

Politechnika Warszawska

W Y D Z I A Ł E L E K T R Y C Z N Y



Instytut Elektrotechniki Teoretycznej i Systemów Informacyjno-Pomiarowych
Zakład Elektrotechniki Teoretycznej i Informatyki Stosowanej

Praca dyplomowa magisterska

na kierunku Informatyka
w specjalności Inżynieria oprogramowania

Analiza wiadomości z serwisów społecznościowych na użytek
rozpoznawania nastrojów.

Jakub Rzepliński

Numer albumu: 233608

promotor
dr inż. Marcin Kołodziej

Warszawa, 2018

Streszczenie

Abstract

Oświadczenie autora (autorów) pracy

Świadom odpowiedzialności prawnej oświadczam, że przedstawiona praca dyplomowa:

- została napisana przeze mnie samodzielnie i nie zawiera treści uzyskanych w sposób niezgodny z obowiązującymi przepisami,
- nie była wcześniej przedmiotem procedur związanych z uzyskaniem tytułu zawodowego lub stopnia naukowego w wyższej uczelni

Oświadczam ponadto, że niniejsza wersja pracy jest identyczna z załączoną wersją elektroniczną.

.....
data

.....
podpis autora (autorów) pracy

Oświadczenie

Wyrażam zgodę / nie wyrażam zgody^{*1} na udostępnianie osobom zainteresowanym mojej pracy dyplomowej. Praca może być udostępniana w pomieszczeniach biblioteki wydziałowej. Zgoda na udostępnienie pracy dyplomowej nie oznacza wyrażenia zgody na jej kopiowanie w całości lub w części.

Brak zgody nie oznacza ograniczenia dostępu do pracy dyplomowej osób:

- reprezentujących władze Politechniki Warszawskiej,
- członków Komisji Akredytacyjnych,
- funkcjonariuszy służb państwowych i innych osób uprawnionych, na mocy odpowiednich przepisów prawnych obowiązujących na terenie Rzeczypospolitej Polskiej,

do swobodnego dostępu do materiałów chronionych międzynarodowymi przepisami o prawach autorskich. Brak zgody nie wyklucza także kontroli tekstu pracy dyplomowej w systemie antyplagiatowym.

.....
data

.....
podpis autora (autorów) pracy

^{*1} - niepotrzebne skreślić

Spis treści

1	Wstęp	1
2	Serwis społecznościowy Twitter	3
3	Analiza sentymentu wypowiedzi	11
4	Wymagania funkcjonalne i нефункционалне	15
5	Wybór narzędzi	19
6	Aplikacja Twitter Analyser	23
7	Badania i wnioski	27
8	Podsumowanie	29

Rozdział 1

Wstęp

Czynnikiem powodującym rozwój świata są informacje i umiejętne ich wykorzystanie. Od jakiegoś czasu źródłem informacji stał się internet, a szczególnie portale społecznościowe, gdzie ludzie wymieniają się informacjami na różne tematy. Wiele globalnych instytucji korzysta z tego typu serwisów, wydając za ich pośrednictwem oficjalne komunikaty. Ilość informacji generowanych przez użytkowników takich serwisów jest tak duża, że wymaga użycia specjalnych narzędzi *Big Data*. Jednak decydującym czynnikiem jest umiejętne wykorzystanie wiedzy zawartej w zgromadzonych informacjach, do czego potrzebny jest czynnik ludzki.

Jednym z ciekawych sposobów wykorzystania informacji z serwisów takich jak np. *Twitter* jest badanie wydźwięku wpisów zamieszczanych przez użytkowników zwanego od angielskiego sformułowania *sentiment analysis* sentymentem. Aby było to możliwe potrzebne jest przetworzenie języka, którym posługują się ludzie czyli *języka naturalnego* na postać, którą mogą obsługiwać systemy NLP (ang. *Natural Language Processing*), a następnie określenie sentymentu za pomocą podejścia słownikowego lub technik maszynowego uczenia się (ang. *machine learning*). Wykorzystanie takich informacji pozwala zbadać opinie ludzi na praktycznie dowolny temat.

Cel pracy

Celem niniejszej pracy jest przedstawienie możliwości wykorzystania analizy sentymentu wypowiedzi użytkowników serwisu społecznościowego Twitter. Stworzone powinno zostać rozwiązanie umożliwiające przetwarzanie w czasie rzeczywistym wiadomości publikowanych w tym serwisie zawierających wybrane słowo kluczowe. Napisany system powinien wyświetlać statystyki wydźwięku opinii użytkowników tego serwisu.

Zakres pracy

Pracę dyplomową można podzielić na trzy części: przegląd najistotniejszych informacji o poruszanych tematach, opis napisanej aplikacji oraz omówienie przeprowadzonych badań i wyciągnięte z nich wnioski.

W pierwszej części w rozdziale 2. opisano zasadę działania serwisu społecznościowego Twitter wraz z udostępnionym przez niego programowalnym interfejsem (ang. *API - Application Programming Interface*). W rozdziale 3. zostały podane informacje na temat analizy sentymentu wypowiedzi w systemach komputerowych.

W drugiej części w rozdziale 4. przedstawiono wymagania funkcjonalne i нефункционалне

postawione przed zbudowaną aplikacją. W rozdziale 5. opisano wybrane do jej implementacji narzędzia, biblioteki oraz języki programowania. W rozdziale 6. znajduje się omówienie stworzonego systemu *Twitter Analyser*.

W ostatniej części w rozdziale 7. omówiono badania sentymentu wypowiedzi przeprowadzone podczas dwóch wydarzeń - meczu piłki nożnej FC Barcelona - Real Madryt oraz Święta Dziękczynienia. W ostatnim rozdziale 8. znajduje się podsumowanie całej pracy dyplomowej.

Rozdział 2

Serwis społecznościowy Twitter

Serwis społecznościowy Twitter jest globalnym serwisem internetowym służącym głównie do zamieszczania wiadomości tzw. *tweet*, które użytkownicy tego serwisu mogą czytać, komentować lub przekazywać dalej. Od kilku lat Twitter jest serwisem, gdzie dochodzi do wymiany zdań na różny temat dotyczących np. polityki, sportu, produktów, wydarzeń społecznych, a profile posiada wiele osób znanych publicznie oraz instytucji.

Serwis ten, nazywany SMS-em internetu, został założony w 2006 r. w Stanach Zjednoczonych przez Jacka Dorsey'a, Ev Williamsa, Noah Glassa oraz Biza Stone'a i od początku powstania sukcesywnie zwiększał swoją popularność poprzez wzrost liczby użytkowników odwiedzających jego witrynę oraz wysyłających wiadomości [1]. W 2012 r. osiągnął ponad 100 milionów użytkowników, którzy zamieszczali łącznie ponad 340 milionów wiadomości dziennie oraz obsługiwał średnio około 1.6 miliarda wyszukujących zapytań dziennie. W 2013 r. Twitter stał się jedną z najczęściej odwiedzanych stron w całym internecie. W tym samym roku inżynierowie Twittera podali informację, że serwis ten obsługuje ok. 143 tys. wiadomości na sekundę. Na początku 2016 r. serwis posiadał ponad 319 milionów użytkowników aktywnych podczas każdego miesiąca. Od listopada 2013 r. akcje Twittera są obecne na nowojorskiej giełdzie.

Tweet, czyli krótka wiadomość tekstowa, była początkowo ograniczona do 140 znaków, ale limit ten został podwojony w 2017 r. dla wszystkich języków oprócz chińskiego, japońskiego i koreańskiego. Użytkownicy mają możliwość wyróżniania wybranych przez siebie tematów przez dodanie do nich znaku '#', co czyni takie wyrażenie tagiem. Inną możliwością oferowaną przez Twittera jest odpowiadanie innym użytkownikom lub zamieszczenie referencji do nich przez dodanie znaku '@' poprzedzającego nazwę profilu innej osoby.

Serwis społecznościowy Twitter opierał się początkowo o typową architekturę trójwarstwową składającą się z warstwy prezentacji, logiki biznesowej oraz warstwy danych. Do napisania tej aplikacji został użyty framework Ruby on Rails wykorzystujący język Ruby, a warstwa bazodanowa opierała się o technologię MySQL. Jednak wraz ze wzrostem ilości przetwarzanych danych inżynierowie Twittera podjęli decyzję w 2011 r. o zmianie technologii na język Scala, który działa na maszynie wirtualnej Javy oraz zrezygnowano z dotychczasowej architektury na rzecz budowy rozproszonych serwisów komunikujących się między sobą. Wraz z przeprowadzonymi zmianami zanotowano ponad 10-krotne polepszenie obsługi tweetów.

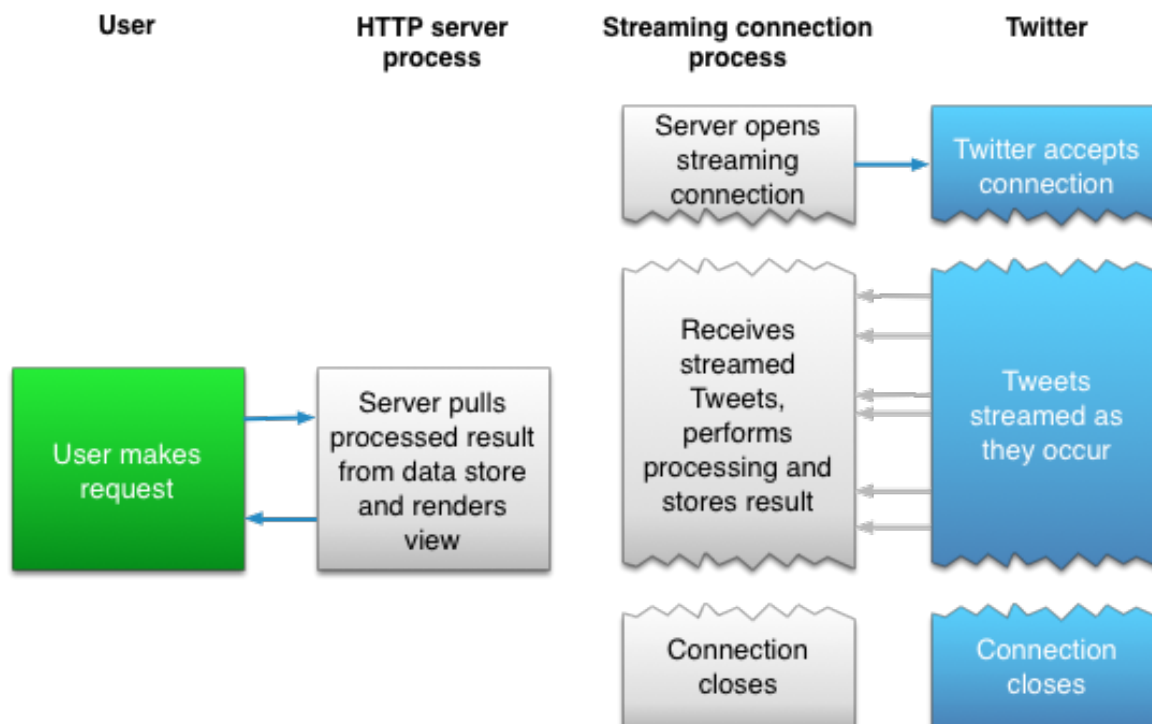
Twitter jest platformą otwartą i udostępnia programowalny interfejs (ang. API - *Application Programming Interface*) w dwóch postaciach: Search API oraz Streaming API.

Programiści korzystający z Search API są w stanie uzyskać dostęp tylko do danych histo-

rycznych, które zostały już wcześniej zamieszczone na łamach serwisu Twitter. Natomiast w przypadku Streaming API dostajemy możliwość śledzenia strumienia danych, które są do naszego dostępu nawet już kilka sekund po zamieszczeniu w serwisie Twitter. Po podłączeniu do takiego strumienia możemy cały czas obserwować nowe wiadomości. Przyjęło się stosować nazywnictwo, że analiza Search API to analiza *back in time*, a Streaming API to śledzenie *real time*.

Obie formy API wymagają wcześniejszej rejestracji na stronie <https://developer.twitter.com/en/apply-for-access> przeznaczonej dla deweloperów zainteresowanych wykorzystywaniem Twitter API. Po pomyślnym przejściu rejestracji dostajemy dane, które po nawiązaniu połączenia z serwisem Twitter umożliwiają mu jednoznacznie określić, że możemy uzyskać dostęp do API.

Search API powstało z wykorzystaniem standardu REST - *Representational State Transfer*. Oba rodzaje API wykorzystują protokół HTTP: do poprawnego działania Streaming API potrzebne jest ciągłe połączenie HTTP, a w przypadku drugiego z nich każda operacja jest wykonywana przy nawiązaniu oddzielnego połączenia. Schemat działania obu wersji API został przedstawiony na rysunku 2.1..



Rysunek 2.1: Schemat działania dwóch rodzajów programistycznego interfejsu API udostępnianego przez serwis społecznościowy Twitter: *Search API* i *Streaming API* [2].

Search API posiada ściśle określone parametry, które mogą być przesłane w żądaniu. W tabeli 2.3 zaprezentowano ich wykaz.

Streaming API nie posiada takich ograniczeń. W języku programowania Java dostępny jest pakiet *twitter4j* zawierający interfejsy *User* oraz *Status*, na które mapowane są przychodzące ze strumienia informacje. W tabelach 2.1 i 2.2 zamieszczono ich dokumentację.

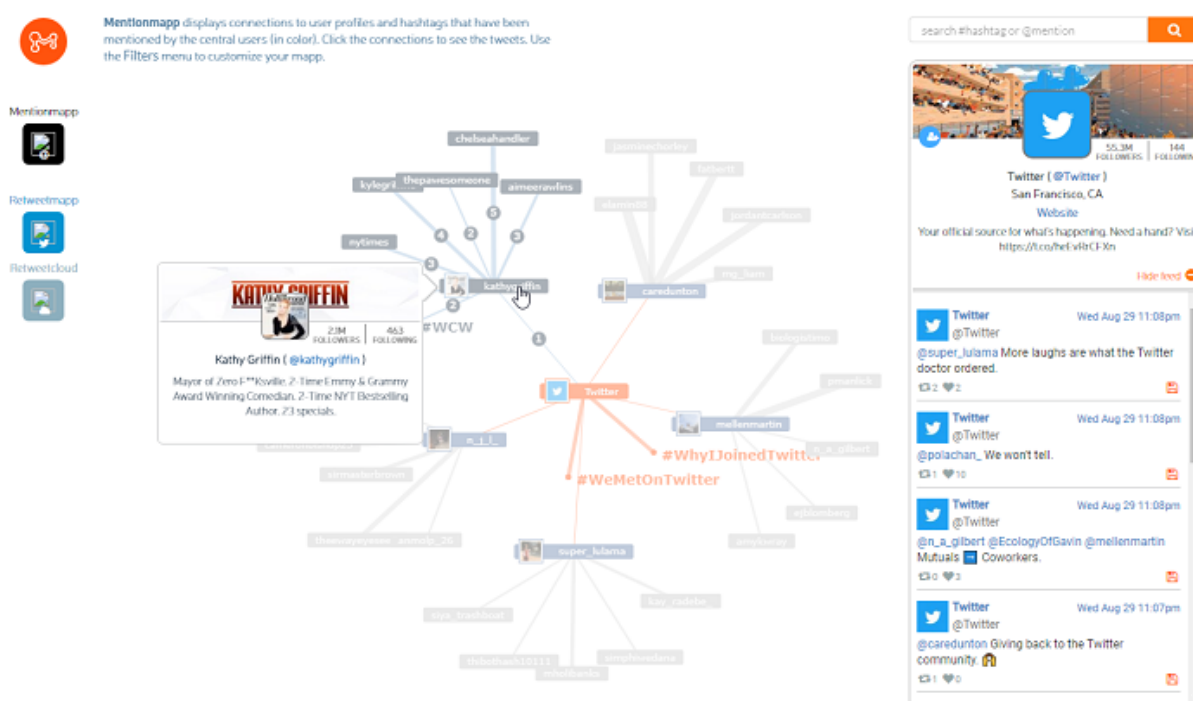
Korzystając z Search API mamy możliwość wysłania 720 zapytań na godzinę, a maksymalna ilość wiadomości jaka może być zwrócona na jedno zapytanie to 100 [4]. Jeśli wykorzystalibyśmy ten limit w maksymalny sposób to daje nam to 72 000 wiadomości na godzinę. W

przypadku Streaming API głównym ograniczeniem jest dostęp do ok. 1 % danych ze strumienia, a maksymalna ilość wiadomości w czasie jednej minuty to 3 000. W przypadku tego API w ciągu godziny możemy uzyskać 180 000 wiadomości na godzinę. Są to ograniczenia, które obowiązują dla rozwiązań typu *open-source*.

Udostępnienie API przez serwis społecznościowy Twitter oraz rosnące znaczenie danych generowanych przez użytkowników tego serwisu spowodowało, że wiele firm oraz instytucji zaczęło przywiązywać dużą wagę do analizy opinii wyrażanych na swój temat lub na tematy pokrewne, zainteresowania pewnymi tematami oraz kształtujących się trendów. Dlatego powstały aplikacje internetowe służące do wyświetlania takich informacji i przeprowadzające wstępną analizę zebranych danych.

Narzędziem wartym omówienia jest *Mentionmap*. Jest to aplikacja, która wyróżnia się spośród innych tym, że rysuje wykres powiązań pomiędzy użytkownikami wchodzącymi w interakcje z danym profilem, a także z ich profilami [5]. Posiada także podstronę umożliwiającą badanie nastrojów społecznych osób zamieszczających wiadomości z konkretnym słowem kluczowym. Nastroje te są przedstawione w postaci chmury nazw kont użytkowników, gdzie kolor nazwy zależy od nastroju prezentowanego przez użytkownika pod kątem słowa kluczowego. Interfejs graficzny narzędzia Mentionmap został przedstawiony na rysunkach 2.2 i 2.3..

Podsumowując warto zauważyć, że nie ma obecnie na rynku aplikacji, która umożliwiałaby śledzenie występowania dowolnego słowa w wiadomościach zamieszczanych w czasie rzeczywistym w serwisie Twitter, rysowałaby wykres zależności pomiędzy użytkownikami tego serwisu, analizowałaby nastroje społeczne użytkowników z wyświetleniem informacji o nastroju wyrażanym w poszczególnych wiadomościach oraz pozwalałaby na analizowanie danych historycznych i zamieszczanych w czasie rzeczywistym. Taka sytuacja pozwala na utworzenie nowej aplikacji, która udostępniałaby te funkcjonalności.



Rysunek 2.2: Narzędzie Mentionmap - wykres powiązań [5].

Tablica 2.1: Dokumentacja interfejsu Status z pakietu twitter4j [7].

Typ zwracany	Nazwa metody	Opis
long[]	getContributors()	
java.util.Date	getCreatedAt()	zwraca datę utworzenia wiadomości
long	getCurrentUserRetweetId()	zwraca id użytkownika, którego wiadomość została podana dalej
int	getDisplayTextRangeEnd()	
int	getDisplayTextRangeStart()	
int	getFavoriteCount()	zwraca informację ile razy została polubiona wiadomość
GeoLocation	getGeoLocation()	zwraca lokalizację użytkownika zamieszczającego tą wiadomość
long	getId()	zwraca id wiadomości
java.lang.String	getInReplyToScreenName()	zwraca nazwę użytkownika, do którego kierowana jest odpowiedź
long	getInReplyToStatusId()	zwraca id wiadomości, do którego kierowana jest odpowiedź
java.lang.String	getLang()	zwraca język zamieszczonej wiadomości
Place	getPlace()	zwraca obiekt Place przypisany do tej wiadomości
Status	getQuotedStatus()	zwraca obiekt Status cytowanej wiadomości
long	getQuotedStatusId()	zwraca id cytowanej wiadomości
URLEntity	getQuotedStatusPermalink()	zwraca obiekt URLEntity reprezentujący bezpośredni odnośnik do cytowanej wiadomości
int	getRetweetCount()	zwraca ile razy wiadomość została podana dalej
Status	getRetweetedStatus()	zwraca oryginalny status, który jest podany dalej w tej wiadomości
Scopes	getScopes()	zwraca obiekt typu Scopes posiadający informację o id miejsc, do których odnosi się ta wiadomość
java.lang.String	getSource()	zwraca źródło wiadomości
java.lang.String	getText()	zwraca tekst wiadomości
User	getUser()	zwraca obiekt typu User powiązany z tą wiadomością
java.lang.String[]	getWithheldInCountries()	zwraca tablicę nazw krajów, w których wiadomość została wstrzymana

boolean	isFavorited()	zwraca informację czy wiadomość została polubiona
boolean	isPossiblySensitive()	zwraca informację czy wiadomość zawiera link do chronionych informacji
boolean	isRetweet()	zwraca informację czy tweet jest podaną dalej wiadomością
boolean	isRetweeted()	informuje czy wiadomość jest podana dalej
boolean	isRetweetedByMe()	informuje czy wiadomość jest podana dalej przez tego użytkownika
boolean	isTruncated()	informuje czy wiadomość jest skrócona (zakończona znakiem "...")

Tablica 2.2: Dokumentacja interfejsu User z pakietu twitter4j [6]. Zamieszczone zostały tylko najważniejsze z metod.

Typ zwracany	Nazwa metody	Opis
java.util.Date	getCreatedAt()	zwraca datę utworzenia profilu użytkownika
java.lang.String	getDescription()	zwraca opis konta użytkownika
java.lang.String	getEmail()	zwraca adres e-mail powiązany z tym kontem
int	getFavouritesCount()	zwraca liczbę wiadomości, którą polubił ten użytkownik
int	getFollowersCount()	podaje ilość użytkowników śledzących profil
int	getFriendsCount()	podaje ilość śledzonych profili
long	getId()	zwraca id użytkownika
java.lang.String	getLang()	zwraca język preferowany przez użytkownika
java.lang.String	getLocation()	zwraca lokalizację użytkownika
java.lang.String	getName()	podaje nazwę użytkownika
java.lang.String	getScreenName()	zwraca nazwę konta
Status	getStatus()	zwraca obiekt typu Status reprezentujący wiadomość wysłaną przez użytkownika
int	getStatusesCount()	podaje ilość wiadomości wysłanych przez użytkownika
java.lang.String	getTimeZone()	podaje strefę czasową użytkownika
java.lang.String	getURL()	zwraca URL do profilu
boolean	isVerified()	podaje informację czy profil jest zweryfikowany

Tablica 2.3: Parametry żądania Twitter Search API [3].

Parametr	Wymagany/Opcjonalny	Opis	Przykład
q	wymagany	zapytanie wyszuki- jące o maksymalnej długości 500 znaków	nasa
geocode	opcjonalny	zwraca wiadomości użytkowników od- dalonych o podany promień od podanej szerokości i długo- ści geograficznej, promień może być podany w milach lub kilometrach	37.781157 -122.398720 1mi
lang	opcjonalny	ogranicza wiadomo- ści do wybranego ję- zyka spośród dostęp- nych kodów ISO 639- 1	pl
locale	opcjonalny	specyfikuje język wy- syłanego zapytania, obecnie tylko <i>ja</i> jest skuteczny	ja
result_type	opcjonalny	określa typ zwraca- nych wiadomości, obecnie dostępne są trzy wartości tego parametru: <i>recent</i> (zwracane są najnowsze wia- domości), <i>popular</i> (zwracane są naj- bardziej popularne wiadomości) <i>mixed</i> (wartość domyślna, zwracane wyniki obejmują najnowsze i najbardziej popu- larne wiadomości)	mixed
count	opcjonalny	specyfikuje ilość zwracanych wiado- mości; maksymalna wartość to 100, a domyślna to 15	100

Rozdział 3

Analiza sentymentu wypowiedzi

Zagadnieniem dobrze ilustrującym postęp technologiczny jest prędkość rozprzestrzeniania się informacji. W starożytności informacja była przekazywana zazwyczaj przez posłańców, od których losów zależało czy i kiedy informacja zostanie dostarczona do nadawcy. Tak było w przypadku posłańca Filippidesa, który według legendy podążał po wygranej bitwie pod Maratonem do Aten, aby uprzedzić Greków przed zbliżającym się atakiem Persów [8]. W dzisiejszych czasach coraz częściej informacja jest dostępna na ekranach urządzeń mobilnych takich jak telefon komórkowy w czasie rzeczywistym zaraz po wystąpieniu jakiegoś zdarzenia. Dlatego bardzo ważne stało się monitorowanie opinii i nastrojów społecznych. Dostępne obecnie narzędzia pozwalają klasyfikować wydzwięk tekstu jako pozytywny, negatywny albo neutralny.

Podstawową operacją języka naturalnego jest ustalenie znaczenia zdania zwane także semantyką. Jego zrozumienie wymaga w systemach informatycznych określenia znaczenia zdania oraz opracowania metody jego zapisu do czego potrzebne są opis faktów danych w zdaniu oraz wyciąganie wniosków z posiadanych już informacji.

Ustalenie znaczenia wyrażenia zdania w języku naturalnym polega na wyznaczeniu występujących w nim obiektów oraz zachodzących między nimi relacji [9]. Do analizy pełnego znaczenia wymagana jest szeroka wiedza o świecie odnosząca się do konkretnego kontekstu. Na tym etapie przetwarzania tekstu pomocne są informacje o powiązaniach między słowami. W zależności od zastosowania przydatna może być na przykład wiedza o tym, że kot jest ssakiem, a rafineria rodzajem przedsiębiorstwa.

Do reprezentacji semantyki języka naturalnego można wykorzystać mechanizmy formalne, które odpowiadają praktycznym potrzebom dokonania interpretacji. Jednym z przykładów takich mechanizmów jest *rachunek predykatów I rzędu*, który pozwala zapisać czy fakt jest prawdziwy lub fałszywy, umożliwia zapisywanie pytań za pomocą użycia zmiennych, a także posiada opracowane metody wnioskowania. Inną metodą reprezentowania znaczenia zdania jest *Teoria Reprezentacji Dyskursu - DRT* (ang. *Discourse Representation Theory*), która polega na przekształcaniu drzewa rozbioru znaczenia zdania na strukturę zwaną *DRS* (ang. *Discourse Representation Structure*). Niestety żaden ze znanych mechanizmów nie jest w stanie odzwierciedlić całej złożoności procesów powiązanych z rozumieniem języka naturalnego. Przyjmując jedną z istniejących metod lub tworząc nową musimy zmierzyć się z wybraniem mechanizmu, który wobec postawionych wymagań najlepiej poradzi sobie z rozległością skali znaczeń jakie chcemy reprezentować, stopniem skomplikowania semantyki oraz kosztem jej uzyskania. Dlatego większość systemów informatycznych nie korzysta z wyrafinowanych reprezentacji znaczenia, ale ogranicza ją do najprostszych scenariuszy.

Analiza tekstu jest złożonym zadaniem. Wymaga zrozumienia zależności pomiędzy występującymi faktami oraz do jego pełnego zrozumienia potrzebna jest wiedza, którą dysponuje człowiek. Dlatego analiza tekstu jest uznawana za problem *AI-zupełny* (ang. *AI - Artificial Intelligence* lub po polsku *SI - Sztuczna Inteligencja*). Zrozumienie kolejnej części tekstu wymaga umiejętnego wyciągania wniosków z poprzedniej części tekstu i powiązania ich z wiedzą nabytą z innych źródeł.

Opis dłuższego tekstu wymaga odtworzenia powiązań pomiędzy kolejnymi zdaniami, ale z powodu ilości możliwych kombinacji sekwencji zdań można opisać tylko wybrane zjawiska oraz wykluczyć niektóre typy powiązań. Kolejną trudnością jest rozstrzyganie do jakich obiektów odnoszą się wyrażenia wskazujące, ponieważ w tekstach stosowane są sposoby opisu tego samego obiektu w różny sposób np. wszystkie frazy odnoszące się do tej samej osoby - *Janek, kolega Maćka, mały chłopiec, ten z prawej, najmłodszy w rodzinie*. Inną ważną kwestią są części wpływające na ciągłość tekstu (nawiązanie do poprzedniego tekstu lub wypowiedzi, ciągłość opisów).

Analiza nastrojów społecznych wypowiedzi jest możliwa po poddaniu tekstu zamianie na reprezentację, którą mogą posługiwać się systemy informatyczne. Poniżej przedstawiono opis kilku takich reprezentacji:

- **bag of words** [10] - najpopularniejszy sposób reprezentacji tekstu w postaci zestawu wyrazów z przyporządkowanymi im liczbami wystąpień w tekście np. zdanie *Jan lubi oglądać filmy, a Maria także lubi je oglądać* można zapisać w notacji *JSON* (ang. *JSON - JavaScript Object Notation*) jako `{"Jan": 1, "lubi": 2, "oglądać": 2, "filmy": 1, "a": 1, "Maria": 1, "także": 1, "je": 1}`; innym przykładem zastosowania takiej reprezentacji jest wykorzystanie jej w mechanizmach odpowiadających za filtrowanie wiadomości e-mail.; wadą takiego podejścia jest utrata informacji o kolejności wyrazów w zdaniu i powiązaniach między nimi.
- **reprezentacja wektorowa** (ang. *vector space model*) [11, 12] - dokument tekstowy reprezentowany jest w postaci wektorów częstości występowania słów; tworzona jest macierz w której terminy (np. wyrazy) odpowiadają wierszom, a dokumenty (np. strony internetowe) kolumnom; wadą tego podejścia jest fakt, że ilość słów może być bardzo duża co skutkuje utworzeniem ogromnej macierzy; ta technika maszynowego uczenia się znajduje swoje zastosowanie m.in. przy określaniu kategorii badanych tekstów.
- **reprezentacja grafowa** [13] - podejście rozszerzające reprezentację wektorową, w której słowa są węzłami połączonymi ze sobą krawędziami jeśli występują razem w tekście; istnieje możliwość zaobrazowania kolejności występowania słów z wykorzystaniem krawędzi skierowanych.

Jako istoty ludzkie posiadamy zdolność do rozpoznawania cudzych emocji po treści wypowiedzi i towarzyszących jej czynników pozawerbalnych, której wydźwięk określa się jako stosunek lub postawa pozostająca w korelacji do pewnego zdarzenia lub sytuacji. Treści publikowane w internecie w formie tekstu niosą ze sobą tylko reprezentację tekstową, dlatego odczytanie nastroju wyrażanego w taki sposób jest trudnym zadaniem. Tym bardziej złożonym procesem wymagającym dokładnego przeanalizowania każdego słowa jest zadanie odczytania wydźwięku wypowiedzi przez napisany system.

Dziedziną zajmującą się analizą wypowiedzi jest przetwarzanie języka naturalnego (ang. *NLP - natural language processing*), gdzie językiem naturalnym określa się język stosowany przez ludzi do komunikacji interpersonalnej. Zadaniem systemów rozumiejących język naturalnych jest przekształcenie go na formę bardziej przyjazną dla komputerów. Natomiast podzbiór

tych systemów służący do określania sentymentu zwraca ogólną informację o wydźwięku w ustalonej skali. Przykładowo analiza sentymentu w uważanej za punkt odniesienia dla takich algorytmów bibliotece NLP, napisanej przez pracowników i studentów amerykańskiego Uniwersytetu Stanforda, zwraca informację o sentymencie w pięciostopniowej skali: bardzo negatywny, negatywny, neutralny, pozytywny, bardzo pozytywny. Jest to ogólna informacja, ale nawet taka w postaci statystyk jest w stanie posłużyć jako cenne źródło informacji.

Algorytmy analizujące sentyment muszą poradzić sobie z następującymi wymaganiami [14, 15]:

- złożoność języka naturalnego,
- niejednoznaczność wypowiedzi np. wieloznaczność słowa zamek: *akcja Zemsty Fredry dzieje się na zamku w Odrzykoniu* lub *zamek w moich drzwiach nie chce się otworzyć* lub syntaktyczna niejednoznaczność gdy w jednej części zdania znajduje się pozytywny wydźwięk, a w kolejnej negatywny: *pogoda była okropna, ale obiad bardzo nam smakował*,
- specyfika stosowanego języka np. w internecie z powodu niejednokrotnie nie zachowanych zasad gramatycznych lub wiadomości nie niosących ze sobą żadnej treści,
- neologizmy np. retweet,
- idiomy np. "urwanie głowy"
- problemy z rozpoznawaniem nazw np. *byliśmy wczoraj na K2* - czy chodzi o film czy o szczyt górski,
- problem wyrwania wypowiedzi z kontekstu,
- tekst o charakterze sarkastycznym.

Techniki badania sentymentu wypowiedzi dzielimy na dwa rodzaje: korzystające ze słowników, które korzystają ze zbudowanych wcześniej list słów kluczowych z określonym sentymentem oraz oparte na klasyfikatorach, które wykorzystują techniki maszynowego uczenia się (ang. *machine learning*). Podstawowe techniki badania sentymentu to:

- **słownik** (ang. *lexicon based approach*) - klasyfikator ten bada tekst na podstawie wystąpień słów przy wykorzystaniu słownika składającego się ze słów z przypisanym pozytywnym lub negatywnym wydźwiękiem; sentyment pojedynczego słowa określa się jako iloraz prawdopodobieństwa wystąpienia z pozytywnym sentymentem do prawdopodobieństwa wystąpienia z negatywnym; wadą takiego podejścia jest brak analizy powiązania między słowami.
- **naiwny klasyfikator Bayesa** (ang. *naive Bayes classifier*) - probabilistyczny klasyfikator wykorzystujący techniki maszynowego uczenia się, opierający się o zasadę niezależności słów oraz wykorzystujący założenie, że teksty o charakterze pozytywnym charakteryzują się określonym słownictwem, a te o charakterze negatywnym charakteryzują się innym; podejście to zakłada także, że tekst, w którym występuje więcej słów o charakterze z kategorii pozytywnej lub negatywnej powinien zostać zaklasyfikowany do tej kategorii;
- **technika maksymalnej entropii** (ang. *maximum entropy technique*) - podejście szacujące rozkład prawdopodobieństwa opierające się na założeniu, że rozkład ma maksymalną entropię, jeśli dane nie są dobrze znane; entropia zwana także miarą niepewności rozkładu jest kolejną metodą wykorzystującą techniki maszynowego uczenia się.

- **metoda wektorów nośnych** (ang. *support vector machines*) - technika *machine learning* opierająca się o ideę hiperpłaszczyzny dzielącej teksty na pozytywne oraz negatywne z jak najmniejszym marginesem; celem jest znalezienie funkcji, dla której błąd sklasyfikowania tekstu będzie najmniejszy; programy wykorzystujące tą metodę dobrze radzą sobie z dużą ilością słów, ale oznaczenie części z nich jako nieistotne może powodować utratę części informacji.

Jak wynika z badań przeprowadzonych na grupie ponad 2000 dorosłych Amerykanów 81% internautów przynajmniej raz szukało opinii o produkcie w internecie, a 20% Amerykanów robi to na co dzień. Około 80% użytkowników internetu spośród wspomnianej grupy badawczej przyznało, że opinie przeczytane w internecie miały znaczny wpływ na dokonane przez nich zakupy. Wyniki tych badań pokazują jak bardzo wydzźwięk informacji, opinii, recenzji oraz poglądów w internecie ma wpływ na zachowania ludzi. Najczęściej spotykane zastosowania analizy sentymentu wypowiedzi zamieszczanych w internecie to:

- analizowanie opinii wyrażanych na temat produktów np. krótko po premierze nowego produktu firmy chcą wiedzieć jak są one odbierane tak aby móc w porę zaplanować wprowadzenie poprawek lub aby prognozować wyniki sprzedaży oraz cenę akcji, firmy próbują także docierać do osób krytykujących aby przekonać ich do zmiany zdania lub żeby krytyków konkurencji zachęcić do zakupu swoich produktów,
- badanie opinii na temat firm np. jako pracodawców lub ich odbiór na rynku,
- analiza opinii na temat ugrupowań politycznych i polityków np. dotycząca obecnej sytuacji lub przyszłych decyzji oraz nastawienia badanych grup społecznych,
- badanie nastrojów społecznych podczas wydarzeń sportowych może posłużyć np. do uzyskania wiedzy jak odbierane są decyzje właścicieli klubów piłkarskich.

Rozdział 4

Wymagania funkcjonalne i niefunkcjonalne

Jak już zostało wspomniane w rozdziale dotyczącym serwisu Twitter przy okazji omówienia dostępnych narzędzi, nie ma obecnie na rynku aplikacji, która umożliwiałaby śledzenie występowania dowolnego słowa w wiadomościach zamieszczanych w czasie rzeczywistym w tym serwisie, rysowałaaby wykres zależności pomiędzy użytkownikami, analizowałaaby nastroje społeczne użytkowników z wyświetleniem informacji o nastroju wyrażanym w poszczególnych wiadomościach oraz pozwalałaby na analizowanie danych historycznych i zamieszczanych w czasie rzeczywistym. Głównym celem tej pracy dyplomowej jest stworzenie aplikacji, która posiadałaby wspomniane funkcjonalności.

Główną funkcjonalnością, którą powinna zapewnić budowana aplikacja jest dostęp do usystematyzowanych danych pochodzących z serwisu Twitter, które będą nieść ze sobą informację o sentymencie. Dane te powinny być także przechowywane w taki sposób, aby móc zapewnić do nich dostęp w dowolnym momencie oraz spoza zaimplementowanego narzędzia.

Wymagania funkcjonalne, które dotyczą budowanego rozwiązania to:

- **przetwarzanie danych z serwisu Twitter** - aplikacja powinna przetwarzać tweety użytkowników serwisu Twitter, do których dostęp można uzyskać przez Streaming API tego serwisu, które zostało szczegółowo omówione w rozdziale 2.;
- **filtrowanie napływających danych po słowie kluczowym** - dane napływające w czasie rzeczywistym powinny być filtrowane pod względem zawartości w treści wiadomości słowa kluczowego, które zostało określone przez użytkownika za pomocą graficznego interfejsu aplikacji;
- **zapis napływających danych** - budowane narzędzie powinno zapisywać przetworzone, uporządkowane i dotyczące wybranego słowa kluczowego dane w lokalnej bazie danych, która umożliwiałaby dostęp do nich przez swój wbudowany pulpit w dowolnym momencie oraz spoza zaimplementowanego narzędzia;
- **duża częstotliwość pobierania danych** - aplikacja powinna pobierać informacje w krótkich odstępach czasu, ponieważ serwis Twitter charakteryzuje duża ilość informacji przesyłanych w każdej sekundzie;
- **prezentowanie danych historycznych** - narzędzie powinno prezentować dane historyczne zgromadzone podczas przetwarzania danych napływających wówczas w czasie

rzeczywistym;

- **prezentowanie podstawowych danych napływających w czasie rzeczywistym** - aplikacja powinna umożliwiać analizowanie podstawowych informacji o wiadomościach i użytkownikach, które będą napływać w czasie rzeczywistym;
- **dostęp do szczegółowej informacji o użytkowniku** - implementowane narzędzie powinno umożliwiać dostęp do informacji o każdym użytkowniku zainteresowanym wybranym słowem kluczowym, którego wiadomość udało się zarejestrować podczas gromadzenia danych z wykorzystaniem Streaming API serwisu Twitter;
- **dostęp do szczegółowej informacji o wiadomości** - aplikacja powinna wyświetlać szczegółową informację, o każdej wiadomości zawierającej wybrane słowo kluczowe i zapisanej podczas gromadzenia danych z wykorzystaniem Streaming API;
- **prezentowanie informacji statystycznej o sentymencie na dany temat** - narzędzie powinno prezentować statystyki sentymentu użytkowników serwisu Twitter, którzy w swoich wiadomościach zawarli wybrane słowo kluczowe i których wiadomości udało się zapisać podczas analizy danych napływających w czasie rzeczywistym;
- **możliwość jednoczesnej analizy danych historycznych i napływających w czasie rzeczywistym** - aplikacja powinna umożliwiać jednoczesną analizę danych historycznych i napływających w czasie rzeczywistym bez konieczności ponownego uruchamiania aplikacji lub zatrzymywania jednej z analiz;
- **przyjazny interfejs graficzny** - narzędzie powinno posiadać wygodny, łatwy do nauczenia oraz prosty interfejs graficzny.

Wymagania niefunkcjonalne, które powinna spełniać przygotowywana aplikacja to:

Spełnienie głównych założeń systemu czasu rzeczywistego

Głównym przypadkiem biznesowym, dla którego tworzona jest wspomniana aplikacja jest sytuacja dużego i globalnego zainteresowania pewnym tematem, które objawia się odnoszeniem się do niego w wiadomościach zamieszczanych przez użytkowników serwisu Twitter. Można stwierdzić, że przygotowywane narzędzie będzie przykładem systemem czasu rzeczywistego jeśli system ten będzie *"urządzeniem technicznym, którego wynik i efekt działania będzie zależny od chwili wypracowania tego wyniku"*. Wspólną cechą definicji takiego systemu jest *"zwrócenie uwagi na równoległość w czasie zmian w środowisku oraz obliczeń realizowanych na podstawie stanu środowiska"*.

Płynna obsługa danych

Budowany system powinien przetwarzać dane w czasie nie większym niż tempo napływania nowych informacji. Koniecznością jest zatem skorzystanie z narzędzi umożliwiających sprawne przetwarzanie danych, ale także ich zapis do bazy w czasie nie większym niż czas trwania określonego okna czasowego.

Możliwość działania aplikacji i analizowania danych historycznych w trybie offline

Dane zapisane podczas sesji korzystania ze strumienia danych serwisu Twitter będą zapisywane w lokalnej bazie danych, dlatego stworzone narzędzie powinno umożliwiać analizowanie danych historycznych bez połączenia z internetem.

Jednorazowe przetwarzanie informacji ze strumienia danych

Aplikacja powinna tylko jeden raz przetwarzać i zapisywać do bazy danych informacje pozyskane ze strumienia. Jeśli w bazie istnieje już informacja o użytkowniku to nowe wiadomości powinny być z nim powiązane.

Niezawodność

System powinien charakteryzować się niezawodnością podczas pracy z dużą ilością danych napływających w krótkich odstępach czasu oraz jak najbliższym prawdy określeniem sentymentu wypowiedzi zawartej w wielu wiadomościach.

System napisany na potrzeby tej pracy dyplomowej powinien spełniać wszystkie z wymienionych wymagań niefunkcjonalnych, ponieważ nie spełnienie nawet jednej z nich może spowodować, że aplikacja będzie nieużyteczna. Postawione wymagania funkcjonalne i niefunkcjonalne, łącząc się z opisem serwisu Twitter, narzędzi Big Data oraz analizą sentymentu wypowiedzi, definiują potrzeby jakie powinny umożliwiać narzędzia wybrane do jej budowy oraz samo narzędzie. Zostanie to omówione w następnych rozdziałach.

Rozdział 5

Wybór narzędzi

Przy pisaniu aplikacji tworzonej wraz z niniejszą pracą dyplomową będzie trzeba się zmierzyć z wymaganiami funkcjonalnymi i нефункциональными określonymi w poprzednim rozdziale. Aby tego dokonać niezbędne jest wybranie odpowiednich narzędzi, które pomogą w ich realizacji. W dalszej części rozdziału zaprezentowano uzasadnienie wyboru konkretnych składowych tworzonego systemu.

Językiem programowania wybranym do napisania tej części aplikacji jest Java. Jest to język, który charakteryzuje niezależność od systemu operacyjnego i procesora. Uniwersalny kod kompilowany jest do kodu bajtowego i następnie wykonywanego przez maszynę wirtualną Javy.

W wersji 8 tego języka, która jest obecnie najczęściej wybierana jako wersja do tworzenia dużych systemów komercyjnych, dodano wiele udogodnień pozwalających na częściowe pisanie kodu w myśl paradygmatu programowania funkcyjnego dzięki wprowadzonym operacjom na strumieniach i wyrażeniom lambda. Java swoją popularność zawdzięcza także wykorzystaniu do tworzenia aplikacji webowych, której przykładem jest aplikacja tworzona na potrzeby tej pracy. Język ten cały czas się rozwija. Przygotowywana jest obecnie 12. wersja Javy, która ma być dostępna w marcu 2019 roku. Przykład zastosowania operacji na strumieniach oraz wyrażenia lambda przedstawiono na rysunku 5.1..

```
private Set<Link> computeInterestedInLinks(Set<TwitterUser> interestedInUsers, String keywordName) {  
    return interestedInUsers  
        .stream()  
        .map(user -> new Link(keywordName, user.getScreenName(), RelationType.INTERESTED_IN))  
        .collect(Collectors.toSet());  
}
```

Rysunek 5.1: Przykład zastosowania operacji na strumieniach oraz wyrażenia lambda w kodzie napisanym w języku Java w wersji 8.

Jedną z bibliotek użytych w aplikacji jest biblioteka Lombok, która jest bardzo pomocna przy pisaniu prostych i powtarzalnych części kodu takich jak np. konstruktory, akcesory, mutatory, czy metody equals i hashCode. Przykładowo żeby napisać konstruktor posiadający jako argumenty wszystkie pola klasy wystarczy dodać na klasie adnotację `@AllArgsConstructor` co zaprezentowano na rysunku 5.2.. Lombok jest biblioteką typu open-source i cały czas się rozwija. Jest bardzo przydatna w podstawowych potrzebach kodu pisanego w Javie.

Językiem programowania wybranym do napisania części prezentacji przygotowywanej aplikacji jest JavaScript. Język ten znajduje swoje zastosowanie w tworzeniu stron internetowych.

```
@AllArgsConstructor
@Getter
public class Link {

    private final String keyword;
    private final String source;
    private final String target;
    private final RelationType type;
}
```

Rysunek 5.2: Przykład zastosowania adnotacji biblioteki Lombok.

Umożliwia reagowanie na zdarzenia wywołane w przeglądarce oraz wyświetlanie efektownych efektów wizualnych.

Dodatkowo w części front-end zostanie użyta biblioteka *React.js*, która jest określana jako *silnik do budowy interfejsów użytkownika*. *React.js* charakteryzuje się wykorzystaniem reaktywnego renderowania, które polega na przechowywaniu w pamięci reprezentacji struktury DOM i ponowne renderowanie interfejsu użytkownika tylko w tych miejscach, które zależą od zmienionego stanu. Jest to biblioteka bardzo wydajna i szybka co będzie bardzo pomocne przy budowie aplikacji.

Frameworkiem, na którym opierać się będzie tworzony system jest *Spring Boot*. *Spring Framework* jest określany jako główna platforma programistyczna, która sprawia, że programowanie w języku Java staje się łatwiejsze i bardziej wydajne. Spring Boot opiera się na zasadzie działania Spring Framework, ale rozwiązuje wiele problemów wynikających z potrzeby konfiguracji projektu przed rozpoczęciem implementowania nowych funkcji. Posiada zbiór wszystkich wymaganych bibliotek oraz konfiguracji zwanych *starterami* oraz wbudowany serwer aplikacyjny Tomcat. Spring jest bardzo dobrze udokumentowanym frameworkiem używanym przez programistów w wielu bardzo zróżnicowanych systemach.

Narzędziem wybranym do przetwarzania danych strumieniowych będzie *Apache Spark Streaming*, które charakteryzuje się łatwością integracji oraz operowaniem na kolekcjach informacji zamiast pojedynczych wiadomościach. Rozwiązanie to opiera się o *RDDs - Resilient Distributed Datasets* reprezentujące kolekcje rekordów, na których można wykonywać operacje znane z paradygmatu programowania funkcyjnego np. *map*, *filter*, *groupByKey* i *join*. Spark Streaming z definicji operuje na każdej porcji danych dokładnie raz. W bardzo rzadkim przypadku wystąpienia błędu może się zdarzyć, że dane nie zostaną przetworzone ani razu. Wówczas zostaną utracone dane z jednego okna czasowego trwającego kilka sekund co w kontekście budowanego rozwiązania wydaje się być dopuszczalne. Spark Streaming dobrze integruje się z językiem Java i frameworkiem Spring. Schemat zasady działania Apache Spark znajduje się na rysunku 5.3..

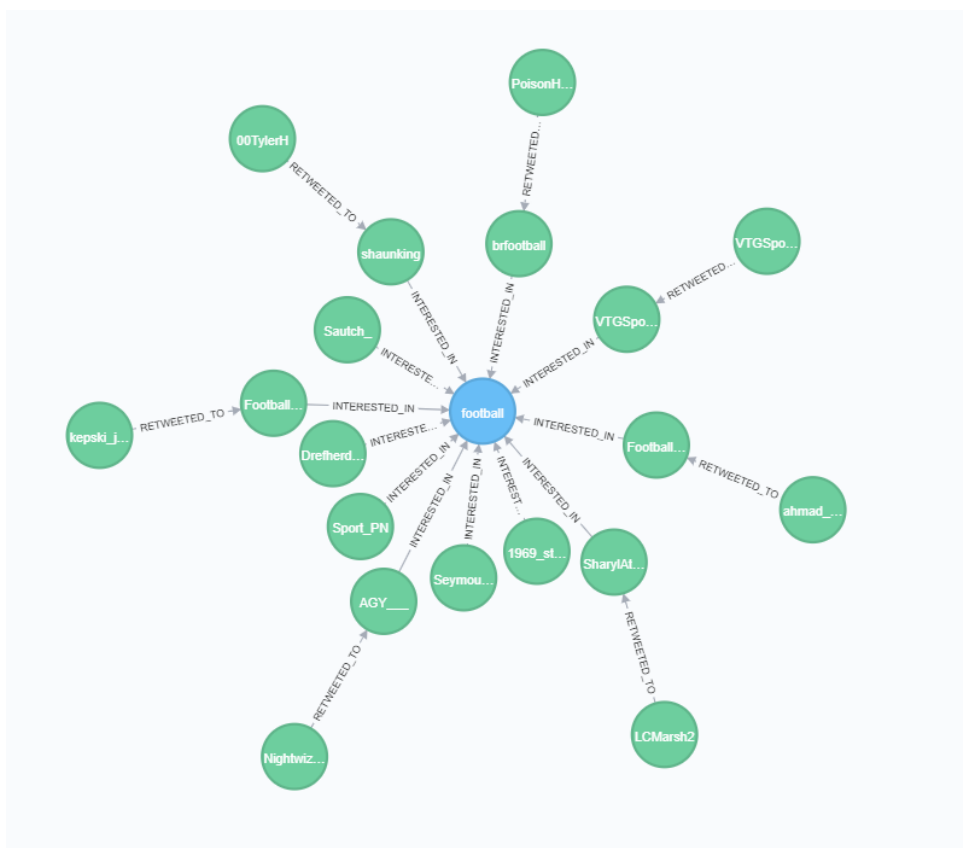
Systemem bazy danych wybranym do budowanego narzędzia jest *neo4j*. Jest to nierelacyjny system bazodanowy, które w odróżnieniu do relacyjnych charakteryzują się brakiem narzuconego schematu modelu informacji. Umożliwia to przechowywanie różnych danych bez potrzeby zapisu ich w postaci relacji, które w przypadku większej ilości danych są bardzo kosz-



Rysunek 5.3: Zasada działania narzędzia *Apache Spark Streaming*.

towne, ponieważ wymagają łączenia informacji z kilku tabel.

Dane w neo4j przechowywane są w postaci grafów, w których węzły odpowiadają wierszom z relacyjnych baz danych, a krawędzie pomiędzy węzłami są odpowiednikiem kluczy obcych. Dzięki takim rozwiązaniom wyszukiwanie zależności między węzłami polega na poruszaniu się po grafie, a nie na kosztownym łączeniu danych. Grafowe bazy danych powstały wraz z rozwojem sieci społecznościowych, dlatego jest naturalne, że tego typu rozwiązanie zostało wybrane do budowy przygotowywanej aplikacji. Neo4j dobrze integruje się z wybranymi już narzędziami. Przykład zastosowania grafowej bazy neo4j zaprezentowano na rysunku 5.4., a zapytania napisanego w języku Cypher przez nią wykorzystywanym na rysunku 5.5..



Rysunek 5.4: Przykład zastosowania grafowej bazy danych Neo4j.

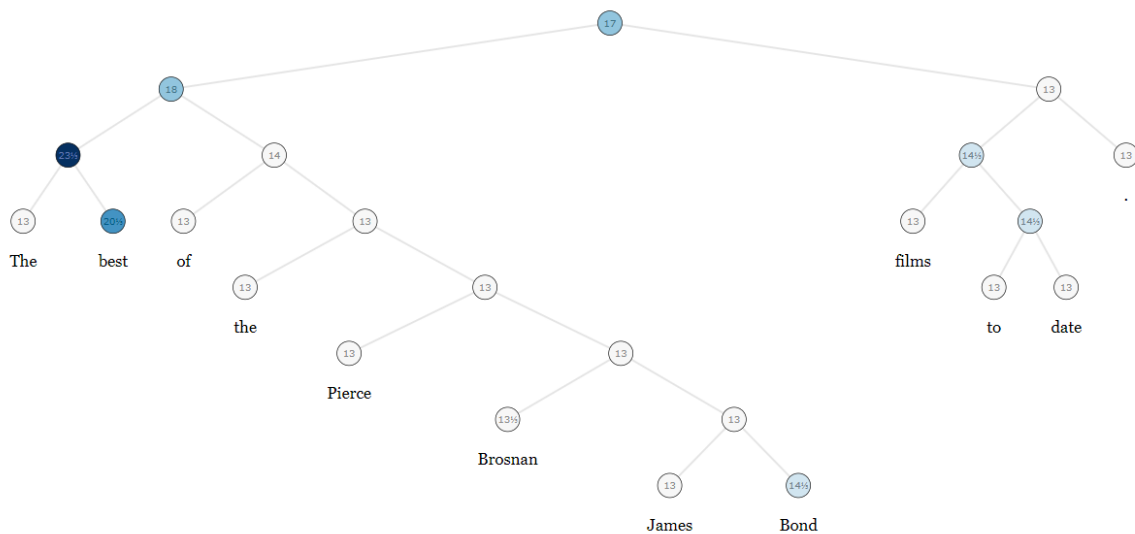
Biblioteką wybraną do analizy sentymentu jest rozwiązanie *Stanford NLP Sentiment Analysis* tworzone i cały czas rozwijane na Uniwersytecie Stanforda. Większość dostępnych rozwiązań określa sentyment tekstu biorąc pod uwagę wydźwięk pojedynczych słów i następnie podając wynik dla całego tekstu. W ten sposób umyka informacja o kolejności wyrazów oraz zostaje zaburzony sentyment. Stanford NLP analizuje każde zdanie tekstu w całości i oblicza

```
@Query("MATCH (keyword:Keyword)<-[relation:INTERESTED_IN]-(user:TwitterUser) " +
      "WHERE keyword.name = {keyword} " +
      "RETURN DISTINCT user")
Stream<TwitterUser> findAllInterestedIn(@Param("keyword") String keyword);
```

Rysunek 5.5: Przykład zapytania napisanego w języku Cypher wykorzystującego Spring Data.

sentymencie na podstawie tego w jaki sposób słowa składają się na znaczenie w dłuższe zwroty.

Narzędzie to korzysta z sieci neuronowej *Recursive Neural Network*, która opiera się na zasadach gramatycznych. Sieć ta została wytrenowana z wykorzystaniem zbioru recenzji filmowych zawierających prawie 12 tysięcy zdań i około 215 tysięcy słów.



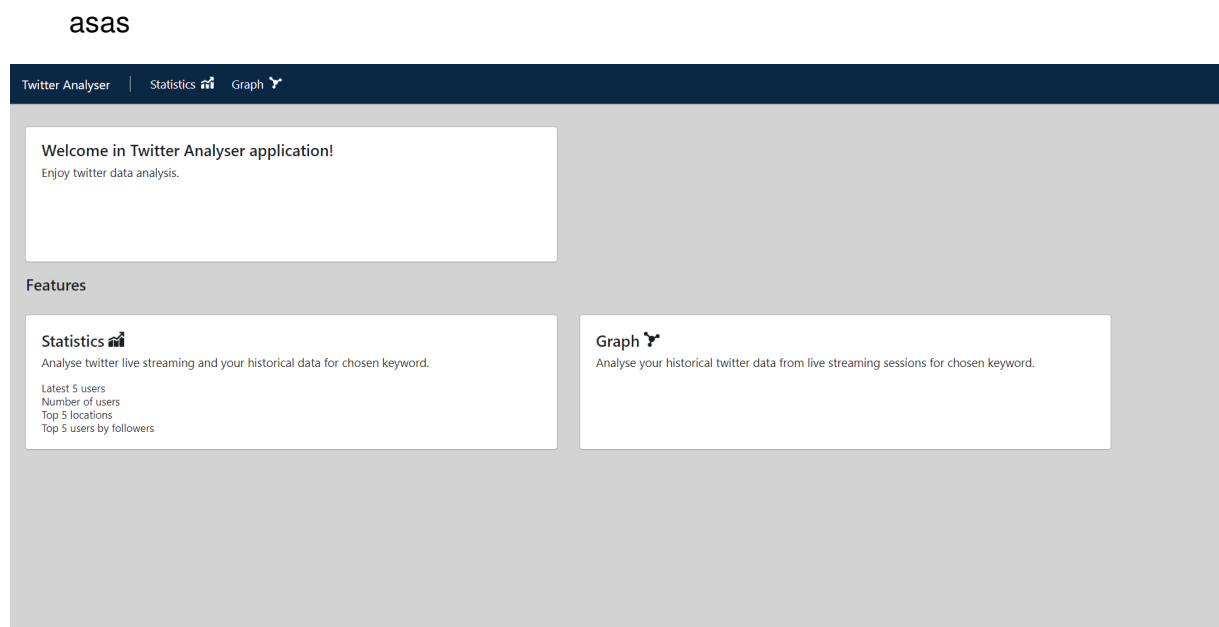
Rysunek 5.6: Przykład rozkładu opinii na temat filmu z serii *James Bond*.

Jednym z pierwszych kroków w analizie wydźwięku tekstu za pomocą tej biblioteki jest podział tekstu na zdania. Stanford NLP w swojej pełnej wersji udostępnia także parser tekstów, który może zostać użyty przy tej okazji. Następnie każde zdanie jest analizowane oddzielnie, a sentyment całego tekstu jest obliczany jako średnia arytmetyczna jego składowych. Wydźwięk podawany jest w pięciostopniowej skali: bardzo negatywny, negatywny, neutralny, pozytywny, bardzo pozytywny. Co ciekawe twórcy tego narzędzia umożliwiają na swojej stronie internetowej zwrócenie im uwagi, że pewne zdanie ze zbioru treningowego źle określa sentyment. Wówczas zostanie to przeanalizowane i poprawione. Przykład rozkładu zdania za pomocą omawianej biblioteki zaprezentowano na rysunku 5.6..

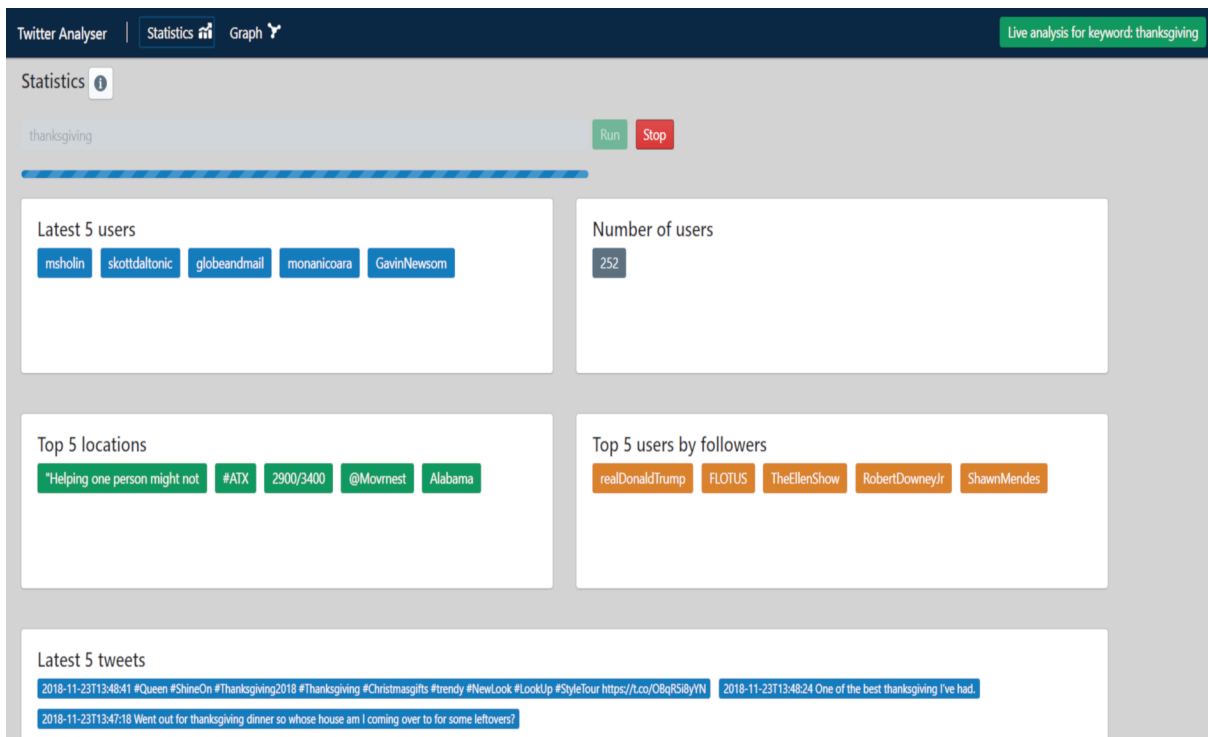
Stanford NLP wspiera obecnie badanie sentymentu tylko w tekstach angielskojęzycznych. Jego twórcy podają, że skuteczność stosowanego przez nich algorytmu wynosi 85.4% oraz że jest to obecnie jedyne rozwiązanie, które wspiera negacje na różnych poziomach budowanych drzew w zwrotach pozytywnych i negatywnych. Narzędzie to dobrze integruje się z wybranymi już wcześniej składowymi budowanym systemu.

Rozdział 6

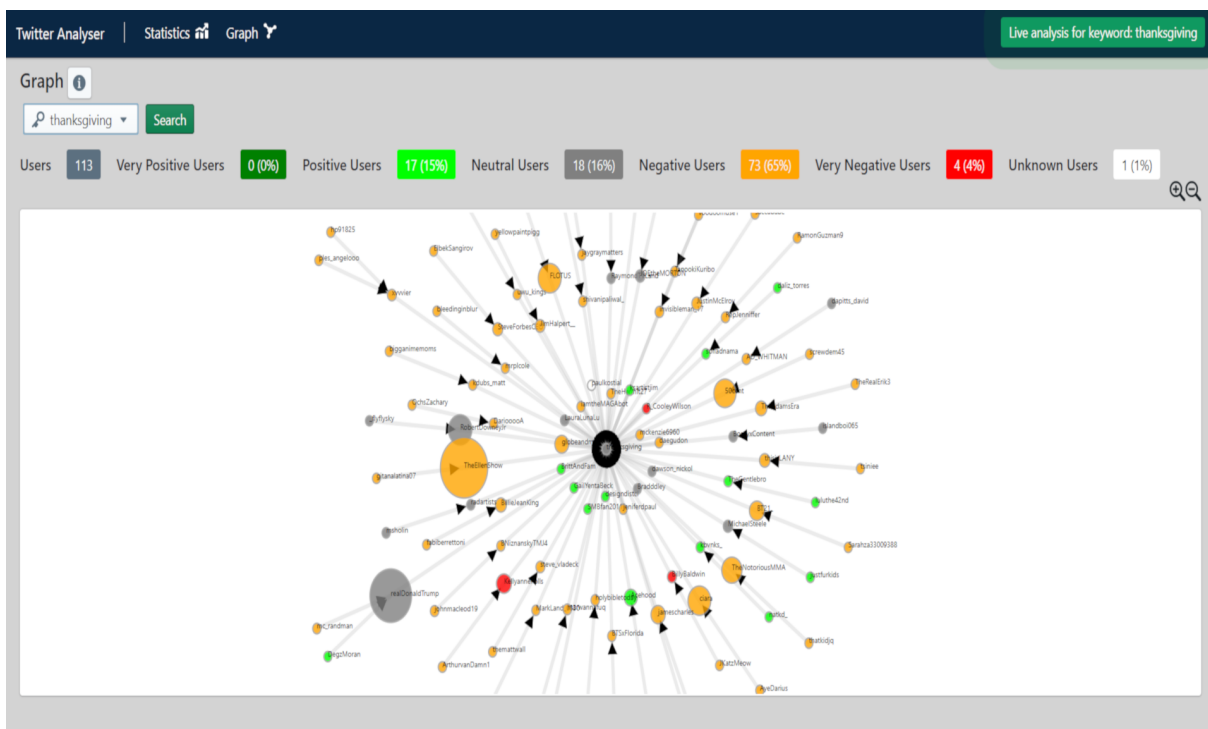
Aplikacja Twitter Analyser



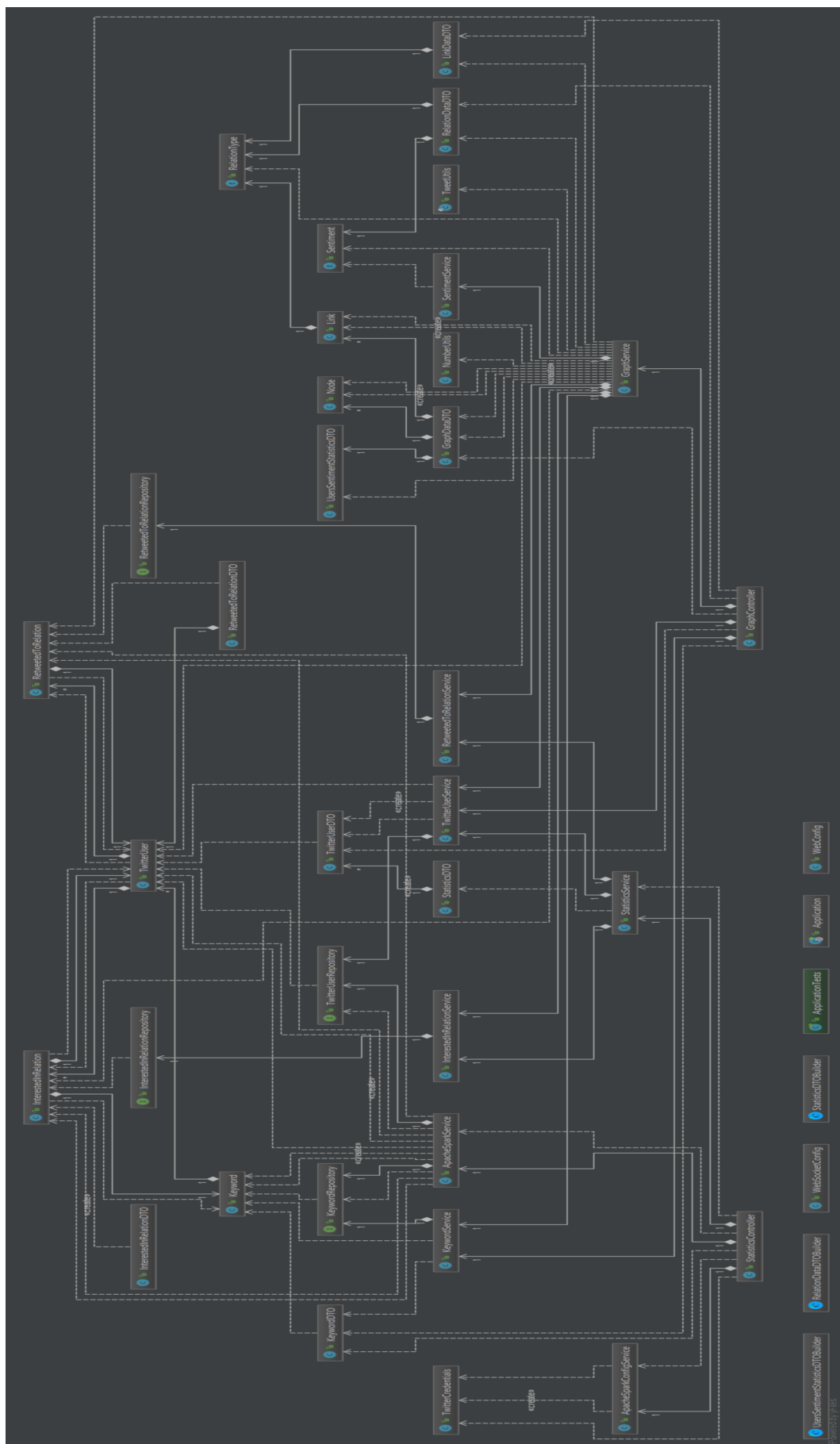
Rysunek 6.1: Przykład zastosowania operacji na strumieniach oraz wyrażenia lambda w kodzie napisanym w języku Java w wersji 8.



Rysunek 6.2: Przykład zastosowania operacji na strumieniach oraz wyrażenia lambda w kodzie napisanym w języku Java w wersji 8.



Rysunek 6.3: Przykład zastosowania operacji na strumieniach oraz wyrażenia lambda w kodzie napisanym w języku Java w wersji 8.



Rysunek 6.4: Przykład zastosowania operacji na strumieniach oraz wyrażenia lambda w kodzie napisanym w języku Java w wersji 8.

Rozdział 7

Badania i wnioski

Rozdział 8

Podsumowanie

Bibliografia

- [1] <https://en.wikipedia.org/wiki/Twitter> [Dostęp 4 listopada 2018]
- [2] <https://medium.com/@ssola/playing-with-twitter-streaming-api-b1f8912e50b0> [Dostęp 11 listopada 2018]
- [3] <https://developer.twitter.com/en/docs/tweets/search/api-reference/get-search-tweets.html> [Dostęp 11 listopada 2018]
- [4] <http://140dev.com/twitter-api-programming-tutorials/aggregating-tweets-search-api-vs-streaming-api> [Dostęp 11 listopada 2018]
- [5] <https://www.talkwalker.com/blog/5-free-twitter-analytics-tools-with-views-from-experts> [Dostęp 19 listopada 2018]
- [6] <http://twitter4j.org/javadoc/twitter4j/User.html> [Dostęp 12 listopada 2018]
- [7] <http://twitter4j.org/javadoc/twitter4j/Status.html> [Dostęp 12 listopada 2018]
- [8] https://pl.wikipedia.org/wiki/Bieg_maratoński [Dostęp 24 listopada 2018]
- [9] A. Mykowiecka, *"Inżynieria lingwistyczna: komputerowe przetwarzanie tekstów w języku naturalnym"*, Wydawnictwo Polsko-Japońskiej Wyższej Szkoły Technik Komputerowych, 2007, s. 127 -170.
- [10] https://en.wikipedia.org/wiki/Bag-of-words_model [Dostęp 2 grudnia 2018]
- [11] https://en.wikipedia.org/wiki/Vector_space_model [Dostęp 11 grudnia 2018]
- [12] V. Dixit, A. Saroliya, *"A semantic Vector Space Model approach for sentiment analysis"*, International Journal of Advanced Research in Computer and Communication Engineering Vol. 2, Issue 8, August 2013.
- [13] F. Zhou, F. Zhang, G. Yang, *"Graph-based text representation model and its realization"*, Natural Language Processing and Knowledge Engineering, 2010.
- [14] <https://blog.brand24.pl/co-to-jest-analiza-sentymentu-oraz-jak-mozesz-ja-wykorzystac/> [Dostęp 24 listopada 2018]
- [15] https://pl.wikipedia.org/wiki/Przetwarzanie_języka_naturalnego [Dostęp 24 listopada 2018]