

Data oddania: _____

Ocena: _____

Julia Szymańska 224441
Przemysław Zdrzałik 224466

Projekt 2. Podsumowania lingwistyczne relacyjnych baz danych

1. Cel

Celem projektu jest stworzenie aplikacji lingwistycznie agregującej zawartości zbioru danych wypadków samochodowych w Zjednoczonych Stanach Ameryki[4]. Jako wynik działania aplikacji zostanie wygenerowany opis wartości danych liczbowych w zbiorze w języku quasi-naturalnym.

2. Charakterystyka podsumowywanej bazy danych

W programie został użyty zbiór danych[4] znajdujący się w pliku CSV, który został przekształcony w bazę danych.

Zbiór danych zawiera informacje o ponad 3 milionach wypadków samochodowych w 49 stanach Zjednoczonych Stanów Ameryki, mających miejsce od lutego 2016 do grudnia 2020. Spośród 47 kolumn znajdujących się w zbiorze danych, wybraliśmy następujące 11 kolumn:

- Dotkliwość - Severity - wpływ wypadku na ruch na drodze, przyjmuje wartości całkowite od 1 do 4 włącznie, gdzie 1 oznacza najmniejszy wpływ na ruch drogowy, natomiast 4 oznacza największy wpływ. Dotkliwość można lingwistycznie opisać jako mały lub duży wpływ na ruch drogowy.
- Czas rozpoczęcia - Start_Time - czas rozpoczęcia się wypadku w lokalnej strefie czasowej, przyjmuje wartości od 8 lutego 2016, do 31 grudnia 2020. Wartość kolumny zostanie zamieniona na wartość całkowitą oznaczającą liczbę sekund od początku 1970 roku.

- Czas zakończenia - End_Time - czas zakończenia się wypadku w lokalnej strefie czasowej, przyjmuje wartości od 8 lutego 2016, do 1 stycznia 2021. Wartość kolumny zostanie zamieniona na wartość całkowitą oznaczającą liczbę sekund od początku 1970 roku.
- Odległość - Distance - długość odcinka ulicy wyrażony w milach, na którego miał wpływ wypadek. Przyjmuje wartości zmiennoprzecinkowe od 0 do 334, gdzie zdecydowana większość danych mieści się w przedziale od 0.00 do 4.00.
- Temperatura - Temperature - temperatura powietrza wyrażona w Fahrenheit'ach, w momencie, gdy zdarzył się wypadek. Przyjmuje wartości zmiennoprzecinkowe od -16.00 do 104.00. Temperature można opisać jako bardzo zimną, zimną, umiarkowaną, ciepłą, bardzo ciepłą. Oczywiście jest to opis subiektywny.
- Temperatura odczuwalna - Wind_Chill - temperatura odczuwalna wyrażona w Fahrenheit'ach, w momencie, gdy zdarzył się wypadek. Przyjmuje wartości zmiennoprzecinkowe od -16.00 do 101.00. Temperaturę odczuwalną można opisać tak samo jak temperaturę.
- Wilgotność - Humidity - wilgotność powietrza wyrażona w procentach w momencie, gdy zdarzył się wypadek. Przyjmuje wartości zmiennoprzecinkowe od 4.00 do 100.00.
- Ciśnienie - Pressure - ciśnienie powietrza wyrażone w inches, w momencie, gdy zdarzył się wypadek. Przyjmuje wartości zmiennoprzecinkowe od 27.00 do 32.00. Ciśnienie można opisać jako wysokie, umiarkowane lub niskie.
- Widoczność - Visibility - widoczność wyrażona w milach, w momencie, gdy zdarzył się wypadek. Przyjmuje wartości zmiennoprzecinkowe od 0.00 do 12.00. Widoczność można opisać jako dobrą, ograniczoną, słabą.
- Prędkość wiatru - Wind_Speed - prędkość wiatru wyrażona w milach na godzinę, w momencie, gdy zdarzył się wypadek. Przyjmuje wartości zmiennoprzecinkowe od 0.00 do 40.00. Wiatr można opisać jako słaby, umiarkowany, silny.
- Ilość opadów - Precipitation - ilość opadów wyrażona w inches, w momencie, gdy zdarzył się wypadek. Jeśli opady nie występowały to kolumna przyjmuje wartość nan. Przyjmuje wartości zmiennoprzecinkowe od 0.00 do 0.50.

Atrybutom nadawane są opisane zwyczajowe wartości lingwistyczne ze względu na zwiększenie przystępności i ułatwienie szybkiego zrozumienia atrybutu przez człowieka, kiedy ten atrybut nie musi być dokładnie opisany. Przykładowo temperatura, mimo że rozumiała dla człowieka w postaci liczbowej, jest łatwiejsza do szybszego zrozumienia w postaci tekstowej, a dla ludzi nie ma dużego znaczenia czy temperatura różni się o 1 czy 2 stopnie, wystarczy opisać ją słownie tak jak wcześniej podaliśmy jako bardzo zimną, zimną, umiarkowaną, ciepłą, bardzo ciepłą.

accidents	
id	bigint
severity	smallint
start_time	timestamp with time zone
end_time	timestamp with time zone
distance	double precision
temperature	double precision
wind_chill	double precision
humidity	double precision
pressure	double precision
visibility	double precision
wind_speed	double precision
precipitation	double precision

Powered by yFiles

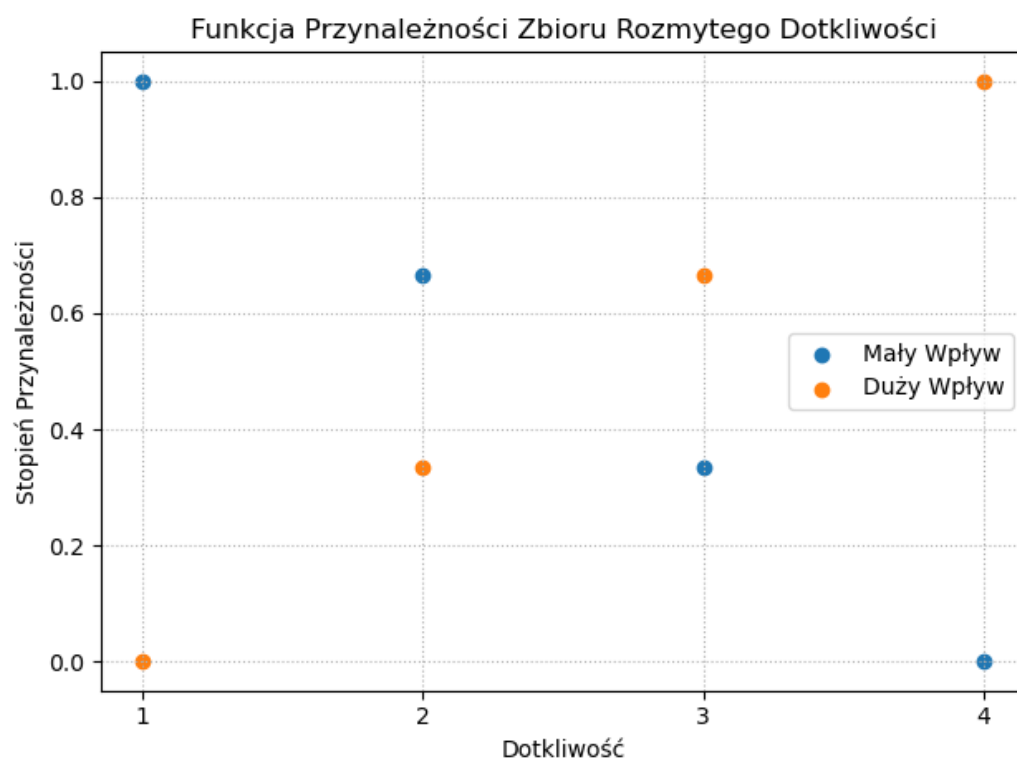
Rysunek 1. Tabela reprezentująca omawiane dane wykonana w DBMS PostgreSQL

3. Atrybuty i liczności obiektów wyrażone zmiennymi lingwistycznymi

Dotkliwość, zmienna lingwistyczna przedstawiona wraz z etykietami w formie wykresu funkcji przynależności oraz wzorów analitycznych:

$$\mu_{dotkliwoscDuza}(x) = \frac{x-1}{3} \quad dla \quad x \in \{1, 2, 3, 4\} \quad (1)$$

$$\mu_{dotkliwoscMala}(x) = \frac{4-x}{3} \quad dla \quad x \in \{1, 2, 3, 4\} \quad (2)$$



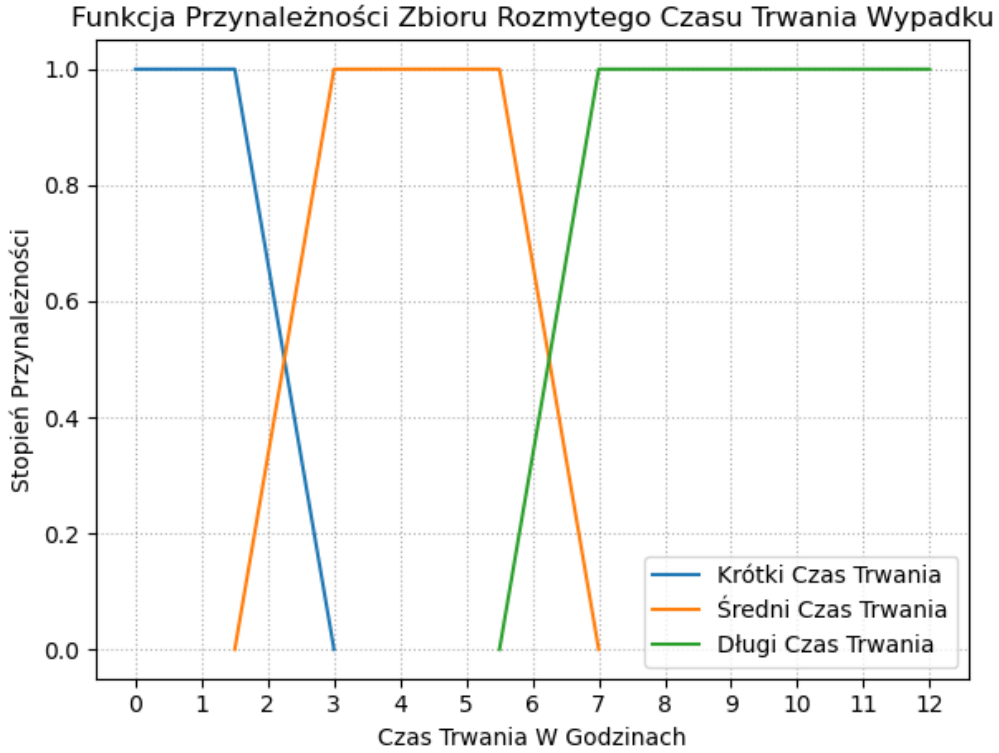
Rysunek 2. Wykres funkcji przynależności zbioru rozmytego dotkliwości.

Czas trwania wypadku, zmienna lingwistyczna przedstawiona wraz z etykietami w formie wykresu funkcji przynależności oraz wzorów analitycznych:

$$\mu_{\text{czasTrwaniaKrotki}}(x) = \begin{cases} 1 & \text{dla } x \leq 1.5 \\ \frac{3-x}{1.5} & \text{dla } 1.5 < x \leq 3 \\ 0 & \text{dla } 3 \leq x \end{cases} \quad (3)$$

$$\mu_{\text{czasTrwaniaSredni}}(x) = \begin{cases} 0 & \text{dla } x \leq 1.5 \\ \frac{x-1.5}{1.5} & \text{dla } 1.5 < x \leq 3 \\ 1 & \text{dla } 3 < x < 5.5 \\ \frac{7-x}{1.5} & \text{dla } 5.5 < x \leq 7 \\ 0 & \text{dla } x \leq 7 \end{cases} \quad (4)$$

$$\mu_{\text{czasTrwaniaDlugi}}(x) = \begin{cases} 0 & \text{dla } x \leq 5.5 \\ \frac{x-5.5}{1.5} & \text{dla } 5.5 < x \leq 7 \\ 1 & \text{dla } 7 \leq x \end{cases} \quad (5)$$



Rysunek 3. Wykres funkcji przynależności zbioru rozmytego czasu trwania wypadku.

Odległość, zmienna lingwistyczna przedstawiona wraz z etykietami w formie wykresu funkcji przynależności oraz wzorów analitycznych:

$$\mu_{OdlegloscKrotki}(x) = \begin{cases} 1 & \text{dla } x \leq 0.5 \\ \frac{1-x}{0.5} & \text{dla } 0.5 < x \leq 1 \\ 0 & \text{dla } 1 \leq x \end{cases} \quad (6)$$

$$\mu_{OdlegloscDlugi}(x) = \begin{cases} 0 & \text{dla } x \leq 0.5 \\ \frac{x-0.5}{0.5} & \text{dla } 0.5 < x \leq 1 \\ 1 & \text{dla } 1 \leq x \end{cases} \quad (7)$$

Temperatura w momencie wypadku, zmienna lingwistyczna przedstawiona wraz z etykietami w formie wykresu funkcji przynależności oraz wzorów analitycznych:

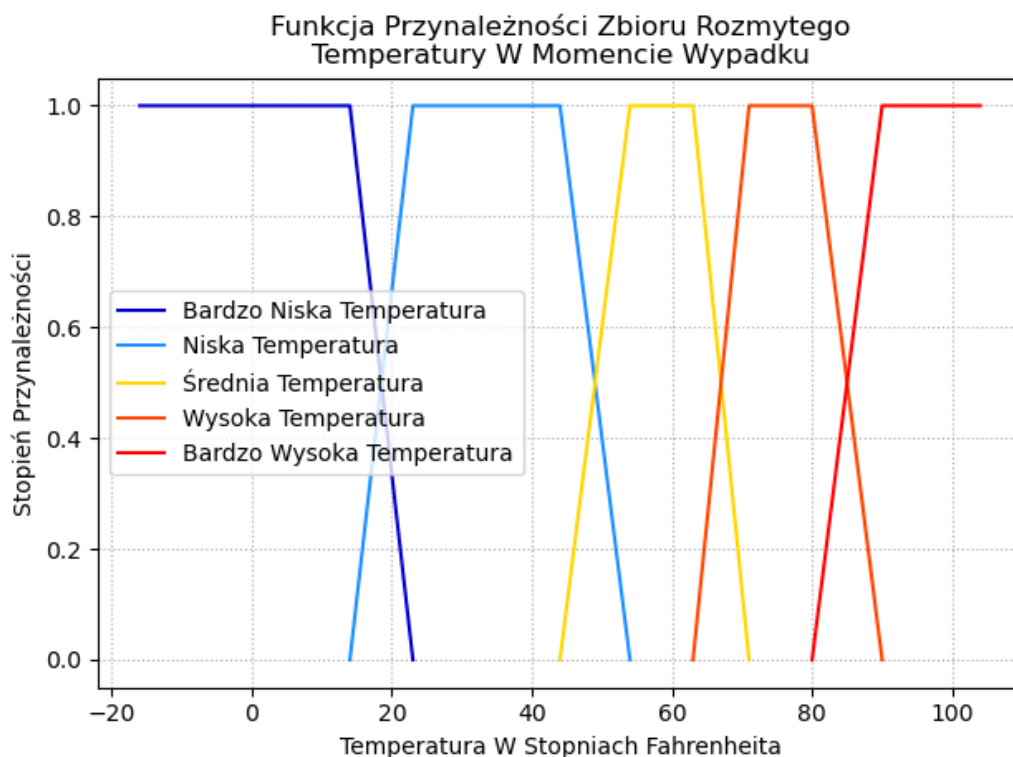
$$\mu_{temperaturaBardzoZimno}(x) = \begin{cases} 1 & \text{dla } x \leq 14 \\ \frac{23-x}{9} & \text{dla } 14 < x \leq 23 \\ 0 & \text{dla } 23 \leq x \end{cases} \quad (8)$$

$$\mu_{temperaturaZimno}(x) = \begin{cases} 0 & \text{dla } x \leq 14 \\ \frac{x-14}{9} & \text{dla } 14 < x \leq 23 \\ 1 & \text{dla } 23 < x < 44 \\ \frac{54-x}{10} & \text{dla } 44 < x \leq 54 \\ 0 & \text{dla } x \leq 54 \end{cases} \quad (9)$$

$$\mu_{temperaturaSrednia}(x) = \begin{cases} 0 & \text{dla } x \leq 44 \\ \frac{x-44}{10} & \text{dla } 44 < x \leq 54 \\ 1 & \text{dla } 54 < x < 63 \\ \frac{71-x}{8} & \text{dla } 63 < x \leq 71 \\ 0 & \text{dla } x \leq 71 \end{cases} \quad (10)$$

$$\mu_{temperaturaCieplo}(x) = \begin{cases} 0 & \text{dla } x \leq 63 \\ \frac{x-63}{8} & \text{dla } 63 < x \leq 71 \\ 1 & \text{dla } 71 < x < 80 \\ \frac{90-x}{10} & \text{dla } 80 < x \leq 90 \\ 0 & \text{dla } x \leq 90 \end{cases} \quad (11)$$

$$\mu_{temperaturaBardzoCieplo}(x) = \begin{cases} 0 & \text{dla } x \leq 80 \\ \frac{x-80}{10} & \text{dla } 80 < x \leq 90 \\ 1 & \text{dla } 90 \leq x \end{cases} \quad (12)$$



Rysunek 4. Wykres funkcji przynależności zbioru rozmytego czasu trwania wypadku.

Zmienne lingwistyczne dla wybranych 10 atrybutów z bazy danych, przedstawione w formie wykresów funkcji przynależności i wzorów analitycznych, wymienione etykiety oraz objaśnione wszystkie symbole ułatwiające czytelnikowi ich zrozumienie [2]. Zbędne jest cytowanie definicji. Konieczne precyzyjnie podane przestrzenie rozważań każdej zmiennej lingwistycznej, wzory i wykresy dla każdej wartości/etykiety.

Jw. kwantyfikatory lingwistyczne – opisane etykietami, wykresami funkcji przynależności i wzorami analitycznymi. Uzasadnione wiedzą dziedzinową wybrane zakresy i etykiety. Precyzyjnie podane przestrzenie rozważań każdego kwantyfikatora lingwistycznego/rozmytego, wzory i wykresy dla każdej wartości/etykiety. Opisy własne z przypisami do literatury, tak by inżynier innej specjalności zrozumiał dalszy opis tego konkretnego ćwiczenia/eksperymentu.

Sekcja uzupełniona jako efekt zadania Tydzień 09 wg Harmonogramu Zajęć na WIKAMP KSR.

4. Narzędzia obliczeniowe: projekt (wybór, implementacja) i diagram UML pakietu obliczeń rozmytych. Diagram UML generatora podsumowań

4.1. Diagram pakietu obliczeń rozmytych

Diagram UML i zwięzły opis pakietu obliczeń rozmytych: źródło pakietu (zewnętrzny/własny/hybrydowy), przypis do literatury. Krótka charakterystyka najważniejszych klas i podstawowych dla zadania ich metod.

Sekcja uzupełniona jako efekt zadania Tydzień 10 wg Harmonogramu Zajęć na WIKAMP KSR.

4.2. Diagram UML generatora podsumowań. Krótka instrukcja użytkownika

Diagram UML generatora podsumowań (warstwy obliczeniowej oraz interfejsu użytkownika). Krótki ilustrowany opis jak użytkownik może korzystać z aplikacji, w szczególności wprowadzać parametry podsumowań, odczytywać wyniki oraz definiować własne etykiety i kwantyfikatory. Wersja JRE i inne wymagania niezbędne do uruchomienia aplikacji przez użytkownika na własnym komputerze.

Sekcja uzupełniona jako efekt zadania Tydzień 11 wg Harmonogramu Zajęć na WIKAMP KSR.

5. Jednopodmiotowe podsumowania lingwistyczne. Miary jakości, podsumowanie optymalne

Wyniki kolejnych eksperymentów wg punktów 2.-4. opisu projektu 2. Listy podsumowań jednopodmiotowych i tabele/rankingi podsumowań dla danych atrybutów obowiązkowe i dokładnie opisane w „captions” (tytułach), konieczny opis kolumn i wierszy tabel. Dla każdego podsumowania podane miary jakości oraz miara jakości podsumowania optymalnego.

Sekcja uzupełniona jako efekt zadania Tydzień 11 wg Harmonogramu Zajęć na WIKAMP KSR.

6. Wielopodmiotowe podsumowania lingwistyczne i ich miary jakości

Wyniki kolejnych eksperymentów wg punktów 2.-4. opisu projektu 2. Uzasadnienie i metoda podziału zbioru danych na rozłączne podmioty. Listy podsumowań wielopodmiotowych i tabele/rankingi podsumowań dla danych atrybutów obowiązkowe i dokładnie opisane w „captions” (tytułach), konieczny opis kolumn i wierszy tabel. Konieczne uwzględnienie wszystkich 4-ch form podsumowań wielopodmiotowych.

** Możliwe sformułowanie zagadnienia wielopodmiotowego podsumowania optymalnego **.

Ewentualne wyniki realizacji punktu „na ocenę 5.0” wg opisu Projektu 2. i ich porównanie do wyników z części obowiązkowej.

Sekcja uzupełniona jako efekt zadania Tydzień 12 wg Harmonogramu Zajęć na WIKAMP KSR.

7. Dyskusja, wnioski

Dokładne interpretacje uzyskanych wyników w zależności od parametrów klasyfikacji opisanych w punktach 3.-4 opisu Projektu 2. Szczególnie istotne są wnioski o charakterze uniwersalnym, istotne dla podobnych zadań. Omówić i wyjaśnić napotkane problemy (jeśli były). Każdy wniosek/problem powinien mieć poparcie w przeprowadzonych eksperymentach (odwołania do konkretnych wyników: tabel i miar jakości). Ocena które wybrane kwantyfikatory, sumaryzatory, kwalifikatory i/lub ich miary jakości mają małe albo duże znaczenie dla wiarygodności i jakości otrzymanych agregacji/podsumowań. Dla końcowej oceny jest to najważniejsza sekcja sprawozdania, gdyż prezentuje poziom zrozumienia rozwiązywanego problemu.

** Możliwości kontynuacji prac w obszarze logiki rozmytej i wnioskowania rozmytego, zwłaszcza w kontekście pracy inżynierskiej, magisterskiej, naukowej, itp. **

Sekcja uzupełniona jako efekt zadań Tydzień 11 i Tydzień 12 wg Harmonogramu Zajęć na WIKAMP KSR.

8. Braki w realizacji projektu 2.

Wymienić wg opisu Projektu 2. wszystkie niezrealizowane obowiązkowe elementy projektu, ewentualnie podać merytoryczne (ale nie czasowe) przyczyny tych braków.

Literatura

- [1] A. Niewiadomski, Zbiory rozmyte typu 2. Zastosowania w reprezentowaniu informacji. Seria „Problemy współczesnej informatyki” pod redakcją L. Rutkowskiego. Akademicka Oficyna Wydawnicza EXIT, Warszawa, 2019.
- [2] S. Zadrozny, Zapytania nieprecyzyjne i lingwistyczne podsumowania baz danych, EXIT, 2006, Warszawa
- [3] A. Niewiadomski, Methods for the Linguistic Summarization of Data: Applications of Fuzzy Sets and Their Extensions, Akademicka Oficyna Wydawnicza EXIT, Warszawa, 2008.
- [4] 2021 Kaggle Inc [internetowa społeczność związana z analizą danych], US Accidents (3 million records – updated) A Countrywide Traffic Accident

Dataset (2016 - 2020) [przełączany 24 kwietnia 2021], Dostępny w: <https://www.kaggle.com/sobhanmoosavi/us-accidents>

Literatura zawiera wyłącznie źródła recenzowane i/lub o potwierdzonej wiarygodności, możliwe do weryfikacji i cytowane w sprawozdaniu.