

PK4: Platforma NLP

Sprawozdanie

Michał Czyż

10 czerwca 2024

1 Wstęp

Projektem przewidzianym do wykonania w tym semestrze była platforma przetwarzania języka naturalnego (*eng. natural language processing*). W ramach projektu zaimplementowane zostało tłumaczenie tekstu, wykrywanie języka tekstu, anonimizacja danych wrażliwych, oraz badanie sentymentu.

Przetwarzanie języka naturalnego to dziedzina nauki łącząca językoznawstwo oraz dziedzinę sztucznej inteligencji. Zajmuje się ona analizą i przetwarzaniem języka ludzkiego przez komputer. Po stworzeniu i wyszkoleniu modelu językowego, program jest w stanie z dużą dokładnością przewidywać oraz wykryć logiczny ciąg zdań języka. Na tej podstawie może być on dalej analizowany i przetwarzany.

Tłumaczenie maszynowe może być realizowane na wiele różnych sposobów. Skupiono się na dwóch rozwiązaniach jakimi są statystyczne tłumaczenie maszynowe oraz duże modele językowe opartych na transformerach. Do tłumaczenia statystycznego wykorzystano pierwszy model opracowany przez IBM w latach 90. Skupia się on w głównej mierze na tłumaczeniu leksykalnym, zatem nie będzie bardzo dokładny przy wychwytywaniu dokładnego kontekstu, ale jest w stanie osiągnąć dostateczne rezultaty. Dodatkowym atutem tego modelu jest jego wysoka wydajność. Model oparty jest o prawdopodobieństwo warunkowe.

2 Analiza tematu

Platforma NLP składa się z trzech usług:

- Core - główny program napisany w języku C++. Jest to główny program skupiający się na komunikacji z warstwą API oraz na przetwarzaniu danych.
- Middleware - serwer API, napisany w Pythonie. Łączący komunikację pomiędzy użytkownikiem, a programem głównym. Middleware udostępnia warstwę programistyczną, którą wykorzystuje frontend.
- Frontend - strona internetowa napisana przy pomocy Next.JS oraz Reacta. Korzysta z API i udostępnia warstwę graficzną użytkownikowi końcowemu.

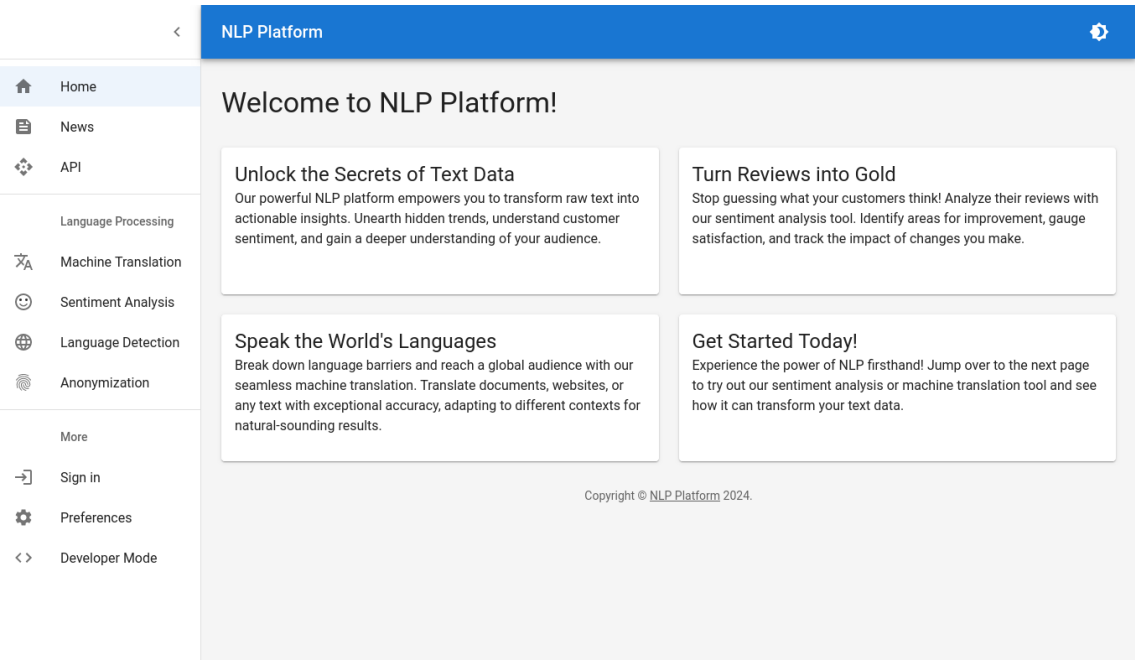
Język C++ wybrano przede wszystkim ze względu na wydajność oraz dostępność wielu elementów programowania obiektowego. Język Python został wybrany, ponieważ ma bardzo rozbudowane biblioteki do nauczania maszynowego oraz posiada bardzo łatwe i wygodną bibliotekę - *Flask*.

Jednym z głównych założeń projektu jest bycie modularnym i łatwym w rozbudowie. W szczególności, iż projekt jest złożony i wymaga wykorzystania wielu języków i bibliotek. Główna struktura programu w C++ wygląda następująco:

- Communication - zajmuje się komunikacją pomiędzy pythonem i programem.
- Engine - zbiór wielu klas obsługujących działanie programu i przetwarzających dane.
- Translate - zbiór klas odpowiedzialnych za obsługę tłumaczenia.
- Sentiment - zbiór klas odpowiedzialnych za obsługę analizy sentymentu.
- AsyncLogger - klasa odpowiedzialna za prawidłowe logowanie stanu programu.

Wykorzystanie biblioteki podczas tworzenia głównego programu w C++ to zeromq, cppmq oraz boost. ZeroMQ jest wykorzystywany do komunikacji pomiędzy programem a interfejsem sie-

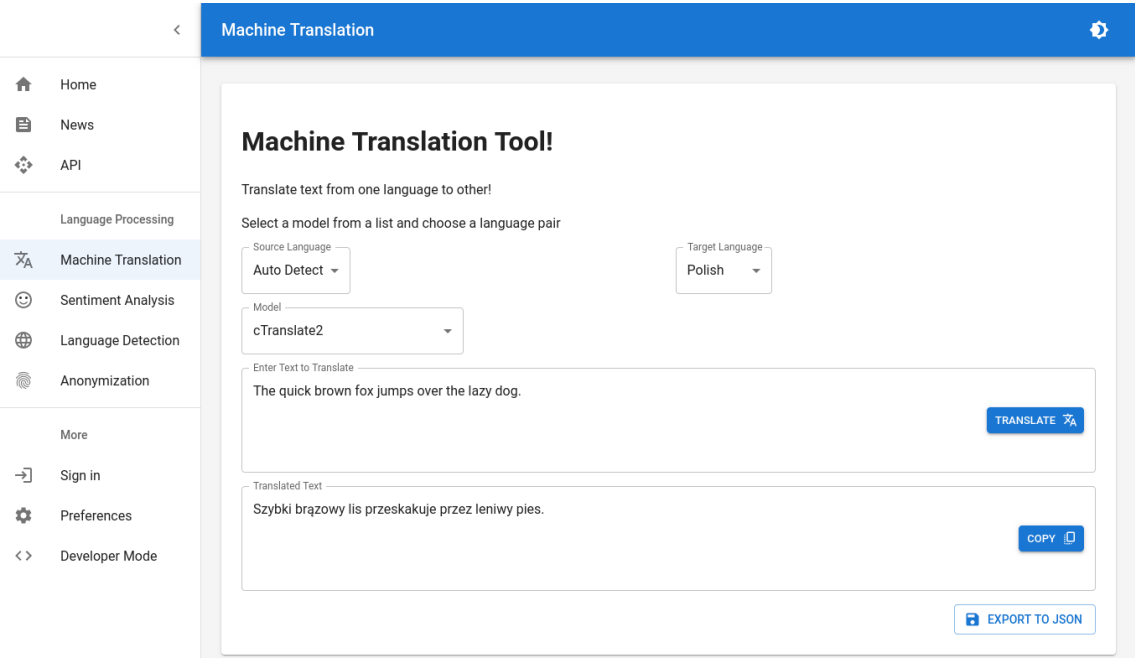
ciowym. Biblioteka cppzmq udostępnia wygodny sposób do inicjalizacji i zarządzania komunikacją. Biblioteka Boost była używana przy łączeniu C++ z Pythonem oraz do komunikacji.



Rysunek 1: Wygląd strony głównej

3 Specyfikacja zewnętrzna

Program posiada warstwę API oraz graficzną co pozwala na interakcję z platformą w wygodny sposób dla użytkownika końcowego jak i także dla programisty, który chciałby wykorzystać funkcjonalność platformy NLP w swoim programie. Serwer API posiada podstawową autoryzację poprzez token API.



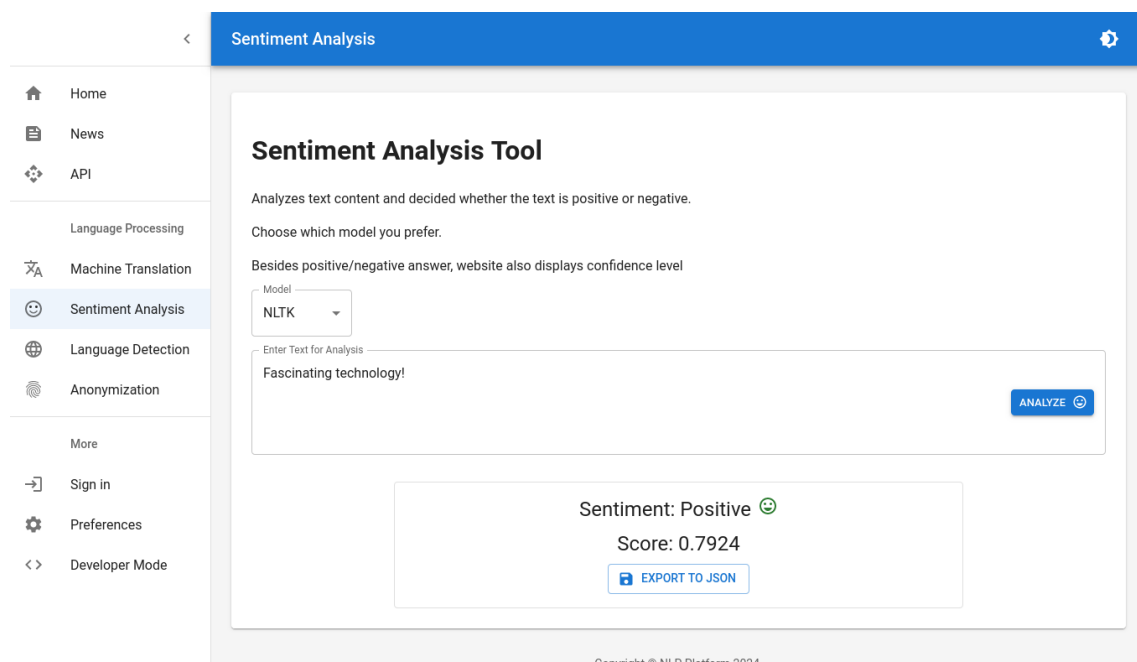
Rysunek 2: Interfejs tłumaczenia maszynowego

3.1 Instrukcja dla użytkownika

Jeżeli usługa jest uruchomiona na komputerze lub port interfejsu webowego jest udostępniony do internetu wystarczy, że użytkownik ma dostęp do internetu oraz posiada zainstalowaną przeglądarkę internetową. Po odwiedzeniu odpowiedniego adresu ukaże się strona internetowa po której może nawigować i wybierać interesujące go funkcje.

3.2 Przykłady działania

Przykłady działania aplikacji zostały przedstawione na rysunkach widocznych w dokumencie. Jest możliwe swobodne nawigowanie po interfejsie webowym, wybranie modelu i pary językowej, która nas interesuje oraz odczytanie i pobranie wyniku w postaci pliku JSON.



Rysunek 3: Interfejs analizy sentymentu

Interfejs Programowania Aplikacji wygląda następująco. Należy wysłać zapytanie na następującą ścieżkę na serwerze:

GET `http://localhost:8080/api/ping`

POST `http://localhost:8080/api/translate`

POST `http://localhost:8080/api/sentiment`

POST `http://localhost:8080/api/detect`

POST `http://localhost:8080/api/anonymize`

POST `http://localhost:8080/api/auth`

POST `http://localhost:8080/api/signin`

Struktura dla przykładowego tłumaczenia przy wykorzystaniu cTranslate2 wygląda następująco:

```
{
  "srcLanguage": "en",
  "dstLanguage": "pl",
  "mode": "ct2",
  "data": "Natural Language Processing"
}
```

4 Specyfikacja wewnętrzna

Diagram klas został przedstawiony na Rysunku 4. Główne klasy to Translate, Communication oraz Sentiment, które odpowiadają za główną funkcjonalność programu. Reszta klas to m.in. Async-

- [AnkiWeb Polish-English Vocabulary](#)
- [European Parliament Languages Dataset](#)

5 Testowanie

Program był testowany w kilku etapach. Na początku testowano interfejs webowy aby sprawdzić czy wszystkie możliwości i kombinacje ustawień modeli są prawidłowe. Sprawdzana była także komunikacja warstwy webowej z API oraz czy wszystkie requesty są poprawnie obsługiwane lub odrzucane.

Później sprawdzano poprawność działania warstwy API z warstwą C++. Komunikacja ta była kluczowa do poprawnego działania programu. Sprawdzono także czy wszystkie przekazywane dane są przekazywane prawidłowo, czy nie można poprzez wysłanie błędnych danych dokonać *remote code execution*.

Ostatecznie sprawdzano warstwę C++, pod kątem bezpieczeństwa wykonywania modeli. A także testowano czy wszystkie klasy zachowują się zgodnie z założeniem.

6 Uruchamianie

Zalecana kolejność uruchamiania usług to warstwa webowa, następnie API i na końcu program w C++.

6.1 Uruchomienie serwera WWW

Wymagane pakiety: `nodejs`, `npm`

Aby uruchomić usługę należy wykonać następujące kroki:

```
$ cd web
$ npm install
$ npm run build
$ npm run start
```

Aplikacja uruchomi się na porcie `3000`.

W przypadku chęci debuggowania aplikacji lub wprowadzania zmian w kodzie można uruchomić serwer deweloperski:

```
$ npm run dev
```

6.2 Uruchomienie serwera API

Wymagane pakiety: `python`, `python-pip`

```
$ cd api
```

Przed uruchomieniem serwera należy aktywować wirtualne środowisko.

```
$ source activate.sh
```

Potem można wywołać serwer API.

```
$ python main.py
```

Aby opuścić wirtualne środowisko, można wykonać polecenie:

```
$ deactivate
```

6.3 Uruchomienie Core Platform

Wymagane pakiety: `base-devel`, `boost`, `cmake`, `python`, `python-pip`, `mold`, `zeromq`, `cppzmq`

Na początku należy aktywować wirtualne środowisko, aby móc wykorzystywać modele.

```
$ cd nlp/models
```

Umożliwia to skrypt `activate.sh`.

```
$ source activate.sh  
$ cd ..
```

Ostatecznie wystarczy skompilować i uruchomić program wykonując:

```
$ make
```

Jeżeli nie chcemy ponownie kompilować programu, można wykorzystać:

```
$ make run
```

Aby opuścić wirtualne środowisko, można wykonać polecenie:

```
$ deactivate
```

6.4 Parametry programu

W C++ można także przekazać parametr `-verbose`, który włącza dodatkowe logowanie do konsoli informacji o działaniu ZeroMQ i dokładnego przebiegu pakietów.

7 Uwagi i wnioski

Podsumowując, platforma NLP okazała się fascynującym i bardzo rozwojowym projektem. Poznałem jak działają modele językowe, jak tworzyć i trenować własne. Nauczyłem się także jak łączyć warstwę sieciową z programami niższego poziomu, jak C++. Mogłem także dobrze rozwinąć swoje umiejętności programowania obiektowego oraz poznać i lepiej zgłębić zagadnienia poznane w ramach laboratorium.