



EduAkademia.pl

prace naukowe na zlecenie

Praca-magisterska-przykład-5

PRACA MAGISTERSKA

AUTOMATYCZNA KLASYFIKACJA, GRUPOWANIE I

ZNAJDOWANIE UOGÓLNIEN DLA J EZYKA POLSKIEGO

OSWIADCZENIE AUTORA PRACY

SWIADCZAM SWIADOMY ODPOWIEDZIALNOSCI KARNEJ ZA POSWIAD
CZENIE NIEPRAWDY,

ZE NINIEJSZĄ PRACĄ, DYPLOMOWĄ, WYKONAŁEM

OSOBISCIE I SAMODZIELNIE,
I NIE KORZYSTAŁEM ZE ŹRÓDEŁ INNYCH

NIZ WYMIENIONE W PRACY

.....

PODPIS

MASTER OF SCIENCE THESIS

AUTOMATIC CLASSIFICATION, CLUSTERING AND
SEARCHING FOR GENERALIZATIONS FOR POLISH
LANGUAGE.

SUPERVISOR:

Adrian Horzyk Ph.D

Krakow 2011

Serdecznie dziękuję promotorowi, rodzinie i przyjaciołom za wsparcie podczas pisania tej pracy.

Spis treści

| | | |
|------|-----------------------------|---|
| 1. | Wstęp..... | |
| | | 6 |
| 1.1. | Cele pracy | |
| | | 6 |
| 1.2. | Zawartość pracy | |
| | | 6 |
| 2. | Wprowadzenie | |
| | | 8 |
| 2.1. | Czym jest lingwistyka?..... | |
| | | 8 |

| | | |
|--------|---|----|
| 2.2. | Lingwistyka komputerowa..... | 8 |
| 2.3. | Historia lingwistyki | 9 |
| 2.4. | Problemy i cele współczesnej lingwistyki komputerowej | 10 |
| 2.5. | Pajak, internetowy..... | 11 |
| 2.6. | Ekstrakcja danych | 14 |
| 3. | Relacje między wyrazami w języku polskim | 17 |
| 3.1. | Części mowy..... | 17 |
| 3.2. | Relacje między wyrazami..... | 18 |
| 4. | Opis własnego rozwiązania, | 21 |
| 4.1. | Założenia projektowe..... | 22 |
| 4.2. | Architektura systemu | 22 |
| 4.2.1. | Interfejs użytkownika' | 23 |

| | |
|--|----|
| 4.2.2. Pajak ₂ internetowy i ekstraktor danych | 25 |
| 4.2.3. Biblioteka CLP..... | 26 |
| 4.2.4. Specjalistyczny graf przyzwyczajen ¹ lingwistycznych | 27 |
| 4.3. Opis działania aplikacji..... | 29 |
| 4.3.1. Tryb online | 29 |
| 4.3.2. Tryb offline..... | 32 |
| 4.3.3. Symulator | 32 |
| 4.4. Generacja specjalistycznego grafu LHG | 33 |
| 5. Prezentacja działania aplikacji | 36 |
| 6. Podsumowanie wyników pracy i wnioski..... | 45 |
| 6.1. Podsumowanie | 45 |
| 6.2. Dalsze możliwości | |

| | |
|---|----|
| sci' rozwoju | 46 |
| 6.3. Porównanie z rozwiązaniami, o podobnej tematyce..... | 47 |
| 7. Zakonczenie' | 50 |

5

1. Wstę

Komunikacja oraz przepływ informacji droga, elektroniczna, sa, nieodzownym elementem współcze-

snego zycia'. W wielu sytuacjach niezwykle uzyteczna' jest automatyzacja tego procesu, nie tylko pod katem, wydajnosci' i bezpieczenstwa, ale takze' poprawnosci' lingwistycznej. O ile w przypadku gotowych szablonów wiadomosci' poprawnos'c' nie jest problemem, to gdy wymagana jest wi,eksza interaktywnos'c', pojawiaja, si,e kłopoty.

Rozmaite inteligentne czatboty posiadaja, wbudowane mechanizmy, które umozliwiaja, im wysoce in-

teraktywna, komunikacj,e z uzytkownikiem, która jednak cz,esto nie jest pozbawiona rozlicznych bł,edów zarówno składniowych, jak i logicznych. Niniejsza praca podejmuje tematyk, automatycznego grupo-wania i klasyfikacji wyrazów w j,ezyku polskim w celu modelowania uogólnien' i przechodniosci' właści'-wosci' obiektów w hierarchii słów.

Zakres pracy obejmuje zbudowanie specjalistycznego grafu lingwistycznego, agregujacego, słowa w

zależności od stopnia uogólnienia pojęć, które reprezentują, na podstawie ich automatycznej klasyfikacji oraz grupowania. Do zbudowania tego grafu powinny być wykorzystane różne korpusy tekstów pozyskiwane z baz danych oraz Internetu. Celem pracy jest zbudowanie takiego grafu oraz pokazanie możliwości przenoszenia właściwości obiektów przez hierarchię klas tego grafu - modelując możliwości uogólniania, jakie zachodzą w ramach języka.

1.1. Cele pracy

Celem pracy jest zbudowanie specjalistycznego grafu przyzwyczajen lingwistycznych (ang. Language Habits Graph, LHG)[9], który ukazuje relacje pomiędzy wyrazami. W szczególności skupia się

on na uogólnianiu i na jego podstawie klasyfikowaniu słów oraz wyróżnianiu ich wspólnych właściwości. Zadaniem było stworzenie pajaka internetowego, który przeszukując wyniki z wyszukiwarki

internetowej buduje bazę wiedzy, wokół określonego przez użytkownika słowa. Na podstawie uzyskanych w ten sposób informacji, program tworzy sieć relacji dla podanego wyrazu. Sieć ta reprezentowana jest na grafie LHG w uproszczonej postaci, z uwzględnieniem siły (częstotliwości występowania) określonych połączeń. Zamodelowana została w ten sposób zhierarchizowana struktura, która obrazuje grupowanie obiektów, jak i przechodniosć cech wspólnych wewnątrz danej klasy.

1.2. Zawartość pracy

1. W rozdziale 1 znajduje się wstęp zawierający krótką informację o tematyce i celu pracy.

2. Rozdział 2 opisuje zagadnienie lingwistyki komputerowej i ukazuje problematykę pracy w szerszej perspektywie.

3. Zagadnienia związane z relacjami występującymi między wyrazami w języku polskim zostały opisane w rozdziale 3.

4. Rozdział 4 zawiera szczegółowy opis własnego rozwiązania.

5. Rozdział 5 przedstawia wyniki działania aplikacji na konkretnych przykładach.

6

1.2. Zawartość pracy

7

6. W rozdziale 6 zamieszczone zostało podsumowanie wyników pracy oraz przemyślenia i wnioski.

7. Rozdział 7 jest zakończeniem, w którym zestawione są wyniki pracy z zamierzonymi celami.

Ł. Wróbel

2. Wprowadzenie

2.1. Czym jest lingwistyka?

Definicja 1 Lingwistyka (językoznawstwo) jest to dział nauk humanistycznych badający, istotę, budowę i rozwój języka[22].

Język mówiony i pisany od wieków był przedmiotem badań i rozważań naukowych. Większość wysiłków koncentrowała się jednak na zgłębianiu jego ewolucji poprzez analizę dzieł literackich, czy też

poszukiwaniu podobieństw pomiędzy różnymi językami i dialektami. Dopiero w XX wieku zaczęło się kształtować bardziej ustrukturalizowane i sformalizowane podejście do tego zagadnienia[22].

Zasadniczo lingwistyka dzieli się na teoretyczną i stosowaną. Językoznawstwo teoretyczne skupia się na badaniu aspektów współczesnej odmiany konkretnego języka i wyszukiwaniu oraz opisywaniu jego ogólnych zasad. Przedmiotem rozważań mogą być też ogólne reguły i aspekty wspólne dla wszyst-

kich języków. Stosowana lingwistyka podejmuje problemy związane z wykorzystaniem wyników tych rozważań w innych dziedzinach. Bada język w szerszej perspektywie zarówno w porównaniu z innymi, jak i jego własnym kształtowaniem się przez lata[22].

Współczesne językoznawstwo korzystające z najnowszej technologii dało początek informatyzacji

tej nauki. Otworzyła ona nie tylko nowe możliwości, ale i liczne wyzwania związane z elektroniczną

komunikacja między ludźmi, a przede wszystkim człowieka z komputerem. Tak narodziła się lingwistyka komputerowa.

2.2. Lingwistyka komputerowa

Definicja 2 Lingwistyka komputerowa to dział lingwistyki używający modeli komputerowych w celu testowania hipotez dotyczących mowy i języka oraz tworzenia programów komputerowych przetwarzających język naturalny[23].

Ta ogólna definicja ukazuje jak bardzo lingwistyka komputerowa zakorzeniona jest w rozmaitych dziedzinach językoznawstwa. Zakres jej zastosowań jest jednak o wiele większy, gdyż obecna jest także między innymi w filozofii, psychologii, informatyce (np. sztuczna inteligencja) czy matematyce (np. rachunek prawdopodobieństwa, teoria automatów, logika matematyczna, statystyka, czy języki formalne)[23].

Dziedzinę tę można podzielić na część teoretyczną i praktyczną. Zagadnienia teoretyczne skupiają się na badaniu możliwości automatyzacji procesów związanych z klasyfikacją podzbiorów językowych, obliczaniu wydajności tych procesów, a także tworzeniu ich formalnych opisów. Praktyczne podejście do lingwistyki komputerowej owocuje natomiast powstawaniem algorytmów i złożonych modeli automatyzujących i przetwarzających, które znajdują zastosowanie w rozmaitych aspektach językowych[23].

Rozwój i cele stojące przed lingwistyką komputerową są ściśle związane z rozwojem technologii,

jak i z rosnącymi potrzebami automatyzacji i formalizacji procesów językowych, które wykraczają poza możliwości człowieka. Coraz nowocześniejsze maszyny z dużą mocą obliczeniową są w stanie realizować nawet bardzo złożone i jeszcze niedawno zbyt kosztowne algorytmy, co z kolei sprawia, że zagadnienia związane ze sztuczną inteligencją są obecnie bardzo popularne.

8

2.3. Historia lingwistyki

9

2.3. Historia lingwistyki

Aby dobrze zrozumieć jak wielki postęp dokonał się w dziedzinie lingwistyki komputerowej warto zagłębić się w historię jej rozwoju. Daje to możliwość zaobserwowania trendów i zagadnień, które były, są i będą przedmiotem badań współczesnego językoznawstwa informatycznego.

Początki lingwistyki jako nauki o języku sięgają czasów starożytnych. Za jej kolebkę uważa się In-

die, gdzie już w pierwszym tysiącleciu p.n.e. dokonywano tłumaczeń starożytnych tekstów religijnych. Za twórcę pierwszej gramatyki uważa się Paniniego, który około V-IV w. p.n.e. stworzył szereg pojęć i definicji oraz około 4 tysiące reguł językowych. Niektóre powstałe wtedy terminy takie jak rozróżnienie

między głoskami, a literami, fleksja, słowotwórstwo, struktura wyrazu, czy klasyfikacja części

mowy używane są do dnia dzisiejszego. Wszelkie wyjątki, które licznie pojawiają się w językach, takich jak łacina, czy greka Panini w swojej gramatyce opisał za pomocą bardzo szczegółowych reguł, które opisywały nawet pojedyncze przypadki[14].

W starożytnej Grecji i Rzymie skupiano się bardziej na ogólnych problemach języka, takich jak jego istota, stosunek do logiki, retoryki i poetyki. Ukształtowały się dwa odmienne punkty widzenia. Ana-logizm, który zakładał, że język ma charakter logiczny oraz anomalizm, który uznawał że język jest zmienny. Z tego okresu wywodzi się pojęcie etymologii (nauka o pochodzeniu i ewolucji znaczenia oraz formy wyrazów), teorii nominacji (nazywania świata za pomocą słów), onomatopei (dźwiękonasładowictwa). Szczególnie rozwinęło się także językoznawstwo opisowe. Powstał podział na części mowy, opisano deklinacje, koniugacje i zasady opisu składni[14].

Dopiero w XVIII w. gramatyka języka łacińskiego przestała być powszechnie uważana za uniwersalną, do czego przyczyniły się badania Abu Hayana at Tauhaidiego, który postulował, że na proces

kształtowania języka mają wpływ kultura i historia narodu. Zapoczątkowało to pojawianie się filologii narodowych: m.in. francuskiej, niemieckiej, angielskiej i słowiańskiej[14].

XIX w. to rozwój językoznawstwa historycznego, który rozpoczął odkrycie przez William'a Jones'a

sanskrytu (starożytnego języka Indii). Jego znaczne podobieństwo do języków europejskich zainicjowało badania nad komparatywistyką (lingwistyką porównawczą). Określone zostały regiony i narody, których formy komunikacji mają wspólne korzenie historyczne, tworząc tak zwane rodziny językowe[20].

Wiek XX to okres strukturalizmu i formalizacji zasad językowych. Zamiast do genezy zaczęto przywiązywać większą wagę do opisu struktury. F. de Saussure uważał, że każdy znak zbudowany

jest z dwóch elementów: znaczącego (jego reprezentacja np. obraz, czy słowo) i znaczonego (poję-

cie z nim związane). Język zdefiniowano więc jako dwuklasowy system semiotyczny, który służy do komunikacji[19]. Powstał także pierwszy matematyczny opis języka, który za pomocą formalnych narzędzi matematycznych takich jak dowody, twierdzenia i rachunki stworzył drzewiastą strukturę wszystkich wyrazów poprawnych w danej gramatyce. N. Chomsky wprowadził gramatykę transformacyjno-generatywną, która jest systemem określającym skończoną ilość reguł. Z ich pomocą rozmówca jest w stanie utworzyć poprawne zdanie, którego nigdy wcześniej nie słyszał, a odbiorca stwierdzić, czy zdanie, które usłyszał jest poprawne w danym języku, czy nie[13].

Formalizacja zasad lingwistycznych stworzyła podwaliny pod powstanie językoznawstwa informa-

tycznego, które wykorzystując ściśle określone reguły umożliwiło automatyczną analizę i przetwarzanie tekstu, wyszukiwanie, a także budowanie poprawnych form językowych przez programy komputerowe. Początkowo w lingwistyce komputerowej dominowały dwa podejścia. Pierwsze oparte było na metodach symbolicznych (teoria generatywna, sztuczna inteligencja), natomiast drugie bazowało na metodach statycznych takich jak analiza tekstu, rozpoznawanie liter itp[6]. Powstały modele statystyczne służące do syntezy mowy, języki programowania oparte o logikę formalną i metody wnioskowania (np. Prolog), a także aplikacje sterowane komendami w języku naturalnym.

Lata 90 to dynamiczny rozwój Internetu oraz mocy obliczeniowej i pojemności pamięciowej komputerów. Pojawiają się złożone aplikacje wykorzystujące zaawansowane algorytmy oparte m.in. o statystykę. W

odpowiedzi na zapotrzebowanie na programy zajmujące się ekstrakcją danych z Internetu,

czy korpusów tekstów oraz narzędzi przetwarzających i analizujących tekst z wykorzystaniem metod sztucznej inteligencji powstało wiele rozwiązań komercyjnych[6].

Ł. Wróbel

2.4. Problemy i cele współczesnej lingwistyki komputerowej

10

2.4. Problemy i cele współczesnej lingwistyki komputerowej

Największym problemem, przed jakim staje współczesne językoznawstwo komputerowe jest wieloznaczność[6]. Dotyczy on fonetyki (homofonia - fonetyczna tożsamość różnych znaków językowych),

morfologii (homonimia - posiadanie różnych znaczeń przez identyczne formy językowe), składni (niejednoznaczność przy określaniu części zdania), semantyki (idiomy, polisemia - posiadanie przez jedno

słowo wielu znaczeń), pragmatyki (metafory, ironia). O ile człowiek korzystając z wyższego poziomu języka może w łatwy sposób rozwiązać te problemy, to już komputer, ma z tym spore kłopoty.

Według Bonnie Webber (2001) dwa podstawowe cele jakie stoją przed lingwistyką informatyczną to[11]:

modelowanie ludzkiego rozumienia i generacji języka naturalnego jako systemu procesów przetwarzających informacje (lingwistyka informatyczna)

wyposażenie komputerów w mechanizmy analizy i generowania języka naturalnego w celu do-

starczenia użytecznej usługi (stosowane przetwarzanie języka naturalnego - ang. applied natural language processing, inżynieria języka naturalnego - ang. natural language engineering, technologia językowa - ang. language technology)

Roland Hausser (2001) uważa, że lingwistyka komputerowa powinna się skupiać na modelowaniu komunikacji człowieka z komputerem, konstruując modele wyjaśniające naturalny przekaz informacji w sposób spójny funkcjonalnie, precyzyjny matematycznie i efektywny obliczeniowo[11].

Z punktu widzenia współczesnych wymagań najintensywniejsze prace prowadzone są nad[11]:

aplikacja umożliwiającą wyszukiwanie specyficznych informacji na podstawie dokładnie opisanych wymagań co do rezultatów

systemem potrafiącym uczyć się ze źródeł przyswajalnych przez człowieka, a następnie potrafiącym

rozwiązać test sprawdzający zdobytą wiedzę, jaki rozwiązałby człowiek

systemem do automatycznego tłumaczenia na inne języki naturalne, a także języki nieistniejące, ale opisane określonym zestawem reguł

aplikacja umożliwiającą komunikację głosową człowieka z komputerem

Niewyczerpanym źródłem wiedzy dla programów wykorzystujących sztuczną inteligencję do wyszukiwania informacji jest Internet. Ogromu zawartych w nim informacji człowiek nie jest w stanie przetworzyć, natomiast posiadająca odpowiedni zestaw reguł oraz sposób wartościowania i ekstrakcji danych maszyna może podjąć się tego zadania. O ile już samo stworzenie narzędzia do inteligentnego pozyskiwania informacji nie jest zadaniem prostym, to dodatkowe trudności pojawiają się podczas przetwarzania danych. System, aby wykorzystać pozyskaną wiedzę w określonym celu, musi prawidłowo rozumieć zarówno posiadane dane, jak i cel ich wykorzystania. W szczególności stanowi to ogromny problem, gdy zadaniem jest przetłumaczenie tekstu na inny język. Znajomość reguł tworzenia poprawnych form nie gwarantuje wierności przekazu. Przykładem może być napisana przed wieloma laty rosyjska aplikacja, która tłumaczyła wyrażenia z języka angielskiego na rosyjski i później znów na angielski. Zamieniła ona

wyrażenie „out of sight, out of mind” (co z oczu, to z serca) na „invisible idiot” (niewidzialny idiota)

-out of mind zostało przetłumaczone jako „niespełna rozumu”, gdyż system nie zrozumiał semantyki całego wyrażenia[1].

Cele bezpośrednio związane z przedmiotem tej pracy to wyszukiwanie i wyodrębnianie informacji, ich analiza i na jej podstawie ukazanie zasad modelowania wybranej dziedziny języka. Następnie podrozdziały to ogólny opis narzędzi i technik wykorzystywanych do rozwiązywania tych problemów wraz z przykładami istniejących implementacji.

Ł. Wróbel

2.5. Pająk internetowy

11

2.5. Pająk internetowy

Definicja 3 Pająk (robot) internetowy (ang. Web crawler, spider bot) to program przeszukujący zawartość sieci Internet w zorganizowany i zautomatyzowany sposób[24].

Rysunek 2.1: Schemat działania pajaka internetowego

Podstawowym zadaniem wszelkiego rodzaju robotów internetowych jest automatyczna nawigacja po

stronach internetowych w sposób określony przez zadane reguły. Każdy pajak ma początkowo określony pulę adresów, od której zaczyna przeglądanie. Odwiedzając kolejne strony przenosi adresy z tej listy na listę adresów już odwiedzonych. Na witrynie, przez którą przechodzi może szukać nowych linków, które dodaje do listy „do odwiedzenia”, pod warunkiem, że nie ma ich jeszcze na żadnej z list. Robot działa dopóki na liście „do odwiedzenia” znajduje się przynajmniej jeden adres. Specyficzne pajaki mogą oczywiście

Ł. Wróbel

2.5. Pajak internetowy

12

wiecie przechodzić przez te same strony wielokrotnie, lub działać w nieskończonej pętli, aby okresowo patrolować określone witryny[24].

Po wejściu na stronę robot sprawdza znajdujący się na serwerze plik robots.txt, który określa do jakiej części serwisu ma dostęp. Zawartość pliku to spis plików i katalogów opisana za pomocą protokołu Robot Exclusion Protocol, niedostępnego dla pajaka. Administrator strony może w ten sposób ukryć część serwisu przed wszelkimi automatami, lub też całkowicie odmówić im dostępu, aby na przykład nie obciążały łącza. Plik ten jest jednak publicznie dostępny, więc zapisywanie w nim informacji na temat niepublicznych danych nie jest dobrym pomysłem. Roboty internetowe (w szczególności wszelkiego rodzaju malware, harvestery i spam boty) mogą po prostu ignorować plik robots.txt, więc nie jest on zabezpieczeniem absolutnym. Informacje dla pajaka mogą znajdować się także w meta tagach, które określają zasady zachowania robota na danej stronie. Najczęściej stosowane tagi to NOINDEX

(uniemożliwia indeksowanie zawartości strony) i NOFOLLOW (uniemożliwia bezpośrednie podążanie za linkami znajdującymi się na stronie)[16].

Pajaki internetowe oprócz odwiedzania witryn służą przede wszystkim do zbierania informacji na nich umieszczonych. Wyszukiwarki internetowe wykorzystują rozmaite boty do indeksowania i sprawdzania aktualnej zawartości stron w sieci, aby wykorzystać uzyskane w ten sposób dane do zwrócenia jak najlepszej listy wyników. Często na potrzeby pobierania danych innych niż linki sprawdzany jest także język strony, czy kodowanie znaków, opisane w nagłówku, aby niepotrzebnie nie czytać całej treści[24]. Kilka przykładowych istniejących implementacji robotów internetowych:

Googlebot

Jest to pajak wykorzystywany przez popularną wyszukiwarkę do indeksowania stron w internecie. Działa on w dwóch wersjach: deep crawl (wędruje po witrynach gromadząc i odwiedzając jak najwięcej znalezionych adresów) oraz fresh crawl (odwiedza zmieniające się często strony, w celu

aktualizacji związanych z nimi informacji). Aby Googlebot mógł odwiedzić daną stronę, musi istnieć do niej link na innej, przez którą już przechodził. Gromadzi on także statystyki takie jak ilość wystąpienia adresu danej strony, która następnie określa jej pozycję w wynikach wyszukiwa-

nia. Największym problemem związanym z Googlebotem jest wysokie wykorzystanie transferu, co może

doprowadzić nawet do czasowego zawieszenia strony z powodu przekroczenia ograniczeń transferowych. Google udostępnia administratorom narzędzie Webmaster Tools, za pomocą którego można ustawić m.in. częstotliwość wizyt pajaka na stronie[22].

WebCrawler

Metawyszukiwarka internetowa agregująca najlepsze wyniki z google, yahoo, ask.com, bing search i innych popularnych wyszukiwarek. Umożliwia wyszukiwanie także obrazków, plików au-

dio i wideo. WebCrawler był pierwszą wyszukiwarką umożliwiającą szukanie w pełnym tekście dokumentów, a nie tylko indeksowanie[25].

Rysunek 2.2: Graficzny interfejs metawyszukiwarki WebCrawler.

Ł. Wróbel

2.5. Pajak internetowy

13

IBM Web Fountain

Web Fountain to jeden z pierwszych projektów, których celem jest skatalogowanie i interpretacja nieustrukturalizowanych lub częściowo ustrukturalizowanych danych tekstowych z Internetu. Jego głównym zadaniem jest poszukiwanie wzorców i zależności na potrzeby statystyczne oraz prezentacji trendów w czasie rzeczywistym[21].

Rysunek 2.3: Graficzny interfejs aplikacji Web Fountain.

WGet

Jest to program pobierający zasoby z serwerów internetowych będący częścią projektu GNU. Wspiera on obecnie transfer za pomocą protokołów HTTP, HTTPS i FTP. Umożliwia pobiera-

nie rekursywne, konwersję linków umożliwiającą przeglądanie lokalnych plików HTML w trybie offline oraz wspiera proxy. WGet był jednym z pierwszych programów wykorzystujących nagłówki HTTP do wznowiania przerwanych transferów. Pajak internetowy wykorzystywany jest przede wszystkim do pobierania rekursywnego, które dokładnie odwzorowuje strukturę zasobów na serwerze. Domyślnie WGet bierze pod uwagę informacje zawarte w pliku robots.txt, ale można go też uruchomić z parametrem -e robots=off. Istnieje także możliwość pobierania tylko plików nowszych niż zapisane lokalnie, a także filtrowania zasobów po nazwach za pomocą wyrażenia regularnych[2].

Ł. Wróbel

2.6. Ekstrakcja danych

14

Rysunek 2.4: Wynik działania programu WGet.

Crawler4j

Jest to prosty, open source'owy interfejs w javie, który pozwala bardzo łatwo i szybko stworzyć własnego pajaka internetowego. Udostępnia dwie metody do przeładowania: shouldVisit (określa kryteria definiujące, czy dany URL powinien zostać odwiedzony, czy nie) i visit (określa co ma zostać zrobione w momencie odwiedzenia strony). Dodatkowo należy zdefiniować kontroler, który

określi początkowy zbiór adresów, miejsce zapisu danych, czy ilość równoległych wątków[3].

2.6. Ekstrakcja danych

Znalezienie strony, na której mogą być zamieszczone poszukiwane informacje, to dopiero połowa procesu ich pozyskiwania. Są one zazwyczaj otoczone kodem html, który ułatwia użytkownikowi ich przeglądanie z użyciem przeglądarki, ale komplikuje proces ekstrakcji właściwego tekstu przez program.

Wydobywanie danych może odbywać się na dwa sposoby. Pierwszy zakłada, że użytkownika interesuje konkretna informacja, której program poszukuje z wykorzystaniem np. słów kluczowych, czy wyrazów

regularnych. Drugi sposób to ekstrakcja całości tekstu w celu zbudowania przez użytkownika bazy wiedzy, z punktu widzenia której istotne jest pozyskanie jak największej ilości danych. Baza ta wykorzystuje między innymi data mining, który na jej podstawie wyszukuje trendy i formułuje predykcje przyszłych zmian.

Aby uzyskane w ten sposób informacje mogły zostać dalej przetwarzane nie mogą zawierać elementów języka html. Niezbędny jest zatem parser, który znajdzie i usunie wszystkie tagi html, czy skrypty

języka JavaScript zostawiając tylko istotne informacje. Teoretycznie możliwe jest wykorzystanie wyra-

żenia regularnego, znajdującego niepozadane elementy, lecz może ono nie uwzględnić czystych błędów w strukturze stron, jak np. znak < wewnątrz taga, czy też fragmenty tekstu przypominające tagi (równa-

nia matematyczne, kod w różnych językach programowania).

Przykłady ekstraktorów danych i parserów html:

Online Data Extractor

Jest to program, który pozwala na wyszukiwanie danych kontaktowych głównie na potrzeby firm, takich jak adresy e-mail, numery telefonu, czy fax, a także adresy URL i meta tagi). Informacje

te znajdują się z pomocą pajaka internetowego przeglądającego określone strony, a także ko-

rzystającego z pomocy popularnych wyszukiwarek. Umożliwia eksport zgromadzonych danych między innymi do plików txt, xls, csv i htm, a także pozwala na zaawansowaną integrację danych do celów data miningu[10].

Ł. Wróbel

2.6. Ekstrakcja danych

Rysunek 2.5: Zrzut ekranu obrazujący działanie aplikacji Online Data Extractor.

HTML Text Extractor

HTML Text Extractor to komercyjne narzędzie działające pod Windowsem, które umożliwia wydobycie kodu html dowolnej strony internetowej, w tym także tych, w których zablokowany jest podgląd źródła, czy też są zaszyfrowane. Ekstrakcja następuje w momencie odwiedzenia strony i daje możliwość wyodrębnienia, zaznaczenia i skopiowania samego tekstu, czy to z całej witryny, czy też jej fragmentu[4].

Rysunek 2.6: Zrzut ekranu obrazujący działanie aplikacji HTML Text Extractor.

Ł. Wróbel

2.6. Ekstrakcja danych

16

Na potrzeby aplikacji zajmujących się przetwarzaniem i analizą języka naturalnego, a także pogłębiania wiedzy na temat kultury i literatury tworzy się korpusy tekstów, które zawierają tekst pokazujący wyszukiwane wyrazy i konstrukcje językowe zastosowane w rzeczywistym kontekście. Często też do

korpusów udostępnionych online dołączone są mniej lub bardziej zaawansowane narzędzia wyszukiwające. Popularne źródła korpusów tekstów to m.in.:

PELCRA

Projekt PELCRA realizowany przez Katedrę Języka Angielskiego Uniwersytetu Łódzkiego przy współpracy z Departamentem Językoznawstwa i Współczesnego Języka Angielskiego Uniwersytetu w Lancaster opracowuje korpusy tekstów w języku polskim i angielskim na potrzeby badań m.in. nad leksykografią, translatoryką i lingwistyką (zarówno ogólną jak i komputerową)[12].

Korpusy wchodzące w skład projektu PELCRA:

– Narodowy Korpus Języka Polskiego - Jeden z największych i najpopularniejszych korpusów języka polskiego współtworzony przez Instytut Podstaw Informatyki PAN, Uniwersytet Łódzki, Wydawnictwo Naukowe PWN, oraz Instytut Języka Polskiego PAN.

– Korpus referencyjny języka polskiego PELCRA - Zawiera on około 93 miliony segmentów słów tekstu. Jego częścią jest dostępny zarówno w wersji tekstowej jak i audio Korpus Polszczyzny Konwersacyjnej, który składa się z nagranych rozmów przypadkowych osób uzupełnionych o dane dotyczące uczestników konwersacji.

– ENHIG - Korpus historyczny zawierający teksty staroangielskie, oznakowane pod względem składniowym próbki poezji, prozy i tłumaczeń stworzony na potrzeby prowadzenia badań nad szykiem zdania.

Rysunek 2.7: Kolokacje wyrazu niebo na podstawie Narodowego Korpusu Języka Polskiego.

IPI PAN

Korpus Polskiej Akademii Nauk zawiera obecnie około 250 milionów segmentów i dostępny jest darmowo. Do wyszukiwania służy dostępny na licencji GPL, dedykowany program Poliqarp[5].

Wydawnictwo Naukowe PWN

Wydawnictwo PWN udostępnia swój korpus językowy odpłatnie. Zawiera on ponad 40 milionów słów z fragmentów 386 różnych ksiązek, 977 numerów 185 różnych gazet i czasopism, 84 nagranych rozmów, 207 stron internetowych oraz kilkuset ulotek reklamowych. Wersja demonstracyjna dostępna darmowo składa się z 7.5 miliona słów[26].

Ł. Wróbel

3. Relacje między wyrazami w języku polskim

3.1. Części mowy

Głównie ze względów praktycznych wyrazy w języku polskim podzielono na klasy językowe zwane potocznie częściami mowy. Podział ten zależny jest od szeregu właściwości fleksyjnych, semantycznych

i składniowych leksemów. Z punktu widzenia morfologii istotne jest, czy przyporządkowane do danej klasy wyrazy są odmienne i w jaki sposób. Klasyfikacja semantyczna wprowadza natomiast podział słów na samoznaczące (autosemantyczne) i nie mające znaczenia (synsemantyczne)[18].

Definicja 4 Części mowy to kategorie gramatyczno-leksykalne, które[18]:

mają te same cechy morfologiczne tzn. każdy leksem należący do danej części mowy ma podobną budowę, podobnie się odmienia np. rzeczowniki przez przypadki i liczby, przymiotniki przez rodzaje, przypadki i liczby,

pełnią identyczną funkcję składniową tzn. występują w roli orzeczenia, podmiotu, przydawki itd., mają podobne znaczenie lub niosą podobną treść ogólną np. nazywają zjawiska atmosferyczne, procesy, wielkości, cechy, liczby, stany itd.,

wchodzą w określone relacje (związki) z innymi wyrazami np. przysówek z czasownikiem tworzą związek przynależności.

Tradycyjny podział wyróżnia dziewięć części mowy[18]:

Rzeczownik Przysówek Przymiotnik Przyimek Liczebnik Czasownik Spójnik

Zaimek Partykuła

Wykrzyknik

Na potrzeby pracy wykorzystane zostały tylko te części mowy, które uznane zostały za niosące najwięcej istotnych informacji oraz pełniące ważne funkcje w zdaniu, a przede wszystkim bezpośrednio wykorzystywane w procesie kojarzenia. Z punktu widzenia kognitywistyki największe znaczenie mają wyrazy, z którymi bezpośrednio można skojarzyć określony obraz oraz takie, za pomocą których można obraz ten opisać.

17

3.2. Relacje między wyrazami

18

3.2. Relacje między wyrazami

Aby lepiej zrozumieć prawa, jakie rządzą relacjami między słowami w języku polskim, należy najpierw wprowadzić definicje pojęć takich jak zdanie i składnik, które są niezbędne z punktu widzenia analizy składniowej.

Definicja 5 Zdanie to odcinek tekstu od wielkiej litery do kropki lub znaku jej równoważnego (pytajnika, wykrzyknika). W odniesieniu do formy mówionej jest to pojedynczy odcinek intonacyjny.[27]

Definicja 6 Składnik to najmniejszy element, jaki da się wyodrębnić w procesie analizy składniowej, mogący pełnić jako całość funkcję składniową i wchodzić w związki z innymi takimi elementami.[27]

Najczęściej mamy do czynienia z sytuacją, w której składnik pokrywa się z wyrazem. Bywają jednak przypadki, gdy ta funkcja składniowa charakteryzuje połączenie wyrazów np. na pewno, niech będzie, Jan Kowalski. Takie połączenie nazywane jest składnikiem złożonym.

Wyróżniamy trzy zasadnicze kategorie grup syntaktycznych: podrzędne, współrzędne i

egzocentryczne.[15] W grupie podrzędnej członem nadrzędnym jest składnik, którego nie można

z niej usunąć, natomiast człon podrzędny niesie dodatkowe informacje, ale nie jest niezbędny z punktu

widzenia poprawności składniowej wyrażenia w zdaniu[15]. Przykładowo dla zdania Babcia upiekła pyszne ciasteczka. fraza pyszne ciasteczka jest grupą podrzędną, gdzie składnikiem nadrzędnym jest wyraz ciasteczka, a podrzędnym pyszne, gdyż zdanie Babcia upiekła ciasteczka. jest poprawne, a Babcia upiekła pyszne. już nie. Ten rodzaj grupy składniowej wykorzystany został w pracy, gdyż niesie on

najwięcej użytecznych informacji z punktu widzenia relatywnego opisywania rzeczywistości.

Grupa współrzędna zwana także szeregiem charakteryzuje się tym, że składniki są równorzędne

i można ją zredukować do jednego składnika (nie będącego spójnikiem)[15]. Przykładowo w zdaniu Babcia upiekła ciasteczka dla mnie i brata. szeregiem będzie fragment mnie i brata, gdzie współrzędnymi składnikami są wyrazy mnie oraz brata, a i jest spójnikiem. Poprawne są więc zdania Babcia upiekła ciasteczka dla mnie. oraz Babcia upiekła ciasteczka dla brata., ale Babcia upiekła ciasteczka dla i. już nie.

Trzecia grupa zwana egzocentryczną składa się z dwóch składników, z których jeden nie jest redukwalny, a składnik nadrzędny wyznacza formę podrzędnego[15]. Przykładem niech będzie zdanie Babcia upiekła dla nas ciasteczka., gdzie grupa egzocentryczną tworzą słowa dla nas, gdyż ani zdanie Babcia upiekła nas ciasteczka., ani Babcia upiekła dla ciasteczka. nie jest poprawne, a wyraz dla definiuje przypadek zaimka my.

Definicja 7 Akomodacja syntaktyczna (łac. accomodare - przystosowywać) oznacza dostosowanie się do siebie jednostek w zdaniu. Dostosowanie to polega najczęściej na przyjmowaniu konkretnych wartości kategorii gramatycznych (takich jak rodzaj, liczba, przypadek) przez jednostkę wchodzącą w konstrukcję z inną jednostką.[27]

Akomodacje syntaktyczne podzielić można na jednostronne, gdy składnik nadrzędny wymaga od

podrzędnego dostosowania formy fleksyjnej, wzajemne, gdy każdy ze składników ma wpływ na formę drugiego, a także międzyfrazowe, w których o formie podrzędnika decyduje nadrzędnik z innej frazy. Reguły zdefiniowane w ramach akomodacji syntaktycznych morfologicznych definiują warunki, które muszą być spełnione przez formy wyrazowe, aby mogły one istnieć wspólnie w jednym zdaniu[27]. Przykładowa taka reguła dla leksemów kot i isć wyglądała by jak na diagramie 3.1.

Wyraz nadrzędny (kot) wymusza na podrzędnym (isć) odpowiednią formę gramatyczną: liczbę, rodzaj i osobę. Pozostałe elementy takie jak czas, czy tryb zależą od źródła pozajęzykowego, jakim jest kontekst, do którego odniesiona jest wypowiedź. Uogólniając powyższy przykład można przedstawić ogólną regułę kształtującą zależność między podmiotem a orzeczeniem w zdaniu jak na diagramie 3.2.

Tego typu zależności akomodacyjnych można zaobserwować dla innych części mowy, jak pokazano na przykładowych diagramach 3.3 i 3.4.

Ł. Wróbel

3.2. Relacje między wyrazami

19

Rysunek 3.1: Przykład oddziaływania akomodacyjnego dla leksemów kot i isć.

Rysunek 3.2: Schemat uogólnienia oddziaływań akomodacyjnych między podmiotem i orzeczeniem.

Rysunek 3.3: Schemat uogólnienia oddziaływan' akomodacyjnych mi'edzy czasownikiem i rzeczowni-kiem.

Rysunek 3.4: Schemat uogólnienia oddziaływan' akomodacyjnych mi'edzy przymiotnikiem i rzeczowni-kiem.

W gramatyce j'ęzyka polskiego wyst'ępuje równiez' akomodacja typu frazy zdaniowej, w której cza-

sownik wymaga zdania podrzędne go poprzedzone go spójnikiem lub zdania pytajno-zależne go np. wyraz sadzi' c w znaczeniu przypuszczać lub wyraża'ć przekonanie wymaga zdania podrzędne go rozpoczynają - cego się spójnikiem że' (Sadz'e, że' gramatyka polska jest trudna.).

Rysunek 3.5: Przykład akomodacji typu frazy zdaniowej.

Zarówno w logice jak i w lingwistyce często wykorzystywane jest pojęcie konotacji.

Definicja 8 W składni j'ęzyka polskiego konotacja nazywamy zjawisko polegające na tym, że' wyst'epujący leksem niejako „zapowiada” pojawienie się innych leksemów lub konstrukcji składniowych.[27]

Przykładem może być wypowiedź Wczoraj mój pies..., która pozbawiona dokonczenia' rodzi pytanie Co się stało z twoim psem?. Naturalnym jest więc, że' bezpośrednio po usłyszeniu takiego niepełne go fragmentu słuchacz oczekiwał będzie czasownika, lub innych uzupełnień które odpowiadają na jego nie-zadane pytanie. Najczęściej konotacja syntaktyczna dotyczy zapowiadania określone go typu konstrukcji,

Ł. Wróbel

3.2. Relacje mi'edzy wyrazami

20

a rzadziej wymusza konkretny leksem, lub określone jego właściwo'ści'. W powyższym przykładzie wyraz „zapowiada” czasownik lub określon'ą grupę czasownikową, o strukturze, która wynika z cech akomodacyjnych i konotacyjnych danego czasownika[27].

Istnieje prosta metoda sprawdzenia, czy dany składnik podlega konotacji. Jeżeli usunięcie rzeczowego składnika sprawi, że konstrukcja stanie się niepoprawna, czy to w sensie składniowym, czy logicznym, oznacza to że jest on konotowany (wymagany)[27]. Jako przykład usunięcie ze zdania Wczoraj mój pies pogryzł moje ulubione pantofle. składnika pogryzł spowoduje otrzymanie niepoprawnej struktury: Wczoraj mój pies moje ulubione pantofle. Wyraz pies jest również konotowany. Zdanie Wczoraj mój pogryzł moje ulubione pantofle. w określonym kontekście (zapewniając odpowiedni podmiot domyślny) mogłoby zostać uznane za poprawne, to jednak bezpośrednie znaczenie zostało zmienione. Niekonotowane są natomiast składniki: wczoraj, mój, moje i ulubione.

Z. Klemensiewicz w swojej interpretacji budowy zdania, wprowadził następującą dychotomiczną

podział związków między składnikami w zdaniu:

związek główny - związek między podmiotem i orzeczeniem
związki poboczne - wszystkie pozostałe związki w zdaniu

Powyższy podział czyni podmiot i orzeczenie najważniejszymi składnikami zdania. Klemensiewicz

przypisywał dominującą rolę podmiotowi, który uznał za „nadrzędnik związku orzekającego”[7]. Orzeczenie niesie informacje o stanie podmiotu, lub wykonywanej przez niego czynności. Wymagane jest jednak uprzednie stwierdzenie istnienia podmiotu. Z punktu widzenia czysto strukturalnego czasownik

uznawany jest za nadrzędnik zdania, gdyż jest on jego elementem koniecznym i wystarczającym[27].
Związki poboczne kategoryzowane są następująco:

związek zgody
związek rzędu

związek przynależności

Związek zgody zachodzi na przykład między rzeczownikiem i przymiotnikiem lub między dwoma

rzeczownikami. Podrzędnik każdorazowo dostosowuje swoją formę fleksyjną do formy nadrzędnika[27].
Przykładowo rzeczownik kot narzuca określony rodzaj, liczbę i przypadek przymiotnikowi perski - kota perskiego, kotu perskiemu, kocie perskim. W przypadku dwóch rzeczowników sytuacja wygląda podobnie: pies bernardyn, psu bernardynowi, psem bernardynem.

W przypadku związku rzędu (głównie z czasownikiem) od rzeczownika wymagany jest ściśle określony przypadek[27]. Przykładem niech będzie wyrażenie: jechać pociągiem, lub karmić psa. Rzeczowniki pociąg i pies muszą wystąpić odpowiednio w narzędniku i bierniku, aby zdanie było składniowo poprawne. Związek rzędu może również zachodzić między rzeczownikiem, a przymiotnikiem: przy drodze, za górami, czy też między dwoma rzeczownikami: szklanka soku, chwila namysłu. Szczególnym przypadkiem związku rzędu jest sytuacja, w której oprócz określonego przypadku rzeczownika wymagany jest również odpowiedni przyimek zespalający go z czasownikiem, na przykład: idziemy na grzyby, czy też czytamy o lingwistyce.

Związek przynależności, który charakteryzuje się tym, że forma wyrazu podrzędnego nie jest wymuszana przez wyraz nadrzędny, występuje głównie między czasownikiem, a nieodmienną częścią mowy, na przykład przysłówkiem[27]. Przykładowo: idę szybko, wpadnę jutro.

Ł. Wróbel

4. Opis własnego rozwiązania

Wraz z narodzinami sztucznej inteligencji powstały nowe możliwości automatycznej analizy tekstu. Czatboty, które do tej pory mogły jedynie wykonywać ściśle zaprogramowane operacje w z góry określonych przypadkach zyskały nowe możliwości, które czyniły je nieporównanie bardziej elastycznymi. Dzięki możliwości wnioskowania, program, któremu dostarczone odpowiednio duże ilości danych był w stanie sam określić reguły stosowane w trakcie rozmowy, a nawet uczyć się ich na bieżąco. Wciąż jednak nie są one niezawodne, a z nowymi możliwościami pojawiły się też nowe wyzwania i problemy. Od idealnego czatbota oczekuje się, żeby pomysłnie przeszedł test Turinga. Oznacza to, że prowadzący z nim konwersację człowiek nie powinien się zorientować, że rozmawia z programem komputerowym.

Aby sprostać temu wyzwaniu czatbot musi spełnić szereg warunków poprawnościowych. Przede wszystkim, budowane przez program zdania muszą mieć prawidłową składnię w języku polskim. Wykorzystywana w tym celu jest tak zwana gramatyka generatywna stworzona przez Noama Chomsky'ego,

która jest skończonym zbiorem reguł umożliwiających zbudowanie wszystkich możliwych, poprawnych zdań. Chatbot jest dzięki niej w stanie zweryfikować, czy dana forma jest poprawna składniowo, a także

taka forma utworzyć. Jest to możliwe nawet, gdy zdanie podlegające weryfikacji zostało napotkane po raz pierwszy. Sama poprawność językowa nie jest jednak wystarczająca, aby prowadzić rozmowę z człowiekiem. Nie mniej ważna, a o wiele bardziej skomplikowana jest warstwa logiczna wypowiedzi. Słowa, które padają podczas rozmowy często mają więcej niż jedno znaczenie, które nie zawsze program komputerowy jest w stanie wywnioskować z kontekstu zdania, czy nawet większej części konwersacji. Porównania, przenoszenie idiomy, czy nawet pospolite przysłowia lub powiedzenia są poważną przeszkodą, która wciąż stanowi wyzwanie dla programistów.

Aplikacja, która jest przedmiotem tej pracy skupia się na sferze semantycznej rozmowy pomiędzy

rzeczonym czatbotem, a użytkownikiem'. Dotyka problemu rozumienia znaczenia słów i relacji pomiędzy nimi. W szczególności modeluje strukturę, która umożliwi programowi zbadanie kontekstu używanych przez użytkownika słów przez pryzmat posiadanej bazy wiedzy, aby lepiej zrozumieć ich znaczenie.

Człowiek, który usłyszy lub przeczyta określony wyraz natychmiast kojarzy go z odpowiadającym

mu obrazem. Dzięki temu od razu jest w stanie połączyć otrzymany obraz z przypisanymi mu (a więc też i słowu, które reprezentuje) cechami. Daje mu to możliwość opisanego co zobaczył, a także porównania z innym słowem (czyli de facto przypisanym mu obrazem). Ta naturalna metoda agregacji jest główną inspiracją przedstawionego dalej rozwiązania. Jakkolwiek program komputerowy nie jest w stanie wyobrazić sobie obrazu w sposób tak plastyczny jak człowiek, to jednak nie jest to przeszkodą, by

za pomocą grupowania słów z nim związanych był go w stanie z porównywalną dokładnością opisać.

Napisana przeze mnie aplikacja wykorzystuje pajaka internetowego, z którego pomocą budowana jest baza językowa będąca podstawowym źródłem wiedzy programu. W zależności od określonych przez

użytkownika słów będących korzeniami drzewa, na podstawie przechowywanego tekstu powstaje lista słów, które potencjalnie mogą być z nimi w relacji. Następnym krokiem jest wyodrębnienie tych wyrazów, które łączą się z korzeniami w poprawny językowo i logicznie zależności. Ostatnim etapem działania programu jest zaprezentowanie zdefiniowanych w ten sposób połączeń na grafie przyzwyczajen lingwistycznych. Szczegółowa realizacja wymienionych kroków, a także ich działanie zostaną opisane w kolejnych podrozdziałach.

21

.

22

4.1. Założenia projektowe

4.1. Założenia projektowe

Niniejszy program zakłada, że podane przez użytkownika słowo (korzeń) jest poprawnym słowem w języku polskim napisanym bez błędów. Poprawa błędów i weryfikacja poprawności nie są w zakresie

wymagań projektu. Ze względu na ograniczenia zakresu pracy aplikacja obsługuje wyłącznie relacje pomiędzy rzeczownikami, czasownikami, przymiotnikami i przysłówkami. Poprawność relacji badana jest w oparciu o wyniki z wyszukiwarki Google, która nakłada dodatkowe ograniczenia dotyczące maksymalnej ilości zapytań. Problem ten opisany został bardziej szczegółowo w dalszej części pracy.

4.2. Architektura systemu

Aplikacja składa się z czterech głównych elementów: pajaka internetowego, modułu odpowiedzial-

nego za tworzenie relacji, specjalistycznego grafu LHG oraz interfejsu użytkownika. Pajak, jest elementem opcjonalnym, gdyż użytkownik ma możliwość stworzenia bazy wiedzy na podstawie innego źródła niż Internet. Uogólniony diagram klas zaprezentowany został na rysunku 4.1.

Rysunek 4.1: Diagram klas

Ł. Wróbel

4.2. Architektura systemu

23

4.2.1. Interfejs użytkownika

Podstawowym zadaniem interfejsu użytkownika jest umożliwienie konfiguracji parametrów aplikacji. Dzięki niemu można zdefiniować wejścia i wyjścia danych, a także sposób pozyskiwania i wyświetla-

nia informacji, zależnie od preferencji użytkownika, lub ograniczeń systemu (na przykład brak dostępu do Internetu).

Bezpośrednio po uruchomieniu programu wyświetlone zostaje okno połączenia z bazą danych pokazane na zrzucie ekranu 4.2. Aby nawiązać łączność należy podać adres serwera, nazwę bazy, a także login i hasło. Połączenie z bazą danych nie jest wymagane i można z niego zrezygnować.

Rysunek 4.2: Okno połączenia z bazą danych

Główne okno programu zaprezentowane na zrzucie ekranu 4.3 jest nieco bardziej skomplikowane.

Najważniejszym elementem jest pole służące do podania słów, wokół których budowany będzie specjalistyczny graf przyzwyczajen lingwistycznych. Każde słowo wpisane po przecinku wykorzystane zostanie podczas wyszukiwania stron internetowych, w oparciu o które powstanie baza wiedzy, a także przy two-

zeniu relacji. Możliwe do wpisania są wyłącznie wyrazy znajdujące się w bazie słownikowej biblioteki CLP opisanej dokładniej w rozdziale 4.2.3.

Rysunek 4.3: Główne okno programu

Posród opcji możliwych do konfiguracji za pomocą głównego okna programu jest także tryb pracy aplikacji. Działanie systemu możliwe jest w jednym z czterech stanów, które definiują sposób pozyskiwania informacji, budowania bazy wiedzy, a także znacząco wpływają na wyświetlony finalnie graf. Tryby pracy aplikacji dzielą się na:

Ł. Wróbel

4.2. Architektura systemu

24

ONLINE - Jest to domyslny tryb działania, w którym aplikacja wykorzystuje połączenie z Internetem w celu zbudowania korpusu tekstu i oceny poprawności utworzonych relacji. Możliwe jest także wykorzystanie plików i bazy danych w celu odczytu i zapisu danych.

OFFLINE - Tryb ten pozwala wyłącznie na zaprezentowanie relacji, które zapisane są w pliku lub bazie danych. Nie wykorzystuje on połączenia z Internetem.

Symulator ONLINE - Symulator pozwala na utworzenie grafu relacji na podstawie korpusu tekstu

bez wykorzystywania Internetu do ich ewaluacji. Umożliwia to zaprezentowanie przykładowego grafu, który jest w stanie zobrazować interesujące informacje, jednocześnie nie korzystając z ograniczonych zasobów wyszukiwarki internetowej Google.

Symulator OFFLINE - Tryb ten umożliwia zasymulowanie pełnego działania programu bez do-

stępu do Internetu z wykorzystaniem korpusów tekstów znajdujących się w plikach lub bazie danych.

W głównym oknie programu zdefiniować można także źródła danych wejściowych. Możliwe jest wczytanie tekstu lub relacji zarówno z plików jak i bazy danych. Okno prezentujące listę plików z danymi do wczytania prezentuje zrzut ekranu 4.4. Jeżeli plik zaczyna się od taga RELATIONS, aplikacja usiłuje odczytać zapisane w nim relacje, w przeciwnym wypadku zostanie on uznany jako korpus tekstu. Informacje te można uzyskać także z bazy danych pod warunkiem, że wcześniej ustanowione zostało

połączenie. Można je zdefiniować ponownie wybierając z menu Plik, a następnie Połącz z bazą danych.

Rysunek 4.4: Wybór plików zawierających relacje lub korpusy tekstów.

W trybie ONLINE możliwe jest określenie liczby linków jaka zostanie przeglądnięta dla każdego słowa kluczowego. Adresy stron odwiedzonych przez pajaka zapisywane są do pliku, aby aplikacja

ponownie nie pobierała informacji z tego samego źródła. Użytkownik ma jednak możliwość ręcznego opróżnienia pliku z linkami.

Ważnymi parametrami z punktu widzenia wyświetlania informacji na grafie są parametry definiujące

maksymalna liczba połączeń, jakie zostaną wyświetlone dla pojedynczego słowa, a także maksymalna liczba połączeń, jakie zostaną wyświetlone dla relacji utworzonych na podstawie pseudo-logicznego rozumowania. Parametry te pozwalają kontrolować ilość informacji na grafie i czynią go bardziej czytelnym. Bardziej szczegółowy opis sposobu przedstawiania danych na grafie LHG znajduje się w rozdziale 4.4.

Ł. Wróbel

4.2. Architektura systemu

25

4.2.2. Pajak internetowy i ekstraktor danych

Zdecydowałem się napisać własną, uproszczoną wersję pajaka internetowego, która wykorzystuje tylko absolutnie podstawowe funkcjonalności, a jednocześnie w zupełności zaspokaja potrzeby aplikacji.

Nie zdecydowałem się na żadne z już istniejących rozwiązań, także aby mieć pełną kontrolę nad wykonywanymi przez pajaka operacjami. Podstawową różnicą w stosunku do większości crawlerów jest fakt, że aplikacja nie podąża za napotkanymi linkami, a tylko przechodzi przez adresy z określonego źródła, co znacznie skraca czas budowania bazy wiedzy i ogranicza ilość pobieranych, choć zbędnych w danym momencie danych.

Pajak rozpoczyna swoje działanie od znalezienia wszystkich linków na stronie startowej i zapisaniu

ich na listę adresów do odwiedzenia. Następnie tworzona jest określona liczba wątków, które po odwiedzeniu usuwają kolejne elementy listy. Pajak wykorzystuje framework Htmlparser, który udostępnia me-

todę pozwalającą na pobranie wyłącznie tekstu z danej witryny pozbawionego tagów HTML. Uzyskany

w ten sposób tekst, który wciąż jeszcze może zawierać wiele niepozadanych elementów przekazany zostanie do ekstraktora, który zajmie się jego dokładnym przetworzeniem. Po przekazaniu tekstu strony do ekstraktora wątek dodaje ją na listę odwiedzonych i przechodzi to przetwarzania kolejnego elementu z

listy oczekujących, lub kończy działanie, gdy taki element nie istnieje.

Ekstraktor danych, którego schemat działania widoczny jest na diagramie 4.5, jest integralną częścią

pajaka internetowego. Jego zadaniem jest usunąć wszelkie niepozadane elementy pozostałe po usunięciu tagów HTML i wyodrębnić tekst, który następnie będzie mógł być wykorzystany w głównej części aplikacji. Pierwszym problemem jaki napotkałem było określenie definicji „czystego tekstu”. Wbrew pozorom pozbycie się elementów języka HTML to dopiero początek procesu jego pozyskiwania. Podstawowym komponentem, z którego zbudowany jest każdy tekst jest forma zdaniowa. Przyjąłem, że naj-

mniej taka forma musi składać się przynajmniej z dwóch wyrazów, aby wnosić jakakolwiek wartość do bazy wiedzy. Z punktu widzenia relacji wyrażenie złożone z pojedynczego słowa jest mało interesujące.

Dla zmniejszenia stopnia skomplikowania problemu oraz zachowania w miarę przewidywalnej bu-

dowy zdań przyjąłem, że powinny się one składać wyłącznie z określonych znaków. Dopuszczalny zbiór

znaków zawierał wyłącznie duże i małe litery alfabetu polskiego, a także białe znaki, kropki, wykrzykniki, pytajniki i przecinki. Wszelkie elementy spoza tego zbioru dyskwalifikują zdanie jako źródło informacji. To dość restrykcyjne ograniczenie uzasadnione jest tym, że znaki takie jak -, :, czy () są w stanie znacząco zmienić strukturę zdania, a także utrudnić automatyczne sprawdzenie jego poprawności. Przykładem może być znak pominięcia fragmentu tekstu podczas cytowania: (...), czy też ewentualne pozostałości po tagach HTML, lub fragmentach skryptów które nie zawsze zostaną dokładnie usunięte.

Największym problemem, z jakim zetknąłem się podczas tworzenia ekstraktora było określenie gra-

nicy między formami zdaniowymi w ciągłym tekście. Teoretycznie zdanie w języku polskim zaczyna się wielką literą, a kończy kropką, znakiem zapytania lub wykrzyknikiem. Dla programu komputerowego taka definicja jest jednak niewystarczająca. W poprawnym zdaniu, mogą bowiem wystąpić wielkie litery,

które nie są początkiem nowego, jak i kropki, które tego zdania nie kończą. Dobrym przykładem może być wypowiedź: "Pan prof. Kowalski udał się na zasłużony urlop." Program szukając nowego zdania znalazłby literę P w słowie Pan, która słusznie uznaliby za jego początek, ale za koniec uznaliby kropkę nastę-

pującą po skrócie prof. Zaproponowane przeze mnie rozwiązanie wyszukuje co prawda zdania zgodnie z tą intuicyjną definicją, ale odrzuca te, które zakończone są skrótem jak na przykład prof., gen., gł.,

np., itp., a wyszukiwanie następnego rozpoczyna się od kolejnej niebędącej elementem żadnego skrótu kropki. Pozyskane w ten sposób zdania zostają zapisane do bazy danych i plików zdefiniowanych przez

użytkownika, a także wykorzystane w procesie budowania grafu relacji.

Ł. Wróbel

4.2. Architektura systemu

26

Rysunek 4.5: Schemat działania ekstraktora danych

4.2.3. Biblioteka CLP

Niezwykle istotnym z punktu widzenia poprawności syntaktycznej elementem aplikacji jest słownik fleksyjny CLP. Został on stworzony przez Katedrę Informatyki Akademii Górniczo-Hutniczej i Katedrę Lingwistyki Komputerowej Uniwersytetu Jagiellońskiego, a jego autorami są Wiesław Lubaszewski, Henryk Wróbel, Marek Gajęcki, Barbara Moskał, Paweł Pietras, Piotr Pisarek i Teresa Rokicka. Dostępny jest on w formie elektronicznej bazy danych SFJP (Słownik Fleksyjny Języka Polskiego) i biblioteki CLP dla języka C.

Podstawowe funkcjonalności biblioteki to:

rozpoznawanie wyrazu na podstawie jego dowolnej formy fleksyjnej

Ł. Wróbel

4.2. Architektura systemu

27

dostarczanie informacji o wyrazie

wygenerowanie form fleksyjnych dla danego wyrazu dostarczanie informacji o formie wyrazu
Każdy wyraz w bazie danych reprezentowany jest przez niepowtarzalny numer, dzięki któremu

można go jednoznacznie zidentyfikować, wraz z jego etykietą fleksyjną. Dla każdego słowa znajdującego się w bazie biblioteka jest w stanie zwrócić jego numer, lub tablicę numerów, jeśli dopasowanie nie

jest jednoznaczne. Na podstawie etykiety fleksyjnej możliwe jest określenie, na przykład części mowy i wzorca odmiany wyrazu. Identyfikator służy do wyznaczenia formy podstawowej danego słowa, a także

wszystkich jego możliwych form fleksyjnych wraz z jego pozycją wśród nich. W bazie słownika znajduje się ponad 120 tysięcy wyrazów. Jest ona na początku ładowana do pamięci RAM, co znacznie

przyspiesza wyszukiwanie. Istniejąca wersja działa wyłącznie pod systemem GNU/Linux.

4.2.4. Specjalistyczny graf przyzwyczajen lingwistycznych

Głównym zadaniem uproszczonej postaci grafu LHG, która jest wynikiem działania aplikacji jest

pokazanie użytkownikowi możliwie jak najbardziej reprezentatywnego fragmentu powstałych relacji. Do wizualizacji wykorzystałem bibliotekę JUNG, gdyż oferuje dość bogate możliwości prezentacji i

spora interakcję z użytkownikiem. Wierzchołki grafu reprezentują formy bazowe słów, które wchodzi w

skład relacji, natomiast krawędzie odpowiadają relacjom między połączonymi przez nie wyrazami. Kolor

każdego wierzchołka zależy od tego, jaką część mowy jest słowo, które reprezentuje. Szczegóły

dotyczące przyporządkowania kolorów do poszczególnych części mowy zostały zaprezentowane w tabeli 4.1.

Tablica 4.1: Kolory poszczególnych części mowy

| Część mowy | Kolor |
|------------|-------|
|------------|-------|

rzeczownik
niebieski
czasownik
zółty
przymiotnik
zielony
przysłówek
turkusowy
liczebnik
pomarańczowy

Na grafie nie są prezentowane wszystkie powstałe i zaimportowane relacje, gdyż stałyby one

nieczytelny. Parametr ustawiony przez użytkownika w głównym oknie programu określa maksymalną liczbę krawędzi, które są wyświetlane dla pojedynczego wierzchołka. Wszystkie mniej lub bardziej po-

prawne połączenia między słowami zaprezentowane są na jednej skierowanej krawędzi. Po kliknięciu

na nią użytkownik może odczytać wszystkie relacje, które reprezentuje, ich typ, a także nadaną im wartość. Przykład wyświetlania szczegółowych informacji o krawędziach został przedstawiony na rzutach

ekranu 4.6 i 4.7. Prezentowane są wyłącznie informacje uznane za najbardziej wartościowe i dające najlepszy obraz działania systemu. Dokładny opis algorytmu decyzyjnego odpowiedzialnego za rysowanie krawędzi opisany został w sekcji 4.4.

Użytkownik ma dodatkowo możliwość wyświetlenia relacji przechodnich. Są one wynikiem pseudo-

logicznego rozumowania modelującego naturalną przechodniosć cech w ramach hierarchii. Nie są one

bezpośrednio uzyskane z tekstu źródłowego, ale powstają poprzez wyciąganie wniosków z otrzymanej struktury, co maksymalizuje ilość uzyskanych informacji i dokładniej oddaje sposób ludzkiego myślenia.

Dodatkowo możliwe jest uzupełnienie otrzymanego grafu nowymi słowami kluczowymi i relacjami przez wybór opcji Dodaj więcej informacji z menu Plik. Należy jednak pamiętać, że zaprezentowany

zostaje zawsze wyłącznie ograniczony wycinek zbudowanej struktury, co sprawia, że niektóre oczeki-

ł. Wróbel

4.2. Architektura systemu

28

wane połączenia mogą nie zostać wyświetlone, a zamiast nich zaprezentowane zostaną te, które zostały najlepiej ocenione.

Jezeli uzytkownik chce zobaczyc wszystkie utworzone i pobrane z bazy wiedzy relacje, ma taka mozliwosc przez wybor opcji Pokaz liste podstawowych relacji z menu Plik. Dodatkowo, jezeli chce zobaczyc relacje, ktore powstaly na skutek wnioskowania, ma mozliwosc wyboru opcji Pokaz liste wszystkich relacji. Przykladowa lista, jaka zostanie zaprezentowana zostala pokazana na zrzucie ekranu 4.8.

Rysunek 4.6: Szczegóły krawędzi - grupowanie, posiadanie

Rysunek 4.7: Szczegóły krawędzi - określenie

Ł. Wróbel

4.3. Opis działania aplikacji

29

Rysunek 4.8: Lista relacji

4.3. Opis działania aplikacji

Celem działania aplikacji jest stworzenie uproszczonej postaci grafu przyzwyczajen¹ lingwistycznych, który ilustruje hierarchię wyrazów wraz z ich określeniami² oraz pokazuje przechodnios³c⁴ niektórych cech w ramach tej struktury. Powstały graf jest wyłącznie⁵ uproszczoną⁶ odmianą⁷ grafu LHG, gdyż⁸

prezentuje wyłącznie⁹ bazowe formy wyrazów, a także¹⁰ wyłącznie¹¹ niektóre typy relacji charakterystycznych dla LHG[9]. Z punktu widzenia czatbota, może¹² to być¹³ cenne źródło¹⁴ wiedzy, które podczas rozmowy można¹⁵ wykorzystać¹⁶ do analizy kontekstu, w jakim użyte¹⁷ zostało¹⁸ jakieś¹⁹ słowo, zadawania pytań²⁰ o określone²¹ cechy wyrazu, czy też²² jego opisywania.

4.3.1. Tryb online

Rysunek 4.9: Wybór trybu online

Po wybraniu trybu online (zrzut ekranu 4.10), do wyszukiwarki Google zostaje wysłane zapytanie

Ł. Wróbel

4.3. Opis działania aplikacji

o podane słowa-korzenie. Ze zwróconej listy wyników pobierana jest określona w parametrach konfi-

guracyjnych przez użytkownika liczba linków, które następnie przekazywane są do kolejnych wątków, pajaka. Po określeniu, czy strona nadaje się jako źródło informacji następuje ekstrakcja tekstu. Powstaje lista zdań, które trafiają do określonego na początku miejsca docelowego bazy wiedzy dla tekstu. Wątki działają dopóki na liście adresów do odwiedzenia znajdują się jakieś elementy. Adresy, które zostały odwiedzone wypisane są w logach: 4.3.1.

2011 03 05

14:34:55

[INFO]

Pobieranie

słów:

<http://eszukam.pl/>

[posokowiec_bawarski_szczenie_pies/debnica_kaszubska/324579/](http://eszukam.pl/posokowiec_bawarski_szczenie_pies/debnica_kaszubska/324579/)

2011 03 05

14:34:55

[INFO]

Pobieranie

zdań:

<http://eszukam.pl/>

[posokowiec_bawarski_szczenie_pies/debnica_kaszubska/324579/](http://eszukam.pl/posokowiec_bawarski_szczenie_pies/debnica_kaszubska/324579/)

2011 03 05

14:34:55

[INFO]

Pobieranie

tekstu:

<http://eszukam.pl/>

[posokowiec_bawarski_szczenie_pies/debnica_kaszubska/324579/](http://eszukam.pl/posokowiec_bawarski_szczenie_pies/debnica_kaszubska/324579/)

2011 03 05

14:34:55

[INFO]

Pobieranie

słów:

<http://tanio.pl/Pies>

[Baskervilleow_Artur_Conan_Doyle_Audiobook1_68234_12277036.html](http://tanio.pl/Baskervilleow_Artur_Conan_Doyle_Audiobook1_68234_12277036.html)

2011 03 05

14:34:55

[INFO]

Pobieranie

zdan':
<http://tanio.pl/Pies>

Baskervilleow Artur Conan Doyle Audiobook1 68234 12277036.html

2011 03 05
14:34:55
[INFO]
Pobieranie
tekstu:
<http://tanio.pl/Pies>

Baskervilleow Artur Conan Doyle Audiobook1 68234 12277036.html

2011 03 05
14:34:55
[INFO]
Pobieranie
słów:
<http://www.>

niepelnosprawni.pl/ledge/x/10686

2011 03 05
14:34:55
[INFO]
Pobieranie
zdan':
<http://www.>

niepelnosprawni.pl/ledge/x/10686

2011 03 05
14:34:55
[INFO]
Pobieranie
tekstu:
<http://www.>

niepelnosprawni.pl/ledge/x/10686

2011 03 05
14:34:55
[INFO]
Pobieranie
słów:
<http://klarka.blog.onet>.

pl/jaki pan taki pies, 2, ID419994422, n

2011 03 05
14:34:55
[INFO]
Pobieranie
zdan':
<http://klarka.blog.onet>.

pl/jaki pan taki pies, 2, ID419994422, n

2011 03 05
14:34:55
[INFO]
Pobieranie
tekstu:
<http://klarka.blog.onet>

.pl/jaki pan taki pies, 2, ID419994422, n

2011 03 05
14:34:55
[INFO]
Pobieranie
słów:
<http://pies.org.pl/>

2011 03 05
14:34:55
[INFO]
Pobieranie
zdan':
<http://pies.org.pl/>

2011 03 05
14:34:55
[INFO]

Pobieranie
tekstu:
<http://pies.org.pl/>

2011 03 05
14:34:55

[INFO]

Pobieranie
słów:
<http://www.filmweb.pl/>

film/Pies+andaluzyjski 1929 906
zdan´:

2011 03 05
14:34:55

[INFO]

Pobieranie

<http://www.filmweb.pl/>

film/Pies+andaluzyjski 1929 906

2011 03 05
14:34:55

[INFO]

Pobieranie
tekstu:
<http://www.filmweb.pl/>

film/Pies+andaluzyjski 1929 906

W następnym kroku wyniki działania pajaka, oraz wybrane przez użytkownika źródła tekstu są ana-

lizowane i wybierane są z nich zdania zawierające, podane przez użytkownika słowa w ich dowolnej formie fleksyjnej. Wyszukiwane są wszystkie zwrócone przez bibliotekę CLP formy dla korzeni. Powstaje w ten sposób nowa lista zdan´. Wykorzystując bibliotekę CLP wszystkie słowa są następnie sprowadzane do formy bazowej, dzięki czemu ich formy fleksyjne nie będą traktowane jako niezależne słowa. Są one

następnie zliczane i sortowane malejąco, w zależności od częstości występowania. Dla każdego korzenia powstaje oddzielna lista, a sam korzeń oczywiście znajduje się na jej szczycie, jako że wystąpił on

w każdym przetworzonym zdaniu. Za nim zazwyczaj umiejscowione są popularne przyimki np. z, do, w, na i spójniki np. i, a, lub. Program filtruje jednak z pomocą biblioteki CLP otrzymane wyniki, tak

aby pozostały wyłącznie określone części mowy. Aplikacja skupia się na najbardziej istotnych i niosących najwięcej informacji, a zatem na rzeczownikach, czasownikach, przymiotnikach, przysłówkach i

liczebnikach. Rozróżnienia, jaka część mowy jest dane słowo dokonuje biblioteka CLP.

Ł. Wróbel

4.3. Opis działania aplikacji

31

Jako że biblioteka CLP dopasowując wyraz porównuje go z wszelkimi możliwymi formami, nie

zawsze zwrócony wynik odpowiada kontekstowi zdania. Przykładowo często występujące słowo kilka identyfikowane jest w pierwszej kolejności jako rzeczownik (jest to nazwa ryby), w drugiej jako przymiotnik, a dopiero w trzeciej kolejności jako zaimek. Podobna sytuacja występuje w przypadku bardzo

często używanego słowa jest, którego forma bazowa być również uznawana jest przez CLP jako rzeczownik (dopełniacz liczby mnogiej rzeczownika bycie), a dopiero kolejny zwrócony wynik dopasowania jest

czasownikiem. Problem ten rozwiązany został przez rozróżnienie na etapie przetwarzania tekstu nie tylko

słów, ale także ich części mowy, co sprawia, że homonimy są rozróżnialne i traktowane jako oddzielne wyrazy.

W tak posortowanej liście wyrazy znajdujące się wysoko mają dużą szansę być w relacji ze słowem kluczowym, jako że często były użyte z nim w jednym zdaniu. Takie statystyczne podejście utrudnia

co prawda stwierdzenie koneksji z rzadziej używanymi słowami, ale omija problemy związane z za-
włościami składni języka polskiego. Nie jest konieczne budowanie skomplikowanego parsera zdań z

wyszukanymi regułami uwzględniającymi niuanse językowe. Związki między wyrazami ustalane są w ramach jednego zdania. Frazy, w których klucz został zamieniony zaimkiem lub wystąpił jako podmiot domyslny w ogóle się tu nie znajdują, gdyż zostaną odfiltrowane jako niezawierające korzenia.

Elementy z posortowanej listy potencjalnych kandydatów do relacji ze słowem kluczowym są ko-

lejno w odpowiedniej formie łączone z korzeniem (także w odpowiedniej formie) z wykorzystaniem szablonów relacji. W zależności od typu relacji do utworzonej w ten sposób formy może zostać dodany

spójnik. Szablony relacji definiują typowe zależności pomiędzy dwoma wyrazami występującymi w języku polskim i określają części mowy jakie wchodzi w ich skład, ich formę fleksyjną oraz kolejność. Nie muszą być one skomplikowane, ani brać pod uwagę innych słów, które mogłyby zaburzyć ustalony schemat. Ważne jest, aby sam układ i forma dwóch składowych wyrazów był poprawny.

Dla każdego schematu utworzona przez niego fraza jest przekazywana do wyszukiwarki Google, po czym

przypisywana jest jej wartość zależna od ilości zwróconych wyników zapytania. Dopóki więc dana relacja między dwoma wyrazami występuje w niezmienionej, dokładnie takiej jak określona formie odpowiednio często na stronach indeksowanych przez wyszukiwarkę, program zakłada, że jest ona

poprawna. Każda para wyrazów (korzeń i słowo z listy wyrazów często z nim występujących w jednym zdaniu) może być w potencjalnie dowolnej zależności, a więc dla odpowiednich części mowy spraw-

dzane są wszystkie możliwe szablony relacji. Te, dla których Google zwróci odpowiednio dużą liczbę wyników zostają uznane za poprawne. Szablony powinny odzwierciedlać najczęściej występujące struktury, gdyż każdy kolejny zwiększa liczbę zapytań o poprawność utworzonej relacji, co bezpośrednio przekłada się na maksymalną ilość uzyskanych wartościowych informacji.

Przykładowo dla korzenia pies i słowa z listy łapa, jako że oba wyrazy są rzeczownikami sprawdzone

zostaną możliwe relacje pomiędzy dwoma rzeczownikami. Tak więc dla relacji posiadania o szablonie MIANOWNIK1 + ma + BIERNIK2 ułożona zostanie fraza pies ma łapę, dla której Google zwróci X wyników. Sprawdzona zostanie też relacja odwrotna, a zatem łapa ma psa. Kolejną potencjalną relacją pomiędzy dwoma rzeczownikami jest grupowanie, czyli przypadek, gdy jeden wyraz należy do podzbioru (czyli jest przykładem) wyrazu drugiego. Szablon tej zależności wygląda następująco: MIANOWNIK1 +

jest + NARZEDNIK2, a zatem pies jest łapą, czy też łapa jest psem nie zwróci zbyt wielu wyników. Dla pary kot i zwierzę fraza kot jest zwierzęciem oceniona zostanie wysoko, a odwrotność zwierzę jest kotem już nie.

Algorytm sprawdza poprawność relacji wysyłając kolejne zapytania do wyszukiwarki do momentu,

aż nie wyczerpią się wyrażenia, które wymagają oceny, lub nie zostanie zwrócony błąd oznaczający przekroczenie maksymalnej ilości zapytań. Wynikiem jego działania jest lista relacji z przyporządkowanymi

im ocenami. Każda uzyskana w ten sposób informacja jest cenna. Otrzymane wyniki dają informację

nie tylko o tym, które wyrażenia są poprawne, ale także o tym, które poprawne nie są. Jako, że nigdy

nie ma absolutnej pewności, czy dana relacja ma sens, czy nie, przyporządkowana jej ocena pozwala

ustalić wyłącznie prawdopodobieństwo. Informacje o utworzonych relacjach wypisywane są w logach aplikacji: 4.3.1.

Ł. Wróbel

4.3. Opis działania aplikacji

32

Listing 4.1: Fragment logów wyświetlanych podczas tworzenia relacji.

2011 03 05

13 : 33 : 28

[INFO]
o s o b a
l u b i :
139000

2011 03 05
13:34:36
[INFO]
o s o b a m a
r a s e , :
0

2011 03 05
13:34:38
[INFO]
c a ł y
j a k
o s o b a : 0

2011 03 05
13:34:40
[INFO]
o s o b a
w i e :
93500

2011 03 05
13:34:42
[INFO]
o s o b a
j e s t
s a , s i a d e m :
0

2011 03 05
13:34:44

[INFO]
o s o b a
t o
s w o j e : 0

2011 03 05
13:34:45
[INFO]
l u d
m a
o s o b e _ :
0

2011 03 05
13:35:04
[INFO]
o s o b a
z r o b i :
19400

2011 03 05
13:35:37
[INFO]
c z ł o w i e k
j e s t
o s o b a _ : 1550000

2011 03 05
13:36:44
[INFO]
.

j a k
c z ł o w i e k :
26800

k a z d y

2011 03 05
13:37:15
[INFO]
swoje
jest
zwierzęciem:
0

2011 03 05
13:37:16
[INFO]
swoje
jest
rasa: 0

2011 03 05
13:37:18
[INFO]
swoje
jest
siedem:
0

2011 03 05
13:37:20
[INFO]
swoje
jest
człowiekiem:
0

2011 03 05
13:37:33
[INFO]
człowiek
wie:
398000

2011 03 05
13:37:47

[INFO]
oko
to
zwierze,: 0

2011 03 05
13:37:49
[INFO]
oko
to
rasa: 0

2011 03 05
13:37:50
[INFO]
oko
to
sa_siad:
0

2011 03 05
13:38:27
[INFO]
nowy
sa_siad: 50500

2011 03 05
13:38:28
[INFO]
nowy
człowiek:
29200

2011 03 05

13:38:30
[INFO]
oko
jest
zwierzęciem:
0

2011 03 05
13:38:32
[INFO]
oko
jest
rasa:
0

2011 03 05
13:38:34
[INFO]
oko
jest
sąsiadem: 0

2011 03 05
13:38:35
[INFO]
oko
jest
człowiekiem:
0

2011 03 05
13:40:17
[INFO]
człowiek
musi: 1280000

2011 03 05
13:41:14
[INFO]
oko
człowieka:

50500

2011 03 05
13:41:16
[INFO]
pasterski
jak
zwierze:

0

2011 03 05
13:41:18
[INFO]
pasterski
jak
rasa:
0

2011 03 05
13:41:48
[INFO]
człowiek
jest
cały:
342000

4.3.2. Tryb offline

Rysunek 4.10: Wybór trybu offline

Główną różnicą między trybem offline, a trybem online jest brak konieczności łączenia się z Internetem. W tym stanie aplikacja nie wykorzystuje pająka internetowego i nie wylicza ocen relacji, a jedynie wyświetla te, które zostały pobrane z bazy danych lub plików. Działanie programu w tym trybie jest

naturalnie o wiele szybsze. Jako, że nie jest możliwe określenie poprawności relacji, użytkownik nie musi

podawac' słowa kluczowego. Graf zostanie zbudowany na podstawie informacji, które juz' zostały zebrane. W trybie tym istnieje także' możliwość' wczytania relacji z plików i ich zapis do bazy danych.

4.3.3. Symulator

Tryb symulatora zaimplementowany został wyłącznie' na potrzeby prezentacji idei działania aplikacji. Nadaje on nowo utworzonym relacjom stałe' oceny poprawności', aby zostały one wyświetlone' na

ł. Wróbel

4.4. Generacja specjalistycznego grafu LHG

33

grafie. Ustalanie poprawności' w ten sposób oczywiście' nie ma nic wspólnego z rzeczywistością', ale pozwala jednak przedstawić' sposób działania programu i wszystkie jego funkcje bez konieczności' korzystania z mechanizmu sprawdzania poprawności'. Zrzut ekranu 4.11 przedstawia wybór symulacji sprawdzania poprawności' odpowiednio dla trybu offline i online.

Rysunek 4.11: Wybór trybu pracy z symulatorem

4.4. Generacja specjalistycznego grafu LHG

Ostatnim etapem działania aplikacji jest wygenerowanie grafu przyzwyczajajen' lingwistycznych w uproszczonej postaci, który zaprezentuje znalezione relacje z przypisanymi do nich wartościami'. Po zakończeniu przetwarzania danych powstaje lista relacji złożona' zarówno z wczytanych przez użytkownik'a, jak i uzyskanych na podstawie analizowanego tekstu. Słowa, które wchodzi' w skład wyrazin' staja' się' węzłami i zapisane sa' w formie bazowej. W zależności' od tego, jaka' część' mowy reprezentuja', przy-

porzadkowany' jest im odpowiedni kolor zgodnie z tabela: 4.1. Słowa z tym samym źródłosłowem,' lecz

będace' różnymi' częściami' mowy sa' też' różnymi' wierzchołkami. Przykład grafu obrazujacego' różnicowanie kolorów w zależności' od części' mowy, jaka' reprezentuje słowo w wierzchołku pokazuje zrzut ekranu 4.12.

Rysunek 4.12: Przykładowy graf z wierzchołkami będącymi różnymi częściami mowy.

Poprzez kliknięcie na wierzchołku grafu użytkownik ma możliwość odczytania, jaka część mowy jest słowo, które dany wierzchołek reprezentuje. Informacja ta nie tylko określona jest za pomocą koloru zgodnego z tabelką 4.1, ale także jest wyświetlana w panelu informacyjnym, jak pokazuje zrzut ekranu 4.13.

Ł. Wróbel

4.4. Generacja specjalistycznego grafu LHG

34

Rysunek 4.13: Informacja o konkretnym wierzchołku wyświetlona na grafie.

Krawędzie grafu definiują połączenia między dwoma wyrazami. Nie odpowiadają one bezpośrednio relacji, ale raczej grupie relacji między wierzchołkami. Zwrot krawędzi określa rolę jaką dany węzeł pełni w ich połączeniu. Zawsze jest on skierowany od nadrzędnego wyrazu w wyrażeniu do podrzędnego.

Między dwoma wierzchołkami może zatem wystąpić co najwyżej dwie krawędzie o różnych zwrotach. Relacje, które uznane zostały za niepoprawne (mają ocenę 0) nie są wyświetlane na grafie, co znacznie poprawia jego czytelność. Decyzja, czy dane relacje tworzą krawędź, która następnie zostanie wyświetlona na grafie zależy od średniej oceny relacji między słowami reprezentowanymi przez wierzchołki. Ich ilość jest ograniczona przez parametry opisane w sekcji 4.2.1. Tylko krawędzie o największej średniej ocenie relacji zostaną wyświetlone.

Relacje utworzone na podstawie analizy hierarchicznej struktury istniejących powiązań są wyświetlane tylko na żądanie użytkownika, aby nie zaciemniać grafu. Powstałe w ten sposób krawędzie oznaczone są kolorem czerwonym, aby łatwo można było je odróżnić i przesledzić ścieżkę rozumowania aplikacji. Liczba dorysowanych w ten sposób krawędzi zależy od zdefiniowanego przez użytkownika

w głównym oknie programu parametru Maksymalna ilość relacji spropagowanych opisanego w sekcji 4.2.1. Graf z widocznymi spropagowanymi relacjami został zaprezentowany na zrzucie ekranu 4.14.

Ł. Wróbel

4.4. Generacja specjalistycznego grafu LHG

35

Rysunek 4.14: Widok pokazujący spropagowane relacje na grafie.

Ł. Wróbel

5. Prezentacja działania aplikacji

Aby dokładnie zaprezentować działanie aplikacji prześledźmy przykład jej wykorzystania w celu wygenerowania struktury relacji dla zadanych słów. Załóżmy, że czatbot rozmawiający, właśnie z użytkownikiem chce przeanalizować otrzymaną wypowiedź z pod kątem określonych słów kluczowych. Jest to konieczne, jeżeli chce nawiązać do cech charakterystycznych tematu rozmowy, lub też sformułować inteligentne pytanie. To, czy czatbot „zrozumie” niektóre elementy rozmowy zależy będzie od tego jakie będą jego skojarzenia, czyli relacje, które posiada w bazie danych.

Przykładowo dla wyrazu gepard dane wejściowe aplikacji mogą wyglądać jak na zrzucie ekranu 5.1.

Oczywiście możliwe jest wykorzystanie większej ilości zasobów (zarówno stron internetowych, jak i korpusów tekstów).

Rysunek 5.1: Przykładowe dane wejściowe dla programu.

Ze 100 wyników zwróconych przez wyszukiwarkę Google dla słowa gepard pobierany jest tekst, który następnie dzielony jest na zdania, a te z kolei na wyrazy, które pojawiają się w kontekście słowa gepard. Najczęściej pojawiające się wyrazy są następnie podstawiane do szablonów relacji wraz ze słowem gepard i ich poprawność jest sprawdzana za pomocą wyszukiwarki Google. Im więcej wyników zostanie zwróconych dla określonej frazy, tym częściej występuje ona w Internecie, co pośrednio przekłada się

także na jej poprawność językową. Gdy do wyszukiwarki trafi zbyt wiele zapytań, tworzenie nowych relacji zostaje przerwane i wyświetlone zostają uzyskane do tej pory wyniki. Sytuacja taka obrazuje wycinek logów 5.

2011 04 09
11:32:56
[INFO]
.
gepard: 17

używany

2011 04 09
11:32:58
[INFO]
wielki
jak
gepard: 0

2011 04 09
11:32:58
[INFO]
Zakończono
tworzenie relacji podstawowych.

36

37

2011 04 09
11:32:59
[INFO]
firma
ma
zdjęcie:
0

2011 04 09
11:33:01
[INFO]
firma
ma
zwierze:
0
zapytań

2011 04 09
11:33:09
[ERROR]
Przekroczono limit

do google.

2011 04 09
11:33:09
[INFO]
człowiek
to firma:
0

2011 04 09
11:33:09
[INFO]
Zakończono tworzenie
relacji
przechodnich.

Jak widać na zrzucie ekranu 5.2 utworzony został graf przyzwyczajen lingwistycznych, który obrazuje wycinek relacji uzyskanych dla wyrazu gepard przy wykorzystaniu zdefiniowanej wcześniej bazy wiedzy. Zgodnie z parametrem, jaki został podany, maksymalnie 10 krawędzi łączy się z pojedynczym wierzchołkiem, co przy ocenach utworzonych relacji pozwala wyświetlić 8 wierzchołków reprezentujących słowa związane relacjami z korzeniem. Popularne czasowniki takie jak być, móc i mieć niebezpiecznie łączy się z wieloma rzeczownikami. Rzeczownik firma pojawia się, gdyż GEPARD TRANS jest nazwą dobrze rozreklamowanej w Internecie polskiej firmy przewozowej. Kolor czarny został przypisany gepardowi, choć w rzeczywistości odnosi się on raczej do cętek niż do samego geparda.

Rysunek 5.2: Graf relacji dla słowa gepard.

Ł. Wróbel

Na pełnej liście utworzonych relacji zaobserwować można relacje, które nie zostały zaprezentowane

na grafie ze względu na zbyt niską wartość krawędzi, takie jak używany gepard (GEpard to nazwa dealera samochodów), czy gepard zostanie. Lista ta widoczna jest na zrzucie ekranu 5.3.

Rysunek 5.3: Przykładowa lista relacji wygenerowanych dla słowa gepard.

Ł. Wróbel

39

Dla wyrazu strus´ baza wiedzy zbudowana została na podstawie 300 wyników z wyszukiwarki Go-ogle, co nieco poprawiło jakość utworzonych relacji. Graf LHG został zaprezentowany na zrzucie ekranu 5.4. Znow pojawiły się na nim czasowniki być i mieć, ale wyświetlone zostały także czasowniki bardziej charakterystyczne dla strusia, jak chować, czy biegać. Warto zauważyć, że biblioteka CLP wyraźnie różni słowa biegać i biec. Wyświetlone przymiotniki, jakie charakteryzują strusia to afrykański, szybki i

polski. Brak rozróżnienia części mowy w momencie analizy zdania i konieczność wykorzystania wszystkich dopasowań zwróconych przez bibliotekę CLP skutkuje tym, że jako poprawna relacja uznane zostało wyrażenie strus´ szybki, gdzie szybki jest dopełniaczem wyrazu szybka. Google zwróciło wiele wyników, pomimo iż wydawałoby się, że kolejność wyrazów nie jest poprawna z gramatycznego punktu widzenia. Okazuje się jednak, że gra Strus´ Szybki Lopez jest na tyle popularna w Internecie, że algorytm sprawdzania poprawności uznaje taką relację za poprawną.

Rysunek 5.4: Graf relacji dla słowa strus'.

Ł. Wróbel

40

Interesujące informacje można także odczytać z pełnej listy relacji zaprezentowanej na zrzucie ekranu 5.5. Mimo, iż relacja mówiąca, że strus jest ptakiem zwróciła kilkaset wyników, to jednak nie została zaprezentowana na grafie LHG ze względu na limit 10 krawędzi dla jednego wierzchołka. Przewidzenie, czy określony limit ograniczy wyświetlenie wartościowych informacji, gdy będzie zbyt mały,

czy pozwoli na wyświetlenie niepotrzebnych relacji, zaciemniających graf, gdy będzie zbyt duży, jest trudne do przewidzenia. Na liście widoczne są także z pozoru niezrozumiałe relacje takie jak kojot strusia, czy pędziwiatr to strus. Nawiązują one jednak do popularnej bajki „Strus Pędziwiatr” studia Warner Bros.

Rysunek 5.5: Przykładowa lista relacji wygenerowanych dla słowa strus'.

Ł. Wróbel

41

Aby wygenerować więcej relacji przechodnich na podstawie relacji już istniejących, wystarczy wczytać te relacje, czy to z plików, czy bazy danych. Nie jest konieczne wykorzystywanie pajaka internetowego, więc ilość linków do pobrania można ustawić na 0. Jeżeli z bazy danych lub z pliku nie zostaną pobrane zdania, na podstawie których powstać mogą nowe relacje, aplikacja stworzy będzie wyłączone

relacje na podstawie analizowania istniejącej hierarchii i wyciągania pseudo-logicznych wniosków.

Przykładowa konfiguracja takich danych wejściowych zaprezentowana została na zrzucie ekranu 5.6.

Rysunek 5.6: Przykładowe dane wejściowe dla generacji dodatkowych relacji przechodnich.

Ł. Wróbel

42

Po pobraniu relacji z bazy danych lub plików, graf przyzwyczajen¹ lingwistycznych rozrasta się, a słowa, które wcześniej charakteryzowały pojedyncze wyrazy łączą się z kilkoma innymi tworząc bardziej rozbudowaną strukturę. Nadal obowiązują jednak reguły dotyczące maksymalnej ilości krawędzi dla jednego wierzchołka, które ustalone zostały w momencie definiowania danych wejściowych¹. Taki bardziej rozbudowany graf zawierający zarówno relacje z dwóch poprzednich przykładów, jak i kilka innych pokazany został na zrzucie ekranu 5.7.

Rysunek 5.7: Przykładowy graf relacji dla kilku słów.

Ł. Wróbel

43

Większa struktura tego typu pozwala na podejmowanie przez aplikację prób utworzenia relacji na podstawie analizy hierarchii. Na zrzucie ekranu 5.8 widac' owoce tego procesu. Wyświetlone' zostały nowe krawędzie, które powstały na podstawie nowych propozycji relacji. Jak widac' na zrzucie ekranu 5.2 poprawna jest relacja gepard jest zwierzęciem. Choc' nie została ona wyświetlona' na bardziej rozbudowanym grafie, to jednak została uwzględniona w procesie wnioskowania. Jako, ze' poprawna jest także' relacja zwierzę lubi, aplikacja uznała łącząc te dwa fakty, ze' rzeczownik gepard i czasownik lubić także'

moga być połączone, relacja. Weryfikacja przeprowadzona na podstawie ilości wyników zwróconych przez wyszukiwarkę Google potwierdziła poprawność takiego połączenia, co zaowocowało powstaniem nowej czerwonej krawędzi na grafie, obrazującej spropagowaną relację. Dokładna lista wszystkich relacji została zaprezentowana na zrzucie ekranu 5.9.

Rysunek 5.8: Przykładowy graf relacji wraz z relacjami przechodnimi dla kilku słów.

Ł. Wróbel

44

Rysunek 5.9: Przykładowa lista relacji wygenerowanych dla kilku różnych słów.

W przykładach zaprezentowanych powyżej starałem się pokazać sposób działania aplikacji i używane przez nią wyniki. Na efekt działania programu silnie wpływają trendy panujące aktualnie w Internecie, a także reklamy rozmaitych produktów, których nazwy często łączą słowa w nieoczekiwany

sposób. Homonimy rozróżniane są poprzez wykorzystanie wszystkich części mowy zwróconych przez bibliotekę CLP i przekazanie ich jako zapytanie do wyszukiwarki w odpowiedniej formie. Dwa identycz-

nie zapisane wyrazy w formie bazowej, będące jednak różnymi częściami mowy, rozróżnione zostaną na podstawie reguł akomodacyjnych.

Ł. Wróbel

6. Podsumowanie wyników pracy i wnioski

6.1. Podsumowanie

Współczesna technologia wciąż jest daleka od skonstruowania inteligentnego robota, który byłby w stanie idealnie komunikować się z ludźmi. Taka maszyna musi znać nie tylko gramatykę danego języka,

ale także rozumieć znaczenie wypowiedzianych słów na poziomie wyższym niż prosta poprawność językowa. Problem badania znaczenia słów, aby następnie poprawnie użyć ich podczas konwersacji nie

może sprowadzać się wyłącznie do analizowania ich kontekstu w pojedynczym zdaniu lub wypowiedzi. Mimo, iż rozszyfrowanie formy frazeologicznej słowa nie jest zadaniem skomplikowanym, to sformułowanie nie tylko poprawnej gramatycznie, ale i logicznie wypowiedzi zależy już od większej liczby czynników. Podobnie jak człowiek, który po raz pierwszy usłyszał dane słowo, jest w stanie wywnio-

skować jego znaczenie z kontekstu zdania i następnie używać go poprawnie, tak i możliwe jest to dla czatbota. Niezależnym jednak elementem, który warunkuje taką umiejętność jest posiadanie odpowied-

niej wiedzy, na podstawie której stosując logiczne rozumowanie możliwe staje się uzyskanie nowych,

poprawnych informacji. Stworzona przeze mnie aplikacja jest propozycją rozwiązania tego problemu. Dostarcza ona danych, które wykorzystane przez czatbota pozwolą mu zarówno na znacznie lepsze rozumienie analizowanych przez niego wypowiedzi, jak i formułowanie własnych logicznie poprawnych wyrażen.

Największym atutem mojego rozwiązania jest fakt, że nie ogranicza się ono do informacji zawartych

wyłącznie w korpusach tekstów. Jakkolwiek są one i długo pozostaną nieocenionym źródłem danych dla czatbotów, to ilość słów zawarta nawet w największych tego typu repozytoriach jest nieporównanie mniejsza niż mnogość form wyrazowych w Internecie. Do osiągnięcia zamierzonego celu niezbędna

była baza wiedzy, która nie ogranicza się wyłącznie do pojedynczych wyrazów, ale także ich możliwych połączeń. Podczas, gdy korpusy tekstów liczą dziesiątki lub setki tysięcy wystąpień wielu popularnych

wyrazów, zasoby oferowane przez Google słowa takie liczą w dziesiątkach, czy nawet setkach milionów. Dodatkowym argumentem przemawiającym na korzyść korzystania bezpośrednio z wiedzy zawartej w Internecie jest fakt ciągłej zmienności umieszczonych tam tekstów. Język ewoluuje, a zabytki pismienictwa, które często są integralną częścią korpusów nie zawierają wielu współcześnie używanych form i wyrażen. Z drugiej strony zawierają formy, które rzadko pojawiają się w mowie potocznej dzisiejszych czasów. Ale nawet znalezienie rzadkich, czy archaicznych sformułowań jest łatwiejsze, gdy ma się do dyspozycji całą wiedzę zgromadzoną w Internecie.

Spore niebezpieczeństwo dla poprawności językowej zgromadzonych danych stanowi sposób ich za-

pisu na stronach internetowych. W szczególnych przypadkach nieporządane znaki, lub fragmenty kodu

moga znaleźć się w tworzonej bazie wiedzy. Był to główny powód, dla którego zdecydowałem się dość restrykcyjnie podejść do filtrowania nieporządkanych informacji. Nie wszystkie zagrożenia dla poprawności przechowywanych danych mają podłoże opierające się o nietypowe znaki lub ich kombinacje. Bardzo wiele portali internetowych zawiera treści niepozbawione błędów. Nie chodzi tylko o literówki w artykułach, czy esejach. Prawdziwą plagą stają się wpisy na rozmaitych forach internetowych i blogach, a także komentarze do wiadomości na portalach informacyjnych. Zarówno poziom ortografii, jak i stylistyki, a nawet logiki wielu wypowiedzi pozostawia wiele do życzenia. Wśród wpisów tego typu częstym zjawiskiem są boty reklamujące, które na wielu forach automatycznie tworzą posty składające się wyłącznie z pisanych jednym ciągiem haseł. Mając do dyspozycji wiele źródeł informacji pojedyncze

45

.

46

6.2. Dalsze możliwości rozwoju

przypadki tego typu byłyby z pewnością niezauważalne. Proceder ten jest jednak powszechny i automa-

tycznie propagowany przy pomocy złośliwych pajaków internetowych na bardzo wiele stron. Ilość tego

typu zdegenerowanych danych może mieć wpływ na proces automatycznego decydowania o poprawności wypowiedzi.

Podstawowa wersja pajaka internetowego, którą zdecydowałem się zaimplementować spełnia swoje zadanie. Wyszukiwarka Google dostarcza wyniki dość dobrej jakości, co bezpośrednio przekłada się na jakość stron odwiedzanych przez robota. Nawet jeżeli pomijane są czasem zasoby, które potencjalnie

mogłyby zostać wykorzystane do budowania bazy wiedzy, to najważniejszym aspektem działania pajaka

jest zebranie reprezentatywnych źródeł. Źródła te mają posłużyć zgromadzeniu informacji o słowach, które jak najczęściej pojawiają się w kontekście wyrazu, który w danej chwili jest obiektem zainteresowania programu.

Ekstrakcja danych i tworzenie relacji przebiega bardzo sprawnie. Waskim gardłem aplikacji jest

zdecydowanie sprawdzanie poprawności relacji. Jako, że Google nie udostępnia API pozwalającego na dostęp do wyników wyszukiwania, konieczne jest pozyskiwanie ich bezpośrednio parsując otrzymaną stronę i wynikową. Częste automatyczne zapytania są jednak bardzo skutecznie rozpoznawane przez algorytmy wyszukiwarki, co skutkuje blokadą możliwości korzystania z niej na kilka godzin, lub konieczności uzupełnienia formularza zwanego CAPTCHA (ang. Completely Automated Public Turing test to tell

Computers and Humans Apart), którego zadaniem jest weryfikacja, czy użytkownik jest człowiekiem. Zaproponowana przeze mnie metoda oceny poprawności relacji nie jest niezawodna. Jej główna niedoskonałość wiąże się z algorytmem wyszukującym Google, który możliwie jak najbardziej poszerza

zadane przez użytkownika kryteria, aby zwrócić jak największą liczbę wyników. Ignorowane są zatem wszelkie znaki interpunkcyjne (zarówno w wyszukiwanej frazie, jak i potencjalnych wynikach), które mogą mieć znaczący wpływ na jakość rezultatu. Prosty przykładem jest sprawdzanie poprawności wyrażenia kot to samochód. Grupowanie takie wydaje się niedorzeczne, ale Google znajdzie to wyrażenie w tekście zawierającym kolejno zdania: Obok przebiegł przerażony kot. To samochód przejeżdżający obok tak go przestraszył. Fakt, że między kot, a to samochód jest kropka zostanie całkowicie zignorowany.

Problem stanowią także tytuły książek, filmów lub artykułów, które często celowo formułowane są w sposób syntaktycznie lub semantycznie niepoprawny. Można by pomyśleć, że nie ma w tym nic złego. O rozmaitych komentarzach pod artykułami na stronach informacyjnych lub wpisach na forach, czy blogach można powiedzieć to samo. Niestety tytuły dzieł publicznie dostępnych są nieporównanie szerzej rozpowszechnione, co sprawia, że Google podczas wyszukiwania natrafi na nie o wiele częściej przeszukując rozmaite strony zawierające reklamy, zapowiedzi, czy recenzje.

6.2. Dalsze możliwości rozwoju

W aplikacji istnieje wciąż wiele obszarów, które można poprawić, aby polepszyć zarówno jakość otrzymywanych wyników, jak i wydajność działania. Główne elementy, które można rozwinąć w przyszłości to:

1. Algorytm sprawdzania poprawności relacji.

Ograniczenia narzucone przez Google są największym mankamentem zaproponowanego przeze mnie rozwiązania. Konieczność oczekiwania pomiędzy zapytaniami, a także możliwość blokady zapytań na kilka godzin znacznie obniża wydajność aplikacji. Niestety inne wypróbowywane przeze mnie wyszukiwarki internetowe zwracają nieporównanie mniej wyników, co w przypadku sprawdzania poprawności metoda statystyczna ma spore znaczenie. Być może w przyszłości znów dostępne będą API, które pozwoli na zautomatyzowane korzystanie z wyszukiwarki Google, co

z pewnością znacznie poprawi wydajność programu. Akceptowalnym rozwiązaniem byłby także znacznie bardziej zaawansowany mechanizm wysyłający zapytania wykorzystujący czyste zmiany

przeglądarek podawanych w nagłówku wysłanego zapytania, a także korzystający z serwerów proxy.

Ł. Wróbel

6.3. Porównanie z rozwiązaniami o podobnej tematyce 47

2. Szablony relacji

Jeżeli problem skuteczniejszego mechanizmu sprawdzania poprawności relacji zostałby rozwiązany, to nie stałoby na przeszkodzie, aby wprowadzić więcej rodzajów relacji, które znacznie wzbogaciłyby wartość lingwistyczną tworzonej bazy wiedzy. Kolejnym krokiem w rozwoju tego elementu systemu mogłoby być wyodrębnienie spójników jako oddzielnych części mowy, a nie jak obecnie traktowanie ich jako „sztywnych” elementów zespalających połączenie, między bardziej

znaczącymi częściami mowy. Na podstawie analizy składniowej tekstu możliwe byłoby dodanie automatycznego tworzenia szablonów relacji na podstawie często spotykanych struktur gramatycznych.

3. Słownik fleksyjny

Kolejną możliwością ulepszenia aplikacji jest wzbogacenie słownika fleksyjnego o nowe formy wyrazowe, a także bardziej przejrzysty sposób określania i definiowania odmiany słów. Możliwość dokładnego określenia, odmiana której cechy wyrazu nas interesuje byłaby znacznie praktyczniejsza niż wybór po indeksie w tablicy. Przykładowo mając wyraz kot w formie bazowej, chcąc uzyskać dopełniacz liczby mnogiej musimy odwołać się do określonego indeksu, który odpowiada zadanej formie. O wiele wygodniejsza byłaby możliwość podania parametrów formy, które nas interesują np. `zwrocForme("kot", DOPELNIACZ, L_MN)`.

6.3. Porównanie z rozwiązaniami o podobnej tematyce

Przykładowa, znaleziona w Internecie implementacja czatbota [17] napisana przez Tomasza Jędrzej-

ewskiego wykorzystuje prosty algorytm tworzenia wypowiedzi. Na podstawie rozmowy z użytkownikiem, zapamiętuje pojawiające się zdania w postaci skierowanego grafu przedstawionego na rysunku 6.1.

Graf ten łączy wierzchołki odpowiadające użytym wyrazom skierowanymi krawędziami, które odpowia-

dają połączeniom między nimi. Wierzchołki \$ i # oznaczają odpowiednio początek i koniec zdania. Generowanie wypowiedzi polega na wyborze losowej ścieżki pomiędzy elementami grafu od wierzchołka \$ do #.

Rysunek 6.1: Graf obrazujący połączenia między wyrazami dla czatbota zaimplementowanego przez Tomasza Jędrzejewskiego.

Algorytm tego typu jest bardzo prosty, a generowane przez czatbota wypowiedzi bardzo często są niepoprawne nie tylko semantycznie, ale również składniowo, gdyż nie wykorzystywany jest żaden słownik morfologiczny. Baza wiedzy, na podstawie której tworzone są zdania jest ograniczona wyłącznie do danych zebranych podczas rozmowy z użytkownikiem. Zebranie dużej ilości danych z wielu konwer-

Ł. Wróbel

6.3. Porównanie z rozwiązaniami o podobnej tematyce 48

sacji nie poprawi jakości działania programu, gdyż losowy wybór ścieżki w grafie może powodować zapętlenie, lub powstawanie długich pozbawionych sensu zdań.

Rozwiązanie, które zaproponowałem korzysta ze zdecydowanie bogatszej bazy wiedzy, co pozwala na generowanie wypowiedzi z wykorzystaniem słów, które nie pojawiły się w trakcie rozmowy. Połączenia między wyrazami są poprawne syntaktycznie, a także odfiltrowywane są te, uznane za nie poprawne z punktu widzenia semantyki. O poprawności decyduje w głównej mierze popularność fraz w Internecie, zatem błędy choć są nieuniknione, pojawiać się będą znacznie rzadziej. Baza wiedzy budowana jest

przez wyszukiwanie informacji związanych z konkretnym tematem. Bezpośrednio przekłada się to na jakość otrzymywanych wyników, które są tym dokładniejsze, im większa jest baza wiedzy, na podstawie której zostały wygenerowane.

Marcin Gadamer w swojej pracy[8] wykorzystał algorytm, który buduje graf przyzwyczajenia lingwistycznych na podstawie korpusów tekstów. Generowanie wypowiedzi odbywa się w interakcji z użytkownikiem, który wskazuje kolejne słowo z listy zaproponowanych możliwości. Przykładowy fragment grafu zaprezentowany został na zrzucie ekranu 6.2. W bazie danych przechowywane są trójki wyrazów, z zachowaniem kolejności, w jakiej występowały w analizowanych zdaniach. Na podstawie dwóch pierwszych form wyrazowych z trójki prezentowane są możliwe poprawne formy kolejnej, wraz z określeniem częstotliwości wystąpienia.

Rysunek 6.2: Fragment grafu LHG z aplikacji Marcina Gadamera prezentujący cztery słowa wypowiedzi i proponowane następniki.

Algorytm ten dobrze sprawdza się podczas automatycznej korekty tekstu, jednak słabiej spisywałby się w przypadku automatycznej generacji zdań. Wrażliwość na zapętlenia, a także powstawanie zbyt

długich wypowiedzi zostały rozwiązane przez interakcję z użytkownikiem. Problemem jest także ilość

Ł. Wróbel

6.3. Porównanie z rozwiązaniami o podobnej tematyce 49

danych przechowywanych w bazie wiedzy. Każde przeanalizowane zdanie rozkładane jest na $n-2$ trójek, gdzie n jest ilością wyrazów w zdaniu.

Generowana przez moją aplikację uproszczona postać grafu LHG skupia się w większym stopniu

na relacjach między konkretnymi parami wyrazów. Każda informacja odnosnie takiego połączenia jest istotna i posiada wartość dla czatbota. Nawet jeżeli do bazy danych zapisywane są niepoprawne relacje, to jest to także ważny element wiedzy programu. Na podstawie zebranych danych nie da się co prawda utworzyć pełnych zdań w języku polskim, ale można je wykorzystać podczas procesu „rozumienia”

sensu wypowiedzi użytkownika, a także generowania odpowiedzi zgodnej z tematem rozmowy.

Ł. Wróbel

7. Zakonczenie'

Cel jaki został postawiony w tej pracy to utworzenie uproszczonej wersji grafu przyzwyczajen' lingwistycznych (LHG), który obrazuje relacje zachodzace, mi'edzy wyrazami w j'ęzyku polskim, w szczeg'olności modelujace, grupowanie wokół okreslonych' cech i budow,e hierarchicznej struktury w obr'ebie danej grupy. Osiagn'iecie tego celu sprowadzone zostało do utworzenia trzech głównych modułów aplikacji oraz interfejsu uzytkownika'.

1.Pajak, internetowy

Aby wykorzystac' wiedz, zgromadzona, w Internecie niezb'edne było wykorzystanie pajaka, inter-netowego, który analizujace, tres'c' odwiedzanych stron pozyskuje informacje, które nast'epnie wy-korzystywane sa, do utworzenia bazy wiedzy. Utworzony przeze mnie robot wykorzystujace, wyniki z wyszukiwarki Google pobiera poprawne zdania w j'ęzyku polskim i zapisuje je w bazie danych stale powi'ekszajace, dost'epne do wykorzystania repozytorium wiedzy. Jego zadaniem jest takze' rozklad zdan' na pojedyncze słowa, które nast'epnie wykorzystywane sa, w mechanizmie tworzenia relacji. Oprócz pajaka, dodałem takze' możliwo's'c' wykorzystania korpusów tekstów znajdujacych, si'e w plikach tekstowych, aby Internet nie był jedynym źródłem' wiedzy, z którego aplikacja moze' korzystac'.

2. Kreator relacji

Zasadniczy mechanizm działania aplikacji miał za zadanie z wyników dostarczonych przez pajaka, utworzyć strukturę, która modeluje relacje między słowami w języku polskim. Zastosowany przeze mnie algorytm zlicza ilość występień wyrazów w kontekście słowa, wokół którego tworzony jest graf relacji. Najczęściej występujące są następnie podstawiane do pre-definiowanych szablonów, które zawierają oba elementy relacji w określonych formach fleksyjnych tworzących poprawne

syntaktycznie wyrażenie. Poprawność takiego wyrażenia jest następnie sprawdzana za pomocą wyszukiwarki Google. Ilość zwróconych wyników określa „popularność” frazy w Internecie, co przekłada się na jej poprawność zarówno syntaktyczną, jak i semantyczną.

Ważną funkcją kreatora relacji jest także proces pseudo-logicznego rozumowania. Wykorzystując znane już poprawne relacje i ich hierarchiczną strukturę sprawdzane są potencjalne połączenia między słowami. Elementy należące do wspólnej grupy zestawiane są z cechami składników tejże grupy. Podobnie wyrazy, które posiadają wspólne cechy sprawdzane są pod kątem przynależności do wspólnej grupy.

3. Specjalistyczny graf przyzwyczajen lingwistycznych

Zadaniem uproszczonego grafu LHG było zilustrowanie modelu relacji w sposób czytelny dla

użytkownika. Zostały na nim pokazane różne rodzaje połączeń między wyrazami oraz wartości

liczbowa określająca ich poprawność. Relacje zapisane zostały jako wyrażenia składające się z

dwóch wyrazów w odpowiedniej formie połączonych ewentualnym spójnikiem i zaprezentowane

na skierowanych krawędziach grafu. Na żądanie użytkownika wyświetlone mogą zostać także relacje uzyskane z procesu tak zwanej „przechodności”, czyli wykonanego wcześniej pseudo-logicznego procesu rozumowania i analizy istniejących połączeń.

50

51

Popularność tematyki sztucznej inteligencji jest obecnie bardzo duża, a zainteresowanie urządzeniami i aplikacjami zdolnymi porozumiewać się z człowiekiem jest szczególnie wysokie. Roboty, które

komunikują się z użytkownikiem i na bieżąco uczą się poszerzając swoją bazę wiedzy przestają być wytworem filmów science fiction, a zaczynają pojawiać się już jako gotowe produkty. Wciąż jeszcze charakteryzowane są raczej jako nowinki techniczne, lub modne gadzety, ale wzrastające zapotrzebowanie może sprawić, że już niedługo napotykac je będziemy w codziennym życiu.

Zaproponowane przeze mnie rozwiązanie jest zaledwie małym krokiem na drodze do stworzenia inteligentnego czatbota. Daje ono jednak w miarę zadowalające rezultaty, które mogą zostać wykorzystane przy tworzeniu nie tylko poprawnych gramatycznie, ale i logicznie zdań. Wizja robota, który nie

tylko formułuje poprawne zdania, ale także rozumie kontekst poszczególnych słów, jakich używa, i jak i tych, które otrzymuje od rozmówcy, staje się coraz bardziej realna. Pomimo wielu problemów, jakie wciąż stoją przed programistami, powstanie czatbota, który pozytywnie przejdzie test Turinga jest tylko kwestią czasu.

Ł. Wróbel

Spis rysunków

| | |
|--|--|
| 2.1 | |
| Schemat działania pajaka internetowego | |
| 11 | |
| 2.2 | |
| Graficzny interfejs metawyszukiwarki WebCrawler. | |
| 12 | |
| 2.3 | |
| Graficzny interfejs aplikacji Web Fountain. | |
| 13 | |
| 2.4 | |
| Wynik działania programu WGet. | |
| 14 | |
| 2.5 | |
| Zrzut ekranu obrazujący działanie aplikacji Online Data Extractor. | |
| 15 | |
| 2.6 | |
| Zrzut ekranu obrazujący działanie aplikacji HTML Text Extractor. | |
| 15 | |
| 2.7 | |
| Kolokacje wyrazu niebo na podstawie Narodowego Korpusu Języka Polskiego. | |
| 16 | |
| 3.1 | |
| Przykład oddziaływania akomodacyjnego dla leksemów kot i is'c. | |
| 19 | |
| 3.2 | |
| Schemat uogólnienia oddziaływań akomodacyjnych między podmiotem i orzeczeniem. . | |

19

3.3

Schemat uogólnienia oddziaływan' akomodacyjnych mi'edzy czasownikiem i rzeczowni-

kiem.

19

3.4

Schemat uogólnienia oddziaływan' akomodacyjnych mi'edzy przymiotnikiem i rzeczow-

nikiem.

19

3.5

Przykład akomodacji typu frazy zdaniowej.

19

4.1

Diagram klas

22

4.2

Okno połączenia
z baza' danych

23

4.3

Główne okno programu

23

4.4

Wybór plików zawierających
relacje lub korpusy tekstów.

24

4.5

Schemat działania ekstraktora danych

26

4.6

Szczegóły kraw'edzi - grupowanie, posiadanie

28

4.7

Szczegóły kraw'edzi - określenie'

28

4.8

Lista relacji

29

4.9

Wybór trybu online

29

4.10

Wybór trybu offline

32

4.11

Wybór trybu pracy z symulatorem

| | |
|---|----|
| 33 | |
| 4.12 | |
| Przykładowy graf z wierzchołkami będącymi różnymi częściami mowy. | 33 |
| 4.13 | |
| Informacja o konkretnym wierzchołku wyświetlona na grafie. | 34 |
| 4.14 | |
| Widok pokazujący spropagowane relacje na grafie. | 35 |
| 5.1 | |
| Przykładowe dane wejściowe dla programu. | 36 |
| 5.2 | |
| Graf relacji dla słowa gepard. | 37 |
| 5.3 | |
| Przykładowa lista relacji wygenerowanych dla słowa gepard. | 38 |
| 5.4 | |
| Graf relacji dla słowa strus'. | 39 |
| 5.5 | |
| Przykładowa lista relacji wygenerowanych dla słowa strus'. | 40 |
| 5.6 | |
| Przykładowe dane wejściowe dla generacji dodatkowych relacji przechodnich. | 41 |
| 5.7 | |
| Przykładowy graf relacji dla kilku słów. | 42 |

52

SPIS RYSUNKÓW

53

| | | |
|-----|--|----|
| 5.8 | Przykładowy graf relacji wraz z relacjami przechodnimi dla kilku słów. | 43 |
| 5.9 | Przykładowa lista relacji wygenerowanych dla kilku różnych słów. | 44 |
| 6.1 | Graf obrazujący połączenia między wyrazami dla czatbota zaimplementowanego przez Tomasza Jędrzejewskiego. | 47 |
| 6.2 | Fragment grafu LHG z aplikacji Marcina Gadamera prezentujący cztery słowa wypowiedzi. | |

| | |
|------------------------------------|----|
| dzi i proponowane następniki. | 48 |
|------------------------------------|----|

Ł. Wróbel

Bibliografia

[1]Geofferey Landis. Out of Sight, Out Of Mind. <http://www.geoffreylandis.com/sight.htm>, 1999.

[2]GNU. WGet. <http://www.gnu.org/software/wget/>, 2010.

[3]Google Code. crawler4j. <http://code.google.com/p/crawler4j/>, 2010.

[4]ICONICO. HTML Text Extractor. <http://www.iconico.com/html extractor/>, 2010. [dostęp: 2010-05-08].

[5]IPI PAN. Korpus IPI PAN. <http://korpus.pl/>, 2008.

[6]M. Junczys-Dowmunt. Zastosowanie automatów skoczonych stanów w przetwarzaniu języków naturalnych, 2008.

[7]Z. Klemensiewicz. Zarys składni polskiej. Warszawa, 1961.

[8]M. Gadamer. Automatyczna kontekstowa korekta tekstów na podstawie Grafu Przyzwyczajen Lingwistycznych (LHG) zbudowanego przez robota internetowego dla języka polskiego, Kraków, 2009.

[9]M. Gadamer, A. Horzyk. Automatyczna kontekstowa korekta tekstów z wykorzystaniem grafu LHG. Computer Science AGH, Vol. 10, pp. 39-57, Kraków, 2010.

- [10]Online Data Extractor. Online data extractor tool. <http://www.onlinedataextractor.com>, 2010.
- [11]M. Piasecki. Cele i zadania lingwistyki informatycznej. Metodologie je_zykoznastwa. Współczesne tendencje i kontrowersje, 2007.
- [12]Piotr Pe_zik, Narodowy Korpus Je_zyka Polskiego . PELCRA. <http://www.nkjp.uni.lodz.pl/>, 2010. [doste_p: 2010-03-10].
- [13]Portal naukowy.pl. Gramatyka generatywna. http://www.naukowy.pl/encyklopedia/Gramatyka_generatywna, 2010.
- [14]Portal naukowy.pl. Historia lingwistyki. http://www.naukowy.pl/encyklopedia/Historia_lingwistyki, 2010.
- [15]Z. Saloni and M. Swidzinski. Składnia współczesnego je_zyka polskiego. PWN, 1985.
- [16]The Web Robots Pages. About /robots.txt. <http://www.robotstxt.org/robotstxt.html>, 2010.
- [17]Tomasz J_edrzejewski. Piszemy czatbota. <http://www.eioba.pl/a/1nrd/piszemy-chatbota>, 2007.
- [18]Wapedia. Klasyfikacja cz_esci´ mowy. http://wapedia.mobi/pl/Klasyfikacja_cz%20C4%99%C5%9Bci_mowy, 2010.

54

BIBLIOGRAFIA

55

- [19]Wikipedia. Ferdinand de saussure. http://pl.wikipedia.org/w/index.php?title=Ferdinand_de_Saussure&oldid=24025783, 2010.
- [20]Wikipedia. Historia je_zykoznastwa. http://pl.wikipedia.org/w/index.php?title=Historia_%20C4%99zykoznawstwa&oldid=23236068, 2010.
- [21]Wikipedia. IBM WebFountain. http://en.wikipedia.org/w/index.php?title=IBM_WebFountain&oldid=391584874, 2010.
- [22]Wikipedia. Je_zykoznastwo. <http://pl.wikipedia.org/w/index.php?title=%20C4%99zykoznawstwo&oldid=23891961>, 2010.
- [23]Wikipedia. Lingwistyka komputerowa. http://pl.wikipedia.org/w/index.php?title=Lingwistyka_komputerowa&oldid=23235735, 2010.
- [24]Wikipedia. Web crawler. http://en.wikipedia.org/w/index.php?title=Web_crawler&oldid=399429057, 2010.
- [25]Wikipedia. Web crawler. http://en.wikipedia.org/w/index.php?title=Web_crawler&oldid=398058393,

2010.

[26] Wydawnictwo Naukowe PWN. Korpus Języka Polskiego Wydawnictwa Naukowego PWN. <http://korpus.pwn.pl/>, 2010.

[27] P. Zmigrodzki. Związki między składnikami w zdaniu polskim. Szkoła języka i kultury polskiej
Uniwersytet Śląski, Katowice.

Ł. Wróbel