

## 3.7.1-bc75-snc-single Lexical Complexity Binary Classification Prediction Transformers Modeling

April 13, 2025

### 0.1 Packages, Library Imports, File Mounts, & Data Imports \*\* Run All \*\*

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

118.3/118.3 kB

8.2 MB/s eta 0:00:00

53.8/53.8 MB

41.4 MB/s eta 0:00:00

```
[2]: !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 https://developer.download.nvidia.com/compute/cuda/repos/ubuntu2204/x86_64
InRelease
Get:3 http://security.ubuntu.com/ubuntu jammy-security InRelease [129 kB]
Hit:4 http://archive.ubuntu.com/ubuntu jammy InRelease
Get:5 http://archive.ubuntu.com/ubuntu jammy-updates InRelease [128 kB]
Get:6 https://r2u.stat.illinois.edu/ubuntu jammy InRelease [6,555 B]
Get:7 https://ppa.launchpadcontent.net/deadsnakes/ppa/ubuntu jammy InRelease
[18.1 kB]
Get:8 https://r2u.stat.illinois.edu/ubuntu jammy/main all Packages [8,833 kB]
Hit:9 https://ppa.launchpadcontent.net/graphics-drivers/ppa/ubuntu jammy
InRelease
Hit:10 https://ppa.launchpadcontent.net/ubuntuugis/ppa/ubuntu jammy InRelease
Hit:11 http://archive.ubuntu.com/ubuntu jammy-backports InRelease
Get:12 http://security.ubuntu.com/ubuntu jammy-security/universe amd64 Packages
[1,243 kB]
```

```

Get:13 https://ppa.launchpadcontent.net/deadsnakes/ppa/ubuntu jammy/main amd64
Packages [34.3 kB]
Get:14 http://archive.ubuntu.com/ubuntu jammy-updates/universe amd64 Packages
[1,542 kB]
Get:15 http://archive.ubuntu.com/ubuntu jammy-updates/main amd64 Packages [3,099
kB]
Get:16 http://security.ubuntu.com/ubuntu jammy-security/main amd64 Packages
[2,788 kB]
Get:17 https://r2u.stat.illinois.edu/ubuntu jammy/main amd64 Packages [2,690 kB]
Fetched 20.5 MB in 4s (5,545 kB/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 32 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 1s (43.7 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 ... 126315 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) ...

```

```

[3]: #@title Imports
import nltk
from nltk.tokenize import RegexpTokenizer
import sentencepiece
import contractions
import spacy

```

```

import evaluate
from datasets import load_dataset, Dataset, DatasetDict

import torch
import torch.nn as nn
from torchinfo import summary

import transformers
from transformers import AutoTokenizer, AutoModel,
    ↳AutoModelForSequenceClassification, TrainingArguments, Trainer, BertConfig,
    ↳BertForSequenceClassification

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

import sklearn
from sklearn.feature_extraction.text import TfidfVectorizer
from sklearn.naive_bayes import MultinomialNB
from sklearn.metrics import classification_report,
    ↳precision_recall_fscore_support, accuracy_score

import json
import datetime
import zoneinfo
from datetime import datetime

```

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

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

Mounted at /content/drive

```
[6]: 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/'
log_filename = "experiment_runs.txt"
log_filepath = os.path.join(dir_results, log_filename)
```

```
[7]: wandbai_api_key = ""
```

```
[8]: !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

```
[9]: !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
```

```
[10]: !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

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

```
[12]: 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}")

# 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())
#         print(df.info())
#         print(df.head())

for df_name, df in loaded_dataframes.items():
    globals()[df_name] = df
    print(f"{df_name} loaded into global namespace.")

```

```

Loaded train_single_df from /content/drive/MyDrive/266-final/data/266-comp-lex-
master/fe-train/train_single_df.csv
Loaded train_multi_df from /content/drive/MyDrive/266-final/data/266-comp-lex-
master/fe-train/train_multi_df.csv
Loaded trial_val_single_df from /content/drive/MyDrive/266-final/data/266-comp-
lex-master/fe-trial-val/trial_val_single_df.csv
Loaded trial_val_multi_df from /content/drive/MyDrive/266-final/data/266-comp-
lex-master/fe-trial-val/trial_val_multi_df.csv
Loaded test_single_df from /content/drive/MyDrive/266-final/data/266-comp-lex-
master/fe-test-labels/test_single_df.csv
Loaded test_multi_df from /content/drive/MyDrive/266-final/data/266-comp-lex-
master/fe-test-labels/test_multi_df.csv
train_single_df loaded into global namespace.
train_multi_df loaded into global namespace.
trial_val_single_df loaded into global namespace.
trial_val_multi_df loaded into global namespace.
test_single_df loaded into global namespace.
test_multi_df loaded into global namespace.

```

- Functional tests pass, we can proceed with Baseline Modeling

## 0.2 Experiments

### 0.2.1 Helper Functions **\*\* Run \*\***

```

[13]: MODEL_LINEAGE = {}

def get_model_and_tokenizer(
    remote_model_name: str = None,
    local_model_path: str = None,
    config=None
):

```

```

"""
Loads the model & tokenizer for classification.
If 'local_model_path' is specified, load from that path.
Otherwise, fall back to 'remote_model_name'.

Optional: 'config' can be a custom BertConfig/AutoConfig object
to override certain configuration parameters.

Records complete traceable lineage in the global MODEL_LINEAGE.
"""
global MODEL_LINEAGE

if local_model_path:
    print(f"Loading from local path: {local_model_path}")
    tokenizer = AutoTokenizer.from_pretrained(local_model_path)

    # If a config object is provided, we pass it to from_pretrained.
    # Otherwise, it just uses the config that is part of local_model_path.
    if config is not None:
        model = AutoModelForSequenceClassification.from_pretrained(
            local_model_path,
            config=config
        )
    else:
        model = AutoModelForSequenceClassification.
↪from_pretrained(local_model_path)

    MODEL_LINEAGE = {
        "type": "offline_checkpoint",
        "path": local_model_path,
        "timestamp": datetime.now().strftime("%Y-%m-%d %H:%M:%S")
    }
elif remote_model_name:
    print(f"Loading from Hugging Face model: {remote_model_name}")
    tokenizer = AutoTokenizer.from_pretrained(remote_model_name)

    if config is not None:
        model = AutoModelForSequenceClassification.from_pretrained(
            remote_model_name,
            config=config
        )
    else:
        model = AutoModelForSequenceClassification.
↪from_pretrained(remote_model_name)

    MODEL_LINEAGE = {
        "type": "huggingface_hub",

```

```

        "path": remote_model_name,
        "timestamp": datetime.now().strftime("%Y-%m-%d %H:%M:%S")
    }
    else:
        raise ValueError("You must provide either a remote_model_name or a local_model_path!")

    return model, tokenizer

```

```

[14]: def freeze_unfreeze_layers(model, layers_to_unfreeze=None):
    """
    Toggles requires_grad = False for all parameters
    except for those whose names contain any string in layers_to_unfreeze.
    By default, always unfreeze classifier/heads.
    """
    if layers_to_unfreeze is None:
        layers_to_unfreeze = ["classifier.", "pooler."]

    for name, param in model.named_parameters():
        if any(substring in name for substring in layers_to_unfreeze):
            param.requires_grad = True
        else:
            param.requires_grad = False

```

```

[15]: def encode_examples(examples, tokenizer, text_col, max_length=256):
    """
    Tokenizes a batch of texts from 'examples[text_col]' using the given
    tokenizer.
    Returns a dict with 'input_ids', 'attention_mask', etc.
    """
    texts = examples[text_col]
    encoded = tokenizer(
        texts,
        truncation=True,
        padding='max_length',
        max_length=max_length
    )
    return encoded

```

```

[16]: def prepare_dataset(df, tokenizer, text_col, label_col, max_length=256):
    """
    Converts a Pandas DataFrame to a Hugging Face Dataset,
    then applies 'encode_examples' to tokenize.
    """
    dataset = Dataset.from_pandas(df)

    dataset = dataset.map(

```



```

        lambda batch: encode_examples(batch, tokenizer, text_col, max_length),
        batched=True
    )

    dataset = dataset.rename_column(label_col, "labels")
    dataset.set_format(type='torch',
                       columns=['input_ids', 'attention_mask', 'labels'])

    return dataset

```

```

[17]: def compute_metrics(eval_pred):
    """
    Computes classification metrics, including accuracy, precision, recall, and
    ↪F1.
    """
    logits, labels = eval_pred
    preds = np.argmax(logits, axis=1)

    metric_accuracy = evaluate.load("accuracy")
    metric_precision = evaluate.load("precision")
    metric_recall = evaluate.load("recall")
    metric_f1 = evaluate.load("f1")

    accuracy_result = metric_accuracy.compute(predictions=preds,
    ↪references=labels)
    precision_result = metric_precision.compute(predictions=preds,
    ↪references=labels, average="binary")
    recall_result = metric_recall.compute(predictions=preds,
    ↪references=labels, average="binary")
    f1_result = metric_f1.compute(predictions=preds, references=labels,
    ↪average="binary")

    return {
        "accuracy" : accuracy_result["accuracy"],
        "precision": precision_result["precision"],
        "recall" : recall_result["recall"],
        "f1" : f1_result["f1"]
    }

```

```

[18]: def gather_config_details(model):
    """
    Enumerates every attribute in model.config
    """
    config_items = {}
    for attr_name, attr_value in vars(model.config).items():
        config_items[attr_name] = attr_value
    return config_items

```

```

def gather_model_details(model):
    """
    Extracts total layers, total params, trainable params, and activation
    ↪function
    from a Transformers model. Adjust logic as needed for different
    ↪architectures.
    """
    details = {}

    try:
        total_params = model.num_parameters()
        trainable_params = model.num_parameters(only_trainable=True)
    except AttributeError:
        all_params = list(model.parameters())
        total_params = sum(p.numel() for p in all_params)
        trainable_params = sum(p.numel() for p in all_params if p.requires_grad)

    details["model_total_params"] = total_params
    details["model_trainable_params"] = trainable_params

    if hasattr(model, "bert") and hasattr(model.bert, "pooler"):
        act_obj = getattr(model.bert.pooler, "activation", None)
        details["pooler_activation_function"] = act_obj.__class__.__name__ if
    ↪act_obj else "N/A"
    else:
        details["pooler_activation_function"] = "N/A"

    details["config_attributes"] = gather_config_details(model)
    return details

def gather_all_run_metrics(trainer, train_dataset=None, val_dataset=None,
    ↪test_dataset=None):
    """
    Gathers final training metrics, final validation metrics, final test
    ↪metrics.
    Instead of only parsing the final train_loss from the log, we also do a full
    trainer.evaluate(train_dataset) to get the same set of metrics that val/
    ↪test have.
    """
    results = {}

    if train_dataset is not None:
        train_metrics = trainer.evaluate(train_dataset)
        for k, v in train_metrics.items():
            results[f"train_{k}"] = v
    else:

```

```

        results["train_metrics"] = "No train dataset provided"

    if val_dataset is not None:
        val_metrics = trainer.evaluate(val_dataset)
        for k, v in val_metrics.items():
            results[f"val_{k}"] = v
    else:
        results["val_metrics"] = "No val dataset provided"

    if test_dataset is not None:
        test_metrics = trainer.evaluate(test_dataset)
        for k, v in test_metrics.items():
            results[f"test_{k}"] = v
    else:
        results["test_metrics"] = "No test dataset provided"

    return results

# def log_experiment_results_json(experiment_meta, model_details, run_metrics,
#                               log_file):
#     """
#     Logs experiment metadata, model details, and metrics to a JSON lines file.
#     Automatically concatenates the 'checkpoint_path' to the 'model_lineage'.
#     """
#     checkpoint_path = model_details.get("checkpoint_path")
#     if checkpoint_path:
#         if "model_lineage" not in model_details:
#             model_details["model_lineage"] = ""
#         if model_details["model_lineage"]:
#             model_details["model_lineage"] += " -> "
#         model_details["model_lineage"] += checkpoint_path

#     record = {
#         "timestamp": str(datetime.datetime.now()),
#         "experiment_meta": experiment_meta,
#         "model_details": model_details,
#         "run_metrics": run_metrics
#     }

#     with open(log_file, "a", encoding="utf-8") as f:
#         json.dump(record, f)
#         f.write("\n")

def log_experiment_results_json(experiment_meta, model_details, run_metrics,
                                log_file):
    """
    Logs experiment metadata, model details, and metrics to a JSON lines file.

```

```

Automatically concatenates the 'checkpoint_path' to the 'model_lineage'
and uses Pacific time for the timestamp.
"""
checkpoint_path = model_details.get("checkpoint_path")
if checkpoint_path:
    if "model_lineage" not in model_details:
        model_details["model_lineage"] = ""
    if model_details["model_lineage"]:
        model_details["model_lineage"] += " -> "
    model_details["model_lineage"] += checkpoint_path

    pacific_time = datetime.now(zoneinfo.ZoneInfo("America/Los_Angeles")) #_
↪update to support pacific time
    timestamp_str = pacific_time.isoformat()

    record = {
        "timestamp": timestamp_str,
        "experiment_meta": experiment_meta,
        "model_details": model_details,
        "run_metrics": run_metrics
    }

    with open(log_file, "a", encoding="utf-8") as f:
        json.dump(record, f)
        f.write("\n")

```

## 0.2.2 Experiment Cohort Design

```

[19]: # Define Experiment Parameters

named_model = "bert-base-cased"
# named_model = "roberta-base"
# named_model = "bert-large"
# named_model = "roberta-large"
# named_model = "" # modern bert

# learning_rate = 1e-3
# learning_rate = 1e-4
learning_rate = 1e-5
# learning_rate = 5e-6
# learning_rate = 5e-7
# learning_rate = 5e-8

# num_epochs = 1
# num_epochs = 3
# num_epochs = 5
num_epochs = 1

```

```

# num_epochs = 15
# num_epochs = 20

# length_max = 128
length_max = 256
# length_max = 348
# length_max = 512

# size_batch = 1
# size_batch = 4
# size_batch = 8
size_batch = 16
# size_batch = 24
# size_batch = 32
# size_batch = 64
# size_batch = 128

# regularization_weight_decay = 0
regularization_weight_decay = 0.1
# regularization_weight_decay = 0.5

y_col = "binary_complexity_75th_split"
# y_col = "binary_complexity"
# y_col = "complexity"

x_task = "single"
# x_task = "multi"

# x_col = "sentence"
x_col = "sentence_no_contractions"
# x_col = "pos_sequence"
# x_col = "dep_sequence"
# x_col = "morph_sequence"

if x_task == "single":
    df_train = train_single_df
    df_val = trial_val_single_df
    df_test = test_single_df
else:
    df_train = train_multi_df
    df_val = trial_val_multi_df
    df_test = test_multi_df

custom_config = BertConfig.from_pretrained("bert-base-cased")

custom_config.hidden_dropout_prob = 0.1
# custom_config.intermediate_size = 3072

```

```

# custom_config.intermediate_size = 6144
# custom_config.num_attention_heads = 12
# custom_config.num_hidden_layers = 12
custom_config.gradient_checkpointing = False
custom_config.attention_probs_dropout_prob = 0.1
# custom_config.max_position_embeddings = 512
# custom_config.type_vocab_size = 2
custom_config.hidden_act = "gelu" # alts: "relu" "silu"
# custom_config.vocab_size = 28996 # must match

# model.bert.pooler.activation = nn.ReLU() # Tanh() replaced as the pooler
↳ layer activation function in side-by-side with 1.1

```

config.json: 0% | 0.00/570 [00:00<?, ?B/s]

```

[20]: def train_transformer_model(
    model,
    tokenizer,
    train_dataset,
    val_dataset,
    output_dir=dir_results,
    num_epochs=num_epochs,
    batch_size=size_batch,
    lr=learning_rate,
    weight_decay=regularization_weight_decay
):
    """
    Sets up a Trainer and trains the model for 'num_epochs' using the given
    ↳ dataset.
    Returns the trained model and the Trainer object for possible re-use on
    ↳ analysis.
    """

    training_args = TrainingArguments(
        output_dir=output_dir,
        num_train_epochs=num_epochs,
        per_device_train_batch_size=batch_size,
        per_device_eval_batch_size=batch_size,
        evaluation_strategy="epoch",
        save_strategy="no",
        logging_strategy="epoch",
        learning_rate=lr,
        weight_decay=weight_decay,
        report_to=["none"], # or "wandb"
        warmup_steps=100
    )

```

```

trainer = Trainer(
    model=model,
    args=training_args,
    train_dataset=train_dataset,
    eval_dataset=val_dataset,
    tokenizer=tokenizer, # optional
    compute_metrics=compute_metrics
)

trainer.train()
return model, trainer

```

---



---

### Model Inspection \*\* Run \*\*

```

[21]: print("model checkpoints:", dir_models)
      !ls /content/drive/MyDrive/266-final/models/

```

```

model checkpoints: /content/drive/MyDrive/266-final/models/
multi_answerdotai
multi_bert-base-cased_binary_complexity_20250408_143322
multi_bert-base-cased_binary_complexity_20250409_175804
multi_bert-base-cased_binary_complexity_20250409_175954
multi_bert-base-cased_binary_complexity_20250409_180139
multi_bert-base-cased_binary_complexity_20250409_185057
multi_bert-base-cased_binary_complexity_20250409_185213
multi_bert-base-cased_binary_complexity_20250409_185333
multi_bert-base-cased_binary_complexity_20250409_234934
multi_bert-base-cased_binary_complexity_20250410_001637
multi_bert-base-cased_binary_complexity_20250410_003117
multi_bert-base-cased_binary_complexity_20250410_004527
multi_bert-base-cased_binary_complexity_20250410_025823
multi_bert-base-cased_binary_complexity_20250410_030623
multi_bert-base-cased_binary_complexity_20250410_031401
multi_bert-base-cased_binary_complexity_20250410_032138
multi_bert-base-cased_binary_complexity_20250410_034203
multi_bert-base-cased_binary_complexity_20250410_034823
multi_bert-base-cased_binary_complexity_20250410_035510
multi_bert-base-cased_binary_complexity_20250410_040140
multi_bert-base-cased_binary_complexity_20250410_174340
multi_bert-base-cased_binary_complexity_20250411_002219
multi_bert-base-cased_binary_complexity_75th_split_20250411_005437
multi_bert-large-cased_binary_complexity_20250411_002650
multi_microsoft
multi_roberta-base_binary_complexity_20250411_002307
multi_roberta-base_binary_complexity_75th_split_20250411_005524

```

multi\_roberta-large\_binary\_complexity\_20250411\_002759  
multi\_xlnet  
single\_answerdotai  
single\_bert-base-cased\_binary\_complexity\_20250408\_043117  
single\_bert-base-cased\_binary\_complexity\_20250408\_043334  
single\_bert-base-cased\_binary\_complexity\_20250408\_043750  
single\_bert-base-cased\_binary\_complexity\_20250409\_175702  
single\_bert-base-cased\_binary\_complexity\_20250409\_175900  
single\_bert-base-cased\_binary\_complexity\_20250409\_180045  
single\_bert-base-cased\_binary\_complexity\_20250409\_185027  
single\_bert-base-cased\_binary\_complexity\_20250409\_185141  
single\_bert-base-cased\_binary\_complexity\_20250409\_185303  
single\_bert-base-cased\_binary\_complexity\_20250409\_234236  
single\_bert-base-cased\_binary\_complexity\_20250410\_000508  
single\_bert-base-cased\_binary\_complexity\_20250410\_002813  
single\_bert-base-cased\_binary\_complexity\_20250410\_004230  
single\_bert-base-cased\_binary\_complexity\_20250410\_025214  
single\_bert-base-cased\_binary\_complexity\_20250410\_030435  
single\_bert-base-cased\_binary\_complexity\_20250410\_031211  
single\_bert-base-cased\_binary\_complexity\_20250410\_031404  
single\_bert-base-cased\_binary\_complexity\_20250410\_031948  
single\_bert-base-cased\_binary\_complexity\_20250410\_034334  
single\_bert-base-cased\_binary\_complexity\_20250410\_035314  
single\_bert-base-cased\_binary\_complexity\_20250410\_035940  
single\_bert-base-cased\_binary\_complexity\_20250410\_173757  
single\_bert-base-cased\_binary\_complexity\_20250410\_173911  
single\_bert-base-cased\_binary\_complexity\_20250410\_174027  
single\_bert-base-cased\_binary\_complexity\_20250410\_175501  
single\_bert-base-cased\_binary\_complexity\_20250410\_210219  
single\_bert-base-cased\_binary\_complexity\_20250410\_213212  
single\_bert-base-cased\_binary\_complexity\_20250410\_214441  
single\_bert-base-cased\_binary\_complexity\_20250410\_214546  
single\_bert-base-cased\_binary\_complexity\_20250410\_214659  
single\_bert-base-cased\_binary\_complexity\_75th\_split\_20250411\_005451  
single\_bert-large-cased\_binary\_complexity\_20250410\_215725  
single\_bert-large-cased\_binary\_complexity\_20250410\_222431  
single\_microsoft  
single\_roberta-base\_binary\_complexity\_20250410\_212304  
single\_roberta-base\_binary\_complexity\_20250410\_212514  
single\_roberta-base\_binary\_complexity\_20250410\_213732  
single\_roberta-base\_binary\_complexity\_20250410\_214805  
single\_roberta-base\_binary\_complexity\_20250410\_221944  
single\_roberta-base\_binary\_complexity\_75th\_split\_20250411\_005603  
single\_roberta-large\_binary\_complexity\_20250410\_221054  
single\_roberta-large\_binary\_complexity\_20250410\_222652  
single\_roberta-large\_binary\_complexity\_20250410\_223030  
single\_roberta-large\_binary\_complexity\_20250410\_223320  
single\_roberta-large\_binary\_complexity\_20250410\_223754



single\_xlnet

```
[22]: # Load Model & Tokenizer
# model, tokenizer = get_model_and_tokenizer(named_model) # deprecated argument
# structure
# model, tokenizer = get_model_and_tokenizer("/content/drive/MyDrive/266-final/
# models/...") # proposed argument usage for checkpointed models

# for name, param in model.named_parameters():
#     print(name)

model, tokenizer = get_model_and_tokenizer(
    remote_model_name="bert-base-cased",
    local_model_path=None,
    config=custom_config
)

# model, tokenizer = get_model_and_tokenizer(
#     local_model_path="my_local_bert_path",
#     config=custom_config
# )

print("=====")
print(named_model, ":")
print("=====")
# print(model)
print("=====")
print(model.config)
print("=====")
print("num_parameters:", model.num_parameters())
print("=====")
print("num_trainable_parameters:", model.num_parameters(only_trainable=True))
```

Loading from Hugging Face model: bert-base-cased

tokenizer\_config.json: 0%| | 0.00/49.0 [00:00<?, ?B/s]

vocab.txt: 0%| | 0.00/213k [00:00<?, ?B/s]

tokenizer.json: 0%| | 0.00/436k [00:00<?, ?B/s]

model.safetensors: 0%| | 0.00/436M [00:00<?, ?B/s]

Some weights of BertForSequenceClassification were not initialized from the model checkpoint at bert-base-cased and are newly initialized:

['classifier.bias', 'classifier.weight']

You should probably TRAIN this model on a down-stream task to be able to use it for predictions and inference.

=====  
bert-base-cased :

```

=====
=====
BertConfig {
  "_attn_implementation_autoset": true,
  "architectures": [
    "BertForMaskedLM"
  ],
  "attention_probs_dropout_prob": 0.1,
  "classifier_dropout": null,
  "gradient_checkpointing": false,
  "hidden_act": "gelu",
  "hidden_dropout_prob": 0.1,
  "hidden_size": 768,
  "initializer_range": 0.02,
  "intermediate_size": 3072,
  "layer_norm_eps": 1e-12,
  "max_position_embeddings": 512,
  "model_type": "bert",
  "num_attention_heads": 12,
  "num_hidden_layers": 12,
  "pad_token_id": 0,
  "position_embedding_type": "absolute",
  "torch_dtype": "float32",
  "transformers_version": "4.50.3",
  "type_vocab_size": 2,
  "use_cache": true,
  "vocab_size": 28996
}

```

```

=====
num_parameters: 108311810
=====
num_trainable_parameters: 108311810

```

### Layer Configuration \*\* Run \*\*

[23]: *# Freeze/Unfreeze Layers & Additional Activation Function Configuration*

```

layers_to_unfreeze = [
  # "bert.embeddings.",
  # "bert.encoder.layer.0.",
  # "bert.encoder.layer.1.",
  # "bert.encoder.layer.8.",
  # "bert.encoder.layer.9.",
  # "bert.encoder.layer.10.",
  "bert.encoder.layer.11.",
  "bert.pooler.",
  "classifier.",

```

```

]

freeze_unfreeze_layers(model, layers_to_unfreeze=layers_to_unfreeze)

for name, param in model.named_parameters():
    print(name, "requires_grad=", param.requires_grad)

print("\nLayers that are 'True' are trainable. 'False' are frozen.")

print("=====")
print(named_model, ":")
print("=====")
# print(model)
print("=====")
print(model.config)
print("=====")
print("num_parameters:", model.num_parameters())
print("=====")
print("num_trainable_parameters:", model.num_parameters(only_trainable=True))

```

```

bert.embeddings.word_embeddings.weight requires_grad= False
bert.embeddings.position_embeddings.weight requires_grad= False
bert.embeddings.token_type_embeddings.weight requires_grad= False
bert.embeddings.LayerNorm.weight requires_grad= False
bert.embeddings.LayerNorm.bias requires_grad= False
bert.encoder.layer.0.attention.self.query.weight requires_grad= False
bert.encoder.layer.0.attention.self.query.bias requires_grad= False
bert.encoder.layer.0.attention.self.key.weight requires_grad= False
bert.encoder.layer.0.attention.self.key.bias requires_grad= False
bert.encoder.layer.0.attention.self.value.weight requires_grad= False
bert.encoder.layer.0.attention.self.value.bias requires_grad= False
bert.encoder.layer.0.attention.output.dense.weight requires_grad= False
bert.encoder.layer.0.attention.output.dense.bias requires_grad= False
bert.encoder.layer.0.attention.output.LayerNorm.weight requires_grad= False
bert.encoder.layer.0.attention.output.LayerNorm.bias requires_grad= False
bert.encoder.layer.0.intermediate.dense.weight requires_grad= False
bert.encoder.layer.0.intermediate.dense.bias requires_grad= False
bert.encoder.layer.0.output.dense.weight requires_grad= False
bert.encoder.layer.0.output.dense.bias requires_grad= False
bert.encoder.layer.0.output.LayerNorm.weight requires_grad= False
bert.encoder.layer.0.output.LayerNorm.bias requires_grad= False
bert.encoder.layer.1.attention.self.query.weight requires_grad= False
bert.encoder.layer.1.attention.self.query.bias requires_grad= False
bert.encoder.layer.1.attention.self.key.weight requires_grad= False
bert.encoder.layer.1.attention.self.key.bias requires_grad= False
bert.encoder.layer.1.attention.self.value.weight requires_grad= False
bert.encoder.layer.1.attention.self.value.bias requires_grad= False
bert.encoder.layer.1.attention.output.dense.weight requires_grad= False

```

[illegible]





```

bert.encoder.layer.10.attention.output.dense.bias requires_grad= False
bert.encoder.layer.10.attention.output.LayerNorm.weight requires_grad= False
bert.encoder.layer.10.attention.output.LayerNorm.bias requires_grad= False
bert.encoder.layer.10.intermediate.dense.weight requires_grad= False
bert.encoder.layer.10.intermediate.dense.bias requires_grad= False
bert.encoder.layer.10.output.dense.weight requires_grad= False
bert.encoder.layer.10.output.dense.bias requires_grad= False
bert.encoder.layer.10.output.LayerNorm.weight requires_grad= False
bert.encoder.layer.10.output.LayerNorm.bias requires_grad= False
bert.encoder.layer.11.attention.self.query.weight requires_grad= True
bert.encoder.layer.11.attention.self.query.bias requires_grad= True
bert.encoder.layer.11.attention.self.key.weight requires_grad= True
bert.encoder.layer.11.attention.self.key.bias requires_grad= True
bert.encoder.layer.11.attention.self.value.weight requires_grad= True
bert.encoder.layer.11.attention.self.value.bias requires_grad= True
bert.encoder.layer.11.attention.output.dense.weight requires_grad= True
bert.encoder.layer.11.attention.output.dense.bias requires_grad= True
bert.encoder.layer.11.attention.output.LayerNorm.weight requires_grad= True
bert.encoder.layer.11.attention.output.LayerNorm.bias requires_grad= True
bert.encoder.layer.11.intermediate.dense.weight requires_grad= True
bert.encoder.layer.11.intermediate.dense.bias requires_grad= True
bert.encoder.layer.11.output.dense.weight requires_grad= True
bert.encoder.layer.11.output.dense.bias requires_grad= True
bert.encoder.layer.11.output.LayerNorm.weight requires_grad= True
bert.encoder.layer.11.output.LayerNorm.bias requires_grad= True
bert.pooler.dense.weight requires_grad= True
bert.pooler.dense.bias requires_grad= True
classifier.weight requires_grad= True
classifier.bias requires_grad= True

```

Layers that are 'True' are trainable. 'False' are frozen.

=====

bert-base-cased :

=====

=====

```

BertConfig {
  "_attn_implementation_autoset": true,
  "architectures": [
    "BertForMaskedLM"
  ],
  "attention_probs_dropout_prob": 0.1,
  "classifier_dropout": null,
  "gradient_checkpointing": false,
  "hidden_act": "gelu",
  "hidden_dropout_prob": 0.1,
  "hidden_size": 768,
  "initializer_range": 0.02,
  "intermediate_size": 3072,

```

```

"layer_norm_eps": 1e-12,
"max_position_embeddings": 512,
"model_type": "bert",
"num_attention_heads": 12,
"num_hidden_layers": 12,
"pad_token_id": 0,
"position_embedding_type": "absolute",
"torch_dtype": "float32",
"transformers_version": "4.50.3",
"type_vocab_size": 2,
"use_cache": true,
"vocab_size": 28996
}

```

```

=====
num_parameters: 108311810
=====
num_trainable_parameters: 7680002

```

#### Dataset Preparation \*\* Run \*\*

```

[24]: # Tokenize & Prepare Datasets

train_data_hf = prepare_dataset(
    df_train,
    tokenizer,
    text_col=x_col,
    label_col=y_col,
    max_length=length_max
)

val_data_hf = prepare_dataset(
    df_val,
    tokenizer,
    text_col=x_col,
    label_col=y_col,
    max_length=length_max
)

test_data_hf = prepare_dataset(
    df_test,
    tokenizer,
    text_col=x_col,
    label_col=y_col,
    max_length=length_max
)

```



```
print("Datasets prepared. Sample from train_data_hf:\n", train_data_hf[10])
# print("Datasets prepared. Sample from train_data_hf:\n", val_data_hf[10])
# print("Datasets prepared. Sample from train_data_hf:\n", test_data_hf[10])
```

```
Map: 0%|          | 0/7662 [00:00<?, ? examples/s]
```

```
Map: 0%|          | 0/421 [00:00<?, ? examples/s]
```

```
Map: 0%|          | 0/917 [00:00<?, ? examples/s]
```

Datasets prepared. Sample from train\_data\_hf:

[illegible]

0.2.3 snc regularization\_weight\_decay = 0.5 learning\_rate = 5e-6 size\_batch = 128  
length\_max = 128 num\_epochs = 1

```
[ ]: # Define Experiment Parameters
named_model = "bert-base-cased"
# named_model = "roberta-base"
# named_model = "bert-large"
# named_model = "roberta-large"
# named_model = "" # modern bert
#####
regularization_weight_decay = 0.5
learning_rate = 5e-6
size_batch = 128
length_max = 128
num_epochs = 1
#####
# x_col = "sentence"
x_col = "sentence_no_contractions"
# x_col = "pos_sequence"
# x_col = "dep_sequence"
# x_col = "morph_sequence"
#####
y_col = "binary_complexity_75th_split"
# y_col = "binary_complexity"
# y_col = "complexity"
#####
x_task = "single"
# x_task = "multi"
if x_task == "single":
    df_train = train_single_df
    df_val = trial_val_single_df
    df_test = test_single_df
else:
    df_train = train_multi_df
    df_val = trial_val_multi_df
    df_test = test_multi_df
#####
# Tokenize & Prepare Datasets
train_data_hf = prepare_dataset(
    df_train,
    tokenizer,
    text_col=x_col,
    label_col=y_col,
    max_length=length_max)
val_data_hf = prepare_dataset(
    df_val,
    tokenizer,
```

```

        text_col=x_col,
        label_col=y_col,
        max_length=length_max)
test_data_hf = prepare_dataset(
    df_test,
    tokenizer,
    text_col=x_col,
    label_col=y_col,
    max_length=length_max)
print("Datasets prepared. Sample from train_data_hf:\n", train_data_hf[10])
# print("Datasets prepared. Sample from train_data_hf:\n", val_data_hf[10])
# print("Datasets prepared. Sample from train_data_hf:\n", test_data_hf[10])
#####
# custom_config = BertConfig.from_pretrained("bert-base-cased")
# custom_config.hidden_act = "gelu" # alts: "relu" "silu"
# custom_config.attention_probs_dropout_prob = 0.1
# custom_config.hidden_dropout_prob = 0.1
# custom_config.gradient_checkpointing = False
#####
model, tokenizer = get_model_and_tokenizer(
    remote_model_name="bert-base-cased",
    local_model_path=None,
    config=None)
#####
# model, tokenizer = get_model_and_tokenizer(
#     remote_model_name=None
#     local_model_path="...CONFIGURE_PATH...",
#     config=custom_config)
print("=====")
print(named_model, ":")
print("=====")
print("num_parameters:", model.num_parameters())
print("num_trainable_parameters at load:", model.
    ↪ num_parameters(only_trainable=True))
print("=====")
print("model lineage:", MODEL_LINEAGE)
print("=====")
#####
layers_to_unfreeze = [
    # "bert.embeddings.",
    # "bert.encoder.layer.0.",
    # "bert.encoder.layer.1.",
    # "bert.encoder.layer.8.",
    # "bert.encoder.layer.9.",
    # "bert.encoder.layer.10.",
    "bert.encoder.layer.11.",
    "bert.pooler.",

```

```

        "classifier.",
    ]
    freeze_unfreeze_layers(model, layers_to_unfreeze=layers_to_unfreeze)
    print(model.config)
    print("=====")
    print("num_parameters:", model.num_parameters())
    print("num_trainable_parameters:", model.num_parameters(only_trainable=True))
    print("=====")
    #####
    print("Experiment configuration used with this experiment:")
    print("model used:", named_model)
    print("learning rate used:", learning_rate)
    print("number of epochs:", num_epochs)
    print("maximum sequence length:", length_max)
    print("batch size used:", size_batch)
    print("regularization value:", regularization_weight_decay)
    print("outcome variable:", y_col)
    print("task:", x_task)
    print("input column:", x_col)

```

```

[ ]: # Train & Evaluate
trained_model, trainer_obj = train_transformer_model(
    model = model,
    tokenizer = tokenizer,
    train_dataset = train_data_hf,
    val_dataset = val_data_hf,
    output_dir = dir_results,
    num_epochs = num_epochs,
    batch_size = size_batch,
    lr = learning_rate,
    weight_decay = regularization_weight_decay)
metrics = trainer_obj.evaluate()
print("Validation metrics:", metrics)
test_metrics = trainer_obj.evaluate(test_data_hf) if test_data_hf else None
print("Test metrics:", test_metrics)

```

```

[ ]: # save model checkpoint
timestamp = datetime.now().strftime("%Y%m%d_%H%M%S")
model_save_path = os.path.join(dir_models,
    f"{x_task}_{named_model}_{y_col}_{timestamp}")
trainer_obj.save_model(model_save_path)
print(f"Model checkpoint saved to: {model_save_path}")
# log experiment results
experiment_info = {
    "model_name": named_model,
    "learning_rate": learning_rate,
    "epochs": num_epochs,

```

```

        "batch_size": size_batch,
        "weight_decay": regularization_weight_decay,
        "x_task": x_task,
        "x_col": x_col,
        "y_col": y_col,
        "layers_to_unfreeze": layers_to_unfreeze}
model_info = gather_model_details(trained_model)
all_run_metrics = gather_all_run_metrics(
    trainer=trainer_obj,
    train_dataset=train_data_hf,
    val_dataset=val_data_hf,
    test_dataset=test_data_hf)
log_experiment_results_json(
    experiment_meta=experiment_info,
    model_details=model_info,
    run_metrics=all_run_metrics,
    log_file=log_filepath)
print(f"EXPERIMENT LOGGED TO: {log_filepath}")

```

0.2.4 snc roberta-base regularization\_weight\_decay = 0.5 learning\_rate = 5e-6  
size\_batch = 128 length\_max = 128 num\_epochs = 1

```

[ ]: # Define Experiment Parameters
# named_model = "bert-base-cased"
named_model = "roberta-base"
# named_model = "bert-large"
# named_model = "roberta-large"
# named_model = "" # modern bert
#####
regularization_weight_decay = 0.5
learning_rate = 5e-6
size_batch = 128
length_max = 128
num_epochs = 1
#####
# x_col = "sentence"
x_col = "sentence_no_contractions"
# x_col = "pos_sequence"
# x_col = "dep_sequence"
# x_col = "morph_sequence"
#####
y_col = "binary_complexity_75th_split"
# y_col = "binary_complexity"
# y_col = "complexity"
#####
x_task = "single"
# x_task = "multi"

```

```

if x_task == "single":
    df_train = train_single_df
    df_val = trial_val_single_df
    df_test = test_single_df
else:
    df_train = train_multi_df
    df_val = trial_val_multi_df
    df_test = test_multi_df
#####
# Tokenize & Prepare Datasets
train_data_hf = prepare_dataset(
    df_train,
    tokenizer,
    text_col=x_col,
    label_col=y_col,
    max_length=length_max)
val_data_hf = prepare_dataset(
    df_val,
    tokenizer,
    text_col=x_col,
    label_col=y_col,
    max_length=length_max)
test_data_hf = prepare_dataset(
    df_test,
    tokenizer,
    text_col=x_col,
    label_col=y_col,
    max_length=length_max)
print("Datasets prepared. Sample from train_data_hf:\n", train_data_hf[10])
# print("Datasets prepared. Sample from train_data_hf:\n", val_data_hf[10])
# print("Datasets prepared. Sample from train_data_hf:\n", test_data_hf[10])
#####
# custom_config = BertConfig.from_pretrained("roberta-base")
# custom_config.hidden_act = "gelu" # alts: "relu" "silu"
# custom_config.attention_probs_dropout_prob = 0.1
# custom_config.hidden_dropout_prob = 0.1
# custom_config.gradient_checkpointing = False
#####
model, tokenizer = get_model_and_tokenizer(
    remote_model_name="roberta-base",
    local_model_path=None,
    config=None)
#####
# model, tokenizer = get_model_and_tokenizer(
#     remote_model_name=None
#     local_model_path="...CONFIGURE_PATH...",
#     config=custom_config)

```

```

print("=====")
print(named_model, ":")
print("=====")
print("num_parameters:", model.num_parameters())
print("num_trainable_parameters at load:", model.
    ↪num_parameters(only_trainable=True))
print("=====")
print("model lineage:", MODEL_LINEAGE)
print("=====")
#####

```

```
[ ]: print(model)
```

```
[ ]: for name, param in model.named_parameters():
    print(name, "requires_grad=", param.requires_grad)
```

```
[ ]: # Inspect the attention_mask tensor for the first few samples
for i in range(5):
    print(train_data_hf[i]['attention_mask'])
```

```
[ ]: layers_to_unfreeze = [
    "roberta.encoder.layer.11.attention.self.query.weight",
    "roberta.encoder.layer.11.attention.self.query.bias",
    "roberta.encoder.layer.11.attention.self.key.weight",
    "roberta.encoder.layer.11.attention.self.key.bias",
    "roberta.encoder.layer.11.attention.self.value.weight",
    "roberta.encoder.layer.11.attention.self.value.bias",
    "roberta.encoder.layer.11.attention.output.dense.weight",
    "roberta.encoder.layer.11.attention.output.dense.bias",
    "roberta.encoder.layer.11.attention.output.LayerNorm.weight",
    "roberta.encoder.layer.11.attention.output.LayerNorm.bias",
    "roberta.encoder.layer.11.intermediate.dense.weight",
    "roberta.encoder.layer.11.intermediate.dense.bias",
    "roberta.encoder.layer.11.output.dense.weight",
    "roberta.encoder.layer.11.output.dense.bias",
    "roberta.encoder.layer.11.output.LayerNorm.weight",
    "roberta.encoder.layer.11.output.LayerNorm.bias",
    "classifier.dense.weight",
    "classifier.dense.bias",
    "classifier.out_proj.weight",
    "classifier.out_proj.bias"
]
freeze_unfreeze_layers(model, layers_to_unfreeze=layers_to_unfreeze)
print(model.config)
print("=====")
print("num_parameters:", model.num_parameters())
print("num_trainable_parameters:", model.num_parameters(only_trainable=True))

```

```

print("=====")
#####
print("Experiment configuration used with this experiment:")
print("model used:", named_model)
print("learning rate used:", learning_rate)
print("number of epochs:", num_epochs)
print("maximum sequence length:", length_max)
print("batch size used:", size_batch)
print("regularization value:", regularization_weight_decay)
print("outcome variable:", y_col)
print("task:", x_task)
print("input column:", x_col)

```

```

[ ]: for name, param in model.named_parameters():
    print(name, "requires_grad=", param.requires_grad)

```

```

[ ]: # Train & Evaluate
trained_model, trainer_obj = train_transformer_model(
    model = model,
    tokenizer = tokenizer,
    train_dataset = train_data_hf,
    val_dataset = val_data_hf,
    output_dir = dir_results,
    num_epochs = num_epochs,
    batch_size = size_batch,
    lr = learning_rate,
    weight_decay = regularization_weight_decay)
metrics = trainer_obj.evaluate()
print("Validation metrics:", metrics)
test_metrics = trainer_obj.evaluate(test_data_hf) if test_data_hf else None
print("Test metrics:", test_metrics)

```

```

[ ]: # save model checkpoint
timestamp = datetime.now().strftime("%Y%m%d_%H%M%S")
model_save_path = os.path.join(dir_models,
    f"{x_task}_{named_model}_{y_col}_{timestamp}")
trainer_obj.save_model(model_save_path)
print(f"Model checkpoint saved to: {model_save_path}")
# log experiment results
experiment_info = {
    "model_name": named_model,
    "learning_rate": learning_rate,
    "epochs": num_epochs,
    "batch_size": size_batch,
    "weight_decay": regularization_weight_decay,
    "x_task": x_task,
    "x_col": x_col,
}

```



```

        "y_col": y_col,
        "layers_to_unfreeze": layers_to_unfreeze}
model_info = gather_model_details(trained_model)
all_run_metrics = gather_all_run_metrics(
    trainer=trainer_obj,
    train_dataset=train_data_hf,
    val_dataset=val_data_hf,
    test_dataset=test_data_hf)
log_experiment_results_json(
    experiment_meta=experiment_info,
    model_details=model_info,
    run_metrics=all_run_metrics,
    log_file=log_filepath)
print(f"EXPERIMENT LOGGED TO: {log_filepath}")

```

0.2.5 snc bert-large-cased regularization\_weight\_decay = 0.5 learning\_rate = 5e-6  
size\_batch = 128 length\_max = 128 num\_epochs = 1

```

[25]: # Define Experiment Parameters
# named_model = "bert-base-cased"
# named_model = "roberta-base"
named_model = "bert-large-cased"
# named_model = "roberta-large"
# named_model = "" # modern bert
#####
regularization_weight_decay = 0.5
learning_rate = 5e-6
size_batch = 128
length_max = 128
num_epochs = 1
#####
# x_col = "sentence"
x_col = "sentence_no_contractions"
# x_col = "pos_sequence"
# x_col = "dep_sequence"
# x_col = "morph_sequence"
#####
y_col = "binary_complexity_75th_split"
# y_col = "binary_complexity"
# y_col = "complexity"
#####
x_task = "single"
# x_task = "multi"
if x_task == "single":
    df_train = train_single_df
    df_val = trial_val_single_df
    df_test = test_single_df

```

```

else:
    df_train = train_multi_df
    df_val = trial_val_multi_df
    df_test = test_multi_df
#####
# Tokenize & Prepare Datasets
train_data_hf = prepare_dataset(
    df_train,
    tokenizer,
    text_col=x_col,
    label_col=y_col,
    max_length=length_max)
val_data_hf = prepare_dataset(
    df_val,
    tokenizer,
    text_col=x_col,
    label_col=y_col,
    max_length=length_max)
test_data_hf = prepare_dataset(
    df_test,
    tokenizer,
    text_col=x_col,
    label_col=y_col,
    max_length=length_max)
print("Datasets prepared. Sample from train_data_hf:\n", train_data_hf[10])
# print("Datasets prepared. Sample from train_data_hf:\n", val_data_hf[10])
# print("Datasets prepared. Sample from train_data_hf:\n", test_data_hf[10])
#####
# custom_config = BertConfig.from_pretrained("roberta-base")
# custom_config.hidden_act = "gelu" # alts: "relu" "silu"
# custom_config.attention_probs_dropout_prob = 0.1
# custom_config.hidden_dropout_prob = 0.1
# custom_config.gradient_checkpointing = False
#####
model, tokenizer = get_model_and_tokenizer(
    remote_model_name="bert-large-cased",
    local_model_path=None,
    config=None)
#####
# model, tokenizer = get_model_and_tokenizer(
#     remote_model_name=None
#     local_model_path="...CONFIGURE_PATH...",
#     config=custom_config)
print("=====")
print(named_model, ":")
print("=====")
print("num_parameters:", model.num_parameters())

```

```

print("num_trainable_parameters at load:", model.
      ↪num_parameters(only_trainable=True))
print("=====")
print("model lineage:", MODEL_LINEAGE)
print("=====")

```

Map: 0%| | 0/7662 [00:00<?, ? examples/s]

Map: 0%| | 0/421 [00:00<?, ? examples/s]

Map: 0%| | 0/917 [00:00<?, ? examples/s]

Datasets prepared. Sample from train\_data\_hf:

```

{'labels': tensor(0), 'input_ids': tensor([ 101, 1252, 1106, 1103, 3824,
1104, 19892, 11220, 1324, 1119,
      1522, 3839, 117, 1272, 1103, 1555, 1104, 1103, 11563, 5609,
1106, 1172, 132, 1152, 2446, 1122, 1113, 1147, 3221, 119,
      102, 0, 0, 0, 0, 0, 0, 0, 0, 0,
      0, 0, 0, 0, 0, 0, 0, 0, 0, 0,
      0, 0, 0, 0, 0, 0, 0, 0, 0, 0,
      0, 0, 0, 0, 0, 0, 0, 0, 0, 0,
      0, 0, 0, 0, 0, 0, 0, 0, 0, 0,
      0, 0, 0, 0, 0, 0, 0, 0, 0, 0,
      0, 0, 0, 0, 0, 0, 0, 0, 0, 0,
      0, 0, 0, 0, 0, 0, 0, 0, 0, 0,
      0, 0, 0, 0, 0, 0, 0, 0]),
'attention_mask': tensor([1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1,
1, 1, 1, 1, 1,
      1, 1, 1, 1, 1, 1, 1, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0,
      0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0,
      0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0,
      0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0,
      0, 0, 0, 0, 0, 0, 0, 0])}

```

Loading from Hugging Face model: bert-large-cased

tokenizer\_config.json: 0%| | 0.00/49.0 [00:00<?, ?B/s]

config.json: 0%| | 0.00/762 [00:00<?, ?B/s]

vocab.txt: 0%| | 0.00/213k [00:00<?, ?B/s]

tokenizer.json: 0%| | 0.00/436k [00:00<?, ?B/s]

model.safetensors: 0%| | 0.00/1.34G [00:00<?, ?B/s]

Some weights of BertForSequenceClassification were not initialized from the model checkpoint at bert-large-cased and are newly initialized:

['classifier.bias', 'classifier.weight']

You should probably TRAIN this model on a down-stream task to be able to use it for predictions and inference.

```

=====
bert-large-cased :
=====
num_parameters: 333581314
num_trainable_parameters at load: 333581314
=====
model lineage: {'type': 'huggingface_hub', 'path': 'bert-large-cased',
'timestamp': '2025-04-11 01:01:59'}
=====

```

```
[26]: print(model)
```

```

BertForSequenceClassification(
  (bert): BertModel(
    (embeddings): BertEmbeddings(
      (word_embeddings): Embedding(28996, 1024, padding_idx=0)
      (position_embeddings): Embedding(512, 1024)
      (token_type_embeddings): Embedding(2, 1024)
      (LayerNorm): LayerNorm((1024,), eps=1e-12, elementwise_affine=True)
      (dropout): Dropout(p=0.1, inplace=False)
    )
    (encoder): BertEncoder(
      (layer): ModuleList(
        (0-23): 24 x BertLayer(
          (attention): BertAttention(
            (self): BertSdpaSelfAttention(
              (query): Linear(in_features=1024, out_features=1024, bias=True)
              (key): Linear(in_features=1024, out_features=1024, bias=True)
              (value): Linear(in_features=1024, out_features=1024, bias=True)
              (dropout): Dropout(p=0.1, inplace=False)
            )
            (output): BertSelfOutput(
              (dense): Linear(in_features=1024, out_features=1024, bias=True)
              (LayerNorm): LayerNorm((1024,), eps=1e-12,
elementwise_affine=True)
              (dropout): Dropout(p=0.1, inplace=False)
            )
          )
          (intermediate): BertIntermediate(
            (dense): Linear(in_features=1024, out_features=4096, bias=True)
            (intermediate_act_fn): GELUActivation()
          )
          (output): BertOutput(
            (dense): Linear(in_features=4096, out_features=1024, bias=True)
            (LayerNorm): LayerNorm((1024,), eps=1e-12, elementwise_affine=True)
            (dropout): Dropout(p=0.1, inplace=False)
          )
        )
      )
    )
  )
)

```

```

    )
)
(pooler): BertPooler(
  (dense): Linear(in_features=1024, out_features=1024, bias=True)
  (activation): Tanh()
)
)
(dropout): Dropout(p=0.1, inplace=False)
(classifier): Linear(in_features=1024, out_features=2, bias=True)
)

```

```

[27]: for name, param in model.named_parameters():
      print(name, "requires_grad=", param.requires_grad)

```

```

bert.embeddings.word_embeddings.weight requires_grad= True
bert.embeddings.position_embeddings.weight requires_grad= True
bert.embeddings.token_type_embeddings.weight requires_grad= True
bert.embeddings.LayerNorm.weight requires_grad= True
bert.embeddings.LayerNorm.bias requires_grad= True
bert.encoder.layer.0.attention.self.query.weight requires_grad= True
bert.encoder.layer.0.attention.self.query.bias requires_grad= True
bert.encoder.layer.0.attention.self.key.weight requires_grad= True
bert.encoder.layer.0.attention.self.key.bias requires_grad= True
bert.encoder.layer.0.attention.self.value.weight requires_grad= True
bert.encoder.layer.0.attention.self.value.bias requires_grad= True
bert.encoder.layer.0.attention.output.dense.weight requires_grad= True
bert.encoder.layer.0.attention.output.dense.bias requires_grad= True
bert.encoder.layer.0.attention.output.LayerNorm.weight requires_grad= True
bert.encoder.layer.0.attention.output.LayerNorm.bias requires_grad= True
bert.encoder.layer.0.intermediate.dense.weight requires_grad= True
bert.encoder.layer.0.intermediate.dense.bias requires_grad= True
bert.encoder.layer.0.output.dense.weight requires_grad= True
bert.encoder.layer.0.output.dense.bias requires_grad= True
bert.encoder.layer.0.output.LayerNorm.weight requires_grad= True
bert.encoder.layer.0.output.LayerNorm.bias requires_grad= True
bert.encoder.layer.1.attention.self.query.weight requires_grad= True
bert.encoder.layer.1.attention.self.query.bias requires_grad= True
bert.encoder.layer.1.attention.self.key.weight requires_grad= True
bert.encoder.layer.1.attention.self.key.bias requires_grad= True
bert.encoder.layer.1.attention.self.value.weight requires_grad= True
bert.encoder.layer.1.attention.self.value.bias requires_grad= True
bert.encoder.layer.1.attention.output.dense.weight requires_grad= True
bert.encoder.layer.1.attention.output.dense.bias requires_grad= True
bert.encoder.layer.1.attention.output.LayerNorm.weight requires_grad= True
bert.encoder.layer.1.attention.output.LayerNorm.bias requires_grad= True
bert.encoder.layer.1.intermediate.dense.weight requires_grad= True
bert.encoder.layer.1.intermediate.dense.bias requires_grad= True
bert.encoder.layer.1.output.dense.weight requires_grad= True

```

[illegible]



bert.encoder.layer.7.output.dense.bias requires\_grad= True  
bert.encoder.layer.7.output.LayerNorm.weight requires\_grad= True  
bert.encoder.layer.7.output.LayerNorm.bias requires\_grad= True  
bert.encoder.layer.8.attention.self.query.weight requires\_grad= True  
bert.encoder.layer.8.attention.self.query.bias requires\_grad= True  
bert.encoder.layer.8.attention.self.key.weight requires\_grad= True  
bert.encoder.layer.8.attention.self.key.bias requires\_grad= True  
bert.encoder.layer.8.attention.self.value.weight requires\_grad= True  
bert.encoder.layer.8.attention.self.value.bias requires\_grad= True  
bert.encoder.layer.8.attention.output.dense.weight requires\_grad= True  
bert.encoder.layer.8.attention.output.dense.bias requires\_grad= True  
bert.encoder.layer.8.attention.output.LayerNorm.weight requires\_grad= True  
bert.encoder.layer.8.attention.output.LayerNorm.bias requires\_grad= True  
bert.encoder.layer.8.intermediate.dense.weight requires\_grad= True  
bert.encoder.layer.8.intermediate.dense.bias requires\_grad= True  
bert.encoder.layer.8.output.dense.weight requires\_grad= True  
bert.encoder.layer.8.output.dense.bias requires\_grad= True  
bert.encoder.layer.8.output.LayerNorm.weight requires\_grad= True  
bert.encoder.layer.8.output.LayerNorm.bias requires\_grad= True  
bert.encoder.layer.9.attention.self.query.weight requires\_grad= True  
bert.encoder.layer.9.attention.self.query.bias requires\_grad= True  
bert.encoder.layer.9.attention.self.key.weight requires\_grad= True  
bert.encoder.layer.9.attention.self.key.bias requires\_grad= True  
bert.encoder.layer.9.attention.self.value.weight requires\_grad= True  
bert.encoder.layer.9.attention.self.value.bias requires\_grad= True  
bert.encoder.layer.9.attention.output.dense.weight requires\_grad= True  
bert.encoder.layer.9.attention.output.dense.bias requires\_grad= True  
bert.encoder.layer.9.attention.output.LayerNorm.weight requires\_grad= True  
bert.encoder.layer.9.attention.output.LayerNorm.bias requires\_grad= True  
bert.encoder.layer.9.intermediate.dense.weight requires\_grad= True  
bert.encoder.layer.9.intermediate.dense.bias requires\_grad= True  
bert.encoder.layer.9.output.dense.weight requires\_grad= True  
bert.encoder.layer.9.output.dense.bias requires\_grad= True  
bert.encoder.layer.9.output.LayerNorm.weight requires\_grad= True  
bert.encoder.layer.9.output.LayerNorm.bias requires\_grad= True  
bert.encoder.layer.10.attention.self.query.weight requires\_grad= True  
bert.encoder.layer.10.attention.self.query.bias requires\_grad= True  
bert.encoder.layer.10.attention.self.key.weight requires\_grad= True  
bert.encoder.layer.10.attention.self.key.bias requires\_grad= True  
bert.encoder.layer.10.attention.self.value.weight requires\_grad= True  
bert.encoder.layer.10.attention.self.value.bias requires\_grad= True  
bert.encoder.layer.10.attention.output.dense.weight requires\_grad= True  
bert.encoder.layer.10.attention.output.dense.bias requires\_grad= True  
bert.encoder.layer.10.attention.output.LayerNorm.weight requires\_grad= True  
bert.encoder.layer.10.attention.output.LayerNorm.bias requires\_grad= True  
bert.encoder.layer.10.intermediate.dense.weight requires\_grad= True  
bert.encoder.layer.10.intermediate.dense.bias requires\_grad= True  
bert.encoder.layer.10.output.dense.weight requires\_grad= True











```

bert.encoder.layer.22.output.dense.bias requires_grad= True
bert.encoder.layer.22.output.LayerNorm.weight requires_grad= True
bert.encoder.layer.22.output.LayerNorm.bias requires_grad= True
bert.encoder.layer.23.attention.self.query.weight requires_grad= True
bert.encoder.layer.23.attention.self.query.bias requires_grad= True
bert.encoder.layer.23.attention.self.key.weight requires_grad= True
bert.encoder.layer.23.attention.self.key.bias requires_grad= True
bert.encoder.layer.23.attention.self.value.weight requires_grad= True
bert.encoder.layer.23.attention.self.value.bias requires_grad= True
bert.encoder.layer.23.attention.output.dense.weight requires_grad= True
bert.encoder.layer.23.attention.output.dense.bias requires_grad= True
bert.encoder.layer.23.attention.output.LayerNorm.weight requires_grad= True
bert.encoder.layer.23.attention.output.LayerNorm.bias requires_grad= True
bert.encoder.layer.23.intermediate.dense.weight requires_grad= True
bert.encoder.layer.23.intermediate.dense.bias requires_grad= True
bert.encoder.layer.23.output.dense.weight requires_grad= True
bert.encoder.layer.23.output.dense.bias requires_grad= True
bert.encoder.layer.23.output.LayerNorm.weight requires_grad= True
bert.encoder.layer.23.output.LayerNorm.bias requires_grad= True
bert.pooler.dense.weight requires_grad= True
bert.pooler.dense.bias requires_grad= True
classifier.weight requires_grad= True
classifier.bias requires_grad= True

```

```

[28]: #####
layers_to_unfreeze = [
    "bert.encoder.layer.23.",
    "bert.pooler.",
    "classifier.",
]
freeze_unfreeze_layers(model, layers_to_unfreeze=layers_to_unfreeze)
print(model.config)
print("=====")
print("num_parameters:", model.num_parameters())
print("num_trainable_parameters:", model.num_parameters(only_trainable=True))
print("=====")
#####
print("Experiment configuration used with this experiment:")
print("model used:", named_model)
print("learning rate used:", learning_rate)
print("number of epochs:", num_epochs)
print("maximum sequence length:", length_max)
print("batch size used:", size_batch)
print("regularization value:", regularization_weight_decay)
print("outcome variable:", y_col)
print("task:", x_task)
print("input column:", x_col)

```

```
#####
print("=====")
print("num_trainable_parameters:", model.num_parameters(only_trainable=True))
```

```
BertConfig {
  "_attn_implementation_autoset": true,
  "architectures": [
    "BertForMaskedLM"
  ],
  "attention_probs_dropout_prob": 0.1,
  "classifier_dropout": null,
  "directionality": "bidi",
  "gradient_checkpointing": false,
  "hidden_act": "gelu",
  "hidden_dropout_prob": 0.1,
  "hidden_size": 1024,
  "initializer_range": 0.02,
  "intermediate_size": 4096,
  "layer_norm_eps": 1e-12,
  "max_position_embeddings": 512,
  "model_type": "bert",
  "num_attention_heads": 16,
  "num_hidden_layers": 24,
  "pad_token_id": 0,
  "pooler_fc_size": 768,
  "pooler_num_attention_heads": 12,
  "pooler_num_fc_layers": 3,
  "pooler_size_per_head": 128,
  "pooler_type": "first_token_transform",
  "position_embedding_type": "absolute",
  "torch_dtype": "float32",
  "transformers_version": "4.50.3",
  "type_vocab_size": 2,
  "use_cache": true,
  "vocab_size": 28996
}
```

=====

```
num_parameters: 333581314
num_trainable_parameters: 13647874
```

=====

Experiment configuration used with this experiment:

```
model used: bert-large-cased
learning rate used: 5e-06
number of epochs: 1
maximum sequence length: 128
batch size used: 128
regularization value: 0.5
```

```
outcome variable: binary_complexity_75th_split
task: single
input column: sentence_no_contractions
=====
num_trainable_parameters: 13647874
```

```
[29]: model.resize_token_embeddings(len(tokenizer))
```

```
[29]: Embedding(28996, 1024, padding_idx=0)
```

```
[30]: for name, param in model.named_parameters():
      print(name, "requires_grad=", param.requires_grad)
```

```
bert.embeddings.word_embeddings.weight requires_grad= False
bert.embeddings.position_embeddings.weight requires_grad= False
bert.embeddings.token_type_embeddings.weight requires_grad= False
bert.embeddings.LayerNorm.weight requires_grad= False
bert.embeddings.LayerNorm.bias requires_grad= False
bert.encoder.layer.0.attention.self.query.weight requires_grad= False
bert.encoder.layer.0.attention.self.query.bias requires_grad= False
bert.encoder.layer.0.attention.self.key.weight requires_grad= False
bert.encoder.layer.0.attention.self.key.bias requires_grad= False
bert.encoder.layer.0.attention.self.value.weight requires_grad= False
bert.encoder.layer.0.attention.self.value.bias requires_grad= False
bert.encoder.layer.0.attention.output.dense.weight requires_grad= False
bert.encoder.layer.0.attention.output.dense.bias requires_grad= False
bert.encoder.layer.0.attention.output.LayerNorm.weight requires_grad= False
bert.encoder.layer.0.attention.output.LayerNorm.bias requires_grad= False
bert.encoder.layer.0.intermediate.dense.weight requires_grad= False
bert.encoder.layer.0.intermediate.dense.bias requires_grad= False
bert.encoder.layer.0.output.dense.weight requires_grad= False
bert.encoder.layer.0.output.dense.bias requires_grad= False
bert.encoder.layer.0.output.LayerNorm.weight requires_grad= False
bert.encoder.layer.0.output.LayerNorm.bias requires_grad= False
bert.encoder.layer.1.attention.self.query.weight requires_grad= False
bert.encoder.layer.1.attention.self.query.bias requires_grad= False
bert.encoder.layer.1.attention.self.key.weight requires_grad= False
bert.encoder.layer.1.attention.self.key.bias requires_grad= False
bert.encoder.layer.1.attention.self.value.weight requires_grad= False
bert.encoder.layer.1.attention.self.value.bias requires_grad= False
bert.encoder.layer.1.attention.output.dense.weight requires_grad= False
bert.encoder.layer.1.attention.output.dense.bias requires_grad= False
bert.encoder.layer.1.attention.output.LayerNorm.weight requires_grad= False
bert.encoder.layer.1.attention.output.LayerNorm.bias requires_grad= False
bert.encoder.layer.1.intermediate.dense.weight requires_grad= False
bert.encoder.layer.1.intermediate.dense.bias requires_grad= False
bert.encoder.layer.1.output.dense.weight requires_grad= False
bert.encoder.layer.1.output.dense.bias requires_grad= False
```







bert.encoder.layer.7.output.LayerNorm.weight requires\_grad= False  
bert.encoder.layer.7.output.LayerNorm.bias requires\_grad= False  
bert.encoder.layer.8.attention.self.query.weight requires\_grad= False  
bert.encoder.layer.8.attention.self.query.bias requires\_grad= False  
bert.encoder.layer.8.attention.self.key.weight requires\_grad= False  
bert.encoder.layer.8.attention.self.key.bias requires\_grad= False  
bert.encoder.layer.8.attention.self.value.weight requires\_grad= False  
bert.encoder.layer.8.attention.self.value.bias requires\_grad= False  
bert.encoder.layer.8.attention.output.dense.weight requires\_grad= False  
bert.encoder.layer.8.attention.output.dense.bias requires\_grad= False  
bert.encoder.layer.8.attention.output.LayerNorm.weight requires\_grad= False  
bert.encoder.layer.8.attention.output.LayerNorm.bias requires\_grad= False  
bert.encoder.layer.8.intermediate.dense.weight requires\_grad= False  
bert.encoder.layer.8.intermediate.dense.bias requires\_grad= False  
bert.encoder.layer.8.output.dense.weight requires\_grad= False  
bert.encoder.layer.8.output.dense.bias requires\_grad= False  
bert.encoder.layer.8.output.LayerNorm.weight requires\_grad= False  
bert.encoder.layer.8.output.LayerNorm.bias requires\_grad= False  
bert.encoder.layer.9.attention.self.query.weight requires\_grad= False  
bert.encoder.layer.9.attention.self.query.bias requires\_grad= False  
bert.encoder.layer.9.attention.self.key.weight requires\_grad= False  
bert.encoder.layer.9.attention.self.key.bias requires\_grad= False  
bert.encoder.layer.9.attention.self.value.weight requires\_grad= False  
bert.encoder.layer.9.attention.self.value.bias requires\_grad= False  
bert.encoder.layer.9.attention.output.dense.weight requires\_grad= False  
bert.encoder.layer.9.attention.output.dense.bias requires\_grad= False  
bert.encoder.layer.9.attention.output.LayerNorm.weight requires\_grad= False  
bert.encoder.layer.9.attention.output.LayerNorm.bias requires\_grad= False  
bert.encoder.layer.9.intermediate.dense.weight requires\_grad= False  
bert.encoder.layer.9.intermediate.dense.bias requires\_grad= False  
bert.encoder.layer.9.output.dense.weight requires\_grad= False  
bert.encoder.layer.9.output.dense.bias requires\_grad= False  
bert.encoder.layer.9.output.LayerNorm.weight requires\_grad= False  
bert.encoder.layer.9.output.LayerNorm.bias requires\_grad= False  
bert.encoder.layer.10.attention.self.query.weight requires\_grad= False  
bert.encoder.layer.10.attention.self.query.bias requires\_grad= False  
bert.encoder.layer.10.attention.self.key.weight requires\_grad= False  
bert.encoder.layer.10.attention.self.key.bias requires\_grad= False  
bert.encoder.layer.10.attention.self.value.weight requires\_grad= False  
bert.encoder.layer.10.attention.self.value.bias requires\_grad= False  
bert.encoder.layer.10.attention.output.dense.weight requires\_grad= False  
bert.encoder.layer.10.attention.output.dense.bias requires\_grad= False  
bert.encoder.layer.10.attention.output.LayerNorm.weight requires\_grad= False  
bert.encoder.layer.10.attention.output.LayerNorm.bias requires\_grad= False  
bert.encoder.layer.10.intermediate.dense.weight requires\_grad= False  
bert.encoder.layer.10.intermediate.dense.bias requires\_grad= False  
bert.encoder.layer.10.output.dense.weight requires\_grad= False  
bert.encoder.layer.10.output.dense.bias requires\_grad= False

[illegible]

[illegible]

[illegible]



```

bert.encoder.layer.22.output.LayerNorm.weight requires_grad= False
bert.encoder.layer.22.output.LayerNorm.bias requires_grad= False
bert.encoder.layer.23.attention.self.query.weight requires_grad= True
bert.encoder.layer.23.attention.self.query.bias requires_grad= True
bert.encoder.layer.23.attention.self.key.weight requires_grad= True
bert.encoder.layer.23.attention.self.key.bias requires_grad= True
bert.encoder.layer.23.attention.self.value.weight requires_grad= True
bert.encoder.layer.23.attention.self.value.bias requires_grad= True
bert.encoder.layer.23.attention.output.dense.weight requires_grad= True
bert.encoder.layer.23.attention.output.dense.bias requires_grad= True
bert.encoder.layer.23.attention.output.LayerNorm.weight requires_grad= True
bert.encoder.layer.23.attention.output.LayerNorm.bias requires_grad= True
bert.encoder.layer.23.intermediate.dense.weight requires_grad= True
bert.encoder.layer.23.intermediate.dense.bias requires_grad= True
bert.encoder.layer.23.output.dense.weight requires_grad= True
bert.encoder.layer.23.output.dense.bias requires_grad= True
bert.encoder.layer.23.output.LayerNorm.weight requires_grad= True
bert.encoder.layer.23.output.LayerNorm.bias requires_grad= True
bert.pooler.dense.weight requires_grad= True
bert.pooler.dense.bias requires_grad= True
classifier.weight requires_grad= True
classifier.bias requires_grad= True

```

```
[31]: model.resize_token_embeddings(len(tokenizer))
```

```
[31]: Embedding(28996, 1024, padding_idx=0)
```

```

[32]: # Train & Evaluate
trained_model, trainer_obj = train_transformer_model(
    model = model,
    tokenizer = tokenizer,
    train_dataset = train_data_hf,
    val_dataset = val_data_hf,
    output_dir = dir_results,
    num_epochs = num_epochs,
    batch_size = size_batch,
    lr = learning_rate,
    weight_decay = regularization_weight_decay)
metrics = trainer_obj.evaluate()
print("Validation metrics:", metrics)
test_metrics = trainer_obj.evaluate(test_data_hf) if test_data_hf else None
print("Test metrics:", test_metrics)

```

```

/usr/local/lib/python3.11/dist-packages/transformers/training_args.py:1611:
FutureWarning: `evaluation_strategy` is deprecated and will be removed in
version 4.46 of Transformers. Use `eval_strategy` instead
    warnings.warn(
<ipython-input-20-c2ee9f934517>:31: FutureWarning: `tokenizer` is deprecated and

```

will be removed in version 5.0.0 for `Trainer.\_\_init\_\_`. Use `processing\_class` instead.

```
trainer = Trainer(

<IPython.core.display.HTML object>

Downloading builder script: 0%|          | 0.00/4.20k [00:00<?, ?B/s]
Downloading builder script: 0%|          | 0.00/7.56k [00:00<?, ?B/s]
Downloading builder script: 0%|          | 0.00/7.38k [00:00<?, ?B/s]
Downloading builder script: 0%|          | 0.00/6.79k [00:00<?, ?B/s]

/usr/local/lib/python3.11/dist-packages/sklearn/metrics/_classification.py:1565:
UndefinedMetricWarning: Precision is ill-defined and being set to 0.0 due to no
predicted samples. Use `zero_division` parameter to control this behavior.
    _warn_prf(average, modifier, f"{metric.capitalize()} is", len(result))

<IPython.core.display.HTML object>

/usr/local/lib/python3.11/dist-packages/sklearn/metrics/_classification.py:1565:
UndefinedMetricWarning: Precision is ill-defined and being set to 0.0 due to no
predicted samples. Use `zero_division` parameter to control this behavior.
    _warn_prf(average, modifier, f"{metric.capitalize()} is", len(result))

Validation metrics: {'eval_loss': 0.565083920955658, 'eval_accuracy':
0.7790973871733967, 'eval_precision': 0.0, 'eval_recall': 0.0, 'eval_f1': 0.0,
'eval_runtime': 7.0908, 'eval_samples_per_second': 59.373,
'eval_steps_per_second': 0.564, 'epoch': 1.0}
Test metrics: {'eval_loss': 0.5492461323738098, 'eval_accuracy':
0.7960741548527808, 'eval_precision': 0.0, 'eval_recall': 0.0, 'eval_f1': 0.0,
'eval_runtime': 8.8918, 'eval_samples_per_second': 103.129,
'eval_steps_per_second': 0.9, 'epoch': 1.0}

/usr/local/lib/python3.11/dist-packages/sklearn/metrics/_classification.py:1565:
UndefinedMetricWarning: Precision is ill-defined and being set to 0.0 due to no
predicted samples. Use `zero_division` parameter to control this behavior.
    _warn_prf(average, modifier, f"{metric.capitalize()} is", len(result))
```

```
[33]: # save model checkpoint
timestamp = datetime.now().strftime("%Y%m%d_%H%M%S")
model_save_path = os.path.join(dir_models,
    f"{x_task}_{named_model}_{y_col}_{timestamp}")
trainer_obj.save_model(model_save_path)
print(f"Model checkpoint saved to: {model_save_path}")
# log experiment results
experiment_info = {
    "model_name": named_model,
    "learning_rate": learning_rate,
    "epochs": num_epochs,
    "batch_size": size_batch,
```



```

    "weight_decay": regularization_weight_decay,
    "x_task": x_task,
    "x_col": x_col,
    "y_col": y_col,
    "layers_to_unfreeze": layers_to_unfreeze}
model_info = gather_model_details(trained_model)
all_run_metrics = gather_all_run_metrics(
    trainer=trainer_obj,
    train_dataset=train_data_hf,
    val_dataset=val_data_hf,
    test_dataset=test_data_hf)
log_experiment_results_json(
    experiment_meta=experiment_info,
    model_details=model_info,
    run_metrics=all_run_metrics,
    log_file=log_filepath)
print(f"EXPERIMENT LOGGED TO: {log_filepath}")

```

Model checkpoint saved to: /content/drive/MyDrive/266-final/models/single\_bert-large-cased\_binary\_complexity\_75th\_split\_20250411\_010303

<IPython.core.display.HTML object>

```

/usr/local/lib/python3.11/dist-packages/sklearn/metrics/_classification.py:1565:
UndefinedMetricWarning: Precision is ill-defined and being set to 0.0 due to no
predicted samples. Use `zero_division` parameter to control this behavior.
    _warn_prf(average, modifier, f"{metric.capitalize()} is", len(result))
/usr/local/lib/python3.11/dist-packages/sklearn/metrics/_classification.py:1565:
UndefinedMetricWarning: Precision is ill-defined and being set to 0.0 due to no
predicted samples. Use `zero_division` parameter to control this behavior.
    _warn_prf(average, modifier, f"{metric.capitalize()} is", len(result))

```

EXPERIMENT LOGGED TO:  
/content/drive/MyDrive/266-final/results/experiment\_runs.txt

**0.2.6 snc roberta-large regularization\_weight\_decay = 0.5 learning\_rate = 5e-6  
size\_batch = 128 length\_max = 128 num\_epochs = 1**

```

[34]: # Define Experiment Parameters
# named_model = "bert-base-cased"
# named_model = "roberta-base"
# named_model = "bert-large-cased"
named_model = "roberta-large"
# named_model = "" # modern bert
#####
regularization_weight_decay = 0.5
learning_rate = 5e-6
size_batch = 128
length_max = 128

```

```

num_epochs = 1
#####
# x_col = "sentence"
x_col = "sentence_no_contractions"
# x_col = "pos_sequence"
# x_col = "dep_sequence"
# x_col = "morph_sequence"
#####
y_col = "binary_complexity_75th_split"
# y_col = "binary_complexity"
# y_col = "complexity"
#####
x_task = "single"
# x_task = "multi"
if x_task == "single":
    df_train = train_single_df
    df_val = trial_val_single_df
    df_test = test_single_df
else:
    df_train = train_multi_df
    df_val = trial_val_multi_df
    df_test = test_multi_df
#####
# Tokenize & Prepare Datasets
train_data_hf = prepare_dataset(
    df_train,
    tokenizer,
    text_col=x_col,
    label_col=y_col,
    max_length=length_max)
val_data_hf = prepare_dataset(
    df_val,
    tokenizer,
    text_col=x_col,
    label_col=y_col,
    max_length=length_max)
test_data_hf = prepare_dataset(
    df_test,
    tokenizer,
    text_col=x_col,
    label_col=y_col,
    max_length=length_max)
print("Datasets prepared. Sample from train_data_hf:\n", train_data_hf[10])
# print("Datasets prepared. Sample from train_data_hf:\n", val_data_hf[10])
# print("Datasets prepared. Sample from train_data_hf:\n", test_data_hf[10])
#####
# custom_config = BertConfig.from_pretrained("roberta-base")

```

```

# custom_config.hidden_act = "gelu" # alts: "relu" "silu"
# custom_config.attention_probs_dropout_prob = 0.1
# custom_config.hidden_dropout_prob = 0.1
# custom_config.gradient_checkpointing = False
#####
model, tokenizer = get_model_and_tokenizer(
    remote_model_name="roberta-large",
    local_model_path=None,
    config=None)
#####
# model, tokenizer = get_model_and_tokenizer(
#     remote_model_name=None
#     local_model_path="...CONFIGURE_PATH...",
#     config=custom_config)
print("=====")
print(named_model, ":")
print("=====")
print("num_parameters:", model.num_parameters())
print("num_trainable_parameters at load:", model.
    ↪num_parameters(only_trainable=True))
print("=====")
print("model lineage:", MODEL_LINEAGE)
print("=====")

```

Map: 0%| | 0/7662 [00:00<?, ? examples/s]

Map: 0%| | 0/421 [00:00<?, ? examples/s]

Map: 0%| | 0/917 [00:00<?, ? examples/s]

Datasets prepared. Sample from train\_data\_hf:

```

{'labels': tensor(0), 'input_ids': tensor([ 101, 1252, 1106, 1103, 3824,
1104, 19892, 11220, 1324, 1119,
1522, 3839, 117, 1272, 1103, 1555, 1104, 1103, 11563, 5609,
1106, 1172, 132, 1152, 2446, 1122, 1113, 1147, 3221, 119,
102, 0, 0, 0, 0, 0, 0, 0, 0, 0,
0, 0, 0, 0, 0, 0, 0, 0, 0, 0,
0, 0, 0, 0, 0, 0, 0, 0, 0, 0,
0, 0, 0, 0, 0, 0, 0, 0, 0, 0,
0, 0, 0, 0, 0, 0, 0, 0, 0, 0,
0, 0, 0, 0, 0, 0, 0, 0, 0, 0,
0, 0, 0, 0, 0, 0, 0, 0, 0, 0,
0, 0, 0, 0, 0, 0, 0, 0, 0, 0,
0, 0, 0, 0, 0, 0, 0, 0, 0, 0,
0, 0, 0, 0, 0, 0, 0, 0]),
'attention_mask': tensor([1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1,
1, 1, 1, 1, 1, 1,
1, 1, 1, 1, 1, 1, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0,
0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0,

```

```

0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0,
0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0,
0, 0, 0, 0, 0, 0, 0, 0, 0]})}

```

Loading from Hugging Face model: roberta-large

```
tokenizer_config.json: 0%|          | 0.00/25.0 [00:00<?, ?B/s]
```

```
config.json: 0%|          | 0.00/482 [00:00<?, ?B/s]
```

```
vocab.json: 0%|          | 0.00/899k [00:00<?, ?B/s]
```

```
merges.txt: 0%|          | 0.00/456k [00:00<?, ?B/s]
```

```
tokenizer.json: 0%|          | 0.00/1.36M [00:00<?, ?B/s]
```

```
model.safetensors: 0%|          | 0.00/1.42G [00:00<?, ?B/s]
```

Some weights of RobertaForSequenceClassification were not initialized from the model checkpoint at roberta-large and are newly initialized:

```
['classifier.dense.bias', 'classifier.dense.weight', 'classifier.out_proj.bias',
'classifier.out_proj.weight']
```

You should probably TRAIN this model on a down-stream task to be able to use it for predictions and inference.

```
=====
```

```
roberta-large :
```

```
=====
```

```
num_parameters: 355361794
```

```
num_trainable_parameters at load: 355361794
```

```
=====
```

```
model lineage: {'type': 'huggingface_hub', 'path': 'roberta-large', 'timestamp':
'2025-04-11 01:04:16'}
```

```
=====
```

```
[35]: print(model)
```

```

RobertaForSequenceClassification(
  (roberta): RobertaModel(
    (embeddings): RobertaEmbeddings(
      (word_embeddings): Embedding(50265, 1024, padding_idx=1)
      (position_embeddings): Embedding(514, 1024, padding_idx=1)
      (token_type_embeddings): Embedding(1, 1024)
      (LayerNorm): LayerNorm((1024,), eps=1e-05, elementwise_affine=True)
      (dropout): Dropout(p=0.1, inplace=False)
    )
    (encoder): RobertaEncoder(
      (layer): ModuleList(
        (0-23): 24 x RobertaLayer(
          (attention): RobertaAttention(
            (self): RobertaSdpaSelfAttention(
              (query): Linear(in_features=1024, out_features=1024, bias=True)
              (key): Linear(in_features=1024, out_features=1024, bias=True)

```

```

        (value): Linear(in_features=1024, out_features=1024, bias=True)
        (dropout): Dropout(p=0.1, inplace=False)
    )
    (output): RobertaSelfOutput(
        (dense): Linear(in_features=1024, out_features=1024, bias=True)
        (LayerNorm): LayerNorm((1024,), eps=1e-05,
elementwise_affine=True)
        (dropout): Dropout(p=0.1, inplace=False)
    )
)
(intermediate): RobertaIntermediate(
    (dense): Linear(in_features=1024, out_features=4096, bias=True)
    (intermediate_act_fn): GELUActivation()
)
(output): RobertaOutput(
    (dense): Linear(in_features=4096, out_features=1024, bias=True)
    (LayerNorm): LayerNorm((1024,), eps=1e-05, elementwise_affine=True)
    (dropout): Dropout(p=0.1, inplace=False)
)
)
)
)
)
(classifier): RobertaClassificationHead(
    (dense): Linear(in_features=1024, out_features=1024, bias=True)
    (dropout): Dropout(p=0.1, inplace=False)
    (out_proj): Linear(in_features=1024, out_features=2, bias=True)
)
)

```

```

[36]: for name, param in model.named_parameters():
        print(name, "requires_grad=", param.requires_grad)

```

```

roberta.embeddings.word_embeddings.weight requires_grad= True
roberta.embeddings.position_embeddings.weight requires_grad= True
roberta.embeddings.token_type_embeddings.weight requires_grad= True
roberta.embeddings.LayerNorm.weight requires_grad= True
roberta.embeddings.LayerNorm.bias requires_grad= True
roberta.encoder.layer.0.attention.self.query.weight requires_grad= True
roberta.encoder.layer.0.attention.self.query.bias requires_grad= True
roberta.encoder.layer.0.attention.self.key.weight requires_grad= True
roberta.encoder.layer.0.attention.self.key.bias requires_grad= True
roberta.encoder.layer.0.attention.self.value.weight requires_grad= True
roberta.encoder.layer.0.attention.self.value.bias requires_grad= True
roberta.encoder.layer.0.attention.output.dense.weight requires_grad= True
roberta.encoder.layer.0.attention.output.dense.bias requires_grad= True
roberta.encoder.layer.0.attention.output.LayerNorm.weight requires_grad= True
roberta.encoder.layer.0.attention.output.LayerNorm.bias requires_grad= True

```











[illegible]





```

roberta.encoder.layer.21.intermediate.dense.weight requires_grad= True
roberta.encoder.layer.21.intermediate.dense.bias requires_grad= True
roberta.encoder.layer.21.output.dense.weight requires_grad= True
roberta.encoder.layer.21.output.dense.bias requires_grad= True
roberta.encoder.layer.21.output.LayerNorm.weight requires_grad= True
roberta.encoder.layer.21.output.LayerNorm.bias requires_grad= True
roberta.encoder.layer.22.attention.self.query.weight requires_grad= True
roberta.encoder.layer.22.attention.self.query.bias requires_grad= True
roberta.encoder.layer.22.attention.self.key.weight requires_grad= True
roberta.encoder.layer.22.attention.self.key.bias requires_grad= True
roberta.encoder.layer.22.attention.self.value.weight requires_grad= True
roberta.encoder.layer.22.attention.self.value.bias requires_grad= True
roberta.encoder.layer.22.attention.output.dense.weight requires_grad= True
roberta.encoder.layer.22.attention.output.dense.bias requires_grad= True
roberta.encoder.layer.22.attention.output.LayerNorm.weight requires_grad= True
roberta.encoder.layer.22.attention.output.LayerNorm.bias requires_grad= True
roberta.encoder.layer.22.intermediate.dense.weight requires_grad= True
roberta.encoder.layer.22.intermediate.dense.bias requires_grad= True
roberta.encoder.layer.22.output.dense.weight requires_grad= True
roberta.encoder.layer.22.output.dense.bias requires_grad= True
roberta.encoder.layer.22.output.LayerNorm.weight requires_grad= True
roberta.encoder.layer.22.output.LayerNorm.bias requires_grad= True
roberta.encoder.layer.23.attention.self.query.weight requires_grad= True
roberta.encoder.layer.23.attention.self.query.bias requires_grad= True
roberta.encoder.layer.23.attention.self.key.weight requires_grad= True
roberta.encoder.layer.23.attention.self.key.bias requires_grad= True
roberta.encoder.layer.23.attention.self.value.weight requires_grad= True
roberta.encoder.layer.23.attention.self.value.bias requires_grad= True
roberta.encoder.layer.23.attention.output.dense.weight requires_grad= True
roberta.encoder.layer.23.attention.output.dense.bias requires_grad= True
roberta.encoder.layer.23.attention.output.LayerNorm.weight requires_grad= True
roberta.encoder.layer.23.attention.output.LayerNorm.bias requires_grad= True
roberta.encoder.layer.23.intermediate.dense.weight requires_grad= True
roberta.encoder.layer.23.intermediate.dense.bias requires_grad= True
roberta.encoder.layer.23.output.dense.weight requires_grad= True
roberta.encoder.layer.23.output.dense.bias requires_grad= True
roberta.encoder.layer.23.output.LayerNorm.weight requires_grad= True
roberta.encoder.layer.23.output.LayerNorm.bias requires_grad= True
classifier.dense.weight requires_grad= True
classifier.dense.bias requires_grad= True
classifier.out_proj.weight requires_grad= True
classifier.out_proj.bias requires_grad= True

```

```
[37]: # Inspect the attention_mask tensor for the first few samples
for i in range(5):
    print(train_data_hf[i]['attention_mask'])
```

[illegible]



```

        "roberta.encoder.layer.23.output.LayerNorm.bias",
        "classifier.dense.weight",
        "classifier.dense.bias",
        "classifier.out_proj.weight",
        "classifier.out_proj.bias",
    ]
    freeze_unfreeze_layers(model, layers_to_unfreeze=layers_to_unfreeze)
    print(model.config)
    print("=====")
    print("num_parameters:", model.num_parameters())
    print("num_trainable_parameters:", model.num_parameters(only_trainable=True))
    print("=====")
    #####
    print("Experiment configuration used with this experiment:")
    print("model used:", named_model)
    print("learning rate used:", learning_rate)
    print("number of epochs:", num_epochs)
    print("maximum sequence length:", length_max)
    print("batch size used:", size_batch)
    print("regularization value:", regularization_weight_decay)
    print("outcome variable:", y_col)
    print("task:", x_task)
    print("input column:", x_col)
    #####
    print("=====")
    print("num_trainable_parameters:", model.num_parameters(only_trainable=True))

```

```

RobertaConfig {
  "_attn_implementation_autoset": true,
  "architectures": [
    "RobertaForMaskedLM"
  ],
  "attention_probs_dropout_prob": 0.1,
  "bos_token_id": 0,
  "classifier_dropout": null,
  "eos_token_id": 2,
  "hidden_act": "gelu",
  "hidden_dropout_prob": 0.1,
  "hidden_size": 1024,
  "initializer_range": 0.02,
  "intermediate_size": 4096,
  "layer_norm_eps": 1e-05,
  "max_position_embeddings": 514,
  "model_type": "roberta",
  "num_attention_heads": 16,
  "num_hidden_layers": 24,
  "pad_token_id": 1,
  "position_embedding_type": "absolute",

```

```

    "torch_dtype": "float32",
    "transformers_version": "4.50.3",
    "type_vocab_size": 1,
    "use_cache": true,
    "vocab_size": 50265
}

```

```

=====

```

```

num_parameters: 355361794
num_trainable_parameters: 13647874

```

```

=====

```

Experiment configuration used with this experiment:

```

model used: roberta-large
learning rate used: 5e-06
number of epochs: 1
maximum sequence length: 128
batch size used: 128
regularization value: 0.5
outcome variable: binary_complexity_75th_split
task: single
input column: sentence_no_contractions

```

```

=====

```

```

num_trainable_parameters: 13647874

```

```

[39]: for name, param in model.named_parameters():
      print(name, "requires_grad=", param.requires_grad)

```

```

roberta.embeddings.word_embeddings.weight requires_grad= False
roberta.embeddings.position_embeddings.weight requires_grad= False
roberta.embeddings.token_type_embeddings.weight requires_grad= False
roberta.embeddings.LayerNorm.weight requires_grad= False
roberta.embeddings.LayerNorm.bias requires_grad= False
roberta.encoder.layer.0.attention.self.query.weight requires_grad= False
roberta.encoder.layer.0.attention.self.query.bias requires_grad= False
roberta.encoder.layer.0.attention.self.key.weight requires_grad= False
roberta.encoder.layer.0.attention.self.key.bias requires_grad= False
roberta.encoder.layer.0.attention.self.value.weight requires_grad= False
roberta.encoder.layer.0.attention.self.value.bias requires_grad= False
roberta.encoder.layer.0.attention.output.dense.weight requires_grad= False
roberta.encoder.layer.0.attention.output.dense.bias requires_grad= False
roberta.encoder.layer.0.attention.output.LayerNorm.weight requires_grad= False
roberta.encoder.layer.0.attention.output.LayerNorm.bias requires_grad= False
roberta.encoder.layer.0.intermediate.dense.weight requires_grad= False
roberta.encoder.layer.0.intermediate.dense.bias requires_grad= False
roberta.encoder.layer.0.output.dense.weight requires_grad= False
roberta.encoder.layer.0.output.dense.bias requires_grad= False
roberta.encoder.layer.0.output.LayerNorm.weight requires_grad= False
roberta.encoder.layer.0.output.LayerNorm.bias requires_grad= False

```



[illegible]













```

roberta.encoder.layer.22.attention.self.query.weight requires_grad= False
roberta.encoder.layer.22.attention.self.query.bias requires_grad= False
roberta.encoder.layer.22.attention.self.key.weight requires_grad= False
roberta.encoder.layer.22.attention.self.key.bias requires_grad= False
roberta.encoder.layer.22.attention.self.value.weight requires_grad= False
roberta.encoder.layer.22.attention.self.value.bias requires_grad= False
roberta.encoder.layer.22.attention.output.dense.weight requires_grad= False
roberta.encoder.layer.22.attention.output.dense.bias requires_grad= False
roberta.encoder.layer.22.attention.output.LayerNorm.weight requires_grad= False
roberta.encoder.layer.22.attention.output.LayerNorm.bias requires_grad= False
roberta.encoder.layer.22.intermediate.dense.weight requires_grad= False
roberta.encoder.layer.22.intermediate.dense.bias requires_grad= False
roberta.encoder.layer.22.output.dense.weight requires_grad= False
roberta.encoder.layer.22.output.dense.bias requires_grad= False
roberta.encoder.layer.22.output.LayerNorm.weight requires_grad= False
roberta.encoder.layer.22.output.LayerNorm.bias requires_grad= False
roberta.encoder.layer.23.attention.self.query.weight requires_grad= True
roberta.encoder.layer.23.attention.self.query.bias requires_grad= True
roberta.encoder.layer.23.attention.self.key.weight requires_grad= True
roberta.encoder.layer.23.attention.self.key.bias requires_grad= True
roberta.encoder.layer.23.attention.self.value.weight requires_grad= True
roberta.encoder.layer.23.attention.self.value.bias requires_grad= True
roberta.encoder.layer.23.attention.output.dense.weight requires_grad= True
roberta.encoder.layer.23.attention.output.dense.bias requires_grad= True
roberta.encoder.layer.23.attention.output.LayerNorm.weight requires_grad= True
roberta.encoder.layer.23.attention.output.LayerNorm.bias requires_grad= True
roberta.encoder.layer.23.intermediate.dense.weight requires_grad= True
roberta.encoder.layer.23.intermediate.dense.bias requires_grad= True
roberta.encoder.layer.23.output.dense.weight requires_grad= True
roberta.encoder.layer.23.output.dense.bias requires_grad= True
roberta.encoder.layer.23.output.LayerNorm.weight requires_grad= True
roberta.encoder.layer.23.output.LayerNorm.bias requires_grad= True
classifier.dense.weight requires_grad= True
classifier.dense.bias requires_grad= True
classifier.out_proj.weight requires_grad= True
classifier.out_proj.bias requires_grad= True

```

```
[40]: model.resize_token_embeddings(len(tokenizer))
```

```
[40]: Embedding(50265, 1024, padding_idx=1)
```

```
[41]: # Train & Evaluate
trained_model, trainer_obj = train_transformer_model(
    model = model,
    tokenizer = tokenizer,
    train_dataset = train_data_hf,
    val_dataset = val_data_hf,
    output_dir = dir_results,
```



```

    num_epochs = num_epochs,
    batch_size = size_batch,
    lr = learning_rate,
    weight_decay = regularization_weight_decay)
metrics = trainer_obj.evaluate()
print("Validation metrics:", metrics)
test_metrics = trainer_obj.evaluate(test_data_hf) if test_data_hf else None
print("Test metrics:", test_metrics)

```

```

/usr/local/lib/python3.11/dist-packages/transformers/training_args.py:1611:
FutureWarning: `evaluation_strategy` is deprecated and will be removed in
version 4.46 of Transformers. Use `eval_strategy` instead
    warnings.warn(
<ipython-input-20-c2ee9f934517>:31: FutureWarning: `tokenizer` is deprecated and
will be removed in version 5.0.0 for `Trainer.__init__`. Use `processing_class`
instead.
    trainer = Trainer(

<IPython.core.display.HTML object>

/usr/local/lib/python3.11/dist-packages/sklearn/metrics/_classification.py:1565:
UndefinedMetricWarning: Precision is ill-defined and being set to 0.0 due to no
predicted samples. Use `zero_division` parameter to control this behavior.
    _warn_prf(average, modifier, f"{metric.capitalize()} is", len(result))

<IPython.core.display.HTML object>

/usr/local/lib/python3.11/dist-packages/sklearn/metrics/_classification.py:1565:
UndefinedMetricWarning: Precision is ill-defined and being set to 0.0 due to no
predicted samples. Use `zero_division` parameter to control this behavior.
    _warn_prf(average, modifier, f"{metric.capitalize()} is", len(result))

Validation metrics: {'eval_loss': 0.6359410881996155, 'eval_accuracy':
0.7790973871733967, 'eval_precision': 0.0, 'eval_recall': 0.0, 'eval_f1': 0.0,
'eval_runtime': 6.7732, 'eval_samples_per_second': 62.156,
'eval_steps_per_second': 0.591, 'epoch': 1.0}
Test metrics: {'eval_loss': 0.6321974992752075, 'eval_accuracy':
0.7960741548527808, 'eval_precision': 0.0, 'eval_recall': 0.0, 'eval_f1': 0.0,
'eval_runtime': 8.8781, 'eval_samples_per_second': 103.288,
'eval_steps_per_second': 0.901, 'epoch': 1.0}

/usr/local/lib/python3.11/dist-packages/sklearn/metrics/_classification.py:1565:
UndefinedMetricWarning: Precision is ill-defined and being set to 0.0 due to no
predicted samples. Use `zero_division` parameter to control this behavior.
    _warn_prf(average, modifier, f"{metric.capitalize()} is", len(result))

```

```

[42]: # save model checkpoint
timestamp = datetime.now().strftime("%Y%m%d_%H%M%S")
model_save_path = os.path.join(dir_models,
    f"{x_task}_{named_model}_{y_col}_{timestamp}")

```

```

trainer_obj.save_model(model_save_path)
print(f"Model checkpoint saved to: {model_save_path}")
# log experiment results
experiment_info = {
    "model_name": named_model,
    "learning_rate": learning_rate,
    "epochs": num_epochs,
    "batch_size": size_batch,
    "weight_decay": regularization_weight_decay,
    "x_task": x_task,
    "x_col": x_col,
    "y_col": y_col,
    "layers_to_unfreeze": layers_to_unfreeze}
model_info = gather_model_details(trained_model)
all_run_metrics = gather_all_run_metrics(
    trainer=trainer_obj,
    train_dataset=train_data_hf,
    val_dataset=val_data_hf,
    test_dataset=test_data_hf)
log_experiment_results_json(
    experiment_meta=experiment_info,
    model_details=model_info,
    run_metrics=all_run_metrics,
    log_file=log_filepath)
print(f"EXPERIMENT LOGGED TO: {log_filepath}")

```

Model checkpoint saved to:  
/content/drive/MyDrive/266-final/models/single\_roberta-  
large\_binary\_complexity\_75th\_split\_20250411\_010518

<IPython.core.display.HTML object>

/usr/local/lib/python3.11/dist-packages/sklearn/metrics/\_classification.py:1565:  
UndefinedMetricWarning: Precision is ill-defined and being set to 0.0 due to no  
predicted samples. Use `zero\_division` parameter to control this behavior.

\_warn\_prf(average, modifier, f"{metric.capitalize()} is", len(result))  
/usr/local/lib/python3.11/dist-packages/sklearn/metrics/\_classification.py:1565:  
UndefinedMetricWarning: Precision is ill-defined and being set to 0.0 due to no  
predicted samples. Use `zero\_division` parameter to control this behavior.  
\_warn\_prf(average, modifier, f"{metric.capitalize()} is", len(result))

EXPERIMENT LOGGED TO:  
/content/drive/MyDrive/266-final/results/experiment\_runs.txt

/usr/local/lib/python3.11/dist-packages/sklearn/metrics/\_classification.py:1565:  
UndefinedMetricWarning: Precision is ill-defined and being set to 0.0 due to no  
predicted samples. Use `zero\_division` parameter to control this behavior.  
\_warn\_prf(average, modifier, f"{metric.capitalize()} is", len(result))

0.2.7 snc answerdotai/ModernBERT-base regularization\_weight\_decay = 0.5 learning\_rate = 5e-6 size\_batch = 128 length\_max = 128 num\_epochs = 1

```
[43]: # Define Experiment Parameters
# named_model = "bert-base-cased"
# named_model = "roberta-base"
# named_model = "bert-large-cased"
# named_model = "roberta-large"
named_model = "answerdotai/ModernBERT-base" # modern bert
#####
regularization_weight_decay = 0.5
learning_rate = 5e-6
size_batch = 128
length_max = 128
num_epochs = 1
#####
# x_col = "sentence"
x_col = "sentence_no_contractions"
# x_col = "pos_sequence"
# x_col = "dep_sequence"
# x_col = "morph_sequence"
#####
y_col = "binary_complexity_75th_split"
# y_col = "binary_complexity"
# y_col = "complexity"
#####
x_task = "single"
# x_task = "multi"
if x_task == "single":
    df_train = train_single_df
    df_val = trial_val_single_df
    df_test = test_single_df
else:
    df_train = train_multi_df
    df_val = trial_val_multi_df
    df_test = test_multi_df
#####
# Tokenize & Prepare Datasets
train_data_hf = prepare_dataset(
    df_train,
    tokenizer,
    text_col=x_col,
    label_col=y_col,
    max_length=length_max)
val_data_hf = prepare_dataset(
    df_val,
    tokenizer,
```

```

        text_col=x_col,
        label_col=y_col,
        max_length=length_max)
test_data_hf = prepare_dataset(
    df_test,
    tokenizer,
    text_col=x_col,
    label_col=y_col,
    max_length=length_max)
print("Datasets prepared. Sample from train_data_hf:\n", train_data_hf[10])
# print("Datasets prepared. Sample from train_data_hf:\n", val_data_hf[10])
# print("Datasets prepared. Sample from train_data_hf:\n", test_data_hf[10])
#####
# custom_config = BertConfig.from_pretrained("roberta-base")
# custom_config.hidden_act = "gelu" # alts: "relu" "silu"
# custom_config.attention_probs_dropout_prob = 0.1
# custom_config.hidden_dropout_prob = 0.1
# custom_config.gradient_checkpointing = False
#####
model, tokenizer = get_model_and_tokenizer(
    remote_model_name="answerdotai/ModernBERT-base",
    local_model_path=None,
    config=None)
#####
# model, tokenizer = get_model_and_tokenizer(
#     remote_model_name=None
#     local_model_path="...CONFIGURE_PATH...",
#     config=custom_config)
print("=====")
print(named_model, ":")
print("=====")
print("num_parameters:", model.num_parameters())
print("num_trainable_parameters at load:", model.
    ↪num_parameters(only_trainable=True))
print("=====")
print("model lineage:", MODEL_LINEAGE)
print("=====")

```

Map: 0%| | 0/7662 [00:00<?, ? examples/s]

Map: 0%| | 0/421 [00:00<?, ? examples/s]

Map: 0%| | 0/917 [00:00<?, ? examples/s]

Datasets prepared. Sample from train\_data\_hf:

```

{'labels': tensor(0), 'input_ids': tensor([ 0, 1708, 7, 5, 7250,
9, 6987, 2681, 37, 851,
4146, 6, 142, 5, 544, 9, 5, 13488, 14574, 7,
106, 131, 51, 2584, 24, 15, 49, 10762, 4, 2,

```



```

        (drop): Dropout(p=0.0, inplace=False)
    )
(layers): ModuleList(
  (0): ModernBertEncoderLayer(
    (attn_norm): Identity()
    (attn): ModernBertAttention(
      (Wqkv): Linear(in_features=768, out_features=2304, bias=False)
      (rotary_emb): ModernBertRotaryEmbedding()
      (Wo): Linear(in_features=768, out_features=768, bias=False)
      (out_drop): Identity()
    )
    (mlp_norm): LayerNorm((768,), eps=1e-05, elementwise_affine=True)
    (mlp): ModernBertMLP(
      (Wi): Linear(in_features=768, out_features=2304, bias=False)
      (act): GELUActivation()
      (drop): Dropout(p=0.0, inplace=False)
      (Wo): Linear(in_features=1152, out_features=768, bias=False)
    )
  )
  (1-21): 21 x ModernBertEncoderLayer(
    (attn_norm): LayerNorm((768,), eps=1e-05, elementwise_affine=True)
    (attn): ModernBertAttention(
      (Wqkv): Linear(in_features=768, out_features=2304, bias=False)
      (rotary_emb): ModernBertRotaryEmbedding()
      (Wo): Linear(in_features=768, out_features=768, bias=False)
      (out_drop): Identity()
    )
    (mlp_norm): LayerNorm((768,), eps=1e-05, elementwise_affine=True)
    (mlp): ModernBertMLP(
      (Wi): Linear(in_features=768, out_features=2304, bias=False)
      (act): GELUActivation()
      (drop): Dropout(p=0.0, inplace=False)
      (Wo): Linear(in_features=1152, out_features=768, bias=False)
    )
  )
)
  (final_norm): LayerNorm((768,), eps=1e-05, elementwise_affine=True)
)
(head): ModernBertPredictionHead(
  (dense): Linear(in_features=768, out_features=768, bias=False)
  (act): GELUActivation()
  (norm): LayerNorm((768,), eps=1e-05, elementwise_affine=True)
)
(drop): Dropout(p=0.0, inplace=False)
(classifier): Linear(in_features=768, out_features=2, bias=True)
)

```

```
[45]: for name, param in model.named_parameters():  
       print(name, "requires_grad=", param.requires_grad)
```

```
model.embeddings.tok_embeddings.weight requires_grad= True  
model.embeddings.norm.weight requires_grad= True  
model.layers.0.attn.Wqkv.weight requires_grad= True  
model.layers.0.attn.Wo.weight requires_grad= True  
model.layers.0.mlp_norm.weight requires_grad= True  
model.layers.0.mlp.Wi.weight requires_grad= True  
model.layers.0.mlp.Wo.weight requires_grad= True  
model.layers.1.attn_norm.weight requires_grad= True  
model.layers.1.attn.Wqkv.weight requires_grad= True  
model.layers.1.attn.Wo.weight requires_grad= True  
model.layers.1.mlp_norm.weight requires_grad= True  
model.layers.1.mlp.Wi.weight requires_grad= True  
model.layers.1.mlp.Wo.weight requires_grad= True  
model.layers.2.attn_norm.weight requires_grad= True  
model.layers.2.attn.Wqkv.weight requires_grad= True  
model.layers.2.attn.Wo.weight requires_grad= True  
model.layers.2.mlp_norm.weight requires_grad= True  
model.layers.2.mlp.Wi.weight requires_grad= True  
model.layers.2.mlp.Wo.weight requires_grad= True  
model.layers.3.attn_norm.weight requires_grad= True  
model.layers.3.attn.Wqkv.weight requires_grad= True  
model.layers.3.attn.Wo.weight requires_grad= True  
model.layers.3.mlp_norm.weight requires_grad= True  
model.layers.3.mlp.Wi.weight requires_grad= True  
model.layers.3.mlp.Wo.weight requires_grad= True  
model.layers.4.attn_norm.weight requires_grad= True  
model.layers.4.attn.Wqkv.weight requires_grad= True  
model.layers.4.attn.Wo.weight requires_grad= True  
model.layers.4.mlp_norm.weight requires_grad= True  
model.layers.4.mlp.Wi.weight requires_grad= True  
model.layers.4.mlp.Wo.weight requires_grad= True  
model.layers.5.attn_norm.weight requires_grad= True  
model.layers.5.attn.Wqkv.weight requires_grad= True  
model.layers.5.attn.Wo.weight requires_grad= True  
model.layers.5.mlp_norm.weight requires_grad= True  
model.layers.5.mlp.Wi.weight requires_grad= True  
model.layers.5.mlp.Wo.weight requires_grad= True  
model.layers.6.attn_norm.weight requires_grad= True  
model.layers.6.attn.Wqkv.weight requires_grad= True  
model.layers.6.attn.Wo.weight requires_grad= True  
model.layers.6.mlp_norm.weight requires_grad= True  
model.layers.6.mlp.Wi.weight requires_grad= True  
model.layers.6.mlp.Wo.weight requires_grad= True  
model.layers.7.attn_norm.weight requires_grad= True  
model.layers.7.attn.Wqkv.weight requires_grad= True
```

model.layers.7.attn.Wo.weight requires\_grad= True  
model.layers.7.mlp\_norm.weight requires\_grad= True  
model.layers.7.mlp.Wi.weight requires\_grad= True  
model.layers.7.mlp.Wo.weight requires\_grad= True  
model.layers.8.attn\_norm.weight requires\_grad= True  
model.layers.8.attn.Wqkv.weight requires\_grad= True  
model.layers.8.attn.Wo.weight requires\_grad= True  
model.layers.8.mlp\_norm.weight requires\_grad= True  
model.layers.8.mlp.Wi.weight requires\_grad= True  
model.layers.8.mlp.Wo.weight requires\_grad= True  
model.layers.9.attn\_norm.weight requires\_grad= True  
model.layers.9.attn.Wqkv.weight requires\_grad= True  
model.layers.9.attn.Wo.weight requires\_grad= True  
model.layers.9.mlp\_norm.weight requires\_grad= True  
model.layers.9.mlp.Wi.weight requires\_grad= True  
model.layers.9.mlp.Wo.weight requires\_grad= True  
model.layers.10.attn\_norm.weight requires\_grad= True  
model.layers.10.attn.Wqkv.weight requires\_grad= True  
model.layers.10.attn.Wo.weight requires\_grad= True  
model.layers.10.mlp\_norm.weight requires\_grad= True  
model.layers.10.mlp.Wi.weight requires\_grad= True  
model.layers.10.mlp.Wo.weight requires\_grad= True  
model.layers.11.attn\_norm.weight requires\_grad= True  
model.layers.11.attn.Wqkv.weight requires\_grad= True  
model.layers.11.attn.Wo.weight requires\_grad= True  
model.layers.11.mlp\_norm.weight requires\_grad= True  
model.layers.11.mlp.Wi.weight requires\_grad= True  
model.layers.11.mlp.Wo.weight requires\_grad= True  
model.layers.12.attn\_norm.weight requires\_grad= True  
model.layers.12.attn.Wqkv.weight requires\_grad= True  
model.layers.12.attn.Wo.weight requires\_grad= True  
model.layers.12.mlp\_norm.weight requires\_grad= True  
model.layers.12.mlp.Wi.weight requires\_grad= True  
model.layers.12.mlp.Wo.weight requires\_grad= True  
model.layers.13.attn\_norm.weight requires\_grad= True  
model.layers.13.attn.Wqkv.weight requires\_grad= True  
model.layers.13.attn.Wo.weight requires\_grad= True  
model.layers.13.mlp\_norm.weight requires\_grad= True  
model.layers.13.mlp.Wi.weight requires\_grad= True  
model.layers.13.mlp.Wo.weight requires\_grad= True  
model.layers.14.attn\_norm.weight requires\_grad= True  
model.layers.14.attn.Wqkv.weight requires\_grad= True  
model.layers.14.attn.Wo.weight requires\_grad= True  
model.layers.14.mlp\_norm.weight requires\_grad= True  
model.layers.14.mlp.Wi.weight requires\_grad= True  
model.layers.14.mlp.Wo.weight requires\_grad= True  
model.layers.15.attn\_norm.weight requires\_grad= True  
model.layers.15.attn.Wqkv.weight requires\_grad= True



```
model.layers.15.attn.Wo.weight requires_grad= True
model.layers.15.mlp_norm.weight requires_grad= True
model.layers.15.mlp.Wi.weight requires_grad= True
model.layers.15.mlp.Wo.weight requires_grad= True
model.layers.16.attn_norm.weight requires_grad= True
model.layers.16.attn.Wqkv.weight requires_grad= True
model.layers.16.attn.Wo.weight requires_grad= True
model.layers.16.mlp_norm.weight requires_grad= True
model.layers.16.mlp.Wi.weight requires_grad= True
model.layers.16.mlp.Wo.weight requires_grad= True
model.layers.17.attn_norm.weight requires_grad= True
model.layers.17.attn.Wqkv.weight requires_grad= True
model.layers.17.attn.Wo.weight requires_grad= True
model.layers.17.mlp_norm.weight requires_grad= True
model.layers.17.mlp.Wi.weight requires_grad= True
model.layers.17.mlp.Wo.weight requires_grad= True
model.layers.18.attn_norm.weight requires_grad= True
model.layers.18.attn.Wqkv.weight requires_grad= True
model.layers.18.attn.Wo.weight requires_grad= True
model.layers.18.mlp_norm.weight requires_grad= True
model.layers.18.mlp.Wi.weight requires_grad= True
model.layers.18.mlp.Wo.weight requires_grad= True
model.layers.19.attn_norm.weight requires_grad= True
model.layers.19.attn.Wqkv.weight requires_grad= True
model.layers.19.attn.Wo.weight requires_grad= True
model.layers.19.mlp_norm.weight requires_grad= True
model.layers.19.mlp.Wi.weight requires_grad= True
model.layers.19.mlp.Wo.weight requires_grad= True
model.layers.20.attn_norm.weight requires_grad= True
model.layers.20.attn.Wqkv.weight requires_grad= True
model.layers.20.attn.Wo.weight requires_grad= True
model.layers.20.mlp_norm.weight requires_grad= True
model.layers.20.mlp.Wi.weight requires_grad= True
model.layers.20.mlp.Wo.weight requires_grad= True
model.layers.21.attn_norm.weight requires_grad= True
model.layers.21.attn.Wqkv.weight requires_grad= True
model.layers.21.attn.Wo.weight requires_grad= True
model.layers.21.mlp_norm.weight requires_grad= True
model.layers.21.mlp.Wi.weight requires_grad= True
model.layers.21.mlp.Wo.weight requires_grad= True
model.final_norm.weight requires_grad= True
head.dense.weight requires_grad= True
head.norm.weight requires_grad= True
classifier.weight requires_grad= True
classifier.bias requires_grad= True
```

```
[46]: # Inspect the attention_mask tensor for the first few samples
for i in range(5):
    print(train_data_hf[i]['attention_mask'])
```

```
tensor([1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1,
        1, 1, 1, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0,
        0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0,
        0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0,
        0, 0, 0, 0, 0, 0, 0, 0, 0])
tensor([1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1,
        1, 1, 1, 1, 1, 1, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0,
        0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0,
        0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0,
        0, 0, 0, 0, 0, 0, 0, 0, 0])
tensor([1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1,
        1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1,
        1, 1, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0,
        0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0,
        0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0,
        0, 0, 0, 0, 0, 0, 0, 0])
tensor([1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1,
        1, 1, 1, 1, 1, 1, 1, 1, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0,
        0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0,
        0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0,
        0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0,
        0, 0, 0, 0, 0, 0, 0, 0])
tensor([1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0,
        0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0,
        0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0,
        0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0,
        0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0,
        0, 0, 0, 0, 0, 0, 0, 0])
```

```
[47]: #####
layers_to_unfreeze = [
    "model.layers.21.attn_norm.weight",
    "model.layers.21.attn.Wqkv.weight",
    "model.layers.21.attn.Wo.weight",
    "model.layers.21.mlp_norm.weight",
    "model.layers.21.mlp.Wi.weight",
    "model.layers.21.mlp.Wo.weight",
    "model.final_norm.weight",
    "head.dense.weight",
    "head.norm.weight",
    "classifier.weight",
    "classifier.bias"]
```

```

freeze_unfreeze_layers(model, layers_to_unfreeze=layers_to_unfreeze)
print(model.config)
print("=====")
print("num_parameters:", model.num_parameters())
print("num_trainable_parameters:", model.num_parameters(only_trainable=True))
print("=====")
#####
print("Experiment configuration used with this experiment:")
print("model used:", named_model)
print("learning rate used:", learning_rate)
print("number of epochs:", num_epochs)
print("maximum sequence length:", length_max)
print("batch size used:", size_batch)
print("regularization value:", regularization_weight_decay)
print("outcome variable:", y_col)
print("task:", x_task)
print("input column:", x_col)
#####
print("=====")
print("num_trainable_parameters:", model.num_parameters(only_trainable=True))

```

```

ModernBertConfig {
  "_attn_implementation_autoset": true,
  "architectures": [
    "ModernBertForMaskedLM"
  ],
  "attention_bias": false,
  "attention_dropout": 0.0,
  "bos_token_id": 50281,
  "classifier_activation": "gelu",
  "classifier_bias": false,
  "classifier_dropout": 0.0,
  "classifier_pooling": "mean",
  "cls_token_id": 50281,
  "decoder_bias": true,
  "deterministic_flash_attn": false,
  "embedding_dropout": 0.0,
  "eos_token_id": 50282,
  "global_attn_every_n_layers": 3,
  "global_rope_theta": 160000.0,
  "gradient_checkpointing": false,
  "hidden_activation": "gelu",
  "hidden_size": 768,
  "initializer_cutoff_factor": 2.0,
  "initializer_range": 0.02,
  "intermediate_size": 1152,
  "layer_norm_eps": 1e-05,
  "local_attention": 128,

```

```

"local_rope_theta": 10000.0,
"max_position_embeddings": 8192,
"mlp_bias": false,
"mlp_dropout": 0.0,
"model_type": "modernbert",
"norm_bias": false,
"norm_eps": 1e-05,
"num_attention_heads": 12,
"num_hidden_layers": 22,
"pad_token_id": 50283,
"position_embedding_type": "absolute",
"reference_compile": null,
"repad_logits_with_grad": false,
"sep_token_id": 50282,
"sparse_pred_ignore_index": -100,
"sparse_prediction": false,
"torch_dtype": "float32",
"transformers_version": "4.50.3",
"vocab_size": 50368
}

```

=====

```

num_parameters: 149606402
num_trainable_parameters: 5607938

```

=====

Experiment configuration used with this experiment:

```

model used: answerdotai/ModernBERT-base
learning rate used: 5e-06
number of epochs: 1
maximum sequence length: 128
batch size used: 128
regularization value: 0.5
outcome variable: binary_complexity_75th_split
task: single
input column: sentence_no_contractions

```

=====

```

num_trainable_parameters: 5607938

```

```

[48]: for name, param in model.named_parameters():
      print(name, "requires_grad=", param.requires_grad)

```

```

model.embeddings.tok_embeddings.weight requires_grad= False
model.embeddings.norm.weight requires_grad= False
model.layers.0.attn.Wqkv.weight requires_grad= False
model.layers.0.attn.Wo.weight requires_grad= False
model.layers.0.mlp_norm.weight requires_grad= False
model.layers.0.mlp.Wi.weight requires_grad= False
model.layers.0.mlp.Wo.weight requires_grad= False

```



model.layers.9.attn\_norm.weight requires\_grad= False  
model.layers.9.attn.Wqkv.weight requires\_grad= False  
model.layers.9.attn.Wo.weight requires\_grad= False  
model.layers.9.mlp\_norm.weight requires\_grad= False  
model.layers.9.mlp.Wi.weight requires\_grad= False  
model.layers.9.mlp.Wo.weight requires\_grad= False  
model.layers.10.attn\_norm.weight requires\_grad= False  
model.layers.10.attn.Wqkv.weight requires\_grad= False  
model.layers.10.attn.Wo.weight requires\_grad= False  
model.layers.10.mlp\_norm.weight requires\_grad= False  
model.layers.10.mlp.Wi.weight requires\_grad= False  
model.layers.10.mlp.Wo.weight requires\_grad= False  
model.layers.11.attn\_norm.weight requires\_grad= False  
model.layers.11.attn.Wqkv.weight requires\_grad= False  
model.layers.11.attn.Wo.weight requires\_grad= False  
model.layers.11.mlp\_norm.weight requires\_grad= False  
model.layers.11.mlp.Wi.weight requires\_grad= False  
model.layers.11.mlp.Wo.weight requires\_grad= False  
model.layers.12.attn\_norm.weight requires\_grad= False  
model.layers.12.attn.Wqkv.weight requires\_grad= False  
model.layers.12.attn.Wo.weight requires\_grad= False  
model.layers.12.mlp\_norm.weight requires\_grad= False  
model.layers.12.mlp.Wi.weight requires\_grad= False  
model.layers.12.mlp.Wo.weight requires\_grad= False  
model.layers.13.attn\_norm.weight requires\_grad= False  
model.layers.13.attn.Wqkv.weight requires\_grad= False  
model.layers.13.attn.Wo.weight requires\_grad= False  
model.layers.13.mlp\_norm.weight requires\_grad= False  
model.layers.13.mlp.Wi.weight requires\_grad= False  
model.layers.13.mlp.Wo.weight requires\_grad= False  
model.layers.14.attn\_norm.weight requires\_grad= False  
model.layers.14.attn.Wqkv.weight requires\_grad= False  
model.layers.14.attn.Wo.weight requires\_grad= False  
model.layers.14.mlp\_norm.weight requires\_grad= False  
model.layers.14.mlp.Wi.weight requires\_grad= False  
model.layers.14.mlp.Wo.weight requires\_grad= False  
model.layers.15.attn\_norm.weight requires\_grad= False  
model.layers.15.attn.Wqkv.weight requires\_grad= False  
model.layers.15.attn.Wo.weight requires\_grad= False  
model.layers.15.mlp\_norm.weight requires\_grad= False  
model.layers.15.mlp.Wi.weight requires\_grad= False  
model.layers.15.mlp.Wo.weight requires\_grad= False  
model.layers.16.attn\_norm.weight requires\_grad= False  
model.layers.16.attn.Wqkv.weight requires\_grad= False  
model.layers.16.attn.Wo.weight requires\_grad= False  
model.layers.16.mlp\_norm.weight requires\_grad= False  
model.layers.16.mlp.Wi.weight requires\_grad= False  
model.layers.16.mlp.Wo.weight requires\_grad= False

```

model.layers.17.attn_norm.weight requires_grad= False
model.layers.17.attn.Wqkv.weight requires_grad= False
model.layers.17.attn.Wo.weight requires_grad= False
model.layers.17.mlp_norm.weight requires_grad= False
model.layers.17.mlp.Wi.weight requires_grad= False
model.layers.17.mlp.Wo.weight requires_grad= False
model.layers.18.attn_norm.weight requires_grad= False
model.layers.18.attn.Wqkv.weight requires_grad= False
model.layers.18.attn.Wo.weight requires_grad= False
model.layers.18.mlp_norm.weight requires_grad= False
model.layers.18.mlp.Wi.weight requires_grad= False
model.layers.18.mlp.Wo.weight requires_grad= False
model.layers.19.attn_norm.weight requires_grad= False
model.layers.19.attn.Wqkv.weight requires_grad= False
model.layers.19.attn.Wo.weight requires_grad= False
model.layers.19.mlp_norm.weight requires_grad= False
model.layers.19.mlp.Wi.weight requires_grad= False
model.layers.19.mlp.Wo.weight requires_grad= False
model.layers.20.attn_norm.weight requires_grad= False
model.layers.20.attn.Wqkv.weight requires_grad= False
model.layers.20.attn.Wo.weight requires_grad= False
model.layers.20.mlp_norm.weight requires_grad= False
model.layers.20.mlp.Wi.weight requires_grad= False
model.layers.20.mlp.Wo.weight requires_grad= False
model.layers.21.attn_norm.weight requires_grad= True
model.layers.21.attn.Wqkv.weight requires_grad= True
model.layers.21.attn.Wo.weight requires_grad= True
model.layers.21.mlp_norm.weight requires_grad= True
model.layers.21.mlp.Wi.weight requires_grad= True
model.layers.21.mlp.Wo.weight requires_grad= True
model.final_norm.weight requires_grad= True
head.dense.weight requires_grad= True
head.norm.weight requires_grad= True
classifier.weight requires_grad= True
classifier.bias requires_grad= True

```

```
[49]: # model.resize_token_embeddings(len(tokenizer))
```

```
[50]: # Train & Evaluate
trained_model, trainer_obj = train_transformer_model(
    model = model,
    tokenizer = tokenizer,
    train_dataset = train_data_hf,
    val_dataset = val_data_hf,
    output_dir = dir_results,
    num_epochs = num_epochs,
    batch_size = size_batch,
```

```

    lr = learning_rate,
    weight_decay = regularization_weight_decay)
metrics = trainer_obj.evaluate()
print("Validation metrics:", metrics)
test_metrics = trainer_obj.evaluate(test_data_hf) if test_data_hf else None
print("Test metrics:", test_metrics)

```

```

/usr/local/lib/python3.11/dist-packages/transformers/training_args.py:1611:
FutureWarning: `evaluation_strategy` is deprecated and will be removed in
version 4.46 of Transformers. Use `eval_strategy` instead
    warnings.warn(
<ipython-input-20-c2ee9f934517>:31: FutureWarning: `tokenizer` is deprecated and
will be removed in version 5.0.0 for `Trainer.__init__`. Use `processing_class`
instead.
    trainer = Trainer(
/usr/local/lib/python3.11/dist-packages/torch/_inductor/compile_fx.py:194:
UserWarning: TensorFloat32 tensor cores for float32 matrix multiplication
available but not enabled. Consider setting
`torch.set_float32_matmul_precision('high')` for better performance.
    warnings.warn(
<IPython.core.display.HTML object>
<IPython.core.display.HTML object>

Validation metrics: {'eval_loss': 2.580798387527466, 'eval_accuracy':
0.22565320665083136, 'eval_precision': 0.21654501216545013, 'eval_recall':
0.956989247311828, 'eval_f1': 0.3531746031746032, 'eval_runtime': 6.0138,
'eval_samples_per_second': 70.006, 'eval_steps_per_second': 0.665, 'epoch': 1.0}
Test metrics: {'eval_loss': 2.600764274597168, 'eval_accuracy':
0.22464558342420937, 'eval_precision': 0.20495495495495494, 'eval_recall':
0.9732620320855615, 'eval_f1': 0.3386046511627907, 'eval_runtime': 6.5903,
'eval_samples_per_second': 139.144, 'eval_steps_per_second': 1.214, 'epoch':
1.0}

```

```

[51]: # save model checkpoint
timestamp = datetime.now().strftime("%Y%m%d_%H%M%S")
model_save_path = os.path.join(dir_models,
    ↪f"{x_task}_{named_model}_{y_col}_{timestamp}")
trainer_obj.save_model(model_save_path)
print(f"Model checkpoint saved to: {model_save_path}")
# log experiment results
experiment_info = {
    "model_name": named_model,
    "learning_rate": learning_rate,
    "epochs": num_epochs,
    "batch_size": size_batch,
    "weight_decay": regularization_weight_decay,
    "x_task": x_task,

```



```

        "x_col": x_col,
        "y_col": y_col,
        "layers_to_unfreeze": layers_to_unfreeze}
model_info = gather_model_details(trained_model)
all_run_metrics = gather_all_run_metrics(
    trainer=trainer_obj,
    train_dataset=train_data_hf,
    val_dataset=val_data_hf,
    test_dataset=test_data_hf)
log_experiment_results_json(
    experiment_meta=experiment_info,
    model_details=model_info,
    run_metrics=all_run_metrics,
    log_file=log_filepath)
print(f"EXPERIMENT LOGGED TO: {log_filepath}")

```

Model checkpoint saved to:

/content/drive/MyDrive/266-final/models/single\_answerdotai/ModernBERT-  
base\_binary\_complexity\_75th\_split\_20250411\_010723

<IPython.core.display.HTML object>

EXPERIMENT LOGGED TO:

/content/drive/MyDrive/266-final/results/experiment\_runs.txt

**0.2.8** snc answerdotai/ModernBERT-large regularization\_weight\_decay = 0.5 learning\_rate = 5e-6 size\_batch = 128 length\_max = 128 num\_epochs = 1

```

[52]: # Define Experiment Parameters
# named_model = "bert-base-cased"
# named_model = "roberta-base"
# named_model = "bert-large-cased"
# named_model = "roberta-large"
named_model = "answerdotai/ModernBERT-large" # modern bert
#####
regularization_weight_decay = 0.5
learning_rate = 5e-6
size_batch = 128
length_max = 128
num_epochs = 1
#####
# x_col = "sentence"
x_col = "sentence_no_contractions"
# x_col = "pos_sequence"
# x_col = "dep_sequence"
# x_col = "morph_sequence"
#####
y_col = "binary_complexity_75th_split"

```

```

# y_col = "binary_complexity"
# y_col = "complexity"
#####
x_task = "single"
# x_task = "multi"
if x_task == "single":
    df_train = train_single_df
    df_val = trial_val_single_df
    df_test = test_single_df
else:
    df_train = train_multi_df
    df_val = trial_val_multi_df
    df_test = test_multi_df
#####
# Tokenize & Prepare Datasets
train_data_hf = prepare_dataset(
    df_train,
    tokenizer,
    text_col=x_col,
    label_col=y_col,
    max_length=length_max)
val_data_hf = prepare_dataset(
    df_val,
    tokenizer,
    text_col=x_col,
    label_col=y_col,
    max_length=length_max)
test_data_hf = prepare_dataset(
    df_test,
    tokenizer,
    text_col=x_col,
    label_col=y_col,
    max_length=length_max)
print("Datasets prepared. Sample from train_data_hf:\n", train_data_hf[10])
# print("Datasets prepared. Sample from train_data_hf:\n", val_data_hf[10])
# print("Datasets prepared. Sample from train_data_hf:\n", test_data_hf[10])
#####
model, tokenizer = get_model_and_tokenizer(
    remote_model_name="answerdotai/ModernBERT-large",
    local_model_path=None,
    config=None)
#####
# model, tokenizer = get_model_and_tokenizer(
#     remote_model_name=None
#     local_model_path="...CONFIGURE_PATH...",
#     config=custom_config)
#####

```



```
['classifier.bias', 'classifier.weight']
```

You should probably TRAIN this model on a down-stream task to be able to use it for predictions and inference.

```
=====
```

```
answerdotai/ModernBERT-large :
```

```
=====
```

```
num_parameters: 395833346
```

```
num_trainable_parameters at load: 395833346
```

```
=====
```

```
model lineage: {'type': 'huggingface_hub', 'path': 'answerdotai/ModernBERT-large', 'timestamp': '2025-04-11 01:08:09'}
```

```
=====
```

```
[53]: print(model)
```

```
ModernBertForSequenceClassification(
  (model): ModernBertModel(
    (embeddings): ModernBertEmbeddings(
      (tok_embeddings): Embedding(50368, 1024, padding_idx=50283)
      (norm): LayerNorm((1024,), eps=1e-05, elementwise_affine=True)
      (drop): Dropout(p=0.0, inplace=False)
    )
    (layers): ModuleList(
      (0): ModernBertEncoderLayer(
        (attn_norm): Identity()
        (attn): ModernBertAttention(
          (Wqkv): Linear(in_features=1024, out_features=3072, bias=False)
          (rotary_emb): ModernBertRotaryEmbedding()
          (Wo): Linear(in_features=1024, out_features=1024, bias=False)
          (out_drop): Identity()
        )
        (mlp_norm): LayerNorm((1024,), eps=1e-05, elementwise_affine=True)
        (mlp): ModernBertMLP(
          (Wi): Linear(in_features=1024, out_features=5248, bias=False)
          (act): GELUActivation()
          (drop): Dropout(p=0.0, inplace=False)
          (Wo): Linear(in_features=2624, out_features=1024, bias=False)
        )
      )
    )
    (1-27): 27 x ModernBertEncoderLayer(
      (attn_norm): LayerNorm((1024,), eps=1e-05, elementwise_affine=True)
      (attn): ModernBertAttention(
        (Wqkv): Linear(in_features=1024, out_features=3072, bias=False)
        (rotary_emb): ModernBertRotaryEmbedding()
        (Wo): Linear(in_features=1024, out_features=1024, bias=False)
        (out_drop): Identity()
      )
      (mlp_norm): LayerNorm((1024,), eps=1e-05, elementwise_affine=True)
```

```

        (mlp): ModernBertMLP(
          (Wi): Linear(in_features=1024, out_features=5248, bias=False)
          (act): GELUActivation()
          (drop): Dropout(p=0.0, inplace=False)
          (Wo): Linear(in_features=2624, out_features=1024, bias=False)
        )
      )
    )
    (final_norm): LayerNorm((1024,), eps=1e-05, elementwise_affine=True)
  )
  (head): ModernBertPredictionHead(
    (dense): Linear(in_features=1024, out_features=1024, bias=False)
    (act): GELUActivation()
    (norm): LayerNorm((1024,), eps=1e-05, elementwise_affine=True)
  )
  (drop): Dropout(p=0.0, inplace=False)
  (classifier): Linear(in_features=1024, out_features=2, bias=True)
)

```

```

[54]: for name, param in model.named_parameters():
      print(name, "requires_grad=", param.requires_grad)

```

```

model.embeddings.tok_embeddings.weight requires_grad= True
model.embeddings.norm.weight requires_grad= True
model.layers.0.attn.Wqkv.weight requires_grad= True
model.layers.0.attn.Wo.weight requires_grad= True
model.layers.0.mlp_norm.weight requires_grad= True
model.layers.0.mlp.Wi.weight requires_grad= True
model.layers.0.mlp.Wo.weight requires_grad= True
model.layers.1.attn_norm.weight requires_grad= True
model.layers.1.attn.Wqkv.weight requires_grad= True
model.layers.1.attn.Wo.weight requires_grad= True
model.layers.1.mlp_norm.weight requires_grad= True
model.layers.1.mlp.Wi.weight requires_grad= True
model.layers.1.mlp.Wo.weight requires_grad= True
model.layers.2.attn_norm.weight requires_grad= True
model.layers.2.attn.Wqkv.weight requires_grad= True
model.layers.2.attn.Wo.weight requires_grad= True
model.layers.2.mlp_norm.weight requires_grad= True
model.layers.2.mlp.Wi.weight requires_grad= True
model.layers.2.mlp.Wo.weight requires_grad= True
model.layers.3.attn_norm.weight requires_grad= True
model.layers.3.attn.Wqkv.weight requires_grad= True
model.layers.3.attn.Wo.weight requires_grad= True
model.layers.3.mlp_norm.weight requires_grad= True
model.layers.3.mlp.Wi.weight requires_grad= True
model.layers.3.mlp.Wo.weight requires_grad= True
model.layers.4.attn_norm.weight requires_grad= True

```

model.layers.4.attn.Wqkv.weight requires\_grad= True  
model.layers.4.attn.Wo.weight requires\_grad= True  
model.layers.4.mlp\_norm.weight requires\_grad= True  
model.layers.4.mlp.Wi.weight requires\_grad= True  
model.layers.4.mlp.Wo.weight requires\_grad= True  
model.layers.5.attn\_norm.weight requires\_grad= True  
model.layers.5.attn.Wqkv.weight requires\_grad= True  
model.layers.5.attn.Wo.weight requires\_grad= True  
model.layers.5.mlp\_norm.weight requires\_grad= True  
model.layers.5.mlp.Wi.weight requires\_grad= True  
model.layers.5.mlp.Wo.weight requires\_grad= True  
model.layers.6.attn\_norm.weight requires\_grad= True  
model.layers.6.attn.Wqkv.weight requires\_grad= True  
model.layers.6.attn.Wo.weight requires\_grad= True  
model.layers.6.mlp\_norm.weight requires\_grad= True  
model.layers.6.mlp.Wi.weight requires\_grad= True  
model.layers.6.mlp.Wo.weight requires\_grad= True  
model.layers.7.attn\_norm.weight requires\_grad= True  
model.layers.7.attn.Wqkv.weight requires\_grad= True  
model.layers.7.attn.Wo.weight requires\_grad= True  
model.layers.7.mlp\_norm.weight requires\_grad= True  
model.layers.7.mlp.Wi.weight requires\_grad= True  
model.layers.7.mlp.Wo.weight requires\_grad= True  
model.layers.8.attn\_norm.weight requires\_grad= True  
model.layers.8.attn.Wqkv.weight requires\_grad= True  
model.layers.8.attn.Wo.weight requires\_grad= True  
model.layers.8.mlp\_norm.weight requires\_grad= True  
model.layers.8.mlp.Wi.weight requires\_grad= True  
model.layers.8.mlp.Wo.weight requires\_grad= True  
model.layers.9.attn\_norm.weight requires\_grad= True  
model.layers.9.attn.Wqkv.weight requires\_grad= True  
model.layers.9.attn.Wo.weight requires\_grad= True  
model.layers.9.mlp\_norm.weight requires\_grad= True  
model.layers.9.mlp.Wi.weight requires\_grad= True  
model.layers.9.mlp.Wo.weight requires\_grad= True  
model.layers.10.attn\_norm.weight requires\_grad= True  
model.layers.10.attn.Wqkv.weight requires\_grad= True  
model.layers.10.attn.Wo.weight requires\_grad= True  
model.layers.10.mlp\_norm.weight requires\_grad= True  
model.layers.10.mlp.Wi.weight requires\_grad= True  
model.layers.10.mlp.Wo.weight requires\_grad= True  
model.layers.11.attn\_norm.weight requires\_grad= True  
model.layers.11.attn.Wqkv.weight requires\_grad= True  
model.layers.11.attn.Wo.weight requires\_grad= True  
model.layers.11.mlp\_norm.weight requires\_grad= True  
model.layers.11.mlp.Wi.weight requires\_grad= True  
model.layers.11.mlp.Wo.weight requires\_grad= True  
model.layers.12.attn\_norm.weight requires\_grad= True

model.layers.12.attn.Wqkv.weight requires\_grad= True  
model.layers.12.attn.Wo.weight requires\_grad= True  
model.layers.12.mlp\_norm.weight requires\_grad= True  
model.layers.12.mlp.Wi.weight requires\_grad= True  
model.layers.12.mlp.Wo.weight requires\_grad= True  
model.layers.13.attn\_norm.weight requires\_grad= True  
model.layers.13.attn.Wqkv.weight requires\_grad= True  
model.layers.13.attn.Wo.weight requires\_grad= True  
model.layers.13.mlp\_norm.weight requires\_grad= True  
model.layers.13.mlp.Wi.weight requires\_grad= True  
model.layers.13.mlp.Wo.weight requires\_grad= True  
model.layers.14.attn\_norm.weight requires\_grad= True  
model.layers.14.attn.Wqkv.weight requires\_grad= True  
model.layers.14.attn.Wo.weight requires\_grad= True  
model.layers.14.mlp\_norm.weight requires\_grad= True  
model.layers.14.mlp.Wi.weight requires\_grad= True  
model.layers.14.mlp.Wo.weight requires\_grad= True  
model.layers.15.attn\_norm.weight requires\_grad= True  
model.layers.15.attn.Wqkv.weight requires\_grad= True  
model.layers.15.attn.Wo.weight requires\_grad= True  
model.layers.15.mlp\_norm.weight requires\_grad= True  
model.layers.15.mlp.Wi.weight requires\_grad= True  
model.layers.15.mlp.Wo.weight requires\_grad= True  
model.layers.16.attn\_norm.weight requires\_grad= True  
model.layers.16.attn.Wqkv.weight requires\_grad= True  
model.layers.16.attn.Wo.weight requires\_grad= True  
model.layers.16.mlp\_norm.weight requires\_grad= True  
model.layers.16.mlp.Wi.weight requires\_grad= True  
model.layers.16.mlp.Wo.weight requires\_grad= True  
model.layers.17.attn\_norm.weight requires\_grad= True  
model.layers.17.attn.Wqkv.weight requires\_grad= True  
model.layers.17.attn.Wo.weight requires\_grad= True  
model.layers.17.mlp\_norm.weight requires\_grad= True  
model.layers.17.mlp.Wi.weight requires\_grad= True  
model.layers.17.mlp.Wo.weight requires\_grad= True  
model.layers.18.attn\_norm.weight requires\_grad= True  
model.layers.18.attn.Wqkv.weight requires\_grad= True  
model.layers.18.attn.Wo.weight requires\_grad= True  
model.layers.18.mlp\_norm.weight requires\_grad= True  
model.layers.18.mlp.Wi.weight requires\_grad= True  
model.layers.18.mlp.Wo.weight requires\_grad= True  
model.layers.19.attn\_norm.weight requires\_grad= True  
model.layers.19.attn.Wqkv.weight requires\_grad= True  
model.layers.19.attn.Wo.weight requires\_grad= True  
model.layers.19.mlp\_norm.weight requires\_grad= True  
model.layers.19.mlp.Wi.weight requires\_grad= True  
model.layers.19.mlp.Wo.weight requires\_grad= True  
model.layers.20.attn\_norm.weight requires\_grad= True

model.layers.20.attn.Wqkv.weight requires\_grad= True  
model.layers.20.attn.Wo.weight requires\_grad= True  
model.layers.20.mlp\_norm.weight requires\_grad= True  
model.layers.20.mlp.Wi.weight requires\_grad= True  
model.layers.20.mlp.Wo.weight requires\_grad= True  
model.layers.21.attn\_norm.weight requires\_grad= True  
model.layers.21.attn.Wqkv.weight requires\_grad= True  
model.layers.21.attn.Wo.weight requires\_grad= True  
model.layers.21.mlp\_norm.weight requires\_grad= True  
model.layers.21.mlp.Wi.weight requires\_grad= True  
model.layers.21.mlp.Wo.weight requires\_grad= True  
model.layers.22.attn\_norm.weight requires\_grad= True  
model.layers.22.attn.Wqkv.weight requires\_grad= True  
model.layers.22.attn.Wo.weight requires\_grad= True  
model.layers.22.mlp\_norm.weight requires\_grad= True  
model.layers.22.mlp.Wi.weight requires\_grad= True  
model.layers.22.mlp.Wo.weight requires\_grad= True  
model.layers.23.attn\_norm.weight requires\_grad= True  
model.layers.23.attn.Wqkv.weight requires\_grad= True  
model.layers.23.attn.Wo.weight requires\_grad= True  
model.layers.23.mlp\_norm.weight requires\_grad= True  
model.layers.23.mlp.Wi.weight requires\_grad= True  
model.layers.23.mlp.Wo.weight requires\_grad= True  
model.layers.24.attn\_norm.weight requires\_grad= True  
model.layers.24.attn.Wqkv.weight requires\_grad= True  
model.layers.24.attn.Wo.weight requires\_grad= True  
model.layers.24.mlp\_norm.weight requires\_grad= True  
model.layers.24.mlp.Wi.weight requires\_grad= True  
model.layers.24.mlp.Wo.weight requires\_grad= True  
model.layers.25.attn\_norm.weight requires\_grad= True  
model.layers.25.attn.Wqkv.weight requires\_grad= True  
model.layers.25.attn.Wo.weight requires\_grad= True  
model.layers.25.mlp\_norm.weight requires\_grad= True  
model.layers.25.mlp.Wi.weight requires\_grad= True  
model.layers.25.mlp.Wo.weight requires\_grad= True  
model.layers.26.attn\_norm.weight requires\_grad= True  
model.layers.26.attn.Wqkv.weight requires\_grad= True  
model.layers.26.attn.Wo.weight requires\_grad= True  
model.layers.26.mlp\_norm.weight requires\_grad= True  
model.layers.26.mlp.Wi.weight requires\_grad= True  
model.layers.26.mlp.Wo.weight requires\_grad= True  
model.layers.27.attn\_norm.weight requires\_grad= True  
model.layers.27.attn.Wqkv.weight requires\_grad= True  
model.layers.27.attn.Wo.weight requires\_grad= True  
model.layers.27.mlp\_norm.weight requires\_grad= True  
model.layers.27.mlp.Wi.weight requires\_grad= True  
model.layers.27.mlp.Wo.weight requires\_grad= True  
model.final\_norm.weight requires\_grad= True



```
head.dense.weight requires_grad= True
head.norm.weight requires_grad= True
classifier.weight requires_grad= True
classifier.bias requires_grad= True
```

```
[55]: # Inspect the attention_mask tensor for the first few samples
```

```
for i in range(5):
    print(train_data_hf[i]['attention_mask'])
```

```
tensor([1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1,
        1, 1, 1, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0,
        0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0,
        0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0,
        0, 0, 0, 0, 0, 0, 0, 0])
tensor([1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1,
        1, 1, 1, 1, 1, 1, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0,
        0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0,
        0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0,
        0, 0, 0, 0, 0, 0, 0, 0])
tensor([1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1,
        1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1,
        1, 1, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0,
        0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0,
        0, 0, 0, 0, 0, 0, 0, 0])
tensor([1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1,
        1, 1, 1, 1, 1, 1, 1, 1, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0,
        0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0,
        0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0,
        0, 0, 0, 0, 0, 0, 0, 0])
tensor([1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0,
        0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0,
        0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0,
        0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0,
        0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0,
        0, 0, 0, 0, 0, 0, 0, 0])
```

```
[56]: #####
```

```
layers_to_unfreeze = [
    "model.layers.27.attn_norm.weight",
    "model.layers.27.attn.Wqkv.weight",
    "model.layers.27.attn.Wo.weight",
    "model.layers.27.mlp_norm.weight",
    "model.layers.27.mlp.Wi.weight",
    "model.layers.27.mlp.Wo.weight",
```

```

        "model.final_norm.weight",
        "head.dense.weight",
        "head.norm.weight",
        "classifier.weight",
        "classifier.bias"
    ]
freeze_unfreeze_layers(model, layers_to_unfreeze=layers_to_unfreeze)
print(model.config)
print("=====")
print("num_parameters:", model.num_parameters())
print("num_trainable_parameters:", model.num_parameters(only_trainable=True))
print("=====")
#####
print("Experiment configuration used with this experiment:")
print("model used:", named_model)
print("learning rate used:", learning_rate)
print("number of epochs:", num_epochs)
print("maximum sequence length:", length_max)
print("batch size used:", size_batch)
print("regularization value:", regularization_weight_decay)
print("outcome variable:", y_col)
print("task:", x_task)
print("input column:", x_col)
#####
print("=====")
print("num_trainable_parameters:", model.num_parameters(only_trainable=True))

```

```

ModernBertConfig {
  "_attn_implementation_autoset": true,
  "architectures": [
    "ModernBertForMaskedLM"
  ],
  "attention_bias": false,
  "attention_dropout": 0.0,
  "bos_token_id": 50281,
  "classifier_activation": "gelu",
  "classifier_bias": false,
  "classifier_dropout": 0.0,
  "classifier_pooling": "mean",
  "cls_token_id": 50281,
  "decoder_bias": true,
  "deterministic_flash_attn": false,
  "embedding_dropout": 0.0,
  "eos_token_id": 50282,
  "global_attn_every_n_layers": 3,
  "global_rope_theta": 160000.0,
  "gradient_checkpointing": false,
  "hidden_activation": "gelu",

```

```

"hidden_size": 1024,
"initializer_cutoff_factor": 2.0,
"initializer_range": 0.02,
"intermediate_size": 2624,
"layer_norm_eps": 1e-05,
"local_attention": 128,
"local_rope_theta": 10000.0,
"max_position_embeddings": 8192,
"mlp_bias": false,
"mlp_dropout": 0.0,
"model_type": "modernbert",
"norm_bias": false,
"norm_eps": 1e-05,
"num_attention_heads": 16,
"num_hidden_layers": 28,
"pad_token_id": 50283,
"position_embedding_type": "absolute",
"reference_compile": null,
"repad_logits_with_grad": false,
"sep_token_id": 50282,
"sparse_pred_ignore_index": -100,
"sparse_prediction": false,
"torch_dtype": "float32",
"transformers_version": "4.50.3",
"vocab_size": 50368
}

```

=====

```

num_parameters: 395833346
num_trainable_parameters: 13309954
=====

```

Experiment configuration used with this experiment:

```

model used: answerdotai/ModernBERT-large
learning rate used: 5e-06
number of epochs: 1
maximum sequence length: 128
batch size used: 128
regularization value: 0.5
outcome variable: binary_complexity_75th_split
task: single
input column: sentence_no_contractions
=====

```

```

num_trainable_parameters: 13309954

```

```

[57]: for name, param in model.named_parameters():
      print(name, "requires_grad=", param.requires_grad)

```

```

model.embeddings.tok_embeddings.weight requires_grad= False

```

model.embeddings.norm.weight requires\_grad= False  
model.layers.0.attn.Wqkv.weight requires\_grad= False  
model.layers.0.attn.Wo.weight requires\_grad= False  
model.layers.0.mlp\_norm.weight requires\_grad= False  
model.layers.0.mlp.Wi.weight requires\_grad= False  
model.layers.0.mlp.Wo.weight requires\_grad= False  
model.layers.1.attn\_norm.weight requires\_grad= False  
model.layers.1.attn.Wqkv.weight requires\_grad= False  
model.layers.1.attn.Wo.weight requires\_grad= False  
model.layers.1.mlp\_norm.weight requires\_grad= False  
model.layers.1.mlp.Wi.weight requires\_grad= False  
model.layers.1.mlp.Wo.weight requires\_grad= False  
model.layers.2.attn\_norm.weight requires\_grad= False  
model.layers.2.attn.Wqkv.weight requires\_grad= False  
model.layers.2.attn.Wo.weight requires\_grad= False  
model.layers.2.mlp\_norm.weight requires\_grad= False  
model.layers.2.mlp.Wi.weight requires\_grad= False  
model.layers.2.mlp.Wo.weight requires\_grad= False  
model.layers.3.attn\_norm.weight requires\_grad= False  
model.layers.3.attn.Wqkv.weight requires\_grad= False  
model.layers.3.attn.Wo.weight requires\_grad= False  
model.layers.3.mlp\_norm.weight requires\_grad= False  
model.layers.3.mlp.Wi.weight requires\_grad= False  
model.layers.3.mlp.Wo.weight requires\_grad= False  
model.layers.4.attn\_norm.weight requires\_grad= False  
model.layers.4.attn.Wqkv.weight requires\_grad= False  
model.layers.4.attn.Wo.weight requires\_grad= False  
model.layers.4.mlp\_norm.weight requires\_grad= False  
model.layers.4.mlp.Wi.weight requires\_grad= False  
model.layers.4.mlp.Wo.weight requires\_grad= False  
model.layers.5.attn\_norm.weight requires\_grad= False  
model.layers.5.attn.Wqkv.weight requires\_grad= False  
model.layers.5.attn.Wo.weight requires\_grad= False  
model.layers.5.mlp\_norm.weight requires\_grad= False  
model.layers.5.mlp.Wi.weight requires\_grad= False  
model.layers.5.mlp.Wo.weight requires\_grad= False  
model.layers.6.attn\_norm.weight requires\_grad= False  
model.layers.6.attn.Wqkv.weight requires\_grad= False  
model.layers.6.attn.Wo.weight requires\_grad= False  
model.layers.6.mlp\_norm.weight requires\_grad= False  
model.layers.6.mlp.Wi.weight requires\_grad= False  
model.layers.6.mlp.Wo.weight requires\_grad= False  
model.layers.7.attn\_norm.weight requires\_grad= False  
model.layers.7.attn.Wqkv.weight requires\_grad= False  
model.layers.7.attn.Wo.weight requires\_grad= False  
model.layers.7.mlp\_norm.weight requires\_grad= False  
model.layers.7.mlp.Wi.weight requires\_grad= False  
model.layers.7.mlp.Wo.weight requires\_grad= False

model.layers.8.attn\_norm.weight requires\_grad= False  
model.layers.8.attn.Wqkv.weight requires\_grad= False  
model.layers.8.attn.Wo.weight requires\_grad= False  
model.layers.8.mlp\_norm.weight requires\_grad= False  
model.layers.8.mlp.Wi.weight requires\_grad= False  
model.layers.8.mlp.Wo.weight requires\_grad= False  
model.layers.9.attn\_norm.weight requires\_grad= False  
model.layers.9.attn.Wqkv.weight requires\_grad= False  
model.layers.9.attn.Wo.weight requires\_grad= False  
model.layers.9.mlp\_norm.weight requires\_grad= False  
model.layers.9.mlp.Wi.weight requires\_grad= False  
model.layers.9.mlp.Wo.weight requires\_grad= False  
model.layers.10.attn\_norm.weight requires\_grad= False  
model.layers.10.attn.Wqkv.weight requires\_grad= False  
model.layers.10.attn.Wo.weight requires\_grad= False  
model.layers.10.mlp\_norm.weight requires\_grad= False  
model.layers.10.mlp.Wi.weight requires\_grad= False  
model.layers.10.mlp.Wo.weight requires\_grad= False  
model.layers.11.attn\_norm.weight requires\_grad= False  
model.layers.11.attn.Wqkv.weight requires\_grad= False  
model.layers.11.attn.Wo.weight requires\_grad= False  
model.layers.11.mlp\_norm.weight requires\_grad= False  
model.layers.11.mlp.Wi.weight requires\_grad= False  
model.layers.11.mlp.Wo.weight requires\_grad= False  
model.layers.12.attn\_norm.weight requires\_grad= False  
model.layers.12.attn.Wqkv.weight requires\_grad= False  
model.layers.12.attn.Wo.weight requires\_grad= False  
model.layers.12.mlp\_norm.weight requires\_grad= False  
model.layers.12.mlp.Wi.weight requires\_grad= False  
model.layers.12.mlp.Wo.weight requires\_grad= False  
model.layers.13.attn\_norm.weight requires\_grad= False  
model.layers.13.attn.Wqkv.weight requires\_grad= False  
model.layers.13.attn.Wo.weight requires\_grad= False  
model.layers.13.mlp\_norm.weight requires\_grad= False  
model.layers.13.mlp.Wi.weight requires\_grad= False  
model.layers.13.mlp.Wo.weight requires\_grad= False  
model.layers.14.attn\_norm.weight requires\_grad= False  
model.layers.14.attn.Wqkv.weight requires\_grad= False  
model.layers.14.attn.Wo.weight requires\_grad= False  
model.layers.14.mlp\_norm.weight requires\_grad= False  
model.layers.14.mlp.Wi.weight requires\_grad= False  
model.layers.14.mlp.Wo.weight requires\_grad= False  
model.layers.15.attn\_norm.weight requires\_grad= False  
model.layers.15.attn.Wqkv.weight requires\_grad= False  
model.layers.15.attn.Wo.weight requires\_grad= False  
model.layers.15.mlp\_norm.weight requires\_grad= False  
model.layers.15.mlp.Wi.weight requires\_grad= False  
model.layers.15.mlp.Wo.weight requires\_grad= False



```

model.layers.24.attn_norm.weight requires_grad= False
model.layers.24.attn.Wqkv.weight requires_grad= False
model.layers.24.attn.Wo.weight requires_grad= False
model.layers.24.mlp_norm.weight requires_grad= False
model.layers.24.mlp.Wi.weight requires_grad= False
model.layers.24.mlp.Wo.weight requires_grad= False
model.layers.25.attn_norm.weight requires_grad= False
model.layers.25.attn.Wqkv.weight requires_grad= False
model.layers.25.attn.Wo.weight requires_grad= False
model.layers.25.mlp_norm.weight requires_grad= False
model.layers.25.mlp.Wi.weight requires_grad= False
model.layers.25.mlp.Wo.weight requires_grad= False
model.layers.26.attn_norm.weight requires_grad= False
model.layers.26.attn.Wqkv.weight requires_grad= False
model.layers.26.attn.Wo.weight requires_grad= False
model.layers.26.mlp_norm.weight requires_grad= False
model.layers.26.mlp.Wi.weight requires_grad= False
model.layers.26.mlp.Wo.weight requires_grad= False
model.layers.27.attn_norm.weight requires_grad= True
model.layers.27.attn.Wqkv.weight requires_grad= True
model.layers.27.attn.Wo.weight requires_grad= True
model.layers.27.mlp_norm.weight requires_grad= True
model.layers.27.mlp.Wi.weight requires_grad= True
model.layers.27.mlp.Wo.weight requires_grad= True
model.final_norm.weight requires_grad= True
head.dense.weight requires_grad= True
head.norm.weight requires_grad= True
classifier.weight requires_grad= True
classifier.bias requires_grad= True

```

```
[58]: # model.resize_token_embeddings(len(tokenizer))
```

```
[59]: # Train & Evaluate
trained_model, trainer_obj = train_transformer_model(
    model = model,
    tokenizer = tokenizer,
    train_dataset = train_data_hf,
    val_dataset = val_data_hf,
    output_dir = dir_results,
    num_epochs = num_epochs,
    batch_size = size_batch,
    lr = learning_rate,
    weight_decay = regularization_weight_decay)
metrics = trainer_obj.evaluate()
print("Validation metrics:", metrics)
test_metrics = trainer_obj.evaluate(test_data_hf) if test_data_hf else None
print("Test metrics:", test_metrics)
```

```

/usr/local/lib/python3.11/dist-packages/transformers/training_args.py:1611:
FutureWarning: `evaluation_strategy` is deprecated and will be removed in
version 4.46 of Transformers. Use `eval_strategy` instead
    warnings.warn(
<ipython-input-20-c2ee9f934517>:31: FutureWarning: `tokenizer` is deprecated and
will be removed in version 5.0.0 for `Trainer.__init__`. Use `processing_class`
instead.
    trainer = Trainer(

<IPython.core.display.HTML object>

W0411 01:08:59.809000 822 torch/_dynamo/convert_frame.py:906] [1/8]
torch._dynamo hit config.cache_size_limit (8)
W0411 01:08:59.809000 822 torch/_dynamo/convert_frame.py:906] [1/8]    function:
'compiled_mlp' (/usr/local/lib/python3.11/dist-
packages/transformers/models/modernbert/modeling_modernbert.py:552)
W0411 01:08:59.809000 822 torch/_dynamo/convert_frame.py:906] [1/8]    last
reason: 1/0: GLOBAL_STATE changed: grad_mode
W0411 01:08:59.809000 822 torch/_dynamo/convert_frame.py:906] [1/8] To log all
recompilation reasons, use TORCH_LOGS="recompiles".
W0411 01:08:59.809000 822 torch/_dynamo/convert_frame.py:906] [1/8] To diagnose
recompilation issues, see
https://pytorch.org/docs/main/torch.compiler\_troubleshooting.html.

<IPython.core.display.HTML object>

Validation metrics: {'eval_loss': 0.6973762512207031, 'eval_accuracy':
0.5415676959619953, 'eval_precision': 0.2222222222222222, 'eval_recall':
0.43010752688172044, 'eval_f1': 0.29304029304029305, 'eval_runtime': 8.4333,
'eval_samples_per_second': 49.921, 'eval_steps_per_second': 0.474, 'epoch': 1.0}
Test metrics: {'eval_loss': 0.7029761672019958, 'eval_accuracy':
0.5267175572519084, 'eval_precision': 0.20803782505910165, 'eval_recall':
0.47058823529411764, 'eval_f1': 0.28852459016393445, 'eval_runtime': 10.0378,
'eval_samples_per_second': 91.355, 'eval_steps_per_second': 0.797, 'epoch': 1.0}

```

```

[60]: # save model checkpoint
timestamp = datetime.now().strftime("%Y%m%d_%H%M%S")
model_save_path = os.path.join(dir_models,
    f"{x_task}_{named_model}_{y_col}_{timestamp}")
trainer_obj.save_model(model_save_path)
print(f"Model checkpoint saved to: {model_save_path}")
# log experiment results
experiment_info = {
    "model_name": named_model,
    "learning_rate": learning_rate,
    "epochs": num_epochs,
    "batch_size": size_batch,
    "weight_decay": regularization_weight_decay,
    "x_task": x_task,

```



```

        "x_col": x_col,
        "y_col": y_col,
        "layers_to_unfreeze": layers_to_unfreeze}
model_info = gather_model_details(trained_model)
all_run_metrics = gather_all_run_metrics(
    trainer=trainer_obj,
    train_dataset=train_data_hf,
    val_dataset=val_data_hf,
    test_dataset=test_data_hf)
log_experiment_results_json(
    experiment_meta=experiment_info,
    model_details=model_info,
    run_metrics=all_run_metrics,
    log_file=log_filepath)
print(f"EXPERIMENT LOGGED TO: {log_filepath}")

```

Model checkpoint saved to:

/content/drive/MyDrive/266-final/models/single\_answerdotai/ModernBERT-large\_binary\_complexity\_75th\_split\_20250411\_010925

<IPython.core.display.HTML object>

EXPERIMENT LOGGED TO:

/content/drive/MyDrive/266-final/results/experiment\_runs.txt

**0.2.9 snc microsoft/deberta-v3-base regularization\_weight\_decay = 0.5 learning\_rate = 5e-6 size\_batch = 128 length\_max = 128 num\_epochs = 1**

```

[61]: # Define Experiment Parameters
# named_model = "bert-base-cased"
# named_model = "roberta-base"
# named_model = "bert-large-cased"
# named_model = "roberta-large"
# named_model = "answerdotai/ModernBERT-base" # modern bert
# named_model = "answerdotai/ModernBERT-large" # modern bert
named_model = "microsoft/deberta-v3-base" # deberta
#####
regularization_weight_decay = 0.5
learning_rate = 5e-6
size_batch = 128
length_max = 128
num_epochs = 1
#####
# x_col = "sentence"
x_col = "sentence_no_contractions"
# x_col = "pos_sequence"
# x_col = "dep_sequence"
# x_col = "morph_sequence"

```

```

#####
y_col = "binary_complexity_75th_split"
# y_col = "binary_complexity"
# y_col = "complexity"
#####
x_task = "single"
# x_task = "multi"
if x_task == "single":
    df_train = train_single_df
    df_val = trial_val_single_df
    df_test = test_single_df
else:
    df_train = train_multi_df
    df_val = trial_val_multi_df
    df_test = test_multi_df
#####
# Tokenize & Prepare Datasets
train_data_hf = prepare_dataset(
    df_train,
    tokenizer,
    text_col=x_col,
    label_col=y_col,
    max_length=length_max)
val_data_hf = prepare_dataset(
    df_val,
    tokenizer,
    text_col=x_col,
    label_col=y_col,
    max_length=length_max)
test_data_hf = prepare_dataset(
    df_test,
    tokenizer,
    text_col=x_col,
    label_col=y_col,
    max_length=length_max)
print("Datasets prepared. Sample from train_data_hf:\n", train_data_hf[10])
# print("Datasets prepared. Sample from train_data_hf:\n", val_data_hf[10])
# print("Datasets prepared. Sample from train_data_hf:\n", test_data_hf[10])
#####
model, tokenizer = get_model_and_tokenizer(
    remote_model_name="microsoft/deberta-v3-base",
    local_model_path=None,
    config=None)
#####
# model, tokenizer = get_model_and_tokenizer(
#     remote_model_name=None
#     local_model_path="...CONFIGURE_PATH...",

```



byte fallback option which is not implemented in the fast tokenizers. In practice this means that the fast version of the tokenizer can produce unknown tokens whereas the sentencepiece version would have converted these unknown tokens into a sequence of byte tokens matching the original piece of text.

```
warnings.warn(

pytorch_model.bin:   0%|                               | 0.00/371M [00:00<?, ?B/s]

Some weights of DebertaV2ForSequenceClassification were not initialized from the
model checkpoint at microsoft/deberta-v3-base and are newly initialized:
['classifier.bias', 'classifier.weight', 'pooler.dense.bias',
'pooler.dense.weight']
You should probably TRAIN this model on a down-stream task to be able to use it
for predictions and inference.

=====
microsoft/deberta-v3-base :
=====
num_parameters: 184423682
num_trainable_parameters at load: 184423682
=====
model lineage: {'type': 'huggingface_hub', 'path': 'microsoft/deberta-v3-base',
'timestamp': '2025-04-11 01:10:40'}
```

```
[62]: print(model)
```

```
DebertaV2ForSequenceClassification(
  (deberta): DebertaV2Model(
    (embeddings): DebertaV2Embeddings(
      (word_embeddings): Embedding(128100, 768, padding_idx=0)
      (LayerNorm): LayerNorm((768,), eps=1e-07, elementwise_affine=True)
      (dropout): Dropout(p=0.1, inplace=False)
    )
    (encoder): DebertaV2Encoder(
      (layer): ModuleList(
        (0-11): 12 x DebertaV2Layer(
          (attention): DebertaV2Attention(
            (self): DisentangledSelfAttention(
              (query_proj): Linear(in_features=768, out_features=768, bias=True)
              (key_proj): Linear(in_features=768, out_features=768, bias=True)
              (value_proj): Linear(in_features=768, out_features=768, bias=True)
              (pos_dropout): Dropout(p=0.1, inplace=False)
              (dropout): Dropout(p=0.1, inplace=False)
            )
          (output): DebertaV2SelfOutput(
            (dense): Linear(in_features=768, out_features=768, bias=True)
            (LayerNorm): LayerNorm((768,), eps=1e-07, elementwise_affine=True)
            (dropout): Dropout(p=0.1, inplace=False)
          )
        )
      )
    )
  )
)
```

```

    )
    (intermediate): DebertaV2Intermediate(
        (dense): Linear(in_features=768, out_features=3072, bias=True)
        (intermediate_act_fn): GELUActivation()
    )
    (output): DebertaV2Output(
        (dense): Linear(in_features=3072, out_features=768, bias=True)
        (LayerNorm): LayerNorm((768,), eps=1e-07, elementwise_affine=True)
        (dropout): Dropout(p=0.1, inplace=False)
    )
)
)
(rel_embeddings): Embedding(512, 768)
(LayerNorm): LayerNorm((768,), eps=1e-07, elementwise_affine=True)
)
)
(pooler): ContextPooler(
    (dense): Linear(in_features=768, out_features=768, bias=True)
    (dropout): Dropout(p=0, inplace=False)
)
(classifier): Linear(in_features=768, out_features=2, bias=True)
(dropout): Dropout(p=0.1, inplace=False)
)

```

```

[63]: for name, param in model.named_parameters():
        print(name, "requires_grad=", param.requires_grad)

```

```

deberta.embeddings.word_embeddings.weight requires_grad= True
deberta.embeddings.LayerNorm.weight requires_grad= True
deberta.embeddings.LayerNorm.bias requires_grad= True
deberta.encoder.layer.0.attention.self.query_proj.weight requires_grad= True
deberta.encoder.layer.0.attention.self.query_proj.bias requires_grad= True
deberta.encoder.layer.0.attention.self.key_proj.weight requires_grad= True
deberta.encoder.layer.0.attention.self.key_proj.bias requires_grad= True
deberta.encoder.layer.0.attention.self.value_proj.weight requires_grad= True
deberta.encoder.layer.0.attention.self.value_proj.bias requires_grad= True
deberta.encoder.layer.0.attention.output.dense.weight requires_grad= True
deberta.encoder.layer.0.attention.output.dense.bias requires_grad= True
deberta.encoder.layer.0.attention.output.LayerNorm.weight requires_grad= True
deberta.encoder.layer.0.attention.output.LayerNorm.bias requires_grad= True
deberta.encoder.layer.0.intermediate.dense.weight requires_grad= True
deberta.encoder.layer.0.intermediate.dense.bias requires_grad= True
deberta.encoder.layer.0.output.dense.weight requires_grad= True
deberta.encoder.layer.0.output.dense.bias requires_grad= True
deberta.encoder.layer.0.output.LayerNorm.weight requires_grad= True
deberta.encoder.layer.0.output.LayerNorm.bias requires_grad= True
deberta.encoder.layer.1.attention.self.query_proj.weight requires_grad= True
deberta.encoder.layer.1.attention.self.query_proj.bias requires_grad= True

```



[illegible]





```

deberta.encoder.layer.10.attention.self.key_proj.weight requires_grad= True
deberta.encoder.layer.10.attention.self.key_proj.bias requires_grad= True
deberta.encoder.layer.10.attention.self.value_proj.weight requires_grad= True
deberta.encoder.layer.10.attention.self.value_proj.bias requires_grad= True
deberta.encoder.layer.10.attention.output.dense.weight requires_grad= True
deberta.encoder.layer.10.attention.output.dense.bias requires_grad= True
deberta.encoder.layer.10.attention.output.LayerNorm.weight requires_grad= True
deberta.encoder.layer.10.attention.output.LayerNorm.bias requires_grad= True
deberta.encoder.layer.10.intermediate.dense.weight requires_grad= True
deberta.encoder.layer.10.intermediate.dense.bias requires_grad= True
deberta.encoder.layer.10.output.dense.weight requires_grad= True
deberta.encoder.layer.10.output.dense.bias requires_grad= True
deberta.encoder.layer.10.output.LayerNorm.weight requires_grad= True
deberta.encoder.layer.10.output.LayerNorm.bias requires_grad= True
deberta.encoder.layer.11.attention.self.query_proj.weight requires_grad= True
deberta.encoder.layer.11.attention.self.query_proj.bias requires_grad= True
deberta.encoder.layer.11.attention.self.key_proj.weight requires_grad= True
deberta.encoder.layer.11.attention.self.key_proj.bias requires_grad= True
deberta.encoder.layer.11.attention.self.value_proj.weight requires_grad= True
deberta.encoder.layer.11.attention.self.value_proj.bias requires_grad= True
deberta.encoder.layer.11.attention.output.dense.weight requires_grad= True
deberta.encoder.layer.11.attention.output.dense.bias requires_grad= True
deberta.encoder.layer.11.attention.output.LayerNorm.weight requires_grad= True
deberta.encoder.layer.11.attention.output.LayerNorm.bias requires_grad= True
deberta.encoder.layer.11.intermediate.dense.weight requires_grad= True
deberta.encoder.layer.11.intermediate.dense.bias requires_grad= True
deberta.encoder.layer.11.output.dense.weight requires_grad= True
deberta.encoder.layer.11.output.dense.bias requires_grad= True
deberta.encoder.layer.11.output.LayerNorm.weight requires_grad= True
deberta.encoder.layer.11.output.LayerNorm.bias requires_grad= True
deberta.encoder.rel_embeddings.weight requires_grad= True
deberta.encoder.LayerNorm.weight requires_grad= True
deberta.encoder.LayerNorm.bias requires_grad= True
pooler.dense.weight requires_grad= True
pooler.dense.bias requires_grad= True
classifier.weight requires_grad= True
classifier.bias requires_grad= True

```

```

[64]: #####
layers_to_unfreeze = [
    "deberta.encoder.layer.11.attention.self.query_proj.weight",
    "deberta.encoder.layer.11.attention.self.query_proj.bias",
    "deberta.encoder.layer.11.attention.self.key_proj.weight",
    "deberta.encoder.layer.11.attention.self.key_proj.bias",
    "deberta.encoder.layer.11.attention.self.value_proj.weight",
    "deberta.encoder.layer.11.attention.self.value_proj.bias",
    "deberta.encoder.layer.11.attention.output.dense.weight",

```

```

"deberta.encoder.layer.11.attention.output.dense.bias",
"deberta.encoder.layer.11.attention.output.LayerNorm.weight",
"deberta.encoder.layer.11.attention.output.LayerNorm.bias",
"deberta.encoder.layer.11.intermediate.dense.weight",
"deberta.encoder.layer.11.intermediate.dense.bias",
"deberta.encoder.layer.11.output.dense.weight",
"deberta.encoder.layer.11.output.dense.bias",
"deberta.encoder.layer.11.output.LayerNorm.weight",
"deberta.encoder.layer.11.output.LayerNorm.bias",
"deberta.encoder.rel_embeddings.weight",
"deberta.encoder.LayerNorm.weight",
"deberta.encoder.LayerNorm.bias",
"pooler.dense.weight",
"pooler.dense.bias",
"classifier.weight",
"classifier.bias"
]
freeze_unfreeze_layers(model, layers_to_unfreeze=layers_to_unfreeze)
print(model.config)
print("=====")
print("num_parameters:", model.num_parameters())
print("num_trainable_parameters:", model.num_parameters(only_trainable=True))
print("=====")
#####
print("Experiment configuration used with this experiment:")
print("model used:", named_model)
print("learning rate used:", learning_rate)
print("number of epochs:", num_epochs)
print("maximum sequence length:", length_max)
print("batch size used:", size_batch)
print("regularization value:", regularization_weight_decay)
print("outcome variable:", y_col)
print("task:", x_task)
print("input column:", x_col)
#####
print("=====")
print("num_trainable_parameters:", model.num_parameters(only_trainable=True))

```

```

DebertaV2Config {
  "_attn_implementation_autoset": true,
  "attention_probs_dropout_prob": 0.1,
  "hidden_act": "gelu",
  "hidden_dropout_prob": 0.1,
  "hidden_size": 768,
  "initializer_range": 0.02,
  "intermediate_size": 3072,
  "layer_norm_eps": 1e-07,
  "legacy": true,

```

```

    "max_position_embeddings": 512,
    "max_relative_positions": -1,
    "model_type": "deberta-v2",
    "norm_rel_ebd": "layer_norm",
    "num_attention_heads": 12,
    "num_hidden_layers": 12,
    "pad_token_id": 0,
    "pooler_dropout": 0,
    "pooler_hidden_act": "gelu",
    "pooler_hidden_size": 768,
    "pos_att_type": [
        "p2c",
        "c2p"
    ],
    "position_biased_input": false,
    "position_buckets": 256,
    "relative_attention": true,
    "share_att_key": true,
    "torch_dtype": "float32",
    "transformers_version": "4.50.3",
    "type_vocab_size": 0,
    "vocab_size": 128100
}

```

=====

```

num_parameters: 184423682
num_trainable_parameters: 8074754

```

=====

Experiment configuration used with this experiment:

```

model used: microsoft/deberta-v3-base
learning rate used: 5e-06
number of epochs: 1
maximum sequence length: 128
batch size used: 128
regularization value: 0.5
outcome variable: binary_complexity_75th_split
task: single
input column: sentence_no_contractions

```

=====

```

num_trainable_parameters: 8074754

```

```

[65]: for name, param in model.named_parameters():
        print(name, "requires_grad=", param.requires_grad)

```

```

deberta.embeddings.word_embeddings.weight requires_grad= False
deberta.embeddings.LayerNorm.weight requires_grad= False
deberta.embeddings.LayerNorm.bias requires_grad= False
deberta.encoder.layer.0.attention.self.query_proj.weight requires_grad= False

```







```

deberta.encoder.layer.9.attention.self.query_proj.bias requires_grad= False
deberta.encoder.layer.9.attention.self.key_proj.weight requires_grad= False
deberta.encoder.layer.9.attention.self.key_proj.bias requires_grad= False
deberta.encoder.layer.9.attention.self.value_proj.weight requires_grad= False
deberta.encoder.layer.9.attention.self.value_proj.bias requires_grad= False
deberta.encoder.layer.9.attention.output.dense.weight requires_grad= False
deberta.encoder.layer.9.attention.output.dense.bias requires_grad= False
deberta.encoder.layer.9.attention.output.LayerNorm.weight requires_grad= False
deberta.encoder.layer.9.attention.output.LayerNorm.bias requires_grad= False
deberta.encoder.layer.9.intermediate.dense.weight requires_grad= False
deberta.encoder.layer.9.intermediate.dense.bias requires_grad= False
deberta.encoder.layer.9.output.dense.weight requires_grad= False
deberta.encoder.layer.9.output.dense.bias requires_grad= False
deberta.encoder.layer.9.output.LayerNorm.weight requires_grad= False
deberta.encoder.layer.9.output.LayerNorm.bias requires_grad= False
deberta.encoder.layer.10.attention.self.query_proj.weight requires_grad= False
deberta.encoder.layer.10.attention.self.query_proj.bias requires_grad= False
deberta.encoder.layer.10.attention.self.key_proj.weight requires_grad= False
deberta.encoder.layer.10.attention.self.key_proj.bias requires_grad= False
deberta.encoder.layer.10.attention.self.value_proj.weight requires_grad= False
deberta.encoder.layer.10.attention.self.value_proj.bias requires_grad= False
deberta.encoder.layer.10.attention.output.dense.weight requires_grad= False
deberta.encoder.layer.10.attention.output.dense.bias requires_grad= False
deberta.encoder.layer.10.attention.output.LayerNorm.weight requires_grad= False
deberta.encoder.layer.10.attention.output.LayerNorm.bias requires_grad= False
deberta.encoder.layer.10.intermediate.dense.weight requires_grad= False
deberta.encoder.layer.10.intermediate.dense.bias requires_grad= False
deberta.encoder.layer.10.output.dense.weight requires_grad= False
deberta.encoder.layer.10.output.dense.bias requires_grad= False
deberta.encoder.layer.10.output.LayerNorm.weight requires_grad= False
deberta.encoder.layer.10.output.LayerNorm.bias requires_grad= False
deberta.encoder.layer.11.attention.self.query_proj.weight requires_grad= True
deberta.encoder.layer.11.attention.self.query_proj.bias requires_grad= True
deberta.encoder.layer.11.attention.self.key_proj.weight requires_grad= True
deberta.encoder.layer.11.attention.self.key_proj.bias requires_grad= True
deberta.encoder.layer.11.attention.self.value_proj.weight requires_grad= True
deberta.encoder.layer.11.attention.self.value_proj.bias requires_grad= True
deberta.encoder.layer.11.attention.output.dense.weight requires_grad= True
deberta.encoder.layer.11.attention.output.dense.bias requires_grad= True
deberta.encoder.layer.11.attention.output.LayerNorm.weight requires_grad= True
deberta.encoder.layer.11.attention.output.LayerNorm.bias requires_grad= True
deberta.encoder.layer.11.intermediate.dense.weight requires_grad= True
deberta.encoder.layer.11.intermediate.dense.bias requires_grad= True
deberta.encoder.layer.11.output.dense.weight requires_grad= True
deberta.encoder.layer.11.output.dense.bias requires_grad= True
deberta.encoder.layer.11.output.LayerNorm.weight requires_grad= True
deberta.encoder.layer.11.output.LayerNorm.bias requires_grad= True
deberta.encoder.rel_embeddings.weight requires_grad= True

```

```

deberta.encoder.LayerNorm.weight requires_grad= True
deberta.encoder.LayerNorm.bias requires_grad= True
pooler.dense.weight requires_grad= True
pooler.dense.bias requires_grad= True
classifier.weight requires_grad= True
classifier.bias requires_grad= True

```

```
[66]: # model.resize_token_embeddings(len(tokenizer))
```

```
[67]: # Train & Evaluate
trained_model, trainer_obj = train_transformer_model(
    model = model,
    tokenizer = tokenizer,
    train_dataset = train_data_hf,
    val_dataset = val_data_hf,
    output_dir = dir_results,
    num_epochs = num_epochs,
    batch_size = size_batch,
    lr = learning_rate,
    weight_decay = regularization_weight_decay)
metrics = trainer_obj.evaluate()
print("Validation metrics:", metrics)
test_metrics = trainer_obj.evaluate(test_data_hf) if test_data_hf else None
print("Test metrics:", test_metrics)
```

```

/usr/local/lib/python3.11/dist-packages/transformers/training_args.py:1611:
FutureWarning: `evaluation_strategy` is deprecated and will be removed in
version 4.46 of Transformers. Use `eval_strategy` instead
  warnings.warn(
<ipython-input-20-c2ee9f934517>:31: FutureWarning: `tokenizer` is deprecated and
will be removed in version 5.0.0 for `Trainer.__init__`. Use `processing_class`
instead.
  trainer = Trainer(
<IPython.core.display.HTML object>
model.safetensors:   0%|          | 0.00/371M [00:00<?, ?B/s]
<IPython.core.display.HTML object>
Validation metrics: {'eval_loss': 0.5420940518379211, 'eval_accuracy':
0.7767220902612827, 'eval_precision': 0.0, 'eval_recall': 0.0, 'eval_f1': 0.0,
'eval_runtime': 6.0467, 'eval_samples_per_second': 69.625,
'eval_steps_per_second': 0.662, 'epoch': 1.0}
Test metrics: {'eval_loss': 0.5194469690322876, 'eval_accuracy':
0.7960741548527808, 'eval_precision': 0.0, 'eval_recall': 0.0, 'eval_f1': 0.0,
'eval_runtime': 6.9785, 'eval_samples_per_second': 131.403,
'eval_steps_per_second': 1.146, 'epoch': 1.0}
/usr/local/lib/python3.11/dist-packages/sklearn/metrics/_classification.py:1565:
UndefinedMetricWarning: Precision is ill-defined and being set to 0.0 due to no

```



predicted samples. Use `zero\_division` parameter to control this behavior.  
\_warn\_prf(average, modifier, f"{metric.capitalize()} is", len(result))

```
[68]: # save model checkpoint
timestamp = datetime.now().strftime("%Y%m%d_%H%M%S")
model_save_path = os.path.join(dir_models,
    ↪f"{x_task}_{named_model}_{y_col}_{timestamp}")
trainer_obj.save_model(model_save_path)
print(f"Model checkpoint saved to: {model_save_path}")
# log experiment results
experiment_info = {
    "model_name": named_model,
    "learning_rate": learning_rate,
    "epochs": num_epochs,
    "batch_size": size_batch,
    "weight_decay": regularization_weight_decay,
    "x_task": x_task,
    "x_col": x_col,
    "y_col": y_col,
    "layers_to_unfreeze": layers_to_unfreeze}
model_info = gather_model_details(trained_model)
all_run_metrics = gather_all_run_metrics(
    trainer=trainer_obj,
    train_dataset=train_data_hf,
    val_dataset=val_data_hf,
    test_dataset=test_data_hf)
log_experiment_results_json(
    experiment_meta=experiment_info,
    model_details=model_info,
    run_metrics=all_run_metrics,
    log_file=log_filepath)
print(f"EXPERIMENT LOGGED TO: {log_filepath}")
```

Model checkpoint saved to: /content/drive/MyDrive/266-final/models/single\_microsoft/deberta-v3-base\_binary\_complexity\_75th\_split\_20250411\_011134

<IPython.core.display.HTML object>

EXPERIMENT LOGGED TO:

/content/drive/MyDrive/266-final/results/experiment\_runs.txt

/usr/local/lib/python3.11/dist-packages/sklearn/metrics/\_classification.py:1565:

UndefinedMetricWarning: Precision is ill-defined and being set to 0.0 due to no predicted samples. Use `zero\_division` parameter to control this behavior.

\_warn\_prf(average, modifier, f"{metric.capitalize()} is", len(result))

0.2.10 snc xlnet/xlnet-base-cased regularization\_weight\_decay = 0.5 learning\_rate = 5e-6 size\_batch = 128 length\_max = 128 num\_epochs = 1

```
[69]: # Define Experiment Parameters
# named_model = "bert-base-cased"
# named_model = "roberta-base"
# named_model = "bert-large-cased"
# named_model = "roberta-large"
# named_model = "answerdotai/ModernBERT-base" # modern bert
# named_model = "answerdotai/ModernBERT-large" # modern bert
# named_model = "microsoft/deberta-v3-base" # deberta
named_model = "xlnet/xlnet-base-cased" #
# named_model = "xlnet/xlnet-large-cased" #
#####
regularization_weight_decay = 0.5
learning_rate = 5e-6
size_batch = 128
length_max = 128
num_epochs = 1
#####
# x_col = "sentence"
x_col = "sentence_no_contractions"
# x_col = "pos_sequence"
# x_col = "dep_sequence"
# x_col = "morph_sequence"
#####
y_col = "binary_complexity_75th_split"
# y_col = "binary_complexity"
# y_col = "complexity"
#####
x_task = "single"
# x_task = "multi"
if x_task == "single":
    df_train = train_single_df
    df_val = trial_val_single_df
    df_test = test_single_df
else:
    df_train = train_multi_df
    df_val = trial_val_multi_df
    df_test = test_multi_df
#####
# Tokenize & Prepare Datasets
train_data_hf = prepare_dataset(
    df_train,
    tokenizer,
    text_col=x_col,
    label_col=y_col,
```

```

        max_length=length_max)
val_data_hf = prepare_dataset(
    df_val,
    tokenizer,
    text_col=x_col,
    label_col=y_col,
    max_length=length_max)
test_data_hf = prepare_dataset(
    df_test,
    tokenizer,
    text_col=x_col,
    label_col=y_col,
    max_length=length_max)
print("Datasets prepared. Sample from train_data_hf:\n", train_data_hf[10])
# print("Datasets prepared. Sample from train_data_hf:\n", val_data_hf[10])
# print("Datasets prepared. Sample from train_data_hf:\n", test_data_hf[10])
#####
model, tokenizer = get_model_and_tokenizer(
    remote_model_name="xlnet/xlnet-base-cased",
    local_model_path=None,
    config=None)
#####
# model, tokenizer = get_model_and_tokenizer(
#     remote_model_name=None
#     local_model_path="...CONFIGURE_PATH...",
#     config=custom_config)
#####
print("=====")
print(named_model, ":")
print("=====")
print("num_parameters:", model.num_parameters())
print("num_trainable_parameters at load:", model.
    ↪ num_parameters(only_trainable=True))
print("=====")
print("model lineage:", MODEL_LINEAGE)
print("=====")

```

Map: 0%| | 0/7662 [00:00<?, ? examples/s]

Map: 0%| | 0/421 [00:00<?, ? examples/s]

Map: 0%| | 0/917 [00:00<?, ? examples/s]

Datasets prepared. Sample from train\_data\_hf:

```

{'labels': tensor(0), 'input_ids': tensor([ 1, 420, 264, 262, 6579,
265, 95435, 474, 313, 1255,
3083, 261, 401, 262, 506, 265, 262, 13083, 13713, 264,
349, 346, 306, 2635, 278, 277, 308, 6868, 260, 2,
0, 0, 0, 0, 0, 0, 0, 0, 0, 0,

```

```

0, 0, 0, 0, 0, 0, 0, 0, 0, 0,
0, 0, 0, 0, 0, 0, 0, 0, 0, 0,
0, 0, 0, 0, 0, 0, 0, 0, 0, 0,
0, 0, 0, 0, 0, 0, 0, 0, 0, 0,
0, 0, 0, 0, 0, 0, 0, 0, 0, 0,
0, 0, 0, 0, 0, 0, 0, 0, 0, 0,
0, 0, 0, 0, 0, 0, 0, 0, 0, 0,
0, 0, 0, 0, 0, 0, 0, 0, 0, 0,
0, 0, 0, 0, 0, 0, 0, 0]),
'attention_mask': tensor([1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1,
1, 1, 1, 1, 1,
1, 1, 1, 1, 1, 1, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0,
0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0,
0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0,
0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0,
0, 0, 0, 0, 0, 0, 0, 0]),
Loading from Hugging Face model: xlnet/xlnet-base-cased
config.json: 0%|          | 0.00/760 [00:00<?, ?B/s]
spiece.model: 0%|          | 0.00/798k [00:00<?, ?B/s]
tokenizer.json: 0%|          | 0.00/1.38M [00:00<?, ?B/s]
pytorch_model.bin: 0%|          | 0.00/467M [00:00<?, ?B/s]

Some weights of XLNetForSequenceClassification were not initialized from the
model checkpoint at xlnet/xlnet-base-cased and are newly initialized:
['logits_proj.bias', 'logits_proj.weight', 'sequence_summary.summary.bias',
'sequence_summary.summary.weight']
You should probably TRAIN this model on a down-stream task to be able to use it
for predictions and inference.

```

```

=====
xlnet/xlnet-base-cased :
=====
num_parameters: 117310466
num_trainable_parameters at load: 117310466
=====
model lineage: {'type': 'huggingface_hub', 'path': 'xlnet/xlnet-base-cased',
'timestamp': '2025-04-11 01:12:18'}
=====

```

```
[70]: print(model)
```

```

XLNetForSequenceClassification(
  (transformer): XLNetModel(
    (word_embedding): Embedding(32000, 768)
    (layer): ModuleList(
      (0-11): 12 x XLNetLayer(
        (rel_attn): XLNetRelativeAttention(

```

```

        (layer_norm): LayerNorm((768,), eps=1e-12, elementwise_affine=True)
        (dropout): Dropout(p=0.1, inplace=False)
    )
    (ff): XLNetFeedForward(
        (layer_norm): LayerNorm((768,), eps=1e-12, elementwise_affine=True)
        (layer_1): Linear(in_features=768, out_features=3072, bias=True)
        (layer_2): Linear(in_features=3072, out_features=768, bias=True)
        (dropout): Dropout(p=0.1, inplace=False)
        (activation_function): GELUActivation()
    )
    (dropout): Dropout(p=0.1, inplace=False)
)
)
(dropout): Dropout(p=0.1, inplace=False)
)
(sequence_summary): SequenceSummary(
    (summary): Linear(in_features=768, out_features=768, bias=True)
    (activation): Tanh()
    (first_dropout): Identity()
    (last_dropout): Dropout(p=0.1, inplace=False)
)
(logits_proj): Linear(in_features=768, out_features=2, bias=True)
)

```

```

[71]: for name, param in model.named_parameters():
        print(name, "requires_grad=", param.requires_grad)

```

```

transformer.mask_emb requires_grad= True
transformer.word_embedding.weight requires_grad= True
transformer.layer.0.rel_attn.q requires_grad= True
transformer.layer.0.rel_attn.k requires_grad= True
transformer.layer.0.rel_attn.v requires_grad= True
transformer.layer.0.rel_attn.o requires_grad= True
transformer.layer.0.rel_attn.r requires_grad= True
transformer.layer.0.rel_attn.r_r_bias requires_grad= True
transformer.layer.0.rel_attn.r_s_bias requires_grad= True
transformer.layer.0.rel_attn.r_w_bias requires_grad= True
transformer.layer.0.rel_attn.seg_embed requires_grad= True
transformer.layer.0.rel_attn.layer_norm.weight requires_grad= True
transformer.layer.0.rel_attn.layer_norm.bias requires_grad= True
transformer.layer.0.ff.layer_norm.weight requires_grad= True
transformer.layer.0.ff.layer_norm.bias requires_grad= True
transformer.layer.0.ff.layer_1.weight requires_grad= True
transformer.layer.0.ff.layer_1.bias requires_grad= True
transformer.layer.0.ff.layer_2.weight requires_grad= True
transformer.layer.0.ff.layer_2.bias requires_grad= True
transformer.layer.1.rel_attn.q requires_grad= True
transformer.layer.1.rel_attn.k requires_grad= True

```









```

transformer.layer.9.rel_attn.layer_norm.bias requires_grad= True
transformer.layer.9.ff.layer_norm.weight requires_grad= True
transformer.layer.9.ff.layer_norm.bias requires_grad= True
transformer.layer.9.ff.layer_1.weight requires_grad= True
transformer.layer.9.ff.layer_1.bias requires_grad= True
transformer.layer.9.ff.layer_2.weight requires_grad= True
transformer.layer.9.ff.layer_2.bias requires_grad= True
transformer.layer.10.rel_attn.q requires_grad= True
transformer.layer.10.rel_attn.k requires_grad= True
transformer.layer.10.rel_attn.v requires_grad= True
transformer.layer.10.rel_attn.o requires_grad= True
transformer.layer.10.rel_attn.r requires_grad= True
transformer.layer.10.rel_attn.r_r_bias requires_grad= True
transformer.layer.10.rel_attn.r_s_bias requires_grad= True
transformer.layer.10.rel_attn.r_w_bias requires_grad= True
transformer.layer.10.rel_attn.seg_embed requires_grad= True
transformer.layer.10.rel_attn.layer_norm.weight requires_grad= True
transformer.layer.10.rel_attn.layer_norm.bias requires_grad= True
transformer.layer.10.ff.layer_norm.weight requires_grad= True
transformer.layer.10.ff.layer_norm.bias requires_grad= True
transformer.layer.10.ff.layer_1.weight requires_grad= True
transformer.layer.10.ff.layer_1.bias requires_grad= True
transformer.layer.10.ff.layer_2.weight requires_grad= True
transformer.layer.10.ff.layer_2.bias requires_grad= True
transformer.layer.11.rel_attn.q requires_grad= True
transformer.layer.11.rel_attn.k requires_grad= True
transformer.layer.11.rel_attn.v requires_grad= True
transformer.layer.11.rel_attn.o requires_grad= True
transformer.layer.11.rel_attn.r requires_grad= True
transformer.layer.11.rel_attn.r_r_bias requires_grad= True
transformer.layer.11.rel_attn.r_s_bias requires_grad= True
transformer.layer.11.rel_attn.r_w_bias requires_grad= True
transformer.layer.11.rel_attn.seg_embed requires_grad= True
transformer.layer.11.rel_attn.layer_norm.weight requires_grad= True
transformer.layer.11.rel_attn.layer_norm.bias requires_grad= True
transformer.layer.11.ff.layer_norm.weight requires_grad= True
transformer.layer.11.ff.layer_norm.bias requires_grad= True
transformer.layer.11.ff.layer_1.weight requires_grad= True
transformer.layer.11.ff.layer_1.bias requires_grad= True
transformer.layer.11.ff.layer_2.weight requires_grad= True
transformer.layer.11.ff.layer_2.bias requires_grad= True
sequence_summary.summary.weight requires_grad= True
sequence_summary.summary.bias requires_grad= True
logits_proj.weight requires_grad= True
logits_proj.bias requires_grad= True

```

```
[72]: #####
layers_to_unfreeze = [
    "transformer.layer.11.rel_attn.q",
    "transformer.layer.11.rel_attn.k",
    "transformer.layer.11.rel_attn.v",
    "transformer.layer.11.rel_attn.o",
    "transformer.layer.11.rel_attn.r",
    "transformer.layer.11.rel_attn.r_r_bias",
    "transformer.layer.11.rel_attn.r_s_bias",
    "transformer.layer.11.rel_attn.r_w_bias",
    "transformer.layer.11.rel_attn.seg_embed",
    "transformer.layer.11.rel_attn.layer_norm.weight",
    "transformer.layer.11.rel_attn.layer_norm.bias",
    "transformer.layer.11.ff.layer_norm.weight",
    "transformer.layer.11.ff.layer_norm.bias",
    "transformer.layer.11.ff.layer_1.weight",
    "transformer.layer.11.ff.layer_1.bias",
    "transformer.layer.11.ff.layer_2.weight",
    "transformer.layer.11.ff.layer_2.bias",
    "sequence_summary.summary.weight",
    "sequence_summary.summary.bias",
    "logits_proj.weight",
    "logits_proj.bias"
]
freeze_unfreeze_layers(model, layers_to_unfreeze=layers_to_unfreeze)
print(model.config)
print("=====")
print("num_parameters:", model.num_parameters())
print("num_trainable_parameters:", model.num_parameters(only_trainable=True))
print("=====")
#####
print("Experiment configuration used with this experiment:")
print("model used:", named_model)
print("learning rate used:", learning_rate)
print("number of epochs:", num_epochs)
print("maximum sequence length:", length_max)
print("batch size used:", size_batch)
print("regularization value:", regularization_weight_decay)
print("outcome variable:", y_col)
print("task:", x_task)
print("input column:", x_col)
#####
print("=====")
print("num_trainable_parameters:", model.num_parameters(only_trainable=True))
```

```
XLNetConfig {
  "_attn_implementation_autoset": true,
  "architectures": [
```

```

    "XLNetLMHeadModel"
],
"attn_type": "bi",
"bi_data": false,
"bos_token_id": 1,
"clamp_len": -1,
"d_head": 64,
"d_inner": 3072,
"d_model": 768,
"dropout": 0.1,
"end_n_top": 5,
"eos_token_id": 2,
"ff_activation": "gelu",
"initializer_range": 0.02,
"layer_norm_eps": 1e-12,
"mem_len": null,
"model_type": "xlnet",
"n_head": 12,
"n_layer": 12,
"pad_token_id": 5,
"reuse_len": null,
"same_length": false,
"start_n_top": 5,
"summary_activation": "tanh",
"summary_last_dropout": 0.1,
"summary_type": "last",
"summary_use_proj": true,
"task_specific_params": {
    "text-generation": {
        "do_sample": true,
        "max_length": 250
    }
},
"torch_dtype": "float32",
"transformers_version": "4.50.3",
"untie_r": true,
"use_mems_eval": true,
"use_mems_train": false,
"vocab_size": 32000
}

```

=====

```

num_parameters: 117310466
num_trainable_parameters: 8270594
=====

```

Experiment configuration used with this experiment:  
model used: xlnet/xlnet-base-cased  
learning rate used: 5e-06

```

number of epochs: 1
maximum sequence length: 128
batch size used: 128
regularization value: 0.5
outcome variable: binary_complexity_75th_split
task: single
input column: sentence_no_contractions
=====
num_trainable_parameters: 8270594

```

```

[73]: for name, param in model.named_parameters():
      print(name, "requires_grad=", param.requires_grad)

```

```

transformer.mask_emb requires_grad= False
transformer.word_embedding.weight requires_grad= False
transformer.layer.0.rel_attn.q requires_grad= False
transformer.layer.0.rel_attn.k requires_grad= False
transformer.layer.0.rel_attn.v requires_grad= False
transformer.layer.0.rel_attn.o requires_grad= False
transformer.layer.0.rel_attn.r requires_grad= False
transformer.layer.0.rel_attn.r_r_bias requires_grad= False
transformer.layer.0.rel_attn.r_s_bias requires_grad= False
transformer.layer.0.rel_attn.r_w_bias requires_grad= False
transformer.layer.0.rel_attn.seg_embed requires_grad= False
transformer.layer.0.rel_attn.layer_norm.weight requires_grad= False
transformer.layer.0.rel_attn.layer_norm.bias requires_grad= False
transformer.layer.0.ff.layer_norm.weight requires_grad= False
transformer.layer.0.ff.layer_norm.bias requires_grad= False
transformer.layer.0.ff.layer_1.weight requires_grad= False
transformer.layer.0.ff.layer_1.bias requires_grad= False
transformer.layer.0.ff.layer_2.weight requires_grad= False
transformer.layer.0.ff.layer_2.bias requires_grad= False
transformer.layer.1.rel_attn.q requires_grad= False
transformer.layer.1.rel_attn.k requires_grad= False
transformer.layer.1.rel_attn.v requires_grad= False
transformer.layer.1.rel_attn.o requires_grad= False
transformer.layer.1.rel_attn.r requires_grad= False
transformer.layer.1.rel_attn.r_r_bias requires_grad= False
transformer.layer.1.rel_attn.r_s_bias requires_grad= False
transformer.layer.1.rel_attn.r_w_bias requires_grad= False
transformer.layer.1.rel_attn.seg_embed requires_grad= False
transformer.layer.1.rel_attn.layer_norm.weight requires_grad= False
transformer.layer.1.rel_attn.layer_norm.bias requires_grad= False
transformer.layer.1.ff.layer_norm.weight requires_grad= False
transformer.layer.1.ff.layer_norm.bias requires_grad= False
transformer.layer.1.ff.layer_1.weight requires_grad= False
transformer.layer.1.ff.layer_1.bias requires_grad= False
transformer.layer.1.ff.layer_2.weight requires_grad= False

```

[illegible]

[illegible]

transformer.layer.7.rel\_attn.layer\_norm.bias requires\_grad= False  
transformer.layer.7.ff.layer\_norm.weight requires\_grad= False  
transformer.layer.7.ff.layer\_norm.bias requires\_grad= False  
transformer.layer.7.ff.layer\_1.weight requires\_grad= False  
transformer.layer.7.ff.layer\_1.bias requires\_grad= False  
transformer.layer.7.ff.layer\_2.weight requires\_grad= False  
transformer.layer.7.ff.layer\_2.bias requires\_grad= False  
transformer.layer.8.rel\_attn.q requires\_grad= False  
transformer.layer.8.rel\_attn.k requires\_grad= False  
transformer.layer.8.rel\_attn.v requires\_grad= False  
transformer.layer.8.rel\_attn.o requires\_grad= False  
transformer.layer.8.rel\_attn.r requires\_grad= False  
transformer.layer.8.rel\_attn.r\_r\_bias requires\_grad= False  
transformer.layer.8.rel\_attn.r\_s\_bias requires\_grad= False  
transformer.layer.8.rel\_attn.r\_w\_bias requires\_grad= False  
transformer.layer.8.rel\_attn.seg\_embed requires\_grad= False  
transformer.layer.8.rel\_attn.layer\_norm.weight requires\_grad= False  
transformer.layer.8.rel\_attn.layer\_norm.bias requires\_grad= False  
transformer.layer.8.ff.layer\_norm.weight requires\_grad= False  
transformer.layer.8.ff.layer\_norm.bias requires\_grad= False  
transformer.layer.8.ff.layer\_1.weight requires\_grad= False  
transformer.layer.8.ff.layer\_1.bias requires\_grad= False  
transformer.layer.8.ff.layer\_2.weight requires\_grad= False  
transformer.layer.8.ff.layer\_2.bias requires\_grad= False  
transformer.layer.9.rel\_attn.q requires\_grad= False  
transformer.layer.9.rel\_attn.k requires\_grad= False  
transformer.layer.9.rel\_attn.v requires\_grad= False  
transformer.layer.9.rel\_attn.o requires\_grad= False  
transformer.layer.9.rel\_attn.r requires\_grad= False  
transformer.layer.9.rel\_attn.r\_r\_bias requires\_grad= False  
transformer.layer.9.rel\_attn.r\_s\_bias requires\_grad= False  
transformer.layer.9.rel\_attn.r\_w\_bias requires\_grad= False  
transformer.layer.9.rel\_attn.seg\_embed requires\_grad= False  
transformer.layer.9.rel\_attn.layer\_norm.weight requires\_grad= False  
transformer.layer.9.rel\_attn.layer\_norm.bias requires\_grad= False  
transformer.layer.9.ff.layer\_norm.weight requires\_grad= False  
transformer.layer.9.ff.layer\_norm.bias requires\_grad= False  
transformer.layer.9.ff.layer\_1.weight requires\_grad= False  
transformer.layer.9.ff.layer\_1.bias requires\_grad= False  
transformer.layer.9.ff.layer\_2.weight requires\_grad= False  
transformer.layer.9.ff.layer\_2.bias requires\_grad= False  
transformer.layer.10.rel\_attn.q requires\_grad= False  
transformer.layer.10.rel\_attn.k requires\_grad= False  
transformer.layer.10.rel\_attn.v requires\_grad= False  
transformer.layer.10.rel\_attn.o requires\_grad= False  
transformer.layer.10.rel\_attn.r requires\_grad= False  
transformer.layer.10.rel\_attn.r\_r\_bias requires\_grad= False  
transformer.layer.10.rel\_attn.r\_s\_bias requires\_grad= False

```

transformer.layer.10.rel_attn.r_w_bias requires_grad= False
transformer.layer.10.rel_attn.seg_embed requires_grad= False
transformer.layer.10.rel_attn.layer_norm.weight requires_grad= False
transformer.layer.10.rel_attn.layer_norm.bias requires_grad= False
transformer.layer.10.ff.layer_norm.weight requires_grad= False
transformer.layer.10.ff.layer_norm.bias requires_grad= False
transformer.layer.10.ff.layer_1.weight requires_grad= False
transformer.layer.10.ff.layer_1.bias requires_grad= False
transformer.layer.10.ff.layer_2.weight requires_grad= False
transformer.layer.10.ff.layer_2.bias requires_grad= False
transformer.layer.11.rel_attn.q requires_grad= True
transformer.layer.11.rel_attn.k requires_grad= True
transformer.layer.11.rel_attn.v requires_grad= True
transformer.layer.11.rel_attn.o requires_grad= True
transformer.layer.11.rel_attn.r requires_grad= True
transformer.layer.11.rel_attn.r_r_bias requires_grad= True
transformer.layer.11.rel_attn.r_s_bias requires_grad= True
transformer.layer.11.rel_attn.r_w_bias requires_grad= True
transformer.layer.11.rel_attn.seg_embed requires_grad= True
transformer.layer.11.rel_attn.layer_norm.weight requires_grad= True
transformer.layer.11.rel_attn.layer_norm.bias requires_grad= True
transformer.layer.11.ff.layer_norm.weight requires_grad= True
transformer.layer.11.ff.layer_norm.bias requires_grad= True
transformer.layer.11.ff.layer_1.weight requires_grad= True
transformer.layer.11.ff.layer_1.bias requires_grad= True
transformer.layer.11.ff.layer_2.weight requires_grad= True
transformer.layer.11.ff.layer_2.bias requires_grad= True
sequence_summary.summary.weight requires_grad= True
sequence_summary.summary.bias requires_grad= True
logits_proj.weight requires_grad= True
logits_proj.bias requires_grad= True

```

```
[74]: # model.resize_token_embeddings(len(tokenizer))
```

```

[75]: # Train & Evaluate
trained_model, trainer_obj = train_transformer_model(
    model = model,
    tokenizer = tokenizer,
    train_dataset = train_data_hf,
    val_dataset = val_data_hf,
    output_dir = dir_results,
    num_epochs = num_epochs,
    batch_size = size_batch,
    lr = learning_rate,
    weight_decay = regularization_weight_decay)
metrics = trainer_obj.evaluate()
print("Validation metrics:", metrics)

```



```
test_metrics = trainer_obj.evaluate(test_data_hf) if test_data_hf else None
print("Test metrics:", test_metrics)
```

/usr/local/lib/python3.11/dist-packages/transformers/training\_args.py:1611:

FutureWarning: `evaluation\_strategy` is deprecated and will be removed in version 4.46 of Transformers. Use `eval\_strategy` instead

```
warnings.warn(
```

<ipython-input-20-c2ee9f934517>:31: FutureWarning: `tokenizer` is deprecated and will be removed in version 5.0.0 for `Trainer.\_\_init\_\_`. Use `processing\_class` instead.

```
trainer = Trainer(
```

```
-----
RuntimeError                                Traceback (most recent call last)
<ipython-input-75-cf30694e0124> in <cell line: 0>()
      1 # Train & Evaluate
----> 2 trained_model, trainer_obj = train_transformer_model(
      3     model = model,
      4     tokenizer = tokenizer,
      5     train_dataset = train_data_hf,

<ipython-input-20-c2ee9f934517> in train_transformer_model(model, tokenizer,
↳ train_dataset, val_dataset, output_dir, num_epochs, batch_size, lr,
↳ weight_decay)
      38     )
      39
---> 40     trainer.train()
      41     return model, trainer

/usr/local/lib/python3.11/dist-packages/transformers/trainer.py in train(self,
↳ resume_from_checkpoint, trial, ignore_keys_for_eval, **kwargs)
    2243         hf_hub_utils.enable_progressBars()
    2244     else:
-> 2245         return inner_training_loop(
    2246             args=args,
    2247             resume_from_checkpoint=resume_from_checkpoint,

/usr/local/lib/python3.11/dist-packages/transformers/trainer.py in
↳ _inner_training_loop(self, batch_size, args, resume_from_checkpoint, trial,
↳ ignore_keys_for_eval)
    2554         )
    2555         with context():
-> 2556             tr_loss_step = self.training_step(model, inputs
↳ num_items_in_batch)
    2557
    2558         if (
```

```

/usr/local/lib/python3.11/dist-packages/transformers/trainer.py in
↳ training_step(self, model, inputs, num_items_in_batch)
    3716
    3717         with self.compute_loss_context_manager():
-> 3718             loss = self.compute_loss(model, inputs,
↳ num_items_in_batch=num_items_in_batch)
    3719
    3720         del inputs

/usr/local/lib/python3.11/dist-packages/transformers/trainer.py in
↳ compute_loss(self, model, inputs, return_outputs, num_items_in_batch)
    3781             loss_kwargs["num_items_in_batch"] = num_items_in_batch
    3782             inputs = {**inputs, **loss_kwargs}
-> 3783             outputs = model(**inputs)
    3784             # Save past state if it exists
    3785             # TODO: this needs to be fixed and made cleaner later.

/usr/local/lib/python3.11/dist-packages/torch/nn/modules/module.py in
↳ _wrapped_call_impl(self, *args, **kwargs)
    1737         return self._compiled_call_impl(*args, **kwargs) # type:
↳ ignore[misc]
    1738         else:
-> 1739         return self._call_impl(*args, **kwargs)
    1740
    1741         # torchrec tests the code consistency with the following code

/usr/local/lib/python3.11/dist-packages/torch/nn/modules/module.py in
↳ _call_impl(self, *args, **kwargs)
    1748             or _global_backward_pre_hooks or _global_backward_hooks
    1749             or _global_forward_hooks or _global_forward_pre_hooks):
-> 1750         return forward_call(*args, **kwargs)
    1751
    1752         result = None

/usr/local/lib/python3.11/dist-packages/transformers/models/xlnet/modeling_xlnet.py in
↳ forward(self, input_ids, attention_mask, mems, perm_mask,
↳ target_mapping, token_type_ids, input_mask, head_mask, inputs_embeds, labels,
↳ use_mems, output_attentions, output_hidden_states, return_dict, **kwargs)
    1541         return_dict = return_dict if return_dict is not None else self.
↳ config.use_return_dict
    1542
-> 1543         transformer_outputs = self.transformer(
    1544             input_ids,
    1545             attention_mask=attention_mask,

/usr/local/lib/python3.11/dist-packages/torch/nn/modules/module.py in
↳ _wrapped_call_impl(self, *args, **kwargs)

```

```

1737         return self._compiled_call_impl(*args, **kwargs) # type:␣
↪ ignore[misc]
1738     else:
-> 1739         return self._call_impl(*args, **kwargs)
1740
1741     # torchrec tests the code consistency with the following code

/usr/local/lib/python3.11/dist-packages/torch/nn/modules/module.py in␣
↪ _call_impl(self, *args, **kwargs)
1748         or _global_backward_pre_hooks or _global_backward_hooks
1749         or _global_forward_hooks or _global_forward_pre_hooks):
-> 1750         return forward_call(*args, **kwargs)
1751
1752     result = None

/usr/local/lib/python3.11/dist-packages/transformers/models/xlnet/modeling_xlnet␣
↪ py in forward(self, input_ids, attention_mask, mems, perm_mask,␣
↪ target_mapping, token_type_ids, input_mask, head_mask, inputs_embeds,␣
↪ use_mems, output_attentions, output_hidden_states, return_dict, **kwargs)
1198     # Positional encoding
1199     pos_emb = self.relative_positional_encoding(qlen, klen, bsz=bsz
-> 1200     pos_emb = pos_emb.to(output_h.device)
1201     pos_emb = self.dropout(pos_emb)
1202

```

**RuntimeError:** CUDA error: device-side assert triggered  
 CUDA kernel errors might be asynchronously reported at some other API call, so␣  
 ↪ the stacktrace below might be incorrect.  
 For debugging consider passing CUDA\_LAUNCH\_BLOCKING=1  
 Compile with `TORCH\_USE\_CUDA\_DSA` to enable device-side assertions.

```

[ ]: # save model checkpoint
timestamp = datetime.now().strftime("%Y%m%d_%H%M%S")
model_save_path = os.path.join(dir_models,␣
↪ f"{x_task}_{named_model}_{y_col}_{timestamp}")
trainer_obj.save_model(model_save_path)
print(f"Model checkpoint saved to: {model_save_path}")
# log experiment results
experiment_info = {
    "model_name": named_model,
    "learning_rate": learning_rate,
    "epochs": num_epochs,
    "batch_size": size_batch,
    "weight_decay": regularization_weight_decay,
    "x_task": x_task,
    "x_col": x_col,

```

```

        "y_col": y_col,
        "layers_to_unfreeze": layers_to_unfreeze}
model_info = gather_model_details(trained_model)
all_run_metrics = gather_all_run_metrics(
    trainer=trainer_obj,
    train_dataset=train_data_hf,
    val_dataset=val_data_hf,
    test_dataset=test_data_hf)
log_experiment_results_json(
    experiment_meta=experiment_info,
    model_details=model_info,
    run_metrics=all_run_metrics,
    log_file=log_filepath)
print(f"EXPERIMENT LOGGED TO: {log_filepath}")

```

**0.2.11** snc xlnet/xlnet-large-cased regularization\_weight\_decay = 0.5 learning\_rate = 5e-6 size\_batch = 128 length\_max = 128 num\_epochs = 1

```

[ ]: # Define Experiment Parameters
# named_model = "bert-base-cased"
# named_model = "roberta-base"
# named_model = "bert-large-cased"
# named_model = "roberta-large"
# named_model = "answerdotai/ModernBERT-base" # modern bert
# named_model = "answerdotai/ModernBERT-large" # modern bert
# named_model = "microsoft/deberta-v3-base" # deberta
# named_model = "xlnet/xlnet-base-cased" #
named_model = "xlnet/xlnet-large-cased" #
#####
regularization_weight_decay = 0.5
learning_rate = 5e-6
size_batch = 128
length_max = 128
num_epochs = 1
#####
# x_col = "sentence"
x_col = "sentence_no_contractions"
# x_col = "pos_sequence"
# x_col = "dep_sequence"
# x_col = "morph_sequence"
#####
y_col = "binary_complexity_75th_split"
# y_col = "binary_complexity"
# y_col = "complexity"
#####
x_task = "single"
# x_task = "multi"

```

```

if x_task == "single":
    df_train = train_single_df
    df_val = trial_val_single_df
    df_test = test_single_df
else:
    df_train = train_multi_df
    df_val = trial_val_multi_df
    df_test = test_multi_df
#####
# Tokenize & Prepare Datasets
train_data_hf = prepare_dataset(
    df_train,
    tokenizer,
    text_col=x_col,
    label_col=y_col,
    max_length=length_max)
val_data_hf = prepare_dataset(
    df_val,
    tokenizer,
    text_col=x_col,
    label_col=y_col,
    max_length=length_max)
test_data_hf = prepare_dataset(
    df_test,
    tokenizer,
    text_col=x_col,
    label_col=y_col,
    max_length=length_max)
print("Datasets prepared. Sample from train_data_hf:\n", train_data_hf[10])
# print("Datasets prepared. Sample from train_data_hf:\n", val_data_hf[10])
# print("Datasets prepared. Sample from train_data_hf:\n", test_data_hf[10])
#####
model, tokenizer = get_model_and_tokenizer(
    remote_model_name="xlnet/xlnet-large-cased",
    local_model_path=None,
    config=None)
#####
# model, tokenizer = get_model_and_tokenizer(
#     remote_model_name=None
#     local_model_path="...CONFIGURE_PATH...",
#     config=custom_config)
#####
print("=====")
print(named_model, ":")
print("=====")
print("num_parameters:", model.num_parameters())

```

```

print("num_trainable_parameters at load:", model.
    ↪ num_parameters(only_trainable=True))
print("=====")
print("model lineage:", MODEL_LINEAGE)
print("=====")

```

```
[ ]: print(model)
```

```
[ ]: for name, param in model.named_parameters():
      print(name, "requires_grad=", param.requires_grad)
```

```
[ ]: #####
layers_to_unfreeze = [
    "transformer.layer.23.rel_attn.q",
    "transformer.layer.23.rel_attn.k",
    "transformer.layer.23.rel_attn.v",
    "transformer.layer.23.rel_attn.o",
    "transformer.layer.23.rel_attn.r",
    "transformer.layer.23.rel_attn.r_r_bias",
    "transformer.layer.23.rel_attn.r_s_bias",
    "transformer.layer.23.rel_attn.r_w_bias",
    "transformer.layer.23.rel_attn.seg_embed",
    "transformer.layer.23.rel_attn.layer_norm.weight",
    "transformer.layer.23.rel_attn.layer_norm.bias",
    "transformer.layer.23.ff.layer_norm.weight",
    "transformer.layer.23.ff.layer_norm.bias",
    "transformer.layer.23.ff.layer_1.weight",
    "transformer.layer.23.ff.layer_1.bias",
    "transformer.layer.23.ff.layer_2.weight",
    "transformer.layer.23.ff.layer_2.bias",
    "sequence_summary.summary.weight",
    "sequence_summary.summary.bias",
    "logits_proj.weight",
    "logits_proj.bias"
]
freeze_unfreeze_layers(model, layers_to_unfreeze=layers_to_unfreeze)
print(model.config)
print("=====")
print("num_parameters:", model.num_parameters())
print("num_trainable_parameters:", model.num_parameters(only_trainable=True))
print("=====")
#####
print("Experiment configuration used with this experiment:")
print("model used:", named_model)
print("learning rate used:", learning_rate)
print("number of epochs:", num_epochs)
print("maximum sequence length:", length_max)

```

```

print("batch size used:", size_batch)
print("regularization value:", regularization_weight_decay)
print("outcome variable:", y_col)
print("task:", x_task)
print("input column:", x_col)
#####
print("=====")
print("num_trainable_parameters:", model.num_parameters(only_trainable=True))

```

```

[ ]: for name, param in model.named_parameters():
    print(name, "requires_grad=", param.requires_grad)

```

```

[ ]: # Train & Evaluate
trained_model, trainer_obj = train_transformer_model(
    model = model,
    tokenizer = tokenizer,
    train_dataset = train_data_hf,
    val_dataset = val_data_hf,
    output_dir = dir_results,
    num_epochs = num_epochs,
    batch_size = size_batch,
    lr = learning_rate,
    weight_decay = regularization_weight_decay)
metrics = trainer_obj.evaluate()
print("Validation metrics:", metrics)
test_metrics = trainer_obj.evaluate(test_data_hf) if test_data_hf else None
print("Test metrics:", test_metrics)

```

```

[ ]: # save model checkpoint
timestamp = datetime.now().strftime("%Y%m%d_%H%M%S")
model_save_path = os.path.join(dir_models,
    f"{x_task}_{named_model}_{y_col}_{timestamp}")
trainer_obj.save_model(model_save_path)
print(f"Model checkpoint saved to: {model_save_path}")
# log experiment results
experiment_info = {
    "model_name": named_model,
    "learning_rate": learning_rate,
    "epochs": num_epochs,
    "batch_size": size_batch,
    "weight_decay": regularization_weight_decay,
    "x_task": x_task,
    "x_col": x_col,
    "y_col": y_col,
    "layers_to_unfreeze": layers_to_unfreeze}
model_info = gather_model_details(trained_model)
all_run_metrics = gather_all_run_metrics(

```

```
    trainer=trainer_obj,  
    train_dataset=train_data_hf,  
    val_dataset=val_data_hf,  
    test_dataset=test_data_hf)  
log_experiment_results_json(  
    experiment_meta=experiment_info,  
    model_details=model_info,  
    run_metrics=all_run_metrics,  
    log_file=log_filepath)  
print(f"EXPERIMENT LOGGED TO: {log_filepath}")
```