

```
In [3]: import matplotlib

import codecs
import re
import copy
import collections

import numpy as np
import pandas as pd
import nltk
from nltk.stem import PorterStemmer
from nltk.tokenize import WordPunctTokenizer
from __future__ import division
```

Precisamos de algumas funções especializadas do NLTK que não estão incluídas por padrão. É possível baixar apenas a parte com as "stopwords", palavras irrelevantes, mas talvez seja mais fácil baixar tudo no NLTK. Observe que é um processo muito demorado; levou mais de 30 minutos em meu computador.

```
In [4]: nltk.download('all')
```

```
[nltk_data] Downloading collection 'all'
[nltk_data] |
[nltk_data] | Downloading package abc to
[nltk_data] | C:\Users\Atelli\AppData\Roaming\nltk_data...
[nltk_data] | Package abc is already up-to-date!
[nltk_data] | Downloading package alpino to
[nltk_data] | C:\Users\Atelli\AppData\Roaming\nltk_data...
[nltk_data] | Package alpino is already up-to-date!
[nltk_data] | Downloading package biocreative_ppi to
[nltk_data] | C:\Users\Atelli\AppData\Roaming\nltk_data...
[nltk_data] | Package biocreative_ppi is already up-to-date!
[nltk_data] | Downloading package brown to
[nltk_data] | C:\Users\Atelli\AppData\Roaming\nltk_data...
[nltk_data] | Package brown is already up-to-date!
[nltk_data] | Downloading package brown_tei to
[nltk_data] | C:\Users\Atelli\AppData\Roaming\nltk_data...
[nltk_data] | Package brown_tei is already up-to-date!
[nltk_data] | Downloading package cess_cat to
[nltk_data] | C:\Users\Atelli\AppData\Roaming\nltk_data...
[nltk_data] | Package cess_cat is already up-to-date!
[nltk_data] | Downloading package cess_esp to
[nltk_data] | C:\Users\Atelli\AppData\Roaming\nltk_data...
[nltk_data] | Package cess_esp is already up-to-date!
[nltk_data] | Downloading package chat80 to
[nltk_data] | C:\Users\Atelli\AppData\Roaming\nltk_data...
[nltk_data] | Package chat80 is already up-to-date!
[nltk_data] | Downloading package city_database to
[nltk_data] | C:\Users\Atelli\AppData\Roaming\nltk_data...
[nltk_data] | Package city_database is already up-to-date!
[nltk_data] | Downloading package cmudict to
[nltk_data] | C:\Users\Atelli\AppData\Roaming\nltk_data...
[nltk_data] | Package cmudict is already up-to-date!
[nltk_data] | Downloading package comparative_sentences to
[nltk_data] | C:\Users\Atelli\AppData\Roaming\nltk_data...
[nltk_data] | Package comparative_sentences is already up-to-
[nltk_data] | date!
[nltk_data] | Downloading package comtrans to
[nltk_data] | C:\Users\Atelli\AppData\Roaming\nltk_data...
[nltk_data] | Package comtrans is already up-to-date!
[nltk_data] | Downloading package conll2000 to
[nltk_data] | C:\Users\Atelli\AppData\Roaming\nltk_data...
[nltk_data] | Package conll2000 is already up-to-date!
[nltk_data] | Downloading package conll2002 to
[nltk_data] | C:\Users\Atelli\AppData\Roaming\nltk_data...
[nltk_data] | Package conll2002 is already up-to-date!
[nltk_data] | Downloading package conll2007 to
[nltk_data] | C:\Users\Atelli\AppData\Roaming\nltk_data...
[nltk_data] | Package conll2007 is already up-to-date!
[nltk_data] | Downloading package crubadan to
[nltk_data] | C:\Users\Atelli\AppData\Roaming\nltk_data...
[nltk_data] | Package crubadan is already up-to-date!
[nltk_data] | Downloading package dependency_treebank to
[nltk_data] | C:\Users\Atelli\AppData\Roaming\nltk_data...
[nltk_data] | Package dependency_treebank is already up-to-date!
[nltk_data] | Downloading package dolch to
[nltk_data] | C:\Users\Atelli\AppData\Roaming\nltk_data...
[nltk_data] | Package dolch is already up-to-date!
[nltk_data] | Downloading package europarl_raw to
[nltk_data] | C:\Users\Atelli\AppData\Roaming\nltk_data...
[nltk_data] | Package europarl_raw is already up-to-date!
[nltk_data] | Downloading package floresta to
```

```
[nltk_data] | C:\Users\Atelli\AppData\Roaming\nltk_data...
[nltk_data] | Package floresta is already up-to-date!
[nltk_data] | Downloading package framenet_v15 to
[nltk_data] | C:\Users\Atelli\AppData\Roaming\nltk_data...
[nltk_data] | Package framenet_v15 is already up-to-date!
[nltk_data] | Downloading package framenet_v17 to
[nltk_data] | C:\Users\Atelli\AppData\Roaming\nltk_data...
[nltk_data] | Package framenet_v17 is already up-to-date!
[nltk_data] | Downloading package gazetteers to
[nltk_data] | C:\Users\Atelli\AppData\Roaming\nltk_data...
[nltk_data] | Package gazetteers is already up-to-date!
[nltk_data] | Downloading package genesis to
[nltk_data] | C:\Users\Atelli\AppData\Roaming\nltk_data...
[nltk_data] | Package genesis is already up-to-date!
[nltk_data] | Downloading package gutenber to
[nltk_data] | C:\Users\Atelli\AppData\Roaming\nltk_data...
[nltk_data] | Package gutenber is already up-to-date!
[nltk_data] | Downloading package ier to
[nltk_data] | C:\Users\Atelli\AppData\Roaming\nltk_data...
[nltk_data] | Package ier is already up-to-date!
[nltk_data] | Downloading package inaugural to
[nltk_data] | C:\Users\Atelli\AppData\Roaming\nltk_data...
[nltk_data] | Package inaugural is already up-to-date!
[nltk_data] | Downloading package indian to
[nltk_data] | C:\Users\Atelli\AppData\Roaming\nltk_data...
[nltk_data] | Package indian is already up-to-date!
[nltk_data] | Downloading package jeita to
[nltk_data] | C:\Users\Atelli\AppData\Roaming\nltk_data...
[nltk_data] | Package jeita is already up-to-date!
[nltk_data] | Downloading package kimmo to
[nltk_data] | C:\Users\Atelli\AppData\Roaming\nltk_data...
[nltk_data] | Package kimmo is already up-to-date!
[nltk_data] | Downloading package knbc to
[nltk_data] | C:\Users\Atelli\AppData\Roaming\nltk_data...
[nltk_data] | Package knbc is already up-to-date!
[nltk_data] | Downloading package lin_thesaurus to
[nltk_data] | C:\Users\Atelli\AppData\Roaming\nltk_data...
[nltk_data] | Package lin_thesaurus is already up-to-date!
[nltk_data] | Downloading package mac_morpho to
[nltk_data] | C:\Users\Atelli\AppData\Roaming\nltk_data...
[nltk_data] | Package mac_morpho is already up-to-date!
[nltk_data] | Downloading package machado to
[nltk_data] | C:\Users\Atelli\AppData\Roaming\nltk_data...
[nltk_data] | Package machado is already up-to-date!
[nltk_data] | Downloading package masc_tagged to
[nltk_data] | C:\Users\Atelli\AppData\Roaming\nltk_data...
[nltk_data] | Package masc_tagged is already up-to-date!
[nltk_data] | Downloading package moes_sample to
[nltk_data] | C:\Users\Atelli\AppData\Roaming\nltk_data...
[nltk_data] | Package moes_sample is already up-to-date!
[nltk_data] | Downloading package movie_reviews to
[nltk_data] | C:\Users\Atelli\AppData\Roaming\nltk_data...
[nltk_data] | Package movie_reviews is already up-to-date!
[nltk_data] | Downloading package names to
[nltk_data] | C:\Users\Atelli\AppData\Roaming\nltk_data...
[nltk_data] | Package names is already up-to-date!
[nltk_data] | Downloading package nombank.1.0 to
[nltk_data] | C:\Users\Atelli\AppData\Roaming\nltk_data...
[nltk_data] | Package nombank.1.0 is already up-to-date!
[nltk_data] | Downloading package nps_chat to
[nltk_data] | C:\Users\Atelli\AppData\Roaming\nltk_data...
[nltk_data] | Package nps_chat is already up-to-date!
[nltk_data] | Downloading package omw to
[nltk_data] | C:\Users\Atelli\AppData\Roaming\nltk_data...
[nltk_data] | Package omw is already up-to-date!
[nltk_data] | Downloading package opinion_lexicon to
[nltk_data] | C:\Users\Atelli\AppData\Roaming\nltk_data...
[nltk_data] | Package opinion_lexicon is already up-to-date!
[nltk_data] | Downloading package paradigms to
[nltk_data] | C:\Users\Atelli\AppData\Roaming\nltk_data...
[nltk_data] | Package paradigms is already up-to-date!
[nltk_data] | Downloading package pil to
[nltk_data] | C:\Users\Atelli\AppData\Roaming\nltk_data...
[nltk_data] | Package pil is already up-to-date!
[nltk_data] | Downloading package pl196x to
[nltk_data] | C:\Users\Atelli\AppData\Roaming\nltk_data...
[nltk_data] | Package pl196x is already up-to-date!
[nltk_data] | Downloading package ppattach to
[nltk_data] | C:\Users\Atelli\AppData\Roaming\nltk_data...
[nltk_data] | Package ppattach is already up-to-date!
[nltk_data] | Downloading package problem_reports to
[nltk_data] | C:\Users\Atelli\AppData\Roaming\nltk_data...
[nltk_data] | Package problem_reports is already up-to-date!
[nltk_data] | Downloading package propbank to
[nltk_data] | C:\Users\Atelli\AppData\Roaming\nltk_data...
[nltk_data] | Package propbank is already up-to-date!
```

```
[nltk_data] | Downloading package ptb to
[nltk_data] | C:\Users\Atelli\AppData\Roaming\nltk_data...
[nltk_data] | Package ptb is already up-to-date!
[nltk_data] | Downloading package product_reviews_1 to
[nltk_data] | C:\Users\Atelli\AppData\Roaming\nltk_data...
[nltk_data] | Package product_reviews_1 is already up-to-date!
[nltk_data] | Downloading package product_reviews_2 to
[nltk_data] | C:\Users\Atelli\AppData\Roaming\nltk_data...
[nltk_data] | Package product_reviews_2 is already up-to-date!
[nltk_data] | Downloading package pros_cons to
[nltk_data] | C:\Users\Atelli\AppData\Roaming\nltk_data...
[nltk_data] | Package pros_cons is already up-to-date!
[nltk_data] | Downloading package qc to
[nltk_data] | C:\Users\Atelli\AppData\Roaming\nltk_data...
[nltk_data] | Package qc is already up-to-date!
[nltk_data] | Downloading package reuters to
[nltk_data] | C:\Users\Atelli\AppData\Roaming\nltk_data...
[nltk_data] | Package reuters is already up-to-date!
[nltk_data] | Downloading package rte to
[nltk_data] | C:\Users\Atelli\AppData\Roaming\nltk_data...
[nltk_data] | Package rte is already up-to-date!
[nltk_data] | Downloading package semcor to
[nltk_data] | C:\Users\Atelli\AppData\Roaming\nltk_data...
[nltk_data] | Package semcor is already up-to-date!
[nltk_data] | Downloading package senseval to
[nltk_data] | C:\Users\Atelli\AppData\Roaming\nltk_data...
[nltk_data] | Package senseval is already up-to-date!
[nltk_data] | Downloading package sentiwordnet to
[nltk_data] | C:\Users\Atelli\AppData\Roaming\nltk_data...
[nltk_data] | Package sentiwordnet is already up-to-date!
[nltk_data] | Downloading package sentence_polarity to
[nltk_data] | C:\Users\Atelli\AppData\Roaming\nltk_data...
[nltk_data] | Package sentence_polarity is already up-to-date!
[nltk_data] | Downloading package shakespeare to
[nltk_data] | C:\Users\Atelli\AppData\Roaming\nltk_data...
[nltk_data] | Package shakespeare is already up-to-date!
[nltk_data] | Downloading package sinica_treebank to
[nltk_data] | C:\Users\Atelli\AppData\Roaming\nltk_data...
[nltk_data] | Package sinica_treebank is already up-to-date!
[nltk_data] | Downloading package smultron to
[nltk_data] | C:\Users\Atelli\AppData\Roaming\nltk_data...
[nltk_data] | Package smultron is already up-to-date!
[nltk_data] | Downloading package state_union to
[nltk_data] | C:\Users\Atelli\AppData\Roaming\nltk_data...
[nltk_data] | Package state_union is already up-to-date!
[nltk_data] | Downloading package stopwords to
[nltk_data] | C:\Users\Atelli\AppData\Roaming\nltk_data...
[nltk_data] | Package stopwords is already up-to-date!
[nltk_data] | Downloading package subjectivity to
[nltk_data] | C:\Users\Atelli\AppData\Roaming\nltk_data...
[nltk_data] | Package subjectivity is already up-to-date!
[nltk_data] | Downloading package swadesh to
[nltk_data] | C:\Users\Atelli\AppData\Roaming\nltk_data...
[nltk_data] | Package swadesh is already up-to-date!
[nltk_data] | Downloading package switchboard to
[nltk_data] | C:\Users\Atelli\AppData\Roaming\nltk_data...
[nltk_data] | Package switchboard is already up-to-date!
[nltk_data] | Downloading package timit to
[nltk_data] | C:\Users\Atelli\AppData\Roaming\nltk_data...
[nltk_data] | Package timit is already up-to-date!
[nltk_data] | Downloading package toolbox to
[nltk_data] | C:\Users\Atelli\AppData\Roaming\nltk_data...
[nltk_data] | Package toolbox is already up-to-date!
[nltk_data] | Downloading package treebank to
[nltk_data] | C:\Users\Atelli\AppData\Roaming\nltk_data...
[nltk_data] | Package treebank is already up-to-date!
[nltk_data] | Downloading package twitter_samples to
[nltk_data] | C:\Users\Atelli\AppData\Roaming\nltk_data...
[nltk_data] | Package twitter_samples is already up-to-date!
[nltk_data] | Downloading package udhr to
[nltk_data] | C:\Users\Atelli\AppData\Roaming\nltk_data...
[nltk_data] | Package udhr is already up-to-date!
[nltk_data] | Downloading package udhr2 to
[nltk_data] | C:\Users\Atelli\AppData\Roaming\nltk_data...
[nltk_data] | Package udhr2 is already up-to-date!
[nltk_data] | Downloading package unicode_samples to
[nltk_data] | C:\Users\Atelli\AppData\Roaming\nltk_data...
[nltk_data] | Package unicode_samples is already up-to-date!
[nltk_data] | Downloading package universal_treebanks_v20 to
[nltk_data] | C:\Users\Atelli\AppData\Roaming\nltk_data...
[nltk_data] | Package universal_treebanks_v20 is already up-to-
[nltk_data] | date!
[nltk_data] | Downloading package verbnet to
[nltk_data] | C:\Users\Atelli\AppData\Roaming\nltk_data...
[nltk_data] | Package verbnet is already up-to-date!
[nltk_data] | Downloading package verbnet3 to
```

```
[nltk_data] | C:\Users\Atelli\AppData\Roaming\nltk_data...
[nltk_data] | Package verbnet3 is already up-to-date!
[nltk_data] | Downloading package webtext to
[nltk_data] | C:\Users\Atelli\AppData\Roaming\nltk_data...
[nltk_data] | Package webtext is already up-to-date!
[nltk_data] | Downloading package wordnet to
[nltk_data] | C:\Users\Atelli\AppData\Roaming\nltk_data...
[nltk_data] | Package wordnet is already up-to-date!
[nltk_data] | Downloading package wordnet_ic to
[nltk_data] | C:\Users\Atelli\AppData\Roaming\nltk_data...
[nltk_data] | Package wordnet_ic is already up-to-date!
[nltk_data] | Downloading package words to
[nltk_data] | C:\Users\Atelli\AppData\Roaming\nltk_data...
[nltk_data] | Package words is already up-to-date!
[nltk_data] | Downloading package ycoe to
[nltk_data] | C:\Users\Atelli\AppData\Roaming\nltk_data...
[nltk_data] | Package ycoe is already up-to-date!
[nltk_data] | Downloading package rslp to
[nltk_data] | C:\Users\Atelli\AppData\Roaming\nltk_data...
[nltk_data] | Package rslp is already up-to-date!
[nltk_data] | Downloading package maxent_treebank_pos_tagger to
[nltk_data] | C:\Users\Atelli\AppData\Roaming\nltk_data...
[nltk_data] | Package maxent_treebank_pos_tagger is already up-
[nltk_data] | to-date!
[nltk_data] | Downloading package universal_tagset to
[nltk_data] | C:\Users\Atelli\AppData\Roaming\nltk_data...
[nltk_data] | Package universal_tagset is already up-to-date!
[nltk_data] | Downloading package maxent_ne_chunker to
[nltk_data] | C:\Users\Atelli\AppData\Roaming\nltk_data...
[nltk_data] | Package maxent_ne_chunker is already up-to-date!
[nltk_data] | Downloading package punkt to
[nltk_data] | C:\Users\Atelli\AppData\Roaming\nltk_data...
[nltk_data] | Package punkt is already up-to-date!
[nltk_data] | Downloading package book_grammars to
[nltk_data] | C:\Users\Atelli\AppData\Roaming\nltk_data...
[nltk_data] | Package book_grammars is already up-to-date!
[nltk_data] | Downloading package sample_grammars to
[nltk_data] | C:\Users\Atelli\AppData\Roaming\nltk_data...
[nltk_data] | Package sample_grammars is already up-to-date!
[nltk_data] | Downloading package spanish_grammars to
[nltk_data] | C:\Users\Atelli\AppData\Roaming\nltk_data...
[nltk_data] | Package spanish_grammars is already up-to-date!
[nltk_data] | Downloading package basque_grammars to
[nltk_data] | C:\Users\Atelli\AppData\Roaming\nltk_data...
[nltk_data] | Package basque_grammars is already up-to-date!
[nltk_data] | Downloading package large_grammars to
[nltk_data] | C:\Users\Atelli\AppData\Roaming\nltk_data...
[nltk_data] | Package large_grammars is already up-to-date!
[nltk_data] | Downloading package tagsets to
[nltk_data] | C:\Users\Atelli\AppData\Roaming\nltk_data...
[nltk_data] | Package tagsets is already up-to-date!
[nltk_data] | Downloading package snowball_data to
[nltk_data] | C:\Users\Atelli\AppData\Roaming\nltk_data...
[nltk_data] | Package snowball_data is already up-to-date!
[nltk_data] | Downloading package bllip_wsj_no_aux to
[nltk_data] | C:\Users\Atelli\AppData\Roaming\nltk_data...
[nltk_data] | Package bllip_wsj_no_aux is already up-to-date!
[nltk_data] | Downloading package word2vec_sample to
[nltk_data] | C:\Users\Atelli\AppData\Roaming\nltk_data...
[nltk_data] | Package word2vec_sample is already up-to-date!
[nltk_data] | Downloading package panlex_swadesh to
[nltk_data] | C:\Users\Atelli\AppData\Roaming\nltk_data...
[nltk_data] | Package panlex_swadesh is already up-to-date!
[nltk_data] | Downloading package mte_teip5 to
[nltk_data] | C:\Users\Atelli\AppData\Roaming\nltk_data...
[nltk_data] | Package mte_teip5 is already up-to-date!
[nltk_data] | Downloading package averaged_perceptron_tagger to
[nltk_data] | C:\Users\Atelli\AppData\Roaming\nltk_data...
[nltk_data] | Package averaged_perceptron_tagger is already up-
[nltk_data] | to-date!
[nltk_data] | Downloading package averaged_perceptron_tagger_ru to
[nltk_data] | C:\Users\Atelli\AppData\Roaming\nltk_data...
[nltk_data] | Package averaged_perceptron_tagger_ru is already
[nltk_data] | up-to-date!
[nltk_data] | Downloading package perluniprops to
[nltk_data] | C:\Users\Atelli\AppData\Roaming\nltk_data...
[nltk_data] | Package perluniprops is already up-to-date!
[nltk_data] | Downloading package nonbreaking_prefixes to
[nltk_data] | C:\Users\Atelli\AppData\Roaming\nltk_data...
[nltk_data] | Package nonbreaking_prefixes is already up-to-date!
[nltk_data] | Downloading package vader_lexicon to
[nltk_data] | C:\Users\Atelli\AppData\Roaming\nltk_data...
[nltk_data] | Package vader_lexicon is already up-to-date!
[nltk_data] | Downloading package porter_test to
[nltk_data] | C:\Users\Atelli\AppData\Roaming\nltk_data...
[nltk_data] | Package porter_test is already up-to-date!
```

```
[nltk_data] | Downloading package wmt15_eval to
[nltk_data] | C:\Users\Atelli\AppData\Roaming\nltk_data...
[nltk_data] | Package wmt15_eval is already up-to-date!
[nltk_data] | Downloading package mwa_ppdb to
[nltk_data] | C:\Users\Atelli\AppData\Roaming\nltk_data...
[nltk_data] | Package mwa_ppdb is already up-to-date!
[nltk_data] | Done downloading collection all
```

Out[4]: True

Baixar o pacote "stopwords" do NLTK.

```
In [5]: from nltk.corpus import stopwords
```

## Ler dados

```
In [27]: with codecs.open("relatorio.txt", "r", encoding="utf-8") as f:
relatorio_1 = f.read()
with codecs.open("coronavac.txt", "r", encoding="utf-8") as f:
coronavac_1 = f.read()
```

## Processar dados

### Verificar palavras irrelevantes em português.

```
In [28]: esw = stopwords.words('portuguese')
esw.append("would")
```

Filtrar tokens (usando expressões regulares).

```
In [29]: word_pattern = re.compile("^\\w+$")
```

Criar função para contagem de tokens.

```
In [30]: def get_text_counter(text):
tokens = WordPunctTokenizer().tokenize(PorterStemmer().stem(text))
tokens = list(map(lambda x: x.lower(), tokens))
tokens = [token for token in tokens if re.match(word_pattern, token) and token not in esw]
return collections.Counter(tokens), len(tokens)
```

Criar função para cálculo da frequência absoluta e da frequência relativa das palavras mais comuns.

```
In [31]: def make_df(counter, size):
abs_freq = np.array([el[1] for el in counter])
rel_freq = abs_freq / size
index = [el[0] for el in counter]
df = pd.DataFrame(data=np.array([abs_freq, rel_freq]).T, index=index, columns=["Frequencia Absoluta", "Frequencia Relativa"])
df.index.name = "Palavras mais comuns"
return df
```

## Analisar textos individuais

Calcular as palavras mais comuns de \_RelatorioAnvisa. Isso demora um pouco. Então, exibir as 10 mais comuns.

```
In [41]: re_counter, re_size = get_text_counter(relatorio_1)
make_df(re_counter.most_common(10), re_size)
```

Out[41]:

	Frequencia Absoluta	Frequencia Relativa
--	---------------------	---------------------

Palavras mais comuns		
vacina	81.0	0.017020
19	67.0	0.014079
covid	56.0	0.011767

	Frequencia Absoluta	Frequencia Relativa
<b>Palavras mais comuns</b>		
<b>2020</b>	52.0	0.010927
<b>uso</b>	51.0	0.010717
<b>vacinas</b>	51.0	0.010717
<b>4</b>	45.0	0.009456
<b>anvisa</b>	45.0	0.009456
<b>eficácia</b>	43.0	0.009036
<b>saúde</b>	42.0	0.008825

Salvar as 1.000 palavras mais comuns de *Relatorio* como CSV.

```
In [42]: re_df = make_df(re_counter.most_common(1000), re_size)
re_df.to_csv("RE_1000.csv")
```

```
In [43]: co_counter, co_size = get_text_counter(coronavac_1)
make_df(co_counter.most_common(10), co_size)
```

Out[43]:

	Frequencia Absoluta	Frequencia Relativa
<b>Palavras mais comuns</b>		
<b>estudo</b>	42.0	0.021298
<b>eficácia</b>	36.0	0.018256
<b>vacina</b>	34.0	0.017241
<b>covid</b>	27.0	0.013692
<b>19</b>	25.0	0.012677
<b>2</b>	21.0	0.010649
<b>anos</b>	19.0	0.009635
<b>dose</b>	18.0	0.009128
<b>0</b>	18.0	0.009128
<b>1</b>	15.0	0.007606

```
In [44]: co_df = make_df(co_counter.most_common(1000), co_size)
co_df.to_csv("co_1000.csv")
```

## Comparar textos

Identificar as palavras mais comuns nos dois documentos.

```
In [45]: all_counter = co_counter + re_counter
all_df = make_df(co_counter.most_common(1000), 1)
most_common_words = all_df.index.values
```

Criar um quadro de dados com as diferenças de frequência das palavras.

```
In [46]: df_data = []
for word in most_common_words:
    re_c = re_counter.get(word, 0) / re_size
    co_c = co_counter.get(word, 0) / co_size
    d = abs(re_c - co_c)
    df_data.append([re_c, co_c, d])
dist_df = pd.DataFrame(data=df_data, index=most_common_words,
                        columns=["Relatorio Anvisa Frequência Relativa", "Apresentação Coronavac Frequência Relativa",
                                "Diferença de Frequência Relativa"])
dist_df.index.name = "Palavras mais comuns"
dist_df.sort_values("Diferença de Frequência Relativa", ascending=False, inplace=True)
```

Exibir as palavras mais distintas.

```
In [47]: dist_df.head(10)
```

Out[47]:

	Relatorio Anvisa	Frequência Relativa	Apresentação Coronavac	Frequência Relativa	Diferença de Frequência Relativa
Palavras mais comuns					
estudo		0.001261		0.021298	0.020037
2020		0.010927		0.001014	0.009912
vacinas		0.010717		0.001014	0.009702
eficácia		0.009036		0.018256	0.009220
anos		0.000841		0.009635	0.008794
0		0.000420		0.009128	0.008708
dose		0.000630		0.009128	0.008497
anvisa		0.009456		0.002028	0.007427
idade		0.000420		0.007099	0.006679
uso		0.010717		0.004057	0.006660

Salvar a lista completa com as palavras distintas como um CSV intitulado "distintas.csv".

In [48]: dist\_df.to\_csv("distintas.csv")

In [ ]: