

✓ Importing Python Libraries and preparing the environment

At this step we will be importing the libraries and modules needed to run our script. Libraries are:

- Pandas
- Pytorch
- Pytorch Utils for Dataset and Dataloader
- Transformers
- DistilBERT Model and Tokenizer

Followed by that we will preapre the device for CUDA execeution. This configuration is needed if you want to leverage on onboard GPU.

```
# Importing the libraries needed
import pandas as pd
import torch
import transformers
from torch.utils.data import Dataset, DataLoader
from transformers import DistilBertModel, DistilBertTokenizer
```

```
# Setting up the device for GPU usage
```

```
from torch import cuda
device = 'cuda' if cuda.is_available() else 'cpu'
```

✓ Importing and Pre-Processing the domain data

We will be working with the data and preparing for fine tuning purposes. *Assuming that the newCorpora.csv is already downloaded in your data folder*

Import the file in a dataframe and give it the headers as per the documentation. Cleaning the file to remove the unwanted columns and create an additional column for training. The final Dataframe will be something like this:

TITLE	CATEGORY	ENCODED_CAT
title_1	Entertainment	1
title_2	Entertainment	1
title_3	Business	2
title_4	Science	3
title_5	Science	3
title_6	Health	4

```
# Import the csv into pandas dataframe and add the headers
df = pd.read_csv('./data/newsCorpora.csv', sep='\t', names=['ID', 'TITLE', 'URL', 'PUBLISHER', 'CATEGORY', 'STORY', 'HOSTNAME', 'TIMESTAMP'])
# df.head()
# # Removing unwanted columns and only leaving title of news and the category which will be the target
df = df[['TITLE', 'CATEGORY']]
# df.head()

# # Converting the codes to appropriate categories using a dictionary
my_dict = {
    'e': 'Entertainment',
    'b': 'Business',
    't': 'Science',
    'm': 'Health'
}

def update_cat(x):
    return my_dict[x]

df['CATEGORY'] = df['CATEGORY'].apply(lambda x: update_cat(x))

encode_dict = {}

def encode_cat(x):
    if x not in encode_dict.keys():
        encode_dict[x] = len(encode_dict)
    return encode_dict[x]

df['ENCODE_CAT'] = df['CATEGORY'].apply(lambda x: encode_cat(x))
```

✓ Preparing the Dataset and Dataloader

We will start with defining few key variables that will be used later during the training/fine tuning stage. Followed by creation of Dataset class - This defines how the text is pre-processed before sending it to the neural network. We will also define the Dataloader that will feed the data in batches to the neural network for suitable training and processing. Dataset and Dataloader are constructs of the PyTorch library for defining and controlling the data pre-processing and its passage to neural network. For further reading into Dataset and Dataloader read the [docs at PyTorch](#)

Triage Dataset Class

- This class is defined to accept the Dataframe as input and generate tokenized output that is used by the DistilBERT model for training.
- We are using the DistilBERT tokenizer to tokenize the data in the `TITLE` column of the dataframe.
- The tokenizer uses the `encode_plus` method to perform tokenization and generate the necessary outputs, namely: `ids`, `attention_mask`
- To read further into the tokenizer, [refer to this document](#)
- `target` is the encoded category on the news headline.
- The *Triage* class is used to create 2 datasets, for training and for validation.
- *Training Dataset* is used to fine tune the model: **80% of the original data**
- *Validation Dataset* is used to evaluate the performance of the model. The model has not seen this data during training.

Dataloader

- Dataloader is used to for creating training and validation dataloader that load data to the neural network in a defined manner. This is needed because all the data from the dataset cannot be loaded to the memory at once, hence the amount of data loaded to the memory and then passed to the neural network needs to be controlled.
- This control is achieved using the parameters such as `batch_size` and `max_len`.
- Training and Validation dataloaders are used in the training and validation part of the flow respectively

```
# Defining some key variables that will be used later on in the training
MAX_LEN = 512
TRAIN_BATCH_SIZE = 4
VALID_BATCH_SIZE = 2
EPOCHS = 1
LEARNING_RATE = 1e-05
tokenizer = DistilBertTokenizer.from_pretrained('distilbert-base-cased')

class Triage(Dataset):
    def __init__(self, dataframe, tokenizer, max_len):
        self.len = len(dataframe)
        self.data = dataframe
        self.tokenizer = tokenizer
        self.max_len = max_len

    def __getitem__(self, index):
        title = str(self.data.TITLE[index])
        title = " ".join(title.split())
        inputs = self.tokenizer.encode_plus(
            title,
            None,
            add_special_tokens=True,
            max_length=self.max_len,
            pad_to_max_length=True,
            return_token_type_ids=True,
            truncation=True
        )
        ids = inputs['input_ids']
        mask = inputs['attention_mask']

        return {
            'ids': torch.tensor(ids, dtype=torch.long),
            'mask': torch.tensor(mask, dtype=torch.long),
            'targets': torch.tensor(self.data.ENCODE_CAT[index], dtype=torch.long)
        }

    def __len__(self):
        return self.len
```

```
# Creating the dataset and dataloader for the neural network

train_size = 0.8
train_dataset=df.sample(frac=train_size,random_state=200)
test_dataset=df.drop(train_dataset.index).reset_index(drop=True)
train_dataset = train_dataset.reset_index(drop=True)

print("FULL Dataset: {}".format(df.shape))
print("TRAIN Dataset: {}".format(train_dataset.shape))
print("TEST Dataset: {}".format(test_dataset.shape))

training_set = Triage(train_dataset, tokenizer, MAX_LEN)
testing_set = Triage(test_dataset, tokenizer, MAX_LEN)

FULL Dataset: (422419, 3)
TRAIN Dataset: (337935, 3)
TEST Dataset: (84484, 3)

train_params = {'batch_size': TRAIN_BATCH_SIZE,
                'shuffle': True,
                'num_workers': 0
               }

test_params = {'batch_size': VALID_BATCH_SIZE,
               'shuffle': True,
               'num_workers': 0
              }

training_loader = DataLoader(training_set, **train_params)
testing_loader = DataLoader(testing_set, **test_params)
```

✓ Creating the Neural Network for Fine Tuning

Neural Network

- We will be creating a neural network with the `DistilBERTClass`.
- This network will have the DistilBERT Language model followed by a `dropout` and finally a `Linear` layer to obtain the final outputs.
- The data will be fed to the DistilBERT Language model as defined in the dataset.
- Final layer outputs is what will be compared to the `encoded category` to determine the accuracy of models prediction.
- We will initiate an instance of the network called `model`. This instance will be used for training and then to save the final trained model for future inference.

Loss Function and Optimizer

- `Loss Function` and `Optimizer` and defined in the next cell.
- The `Loss Function` is used to calculate the difference in the output created by the model and the actual output.
- `Optimizer` is used to update the weights of the neural network to improve its performance.

Further Reading

- You can refer to my [Pytorch Tutorials](#) to get an intuition of Loss Function and Optimizer.
- [Pytorch Documentation for Loss Function](#)
- [Pytorch Documentation for Optimizer](#)
- Refer to the links provided on the top of the notebook to read more about DistilBERT.

Creating the customized model, by adding a drop out and a dense layer on top of distil bert to get the final output for the model.

```
class DistilBERTClass(torch.nn.Module):
    def __init__(self):
        super(DistilBERTClass, self).__init__()
        self.l1 = DistilBertModel.from_pretrained("distilbert-base-uncased")
        self.pre_classifier = torch.nn.Linear(768, 768)
        self.dropout = torch.nn.Dropout(0.3)
        self.classifier = torch.nn.Linear(768, 4)

    def forward(self, input_ids, attention_mask):
        output_1 = self.l1(input_ids=input_ids, attention_mask=attention_mask)
        hidden_state = output_1[0]
        pooler = hidden_state[:, 0]
        pooler = self.pre_classifier(pooler)
        pooler = torch.nn.ReLU()(pooler)
        pooler = self.dropout(pooler)
        output = self.classifier(pooler)
        return output
```

```
model = DistillBERTClass()
model.to(device)
```

[Show hidden output](#)

```
# Creating the loss function and optimizer
loss_function = torch.nn.CrossEntropyLoss()
optimizer = torch.optim.Adam(params = model.parameters(), lr=LEARNING_RATE)
```

✓ Fine Tuning the Model

After all the effort of loading and preparing the data and datasets, creating the model and defining its loss and optimizer. This is probably the easier steps in the process.

Here we define a training function that trains the model on the training dataset created above, specified number of times (EPOCH), An epoch defines how many times the complete data will be passed through the network.

Following events happen in this function to fine tune the neural network:

- The dataloader passes data to the model based on the batch size.
- Subsequent output from the model and the actual category are compared to calculate the loss.
- Loss value is used to optimize the weights of the neurons in the network.
- After every 5000 steps the loss value is printed in the console.

As you can see just in 1 epoch by the final step the model was working with a miniscule loss of 0.0002485 i.e. the output is extremely close to the actual output.

```
# Function to calculate the accuracy of the model

def calculate_accu(big_idx, targets):
    n_correct = (big_idx==targets).sum().item()
    return n_correct

# Defining the training function on the 80% of the dataset for tuning the distilbert model

def train(epoch):
    tr_loss = 0
    n_correct = 0
    nb_tr_steps = 0
    nb_tr_examples = 0
    model.train()
    for _,data in enumerate(training_loader, 0):
        ids = data['ids'].to(device, dtype = torch.long)
        mask = data['mask'].to(device, dtype = torch.long)
        targets = data['targets'].to(device, dtype = torch.long)

        outputs = model(ids, mask)
        loss = loss_function(outputs, targets)
        tr_loss += loss.item()
        big_val, big_idx = torch.max(outputs.data, dim=1)
        n_correct += calculate_accu(big_idx, targets)

        nb_tr_steps += 1
        nb_tr_examples+=targets.size(0)

        if _%5000==0:
            loss_step = tr_loss/nb_tr_steps
            accu_step = (n_correct*100)/nb_tr_examples
            print(f"Training Loss per 5000 steps: {loss_step}")
            print(f"Training Accuracy per 5000 steps: {accu_step}")

            optimizer.zero_grad()
            loss.backward()
            ## When using GPU
            optimizer.step()

    print(f'The Total Accuracy for Epoch {epoch}: {(n_correct*100)/nb_tr_examples}')
    epoch_loss = tr_loss/nb_tr_steps
    epoch_accu = (n_correct*100)/nb_tr_examples
    print(f"Training Loss Epoch: {epoch_loss}")
    print(f"Training Accuracy Epoch: {epoch_accu}")

    return

for epoch in range(EPOCHS):
    train(epoch)
```

```

Epoch: 0, Loss: 6.332988739013672
Epoch: 0, Loss: 0.0013066530227661133
Epoch: 0, Loss: 0.0029534101486206055
Epoch: 0, Loss: 0.005258679389953613
Epoch: 0, Loss: 0.0020235776901245117
Epoch: 0, Loss: 0.0023298263549804688
Epoch: 0, Loss: 0.0034378767013549805
Epoch: 0, Loss: 0.004993081092834473
Epoch: 0, Loss: 0.008559942245483398
Epoch: 0, Loss: 0.0014510154724121094
Epoch: 0, Loss: 0.0028634071350097656
Epoch: 0, Loss: 0.0006411075592041016
Epoch: 0, Loss: 0.0012137889862060547
Epoch: 0, Loss: 0.002307891845703125
Epoch: 0, Loss: 0.00028586387634277344
Epoch: 0, Loss: 0.0029143095016479492
Epoch: 0, Loss: 0.0002485513687133789

```

✓ Validating the Model

During the validation stage we pass the unseen data (Testing Dataset) to the model. This step determines how good the model performs on the unseen data.

This unseen data is the 20% of `newscorpora.csv` which was separated during the Dataset creation stage. During the validation stage the weights of the model are not updated. Only the final output is compared to the actual value. This comparison is then used to calculate the accuracy of the model.

As you can see the model is predicting the correct category of a given headline to a 99.9% accuracy.

```

def valid(model, testing_loader):
    model.eval()
    n_correct = 0; n_wrong = 0; total = 0
    with torch.no_grad():
        for _, data in enumerate(testing_loader, 0):
            ids = data['ids'].to(device, dtype = torch.long)
            mask = data['mask'].to(device, dtype = torch.long)
            targets = data['targets'].to(device, dtype = torch.long)
            outputs = model(ids, mask).squeeze()
            loss = loss_function(outputs, targets)
            tr_loss += loss.item()
            big_val, big_idx = torch.max(outputs.data, dim=1)
            n_correct += calculate_acc(big_idx, targets)

        nb_tr_steps += 1
        nb_tr_examples += targets.size(0)

        if _%5000==0:
            loss_step = tr_loss/nb_tr_steps
            accu_step = (n_correct*100)/nb_tr_examples
            print(f"Validation Loss per 100 steps: {loss_step}")
            print(f"Validation Accuracy per 100 steps: {accu_step}")
    epoch_loss = tr_loss/nb_tr_steps
    epoch_accu = (n_correct*100)/nb_tr_examples
    print(f"Validation Loss Epoch: {epoch_loss}")
    print(f"Validation Accuracy Epoch: {epoch_accu}")

    return epoch_accu

print('This is the validation section to print the accuracy and see how it performs')
print('Here we are leveraging on the dataloader created for the validation dataset, the approach is using more of pytorch')

acc = valid(model, testing_loader)
print("Accuracy on test data = %.2f%%" % acc)

```

```

This is the validation section to print the accuracy and see how it performs
Here we are leveraging on the dataloader created for the validation dataset, the approach is using more of pytorch
Accuracy on test data = 99.99%

```

✓ Saving the Trained Model Artifacts for inference

This is the final step in the process of fine tuning the model.

The model and its vocabulary are saved locally. These files are then used in the future to make inference on new inputs of news headlines.

Please remember that a trained neural network is only useful when used in actual inference after its training.

In the lifecycle of an ML projects this is only half the job done. We will leave the inference of these models for some other day.

```
# Saving the files for re-use
```

```
output_model_file = './models/pytorch_distilbert_news.bin'  
output_vocab_file = './models/vocab_distilbert_news.bin'
```

```
model_to_save = model  
torch.save(model_to_save, output_model_file)  
tokenizer.save_vocabulary(output_vocab_file)
```

```
print('All files saved')  
print('This tutorial is completed')
```

```
All files saved  
This tutorial is completed
```