

2_0_Lexical_Complexity_Binary_Classification_Prediction_Data_Preparat

April 5, 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
```

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

```
Hit:1 https://cloud.r-project.org/bin/linux/ubuntu jammy-cran40/ InRelease
Hit:2 https://developer.download.nvidia.com/compute/cuda/repos/ubuntu2204/x86_64
InRelease
Hit:3 http://archive.ubuntu.com/ubuntu jammy InRelease
Hit:4 http://security.ubuntu.com/ubuntu jammy-security InRelease
Hit:5 https://r2u.stat.illinois.edu/ubuntu jammy InRelease
Hit:6 http://archive.ubuntu.com/ubuntu jammy-updates InRelease
Hit:7 http://archive.ubuntu.com/ubuntu jammy-backports InRelease
Hit:8 https://ppa.launchpadcontent.net/deadsnakes/ppa/ubuntu jammy InRelease
Hit:9 https://ppa.launchpadcontent.net/graphics-drivers/ppa/ubuntu jammy
InRelease
Hit:10 https://ppa.launchpadcontent.net/ubuntugis/ppa/ubuntu jammy InRelease
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
tree is already the newest version (2.0.2-1).
0 upgraded, 0 newly installed, 0 to remove and 42 not upgraded.
```

```
[ ]: #@title Imports
```

```
import transformers
```

```

import evaluate

import nltk

import contractions

from datasets import load_dataset
from torchinfo import summary

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

import os
import pandas as pd
import numpy as np

import sklearn

```

```
[ ]: # @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
  fe-train
  fe-trial-val
  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, 6 files

```
[19]: !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:

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

/content/drive/MyDrive/266-final/data/266-comp-lex-master/fe-trial-val:

/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
```

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

```
[35]: # train_single_df = pd.read_csv(os.path.join(dir_data, "train",
↳ "lcp_single_train.tsv"), sep="\t")
# train_multi_df = pd.read_csv(os.path.join(dir_data, "train", "lcp_multi_train.
↳ tsv"), sep="\t")

# trail_val_single_df = pd.read_csv(os.path.join(dir_data, "trial",
↳ "lcp_single_trial.tsv"), sep="\t")
# trail_val_multi_df = pd.read_csv(os.path.join(dir_data, "trial",
↳ "lcp_multi_trial.tsv"), sep="\t")

# test_single_df = pd.read_csv(os.path.join(dir_data, "test-labels",
↳ "lcp_single_test.tsv"), sep="\t")
# test_multi_df = pd.read_csv(os.path.join(dir_data, "test-labels",
↳ "lcp_multi_test.tsv"), sep="\t")
```

```
[ ]: ## Try to load the files containing unterminated strings
# try:
#     # Approach 1: Try with the C engine but with error handling
#     multi_test_df = pd.read_csv(
#         os.path.join(dir_data, "test", "lcp_multi_test.tsv"),
#         sep="\t",
#         on_bad_lines='skip' # Skip bad lines
#     )
```

```

#     print("Loaded with skipping bad lines")
# except Exception as e:
#     print(f"First approach failed: {e}")
#     try:
#         # Approach 2: Try with the Python engine which might be more forgiving
#         multi_test_df = pd.read_csv(
#             os.path.join(dir_data, "test", "lcp_multi_test.tsv"),
#             sep="\t",
#             engine="python",
#             quoting=3 # QUOTE_NONE
#         )
#     print("Loaded with Python engine")

```

```

[28]: # Load train data into train*_df
train_single_df = pd.read_csv(
    os.path.join(dir_data, "train", "lcp_single_train.tsv"),
    sep = "\t",
    engine = "python",
    quoting = 3
)
train_multi_df = pd.read_csv(
    os.path.join(dir_data, "train", "lcp_multi_train.tsv"),
    sep = "\t",
    engine = "python",
    quoting = 3
)

# Load trial data into trial_val*_df
trial_val_single_df = pd.read_csv(
    os.path.join(dir_data, "trial", "lcp_single_trial.tsv"),
    sep = "\t",
    engine = "python",
    quoting = 3
)
trial_val_multi_df = pd.read_csv(
    os.path.join(dir_data, "trial", "lcp_multi_trial.tsv"),
    sep = "\t",
    engine = "python",
    quoting = 3
)

# Load test data (with labels) into test*_df
test_single_df = pd.read_csv(
    os.path.join(dir_data, "test-labels", "lcp_single_test.tsv"),
    sep = "\t",
    engine = "python",
    quoting = 3
)

```

```

)
test_multi_df = pd.read_csv(
    os.path.join(dir_data, "test-labels", "lcp_multi_test.tsv"),
    sep = "\t",
    engine = "python",
    quoting = 3
)

print("Data successfully loaded into train, trial-val, and test variables")

```

Data successfully loaded into train, trial-val, and test variables

[36]: *#@title EDA*

```

[30]: import numpy as np

def print_dataframe_summary(df_name, df):
    # Print section header
    print(f"===== {df_name} =====")

    # Shape and Columns
    print(f"Shape: {df.shape}")
    print(f"Columns: {list(df.columns)}\n")

    # Data Types
    print("Data Types:")
    print(df.dtypes)
    print()

    # Missing Values
    print("Missing Values (by column):")
    print(df.isna().sum())
    print()

    # 'complexity' column stats
    desc = df['complexity'].describe() # count, mean, std, min, 25%, 50%, 75%,
↪max
    print("'complexity' Column Stats (incl. quartiles and median):")
    print(desc)

    # Calculate frequency counts for each quartile range
    q1 = desc['25%']
    q2 = desc['50%'] # This is the median
    q3 = desc['75%']
    q_max = desc['max']

    # Note: We'll define the ranges as:

```

```

# <= Q1
# > Q1 and <= Q2
# > Q2 and <= Q3
# > Q3

freq_q1 = np.sum(df['complexity'] <= q1)
freq_q2 = np.sum((df['complexity'] > q1) & (df['complexity'] <= q2))
freq_q3 = np.sum((df['complexity'] > q2) & (df['complexity'] <= q3))
freq_q4 = np.sum(df['complexity'] > q3)

print()
print("Quartile Frequency Counts (tab-separated next to each quartile):")
print(f"25%: {q1}\tCount (<= Q1): {freq_q1}")
print(f"50% (Median): {q2}\tCount (Q1 < x <= Q2): {freq_q2}")
print(f"75%: {q3}\tCount (Q2 < x <= Q3): {freq_q3}")
print(f"100% (Max): {q_max}\tCount (Q3 < x <= Max): {freq_q4}")

print("=====\n")

# Now we call this for each of our dataframes
print_dataframe_summary("train_single_df", train_single_df)
print_dataframe_summary("train_multi_df", train_multi_df)
print_dataframe_summary("trial_val_single_df", trial_val_single_df)
print_dataframe_summary("trial_val_multi_df", trial_val_multi_df)
print_dataframe_summary("test_single_df", test_single_df)
print_dataframe_summary("test_multi_df", test_multi_df)

```

===== train_single_df =====

Shape: (7662, 5)

Columns: ['id', 'corpus', 'sentence', 'token', 'complexity']

Data Types:

```

id          object
corpus      object
sentence    object
token       object
complexity  float64
dtype: object

```

Missing Values (by column):

```

id          0
corpus      0
sentence    0
token       7
complexity  0
dtype: int64

```

'complexity' Column Stats (incl. quartiles and median):

```

count      7662.000000
mean       0.302288
std        0.132977
min        0.000000
25%        0.211538
50%        0.279412
75%        0.375000
max        0.861111
Name: complexity, dtype: float64

```

```

Quartile Frequency Counts (tab-separated next to each quartile):
25%: 0.2115384615384615 Count (<= Q1): 1928
50% (Median): 0.2794117647058823 Count (Q1 < x <= Q2): 1937
75%: 0.375 Count (Q2 < x <= Q3): 1984
100% (Max): 0.8611111111111112 Count (Q3 < x <= Max): 1813
=====

```

```

===== train_multi_df =====
Shape: (1517, 5)
Columns: ['id', 'corpus', 'sentence', 'token', 'complexity']

```

```

Data Types:
id          object
corpus      object
sentence    object
token       object
complexity  float64
dtype: object

```

```

Missing Values (by column):
id          0
corpus      0
sentence    0
token       0
complexity  0
dtype: int64

```

```

'complexity' Column Stats (incl. quartiles and median):
count      1517.000000
mean       0.418362
std        0.155536
min        0.027778
25%        0.302632
50%        0.409091
75%        0.529412
max        0.975000
Name: complexity, dtype: float64

```

Quartile Frequency Counts (tab-separated next to each quartile):

25%: 0.3026315789473685 Count (<= Q1): 382

50% (Median): 0.409090909090909 Count (Q1 < x <= Q2): 377

75%: 0.5294117647058824 Count (Q2 < x <= Q3): 380

100% (Max): 0.975 Count (Q3 < x <= Max): 378

=====

===== trial_val_single_df =====

Shape: (421, 5)

Columns: ['id', 'subcorpus', 'sentence', 'token', 'complexity']

Data Types:

id object

subcorpus object

sentence object

token object

complexity float64

dtype: object

Missing Values (by column):

id 0

subcorpus 0

sentence 0

token 0

complexity 0

dtype: int64

'complexity' Column Stats (incl. quartiles and median):

count 421.000000

mean 0.298631

std 0.137619

min 0.000000

25% 0.214286

50% 0.266667

75% 0.359375

max 0.875000

Name: complexity, dtype: float64

Quartile Frequency Counts (tab-separated next to each quartile):

25%: 0.2142857142857143 Count (<= Q1): 106

50% (Median): 0.2666666666666667 Count (Q1 < x <= Q2): 107

75%: 0.359375 Count (Q2 < x <= Q3): 103

100% (Max): 0.875 Count (Q3 < x <= Max): 105

=====

===== trial_val_multi_df =====

Shape: (99, 5)

Columns: ['id', 'subcorpus', 'sentence', 'token', 'complexity']

Data Types:

```
id          object
subcorpus   object
sentence    object
token       object
complexity   float64
dtype: object
```

Missing Values (by column):

```
id          0
subcorpus   0
sentence    0
token       0
complexity   0
dtype: int64
```

'complexity' Column Stats (incl. quartiles and median):

```
count      99.000000
mean        0.417961
std         0.153752
min         0.000000
25%         0.309028
50%         0.421875
75%         0.513932
max         0.825000
```

Name: complexity, dtype: float64

Quartile Frequency Counts (tab-separated next to each quartile):

```
25%: 0.3090277777777778 Count (<= Q1): 25
50% (Median): 0.421875 Count (Q1 < x <= Q2): 25
75%: 0.5139318885448916 Count (Q2 < x <= Q3): 24
100% (Max): 0.825 Count (Q3 < x <= Max): 25
```

=====

===== test_single_df =====

Shape: (917, 5)

Columns: ['id', 'corpus', 'sentence', 'token', 'complexity']

Data Types:

```
id          object
corpus       object
sentence     object
token        object
complexity   float64
dtype: object
```

Missing Values (by column):

```
id          0
corpus      0
sentence    0
token       0
complexity  0
dtype: int64
```

'complexity' Column Stats (incl. quartiles and median):

```
count      917.000000
mean        0.296362
std         0.127290
min         0.000000
25%         0.214286
50%         0.276316
75%         0.357143
max         0.777778
```

Name: complexity, dtype: float64

Quartile Frequency Counts (tab-separated next to each quartile):

```
25%: 0.2142857142857143 Count (<= Q1): 237
50% (Median): 0.2763157894736842      Count (Q1 < x <= Q2): 224
75%: 0.3571428571428571 Count (Q2 < x <= Q3): 229
100% (Max): 0.7777777777777777 Count (Q3 < x <= Max): 227
```

=====

===== test_multi_df =====

Shape: (184, 5)

Columns: ['id', 'corpus', 'sentence', 'token', 'complexity']

Data Types:

```
id          object
corpus      object
sentence    object
token       object
complexity  float64
dtype: object
```

Missing Values (by column):

```
id          0
corpus      0
sentence    0
token       0
complexity  0
dtype: int64
```

'complexity' Column Stats (incl. quartiles and median):

```
count      184.000000
mean        0.422312
```

```
std      0.155785
min      0.000000
25%      0.316667
50%      0.428571
75%      0.527778
max      0.800000
Name: complexity, dtype: float64
```

```
Quartile Frequency Counts (tab-separated next to each quartile):
25%: 0.31666666666666666 Count (<= Q1): 47
50% (Median): 0.4285714285714286      Count (Q1 < x <= Q2): 46
75%: 0.5277777777777778 Count (Q2 < x <= Q3): 46
100% (Max): 0.8 Count (Q3 < x <= Max): 45
=====
```

```
[39]: print(train_single_df.head())
```

```

              id corpus \
0  3ZLW647WALVGE8EBR50EGUBPU4P32A  bible
1  34ROBODSP1ZBN3DVY8J8XSIY551E5C  bible
2  3S1WOPCJFGTJU2SGNAN2Y213N6WJE3  bible
3  3BFNCI9LYKQN09BHXHH9CLSX5KP738  bible
4  3G5RUKN2EC3YIWSKUXZ8ZVH95R49N2  bible

              sentence      token  complexity
0  Behold, there came up out of the river seven c...    river    0.000000
1  I am a fellow bondservant with you and with yo...  brothers    0.000000
2  The man, the lord of the land, said to us, 'By...  brothers    0.050000
3  Shimei had sixteen sons and six daughters; but...  brothers    0.150000
4              "He has put my brothers far from me.  brothers    0.263889
```

```
[40]: print(train_multi_df.head())
```

```

              id corpus \
0  3S37Y8CWI80N8KVM53U4E6JKCDC4WE  bible
1  3WGCNLZJKF877FYC1Q6COKNWDWD11  bible
2  3UQMW19E6D6WQ5TH2HDD74IVKTP5CB  bible
3  36JW4WBR06KF9AXMUL4N476OMF8FHD  bible
4  3HRWUH63QU2FH9Q8R7MRNFC7JX2N5A  bible

              sentence      token \
0  but the seventh day is a Sabbath to Yahweh you...    seventh day
1  But let each man test his own work, and then h...        own work
2  To him who by understanding made the heavens; ...  loving kindness
3  Remember to me, my God, this also, and spare m...  loving kindness
4  Because your loving kindness is better than li...  loving kindness

complexity
```

0	0.027778
1	0.050000
2	0.050000
3	0.050000
4	0.075000

[]:	<input type="text"/>
[]:	<input type="text"/>
[]:	<input type="text"/>
[]:	<input type="text"/>
[]:	<input type="text"/>
[]:	<input type="text"/>
[]:	<input type="text"/>
[]:	<input type="text"/>