

# **MI201 – Project Report**

## Development of Text Emotion Classifiers using Machine Learning Techniques

Gabriel DUPUIS, Louis MARCHAL

February 7, 2025

**ENSTA**



**IP PARIS**

# Contents

Introduction . . . . .	2
Data Analysis . . . . .	3
Classical Machine Learning Predictors . . . . .	18
MLP for Prediction . . . . .	21
Performance Analysis . . . . .	31
Comparison of Performance with an LLM . . . . .	33
Justification of the Embeddings . . . . .	35
How BERT Works . . . . .	36
Fine-tuning for a Classifier . . . . .	38

## Introduction

In the digital age, Internet moderation has become crucial to ensure a safe and respectful communication space, helping to limit the spread of abusive or inappropriate content.

This project therefore aims to develop a sentiment classification system based on text embeddings, in the context of content moderation on Twitter.

By leveraging the contextual representations generated by pre-trained models such as BERT, we seek to efficiently identify the tone of different messages (positive, neutral, or negative).

The main challenge is to design a moderation bot capable of automatically detecting and flagging potentially problematic messages, relying on classical machine learning techniques as well as more complex architectures like multi-layer perceptron (MLP) neural networks or fine-tuning of language models (LLM).

# Q0: Data Analysis using Classical Unsupervised Machine Learning Methods

In this section, we aim to understand the data we are working with and analyze it using various classical unsupervised machine learning methods (PCA, TSNE, KMeans).

## Data Format

```
1 train_dataset = path+'/train.csv'  
2 test_dataset = path+'/test.csv'  
3  
4 # Check if the path exists  
5 print (os.path.exists(train_dataset))  
6 print (os.path.exists(test_dataset))  
7  
8 # Load the CSV file into a DataFrame  
9 train_df = pd.read_csv(train_dataset, encoding='ISO-8859-1')  
10 test_df = pd.read_csv(test_dataset, encoding='ISO-8859-1')
```

First, we examine what the texts we will work on look like and their distribution by class. We have chosen texts from Twitter to work on real messages.

```
1 # Distribution of sentiment labels  
2 sns.countplot(x='sentiment', data=train_df)  
3 plt.title("Distribution of Sentiment Labels")  
4 plt.show()
```

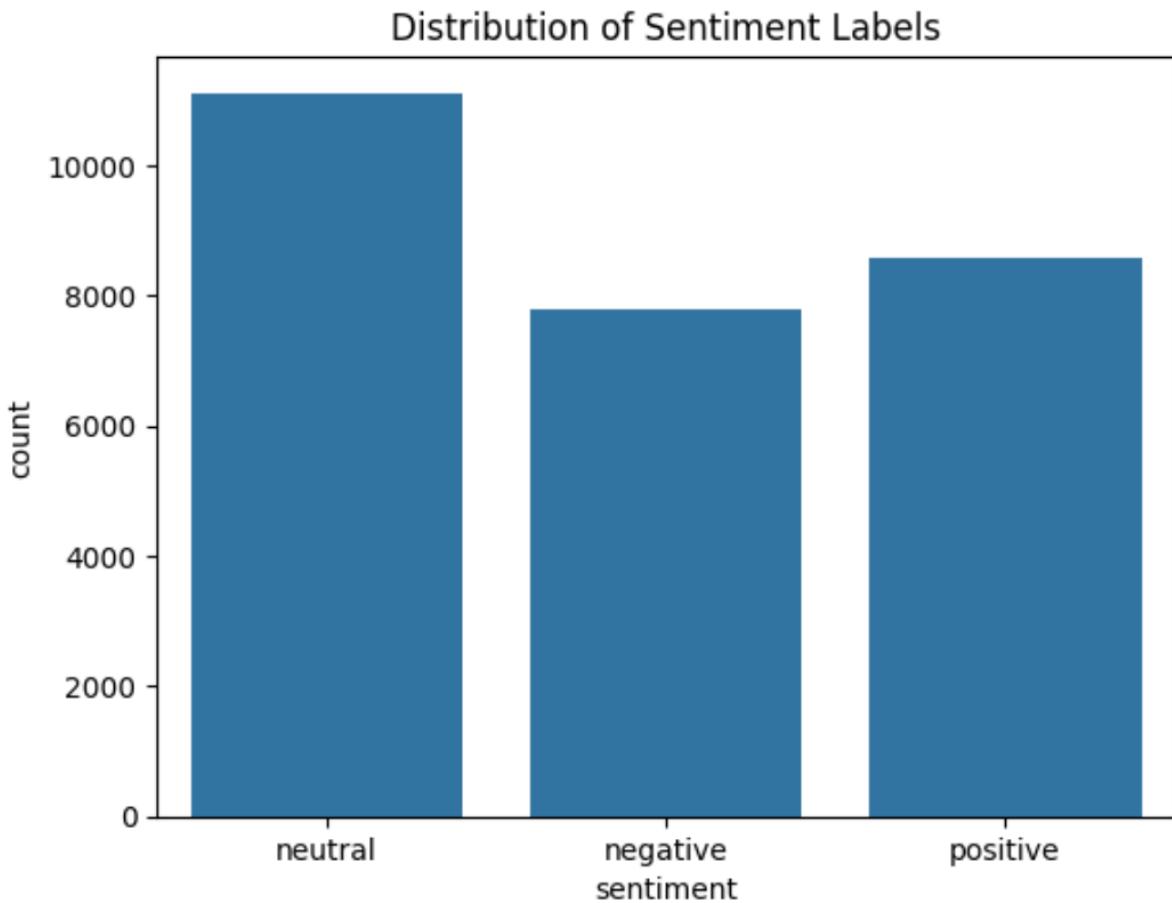


Figure 1: Histogram by Class

We apply preprocessing on the text to remove all stop-words, leaving only the words that are of interest.

**Note:** Since we are working with texts from messages that are not in perfect English, and which include emojis and special characters, this step is very important to isolate the truly important content and improve the quality of our predictions.

```
1 # Text preprocessing function that removes stopwords and converts text to
2   ↵ lowercase
3
4 def preprocess_text(text):
5     # Convert to lowercase
6     text = text.lower()
7     # Remove stopwords
8     stop_words = set(stopwords.words('english'))
9     text = " ".join([word for word in text.split() if word not in
10       ↵ stop_words])
11
12     return text
```

```

1 # Apply preprocessing on train dataset
2 train_df['processed_text'] = train_df['text'].apply(preprocess_text)
3 train_df.head()

```

	textID	text	selected_text	sentiment	Time of Tweet	Age of User	Country	Population -2020	Land Area (Km²)	Density (P/Km²)	sentiment_class	processed_text
0	cb774db0d1	I'd have responded, if I were going	I'd have responded, if I were going	neutral	morning	0-20	Afghanistan	38928346	652860.0	60	1	i'd responded, going
1	549e992a42	Sooo SAD I will miss you here in San Diego!!!	Sooo SAD	negative	noon	21-30	Albania	2877797	27400.0	105	0	soooo sad miss san diego!!!
2	088c60f138	my boss is bullying me...	bullying me	negative	night	31-45	Algeria	43851044	2381740.0	18	0	boss bullying me...
3	9642c003ef	what interview leave me alone	leave me alone	negative	morning	46-60	Andorra	77265	470.0	164	0	interview! leave alone
4	358b9e861	Sons of ****, why couldn't they put them on t...	Sons of ****,	negative	noon	60-70	Angola	32866272	1246700.0	26	0	sons ****, couldn't put releases already bought

Figure 2: Data Preview

We observe that the preprocessing worked, and we will now work solely with the pre-processed text, disregarding the other features.

```

1 # 1. Transform sentiment into 3 classes
2 # Example mapping: positive -> 2, neutral -> 1, negative -> 0
3 sentiment_mapping = {"positive": 2, "neutral": 1, "negative": 0}
4 train_df["sentiment_class"] = train_df["sentiment"].map(sentiment_mapping)
5 test_df["sentiment_class"] = test_df["sentiment"].map(sentiment_mapping)

6
7 # 2. Extract all the values from the 'processed_text' column into a list
8 trainval_x = train_df["processed_text"].tolist()
9 trainval_y = train_df["sentiment_class"].tolist()

10
11 train_x, val_x, train_y, val_y = train_test_split(trainval_x, trainval_y,
12   ↳ test_size=0.25, random_state=42)

13 test_x = test_df["processed_text"].tolist()
14 test_y = test_df["sentiment_class"].tolist()

```

We then extract the embeddings using BERT, which we will use throughout this project.

## Visualization using PCA

We know that the dimensionality of the data is 768. However, we would like to visualize the data in 1D, 2D, or 3D.

To achieve this, we will use the PCA method (known in French as ACP), which allows us to project our data onto the principal components of the covariance matrix for the most meaningful visualization.

### Variance Analysis

First, we want to examine how the variance evolves as a function of the number of components chosen. In particular, we wish to determine the number of components needed to capture 95% of the variance.

```
1 # PCA with all components
2 pca = PCA()
3 train_embeddings_pca = pca.fit_transform(train_embeddings.numpy())
4 # Variance ratio per component
5 explained_variance_ratio = pca.explained_variance_ratio_
6
7 cumul_var = 0
8 i = 0
9 # Retrieve the components whose cumulative sum constitutes 95% of the
10 # variance
11 while (cumul_var < 0.95):
12     cumul_var += explained_variance_ratio[i]
13     i += 1
14
15 n_components = i+1
16 print("We should have {} components to capture {}% of the
17 # variance.".format(i+1, int(cumul_var*100)))
```

1 We should have 23 components to capture 95% of the variance.

We then display the raw values and the variance histogram.

```
1 pca = PCA(n_components)
2 train_embeddings_pca = pca.fit_transform(train_embeddings.numpy())
3
4 explained_variance_ratio = pca.explained_variance_ratio_
5 explained_variance = pca.explained_variance_
6 print("Explained Variance Ratio per Component:")
7 print(explained_variance_ratio)
8 print("\nExplained Variance per Component:")
9 print(explained_variance)
10 print(f"Cumulative explained variance for 3 components:
11   {sum(explained_variance_ratio):.2f}")
12
13 # Display the variance histogram
14 plt.figure(figsize=(8, 6))
15 components = np.arange(1, n_components + 1)
16 plt.bar(components, explained_variance_ratio, color='skyblue',
17         edgecolor='black')
18 plt.xlabel('Principal Component')
19 plt.ylabel('Explained Variance Ratio')
20 plt.title('Explained Variance Ratio per Principal Component')
21 plt.xticks(components)
22 plt.show()
```

```
1 Explained Variance Ratio per Component:
2 [0.69273156 0.10438945 0.0352178  0.02309251 0.02007595 0.01203258
3  0.01042923 0.00811969 0.00682142 0.00531265 0.00430818 0.00391379
4  0.00334314 0.00318951 0.00295682 0.00272996 0.00255025 0.00215913
5  0.00208065 0.00204379 0.00179613 0.00175866 0.00152435]
6
7 Explained Variance per Component:
8 [20.593342    3.1032622   1.0469453   0.686488    0.5968127   0.35770142
9   0.31003737   0.24138017   0.20278552   0.15793316   0.12807232   0.11634817
10  0.09938393   0.09481703   0.08789954   0.08115551   0.07581323   0.06418614
11  0.06185301   0.06075718   0.0533948   0.05228084   0.04531556]
12 Cumulative explained variance for 3 components: 0.95
```

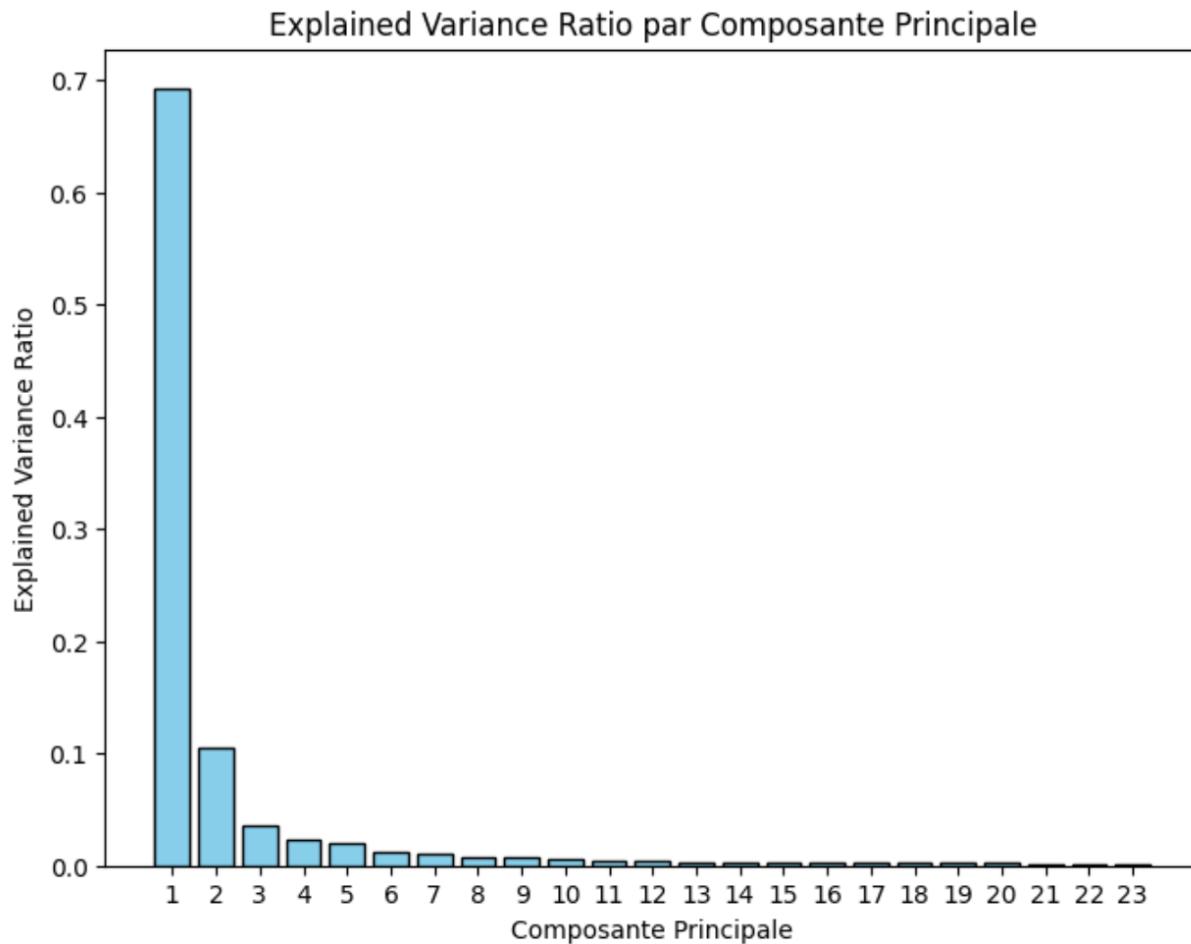


Figure 3: Histogram of Explained Variance Ratio per Principal Component

## Standard PCA

We now perform PCA for each dimension to visualize the classification.

In 1D:

```
1  from sklearn.decomposition import PCA
2  import matplotlib.pyplot as plt
3
4  n_components = 1
5  pca = PCA(n_components)
6  train_embeddings_1d = pca.fit_transform(train_embeddings.numpy())
7
8  # Histogram visualization
9  plt.figure(figsize=(10, 8))
10 plt.hist(train_embeddings_1d[:, 0], bins=30, color='skyblue',
11          edgecolor='black')
12 plt.title("1D Histogram Visualization of PCA on the Embeddings (Training
13          Set)")
```

```
12 plt.xlabel("Principal Component 1")
13 plt.ylabel("Frequency")
14 plt.show()
```

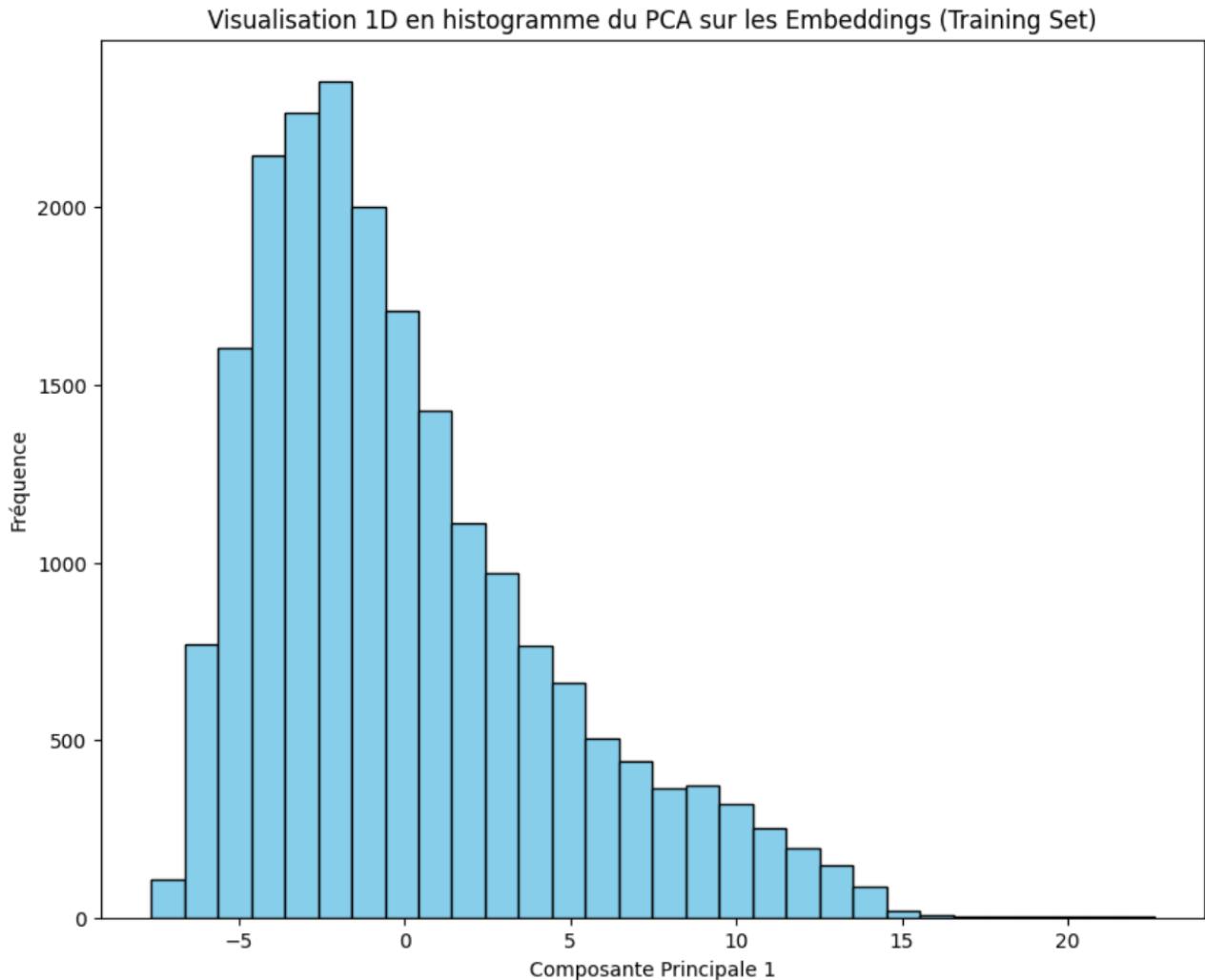


Figure 4: 1D PCA Histogram

We note that the majority of the points are located before zero, and

In 2D:

```
1 from sklearn.decomposition import PCA
2 from sklearn.decomposition import KernelPCA
3
4 n_components = 2
5 pca = PCA(n_components)
6 # pca = KernelPCA(n_components=2, kernel='rbf', gamma=0.1)
7 # pca = KernelPCA(n_components=2, kernel='poly', degree=2)
8 # pca = KernelPCA(n_components=2, kernel='rbf', gamma=0.1)
```

```
9 # pca = KernelPCA(n_components=2, kernel='rbf', gamma=0.1)
10
11 train_embeddings_2d = pca.fit_transform(train_embeddings.numpy())
12
13 plt.figure(figsize=(10, 8))
14 scatter = plt.scatter(train_embeddings_2d[:, 0], train_embeddings_2d[:, 1],
15                      c=train_y, cmap='viridis', alpha=0.5)
16 plt.colorbar(scatter, label='Sentiment Class')
17 plt.title("2D Visualization of PCA on the Embeddings")
18 plt.xlabel("Principal Component 1")
19 plt.ylabel("Principal Component 2")
20 plt.show()
```

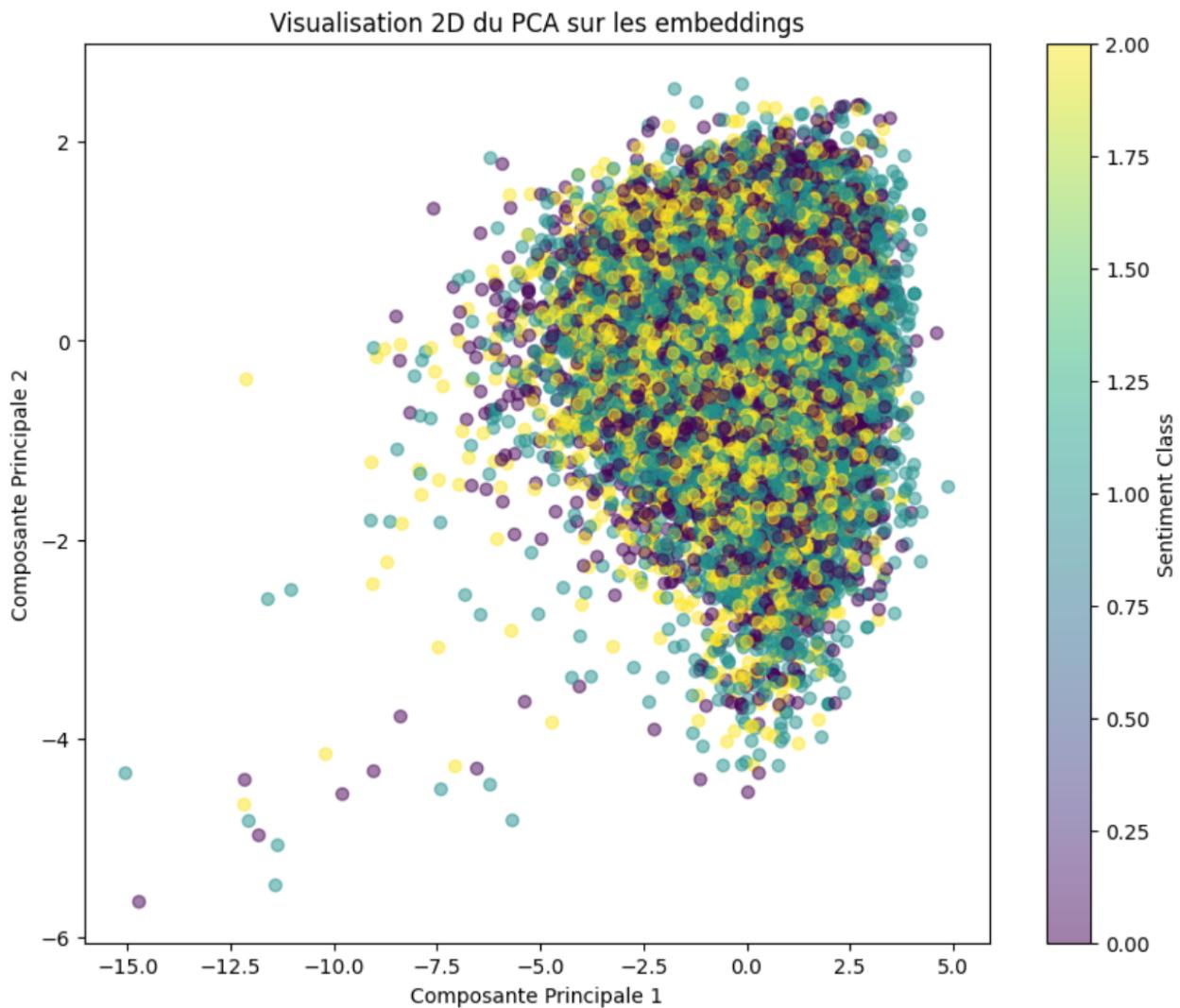


Figure 5: 2D Visualization of PCA

We observe a strong overlap of all our data, which is normal since we have projected ap-

proximately 800 dimensions onto only 2; it would have been surprising to observe a linear relationship. Therefore, it might be worthwhile to visualize the data in 3D.

### In 3D:

```
1 from sklearn.decomposition import KernelPCA
2
3 pca = PCA(3)
4 train_embeddings_3d = pca.fit_transform(train_embeddings.numpy())
5
6
7 fig = plt.figure(1, figsize=(8, 6))
8 ax = fig.add_subplot(111, projection="3d", elev=-150, azim=110)
9
10 scatter = ax.scatter(
11     train_embeddings_3d[:, 0],
12     train_embeddings_3d[:, 1],
13     train_embeddings_3d[:, 2],
14     c=train_y,
15     cmap='Spectral',          # Colormap to better distinguish classes
16     s=40,                     # Size of the points
17     alpha=0.8                 # Slight transparency
18 )
19
20 # Axis configuration
21 ax.set(
22     title="Projection of the First Three PCA Dimensions",
23     xlabel="Principal Component 1",
24     ylabel="Principal Component 2",
25     zlabel="Principal Component 3",
26 )
27
28 # Removal of tick labels for clarity
29 ax.xaxis.set_ticklabels([])
30 ax.yaxis.set_ticklabels([])
31 ax.zaxis.set_ticklabels([])
32
33 # Add a colorbar to represent the classes
34 colorbar = fig.colorbar(scatter, ax=ax, pad=0.1, aspect=10)
35 colorbar.set_label('Classes')
36
```

```
37 # Display the plot  
38 plt.show()
```

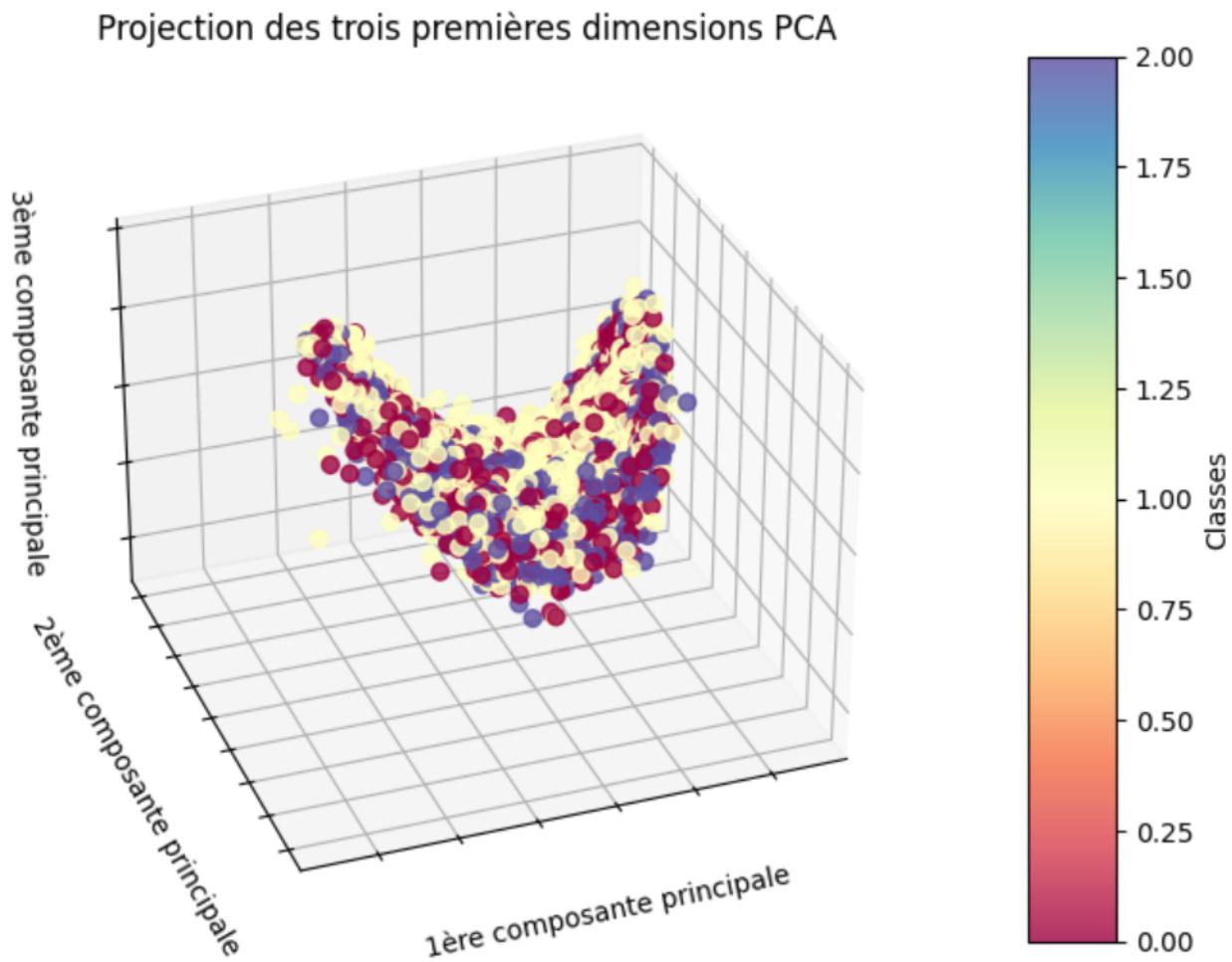


Figure 6: 3D Visualization of the PCA

It can be observed that the 3D shape is a "chip", meaning that it is concave along one axis and convex along the other. It might therefore be interesting to study the data with a PCA using a doubly polynomial kernel of order 2.

## Kernel PCA

We replace the PCA definition line with:

```
1 pca = KernelPCA(n_components=3, kernel='poly', degree=2)
```

We then obtain:

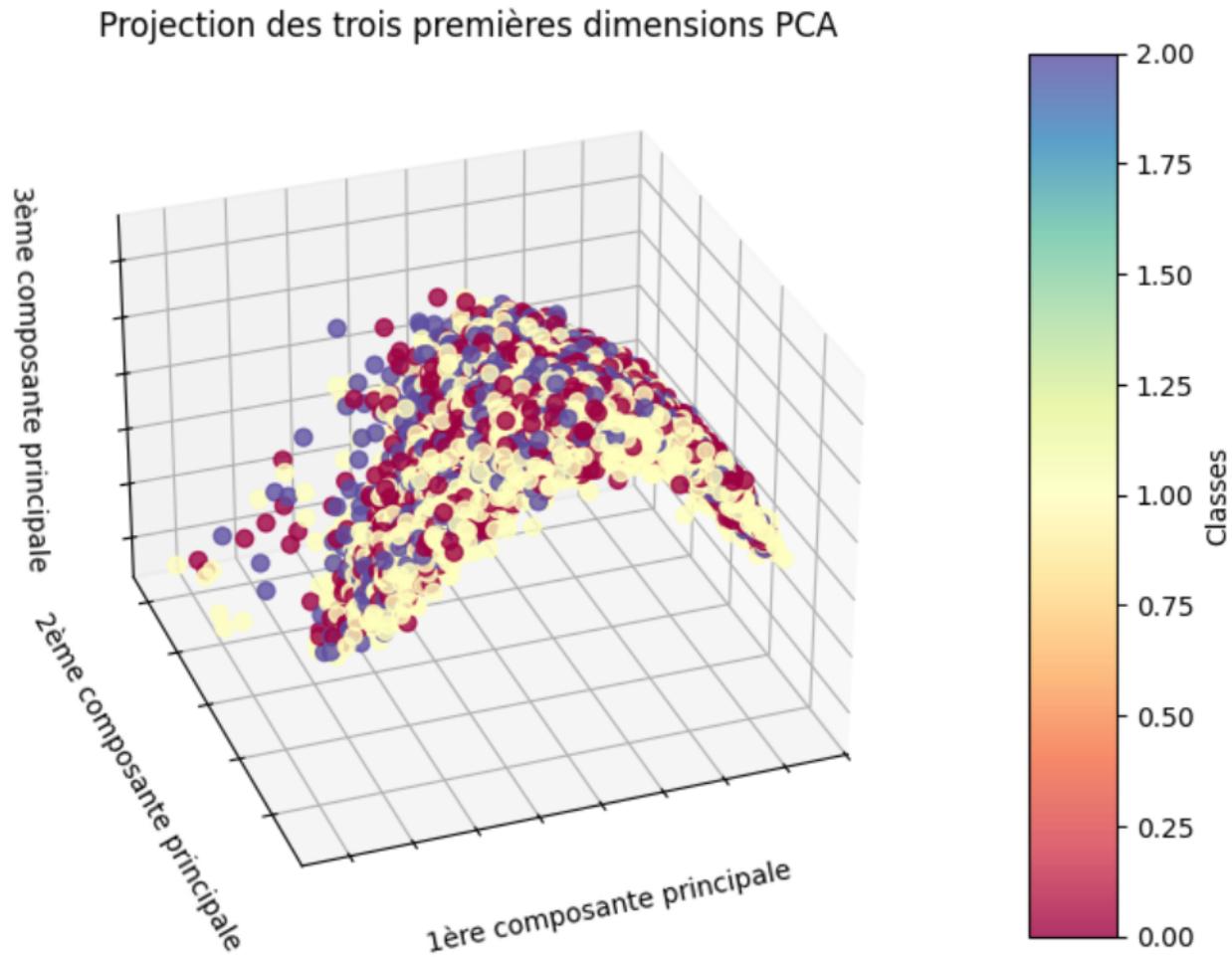


Figure 7: 3D Visualization of the KernelPCA

The points still remain very overlapped but a clearer distinction appears between the points of class 1 and the other classes.

## Visualization using T-SNE

Using T-SNE, here is what we obtain.

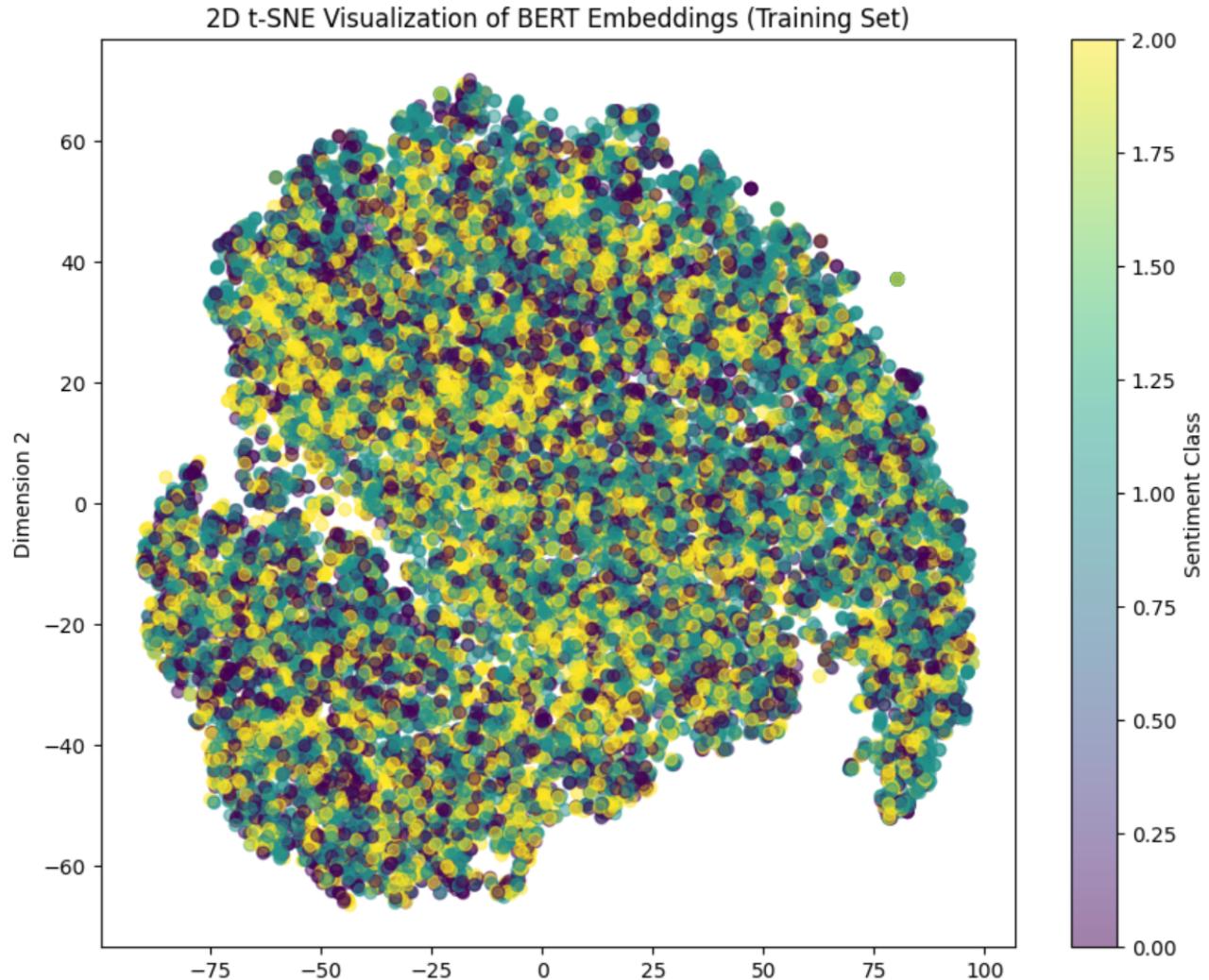


Figure 8: 2D Visualization of the TSNE

We cannot deduce anything from this figure, and the computation time is extremely long compared to PCA.

## Classification using KMEANS

Starting from the PCA, we can try to run the 3-means algorithm on the processed data to attempt unsupervised classification. Using  $n_{components} = 23$  to capture 95% of the variance, here are the 3-means results obtained.

```
1  from sklearn.decomposition import PCA
2  from sklearn.cluster import KMeans, DBSCAN
3
4  pca = PCA(n_components)
5  train_embeddings_pca = pca.fit_transform(train_embeddings.numpy())
6
7  k = 3
8  kmeans = KMeans(n_clusters=k, random_state=42)
9  clusters = kmeans.fit_predict(train_embeddings_pca)
10
11 # 3. Visualize the result
12 plt.figure(figsize=(8,6))
13 sns.scatterplot(x=train_embeddings_pca[:,0], y=train_embeddings_pca[:,1],
14                  hue=clusters, palette="viridis", s=60)
15 plt.title("K-means after PCA reduction")
16 plt.xlabel("Principal Component 1")
17 plt.ylabel("Principal Component 2")
18 plt.legend(title="Clusters")
19 plt.show()
```

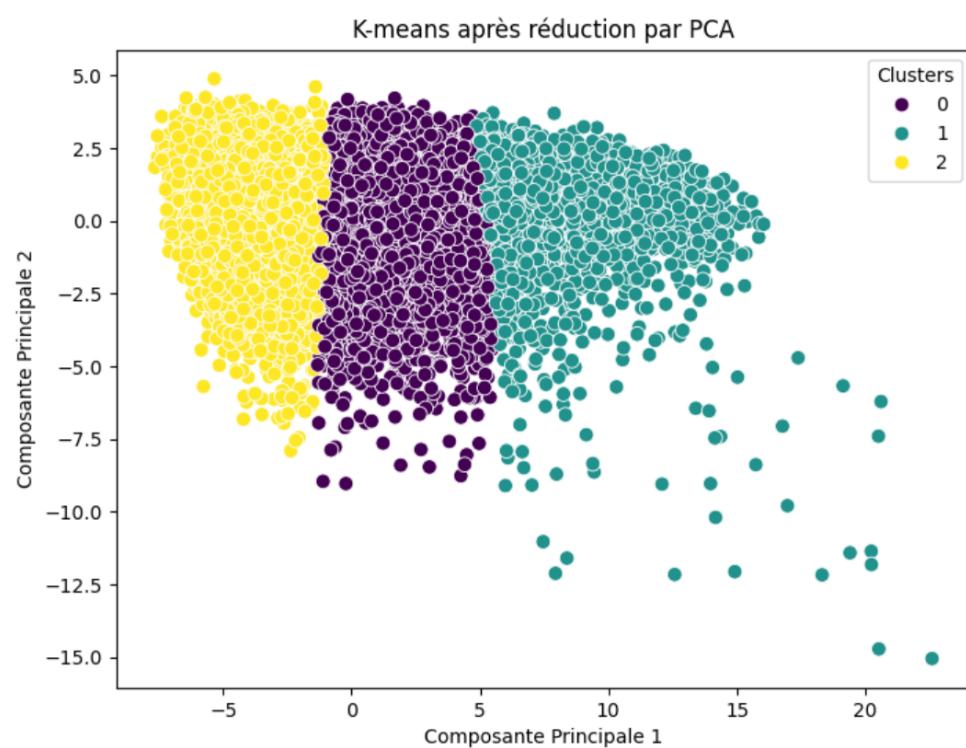


Figure 9: Visualization of the 3-means

We observe that this does not correspond at all to the classification we observed in the simple 2D PCA figure and we will verify it by calculating the accuracy.

```
1 def kmeans_accuracy(y_true, y_pred):
2     """
3         Function that calculates the accuracy of the predictions without taking
4             the label into account.
5         For each cluster, the majority class is chosen to determine the
6             prediction.
7         It is then compared with the true labels.
8     """
9
10    n = len(y_true)
11    y_true = np.array(y_true)
12    y_pred = np.array(y_pred)
13
14    # labels[i] represents cluster i. labels[i][j] is the number of times
15    # the points in cluster i are classified as j.
16    labels = [[0,0,0] for _ in range(3)]
17    for i in range(n):
18        labels[y_pred[i]][y_true[i]] += 1
```

```
17     mapping = [np.argmax(labels[i]) for i in range(len(labels))]
18     print(mapping)
19     print("The mapping is: {}".format(mapping[0] != mapping[1] !=
20         ↪ mapping[2]))
21
22     score = 0
23     for i in range(n):
24         score += (y_true[i] == mapping[y_pred[i]])
25
26     return score/n
27
28 kmeans_accuracy(train_y, clusters)
```

We then obtain:

```
1 [1, 1, 1]
2 The mapping is: False
3 0.4041727316836487
```

All points are predicted as 1, which corresponds to the neutral class that dominates the dataset and is therefore uniformly distributed in each cluster. This means that our **k-means** is not capable of finding patterns and only gives us the ratio of neutral elements in the dataset (0.40) ... By trying with a KernelPCA, we obtain exactly the same results. Therefore, PCA and clustering methods are not useful here and more complex methods will be needed to predict the points in the space.

# Q1: Use of Classical Machine Learning Techniques

## Data Preprocessing

The sentence embeddings extracted via BERT are used as input features for sentiment classification. Let  $X$  be the embeddings matrix and  $Y$  the associated sentiment labels. After extraction, the data is split into a training set and a test set:

```
1 X_train = train_embeddings.cpu().numpy()
2 X_test = test_embeddings.cpu().numpy()
3
4 y_train = np.array(train_y)
5 y_test = np.array(test_y)
```

Since machine learning models are sensitive to the scale of the data, the embeddings are normalized using a StandardScaler to improve algorithm convergence:

```
1 scaler = StandardScaler()
2 X_train = scaler.fit_transform(X_train)
3 X_test = scaler.transform(X_test)
```

## Training and Evaluation of the Models

A class ModelTrainer is defined to manage the training and evaluation of the classification models. It includes:

- Training of the model with timing.
- Evaluation of performance via accuracy and a detailed report.
- Display of a confusion matrix to visualize common errors.

```
1 class ModelTrainer:
2     def __init__(self, model, model_name, **hyperparams):
3         self.model_name = model_name
4         self.model = model(**hyperparams)
5         self.training_time = None
6         self.results = []
7
8     def train(self, X_train, y_train):
9         start_time = time.time()
10        self.model.fit(X_train, y_train)
```

```
11     self.training_time = time.time() - start_time
12
13     def evaluate(self, X_test, y_test):
14         y_pred = self.model.predict(X_test)
15         acc = accuracy_score(y_test, y_pred)
16         report = classification_report(y_test, y_pred)
17         cm = confusion_matrix(y_test, y_pred)
18
19         self.results = {
20             'accuracy': acc,
21             'report': report,
22             'training_time': self.training_time,
23             'confusion_matrix': cm,
24             'predictions': y_pred
25         }
26
27         print(f"--- {self.model_name} ---")
28         print(f"Accuracy: {acc:.4f}")
29         print(report)
30         print(f"Training Time: {self.training_time:.2f} seconds")
31
32         # Display the confusion matrix
33         plt.figure(figsize=(6, 4))
34         sns.heatmap(cm, annot=True, fmt='d', cmap='Blues',
35             xticklabels=set(y_test), yticklabels=set(y_test))
36         plt.xlabel('Predicted Label')
37         plt.ylabel('True Label')
38         plt.title(f'Confusion Matrix - {self.model_name}')
39         plt.show()
40
41
42         return self.results
```

## Experimentation with Different Models

Different classification models are tested to compare their performance on the BERT embeddings.

## Logistic Regression

Different values for the parameter `max_iter` are tested to analyze their impact on convergence:

```
1 parameter = [{'max_iter':50}, {'max_iter':100}, {'max_iter':200}]
2 for params in parameter:
3     trainer = ModelTrainer(LogisticRegression, f"LogisticRegression with
4         ↪ max_iter={params['max_iter']}", **params)
5     trainer.train(X_train, y_train)
6     trainer.evaluate(X_test, y_test)
```

## Support Vector Machine (SVM)

Several kernels (`kernel`) are explored to observe their influence on classification:

```
1 parameter = [
2     {"kernel": "linear"},
3     {"kernel": "poly", "degree": 3, "gamma": "scale", "coef0": 1},
4     {"kernel": "rbf", "gamma": "scale"},
5     {"kernel": "sigmoid", "gamma": "scale", "coef0": 0}
6 ]
7
8 for params in parameter:
9     trainer = ModelTrainer(SVC, f"SVM with kernel={params['kernel']}", 
10         ↪ **params)
11     trainer.train(X_train, y_train)
12     trainer.evaluate(X_test, y_test)
```

## Random Forest

We tested several ensemble sizes (`n_estimators`):

```
1 parameter = [{'n_estimators':50}, {'n_estimators':100}, {'n_estimators':150},
2     ↪ {'n_estimators':200}]
3 for params in parameter:
4     trainer = ModelTrainer(RandomForestClassifier, f"RandomForestClassifier
5         ↪ with n_estimators={params['n_estimators']}", **params)
6     trainer.train(X_train, y_train)
7     trainer.evaluate(X_test, y_test)
```

## Q2: Use of an MLP for Prediction

### Definition and Code

First, before building our MLP, we must extract our embeddings using BERT to be able to process our text as a 768-dimensional vector.

```
1 # Hyperparameters
2 PRETRAINED_MODEL = "bert-base-uncased"
3 MAX_LENGTH = 128
4 BATCH_SIZE = 64
5 NUM_CLASSES = 3
6 LEARNING_RATE = 2e-5
7 EPOCHS = 10
8 device = torch.device("cuda" if torch.cuda.is_available() else "cpu")
9
10 # Initialize tokenizer, dataset, and dataloader
11 tokenizer = BertTokenizer.from_pretrained(PRETRAINED_MODEL)
12 traindataset = TextDataset(train_x, train_y, tokenizer, MAX_LENGTH)
13 trainloader = DataLoader(traindataset, batch_size=BATCH_SIZE, shuffle=False)
14 valdataset = TextDataset(val_x, val_y, tokenizer, MAX_LENGTH)
15 valloader = DataLoader(valdataset, batch_size=BATCH_SIZE, shuffle=False)
16
17 bert = BertModel.from_pretrained(PRETRAINED_MODEL).to(device)
18 # Extract embeddings
19 train_embeddings = extract_embeddings(bert, trainloader, device)
20 train_embeddings = train_embeddings.cpu()
21
22 val_embeddings = extract_embeddings(bert, valloader, device)
23 val_embeddings = val_embeddings.cpu()
24
25 # Test embeddings
26 testdataset = TextDataset(test_x, test_y, tokenizer, MAX_LENGTH)
27 testloader = DataLoader(testdataset, batch_size=BATCH_SIZE, shuffle=False)
28
29 bert = BertModel.from_pretrained(PRETRAINED_MODEL).to(device)
30 # Extract embeddings
31
32 test_embeddings = extract_embeddings(bert, testloader, device)
33 test_embeddings = test_embeddings.cpu()
```

All embeddings are respectively stored in `train_embeddings` or `test_embeddings`. We now define our model.

```
1 criterion = nn.CrossEntropyLoss()
2 # Model Definition
3 class SentimentClassifier2(nn.Module):
4     def __init__(self, hidden_sizes, output_size):
5         super(SentimentClassifier2, self).__init__()
6
6         n = len(hidden_sizes)
7         self.fc_layers = nn.ModuleList()
8         self.bn_layers = nn.ModuleList()
9         # Fully connected layer for classification
10        for i in range(n-1):
11            self.fc_layers.append(nn.Linear(hidden_sizes[i],
12                → hidden_sizes[i+1]))
13            self.bn_layers.append(nn.BatchNorm1d(hidden_sizes[i + 1]))
14
15        self.fc_layers.append(nn.Linear(hidden_sizes[-1], output_size))
16        self.dropout = nn.Dropout(0.2)
17
18    def forward(self, x):
19        # Pass through the fully connected layers
20        n = len(self.fc_layers)
21        for fc, bn in zip(self.fc_layers, self.bn_layers):
22            x = fc(x)
23            x = bn(x)
24            x = F.gelu(x)
25            x = self.dropout(x)
26
27        x = self.fc_layers[-1](x)
28        return x
29
30 # Training script
31
32 def train(model, train_loader, optimizer, epoch, log_interval=50):
33     model.train()
34     loss_cpu = 0
35     correct = 0
36     total = 0
```

```
37     for batch_idx, data in enumerate(train_loader, 0):
38         # Get the inputs; data is a list of [inputs, labels]
39         inputs, target = data['input_ids'], data['label']
40         inputs, target = inputs.cuda(), target.cuda()
41         # inputs = inputs.detach()
42         optimizer.zero_grad()
43
44         outputs = model(inputs)
45         loss = criterion(outputs, target)
46
47         loss.backward()
48         optimizer.step()
49         _, predicted = torch.max(outputs.data, 1)
50         loss_cpu += loss.item()
51         total += target.size(0)
52         correct += predicted.eq(target.data).cpu().sum()
53
54     if batch_idx % log_interval == 0:
55         print(' | Epoch [%3d/%3d] Iter[%3d/%3d]\t\tLoss: %.4f Acc@1:
56             %.3f%%'
57             % (epoch, EPOCHS, batch_idx+1,
58                 (len(train_loader)//BATCH_SIZE)+1, loss.item(),
59                 → 100.*correct/total))
60         # n_iter = epoch * len(train_loader) + batch_idx
61
62     return loss_cpu/len(train_loader)
63
64
65
66
67
68     # Initialize F1 metric
69     f1_metric = MulticlassF1Score(num_classes=3, average='macro').to(device)
70
71     with torch.no_grad():
72         all_preds = []
73         all_labels = []
74
```

```
75     for batch_idx, data in enumerate(test_loader, 0):
76         inputs, target = data['input_ids'], data['label']
77         inputs, target = inputs.cuda(), target.cuda()
78
79         outputs = model(inputs)
80         loss = criterion(outputs, target)
81         test_loss_MSE += loss.item()
82
83         _, predicted = torch.max(outputs.data, 1)
84         total += target.size(0)
85         correct += predicted.eq(target.data).cpu().sum()
86
87         # Collect predictions and labels
88         all_preds.append(predicted.cpu())
89         all_labels.append(target.cpu())
90
91     # Concatenate predictions and labels
92     all_preds = torch.cat(all_preds).numpy()
93     all_labels = torch.cat(all_labels).numpy()
94
95     # Calculate F1-score
96     f1_score = f1_metric(torch.tensor(all_preds).to(device),
97                          torch.tensor(all_labels).to(device))
98
99     # Plot confusion matrix if required
100    if plot:
101        cm = confusion_matrix(all_labels, all_preds)
102        plt.figure(figsize=(6, 4))
103        sns.heatmap(cm, annot=True, fmt='d', cmap='Blues',
104                    xticklabels=["negative", "neutral", "positive"],
105                    yticklabels=["negative", "neutral", "positive"])
106        plt.xlabel('Predicted Label')
107        plt.ylabel('True Label')
108        plt.title('Confusion Matrix - MLP')
109        plt.show()
110
111    # Final metrics
112    test_loss_MSE = test_loss_MSE / len(test_loader)
113    accuracy = 100. * correct / total
```

```
111     print("\n| Validation Epoch #%-d\t\tLoss: %.4f Acc@1: %.2f%% F1: %.2f" %
112         → (epoch, test_loss_MSE, accuracy, f1_score.item()))
113
114     return test_loss_MSE, accuracy, f1_score.item()
```

Compared to the initial code provided, here are the different choices we made.

1. We chose to have multiple layers whose topology is given by the array `hidden_size`, the number of neurons in layer  $i$  is given by `hidden_sizes[i]`.
2. After each hidden layer, the activation function is `gelu`. It is defined as  $GELU(x) = x\phi(x)$  with  $\phi$  being a normal distribution. This function preserves the properties of ReLU while retaining relationships in the embedding that are lost by ReLU.
3. After each layer, we apply a **dropout of 0.2**, which disables a portion of the neurons at each step and helps prevent loss plateaus and overfitting.
4. We use cross-entropy because it is the classical loss function for classification problems.
5. We use the F1-score. This metric quantifies the number of false positives and false negatives through precision and recall. It is important when trying to prevent negative texts from being classified as positive for moderation. A good value is around **0.65**. In addition, we also display the confusion matrix.
6. We normalize the data at each layer's output.

The training code remains the same as initially provided with the few changes mentioned.

```
1 print ("Let us Train.")
2
3 EPOCHS = 40
4 input_size = len(train_embeddings[0])
5 output_size = 3
6 hidden_sizes = [input_size, 512, 128, 16]
7
8 model = SentimentClassifier2(hidden_sizes, output_size).to(device)
9 model_test = SentimentClassifier2(hidden_sizes, output_size).to(device)
10
11 best_error = float('inf')
12 LEARNING_RATE = 1e-4
13
14 from torch.optim.lr_scheduler import ReduceLROnPlateau, StepLR
15
16 traindataset = EmbeddingDataset(train_embeddings, train_y)
```

```
17 trainloader = DataLoader(traindataset, batch_size=BATCH_SIZE, shuffle=True)
18 valdataset = EmbeddingDataset(val_embeddings, val_y)
19 valloader = DataLoader(valdataset, batch_size=BATCH_SIZE, shuffle=False)
20
21 optimizer = torch.optim.AdamW(model.parameters(), lr=LEARNING_RATE,
22   ↪ weight_decay=1e-4)
23 # lr_scheduler = StepLR(optimizer, step_size=20, gamma=0.1)
24 lr_scheduler = ReduceLROnPlateau(optimizer, 'min', patience=7)
25
26 best_error = float('inf')
27 train_history = []
28 val_history = []
29 print('-----')
30 print('3/ training model 1')
31 print('-----')
32
33 iter = 0
34 for epoch in range(EPOCHS):
35     loss = train(model, trainloader, optimizer, epoch)
36     train_history.append(loss)
37
38     val_loss, acc, _ = test(model, valloader, epoch)
39     val_history.append(val_loss)
40
41     lr_scheduler.step(val_loss)
42     print('lr =', get_lr(optimizer))
43
44     if val_loss < best_error:
45         best_error = val_loss
46         print('Best RMSE is of : ' + str(best_error), 'for epoch :', epoch,
47           ↪ 'ERROR TEST = ', val_loss)
48         # model_test.parameters() = model.state_dict()
49         model_test.load_state_dict(model.state_dict(), strict=True)
50         torch.save(model_test.state_dict(), "best_model.pth") # Save the
51           ↪ model
52
53 print ("Training Done!")
```

“latex

## Hyperparameter Choices and Performance

With all the choices mentioned above, we observed several phenomena to be mitigated.

1. Adding new layers marginally improves the model's effectiveness. To have an effective contribution, we chose that **each layer would be a power of 2** and **the batch size would be 32**. `hidden_sizes = [input_size, 512, 128, 16]`
2. The loss function reaches a plateau after N=25 epochs. To mitigate this, we added a **dropout of 0.2**, we **used a ReduceLROnPlateau scheduler**, we **normalized the data at each layer's output** and we chose a learning rate of  $\eta = 10^{-4}$  to favor convergence toward an optimum. Additionally, we chose EPOCHS = 50 to avoid overfitting.
3. The validation loss function differs from the training loss, which is characteristic of overfitting. The dropout once again helps improve performance in this case.

## Results

We then obtain for the previous choices:

```
1 Let us Train.  
2 -----  
3 / training model 1  
4 -----  
5 | Epoch [ 0/ 40] Iter[ 1/ 6] Loss: 1.0953 Acc@1: 34.375%  
6 | Epoch [ 0/ 40] Iter[ 51/ 6] Loss: 1.1323 Acc@1: 35.049%  
7 | Epoch [ 0/ 40] Iter[101/ 6] Loss: 1.0297 Acc@1: 39.851%  
8 | Epoch [ 0/ 40] Iter[151/ 6] Loss: 0.9967 Acc@1: 42.074%  
9 | Epoch [ 0/ 40] Iter[201/ 6] Loss: 1.0124 Acc@1: 44.364%  
10 | Epoch [ 0/ 40] Iter[251/ 6] Loss: 0.9834 Acc@1: 45.667%  
11 | Epoch [ 0/ 40] Iter[301/ 6] Loss: 0.9906 Acc@1: 46.865%  
12  
13 | Validation Epoch #0 Loss: 1.0187 Acc@1: 46.10% F1:  
14   ↳ 0.45  
15 lr = 0.0001  
16 Best RMSE is of : 1.0187300780305155 for epoch : 0 ERROR TEST =  
17   ↳ 1.0187300780305155  
18 | Epoch [ 1/ 40] Iter[ 1/ 6] Loss: 0.9204 Acc@1: 60.938%  
19 | Epoch [ 1/ 40] Iter[ 51/ 6] Loss: 1.0690 Acc@1: 55.300%  
20 | Epoch [ 1/ 40] Iter[101/ 6] Loss: 0.9006 Acc@1: 55.554%  
21 | Epoch [ 1/ 40] Iter[151/ 6] Loss: 0.9869 Acc@1: 55.867%  
22 | Epoch [ 1/ 40] Iter[201/ 6] Loss: 1.0376 Acc@1: 56.172%  
23 | Epoch [ 1/ 40] Iter[251/ 6] Loss: 0.9264 Acc@1: 56.431%
```

```
22 | Epoch [ 1/ 40] Iter[301/  6]           Loss: 0.8589 Acc@1: 56.785%
23
24 ...
25
26 | Validation Epoch #39                  Loss: 0.7976 Acc@1: 64.32% F1:
27   ↵ 0.64
28 lr = 1.0000000000000002e-06
Training Done!
```

The evolution of the Loss curve is given below.

```
1 plt.figure(figsize=(8, 6))
2 plt.plot(train_history[4:], label='Train Loss', marker='o')
3 plt.plot(val_history[4:], label='Validation Loss', marker='o')
4 plt.xlabel("Epoch")
5 plt.ylabel("Loss")
6 plt.title("Training and Validation Loss over Epochs")
7 plt.legend()
8 plt.grid(True)
9 plt.show()
```

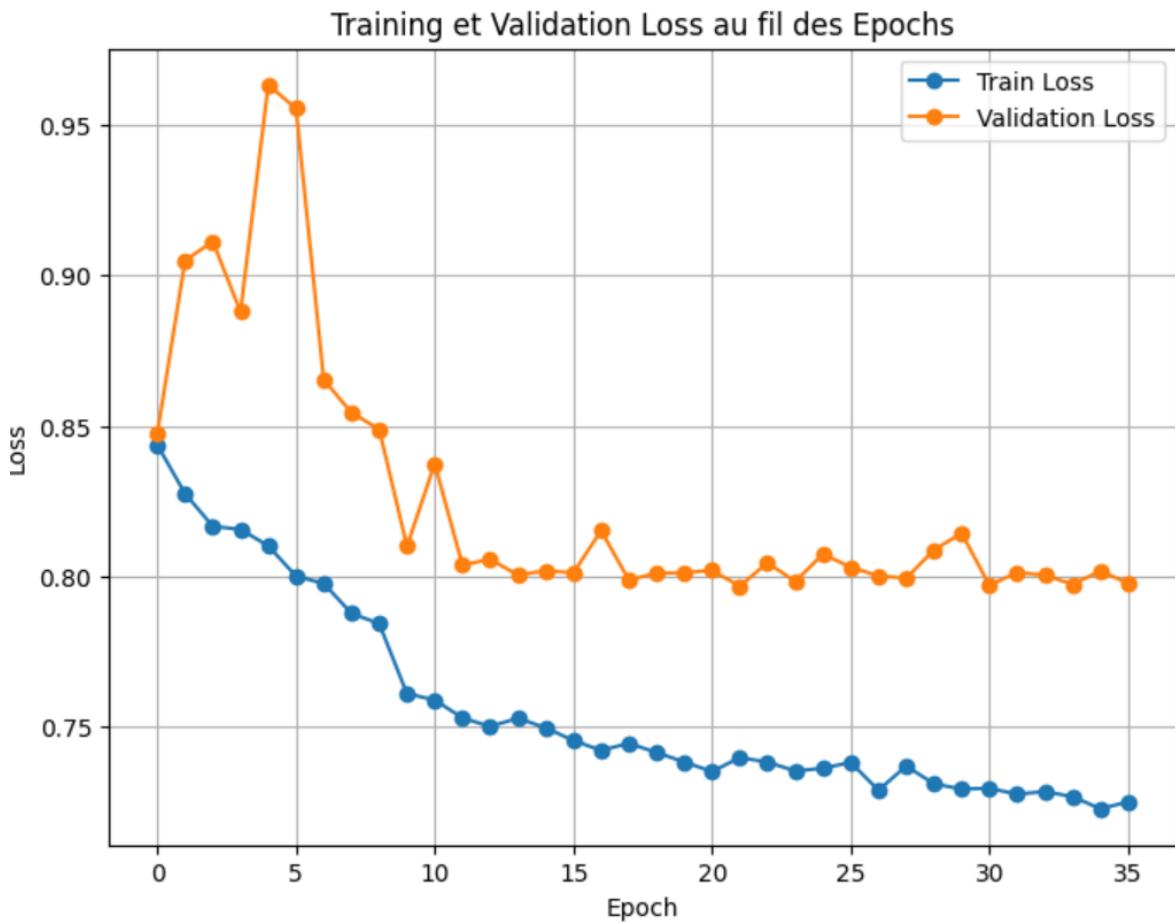


Figure 10: Evolution of the loss curve over Epochs (40 epochs)

We then test the model on the validation data.

```
1 testdataset = TextDataset(test_x,test_y, tokenizer, MAX_LENGTH)
2 testloader = DataLoader(testdataset, batch_size=BATCH_SIZE, shuffle=False)
3
4 bert = BertModel.from_pretrained(PRETRAINED_MODEL).to(device)
5 # Extract embeddings
6
7 test_embeddings = extract_embeddings(bert, testloader, device)
8 test_embeddings = test_embeddings.cpu()
9
10 print('===== ')
11 print('Test set = ')
12 testdataset = EmbeddingDataset(test_embeddings, test_y)
13 testdataset = DataLoader(testdataset, batch_size=BATCH_SIZE, shuffle=True)
14 test_loss, accu, f1 = test(model, testdataset, epoch, True)
15 print('===== ')
```

```
16 print("Test_loss: {}, Accu: {}, F1-score: {}, Taux d'erreur critique:  
    ↴ ??".format(test_loss, accu, f1))
```

We then obtain these indicators:

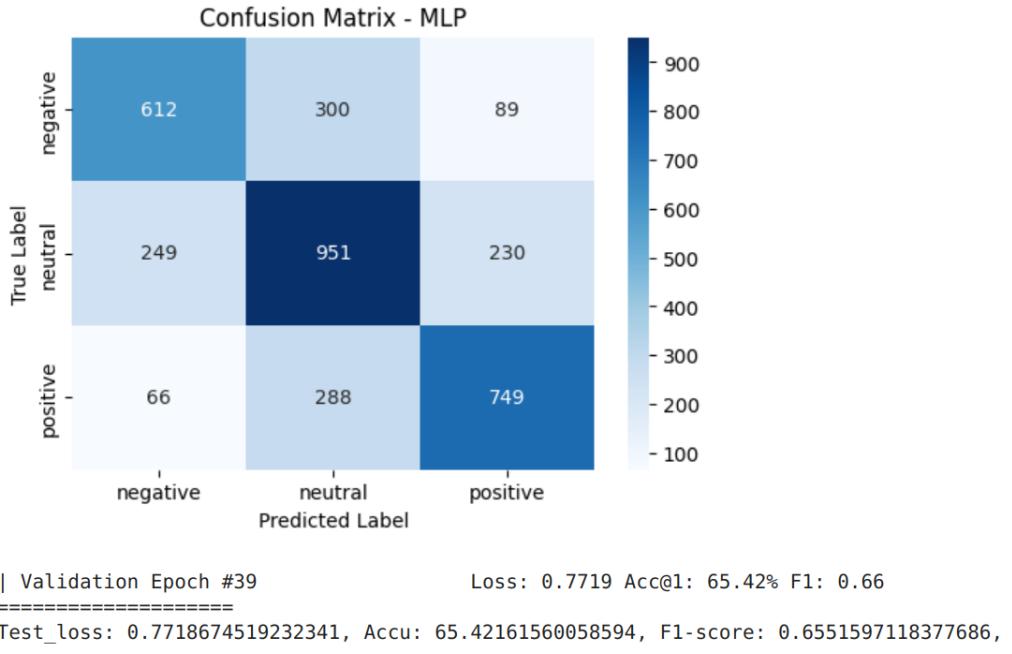


Figure 11: MLP Results on the validation data

We have a fairly good F1-score of **F1-score = 0.65**. However, despite all our trials, the maximum accuracy is **Acc = 65.42%**. This is rather good since it approaches the accuracy of the best models, as we will see in the following sections. Finally, we note that 89 negative texts are classified as positive, which is low compared to the amount of data but too high to use this predictor as a moderation tool.

**Note:** We obtained an accuracy of 65.7% after several training trials by changing the random seed, and we consider the MLP to have this accuracy for the remainder.

## Q3: Performance Analysis and Selection of the Best Model

We analyze the performance of the sentiment classification models using three main metrics:

- **Accuracy:** percentage of correct classifications.
- **Average F1-score:** balance between precision and recall.
- **Critical Error Rate:** proportion of completely inverted predictions (e.g., positive classified as negative and vice versa). This score, which may initially seem exotic, is considered important because it is thought to be much less serious to classify a text as neutral knowing that it is positive or negative (or vice versa) than to classify a text as positive when it is negative (or vice versa).

### Model Performance

The results are presented in the table below:

Model	Accuracy	F1-score	Critical Error Rate	Time (s)
Logistic Regression (50)	63.6%	0.64	4%	2.01
Logistic Regression (100)	65.3%	0.65	4%	4.64
Logistic Regression (200)	65.6%	0.66	4%	5.48
SVM (linear)	64.2%	0.64	5%	2806.39
SVM (poly)	64.8%	0.65	4%	236.24
SVM (RBF)	63.5%	0.63	4%	212.49
SVM (sigmoid)	34.0%	0.34	25%	153.29
Random Forest (50 trees)	53.7%	0.53	6%	21.06
Random Forest (100 trees)	54.2%	0.53	5%	41.76
Random Forest (150 trees)	55.6%	0.54	5%	62.47
Random Forest (200 trees)	55.0%	0.54	5%	84.14
MLP	65.7%	0.66	4.75%	N/A

Table 1: Comparison of model performance

### Comparative Analysis

**Accuracy** The MLP model achieves the highest accuracy with **65.7%**, slightly surpassing logistic regression (**65.6%**). The SVM and random forest models show lower performance, except for the poly kernel which reaches **64.8%**.

**F1-score** The evolution of the F1-score follows that of accuracy, with the MLP and logistic regression at the top (**0.66**), and SVM poly just behind (**0.65**). The random forests remain behind (**0.54**).

**Critical Error Rate** The SVM sigmoid remains the worst choice with a critical error rate of **25%**. The MLP displays a slightly higher rate than logistic regression (**4.75%** vs. **4%**), which is still within the same range of performance.

**Training Time** The training time of the MLP is not provided here, but it strongly depends on the network size and available computational power. In contrast, the linear and poly SVMs are the most time-consuming (**2806s** and **236.24s** respectively), while logistic regression remains very fast.

## Selection of the Optimal Model

Based on these results, two models stand out:

- **MLP**, which offers the highest accuracy (**65.7%**) and an F1-score identical to logistic regression (**0.66**).
- **Logistic Regression with max\_iter=200**, which has a very close accuracy (**65.6%**), a lower critical error rate (**4%**) and a very short training time.

Thus, the MLP is a promising alternative, although it may require more computational resources depending on its configuration.

## Conclusion

Logistic regression and the MLP are the two best candidates for this task. If computational constraints are limited, the MLP can be an interesting choice to obtain a slight improvement in performance. However, logistic regression remains an excellent compromise between accuracy, speed, and robustness.

## Q4: Comparison with the Performance of an LLM

We aim to compare the performance of our best model with that of an LLM. We therefore import the **Roberta** model from Hugging Face, which is specialized in sentiment classification, in order to compare its performance on the test data with our model.

```
1  from transformers import pipeline
2  import numpy as np
3  from sklearn.metrics import accuracy_score, classification_report
4
5  label_mapping = {"negative": 0, "neutral": 1, "positive": 2}
6
7  # Initialization of an LLM specialized in sentiment detection via a
8  # Hugging Face pipeline. Here we load Roberta
9  llm_classifier = pipeline(
10      "sentiment-analysis",
11      model="cardiffnlp/twitter-roberta-base-sentiment-latest",
12      tokenizer="cardiffnlp/twitter-roberta-base-sentiment-latest"
13  )
14
15  # Evaluation of the LLM on the preprocessed texts
16  llm_preds = []
17  for text in test_x:
18      result = llm_classifier(text)
19      predicted_label = result[0]["label"].lower()
20      llm_preds.append(label_mapping[predicted_label])
21
22  llm_preds = np.array(llm_preds)
23  y_test_np = np.array(test_y)
24
25  # Calculate the model's accuracy
26  llm_accuracy = accuracy_score(y_test_np, llm_preds)
27  llm_report = classification_report(
28      y_test_np, llm_preds, target_names=["negative", "neutral", "positive"]
29  )
30
31  print("== LLM Sentiment Classifier Performance ==")
32  print("LLM Accuracy: {:.4f}".format(llm_accuracy))
33  print("\nLLM Classification Report:")
34  print(llm_report)
```

We then obtain:

```
1     === LLM Sentiment Classifier Performance ===  
2     LLM Accuracy: 0.7043  
3  
4     LLM Classification Report:  
5             precision    recall    f1-score   support  
6  
7     negative        0.71      0.72      0.71      1001  
8     neutral         0.71      0.59      0.65      1430  
9     positive        0.70      0.84      0.76      1103  
10  
11    accuracy           0.70      0.70      0.70      3534  
12    macro avg       0.70      0.72      0.71      3534  
13    weighted avg    0.70      0.70      0.70      3534
```

It is then very clear that **the MLP's accuracy is very close to that of the LLM** since if the LLM represents the upper limit, **the maximum accuracy would be 70%** whereas the MLP achieves **65.7%**. This maximum limit therefore seems to be determined by the data preprocessing, and reflects the low variance distribution during the principal component decomposition.

Model	Accuracy	F1-score	Critical Error Rate	Time (s)
Roberta-LLM	70%	0.71	N/A %	N/A

Table 2: Performance of the LLM on the test set

## Q5: Use of BERT

Using BERT embeddings allows the transformation of text into representations that are practical to use. The objective is to obtain vectors that incorporate the meaning of words based on their context, thus facilitating machine learning on natural language processing tasks.

BERT generates contextualized embeddings, which means that the representation of a word depends on the context in which it appears. Unlike classical approaches such as Word2Vec or TF-IDF, this allows for better handling of polysemy and the subtleties of language.

The model is pre-trained on large amounts of data and thus integrates advanced linguistic and semantic knowledge, which is very useful for handling complex cases of language form such as double negation.

Finally, the use of embeddings also reduces dimensionality and noise compared to raw textual representations. Classical methods, such as one-hot encoding, produce very long and sparse vectors, whereas BERT embeddings provide dense and compact vectors, which are more effective for machine learning.

## Q6: How BERT Works

### Input Representation

The input is composed of three types of embeddings:

- **Token Embeddings:** Each word or sub-word (via WordPiece) is converted into a learned vector.
- **Segment Embeddings:** When the input contains two sentences (for the NSP task), an embedding indicates which sentence each token belongs to.
- **Positional Embeddings:** Since the Transformer does not understand the order of words, these embeddings add information about the position of each token.

These three embeddings are summed before being fed into the model.

### BERT Architecture

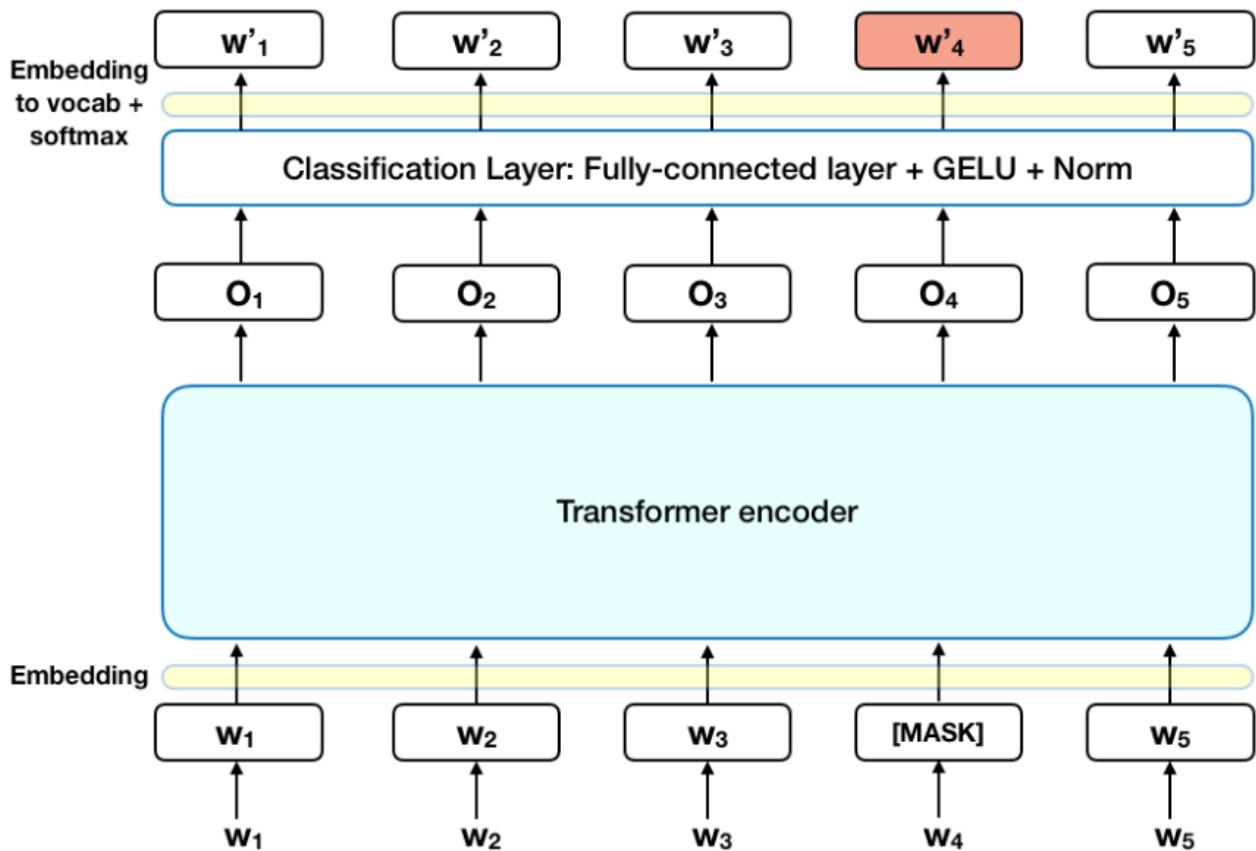


Figure 12: BERT Architecture

BERT consists of several layers of Transformer encoders (12 for BERT-Base, 24 for BERT-Large). Each layer includes:

- **Multi-Head Self-Attention:** Each word can interact with all others in the sequence to capture the global context. Multiple attention heads are used to learn different relationships.
- **Feed-Forward Network:** After self-attention, each word passes through a simple neural network (two linear layers with GELU activation).
- **Layer Normalization and Residual Connections:** Each layer is stabilized by normalization and residual connections to improve learning.

## BERT Pre-training

BERT learns from two main tasks:

- **Masked Language Modeling (MLM):** 15% of the words are masked, and the model must predict them using their context.
- **Next Sentence Prediction (NSP):** The model receives two sentences and must guess whether they logically follow one another.

## Bidirectionality

Unlike models that read text in one direction, BERT simultaneously captures both left and right contexts, allowing it to better understand the overall meaning of a text.

## Q7: Fine-tuning RoBERTa with LoRA for Sentiment Classification

In this section, we train a sentiment classification model based on **RoBERTa**, a pre-trained language model, using the **LoRA** (Low-Rank Adaptation) technique. LoRA allows for efficient model adaptation without requiring an update of all of the model's weights, thereby reducing computational complexity.

### Model Setup

We use the roberta-base model from the transformers library, which is specialized in natural language processing. We apply LoRA to limit the number of parameters adjusted during training.

```
1 import torch
2 import numpy as np
3 from transformers import AutoTokenizer, AutoModelForSequenceClassification
4 from peft import LoraConfig, get_peft_model, TaskType
5
6 # Load the model and tokenizer
7 MODEL_NAME = "roberta-base"
8 tokenizer = AutoTokenizer.from_pretrained(MODEL_NAME)
9 model = AutoModelForSequenceClassification.from_pretrained(MODEL_NAME,
10   num_labels=3)
11
12 # Apply LoRA
13 lora_config = LoraConfig(
14     task_type=TaskType.SEQ_CLS,
15     inference_mode=False,
16     r=8,
17     lora_alpha=32,
18     lora_dropout=0.1
19 )
20 model = get_peft_model(model, lora_config)
21 model.print_trainable_parameters()
```

### Data Preparation

We use the same dataset as we used previously for the MLP

```
1 class SentimentDataset(Dataset):
2     def __init__(self, texts, labels, tokenizer, max_length=128):
3         self.texts = texts
4         self.labels = labels
5         self.tokenizer = tokenizer
6         self.max_length = max_length
7
8     def __len__(self):
9         return len(self.texts)
10
11    def __getitem__(self, idx):
12        encoded = self.tokenizer(
13            self.texts[idx],
14            truncation=True,
15            padding="max_length",
16            max_length=self.max_length,
17            return_tensors="pt"
18        )
19        return {key: val.squeeze(0) for key, val in encoded.items()} |
→ {"labels": torch.tensor(self.labels[idx], dtype=torch.long)}
```

We define the training, validation, and test sets:

```
1 train_dataset = SentimentDataset(train_x, train_y, tokenizer)
2 val_dataset = SentimentDataset(val_x, val_y, tokenizer)
3 test_dataset = SentimentDataset(test_x, test_y, tokenizer)
```

## Model Training

We use Trainer from the transformers library to train our model with the following arguments:

```
1 from transformers import TrainingArguments, Trainer
2 from sklearn.metrics import accuracy_score
3
4 # Define the metrics
5 def compute_metrics(eval_pred):
6     logits, labels = eval_pred
7     preds = np.argmax(logits, axis=-1)
```

```
8     return {"accuracy": accuracy_score(labels, preds)}
```

```
9
```

```
10    # Hyperparameter configuration
```

```
11    training_args = TrainingArguments(
```

```
12        output_dir=".lora_roberta_sentiment",
```

```
13        evaluation_strategy="epoch",
```

```
14        save_strategy="epoch",
```

```
15        learning_rate=2e-5,
```

```
16        per_device_train_batch_size=8,
```

```
17        per_device_eval_batch_size=8,
```

```
18        num_train_epochs=3,
```

```
19        weight_decay=0.01,
```

```
20        logging_steps=10,
```

```
21        load_best_model_at_end=True,
```

```
22    )
```

```
23
```

```
24    # Initialize the Trainer
```

```
25    trainer = Trainer(
```

```
26        model=model,
```

```
27        args=training_args,
```

```
28        train_dataset=train_dataset,
```

```
29        eval_dataset=val_dataset,
```

```
30        compute_metrics=compute_metrics,
```

```
31        tokenizer=tokenizer,
```

```
32    )
```

```
33
```

```
34    # Train the model
```

```
35    import time
```

```
36    start_time = time.time()
```

```
37    trainer.train()
```

```
38    training_time = time.time() - start_time
```

```
39    print(f"\nTraining completed in {training_time:.2f} seconds")
```

## Evaluation and Results

We evaluate the model on the validation and test sets by analyzing the classification metrics.

```
1 from sklearn.metrics import classification_report, confusion_matrix
```

```
2 import seaborn as sns
```

```
3 import matplotlib.pyplot as plt
```

```
4
5 # Evaluation on the validation set
6 eval_results = trainer.evaluate()
7 print("\n--- Evaluation Results on Validation ---")
8 print(eval_results)

9
10 # Predictions on the test set
11 preds_output = trainer.predict(test_dataset)
12 predictions = np.argmax(preds_output.predictions, axis=-1)

13
14 accuracy = accuracy_score(test_y, predictions)
15 report = classification_report(test_y, predictions, target_names=["negative",
16   "neutral", "positive"])
17 cm = confusion_matrix(test_y, predictions)

18 print("\n--- Classification Report on Test ---")
19 print(report)
20 print(f"Accuracy: {accuracy:.4f}")

21
22 plt.figure(figsize=(6, 4))
23 sns.heatmap(cm, annot=True, fmt='d', cmap='Blues',
24   xticklabels=["negative", "neutral", "positive"],
25   yticklabels=["negative", "neutral", "positive"])
26 plt.xlabel("Prediction")
27 plt.ylabel("Actual")
28 plt.title("Confusion Matrix - RoBERTa with LoRA (Test)")
29 plt.show()
```

Model	Accuracy	F1-score	Critical Error Rate	Time (s)
Fine-Tuning	74.22%	0.74	2.96%	1134

Table 3: Performance of the Fine-Tuning on the test set

## Conclusion

Thus, we observed that the best performance is achieved through fine-tuning the language model, allowing for better adaptation to the sentiment classification task. However, accuracy remains improvable and could be further optimized by exploring advanced text preprocessing strategies, such as improving data cleaning, expanding contractions or integrating specific linguistic features.