

Brakuje konkretnej literatury w języku polskim

- D. Spinczyk, M. Dzieciątko, *Text mining. Metody, narzędzia, zastosowania*, PWN (2016),

Polecam również poniższe pozycje w jęz. angielskim:

- Ch. Aggarwal, Ch-X Zhai, C. O'Neil *Mining Text Data*, Springer (2012).
- D. Robinson, J. Silge, *Text Mining with R*, O'Reilly (2017)

## Text mining wg Wikipedii (ang.)

**Text mining**, also referred to as **text data mining**, roughly equivalent to **text analytics**, is the process of deriving high-quality information from text.

### Text mining wg Wikipedii (ang.)

**Text mining**, also referred to as **text data mining**, roughly equivalent to **text analytics**, is the process of deriving high-quality information from text.

### Text mining wg Wikipedii (pol.)

Text mining (eksploracja tekstu) — ogólna nazwa metod eksploracji danych służących do wydobywania danych z tekstu i ich późniejszej obróbki.

## Text mining wg Wikipedii (ang.)

**Text mining**, also referred to as **text data mining**, roughly equivalent to **text analytics**, is the process of deriving high-quality information from text.

## Text mining wg Wikipedii (pol.)

Text mining (eksploracja tekstu) — ogólna nazwa metod eksploracji danych służących do wydobywania danych z tekstu i ich późniejszej obróbki.

## Text mining wg Marti Hearst

Another way to view text data mining is as a process of exploratory data analysis that leads to heretofore unknown information, or to answers for questions for which the answer is not currently known.



[Grafika pobrana z:  
<https://www.ischool.berkeley.edu>]

## Po co text mining?

Z drugiej strony, warto zadać sobie pytanie **po co potrzebujemy eksploracji tekstu?** lub **jakie jest zadanie eksploracji tekstu?**. Ogólną odpowiedzią jest oczywiście: **aby (w automatyczny sposób) zrozumieć zawartość danego tekstu...**

## Po co text mining?

Z drugiej strony, warto zadać sobie pytanie **po co potrzebujemy eksploracji tekstu?** lub **jakie jest zadanie eksploracji tekstu?**. Ogólną odpowiedzią jest oczywiście: **aby (w automatyczny sposób) zrozumieć zawartość danego tekstu...**

## Po co text mining?

... niestety to założenie wydaje się być zbyt trudne. Dlatego skupiamy się raczej pomniejszych zdaniach.

## Dlaczego analiza tekstu jest **trudna**?

Cieężko jest oddać abstrakcyjne pojęcia w postaci innych, dobrze zdefiniowanych pojęć



## Dlaczego analiza tekstu jest **trudna**?

Cieężko jest oddać abstrakcyjne pojęcia w postaci innych, dobrze zdefiniowanych pojęć



**Time** **flies** like an  
arrow.

Niezliczone kombinacje subtelnych i abstrakcyjnych relacji pomiędzy pojeciami



## Dlaczego analiza tekstu jest **trudna**?

Cieężko jest oddać abstrakcyjne pojęcia w postaci innych, dobrze zdefiniowanych pojęć



**Time** **flies** like an arrow.



Niezliczone kombinacje subtelnych i abstrakcyjnych relacji pomiędzy pojeciami

Wiele sposobów opisywania tych samych pojęć

## Dlaczego analiza tekstu jest **trudna**?

Cieężko jest oddać abstrakcyjne pojęcia w postaci innych, dobrze zdefiniowanych pojęć



**Time** **flies** like an arrow.



Niezliczone kombinacje subtelnych i abstrakcyjnych relacji pomiędzy pojeciami

Wiele sposobów opisywania tych samych pojęć

Wysoka wymiarowość problemu



## Dlaczego analiza tekstu jest **trudna**?

Cieężko jest oddać abstrakcyjne pojęcia w postaci innych, dobrze zdefiniowanych pojęć



Time flies like an arrow.



## Niezliczone kombinacje subtelnych i abstrakcyjnych relacji pomiędzy pojęciami

## Wiele sposobów opisywania tych samych pojęć

## Wysoka wymiarowość problemu



Type of Feature	Genre	Plot	Character	Setting	Theme
Activity book?	✓	✓	✓		
Reference book?		✓	✓		
Text book?				✓	
Graphic novel?					✓
Classical novel?	✓				

Bardzo wiele cech (features)

Dlaczego analiza tekstu może być **łatwa**?

## Dlaczego analiza tekstu może być **łatwa**?

W tekście zwykle jest spora ilość nadmiarowych lub powtarzających się informacji.

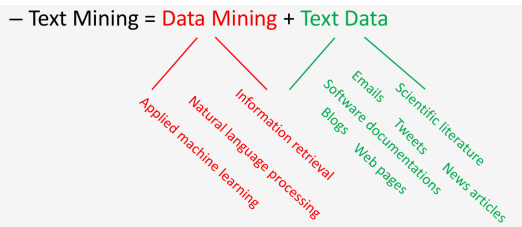
## Dlaczego analiza tekstu może być **łatwa**?

W tekście zwykle jest spora ilość nadmiarowych lub powtarzających się informacji.

W zasadzie większość prostych algorytmów może osiągnąć całkiem dobre wyniki przy wykonywaniu w następujących nieskomplikowanych zadań:

- wydobać "istotne" wyrażenia,
- znaleźć istotnie powiązane słowa,
- stwórz pewnego rodzaju podsumowanie dokumentów

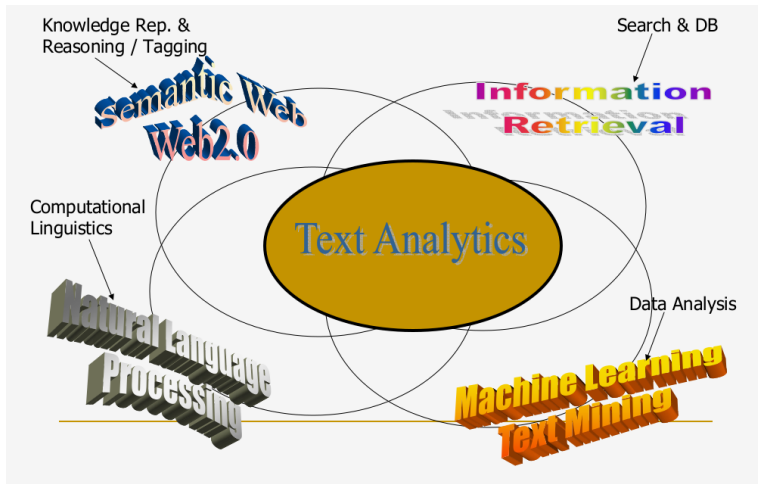
Można również próbować zilustrować powiązania pomiędzy eksploracją tekstu a innymi dziedzinami:



	Finding Patterns	Finding "Nuggets"	
		Novel	Non-Novel
Non-textual data	General data-mining	Exploratory analysis	Database queries
Textual data	Comp Ling		Information retrieval

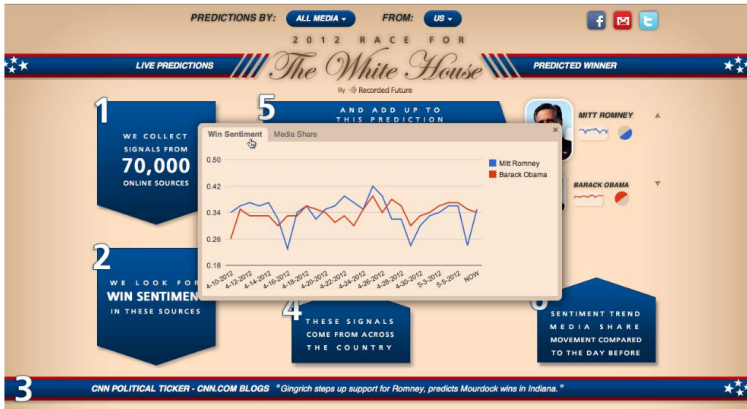
**Text Mining**

Można również próbować zilustrować powiązania pomiędzy eksploracją tekstu a innymi dziedzinami:





## Przykłady: analiza sentymentu – wybory



## Przykłady: podsumowywanie dokumentów



## Przykłady: systemy rekomendujące

### FOREIGN SUGGESTIONS (about 104) [See all >](#)



#### Tell No One

Because you enjoyed:  
Memento  
Syriana  
Children of Men



#### Let the Right One In

Because you enjoyed:  
Seven Samurai  
This Is Spinal Tap  
The Big Lebowski



#### I've Loved You So Long

Because you enjoyed:  
The Queen  
Syriana  
Good Night, and Good Luck



#### Downfall

Because you enjoyed:  
Das Boot  
The Killing Fields  
Seven Samurai



### DRAMA SUGGESTIONS (about 82) [See all >](#)



#### The Wrestler

Because you enjoyed:  
Sin City  
Reservoir Dogs  
The Big Lebowski



#### The Visitor

Because you enjoyed:  
Gandhi  
The Motorcycle Diaries  
The Queen



#### Brick

Because you enjoyed:  
The Big Lebowski  
Rushmore  
Fight Club

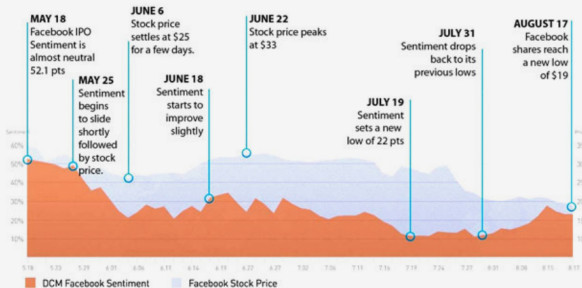


#### The Pianist

Because you enjoyed:  
Amadeus  
The Killing Fields  
Empire of the Sun



## Przykłady: analiza tekstu w serwisach finansowych

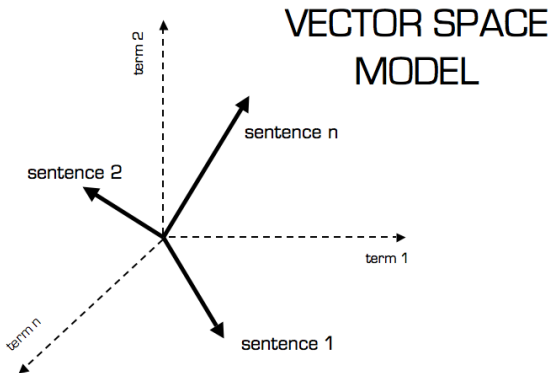




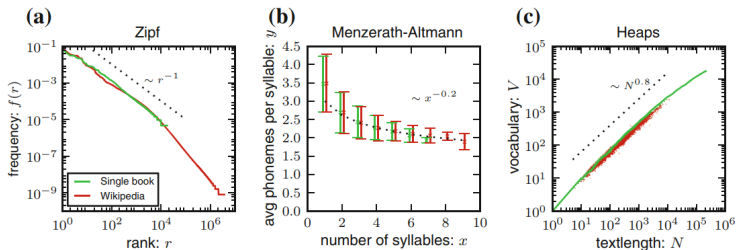
## Ogólny plan wykładu

- 1 reprezentacja tekstu
- 2 prawo Zipfa
- 3 przetwarzanie języka naturalnego (NLP)
- 4 analiza sentymentu
- 5 topic modeling
- 6 analiza mediów społecznościowych

## 2 Reprezentacja tekstu...



### 3 Prawo Zipfa i pokrewne...



[Altmann, Gerlach, Statistical laws in Linguistics, Creativity and Universality in Language, Springer (2017)]



## 4 przetwarzanie języka naturalnego (NLP)

### Part of speech:

NP NP RB VBD IN NP NP CC PRF VBZ RB VBG PRP IN PRP .  
Mrs. Clinton previously worked for Mr. Obama, but she is now distancing herself from him .

### Named entity recognition:

Person Date Person Date  
Mrs. Clinton previously worked for Mr. Obama, but she is now distancing herself from him.

### Co-reference:

Mention Ment M Mention M  
Mrs. Clinton previously worked for Mr. Obama, but she is now distancing herself from him.

### Basic dependencies:

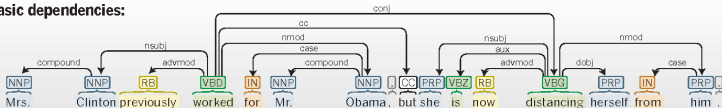
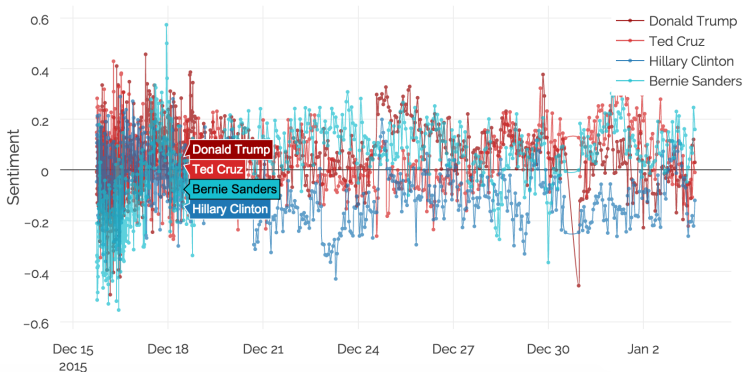


Fig. 1. Many language technology tools start by doing linguistic structure analysis. Here we show output from Stanford CoreNLP. As shown from top to

[Hirschberg, Manning, Advances in natural language processing, Science 349, 261 (2015)]

## 5 Analiza sentymentu

### How Twitter Feels About the 2016 Election Candidates



## 5 Analiza sentymentu: klasyfikatory słownikowe vs uczenie pod nadzorem

### Dictionary-based Approach

Create lists of **positive/negative** words (phrases).

Negative

suck  
terrible  
awful  
unwatchable  
hideous

Positive

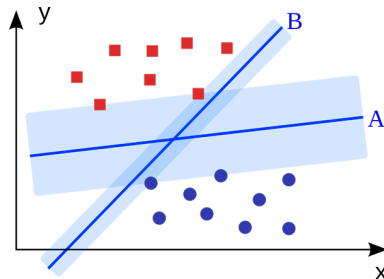
dazzling  
brilliant  
phenomenal  
excellent  
fantastic

Sentiment = |Positive words| - |Negative words|

Around 65% accuracy!

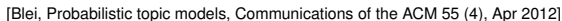


5

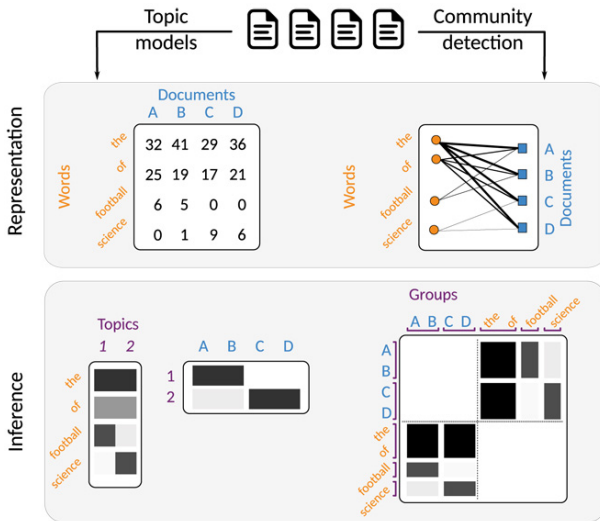


[<https://www.slideshare.net/jchoi7s/cs571-sentiment-analysis>]

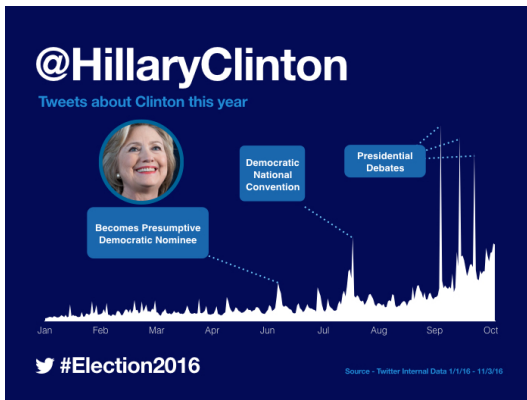
[<https://medium.com/nlpython/sentiment-analysis-analysis-part-2-support-vector-machines-31f78baeee09>]



## 6 Topic modelling

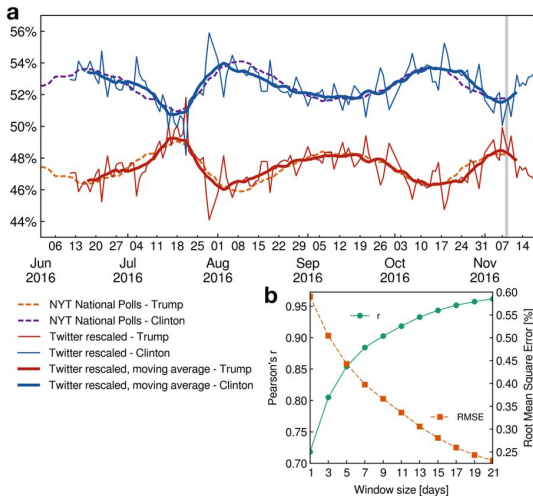


## 7 Analiza mediów społecznościowych



[Bovet, Morone, Makse, Validation of Twitter opinion trends with national polling aggregates: Hillary Clinton vs Donald Trump, Scientific Reports (2018)]

## 7 Analiza mediów społecznościowych



[Bovet, Morone, Makse, Validation of Twitter opinion trends with national polling aggregates: Hillary Clinton vs Donald Trump, Scientific Reports (2018)]