

GYMNASIUM JANA KEPLERA

Parléřova 2/118, 169 00 Praha 6



Komparativní analýza Bayesiánského a klasického přístupu k třídění a vyhledávání informací

Maturitní práce

Autor: Matěj Dvořák

Třída: R8.A

Školní rok: 2020/2021

Předmět: Informatika

Vedoucí práce: Šimon Schierreich

Praha, 2021



GYMNASIUM JANA KEPLERA
Kabinet informatiky

ZADÁNÍ MATURITNÍ PRÁCE

Student: Matěj Dvořák

Třída: R8.A

Školní rok: 2020/2021

Platnost zadání: 30. 9. 2021

Vedoucí práce: Šimon Schierreich

Název práce: Komparativní analýza Bayesiánského a klasického přístupu k třídění a vyhledávání informací

Pokyny pro vypracování:

Cílem práce je vytvořit systém pro automatickou klasifikaci textu v anglickém jazyce. Systém bude umožňovat a) sestavení vhodného hierarchického systému a jeho ruční obohacování o texty z různých zdrojů, b) analýzu pomocí Bayesiánského klasifikátoru a určení pravděpodobnosti příslušnosti souboru do dané složky a c) jednoduchý crawler vyhledávající stránky spadající do zadané kategorie.

Doporučená literatura:

[1] SRIVASTAVA, Ashok N. a Mehran SAHAMI, ed. Text Mining: Classification, Clustering, and Applications. London: Chapman and Hall/CRC, 2009. ISBN 978-1-4200-5940-3.

[2] RASCHKA, Sebastian. Naive Bayes and Text Classification I - Introduction and Theory [online]. 2017.

Dostupné z: <https://arxiv.org/abs/1410.5329>

URL repozitáře:

<https://github.com/xdvom03/klaus>

vedoucí práce

student

V Praze dne 29. 10. 2020

Prohlášení

Prohlašuji, že jsem svou práci vypracoval samostatně a použil jsem pouze prameny a literaturu uvedené v seznamu bibliografických záznamů. Nemám žádné námitky proti zpřístupňování této práce v souladu se zákonem č. 121/2000 Sb. o právu autorském, o právech souvisejících s právem autorským a o změně některých zákonů (autorský zákon) ve znění pozdějších předpisů.

V Praze dne 8. dubna 2021

Matěj Dvořák

Poděkování

Děkuji vedoucímu práce Šimonu Schierreichovi za poskytnutí materiálů a konzultace. Děkuji spolužákům (Šimon Brecher, Benjamín Verner) za konzultace struktury dat a algoritmů.

Abstrakt

Práce se zabývá tvorbou programu pro zpracování textu pomocí systému Naive Bayes. Tuto myšlenku zobecňuje pro hierarchický systém kategorií. Výsledný klasifikátor využívá pro tvorbu zaměřeného crawleru umožňujícího vyhledávání stránek do dané kategorie. Na závěr práce rozebírá dva příklady spuštění takového crawleru na příkladných datech.

Klíčová slova

zpracování přirozeného jazyka, naivně Bayesiánská klasifikace, webový adresář, sémantické vyhledávání, směrovaný crawler

Abstract

The project is concerned with generalising the Naive Bayes text classification algorithm to a hierarchical class system. The resulting automatic classifier is used to create a focused crawler searching for World Wide Web sites within a sought class. The work concludes with discussing two example runs of such a crawler on example data.

Keywords

natural language processing, naive Bayes, web directory, semantic search, focused crawler

Obsah

1	Teoretická část	3
1.1	Existující technologie	3
2	Implementační část	5
2.1	Klasifikace	5
2.1.1	Dvě kategorie	5
2.1.2	Zobecnění pro N kategorií	5
2.1.3	Hierarchický systém	7
2.1.4	Benchmark klasifikace	8
2.2	Crawler	8
2.2.1	robots.txt	8
2.2.2	Skóre	8
2.2.3	Algoritmus crawleru	9
2.3	Zdrojová data	9
2.3.1	Tokenizace	9
2.3.2	Opakované tokeny	9
2.3.3	Boilerplate	10
2.4	Uživatelské prostředí	10
2.5	Přiložená klasifikace	10
2.6	Design	11
2.6.1	Použité technologie	11
2.6.2	Ukládání dat	12
2.6.3	Databáze?	12
2.6.4	Rychlost	12
2.6.5	Slepé uličky	13
2.6.6	Možná vylepšení	14
3	Technická dokumentace	15
3.1	Požadavky spuštění	15
3.2	Release	15
3.3	Kompilace ze zdrojového kódu	15
3.4	Navigace	15
3.5	Import/export	17
4	Závěr	19
	Závěr	19
4.1	Tvorba hierarchického systému	19
4.2	Klasifikace daného textu	19
4.3	Crawler	19
4.3.1	Test crawleru	19
4.4	Shrnutí	20
	Seznam obrázků	21
	Seznam tabulek	22

1. Teoretická část

Před existencí internetového vyhledávače Google byly vyhledávače často neefektivní: nacházely spam a nenacházely relevantní stránky s gramaticky jinými formami hledaných klíčových slov. Proto se na ně uživatelé nespolehali jako na jediný způsob, jak nacházet informace na internetu. Jednou z alternativ byly webové adresáře (web directories), veřejně dostupné indexy odkazů tříděné podle tématu. Mezi takové adresáře patřilo například Yahoo! Directory nebo DMOZ.

Adresáře byly udržovány ručně na základě podnětů od uživatelů. Tento systém se s růstem počtu dokumentů na internetu ukázal neúnosným, zatímco vyhledávače šlo škálovat. Ve výsledku tak dnes DMOZ ani Yahoo! Directory oficiálně neexistují a dobrovolníky vedené Curlie (vycházející z DMOZ)¹ je těžce zastaralé, značná část odkazů už ani nefunguje. Výsledkem je dnešní stav, kde má na vyhledávání téměř monopol Google.

Vyhledávače na základě klíčových slov ale také nejsou ideální. Dobře si poradí s konkrétními požadavky (např. chybová hláška nebo název restaurace). Necháme-li stranou problematiku jejich monetizace a SEO, stále trpí omezenou schopností pracovat do hloubky se širokým dotazem. Požadavku o několika slovech odpovídá řada dokumentů, výběr z nich ale probíhá arbitrárně nebo na základě řady faktorů, jako jsou stáří odkazu, další klíčová slova, nebo informace o vyhledávajícím uživateli (např. lokace). Tyto faktory jsou těžko ovlivnitelné, prozkoumat vyhledávaný široký dotaz ze všech úhlů pohledu proto obvykle nejde. Někdy navíc uživatel nemusí správné klíčové slovo znát.

Adresáře místo toho fungují na základě tématických okruhů. Je-li adresář dobře veden, každý okruh obsahuje řadu podokruhů korespondujících k různým podtématům. Adresář tak slouží jako encyklopedie a počáteční bod pro ruční vyhledávání informací. Je tak kromě vyhledávače i pedagogickým nástrojem.

Problém s udržováním adresářů je ale stále platný. Tento projekt si klade za cíl jej částečně automatizovat - vyhledávat texty podobné již přidaným stránkám. Uživatel, manažer adresáře, by pak musel jen adresář prohlubovat a složky dále dělit.

1.1 Existující technologie

Zpracování přirozeného jazyka pomocí metody Naive Bayes má široké uplatnění, spíše je ale využíváno pro konkrétnější účely, jako klasifikace uživatelských recenzí². Existuje projekt Klassify³, umožňující obecně natrénovat klasifikátor, ale jde spíše o demonstraci a chybí základní možnosti úprav databáze, mazání/přesouvání souborů, apod. Našel jsem jeden pokus o obecnou klasifikaci⁴ vycházející z knihovnického systému, který není dobře přenositelný na internetové stránky. Ukázala se na něm ale platnost principu.

¹<https://curlie.org/>

²<https://arxiv.org/pdf/1610.09982.pdf>

³<https://github.com/fatiharikli/klassify>

⁴https://www.researchgate.net/publication/4046602_Classification_of_Web_Documents_Using_a_Naive_Bayes_Method

2. Implementační část

2.1 Klasifikace

Program funguje na základě statistické metody Naive Bayes, kterou se řídí při rozhodování, kam daný dokument zařadit. Nejjednodušším případem je rozhodování mezi dvěma kategoriemi.

2.1.1 Dvě kategorie

Dokument modelujeme jako seznam slov v něm obsažených (viz část Tokenizace). U kategorií vybudujeme součet pro všechny dokumenty v nich obsažené korpus, v němž je uveden počet výskytů každého slova. Při porovnávání dvou kategorií pracujeme s normalizovanými korpusy. U většího z korpusů podělíme výskyty všech slov stejnou konstantou tak, aby měly oba výsledné korpusy dohromady stejný počet slov. U každého slova určíme počet jeho výskytů v obou korpusech. Označme výskyt slova v kategorii A a , v kategorii B b . Pravděpodobnost pro dané slovo, že dokument jej obsahující patří do kategorie A, určíme jako:

$$P_A = \frac{a + k}{a + b + 2k}$$

Ve výpočtu zahrnujeme konstantu k , takzvaný “smoothing parameter”, podstatnou zejména v případě, že a nebo b je velmi malé. Program využívá konzervativní Laplace smoothing, kde $k = 1$ (Lidstone smoothing, nižší hodnoty k , by v takových případech vedlo k vyšší váze takových slov). Takto určíme pravděpodobnosti pro všechna slova v dokumentu. Dále vybereme relevantní klíčová slova.

Nevyužíváme všechna slova v dokumentu, protože bychom mohli zesílit šum. Pro každou z kategorií uvažujeme jen X nejsilnějších slov proti kategorii, kde X je polovina počtu slov nad stanovenou “zajímavostí” (konkrétně slova s pravděpodobností pro jednu z kategorií pod 0.2). X má nastavenou minimální hodnotu, takže vždy uvažujeme alespoň 12 slov, ale většinou je zajímavých slov více. Smysl tohoto výběru je snaha o nalezení stejného množství důkazů proti oběma stranám a porovnání jejich celkové síly.

Skóre pro tato klíčová slova vynásobíme. Tento proces provedeme stejně pro kategorie v opačném pořadí (vyměníme a a b). Výsledné pravděpodobnosti pro obě kategorie vynásobíme společnou konstantou tak, aby se sečetly na 1 (tomuto dále v textu budeme říkat “normalizace pravděpodobnosti”, NP). Tento krok je primárně řízen přehledností a snahou vytvořit čisté pravděpodobnosti, ne dvě velmi nízká skóre.

2.1.2 Zobecnění pro N kategorií

Problém zařazení do jedné z několika kategorií má několik řešení. Při volbě modelu jsem se zejména řídil požadavkem, aby obsah jedné kategorie neovlivnil rozhodování mezi ostatními kategoriemi. Proto například nevolím metodu, kde se každá kategorie srovná s celým vzorkem a vyhraje ta s

nejlepším výsledkem¹. Každá kategorie by byla v takovém případě součástí všech rozhodování. Pokud by například jedna z kategorií byla “úplně mimo”, všechna ostatní skóre by mohla vyjít poblíž jedné a rozlišovací schopnost by se ztratila. Toto je případ např. nejvyšší úrovně přiloženého rozdělení, kde by takto působily legální dokumenty².

Namísto toho začínám nalezením pravděpodobností pro každou dvojici kategorií. Dále vycházím z předpokladu, že správná kategorie by měla porazit všechny ostatní. Jakákoli kategorie tedy může být sama o sobě kontrolní pro jakoukoli jinou. U každé kategorie získáme výslednou pravděpodobnost příslušnosti vynásobením pravděpodobností proti všem ostatním kategoriím. Pro prezentaci provedeme NP všech výsledků.

Obrázek 2.1: Vysvětlení klasifikace

Table explainer							
	articles	fiction	legal	news	non-english	utility	
articles		0.0	0.0	-0.018	0.0	0.0	-0.018
fiction	-35.068		0.0	-31.381	0.0	-9.394	-75.843
legal	-97.119	-63.411		-61.916	0.0	-23.784	-246.23
news	-4.008	0.0	0.0		0.0	-0.428	-4.436
non-english	-442.527	-395.169	-297.29	-383.392		-333.856	-1852.235
utility	-10.242	0.0	0.0	-1.054	0.0		-11.297

Na této ukázce je vidět klasifikace Wikipedia Terms of Use. Čísla v tabulce jsou logaritmy jednotlivých pravděpodobností, čísla vpravo od tabulky jejich součty. Jak je vidět, kategorie /articles/ porazila všechny ostatní kategorie a jejich vzájemná skóre už nejsou podstatná.

Další výhodou tohoto systému je fakt, že výstupem je stále pravděpodobnost, nikoli pouze pořadí.

Výstup vícekategoričké klasifikace

Pravděpodobnost klasifikace do dané kategorie vrací funkce (prob). Na vyšší úrovni je funkce (place), jež najde nejpravděpodobnější konečné umístění stránky. To získáme klasifikační úroveň po úrovni. Na každé úrovni sledujeme nejlepší možnost, pokud jí přísluší pravděpodobnost překračující danou mezní hodnotu. Zde ale nevyužíváme NP. Logika je taková, že pokud je NP podstatný efekt, znamená to, že všechny kategorie mají malé skóre, což znamená, že proti všem kategoriím existují důkazy. To může nastat v případě, že je klasifikace založená na mylném předpokladu: stránka pravděpodobně nepatří do žádné z možností. To může mít tři hlavní důvody: Bud' nepatří do kategorie, v rámci níž klasifikujeme (tj. někde předtím nastala chyba v klasifikaci), nebo daná podkategorie neexistuje, nebo nepatří do podkategorie, ale jen do celé kategorie.

¹Tento předpoklad s vznešeným názvem *Maximum-likelihood estimate* je ve většině zdrojů, např. v <https://arxiv.org/pdf/1410.5329.pdf>, nebo v <https://web.archive.org/web/20200708084909/https://web.stanford.edu/class/cs124/lec/naivebayes.pdf> (slide 55). Pokud vím, mou alternativu nikdo jiný zatím netestoval, možná také proto, že její časová náročnost roste s druhou mocninou počtu kategorií, takže pro více kategorií prakticky vynucuje hierarchický systém.

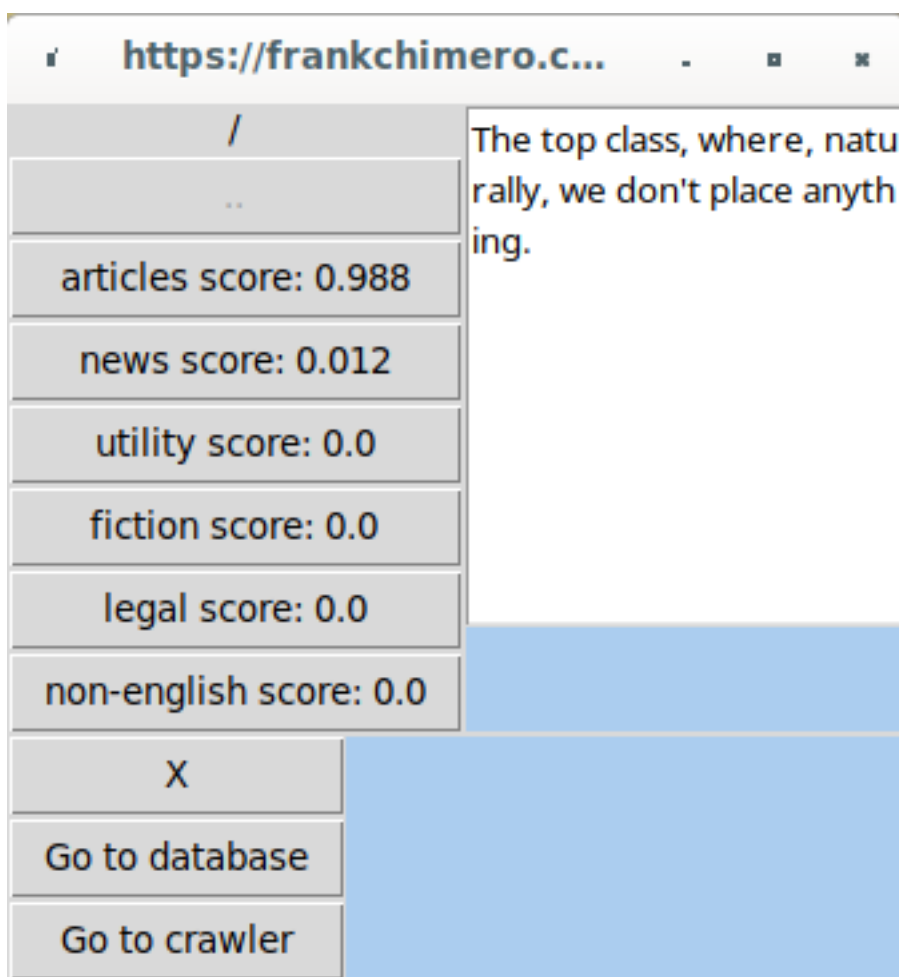
²Jednou z prvních kategorií, kterou jsem testoval, jsou dokumenty typu “Terms of Use” a “Privacy Policy”, pro jejich snadnou odlišitelnost (a velký výskyt, pro crawler jsou užitečné, aby se jich dokázal vyvarovat)

Logaritmický formát čísel

Ježto se s dostatkem slov můžeme dostat na pravděpodobnosti velice blízké nule, hrozí ztráta přesnosti, v horším případě zaokrouhlení na nulu. Proto všude namísto s pravděpodobnostmi počítáme s jejich přirozenými logaritmy. Ty také zobrazujeme ve vysvětlovacím okně (viz manuál).

2.1.3 Hierarchický systém

Obrázek 2.2: Možnost výběru další úrovně



Kategorie může obsahovat podkategorie. Kategorie může najednou obsahovat soubory i podkategorie. Při rozhodování se berou v úvahu všechny dokumenty v kategorii a všech úrovních jejích podkategorií. Klasifikace pak může být vícezkroková. Na obrázku je vidět okno při vysvětlení klasifikace, kde je možné kliknutím na příslušné tlačítko pokračovat v klasifikaci v libovolné podkategorii.

2.1.4 Benchmark klasifikace

Tento systém jsem testoval na často užívaném datasetu 20 Newsgroups³. Hierarchie jsem sestavil dle doporučení autora datasetu⁴. Výsledná přesnost byla 6093 z 7528 dokumentů klasifikovaných správně, tj. zhruba 81%, což je srovnatelné s alternativními technologiemi. Naive Bayes s dodatečnými “vylepšeními” jako stemming (zjednodušení na základní gramatické tvary slov) a odstranění stop words (velmi častých slov bez významového obsahu) získal (dle prvního zdroje) výsledek 6100 z 7528, z čehož usuzuji, že efekt těchto technologií je zanedbatelný v rámci chyby měření. K podobnému výsledku dochází i většina literatury, např. zde⁵ dochází autoři k výsledku, že v angličtině různé lematizační algoritmy vedou k menšímu korpusu, ale vždy za cenu poklesu efektivity. Na škále, na které pracuje tento projekt, nemá zmenšování korpusu prioritu.

2.2 Crawler

Crawler vychází z předpokladu, že stránky budou spíše odkazovat na stránky s podobným tématem. Snaží se tedy stále sledovat co nejlépe odpovídající stránky a odkazy z nich. Pro specifické požadavky na několikáté úrovni klasifikace to činí tak, že (obvykle) najde nejprve obecnější požadavek a až z něj požadavek konkrétní. Například pro /news/sports/ se pravděpodobně nejprve dostane crawler na stránku novin a až poté najde sportovní sekci. Tomuto algoritmu by měl být přizpůsoben hierarchický systém.

2.2.1 robots.txt

Bot splňuje robots.txt podle původní specifikace⁶. Je identifikován v hlavičce požadavku na stránku, kde je také uvedena adresa projektu a má emailová adresa pro případ, že by někde způsobil problémy. Na stránky neposílá přehnaně mnoho požadavků - teoretické maximum během jednoho crawlu je zhruba 20*(počet stránek na doménu) a vzhledem ke klasifikaci je posílá postupně a se zdržením.

2.2.2 Skóre

Crawler má za cíl najít stránky v dané kategorii, vychází tedy ze skóre určujícího, jak moc se daná stránka cílové kategorii přiblížila. Skóre vychází z nejpravděpodobnějšího umístění kategorie. Je součtem počtu úrovní (od počáteční), ve kterých se shoduje cílové a skutečné umístění, a pravděpodobnosti pro první neshodující se kategorii. Například: hledá-li crawler /news/sports/tennis/ a najde /news/news-reports/, je skóre 1 plus pravděpodobnost, že namísto do /news/news-reports/ patří stránka do /news/sports/.

³<https://ana.cachopo.org/datasets-for-single-label-text-categorization>

⁴<http://qwone.com/~jason/20Newsgroups/>

⁵Toman, Michal & Tesar, Roman & Jezek, Karel. (2006). Influence of Word Normalization on Text Classification.

⁶<http://www.robotstxt.org/orig.html>

2.2.3 Algoritmus crawleru

Vstupem je počáteční URL adresa a počet kroků. Algoritmus crawleru:

1. Ze zatím prošlých URL vybereme to s nejlepším skóre, které jsme zatím takto nevybrali (tj. po vybrání odstraníme URL z možností výběru).
2. Najdeme všechny odkazy vedoucí z tohoto URL. Odstraníme duplikátní odkazy⁷.
3. Vyloučíme všechny nefunkční odkazy⁸ a všechny odkazy vedoucí na již nám známé domény. U vybraných odkazů se podíváme, ze které adresy přišla odpověď (odkaz totiž mohl vést na přesměrování), a případně tuto adresu odstraníme, pokud je duplikátní nebo nepovolená.
4. Stáhneme všechny zbylé cílové stránky, oskórujeme. Vybereme z nich tu s nejlepším skóre. Skóre měříme v náhodném pořadí, je-li odkazů moc, můžeme po daném počtu přestat.
5. Na nové doméně několikrát (počet je nastavitelný) sledujeme odkazy zůstávající uvnitř domény, vždy opět na ten s nejlepším skóre. Všechna URL po cestě ukládáme do seznamu pro bod 1.
6. Opakujeme kroky 1-5 tolikrát, kolik domén je specifikováno projít.

Pokud v kroku 4 nemáme žádné funkční odkazy, doménu pouze odstraníme a zkusíme jinou. Pokud dojdou domény, systém zavolá chybu, která se uživateli zobrazí. Pokud nenajdeme funkční odkazy v kroku 5, zkusíme to znovu, ale neplatný pokus počítáme jako pokus (abychom se mohli dostat pryč z domény, která neodkazuje na sebe samu).

Při výpočtu skóre automaticky uložíme nalezené umístění URL do dat, která periodicky ukládáme do souboru /DATA/discovered. Tento soubor pak můžeme procházet (viz manuál).

2.3 Zdrojová data

2.3.1 Tokenizace

Většina dokumentů je získána stažením HTML z internetu. V takovém případě nejprve odstraníme všechn text v tazích <script>, <style>, nebo <head> a následně všechny HTML tagy. Dále text zbavíme diakritiky a všech ne-alfanumerických znaků. Všechn odstraněný text nahrazujeme mezerami. Výsledek rozdělíme na tokeny podle mezer.

2.3.2 Opakované tokeny

Během vývoje jsem zkoušel dva extrémní způsoby řešení opakování tokenů. Bud' opakování nepovažujeme za problém a tokeny počítáme vícekrát, nebo počítáme každý token pouze jednou. V prvním případě může hromadné opakování tokenů na jedné stránce pokřivit vnímání tokenů (což

⁷Za duplikátní považujeme odkazy lišící se jen rozdílem http/s, velikostí písmen v doméně, fragmentem (částí za # v odkazu).

⁸Odkazy na nefunkční nebo nedostupné stránky, příliš krátké stránky, ale také stránky zakazující funkci crawleru souborem robots.txt

se například stalo s malými čísly, která disproporčně naznačovala matematiku, když přitom může jít např. o nadpisy kapitol), ve druhém případě nefungují častá slova ve dlouhých textech (jako je například kategorie /fiction/). Proto text rozdělujeme na úseky o dané délce (několika stovek slov), ve kterých tokeny neopakujeme, mezi nimi ale ano. Analogicky vybíráme tokeny v klasifikovaném textu.

2.3.3 Boilerplate

Na doméně se často nachází text nesouvisející s konkrétní stránkou, určený např. pro navigaci na stránce, nebo pro splnění legálních požadavků (disclaimer psaný malým písmem ve spodní části stránky). Takový text se snažíme u trénovacích dat odfiltrovat, aby neovlivňoval povahu kategorie. Proto u každého dokumentu poté, co ho stáhneme, během výstavby dat nejprve vytvoříme jeho verzi s co největší částí takového textu odstraněnou. Pro každou doménu v systému, od které máme alespoň dva soubory, nalezneme překryv všech jejich textů. Překryvem je jakýkoli úsek o daném minimálním počtu slov, který se vyskytuje ve všech souborech. Takové úseky pak ze souborů odstraníme, čímž vzniknou texty očištěné, které ukládáme ve složce DATA/core/.

2.4 Uživatelské prostředí

Grafické prostředí je limitováno schopnostmi knihovny LTK, v rámci nichž funguje analogicky k procházení souborového systému. Ježto nemusí uživateli zobrazovat dokumenty ani velké množství kategorií najednou, je na obrazovce hodně volného místa, takže jsou všechny možnosti zobrazeny najednou namísto ukrytí do submenu.

Hlavní netradiční částí UI je systém přesouvání souborů a kategorií. Ježto v knihovně není jednoduché vytvořit click-and-drag, systém běžně používaný pro přesouvání, inspiroval jsem se na C2 Wiki⁹ a využil systém interně nazývaný Bucket: soubory nebo složky může uživatel zkopírovat do části Bucket, ve které zůstávají při pohybu skrz souborový systém, a ze které mohou být přesunuty do jiného místa. Umožňuje to jednodušší přesuny mezi vzdálenými kategoriemi. Podobně přesouváme i celé kategorie. Tlačítka s módy umožňují výběr mezi přesunem, smazáním, nebo jen odstraněním z části Bucket. Pro odlišení od tlačítek provádějících nějakou akci je na nich text psán velkými písmeny.

2.5 Přiložená klasifikace

Souběžně s vývojem prostředí jsem vyvíjel i vlastní soubor trénovacích dat kvůli testování systému. Mimo jiné jsem díky němu zjistil, že užitečnější systém uvažování, než ontologie (ukládání textů podle toho, o čem *jsou*), je ukládání podle principu vyhledávacího (požadavek, pro který by nalezení daného dokumentu dávalo smysl). Tím se vyhneme často filozofickým otázkám o tom, co je *doopravdy* např. vzdělávání.

⁹<http://wiki.c2.com/?HandVsPointer>

2.6 Design

2.6.1 Použité technologie

Jazyk

Program jsem se rozhodl vytvořit v jazyku Common Lisp. Roli při výběru hrály: kombinace funkcionálního a imperativního paradigmatu, rychlost výpočtů, jednoduchá prefixová syntaxe vhodná pro zápis matematických vzorců na více řádků (což usnadňuje práci s výpočty), předchozí zkušenost se souvisejícími problémy v tomto jazyku, výběr v inspirujícím spamovém filtru (viz níže). Klasifikace je naprogramována převážně funkcionálně, crawler převážně imperativně. Uživatelské rozhraní je řízeno na bázi událostí (stisknutí tlačítek, zavření okna).

Knihovny

LTK¹⁰: Grafické rozhraní. Jde o port Tcl/Tk pro Common Lisp a koncepčně nejjednodušší řešení, ale některé funkce v knihovně chybí (zejména volba barvy tlačítka a okno pro výběr souboru - které jsem ale vytvořil vlastní a snad funguje).

Drakma¹¹: Stahování HTML. Nejprve jsem používal knihovnu Dexador¹², rozšíření Drakma, ale ukázalo se, že trpí chybou, která způsobuje rostoucí využití paměti. Tuto chybu jsem nahlásil¹³, ale knihovna není udržovaná a nepřišla žádná reakce.

QURI¹⁴ umožňuje jednak úpravy znaků v názvech (např. `example.com/B%C3%A1g%C5%99i` -> `example.com/Bágři`), jednak práci s odkazy - najde cíl odkazu i pro řadu speciálních případů (kromě pasti, kterou je dvojí lomítko pro relativní protokol) a dokáže obecně najít klíčovou část domény pro odkazy (mimo jiné pozná, že `.co.uk` je jen jedna koncovka).

Plump¹⁵ dokáže rozklíčovat escaped znaky v HTML (`’` -> `'`).

Trivial-timeout¹⁶ zastaví načítání nereagující stránky.

Cl-strings¹⁷ obsahuje řadu utilit pro práci s textem.

Alexandria¹⁸ umožňuje kopírování struktury hash table.

¹⁰<http://www.peter-herth.de/ltk/index.html>

¹¹<https://github.com/edicl/drakma>

¹²<https://github.com/fukamachi/dexador>

¹³<https://github.com/fukamachi/dexador/issues/97>

¹⁴<https://github.com/fukamachi/quri>

¹⁵<https://github.com/Shinmera/plump>

¹⁶<https://common-lisp.net/project/trivial-timeout/>

¹⁷<https://github.com/diogoalexandrefranco/cl-strings>

¹⁸<https://gitlab.common-lisp.net/alexandria/alexandria>

2.6.2 Ukládání dat

Program má tři druhy dat: Korpusy, konfigurační data ke složkám a zdrojové dokumenty. Zdrojové dokumenty jsou z nich zdaleka největší, proto je každý dokument uložen v samostatném souboru. Při vložení souboru se vytvoří jeho tři verze - hrubý text tak, jak byl vložen, text očištěný dle tokenizace a text s odstraněnými duplikátními prvky (viz Boilerplate). Tyto zpracované verze textu by bylo lze vytvořit znovu, např. v případě změny povolených znaků.

Ostatní data nepřesahují pro typické používání jednotky MB, jsou tedy uloženy pro všechny kategorie dohromady v jednom souboru. To má výhodu jednoduchého exportu, viz manuál.

2.6.3 Databáze?

Data program při spuštění načte ze souborů a opět je do nich uloží, pokud uživatel uloží změny. Jinak má celý korpus v pracovní paměti, což zrychluje klasifikaci.

V současném stavu má korpus 8.6 MB. Obsahuje kolem 100000 různých slov a 200 kategorií. Počet různých slov už řádově nevzroste, např. Oxford English Dictionary obsahuje kolem 600000 slov¹⁹, ale řada z nich je reálně nevyužívaná. Počet kategorií může růst s vývojem. Spíše ale porostou počty slov, což na velikost dat prakticky nebude mít vliv. Na problémy by mohl tento systém narazit při řádově větším počtu kategorií.

2.6.4 Rychlost

Úpravy databáze

Program se snažíme zrychlit zejména dvěma technikami: prací v paměti, kdy omezuje práci se soubory na načtení a uložení, a omezením redundantní práce. U časově náročných činností, tedy odstranění boilerplate a výstavba korpusu, ukládáme vstupní data, ze kterých vznikla předchozí verze takto zpracovaných dat. Pokud zjistíme, že uživatel nezpůsobil žádné změny, necháme data tak, jak jsou. Tím udržujeme prostředí pro úpravu databáze použitelným.

Klasifikace

Většinu doby běhu klasifikace stráví program požadavky na vyhledání v hash table s počty slov v jednotlivých korpusech. Nic výrazně rychlejšího, než hash table, už ale nevymyslíme a všechny požadavky jsou nutné pro výpočet skóre. Teoreticky by se dal program zrychlit předpočítáním všech skóre všech slov, ale prakticky by to nepřiměřeně zpomalilo tvorbu dat.

¹⁹<https://www.oed.com/>

Crawler

Crawler na rozdíl od práce s databází data stránek neukládá. S HTML pracuje v paměti a zahodí je, jakmile už není potřeba. Dále má cache všech pravidel robots.txt, na která cestou narazil (ne souborů robots.txt, ale už zpracovaných pravidel), aby je nemusel stahovat a zpracovávat stále dokola. Ve výsledku tvoří většinu času běhu crawleru stahování dat ze stránek, což na straně programu vylepšit nemůžeme.

2.6.5 Slepé uličky

Nejdynamičtější byl vývoj uživatelského prostředí, které nejdřív mělo řadu zbytečných funkcí, jako historie vložených URL, a do databáze se musel uživatel prokliknout s jedním URL pro jeho zařazení (takže přidávání víc URL vyžadovalo pokaždé znovu najít cílovou kategorii).

Skórování mezi více kategoriemi nejdřív fungovalo na bázi algoritmu podobného PageRanku²⁰, kde se váha každé kategorie jako důkazu proti ostatním kategoriím určovala na základě výsledného skóre dané kategorie. Do věci to vnášelo složité další výpočty, které neměly dobré logické opodstatnění (na rozdíl od samotného PageRanku, který koresponduje s určitým modelem chování uživatele).

Hodně změn nastalo v souborovém systému. Před současným uložením do jednoho souboru jsem metaforu souborového systému v klasifikaci zakládal na skutečném souborovém systému, takže kategoriím odpovídaly složky. To bylo časem zbytečně nepřehledné.

Zbytečně dlouho jsem razil myšlenku informační konzole v GUI, ve které by šlo vidět většinu činností programu. Oproti běžné textové konzoli neměla žádné výhody (chtěl jsem ji rozdělit podle skupin oznámení, ale nikdy jsem to neudělal) a uživatelské prostředí značně komplikovala.

Prvotní inspirace programu pochází z Bayesiánského spamového filtru, který vyvinul Paul Graham²¹. Chtěl jsem jeho myšlenku zobecnit. Bohužel, detekce spamu je jako klasifikace textu značně specifický problém. Graham se v původním filtru dopustil řady zjednodušení, která fungovala konkrétně pro spam (nebo alespoň neškodila, většina spamu je odlišitelná snadno, tedy i s částečně chybným systémem). Musel jsem jednu po druhé vložit řadu oprav: kontrolu pro počet dokumentů v kategorii, kontrolu pro celkový počet slov v kategorii, smoothing, balancovaný výběr slov. Program se prakticky od září až do prosince, kdy padla poslední z (těchto) chyb, choval většinou absurdně.

Crawler měl řadu verzí, které uvažovaly pouze omezenou frontu možností pokračování, řazenou podle skóre, s různými variacemi. Výsledkem typicky bylo, že všechny možnosti měly nulové skóre a crawler se náhodně motal, protože když něco náhodou našel, nedal do fronty dost různých odkazů, takže mu brzy zase došly. Samotné skóre také mělo řadu pochybných variant, například prostě pravděpodobnost, že patří stránka přesně do dané kategorie (bez bodů za částečný překryv), což prakticky znamenalo, že pokud nenašel crawler přesně to, co hledal, všechny možnosti byly blízko nuly a rozhodoval spíše počet slov, podle kterých se klasifikace mohla řídit.

²⁰<http://infolab.stanford.edu/backrub/google.html>

²¹<http://paulgraham.com/spam.html>

2.6.6 Možná vylepšení

Technické nedostatky

U stránek pouze stahujeme HTML, což se může někdy značně lišit od toho, co by viděl uživatel. Nejdestruktivnější případy (kdy celou stránku řídí Javascript a není na ní vůbec statický text) dokážeme odchytit varováním o krátkém textu. Lepší by ale bylo využít například program Selenium²² pro napodobení webového prohlížeče. Na přiloženém korpusu se to projevuje mimo jiné nesprávným zpracováním stránek, na nichž je LaTeX. Dále program nevyužívá některých částí HTML standardu, zejména tagu `<base>` v hlavičce, což někde může způsobit, že crawler nedokáže sledovat odkazy.

Kategorizaci by mohla vylepšit opatrná verze technologie stemming.

Další funkce

U crawleru by se hodilo zobrazit více statistik o jeho cestě, možnost např. stromového zobrazení, ze kterých stránek se kam dostal. Algoritmus crawleru také není optimální, všechny domény prochází stejně dlouho. Určitě existuje lepší způsob vyvážení snahy o nalezení dobrého skóre a snahy o variabilitu výsledků.

Grafická knihovna LTK byla možná retrospektivně špatným rozhodnutím. Výsledné prostředí není příliš přehledné ani ovladatelné. Po úpravě čehokoli v okně musíme čekat na jeho aktualizaci, která se řídí mnou zatím nepochopenými pravidly. Text na tlačítkách se chová nerozumně, dlouhé odkazy způsobují natahování celého prostředí. Není dobrý způsob zrušení tohoto chování (a např. sbalení textu do tří teček, přičemž celý odkaz by byl vidět v tooltipu při najetí myší).

²²<https://www.selenium.dev/>

3. Technická dokumentace

3.1 Požadavky spuštění

Pro GUI je nutný grafický systém Wish z Tk. Instalovat lze např. pomocí APT: “sudo apt install tk”. Pro kompilaci ze zdroje je také potřeba stejně instalovat kompilátor SBCL.

Na systému Windows je nutno instalovat Tcl/Tk, např. IronTcl¹, a před spuštěním (main) zavolat (setf ltk:*wish-pathname* “cesta k souboru wish”) - pozor na potřebu dvou zpětných lomítek (escaped backslash) v adrese, aby v adrese zůstalo jedno.

3.2 Release

Nejjednodušší cestou pro spuštění programu je rozbalit jeden ze souborů ve složce /releases/ (značeny jsou dnem vytvoření) a spustit spustitelný soubor v rozbalené složce. Příložená data jsou vedena jako importy, oproti kompilaci ze souboru tedy není nutno nic stahovat.

Release jsem testoval na nově instalovaném systému Ubuntu 20.04 LTS ve Virtual Machine, kde ho šlo spustit a používat jen s požadavky výše uvedenými.

Bohužel se mi nepodařilo vytvořit funkční release pro systém Windows, ale kompilace ze zdrojového kódu alespoň na Windows 7 funguje.

3.3 Kompilace ze zdrojového kódu

Kvůli knihovnám je nutné nainstalovat manager knihoven pro Common Lisp, Quicklisp, tak, aby existoval soubor /portacle/all/quicklisp/setup.lisp. Doporučuji stáhnout systém Portacle² a sloučit složky portacle. Spusťte buďto instalované Portacle, nebo SBCL ve složce /portacle/. Pak stačí zkompileovat soubor portacle/load.lisp a do REPL³ zadat (main). Po chvíli načítání se objeví okno grafického rozhraní.

3.4 Navigace

Po spuštění se uživatel dostane do klasifikačního okna, kde může zjistit klasifikaci vloženého dokumentu (tlačítka “place”), případně si ji nechat vysvětlit (zobrazit pravděpodobnostní tabulky a klíčová slova, tlačítka “open database”). Může načíst text ze souboru⁴.

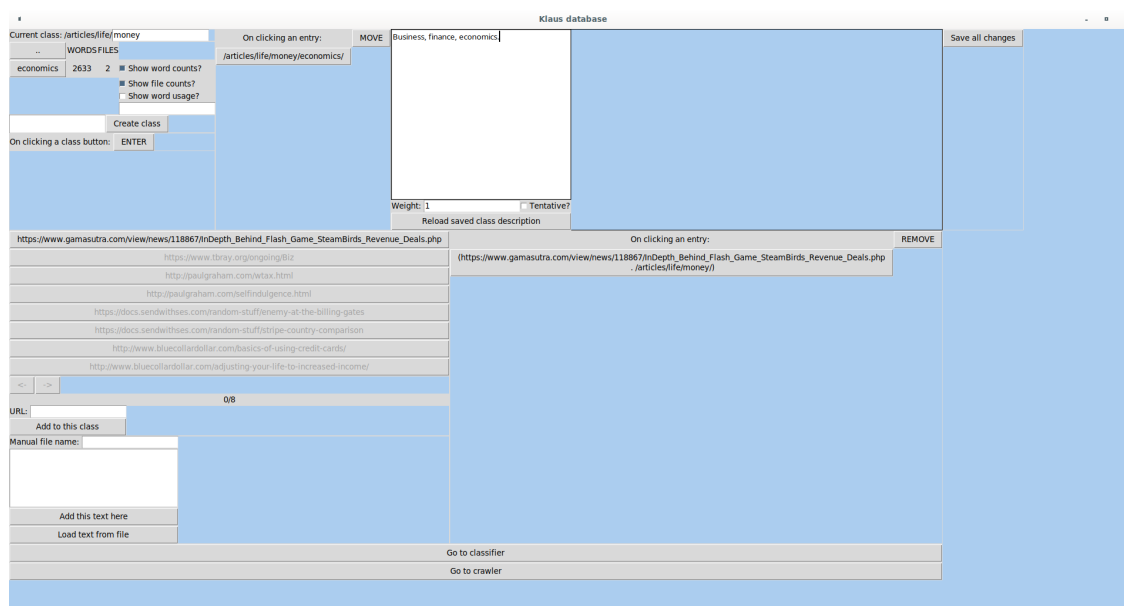
¹<https://www.iron-tcl.com/>

²<https://portacle.github.io>

³Read Eval Print Loop, konzole pro Common Lisp

⁴Platí tu též problém pro vadné soubory.

Obrázek 3.1: Databáze



Druhé prostředí umožňuje úpravy databáze trénovacích dat. Uživatel se může pohybovat mezi jednotlivými kategoriemi, zobrazit v nich obsažené soubory, komentáře k hierarchickému systému, procházet korpus. Přesouvá data kliknutím na tlačítko příslušného souboru, což jej přesune do části vpravo (pod On clicking on entry: REMOVE), kde může provést se souborem akci.

Může soubory přidávat (buď pomocí URL odkazu, nebo vepsáním do okna pro manuální vstup (je také potřeba dát textu název), případně lze do manuálního okna text načíst ze souboru⁵). Může také přesouvat, přejmenovávat, a mazat kategorie. Většinu dat (název momentální kategorie, komentář, váha, status “tentative” znamenající, že klasifikátor se nebude do dané kategorie nic pokoušet zařadit) může uživatel upravit prostě jejich přepisem a přesunem do jiné kategorie, což automaticky změny zaznamená, a následným uložením buď přes tlačítko, nebo při zavírání okna. Pozor, změny jsou zaznamenány až při přechodu do jiné kategorie.

Standardně se uživatel kliknutím na kategorii do dané kategorie přesune, změní-li ale ENTER na EDIT, chovají se kategorie analogicky jako soubory, se sekci Bucket vpravo od sekce kategorií. Vše je uloženo najednou tlačítkem “uložit změny” (Save all changes), při zavření okna vyskočí možnost uložit případné neuložené změny.

Třetí prostředí umožňuje spuštění crawleru a prezentaci jeho výsledků. Zatímco crawler běží, tiskne svůj momentální stav do konzole, je-li spuštěná, a posouvá se progress bar. Po proběhnutí crawleru se objeví databáze všech prošlých URL, kterou lze standardně procházet a editovat. Po zavření se data z crawleru nachází ve dvou souborech: DATA/discovered obsahuje všechna klasifikovaná URL a DATA/scores skóre pro URL, ze kterých crawler pokračoval dále.

Mezi prostředími lze vždy přejít tlačítkem ve spodní části okna. Při přechodu z databáze do crawleru se předvyplní momentální pozice v databázi jako cíl crawleru.

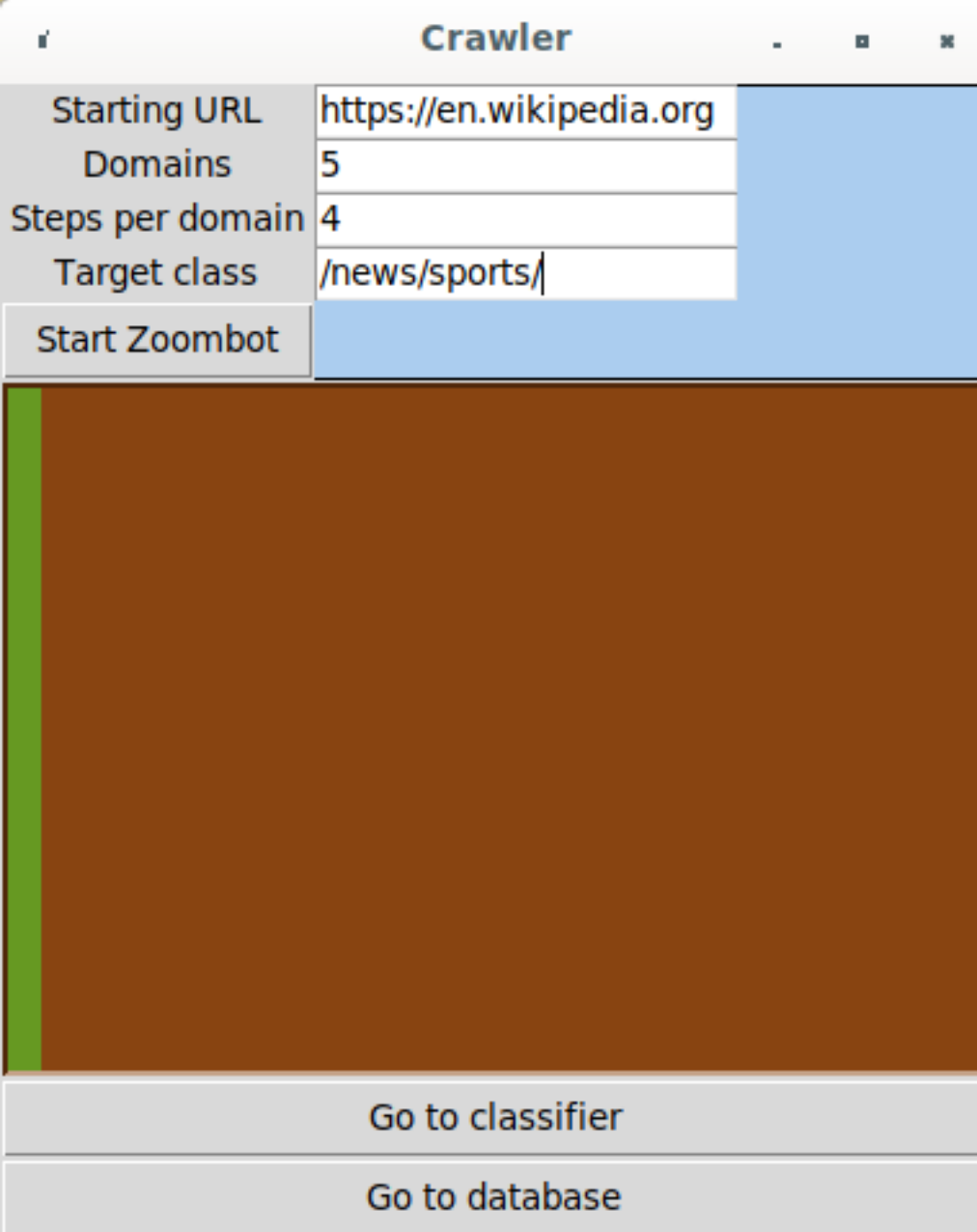
⁵Pokus o vložení ne-textového souboru způsobí chybu, která z nejasných důvodů nejde zachytit běžným handler-case (jak je zachytávána většina vadných vstupů; uživateli se pak pouze zobrazí chybová hláška).

3.5 Import/export

Zpracovaná data lze snadno exportovat, čímž odpadne návaznost na zdrojové dokumenty. Data jsou uložena ve dvou souborech: DATA/urls a DATA/corpora. Pokud tyto dva soubory dáme do společné složky, kterou umístíme do složky DATA/imports/, bude je systém brát jako pevně daný dodatečný korpus. Ten již nelze dále upravovat (zobrazuje se deaktivovanými tlačítky), zato se rychleji načte. Importových složek lze mít ve složce DATA/imports/ libovolné množství.

Výsledek crawleru lze použít jako databázi. Pokud přejmenuje uživatel soubor discovered na urls, program při spuštění všechny soubory v něm stáhne. Uživatel může crawler opravit, pokud něco umístil špatně, vytvořit export a vrátit se ke svému původnímu datasetu (který musí před přejmenováním zálohovat).

Obrázek 3.2: Crawler



Starting URL	https://en.wikipedia.org
Domains	5
Steps per domain	4
Target class	/news/sports/
Start Zoombot	

Go to classifier

Go to database

4. Závěr

Zadáním práce bylo vytvořit jednotný systém umožňující tři základní funkce:

4.1 Tvorba hierarchického systému

Soubory je *možné* vkládat z různých zdrojů, ačkoliv v praxi jsem uplatnil jen stahování z internetu, protože systém je určen pro klasifikaci stránek, nikoli náhodných souborů.

4.2 Klasifikace daného textu

Systém umožňuje jak přímo zjistit umístění textu, tak získat podrobná data o důvodech, z jakých byl text umístěn, kam byl.

4.3 Crawler

Myšlenku crawleru jsem v zadání úmyslně zlehčil ("jednoduchý crawler"), protože jsem neměl předem představu o tom, jak dobře může crawler dopadnout. Formálně vzato systém crawleru existuje, zajímavější je ale otázka, jestli je efektivní v praktickém účelu.

4.3.1 Test crawleru

Abych to zjistil, nechal jsem crawler běžet na základě stejných dat za stejných podmínek ze stejné počáteční stránky s dvěma různými cíli. Počáteční data jsou v release "5-4-2021", počáteční stránkou je "https://en.wikipedia.org", crawler hledá 160 domén, 4 stránky na doménu¹.

První crawl

Skóre stránek, na kterých cestou byl, je rozloženo takto:

¹Tato hodnota se osvědčila: více stránek je zbytečně pomalé, pro méně stránek hrozí buď to, že dojde fronta, nebo to, že se crawler nedostane na zajímavá místa uvnitř domény. Každopádně doporučuji rozmezí 2-5.

Počet správných úrovní	0	1	2	3	4	5
Počet stránek	130	20	225	137	33	10

Tabulka 4.1: Skóre s cílem /articles/stem/science/physics/thermodynamics/

Počet správných úrovní	0	1	2	3	4	5
Počet stránek	178	92	52	11	4	245

Tabulka 4.2: Skóre s cílem /articles/life/recreation/art/fiction/

Crawl trval 6 hodin a 50 minut. 17 nalezených stránek bylo duplikátních, jen s rozdílným query v adrese, což obecně může znamenat zcela jiný obsah, ačkoliv zde šlo o irelevantní požadavky.

V tabulce je vidět, že se crawler většinu času dokázal udržet kolem /articles/stem/science/, ale pak měl problémy. Z toho vyplývá, že rozdělení science je špatně navrženo. S tím nelze nesouhlasit, protože kategorie nature obsahuje pouze podkategorie math a physics. Chybí jí vlastní povaha, jde jen o nádobu pro dvě jiné kategorie. Pro další využití jde o signál, že tato část potřebuje jiné řešení.

V příloženém souboru crawl.ods na listu 1 je seznam všech stránek a jejich skóre, aby si mohl čtenář dokumentace udělat představu o výsledcích.

Druhý crawl

Crawl trval 9 hodin.

38 stránek bylo duplikátních. Tento crawl byl mnohem úspěšnější, až 42 % prošlých stránek bylo podle crawleru relevantních. Je to také proto, že jde o mnohem širší kategorii a zároveň o kategorii širěji zastoupenou v osobních stránkách, kterých je na internetu nevyčerpatelná zásoba. Naznačuje to, že by se měla tato kategorie prohloubit a zaplnit různými podkategoriemi, které se nabízejí podle žánrů² nebo podle typu obsahu (recenze, fanpage, diskuze o spisovatelství...). Všechny výsledky jsou v příloženém souboru crawl.ods na listu 2.

4.4 Shrnutí

Nejobecnějším ponaučením z práce na projektu je obrovská užitečnost databáze issues, do které jsem ukládal všechny chyby a všechna plánovaná a probíhající vylepšení programu. Github umožňuje issues uzavřít, ale kdykoli se k nim vrátit, což je užitečné, pokud se chyba vrátí v jiném kontextu. Seznam také pomohl při tvorbě dokumentace. Konkrétněji pro projekt samotný: domnívám se, že problém je na správné straně hranice nemožnosti, a že crawler to doložil, i když dataset, který jsem vytvořil, není zdaleka všepokrývací, spíše jen demonstrační.

²I přesto, že trénovací data byla pouze literární, do kategorie spadla řada filmů. To vzhledem k obecnější kategorii fikce nevadí, ale je to zajímavé zobecnění. Klasifikace někdy v případě, že cílová kategorie neexistuje, uvažuje metaforicky: přiřadí astrologii k historii náboženství nebo k vesmíru, venkovní sporty zařadí pod počasí. Nedokáže totiž rozlišit mezi kategoriemi, pro které nemá trénovací data.

Seznam obrázků

2.1	Vysvětlení klasifikace	6
2.2	Možnost výběru další úrovně	7
3.1	Databáze	16
3.2	Crawler	18

Seznam tabulek

4.1	Skóre s cílem /articles/stem/nature/physics/thermodynamics/	19
4.2	Skóre s cílem /articles/life/recreation/art/fiction/	20