

Introdução ao Processamento de Linguagem Natural usando Python

# O QUE É PYLADIES?

PyLadies é um grupo internacional de mentoria com foco em ajudar mais mulheres a tornarem-se participantes ativas e líderes da comunidade Python.



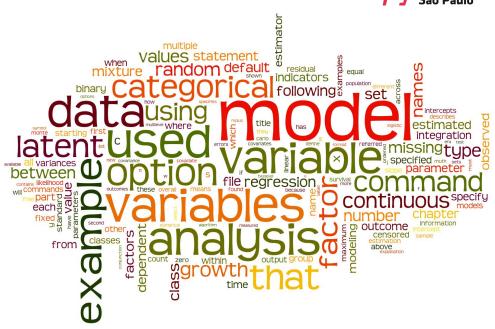


#souPyLadiesSP

# Roteiro de hoje

pyladies. São Paulo

- Introdução
- Objetivo
- Ferramentas
- Relembrando...
- Fundamentos de Processamento de Linguagem Natural
- Pré Processamento dos dados
- Análise na Prática
- Um pouco de Deep Learning
- Quer praticar mais?
- Referências Bibliográficas



# Introdução - O que é?



### Processamento de língua natural (PLN)

É uma subárea da ciência da computação, **inteligência artificial** e da linguística que estuda os problemas da **geração e compreensão automática de línguas humanas naturais**. Sistemas de geração de língua natural convertem informação de **bancos de dados** de computadores em linguagem compreensível ao ser humano e sistemas de compreensão de língua natural convertem ocorrências de linguagem humana em representações mais formais, mais **facilmente manipuláveis** por programas de computador.

# Ciência da Computação e Linguística



Processamento de linguagem natural é junção entre duas áreas: Ciência da computação e a Linguística.



# Tipos de Linguagem



#### Linguagem Estruturada

Quando a linguagem é estruturada e independe de interpretação.

É fácil de ser processada pelo computador pois ela é definida por um conjunto restrito de regras ou gramaticais

#### **Exemplo:**

Lógica 
$$\rightarrow$$
 (A + B) & (A + C)

Programação → Select nome from tabela;

#### Linguagem Não Estruturada

Qualquer texto que não siga uma estrutura pré-definida.

Possui regras gramaticais e algumas frases podem ter uma estrutura bem simples, mas na maior parte do tempo, a linguagem natural não é estruturada e ambígua.

# **Tipos de Linguagem**



E como os computadores podem processar linguagem não estruturada?

Keywords, Parts of Speech, Names Entities, Datas, Quantidades, Contagem de Palavras, Bag of Words, Estatística....etc....

Ou seja, o computador primeiro precisa fazer a extração da informação útil daquele texto (Parse) de uma sentença antes de processá-la .

# Linguagens



#### **Línguas Naturais**

Usadas no dia a dia e produzida por humanos.

Exemplo: Português, Inglês, Espanhol, Alemão.

### **Línguas Artificiais**

Linguagens de programação e notações matemáticas. \



PLN pode ser definido como uma forma de descobrir quem faz o quê, a quem, quando, onde, como e porquê.

# Estágios da análise de NLP



contexto.

sintática



variantes das

palavras, como as inflexões verbais

# **Aplicações**



casandra decision tree communication skills pig hive predictive communication skills productive mattlab predictive computer SQL R hadoop machine learning forecasting regression python sas mathematical computer science stata computer science computer science computer science classification classificatio

WorldCloud

#### **Análise de Sentimentos**



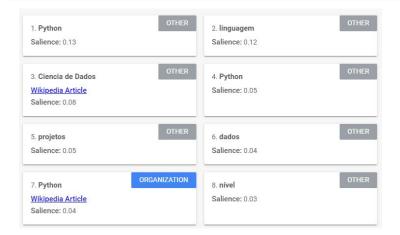
# **Aplicações**

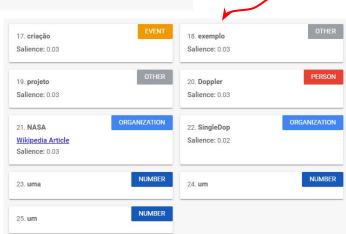


Google Cloud Platform
Natural Language API

### Identificação de entidades

 $\label{eq:continuous} $$\langle Python \rangle_1$ não é somente $\langle uma \rangle_{23}$ $\langle linguagem \rangle_2$ mais comumente conhecida para $\langle Ciencia de Dados \rangle_3$ ! $\langle Python \rangle_4$ também é utilizada em $\langle projetos \rangle_5$ grandes como por $\langle exemplo \rangle_{15}$ a $\langle NASA \rangle_{22}$ ! A $\langle NASA \rangle_{22}$ utiliza Python em diversos de seus $\langle projeto \rangle_{19}$, como por $\langle exemplo \rangle_{18}$ o $\langle SingleDop \rangle_{20}$, $\langle um \rangle_{24}$ $\langle toolkit \rangle_{11}$ que recuperara $\langle ventos \rangle_{10}$ bidimensionais de baixo $\langle nível \rangle_8$ a partir de $\langle dados \rangle_6$ de $\langle radar \rangle_9$ $\langle Doppler \rangle_{21}$ reais ou $\langle simulados \rangle_{16}$, e mais recentemente $\langle Python \rangle_7$ foi utilizado na $\langle criação \rangle_{17}$ da primeira $\langle imagem \rangle_{12}$ de $\langle um \rangle_{25}$ $\langle buraco negro \rangle_{13}$ da $\langle historia \rangle_{14}$!}$ 





# Objetivo de hoje



# Análise de dados do Reclame Aqui



# Ferramentas - O que vamos usar?

## **Bibliotecas e Ferramentas**















## Relembrando...



#### Variáveis - Int e Strings

- <u>String</u> representa um conjunto de caracteres disposto numa determinada ordem. Sempre que falarmos o termo String, estaremos nos referindo a um conjunto de caracteres.
- Int são os dados compostos por caracteres numéricos.

#### Listas

Uma lista no Python armazena valores separados por vírgulas.

#### Expressão Regular/Regex

 São usadas para identificar se um padrão existe em uma determinada sequência de caracteres (string) ou não.

#### Função

É uma sequência de comandos que executa alguma tarefa e que tem um nome definido por nós.

#### List comprehensions

 É uma estrutura importantíssima para se trabalhar com grandes conjuntos de dados, tendo uma performance superior a outras estruturas em python e simplifica a escrita do código.

# **List Comprehesion**





```
Ist = []
for x in 'PyLadies':
    Ist.append(x)
```

# **Exercícios**

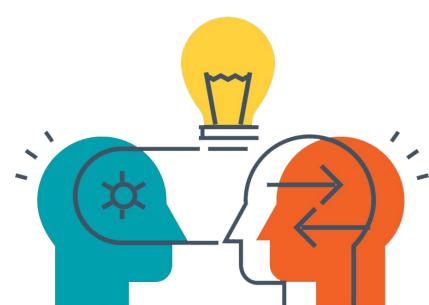




Notebook 1 - Relembrando Python ... com exercícios



Fundamentos de Processamento de Linguagem **Natural** 



# **Exercícios**





Notebook 2 - Intro NLP ... com exercícios também



# Pipeline de NLP





# Pipeline de NLP



#### O que vamos ver hoje!



Aplicar ao nosso dataset (corpus) algum desses tipos de análise e realiza transformações no dataset:

- Limpeza
- Normalização
- Tokenização
- Stemming
- Lemmatization
- outros

Só terá resultados úteis se realizar o processamento de texto de forma bem-feita

- Aplicar técnica estatística
- Coletar informações resumidas
- Outros tipos de análise

Aplicação de técnica de Machine Learning ou Deep Learning para automatizar o processo e entregar o resultado final

# Corpus



#### O que é um Corpus?

Corpus é o conjunto de textos escritos e registros orais em uma determinada língua e que serve como base de análise.

Corpora é o plural de Corpus.

O termo dataset também é usado quando falamos de Corpus



# Corpus

# pyladies São Paulo

#### Corpus machado, nativo do nltk

O nltk contém alguns outros corpus em português:

- Memórias Póstumas de Brás Cubas (1881)
- Dom Casmurro (1899)
- Gênesis
- Folha de São Paulo (1994)

#### 1. Exemplo - Corpus

In [13]: from nltk.corpus import machado

# Verificando o conjunto de textos contido no Corpus Machado
print(machado.fileids())

# Cada arquivo corresponde a uma das obras de Machado de Assis.

['contos/macn001.txt', 'contos/macn002.txt', 'contos/macn003.txt', 'contos/macn004.txt', 'contos/macn005.txt', 'contos/macn005.txt', 'contos/macn006.txt', 'contos/macn06.txt', 'contos/m cn006.txt', 'contos/macn007.txt', 'contos/macn008.txt', 'contos/macn010.txt', 'contos/macn011.txt', 'contos/macn012.txt', 'contos/macn013.txt', 'contos/macn014.txt', 'contos/macn015.txt', 'contos/macn016.txt', n017.txt', 'contos/macn018.txt', 'contos/macn019.txt', 'contos/macn020.txt', 'contos/macn021.txt', 'contos/mac 'contos/macn023.txt', 'contos/macn024.txt', 'contos/macn025.txt', 'contos/macn026.txt', 'contos/macn027.txt', 'contos/macn027.txt', 'contos/macn026.txt', 'contos/macn027.txt', 'contos/macn026.txt', 'contos/macn027.txt', 'contos/macn026.txt', 'contos/macn027.txt', 'contos/macn027.txt', 'contos/macn027.txt', 'contos/macn026.txt', 'contos/macn027.txt', n028.txt', 'contos/macn029.txt', 'contos/macn030.txt', 'contos/macn031.txt', 'contos/macn032.txt', 'contos/macn033.txt', 'contos/macn034.txt', 'contos/macn035.txt', 'contos/macn036.txt', 'contos/macn037.txt', 'contos/macn038.txt', 'contos/mac n039.txt', 'contos/macn040.txt', 'contos/macn041.txt', 'contos/macn042.txt', 'contos/macn043.txt', 'contos/macn044.txt', 'contos/macn045.txt', 'contos/macn046.txt', 'contos/macn047.txt', 'contos/macn048.txt', 'contos/macn049.txt', 'contos/macn049.txt', 'contos/macn046.txt', n050.txt', 'contos/macn051.txt', 'contos/macn052.txt', 'contos/macn053.txt', 'contos/macn054.txt', 'contos/macn055.txt', 'contos/macn056.txt', 'contos/macn057.txt', 'contos/macn058.txt', 'contos/macn059.txt', 'contos/macn060.txt', 'contos/macn060.txt', 'contos/macn059.txt', n061.txt', 'contos/macn062.txt', 'contos/macn063.txt', 'contos/macn064.txt', 'contos/macn065.txt', 'contos/macn066.txt', 'contos/macn067.txt', 'contos/macn068.txt', 'contos/macn069.txt', 'contos/macn070.txt', 'contos/macn071.txt', 'contos/macn070. n072.txt', 'contos/macn073.txt', 'contos/macn074.txt', 'contos/macn075.txt', 'contos/macn076.txt', 'contos/macn077.txt', 'contos/macn078.txt', 'contos/macn079.txt', 'contos/macn080.txt', 'contos/macn081.txt', 'contos/macn082.txt', 'contos/macn080.txt', n083.txt', 'contos/macn084.txt', 'contos/macn085.txt', 'contos/macn086.txt', 'contos/macn087.txt', 'contos/macn088.txt', 'contos/macn089.txt', 'contos/macn090.txt', 'contos/macn091.txt', 'contos/macn092.txt', 'contos/macn093.txt', 'contos/macn093.txt', 'contos/macn092.txt', 'contos/macn093.txt', n094.txt', 'contos/macn095.txt', 'contos/macn096.txt', 'contos/macn097.txt', 'contos/macn098.txt', 'contos/macn099.txt', 'contos/macn100.txt', 'contos/macn101.txt', 'contos/macn102.txt', 'contos/macn103.txt', 'contos/macn104.txt', 'contos/macn104.txt', 'contos/macn105.txt', 'contos/macn105.txt',

# **Corpus**



Entre no notebook e execute as células do "1. Exemplo - Corpus"

O que mais podemos fazer com esse corpus?



**NLTK** 





#### **Dataset**



#### Conjunto de Reclamações do Reclame Aqui

Nosso conjunto de dados é uma pequena amostra de reclamações retiradas do famoso site reclame aqui.

Vamos usar pandas para importar o csv

```
# Importando Pandas
import pandas as pd

# Vamos agora importar os dados que vamos trabalhar!
reclamacoes = pd.read_csv('reclamacoes.csv', sep=';')
print(reclamacoes.shape)
reclamacoes.head()
```



# Atenção!



Os nomes das lojas foram ocultados para mantermos a segurança da loja. A empresa foi substituído por códigos para representar cada uma individualmente.

Dentro do corpo da reclamação, o nome da loja foi substituído por "LOJA"

Nomes e emails também foram removidos.



# **Dataset Reclame Aqui**



No notebook e execute as células do

"Dados Reclame Aqui - 1. Importando os Dados"

Analise os dados, olhe o que temos em cada coluna, e como o dado está









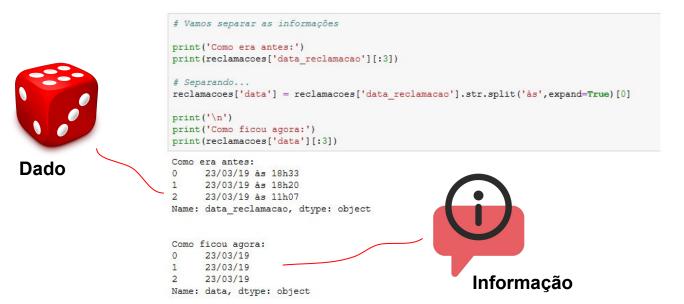


Antes de qualquer coisa, precisamos limpar nossos dados!

#### 2.1 Separando os dados da coluna data\_reclamacao

Repare bem na coluna data\_reclamacao

Ela tem duas informações em um único dado... vamos separar!





#### 2.1 Separando os dados da coluna data\_reclamacao



Sua vez!

Faça a mesma coisa para a informação hora.

Você consegue identificar algo que está faltando ao utilizar essa nossa técnica?

Dica: utilize reclamacoes.data[0] para identificar

Tente resolver!





# pyladies São Paulo

#### 2.2 Quebrando a coluna local em Cidade e Estado

Sua vez...de novo!

Da mesma forma que você separou os dados de Data/Hora, separe agora os dados de Cidade/Estado da variável local

Não se esqueça de resolver aquele problema que identificamos no slide anterior!





#### Como ficou nossos dados

pyladies São Paulo

Depois de todas essas alterações como estão nossos dados?

```
# Visualizando as alterações que fizemos
reclamacoes[['data_reclamacao', 'data', 'hora', 'local', 'cidade', 'estado']].head()
```



|   | data_reclamacao   | data     | hora  | local                    | cidade              | estado |
|---|-------------------|----------|-------|--------------------------|---------------------|--------|
| 0 | 23/03/19 às 18h33 | 23/03/19 | 18h33 | Guarulhos - SP           | Guarulhos           | SP     |
| 1 | 23/03/19 às 18h20 | 23/03/19 | 18h20 | Taubaté - SP             | Taubaté             | SP     |
| 2 | 23/03/19 às 11h07 | 23/03/19 | 11h07 | Franco da Rocha - SP     | Franco da Rocha     | SP     |
| 3 | 23/03/19 às 10h57 | 23/03/19 | 10h57 | Teresina - Pl            | Teresina            | PI     |
| 4 | 22/03/19 às 19h49 | 22/03/19 | 19h49 | São Gonçalo do Pará - MG | São Gonçalo do Pará | MG     |



# pyladies. São Paulo

#### 2.3 Alteração dos tipos das variáveis

Você percebeu algo estranho quando importou o seu dataset ou quando separou a data da coluna data\_reclamacao?

```
# Verificando o tipo de dados
reclamacoes.info()
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 28 entries, 0 to 27
Data columns (total 9 columns):
                     28 non-null object
empresa
                     28 non-null object
data reclamacao
local
                     28 non-null object
                     28 non-null object
titulo reclamacao
                     28 non-null object
corpo reclamacao
tags
                     27 non-null object
                     28 non-null int64
teve resposta
                     28 non-null object
data
                     28 non-null object
hora
dtypes: int64(1), object(8)
memory usage: 2.0+ KB
```





dtypes: category(1), datetime64[ns](1), object(7)

memory usage: 1.9+ KB

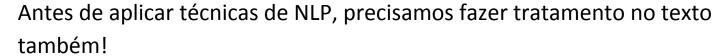


# Alteração dos tipos das variáveis Resultado:

```
# Colunas que são categoricas
reclamacoes['teve resposta'] = reclamacoes['teve resposta'].astype('category')
# Colunas que são datetime
reclamacoes['data'] = pd.to datetime(reclamacoes['data'])
reclamacoes.info()
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 28 entries, 0 to 27
Data columns (total 9 columns):
                    28 non-null object
empresa
data reclamacao
                  28 non-null object
local
                    28 non-null object
titulo reclamacao
                   28 non-null object
corpo reclamacao
                    28 non-null object
                    27 non-null object
tags
                    28 non-null category
teve resposta
data
                     28 non-null datetime64[ns]
                    28 non-null object
```

Por que isso é importante?

#### 2.4 Lower Case



Uma delas é deixar todo o texto em caixa alta ou caixa baixa.

```
print('Antes:')
print(reclamacoes['corpo reclamacao'].head())
Antes:
    Nunca mais compro nessa loja pelo fato de que ...
   Eu a LOJA comprar uma luva de musculação, na q...
   Estive na LOJA da Marginal Tiete no dia 15 de ...
    Comprei 3 produtos no dia 13.03 e recebi email...
     comprei um tenis esportivo e ao receber o avis...
Name: corpo reclamacao, dtype: object
# Aplicando Lover Case
reclamacoes['corpo reclamacao'] = [str(token).lower() for token in reclamacoes['corpo reclamacao']]
print('Depois:')
reclamacoes.corpo reclamacao.head()
Depois:
    nunca mais compro nessa loja pelo fato de que ...
   eu a loja comprar uma luva de musculação, na g...
   estive na loja da marginal tiete no dia 15 de ...
     comprei 3 produtos no dia 13.03 e recebi email...
     comprei um tenis esportivo e ao receber o avis...
Name: corpo reclamacao, dtype: object
```





### 2.5 Tokenização



Nada mais que uma segmentação de Palavras ou quebra a sequência de caracteres

Existem 2 formas de tokenizar um texto:

- Por Palavra/Tokens
- Por Sentença

#### 2.5 Tokenização

pyladies. São Paulo

13 Tokens

Por Palavra/Tokens

Alhistória do NLP começou na década de 1950 com Alan Turing

Por Sentença

O Processamento de Linguagem Natural (PLN) é a subárea da Inteligência Artificial (IA) que estuda a capacidade e as limitações de uma máquina em entender a linguagem dos seres humanos. O objetivo do PLN é fornecer aos computadores a capacidade de entender e compor textos. 'Entender" um texto significa reconhecer o contexto, fazer análise sintática, semântica, lexical e morfológica, criar resumos, extrair informação, interpretar os sentidos, analisar sentimentos e até aprender conceitos com os textos processados.

3 Sentenças

#### 2.5 Tokenização

Sua vez!

Separe os nossos textos em tokens e coloque em uma nova coluna chamada corpo\_reclamacao\_tokens





OBS: Não esqueça de colocar o resultado dentro de uma nova coluna do dataset para não comprometer os nossos dados originais ok?

Dica: Use list Comprehension ou For



#### 2.6 Stopwords

Um detalhe **muito importante** no processamento de linguagem natural é identificar as chamadas **stopwords do idioma**.

Stopword nada mais é que uma palavra que possui **apenas significado sintático** dentro da sentença, porém **não traz informações relevantes sobre o seu sentido**.

Caso contrário, os algoritmos de Machine Learning podem dar **importância para palavras como: "e", "ou", "para"**....e isso certamente atrapalha a análise.

**OBS:** O processo de "tokenização" do NLTK considera as pontuações do texto como tokens, por isso não podemos deixar de retirá-los também.





# 2.6 Stopwords Vamos para o Notebook!

```
# Removendo StopWords de todas as reclamações
# Percorre a lista de reclamações e cria uma coluna nova com o texto sem stopWords
for idx, text in enumerate (reclamacoes.corpo reclamacao):
    print('Removendo StopWords do index {}'.format(idx))
    reclamacoes.at[idx, 'corpo reclamacao semStopWords'] = remove stopwords(text, portuguese stopswords)
    print('---'*20)
Removendo StopWords do index 0
Tamanho do texto original 467
Tamanho do texto sem stopwords 223
Foram removidas 244 stopwords
Removendo StopWords do index 1
Tamanho do texto original 159
Tamanho do texto sem stopwords 78
Foram removidas 81 stopwords
Removendo StopWords do index 2
Tamanho do texto original 160
Tamanho do texto sem stopwords 85
Foram removidas 75 stopwords
Removendo StopWords do index 3
Tamanho do texto original 288
Tamanho do texto sem stopwords 163
Foram removidas 125 stopwords
```







#### 2.6 Stopwords

Tokens como "de", "estive na", "na", "da" não trazem valor a análise





Repare que foram removidas os tokens desnecessários!



#### 2.7 Normalização das palavras - Stemming

**Stemming** é uma técnica de remover prefixos e sufixos de uma palavra, chamada stem. Por exemplo, o stem da palavra reclamação é reclam. Essa técnica é muito usada em mecanismos de buscas para indexação de palavras. Pois, ao invés de armazenar todas as formas de uma palavra, o mecanismo de busca armazena apenas o stem da palavra, reduzindo o tamanho do índice e aumentando a performance do processo de busca.



```
nltk.download('rslp')
stemmer = nltk.stem.RSLPStemmer()

palavras = ['reclamação', 'reclamei', 'reclamando']

for w in palavras:
    print(stemmer.stem(w))

[nltk data] Downloading package rslp to /home/nbuser/nltk data...
```

Package rslp is already up-to-date!

[nltk data]

reclam reclam

reclam

#### 2.7 Normalização das palavras - Lemmatization



**Lemmatization** consiste em aplicar técnicas para deflexionar as palavras, retirando a conjugação verbal, caso seja um verbo, e altera os substantivos e os adjetivos para o singular masculino, de maneira a reduzir a palavra até sua forma de dicionário.

```
from nltk.stem import WordNetLemmatizer
nltk.download('wordnet')
lemmatizer = WordNetLemmatizer()
```

Exemplo de lematização, porém não existe uma biblioteca em português apenas em inglês no momento.

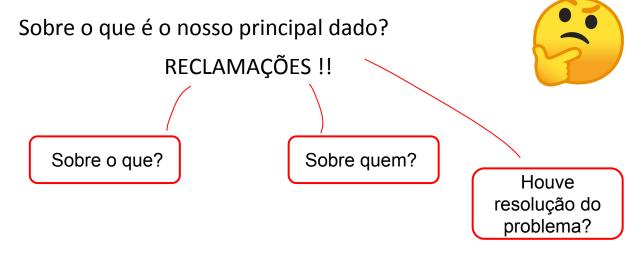
```
palavras = ['jumps', 'ladies', 'oranges']

for w in palavras:
    print(lemmatizer.lemmatize(w))
```

# Análise na prática

#### 3. Análise dos Dados







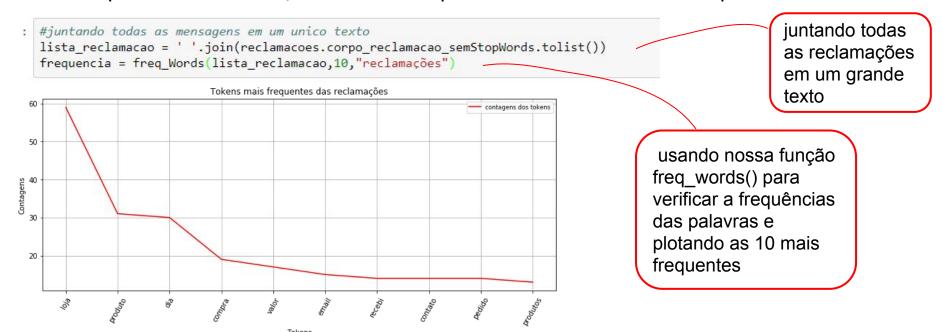


Vamos para o Notebook!

# Análise na prática



FreqDist - é usada para codificar "distribuições de frequência", que conta o número de vezes que cada resultado, no nosso caso palavras ocorre no nosso corpus



#### WordCloud



# Representação visual dos dados de texto





Na nuvem, o tamanho da palavra mostra a frequência com que ela aparece no texto, quanto maior, mais ela aparece.

# WordCloud



## **Exercício!**

Como você faria para criar um WorldCloud de uma única empresa?



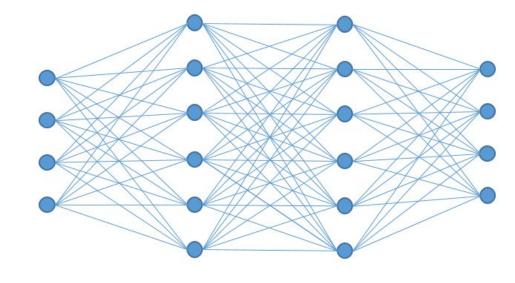
# WordCloud

Uma outra forma de fazer uma WorldCloud é utilizando um pyladies template



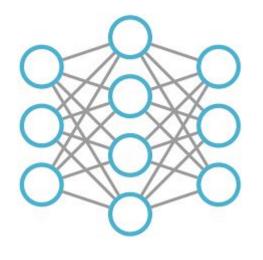


# Um pouco de Deep Learning!



# **Um pouco de Deep Learning!**





Usando **spaCy** para coletar mais informações

# spaCy



Biblioteca de software de código aberto para processamento avançado de linguagem natural

Ao contrário do *NLTK*, que é amplamente utilizado para ensino e pesquisa, o spaCy se concentra no fornecimento de software para uso em produção

Vantagem do Spacy: Possui modelos de Deep Learning pré-treinados em Português!





# spaCy - NER (Named Entity Recognition)



É a extração de informações que procura localizar e classificar menções de entidades nomeadas em texto não estruturado em categorias predefinidas, como nomes de pessoas e organizações, locais, códigos médicos, expressões de tempo, quantidades, valores monetários, percentagens, etc.

#### **Exemplo:**

 $[Jim]_{Person}$  comprou 300 ações da  $[Acme Corp.]_{Organization}$  em  $[2006]_{Time}$  .

```
# Visualizando de uma forma mais bonita!

from spacy import displacy
displacy.render(doc, style="ent")

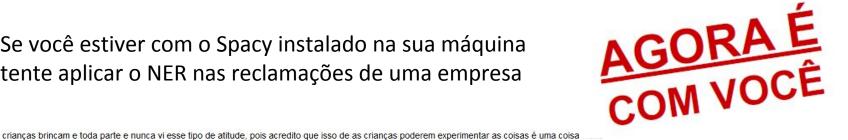
Apesar da Maria PER morar em São Paulo Loc , ela me disse que seu sonho era morar em Nova York Loc
```

# spaCy - NER (Named Entity Recognition)



#### Sua vez!

Se você estiver com o Spacy instalado na sua máquina tente aplicar o NER nas reclamações de uma empresa



legal, e meu filho ser chamado a atenção eu de verdade não entendo, e em segundo lugar e não menos importante uma pessoa se passar por gerente de loja, por que sei que toda loja tem um gerente geral e querer justificar que tinha placa se esse nem era o problema. NUNCA MISC MAIS COMPRO EM NENHUMA LOJA MISC . E QUEM EU PUDER DIZER E CONVENCER ORG DE NÃO MISC COMPRAR Eu a LOJA misc comprar uma luva de musculação, na qual a gondula está descrito melhor custo benefício, os produtos estão todos misturados na gondula sem gualquer identificação para qual seria o produto de 24.99, guando fui ao caixa com o vale troca e um produto de mesma gondula, o caixa informou que haveria uma diferenca a pagar mas o produto bfoi retirado desta gondula com varios tipos de luva de musculação, sendo está preta com polegar. Tirei Loc uma foto e chamei um atendente no qual disse que a loja não tem culpa dos produtos estarem misturados pois a loja estava muito cheia, e não aceitou trocar o produto no qual havia comprado no dia anterior. A diferença era se 5.00 ,mas o que

# spaCy - POS Tagging (Part of Speech Tagging)

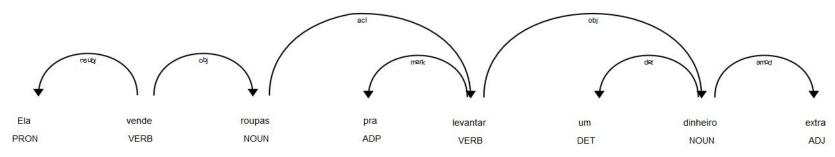
pyladies São Paulo

É a análise das classes gramaticais de um texto/frase.

Com ela é possível identificar os verbos, substantivos, adjetivos de uma frase.

Muito utilizada quando queremos gerar tradução automática de textos ou prever a próxima palavra (já que precisamos saber quais foram as últimas palavras antes da próxima e quais o contextos que elas estavam.

#### **Exemplo:**



# Quer praticar mais?



List Comprehensions - <a href="http://twixar.me/XM6K">http://twixar.me/XM6K</a>
Sumarização de Textos - <a href="https://bit.ly/2SMmVi4">https://spacy.io/usage/linguistic-features</a>
POS - <a href="http://twixar.me/qPCK">http://twixar.me/qPCK</a>

# Quem fez esse curso acontecer



Organizadoras

Deborah Froni

Linkedin - in/deborah foroni/

Jessica Cabral
Linkedin - in/jessica-cabral-carvalho/

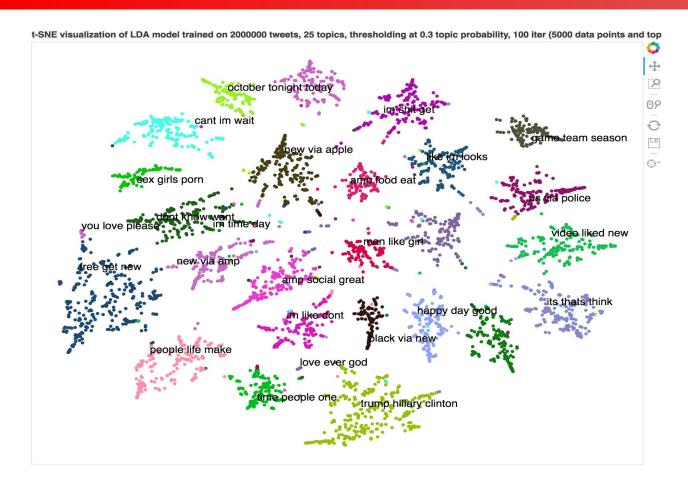
Juliana Neves
Linkedin - in/juliana-neves/

# Referências bibliográficas



- 1. <a href="http://www.eripi.com.br/2017/images/anais/minicursos/5.pdf">http://www.eripi.com.br/2017/images/anais/minicursos/5.pdf</a>
- 2. <a href="https://www.datacamp.com/community/tutorials/wordcloud-python">https://www.datacamp.com/community/tutorials/wordcloud-python</a>
- 3. <a href="https://github.com/bsacash/Introduction-to-NLP/tree/master/1.%20Quick%20Pytho">https://github.com/bsacash/Introduction-to-NLP/tree/master/1.%20Quick%20Pytho</a> <a href="master/1.%20Quick%20Pytho">n%20Refresher</a>
- 4. <a href="https://en.wikipedia.org/wiki/Named-entity\_recognition">https://en.wikipedia.org/wiki/Named-entity\_recognition</a>
- 5. <a href="https://pt.wikipedia.org/wiki/Processamento\_de\_linguagem\_natural">https://pt.wikipedia.org/wiki/Processamento\_de\_linguagem\_natural</a>
- 6. <a href="https://code.nasa.gov/?q=python">https://code.nasa.gov/?q=python</a>
- 7. <a href="https://www.datacamp.com/community/tutorials/stemming-lemmatization-python">https://www.datacamp.com/community/tutorials/stemming-lemmatization-python</a>
- 8. <a href="https://towardsdatascience.com/python-list-comprehensions-in-5-minutes-40a68cbe4561">https://towardsdatascience.com/python-list-comprehensions-in-5-minutes-40a68cbe4561</a>

# Material de Apoio - Outras aplicações





# Segmentação de Palavras/Documentos

# Material de Apoio - Outras aplicações



#### Entendimento da sintaxe do texto

number=SINGULAR

proper=PROPER

number=SINGULAR

proper=PROPER

mood=INDICATIVE

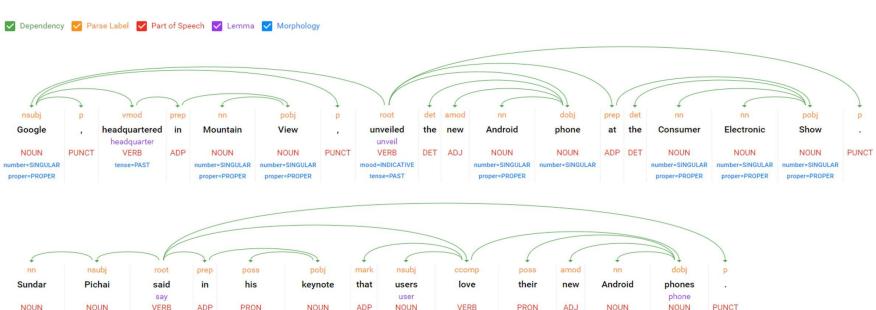
tense=PAST

case=GENITIVE

gender=MASCULINE number=SINGULAR

person=THIRD

number=SINGULAR



number=PLURAL

mood=INDICATIVE

tense=PRESENT

case=GENITIVE

number=PLURAL

person=THIRD

number=SINGULAR

proper=PROPER

number=PLURAL