

## Laboratorium 2. Klasyfikacja danych tekstowych

Mając dane oznaczone etykietami, można wykonać prosty eksperyment klasyfikacji. Jednakże, surowy tekst napisany w języku naturalnym nie nadaje się do trenowania modeli. Aby było to możliwe konieczna jest wektoryzacja danych. Rozwiązanie tego laboratorium pozwoli poznać przykładowy sposób transformacji surowych danych tekstowych z etykietami do prostego modelu zdolnego do klasyfikacji.

### Zadanie 2.1. Przygotowanie danych

Wczytaj dane z pliku `imdb.csv` i dla treści każdego dokumentu:

1. Zamień w tekście wszystkie duże litery na małe
2. Dokonaj tokenizacji z użyciem `TreebankWordTokenizer`<sup>1</sup>
3. Usuń z tokenów stopwords<sup>2</sup> (możesz też usunąć inne tokeny, które nie mają znaczenia np. "<br />")
4. Wykonaj stemming na tokenach z poprzedniego punktu z użyciem `PorterStemmer`<sup>3</sup>
5. Wykonaj lematyzację na tokenach z punktu 3. z użyciem `WordNetLemmatizer`<sup>4</sup>
6. Zapisz do tablic cztery zestawy tokenów:
  - Tokeny oryginalne
  - Tokeny bez stop words
  - Tokeny po stemmingu
  - Tokeny po lematyzacji

### Zadanie 2.2. Klasyfikacja

Wczytaj zestawy tokenów z pierwszego zadania i dla każdego zestawu:

1. Podziel dane na treningowe (70%) oraz testowe (30%) z użyciem `train_test_split`<sup>5</sup>
2. Dokonaj wektoryzacji za pomocą metody `CountVectorizer`<sup>6</sup>
3. Wytrenuj model z użyciem klasyfikatora `MultinomialNB`<sup>7</sup> na danych treningowych (funkcja `fit()`)
4. Dokonaj predykcji modelu na danych testowych i zapisz wynik (funkcja `predict()`)
5. Wyświetl dokładność wyrażoną za pomocą metryki `accuracy_score`<sup>8</sup>
6. Porównaj uzyskane wyniki i zapisz w komentarzach spostrzeżenia

### Zadanie 2.3. Eksperyment

Wczytaj zestawy tokenów z pierwszego zadania i dla każdego zestawu:

1. Dokonaj wektoryzacji z użyciem `TfidfVectorizer`<sup>9</sup>
2. Dane podziel na 5 foldów z użyciem funkcji `StratifiedKFold`<sup>10</sup>
3. Utwórz odpowiednie tablice na wyniki
4. Dla każdego foldu:
  - Wytrenuj model z użyciem klasyfikatora `MLPClassifier`<sup>11</sup> na danych **treningowych** (funkcja `fit()`)
  - Dokonaj predykcji modelu na danych **testowych** (funkcja `predict()`)
  - Zapisz w tablicy z punktu 3 dokładność wyrażoną za pomocą metryki `accuracy_score`
5. Wyznacz średnią dokładność oraz odchylenie standardowe i wyświetl wyniki

<sup>1</sup><https://www.nltk.org/api/nltk.tokenize.treebank.html>

<sup>2</sup><https://www.nltk.org/book/ch02.html>

<sup>3</sup><https://www.nltk.org/api/nltk.stem.porter.html>

<sup>4</sup><https://www.nltk.org/api/nltk.stem.wordnet.html>

<sup>5</sup>[https://scikit-learn.org/stable/modules/generated/sklearn.model\\_selection.train\\_test\\_split.html](https://scikit-learn.org/stable/modules/generated/sklearn.model_selection.train_test_split.html)

<sup>6</sup>[https://scikit-learn.org/stable/modules/generated/sklearn.feature\\_extraction.text.CountVectorizer.html](https://scikit-learn.org/stable/modules/generated/sklearn.feature_extraction.text.CountVectorizer.html)

<sup>7</sup>[https://scikit-learn.org/stable/modules/generated/sklearn.naive\\_bayes.MultinomialNB.html](https://scikit-learn.org/stable/modules/generated/sklearn.naive_bayes.MultinomialNB.html)

<sup>8</sup>[https://scikit-learn.org/stable/modules/generated/sklearn.metrics.accuracy\\_score.html](https://scikit-learn.org/stable/modules/generated/sklearn.metrics.accuracy_score.html)

<sup>9</sup>[https://scikit-learn.org/stable/modules/generated/sklearn.feature\\_extraction.text.TfidfVectorizer.html](https://scikit-learn.org/stable/modules/generated/sklearn.feature_extraction.text.TfidfVectorizer.html)

<sup>10</sup>[https://scikit-learn.org/stable/modules/generated/sklearn.model\\_selection.StratifiedKFold.html](https://scikit-learn.org/stable/modules/generated/sklearn.model_selection.StratifiedKFold.html)

<sup>11</sup>[https://scikit-learn.org/stable/modules/generated/sklearn.neural\\_network.MLPClassifier.html](https://scikit-learn.org/stable/modules/generated/sklearn.neural_network.MLPClassifier.html)