

COMO CONVENCER QUE SEUS DADOS SÃO INTERESSANTES

Processamento de Linguagem Natural e Storytelling na prática

MATHEUS EDUARDO RODRIGUES FREITAG

CIENTISTA DE DADOS

SUMÁRIO

1 NLP?

2 IDENTIFICANDO

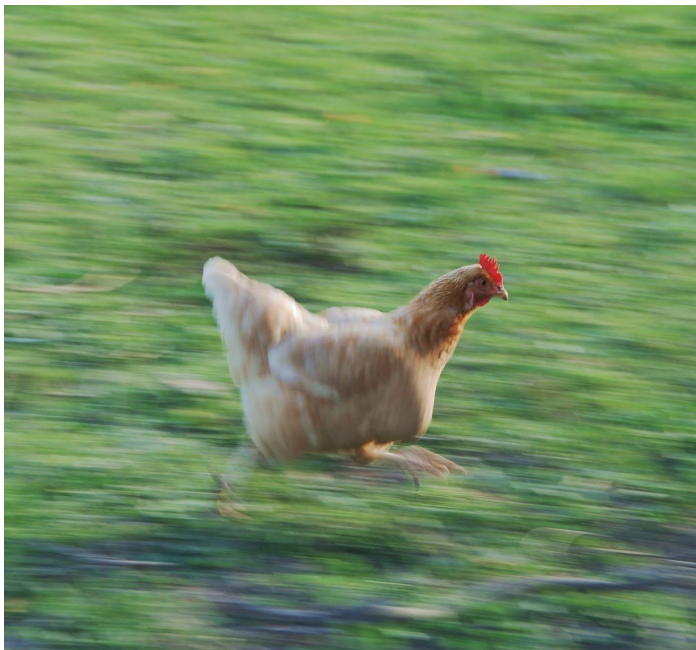
3 ENTENDENDO

4 COMPARANDO

5 SENTINDO

6 EXPANDINDO

7 FUTURO



NLP?

O que é isso, de onde veio, e porque eu deveria me importar?

IDENTIFICANDO

Vamos imaginar que o G1 contratou a Indeorum para categorizar seus artigos

MENU

G1

RIO GRANDE DO SUL

rs24

BUSCAR

Rio Grande do Sul recebe novos grupos de imigrantes venezuelanos

No total, 130 imigrantes vieram para o estado no fim de semana, conforme o Exército, para cidades do interior e da Região Metropolitana de Porto Alegre. Ministério confirma vinda de mais 10 venezuelanos na terça-feira (20).

Por G1 RS

19/08/2019 16h51 - Atualizado há 6 dias





Rio Grande do Sul recebe novos grupos de imigrantes venezuelanos

Quase um ano após a chegada dos primeiros imigrantes venezuelanos ao Rio Grande do Sul, mais dois grupos chegaram ao estado no último fim de semana, trazidos de Roraima. Conforme o Exército, 41 chegaram na sexta-feira (16) e mais 89 vieram no domingo (18).

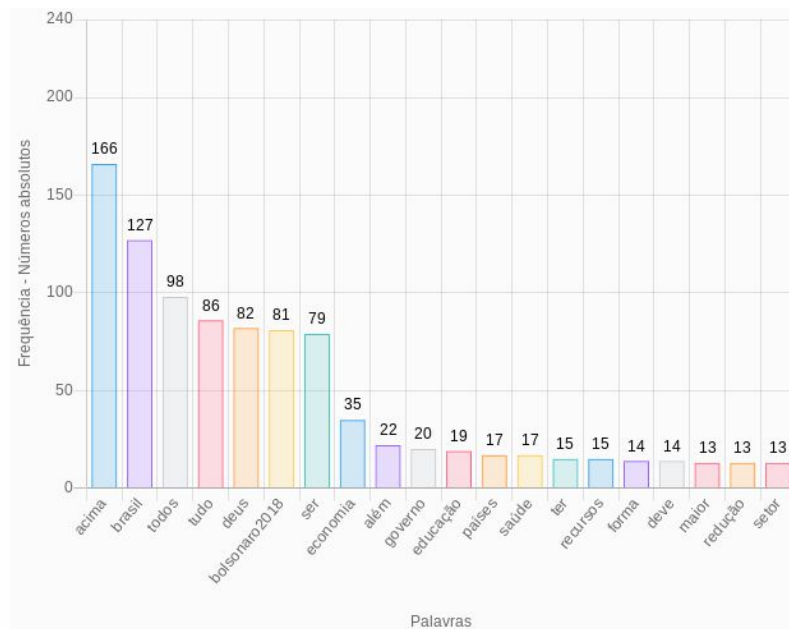
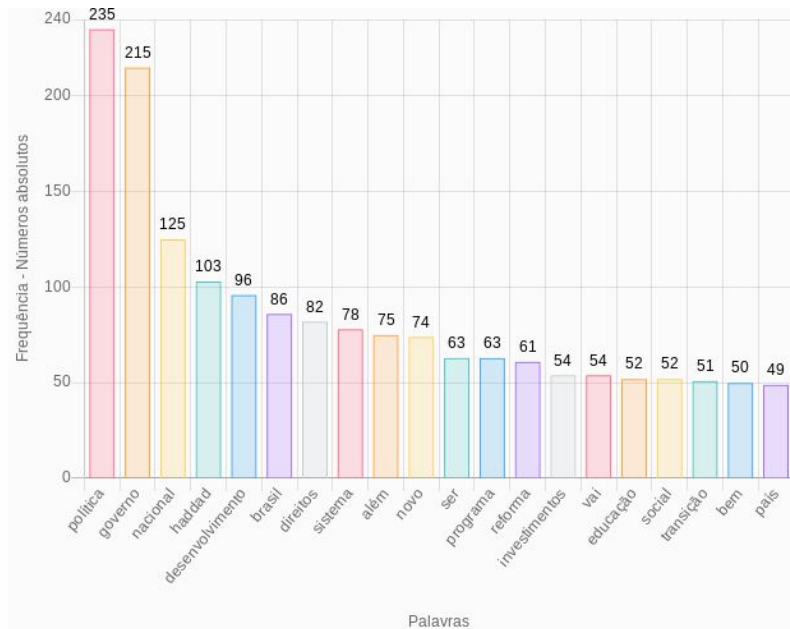
Mais 10 serão trazidos na terça-feira (20), para Canoas, segundo informações do Ministério do Desenvolvimento Social.

Os imigrantes que já chegaram ao estado foram recebidos em Porto Alegre



ENTENDENDO

Vamos imaginar um país fictício que esteja em época de eleições.



TF-IDF

TF: Term Frequency

Quantas vezes um termo aparece em um texto

$$\text{TF}(t) = (\# \text{ de vezes que o termo aparece}) / (\# \text{ de termos no texto})$$

IDF: Inverse Document Frequency

Quão especial é um termo

$$\text{IDF}(t) = \log(\# \text{ de documentos}) / (\# \text{ de documentos que } t \text{ aparece})$$

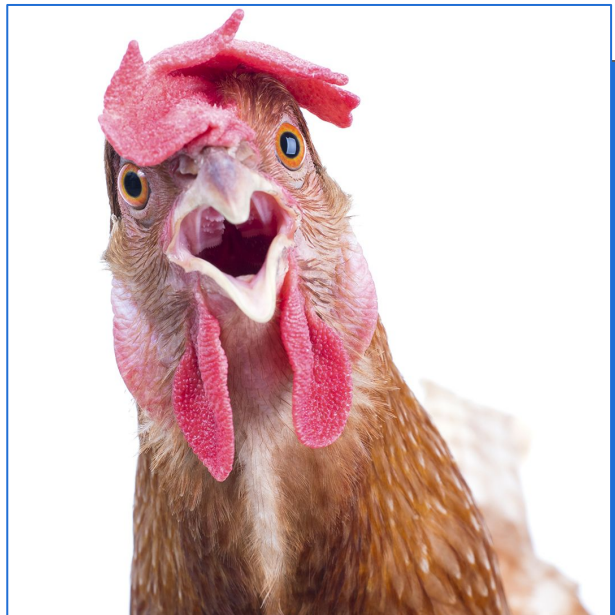
Um documento com 100 palavras, onde a palavra "gato" aparece 3 vezes possui TF igual a $(3 / 100) = 0.03$

Agora considere que temos 10M de documentos e que "gato" aparece em 1K deles, portanto o IDF é $\log(10,000,000 / 1,000) = 4$.

Então o TF-IDF se torna o produto desses dois valores: $0.03 * 4 = 0.12$

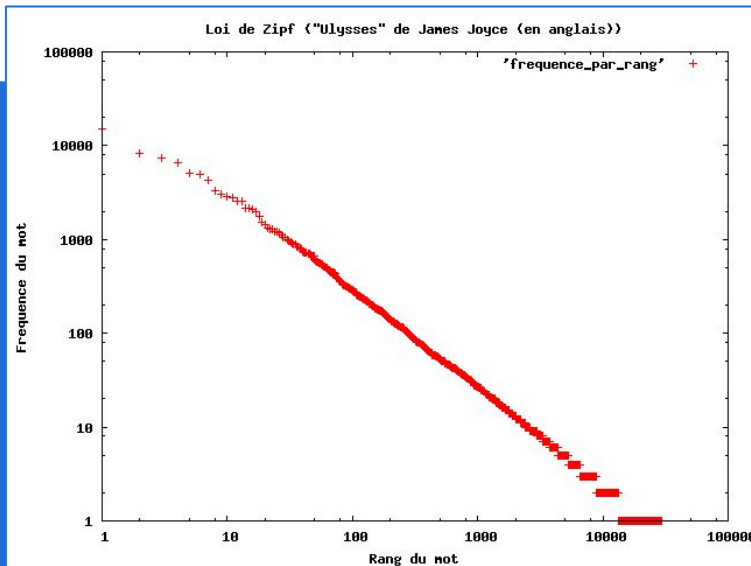


MATEMÁTICA????



- Se uma palavra aparece muito em um documento, ela é importante. Score sobe.
- Mas se essa palavra aparece em muitos documentos, ela não serve de identificador único. O Score cai.

Então palavras sem muito significado e que aparecem muito, como artigos e preposições, recebem score baixo, e palavras que aparecem muito em um único documento recebem score alto



LEI DE ZIPF

Análise Presidencial

~18k palavras

~11k palavras (sem stopwords)

~4.5k palavras únicas

Em Ulysses:

- A palavra mais comum surgiu 8000 vezes
- A décima, 800 vezes
- A centésima, 80 vezes
- A milésima, 8 vezes.

Manuscrito Voynich



COMPARANDO

Como fazemos para comparar dois corpus de texto? E como quantificar as semelhanças e diferenças?



QUÃO PARECIDOS SÃO ESSES TEXTOS?

“[...] A limpeza e secagem correta dos pés após o banho é um fator decisivo para evitar a proliferação de fungos que podem causar odores fortes[...]

“[...] A limpeza e secagem dos pés é um fator decisivo para evitar a proliferação de fungos que causam odores fortes[...]

“[...] A limpeza e secagem correta dos pés após o banho auxilia no combate a proliferação de fungos e do mau odor nos pés [...]

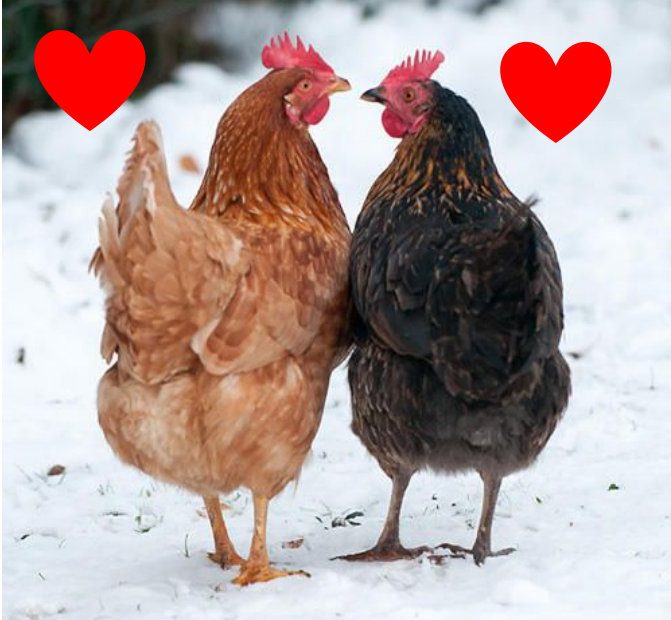
QUÃO PARECIDOS SÃO ESSES TEXTOS?

Matheus Freitag

Mateus Freitag

Mateus Freitas

DISTÂNCIA DE LEVENSHTein



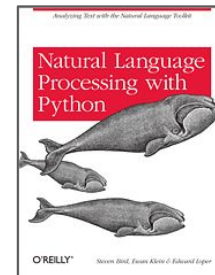
SENTINDO

Como máquinas podem
compreender emoções
inerentemente humanas?



EXPANDINDO

- **Applied Data Science with Python Specialization** <https://bit.ly/2danP4n>
- **Natural Language Processing with Python**
- **Recurrent Neural Networks**
- **Kaggle Competitions**



FUTURO

Palavras sofrem muitas variações linguísticas (especialmente línguas românticas)



STEMMING



LEMMATIZAÇÃO



OBRIGADO!

Alguma pergunta?

Matheus Freitag