CIÊNCIA DE DADOS (BIG DATA PROCESSING AND ANALYTICS)

RECUPERAÇÃO DA INFORMAÇÃO NA WEB E EM REDES SOCIAIS

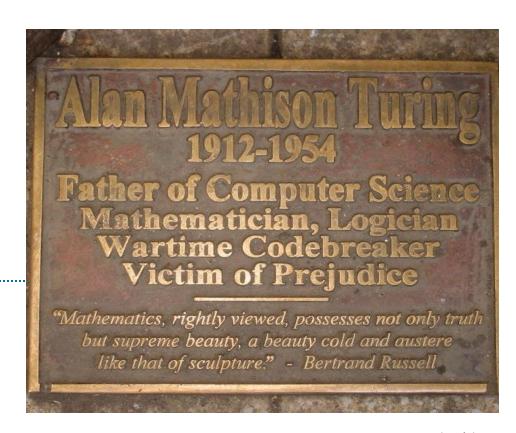




TRILHA 4 PROCESSAMENTO DE LINGUAGEM NATURAL: EXTRAÇÃO DE PALAVRAS-CHAVE

Processamento de Linguagem Natural (PLN)

- Área da ciência da computação, em particular da inteligência artificial.
- Habilitar os computadores a compreender a linguagem natural dos humanos.
- Combinar **linguísticas** com regras, **estatísticas** e modelos de **aprendizado de máquinas**
- Em 1950, **Alan Turing** publicou as bases para avaliação da linguagem natural em seu artigo "Computing Machinery and Intelligence".
- Rápida evolução a partir da década de 1980.
 Progressos recentes culminaram na criação de assistentes virtuais como Waton, Siri e Alexa.



Fonte: Wikipédia.

Casos de uso em Processamento de Linguagem Natural (PLN)

- Auto completar de formulários
- Avaliação de CV (Recrutamento)
- Sumarização de textos
- Correção de textos
- Traduções
- Chatbots
- Análise de sentimentos
- Detecção de SPAM
- Buscas de documentos
- Geração de textos



Fonte: GettyImages.

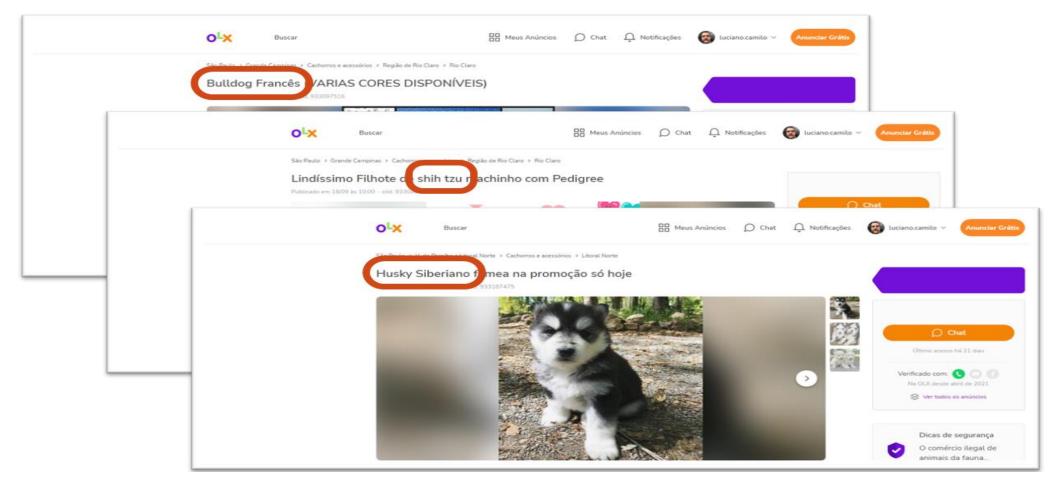
Fluxograma geral de PLN

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 ... ••• ...

Fonte: Elaborada pelo autor.

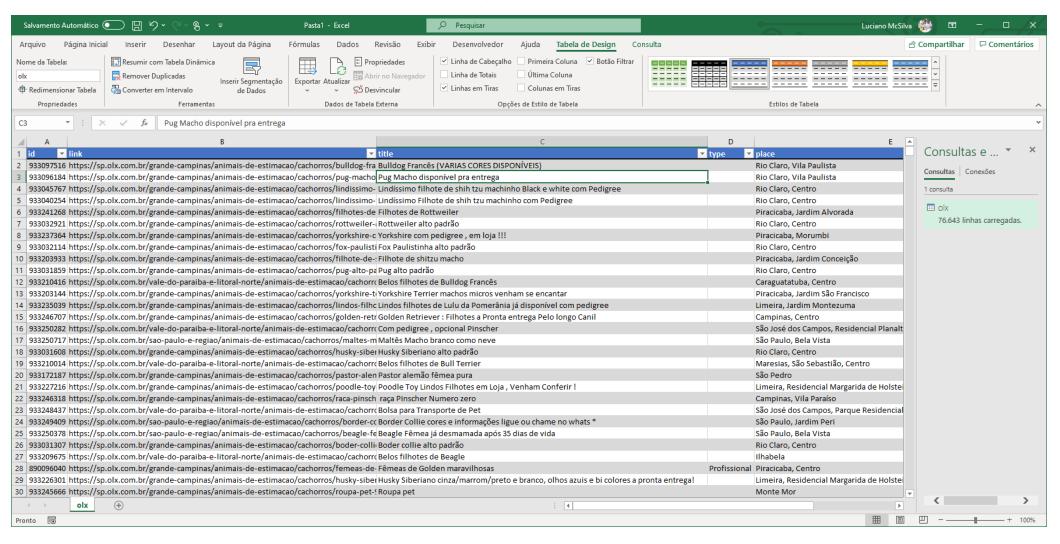
O caso de uso

O caso de uso

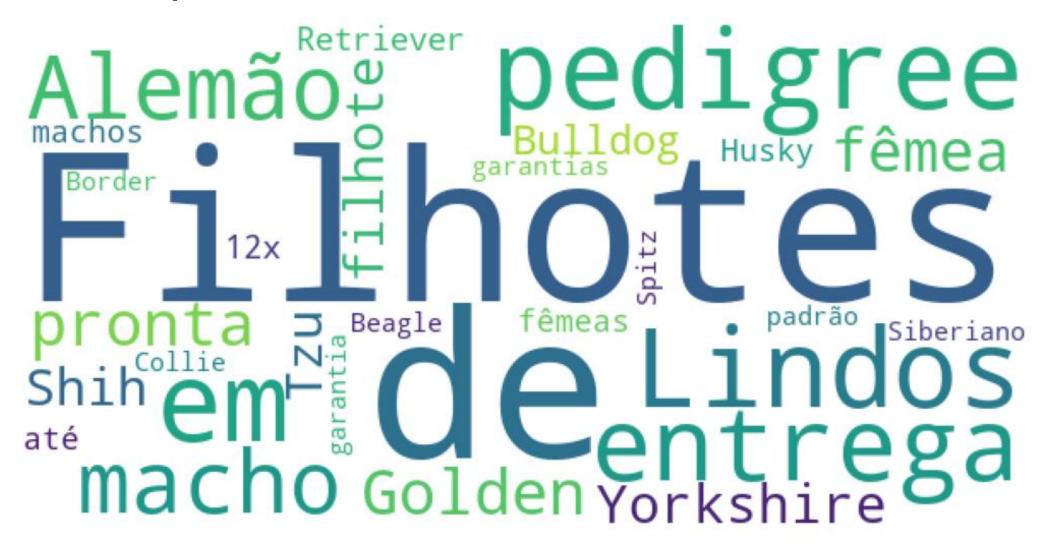


Fonte: http://www.olx.com.br. Acesso em: 29 set. 2021.

76.643 anúncios classificados organizados

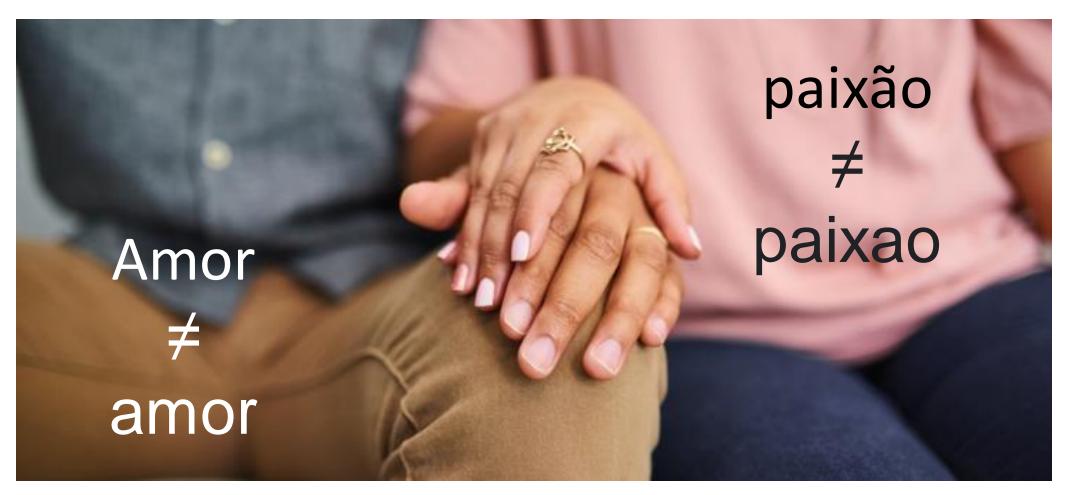


Nuvem de palavras sem nenhum tratamento





Pré-tratamento – Texto insensível e remoção de caracteres especiais



Fonte: Gettylmages

Pré-tratamento – Remoção de plural

Palavras terminadas em	Regra	Exemplos
'a', 'e', 'i', 'o', 'u', 'n'	Acrescenta-se a letra 's'	carro – carros, cadeira – cadeiras, lúmen – lúmens
'm'	Troca-se o 'm' por 'ns'	viagem – viagens, som – sons
'ão'	Utiliza-se: 'aõs', 'ões' ou 'ães'	pão – pães, cidadão – cidadãos, avião – aviões
'r', 's' ou 'z'	Acrescenta-se 'es'	cor – cores, holandês – holandeses, raiz – raízes
'al', 'el', 'ol' ou 'ul'	Troca-se o 'l' por 'is'	farol – faróis, anzol – anzóis
ʻil'	Troca-se o 'il' por 'is'	míssil – misseis, cantil – cantis
ʻx'	Não se altera	durex – durex, xerox – xerox

Pré-tratamento – Palavras de Parada



Fonte: Gettylmages

а	às	num	pelo	isto	nossa	tiver	seriam	tiverem	estivesse
0	tu	nem	pela	está	delas	terei	tinham	teremos	estiverem
е	te	meu	isso	haja	estes	terão	tivera	estivera	houvessem
é	há	nós	seus	eram	estas	teria	tenham	houvemos	houvermos
à	que	lhe	quem	fora	estou	quando	teriam	houveram	houveriam
de	com	vos	esse	seja	estão	também	aqueles	houvesse	tivéramos
do	não	teu	eles	será	houve	depois	aquelas	houverem	estivessem
da	uma	tua	você	tive	hajam	minhas	estamos	houverei	estivermos
em	por	hei	essa	teve	somos	nossos	estavam	houverão	houvéramos
um	dos	hão	suas	terá	fomos	nossas	estejam	houveria	houveremos
os	mas	sou	numa	muito	foram	aquele	estiver	fôssemos	tivéssemos
no	ele	são	elas	entre	sejam	aquela	havemos	seríamos	estivéramos
se	das	era	qual	mesmo	fosse	aquilo	houvera	tínhamos	houvéssemos
na	seu	fui	este	minha	forem	estive	hajamos	tenhamos	houveríamos
as	sua	foi	dele	pelos	serei	esteve	houverá	tivessem	estivéssemos
ao	nos	for	lhes	deles	serão	estava	fôramos	tivermos	
ou	até	tem	meus	essas	seria	esteja	sejamos	teríamos	
já	ela	tém	teus	esses	tenho	houver	seremos	estivemos	
eu	sem	para	tuas	pelas	temos	éramos	tivemos	estiveram	
só	aos	mais	dela	vocês	tinha	fossem	tiveram	estávamos	
me	nas	como	esta	nosso	tenha	formos	tivesse	estejamos	

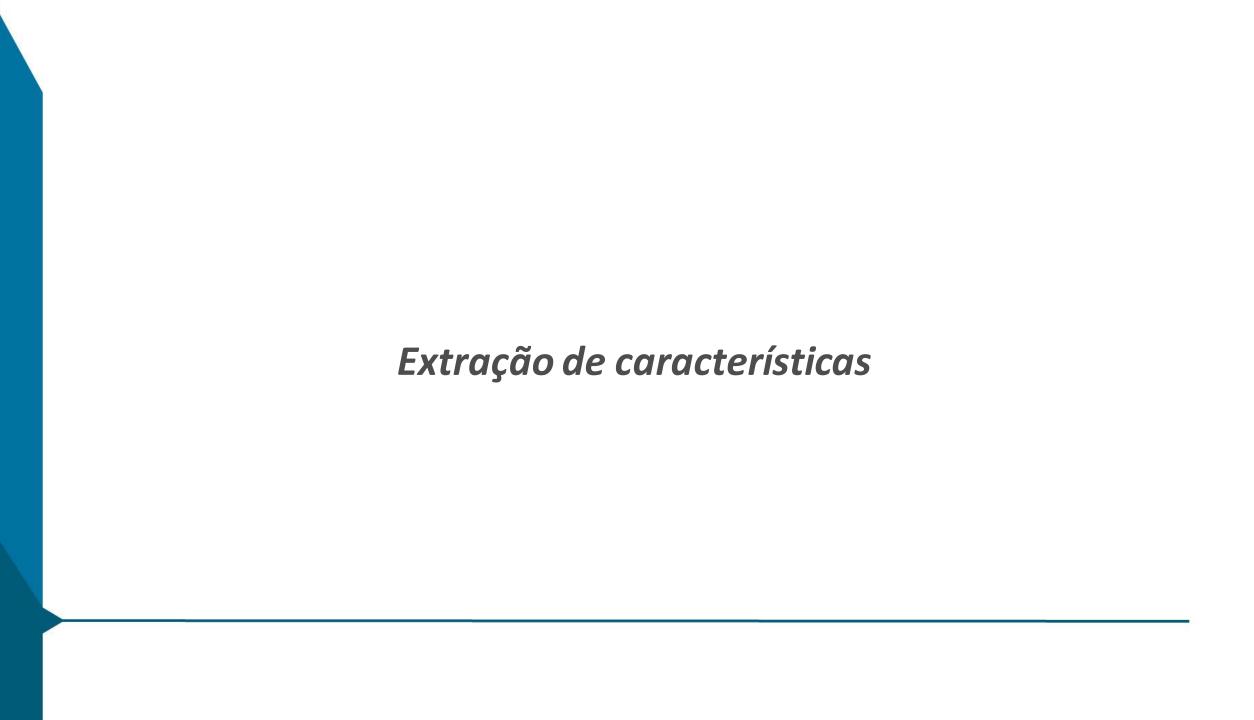
Fonte: "Natural Language Toolkit" por NLTK team licenciado sob Apache 2.0

Pré-tratamento – Palavras a serem ignoradas (dado o contexto)

С	dia	hora	total	ideal	pronta	fisica	vendese	adquirir	incriveis	lindissimo	veterinaria
r	via	leia	todas	reais	pronto	cabeca	reserva	vacinada	garantias	disponivel	assistencia
X	vai	veja	todos	puros	lindas	adulto	reserve	adoravel	microchip	entregamos	vermifugado
ja	chip	veje	todas	sonho	lindos	vacina	pelagem	condicao	vacinados	entregamos	apartamento
so	hoje	casa	venha	preco	padrao	azuis]	tamanho	amorosos	saudaveis	disposicao	companheiro
sp	raca	apto	belos	unico	compro	prontos	conosco	imediata	vacinados	parcelamos	filhotinhos
hs	loja	belo	lojas	unica	compra	entrega	filhote	conhecer	adoraveis	belissimos	maravilhosos
cm	dias	azul	porte	sonho	chamar	suporte	amarelo	namorada	seguranca	felicidade	oportunidade
ate	voce	macho	amigo	otimo	melhor	confira	machinho	conferir	condicoes	transporte	parcelamento
pra	todo	femea	super	otima	cartao	adquira	pedigree	proprias	polegadas	oferecemos	vermifugados
vet	toda	lindo	hiper	chama	varias	tamanho	lindinho	servicos	perfeitos	brincalhao	olhocaramelo
vet	vida	linda	meses	cinza	visita	vacinas	lindinha	retirada	descricao	exclusivas	exclusividade
top	pura	fofos	vezes	preto	unicos	garanta	fofinhos	cachorro	alexandre	exclusivos	exclusividades
bem	casa	saude	docil	preta	unicas	amoroso	garantia	legitimo	companhia	beneficios	brancochocolate
vez	capa	racas	bebes	canil	otimos	contato	contrato	tamanhos	qualidade	filhotinho	vacinado
sim	amor	vendo	chame	olhos	otimas	contato	promocao	linhagem	companhia	merlepreto	maravilhoso
nao	face	vende	horas	white	gratis	procuro	whatsapp	filhotes	excelente	lindissimos	
hrs	aqui	venda	ainda	black	fofura	alegria	gratuito	vermelho	exclusiva	disponiveis	
ver	novo	juros	unica	merle	branco	procura	saudavel	merlered	exclusivo	veterinario	
lar	bebe	whats	feliz	machos	cabeca	visitar	clinicas	lindinhos	pagamento	informacoes	
apt	info	ligue	fotos	femeas	ultimo	conheca	melhores	lindinhas	chocolate	procedencia	

Pré-tratamento – Código utilizado

```
import nltk
2. from nltk.corpus import stopwords
   from unidecode import unidecode
   import string
5.
   # download das palavras de parada em português
nltk.download('stopwords')
   stop = stopwords.words('portuguese')
10. #Transforma em minúscula e remoção de acentos
11. titles['title trd'] = titles['title'].str.lower().apply(lambda x: unidecode(x))
12.
13. # remoção de pontuação
14. titles['title trd'] = titles['title trd'].str.replace('[{}]'.format(string.punctuation), '')
15.
16. # remoção de números
17. titles['title trd'] = titles['title trd'].str.replace('[{}]'.format(string.digits), '')
18.
19. # remoção de palavras de parada
20. titles['title trd'] = titles['title trd'].str.apply(lambda x: ' '.join([w for w in x.split() if w not in (stop)]))
21.
22. # remoção plural
23. titles['title trd'] = titles['title trd'].apply(lambda x: ' '.join([singularizar(word) for word in x.split()]))
24.
25. # remoção palavras a serem ignoradas
26. titles['title_treated_ignored'] = titles['title_trd'].str.apply(lambda x: ' '.join([w for w in x.split() if w not in (ignorewords)]))
```



(falta de) Padrão

11

O bom de usar padrões de mercado é que existem muitos diferentes para você escolher.

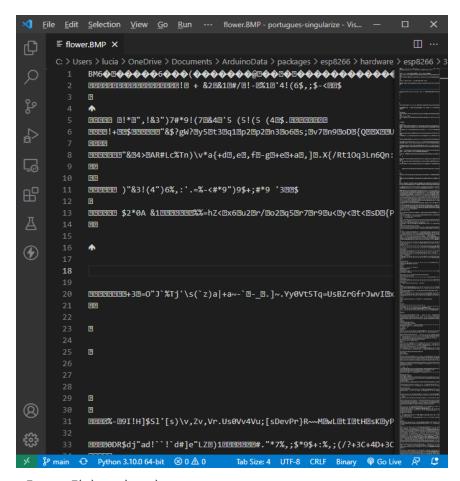


Andrew S. Tanenbaum



Fonte: Wikipédia

Arquivos

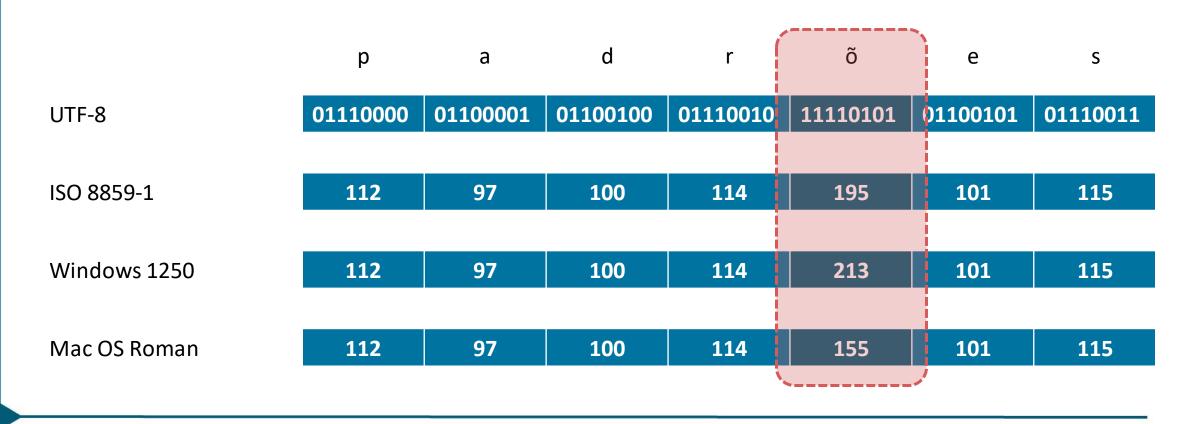


Fonte: Elaborada pelo autor.

- Arquivos computacionais são meras representações de 0s e 1s e metadados.
- Comumente separados em Binários e Textos.
- Arquivos textos não demandam metadados complexos para sua leitura, somente codificação dos bytes (*charset*).
- Há milhares de codificação possível para os bits de um texto. Os mais comuns em português são UTF-8, Windows 1250 e ISSO 8859-1.

Codificação de palavras

Diversos padrões de codificação



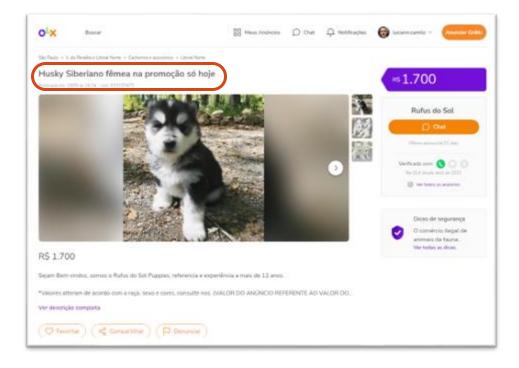
Codificação de palavras

Diversos padrões de codificação

Padrões		Decodificação							
		UTF-8	ISO 8859-1	Windows 1250	Mac OS Roman				
	UTF-8	padrões	padrões	padrões	padrões				
codificação	ISO 8859-1	padr � es	padrões	padrőes	padries				
	Windows 1250	padrões	padroes	padrões	padroes				
	Mac OS Roman	ac OS Roman padr es		padr>es	padrões				

Transformação feita pelo site https://string-functions.com/encodedecode.aspx

Husky Siberiano fêmea na promoção só hoje

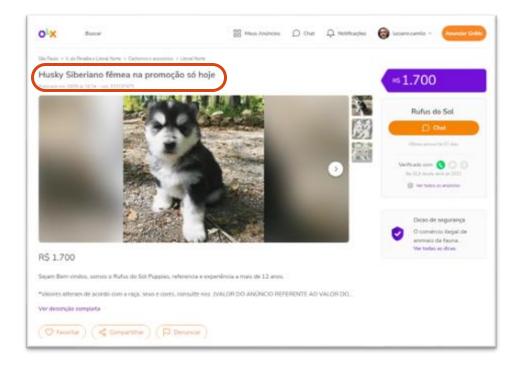


Separar cada palavra como um item próprio

Anúncio Classificado #933187475

- Husky
- Siberiano
- fêmea
- na
- promoção
- só
- hoje

husky siberiano fêmea na promoção só hoje

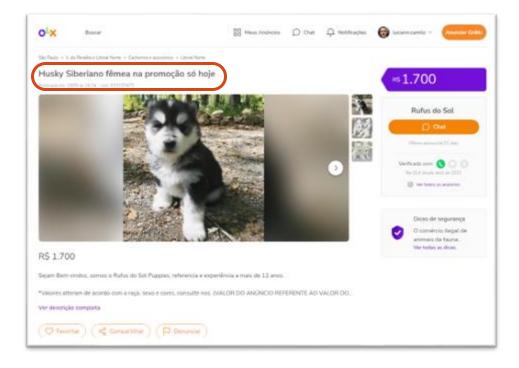


Todas em minúscula

Anúncio Classificado #933187475

- husky
- siberiano
- fêmea
- na
- promoção
- só
- hoje

husky siberiano femea na promocao so hoje

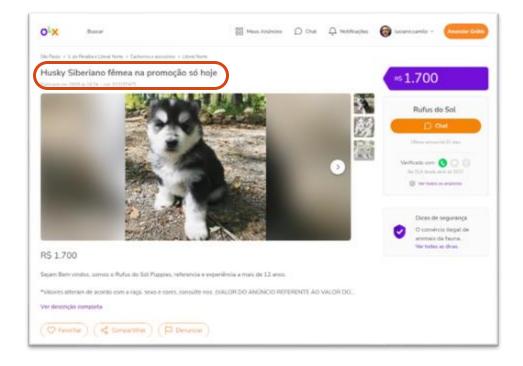


Sem caracteres especiais

Anúncio Classificado #933187475

- husky
- siberiano
- femea
- na
- promocao
- so
- hoje

husky siberiano femea na promocao so hoje

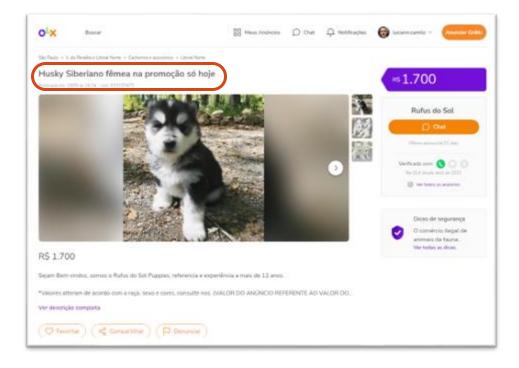


Removendo palavras de parada

Anúncio Classificado #933187475

- husky
- siberiano
- femea
- na
- promocao
- SO
- hoje

husky siberiano femea na promocao so hoje



Removendo palavras a serem ignoradas

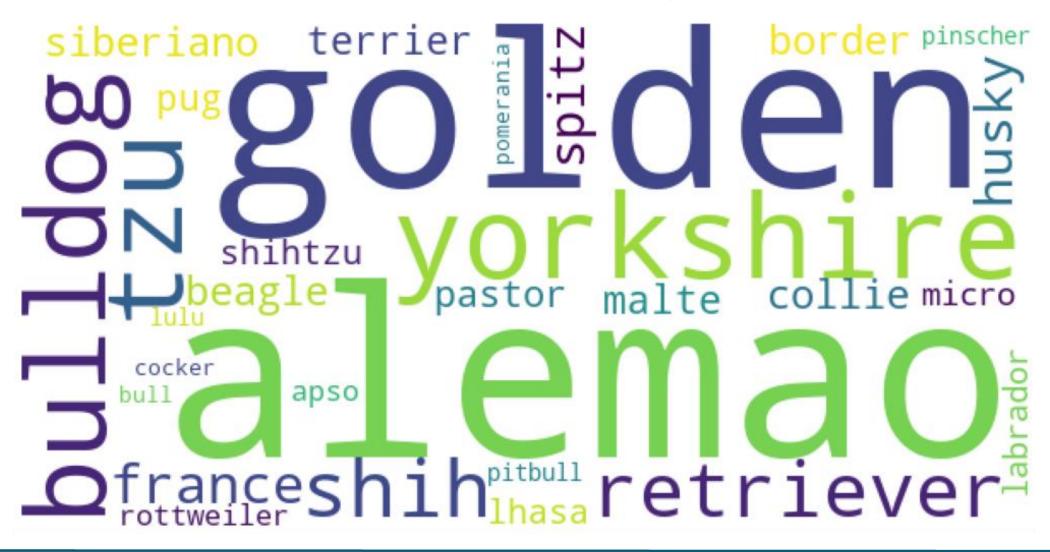
Anúncio Classificado #933187475

- husky
- siberiano
- femea
- na
- promocao
- 50
- hoje

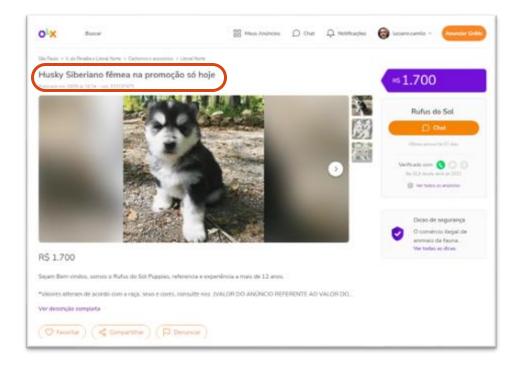
Unigramas – Código utilizado

```
from nltk.tokenize import word tokenize
2.
   from nltk import FreqDist
3.
   #cria uma lista com todos os anúncios concatenados
   treated text = ' '.join(titles['title treated ignored'])
6.
   # cria o token quebrando por palavras
   tokens = word tokenize(treated text)
9.
10. # calcula a frequência, ou seja, quantas vezes cada palavra apareceu
11. unigram fd = nltk.FreqDist(tokens)
12.
13. # ordena a frequência do maior para o menor
14. sorted ugm = sorted(unigram fd.items(), key=lambda x: x[1], reverse=True)
15.
16. # Carrega na biblioteca pandas para visualizar melhor
17. pd unigram df = pd.DataFrame.from dict(sorted ugm)
18.
19. # Apresenta as 30 palavras mais frequentes (desenhada na nuvem de palavras)
20. pd unigram df.head(30)
```

Nuvem de palavras após pré-tratamento (unigrama)



Husky Siberiano fêmea na promoção só hoje

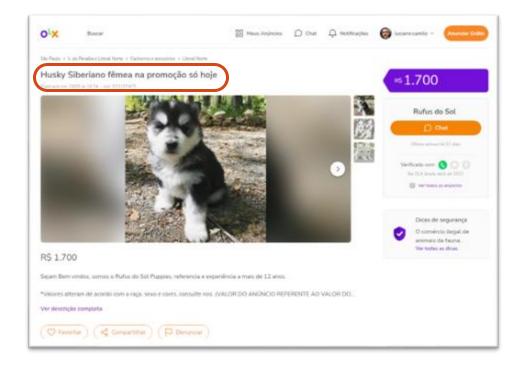


Separar cada palavra como um item próprio

Anúncio Classificado #933187475

- Husky Siberiano
- Siberiano fêmea
- fêmea na
- na promoção
- promoção só
- só hoje

husky siberiano fêmea na promoção só hoje

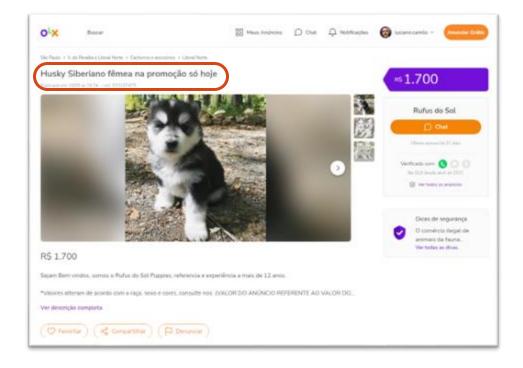


Todas em minúscula

Anúncio Classificado #933187475

- husky siberiano
- siberiano fêmea
- fêmea na
- na promoção
- promoção só
- só hoje

husky siberiano femea na promocao so hoje

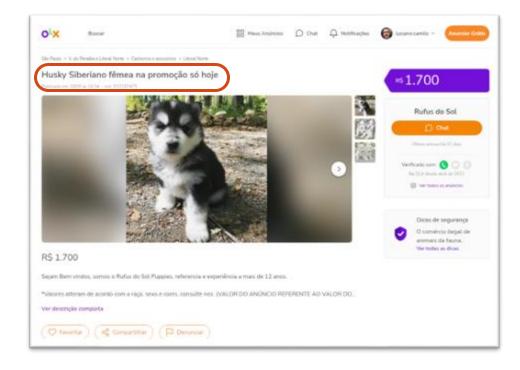


Sem caracteres especiais

Anúncio Classificado #933187475

- husky siberiano
- siberiano femea
- · femea na
- na promocao
- promocao so
- so hoje

husky siberiano femea na promocao so hoje

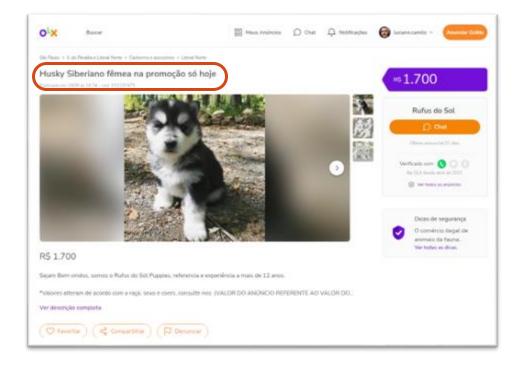


Removendo palavras de parada

Anúncio Classificado #933187475

- husky siberiano
- siberiano femea
- femea na
- na promocao
- femea promoção
- promoção so
- so hoje

husky siberiano femea na promocao so hoje



Removendo palavras a serem ignoradas

Anúncio Classificado #933187475

- husky siberiano
- siberiano femea
- femea na
- na promocao
- femea promoção
- promoção so
- so hoje

Bigramas – Código utilizado

```
from nltk.tokenize import word tokenize
   from nltk import FreqDist
   from nltk.util import ngrams
4.
   #cria uma lista com todos os anúncios concatenados
   treated text = ' '.join(titles['title treated ignored'])
7.
8. # cria o token quebrando por palavras
9. tokens = word tokenize(treated text)
10.
11. # cria o token por bigrama
12. bigram tokens = nltk.bigrams(tokens)
13.
14. # calcula a frequência, ou seja, quantas vezes cada palavra apareceu
15. bigram fd = nltk.FreqDist(bigram tokens)
16.
17. # ordena a frequência do maior para o menor
18. sorted bgm = sorted(bigram fd.items(), key=lambda x: x[1], reverse=True)
19.
20. # Carrega na biblioteca pandas para visualizar melhor
21. pd bigram df = pd.DataFrame.from dict(sorted bgm)
22.
23. # Apresenta as 30 palavras mais frequentes (desenhada na nuvem de palavras)
24. pd bigram df.head(30)
```

Nuvem de palavras após seleção dos 30 mais frequentes bigramas

