CIÊNCIA DE DADOS (BIG DATA PROCESSING AND ANALYTICS)

RECUPERAÇÃO DA INFORMAÇÃO NA WEB E EM REDES SOCIAIS





TRILHA 5 PROCESSAMENTO DE LINGUAGEM NATURAL IDENTIFICAÇÃO DE PADRÕES DE TEXTO EM DOCUMENTOS

Processamento de Linguagem Natural (PLN) – Palavras

Palavra

Morfologia

por exemplo: radical, lema, afixos

Flexões

por exemplo: gênero, número, grau

Combinações

por exemplo: palavras compostos

Wikcionário em outras línguas

Mais de 1.000.000 de entradas

inglês (English) • francês (français) • malgaxe (malagasy)

Mais de 500.000 entradas

alemão (Deutsch) • chinês (中文) • espanhol (español) • holandês (nederlands) • lituano (lietuvių) • polonês / polaco (polski) • russo (русский) • servo-croata (srpskohrvatski / српскохрватски) • sueco (svenska)

Mais de 200.000 entradas

canarês (ජියූයි) • catalão (català) • coreano (한국어) • curdo (kurdî

/ • (دىم) finlandês (suomi) • grego (ελληνικά) • húngaro
(magyar) • ido • italiano • português • tâmil (தமிழ்) • turco
(türkçe) • vietnamita (tiếng việt)

Fonte: Wikicionário – CC BY-SA 3.0

Pré-tratamento – Lematização, steming e outras reduções

- Lematização consiste em reduzir a palavra para sua ideia principal.
 - Conhecimento de muitas palavras.
 - Complexo e demorado.

- Martin Porter publicou algoritmo de *steming* em 1980.
 - Baseado em regras.
 - Rápido e com bons resultados.

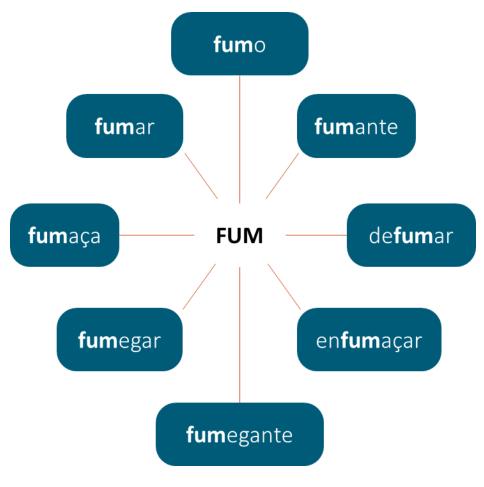


Ilustração do professor, baseada na apostila português IV da fundação CECIERJ.

- 1. João adora as aulas de recuperação de informação. Maria adora também.
- 2. Maria também adora as aulas de machine learning.

1. João adora as aulas de recuperação de informação. Maria adora também.

2. Maria também adora as aulas de machine learning.

	joao	adora	aulas	recuperacao	informacao	maria	tambem	machine	learning
1									
2									

1. João adora as aulas de recuperação de informação. Maria adora também.

{joao: 1, adora: 2, aulas: 1, recuperacao: 1, informacao: 1, maria: 1, tambem: 1, machine: 0, learning: 0}

2. Maria também adora as aulas de machine learning.

		joao	adora	aulas	recuperacao	informacao	maria	tambem	machine	learning
1	.	1	2	1	1	1	1	1	0	0
2										

1. João adora as aulas de recuperação de informação. Maria adora também.

{joao: 1, adora: 2, aulas: 1, recuperacao: 1, informacao: 1, maria: 1, tambem: 1, machine: 0, learning: 0}

2. Maria também adora as aulas de machine learning.

{joao: 0, adora: 1, aulas:0, recuperacao: 0, informacao: 0, maria: 1, tambem: 1, machine: 1, learning: 1}

	joao	adora	aulas	recuperacao	informacao	maria	tambem	machine	learning		
1	1	2	1	1	1	1	1	0	0		
2	0	1	0	0	0	1	1	1	1		
Т	1	3	1	1	1	2	2	1	1		

TF-IDF

 Contagem de palavras ignora tamanho do corpus e relevância da palavra no contexto (corpora).

• *Text Frequency (TF)* mensura a importância relativa da palavra em um documento.

$$tf(t_n, d) = \frac{|\{t_n \in d\}|}{|\{t \in d\}|}$$

 Inverse Document Frequency (IDF) mensura a importância relativa da palavra em um documento.

$$idf(t_n, D) = log \frac{|\{d \in D\}|}{|\{d \in D : t_n \in d\}|}$$

Lorem Ipsum

"Neque porro quisquam est qui dolorem ipsum quia dolor sit amet, consectetur, adipisci velit..."

"Não há quem goste de dor, que a procure e a queira ter, simplesmente porque é dor..."

Lorem ipsum dolor sit amet, consectetur adipiscing elit. Donec ipsum magna, efficitur vitae scelerisque at, finibus ut metus. Nulla fermentum erat felis, in condimentum purus varius vel. Quisque fringilla, mauris in pretium finibus, ex nisl lobortis velit. Vel molestie velit diam ac turpis. Nunc eu sodales dolor. Maecenas vel enim imperdiet, malesuada massa et, placerat ex. In et sapien sed sapien pharetra tempor. Integer commodo lacus eget nibh ultrices vehicula. Mauris nec scelerisque tellus. Vivamus et urna venenatis, lacinia turpis in, blandit turpis. Ut interdum luctus vestibulum. Nulla ac porta nibh. Cras faucibus pellentesque quam a scelerisque.

Duis mattis eu tellus sit amet vehicula. Praesent pretium magna id nisl consectetur scelerisque. Cras volutpat consectetur mi eu commodo. Interdum et malesuada fames ac ante ipsum primis in faucibus. Maecenas efficitur eu nibh ac tristique. Sed efficitur posuere feugiat. Morbi id auctor sapien, a molestie velit. Aliquam vel turpis egestas, fringilla lacus quis, ultrices nibh. Morbi eget mauris semper metus vehicula rhoncus sit amet et massa. Donec eu magna lectus.

Done suscipit fermentum tempor. Praesent vitae massa odio. Sed sed fermentum ex. in aliquam nibh. Nulla sed lectus lectus. Aliquam quis dictum mi. Etiam semper tellus a efficitur facilisis. Cras vitae nisl vitae est pharetra consectetur. Ut ornare quis metus sit amet rhoncus.

Quisque facilisis malesuada ornare. Donec sit amet vehicula magna. Donec id vulputate ligula. Nunc sollicitudin neque eu sem portitior, quis porta libero mattis. Donec ultricies, dolor mollis ultrices convallis, sem turpis eleifend elit, eu viverra nibh massa vel odio. Sed non lobortis turpis. Quisque iaculis, ligula non elementum luctus, ligula purus luctus tortor, sed congue felis mace i psum. Sed efficitur cursus mi vitae condimentum. Nullam sit amet arcu hendrerit, sagittis libero ac, iaculis orci. Ut hendrerit dolor ac posuere pulvinar. Phasellus vel fermentum leo, a eleifend dolor. Morbi di augue sit amet tortor mollis sagittis sit amet sed sem.

Vivamus sollicitudin elit lorem, eget ultricies odio pretium a. Sed ornare, tellus ultrices iaculis maximus, nisi erat eleifend diam, sed dictum est libero vitae purus. Vivamus metus nisi, consectetur vel tincidunt convallis, semper eu libero. Etiam nec pellentesque sem. Suspendisse scelerisque lacinia dui in varius. Aliquam turpis fells, aliquet et tincidunt sit amet, tincidunt ac elit. Integer nec tristique nibh. Nunc fringilla nisi enim, eget faucibus nisi pellentesque nec. In eget semper tortor. Mauris gravida dui in ante efficitur dictum id ac velit. Sed aliquam mollis risus eget efficitur. Vivamus leo ligula, consequat at ullamcorper at, ultrices ut nibh.

Interdum et malesuada fames ac ante ipsum primis in faucibus. Nulla dictum bibendum fells eget semper. Sed ipsum ante, tempor volutpat interdum faucibus, aliquet scelerisque neque. Quisque vehicula, nunc a consectetur facilisis, massa metus maximus metus, non consequat fells nibh dapibus nunc. Suspendisse euismod risus eget justo lobortis, id ultricies turpis posuere. Suspendisse faucibus est vel gravida posuere. Sed vitae fermentum est, sed maximus risus. Etiam suscipit, enim ac porta ultricies, lectus mauris facilisis purus, eget imperdiet felis purus eget metus. Curabitur quis augue odio. Sed pulvinar, odio a consectetur cursus, metus augue aliquam turpis, quis posuere leo nisl sed turpis. Duis iaculis lacus vel turpis varius, sit amet cursus libero bibendum. Vestibulum ultricies est sit amet cursus facilisis. Proin auctor purus a convallis blandit. Nulla hendrerit tincidunt ipsum ac vestibulum. Nulla ut tempor tellus, at rutrum nunc. Nunc quis magna dolor.

Ut vulputate nisi in turpis egestas consequat. Donec at iaculis dui. Proin lacinia nibh ac auctor posuere. Nullam ultricies, dui ac imperdiet mollis, purus felis ailquam nunc. a pellentesque libero dolor ac nibh. Nullam ac libero lacinia risus faucibus mollis. Praesent at felis sodales, vulputate mauris et, viverra est. Nulla mattis ut augue vitae portitor. Donec consequat dolor et nibh efficitur, nec gravida dui tincidunt. Sed malesuada dui eu est sagittis elefend. Proin portitor massa et porta laoreet. Phasellus ut congue orci. Donec aliquam metus sollicitudin nisi interdum convallis. Pellentesque in eros sagittis, fringilla diam tincidunt, rhoncus libero. Nullam ac lacus enim. Cras sagittis pures geet quam imperdiet. vitae portitor leo fermentum, Curabitur sed ornare nisi.

Suspendisse ex ante, consequat sit amet sem non, bibendum mollis nisl. Aliquam elementum vellt tortor, vitae tincidunt mi interdum ac. In tristique arcu id nulla porta, eget gravida purus vestibulum. Phasellus vehicula aliquam bibendum. Sed ac risus placerat. molestie odio egestas, convallis mi. Fusce rhoncus, justo at malesuada scelerisque, nunc mauris fringilla metus, in feuglat nisi felia sa cex. Nulla scelerisque elit eget ipsum facilisis, nec condimentum ipsum feuglat. Maecenas dignissim risus arcu, nec imperdiet elit sodales non.

Fonte: < https://pt.lipsum.com>.

TF-IDF

- 1. João adora as aulas de recuperação de informação. Maria adora também.
- 2. Maria também adora as aulas de machine learning.

BoW

	1	2
joao	1	0
adora	2	1
aulas	1	0
recuperacao	1	0
informacao	1	0
maria	1	1
tambem	1	1
machine	0	1
learning	0	1

Exercício Prático

TRILHA 3 Recuperação de informação por raspagem – robôs

Arnaldo Jabor é um cineasta, roteirista, diretor de cinema e TV, produtor cinematográfico, dramaturgo, crítico, jornalista e escritor brasileiro.

Durante 26 anos escreveu mais de 1500 artigos para jornais e revistas. E teve seu nome associado a muito mais textos vistos na internet!

O exercício consiste em puxar e formatar todos os artigos de Arnaldo Jabor, disponível no site do jornal O Tempo.



Fonte: https://www.otempo.com.br/opiniao/arnaldo-jabor

Textos atribuídos a Arnaldo Jabor

Apesar de ter mais estudo do que qualquer professor de humanas da geração Paulo Freire e mais inteligência emocional do que qualquer outro político brasileiro, não há polidez em suas palavras e tão pouco elegância em seu comportamento.

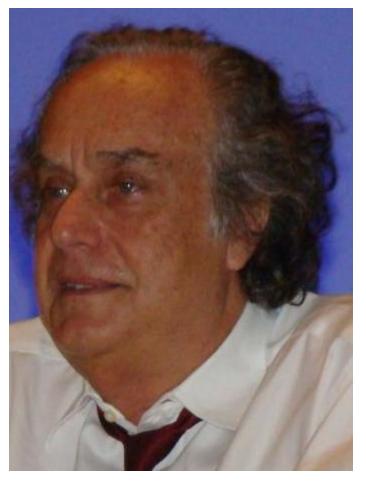
Trecho do texto atribuído a Arnaldo Jabor

Brasileiro é um povo alegre. Mentira. Brasileiro é bobalhão. Fazer piadinha com as imundices que acompanhamos todo dia é o mesmo que tomar bofetada na cara e dar risada. Depois de um massacre que durou quatro dias em São Paulo.

Trecho do texto atribuído a Arnaldo Jabor

Eu vi também a entrevista que a Marina Silva deu para outro grupo de jornalistas no dia seguinte. Meu Deus! Que coisa mais frágil! Que coisa mais pobre de ideias. Que coisa triste ver aquela boa mulher de ótimo caráter sem dúvida, mas sent adinha ali com sua vozinha falando em ética.

Trecho da transcrição da coluna à rádio CBN em 2 ago. 2018

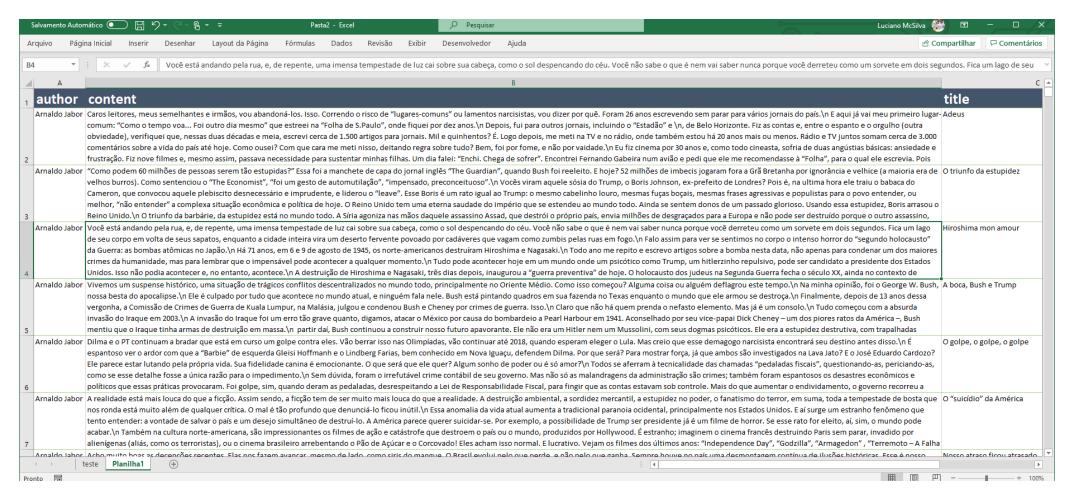


Fonte: Wikipédia

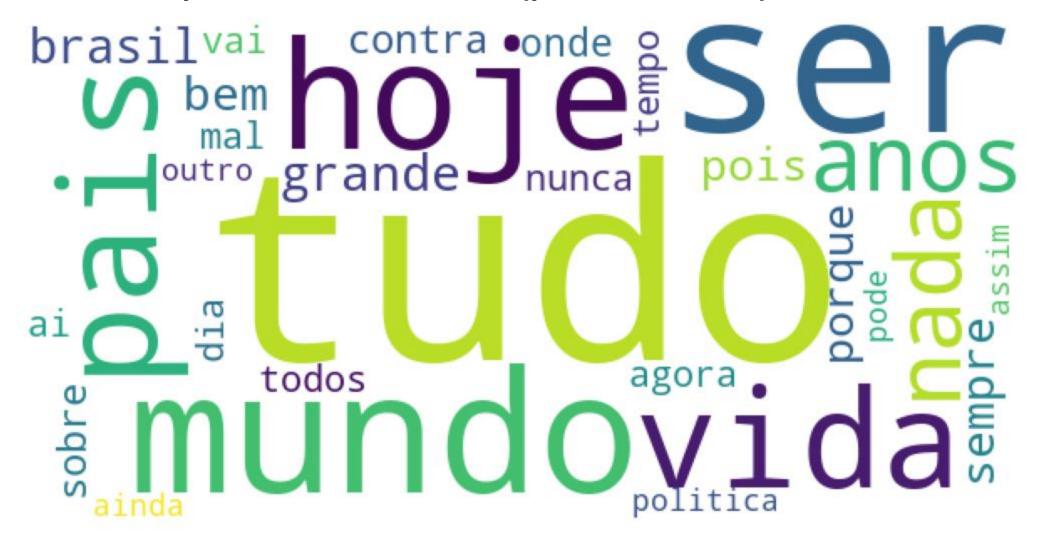
Fluxograma geral de PLN

Coleção de documentos de Origem dos Usuário informação Recuperação documentos Coleta (Ex.: web spiders) Processamento de linguagem natural Modelagem / Caso de Uso Extração de características Pré-tratamento Sumarização Unigramas Remoção de caracteres especiais Análise de Sentimento Remoção de palavras de parada Bigrama Tradução Tf-idf Remoção de Plural ... ••• ...

Arquivo com o conteúdo



Nuvem de palavras Arnaldo Jabor (pós-tratamento)

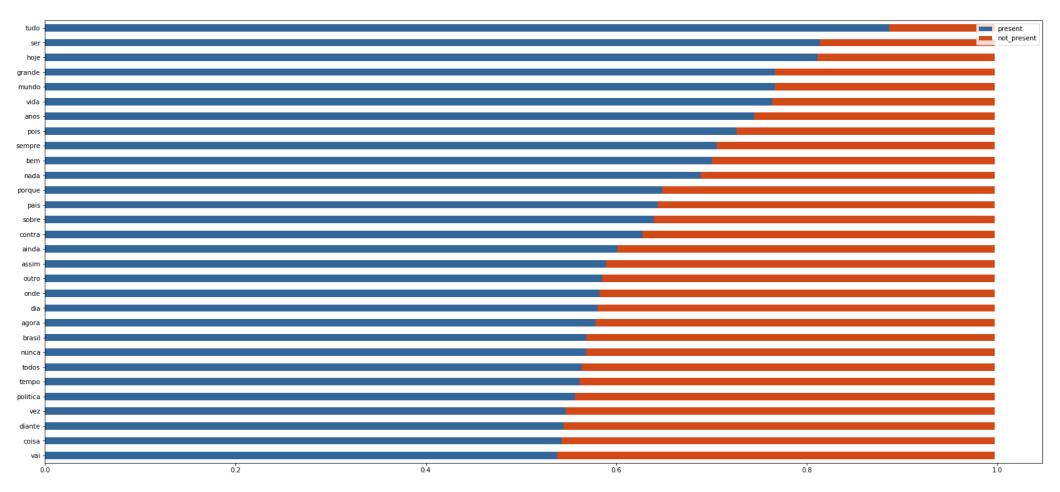


Saco de Palavras – Arnaldo Jabor

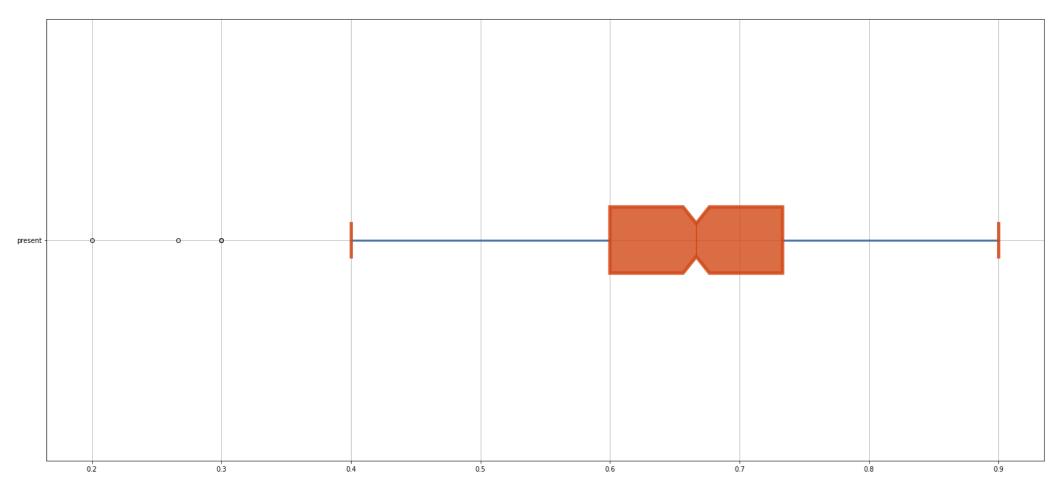
	tudo	ser	hoje	grande	brasileiro	mundo	vida	morte	morrer	otario	ahh	anos	nada	futuro
0	2	3	4	0	0	1	0	1	0	0	0	2	0	
1	4	0	3	1	0	0	2	0	0	0	0	5	3	
2	2	4	6	2	0	10	2	0	0	0	0	0	1	
3	1	3	2	2	1		4	0	0	0	0	6	1	
4	2	4	3	2	0	7) 1	0	0	0	1	1	2	
				•••					•••					
419	3	0	4	2	0	10	4	1	0	0	0	3	2	
420	4	2	0	4	2	1	4	1	0	0	0	0	4	
421	0	6	1	1	0	5	0	0	0	0	0	0	2	
422	4	2	2	0	0	0	1	0	1	0	0	5	11	
423	4	5	3	1	0	8	0	1	0	0	0	1	0	

424 rows x 28080 columns

Frequência das 30 palavras mais utilizadas – Arnaldo Jabor



Presença das 30 palavras mais utilizadas – Arnaldo Jabor



Presença das 30 palavras mais utilizadas – Arnaldo Jabor

	tudo	ser	hoje	grande	opunm	vida	anos	pois	sembre	pem	nada	porque	pais	sobre	contra	ainda	assim	outro	onde	dia	agora	brasil	nunca	todos	tempo	politica	vez	diante	coisa	vai	osn %
Texto 1	1	0	0	0	0	0	0	0	0	0	0	0	0	0	1	0	0	1	0	1	1	0	0	0	0	0	0	0	0	0	17%
Texto 2	0	1	0	0	0	1	0	0	0	0	0	1	0	0	0	1	0	0	0	1	0	0	0	0	0	0	1	0	2	0	23%
CBN	0	1	0	0	1	0	0	2	0	1	1	0	2	3	1	0	0	1	0	1	0	4	0	1	1	0	0	0	3	2	50%
Esperado	89%	81%	81%	77%	77%	76%	75%	73%	71%	70%	69%	65%	64%	64%	63%	60%	59%	58%	58%	58%	58%	57%	57%	56%	56%	56%	55%	54%	54%	54%	62%

