

Data oddania: \_\_\_\_\_

Ocena: \_\_\_\_\_

Radosław Grela 216769

Jakub Wachała 216914

## Zadanie 1: ekstrakcja cech, miary podobieństwa, klasyfikacja

### 1. Cel

Celem naszego zadania było stworzenie aplikacji do klasyfikacji tekstów za pomocą metody  $k$ -NN ( $k$  najbliższych sąsiadów) oraz różnych metryk i miar podobieństwa, a następnie porównać kategorie z tymi wygenerowanymi przez aplikację.

### 2. Wprowadzenie

Głównym zagadnieniem projektowym, z którym mieliśmy do czynienia w ramach zadania 1 była klasyfikacja statystyczna tekstów na podstawie wektora wyekstrahowanych cech. Do przeprowadzenia eksperymentu zaimplementowaliśmy algorytm *k-najbliższych sąsiadów*.

Algorytm  $k$ -najbliższych sąsiadów (*k-NN - k-nearest neighbors*) to jeden z algorytmów zaliczanych do grupy algorytmów leniwych. Jest to taka grupa algorytmów, która szuka rozwiązania dopiero, gdy pojawia się wzorec testujący. Przechowuje wzorce uczące, a dopiero później wyznacza się odległość wzorca testowego względem wzorców treningowych. [9]

Algorytm ten działa w taki sposób, że dla każdego wzorca testowego obliczana jest odległość za pomocą wybranej metryki względem wzorców treningowych, a następnie wybierana jest  $k$  najbliższych wzorców treningowych. Wynik wyznaczony jest jako najczęstszy element wśród nich. W naszym zadaniu odległość ta jest równa skali podobieństwa tekstów, a im ta odległość jest mniejsza, tym lepiej.

## 2.1. Ekstrakcja cech

Do ekstrakcji cech charakterystycznych tekstu utworzyliśmy wektor cech, który opisuje tekst za pomocą 11 cech. Liczba słów zawsze jest liczona po zastosowaniu stop-listy oraz stemizacji, bez znaków przestankowych.

- $C_1$  - Stosunek słów kluczowych do wszystkich słów w pierwszych 10% tekstu. Obliczona jest za pomocą wzoru:

$$S_{k10} = S_{10} \cap S_k \quad (1)$$

$$C_1 = \frac{|S_{k10}|}{|S_{10}|} \quad (2)$$

gdzie

$S_k$  - zbiór wszystkich słów kluczowych,

$S_{10}$  - zbiór słów w pierwszych 10% tekstu,

$S_{k10}$  - zbiór słów kluczowych w pierwszych 10% tekstu,

$|S_{k10}|$  - liczba słów kluczowych w pierwszych 10% tekstu,

$|S_{10}|$  - liczba wszystkich słów w dokumencie w pierwszych 10% tekstu.

- $C_2$  - Stosunek słów kluczowych do wszystkich słów w ostatnich 10% tekstu. Obliczona jest za pomocą wzoru:

$$S_{k90} = S_{90} \cap S_k \quad (3)$$

$$C_2 = \frac{|S_{k90}|}{|S_{90}|} \quad (4)$$

gdzie

$S_{90}$  - zbiór słów w ostatnich 10% tekstu,

$S_{k90}$  - zbiór słów kluczowych w ostatnich 10% tekstu,

$|S_{k90}|$  - liczba słów kluczowych w ostatnich 10% tekstu,

$|S_{90}|$  - liczba wszystkich słów w dokumencie w ostatnich 10% tekstu.

- $C_3$  - Stosunek słów kluczowych do wszystkich słów w dokumencie. Obliczona jest za pomocą wzoru:

$$S_{k3} = S \cap S_k \quad (5)$$

$$C_3 = \frac{|S_{k3}|}{|S|} \quad (6)$$

gdzie

$S_{k3}$  - zbiór wszystkich słów kluczowych znajdujących się w dokumencie,

$|S_{k3}|$  - liczba wszystkich słów kluczowych znajdujących się w dokumencie,

$S$  - zbiór wszystkich słów w dokumencie,

$|S|$  - liczba wszystkich słów w dokumencie.

- $C_4$  - Stosunek słów kluczowych, których ilość liter  $\in (0,4]$  do wszystkich słów w dokumencie. Obliczona jest za pomocą wzoru:

$$S_{k_4} = S \cap S_{k_{04}} \quad (7)$$

$$C_4 = \frac{|S_{k_4}|}{|S|} \quad (8)$$

gdzie

$S_{k_{04}}$  - zbiór słów kluczowych, których ilość liter  $\in (0,4]$ ,

$S_{k_4}$  - zbiór słów kluczowych znajdujących się w dokumencie, których ilość liter  $\in (0,4]$ ,

$|S_{k_4}|$  - liczba słów kluczowych znajdujących się w dokumencie, których ilość liter  $\in (0,4]$ .

- $C_5$  - Stosunek słów kluczowych, których ilość liter jest  $\geq 8$  do wszystkich słów w dokumencie. Obliczona jest za pomocą wzoru:

$$S_{k_5} = S \cap S_{k_{8+}} \quad (9)$$

$$C_5 = \frac{|S_{k_5}|}{|S|} \quad (10)$$

gdzie

$S_{k_{8+}}$  - zbiór słów kluczowych, których ilość liter  $\geq 8$ ,

$S_{k_5}$  - zbiór słów kluczowych znajdujących się w dokumencie, których ilość liter  $\geq 8$ ,

$|S_{k_5}|$  - liczba słów kluczowych znajdujących się w dokumencie, których ilość liter  $\geq 8$ .

- $C_6$  - Stosunek linii do ilości akapitów. Obliczona jest za pomocą wzoru:

$$C_6 = \frac{|l|}{|a|} \quad (11)$$

gdzie

$|l|$  - liczba linii,

$|a|$  - liczba akapitów.

Na rysunku 1 został zamieszczony przykładowy artykuł. Dla tego przykładu,  $|l| = 24$ , a  $|a| = 8$ .

- $C_7$  - Stosunek słów, których ilość liter jest większa niż 6 do wszystkich słów w dokumencie. Obliczona jest za pomocą wzoru:

$$C_7 = \frac{|S_{k_{6+}}|}{|S|} \quad (12)$$

gdzie

$S_{k_{6+}}$  - zbiór słów w dokumencie, których ilość liter jest większa niż 6,

$|S_{k_{6+}}|$  - liczba słów w dokumencie, których ilość liter jest większa niż 6.

Inco Ltd said it did not expect its earlier reported removal from the Dow Jones industrial index to make a major impact on the company's stock.

"We don't think that individuals or institutions buy our shares because we were one of the Dow Jones industrials," spokesman Ken Cherney said in reply to a query.

Inco closed 1-3/8 lower at 19-3/8 in second most active trading on the Toronto Stock Exchange.

The Wall Street Journal, which selects the index, said Inco was dropped to make the index more representative of the market. Inco, the non-Communist world's largest nickel producer, was a member of the index since 1928.

Replacing Inco and Owens-Illinois Inc will be Coca-Cola Co and Boeing Co, effective tomorrow.

Nickel analyst Ilmar Martens at Walwyn Stodgell Cochran Murray Ltd said Inco's removal from the index would likely spark short-term selling pressure on the stock.

"Some investors who have Inco may suddenly say, 'well, because it's not now a Dow stock, we should eliminate that investment,'" said Martens, although he added the move was unlikely to have a serious long-term impact on Inco stock.

Inco has struggled in recent years against sharply lower nickel prices. Its net earnings fell to 200,000 U.S. dlrs in 1986 from 52.2 mln dlrs the previous year.

Rysunek 1. Przykładowy artykuł.

- $C_8$  - Stosunek słów, których ilość liter jest  $\leq 6$  do wszystkich słów w dokumencie. Obliczona jest za pomocą wzoru:

$$C_8 = \frac{|S_{k_{06}}|}{|S|} \quad (13)$$

gdzie

$S_{k_{06}}$  - zbiór słów w dokumencie, których ilość liter jest  $\leq 6$ ,

$|S_{k_{06}}|$  - liczba słów w dokumencie, których ilość liter jest  $\leq 6$ .

- $C_9$  - Ilość słów unikalnych. Jest to liczba słów, które wystąpiły w tekście co najmniej raz. Przykładowo, dla zdania „*Być albo nie być*” ilość słów unikalnych jest równa 3 (*być, albo, nie*).
- $C_{10}$  - Ilość słów w dokumencie, których ilość liter  $\in [5,8]$ . Pseudokod obliczający wartość cechy  $C_{10}$ :
  - $C_{10}=0$
  - Dla każdego słowa w artykule:
    - Jeżeli długość słowa  $\geq 5$  i długość słowa  $\leq 8$ :
      - $C_{10}++$ ;
  - Zwróć  $C_{10}$
- $C_{11}$  - Najczęściej występujące słowo kluczowe. Przykładowo, dla tekstu na rysunku 1 i zbioru słów kluczowych  $\{Inco, pressure, against, year\}$  najczęściej występującym słowem kluczowych jest słowo *Inco*. Jest to cecha tekstowa, której podobieństwo z innym słowem mierzymy jedną z dwóch miar podobieństwa ciągów znaków opisanych w sekcji *Metryki i miary podobieństwa*.

## 2.2. Wyznaczanie słów kluczowych

Wyznaczenie słów kluczowych przebiega w następujący sposób: na początek za pomocą klasy WordCounter zliczane są wszystkie słowa w artykułach oraz jednocześnie dodawane do odpowiednich list w tej klasie. Każda zmienna

jest listą stringów o nazwie wordCountDictionary + nazwa kraju. Dodatkowo, przechowywany jest słownik typu  $\langle string, int \rangle$ , którego kluczem jest słowo, a wartością to ilość wystąpień tego słowa we wszystkich artykułach. Po podliczeniu wszystkich słów oraz przydzieleniu do odpowiednich list wybieramy po 18 najpopularniejszych słów dla każdego kraju, które występują tylko w tym jednym konkretnym kraju. Na koniec  $18 * 6 = 108$  słów zostaje słowami kluczowymi. Cały proces wyznaczania słów kluczowych jest dokonywany po zastosowaniu stop-listy oraz po stemizacji. Ponadto, proces wybierania słów kluczowych pomija 20% wszystkich podliczonych słów, aby proces dopasowywania słów kluczowych do krajów nie trwał zbyt długo.

### 2.3. Metryki i miary podobieństwa

Do liczenia odległości pomiędzy artykułami oraz obliczenia miary podobieństwa używaliśmy 3 metryk i 2 miar podobieństwa ciągów tekstowych.

1. Metryka Euklidesowa - aby obliczyć odległość  $d_e(x, y)$  między wektorami  $x$  i  $y$  należy obliczyć pierwiastek kwadratowy z sumy kwadratów różnic wartości współrzędnych wektora o tych samych indeksach. Wzór jest następujący [6]:

$$d_e(x, y) = \sqrt{(x_1 - y_1)^2 + \dots + (x_n - y_n)^2} \quad (14)$$

gdzie  $x_i$  i  $y_i$  to cechy wektora.

2. Metryka uliczna - odległość  $d_m(x, y)$  jest równa sumie wartości bezwzględnych z różnic wartości współrzędnych wektora o tych samych indeksach [4]:

$$d_m(x, y) = \sum_{n=1}^N |x_n - y_n| \quad (15)$$

gdzie  $x_i$  i  $y_i$  to cechy wektora.

3. Metryka Czebyszewa - odległość  $d_c(x, y)$  w tej metryce jest równa maksymalnej wartości bezwzględnych różnic współrzędnych punktów  $x$  oraz  $y$ , zgodnie ze wzorem [5]:

$$d_c(x, y) = \max_i |x_i - y_i| \quad (16)$$

gdzie  $x_i$  i  $y_i$  to cechy wektora.

4. Miara  $n$ -gramów - metoda ta określa podobieństwo łańcuchów tekstowych  $s_1, s_2$  w oparciu o ilość wspólnych podciągów  $n$ -elementowych, czyli  $n$ -gramów [3]:

$$sim_n(s_1, s_2) = \frac{1}{N - n + 1} \sum_{i=1}^{N-n+1} h(i) \quad (17)$$

gdzie

$h(i) = 1$ , jeśli  $n$ -elementowy podciąg zaczynający się od  $i$ -tej pozycji w  $s_1$  występuje co najmniej raz w  $s_2$ , w przeciwnym razie  $h(i) = 0$

$N - n + 1$  - ilość możliwych  $n$ -elementowych podciągów w  $s_1$ .

W naszym programie  $n$  jest stałe i wynosi 3.

5. Uogólniona miara *n-gramów* (Miara Niewiadomskiego) - ta miara jest ulepszoną wersją miary *n-gramów*. Bada ona podobieństwo poprzez sprawdzenie podciągów różnej długości od jedno- do *N*-elementowych, gdzie *N* jest długością słowa [3]:

$$\mu_N(s_1, s_2) = \frac{2}{N^2 + N} \sum_{i=1}^{N(s_1)} \sum_{j=1}^{N(s_1)-i+1} h(i, j) \quad (18)$$

gdzie

$h(i, j) = 1$ , jeśli *i*-elementowy podciąg w słowie  $s_1$  zaczynający się od *j*-tej pozycji w słowie  $s_1$  pojawia się co najmniej raz w słowie  $s_2$ , w przeciwnym razie  $h(i, j) = 0$ ,

$N(s_1), N(s_2)$  – ilość liter w słowach  $s_1$  i  $s_2$ ,

$N = \max\{N(s_1), N(s_2)\}$ ,

$\frac{N^2+N}{2}$  - ilość możliwych podciągów od 1-elementowych do *N*-elementowych w słowie o długości *N*.

Aby porównać wektory za pomocą metryk w algorytmie *k*-NN, najpierw wyznaczamy miarę podobieństwa z ostatniej, 11 cechy, która jest cechą tekstową. Wyznaczamy ją za pomocą jednej z dwóch miar. Ponieważ w tych miarach im bliżej 1, tym lepiej, odejmujemy tą liczbę od 1, a następnie używamy jej w metryce.

#### 2.4. Miary jakości

W wynikach klasyfikacji używamy następujących miar jakości [10]:

- Accuracy:

$$ACC = \frac{TP + TN}{TP + TN + FP + FN} \quad (19)$$

- Precision

$$PPV = \frac{TP}{TP + FP} \quad (20)$$

- Recall

$$TPR = \frac{TP}{TP + FN} \quad (21)$$

Oznaczenia użytych symboli:

TP - miara prawdziwie pozytywna (*true positive*), czyli poprawnie zaklasyfikowane artykuły dla badanej etykiety

TN - miara prawdziwie negatywna (*true negative*), czyli poprawnie zaklasyfikowane artykuły dla pozostałych etykiet

FP - miara fałszywie pozytywna (*false positive*), czyli niepoprawnie zaklasyfikowane artykuły dla badanej etykiety

FN - miara fałszywie negatywna (*false negative*), czyli niepoprawnie zaklasyfikowane artykuły dla pozostałych etykiet

#### 2.5. Normalizacja wektorów cech

Wszystkie wartości wektora cech zostaną znormalizowane za pomocą wzoru (22) [2]:

$$x_{scaled} = \frac{x - x_{min}}{x_{max} - x_{min}} \quad (22)$$

Dzięki temu, wektor będzie zawierał tylko wartości z przedziału  $[0, 1]$ .

### 3. Opis implementacji

Nasza aplikacja została utworzona w języku C# i jest to aplikacja konsolowa. Poniżej opisane zostały wszystkie klasy oraz dane zawarte w naszym projekcie:

- Klasa Program to klasa główna naszego programu. Jest swego rodzaju kontrolerem dla pozostałych klas. Znajduje się tutaj funkcja *main*, która rozpoczyna wykonywanie programu.
- W katalogu *dane* znajdują się wszystkie pliki z artykułami, które są wykorzystywane do badań.
- Klasa Metric jest klasą abstrakcyjną. Odpowiada za obliczenia odległości tekstów. Po tej klasie dziedziczą klasy: EuclideanMetric, ChebyshevMetric oraz ManhattanMetric.
- Klasa Measure jest klasą abstrakcyjną. Po niej dziedziczą klasy *GeneralizedNGramsMeasure* i *NGramsMeasure*, które odpowiadają za obliczanie miar podobieństwa łańcuchów tekstowych.
- Klasa Feature jest klasą abstrakcyjną. Po niej dziedziczy 10 klas: Feature 1-10, które reprezentują każdą z 10 wyekstrahowanych przez nas cech.
- Klasa Stemmer to klasa, która odpowiada za stemizację tekstów. Została ona zapożyczona z [7]
- Klasa StopwordTool jest klasą odpowiedzialną za usuwanie słów znajdujących się na stopliście. Również została znaleziona i zapożyczona z Internetu ze strony [8]
- WordCounter jest używany do zliczania słów wszystkich artykułów i podania ich liczności. Potrzebny głównie do wyznaczenia słów kluczowych.
- Klasa KeyWords odpowiada za wyznaczenie 100 słów kluczowych. Metoda wyznaczania słów kluczowych została opisana w sekcji 2.
- Klasa FileReader odpowiada za otwieranie każdego pliku z artykułami
- FileParser to klasa odpowiedzialna za parsowanie danych z konkretnego pliku.
- Article to klasa reprezentująca artykuł. Zawiera takie cechy jak: tekst oryginalny, tekst przetworzony, *place*, *classifiedPlace*, wektor cech.
- Klasa Neighbor to klasa, która przechowuje artykuł oraz obliczoną wartość algorytmu k-NN dla konkretnego, obecnie sprawdzanego artykułu w algorytmie. Wykorzystujemy ją, aby znaleźć najbliższych k sąsiadów.
- KNN to klasa odpowiedzialna za algorytm k najbliższych sąsiadów.

Na rysunku 2 przedstawiony został wynik z konsoli po przykładowym uruchomieniu programu, natomiast na rysunku 3 przedstawiony został diagram UML naszego programu.

```

Settings: k=2 , training=30%, test=70%, metric=EuclideanMetric

Place Precision Recall
usa 80,892 81,183
canada 8,921 11,524
france 0,621 0,595
japan 7,925 5,949
west-germany 9,459 5,578
uk 12,559 12,841
Accuracy: 66,415

```

Rysunek 2. Wynik z przykładowego uruchomienia programu.

## 4. Materiały i metody

Wykonana przez nas klasyfikacja została wykonana za pomocą wszystkich trzech metryk oraz dwóch miar podobieństwa. Każdy przypadek testowy był klasyfikowany dla dziesięciu różnych wartości  $k$  najbliższych sąsiadów: 2, 3, 4, 5, 7, 10, 13, 15, 20, 25.

Klasyfikacji dokonywaliśmy tylko na tych tekstach, które miały jedną z etykiet: *west-germany*, *usa*, *france*, *uk*, *canada*, *japan* i były to ich jedyne etykiety.

Dokonaliśmy pięciu różnych podziałów na dane testowe oraz treningowe:

- 30% dane treningowe, 70% dane testowe
- 50% dane treningowe, 50% dane testowe
- 70% dane treningowe, 30% dane testowe
- 80% dane treningowe, 20% dane testowe
- 85% dane treningowe, 15% dane testowe

Poniżej zostały opisane 4 wykonane przez nas eksperymenty.

### 4.1. Badanie zależności Accuracy od parametru $k$

W tym eksperymencie badaliśmy wpływ doboru parametru  $k$  na Accuracy. Program został uruchomiony dla 10 różnych wartości  $k \in \{2, 3, 4, 5, 7, 10, 13, 15, 20, 25\}$ .

Klasyfikacja tekstów została wykonana dla stałej wartości podziału zbioru cech na testowe i treningowe. Był to podział 50% dane treningowe, 50% dane testowe.

Metryką, jakiej użyliśmy była metryka euklidesowa.

### 4.2. Badanie wyników klasyfikacji w zależności od wartości proporcji podziału zbioru

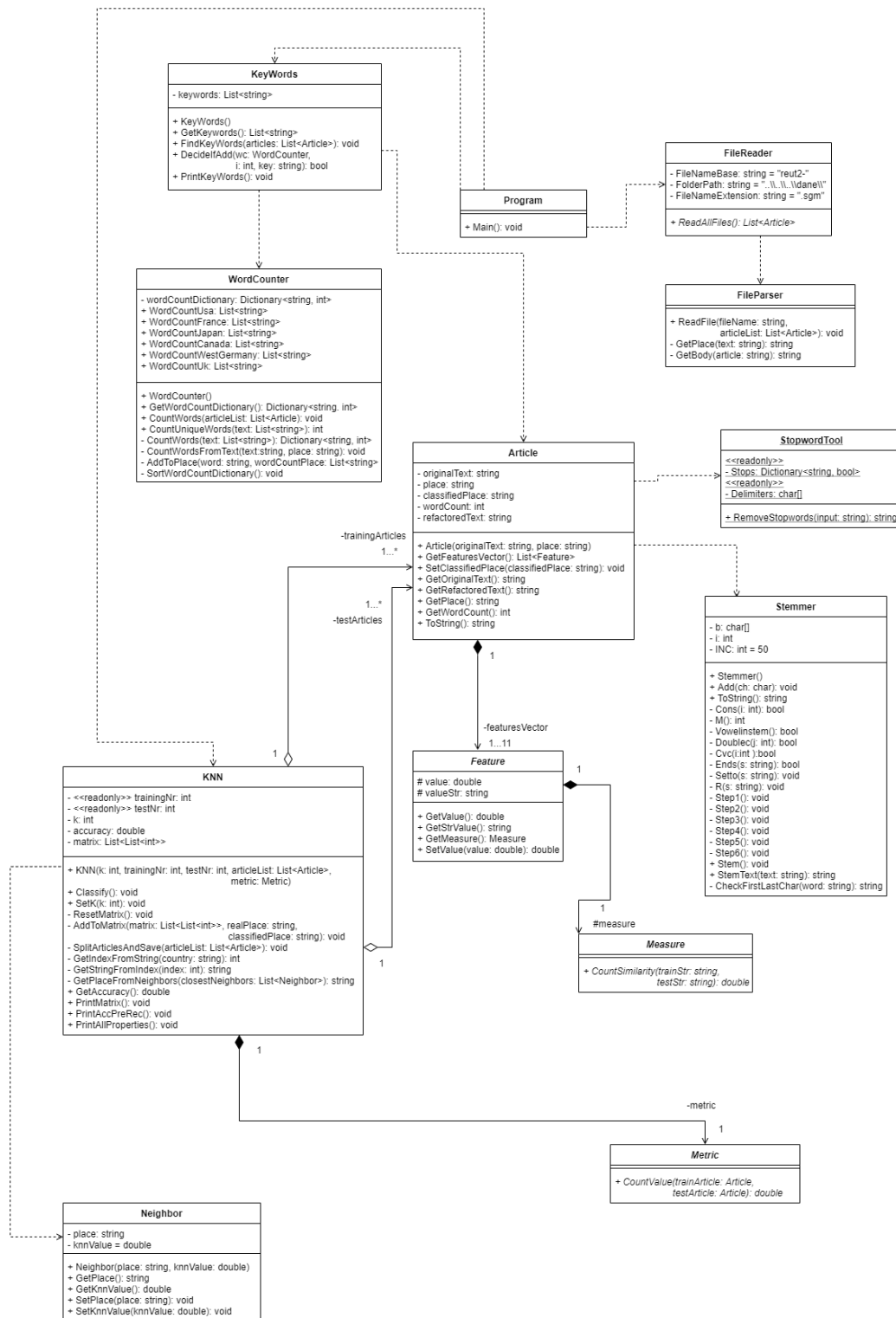
W tym eksperymencie badaliśmy wpływ wartości proporcji podziału zbioru na Accuracy. Program został uruchomiony dla  $k=10$ .

Badane podziały były następujące:

- 30% dane treningowe, 70% dane testowe
- 50% dane treningowe, 50% dane testowe
- 70% dane treningowe, 30% dane testowe
- 80% dane treningowe, 20% dane testowe
- 85% dane treningowe, 15% dane testowe

Metryką, jakiej użyliśmy była metryka uliczna.





Rysunek 3. Diagram UML.

#### 4.3. Badanie zależności Accuracy od wyboru metryki

W tym eksperymencie badaliśmy zależność Accuracy od wyboru metryki. Program został uruchomiony dla  $k=13$ . Podział na dane treningowe i testowe był stały i wynosił 70% treningowe i 30% testowe.

#### 4.4. Badanie zależności Accuracy od wyboru podzbioru cech

W tym eksperymencie badaliśmy zależność Accuracy od wyboru podzbioru cech. Program został uruchomiony dla  $k=20$ . Metryka, jakiej użyliśmy to metryka Czebyszewa. Podział na dane treningowe i testowe był stały i wynosił 50% treningowe i 50% testowe. Podzbiory cech jakie badaliśmy były następujące:

- Wszystkie cechy
- $C_1, C_2, C_3, C_4, C_5, C_{11}$
- $C_6, C_7, C_8, C_9, C_{10}$
- $C_1, C_2, C_3, C_8, C_9, C_{10}$
- $C_4, C_5, C_6, C_7, C_{11}$

### 5. Wyniki

#### 5.1. Wartości wektorów cech przed normalizacją

- Cecha  $C_1$  zawierała się w wartościach  $\in [0, 1]$
- Cecha  $C_2$  zawierała się w wartościach  $\in [0, 0.5]$
- Cecha  $C_3$  zawierała się w wartościach  $\in [0, 0.155]$
- Cecha  $C_4$  zawierała się w wartościach  $\in [0, 0.075]$
- Cecha  $C_5$  zawierała się w wartościach  $\in [0, 0.1]$
- Cecha  $C_6$  zawierała się w wartościach  $\in [1, 14]$
- Cecha  $C_7$  zawierała się w wartościach  $\in [0, 0.591]$
- Cecha  $C_8$  zawierała się w wartościach  $\in [0.409, 1]$
- Cecha  $C_9$  zawierała się w wartościach  $\in [1, 420]$
- Cecha  $C_{10}$  zawierała się w wartościach  $\in [1, 574]$
- Cecha  $C_{11}$  jest cechą tekstową i nie trzeba jej było normalizować, gdyż po użyciu miary  $n$ -gramów przyjmuje wartości  $[0, 1]$ .

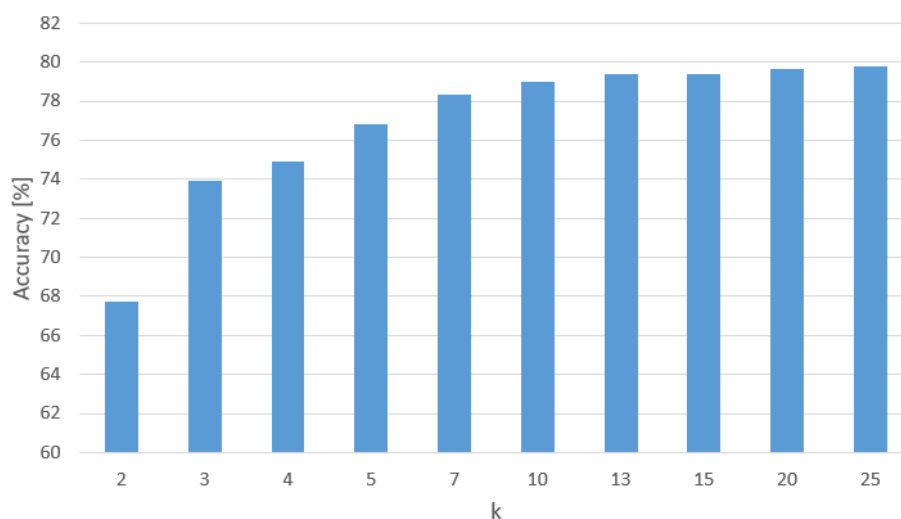
## 5.2. Badanie wyników klasyfikacji w zależności od parametru k

Parametr k	Accuracy [%]	Place	Precision [%]	Recall [%]
2	67,705	USA	81,637	81,394
		CAN	8,952	10,733
		FRA	6,400	6,957
		JAP	11,828	7,801
		W-G	13,821	9,714
		UK	14,021	16,708
3	73,921	USA	81,271	90,093
		CAN	8,108	4,712
		FRA	11,321	5,217
		JAP	13,333	6,383
		W-G	17,284	8,000
		UK	17,483	12,285
4	74,900	USA	80,895	92,398
		CAN	13,939	6,021
		FRA	14,706	4,348
		JAP	8,000	3,546
		W-G	20,000	4,571
		UK	15,086	8,600
5	76,799	USA	80,703	94,684
		CAN	13,265	3,403
		FRA	15,385	3,478
		JAP	8,642	2,482
		W-G	34,211	7,429
		UK	17,742	8,108
7	78,312	USA	80,236	97,342
		CAN	15,686	2,094
		FRA	50,000	1,739
		JAP	10,000	1,064
		W-G	18,750	1,714
		UK	23,009	6,388
10	79,024	USA	79,897	98,401
		CAN	16,129	1,309
		FRA	100,000	0,870
		JAP	5,000	0,355
		W-G	0,000	0,000
		UK	7,937	1,229

Tabela 1. Zależność Accuracy od wartości k.

Parametr k	Accuracy [%]	Place	Precision [%]	Recall [%]
13	79,365	USA	79,844	99,257
		CAN	19,048	1,047
		FRA	0,000	0,000
		JAP	0,000	0,000
		W-G	0,000	0,000
		UK	16,000	0,983
15	79,410	USA	79,899	99,517
		CAN	22,222	1,047
		FRA	0,000	0,000
		JAP	0,000	0,000
		W-G	0,000	0,000
		UK	27,778	1,229
20	79,632	USA	79,807	99,684
		CAN	18,750	0,785
		FRA	0,000	0,000
		JAP	0,000	0,000
		W-G	0,000	0,000
		UK	0,000	0,000
25	79,795	USA	79,828	99,814
		CAN	25,000	0,785
		FRA	0,000	0,000
		JAP	0,000	0,000
		W-G	0,000	0,000
		UK	0,000	0,000

Tabela 2. Zależność Accuracy od wartości k (cd.).

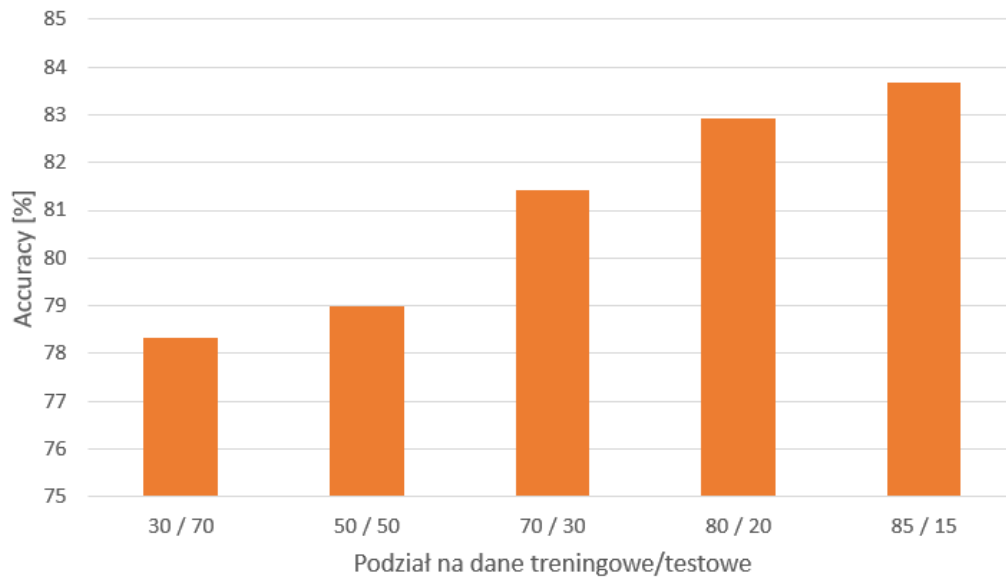


Rysunek 4. Wykres przedstawiający zależność Accuracy od wartości k (dane treningowe/testowe 50%/50%, Metryka euklidesowa).

### 5.3. Badanie wyników klasyfikacji w zależności od podziału na dane treningowe i testowe

Dane tren./test.	Accuracy [%]	Place	Precision [%]	Recall [%]
30/70	78,325	USA	79,525	98,918
		CAN	10,526	0,743
		FRA	0,000	0,000
		JAP	0,000	0,000
		W-G	0,000	0,000
		UK	12,500	0,963
50/50	78,979	USA	79,958	98,773
		CAN	17,241	1,309
		FRA	0,000	0,000
		JAP	7,692	0,355
		W-G	0,000	0,000
		UK	13,462	1,720
70/30	81,409	USA	81,839	98,522
		CAN	10,000	0,439
		FRA	0,000	0,000
		JAP	9,091	0,599
		W-G	50,000	1,042
		UK	13,333	3,030
80/20	82,911	USA	83,570	98,764
		CAN	20,000	0,649
		FRA	0,000	0,000
		JAP	0,000	0,000
		W-G	100,000	1,574
		UK	23,333	6,140
85/15	83,667	USA	84,634	98,846
		CAN	25,000	0,952
		FRA	0,000	0,000
		JAP	0,000	0,000
		W-G	100,000	2,500
		UK	11,111	2,740

Tabela 3. Zależność Accuracy od pięciu wartości proporcji podziału zbioru dla  $k=10$ , metryka uliczna.

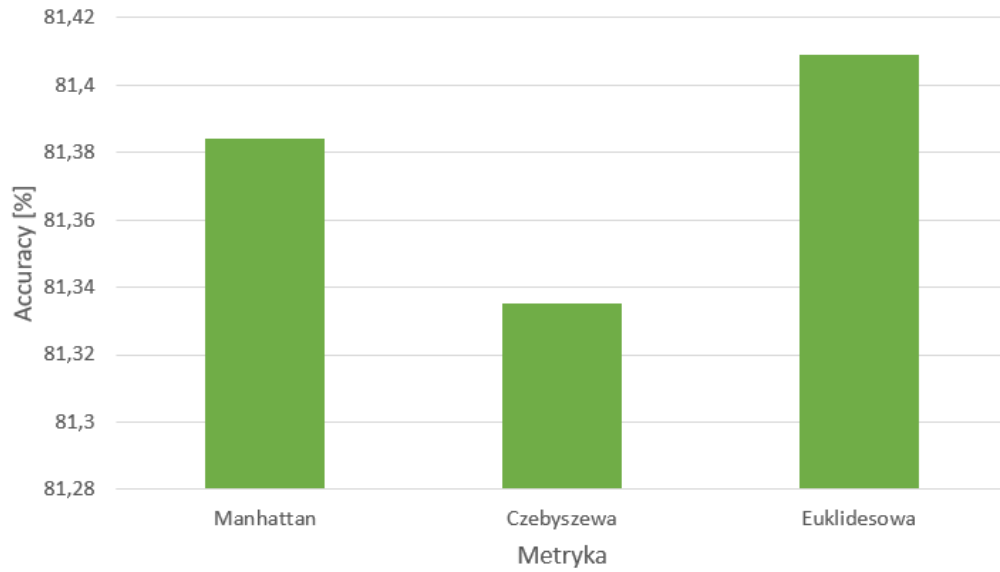


Rysunek 5. Wykres przedstawiający zależność Accuracy od pięciu wartości proporcji podziału zbioru,  $k=10$ , metryka uliczna.

#### 5.4. Badanie zależności Accuracy od wyboru metryki

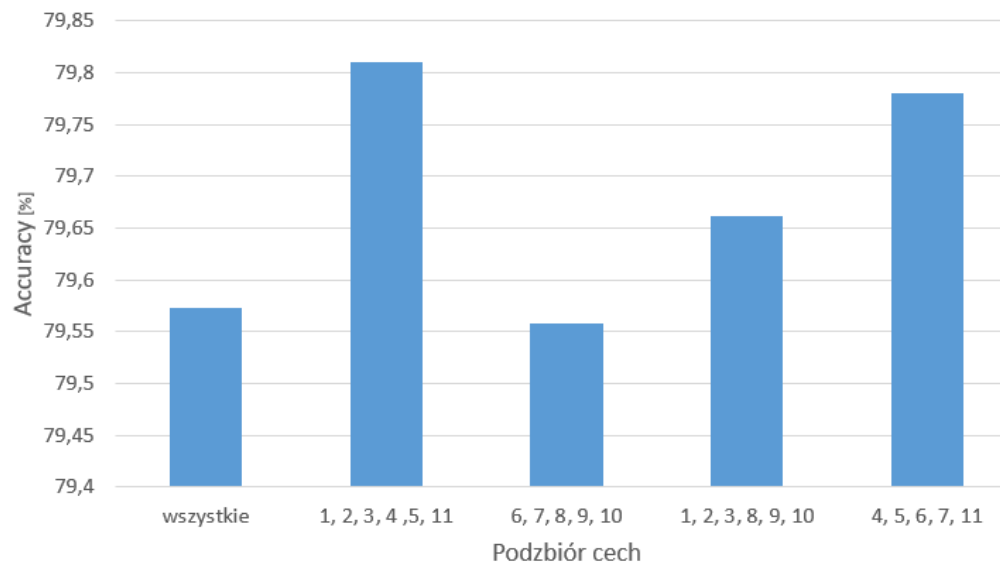
Metryka	Accuracy [%]	Place	Precision [%]	Recall [%]
euklidesowa	81,409	USA	81,918	99,457
		CAN	25,000	0,439
		FRA	0,000	0,000
		JAP	0,000	0,000
		W-G	0,000	0,000
		UK	25,000	3,535
Czebyszewa	81,335	USA	81,900	99,337
		CAN	40,000	0,877
		FRA	0,000	0,000
		JAP	0,000	0,000
		W-G	0,000	0,000
		UK	25,000	3,535
uliczna	81,384	USA	81,816	99,457
		CAN	100,000	0,877
		FRA	0,000	0,000
		JAP	33,333	0,599
		W-G	0,000	0,000
		UK	16,000	2,020

Tabela 4. Zależność Accuracy od wyboru metryki dla  $k=13$  i podziału 70/30.



Rysunek 6. Wykres przedstawiający zależność Accuracy od wyboru metryki dla  $k=13$  i podziału 70/30.

### 5.5. Badanie różnic w wyborze podzbioru cech



Rysunek 7. Wykres przedstawiający zależność Accuracy od wyboru podzbioru cech,  $k=20$ , podział 50/50, metryka Czebyszewa.

Podzbiór cech	Accuracy [%]	Place	Precision [%]	Recall [%]
Wszystkie cechy	79,573	USA	79,816	99,740
		CAN	33,333	1,047
		FRA	0,000	0,000
		JAP	0,000	0,000
		W-G	0,000	0,000
		UK	0,000	0,000
$C_1, C_2, C_3, C_4, C_5, C_{11}$	79,81	USA	79,81	100,000
		CAN	0,000	0,000
		FRA	0,000	0,000
		JAP	0,000	0,000
		W-G	0,000	0,000
		UK	0,000	0,000
$C_6, C_7, C_8, C_9, C_{10}$	79,558	USA	79,863	99,591
		CAN	22,222	0,524
		FRA	0,000	0,000
		JAP	0,000	0,000
		W-G	0,000	0,000
		UK	14,286	0,737
$C_1, C_2, C_3, C_8, C_9, C_{10}$	79,662	USA	80,357	99,461
		CAN	16,667	0,524
		FRA	20,690	5,217
		JAP	22,222	0,709
		W-G	50,000	1,143
		UK	25,000	1,720
$C_4, C_5, C_6, C_7, C_{11}$	79,78	USA	79,813	99,944
		CAN	0,000	0,000
		FRA	0,000	0,000
		JAP	0,000	0,000
		W-G	0,000	0,000
		UK	25,000	0,246

Tabela 5. Zależność Accuracy od wyboru podzbioru cech dla  $k=20$ , podziału 50/50 i metryki Czebyszewa.

## 6. Dyskusja

### 6.1. Wyniki klasyfikacji w zależności od parametru $k$

Dobór odpowiedniej wartości parametru  $k$  ma duży wpływ na wynik klasyfikacji. Im większa wartość  $k$ , tym większa skuteczność klasyfikacji. Można to zauważyć na rysunku 4, jak i również tabeli 1. Największy skok widać między  $k = 2$  a  $k = 3$ . Natomiast przy wartościach wyższych niż 10 różnica jest bardzo niewielka. Najlepsze wyniki klasyfikacji były osiągane, gdy  $k = 25$  - ok. 80%, natomiast najgorsze, gdy  $k = 2$  - ok. 68%.



## 6.2. Wyniki klasyfikacji w zależności od podziału na dane treningowe i testowe

W przypadku podziału na dane testowe i treningowe, dla większej ilości danych treningowych algorytm był w stanie lepiej zaklasyfikować teksty i skuteczność klasyfikacji była wyższa. Dla badanych przez nas podziałów największą skuteczność osiągnął podział 85% dane treningowe i 15% dane testowe. Najgorsza jakość klasyfikacji została osiągnięta dla podziału 30% treningowe / 70% testowe.

Analizując rysunek 5 oraz tabelę 3 największy skok wystąpił między podziałami 50/50 a 70/30. W przypadku tego drugiego podziału algorytm był w stanie się lepiej nauczyć.

## 6.3. Zależność Accuracy od wyboru metryki

Wybór metryki w przypadku klasyfikowania artykułów wg kategorii *places* miał bardzo mały wpływ na wyniki klasyfikacji. Widać to na rysunku 6, a najlepiej w tabeli 4. Pomiedzy najlepszym wynikiem dla metryki euklidesowej, a najgorszym dla metryki Czebyszewa różnica była mniejsza niż 0.1 p.p. Możemy więc powiedzieć, że w przypadku algorytmu k-NN wybór metryki ma minimalny wpływ na wyniki klasyfikacji, jednak najlepszą metryką jest metryka euklidesowa.

## 6.4. Różnice w wyborze podzbioru cech

Jak można zauważyć na rysunku 7 oraz tabeli 5, najlepszą wartość klasyfikacji (79.81%) osiągnął podzbiór składający się z cech:  $C_1, C_2, C_3, C_4, C_5, C_{11}$ . Są to cechy zależne od słów kluczowych. Natomiast cechy niezależne od słów kluczowych, tj.  $C_6, C_7, C_8, C_9, C_{10}$  osiągnęły wynik najgorszy (79.558%). Dla wszystkich cech ten wynik był podobny do wyniku dla cech 6-10. Różnice w skuteczności są niewielkie, lecz zauważalne.

## 7. Wnioski

- Liczba k sąsiadów ma spory wpływ na skuteczność klasyfikacji. Jednakże, zmiana metryki, bądź podziału na dane testowe i treningowe również ma wpływ na wynik klasyfikacji.
- Również istotny jest podział artykułów na dane testowe i treningowe.
- W przypadku zbyt małej ilości danych treningowych wystąpi zjawisko niedouczenia.
- W przypadku zbyt dużej ilości danych treningowych wystąpi zjawisko przeuczenia.
- Wybór metryki ma minimalny, bliski zeru wpływ na wyniki klasyfikacji.
- Najlepszymi cechami do klasyfikacji tekstów będą te cechy, które zależą od słów kluczowych, ale tylko wtedy, gdy te zostaną wybrane równomiernie dla każdej kategorii.

## Literatura

- [1] Niewiadomski, Adam. *Methods for the Linguistic Summarization of Data: Applications of Fuzzy Sets and Their Extensions*. Akademicka Oficyna Wydawnicza EXIT. Warszawa, 2008. ISBN 978-83-60434-40-6
- [2] Ciaburro, Giuseppe. *Regression Analysis with R*. Packt Publishing, 2018. ISBN 978-1-78862-730-6
- [3] [https://ftims.edu.p.lodz.pl/pluginfile.php/132368/mod\\_folder/content/0/ksr-wyklad-2009.pdf?forcedownload=1](https://ftims.edu.p.lodz.pl/pluginfile.php/132368/mod_folder/content/0/ksr-wyklad-2009.pdf?forcedownload=1) [dostęp 22.03.2020]
- [4] [https://en.wikipedia.org/wiki/Taxicab\\_geometry](https://en.wikipedia.org/wiki/Taxicab_geometry) [dostęp 01.04.2020]
- [5] [https://en.wikipedia.org/wiki/Chebyshev\\_distance](https://en.wikipedia.org/wiki/Chebyshev_distance) [dostęp 01.04.2020]
- [6] [https://en.wikipedia.org/wiki/Euclidean\\_distance](https://en.wikipedia.org/wiki/Euclidean_distance) [dostęp 01.04.2020]
- [7] <https://tartarus.org/martin/PorterStemmer/csharp.txt> [dostęp 22.03.2020]
- [8] <https://www.dotnetperls.com/stopword-dictionary> [dostęp 22.03.2020]
- [9] <http://home.agh.edu.pl/horzyk/lectures/miw/KNN.pdf> [dostęp 22.03.2020]
- [10] [https://pl.wikipedia.org/wiki/Tablica\\_pomy%C5%82ek](https://pl.wikipedia.org/wiki/Tablica_pomy%C5%82ek) [dostęp 01.04.2020]