

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

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

Justyna Hubert 210200

Karol Podlewski 210294

## Zadanie 1: Ekstrakcja cech, miary podobieństwa, klasyfikacja\*

### 1. Cel

Celem zadania było stworzenie aplikacji do klasyfikacji metodą k-NN, korzystając z różnych sposobów ekstrakcji wektorów cech oraz istniejących miar podobieństwa porównać kategorie tekstów do tych przypisanych przez aplikację.

### 2. Wprowadzenie

term frequency:

$$tf_{i,j} = \frac{n_{i,j}}{\sum_k n_{k,j}}$$

inverse document frequency:

$$idf_i = \log \frac{|D|}{|\{d : t_i \in d\}|}$$

metryka euklidesowa:

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

metryka uliczna:

$$d_m(x, y) = \sum_{k=1}^n |x_k - y_k|$$

---

\* SVN: <https://github.com/hubjust/KSR>

metryka czebyszewa:

$$d_{ch}(x, y) = \max_i |x_i - y_i|$$

We wprowadzeniu należy zaprezentować całą teorię potrzebną do realizacji zadania (przy czym należy tu ograniczyć się wyłącznie do tego, co było wykorzystane) tak aby osoba, która nigdy wcześniej nie zetknęła się z tą tematyką, potrafiła zrozumieć dalszy opis. Część ta powinna wprowadzać wszystkie wykorzystywane wzory, oznaczenia itp., do których należy się odwoływać w dalszej części niniejszego sprawozdania. Zamieszczony tu własny opis teorii (a nie skopiowany!) należy poprzeć odwołaniami bibliograficznymi do literatury zamieszczonej na końcu.

### 3. Opis implementacji

Program został stworzony w języku C#. Graficzny interfejs użytkownika został stworzony przy wykorzystaniu Windows Presentation Foundation. Logika aplikacji została odseparowana od GUI, w zgodzie ze wzorcem projektowym Model-view-viewmodel (MVVM), poprzez implementacje trzech projektów (Logic, ViewModel i GUI).

Należy tu zamieścić krótki i zwięzły opis zaprojektowanych klas oraz powiązań między nimi. Powinien się tu również znaleźć diagram UML (diagram klas) prezentujący najistotniejsze elementy stworzonej aplikacji. Należy także podać, w jakim języku programowania została stworzona aplikacja.

### 4. Materiały i metody

W tym miejscu należy opisać, jak przeprowadzone zostały wszystkie badania, których wyniki i dyskusja zamieszczane są w dalszych sekcjach. Opis ten powinien być na tyle dokładny, aby osoba czytająca go potrafiła wszystkie przeprowadzone badania samodzielnie powtórzyć w celu zweryfikowania ich poprawności (a zatem m.in. należy zamieścić tu opis architektury sieci, wartości współczynników użytych w kolejnych eksperymentach, sposób inicjalizacji wag, metodę uczenia itp. oraz informacje o danych, na których prowadzone były badania). Przy opisie należy odwoływać się i stosować do opisanych w sekcji drugiej wzorów i oznaczeń, a także w jasny sposób opisać cel konkretnego testu. Najlepiej byłoby wyraźnie wyszczególnić (ponumerować) poszczególne eksperymenty tak, aby łatwo było się do nich odwoływać dalej.

### 5. Wyniki

W tej sekcji należy zaprezentować, dla każdego przeprowadzonego eksperymentu, kompletny zestaw wyników w postaci tabel, wykresów itp. Powinny być one tak ponazywane, aby było wiadomo, do czego się odnoszą. Wszystkie tabele i wykresy należy oczywiście opisać (opisać co jest na osiach, w

kolumnach itd.) stosując się do przyjętych wcześniej oznaczeń. Nie należy tu komentować i interpretować wyników, gdyż miejsce na to jest w kolejnej sekcji. Tu również dobrze jest wprowadzić oznaczenia (tabel, wykresów) aby móc się do nich odwoływać poniżej.

## **6. Dyskusja**

Sekcja ta powinna zawierać dokładną interpretację uzyskanych wyników eksperymentów wraz ze szczegółowymi wnioskami z nich płynącymi. Najcenniejsze są, rzecz jasna, wnioski o charakterze uniwersalnym, które mogą być istotne przy innych, podobnych zadaniach. Należy również omówić i wyjaśnić wszystkie napotkane problemy (jeśli takie były). Każdy wniosek powinien mieć poparcie we wcześniej przeprowadzonych eksperymentach (odwołania do konkretnych wyników). Jest to jedna z najważniejszych sekcji tego sprawozdania, gdyż prezentuje poziom zrozumienia badanego problemu.

## **7. Wnioski**

W tej, przedostatniej, sekcji należy zamieścić podsumowanie najważniejszych wniosków z sekcji poprzedniej. Najlepiej jest je po prostu wypunktować. Znow, tak jak poprzednio, najistotniejsze są wnioski o charakterze uniwersalnym.

## **Literatura**

- [1] Methods for the linguistic summarization of data - applications of fuzzy sets and their extensions, Adam Niewiadomski, Akademicka Oficyna Wydawnicza EXIT, Warszawa 2008