```
!pip install transformers
!pip install sentencepiece
import sentencepiece
import torch
import torch.nn as nn
import numpy as np
import matplotlib.pyplot as plt
import seaborn as sns
from sklearn.metrics import confusion_matrix
from datetime import datetime
from pathlib import Path
import pandas as pd
import torchtext.data as ttd
```

```
    Requirement already satisfied: transformers in /usr/local/lib/python3.6/dist-packages
    Requirement already satisfied: requests in /usr/local/lib/python3.6/dist-packages (fr
    Requirement already satisfied: tqdm>=4.27 in /usr/local/lib/python3.6/dist-packages (
    Requirement already satisfied: numpy in /usr/local/lib/python3.6/dist-packages (from
    Requirement already satisfied: tokenizers==0.9.4 in /usr/local/lib/python3.6/dist-pac
    Requirement already satisfied: sacremoses in /usr/local/lib/python3.6/dist-packages (
    Requirement already satisfied: dataclasses; python_version < "3.7" in /usr/local/lib/
    Requirement already satisfied: filelock in /usr/local/lib/python3.6/dist-packages (fr
    Requirement already satisfied: regex!=2019.12.17 in /usr/local/lib/python3.6/dist-pac
    Requirement already satisfied: packaging in /usr/local/lib/python3.6/dist-packages (1
    Requirement already satisfied: urllib3!=1.25.0,!=1.25.1,<1.26,>=1.21.1 in /usr/local/
    Requirement already satisfied: chardet<4,>=3.0.2 in /usr/local/lib/python3.6/dist-pac
    Requirement already satisfied: idna<3,>=2.5 in /usr/local/lib/python3.6/dist-packages
    Requirement already satisfied: certifi>=2017.4.17 in /usr/local/lib/python3.6/dist-pa
    Requirement already satisfied: joblib in /usr/local/lib/python3.6/dist-packages (from
    Requirement already satisfied: six in /usr/local/lib/python3.6/dist-packages (from sa
    Requirement already satisfied: click in /usr/local/lib/python3.6/dist-packages (from
    Requirement already satisfied: pyparsing>=2.0.2 in /usr/local/lib/python3.6/dist-pack
    Requirement already satisfied: sentencepiece in /usr/local/lib/python3.6/dist-package
```

# ▾ Loading Dataset

We will use The 20 Newsgroups dataset Dataset homepage:

Scikit-learn includes some nice helper functions for retrieving the 20 Newsgroups dataset--
https://scikit-learn.org/stable/modules/generated/sklearn.datasets.fetch_20newsgroups.html.
We'll use them below to retrieve the dataset.

Also look at results fron non- neural net models here : https://scikit-
learn.org/stable/auto_examples/text/plot_document_classification_20newsgroups.html#sphx-
glr-auto-examples-text-plot-document-classification-20newsgroups-py

```
gpu_info = !nvidia-smi
gpu_info = '\n'.join(gpu_info)
if gpu_info.find('failed') >= 0:
  print('Select the Runtime > "Change runtime type" menu to enable a GPU accelera
  print('and then re-execute this cell.')
else:
```

```
print(gpu_info)
```

```
    Wed Dec  2 07:46:54 2020
    +-----------------------------------------------------------------------------+
    | NVIDIA-SMI 455.38       Driver Version: 418.67       CUDA Version: 10.1     |
    |-------------------------------+----------------------+----------------------+
    | GPU  Name        Persistence-M| Bus-Id        Disp.A | Volatile Uncorr. ECC |
    | Fan  Temp  Perf  Pwr:Usage/Cap|         Memory-Usage | GPU-Util  Compute M. |
    |                               |                      |               MIG M. |
    |===============================+======================+======================|
    |   0  Tesla V100-SXM2...  Off  | 00000000:00:04.0 Off |                    0 |
    | N/A   33C    P0    24W / 300W |      0MiB / 16130MiB |      0%      Default |
    |                               |                      |                 ERR! |
    +-------------------------------+----------------------+----------------------+

    +-----------------------------------------------------------------------------+
    | Processes:                                                                  |
    |  GPU   GI   CI        PID   Type   Process name                  GPU Memory |
    |        ID   ID                                                   Usage      |
    |=============================================================================|
    |  No running processes found                                                 |
    +-----------------------------------------------------------------------------+
```

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

```
    cuda:0
```

```python
from sklearn.datasets import fetch_20newsgroups

train = fetch_20newsgroups(subset='train',
                           remove=('headers', 'footers', 'quotes'))

test = fetch_20newsgroups(subset='test',
                          remove=('headers', 'footers', 'quotes'))
```

```python
print(train.data[0])
```

```
    I was wondering if anyone out there could enlighten me on this car I saw
    the other day. It was a 2-door sports car, looked to be from the late 60s/
    early 70s. It was called a Bricklin. The doors were really small. In addition,
    the front bumper was separate from the rest of the body. This is
    all I know. If anyone can tellme a model name, engine specs, years
    of production, where this car is made, history, or whatever info you
    have on this funky looking car, please e-mail.
```

```python
print(train.target[0])
```

```
    7
```

```python
train.target_names
```

```
    ['alt.atheism',
     'comp.graphics',
```
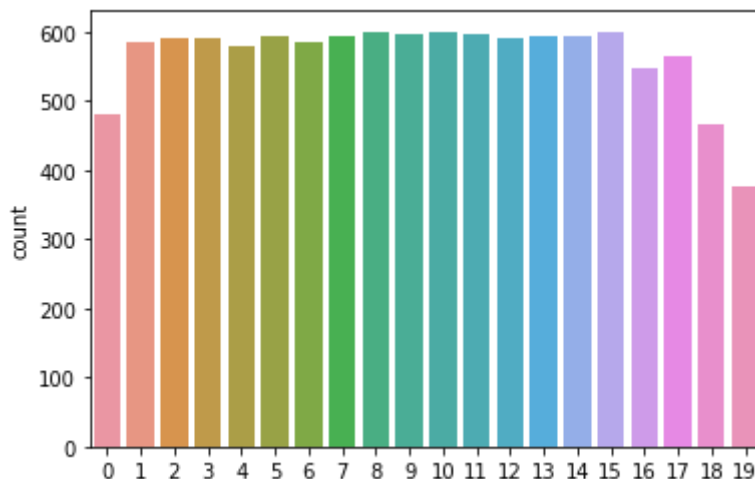
```
        'comp.os.ms-windows.misc',
        'comp.sys.ibm.pc.hardware',
        'comp.sys.mac.hardware',
        'comp.windows.x',
        'misc.forsale',
        'rec.autos',
        'rec.motorcycles',
        'rec.sport.baseball',
        'rec.sport.hockey',
        'sci.crypt',
        'sci.electronics',
        'sci.med',
        'sci.space',
        'soc.religion.christian',
        'talk.politics.guns',
        'talk.politics.mideast',
        'talk.politics.misc',
        'talk.religion.misc']
```

```
import seaborn as sns

# Plot the number of tokens of each length.
sns.countplot(train.target);
```

/usr/local/lib/python3.6/dist-packages/seaborn/_decorators.py:43: FutureWarning: Pass
  FutureWarning



## XLNet with 128 features and truncating at end

```
from transformers import XLNetTokenizer

# Load the BERT tokenizer.
print('Loading XLNet tokenizer...')
tokenizer = XLNetTokenizer.from_pretrained('xlnet-base-cased', do_lower_case=True
```

Loading XLNet tokenizer...

Downloading: 100%                          798k/798k [00:00<00:00, 1.77MB/s]

```python
# Tokenize all of the sentences and map the tokens to thier word IDs.
input_ids = []
attention_masks = []

# For every sentence...
for sent in train.data:
    # `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(
                        sent,                      # Sentence to encode.
                        add_special_tokens = True, # Add '[CLS]' and '[SEP]'
                        truncation=True, #Truncate the sentences
                        max_length = 128,          # Pad & truncate all sentence
                        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(train.target)

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

```
    /usr/local/lib/python3.6/dist-packages/transformers/tokenization_utils_base.py:2142:
      FutureWarning,
    Original:  I was wondering if anyone out there could enlighten me on this car I saw
    the other day. It was a 2-door sports car, looked to be from the late 60s/
    early 70s. It was called a Bricklin. The doors were really small. In addition,
    the front bumper was separate from the rest of the body. This is
    all I know. If anyone can tellme a model name, engine specs, years
    of production, where this car is made, history, or whatever info you
    have on this funky looking car, please e-mail.
    Token IDs: tensor([    5,     5,     5,    17,   150,    30,  7083,   108,  1216,
               105,   121, 22531,   110,    31,    52,   398,    17,   150,   685,
                18,    86,   191,     9,    36,    30,    24,   159,    13,  9381,
              1721,   398,    19,   719,    22,    39,    40,    18,   471,  1639,
                23,   167,   319,  2415,    23,     9,    36,    30,   271,    24,
              5989,  1554,     9,    18,  3965,    55,   343,   316,     9,    25,
               864,    19,    18,   605, 19990,    30,  1731,    40,    18,   904,
                20,    18,   458,     9,    52,    27,    71,    17,   150,   175,
                 9,   108,  1216,    64,   759,  1088,    24,  1342,   304,    19,
```

```
        2012,    17,    23, 10112,    23,    19,   123,    20,   845,    19,
         131,    52,   398,    27,   140,    19,   614,    19,    49,  2636,
        7549,    44,    47,    31,    52,  1572,  3531,   589,   398,    19,
        1282,    17,    93,    13,  1635,     9,     4,     3])
```

```python
test_input_ids = []
test_attention_masks = []

# For every sentence...
for sent in test.data:
    # `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(
                        sent,                      # Sentence to encode.
                        add_special_tokens = True, # Add '[CLS]' and '[SEP]'
                        truncation=True, #Truncate the sentences
                        max_length = 128,          # Pad & truncate all sentence
                        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.
    test_input_ids.append(encoded_dict['input_ids'])

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

# Convert the lists into tensors.
test_input_ids = torch.cat(test_input_ids, dim=0)
test_attention_masks = torch.cat(test_attention_masks, dim=0)
test_labels = torch.tensor(test.target)

# Print sentence 0, now as a list of IDs.
print('Original: ', test.data[0])
print('Token IDs:', test_input_ids[0])
```

```
    /usr/local/lib/python3.6/dist-packages/transformers/tokenization_utils_base.py:2142:
      FutureWarning,
    Original:  I am a little confused on all of the models of the 88-89 bonnevilles.
    I have heard of the LE SE LSE SSE SSEI. Could someone tell me the
    differences are far as features or performance. I am also curious to
    know what the book value is for prefereably the 89 model. And how much
    less than book value can you usually get them for. In other words how
    much are they in demand this time of year. I have heard that the mid-spring
    early summer is the best time to buy.
    Token IDs: tensor([    5,     5,     5,     5,    17,   150,   569,    24,   293,   68
              31,    71,    20,    18,  2626,    20,    18, 10227,    13,  4406,
              17,  4769,   667,  1948,    23,     9,    17,   150,    47,  1133,
              20,    18,    17,   529,    17,  1022,    17,   368,  1022,    17,
```

```
        23,  1022,    17,    23,    23,  5730,     9,   121,   886,   759,
       110,    18,  3589,    41,   420,    34,  1091,    49,   922,     9,
        17,   150,   569,    77,  8595,    22,   175,   113,    18,   522,
       991,    27,    28,  3948,    93,  3513,    18, 11903,  1342,     9,
        21,   160,   178,   486,   100,   522,   991,    64,    44,  1044,
       133,   107,    28,     9,    25,    86,  1006,   160,   178,    41,
        63,    25,  1480,    52,    92,    20,   119,     9,    17,   150,
        47,  1133,    29,    18,  1359,    13, 20343,   319,  1148,    27,
        18,   252,    92,    22,   971,     9,     4,     3])
```

```python
from torch.utils.data import TensorDataset, random_split

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

# Create a 90-10 train-validation split.

# Calculate the number of samples to include in each set.
train_size = int(0.9 * 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))
print('{:>5,} test samples'.format(len(test_dataset)))
```

```
    10,182 training samples
    1,132 validation samples
    7,532 test samples
```

```python
from torch.utils.data import DataLoader, RandomSampler, SequentialSampler

# 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 = 8

# 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 sequenti
            batch_size = batch_size # Evaluate with this batch size.
        )
```

```
test_dataloader = DataLoader(
            test_dataset,  # The training samples.
            sampler = RandomSampler(test_dataset), # Select batches randomly
            batch_size = batch_size # Trains with this batch size.
        )
```

# Training just the Fully Connected Classifier layer by freezing the XLNet model weights

```
from transformers import  XLNetModel

bert_model = XLNetModel.from_pretrained('xlnet-base-cased')
```

Downloading: 100%                                760/760 [00:00<00:00, 1.95kB/s]


Downloading: 100%                                467M/467M [00:06<00:00, 68.2MB/s]


```
bert_model
```

```
            (dropout): Dropout(p=0.1, inplace=False)
          )
          (dropout): Dropout(p=0.1, inplace=False)
        )
        (8): 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)
          )
          (dropout): Dropout(p=0.1, inplace=False)
        )
        (9): 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)
          )
          (dropout): Dropout(p=0.1, inplace=False)
        )
        (10): 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)
          )
          (dropout): Dropout(p=0.1, inplace=False)
        )
        (11): 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)
          )
          (dropout): Dropout(p=0.1, inplace=False)
        )
      )
      (dropout): Dropout(p=0.1, inplace=False)
    )
```

```python
# Define the model
class linear(nn.Module):

  def __init__(self, bert_model, n_outputs, dropout_rate):

    super(linear, self).__init__()

    #self.D = bert_model.config.to_dict()['hidden_size']
    self.bert_model = bert_model
    self.K = n_outputs
    self.dropout_rate=dropout_rate

    # embedding layer
    #self.embed = nn.Embedding(self.V, self.D)


    # dense layer
    self.fc = nn.Linear(768 , self.K)

    # dropout layer
    self.dropout= nn.Dropout(self.dropout_rate)

  def forward(self, X):

    with torch.no_grad():
      embedding = self.bert_model(X)[0][:,0,:]

    #embedding= self.dropout(embedding)
```

```
    output = self.fc(embedding)
    output= self.dropout(output)

    return output


n_outputs = 20
dropout_rate = 0.5


#model = RNN(n_vocab, embed_dim, n_hidden, n_rnnlayers, n_outputs, bidirectional,
model = linear(bert_model, n_outputs, dropout_rate)
model.to(device)
```

```
            )
            (8): 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)
              )
              (dropout): Dropout(p=0.1, inplace=False)
            )
            (9): 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)
              )
              (dropout): Dropout(p=0.1, inplace=False)
            )
            (10): 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)
              )
              (dropout): Dropout(p=0.1, inplace=False)
            )
            (11): 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_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)
      )
      (dropout): Dropout(p=0.1, inplace=False)
    )
  )
  (dropout): Dropout(p=0.1, inplace=False)
)
(fc): Linear(in_features=768, out_features=20, bias=True)
(dropout): Dropout(p=0.5, inplace=False)
)
```

```
print(model)
```

```
        (dropout): Dropout(p=0.1, inplace=False)
      )
      (8): 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)
        )
        (dropout): Dropout(p=0.1, inplace=False)
      )
      (9): 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)
        )
        (dropout): Dropout(p=0.1, inplace=False)
      )
      (10): 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)
        )
        (dropout): Dropout(p=0.1, inplace=False)
      )
      (11): 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)
          )
          (dropout): Dropout(p=0.1, inplace=False)
        )
      )
      (dropout): Dropout(p=0.1, inplace=False)
    )
    (fc): Linear(in_features=768, out_features=20, bias=True)
    (dropout): Dropout(p=0.5, inplace=False)
  )
```

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

```
    bert_model.layer.8.ff.layer_norm.weight torch.Size([768])
    bert_model.layer.8.ff.layer_norm.bias torch.Size([768])
    bert_model.layer.8.ff.layer_1.weight torch.Size([3072, 768])
    bert_model.layer.8.ff.layer_1.bias torch.Size([3072])
    bert_model.layer.8.ff.layer_2.weight torch.Size([768, 3072])
    bert_model.layer.8.ff.layer_2.bias torch.Size([768])
    bert_model.layer.9.rel_attn.q torch.Size([768, 12, 64])
    bert_model.layer.9.rel_attn.k torch.Size([768, 12, 64])
    bert_model.layer.9.rel_attn.v torch.Size([768, 12, 64])
    bert_model.layer.9.rel_attn.o torch.Size([768, 12, 64])
    bert_model.layer.9.rel_attn.r torch.Size([768, 12, 64])
    bert_model.layer.9.rel_attn.r_r_bias torch.Size([12, 64])
    bert_model.layer.9.rel_attn.r_s_bias torch.Size([12, 64])
    bert_model.layer.9.rel_attn.r_w_bias torch.Size([12, 64])
    bert_model.layer.9.rel_attn.seg_embed torch.Size([2, 12, 64])

    bert_model.layer.9.rel_attn.layer_norm.weight torch.Size([768])
    bert_model.layer.9.rel_attn.layer_norm.bias torch.Size([768])
    bert_model.layer.9.ff.layer_norm.weight torch.Size([768])
    bert_model.layer.9.ff.layer_norm.bias torch.Size([768])
    bert_model.layer.9.ff.layer_1.weight torch.Size([3072, 768])
    bert_model.layer.9.ff.layer_1.bias torch.Size([3072])
    bert_model.layer.9.ff.layer_2.weight torch.Size([768, 3072])
    bert_model.layer.9.ff.layer_2.bias torch.Size([768])
    bert_model.layer.10.rel_attn.q torch.Size([768, 12, 64])
    bert_model.layer.10.rel_attn.k torch.Size([768, 12, 64])
    bert_model.layer.10.rel_attn.v torch.Size([768, 12, 64])
    bert_model.layer.10.rel_attn.o torch.Size([768, 12, 64])
    bert_model.layer.10.rel_attn.r torch.Size([768, 12, 64])
    bert_model.layer.10.rel_attn.r_r_bias torch.Size([12, 64])
    bert_model.layer.10.rel_attn.r_s_bias torch.Size([12, 64])
    bert_model.layer.10.rel_attn.r_w_bias torch.Size([12, 64])
    bert_model.layer.10.rel_attn.seg_embed torch.Size([2, 12, 64])
    bert_model.layer.10.rel_attn.layer_norm.weight torch.Size([768])
    bert_model.layer.10.rel_attn.layer_norm.bias torch.Size([768])
    bert_model.layer.10.ff.layer_norm.weight torch.Size([768])
    bert_model.layer.10.ff.layer_norm.bias torch.Size([768])
    bert_model.layer.10.ff.layer_1.weight torch.Size([3072, 768])
    bert_model.layer.10.ff.layer_1.bias torch.Size([3072])
    bert_model.layer.10.ff.layer_2.weight torch.Size([768, 3072])
    bert_model.layer.10.ff.layer_2.bias torch.Size([768])
    bert_model.layer.11.rel_attn.q torch.Size([768, 12, 64])
    bert_model.layer.11.rel_attn.k torch.Size([768, 12, 64])
    bert_model.layer.11.rel_attn.v torch.Size([768, 12, 64])
```

```
            bert_model.layer.11.rel_attn.o torch.Size([768, 12, 64])
            bert_model.layer.11.rel_attn.r torch.Size([768, 12, 64])
            bert_model.layer.11.rel_attn.r_r_bias torch.Size([12, 64])
            bert_model.layer.11.rel_attn.r_s_bias torch.Size([12, 64])
            bert_model.layer.11.rel_attn.r_w_bias torch.Size([12, 64])
            bert_model.layer.11.rel_attn.seg_embed torch.Size([2, 12, 64])
            bert_model.layer.11.rel_attn.layer_norm.weight torch.Size([768])
            bert_model.layer.11.rel_attn.layer_norm.bias torch.Size([768])
            bert_model.layer.11.ff.layer_norm.weight torch.Size([768])
            bert_model.layer.11.ff.layer_norm.bias torch.Size([768])
            bert_model.layer.11.ff.layer_1.weight torch.Size([3072, 768])
            bert_model.layer.11.ff.layer_1.bias torch.Size([3072])
            bert_model.layer.11.ff.layer_2.weight torch.Size([768, 3072])
            bert_model.layer.11.ff.layer_2.bias torch.Size([768])
            fc.weight torch.Size([20, 768])
            fc.bias torch.Size([20])
```

```python
import random

seed = 123

random.seed(seed)
np.random.seed(seed)
torch.manual_seed(seed)
torch.cuda.manual_seed_all(seed)

learning_rate = 0.001
epochs=10

# STEP 5: INSTANTIATE LOSS CLASS
criterion = nn.CrossEntropyLoss()

# STEP 6: INSTANTIATE OPTIMIZER CLASS

optimizer = torch.optim.Adam(model.parameters(), lr=learning_rate)

# Freeze embedding Layer

#freeze embeddings
#model.embed.weight.requires_grad  = False

# STEP 7: TRAIN THE MODEL

train_losses= np.zeros(epochs)
valid_losses= np.zeros(epochs)


for epoch in range(epochs):

  t0= datetime.now()
  train_loss=[]

  model.train()
  for batch in train_dataloader:

    # forward pass
```

```
    output= model(batch[0].to(device))
    loss=criterion(output,batch[2].to(device))

    # set gradients to zero
    optimizer.zero_grad()

    # backward pass
    loss.backward()
    optimizer.step()
    train_loss.append(loss.item())

  train_loss=np.mean(train_loss)

  valid_loss=[]
  model.eval()
  with torch.no_grad():
    for batch in validation_dataloader:

      # forward pass
      output= model(batch[0].to(device))
      loss=criterion(output,batch[2].to(device))

      valid_loss.append(loss.item())

    valid_loss=np.mean(valid_loss)

  # save Losses
  train_losses[epoch]= train_loss
  valid_losses[epoch]= valid_loss
  dt= datetime.now()-t0
  print(f'Epoch {epoch+1}/{epochs}, Train Loss: {train_loss:.4f}    Valid Loss: {
```

```
    Epoch 1/10, Train Loss: 3.5173     Valid Loss: 2.8577, Duration: 0:00:42.563540
    Epoch 2/10, Train Loss: 3.3863     Valid Loss: 2.7111, Duration: 0:00:43.053556
    Epoch 3/10, Train Loss: 3.3450     Valid Loss: 2.7488, Duration: 0:00:43.140066
    Epoch 4/10, Train Loss: 3.3610     Valid Loss: 2.7833, Duration: 0:00:43.226663
    Epoch 5/10, Train Loss: 3.3945     Valid Loss: 2.8700, Duration: 0:00:43.207317
    Epoch 6/10, Train Loss: 3.3625     Valid Loss: 2.8428, Duration: 0:00:43.233899
    Epoch 7/10, Train Loss: 3.3979     Valid Loss: 2.7324, Duration: 0:00:43.238347
    Epoch 8/10, Train Loss: 3.4011     Valid Loss: 2.6344, Duration: 0:00:43.202970
    Epoch 9/10, Train Loss: 3.4075     Valid Loss: 2.7247, Duration: 0:00:43.184335
    Epoch 10/10, Train Loss: 3.4261     Valid Loss: 2.6287, Duration: 0:00:43.196169
```

```
  # Accuracy- write a function to get accuracy
  # use this function to get accuracy and print accuracy
  def get_accuracy(data_iter, model):
    model.eval()
    with torch.no_grad():
      correct =0
      total =0

      for batch in data_iter:

        output=model(batch[0].to(device))
        _,indices = torch.max(output,dim=1)
```

```
        correct+= (batch[2].to(device)==indices).sum().item()
        total += batch[2].shape[0]

    acc= correct/total

    return acc


  train_acc = get_accuracy(train_dataloader, model)
  valid_acc = get_accuracy(validation_dataloader, model)
  test_acc = get_accuracy(test_dataloader ,model)
  print(f'Train acc: {train_acc:.4f},\t Valid acc: {valid_acc:.4f},\t Test acc: {te
```

```
      Train acc: 0.3389,        Valid acc: 0.2686,        Test acc: 0.2593
```

```
  # Write a function to get predictions

  def get_predictions(test_iter, model):
    model.eval()
    with torch.no_grad():
      predictions= np.array([])
      y_test= np.array([])

      for batch in test_iter:

        output=model(batch[0].to(device))
        _,indices = torch.max(output,dim=1)
        predictions=np.concatenate((predictions,indices.cpu().numpy()))
        y_test = np.concatenate((y_test,batch[2].numpy()))

    return y_test, predictions


  y_test, predictions=get_predictions(test_dataloader, model)


  # Confusion Matrix
  cm=confusion_matrix(y_test,predictions)
  cm
```

```
      array([[ 19,    3,   10,    3,   42,   13,   10,   20,   10,    0,    9,   22,    2,
              29,    2,   54,   34,   15,   13,    9],
             [  0,   26,   64,   29,   58,   60,   42,   13,    4,    0,    0,   21,   30,
              11,    5,    6,    9,    6,    3,    2],
             [  1,   14,   84,   50,   63,   66,   21,   21,    8,    2,    0,   19,   15,
               5,    2,    6,    9,    2,    5,    1],
             [  0,    5,   52,   97,   73,   22,   45,   14,    7,    1,    2,   15,   32,
               2,    1,    2,   11,    0,    4,    7],
             [  0,    5,   51,   73,   98,   17,   48,   20,    0,    1,    1,   13,   24,
              12,    3,    0,   13,    0,    5,    1],
             [  1,   18,   62,   26,   48,  133,   33,    6,    5,    0,    4,   22,   14,
               2,    1,    3,    7,    3,    6,    1],
             [  0,    4,   19,   25,   30,    9,  253,   15,    3,    1,    5,    3,   11,
               3,    0,    1,    5,    0,    2,    1],
             [  2,    4,   27,   16,   42,    9,   43,  121,   39,    3,    8,    5,   23,
              13,    0,    2,   31,    3,    4,    1],
             [  1,    5,   27,   20,   63,   12,   42,   48,   81,    5,    3,   13,   10,
```

```
                16,    2,    4,   33,    0,   12,    1],
            [   1,    1,   27,    5,   36,   13,   34,   31,   29,   66,   38,   18,    6,
                19,    3,   12,   44,    3,    7,    4],
            [   1,    2,   22,    2,   41,   10,   27,   31,   25,   20,  129,    7,    7,
                17,    2,    4,   35,    9,    6,    2],
            [   1,    6,   23,   13,   55,   29,   15,   20,    8,    5,    3,  123,   10,
                16,    1,   13,   31,   12,   11,    1],
            [   1,    6,   46,   24,   68,   19,   35,   31,    8,    2,    4,   30,   75,
                21,    3,    1,   14,    0,    5,    0],
            [   3,   11,   13,    6,   32,    8,   18,   28,   32,    1,    3,   10,    9,
               152,    2,   22,   31,    5,    6,    4],
            [   1,    8,   22,    9,   58,   10,   25,   28,   15,    4,    4,   35,   33,
                34,   58,   11,   27,    5,    5,    2],
            [   6,    0,   10,    6,   42,    9,    5,   18,   10,    1,    5,   15,    4,
                30,    1,  176,   35,    7,    5,   13],
            [   1,    3,   19,    5,   47,    9,   15,   25,   23,    3,    5,   25,    6,
                34,    1,    9,  105,   15,    8,    6],
            [   2,    4,   11,    6,   32,    5,   10,   15,   14,    5,    3,   14,    4,
                31,    3,   26,   53,  120,    5,   13],
            [   5,    0,   14,    3,   22,   10,    9,   18,   11,    3,    4,   22,    2,
                28,    1,    7,   99,   22,   23,    7],
            [   5,    4,   17,    5,   23,   11,    9,   14,   12,    3,    1,   11,    2,
                18,    2,   51,   33,    7,    9,   14]])


# Write a function to print confusion matrix
# plot confusion matrix
# need to import confusion_matrix from sklearn for this function to work
# need to import seaborn as sns
# import seaborn as sns
# import matplotlib.pyplot as plt
# from sklearn.metrics import confusion_matrix

def plot_confusion_matrix(y_true,y_pred,normalize=None):
  cm=confusion_matrix(y_true,y_pred,normalize=normalize)
  fig, ax = plt.subplots(figsize=(6,5))
  if normalize == None:
    fmt='d'
    fig.suptitle('Confusion matrix without Normalization', fontsize=12)

  else :
    fmt='0.2f'
    fig.suptitle('Normalized confusion matrix', fontsize=12)

  ax=sns.heatmap(cm,cmap=plt.cm.Blues,annot=True,fmt=fmt)
  ax.axhline(y=0, color='k',linewidth=1)
  ax.axhline(y=cm.shape[1], color='k',linewidth=2)
  ax.axvline(x=0, color='k',linewidth=1)
  ax.axvline(x=cm.shape[0], color='k',linewidth=2)

  ax.set_xlabel('Predicted label', fontsize=12)
  ax.set_ylabel('True label', fontsize=12)


plot_confusion_matrix(y_test,predictions)
```
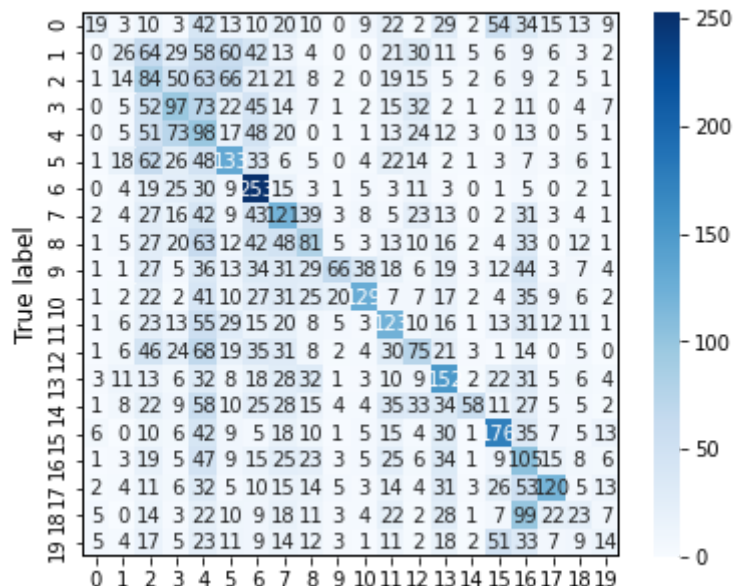
Confusion matrix without Normalization



## Finetuning the pre-trained XLNet Model

```
from transformers import XLNetForSequenceClassification, AdamW, XLNetConfig

# Load BertForSequenceClassification, the pretrained BERT model with a single
# linear classification layer on top.
model = XLNetForSequenceClassification.from_pretrained(
    "xlnet-base-cased", # Use the 12-layer Robera model, with an uncased vocab.
    num_labels = 20, # The number of output labels
    output_attentions = False, # Whether the model returns attentions weights.
    output_hidden_states = False, # Whether the model returns all hidden-states.
)
```

```
    Some weights of the model checkpoint at xlnet-base-cased were not used when initializ
    - This IS expected if you are initializing XLNetForSequenceClassification from the ch
    - This IS NOT expected if you are initializing XLNetForSequenceClassification from th
    Some weights of XLNetForSequenceClassification were not initialized from the model ch
    You should probably TRAIN this model on a down-stream task to be able to use it for p
```

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

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

```
# Tell pytorch to run this model on the GPU.
model.to(device)
```

```
          (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)
          )
          (dropout): Dropout(p=0.1, inplace=False)
        )
        (9): 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)
          )
          (dropout): Dropout(p=0.1, inplace=False)
        )
        (10): 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)
          )
          (dropout): Dropout(p=0.1, inplace=False)
        )
        (11): 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)
          )
          (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)
      (first_dropout): Identity()
      (last_dropout): Dropout(p=0.1, inplace=False)
    )
    (logits_proj): Linear(in_features=768, out_features=20, bias=True)
  )
```

Just for curiosity's sake, we can browse all of the model's parameters by name here.

In the below cell, I've printed out the names and dimensions of the weights for:

1. The embedding layer.
2. The first of the twelve transformers.
3. The output layer.

```
# Get all of the model's parameters as a list of tuples.
params = list(model.named_parameters())

print('The BERT model has {:} different named parameters.\n'.format(len(params)))

print('==== Embedding Layer ====\n')

for p in params[0:5]:
    print("{:<55} {:>12}".format(p[0], str(tuple(p[1].size()))))

print('\n==== First Transformer ====\n')

for p in params[5:21]:
    print("{:<55} {:>12}".format(p[0], str(tuple(p[1].size()))))

print('\n==== Output Layer ====\n')

for p in params[-4:]:
    print("{:<55} {:>12}".format(p[0], str(tuple(p[1].size()))))
```

```
The BERT model has 210 different named parameters.

==== Embedding Layer ====

transformer.mask_emb                                         (1, 1, 768)
transformer.word_embedding.weight                          (32000, 768)
transformer.layer.0.rel_attn.q                              (768, 12, 64)
transformer.layer.0.rel_attn.k                              (768, 12, 64)
transformer.layer.0.rel_attn.v                              (768, 12, 64)

==== First Transformer ====

transformer.layer.0.rel_attn.o                              (768, 12, 64)
transformer.layer.0.rel_attn.r                              (768, 12, 64)
transformer.layer.0.rel_attn.r_r_bias                           (12, 64)
transformer.layer.0.rel_attn.r_s_bias                           (12, 64)
transformer.layer.0.rel_attn.r_w_bias                           (12, 64)
transformer.layer.0.rel_attn.seg_embed                       (2, 12, 64)
transformer.layer.0.rel_attn.layer_norm.weight                    (768,)
transformer.layer.0.rel_attn.layer_norm.bias                      (768,)
transformer.layer.0.ff.layer_norm.weight                          (768,)
transformer.layer.0.ff.layer_norm.bias                            (768,)
transformer.layer.0.ff.layer_1.weight                         (3072, 768)
transformer.layer.0.ff.layer_1.bias                              (3072,)
transformer.layer.0.ff.layer_2.weight                         (768, 3072)
transformer.layer.0.ff.layer_2.bias                               (768,)
transformer.layer.1.rel_attn.q                              (768, 12, 64)
transformer.layer.1.rel_attn.k                              (768, 12, 64)

==== Output Layer ====

sequence_summary.summary.weight                              (768, 768)
```

```
sequence_summary.summary.bias                              (768,)
logits_proj.weight                                       (20, 768)
logits_proj.bias                                          (20,)
```

## ▾ 4.2. Optimizer & Learning Rate Scheduler

Now that we have our model loaded we need to grab the training hyperparameters from within the stored model.

For the purposes of fine-tuning, the authors recommend choosing from the following values (from Appendix A.3 of the [BERT paper](#)):

> - **Batch size:** 16, 32
> - **Learning rate (Adam):** 5e-5, 3e-5, 2e-5
> - **Number of epochs:** 2, 3, 4

We chose:

- Batch size: 32 (set when creating our DataLoaders)
- Learning rate: 2e-5
- Epochs: 4 (we'll see that this is probably too many...)

The epsilon parameter `eps = 1e-8` is "a very small number to prevent any division by zero in the implementation" (from [here](#)).

You can find the creation of the AdamW optimizer in `run_glue.py` [here](#).

## ▾ 4.3. Training Loop

Define a helper function for calculating accuracy.

Helper function for formatting elapsed times as `hh:mm:ss`

We're ready to kick off the training!

```
# STEP 6: INSTANTIATE OPTIMIZER CLASS
epochs = 2
no_decay = ['bias', 'LayerNorm.weight']
optimizer_grouped_parameters = [
        {'params': [p for n, p in model.named_parameters()
          if not any(nd in n for nd in no_decay)],
         'weight_decay': 0.5},

        {'params': [p for n, p in model.named_parameters()
         if any(nd in n for nd in no_decay)],
          'weight_decay': 0.0}
```

https://colab.research.google.com/drive/182n0Fr1iBf9ELebBFzlbYUOSQ665DM6F?authuser=2#scrollTo=eHqvmICSfnPZ&printMode=true          19/25

```python
        ]

optimizer = AdamW(optimizer_grouped_parameters,
                  lr = 5e-5,
                  eps = 1e-8
                )

no_decay = ['bias', 'LayerNorm.weight']



from transformers import get_linear_schedule_with_warmup

# Total number of training steps is [number of batches] x [number of epochs].
total_steps = len(train_dataloader) * epochs

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


import random
from datetime import datetime


seed = 123

random.seed(seed)
np.random.seed(seed)
torch.manual_seed(seed)
torch.cuda.manual_seed_all(seed)


epochs = 2

# STEP 7: TRAIN THE MODEL

train_losses= np.zeros(epochs)
valid_losses= np.zeros(epochs)


for epoch in range(epochs):

  t0= datetime.now()
  train_loss=[]

  model.train()
  for batch in train_dataloader:
    b_input_ids = batch[0]
    b_input_mask = batch[1]
    b_labels = batch[2]
    # forward pass

    outputs = model(b_input_ids.to(device),
                    token_type_ids=None,
                    attention_mask=b_input_mask.to(device),
```

```
                       labels=b_labels.to(device))


    # set gradients to zero
    optimizer.zero_grad()
    # backward pass
    outputs.loss.backward()
    torch.nn.utils.clip_grad_norm_(model.parameters(), 1.0)
    optimizer.step()
    scheduler.step()
    train_loss.append(outputs.loss.item())

  train_loss=np.mean(train_loss)

  valid_loss=[]
  model.eval()
  with torch.no_grad():
    for batch in validation_dataloader:

      # forward pass
      b_input_ids = batch[0].to(device)
      b_input_mask = batch[1].to(device)
      b_labels = batch[2].to(device)
    # forward pass

      outputs = model(b_input_ids,
                      token_type_ids=None,
                      attention_mask=b_input_mask,
                      labels=b_labels)

      valid_loss.append(outputs.loss.item())

    valid_loss=np.mean(valid_loss)

  # save Losses
  train_losses[epoch]= train_loss
  valid_losses[epoch]= valid_loss
  dt= datetime.now()-t0
  print(f'Epoch {epoch+1}/{epochs}, Train Loss: {train_loss:.4f}    Valid Loss: {
```

```
    Epoch 1/2, Train Loss: 1.4028    Valid Loss: 1.1523, Duration: 0:02:25.167743
    Epoch 2/2, Train Loss: 0.7679    Valid Loss: 1.0338, Duration: 0:02:25.047346
```

```
# Accuracy- write a function to get accuracy
# use this function to get accuracy and print accuracy
def get_accuracy(data_iter, model):
  model.eval()
  with torch.no_grad():
    correct =0
    total =0

    for batch in data_iter:

      b_input_ids = batch[0].to(device)
      b_input_mask = batch[1].to(device)
```

```python
        b_labels = batch[2].to(device)
      # forward pass

      outputs = model(b_input_ids,
                          token_type_ids=None,
                          attention_mask=b_input_mask,
                          labels=b_labels)


      _,indices = torch.max(outputs.logits,dim=1)
      correct+= (b_labels==indices).sum().item()
      total += b_labels.shape[0]

    acc= correct/total

    return acc


train_acc = get_accuracy(train_dataloader, model)
valid_acc = get_accuracy(validation_dataloader, model)
test_acc = get_accuracy(test_dataloader, model)


print(f'Train acc: {train_acc:.4f},\t Valid acc: {valid_acc:.4f},\t Test acc: {te
```

```
    Train acc: 0.8517,       Valid acc: 0.6935,       Test acc: 0.6771
```

```python
# Write a function to get predictions
def get_predictions(data_iter, model):
  model.eval()
  with torch.no_grad():
    predictions= np.array([])
    y_test= np.array([])

    for batch in data_iter:

      b_input_ids = batch[0].to(device)
      b_input_mask = batch[1].to(device)
      b_labels = batch[2].to(device)
    # forward pass

      outputs = model(b_input_ids,
                          token_type_ids=None,
                          attention_mask=b_input_mask,
                          labels=b_labels)


      _,indices = torch.max(outputs.logits,dim=1)
      predictions=np.concatenate((predictions,indices.cpu().numpy()))
      y_test = np.concatenate((y_test,b_labels.cpu().numpy()))

  return y_test, predictions


y_valid, predictions=get_predictions(validation_dataloader, model)


predictions.max()
```

```
     19.0
```

```
# Confusion Matrix
from sklearn.metrics import confusion_matrix
cm=confusion_matrix(y_valid,predictions)
cm
```

```
     array([[32,  0,  0,  0,  0,  0,  0,  4,  0,  0,  1,  1,  0,  0,  1,  7,
              0,  0,  1,  1],
            [ 0, 36,  7,  1,  0,  4,  0,  0,  0,  1,  0,  0,  2,  0,  3,  0,
              0,  1,  0,  0],
            [ 0,  1, 33,  6,  0,  4,  1,  2,  1,  1,  0,  2,  2,  1,  0,  1,
              0,  0,  0,  0],
            [ 0,  1,  7, 52,  3,  0,  2,  4,  0,  0,  0,  0,  4,  0,  0,  1,
              0,  0,  0,  0],
            [ 1,  3,  2,  7, 31,  0,  3,  4,  0,  0,  0,  0,  6,  0,  0,  1,
              0,  0,  0,  0],
            [ 0,  3,  6,  0,  0, 52,  0,  0,  0,  0,  0,  1,  0,  0,  0,  0,
              0,  2,  0,  1],
            [ 0,  1,  1,  1,  0,  0, 49,  5,  1,  0,  1,  0,  0,  0,  0,  0,
              0,  0,  0,  0],
            [ 0,  0,  0,  1,  0,  0,  0, 42,  9,  0,  0,  0,  1,  0,  0,  0,
              4,  0,  3,  0],
            [ 0,  0,  0,  0,  0,  0,  1,  7, 43,  1,  0,  1,  2,  1,  2,  0,
              3,  1,  0,  0],
            [ 0,  0,  0,  0,  0,  1,  0,  4,  1, 44,  6,  0,  1,  0,  1,  0,
              0,  1,  0,  0],
            [ 0,  0,  0,  0,  0,  1,  0,  0,  3,  3, 50,  1,  0,  0,  0,  0,
              0,  0,  0,  1],
            [ 0,  1,  1,  0,  0,  1,  0,  0,  0,  1,  0, 44,  0,  1,  1,  1,
              1,  2,  6,  0],
            [ 1,  1,  0,  3,  0,  2,  3,  0,  2,  0,  0,  1, 28,  1,  0,  0,
              1,  0,  0,  0],
            [ 0,  1,  1,  1,  0,  0,  0,  1,  1,  0,  0,  0,  0, 48,  2,  0,
              0,  0,  1,  0],
            [ 0,  1,  0,  1,  0,  0,  0,  1,  3,  2,  0,  2,  1,  1, 32,  0,
              0,  0,  0,  0],
            [ 8,  0,  0,  0,  0,  0,  0,  0,  0,  0,  1,  0,  0,  2,  0, 56,
              0,  0,  2,  2],
            [ 1,  1,  0,  0,  0,  0,  1,  1,  0,  1,  0,  1,  0,  0,  1,  0,
             39,  1,  9,  1],
            [ 2,  1,  0,  0,  0,  0,  0,  1,  1,  0,  1,  0,  0,  0,  0,  1,
              0, 41,  2,  1],
            [ 4,  0,  0,  0,  0,  0,  0,  1,  2,  0,  2,  2,  0,  2,  1,  2,
              2,  6, 30,  0],
            [11,  0,  1,  0,  0,  0,  0,  2,  1,  0,  1,  1,  0,  1,  2, 10,
              6,  1,  2,  3]])
```

```
import seaborn as sns
import matplotlib.pyplot as plt
from sklearn.metrics import confusion_matrix

def plot_confusion_matrix(y_true,y_pred,normalize=None):
  cm=confusion_matrix(y_true,y_pred,normalize=normalize)
  fig, ax = plt.subplots(figsize=(6,5))
  if normalize == None:
    fmt='d'
```

```
    fig.suptitle('Confusion matrix without Normalization', fontsize=12)

  else :
    fmt='0.2f'
    fig.suptitle('Normalized confusion matrix', fontsize=12)

  ax=sns.heatmap(cm,cmap=plt.cm.Blues,annot=True,fmt=fmt)
  ax.axhline(y=0, color='k',linewidth=1)
  ax.axhline(y=cm.shape[1], color='k',linewidth=2)
  ax.axvline(x=0, color='k',linewidth=1)
  ax.axvline(x=cm.shape[0], color='k',linewidth=2)

  ax.set_xlabel('Predicted label', fontsize=12)
  ax.set_ylabel('True label', fontsize=12)


plot_confusion_matrix(y_valid,predictions)
```

Confusion matrix without Normalization