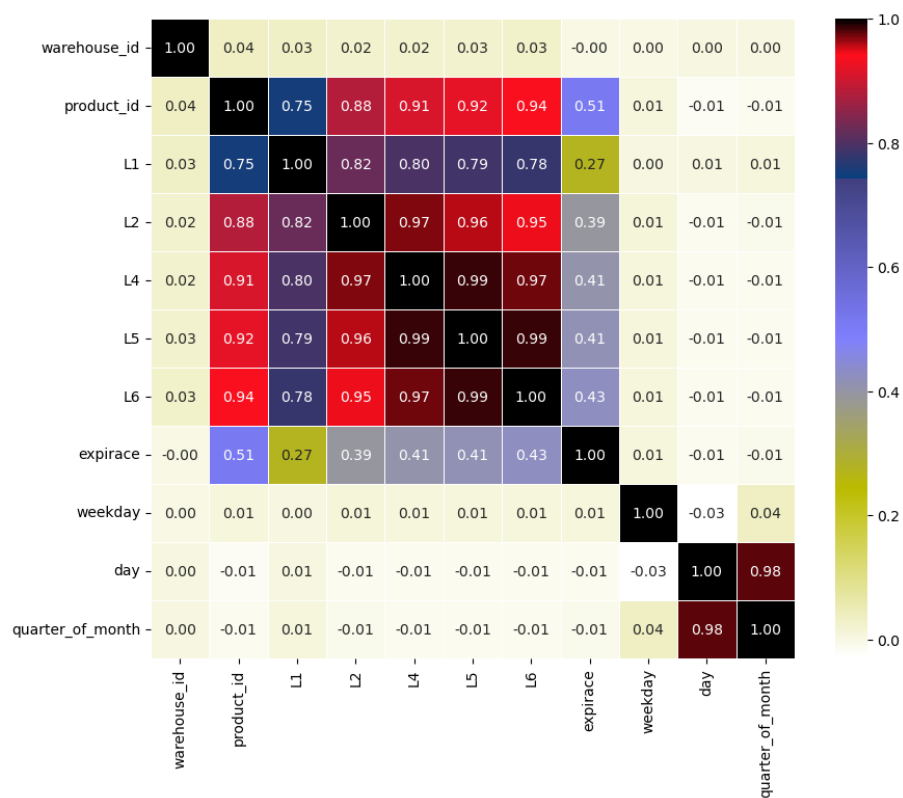
dne v týdnu. Je tedy možné považovat výsledky (kromě těchto dvou výjimek) za statisticky významné.

Dále jsem použila výpočet koeficientů vzájemné informace², který říká, jaká je podobnost mezi dvěma proměnnými [20]. Matice vypočítaných koeficientů je na obr. 4.7, jelikož se jedná o symetrickou vlastnost, jsou vynechány hodnoty pod vedlejší diagonálou. Z výsledků je opět vidět, že ID produktu sdílí informaci s úrovněmi kategorizace tím více, čím je kategorizace jemnější.

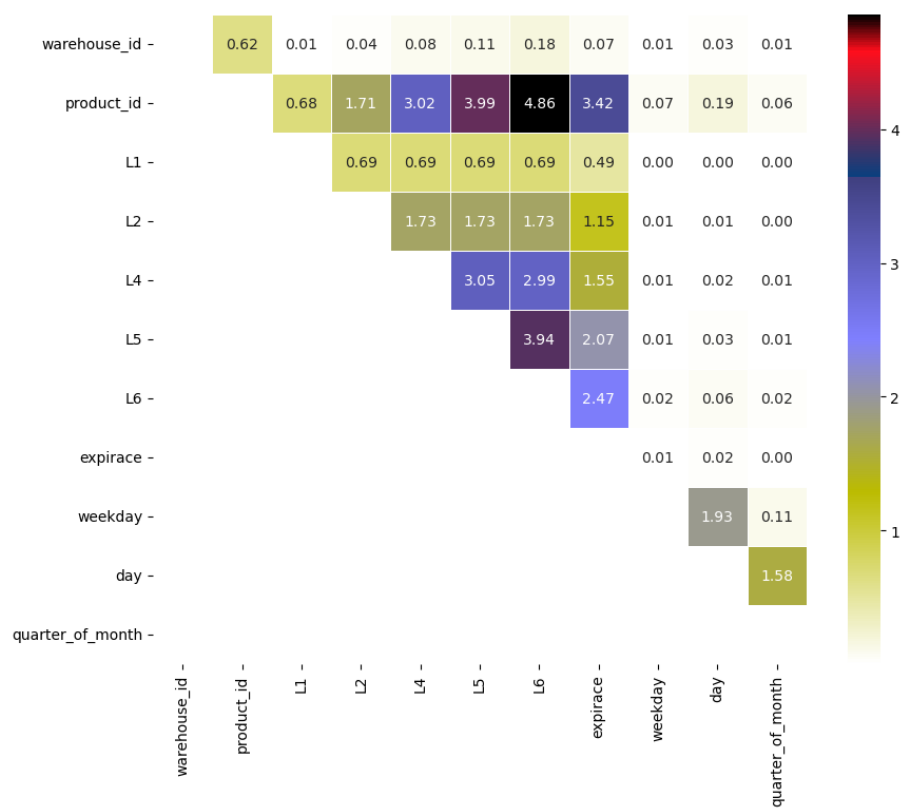
Dále jsem pro znázornění vztahu mezi proměnnými použila koeficient Cramerovo V. Koeficient jsem postupně počítala pro každou dvojici příznaků. Koeficient nabývá hodnot mezi 0 a 1. Číslo přibližné nule indikuje, že mezi proměnnými není asociace, číslo blízké jedničce vysokou závislost [22]. Na obr. 4.8 lze vidět, že pro kategorie L1 až L6 je hodnota koeficientu po zaokrouhlení rovna jedné. Vysoká závislost je pak i mezi příznakem expirace a ID produktu a kategorií L1. Dále logicky mezi číslem dne a dnem v týdnu a obdobím v měsíci.

Další statistikou spočtenou na datech je Theilovo U (neboli koeficient nejistoty), který opět nabývá hodnot mezi 0 a 1 a měří vztah mezi dvěma proměnnými. Na rozdíl od předchozích statistik tento koeficient není symetrický a z výsledků lze

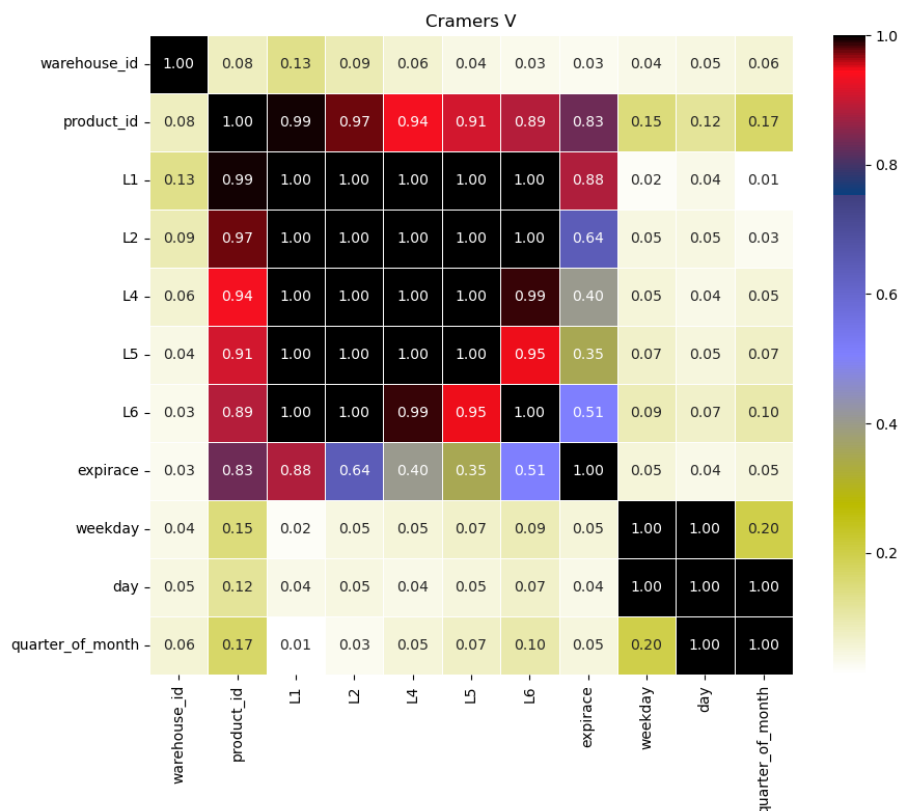
²*mutual information*



Obrázek 4.6: Matice korelačních koeficientů mezi příznaky.



Obrázek 4.7: Matice koeficientů vzájemné informace mezi příznaky.



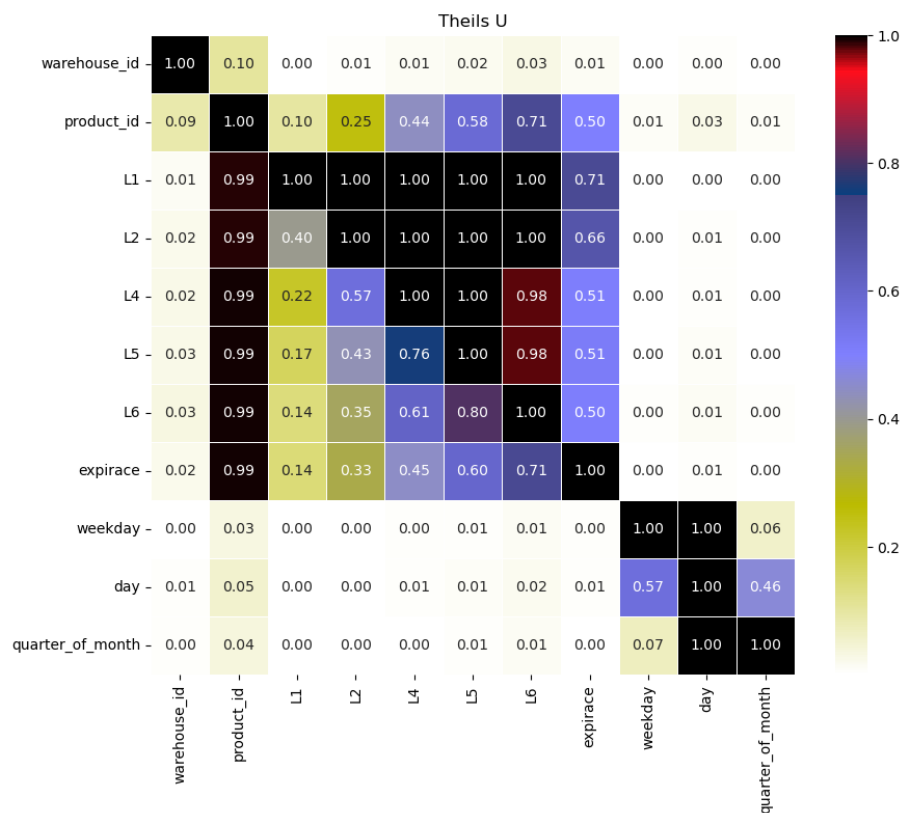
Obrázek 4.8: Matice koeficientů Cramerovo V mezi příznaky.

vyvodit, ze které proměnné ze dvou zkoumaných můžeme vyvodit informaci o druhé proměnné [23]. Z výsledků zobrazených v matici na obr. 4.9 plyne, že z ID produktu lze vyvodit část informace o kategoriích a expiraci. Zatímco úrovně L1 a L2 o ID produktu mnoho informace nenesou. Jak bylo ukázáno i v předchozích statistikách a jak vyplývá z logiky pro získání dne v týdnu a období měsíce, číslo dne nese informaci o těchto dvou příznacích.

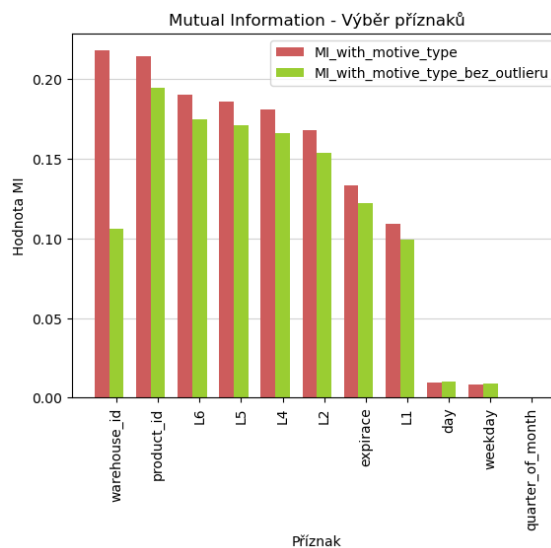
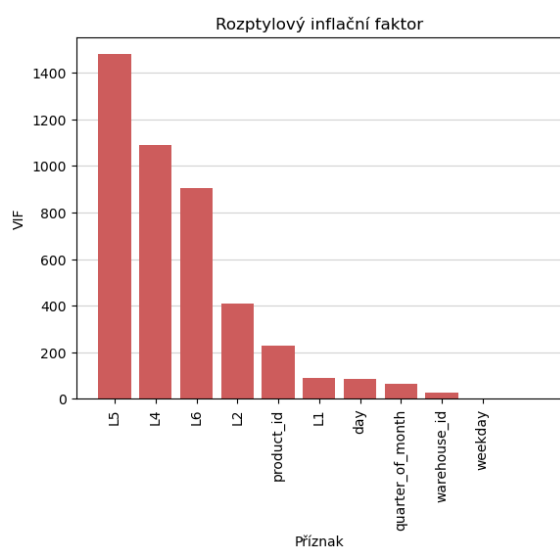
Z vypočítaných statistik na datasetu je patrné, že některé příznaky jsou významně závislé, a proto je třeba je z dat odstranit. Kandidáti na vynechání jsou kategorie L2, L4, L6 a číslo dne. V dalších testech budou také vybráni kandidáti a v závěru vyhodnotím, které příznaky byly podle aplikovaných metod vybrány jako vhodné k vynechání a které nikoli.

V dalším testu jsem otestovala multikolinearitu dat pomocí rozptylového inflačního faktoru (VIF). Jako hraniční faktor jsem zvolila hodnotu 40 VIF. Vysvětlující proměnné jsem odebírala z datasetu postupně a odebírání jsem ukončila až, když hodnota VIF nebyla nižší než hraniční. Tímto došlo k redukci příznaků z jedenácti na pět, a to na kategorii L1, číslo dne, období měsíce, ID prodejny a den v týdnu. Hodnoty koeficientu VIF na datech jsou na obr. 4.10.

Jako další metodu po výběr příznaků jsem vypočítala hodnotu koeficientů vzájemné informace mezi všemi příznaky s cílovým sloupcem - ID shrinku. Na obrázku 4.11 lze vidět, jak jednotlivé proměnné souvisí s cílovým sloupcem. Pro výpočet tohoto koeficientu jsem použila jak data bez outlierů, tak tenkrát data před jejich od-



Obrázek 4.9: Matice koeficientů Theilovo U mezi příznaky.

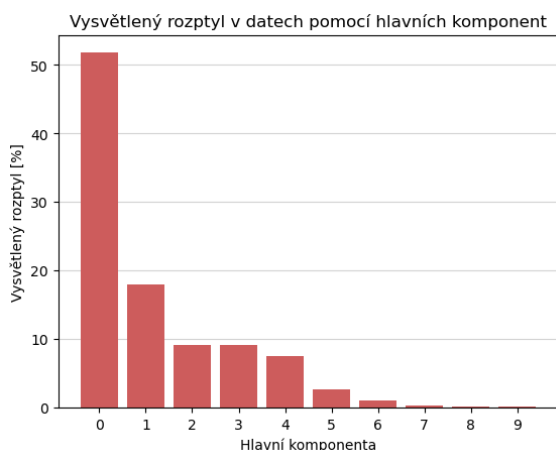


Obrázek 4.10: Rozptylový inflační faktor. Obrázek 4.11: Koeficienty vzájemné informace mezi příznaky a cílovým sloupcem typ shrinku.

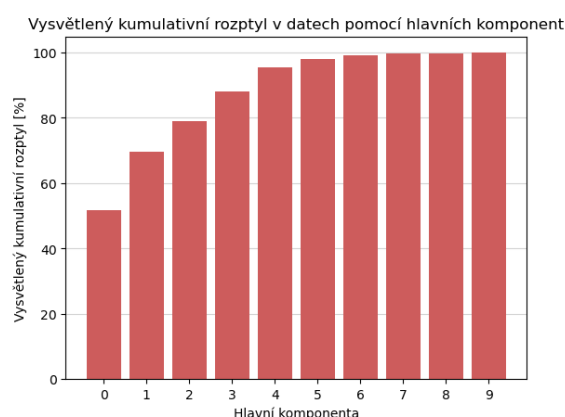
straněním. Zde můžeme vidět, že významnost příznaku ID prodejny klesla o téměř polovinu. Nejvíce informace je sdílena s ID produktu, kategorií L6, dále L5, L4, L2 a expirace. Příznaky související s časovými údaji podle tohoto kritéria nenesou mnoho společné informace.

Jako hlavní metodu pro výběr proměnných jsem se rozhodla použít metodu PCA, tuto metodu je možné použít protože kategorické proměnné jsem převedla na číselné hodnoty v předchozích krocích. Alternativou by bylo použití metody MCA, která se používá pro kategorické datasety, viz dále. Ve své práci jsem využila implementaci PCA v knihovně *Prince* v jazyce Python. Předtím než jsem metodu aplikovala jsem otestovala předpoklad homoskedasticity, tedy shodnost rozptylů v datech, pomocí Bartlettova testu implementovaného v knihovně *factor_analyzer*. Nulová hypotéza o shodnosti rozptylů nebyla vyvrácena (p-hodnota vyšla nulová). Metodu PCA je proto možné použít.

Na obrázcích 4.12 a 4.13 je znázorněno prvních deset komponent a rozptyl který v datech vysvětlují. Na základě hodnot jsem vybrala prvních pět komponent. Již pátá komponenta (označená č. 4) spolu s předchozími vysvětluje více jak 95 % variability dat. V dalším kroku jsem vypočítala příspěvky příznaků k těmto pěti komponentám a vybrala jsem ty příznaky, které přispívají nejvíce k prvním pěti komponentám. Jejich příspěvek je znázorněný na obr. 4.14. Na základě výsledků analýzy hlavních komponent byly vybrány jako vhodné příznaky pro další práci s daty tyto příznaky - ++ ID prodejny, den v týdnu, expirace, den a období v měsíci.

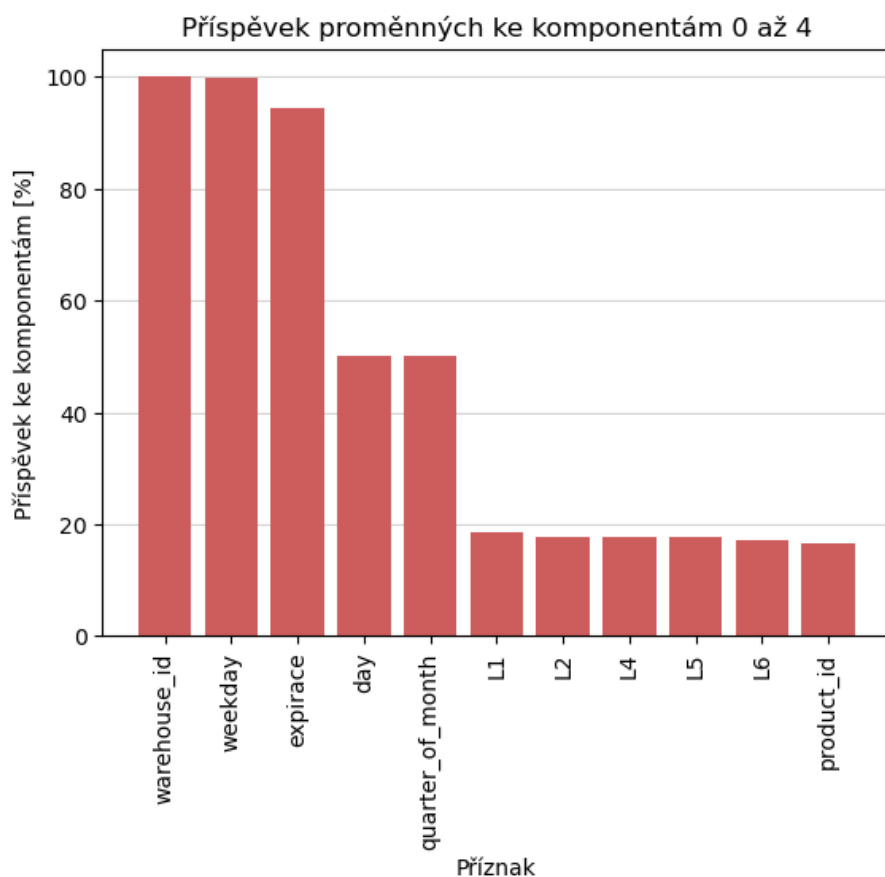


Obrázek 4.12: PCA - vysvětlený rozptyl hlavních komponent.



Obrázek 4.13: PCA - kumulativní vysvětlený rozptyl hlavních komponent.

Jak již bylo zmíněno pro redukci dimenzionality u kategorických dat lze použít metodu MCA, opět jsem využila implementaci z knihovny *Prince*. V této implementaci jsou nominální kategorické hodnoty kódovány tak, že narůstá počet sloupců, a proto bylo nutné, vzhledem k nárokům na paměť k uložení matice, omezit množství dat. Vybrala jsem náhodný 20% vzorek dat, na které jsem MCA aplikovala. Vypočítala jsem prvních pět komponent, které dohromady popisují 79 % variability dat. Jelikož byla každá kategorie chápána jako samostatná proměnná příspěvky jednotlivých příznaků ke komponentám byly rozmístěny mezi všechny kategorie, nikoli k jednotlivým příznakům. Po agregaci podle původních příznaků největší příspěvek mělo ID



Obrázek 4.14: Příspěvek proměnných ke komponentám 0 až 4.

produktu, kategorie L6, L5, L4, zatímco nejmenší ID prodejny, L1, den v týdnu a období měsíce. Tyto výsledky je třeba brát se zvážením neboť výpočty probíhali na řádově menším vzorku než u předchozích metod.

Shrnutí pro výběr dat

Na základě předchozích metod byly původní příznaky datasetu zredukovány na menší počet. Vzhledem k tomu, že různé metody vybraly různé příznaky, bylo stanoveno více možných výběrů.

Korelované jsou hodnoty ID produktu, L6, L5, L4 a expirace. Dále také z označení dne lze určit období měsíce. Ze zmíněných korelovaných příznaků není proto vhodné začlenit více než jeden příznak. Pokud je tato myšlenka aplikována na výsledky metod PCA a MCA a výsledků zjištěných pomocí hraniční hodnoty VIF.

1. Následující sloupce byly získány podle hodnoty rozptylového inflačního faktoru. Touto metodou byl navržen i sloupec s číslem dne, ten však z důvodů korelace nebyl zahrnutý

1.1. L1, období měsíce, ID prodejny a den v týdnu.

K této variantě existují i dvě alternativy, ve kterých je obměněna úroveň kategorizace produktu:

1.2. L5, období měsíce, ID prodejny a den v týdnu

1.3. L4, období měsíce, ID prodejny a den v týdnu

2. Metodou PCA bylo zjištěno, které příznaky nejvíce přispívají ke komponentám, které popisují téměř 96 % rozptylu v původních datech - jedná se o příznaky ID prodejny, den v týdnu, expirace, období v měsíci a číslo dne. Naopak metoda MCA vybrala kategorie L4 až L6 jako důležité. Sloučením a přihlédnutím ke korelačním koeficientům byly vybrány pět příznaků

2.1. ID prodejny, den v týdnu, expirace, období v měsíci, L5.

Tato varianta příznaků byla ještě rozšířena o příznaky, které se týkají produktů. Přidané příznaky jsou spolu korelované, přesto

2.2. ID prodejny, den v týdnu, expirace, období v měsíci, L5, L2

2.3. ID prodejny, den v týdnu, expirace, období v měsíci, L5, L2, ID produktu

2.4. ID prodejny, den v týdnu, období v měsíci, L2, ID produktu

L2, L5, období v měsíci, ID prodejny, den v týdnu, ID produktu, target

Všech sedm možných výběrů bylo otestováno metodou gradient boosting. Pro další výpočty byla použita pouze varianta, která vykazovala nejlepší přesnost. Tabulka 4.3 uvádí získané přesnosti.

Tabulka 4.3: Tabulka dosažených přesností dosažených metodou gradient boosting pro varianty výběru příznaků.

Varianta	Přesnost [%]	
	Trénovací data	Testovací data
1.1.	79,27	79,15
1.2.	82,88	82,77
1.3.	82,80	82,67
2.1.	83,21	83,07
2.2.	83,38	83,30
2.3.	83,67	83,54
2.4.	83,44	83,33
2.5	83,60	83,45

4.2.2 Klasifikace dat

Tato část se věnuje předpovědi typu shrinku z dostupných dat. V předchozích sekcích bylo popsáno předzpracování dat a výběr vhodných příznaků pro úlohu klasifikace. Byly navrženy dvě skupiny příznaků, na kterých budou provedeny výpočty. K obě variantám se bude přistupovat stejným postupem a následně budou porovnány dosažené výsledky.

Cílový sloupec, který je předpovídan, je pouze jeden. Jedná se o ID shrinku. To obsahuje pět různých kategorií (označené číslicemi od 0 do 4). Proto úlohu můžeme označit jako klasifikační úlohu pro více tříd (neboli *multiclass classification*). Jazyk Python nabízí v knihovně *scikit-learn* řadu metod, které podporují klasifikace do více tříd.[21]

Vybrala jsem následující metody pro klasifikaci ID shrinku:

- logistická regrese OVR,
- multinomická logistická regrese,
- random forest klasifikátor,
- gradient boosting klasifikátor.

Logistickou regresi jsem použila jako základní metodu pro klasifikaci v případě, že vstupní dataset se skládá z kategorických proměnných [25]. Balíček *scikit-learn* umožňuje klasifikaci do více tříd spočítat dvěma způsoby, které se liší v přístupu provedení klasifikace. První přístup využívá schématu OVR (*One-vs-Rest* neboli jeden proti všem). Při použití OVR se každá třída trénuje samostatně. Pro každou třídu je úloha převedena na binární klasifikaci, kdy zkoumaná třída je označena jako pozitivní a všechny zbylé jako negativní. Pokud máme N tříd, pak je vyhodnoceno N binárních logistických regresí. Naproti tomu multinomická log. regrese nevyhodnocuje třídy odděleně, ale používá funkci softmax. Ta predikuje zda, daný bod náleží do jedné z tříd.[26]

Další zvolenou metodu je klasifikátor implementující random forest algoritmus. Tento algoritmus jsem zvolila vzhledem ke skutečnosti, že je úspěšně využíván pro problémy z reálného světa a dovede pracovat s velkým objemem dat, které tyto úlohy obvykle zahrnují [27]. Zároveň volba parametrů pro tuto metodu je intuitivní. Poslední zvolenou metodou je klasifikátor, který využívá gradient boosting. Tento klasifikátor také vytváří rozhodovací stromy jako random forest. Narozdíl od zmíněného klasifikátoru, jsou ale stromy vytvářeny postupně v závislosti na naposledy vytvořeném stromu. Stromy jsou k sobě agregovány během procesu trénování. Zatímco random forest vytváří stromy nezávisle a agreguje je až na konci procesu.[28]

Výsledky

Nejprve jsem pracovala s vybranými příznaky - L1, období měsíce, ID prodejny a den v týdnu. Data jsem rozdělila na tři skupiny - data pro trénování, validaci a testování v poměru 8:1:1.

Naimplementovala jsem metodu `perform_classification` umožňuje spustit vybraný model s požadovanými parametry z knihovny *scikit-learn*, provede *k*-fold crossvalidaci, nafituje model na trénovací data a poté ověří přesnost na trénovacích datech. V případě, že jsou předány i parametry pro ladění, na validačních je model doladěn a pak na nejlepších parametrech opět otestován. Celá metoda se nachází v příloze této práce.

```
01 def perform_classification(  
02     model,  
03     parameters,  
04     train_x, train_y, valid_x, valid_y, test_x, test_y,  
05     tuning_parameters,  
06     k-fold=5)
```

V tabulce 4.4 jsou uvedeny přesnosti pro čtyři vybrané klasifikační metody. Uvedena je jak přesnost na trénovacích datech, tak na testovacích datech. Ve všech metodách byla použita metoda křížové validace, kdy data byla rozdělena do pěti skupin. Výsledná přesnost je pak průměrem dílčích přesností.

Tabulka 4.4: Tabulka dosažených přesností pro čtyři vybrané klasifikační metody na datech se shrinky typu damages s vybranými příznaky podle varianty 1.

Metoda	Přesnost [%]	
	Trénovací data	Testovací data
Logistická regrese OVR	77,04	76,95
Multinomická logistická regrese	77,33	77,27
Random forest	82,54	82,12
Gradient boosting	83,80	83,78

Nejlepších výsledků dosahuje klasifikátor gradient boosting. Přesnost na testovacích datech je téměř 84 %. V obou metodách, které využívají rozhodovací stromy, jsem implementovala ladění parametrů. Na následujících obrázcích ?? je znázorněna závislost mezi jednotlivými hodnotami parametrů a dosažené přesnosti.

V případě random forest klasifikátoru se jedná o parametry, které určují počet stromů, hloubku stromu, minimální počet vzorků, který má obsahovat list a minimální počet vzorků, kdy se může rozdělit uzel stromu. Pro klasifikátor gradient boosting byly také laděny parametry pro počet stromů, hloubku a dále míru učení.

Metoda random forest

Kapitola 5

Korelační analýza

Vypočítala jsem korelaci mezi hodnotou shrinkem a tržbami. Jedno pozorování je na agregované na produkt, prodejnu a den záznamu. V rámci analýzy jsem srovnávala mezi pouze záznamy produktů, které se vyskytují ve stejné kategorii. Je možné specifikovat úroveň kategorie, na které se analýza spočítá. Dále je třeba určit konkrétní název kategorie.

5.1 Postup

Hodnotu shrinku jsem porovnávala s následujícími ukazateli.

- Tržby daného produktu.
- Tržby daného produktu, které byly v daný den v promoakci - ukázalo se, že takové, až na výjimky nejsou.
- Součet tržeb všech ostatních produktů v kategorii.
- Součet tržeb všech ostatních produktů v kategorii, které byly v daný den v promoakci.
- Součet tržeb všech ostatních produktů v kategorii, které byly v daný den v promoakci nebo byly v rozmezí jednoho týdne po promoakci.

Ke každému ukazateli, jsem ještě vytvořila analogický ukazatel, který uvažoval zpoždění shrinku. V takovém ukazateli, se nebrala hodnota prodeje ze stejného dne, jako byl den záznamu shrinku, ale hodnota z předchozího dne. Důvodem pro vytvoření takových ukazatelů byla hypotéza, že shrink se může projevit až další den po uskutečněných tržbách. Důvodem může být to, že

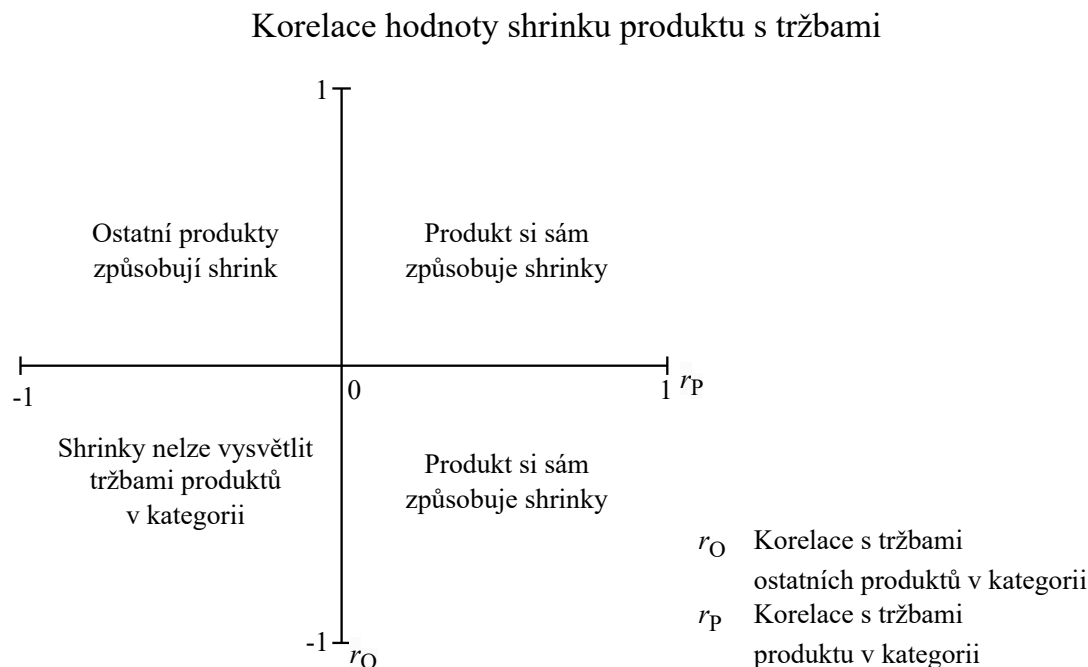
Na základě korelační analýzy je možné roztrždit produkty v kategorii do tří skupin:

Kategorie P - Produkty, které si samy způsobují shrink.

Kategorie O - Produkty, jejichž shrink je způsoben tím, že ostatní produkty v kategorii jsou v promoakci.

Kategorie X - Produkty, jejichž shrink se nepodařilo vysvětlit pomocí korelační analýzy.

Na obrázku 5.1 je znázorněno rozdělení produktů vzhledem ke korelačnímu koeficientu.



Obrázek 5.1: Kategorizace produktů podle korelace hodnoty shrinku produktu s tržbami.

Hypotéza pro zařazení do kategorie P je následující:

Pokud je korelace zaznamenaného shrinku s tržbami téhož produktu kladná, produkt si způsobuje shrinky sám. Abych mohla tuto hypotézu potvrdit, nebo vyvrátit, je třeba statisticky otestovat významnost korelačního koeficientu. Formulovala jsem nulovou hypotézu H_0 a alternativní hypotézu H_A pro koeficient r_P , který měří korelaci mezi hodnotou shrinku a tržbami produktu.

$$H_0 : r_P = 0 \quad \text{Výběry nejsou korelované.}$$

$$H_A : r_P \neq 0 \quad \text{Výběry jsou korelované.}$$

Hypotéza pro zařazení do kategorie O je následující:

Pokud jsou kladně korelované hodnoty zaznamenaného shrinku a tržby ostatních produktů a zároveň korelace shrinků produktu s vlastními tržbami je záporná, potom lze vyslovit hypotézu, že shrinky na produktu jsou způsobené ostatními produkty v promoakci. Pro toto tvrzení je opět nutné statisticky otestovat koeficienty korelace. Pro koeficient r_P je statistický test stejný jako v předchozím případě. Pro koeficient r_O měřící, jak jsou korelované shrinky a tržby ostatních produktů, je třeba otestovat následující hypotézy.

$$H_0 : r_O = 0 \quad \text{Výběry nejsou korelované.}$$

$$H_A : r_O \neq 0 \quad \text{Výběry jsou korelované.}$$

Pokud na zvolené hladině významnosti zamítneme nulovou hypotézu pro zkoumané korelační koeficienty, můžeme tvrdit že s danou pravděpodobností je koeficient statisticky významný. Na základě hodnoty korelace lze pak produkt zařadit do příslušné kategorie. Produkty, u kterých nelze zamítnout, není možné zařadit do tří uvedených kategorií.

Pro výpočet korelačního koeficientu je ještě třeba ověřit předpoklady. Pro Pearsonův korelační koeficient se jedná o předpoklad normality dat, shodnost rozptylů a nezávislost dat. Pro Spearmanův korelační koeficient jsou předpokládány nezávislé stejně rozdělené náhodné veličiny.

5.2 Implementace

V této části je uveden přesný postup pro získání kategorizace produktů. Kód je napsaný v jazyce Python. Součástí kódu je výběr kategorií, které jsou zkoumány, propojení dat shrinků, prodejů a promoakcí, výpočet korelace a ověření předpokladů, statistické testování a rozřazení produktů.

UML diagram pro analýzu.

5.2.1 Vstupy a výstupy

Pro korelační analýzu zaznamenaných shrinků s tržbami dalších produktů je třeba zajistit data, které se týkají zaznamenaných prodejů, produktů a prodejen. V následující části jsou popsány tabulková data, která jsou nezbytná pro správné spuštění analýzy. Dále jsou definované i vstupy, které musí definovat uživatel pro specifikování názvů konkrétních sloupců v souborech a parametry pro analýzu.

Celkem je požadováno pět vstupních tabulek - záznamy shrinků, záznamy prodejů, záznamy o promoakcích, číselník produktů s rozdělením produktové hierarchie. Tabulka se zaznamenanými shrinky musí obsahovat sloupec s datem záznamu, ID produktu, ID prodejny, hodnotu zaznamenaného shrinku. Tabulka s prodeji potřebuje stejné sloupce jako tabulka se shrinky s výjimkou že hodnota prodejů je celková prodaná částka, která byla zaznamenaná na dané prodejně v jeden den u daného produktu. Tabulka s údaji o promoakcích by měla obsahovat ID produktu, kterého se promoakce týká, začáteční a koncové datum promoakce a ID prodejny, pro kterou promoakce platí. Všechny záznamové tabulky musí pokrývat stejné časové období. Období může být libovolně dlouhé. Tabulka produktové hierarchie obsahuje ID produktu, jeho název a libovolně hluboký strom hierarchií. Každá úroveň stromu má vlastní sloupec. Všechny úrovně jsou vyplněné pro každý produkt, tato podmínka je nutná jen pro kategorie, které bude chtít uživatel využít při analýze. Tabulka s hierarchie produktů slouží k tomu, aby mohla být napojena na ostatní tabulky a data se pak mohla vyfiltrovat pouze na záznamy týkající se vybrané kategorie.

Před spuštěním hlavní výpočetní části musí uživatel vypsát konkrétní pojmenování sloupců v tabulce do proměnných. Sloupce, které v různých tabulkách označují tytéž hodnoty, musí mít stejný název. V následujícím kódu 5.1 je ukázka zadání.

V komentářích je slovní popis o jaký sloupec se jedná. Sloupec by však měl být jasný přímo z názvu proměnné.

```

01 product_col      = 'product_id'          # Product ID column
02 product_name_col = 'name'                # Product name
   column
03 whs_id_col       = 'warehouse_id'        # Store ID column
04 date_col         = 'date_of_transaction'  # Date of
   transactions column - for sales and shrinks tables
05 value_col_shrink = 'cost_value'          # Column with value
   of shrinks (shrink table)
06 value_col_sales  = 'cost_value'          # Column with value
   of total sales (sales table)
07 promo_col_from   = 'promotion_date_from' # Starting date of
   promotion (promotion table)
08 promo_col_to     = 'promotion_date_to'   # Starting date of
   promotion (promotion table)
09 categories       = ['L3', 'L4', 'L5', 'L6', 'name'] # Categories
   that we want to map to product ID (product hierarchy)

```

Kód 5.1: Definice konkrétních názvů sloupců.

Uživatel dále (viz ukázka kódu 5.2) zadefinuje formát data, který se používá v datumových sloupcích, aby se tyto sloupce mohly převést z textového řetězce na typ `datetime`. V proměnné `category_column` je třeba vybrat jednu kategorii (název sloupce). Na této úrovni se poté budou procházet jednotlivé kategorie, v rámci každé z nich se pak budou porovnávat a třídit produkty. V dalších proměnných může uživatel změnit umístění tj. název složky, kam se ukládají výsledky kategorizace a grafy. Složky s těmito názvy se vytvoří jako podsložky aktuální cesty.

```

01 date_format      = '%Y-%m-%d'            # Format of date
   columns
02 category_column  = 'L6'                  # Level od product
   hierarchy, on which to compare products (have to be in
   categories)

```

Kód 5.2: Definice parametrů.

5.2.2 Popis funkcí a struktura kódu

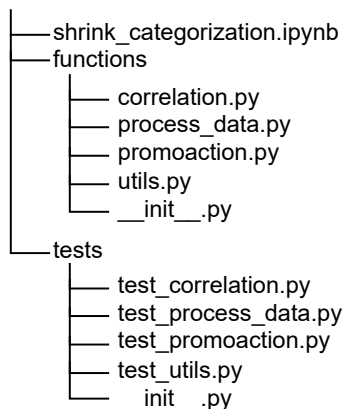
Kód pro korelační analýzu je umístěn ve složce `shrink_categorization`, struktura složky je vidět na obrázku 5.2.

Funkce jsou rozčleněny do modulů podle toho, na jaký výpočet jsou zaměřené. Každá funkce má je zdokumentovaná pomocí docstring obsaženého ve své definici. Dokumentace funkce se skládá ze stručného popisu, co funkce dělá, jaké má vstupní parametry a jaký je jejich význam a co funkce vrací. Funkce jsou otestované pomocí unit testů.

Pro práci s tabulkovými daty, které jsou hlavním vstupem, jsem použila balíček *pandas* jazyka Python.

V závorce je nástin toho, co budu popisovat (TBD)

- Výběr kategorie (Aby se nemuselo zadávat, přesné názvy kategorií systém vybere prvních n nejsilnějších kategorií z pohledu shrinků.)



Obrázek 5.2: Struktura souborů pro kód zpracovávající korelační analýzu.

- Propojení dat
- Korelace (Pearsonův, Spearmanův korelační koeficient) (ověření předpokladů (Kolmogorov-Smirnov test pro IID), testování statistické významnosti)
- Kategorizace (Nastavení prahu pro velikost korelačního koeficientu)
- Pomocné funkce

Funkce pro přiřazení kategorií k produktům

Jak je uvedeno na začátku sekce 5.2.1, uživatel musí specifikovat názvy sloupců kategorií, které bude v analýze používat. Seznam těchto kategorií je pak parametrem pro funkci `assign_levels`. Další parametry jsou `DataFrame`, kam se mají kategorie napojit a `DataFrame` odkud se kategorie napojují. Tyto `DataFrame`y musí mít společný sloupec, podle kterého se napojení provede. Defaultně se jedná o sloupec s ID produktu. Defaultně se provádí *left join*, aby nedošlo ke ztrátě dat, kdyby nějaký produkt neměl v `DataFrame`u kategorií zastoupení. Funkci je také možné předat další argumenty, které se dají volat ve funkci `merge` knihovny *pandas*. V analýze shrinků jeden řádek dat odpovídá transakci jednoho produktu, proto byl zvoleno ID produktu jako propojovací sloupec.

Funkce pro vytipování rizikových kategorií

Funkce `define_risk_categories` vybere prvních n kategorií v dané produktové hierarchii, kde suma hodnot v dané kategorii, je nejvyšší, resp. nejnižší. Funkce vrátí seznam těchto kategorií. Prvním vstupním parametrem je `DataFrame`, který obsahuje minimálně tři sloupce. Tyto sloupce je třeba definovat jako další parametry funkce. Jedná se o sloupec `value_column`, ve kterém jsou hodnoty, které ohodnocují řádky `DataFrame`u a kategorie. Další sloupec je jedna z úrovní produktové hierarchie, ve sloupci se nachází názvy, nebo jiné označení, kategorií. Posledním povinným parametrem je počet kategorií, které má funkce vrátit. Pokud je zadán tento počet tak, že je větší než je počet unikátních kategorií, vrátí se všechny kategorie seřazené

od nejrizikovější. Dale je funkci možné předat keyword argumenty, které se předají funkci `sort_values` z knihovny *pandas*. Jedná se např. o parametr pro vzestupné, nebo sestupné řazení. Defaultní řazení je vzestupné, což znamená, že se vezmou kategorie s nejnižší hodnotou. V této analýze sledujeme vyhozené množství, resp. peníze. Tento ukazatel je záporný, tedy vzestupné řazení vybere ty kategorie, jejichž ztráta byla nejvyšší. Vrácený seznam kategorií je tedy seřazen od nejrizikovější kategorie.

Funkce pro výběr pouze dané kategorie ze všech záznamů

Ve funkci `select_category` jsou vstupem `DataFrame`, název kategorie a úroveň, ve které se daná kategorie nachází. Funkce vrátí `DataFrame` pouze s těmi řádky, kde je obsažena jmenovaná kategorie. V případě, že tato kategorie v datech není, je vrácen prázdný `DataFrame`.

Tato funkce je volána ve funkci `process_dataframes`. Vstupy jsou totožné, avšak namísto jednoho `DataFrame` je možné jich zadat více jako samostatné parametry. Tělo funkce je v kódu 5.3. Funkce vrátí seznam všech vstupních `DataFrame`ů, a to pouze řádky, které obsahují zadanou kategorii.

```
01 def process_dataframes(category: str, category_column: str, *
02     dataframes) -> list[pd.DataFrame]:
03     result_dataframes = []
04
05     print('Sizes of dataframes: ')
06
07     for df in dataframes:
08         result_df = select_category(df, category,
09             category_column)
10         print("Shape of original df: " , df.shape, "Shape of new
11             df: " , result_df.shape)
12         result_dataframes.append(result_df)
13
14     return result_dataframes
```

Kód 5.3: Funkce pro výběr pouze dané kategorie z více `DataFrame`ů.

Funkce pro přiřazení promoakcí

V rámci korelační analýzy bylo potřeba přiřadit k jednotlivým zaznamenaným transakcím, zda byl produkt v den záznamu v promoakci nebo nikoli. V ideálním případě by tento příznak mohl být již uvedený u každého záznamu. Pokud tomu tak, ale není, je nutné data o promocích provázat na základě datumu, produktu a prod-jeny podle číselníku promoakcí. Data vybrané společnosti, na jejichž datech analýza probíhá, nemají promoakce přímo napojené na již proběhlé, zaznamenané transakce. Data o promoacích jsou uložena v číselníku promoakcí. Ten obsahuje ID produktu, prodejny, začátek a konec promoakce a prioritu promoakce. V jeden den může být více promoakcí, v takovém případě platí ta s nejvyšší prioritou.

Pomocné funkce

Funkce `create_folder` vytvoří složku se zadaným jménem v aktuální cestě, pouze pokud již taková složka neexistuje.

Pro základní vizualizaci jsem vytvořila funkci, která pomocí knihovny *matplotlib* vytváří bodový graf dvou proměnných. Graf je buď uložen nebo zobrazený při spuštění funkce. Funkci lze předat DataFrame a názvy dvou sloupců, které reprezentují vstupy pro osy x a y grafu. Další vstupy jsou názvy os a grafu, případně název souboru, pokud uživatel graf uložit.

5.2.3 Testování

Pro testování funkcí jsem použila knihovny *pytest* jazyka Python. Testy lze spustit příkazem `python -m pytest tests` v kořenovém adresáři projektu - `shrink_categorization`.

5.3 Výsledky

Zvolila jsem 5% hladinu významnosti pro testování statistické významnosti koeficientů korelace r_P a r_O .

Nejprve jsem se zaměřila na kategorie z úrovně 4, a to prvních deset kategorií s nejvyšší hodnotou shrinků (tj. s nejvyšší zaznamenanou ztrátou)

V tabulce č. 5.1 (*TBD: Bude tam víc kategorií L4.*) jsou porovnání výsledků kategorizace pro kategorii Masné výrobky a Slané pečivo. Měřila jsem postupně korelaci velikosti shrinku s různými ukazateli pro celkové tržby ostatních produktů. Pro určení míry korelace jsem zvolila Spearmanův korelační koeficient, jelikož data nesplňují předpoklady, které jsou nutné pro použití Pearsonova korelačního koeficientu - data nejsou nezávislá a stejně rozdělená. Z uvedených počtů produktů u jednotlivých kategoriích pro různé ukazatele, je patrné, že výsledky se příliš neliší. Pokud bychom se ale zaměřovali na celkové prodeje, nikoli promoční, tak získáváme větší množství produktů, u nichž nebylo možné vysvětlit shrink pomocí korelace.

Masné výrobky - pultový prodej

Shrink byl zaznamenaný u 111 produktů v této kategorii úrovně 4. 96 produktů bylo klasifikováno jako kategorie P, dva jako kategorie O, další dva jako X. U zbylých produktů nebyl koeficient korelace statisticky významný a proto nejde u těchto produktů vyslovit hypotézu pro jejich zařazení.

Produkty, které patří do kategorie O: Velikonoční klobása a Velikonoční česnekový šál - jedná se zcela jistě o sezónní výrobky. *TBD: vyjmenovat i další produkty, které si způsobují samy, ale je to většina produktů a salámy, šunky, z kategorie...*

Dále jsem zkoumala podkategorie Masných výrobků. Porovnávala jsem prodeje v rámci kategorií na šesté úrovni produktové hierarchie. V podkategorii Salámy s krátkou dobou spotřeby se kategorizace potvrdila. Pro kategorii, do níž patří

Tabulka 5.1: Počet produktů v kategoriích v závislosti na různých ukazatelích.
Kategorizace proběhla na základě Spearmanova korelačního koeficientu

Korelace s tržbami ostatních produktů						
Kategorie	Všechny prodeje		Prodeje v promoakci		Prodeje v a po promoakci	
	Stejný den	Další den	Stejný den	Další den	Stejný den	Další den
Masné výrobky - pultový prodej						
P	96	96	96	96	96	96
O	1	1	2	2	2	2
X	3	3	2	2	2	2
Nevýzn.	11	11	11	11	11	11
Slané pečivo						
P	131	131	131	131	131	131
O	9	13	24	24	24	24
X	16	12	1	1	1	1
Nevýzn.	90	90	90	90	90	90

sezónní výrobky - Netučné masné výrobky, nově z této podkategorie byl jako kategorie O označen i produkt Kladenská pečeně.

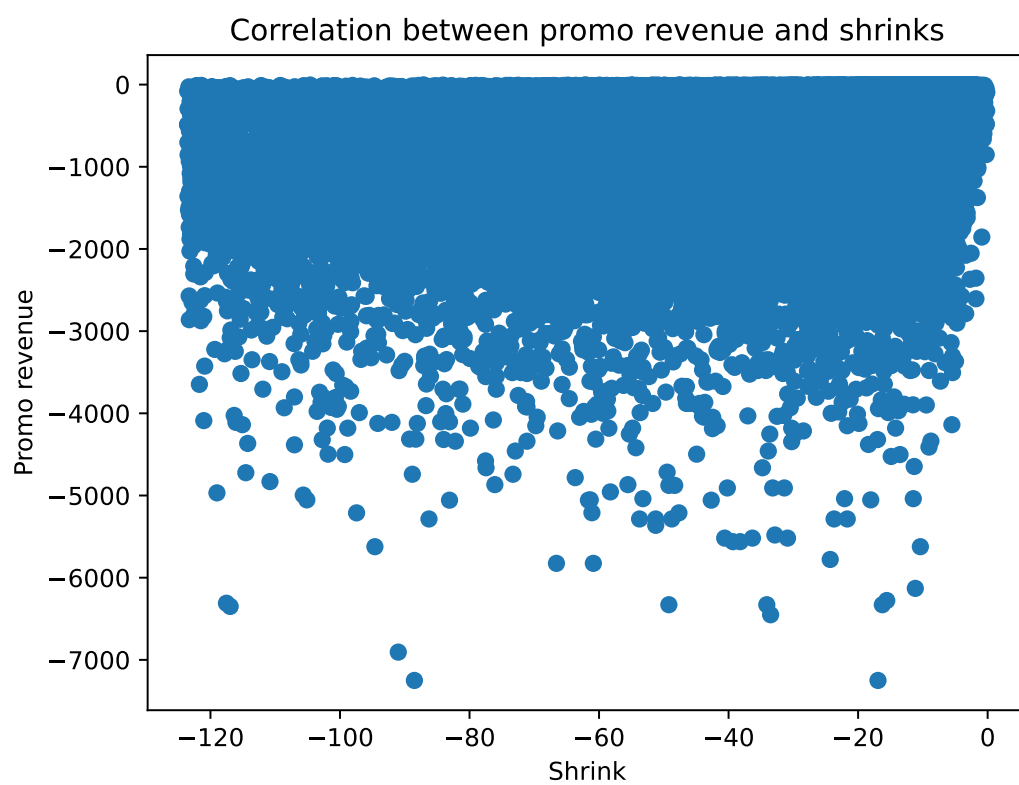
Slané pečivo

Produkty, které patří do kategorie O: pletýnka adélka, pletýnka malena, pletýnka sypaná mákem, pletýnka sypaná solí a km., bramborové pečivo s cibulí, veka cupeko kb, rohlík grahamový aspec, rohlík pivní ora, rohlík staročeský, rohlík obilný mam, rohlík na strouh.karlova, rohl.n str.penam, houska mašek, houska raženka, bageta rust.poh., bageta s grahamem, bageta chlebová, banketka cereální, dalamánek, kostka cereál malena, ciabatta mini natural, twistr se sýrem a špenatem, anglický rohlík,

Nevidím v této kategorii pravidlo v tom, co je v produktech, kterým shrink způsobují ostatní produkty, a které si za to mohou samy. TBD: názvy zobecním...

Sladké pečivo

čokorolka, závin mák , skořicový vrut, donut bílý se sušnkami, závitek cereální nugátový, rohlík lístový s ořech. náp, kobliha vanilková šiška, kobliha vanilková v koš, kobliha s jablky a skořicí, kobliha s lísko.náp.a pol, koblih kapsa s jablky péce, koláč wellartův, koláč tlač. vícezrnny, koláč šátek makový, koláč s ovocnou nápl, koláček švestkový , koláček meruňkový , koláček borůvkový, koláč rohový tři náplně , koláč s makovou náplní, koláč s tvarohovo nápl, máslový koláč tvaroh, šátek kyn. tvarohová nápl, šáteček s náplní višňovou , šáteček s tvaroh.náplní, šátek makový, loupák o., loupák v., loupák m., závin tvaroh , závin kvásk. makový, bavorská hvězdice, makovka pletená mašek, croissant máslový



Obrázek 5.3: Závislost mezi tržbami produktu a tržbami ostatních produktů v kategorii během promoakce (Masné výrobky - pultový prodej).

Plodová zelenina

Shrink byl zaznamenaný u 28 produktů v této kategorii. 20 produktů bylo klasifikováno jako kategorie P, u zbylých produktů nebyl koeficient korelace statisticky významný a proto nejde u těchto produktů vyslovit hypotézu pro jejich zařazení.

Závěr

Cílem práce bylo

Literatura

- [1] BAUDIN, Michel. *Lean Logistics: The Nuts and Bolts of Delivering Materials and Goods*. New York: Productivity Press, 2005. ISBN 978-1563272967.
- [2] CHRISTOPHER, Martin. *Logistics & Supply Chain Management*. 5th ed. Harlow: Pearson Education Limited, 2016. ISBN 9781292083797.
- [3] What is the difference between Logistics and Supply Chain Management. In: *IIM Udaipur Chronicles* [online]. 11. 10. 2019. [cit. 2022-11-07] Dostupné z: <https://www.iimu.ac.in/blog/what-is-the-difference-between-logistics-and-supply-chain-management/>
- [4] JIRSÁK, Petr, MERVART, Michal, VINŠ, Marek. *Logistika pro ekonomy – vstupní logistika*. 1. vydání. Praha: Wolters Kluwer ČR, 2012.
- [5] JONES, Daniel T., Peter HINES a Nick RICH. Lean logistics. *International Journal of Physical Distribution & Logistics Management*. 1997, **27**(3/4), 153-173. ISSN 0960-0035. Dostupné z: doi:10.1108/09600039710170557
- [6] KURITA, Takio. *Principal component analysis (PCA). Computer Vision: A Reference Guide*. 2019, 1-4. [cit. 2022-11-07] Dostupné z: https://link.springer.com/content/pdf/10.1007/978-3-030-03243-2_649-1.pdf
- [7] Analýza hlavních komponent. In: *Wikipedia: the free encyclopedia* [online]. San Francisco (CA): Wikimedia Foundation, 2022 [cit. 2023-03-04]. Dostupné z: https://cs.wikipedia.org/wiki/Analýza_hlavních_komponent
- [8] DI FRANCO, Giovanni. *Multiple correspondence analysis: one only or several techniques?*. *Quality & Quantity*, 2016, 50.3: 1299-1315. [cit. 2023-03-05]. Dostupné z: doi:10.1007/s11135-015-0206-0
- [9] ABDI, Hervé, VALENTIN, Dominique. *Multiple correspondence analysis*. In: *Encyclopedia of measurement and statistics*. 2007, 2.4: 651-657. [cit. 2023-03-05]. Dostupné z: <https://personal.utdallas.edu/~Herve/Abdi-MCA2007-pretty.pdf>
- [10] GREENACRE, Michael. *Correspondence analysis in practice*. Chapman and Hall/CRC, 2017. [cit. 2023-03-05].
- [11] Correspondence analysis. In *Wikiwand* [online], [cit. 2023-03-06]. Dostupné z: https://www.wikiwand.com/en/Correspondence_analysis

- [12] WRONKA, Anna. LEAN LOGISTICS. *Journal of Positive Management*. 2017, 7(2), 55-63. ISSN 2392-1412. Dostupné z: doi:10.12775/JPM.2016.012
- [13] SUTHERLAND Joel, BENNETT Bob. *The Seven Deadly Wastes of Logistics: Applying Toyota Production System Principles to Create Logistics Value*. Bethlehem, PA: Lehigh University, 2007. Dostupné z: <https://www.researchgate.net/publication/265356600>
- [14] SKHMOT, Nawras. *The Lean Way Blog: The 8 Wastes of Lean*. *The Lean Way* [online]. 5. 8. 2017 [cit. 2022-11-17]. Dostupné z: <https://theleanway.net/The-8-Wastes-of-Lean>
- [15] SKHMOT, Nawras. *The Lean Way Blog: What is Lean?*. *The Lean Way* [online]. 5. 8. 2017 [cit. 2022-11-17]. Dostupné z: <https://theleanway.net/what-is-lean>
- [16] SKHMOT, Nawras. *The Lean Way Blog: What is Muda, Mura, and Muri?*. *The Lean Way* [online]. 5. 8. 2017 [cit. 2022-11-17]. Dostupné z: <https://theleanway.net/muda-mura-muri>
- [17] Learning the Lingo: 3 definitions related to unsold food inventory. In: *Blog - Spoiler Alert* [online]. 17. 06. 2019. [cit. 2022-02-07] Dostupné z: <https://blog.spoileralert.com/3-definitions-unsold-food-inventory>
- [18] HALFORD, M. Prince., *Prince* [online]. [cit. 2023-03-08]. Dostupné z: <https://github.com/MaxHalford/prince>
- [19] BAIJAYANTA, Roy., *All about Categorical Variable Encoding* [online]. [cit. 2023-03-13]. Dostupné z: <https://towardsdatascience.com/all-about-categorical-variable-encoding-305f3361fd02>
- [20] scikit-learn Machine Learning in Python [online]. 2023 [cit. 2023-03-21]. Dostupné z: <https://scikit-learn.org/stable/>
- [21] Multiclass and multioutput algorithms. scikit [online]. [cit. 2023-04-10]. Dostupné z: <https://scikit-learn.org/stable/modules/multiclass.html>
- [22] BOBBITT, Zach. *How to Interpret Cramer's V* In: *Statology* [online]. 2021 [cit. 2023-03-21]. Dostupné z: <https://www.statology.org/interpret-cramers-v/>
- [23] ZYCHLINSKI, Shaked. *The Search for Categorical Correlation* [online]. 2018 [cit. 2023-03-14]. Dostupné z: <https://towardsdatascience.com/the-search-for-categorical-correlation-a1cf7f1888c9>
- [24] HOLČÍK, Jiří, KOMENDA, Martin (eds.) a kol. *Matematická biologie: e-learningová učebnice* [online]. 1. vydání. Brno: Masarykova univerzita, 2015. [cit. 2023-03-14]. ISBN 978-80-210-8095-9.
- [25] FROST, Jim. *Choosing the correct type of regression analysis*. [online]. [cit. 2023-04-10]. Dostupné z: <https://statisticsbyjim.com/regression/choosing-regression-analysis/>

- [26] Multiclass logistic regression. In *Refactored* [online]. [cit. 2023-04-10]. Dostupné z: <https://refactored.ai/>
- [27] BIAU, Gérard; SCORNET, Erwan. *A random forest guided tour*. In: *Test* [online]. 2016, 25: 197-227. [cit. 2023-04-10].
- [28] LOK, Leon. *Decision trees, random forests and gradient boosting: What's the difference?* [online]. 5. 1. 2022. [cit. 2023-04-10]. Dostupné z: <https://leonlok.co.uk/blog/decision-trees-random-forests-gradient-boosting-whats-the-difference/>

Příloha A

Obsah přiloženého CD

BP_Gruberova.pdf – soubor s elektronickou verzí této bakalářské práce.

PreprocessData – složka se zdrojovými kódy naimplementovaného balíku
PreprocessData.jl.