Project Name: Sentiment Analysis of comments. Goal: To predict sentiment of comments using Natural Language Processing **Dataset**: IIT - KGP Kshitij Al Hackathon dataset **About the Project**: With the increasing reliance on online platforms for reviews, feedback, and social media interactions, understanding customer sentiment has become a critical task for businesses across industries. Sentiment analysis allows companies to gain deeper insights into customer opinions, which can help improve products, services, and customer satisfaction. However, interpreting the true sentiment behind vast amounts of text data, such as reviews, comments, and social media posts, can be challenging due to the nuances of language, tone, and context. To address this, we are developing a Natural Language Processing (NLP) model that analyzes customer feedback and classifies the sentiment expressed in the text. By training the model on a variety of text data with labeled sentiments, it will be able to accurately predict whether a given review or comment is positive, negative, or neutral. This tool will help businesses quickly identify areas of concern and opportunities for improvement by analyzing large volumes of customer feedback. The ultimate goal is to provide actionable insights that can enhance decision-making, improve customer relationships, and drive business success. **Section 1: Data Collection** In [1]: import pandas as pd import numpy as np import warnings warnings.filterwarnings("ignore") In [2]: url = r"C:\Users\user\Desktop\Code Fragments\Machine Learning\Datasets\AI-Hackathon-test-data-set.csv" data = pd.read\_csv(url) data Out[2]: "Comments, Sentiment" **0** "I love the new features in the update.,Positive" "Fast and efficient customer support., Positive" "The product quality is amazing!, Positive" 2 "Fast and efficient customer support., Positive" am disappointed with the delivery., Neutral" ••• "Fast and efficient customer support., Positive" 19998 "The room was clean but noisy., Neutral" 19999 "The product quality is amazing!, Positive" 20000 "I am disappointed with the delivery., Negative" 20002 "The food was cold and tasteless., Negative" 20003 rows × 1 columns We have ~ 20000 rows in our dataset, but the dataset hasn't be converted into our desired dataframe format. We'll have to manipulate it manually. data["Sentiment"] = data['"Comments, Sentiment"'].apply(lambda x: x[-9:-1]) data["Comments"] = data['"Comments, Sentiment"'].apply(lambda x: x[1:-10]) data Out[5]: "Comments, Sentiment" Sentiment **Comments 0** "I love the new features in the update.,Positive" Positive I love the new features in the update. "Fast and efficient customer support., Positive" Fast and efficient customer support. Positive "The product quality is amazing!, Positive" The product quality is amazing! 2 Positive "Fast and efficient customer support., Positive" Fast and efficient customer support. Positive am disappointed with the delivery., Neutral" I am disappointed with the delivery ,Neutral ••• 19998 "Fast and efficient customer support., Positive" Positive Fast and efficient customer support. 19999 "The room was clean but noisy., Neutral" ,Neutral The room was clean but noisy "The product quality is amazing!, Positive" The product quality is amazing! 20000 Positive I am disappointed with the delivery. "I am disappointed with the delivery., Negative" Negative "The food was cold and tasteless., Negative" The food was cold and tasteless. 20002 Negative 20003 rows × 3 columns Now that we have both the comments & sentiments separated, we can drop the attached column. data.drop(columns = ['"Comments, Sentiment"'], inplace = True) data Out[7]: **Sentiment Comments** Positive I love the new features in the update. Positive Fast and efficient customer support. The product quality is amazing! Positive Fast and efficient customer support. I am disappointed with the delivery ,Neutral ••• Fast and efficient customer support. 19998 Positive The room was clean but noisy 19999 ,Neutral Positive The product quality is amazing! 20000 I am disappointed with the delivery. 20001 Negative The food was cold and tasteless. 20002 Negative 20003 rows × 2 columns Since our dataframe is ready, lets move on to our data manipulation phase. Section 2: Data Cleaning / Manipulation Initially, we can encode our sentiment values. In  $[9]: d = {$ "Positive": 1, "Negative": -1, ",Neutral": 0 data["Sentiment"] = data["Sentiment"].map(d) Lets start with dealing with null values. In [11]: print(data[data["Comments"].isnull()]) print(data[data["Sentiment"].isnull()]) print(data[data['Comments'] == '']) Empty DataFrame Columns: [Sentiment, Comments] Index: [] Empty DataFrame Columns: [Sentiment, Comments] Index: [] Sentiment Comments 86 **153** 180 181 202 . . . 19854 19903 19941 19952 19993 [1001 rows x 2 columns] In [13]: data = data[data["Comments"]!=''] As observed, our dataset has no null values. One thing commonly observed in textual datasets is redundancy, lets deal with redundancy now. In [15]: data["Sentiment"].value\_counts() Out[15]: Sentiment 7461 7358 4183 Name: count, dtype: int64 As we can see, there is an deficit of neutral values. This may lead to underfitting / bias. We have to deal with this inbalance. data["Comments"].value\_counts() Out[17]: Comments I love the new features in the update. 1827 The product quality is amazing! 1811 Fast and efficient customer support. 1805 The app interface is confusing. 1801 I am disappointed with the delivery. 1799 Excellent ambiance and friendly staff. 1796 The food was cold and tasteless. 1793 Very poor customer service experience. 1793 The service was okay, not great 1712 The room was clean but noisy 1711 The room was clean but noisy. 200 The service was okay, not great. 193 I am disappointed with the delivery 106 Excellent ambiance and friendly staff 104 Very poor customer service experience 100 The food was cold and tasteless 95 I love the new features in the update 91 Fast and efficient customer support 90 The app interface is confusing The product quality is amazing 86 Blah blah blah ###random-text-12 Name: count, dtype: int64 As we can see, our dataset is extremely poor in terms of redundancy & annotations. With the occurence frequency in mind, the Countering the problems -• redundancy: we do not try to solve the redundancy since dropping duplicates greatly reduces data size • annotations: for same data being annotated differently, we can adpot the annotation assigned to majority of the sub-data. • imbalance: we train using 4k values each, and then fine tune while adding left over data. **Section 2.1: Addressing Annotation Errors** In [19]: unique\_coms = data['Comments'].unique() unique\_coms Out[19]: array(['I love the new features in the update.', 'Fast and efficient customer support.', 'The product quality is amazing!', 'I am disappointed with the delivery', 'The service was okay, not great', 'The food was cold and tasteless.', 'I am disappointed with the delivery.', 'Excellent ambiance and friendly staff.', 'Very poor customer service experience.', 'The room was clean but noisy.', 'The app interface is confusing.', 'The room was clean but noisy', 'Excellent ambiance and friendly staff', 'Very poor customer service experience', 'The product quality is amazing', 'The food was cold and tasteless', 'The service was okay, not great.', 'Fast and efficient customer support', 'The app interface is confusing', 'I love the new features in the update', 'Blah blah blah', '###random-text-12'], dtype=object) We will proceed with removing the annotation inconsistency, by selecting the one with higher frequency. In [21]: for i in unique\_coms: f1 = len(data[data["Comments"] == i]) f2 = len(data[data["Comments"] == i+"."]) **if** f1>f2: data.drop(data[data['Comments'] == i+"."].index, inplace=True) else: data.drop(data[data['Comments'] == i].index, inplace=True) In [23]: data['Comments'].unique() Out[23]: array(['I love the new features in the update.', 'Fast and efficient customer support.', 'The product quality is amazing!', 'The service was okay, not great', 'The food was cold and tasteless.', 'I am disappointed with the delivery.', 'Excellent ambiance and friendly staff.', 'Very poor customer service experience.', 'The app interface is confusing.', 'The room was clean but noisy', 'The product quality is amazing', 'Blah blah blah', '###random-text-12'], dtype=object) Now that our annotation issues are clear, we can move on data manipulation techniques for textual data. We do the following: remove stop words lemmatize basic formatting In [25]: import re **from** nltk.corpus **import** stopwords from nltk.stem import WordNetLemmatizer stop\_words = set(stopwords.words('english')) lemmatizer = WordNetLemmatizer() def process(text): text = text.lower() text = re.sub(r'[^\w\s]', '', text) words = text.split() words = [lemmatizer.lemmatize(word) for word in words if word not in stop\_words] return ' '.join(words) data["Comments"] = data["Comments"].apply(lambda x: process(x)) we create three parts of the dataset with different sentiments In [27]: pos\_data = data[data["Sentiment"] == 1] neg\_data = data[data["Sentiment"] == -1] net\_data = data[data["Sentiment"] == 0] We merge the three data parts into one, while leaving out parts of positive & negative data to ensure equal representation to neutral data In [29]: use\_pos = pos\_data[:4000] use\_neg = neg\_data[:4000] data1 = pd.concat([use\_pos, use\_neg, net\_data]) We shuffle the data to maintain equal representation of data in all parts [ Concept of Stratified Shuffle Split ] In [31]: from sklearn.utils import shuffle data1 = shuffle(data1, random state = 42) data1 Out[31]: **Sentiment Comments** app interface confusing -1 1785 **7532** 1 fast efficient customer support 9761 disappointed delivery -1 product quality amazing 9807 0 11814 room clean noisy 0 18833 room clean noisy 0 app interface confusing 3345 -1 3865 -1 food cold tasteless 2255 app interface confusing 9132 food cold tasteless -1 11510 rows  $\times$  2 columns In [33]: x = data1["Comments"]y = data1["Sentiment"]  $x_{train} = x[:int(len(x)*0.75)]$  $y_{train} = y[:int(len(y)*0.75)]$  $x_{\text{test}} = x[int(len(x)*0.75):]$  $y_{test} = y[int(len(y)*0.75):]$ In [35]: x\_train[3] Out[35]: 'fast efficient customer support' Our dataset is now ready for training phase 1. Section 3: Model selection Google's bert is a very popular & relaible model. We will procees with it's use for our goal. In [37]: **from** transformers **import** BertTokenizer **as** BT, TFBertForSequenceClassification **as** tbsc from sklearn.metrics import classification\_report as cr WARNING:tensorflow:From C:\Users\user\anaconda3\Lib\site-packages\tf\_keras\src\losses.py:2976: The name tf.losses.sparse\_softmax\_cross\_entropy is deprecate d. Please use tf.compat.v1.losses.sparse\_softmax\_cross\_entropy instead. We use TFBertSC for our classification task and BertTokenizer of tokenizatiion of text. In [38]: tokenizer = BT.from\_pretrained('bert-base-uncased', do\_lower\_case=True) x\_train\_en = tokenizer.batch\_encode\_plus(x\_train.astype(str).tolist(), padding=True, truncation=True, max length = 128, return\_tensors='tf') x\_test\_en= tokenizer.batch\_encode\_plus(x\_test.astype(str).tolist(), padding=True, truncation=True, max\_length = 128, return\_tensors='tf') In [41]: model = tbsc.from\_pretrained('bert-base-uncased', num\_labels=3) WARNING:tensorflow:From C:\Users\user\anaconda3\Lib\site-packages\tf\_keras\src\backend.py:873: The name tf.get\_default\_graph is deprecated. Please use tf.co mpat.v1.get\_default\_graph instead. All PyTorch model weights were used when initializing TFBertForSequenceClassification. Some weights or buffers of the TF 2.0 model TFBertForSequenceClassification were not initialized from the PyTorch model and are newly initialized: ['classif ier.weight', 'classifier.bias'] You should probably TRAIN this model on a down-stream task to be able to use it for predictions and inference. As we have three classifications, we set num\_labels to 3. This allows the model to predict values between [0-3) To avoid errors, we add 1 to all values of y\_train & y\_test (since minimum value is -1) In [43]: y\_train = y\_train + 1 y\_test = y\_test + 1 Now, we compile our model using an optimizer, specifying criteria of loss & accuracy. In [45]: **import** tensorflow **as** tf optimizer = tf.keras.optimizers.Adam(learning\_rate=2e-5) loss = tf.keras.losses.SparseCategoricalCrossentropy(from\_logits=True) metric = tf.keras.metrics.SparseCategoricalAccuracy('accuracy') model.compile(optimizer=optimizer, loss=loss, metrics=[metric]) We finally fit our data into the model. In [47]: ''' model.fit( [x\_train\_en['input\_ids'], x\_train\_en['token\_type\_ids'], x\_train\_en['attention\_mask']], y\_train, validation\_data=( [x\_test\_en['input\_ids'], x\_test\_en['token\_type\_ids'], x\_test\_en['attention\_mask']],y\_test), batch\_size=32, epochs=3 1 1 1 [x\_train\_en['input\_ids'], x\_train\_en['token\_type\_ids'], x\_train\_en['attention\_mask']],\n "\nmodel.fit(\n y\_train,\n validation\_data=(\n Out[47]: [x\_t est\_en['input\_ids'], x\_test\_en['token\_type\_ids'], x\_test\_en['attention\_mask']],y\_test),\n batch\_size=32,\n epochs=3\n)\n" Since the training time is large, and no availability of an accelerator in this environment, the model has been trained using a GPU in a kaggle notebook (https://www.kaggle.com/code/gouravjana/notebookd86c597d7a/notebook?scriptVersionId=218325013) as an continuation to this project.

The next few code blocks would just be repetition of a few steps from the original notebook, Please skip to continue from the new additions

```
In [1]: import pandas as pd
        import numpy as np
        data = pd.read_csv("/kaggle/input/aihsentiment/AI-Hackathon-test-data-set-3.csv")
        pos_data = data[data["Sentiment"] == 1]
        neg_data = data[data["Sentiment"] == -1]
        net_data = data[data["Sentiment"] == 0]
        use_pos = pos_data[:4000]
        use_neg = neg_data[:4000]
        data1 = pd.concat([use_pos, use_neg, net_data])
        from sklearn.utils import shuffle
        data1 = shuffle(data1, random_state = 42)
        x = data1["Comments"]
        y = data1["Sentiment"]
        x_{train} = x[:int(len(x)*0.75)]
        y_{train} = y[:int(len(y)*0.75)]
        x_{test} = x[int(len(x)*0.75):]
        y_{test} = y[int(len(y)*0.75):]
        from transformers import BertTokenizer as BT, TFBertForSequenceClassification as tbsc
        from sklearn.metrics import classification report as cr
        tokenizer = BT.from_pretrained('bert-base-uncased', do_lower_case=True)
        x_train_en = tokenizer.batch_encode_plus(x_train.astype(str).tolist(),
                                                       padding=True,
                                                       truncation=True,
                                                       max_length = 128,
                                                       return tensors='tf')
        x_test_en= tokenizer.batch_encode_plus(x_test.astype(str).tolist(),
                                                       padding=True,
                                                       truncation=True,
                                                       max_length = 128,
                                                       return tensors='tf')
        model = tbsc.from_pretrained('bert-base-uncased', num_labels=3)
        y_{train} = y_{train} + 1
        y_{test} = y_{test} + 1
        import tensorflow as tf
        optimizer = tf.keras.optimizers.Adam(learning_rate=2e-5)
        loss = tf.keras.losses.SparseCategoricalCrossentropy(from_logits=True)
        metric = tf.keras.metrics.SparseCategoricalAccuracy('accuracy')
        model.compile(optimizer=optimizer, loss=loss, metrics=[metric])
       tokenizer_config.json:
                                0%
                                               0.00/48.0 [00:00<?, ?B/s]
       vocab.txt: 0%
                                   0.00/232k [00:00<?, ?B/s]
       tokenizer.json: 0%
                                        0.00/466k [00:00<?, ?B/s]
```

This is where the continuation of the notebook starts from

0.00/570 [00:00<?, ?B/s]

All PyTorch model weights were used when initializing TFBertForSequenceClassification.

0.00/440M [00:00<?, ?B/s]

You should probably TRAIN this model on a down-stream task to be able to use it for predictions and inference.

We start with the model training

ier.weight', 'classifier.bias']

config.json: 0%

model.safetensors: 0%

```
In [2]: model.fit(
     [x_train_en['input_ids'], x_train_en['token_type_ids'], x_train_en['attention_mask']],
     y_train,
     validation_data=(
      [x_test_en['input_ids'], x_test_en['token_type_ids'], x_test_en['attention_mask']],y_test),
     batch_size=32,
     epochs=3
   Epoch 1/3
   Epoch 2/3
   Epoch 3/3
   Out[2]: <tf_keras.src.callbacks.History at 0x7915b8aa3e50>
```

Some weights or buffers of the TF 2.0 model TFBertForSequenceClassification were not initialized from the PyTorch model and are newly initialized: ['classif

As the evaluation shows, our model gains an accuracy of about ~96%.

Lets evaluate our model on the test data.

```
In [3]: test_loss, test_accuracy = model.evaluate(
            [x_test_en['input_ids'], x_test_en['token_type_ids'], x_test_en['attention_mask']],
            y_test
        print(f'Test loss: {test_loss}, Test accuracy: {test_accuracy}')
```

Test loss: 0.1621810644865036, Test accuracy: 0.9523965716362

As we can see, our model performs well on the test data as well. We can proceed with the use of this model for our goal.

## Section 4: User interactive space

dataset itself.

Here, we allow the model to be used for different data.

**Note**: The performance of this model on real world data is extremely poor,

due to the source dataset having effectively only 22 data points. As the motive of the hackathon's organizers was to test the preprocessing abilities of the participants, this model's performance is relevant in the scope of evaluating against the

```
In [4]: while True:
            print("Enter the array of reviews you want to predict.")
            arr = list(input().rstrip().split(", "))
            df = pd.DataFrame(arr, columns=['Comments'])
            df en = tokenizer.batch_encode_plus(x_train.astype(str).tolist(),
                                                 padding=True,
                                                truncation=True,
                                                max_length = 128,
                                                return_tensors='tf')
            pred = model.predict(
                [df_en['input_ids'],df_en['token_type_ids'], df_en['attention_mask']])
            logits = pred.logits
            pred_labels = tf.argmax(logits, axis=1)
            pred_labels = pred_labels.numpy()
            label = {
                2: 'Positive',
                0: 'Negative',
                1: 'Neutral'
            pred_labels = [label[i] for i in pred_labels]
            print("Predicted sentiments are - \n")
            for i in range(len(arr)):
                print(arr[i],":",pred_labels[i])
            print("Do you wish to continue?\n1 - Yes\n2 - No")
            if int(input()) == 2:
                break
        print("Thank you.")
```

```
Enter the array of reviews you want to predict.
286/286 [============== ] - 12s 33ms/step
Predicted sentiments are -
Good product : Positive
Do you wish to continue?
2 - No
Thank you.
```

This marks the end of our project.

## Conclusion

This project is aimed at determining the sentiment of reviews of products accurately.

The true core challenge of the project was to counter the challenges faced by having a poor dataset.

Challenges faces in the dataset -

- Heavy redundancy (>98%), causing no real input value in the dataset (the dataset had 22 rows post dropping duplicate values, which originally had 20000 rows)
- Wrong, with low variety in data in the first place, different instances of identical features having different label added to the loss of value in the data

Due to the above reasons (very low variety in data), the model trained on the same is essentially not of practical use.

With the challenge being to make the model accurate for data within the dataset itself,

our model gained an accuracy of ~96%, which originally hit an bottleneck accuracy of 60% prior to all preprocessing ideas being implemented. With that, the goal has been met & the project is complete.

What makes this project different (& better)? -> Accuracy of ~95% compared to the bottleneck average of 60%

What was the main factor behind the accuracy increase? Resolving massive data annotation errors in the original dataset.

Due to the inconsistent annotation, fitting the model with the entire dataset couldn't reproduce a high accuracy. The reason being the same data being present multiple times, with different annotations.