



SKAZANI NA PASTE BIN

NLP KONTRA ŚMIETNIK INTERNETU

MARTA ZAJKOWSKA & MICHAŁ JADCZUK



НАГК



УЕН



Zespół



Marta Zajkowska

UX/UI Designer



Maciej Sawicki

Full Stack Developer



Michał Jadczyk

Data Scientist

EXATEL

“

Amount of available data increases every day. Finding useful information – *information we can work with* – is possible, but quite tedious and difficult without spending lots of resources.

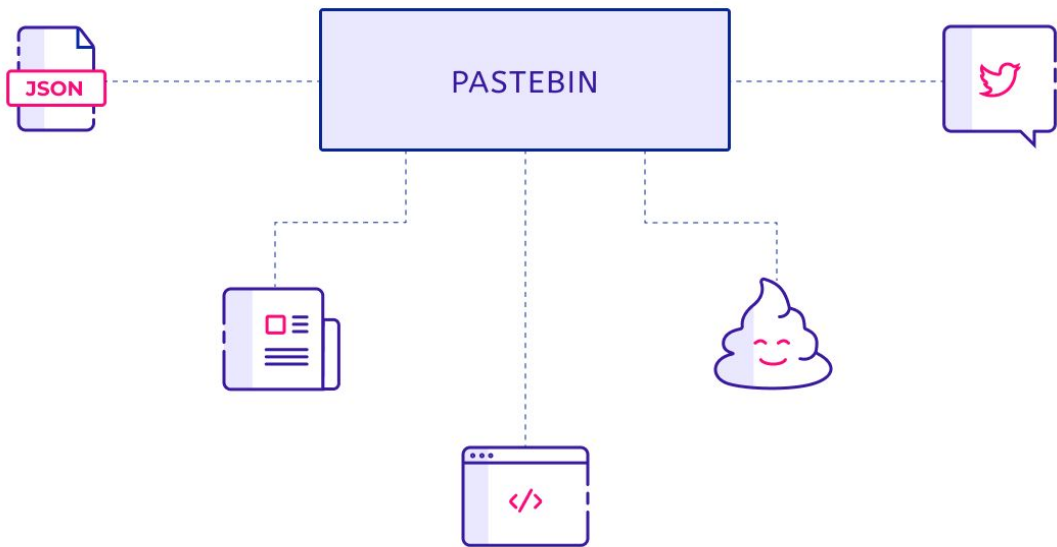
Can you write the software which, given a lot of uncategorized data (say – *tweets, blog posts, application logs, pcaps*), groups them by **the similarity of the discussed topic**?

The goal is that certain groups can be ignored as uninteresting and other browsed manually. Could you **score and sort** the information in a useful way?

If you've considered trying your skills with **unsupervised machine learning**, this challenge is for you!

”

Problem



Przypadki użycia



Użytkownik nie zna zbioru danych
i **nie wie**, czego szuka

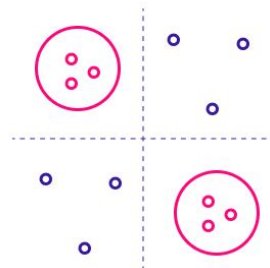


Użytkownik nie zna zbioru danych,
ale **wie**, czego szuka

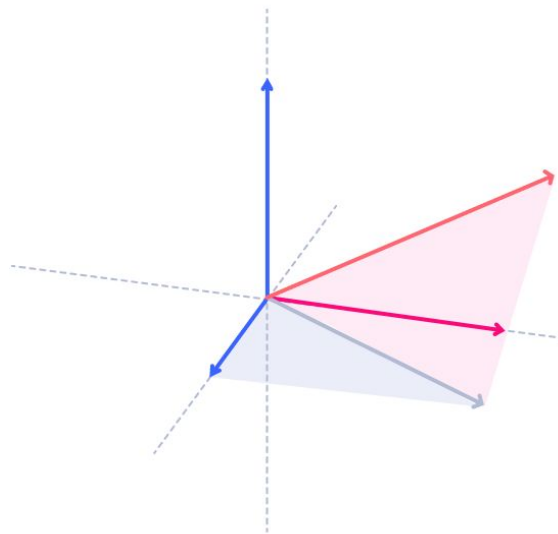
Funkcjonalności



**Wyszukiwanie
kontekstowe**



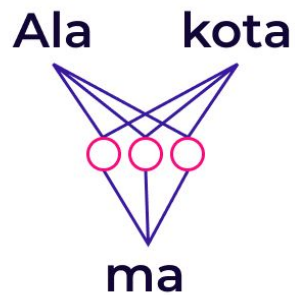
**Grupowanie
danych**



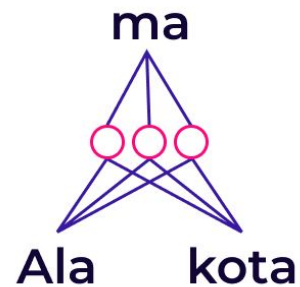
Przykład

FRAZA	ALA	MIEĆ	KOT	SZCZĘŚCIE
Ala ma kota	1	1	1	0
Kot ma Alę	1	1	1	0
Mamy kota, mamy szczęście	0	2	1	1

word2vec

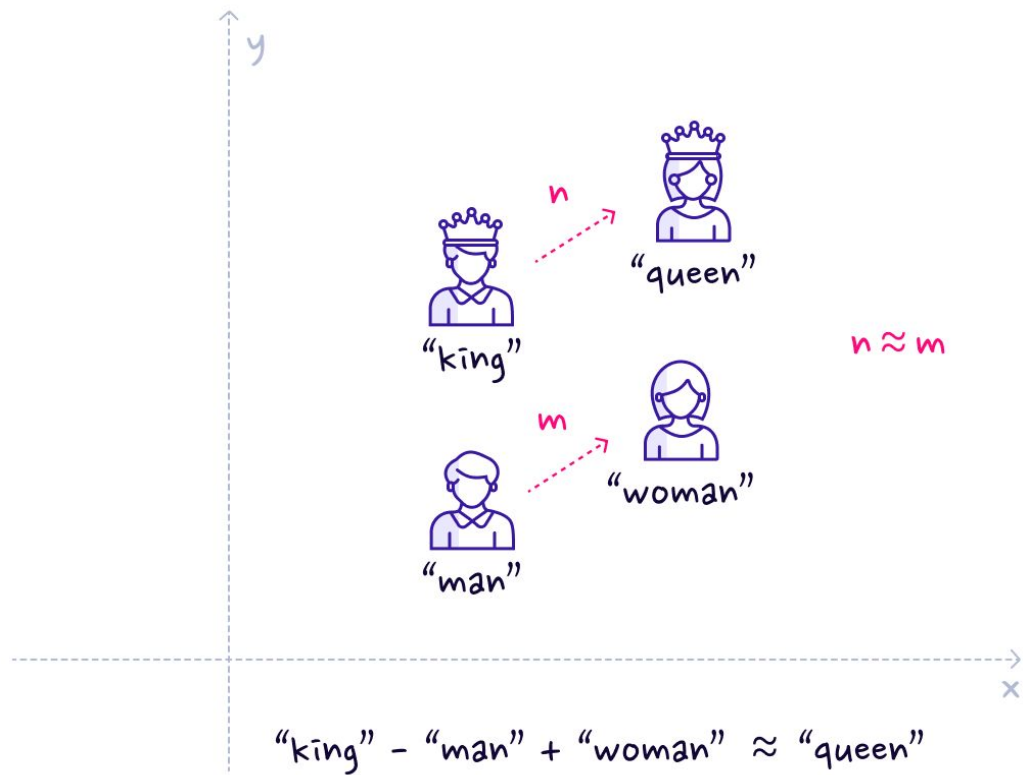


Continuous bag-of-words

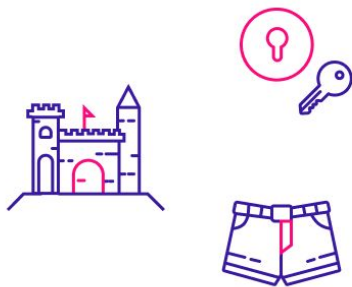


Skip-gram

word2vec



word2vec: wady



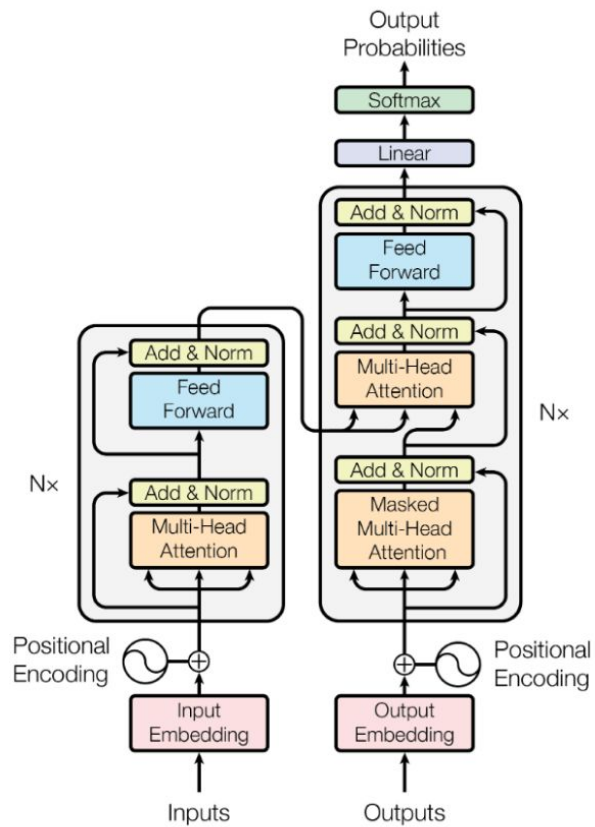
Wektory niezależne
od kontekstu

→
word

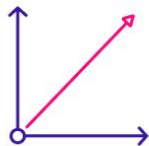
Wektory
na poziomie słów

**ATTENTION IS ALL
YOU NEED!**

Architektura



Transformer: zalety



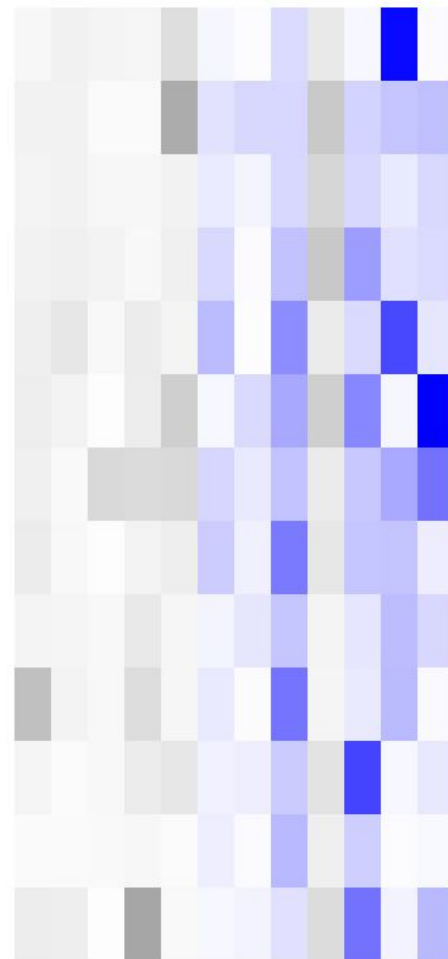
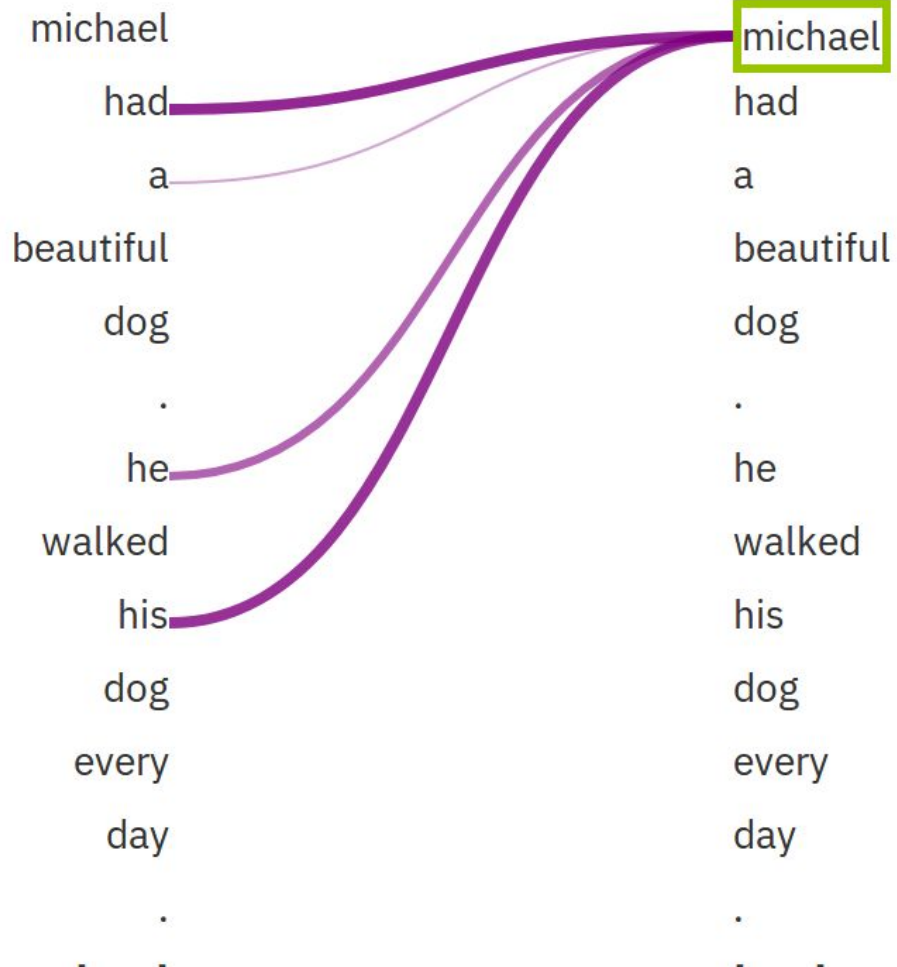
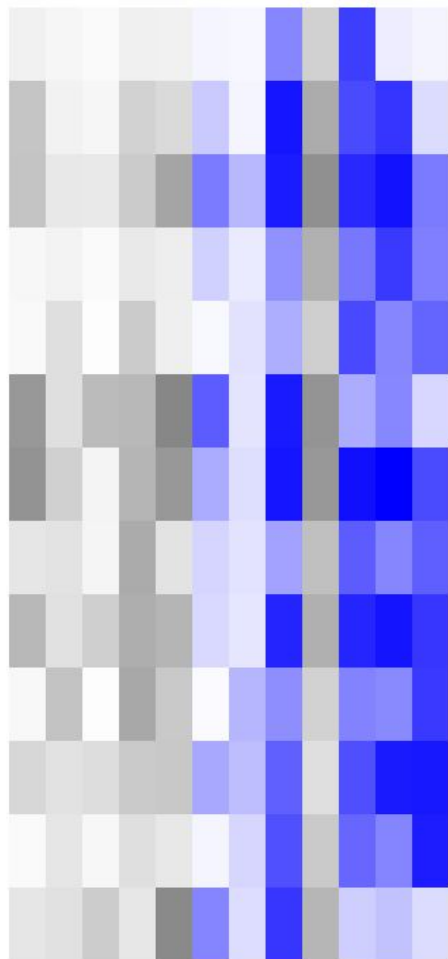
Skalowalność



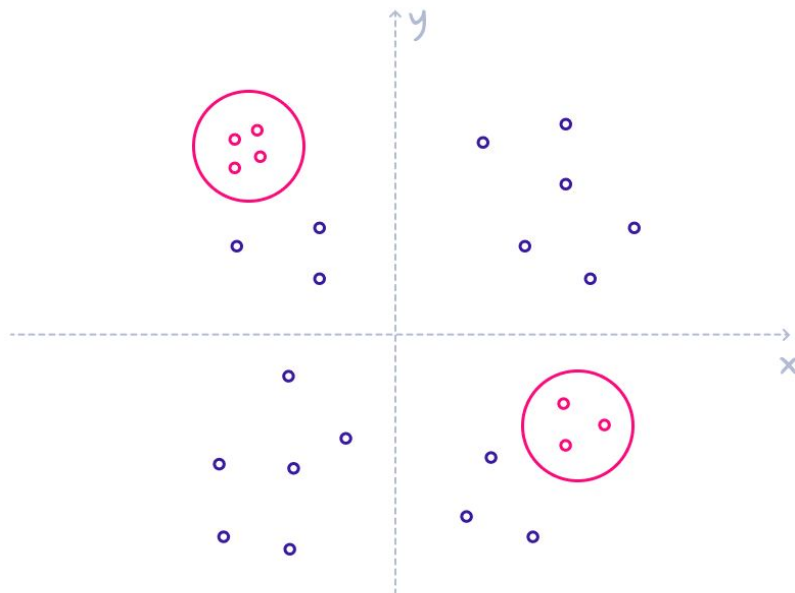
Analiza w jednym kroku



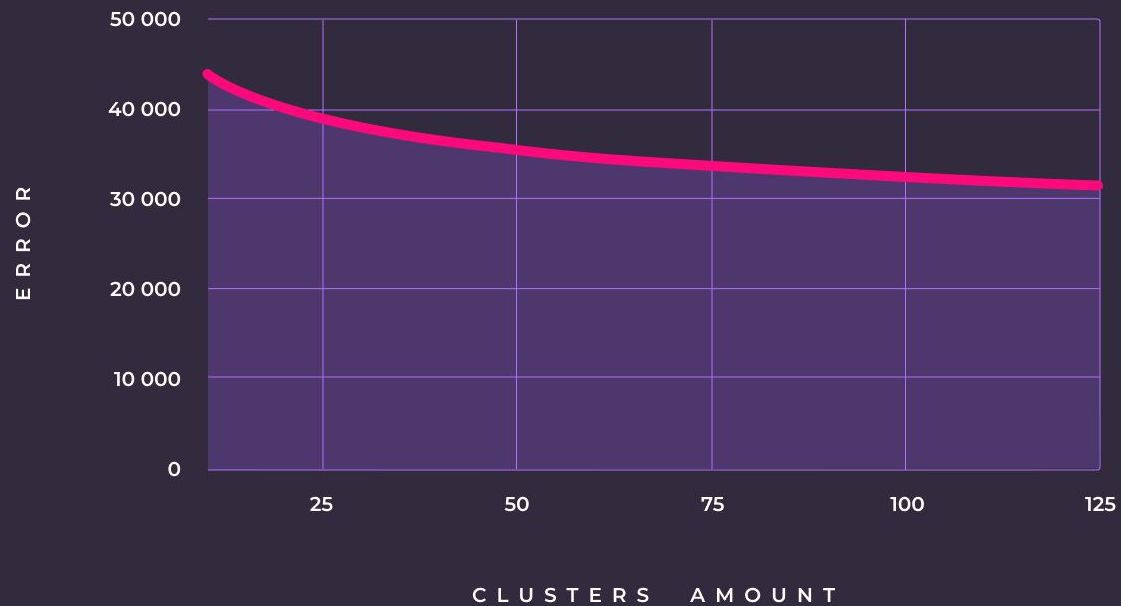
Warstwy uwagi



Grupowanie



Grupowanie



Categories 36



■ All

- ☐ Cars
- ☐ Chat logs
- ☐ Code (Android)
- ☐ Code (C++)
- ☐ Code (CSS, HTML)
- ☐ Code (data science)
- ☐ Code (math, plots)
- ☐ Code (ORM, SQL)
- ☐ Code (PHP)
- ☐ Code (Roblox)
- ☐ Code (snippets, configs)
- ☐ Code (SQL)
- ☐ Courses, Job Postings
- ☐ Crash logs
- ☐ Credit cards
- ☐ Donald Trump

Documents > All

1,427,833 results



60e/60e2e857f8a38eab61cf7afcd1972e1f28414675



```
1 <!DOCTYPE html>
2 <html>
3 <head>
4 <title>
5 <meta charset="utf-8">
6 <style>
7   #dnd-area {
8     position: relative;
```



5ce/5ce28798e1a51698acdd33d6cdb2e14a6915c47a



```
1 Hey Guys! :
2 Here is a amazing scripter from youtube!
3
4 Link: https://youtube.com/c/RawNetworksYouTube
```

Jak żyć na hackathonie?



Interdyscyplinarny
zespół



Gotowy produkt



Dobra prezentacja



Współpraca
z mentorami

 pystok.slack.com

#hackathons

HACK
YEAH

DZIĘKUJEMY ZA UWAGĘ

EXATEL

10 000 PLN

DATA CATEGORIZING SOFTWARE

HACK
YEAH

2019/9/15 16:06