



Majkowska Agata

Migdał-Najman Kamila

Najman Krzysztof

Raca Katarzyna

Identyfikacja grup wiekowych użytkowników portalu X

www.ug.edu.pl



Plan prezentacji



Problem i cel badawczy

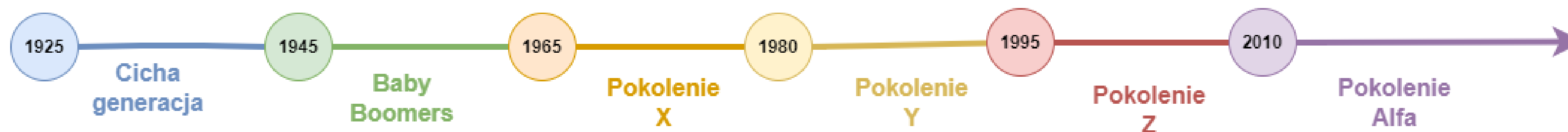
Z badań językowych wynika, że przedstawiciele różnych grup wiekowych posługują się **odmiennym słownictwem** i formami gramatycznymi. Wydaje się, że różnicują je także powszechnie stosowane w wypowiedziach **znaki graficzne**, takie jak emotikony, emoji, piktogramy i inne znaki graficzne.

Celem prezentowanych badań jest próba identyfikacji wieku użytkownika portalu X na podstawie **charakterystycznych elementów w tekście takich jak słowa czy emoji**. Do realizacji badania wykorzystano narzędzia Text Mining oraz naiwny klasyfikator Bayesa.

Pokolenie

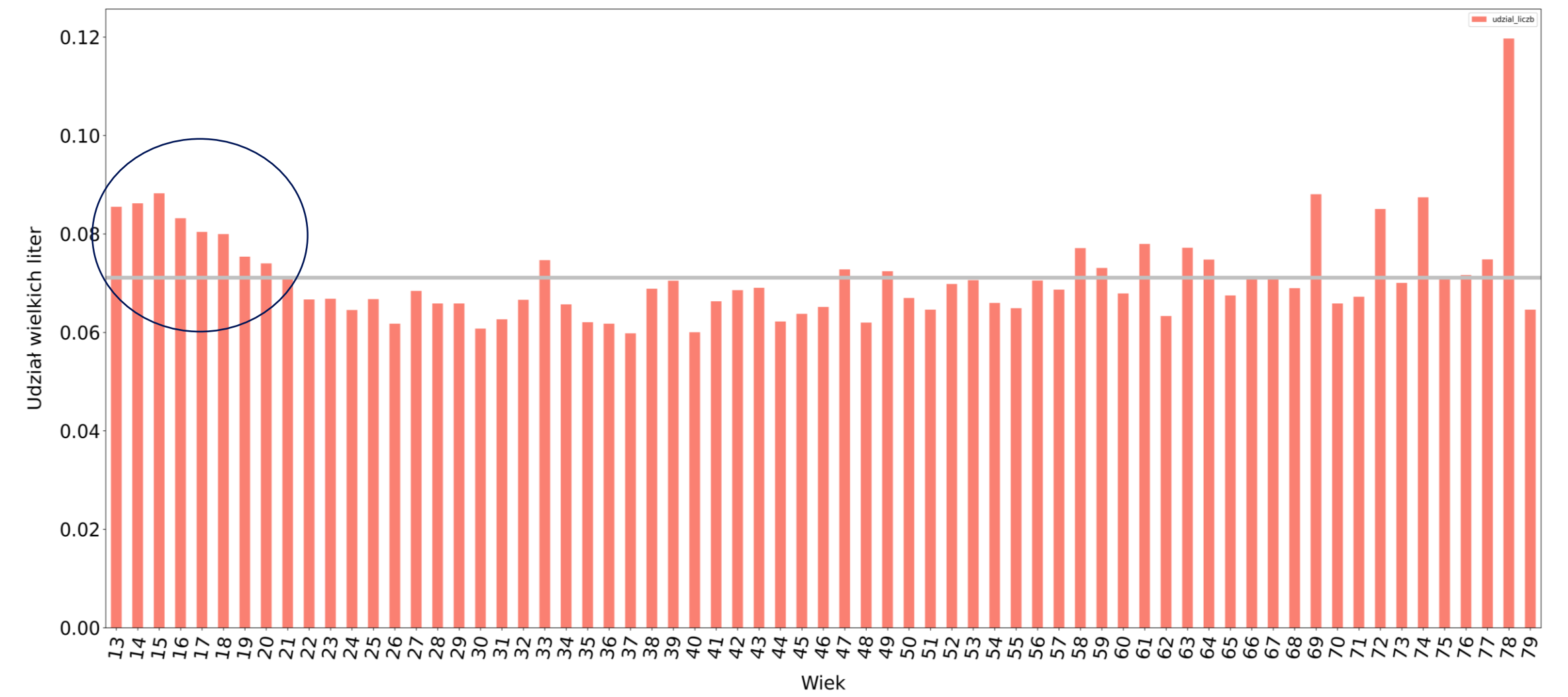
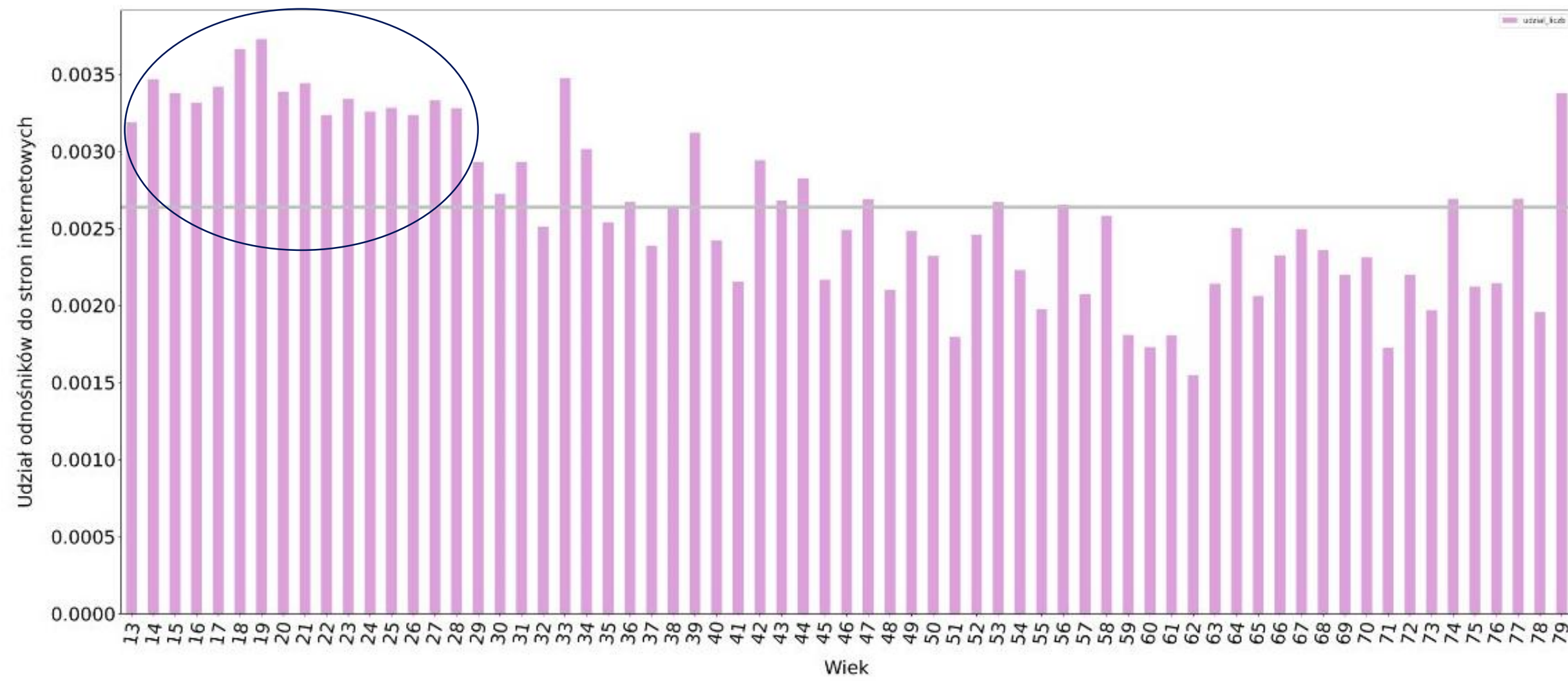
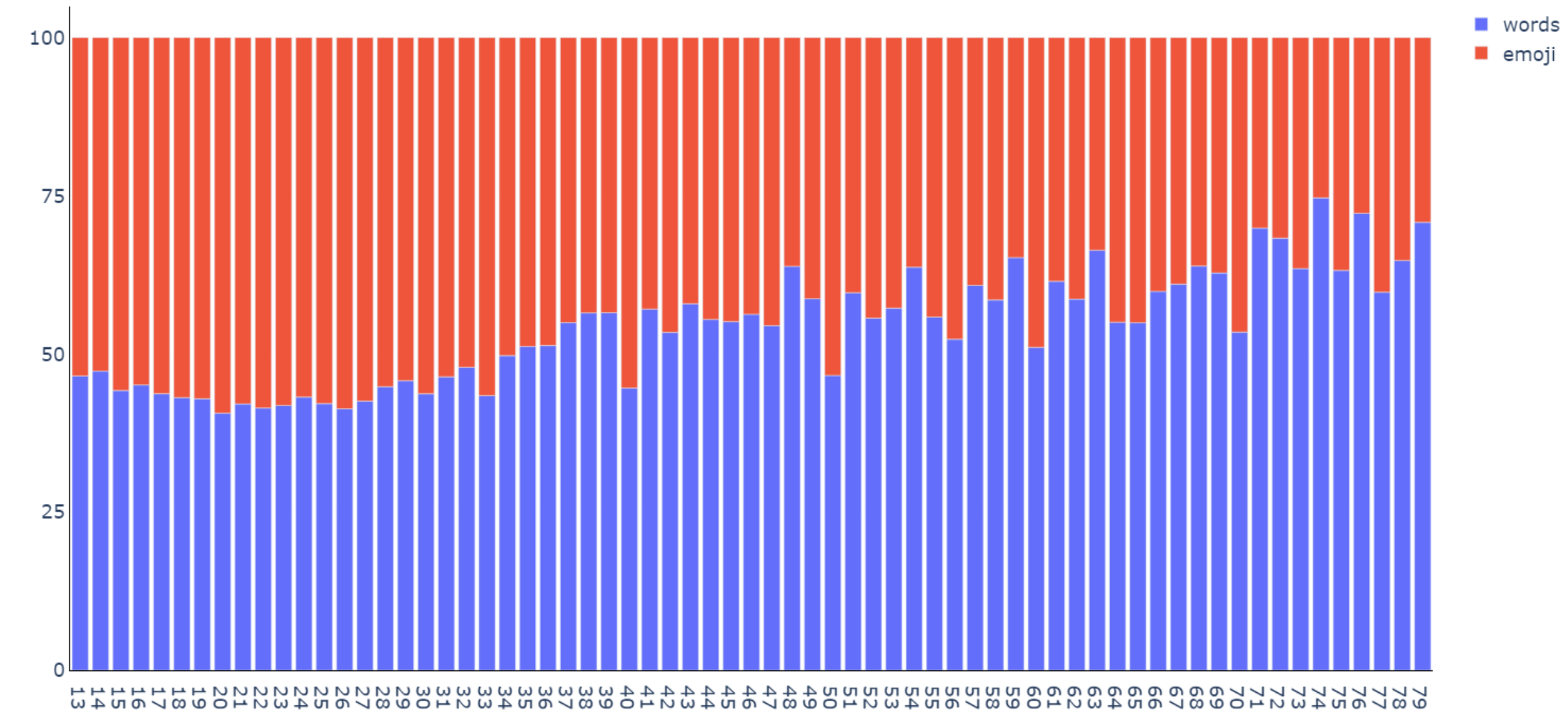
Grupa, która ma wspólne dzieciństwo, wspólny wiek młodości i u których etap dojrzałości przypada na ten sam okres. Generację tworzą jednostki, które uczestniczyły w tych samych przemianach i wydarzeniach, jakie miały miejsce w okresie ich pobudliwości. Źródłem tej odmienności są różne doświadczenia i perspektywy życiowe. [Dilthey, 1924]

Grupa społeczna, która nawiązuje ze sobą wspólną więź i wspólnie działa. Charakteryzuje się specyficznymi wzorcami zachowań, widocznie różnymi od innych grup postawami, poglądami, uznawanymi wartościami, aspiracjami, sposobem życia czy na swój sposób odmiennym pojmowaniem rzeczywistości. W definicji tej zwraca się uwagę na aspekt czasowy i jakościowy. [Sztompka, 2003]

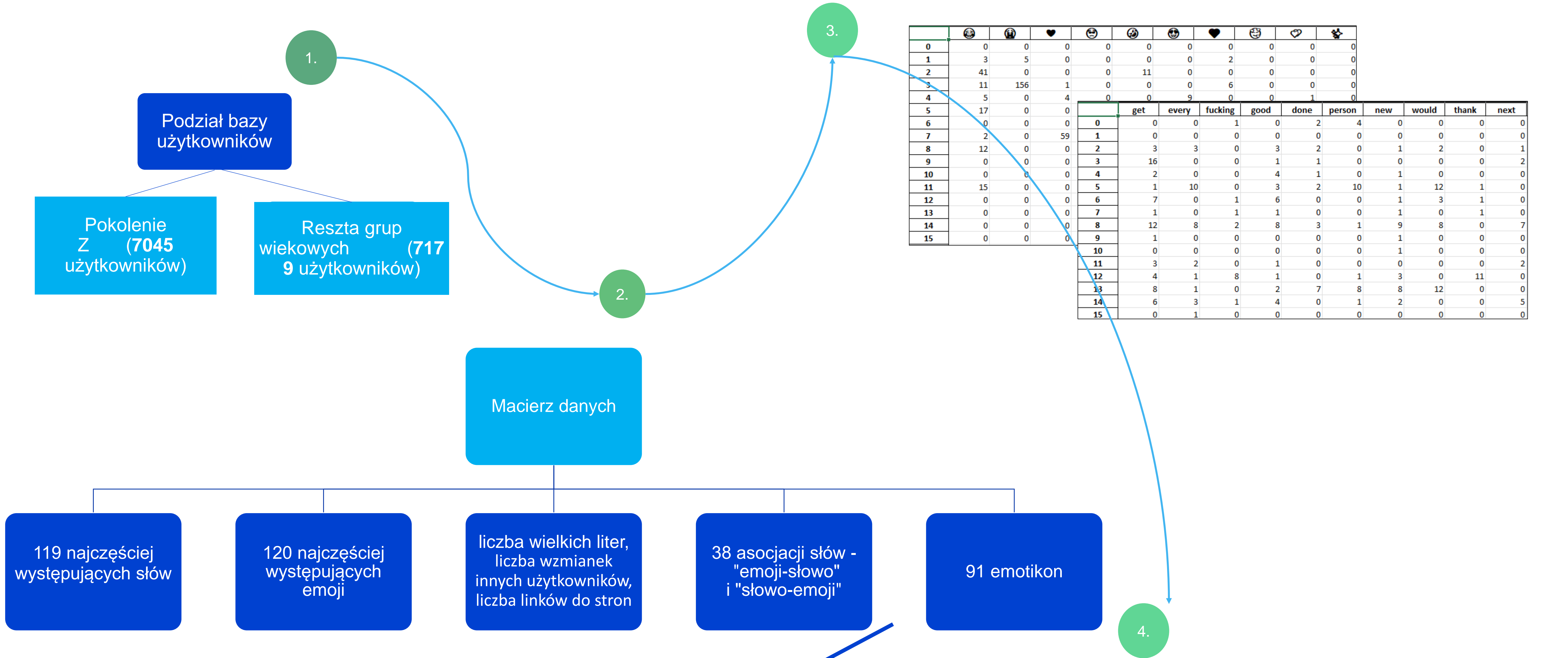


Uzyskane grupowania we wcześniejszych badaniach

Metoda grupowania	Liczba powstałych grup	Cechy wyodrębnione do analizy
grupowanie aglomeracyjne	4	słowa
sieć neuronowa	5	emoji
grupowanie aglomeracyjne	4	emotikony, stopwords, interpunkcja, liczba wzmianek o użytkownikach (@użytkownik), udostępnianie linków do stron, liczba wielkich liter
grupowanie aglomeracyjne	4	słowa i emoji



Przygotowanie danych i zastosowana metoda klasyfikacji



	🤔	🤩	❤️	😬	🤪	😄	❤️	😁	🤍	👑
0	0	0	0	0	0	0	0	0	0	0
1	3	5	0	0	0	0	2	0	0	0
2	41	0	0	0	11	0	0	0	0	0
3	11	156	1	0	0	0	6	0	0	0
4	5	0	4	0	0	9	0	0	1	0
5	17	0	0	0	0	0	0	0	0	0
6	0	0	0	0	0	0	0	0	0	0
7	2	0	59	0	0	0	0	0	0	0
8	12	0	0	0	3	3	2	0	1	2
9	0	0	0	0	3	16	0	1	1	0
10	0	0	0	0	4	2	0	0	1	0
11	15	0	0	0	5	1	10	0	3	2
12	0	0	0	0	6	7	0	1	6	0
13	0	0	0	0	7	1	0	1	1	0
14	0	0	0	0	8	12	8	2	8	3
15	0	0	0	0	9	1	0	0	0	0

	get	every	fucking	good	done	person	new	would	thank	next
0	0	0	1	0	2	4	0	0	0	0
1	0	0	0	0	0	0	0	0	0	0
2	3	3	0	3	2	0	1	2	0	1
3	16	0	0	1	1	0	0	0	0	2
4	2	0	0	4	1	0	1	0	0	0
5	1	10	0	3	2	10	1	12	1	0
6	7	0	1	6	0	0	1	3	1	0
7	1	0	1	1	0	0	1	0	1	0
8	12	8	2	8	3	1	9	8	0	7
9	1	0	0	0	0	0	1	0	0	0
10	0	0	0	0	0	0	1	0	0	0
11	3	2	0	1	0	0	0	0	0	2
12	4	1	8	1	0	1	3	0	11	0
13	8	1	0	2	7	8	8	12	0	0
14	6	3	1	4	0	1	2	0	0	5
15	0	1	0	0	0	0	0	0	0	0

["😬 love", "🤩 omg", "❤️ love", "😬 thank", "🤪 lmao", "😬 love", "🤍 bts", "🤩 thank", "❤️ thank", "👑 brunette", "😬 lol", "🤪 lol", "🤩 thank", "❤️ love", "😬 love", "🤍 thank", "🔥 fire", "😬 thank", "🙏 thank", "love 😬", "omg 🤩", "love ❤️", "thank 😬", "lmao 🤪", "love 😬", "bts 🤍", "thank 🤩", "thank ❤️", "brunette 👑", "lol 😬", "lol 🤪", "thank 🤩", "love ❤️", "love 😬", "thank 🤍", "fire 🔥", "thank 🤩", "thank 🙏"]

NAIWNY KLASYFIKATOR BAYESA

$$P(A|B) = \frac{P(A \cap B)}{P(B)} = \frac{P(B|A) \cdot P(A)}{P(B)}$$

Ocena podobieństwa uzyskanych klasyfikacji – macierz błędów

FAKTYCZNE DANE

PROGNOZA

	Pozytywna	Negatywna
Pozytywna	TP	FP
Negatywna	FN	TN

Precyzja – wskazuje udział liczby prawidłowo pozytywnie przydzielonych względem zaprognozowanych pozytywnie (TP/ TP+FP))

Czułość – wskazuje jaka część obserwacji jest prawidłowo pozytywnie przydzielona względem wszystkich pozytywnych tj. (TP/(TP+FN))

Dokładność – wskazuje jaka część obserwacji została zaklasyfikowana poprawnie (TP+TN)/N

F1-score – średnia harmoniczna z precyzji oraz czułości.(2TP/2TP+FP+FN)

Porównanie wyników klasyfikacji dla różnych zbiorów danych

Odsetek poprawnie zaklasyfikowanych użytkowników

Model	Klasa	Słowa	Emoji	Asocjacje	Emotikony	Wszystko
Naiwny klasyfikator Bayesa	Pokolenie Z	91%	77%	16%	39%	70%
	Reszta grup wiekowych	59%	63%	80%	72%	65%

Porównanie wyników klasyfikacji dla różnych zbiorów danych

Model	Słowa	Emoji	Asocjacje	Emotikony	Wszystko
Dokładność	72%	69%	54%	58%	67%
Precyzja	74%	67%	52%	56%	65%
Czułość	70%	74%	94%	85%	76%
F1-score	72%	71%	67%	67%	70%

Wnioski końcowe

- W identyfikacji pokolenia Z najwyższą precyzją (74%) i dokładnością (72%) charakteryzuje się model wykorzystujący jedynie słowa.
- Uwzględnienie w zbiorze danych innych elementów tekstu nie daje poprawy klasyfikacji.
- W identyfikacji pozostałych pokoleń najskuteczniejszy (80%) jest zbiór danych obejmujący asocjacje emoji - słowo.

Dalsze kierunki badań

- Zastosowanie innych modeli klasyfikacyjnych w tym regresji logistycznej, drzew klasyfikacyjnych i modeli AI.
- Walidacja modeli na nowych danych.
- Próba stworzenia modelu pozwalającego na klasyfikację użytkownika do pokolenia Z na podstawie pojedynczego wpisu.
- Próba stworzenia modelu pozwalającego na identyfikację wieku autora wpisu na podstawie wiadomości tekstowej napisanej w innym języku niż angielski.



Majkowska Agata

Migdał-Najman Kamila

Najman Krzysztof

Raca Katarzyna

Identyfikacja grup wiekowych użytkowników portalu X

www.ug.edu.pl

