

# ml\_imdb\_random\_forest

March 12, 2021

```
[4]: from google.colab import drive
drive.mount('/content/drive')
```

Drive already mounted at /content/drive; to attempt to forcibly remount, call drive.mount("/content/drive", force\_remount=True).

```
[5]: import sys, os
sys.path.insert(0, '/content/drive/My Drive/')
```

```
[8]: !unzip /content/drive/MyDrive/imdb_data/aclImdb.zip
```

Archive: /content/drive/MyDrive/imdb\_data/aclImdb.zip  
replace \_\_MACOSX/.\_aclImdb? [y]es, [n]o, [A]ll, [N]one, [r]ename: N

```
[1]: %matplotlib inline

## Install some stuff you likely don't have
!pip3 -q install --upgrade pip
!pip3 -q install pandas_path
!pip3 -q install fasttext
!pip3 -q install umap-learn[plot]
!pip3 -q install ipywidgets --user

import json
import logging
from pathlib import Path
import random
import tarfile
import tempfile
import warnings

import matplotlib.pyplot as plt
import numpy as np
import pandas as pd
import pandas_path # Path style access for pandas
from tqdm import tqdm
```

```
import os

from PIL import Image
import umap
import umap.plot
```

```
[2]: import seaborn as sns
from pylab import rcParams
import matplotlib.pyplot as plt
from matplotlib import rc
```

```
[3]: import re

def preprocess_text(sen):
    # Removing html tags
    sentence = remove_tags(sen)
    # Remove punctuations and numbers
    sentence = re.sub('[^a-zA-Z]', ' ', sentence)
    # Single character removal
    sentence = re.sub(r"\s+[a-zA-Z]\s+", ' ', sentence)
    # Removing multiple spaces
    sentence = re.sub(r'\s+', ' ', sentence)
    return sentence
TAG_RE = re.compile(r'<[>]+>')
def remove_tags(text):
    return TAG_RE.sub('', text)
```

```
[4]: def path_to_dataframe(path):
    path1 = path+'neg/'
    all_files = os.listdir(path1)
    all_files.sort(key=lambda x:int(x[:-4]))
    dictionary_list = []
    for file in all_files:
        with open(os.path.join(path1, file), 'r', encoding="utf8") as f:
            text = preprocess_text(f.read())
            dictionary_list.append({'text':text, 'label':0})
    path2 = path+'pos/'
    all_files = os.listdir(path2)
    all_files.sort(key=lambda x:int(x[:-4]))
    for file in all_files:
        with open(os.path.join(path2, file), 'r', encoding="utf8") as f:
            text = preprocess_text(f.read())
            dictionary_list.append({'text':text, 'label':1})
    random.shuffle(dictionary_list)
    df = pd.DataFrame.from_dict(dictionary_list)
    return df
```

```
[5]: train_path = '/content/aclImdb/train/'
test_path = '/content/aclImdb/test/'
```

```
original_set = path_to_dataframe(train_path)
test_set = path_to_dataframe(test_path)

print(original_set.info())
```

```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 25000 entries, 0 to 24999
Data columns (total 2 columns):
 #   Column  Non-Null Count  Dtype
---  -
 0   text    25000 non-null   object
 1   label   25000 non-null   int64
dtypes: int64(1), object(1)
memory usage: 390.8+ KB
None
```

```
[6]: %matplotlib inline
      %config InlineBackend.figure_format='retina'

      sns.set(style='whitegrid', palette='muted', font_scale=1.2)

      HAPPY_COLORS_PALETTE = ["#01BEFE", "#FFDD00", "#FF7D00", "#FF006D", "#ADFF02",
                              ↪ "#8F00FF"]

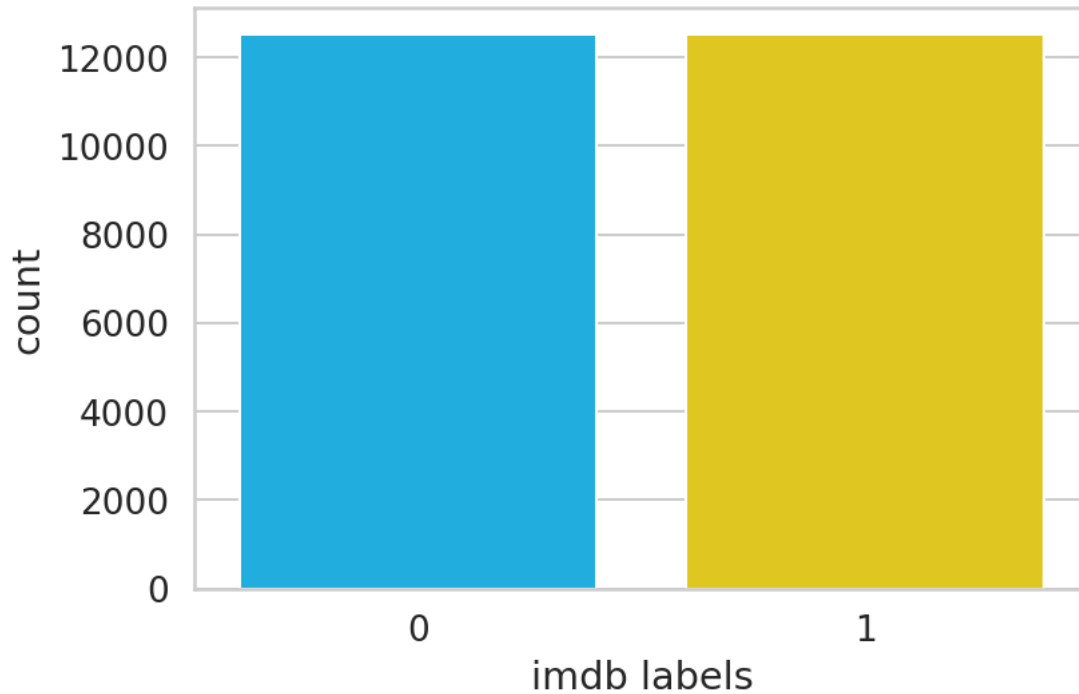
      sns.set_palette(sns.color_palette(HAPPY_COLORS_PALETTE))

[7]: sns.countplot(original_set.label)
      plt.xlabel('imdb labels')
```

```
/usr/local/lib/python3.7/dist-packages/seaborn/_decorators.py:43: FutureWarning:
Pass the following variable as a keyword arg: x. From version 0.12, the only
valid positional argument will be `data`, and passing other arguments without an
explicit keyword will result in an error or misinterpretation.
```

```
FutureWarning
```

```
[7]: Text(0.5, 0, 'imdb labels')
```



```
[8]: import re
import nltk
from sklearn.datasets import load_files
nltk.download('stopwords')
nltk.download('wordnet')
import pickle
from nltk.corpus import stopwords
```

```
[nltk_data] Downloading package stopwords to /root/nltk_data...
[nltk_data] Package stopwords is already up-to-date!
[nltk_data] Downloading package wordnet to /root/nltk_data...
[nltk_data] Package wordnet is already up-to-date!
```

```
[33]: preprocessed_train_x, y_train = original_set.text[:15000], original_set.label.
      ↪array[:15000]
preprocessed_test_x, y_test = test_set.text, test_set.label.array
```

```
[10]: train_piles = []

from nltk.stem import WordNetLemmatizer

stemmer = WordNetLemmatizer()

for sen in range(0, len(preprocessed_train_x)):
    # Remove all the special characters
```

```

train_pile = re.sub(r'\W', ' ', str(preprocessed_train_x[sen]))

# remove all single characters
train_pile = re.sub(r'\s+[a-zA-Z]\s+', ' ', train_pile)

# Remove single characters from the start
train_pile = re.sub(r'\^[a-zA-Z]\s+', ' ', train_pile)

# Substituting multiple spaces with single space
train_pile = re.sub(r'\s+', ' ', train_pile, flags=re.I)

# Removing prefixed 'b'
train_pile = re.sub(r'^b\s+', '', train_pile)

# Converting to Lowercase
train_pile = train_pile.lower()

# Lemmatization
train_pile = train_pile.split()

train_pile = [stemmer.lemmatize(word) for word in train_pile]
train_pile = ' '.join(train_pile)

train_piles.append(train_pile)

```

```
[21]: print(train_piles[2])
```

it should be against the law not to experience this extremely funny stand up show with eddie murphy have never seen anything like it murphy go on for almost minute about dick pussy tit and insults so many famous people including his own family among the people who get it by murphy are elvis mr michael jackson stevie wonder mick jagger luther vandross and james brown have seriously never laughed so hard of anything my entire life mean when person doesn know who mr is but still laugh so hard of murphy a mr there something about it at the time saw the show couldn remember who mr wa but still laughed now know who he is and that just make it so much more funny because that what eddie do he can make those impression so good that it don matter who the hell he trying to do it still hilarious and on top of that we learn that murphy actually is very good singer please watch it

```
[11]: test_piles = []

for sen in range(0, len(preprocessed_test_x)):
    # Remove all the special characters
    test_pile = re.sub(r'\W', ' ', str(preprocessed_test_x[sen]))

    # remove all single characters
    test_pile = re.sub(r'\s+[a-zA-Z]\s+', ' ', test_pile)

```

```

# Remove single characters from the start
test_pile = re.sub(r'\^[a-zA-Z]\s+', ' ', test_pile)

# Substituting multiple spaces with single space
test_pile = re.sub(r'\s+', ' ', test_pile, flags=re.I)

# Removing prefixed 'b'
test_pile = re.sub(r'^b\s+', '', test_pile)

# Converting to Lowercase
test_pile = test_pile.lower()

# Lemmatization
test_pile = test_pile.split()

test_pile = [stemmer.lemmatize(word) for word in test_pile]
test_pile = ' '.join(test_pile)

test_piles.append(test_pile)

```

```

[44]: from sklearn.feature_extraction.text import CountVectorizer
vectorizer = CountVectorizer(max_features=1500, min_df=5, max_df=0.7,
    ↳stop_words=stopwords.words('english'))
x_train = vectorizer.fit_transform(train_piles).toarray()
x_test = vectorizer.fit_transform(test_piles).toarray()

```

```

[45]: print(x_train[100])

```

```

[1 0 0 ... 0 0 0]

```

```

[46]: from sklearn.feature_extraction.text import TfidfTransformer
tfidfconverter = TfidfTransformer()
x_train = tfidfconverter.fit_transform(x_train).toarray()
x_test = tfidfconverter.fit_transform(x_test).toarray()

```

```

[47]: print(x_train[100])

```

```

[0.0815068 0.          0.          ... 0.          0.          0.          ]

```

```

[48]: print(y_train)

```

```

<PandasArray>
[0, 0, 1, 1, 0, 1, 0, 1, 1, 1,
 ...
 0, 1, 0, 1, 1, 1, 0, 1, 1, 0]
Length: 15000, dtype: int64

```

```
[49]: from sklearn.ensemble import RandomForestClassifier
      from sklearn.metrics import accuracy_score

[64]: rf = RandomForestClassifier(n_estimators=1000, max_depth=None, random_state=42,
      ↪n_jobs=-1)
      rf_model = rf.fit(x_train,y_train)

[65]: test_predictions = rf_model.predict(x_test)
      print(test_predictions)
      print(accuracy_score(y_test,test_predictions))
```

```
[0 1 1 ... 0 1 1]
0.561
```

```
[66]: predictions_df = pd.DataFrame()
      predictions_df['y_pred'] = test_predictions
      predictions_df['y_truth'] = test_set['label']
      predictions_df['y_pred'] = predictions_df['y_pred'].astype(int)

      predictions_df.head()
```

```
[66]:   y_pred  y_truth
0         0         0
1         1         1
2         1         1
3         1         0
4         1         1
```

```
[67]: predictions_df.to_csv('predictions_rf.csv')
```

```
[68]: from google.colab import files
      files.download('/content/predictions_rf.csv')
```

<IPython.core.display.Javascript object>

<IPython.core.display.Javascript object>

```
[ ]: max_depths = [5, 300, 500, 700, 1000, 2000]
      depth_accuracy = []

      for depth in max_depths:
          rf = RandomForestClassifier(n_estimators=1000,max_depth=depth,n_jobs=-1)
          rf_model = rf.fit(x_train,y_train)
          test_predictions = rf_model.predict(x_test)
          depth_accuracy.append(accuracy_score(y_test,test_predictions))
```

```
[76]: estimators_array = [100, 300, 400, 800, 1000, 200, 3000]
      estimator_accuracy = []
```

```

for n in estimators_array:
    rf = RandomForestClassifier(n_estimators=n,max_depth=None,n_jobs=-1)
    rf_model = rf.fit(x_train,y_train)
    test_predictions = rf_model.predict(x_test)
    estimator_accuracy.append(accuracy_score(y_test,test_predictions))

```

```

[77]: print(depth_accuracy)
      print(estimator_accuracy)

```

```

[0.52016, 0.56012, 0.5576, 0.56032, 0.5616, 0.55844]
[0.55608, 0.56048, 0.56168, 0.55768, 0.56092, 0.559, 0.56012]

```

```

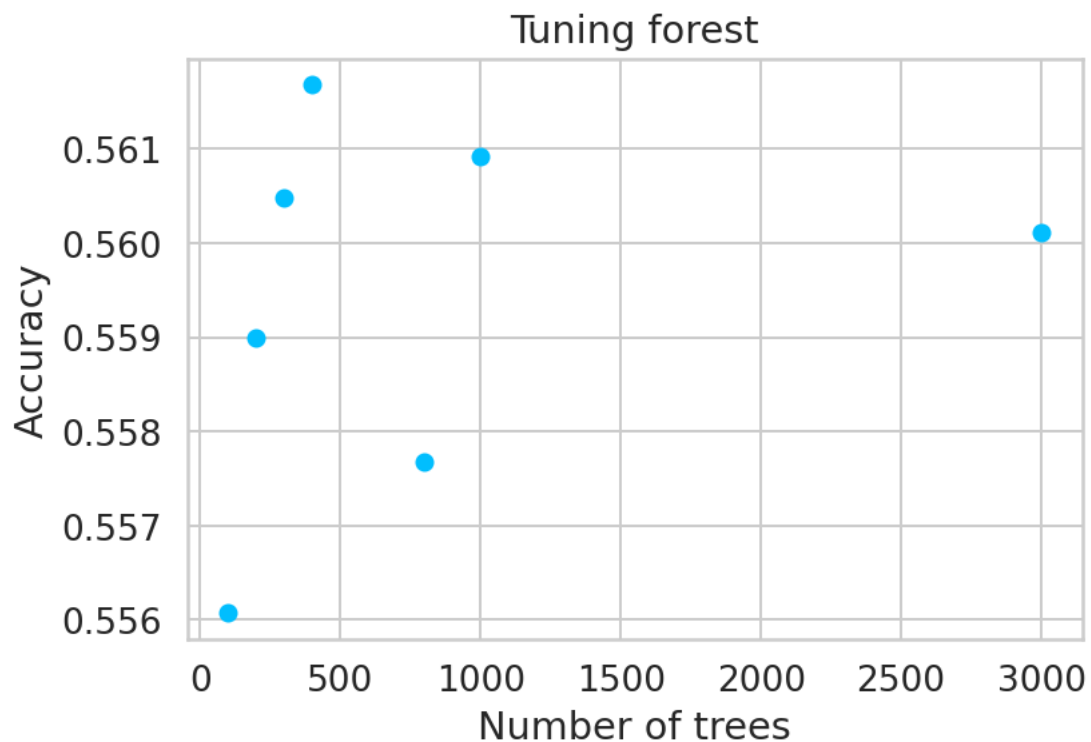
[83]: plt.scatter(estimators_array, estimator_accuracy)
      plt.title('Tuning forest')
      plt.ylabel('Accuracy')
      plt.xlabel('Number of trees')

```

```

[83]: Text(0.5, 0, 'Number of trees')

```



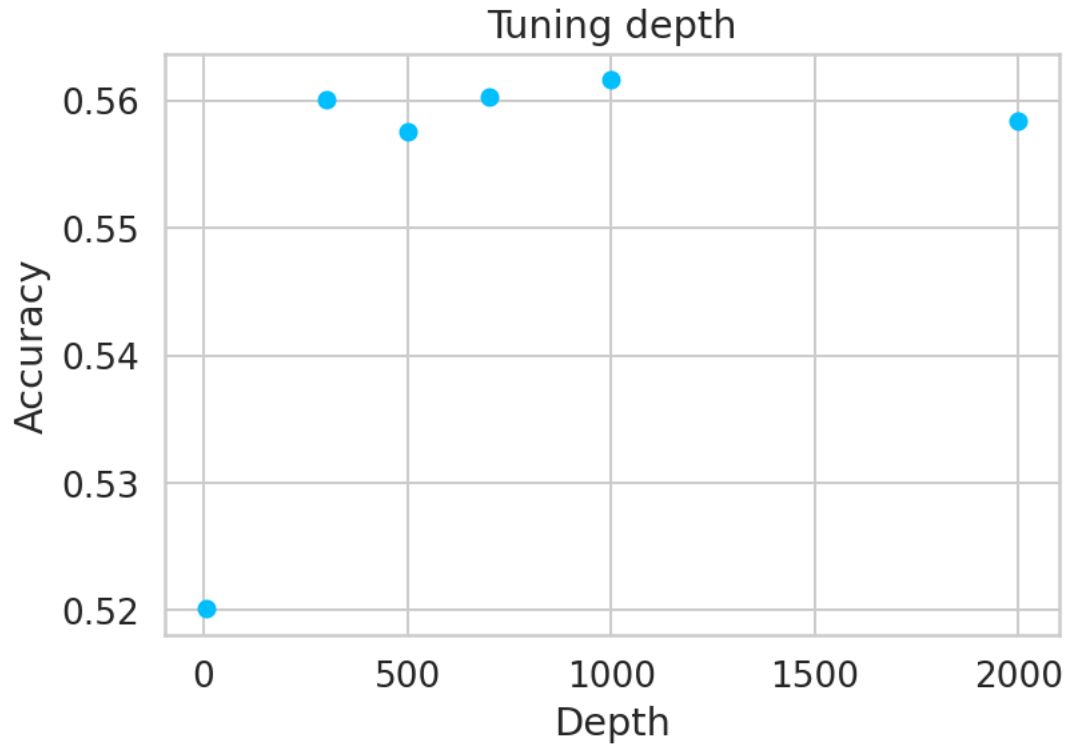
```

[79]: plt.scatter(max_depths, depth_accuracy)
      plt.title('Tuning depth')
      plt.ylabel('Accuracy')
      plt.xlabel('Depth')

```



[79]: Text(0.5, 0, 'Depth')



```
[80]: !apt-get install texlive texlive-xetex texlive-latex-extra pandoc
      !pip install py pandoc
```

```
Reading package lists... Done
Building dependency tree
Reading state information... Done
pandoc is already the newest version (1.19.2.4~dfsg-1build4).
pandoc set to manually installed.
The following additional packages will be installed:
  fonts-droid-fallback fonts-lato fonts-lmodern fonts-noto-mono fonts-texgyre
  javascript-common libcupsfilters1 libcupsimage2 libgs9 libgs9-common
  libijs-0.35 libjbig2dec0 libjs-jquery libkpathsea6 libpotrace0 libptexenc1
  libruby2.5 libsynchronetex1 libtexlua52 libtexlua52 libzzip-0-13 lmodern
  poppler-data preview-latex-style rake ruby ruby-did-you-mean ruby-minitest
  ruby-net-telnet ruby-power-assert ruby-test-unit ruby2.5
  rubygems-integration tlutils tex-common tex-gyre texlive-base
  texlive-binaries texlive-fonts-recommended texlive-latex-base
  texlive-latex-recommended texlive-pictures texlive-plain-generic tipa
Suggested packages:
  fonts-noto apatch | lighttpd | httpd poppler-utils ghostscript
  fonts-japanese-mincho | fonts-ipafont-mincho fonts-japanese-gothic
```

```
| fonts-ipafont-gothic fonts-arphic-ukai fonts-arphic-uming fonts-nanum ri
ruby-dev bundler debhelper gv | postscript-viewer perl-tk xpdf-reader
| pdf-viewer texlive-fonts-recommended-doc texlive-latex-base-doc
python-pygments icc-profiles libfile-which-perl
libspreadsheet-parseexcel-perl texlive-latex-extra-doc
texlive-latex-recommended-doc texlive-pstricks dot2tex prerex ruby-tcltk
| libtcltk-ruby texlive-pictures-doc vprerex
```

The following NEW packages will be installed:

```
fonts-droid-fallback fonts-lato fonts-lmodern fonts-noto-mono fonts-texgyre
javascript-common libcupsfilters1 libcupsimage2 libgs9 libgs9-common
libijs-0.35 libjbig2dec0 libjs-jquery libkpathsea6 libpotrace0 libptexenc1
libruby2.5 libsynchronet1 libtexlua52 libtexlua52 libzzip-0-13 lmodern
poppler-data preview-latex-style rake ruby ruby-did-you-mean ruby-minitest
ruby-net-telnet ruby-power-assert ruby-test-unit ruby2.5
rubygems-integration tlutils tex-common tex-gyre texlive texlive-base
texlive-binaries texlive-fonts-recommended texlive-latex-base
texlive-latex-extra texlive-latex-recommended texlive-pictures
texlive-plain-generic texlive-xetex tipa
```

0 upgraded, 47 newly installed, 0 to remove and 29 not upgraded.

Need to get 146 MB of archives.

After this operation, 460 MB of additional disk space will be used.

Get:1 <http://archive.ubuntu.com/ubuntu bionic/main amd64 fonts-droid-fallback>  
all 1:6.0.1r16-1.1 [1,805 kB]

Get:2 <http://archive.ubuntu.com/ubuntu bionic/main amd64 fonts-lato> all 2.0-2  
[2,698 kB]

Get:3 <http://archive.ubuntu.com/ubuntu bionic/main amd64 poppler-data> all  
0.4.8-2 [1,479 kB]

Get:4 <http://archive.ubuntu.com/ubuntu bionic/main amd64 tex-common> all 6.09  
[33.0 kB]

Get:5 <http://archive.ubuntu.com/ubuntu bionic/main amd64 fonts-lmodern> all  
2.004.5-3 [4,551 kB]

Get:6 <http://archive.ubuntu.com/ubuntu bionic/main amd64 fonts-noto-mono> all  
20171026-2 [75.5 kB]

Get:7 <http://archive.ubuntu.com/ubuntu bionic/universe amd64 fonts-texgyre> all  
20160520-1 [8,761 kB]

Get:8 <http://archive.ubuntu.com/ubuntu bionic/main amd64 javascript-common> all  
11 [6,066 B]

Get:9 <http://archive.ubuntu.com/ubuntu bionic-updates/main amd64 libcupsfilters1>  
amd64 1.20.2-0ubuntu3.1 [108 kB]

Get:10 <http://archive.ubuntu.com/ubuntu bionic-updates/main amd64 libcupsimage2>  
amd64 2.2.7-1ubuntu2.8 [18.6 kB]

Get:11 <http://archive.ubuntu.com/ubuntu bionic/main amd64 libijs-0.35> amd64  
0.35-13 [15.5 kB]

Get:12 <http://archive.ubuntu.com/ubuntu bionic/main amd64 libjbig2dec0> amd64  
0.13-6 [55.9 kB]

Get:13 <http://archive.ubuntu.com/ubuntu bionic-updates/main amd64 libgs9-common>  
all 9.26~dfsg+0-0ubuntu0.18.04.14 [5,092 kB]

Get:14 <http://archive.ubuntu.com/ubuntu bionic-updates/main amd64 libgs9> amd64

9.26~dfsg+0-0ubuntu0.18.04.14 [2,265 kB]  
Get:15 <http://archive.ubuntu.com/ubuntu> bionic/main amd64 libjs-jquery all 3.2.1-1 [152 kB]  
Get:16 <http://archive.ubuntu.com/ubuntu> bionic-updates/main amd64 libkpathsea6 amd64 2017.20170613.44572-8ubuntu0.1 [54.9 kB]  
Get:17 <http://archive.ubuntu.com/ubuntu> bionic/main amd64 libpotrace0 amd64 1.14-2 [17.4 kB]  
Get:18 <http://archive.ubuntu.com/ubuntu> bionic-updates/main amd64 libptexenc1 amd64 2017.20170613.44572-8ubuntu0.1 [34.5 kB]  
Get:19 <http://archive.ubuntu.com/ubuntu> bionic/main amd64 rubygems-integration all 1.11 [4,994 B]  
Get:20 <http://archive.ubuntu.com/ubuntu> bionic-updates/main amd64 ruby2.5 amd64 2.5.1-1ubuntu1.7 [48.6 kB]  
Get:21 <http://archive.ubuntu.com/ubuntu> bionic/main amd64 ruby amd64 1:2.5.1 [5,712 B]  
Get:22 <http://archive.ubuntu.com/ubuntu> bionic-updates/main amd64 rake all 12.3.1-1ubuntu0.1 [44.9 kB]  
Get:23 <http://archive.ubuntu.com/ubuntu> bionic/main amd64 ruby-did-you-mean all 1.2.0-2 [9,700 B]  
Get:24 <http://archive.ubuntu.com/ubuntu> bionic/main amd64 ruby-minitest all 5.10.3-1 [38.6 kB]  
Get:25 <http://archive.ubuntu.com/ubuntu> bionic/main amd64 ruby-net-telnet all 0.1.1-2 [12.6 kB]  
Get:26 <http://archive.ubuntu.com/ubuntu> bionic/main amd64 ruby-power-assert all 0.3.0-1 [7,952 B]  
Get:27 <http://archive.ubuntu.com/ubuntu> bionic/main amd64 ruby-test-unit all 3.2.5-1 [61.1 kB]  
Get:28 <http://archive.ubuntu.com/ubuntu> bionic-updates/main amd64 libruby2.5 amd64 2.5.1-1ubuntu1.7 [3,068 kB]  
Get:29 <http://archive.ubuntu.com/ubuntu> bionic-updates/main amd64 libsyntax1 amd64 2017.20170613.44572-8ubuntu0.1 [41.4 kB]  
Get:30 <http://archive.ubuntu.com/ubuntu> bionic-updates/main amd64 libtexlua52 amd64 2017.20170613.44572-8ubuntu0.1 [91.2 kB]  
Get:31 <http://archive.ubuntu.com/ubuntu> bionic-updates/main amd64 libtexluajit2 amd64 2017.20170613.44572-8ubuntu0.1 [230 kB]  
Get:32 <http://archive.ubuntu.com/ubuntu> bionic-updates/main amd64 libzip-0-13 amd64 0.13.62-3.1ubuntu0.18.04.1 [26.0 kB]  
Get:33 <http://archive.ubuntu.com/ubuntu> bionic/main amd64 lmodern all 2.004.5-3 [9,631 kB]  
Get:34 <http://archive.ubuntu.com/ubuntu> bionic/main amd64 preview-latex-style all 11.91-1ubuntu1 [185 kB]  
Get:35 <http://archive.ubuntu.com/ubuntu> bionic/main amd64 t1utils amd64 1.41-2 [56.0 kB]  
Get:36 <http://archive.ubuntu.com/ubuntu> bionic/universe amd64 tex-gyre all 20160520-1 [4,998 kB]  
Get:37 <http://archive.ubuntu.com/ubuntu> bionic-updates/main amd64 texlive-binaries amd64 2017.20170613.44572-8ubuntu0.1 [8,179 kB]  
Get:38 <http://archive.ubuntu.com/ubuntu> bionic/main amd64 texlive-base all

```

2017.20180305-1 [18.7 MB]
Get:39 http://archive.ubuntu.com/ubuntu bionic/universe amd64 texlive-fonts-
recommended all 2017.20180305-1 [5,262 kB]
Get:40 http://archive.ubuntu.com/ubuntu bionic/main amd64 texlive-latex-base all
2017.20180305-1 [951 kB]
Get:41 http://archive.ubuntu.com/ubuntu bionic/main amd64 texlive-latex-
recommended all 2017.20180305-1 [14.9 MB]
Get:42 http://archive.ubuntu.com/ubuntu bionic/universe amd64 texlive all
2017.20180305-1 [14.4 kB]
Get:43 http://archive.ubuntu.com/ubuntu bionic/universe amd64 texlive-pictures
all 2017.20180305-1 [4,026 kB]
Get:44 http://archive.ubuntu.com/ubuntu bionic/universe amd64 texlive-latex-
extra all 2017.20180305-2 [10.6 MB]
Get:45 http://archive.ubuntu.com/ubuntu bionic/universe amd64 texlive-plain-
generic all 2017.20180305-2 [23.6 MB]
Get:46 http://archive.ubuntu.com/ubuntu bionic/universe amd64 tipa all 2:1.3-20
[2,978 kB]
Get:47 http://archive.ubuntu.com/ubuntu bionic/universe amd64 texlive-xetex all
2017.20180305-1 [10.7 MB]
Fetched 146 MB in 5s (26.5 MB/s)
Extracting templates from packages: 100%
Preconfiguring packages ...
Selecting previously unselected package fonts-droid-fallback.
(Reading database ... 160975 files and directories currently installed.)
Preparing to unpack .../00-fonts-droid-fallback_1%3a6.0.1r16-1.1_all.deb ...
Unpacking fonts-droid-fallback (1:6.0.1r16-1.1) ...
Selecting previously unselected package fonts-lato.
Preparing to unpack .../01-fonts-lato_2.0-2_all.deb ...
Unpacking fonts-lato (2.0-2) ...
Selecting previously unselected package poppler-data.
Preparing to unpack .../02-poppler-data_0.4.8-2_all.deb ...
Unpacking poppler-data (0.4.8-2) ...
Selecting previously unselected package tex-common.
Preparing to unpack .../03-tex-common_6.09_all.deb ...
Unpacking tex-common (6.09) ...
Selecting previously unselected package fonts-lmodern.
Preparing to unpack .../04-fonts-lmodern_2.004.5-3_all.deb ...
Unpacking fonts-lmodern (2.004.5-3) ...
Selecting previously unselected package fonts-noto-mono.
Preparing to unpack .../05-fonts-noto-mono_20171026-2_all.deb ...
Unpacking fonts-noto-mono (20171026-2) ...
Selecting previously unselected package fonts-texgyre.
Preparing to unpack .../06-fonts-texgyre_20160520-1_all.deb ...
Unpacking fonts-texgyre (20160520-1) ...
Selecting previously unselected package javascript-common.
Preparing to unpack .../07-javascript-common_11_all.deb ...
Unpacking javascript-common (11) ...
Selecting previously unselected package libcupsfilters1:amd64.

```

```

Preparing to unpack .../08-libcupsfilters1_1.20.2-0ubuntu3.1_amd64.deb ...
Unpacking libcupsfilters1:amd64 (1.20.2-0ubuntu3.1) ...
Selecting previously unselected package libcupsimage2:amd64.
Preparing to unpack .../09-libcupsimage2_2.2.7-1ubuntu2.8_amd64.deb ...
Unpacking libcupsimage2:amd64 (2.2.7-1ubuntu2.8) ...
Selecting previously unselected package libijs-0.35:amd64.
Preparing to unpack .../10-libijs-0.35_0.35-13_amd64.deb ...
Unpacking libijs-0.35:amd64 (0.35-13) ...
Selecting previously unselected package libjbig2dec0:amd64.
Preparing to unpack .../11-libjbig2dec0_0.13-6_amd64.deb ...
Unpacking libjbig2dec0:amd64 (0.13-6) ...
Selecting previously unselected package libgs9-common.
Preparing to unpack .../12-libgs9-common_9.26~dfsg+0-0ubuntu0.18.04.14_all.deb
...
Unpacking libgs9-common (9.26~dfsg+0-0ubuntu0.18.04.14) ...
Selecting previously unselected package libgs9:amd64.
Preparing to unpack .../13-libgs9_9.26~dfsg+0-0ubuntu0.18.04.14_amd64.deb ...
Unpacking libgs9:amd64 (9.26~dfsg+0-0ubuntu0.18.04.14) ...
Selecting previously unselected package libjs-jquery.
Preparing to unpack .../14-libjs-jquery_3.2.1-1_all.deb ...
Unpacking libjs-jquery (3.2.1-1) ...
Selecting previously unselected package libkpathsea6:amd64.
Preparing to unpack .../15-libkpathsea6_2017.20170613.44572-8ubuntu0.1_amd64.deb
...
Unpacking libkpathsea6:amd64 (2017.20170613.44572-8ubuntu0.1) ...
Selecting previously unselected package libpotrace0.
Preparing to unpack .../16-libpotrace0_1.14-2_amd64.deb ...
Unpacking libpotrace0 (1.14-2) ...
Selecting previously unselected package libptexenc1:amd64.
Preparing to unpack .../17-libptexenc1_2017.20170613.44572-8ubuntu0.1_amd64.deb
...
Unpacking libptexenc1:amd64 (2017.20170613.44572-8ubuntu0.1) ...
Selecting previously unselected package rubygems-integration.
Preparing to unpack .../18-rubygems-integration_1.11_all.deb ...
Unpacking rubygems-integration (1.11) ...
Selecting previously unselected package ruby2.5.
Preparing to unpack .../19-ruby2.5_2.5.1-1ubuntu1.7_amd64.deb ...
Unpacking ruby2.5 (2.5.1-1ubuntu1.7) ...
Selecting previously unselected package ruby.
Preparing to unpack .../20-ruby_1%3a2.5.1_amd64.deb ...
Unpacking ruby (1:2.5.1) ...
Selecting previously unselected package rake.
Preparing to unpack .../21-rake_12.3.1-1ubuntu0.1_all.deb ...
Unpacking rake (12.3.1-1ubuntu0.1) ...
Selecting previously unselected package ruby-did-you-mean.
Preparing to unpack .../22-ruby-did-you-mean_1.2.0-2_all.deb ...
Unpacking ruby-did-you-mean (1.2.0-2) ...
Selecting previously unselected package ruby-minitest.

```

```

Preparing to unpack .../23-ruby-minitest_5.10.3-1_all.deb ...
Unpacking ruby-minitest (5.10.3-1) ...
Selecting previously unselected package ruby-net-telnet.
Preparing to unpack .../24-ruby-net-telnet_0.1.1-2_all.deb ...
Unpacking ruby-net-telnet (0.1.1-2) ...
Selecting previously unselected package ruby-power-assert.
Preparing to unpack .../25-ruby-power-assert_0.3.0-1_all.deb ...
Unpacking ruby-power-assert (0.3.0-1) ...
Selecting previously unselected package ruby-test-unit.
Preparing to unpack .../26-ruby-test-unit_3.2.5-1_all.deb ...
Unpacking ruby-test-unit (3.2.5-1) ...
Selecting previously unselected package libruby2.5:amd64.
Preparing to unpack .../27-libruby2.5_2.5.1-1ubuntu1.7_amd64.deb ...
Unpacking libruby2.5:amd64 (2.5.1-1ubuntu1.7) ...
Selecting previously unselected package libsyntax1:amd64.
Preparing to unpack .../28-libsyntax1_2017.20170613.44572-8ubuntu0.1_amd64.deb
...
Unpacking libsyntax1:amd64 (2017.20170613.44572-8ubuntu0.1) ...
Selecting previously unselected package libtexlua52:amd64.
Preparing to unpack .../29-libtexlua52_2017.20170613.44572-8ubuntu0.1_amd64.deb
...
Unpacking libtexlua52:amd64 (2017.20170613.44572-8ubuntu0.1) ...
Selecting previously unselected package libtexluaajit2:amd64.
Preparing to unpack
.../30-libtexluaajit2_2017.20170613.44572-8ubuntu0.1_amd64.deb ...
Unpacking libtexluaajit2:amd64 (2017.20170613.44572-8ubuntu0.1) ...
Selecting previously unselected package libzip-0-13:amd64.
Preparing to unpack .../31-libzip-0-13_0.13.62-3.1ubuntu0.18.04.1_amd64.deb ...
Unpacking libzip-0-13:amd64 (0.13.62-3.1ubuntu0.18.04.1) ...
Selecting previously unselected package lmodern.
Preparing to unpack .../32-lmodern_2.004.5-3_all.deb ...
Unpacking lmodern (2.004.5-3) ...
Selecting previously unselected package preview-latex-style.
Preparing to unpack .../33-preview-latex-style_11.91-1ubuntu1_all.deb ...
Unpacking preview-latex-style (11.91-1ubuntu1) ...
Selecting previously unselected package tlutils.
Preparing to unpack .../34-tlutils_1.41-2_amd64.deb ...
Unpacking tlutils (1.41-2) ...
Selecting previously unselected package tex-gyre.
Preparing to unpack .../35-tex-gyre_20160520-1_all.deb ...
Unpacking tex-gyre (20160520-1) ...
Selecting previously unselected package texlive-binaries.
Preparing to unpack .../36-texlive-
binaries_2017.20170613.44572-8ubuntu0.1_amd64.deb ...
Unpacking texlive-binaries (2017.20170613.44572-8ubuntu0.1) ...
Selecting previously unselected package texlive-base.
Preparing to unpack .../37-texlive-base_2017.20180305-1_all.deb ...
Unpacking texlive-base (2017.20180305-1) ...

```

```

Selecting previously unselected package texlive-fonts-recommended.
Preparing to unpack .../38-texlive-fonts-recommended_2017.20180305-1_all.deb ...
Unpacking texlive-fonts-recommended (2017.20180305-1) ...
Selecting previously unselected package texlive-latex-base.
Preparing to unpack .../39-texlive-latex-base_2017.20180305-1_all.deb ...
Unpacking texlive-latex-base (2017.20180305-1) ...
Selecting previously unselected package texlive-latex-recommended.
Preparing to unpack .../40-texlive-latex-recommended_2017.20180305-1_all.deb ...
Unpacking texlive-latex-recommended (2017.20180305-1) ...
Selecting previously unselected package texlive.
Preparing to unpack .../41-texlive_2017.20180305-1_all.deb ...
Unpacking texlive (2017.20180305-1) ...
Selecting previously unselected package texlive-pictures.
Preparing to unpack .../42-texlive-pictures_2017.20180305-1_all.deb ...
Unpacking texlive-pictures (2017.20180305-1) ...
Selecting previously unselected package texlive-latex-extra.
Preparing to unpack .../43-texlive-latex-extra_2017.20180305-2_all.deb ...
Unpacking texlive-latex-extra (2017.20180305-2) ...
Selecting previously unselected package texlive-plain-generic.
Preparing to unpack .../44-texlive-plain-generic_2017.20180305-2_all.deb ...
Unpacking texlive-plain-generic (2017.20180305-2) ...
Selecting previously unselected package tipa.
Preparing to unpack .../45-tipa_2%3a1.3-20_all.deb ...
Unpacking tipa (2:1.3-20) ...
Selecting previously unselected package texlive-xetex.
Preparing to unpack .../46-texlive-xetex_2017.20180305-1_all.deb ...
Unpacking texlive-xetex (2017.20180305-1) ...
Setting up libgs9-common (9.26~dfsg+0-0ubuntu0.18.04.14) ...
Setting up libkpathsea6:amd64 (2017.20170613.44572-8ubuntu0.1) ...
Setting up libjs-jquery (3.2.1-1) ...
Setting up libtexlua52:amd64 (2017.20170613.44572-8ubuntu0.1) ...
Setting up fonts-droid-fallback (1:6.0.1r16-1.1) ...
Setting up libsynchronet1:amd64 (2017.20170613.44572-8ubuntu0.1) ...
Setting up libptexenc1:amd64 (2017.20170613.44572-8ubuntu0.1) ...
Setting up tex-common (6.09) ...
update-language: texlive-base not installed and configured, doing nothing!
Setting up poppler-data (0.4.8-2) ...
Setting up tex-gyre (20160520-1) ...
Setting up preview-latex-style (11.91-1ubuntu1) ...
Setting up fonts-texgyre (20160520-1) ...
Setting up fonts-noto-mono (20171026-2) ...
Setting up fonts-lato (2.0-2) ...
Setting up libcupsfilters1:amd64 (1.20.2-0ubuntu3.1) ...
Setting up libcupsimage2:amd64 (2.2.7-1ubuntu2.8) ...
Setting up libjbig2dec0:amd64 (0.13-6) ...
Setting up ruby-did-you-mean (1.2.0-2) ...
Setting up tlutils (1.41-2) ...
Setting up ruby-net-telnet (0.1.1-2) ...

```

```

Setting up libijs-0.35:amd64 (0.35-13) ...
Setting up rubygems-integration (1.11) ...
Setting up libpotrace0 (1.14-2) ...
Setting up javascript-common (11) ...
Setting up ruby-minitest (5.10.3-1) ...
Setting up libzip-0-13:amd64 (0.13.62-3.1ubuntu0.18.04.1) ...
Setting up libgs9:amd64 (9.26~dfsg+0-0ubuntu0.18.04.14) ...
Setting up libtexluaajit2:amd64 (2017.20170613.44572-8ubuntu0.1) ...
Setting up fonts-lmodern (2.004.5-3) ...
Setting up ruby-power-assert (0.3.0-1) ...
Setting up texlive-binaries (2017.20170613.44572-8ubuntu0.1) ...
update-alternatives: using /usr/bin/xdvi-xaw to provide /usr/bin/xdvi.bin
(xdvi.bin) in auto mode
update-alternatives: using /usr/bin/bibtex.original to provide /usr/bin/bibtex
(bibtex) in auto mode
Setting up texlive-base (2017.20180305-1) ...
mktexlsr: Updating /var/lib/texmf/ls-R-TEXLIVEDIST...
mktexlsr: Updating /var/lib/texmf/ls-R-TEXMFMAIN...
mktexlsr: Updating /var/lib/texmf/ls-R...
mktexlsr: Done.
tl-paper: setting paper size for dvips to a4: /var/lib/texmf/dvips/config
/config-paper.ps
tl-paper: setting paper size for dvipdfmx to a4: /var/lib/texmf/dvipdfmx
/dvipdfmx-paper.cfg
tl-paper: setting paper size for xdvi to a4: /var/lib/texmf/xdvi/XDvi-paper
tl-paper: setting paper size for pdftex to a4:
/var/lib/texmf/tex/generic/config/pdftexconfig.tex
Setting up texlive-fonts-recommended (2017.20180305-1) ...
Setting up texlive-plain-generic (2017.20180305-2) ...
Setting up texlive-latex-base (2017.20180305-1) ...
Setting up lmodern (2.004.5-3) ...
Setting up texlive-latex-recommended (2017.20180305-1) ...
Setting up texlive-pictures (2017.20180305-1) ...
Setting up tipa (2:1.3-20) ...
Regenerating '/var/lib/texmf/fmtutil.cnf-DEBIAN'... done.
Regenerating '/var/lib/texmf/fmtutil.cnf-TEXLIVEDIST'... done.
update-fmtutil has updated the following file(s):
    /var/lib/texmf/fmtutil.cnf-DEBIAN
    /var/lib/texmf/fmtutil.cnf-TEXLIVEDIST
If you want to activate the changes in the above file(s),
you should run fmtutil-sys or fmtutil.
Setting up texlive (2017.20180305-1) ...
Setting up texlive-latex-extra (2017.20180305-2) ...
Setting up texlive-xetex (2017.20180305-1) ...
Setting up ruby2.5 (2.5.1-1ubuntu1.7) ...
Setting up ruby (1:2.5.1) ...
Setting up ruby-test-unit (3.2.5-1) ...
Setting up rake (12.3.1-1ubuntu0.1) ...

```



```

Setting up libruby2.5:amd64 (2.5.1-1ubuntu1.7) ...
Processing triggers for mime-support (3.60ubuntu1) ...
Processing triggers for libc-bin (2.27-3ubuntu1.2) ...
/sbin/ldconfig.real: /usr/local/lib/python3.7/dist-
packages/ideep4py/lib/libmkldnn.so.0 is not a symbolic link

Processing triggers for man-db (2.8.3-2ubuntu0.1) ...
Processing triggers for fontconfig (2.12.6-0ubuntu2) ...
Processing triggers for tex-common (6.09) ...
Running updpmap-sys. This may take some time... done.
Running mktexlsr /var/lib/texmf ... done.
Building format(s) --all.
    This may take some time... done.
Collecting py pandoc
  Downloading py pandoc-1.5.tar.gz (26 kB)
Requirement already satisfied: setuptools in /usr/local/lib/python3.7/dist-
packages (from py pandoc) (54.0.0)
Requirement already satisfied: pip>=8.1.0 in /usr/local/lib/python3.7/dist-
packages (from py pandoc) (21.0.1)
Requirement already satisfied: wheel>=0.25.0 in /usr/local/lib/python3.7/dist-
packages (from py pandoc) (0.36.2)
Building wheels for collected packages: py pandoc
  Building wheel for py pandoc (setup.py) ... done
  Created wheel for py pandoc: filename=py pandoc-1.5-py3-none-any.whl size=17037
sha256=58a2f3c37c05f4ae8ea6731bf8ff9ba8272bcb58e601253a63487fb9ff59645a
  Stored in directory: /root/.cache/pip/wheels/06/a7/1d/3259fbf0089b15d491d9166b
97de3a23d1f915bf7591abafc7
Successfully built py pandoc
Installing collected packages: py pandoc
Successfully installed py pandoc-1.5

```

```

[81]: !cp "/content/drive/MyDrive/Colab Notebooks/ml_imdb_random_forest.ipynb" ./
!jupyter nbconvert --to PDF "ml_imdb_random_forest.ipynb"

```

```

[NbConvertApp] Converting notebook ml_imdb_random_forest.ipynb to PDF
[NbConvertApp] Support files will be in ml_imdb_random_forest_files/
[NbConvertApp] Making directory ./ml_imdb_random_forest_files
[NbConvertApp] Making directory ./ml_imdb_random_forest_files
[NbConvertApp] Making directory ./ml_imdb_random_forest_files
[NbConvertApp] Writing 76905 bytes to ./notebook.tex
[NbConvertApp] Building PDF
[NbConvertApp] Running xelatex 3 times: [u'xelatex', u'./notebook.tex',
'-quiet']
[NbConvertApp] Running bibtex 1 time: [u'bibtex', u'./notebook']
[NbConvertApp] WARNING | bibtex had problems, most likely because there were no
citations
[NbConvertApp] PDF successfully created
[NbConvertApp] Writing 140939 bytes to ml_imdb_random_forest.pdf

```

```
[82]: files.download('/content/ml_imdb_random_forest.ipynb')  
files.download('/content/ml_imdb_random_forest.pdf')
```

<IPython.core.display.Javascript object>

<IPython.core.display.Javascript object>

<IPython.core.display.Javascript object>

<IPython.core.display.Javascript object>

Statement of Colaboration: Tokenization and Vectorizing techniques were taken from [this link](#).