

Data Mining HW2 Report

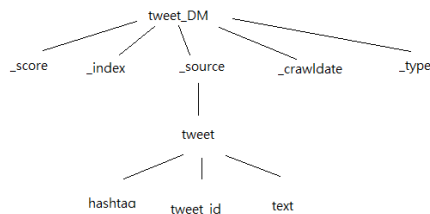
In this homework, I chose to fine-tune pre-trained BERT model for the competition. BERT is a transformer-based machine learning model introduced By Google. It's designed for NLP tasks and has become a foundational model for applications like text classification, question answering, and sentiment analysis.

The following are the steps I walked through:

First comes with some preprocessing step, I use google colab environment to run my code , and I have upload data of competition to my google drive then mount my drive to access these data.

```
1. from google.colab import drive
2. drive.mount('/content/drive')
```

After reading all the data needed as the dataframes, I observe the structure of the json dataset. The structure is like



Thus, I extract the columns under tweet to tweet_dm

```
1. # extract the columns under "tweet" to tweets_dm for simplification
2. tweets_dm["hashtags"] = tweets_dm['_source'].apply(lambda x: x['tweet']['hashtags'])
3. tweets_dm["tweet_id"] = tweets_dm['_source'].apply(lambda x: x['tweet']['tweet_id'])
4. tweets_dm["text"] = tweets_dm['_source'].apply(lambda x: x['tweet']['text'])
```

Then, apply identification tag by merging dataframe of tweet_dm and identification.

```

1. # merge tweets_dm and identification(append identification label to tweets_dm)
2. df_merged = tweets_dm.merge(identification,how="left",on="tweet_id")

```

Now we can split the dataset in training and testing by the identification tag.

```

1. # split the dataset
2. df_train = df_merged[df_merged["identification"]=="train"]
3. df_test = df_merged[df_merged["identification"]=="test"]
4. print(df_train.shape)
5. print(df_test.shape)

```

And apply emotion tag on the training set by merging dataframe of df_train and emotion

```

1. df_train = df_train.merge(emotion,how="left",on="tweet_id")

```

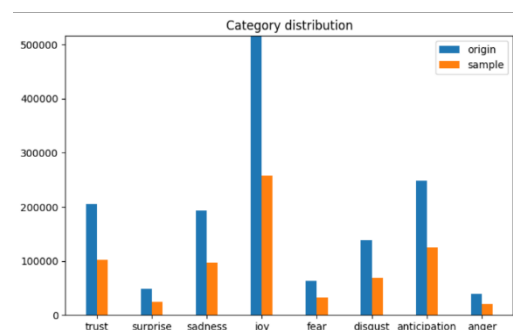
Here, I sample 50% of the training data because it's too big for my computation resource.

```

1. #sample half of data
2. df_train = df_train.sample(frac=0.5, random_state=42)

```

Plot the distribution of all categories before and after sampling.



The hyperparameter I define for machine learning(but only train for one epoch at the end)

```

1. # define the hyperparameters
2. lr = 3e-5
3. epochs = 3
4. batch_size = 64

```

There are 8 categories in total and I also use 2 dictionaries for the conversion between emotion(e.x. “fear”) and label(e.x. 3) since we need numeric labels for machine learning.

```

1. category = ["anger", "anticipation", "disgust", "fear", "sadness", "surprise", "trust", "joy"]
2. emotion_to_label = {c: i for i, c in enumerate(category)}
3. label_to_emotion = {i: c for i, c in enumerate(category)}

```

I load the bert-base-uncased tokenizer here for tokenization of text.

```

1. tokenizer = T.BertTokenizer.from_pretrained("google-bert/bert-base-uncased", cache_dir="./cache/")

```

Definition of dataset and collate function here.They’re help for machine learning by stack a batch of data and package them to pytensor type.Also, I do the tokenization during collate function.

Another worth noting thing is we only use three columns of the data.They’re tweet_id(identify a data point), text(token for input data), and emotion(label for prediction)

```

1. import torch
2. from torch.utils.data import Dataset, DataLoader
3.
4. # define dataset for machine learning
5. class Dataset(Dataset):
6.     def __init__(self, dataframe):
7.         self.dataframe = dataframe

```

```

8.
9.     def __len__(self):
10.         return len(self.dataframe)
11.
12.     def __getitem__(self, index):
13.         id = self.dataframe.iloc[index,0]
14.         x = self.dataframe.iloc[index,1]
15.         y = self.dataframe.iloc[index,2]
16.         return id, x, y
17.
18. # define collate function for machine learning
19. def collate_fn(batch):
20.     id = [data[0] for data in batch]
21.     # tokenize and package into PyTorch
22.     sentence = [data[1] for data in batch]
23.     token = tokenizer(sentence, truncation=True, padding=True, return_tensors="pt").to(device)
24.     # convert emotions to labels here also package into PyTorch
25.     label = torch.tensor([emotion_to_label[data[2]] for data in batch]).to(device)
26.
27.     return id, token, label
28.
29. # build dataset of training set
30. ds_train = Dataset(df_train[["tweet_id", "text", "emotion"]])
31.
32. # build dataloader of training set
33. dl_train = DataLoader(ds_train, batch_size=batch_size, shuffle=True, collate_fn=collate_fn)

```

Definition of the model, I use the bert-base-uncased pre-trained model for this task and add a linear layer for classification.

During forward, model take the tokens as input, and we use pooler_output to feed in linear layer then produce prediction(score for every categories)

```

1. # construct the model
2. class MyModel(torch.nn.Module):

```

```

3.     def __init__(self, *args, **kwargs):
4.         super().__init__(*args, **kwargs)
5.         # define the modules used in my model
6.         # use pretrained BERT model here
7.         self.bert = T.BertModel.from_pretrained("google-bert/bert-base-uncased", cache_
            dir="./cache/")
8.         hidden_size = self.bert.config.hidden_size
9.         # apply a linear layer for classification
10.        self.classifier = torch.nn.Linear(hidden_size, len(category))
11.
12.    def forward(self, **kwargs):
13.        # forward pass
14.        outputs = self.bert(input_ids = kwargs["input_ids"], attention_mask=kwargs["att
            ention_mask"])
15.        # use pooled_output for feature of whole sentence instead of single word
16.        pooled_output = outputs.pooler_output
17.
18.        # predictions
19.        classification_output = self.classifier(pooled_output)
20.
21.        return classification_output
22.
23. model = MyModel().to(device)

```

I use Adam as optimizer and cross entropy as loss function.

```

1.  # define the optimizer
2.  optimizer = torch.optim.Adam(model.parameters(), lr=1e-3)
3.
4.  # define the loss functions
5.  criterion_classification = torch.nn.CrossEntropyLoss() # for classification

```

It takes about 2 hours for training one epoch using sampled data(50% of training set). Since colab has limited free GPU resource, final model is trained for one epoch.

For testing, I apply a dummy column for df_test dataframe to pass through the dataset and dataloader.

```
1. df_test["dummy"] = "fear"
```

Argmax() will find the index of max value in the prediction of all categories, and I convert them to origin emotion by label_to_emotion dictionary. Then write the result in csv file.

```
1. import csv
2.
3. with open('output.csv', 'w', newline='') as csvfile:
4.
5.     writer = csv.writer(csvfile)
6.
7.     writer.writerow(["id", "emotion"])
8.
9.     pbar = tqdm(dl_test)
10.    eval_model.eval()
11.    with torch.no_grad():
12.        for ids, batch_tokens, batch_labels in pbar:
13.            # predict
14.            batch_preds = eval_model(**batch_tokens)
15.
16.            # convert label to emotion
17.            label_preds = torch.argmax(batch_preds, dim=1)
18.
19.            for item in zip(ids, label_preds):
20.                # write output csv file
21.                id = item[0]
22.                emotion = label_to_emotion[item[1].item()]
23.
24.                writer.writerow([id, emotion])
```