

# STYLOWE KALKULACJE

(ZAJAWKA O ATRYBUCJI AUTORSTWA)

ATRYBUCJA AUTORSTWA

- **ustalenie autora** dzieła anonimowego lub o niepewnym autorstwie na podstawie właściwości samego tekstu
- w oparciu o wyniki **analizy** treści, **cech językowo-stylistycznych tekstu** oraz wskaźników historyczno-bibliograficznych, edytorskich i drukarskich
- zakłada istnienie **stylistycznego odcisku palca** - nieuświadomionych nawyków stylistycznych, do badania których wykorzystuje **stylometrię**

# STYLOMETRIA

- bazuje na **badaniach ilościowych** (wspomaganych komputerowo, chociaż przez długi czas szło to ręcznie)
- pomija cechy stylistyczne w ujęciu tradycyjnym - metafory, epitety, stylizację, ironię...
- **analiza porównawcza** - konieczny i bardzo istotny kontekst

- rozkład części mowy
- rozkład form czasownikowych
- połączenia międzywyrazowe
- zbitki literowe
- bogactwo słownictwa
- typowe konstrukcje zdaniowe
- relacje składniowe
- frekwencja wyrazów

DELTA BURROWSA (NIEPEŁNA)

T1: “Cupcake ipsum dolor sit amet.”

T2: “Cupcake ipsum dolor sit amet.”

T3: “Muffin ipsum dolor sit amet.”

T4: “Muffin cupcake cupcake sit amet.”



t\_1 = { cupcake: 1, ipsum: 1, dolor: 1, sit: 1, amet: 1 }

t\_2 = { cupcake: 1, ipsum: 1, dolor: 1, sit: 1, amet: 1 }

t\_3 = { muffin: 1, ipsum: 1, dolor: 1, sit: 1, amet: 1 }

t\_4 = { muffin: 1, cupcake: 2, sit: 1, amet: 1 }

```
def diff_between(text_1, text_2)
  diff = text_1.merge(text_2) do |_token, freq_1, freq_2|
    (freq_1-freq_2).abs
  end
  diff.values.sum
end
```

t\_1:t\_2 {cupcake: 0, ipsum: 0, dolor: 0, sit: 0, amet: 0}

t\_1:t\_3 {cupcake: 1, ipsum: 0, dolor: 0, sit: 0, amet: 0, muffin: 1}

t\_1:t\_4 {cupcake: 1, ipsum: 1, dolor: 1, amet: 0, muffin: 1}

diff\_by\_token\_1\_2 = 0

diff\_by\_token\_1\_3 = 2

diff\_by\_token\_1\_4 = 4

“Cupcake ipsum dolor sit amet.”

“Muffin ipsum dolor sit amet.”

“Muffin cupcake cupcake sit amet.”

POKA PRZYKŁAD!

KORPUS

**102 powieści angielskie** - trzy powieści Jane Austen z Project Gutenberg +  
99 powieści (też z Gutenberg Project, ale pobranych z repo:

[https://github.com/computationalstylistics/100\\_english\\_novels](https://github.com/computationalstylistics/100_english_novels))

łącznie: **14 597 120** wyrazów

the, and, to, of, a, i, in, he, was, that, it, her, you, his, she, had, with, as, for, not, s, at,  
but, be, on, is, him, my, have, me, said, all, so, which, by, this, from, they, no, if,  
would, were, there, what, one, t, when, been, an, or, we, who, could, do, out, up, are,  
very, your, will, them, mr, now, more, man, then, little, like, their, into, about, some,  
did, than, know, can, see, well, should, any, time, good, never, only, come, upon, has,  
how, before, down, must, old, much, own, say, over, think, go, am, such, might, other,  
after, made, again, don, thought, himself, our, came, great, mrs, day, two, too, here,  
us, way, back, lady, face, young, eyes, went, sir, life, where, shall, away, long, hand,  
may, looked, first, house, miss, even, its, yet, nothing, still, make, room, last, though,  
father, look, men, just, these

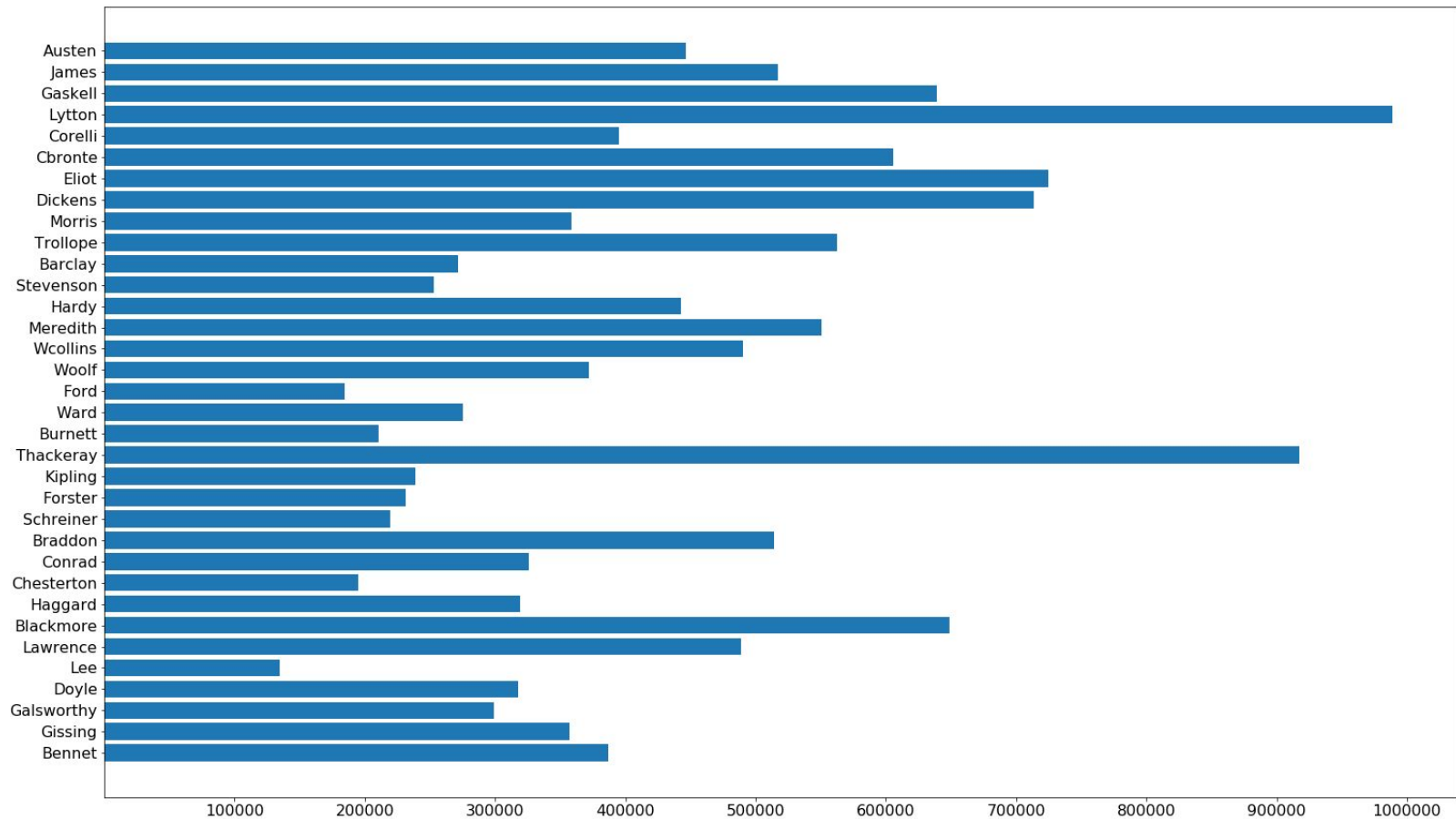


GRUPY

**34 autorów**, po 3 teksty każdego

długość od **134 367** do **988 525** wyrazów w grupie

w każdej grupie autorskiej zliczamy wyrazy, tworząc **frekwencyjny słownik odniesienia** dla danego autora

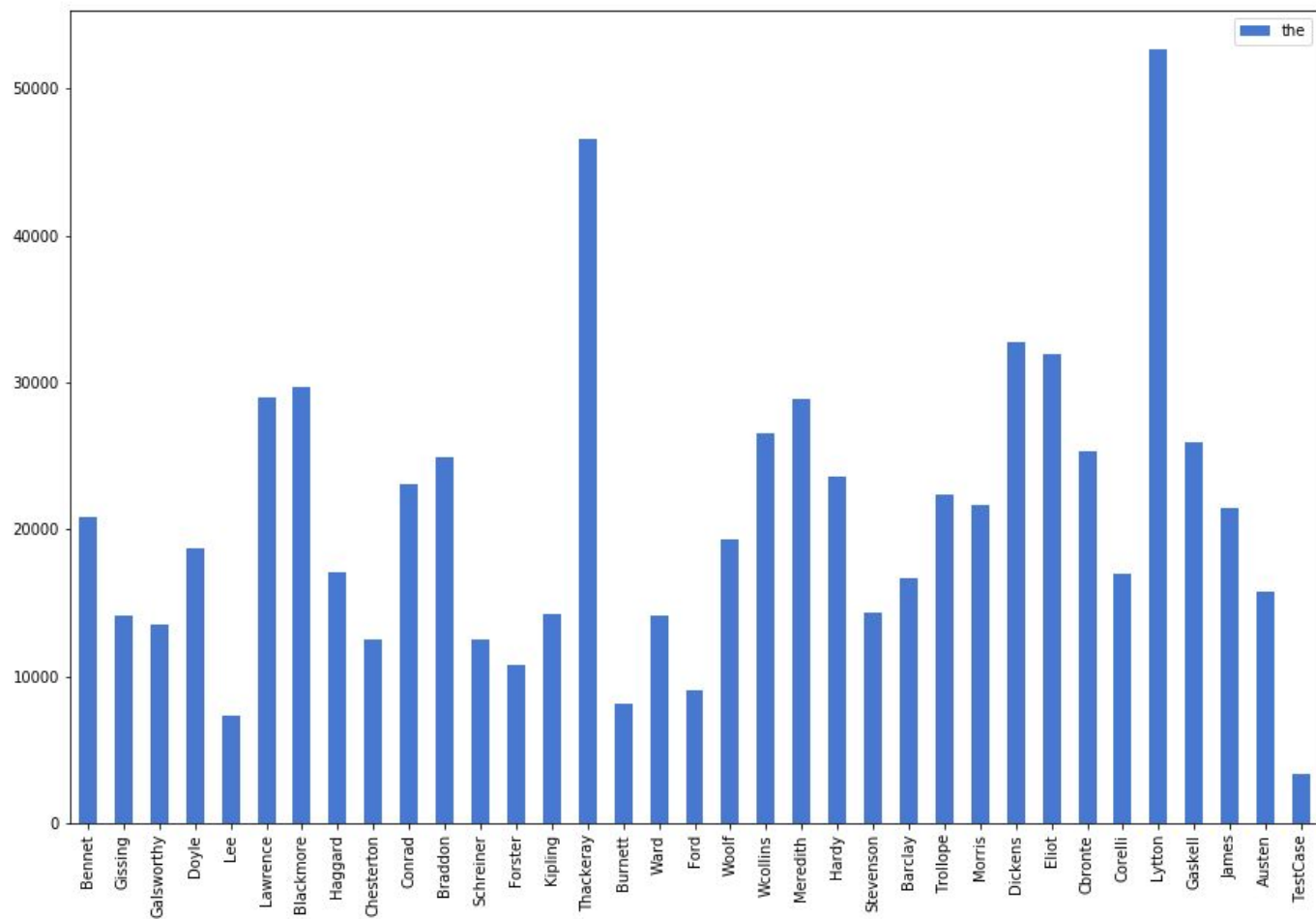


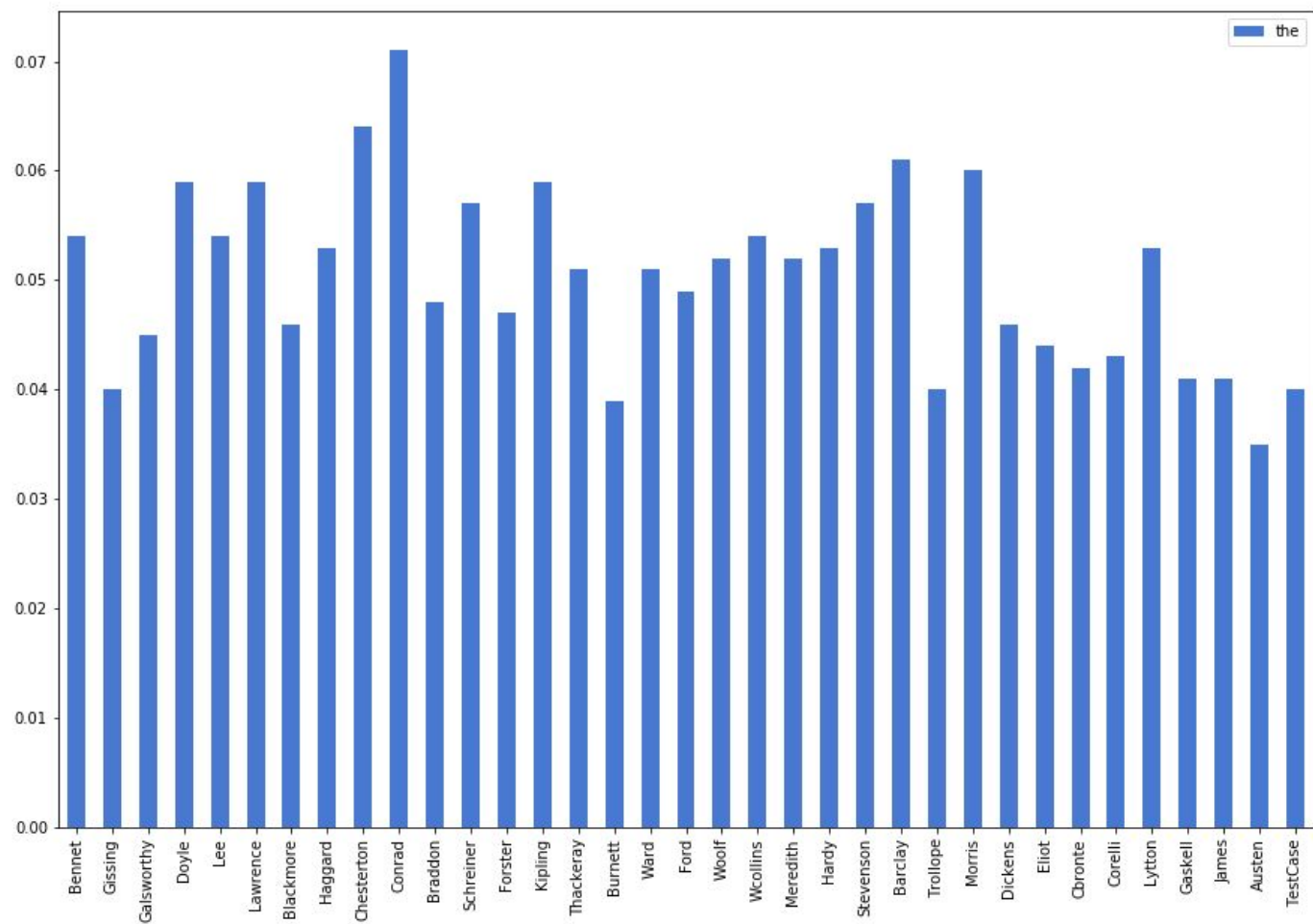
TEKST ZE STRYCHU

będzie go udawać powieść **Jane Austen - “Persuasion”**

**84092** wyrazy, które również zliczamy

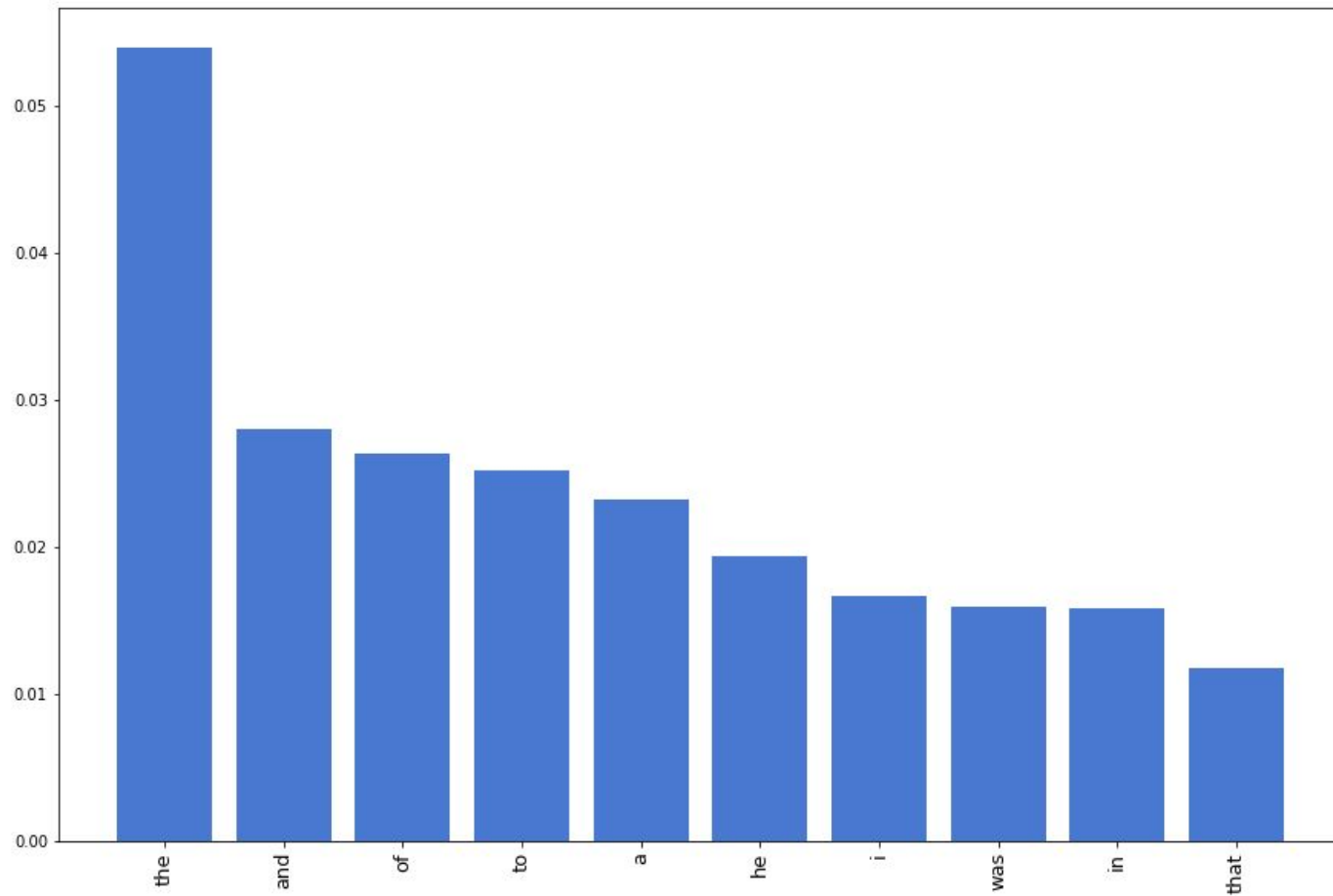
PORÓWNAJMY SOBIE







10 mfw - Bennet



Z-STANDARYZACJA

$$\mathbf{Z\_SCORE} = \frac{\text{rzeczywista wartość zmiennej} - \text{średnia}}{\text{odchylenie standardowe}}$$

**Dla każdego wyrazu W w słownikach frekwencyjnych:**

$$\mathbf{Z\_SCORE} = \frac{\text{częstość względna W} - \text{średnia częstość względna W w korpusie}}{\text{odchylenie standardowe częstości względnej W w korpusie}}$$

```
most_frequent_terms.each do |term|  
  mean = 0  
  groups.each do |group|  
    mean += group.relative_frequency_for(term)  
  end  
  mean /= groups.count  
end
```

```
most_frequent_terms.each do |t|  
  std = 0  
  groups.each do |g|  
    diff = g.relative_frequency_for(t) - corpus.mean_for(t)  
    std += diff**2  
  end  
  std = Math.sqrt(std / (groups.count - 1))  
end
```

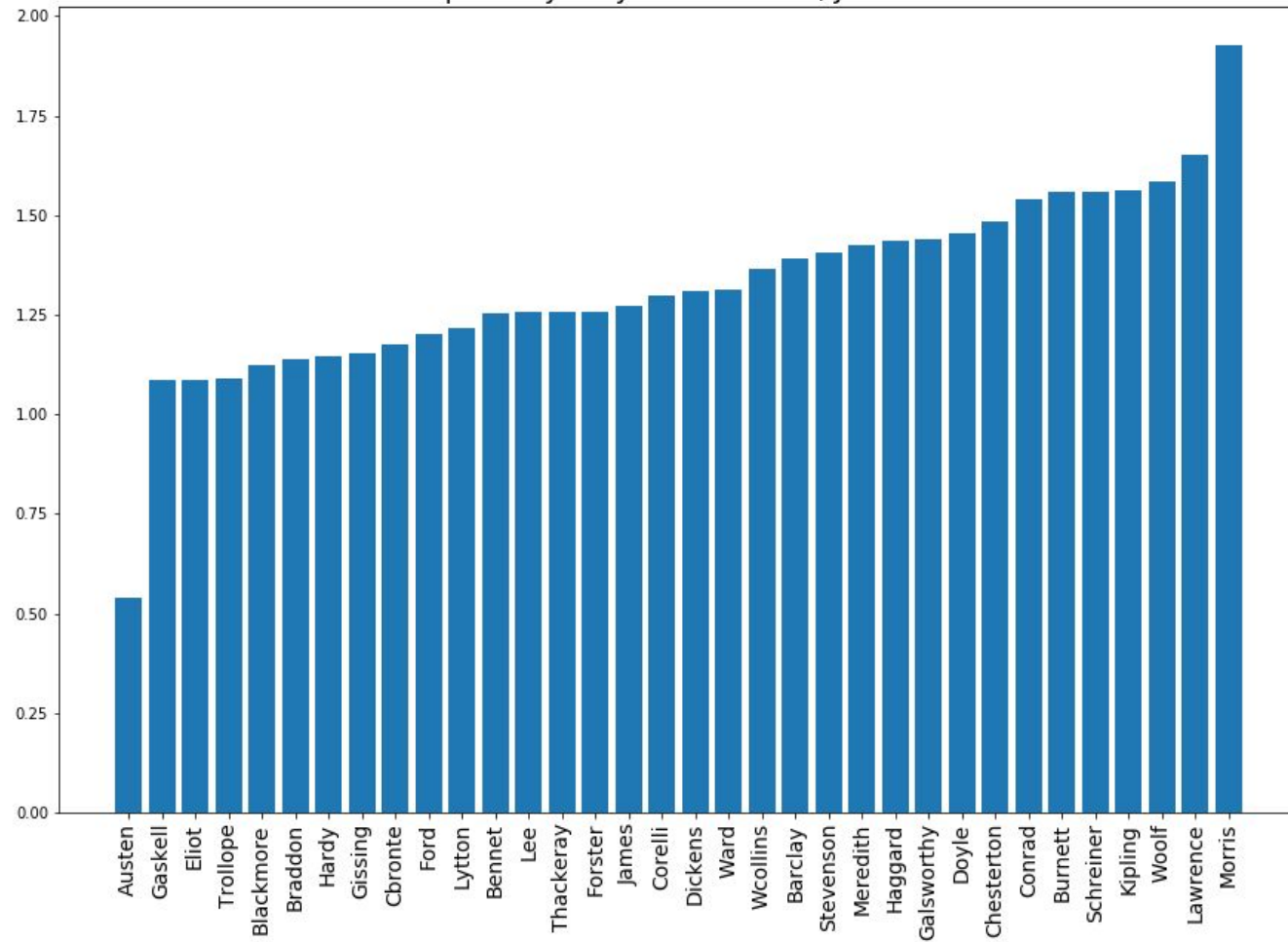
```
def diff_between(z_scores1, z_scores2)
  diff = z_scores1.merge(z_scores2) do |_token, z1, z2|
    (z1 - z2).abs
  end
  diff.values.sum
end
```

```
deltas = groups.map do |group|  
  diff_between(test.z_scores, group.z_scores)  
end
```



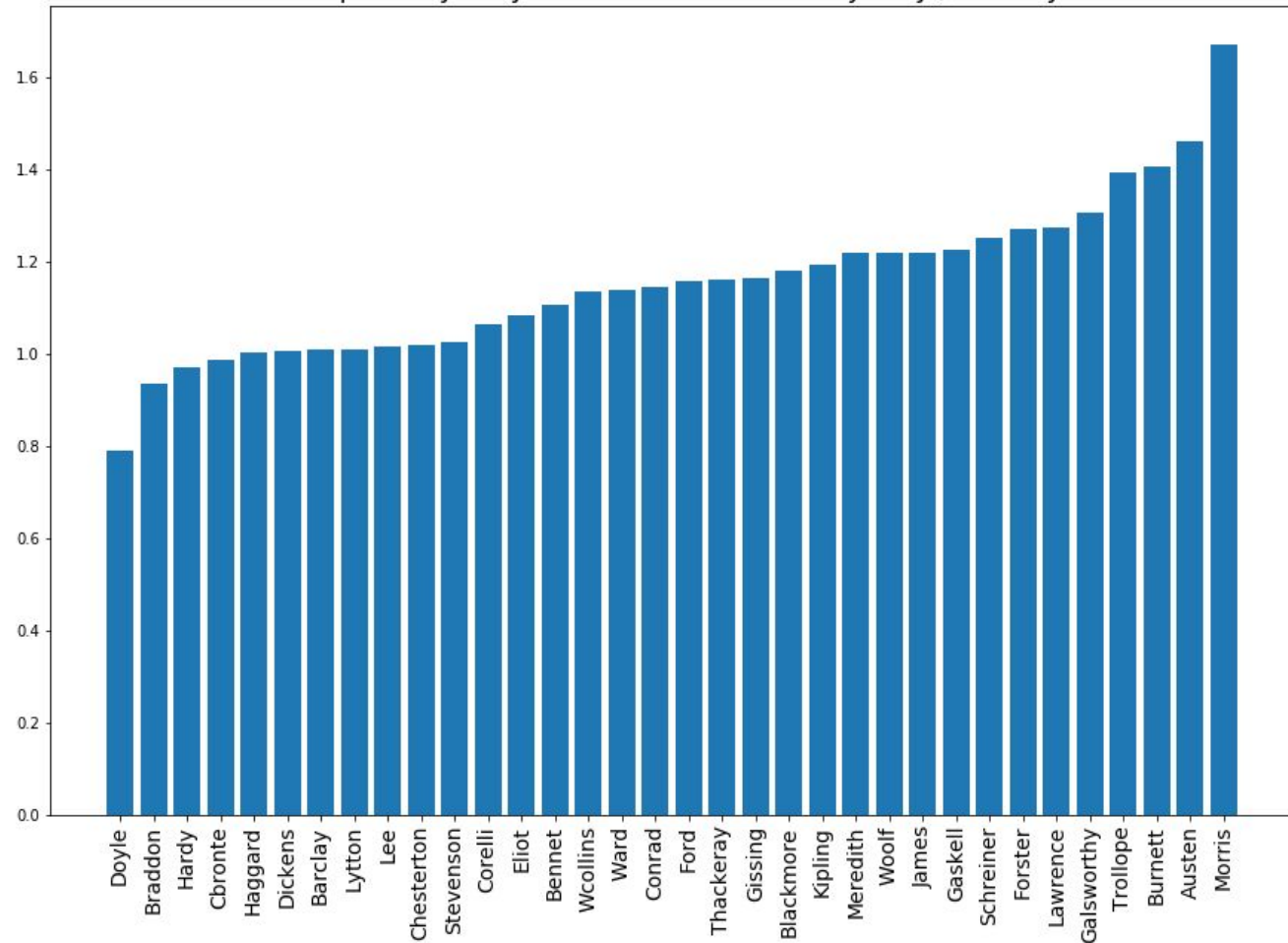
WYNIK?

Tekst porównywany - 'Persuasion', Jane Austen



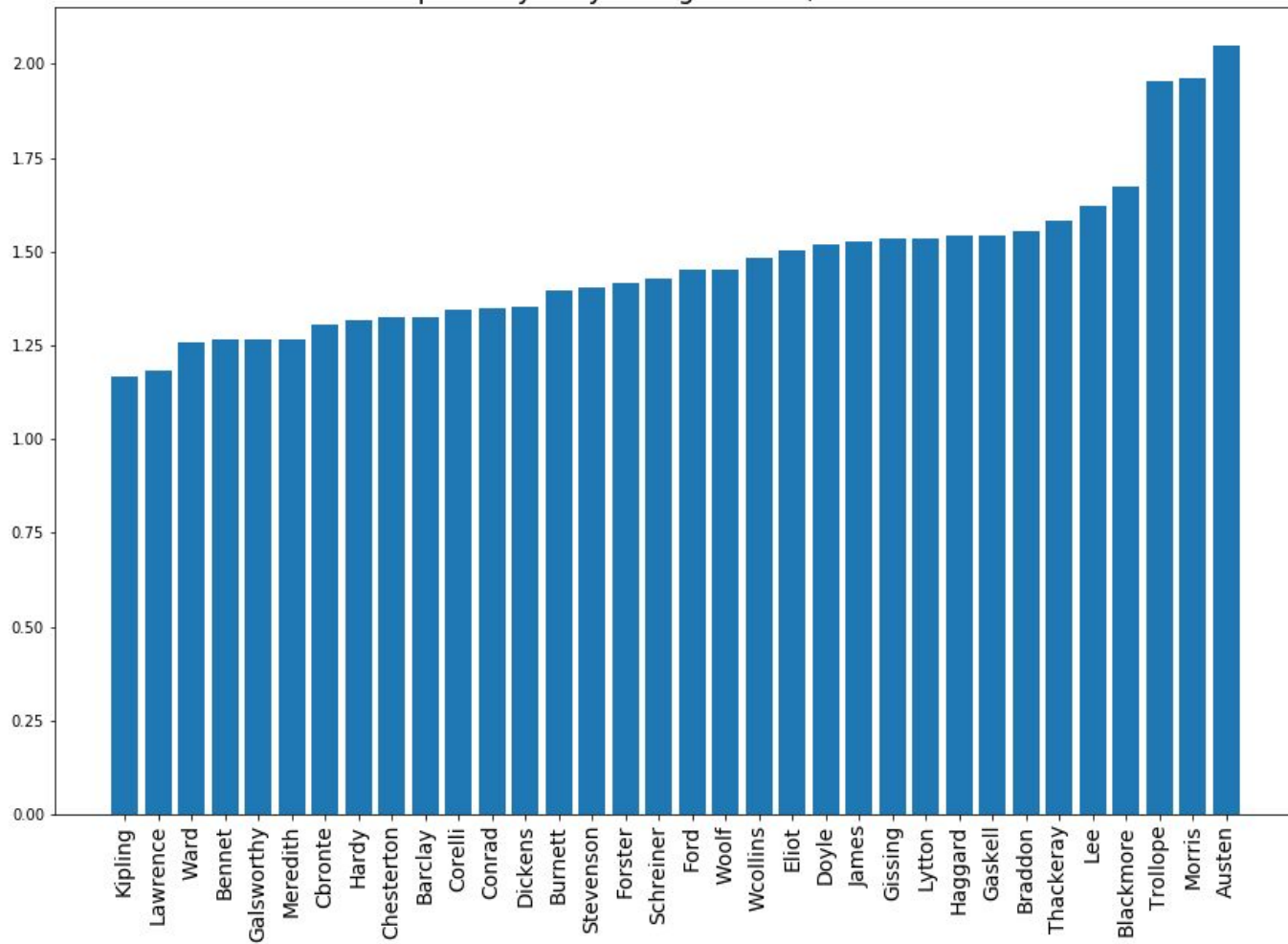
TO MOŻE INNY TEKST?

Tekst porównywany - 'Tales of Terror and Mystery', A.C.Doyle



A JEŚLI AUTOR Z ZUPEŁNIE INNEJ BECZKI?

Tekst porównywany - 'Magis Shifts', Ilona Andrews



TO SKĄD MAM WIEDZIEĆ?

## **Arthur Conan Doyle:**

- "The Sign of the Four",
- "Tales of Terror and Mystery",
- "The White Company",
- "Memoirs of Sherlock Holmes"



## **Jane Austen:**

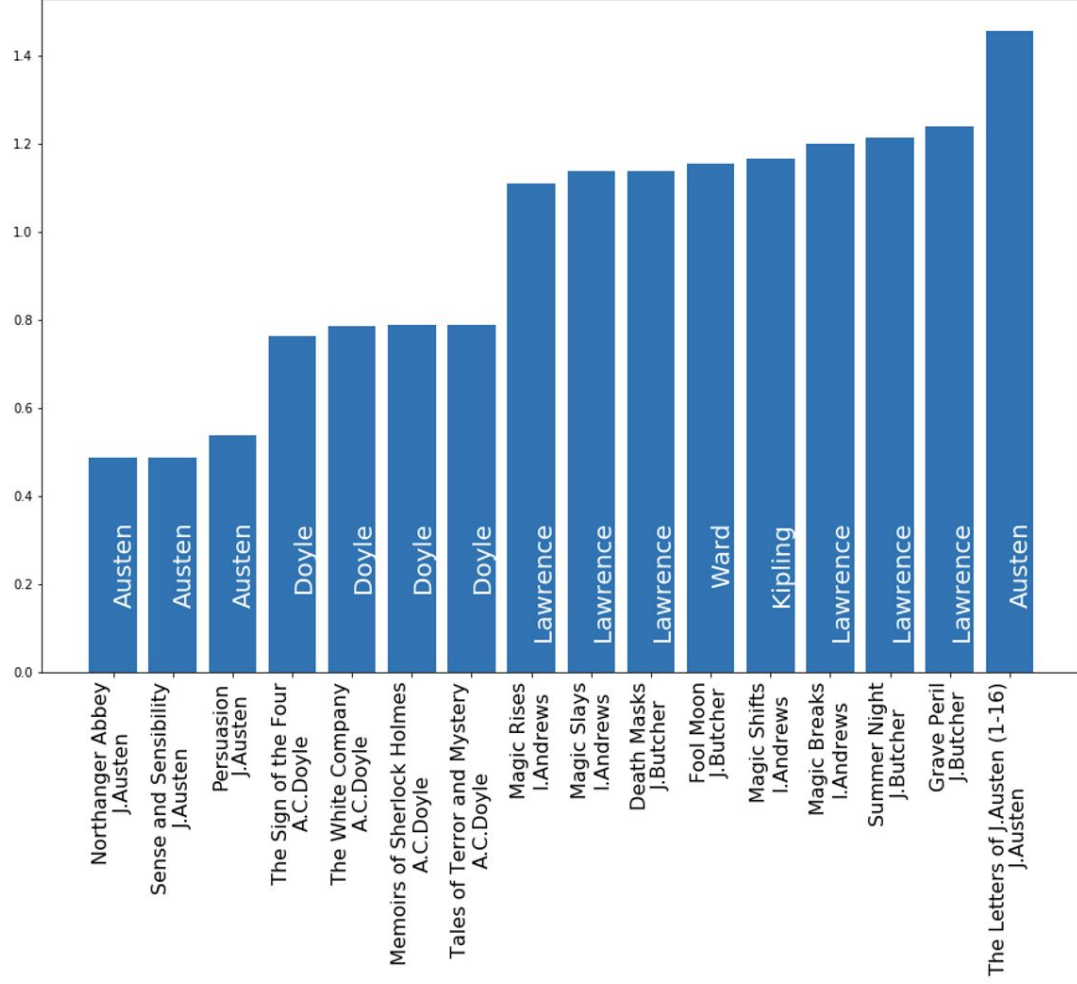
- "The Letters of Jane Austen (1-16)",
- "Northanger Abbey",
- "Persuasion",
- "Sense and Sensibility"

## **Jim Butcher:**

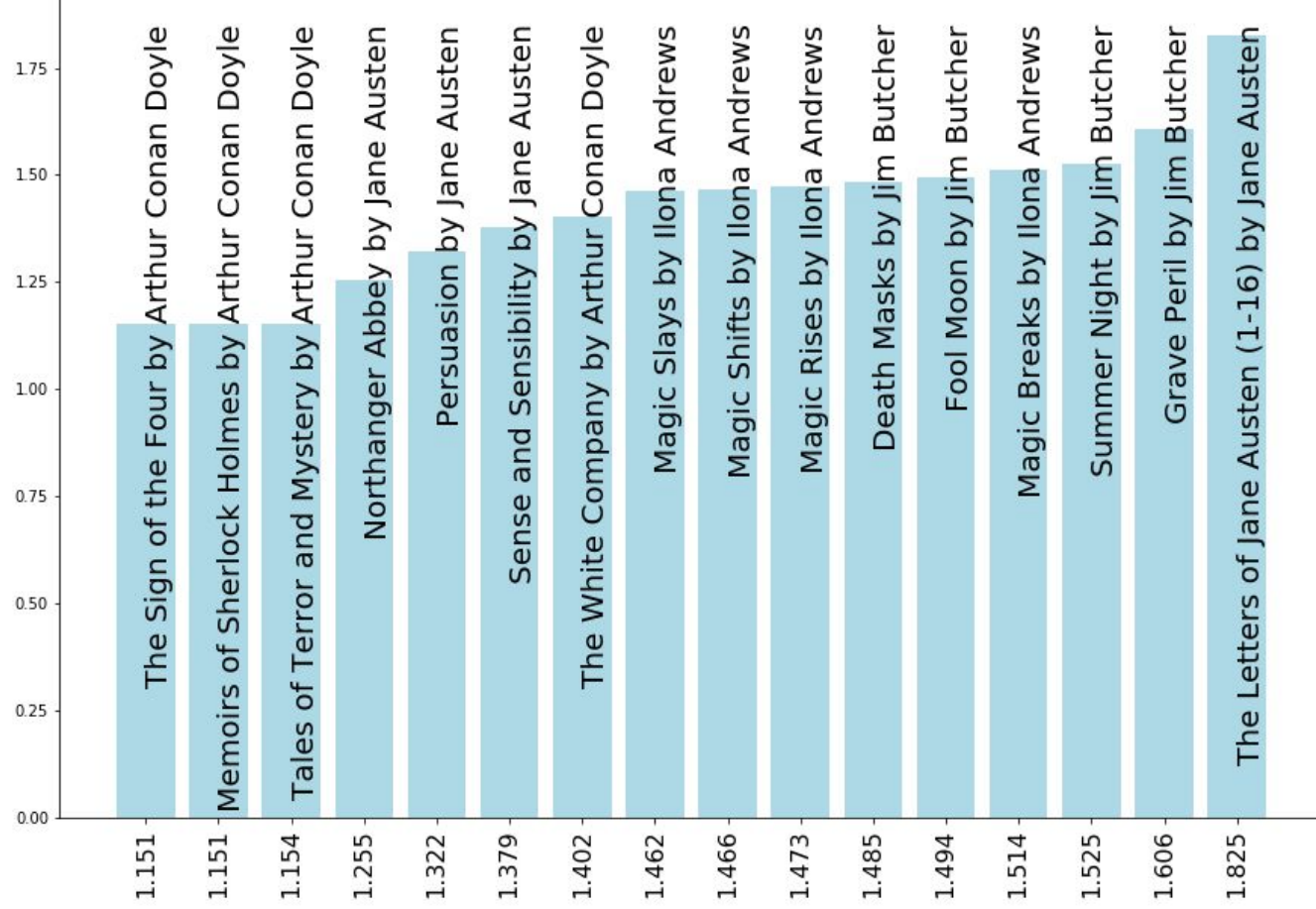
- "Fool Moon",
- "Grave Peril",
- "Summer Night",
- "Death Masks"

## **Ilona Andrews:**

- "Magic Slays",
- "Magic Rises",
- "Magic Breaks",
- "Magic Shifts"



# Srednie delty



O CO CHODZI Z LISTAMI?

PODSUMOWANIE

- **poszlakowe** potwierdzenie autorstwa
- mocno osadzona w **kontekście** - ważny jest odpowiedni dobór korpusu
- wyliczenia są dość proste, pogubić się można tylko w natłoku tekstów, grup, częstości i wyrazów



**formalnie:** Burrows oryginalnie nie bazował na grupach autorskich, tj. porównywał tekst testowy z każdym innym tekstem, a dopiero potem wyniki łączył po autorach, jednak dla uproszczenia zastosowałam swoją wariację grupową

**statystycznie:** delta Burrowsa bazuje na standaryzacji Z (przewidzianej dla rozkładów normalnych) i odległości miejskiej (Manhattan distance), która lepiej działa dla rozkładów Laplace'a - dla standaryzacji Z lepsza byłaby odległość euklidesowa, dla odległości miejskiej lepsza byłaby standaryzacja medianą i rozrzutem

**lista najczęstszych słów w korpusie** - możemy zmieniać jej warianty:  
wybrać tylko te słowa, które są obecne we wszystkich tekstach (lub we wszystkich grupach autorskich) lub w jakimś ich procencie, najczęstsze słowa w danym języku w ogóle (w oparciu o statystyki językoznawców)

**inne cele:** ponieważ metoda ta mierzy różnice między tekstami, możemy zamiast grup autorskich sprawdzać znak emocjonalny tekstu, płeć autora, czas powstania tekstu i wiele innych rzeczy

ŹRÓDŁA

- *Questions of Attribution: Attribution and Beyond*, J.Burrows
- *Delta: A Measure of Stylistic Difference and Guide to Likely Authorship*, J.Burrows
- *Understanding and Explaining Delta measures for Authorship Attribution*, S.Evert i in.
- *Interpreting Burrows' Delta: Geometric and Probabilistic Foundations*, S.Argamon
- *Authorship Attribution*, P.Joula
- *An Open Stylometric System Based on Multilevel Text Analysis*, M.Eder, M.Piasecki, T.Walkowiak
- *Stylometria komputerowa w służbie tłumacza*, J.Rybicki
- *Metody ścisłe w literaturoznawstwie i pułapki pozornego obiektywizmu - przykład stylometrii*, M.Eder