

SWAT - PyNER Intercompatibility

Konrad Delong Antek Piechnik

6 lutego 2009

Spis treści

1	Wstęp	2
2	Struktura projektu	2
2.1	Tagger tekstów do testów	2
2.2	Generator reguł n-gramowych	2
2.3	Rdzeń	2
3	Zastosowane technologie	3
3.1	Tagger	3
3.2	Rdzeń	3
4	Zastosowane narzędzia	3
4.1	Jython	3
4.2	Wykorzystane biblioteki	3
5	Opis powiązania aplikacji z Interfejsem IQuantumDetector z projektem SWAT	4
5.1	Header w pliku Java	4
5.2	Kod aplikacji w Javie	4
5.3	Przykładowe importy w aplikacji pythonowskiej	4
5.4	Przykładowe użycie klas Javy w pythonie (za pomocą jythona	4

1 Wstep

W ponizszym dokumencie opisane zostaly narzedzia (wraz z przykladowym kodem) z których korzystaliśmy podczas tworzenia projektu oraz przy tworzeniu pełnej komunikacji kodu w Pythonie z klasami oraz metodami Javy stosowanymi w projekcie SWAT.

2 Struktura projektu

Projekt sklada sie z nastepujacych czesci:

2.1 Tagger tekstow do testow

Aplikacja majaca na celu skanowanie przykladowych tekstow a nastepnie generowania specyficznych raportow odnosnie ilosci wystapien poszczegolnych encji w tekscie. Raporty generowane będą w formie sprzyjającej reszcie systemu efektywne ich użytkowanie. Wykorzystywane one będą następnie do badania sprawności wykorzystywanego algorytmu ewolucyjnego oraz trafności pojawiających się nowych reguł. Tagger bazować będzie na liście danych encji, z której nie będziemy korzystać w pracy Rdzenia aplikacji.

2.2 Generator reguł n-gramowych

W związku z chęcią implementacji reguł bazujących na najpopularniejszych n-gramach utworzony zostanie generator reguł który zbierze najczęściej występujące n 3,4,5 gramy występujące w bazie nazwisk przygotowanej wcześniej (do celów tagowania tekstow testowych).

Generator, w celu uniknięcia zbierania nic nie znaczących n-gramów porównywał będzie nie tylko jego zawartość ale również pozycję w słowie.

2.3 Rdzen

W tej części aplikacji znajdować się będzie zaimplementowany algorytm który przekazywał będzie zestawy reguł, otrzymywał raporty dotyczace ich sprawnosci a następnie tworzył kolejne zestawy reguł na podstawie otrzymanych wyników.

W celu zarówno uproszczenia implementacji jak i ulatwienia dalszej rozbudowy systemu chcemy odizolowac rdzeń aplikacji od całej reszty i nie powierzać mu zadań takich jak skanowanie tekstu. Pozwoli to zarazem na większe pole manewru przy

wyborze języka implementacji.

Rdzeń aplikacji uruchamiany będzie w celu przygotowania optymalnego zestawu reguł, a więc za każdym razem gdy dodane zostaną nowe reguły tudzież gdy dodane zostaną do systemu nowe teksty testowe, mające zwiększyć odporność systemu na coraz bardziej różnorodne dane wejściowe.

Skaner będzie starał się dopasowywać kombinacje słów w tekście do przedstawionych reguł a następnie starał się zapisać wyniki dla poszczególnych reguł (tak, aby algorytm ewolucyjny mógł sprawnie dobierać nowe reguły na podstawie wyników indywidualnych kombinacji zastosowanych wcześniej).

3 Zastosowane technologie

Do implementacji znacznej części aplikacji pragniemy wykorzystać język Python, ze względu na jego perfekcyjne przystosowanie do prac nad przetwarzaniem języka naturalnego. Oczywiście za pomocą pewnych narzędzi będziemy się starali powiązać go z obecną architekturą systemu, napisanego w dużej części w języku Java.

3.1 Tagger

Wykonano wstępny prototyp taggera w Pythonie na potrzeby testowe.

3.2 Rdzeń

Rozpoczęto prace w Pythonie. Dodatkowo będziemy jak najszybciej starali się skończyć z Jythonem.

4 Zastosowane narzędzia

4.1 Jython

Bridge pozwalający na transparentną komunikację między Pythonem a klasami Javy. Będzie niezmiernie przydatny do komunikacji Rdzenia aplikacji napisanego w Pythonie z obecną architekturą systemu.

4.2 Wykorzystane biblioteki

- Egothor

- Morfologik - Stemming

5 Opis powiazania aplikacji z Interfejsem IQuantumDetector z projektem SWAT

5.1 Header w pliku Java

```
import jyinterface.factory.JythonFactory;
```

5.2 Kod aplikacji w Javie

```
String shortName =  
    "org.ppbw.agh.swat.hoover.smith.quantum.detection.IQuantumDetector";  
Object obj = JythonFactory.getJythonObject(shortName,  
                                            "pyner/ngrams_detector.py",  
                                            "NgramsDetector");
```

5.3 Przykładowe importy w aplikacji pythonowskiej

```
from org.ppbw.agh.swat.hoover.smith.quantum import QuantumType  
...  
from org.ppbw.agh.swat.hoover.smith.stemmer import StemmerPL
```

5.4 Przykładowe użycie klas Javy w pythonie (za pomoca jythona

```
if ngram in prev_word:  
    dq[word_id] =  
        DetectedQuantum(leafSegment.getWordToken(word_id), QuantumType.SURNAME)
```