

```
import numpy as np
import pandas as pd
import time
import datetime
import gc
import random
from nltk.corpus import stopwords
import re
```

```
import torch
import torch.nn as nn
from torch.utils.data import TensorDataset, DataLoader, RandomSampler, SequentialSampler,random_split
from sklearn.model_selection import train_test_split
from sklearn.metrics import classification_report
```

```
!pip install transformers
```

```
Looking in indexes: https://pypi.org/simple, https://us-python.pkg.dev/colab-wheels/public/simple/
Collecting transformers
  Downloading transformers-4.26.0-py3-none-any.whl (6.3 MB)
    6.3/6.3 MB 95.0 MB/s eta 0:00:00
Requirement already satisfied: numpy>=1.17 in /usr/local/lib/python3.8/dist-packages (from transformers) (1.21.6)
Collecting tokenizers!=0.11.3,<0.14,>=0.11.1
  Downloading tokenizers-0.13.2-cp38-cp38-manylinux_2_17_x86_64.manylinux2014_x86_64.whl (7.6 MB)
    7.6/7.6 MB 102.0 MB/s eta 0:00:00
Requirement already satisfied: requests in /usr/local/lib/python3.8/dist-packages (from transformers) (2.25.1)
Requirement already satisfied: packaging>=20.0 in /usr/local/lib/python3.8/dist-packages (from transformers) (21.3)
Requirement already satisfied: regex!=2019.12.17 in /usr/local/lib/python3.8/dist-packages (from transformers) (2022.6.2)
Requirement already satisfied: filelock in /usr/local/lib/python3.8/dist-packages (from transformers) (3.9.0)
Requirement already satisfied: tqdm>=4.27 in /usr/local/lib/python3.8/dist-packages (from transformers) (4.64.1)
Collecting huggingface-hub<1.0,>=0.11.0
  Downloading huggingface_hub-0.12.0-py3-none-any.whl (190 kB)
    190.3/190.3 KB 24.7 MB/s eta 0:00:00
Requirement already satisfied: pyyaml>=5.1 in /usr/local/lib/python3.8/dist-packages (from transformers) (6.0)
Requirement already satisfied: typing-extensions>=3.7.4.3 in /usr/local/lib/python3.8/dist-packages (from huggingface-hub<1.0,>=0.11.0->transformers) (4.5.0)
Requirement already satisfied: pyparsing!=3.0.5,>=2.0.2 in /usr/local/lib/python3.8/dist-packages (from packaging>=20.0->transformers) (3.1.0)
Requirement already satisfied: idna<3,>=2.5 in /usr/local/lib/python3.8/dist-packages (from requests->transformers) (2.10)
Requirement already satisfied: certifi>=2017.4.17 in /usr/local/lib/python3.8/dist-packages (from requests->transformers) (2022.12.7)
Requirement already satisfied: urllib3<1.27,>=1.21.1 in /usr/local/lib/python3.8/dist-packages (from requests->transformers) (1.24.3)
Requirement already satisfied: chardet<5,>=3.0.2 in /usr/local/lib/python3.8/dist-packages (from requests->transformers) (4.0.0)
Installing collected packages: tokenizers, huggingface-hub, transformers
Successfully installed huggingface-hub-0.12.0 tokenizers-0.13.2 transformers-4.26.0
```

```
import transformers
from transformers import BertForSequenceClassification, AdamW, BertConfig,BertTokenizer,get_linear_schedule_with_warmup
```

```
device = torch.device("cuda:0" if torch.cuda.is_available() else "cpu")
device
```

```
device(type='cuda', index=0)
```

```
df = pd.read_csv('/content/train_all_tasks.csv')
```

```
df.head()
```

|   | rewire_id               | text   | label_sexist | label_category | label_vector |
|---|-------------------------|--|--------------|----------------|--------------|
| 0 | sexism2022_english-7358 | Damn, this writing was pretty chaotic                                | not sexist   | none           | none         |
| 1 | sexism2022_english-2367 | Yeah, and apparently a bunch of misogynistic v...<br>How the FUCK is | not sexist   | none           | none         |

```
!pip install nltk
```

```
Looking in indexes: https://pypi.org/simple, https://us-python.pkg.dev/colab-wheels/public/simple/
Requirement already satisfied: nltk in /usr/local/lib/python3.8/dist-packages (3.7)
Requirement already satisfied: joblib in /usr/local/lib/python3.8/dist-packages (from nltk) (1.2.0)
Requirement already satisfied: tqdm in /usr/local/lib/python3.8/dist-packages (from nltk) (4.64.1)
Requirement already satisfied: click in /usr/local/lib/python3.8/dist-packages (from nltk) (7.1.2)
Requirement already satisfied: regex>=2021.8.3 in /usr/local/lib/python3.8/dist-packages (from nltk) (2022.6.2)
```

```
import nltk
nltk.download('stopwords')
```

```
[nltk_data] Downloading package stopwords to /root/nltk_data...
[nltk_data] Unzipping corpora/stopwords.zip.
True
```

```
sw = stopwords.words('english')

def clean_text(text):

    text = text.lower()

    text = re.sub(r"^[^a-zA-Z?!.!;]+", " ", text) # replacing everything with space except (a-z, A-Z, ".", "?", "!", ";")

    text = re.sub(r"http\S+", "",text) #Removing URLs
    #text = re.sub(r"http", "",text)

    html=re.compile(r'<.*?>')

    text = html.sub(r'',text) #Removing html tags

    punctuations = '@#!?+&*[]-~.,:;/()$%^<|{}^' + "`" + " _"
    for p in punctuations:
        text = text.replace(p, '') #Removing punctuations

    text = [word.lower() for word in text.split() if word.lower() not in sw]

    text = " ".join(text) #removing stopwords

    emoji_pattern = re.compile("[
        u\"\\U0001F600-\\U0001F64F\" # emoticons
        u\"\\U0001F300-\\U0001F5FF\" # symbols & pictographs
        u\"\\U0001F680-\\U0001F6FF\" # transport & map symbols
        u\"\\U0001F1E0-\\U0001F1FF\" # flags (iOS)
        u\"\\U00002702-\\U000027B0\"
        u\"\\U000024C2-\\U0001F251\"
        ]+", flags=re.UNICODE)
    text = emoji_pattern.sub(r'', text) #Removing emojis

    return text
```

```
df['text'] = df['text'].apply(lambda x: clean_text(x))
```

```
tweets = df.text.values
labels = df.label_sexist.values
```

```
labels
```

```
array(['not sexist', 'not sexist', 'not sexist', ..., 'not sexist',
       'not sexist', 'not sexist'], dtype=object)
```

```
#labels = pd.get_dummies(df['label_sexist'],drop_first = True)
from sklearn.preprocessing import LabelEncoder
```

```
le = LabelEncoder()
```

```
labels = le.fit_transform(labels)
```

```
labels
```

```
array([0, 0, 0, ..., 0, 0, 0])
```

tweets

```
array(['damn, writing pretty chaotic',
      'yeah, apparently bunch misogynistic virgins ones turn gay woman straight, lol',
      'fuck woman still mp', ...,
      'lol imagine would sims like deserve ugly guy jfl',
      'not, girls go dates kiss first date, text back first date',
      'incel girlfriend fuck anyone says kissed legit girl fuck looking fag'],
      dtype=object)
```

labels

```
array([0, 0, 0, ..., 0, 0, 0])
```

# Load the BERT tokenizer

```
tokenizer = BertTokenizer.from_pretrained('bert-base-uncased', do_lower_case=True)
```

```
Downloading (...)solve/main/vocab.txt: 100%          232k/232k [00:00<00:00, 261kB/s]
Downloading (...)okenizer_config.json: 100%          28.0/28.0 [00:00<00:00, 1.63kB/s]
Downloading (...)lve/main/config.json: 100%          570/570 [00:00<00:00, 27.4kB/s]
```

```
print(' Original: ', tweets[0])
```

# Print the sentence split into tokens.

```
print('Tokenized: ', tokenizer.tokenize(tweets[0]))
```

# Print the sentence mapped to token ids.

```
print('Token IDs: ', tokenizer.convert_tokens_to_ids(tokenizer.tokenize(tweets[0])))
```

```
Original:  damn, writing pretty chaotic
Tokenized: ['damn', ' ', 'writing', 'pretty', 'chaotic']
Token IDs: [4365, 1010, 3015, 3492, 19633]
```

```
max_len = 0
```

# For every sentence...

```
for sent in tweets:
```

```
    # Tokenize the text and add `[CLS]` and `[SEP]` tokens.
```

```
    input_ids = tokenizer.encode(sent, add_special_tokens=True)
```

```
    # Update the maximum sentence length.
```

```
    max_len = max(max_len, len(input_ids))
```

```
print('Max sentence length: ', max_len)
```

```
Max sentence length: 59
```

```
input_ids = []
```

```
attention_masks = []
```

# For every tweet...

```
for tweet in tweets:
```

```
    # `encode_plus` will:
```

```
    # (1) Tokenize the sentence.
```

```
    # (2) Prepend the `[CLS]` token to the start.
```

```
    # (3) Append the `[SEP]` token to the end.
```

```
    # (4) Map tokens to their IDs.
```

```
    # (5) Pad or truncate the sentence to `max_length`
```

```
    # (6) Create attention masks for [PAD] tokens.
```

```
    encoded_dict = tokenizer.encode_plus(
```

```
        tweet,                # Sentence to encode.
```

```
        add_special_tokens = True, # Add `[CLS]` and `[SEP]`
```

```
        max_length = max_len,    # Pad & truncate all sentences.
```

```
        pad_to_max_length = True,
```

```
        return_attention_mask = True, # Construct attn. masks.
```

```
        return_tensors = 'pt',    # Return pytorch tensors.
```

```
)
```

```

# Add the encoded sentence to the list.
input_ids.append(encoded_dict['input_ids'])

# And its attention mask (simply differentiates padding from non-padding).
attention_masks.append(encoded_dict['attention_mask'])

# Convert the lists into tensors.
input_ids = torch.cat(input_ids, dim=0)
attention_masks = torch.cat(attention_masks, dim=0)
labels = torch.tensor(labels)

## Print sentence 0, now as a list of IDs.
print('Original: ', tweets[0])
print('Token IDs:', input_ids[0])

Truncation was not explicitly activated but `max_length` is provided a specific value, please use `truncation=True` to explicitly truncate inputs to the maximum length.
/usr/local/lib/python3.8/dist-packages/transformers/tokenization_utils_base.py:2339: FutureWarning: The `pad_to_max_length` argument is deprecated and will be removed in a future version, using the `padding` model argument instead.
warnings.warn(
Original: damn, writing pretty chaotic
Token IDs: tensor([ 101, 4365, 1010, 3015, 3492, 19633, 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])

# Combine the training inputs into a TensorDataset.
dataset = TensorDataset(input_ids, attention_masks, labels)

# Create a 90-10 train-validation split.

# Calculate the number of samples to include in each set.
train_size = int(0.8 * len(dataset))
#val_size = int(0.2 * len(dataset))
val_size = len(dataset) - train_size

# Divide the dataset by randomly selecting samples.
train_dataset, val_dataset = random_split(dataset, [train_size, val_size])

print('{:>5,} training samples'.format(train_size))
print('{:>5,} validation samples'.format(val_size))

11,200 training samples
2,800 validation samples

# The DataLoader needs to know our batch size for training, so we specify it
# here. For fine-tuning BERT on a specific task, the authors recommend a batch
# size of 16 or 32.
batch_size = 32

# Create the DataLoaders for our training and validation sets.
# We'll take training samples in random order.
train_dataloader = DataLoader(
    train_dataset, # The training samples.
    sampler = RandomSampler(train_dataset), # Select batches randomly
    batch_size = batch_size # Trains with this batch size.
)

# For validation the order doesn't matter, so we'll just read them sequentially.
validation_dataloader = DataLoader(
    val_dataset, # The validation samples.
    sampler = SequentialSampler(val_dataset), # Pull out batches sequentially.
    batch_size = batch_size # Evaluate with this batch size.
)

# Load BertForSequenceClassification, the pretrained BERT model with a single
# linear classification layer on top.
model = BertForSequenceClassification.from_pretrained(
    "bert-base-uncased", # Use the 12-layer BERT model, with an uncased vocab.
    num_labels = 2, # The number of output labels--2 for binary classification.
    # You can increase this for multi-class tasks.
    output_attentions = False, # Whether the model returns attentions weights.
    output_hidden_states = False, # Whether the model returns all hidden-states.
)

```

```
# if device == "cuda:0":
# # Tell pytorch to run this model on the GPU.
#     model = model.cuda()
model = model.to(device)
```

Downloading (...) "pytorch\_model.bin": 100%

440M/440M [00:02<00:00, 146MB/s]

Some weights of the model checkpoint at bert-base-uncased were not used when initializing BertForSequenceClassification: ['cls.predictio  
- This IS expected if you are initializing BertForSequenceClassification from the checkpoint of a model trained on another task or with  
- This IS NOT expected if you are initializing BertForSequenceClassification from the checkpoint of a model that you expect to be exactl  
Some weights of BertForSequenceClassification were not initialized from the model checkpoint at bert-base-uncased and are newly initiali  
You should probably TRAIN this model on a down-stream task to be able to use it for predictions and inference.

```
optimizer = AdamW(model.parameters()),
                lr = 2e-5, # args.learning_rate - default is 5e-5, our notebook had 2e-5
                eps = 1e-8 # args.adam_epsilon - default is 1e-8.
            )
```

/usr/local/lib/python3.8/dist-packages/transformers/optimization.py:306: FutureWarning: This implementation of AdamW is deprecated and w  
warnings.warn(

```
# Number of training epochs. The BERT authors recommend between 2 and 4.
# We chose to run for 4, but we'll see later that this may be over-fitting the
# training data.
epochs = 5
```

```
# Total number of training steps is [number of batches] x [number of epochs].
# (Note that this is not the same as the number of training samples).
total_steps = len(train_dataloader) * epochs
```

```
# Create the learning rate scheduler.
scheduler = get_linear_schedule_with_warmup(optimizer,
                                             num_warmup_steps = 0, # Default value in run_glue.py
                                             num_training_steps = total_steps)
```

```
# Function to calculate the accuracy of our predictions vs labels
def flat_accuracy(preds, labels):
    pred_flat = np.argmax(preds, axis=1).flatten()
    labels_flat = labels.flatten()
    return np.sum(pred_flat == labels_flat) / len(labels_flat)
```

```
def format_time(elapsed):
    """
    Takes a time in seconds and returns a string hh:mm:ss
    """
    # Round to the nearest second.
    elapsed_rounded = int(round((elapsed)))
    # Format as hh:mm:ss
    return str(datetime.timedelta(seconds=elapsed_rounded))
```

```
seed_val = 42
random.seed(seed_val)
np.random.seed(seed_val)
torch.manual_seed(seed_val)
torch.cuda.manual_seed_all(seed_val)
training_stats = []
```

```
# Measure the total training time for the whole run.
total_t0 = time.time()
```

```
# For each epoch...
for epoch_i in range(0, epochs):
```

```
    # =====
    #             Training
    # =====
    # Perform one full pass over the training set.
    print("")
    print('==== Epoch {:} / {:} ====='.format(epoch_i + 1, epochs))
    print('Training...')
    # Measure how long the training epoch takes.
```

```

t0 = time.time()
total_train_loss = 0
model.train()
for step, batch in enumerate(train_dataloader):
    # Unpack this training batch from our dataloader.
    #
    # As we unpack the batch, we'll also copy each tensor to the device using the
    # `to` method.
    #
    # `batch` contains three pytorch tensors:
    # [0]: input ids
    # [1]: attention masks
    # [2]: labels
    b_input_ids = batch[0].to(device)
    b_input_mask = batch[1].to(device)
    b_labels = batch[2].to(device)
    optimizer.zero_grad()
    output = model(b_input_ids,
                    token_type_ids=None,
                    attention_mask=b_input_mask,
                    labels=b_labels)

    loss = output.loss
    total_train_loss += loss.item()
    # Perform a backward pass to calculate the gradients.
    loss.backward()
    # Clip the norm of the gradients to 1.0.
    # This is to help prevent the "exploding gradients" problem.
    torch.nn.utils.clip_grad_norm_(model.parameters(), 1.0)
    # Update parameters and take a step using the computed gradient.
    # The optimizer dictates the "update rule"--how the parameters are
    # modified based on their gradients, the learning rate, etc.
    optimizer.step()
    # Update the learning rate.
    scheduler.step()

# Calculate the average loss over all of the batches.
avg_train_loss = total_train_loss / len(train_dataloader)

# Measure how long this epoch took.
training_time = format_time(time.time() - t0)
print("")
print(" Average training loss: {:.2f}".format(avg_train_loss))
print(" Training epoch took: {}".format(training_time))
# =====
# Validation
# =====
# After the completion of each training epoch, measure our performance on
# our validation set.
print("")
print("Running Validation...")
t0 = time.time()
# Put the model in evaluation mode--the dropout layers behave differently
# during evaluation.
model.eval()
# Tracking variables
total_eval_accuracy = 0
best_eval_accuracy = 0
total_eval_loss = 0
nb_eval_steps = 0
# Evaluate data for one epoch
for batch in validation_dataloader:
    b_input_ids = batch[0].to(device)
    b_input_mask = batch[1].to(device)
    b_labels = batch[2].to(device)
    # Tell pytorch not to bother with constructing the compute graph during
    # the forward pass, since this is only needed for backprop (training).
    with torch.no_grad():
        output= model(b_input_ids,
                       token_type_ids=None,
                       attention_mask=b_input_mask,
                       labels=b_labels)

    loss = output.loss
    total_eval_loss += loss.item()
    # Move logits and labels to CPU if we are using GPU
    logits = output.logits
    logits = logits.detach().cpu().numpy()
    label_ids = b_labels.to('cpu').numpy()

```

```

    # Calculate the accuracy for this batch of test sentences, and
    # accumulate it over all batches.
    total_eval_accuracy += flat_accuracy(logits, label_ids)
# Report the final accuracy for this validation run.
avg_val_accuracy = total_eval_accuracy / len(validation_dataloader)
print(" Accuracy: {:.2f}".format(avg_val_accuracy))
# Calculate the average loss over all of the batches.
avg_val_loss = total_eval_loss / len(validation_dataloader)
# Measure how long the validation run took.
validation_time = format_time(time.time() - t0)
if avg_val_accuracy > best_eval_accuracy:
    torch.save(model, 'bert_model')
    best_eval_accuracy = avg_val_accuracy
#print(" Validation Loss: {:.2f}".format(avg_val_loss))
#print(" Validation took: {}".format(validation_time))
# Record all statistics from this epoch.
training_stats.append(
    {
        'epoch': epoch_i + 1,
        'Training Loss': avg_train_loss,
        'Valid. Loss': avg_val_loss,
        'Valid. Accur.': avg_val_accuracy,
        'Training Time': training_time,
        'Validation Time': validation_time
    }
)
print("")
print("Training complete!")

print("Total training took {} (h:mm:ss)".format(format_time(time.time()-total_t0)))

```

```

===== Epoch 1 / 5 =====
Training...

```

```

Average training loss: 0.18
Training epoch took: 0:01:55

```

```

Running Validation...
Accuracy: 0.85

```

```

===== Epoch 2 / 5 =====
Training...

```

```

Average training loss: 0.22
Training epoch took: 0:01:54

```

```

Running Validation...
Accuracy: 0.85

```

```

===== Epoch 3 / 5 =====
Training...

```

```

Average training loss: 0.13
Training epoch took: 0:01:55

```

```

Running Validation...
Accuracy: 0.85

```

```

===== Epoch 4 / 5 =====
Training...

```

```

Average training loss: 0.09
Training epoch took: 0:01:55

```

```

Running Validation...
Accuracy: 0.82

```

```

===== Epoch 5 / 5 =====
Training...

```

```

Average training loss: 0.06
Training epoch took: 0:01:55

```

```

Running Validation...
Accuracy: 0.82

```

```

Training complete!
Total training took 0:10:29 (h:mm:ss)

```

```
model = torch.load('bert_model')
```

```
df_test = pd.read_csv('/content/test_task_a_entries.csv')
df_test['text'] = df_test['text'].apply(lambda x:clean_text(x))
test_tweets = df_test['text'].values
```

```
df_test.head()
```

|   | rewire_id                | text  |
|---|--------------------------|---|
| 0 | sexism2022_english-7207  | oregon coast snow colder witch titty though       |
| 1 | sexism2022_english-10731 | tall man must certaily better women large musc... |
| 2 | sexism2022_english-11374 | hit rode chode, subreddit                         |
| 3 | sexism2022_english-7356  | lawyer chick shoot mcdonald guy, every time, m... |
| 4 | sexism2022_english-11976 | true, totally hating females want female ish      |

```
test_input_ids = []
test_attention_masks = []
for tweet in test_tweets:
    encoded_dict = tokenizer.encode_plus(
        tweet,
        add_special_tokens = True,
        max_length = max_len,
        pad_to_max_length = True,
        return_attention_mask = True,
        return_tensors = 'pt',
    )
    test_input_ids.append(encoded_dict['input_ids'])
    test_attention_masks.append(encoded_dict['attention_mask'])
test_input_ids = torch.cat(test_input_ids, dim=0)
test_attention_masks = torch.cat(test_attention_masks, dim=0)
```

```
/usr/local/lib/python3.8/dist-packages/transformers/tokenization_utils_base.py:2339: FutureWarning: The `pad_to_max_length` argument is
warnings.warn(
```

```
test_dataset = TensorDataset(test_input_ids, test_attention_masks)
test_dataloader = DataLoader(
    test_dataset, # The validation samples.
    sampler = SequentialSampler(test_dataset), # Pull out batches sequentially.
    batch_size = batch_size # Evaluate with this batch size.
)
```

```
predictions = []
for batch in test_dataloader:
    b_input_ids = batch[0].to(device)
    b_input_mask = batch[1].to(device)
    with torch.no_grad():
        output= model(b_input_ids,
                       token_type_ids=None,
                       attention_mask=b_input_mask)

        logits = output.logits
        logits = logits.detach().cpu().numpy()
        pred_flat = np.argmax(logits, axis=1).flatten()

        predictions.extend(list(pred_flat))
```

```
df_output = pd.DataFrame()
df_output['rewire_id'] = df_test['rewire_id']
df_output['label_pred'] =predictions
df_output.to_csv('submission2.csv',index=False)
```

```
df1 = pd.read_csv('/content/submission2.csv')
```

```
df1.head()
```



|   | rewire_id                | label_pred |
|---|--------------------------|------------|
| 0 | sexism2022_english-7207  | 1          |
| 1 | sexism2022_english-10731 | 0          |
| 2 | sexism2022_english-11374 | 0          |
| 3 | sexism2022_english-7356  | 0          |
| 4 | sexism2022_english-11976 | 0          |

```
df1.columns
```

```
Index(['rewire_id', 'label_pred'], dtype='object')
```

```
df1['label_pred'].value_counts()
```

```
0    3047
1     953
Name: label_pred, dtype: int64
```

```
df1['label_pred'] = df1['label_pred'].replace(1,'sexist')
```

```
df1['label_pred'] = df1['label_pred'].replace(0,'not sexist')
```

```
df1.head()
```

|   | rewire_id                | label_pred |
|---|--------------------------|------------|
| 0 | sexism2022_english-7207  | sexist     |
| 1 | sexism2022_english-10731 | not sexist |
| 2 | sexism2022_english-11374 | not sexist |
| 3 | sexism2022_english-7356  | not sexist |
| 4 | sexism2022_english-11976 | not sexist |

```
df_output = pd.DataFrame()
df_output['rewire_id'] = df_test['rewire_id']
df_output['label_pred'] = df1['label_pred']
df_output.to_csv('submission2.csv', index=False)
```

```
data = pd.read_csv('/content/submission2.csv')
```

```
data.columns
```

```
data.shape
```

```
data['label_pred'].value_counts()
```

✓ 0s completed at 9:10 PM

● ×