

**Masarykova univerzita
Přírodovědecká fakulta**

Data mining I

Projekt – Internet

Tomancová Nikola, 484549
Finanční a pojistná matematika
22.05.2021

Úvod

Tento semestrální projekt pracuje s datovým souborem Internet, který vznikl na základě výzkumu populace České republiky. Průzkum proběhl přibližně na přelomu 20. a 21. století a zaměřil se na okolnosti, které by mohly ovlivnit například množství času stráveného na internetu nebo způsoby, jakým obvykle dotázaní uživatelé tyto chvíle trávili, tedy jaké internetové služby nejčastěji využívali.

Mým úkolem bylo na základě těchto dostupných dat zjistit, jaká byla situace v ČR z hlediska využívání internetu právě na přelomu 20. a 21. století. Dále se můj projekt věnuje hledání v tehdejší době nejčastěji využívaných služeb na internetu, a nakonec zjišťování služeb, které byly obvykle využívány současně.

1 Příprava dat

Seznámení s daty

Data Internet obsahují celkem 32 proměnných a 2 889 případů, jedná se o reprezentativní vzorek občanů České republiky, kteří byli v té době v produktivním věku (tj. 15 – 65 let). Proměnné můžeme rozdělit do dvou skupin, první jsou proměnné demografické, které se týkají místa bydliště respondentů, jejich pohlaví, věku, vzdělání nebo například jejich finanční situace. Podrobněji se popisu jednotlivých demografických proměnných věnuje Tabulka 1.

Ve druhé skupině figurují především dummy proměnné, tedy proměnné přiřazující daným pozorováním pouze hodnoty 1 nebo 0. Většina otázek se týkala využívání základních internetových služeb, a proto možné odpovědi na takové otázky byly jen „ano“ (zakódované v datech jako hodnota 1) nebo „ne“ (hodnota 0 v našich datech). Za dummy proměnnými pak následují ještě tři proměnné týkající se času stráveného na internetu. Nejprve je to ordinální proměnná s kategoriemi 1 až 6 a za ní následují dvě numerické proměnné udávající množství času v minutách. Tuto skupinu proměnných budeme v tomto projektu nazývat „internetové proměnné“ a jsou shrnuty v Tabulce 2.

Vidíme, že proměnné v našem datovém souboru jsou převážně kvalitativní, jediné kvantitativní jsou proměnné *id*, *hmem*, *q75_1* a *q78_1*. Vzhledem k tomu, že ne všechny proměnné byly v datovém souboru přesně specifikovány, bylo nutné si popisy některých proměnných ještě podrobněji dodefinovat. Především se to týkalo proměnných *q75_1* a *q78_1* označující počet minut strávených na internetu. Budeme předpokládat, že v případě proměnné *q75_1* měli respondenti sdělit počet minut strávených na internetu průměrně ve všední dny a naopak u proměnné *q78_1* udávali průměrný počet minut strávených na internetu během víkendového dne.

Vzhledem k tomu, že ordinální proměnné *agecat* a *q40_23* (označující věkové kategorie a množství času stráveného surfováním po internetu) budeme v následujících kapitolách ještě podrobněji rozebírat, bude vhodné zde uvést ještě tabulky vysvětlující význam jednotlivých kategorií těchto proměnných, viz Tabulky 3 a 4.

Název	Popis	Typ
id	Pořadí respondenta	Numerická
reg	Kraj (Středočeský apod.)	Nominální
vb	Velikost města	Ordinální
hinc	Celkový čistý příjem v domácnosti	Ordinální
hmem	Počet členů v domácnosti	Numerická
hint	Jak často jsou připojeni k internetu	Ordinální
fixline	Pevná linka v domácnosti	Nominální
hcar	Počet aut v domácnosti	Ordinální
sex	Pohlaví	Nominální
agecat	Věkové kategorie	Ordinální
marits	Rodinný stav	Nominální
red	Vzdělání	Ordinální
empst	Zaměstnanecké postavení	Nominální
hhead	Hlava domácnosti	Nominální
pinc	Celkový čistý příjem na osobu	Ordinální
sim	Přítomnost SIM karet v domácnosti	Nominální
seg_ls	Typ osobnosti	Nominální

Tabulka 1: Demografické proměnné

Čištění dat, transformace proměnných

Ještě před samotnou analýzou dat bylo nutné data zkontrolovat a ujistit se, že proměnné neobsahují žádné zřejmě chybné hodnoty. V případě kategoriálních proměnných nebyl z tohoto hlediska nalezen žádný problém. Jak je znázorněno v tabulce na Obrázku 1, po vypsání základních popisných statistik těchto proměnných se ukázalo, že všechny kategoriální proměnné se pohybují ve svém oboru hodnot. Jediný problém u kategoriálních veličin by mohl být v proměnné *pinc* znázorňující osobní měsíční příjem, jelikož obsahuje velmi vysoký počet chybějících pozorování. Ovšem většina respondentů, kteří na otázku týkající se osobního příjmu neodpověděli, patří do věkové kategorie 1, případně 2, tzn. jsou ve věku od 15 do 24 let. Je proto velmi pravděpodobné, že se jedná o studenty bez vlastního příjmu a bylo by možné tyto chybějící hodnoty nahradit nulou.

Tabulka dále ukázala podezřelé hodnoty u numerických proměnných *q75_1* a *q78_1*. Maximum těchto proměnných je totiž v obou případech 1 440 minut, což se rovná v přepočtu 24 hodinám. Jedná se tudíž evidentně o nesprávné hodnoty, které by ale mohly náš další výzkum významně ovlivnit. Z toho důvodu byla tato pozorování vyjmuta z datového souboru, stejně jako pozorování s dalšími vysokými hodnotami, které se objevily v tabulkách na

Název	Popis	Typ
q74_1	Elektronická pošta (e-mail)	Nominální
q74_2	Hledání informací, kontaktů (obecně)	Nominální
q74_3	Vyhledávání, objednání nebo nákup zboží či služeb	Nominální
q74_4	Hraní her na internetu (on-line)	Nominální
q74_5	Obsluha bankovního konta z domova	Nominální
q74_6	Telefonování přes internet	Nominální
q74_7	Chat, diskusní skupiny	Nominální
q74_8	Využívání služeb okamžité komunikace (ICQ apod.)	Nominální
q74_9	Poslouchání hudby nebo sledování hud. klipů přes internet	Nominální
q65_13	Internetové připojení	Nominální
q65_17	Stolní počítač	Nominální
q65_18	Notebook	Nominální
q40_23	Surfování na internetu	Ordinální
q75_1	Čas na internetu – denně	Numerická
q78_1	Čas na internetu – o víkendu	Numerická

Tabulka 2: Proměnné - internet

Obrázku 2. Celkem bylo odstraněno 20 pozorování.

S ohledem na analýzu v následující kapitole byla vhodná ještě jednoduchá transformace části „internetových“ proměnných. Byla vytvořena nová proměnná *internet*, která vznikla jako součet hodnot u proměnných q74_1 až q74_9. Jinými slovy, proměnná *internet* označuje počet internetových služeb, které daný respondent využívá, jedná se tedy o ordinální proměnnou, která nabývá hodnot od 0 do 9. Relativní četnosti této proměnné znázorňuje sloupcový graf na Obrázku 3.

2 Češi a internet

Nyní se zaměříme na některé zajímavé vztahy mezi proměnnými. V této kapitole nás bude zajímat vztah některých demografických proměnných k využívání internetových služeb a množství času stráveného na internetu.

Nechali jsme si vypsát k jednotlivým proměnným vždy tři další nejvíce korelované proměnné. Z výsledků bylo patrné, že s „internetovými“ proměnnými je ve vyšší míře korelovaná nejčastěji demografická proměnná *agecat*, největší závislost má věk na poslech hudby, chatování, surfování po internetu a trávení času na internetu během víkendu. Nicméně takový výsledek, že způsob a množství využívání internetu závisí na věku uživatele, byl vcelku

Proměnná <i>agecat</i>	
Kategorie	Věk od - do
1	15 - 20
2	21 - 24
3	25 - 30
4	31 - 40
5	41 - 50
6	51 - 65

Tabulka 3: Proměnná *agecat*

Proměnná <i>q40_23</i>	
Kategorie	Frekvence surfování
1	denně nebo téměř denně
2	přibližně dvakrát až třikrát týdně
3	přibližně jednou týdně
4	přibližně jednou až třikrát měsíčně
5	méně často
6	nikdy nebo téměř nikdy

Tabulka 4: Proměnná *q40_23*

očekávatelný. Naopak poměrně zajímavá je korelace mezi rodinným stavem a poslechem hudby na internetu nebo například mezi vzděláním a využíváním internetového bankovníctví, vyhledávání informací na internetu nebo vlastnictvím notebooku. Ostatní demografické proměnné už s těmi „internetovými“ příliš velkou korelaci neměly.

Jako názorná ukázka korelace mezi věkem a využíváním internetu slouží tabulka na Obrázku 4. Lze pozorovat, že s přibývajícím věkem klesá množství využívaných služeb, množství času stráveného surfováním na internetu i čas strávený na internetu během víkendů.

Následující tabulka (Obrázek 5) pak ještě přibližuje vztah proměnných *agecat* a *internet*. Můžeme vidět, že skupina lidí nevyužívající žádné z nabízených internetových služeb patří nejčastěji do nejstarší věkové kategorie, jsou tedy ve věku 51 až 65 let. Naopak všechny služby, ze kterých bylo možné vybírat, využívá nejvíce druhá věková kategorie (věk 21 až 24 let). Jsou to totiž lidé, kteří jsou sice dost mladí na to, aby na internetu chatovali, poslouchali hudbu či hráli hry, ale zároveň i dostatečně staří na to, aby již byli samostatní a měli tedy například i své vlastní internetové bankovníctví (na rozdíl od respondentů ve věku 15 až 20 let).

Podrobněji se na korelace mezi některými proměnnými zaměřuje ještě korelogram na Obrázku 6. Korelogram se dělí na 2 části – horní a dolní trojúhelník. V horním trojúhelníku se nacházejí korelační koeficienty proměnných (některé z nich jsou vypsány už i v tabulce na Obrázku 4), vypsány modrou barvou v případě kladné korelace a červenou v případě korelace záporné. Míra korelace je pak naznačena i v dolním trojúhelníku korelogramu pomocí

The MEANS Procedure

Variable	Label	N	N Miss	Minimum	Maximum
id		2889	0	1.0000000	2889.00
reg		2889	0	11.0000000	81.0000000
vb		2889	0	1.0000000	6.0000000
hinc		2887	2	11.0000000	31.0000000
hmem		2885	4	1.0000000	11.0000000
hint		2864	25	1.0000000	5.0000000
fixline		2887	2	1.0000000	2.0000000
hcar		2887	2	1.0000000	3.0000000
sex		2889	0	1.0000000	2.0000000
agecat		2889	0	1.0000000	6.0000000
marits		2889	0	1.0000000	5.0000000
red		2889	0	1.0000000	4.0000000
empst		2889	0	1.0000000	9.0000000
hhead		2889	0	1.0000000	2.0000000
pinc		2403	486	11.0000000	31.0000000
sim		2889	0	1.0000000	2.0000000
seg_ls		2889	0	1.0000000	8.0000000
q74_1	q74#1	2889	0	0	1.0000000
q74_2	q74#2	2889	0	0	1.0000000
q74_3	q74#3	2889	0	0	1.0000000
q74_4	q74#4	2889	0	0	1.0000000
q74_5	q74#5	2889	0	0	1.0000000
q74_6	q74#6	2889	0	0	1.0000000
q74_7	q74#7	2889	0	0	1.0000000
q74_8	q74#8	2889	0	0	1.0000000
q74_9	q74#9	2889	0	0	1.0000000
q65_13	q65#13	2889	0	0	1.0000000
q65_17	q65#17	2889	0	0	1.0000000
q65_18	q65#18	2889	0	0	1.0000000
q40_23		2889	0	1.0000000	6.0000000
q75_1		2889	0	0	1440.00
q78_1		2889	0	0	1440.00

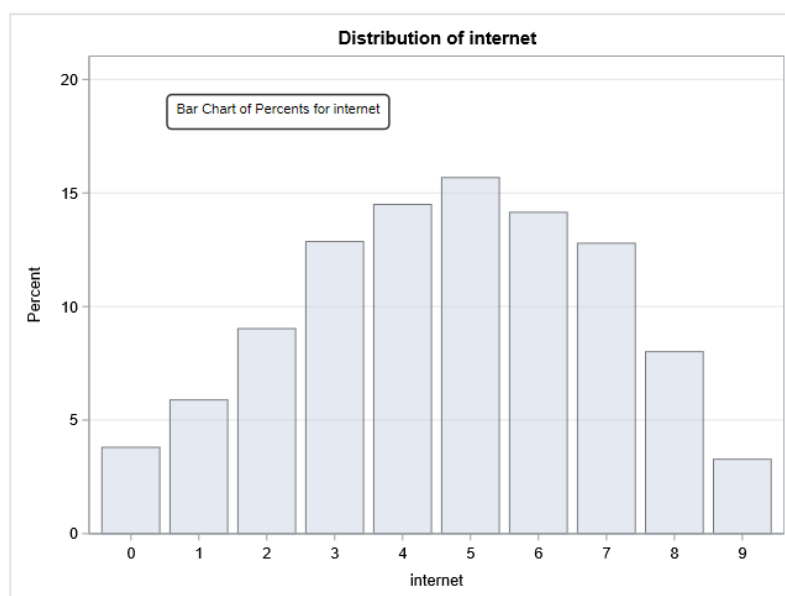
Obrázek 1: Chybná pozorování, minima a maxima

odstínu barev (čím tmavší odstín, tím silnější korelace, tj. čím vyšší korelační koeficient v absolutní hodnotě).

Korelogram potvrdil naše předchozí závěry týkající se proměnné *agecat*. Čím byli respondenti starší, tím méně času trávili u internetu (platí pro všední dny i pro víkendy) a tím méně internetových služeb používali. U proměnné *red* se ukázalo, že úroveň maximálního dosaženého vzdělání nemá příliš velký vliv na celkové využívání internetu. Množství služeb má samozřejmě kladnou korelaci s časem stráveným u internetu. Nakonec jsme se pokusili porovnat čas, který respondenti tráví na internetu ve všední dny a o víkendu, protože by mohlo platit například pravidlo, že ti, kteří tráví na internetu svůj čas v práci během všedních dnů, si pak chtějí o víkendu od počítačů naopak odpočinout a chodí na něj méně. Výsledky korelogramu ovšem tuto hypotézu popírají – platí naopak pravidlo, že respondenti, kteří tráví

Extreme Values					
Lowest			Highest		
Order	Value	Freq	Order	Value	Freq
1	0	136	78	1000	1
2	1	8	79	1020	1
3	2	14	80	1200	2
4	3	13	81	1400	1
5	4	7	82	1440	2

Extreme Values					
Lowest			Highest		
Order	Value	Freq	Order	Value	Freq
1	0	283	68	1000	4
2	1	7	69	1020	1
3	2	16	70	1200	7
4	3	6	71	1400	2
5	4	4	72	1440	5

Obrázek 2: Extrémní hodnoty proměnných *q75_1* a *q78_1*Obrázek 3: Relativní četnosti proměnné *internet*

více času na internetu v týdnu, pak tráví na internetu více času i o víkendu. Toto pravidlo dokonce platí napříč všemi věkovými kategoriemi, i když bychom vzali jednotlivé věkové kategorie zvlášť, korelační koeficient u proměnných *q75_1* a *q78_1* vychází vždy vyšší než 0,5.

3 Analýza dat

V této kapitole se budeme zabývat dalšími otázkami, které se týkají datového souboru Internet. Budeme zkoumat, jaké služby nejčastěji lidé využívají nebo které služby jsou používány současně. Budeme tedy hledat nějaká nejčastější očekávaná i překvapivá pravidla ve využívání jednotlivých internetových služeb.

Právě pro tyto účely je určená analýza nákupního košíku. Tato metoda se totiž snaží najít

Pearson Correlation Coefficients, N = 2869				
	agecat	internet	q40_23	q78_1
agecat	1.00000	-0.39981	0.42249	-0.32961
internet	-0.39981	1.00000	-0.53388	0.41103
q40_23	0.42249	-0.53388	1.00000	-0.37723
q78_1	-0.32961	0.41103	-0.37723	1.00000

Obrázek 4: Korelace

mezi daty nějaké často se vyskytující vztahy a objevit v nich tímto způsobem zajímavé vzory či pravidla. V případě našich dat to znamená, že se pokusíme naleznout určité vzory chování uživatelů internetu, zjistit, které služby jsou obvykle využívány současně a které naopak nikoliv.

Transakční data

Aby bylo možné provést analýzu nákupního košíku, bude nejprve nutná úprava původních dat. Budeme nyní brát v potaz pouze „internetové“ dummy proměnné, vytvoříme proto nejprve tabulku tvořenou pouze vyhovujícími proměnnými.

Ani nově vzniklá tabulka ovšem stále není dostačující. Dále je nutné tato „raw“ data převést na transakční. Transakční data obsahují pouze dvě proměnné, jednou z nich je proměnná *ID* označující jednotlivé respondenty, zákazníky apod. Druhá proměnná bývá obvykle pojmenována jako Položka (*Item*) nebo například Transakce (*Transaction*). Jak již název napovídá, do této proměnné uložíme internetové služby, které daný respondent používá.

Četnost služeb

Nově vzniklá data využijeme nejprve k tomu, abychom zjistili, jaké služby jsou využívány nejčastěji. Jak ukazuje graf relativních četností na Obrázku 7, nejrozšířenější internetovou službou byl e-mail, avšak poměrně často byl internet využíván samozřejmě i za účelem vyhledávání informací. Třetí největší zastoupení mezi službami mělo nakupování přes internet. Tabulka relativních četností tedy nepřinesla příliš překvapivé výsledky.

Pro zajímavost ještě uvádíme graf relativních četností pro vlastnictví počítače či notebooku, viz Obrázek 8. Notebooky výrazně převyšují stolní počítače. Lze předpokládat, že vzhledem k technologickému pokroku a potřeby počítačů k práci či výuce z domova by toto porovnání v dnešní době již vyšlo velmi odlišně.

Analýza nákupního košíku

Nyní bude následovat přímo samotná analýza nákupního košíku prováděná v programu SAS Enterprise Miner. Hladinu podpory (*support*) jsme stanovili na 15 % a s mírou spolehlivosti (*confidence*) jsme se pohybovali v hodnotách 80 % a více. Jako maximální počet položek v

The FREQ Procedure

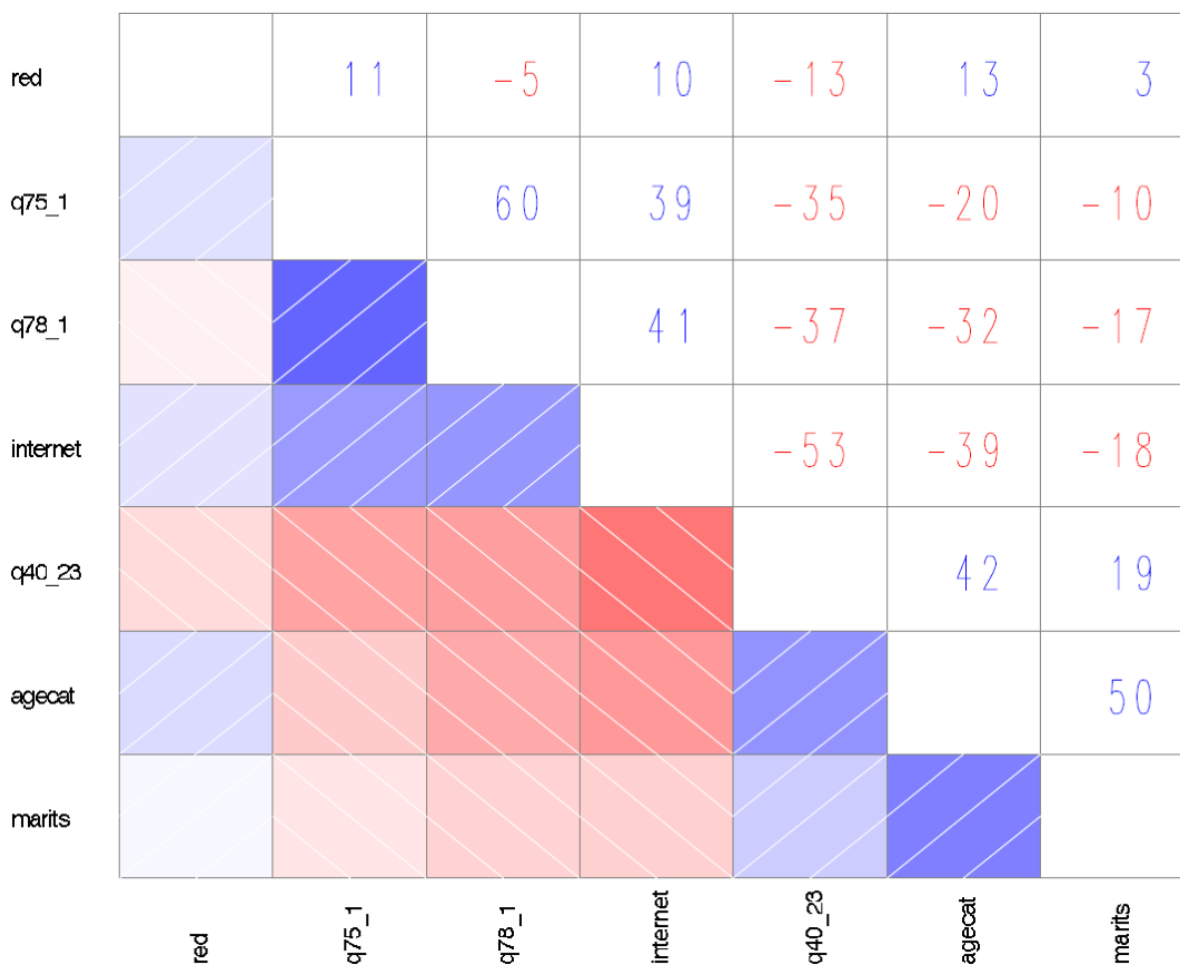
Table of agecat by internet											
agecat	internet										Total
	0	1	2	3	4	5	6	7	8	9	
1	1	11	18	25	42	63	52	96	67	13	388
	0.03	0.38	0.63	0.87	1.46	2.20	1.81	3.35	2.34	0.45	13.52
	0.26	2.84	4.64	6.44	10.82	16.24	13.40	24.74	17.27	3.35	
	0.92	6.51	6.95	6.78	10.10	14.00	12.81	26.16	29.13	13.83	
2	4	14	16	42	41	57	74	63	55	40	406
	0.14	0.49	0.56	1.46	1.43	1.99	2.58	2.20	1.92	1.39	14.15
	0.99	3.45	3.94	10.34	10.10	14.04	18.23	15.52	13.55	9.85	
	3.67	8.28	6.18	11.38	9.86	12.67	18.23	17.17	23.91	42.55	
3	2	6	14	44	40	72	68	59	41	15	361
	0.07	0.21	0.49	1.53	1.39	2.51	2.37	2.06	1.43	0.52	12.58
	0.55	1.66	3.88	12.19	11.08	19.94	18.84	16.34	11.36	4.16	
	1.83	3.55	5.41	11.92	9.62	16.00	16.75	16.08	17.83	15.96	
4	3	17	40	64	96	95	80	76	33	11	515
	0.10	0.59	1.39	2.23	3.35	3.31	2.79	2.65	1.15	0.38	17.95
	0.58	3.30	7.77	12.43	18.64	18.45	15.53	14.76	6.41	2.14	
	2.75	10.06	15.44	17.34	23.08	21.11	19.70	20.71	14.35	11.70	
5	21	47	76	81	83	71	64	43	18	10	514
	0.73	1.64	2.65	2.82	2.89	2.47	2.23	1.50	0.63	0.35	17.92
	4.09	9.14	14.79	15.76	16.15	13.81	12.45	8.37	3.50	1.95	
	19.27	27.81	29.34	21.95	19.95	15.78	15.76	11.72	7.83	10.64	
6	78	74	95	113	114	92	68	30	16	5	685
	2.72	2.58	3.31	3.94	3.97	3.21	2.37	1.05	0.56	0.17	23.88
	11.39	10.80	13.87	16.50	16.64	13.43	9.93	4.38	2.34	0.73	
	71.56	43.79	36.68	30.62	27.40	20.44	16.75	8.17	6.96	5.32	
Total	109	169	259	369	416	450	406	367	230	94	2869
	3.80	5.89	9.03	12.86	14.50	15.68	14.15	12.79	8.02	3.28	100.00

Obrázek 5: Internet a věk

jednom pravidle byla nejprve nastavena hodnota 2. Zahrnutí proměnných *q65_13*, *q65_17* a *q65_18* žádné zajímavé závěry nepřineslo, rozhodli jsme se je proto z dat určených k analýze odebrat a nadále ponechat pouze položky týkající se služeb samotných.

Analýza dvouprvkových pravidel především potvrdila důležitost e-mailu a vyhledávání na internetu, protože všechna pravidla nacházející se v grafu spolehlivosti a podpory vpravo (tzn. ty s největší spolehlivostí) byla vytvořena ve formě „nějaká služba \Rightarrow e-mail“ (resp. „služba \Rightarrow hledání“ nebo „služba \Rightarrow nákupy“). Vzhledem k tomu, že e-mail, vyhledávání a nákupy byly nejčastěji používanými službami (viz Obrázek 7), tak není příliš překvapivé, že měla právě pravidla s těmito službami tak velkou spolehlivost. Pro podrobnější výsledky uvádíme na Obrázku 9 několik prvních řádků z Tabulky pravidel (*Rules Table*) seřazených podle spolehlivosti.

Při seřazení tabulky podle sloupce *Lift* se na první příčku dostala dvojice hudba a hraní her (služby se spolu vyskytovali současně 1,65x častěji oproti náhodě). Může to být způsobeno například tím, že hráči online her si často pustí hudbu přímo k hraní jako kulisu. Po seřazení

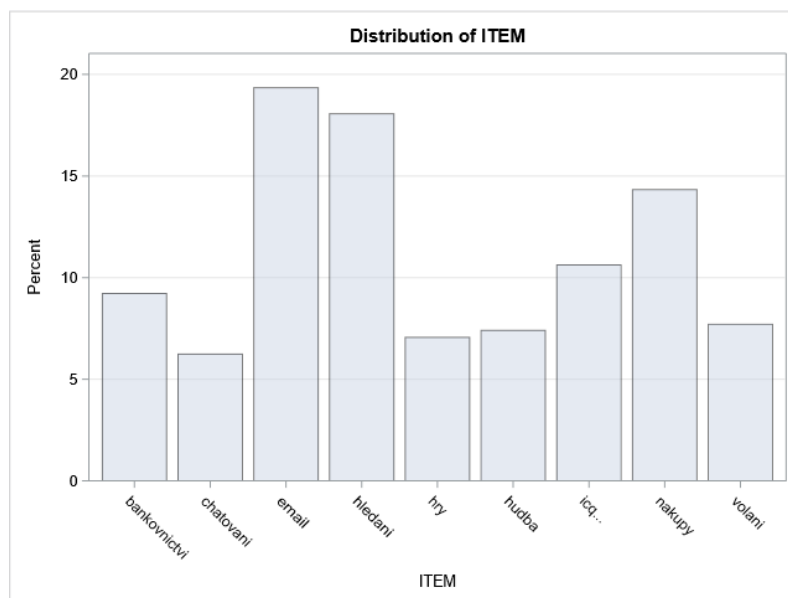


Obrázek 6: Korelogram vybraných proměnných

tabulky podle podpory lze vidět, jaké služby jsou nejčastěji využívány současně. Takovými dvojicemi služeb byly především kombinace služeb s největší relativní četností (např. email a vyhledávání: 83,7 %, nákupy a email 67,54 % apod.).

Pro dokreslení uvádíme ještě Link Graph s orientovanými šipkami (Obrázek 10). I z něj jsou patrná ta nejsilnější pravidla (největší podporu mají ta největší a nejčervenější kolečka, sílu pravidel pak znázorňuje tloušťka a barva šipek).

Dále jsme vyzkoušeli variantu analýzy pro až čtyřprvková pravidla. Výsledky této analýzy zpracované graficky se nachází na Obrázku 12. Čtverečky grafu jsou u vysoké spolehlivosti v naprosté většině ve formě „služba + služba + služba \Rightarrow e-mail“, například „hudba + nákupy + bankovníctví \Rightarrow e-mail“ (podpora 17,13 %, spolehlivost 99,79 %). Pokud tedy respondent využívá některou ze zjištěných trojic služeb, pak téměř se stoprocentní jistotou lze říci, že bude mít i e-mail. Po kombinacích se službou e-mail následují obdobná pravidla, ve kterých

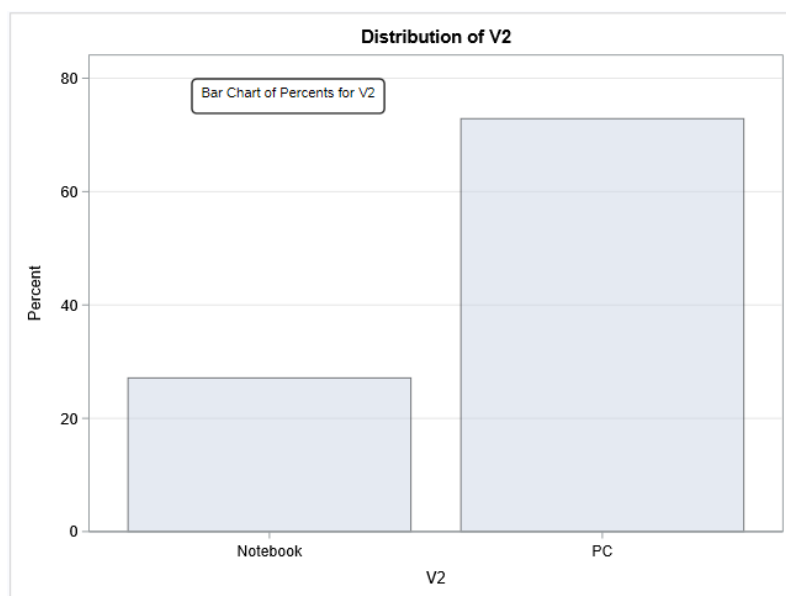


Obrázek 7: Relativní četnosti jednotlivých služeb

figuruje na pravé straně vyhledávání nebo vyhledávání + e-mail (resp. nákupy, nákupy + e-mail nebo nákupy + hledání), poté se na pravou stranu přidávají i služby okamžité komunikace (dále IM podle anglického překladu *Instant Messaging*) atd. V případě IM se už ale spolehlivost pohybuje v nižších hodnotách, mezi 80 až 90 %, například „hudba + hledání + chatování \Rightarrow IM“ (podpora 14,86 %, spolehlivost 87,5 %), „hudba + chatování \Rightarrow IM + e-mail“ (podpora 15,11 %, spolehlivost 85,37 %).

Na zajímavý bod, který ukázala matice pravidel (*Rule Matrix*) poukazuje Obrázek 11. Jedná se o pravidlo „nákupy + hledání + e-mail \Rightarrow bankovníctví“ (spolehlivost 80,51 %), které je svou strukturou velmi odlišné od ostatních okolních pravidel. Vyplývá z něj vztah mezi nakupováním po internetu a vlastnictvím internetového bankovníctví. Tato vazba je poměrně logická, jelikož nákupy přes internet jsou často hrazeny online formou právě prostřednictvím internetového bankovníctví.

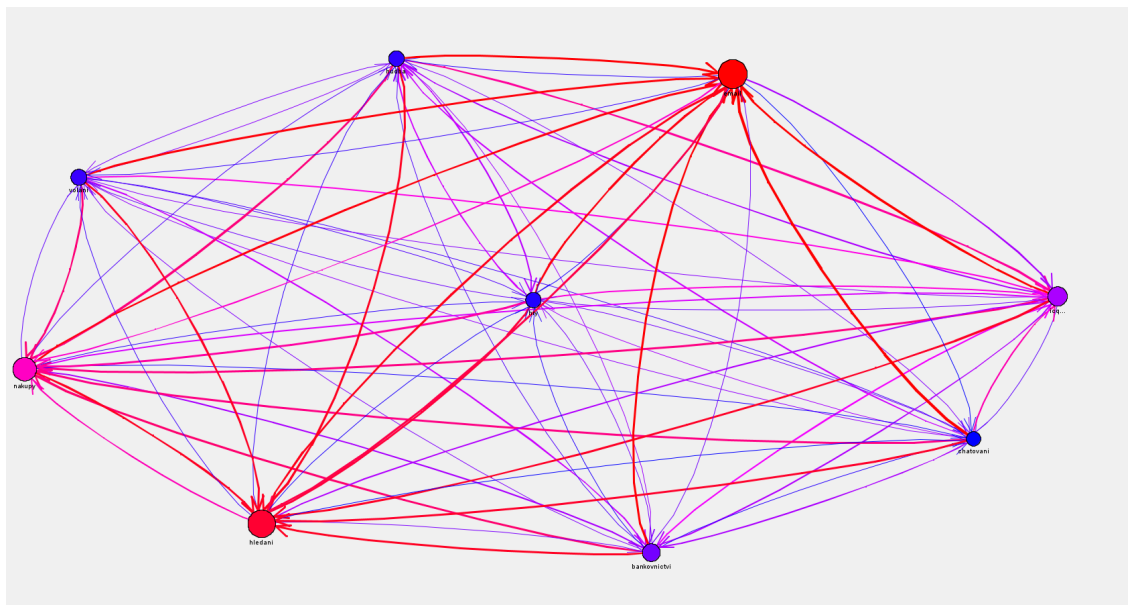
Nakonec jsme pro zajímavost vyzkoušeli přidat do transakčních dat i informaci, do jaké věkové kategorie daní respondenti patří. Nalezená pravidla ovšem opět pouze poukázala na významnost služeb e-mail a internetového vyhledávání. Největší spolehlivost pro využívání e-mailu měla 2. věková kategorie (tj. věk 21 – 24 let), a to ve výši 97 %, takže většina respondentů ve věku 21 až 24 let měla a používala e-mail.



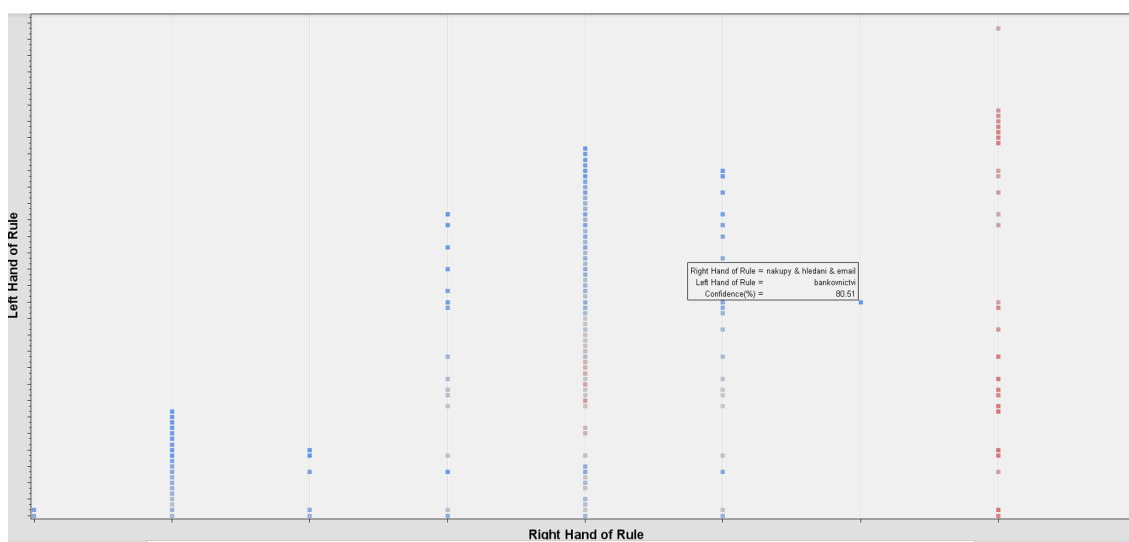
Obrázek 8: Notebooky vs. počítače

Confidence(%) ▼	Support(%)	Lift	Transaction Count	Rule	Left Hand of Rule	Right Hand of Rule	Rule Item 1	Rule Item 2	Rule Item 3
98.58	29.90	1.05	831.00	chatovani... chatovani	email	chatovani	=====	email	
98.40	50.88	1.05	1414.0	icq... ==>... icq...	email	icq...	=====	email	
98.00	43.97	1.04	1222.0	bankovni... bankovni...	email	bankovni...	=====	email	
97.90	35.26	1.04	980.00	hudba ==... hudba	email	hudba	=====	email	
97.31	36.49	1.03	1014.0	volani ==... volani	email	volani	=====	email	
96.85	67.54	1.03	1877.0	nakupy ==... nakupy	email	nakupy	=====	email	
96.65	33.21	1.03	923.00	hry ==>... hry	email	hry	=====	email	
95.29	83.77	1.01	2328.0	hledani ==... hledani	email	hledani	=====	email	
94.41	34.01	1.07	945.00	hudba ==... hudba	hledani	hudba	=====	hledani	
94.15	42.25	1.07	1174.0	bankovni... bankovni...	hledani	bankovni...	=====	hledani	
93.96	65.53	1.07	1821.0	nakupy ==... nakupy	hledani	nakupy	=====	hledani	
93.86	35.19	1.07	978.00	volani ==... volani	hledani	volani	=====	hledani	
93.60	48.40	1.06	1345.0	icq... ==>... icq...	hledani	icq...	=====	hledani	
93.00	28.21	1.06	784.00	chatovani... chatovani	hledani	chatovani	=====	hledani	
90.79	31.20	1.03	867.00	hry ==>... hry	hledani	hry	=====	hledani	
88.99	83.77	1.01	2328.0	email ==... email	hledani	email	=====	hledani	
84.20	37.78	1.21	1050.0	bankovni... bankovni...	nakupy	bankovni...	=====	nakupy	
83.01	31.13	1.19	865.00	volani ==... volani	nakupy	volani	=====	nakupy	
82.52	29.72	1.18	826.00	hudba ==... hudba	nakupy	hudba	=====	nakupy	
81.35	42.07	1.17	1169.0	icq... ==>... icq...	nakupy	icq...	=====	nakupy	
81.26	24.65	1.17	685.00	chatovani... chatovani	nakupy	chatovani	=====	nakupy	
78.32	28.21	1.51	784.00	hudba ==... hudba	icq...	hudba	=====	icq...	
77.28	26.56	1.11	738.00	hry ==>... hry	nakupy	hry	=====	nakupy	
75.33	22.85	1.46	635.00	chatovani... chatovani	icq...	chatovani	=====	icq...	

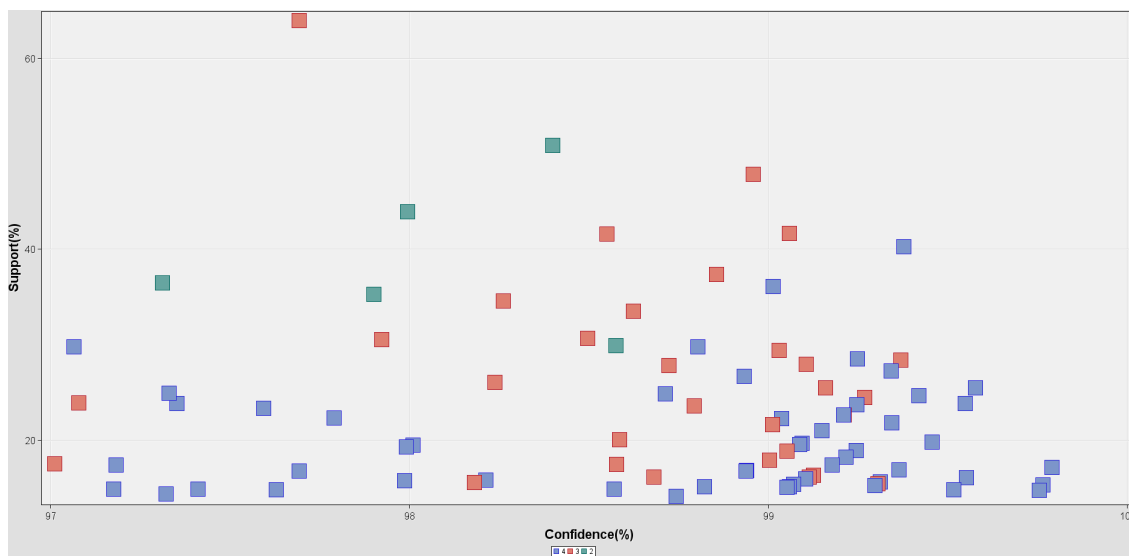
Obrázek 9: Pravidla podle spolehlivosti



Obrázek 10: Link Graph



Obrázek 11: Rule Matrix, max. čtyřprvková pravidla



Obrázek 12: Dvou, tří a čtyřprvková pravidla

Závěr

V tomto projektu jsme hledali odpověď postupně na tři položené otázky. Nejprve jsme se zabývali obecnou otázkou, jaká byla situace v ČR z hlediska využívání internetu. Zjistili jsme, že věk respondenta výrazně ovlivňuje jeho míru využívání internetu. Mladší lidé trávili na internetu více času a používali více služeb než starší ročníky. Pro všechny věkové kategorie je také kladná závislost mezi množstvím času stráveného u internetu během všedního dne a víkendu.

Účelem další části projektu bylo zjistit, jaké internetové služby respondenti využívali nejčastěji. Necelá pětina respondentů udala, že aktivně využívá e-mail. Dalšími hojně se vyskytujícími službami v našich datech bylo vyhledávání na internetu, nákupy přes internet nebo například využívání služeb okamžité komunikace.

Nakonec jsme zjišťovali, jaké služby byly v té době nejčastěji používány současně. Vzhledem k vysoké četnosti výskytu kladné odpovědi u služeb e-mail a vyhledávání byly nejčastěji současně využívané služby právě různé kombinace služeb obsahující e-mail a vyhledávání. Narazili jsme ale i na překvapivější souvislosti, například na vztah mezi hraním her a poslechem hudby nebo nakupováním přes internet a vlastnictvím internetového bankovníctví.

Možným nedostatkem této analýzy by mohl být fakt, že kvůli přílišnému výskytu pravidel „služba \Rightarrow e-mail nebo vyhledávání“ jsme mohli přehlédnout nějaká sice méně spolehlivá, ale zajímavá a méně očekávaná pravidla.

Bylo by nepochybně zajímavé podobnou analýzu provést i na datech pocházejících z dnešní doby a získané výsledky porovnat s výsledky z tohoto projektu. Bylo by možné pozorovat změny četnosti využívání jednotlivých služeb, které by byly jistě podstatně vyšší než na přelomu tisíciletí.