

Laboratorium 1

Część I.

Rozpoczynamy od instalacji Pythona, PyCharm i biblioteki NLTK. Ewentualnie, można korzystać z Anaconda.

Część II.

Korzystamy z biblioteki NLTK. Jako pierwszy należy wykonać skrypt służący do załadowania biblioteki oraz danych z NLTK Book:

```
import nltk
nltk.download()
from nltk.book import *
```

Część III.

Poniżej przeanalizujemy przykłady ilustrujące działanie wybranych użytecznych funkcji:

- funkcja *len()* zwraca liczbę wystąpień tokenów, czyli ciągów znaków (słowa, znaki specjalne np. emotikony); złożenie *len(set())* zwraca liczbę wystąpień bez powtórzeń,
- funkcja *count()* zwraca liczbę wystąpień ustalonego słowa,
- metoda *similar()* z parametrem, będącym słowem wykonana na tekście zwraca słowa podobne,
- metoda *concordance()* z parametrem, będącym słowem zwraca występowania słowa wraz z kontekstem,
- metoda *common_contexts()* dla listy słów zwraca konteksty występowania słów z listy

Przykłady do przetestowania:

```
print(len(text3))
print(len(set(text3)))
text1.count("monstrous")
text1.similar("monstrous")
text1.concordance("monstrous")
text1.common_contexts(["stories", "pictures"])
```

Część IV.

Zadania.

1. Dla każdego z tekstów z NLTK Book należy wybrać dwa słowa i napisać kod, który zwraca wspólne konteksty występowania tych słów.

- Korzystając z funkcji `lexical_diversity` ($\text{len}(\text{text}) / \text{len}(\text{set}(\text{text}))$) - liczba liter w tekście podzielona przez liczbę różnych liter) uzupełnić poniższą tabelkę dla danych z `nlk.book`

| tekst | liczba słów | słowa różne | lexical_diversity |
|-------|-------------|-------------|-------------------|
| text1 | 260819 | 19317 | 13.50 |
| text2 | | | |
| text3 | | | |
| text4 | | | |
| text5 | | | |

- Stwórz listę wszystkich słów 4-literowych z `text1`. Ile ich jest?
- W `text1` znajdź wszystkie wystąpienia słów długości większej niż 17.
- Korzystając z funkcji `set` i `sorted` wyznacz słownik dla każdego ze zdań `sent1`, ..., `sent8` oraz wspólny słownik dla wszystkich wymienionych zdań.
- Podaj definicję funkcji `VocabSize()`, która dla tekstu zwraca rozmiar słownika - czyli wylicza ile jest słów różnych. Zastosuj do każdego z tekstów z `nlk.book`.
- Wyznacz 10 najczęściej występujących słów w `text1`.
- Sprawdź jakie są najdłuższe słowa w każdym z tekstów `text1`, ..., `text6`.

Rozwiązania część IV zad 1-8 proszę przesłać do Portalu Edukacyjnego w pliku z rozszerzeniem `.py` i nazwą:
nazwiskoLab1.py