

3_0_Lexical_Complexity_Binary_Classification_Prediction_Baseline_Model

April 6, 2025

```
[ ]: #@title Install Packages
```

```
[ ]: !pip install -q transformers  
!pip install -q torchinfo  
!pip install -q datasets  
!pip install -q evaluate  
!pip install -q nltk  
!pip install -q contractions
```

491.2/491.2 kB

10.2 MB/s eta 0:00:00

116.3/116.3 kB

10.5 MB/s eta 0:00:00

183.9/183.9 kB

16.6 MB/s eta 0:00:00

143.5/143.5 kB

12.4 MB/s eta 0:00:00

194.8/194.8 kB

17.4 MB/s eta 0:00:00

ERROR: pip's dependency resolver does not currently take into account all the packages that are installed. This behaviour is the source of the following dependency conflicts.

torch 2.6.0+cu124 requires nvidia-cublas-cu12==12.4.5.8; platform_system == "Linux" and platform_machine == "x86_64", but you have nvidia-cublas-cu12 12.5.3.2 which is incompatible.

torch 2.6.0+cu124 requires nvidia-cuda-cupti-cu12==12.4.127; platform_system == "Linux" and platform_machine == "x86_64", but you have nvidia-cuda-cupti-cu12 12.5.82 which is incompatible.

torch 2.6.0+cu124 requires nvidia-cuda-nvrtc-cu12==12.4.127; platform_system == "Linux" and platform_machine == "x86_64", but you have nvidia-cuda-nvrtc-cu12 12.5.82 which is incompatible.

torch 2.6.0+cu124 requires nvidia-cuda-runtime-cu12==12.4.127; platform_system == "Linux" and platform_machine == "x86_64", but you have nvidia-cuda-runtime-cu12 12.5.82 which is incompatible.

torch 2.6.0+cu124 requires nvidia-cudnn-cu12==9.1.0.70; platform_system == "Linux" and platform_machine == "x86_64", but you have nvidia-cudnn-cu12 9.3.0.75 which is incompatible.

torch 2.6.0+cu124 requires nvidia-cufft-cu12==11.2.1.3; platform_system == "Linux" and platform_machine == "x86_64", but you have nvidia-cufft-cu12 11.2.3.61 which is incompatible.

torch 2.6.0+cu124 requires nvidia-curand-cu12==10.3.5.147; platform_system == "Linux" and platform_machine == "x86_64", but you have nvidia-curand-cu12 10.3.6.82 which is incompatible.

torch 2.6.0+cu124 requires nvidia-cusolver-cu12==11.6.1.9; platform_system == "Linux" and platform_machine == "x86_64", but you have nvidia-cusolver-cu12 11.6.3.83 which is incompatible.

torch 2.6.0+cu124 requires nvidia-cuspars-cu12==12.3.1.170; platform_system == "Linux" and platform_machine == "x86_64", but you have nvidia-cuspars-cu12 12.5.1.3 which is incompatible.

torch 2.6.0+cu124 requires nvidia-nvjitlink-cu12==12.4.127; platform_system == "Linux" and platform_machine == "x86_64", but you have nvidia-nvjitlink-cu12 12.5.82 which is incompatible.

gcsfs 2025.3.2 requires fsspec==2025.3.2,² but you have fsspec 2024.12.0 which is incompatible.

2.6 MB/s eta 0:00:00

289.9/289.9 kB

7.0 MB/s eta 0:00:00

118.3/118.3 kB

10.7 MB/s eta 0:00:00

```
[ ]: !sudo apt-get update
      !sudo apt-get install tree
```

```
Get:1 https://cloud.r-project.org/bin/linux/ubuntu jammy-cran40/ InRelease
[3,632 B]
Hit:2 http://archive.ubuntu.com/ubuntu jammy InRelease
Get:3 http://archive.ubuntu.com/ubuntu jammy-updates InRelease [128 kB]
Get:4 https://developer.download.nvidia.com/compute/cuda/repos/ubuntu2204/x86_64
InRelease [1,581 B]
Get:5 https://r2u.stat.illinois.edu/ubuntu jammy InRelease [6,555 B]
Get:6 http://security.ubuntu.com/ubuntu jammy-security InRelease [129 kB]
Get:7 http://archive.ubuntu.com/ubuntu jammy-backports InRelease [127 kB]
Get:8 https://developer.download.nvidia.com/compute/cuda/repos/ubuntu2204/x86_64
Packages [1,381 kB]
Hit:9 https://ppa.launchpadcontent.net/deadsnakes/ppa/ubuntu jammy InRelease
Get:10 https://r2u.stat.illinois.edu/ubuntu jammy/main all Packages [8,804 kB]
Hit:11 https://ppa.launchpadcontent.net/graphics-drivers/ppa/ubuntu jammy
InRelease
Get:12 http://archive.ubuntu.com/ubuntu jammy-updates/main amd64 Packages [3,092
kB]
Hit:13 https://ppa.launchpadcontent.net/ubuntugis/ppa/ubuntu jammy InRelease
Get:14 http://archive.ubuntu.com/ubuntu jammy-updates/restricted amd64 Packages
[4,148 kB]
Get:15 https://r2u.stat.illinois.edu/ubuntu jammy/main amd64 Packages [2,683 kB]
Get:16 http://archive.ubuntu.com/ubuntu jammy-updates/universe amd64 Packages
[1,540 kB]
Get:17 http://security.ubuntu.com/ubuntu jammy-security/universe amd64 Packages
[1,241 kB]
Get:18 http://security.ubuntu.com/ubuntu jammy-security/main amd64 Packages
[2,775 kB]
Get:19 http://security.ubuntu.com/ubuntu jammy-security/restricted amd64
Packages [3,978 kB]
Fetched 30.0 MB in 2s (12.7 MB/s)
Reading package lists... Done
W: Skipping acquire of configured file 'main/source/Sources' as repository
'https://r2u.stat.illinois.edu/ubuntu jammy InRelease' does not seem to provide
it (sources.list entry misspelt?)
Reading package lists... Done
Building dependency tree... Done
Reading state information... Done
The following NEW packages will be installed:
```

```

tree
0 upgraded, 1 newly installed, 0 to remove and 45 not upgraded.
Need to get 47.9 kB of archives.
After this operation, 116 kB of additional disk space will be used.
Get:1 http://archive.ubuntu.com/ubuntu jammy/universe amd64 tree amd64 2.0.2-1
[47.9 kB]
Fetched 47.9 kB in 0s (360 kB/s)
debconf: unable to initialize frontend: Dialog
debconf: (No usable dialog-like program is installed, so the dialog based
frontend cannot be used. at /usr/share/perl5/Debconf/FrontEnd/Dialog.pm line 78,
<> line 1.)
debconf: falling back to frontend: Readline
debconf: unable to initialize frontend: Readline
debconf: (This frontend requires a controlling tty.)
debconf: falling back to frontend: Teletype
dpkg-preconfigure: unable to re-open stdin:
Selecting previously unselected package tree.
(Reading database ... 126213 files and directories currently installed.)
Preparing to unpack ../tree_2.0.2-1_amd64.deb ...
Unpacking tree (2.0.2-1) ...
Setting up tree (2.0.2-1) ...
Processing triggers for man-db (2.10.2-1) ...

```

```

[ ]: #@title Imports
import nltk
from nltk.tokenize import RegexpTokenizer

import evaluate
import transformers

import contractions

from torchinfo import summary
from datasets import load_dataset

from transformers import AutoTokenizer, AutoModel, \
    AutoModelForSequenceClassification
from transformers import TrainingArguments, Trainer

import os
import pandas as pd
import numpy as np
import matplotlib.pyplot as plt
import seaborn as sns

import sklearn

```

```
import spacy
```

```
[ ]: # @title Mount Google Drive
```

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

Mounted at /content/drive

```
[ ]: dir_root = '/content/drive/MyDrive/266-final/'  
# dir_data = '/content/drive/MyDrive/266-final/data/'  
# dir_data = '/content/drive/MyDrive/266-final/data/se21-t1-comp-lex-master/'  
dir_data = '/content/drive/MyDrive/266-final/data/266-comp-lex-master'  
dir_models = '/content/drive/MyDrive/266-final/models/'  
dir_results = '/content/drive/MyDrive/266-final/results/'
```

```
[ ]: !tree /content/drive/MyDrive/266-final/data/266-comp-lex-master/
```

```
/content/drive/MyDrive/266-final/data/266-comp-lex-master/  
  fe-test-labels  
    test_multi_df.csv  
    test_single_df.csv  
  fe-train  
    train_multi_df.csv  
    train_single_df.csv  
  fe-trial-val  
    trial_val_multi_df.csv  
    trial_val_single_df.csv  
  test-labels  
    lcp_multi_test.tsv  
    lcp_single_test.tsv  
  train  
    lcp_multi_train.tsv  
    lcp_single_train.tsv  
  trial  
    lcp_multi_trial.tsv  
    lcp_single_trial.tsv
```

6 directories, 12 files

```
[ ]: !ls -R /content/drive/MyDrive/266-final/data/266-comp-lex-master/
```

```
/content/drive/MyDrive/266-final/data/266-comp-lex-master/:  
fe-test-labels  fe-train  fe-trial-val  test-labels  train  trial
```

```
/content/drive/MyDrive/266-final/data/266-comp-lex-master/fe-test-labels:  
test_multi_df.csv  test_single_df.csv
```

```
/content/drive/MyDrive/266-final/data/266-comp-lex-master/fe-train:
```

```
train_multi_df.csv  train_single_df.csv
```

```
/content/drive/MyDrive/266-final/data/266-comp-lex-master/fe-trial-val:  
trial_val_multi_df.csv  trial_val_single_df.csv
```

```
/content/drive/MyDrive/266-final/data/266-comp-lex-master/test-labels:  
lcp_multi_test.tsv  lcp_single_test.tsv
```

```
/content/drive/MyDrive/266-final/data/266-comp-lex-master/train:  
lcp_multi_train.tsv  lcp_single_train.tsv
```

```
/content/drive/MyDrive/266-final/data/266-comp-lex-master/trial:  
lcp_multi_trial.tsv  lcp_single_trial.tsv
```

```
[ ]: #@title Import Data
```

```
[ ]: # data_dir = "/content/drive/MyDrive/266-final/data/266-comp-lex-master/"
```

```
df_names = [  
    "train_single_df",  
    "train_multi_df",  
    "trial_val_single_df",  
    "trial_val_multi_df",  
    "test_single_df",  
    "test_multi_df"  
]  
  
loaded_dataframes = {}  
  
for df_name in df_names:  
    if "train" in df_name:  
        subdir = "fe-train"  
    elif "trial_val" in df_name:  
        subdir = "fe-trial-val"  
    elif "test" in df_name:  
        subdir = "fe-test-labels"  
    else:  
        subdir = None  
  
    if subdir:  
        read_path = os.path.join(dir_data, subdir, f"{df_name}.csv")  
        loaded_df = pd.read_csv(read_path)  
        loaded_dataframes[df_name] = loaded_df  
        print(f"Loaded {df_name} from {read_path}")  
  
# Optional: quick check of loaded data  
for df_name, df in loaded_dataframes.items():
```

```
print(f"\n>>> {df_name} shape: {df.shape}")
if 'binary_complexity' in df.columns:
    print(df['binary_complexity'].value_counts())
```

```
[ ]: #@title Baseline Modeling Preparation
```

0.1 Model & Vector Designs

```
[ ]:
```