# AMAZON SENTIMENT ANALYSIS

**UEH** | College of Technology and Design
**UNIVERSITY**

**Học phần: NLP**
**Mã lớp: 23C1INF50907601**
**Giảng viên: TS. Đặng Ngọc Hoàng Thành**

# PREPROCESS & EDA

# DATA FORMAT

```
__label__1 Batteries died within a year ...: I bought this charger in Jul 2003 and it worked OK for a while. The design is nice and
convenient. However, after about a year, the batteries would not hold a charge. Might as well just get alkaline disposables, or look
elsewhere for a charger that comes with batteries that have better staying power.
__label__2 works fine, but Maha Energy is better: Check out Maha Energy's website. Their Powerex MH-C204F charger works in 100 minutes for
rapid charge, with option for slower charge (better for batteries). And they have 2200 mAh batteries.
```

**__label__1: negative**

**__label__2: postive**

# PREPROCESS

**A. TIỀN XỬ LÝ**

- Bước 1: Loại bỏ các đường link URL
- Bước 2: Loại bỏ các thẻ HTML
- Bước 3: Mở rộng các từ viết tắt
- Bước 4: Loại bỏ dấu chấm câu và chữ số
- Bước 5: Loại bỏ emoji
- Bước 6: Chuyển sang chữ thường
- Bước 7: Loại bỏ stopwords khỏi các câu

# PREPROCESS

**A. TIỀN XỬ LÝ**

| | text | label | preprocess_sentence |
|---|---|---|---|
| **0** | dare you to finish this book: I dare you to fi... | 0 | dare finish book dare finish book simple inter... |
| **1** | Amazon should not be promoting "hate" stuff li... | 0 | amazon promoting hate stuff like hate jews hom... |
| **2** | Herbie plays it easy-listening: As an HH great... | 0 | herbie plays easy listening hh great fan disap... |
| **3** | Freddie Will Live 4Ever: As long as there are ... | 1 | freddie live 4ever long humans electricity con... |
| **4** | Cord description 100% incorrect: New in origin... | 0 | cord description 100 incorrect new original pa... |
| **...** | ... | ... | ... |
| **39995** | Powdered wax is a bad idea: My daughter receiv... | 0 | powdered wax bad idea daughter received one tw... |
| **39996** | WHAT was the Network Thinking??: Answer...they... | 1 | network thinking answer terrific show yanked 4... |
| **39997** | After Shave Cream: This product is a great fol... | 1 | shave cream product great follow close shave v... |
| **39998** | Zzzzzzzzzzzzzzzzzzzz: I thought I was gonna fa... | 0 | zzzzzzzzzzzzzzzzzzzz thought going fall asleep... |
| **39999** | Put together by a blind man: I received the bo... | 0 | put together blind man received book quickly h... |

40000 rows × 3 columns

Train dataframe của nhóm

# PREPROCESS

**A. TIỀN XỬ LÝ**

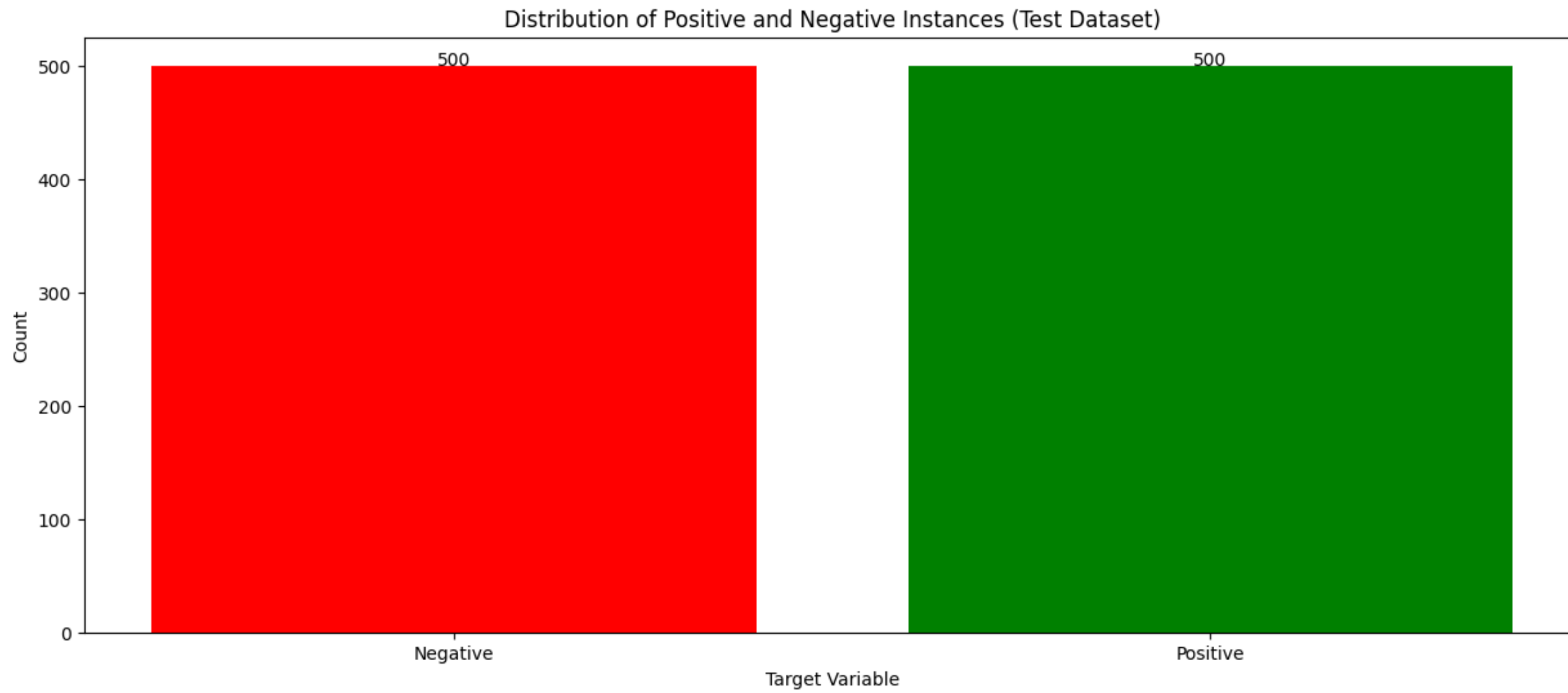| | text | label | preprocess_sentence |
|---|---|---|---|
| **0** | Man, this is sick stuff!!!!: I set out on a mi... | 0 | man sick stuff set mission seek shocking films... |
| **1** | Still waiting for hardware: Ordered (1) chair ... | 0 | still waiting hardware ordered 1 chair 1 loves... |
| **2** | OK Idea, Poorly Executed: Julia was good. Brad... | 0 | ok idea poorly executed julia good brad usual ... |
| **3** | Don't bother: This is a total waste of paper. ... | 0 | bother total waste paper save forrest print bo... |
| **4** | Good, but disappointing: This was a good movie... | 0 | good disappointing good movie general however ... |
| **...** | ... | ... | ... |
| **995** | no problemo: Our school has 100 of the LaCie P... | 1 | problemo school 100 lacie p3 xp never problem ... |
| **996** | Thank you Decapitated!: This album ironically ... | 1 | thank decapitated album ironically makes hope ... |
| **997** | Super item!: This vest came in handy for a my ... | 1 | super item vest came handy son birthday party ... |
| **998** | Great for the price: The car charger works as ... | 1 | great price car charger works expected ear bud... |
| **999** | Mosher Wows Us Again: Howard Frank Mosher sets... | 1 | mosher wows us howard frank mosher sets stage ... |

1000 rows × 3 columns

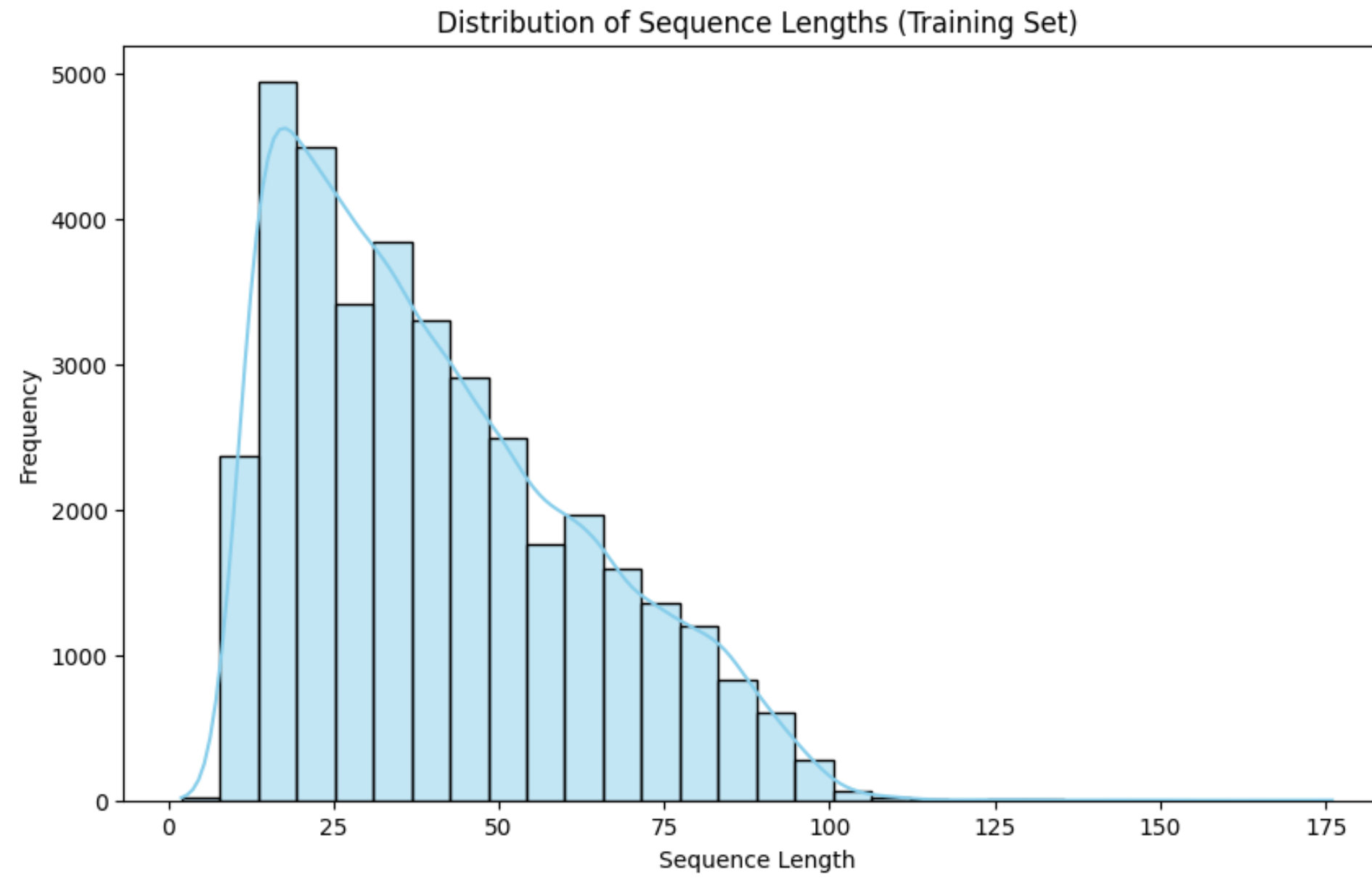Test dataframe của nhóm

# PREPROCESS

**B. LOẠI BỎ GÁNH NẶNG**

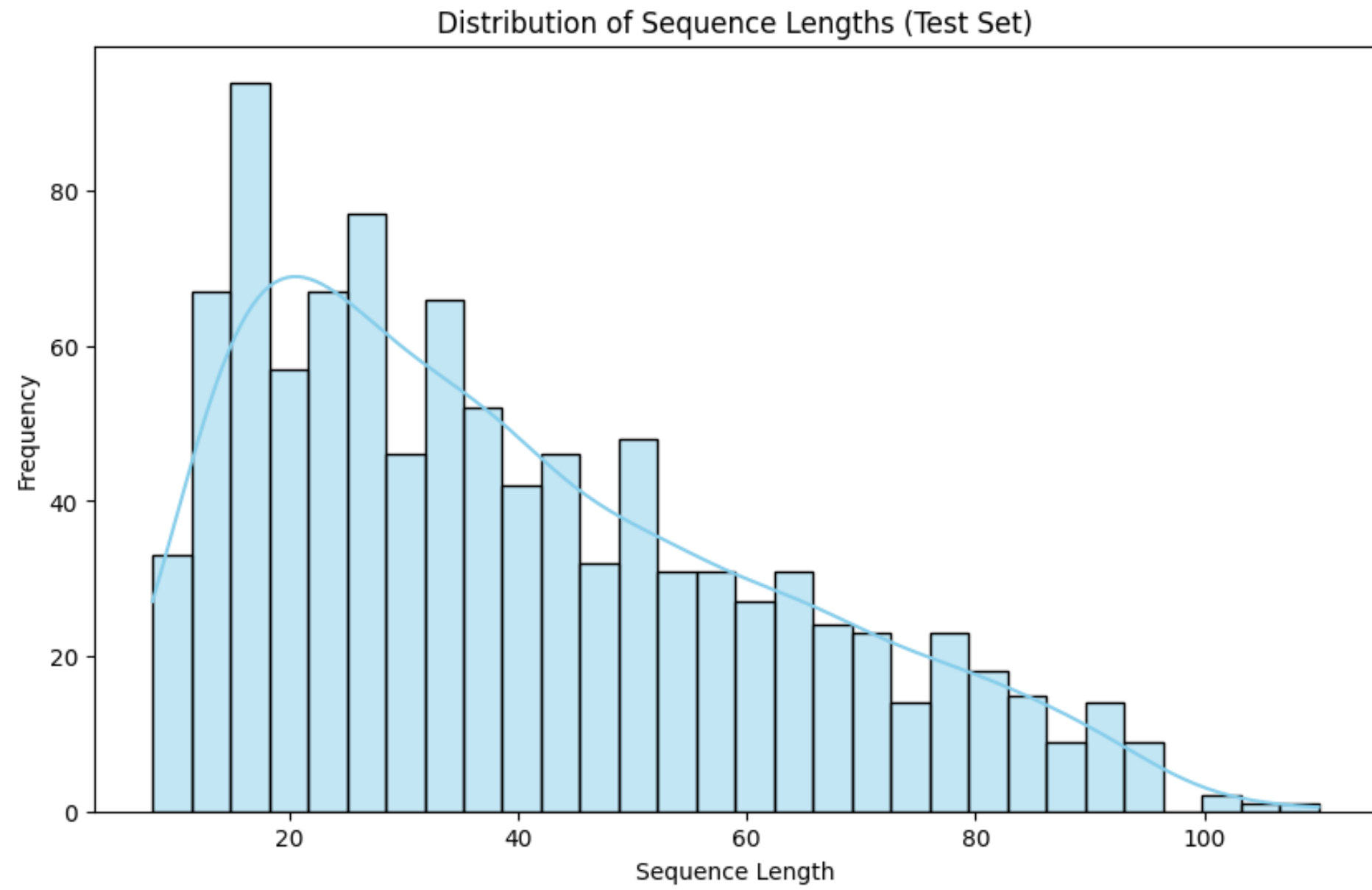| | text | label | preprocess_sentence |
|---|---|---|---|
| **0** | dare finish book dare finish book simple inter... | 0 | dare finish book dare finish book simple inter... |
| **1** | amazon promoting hate stuff like hate jews hom... | 0 | amazon promoting hate stuff like hate jews hom... |
| **2** | herbie plays easy listening hh great fan disap... | 0 | herbie plays easy listening hh great fan disap... |
| **3** | freddie live 4ever long humans electricity con... | 1 | freddie live 4ever long humans electricity con... |
| **4** | cord description 100 incorrect new original pa... | 0 | cord description 100 incorrect new original pa... |
| **...** | ... | ... | ... |
| **39429** | powdered wax bad idea daughter received one tw... | 0 | powdered wax bad idea daughter received one tw... |
| **39430** | network thinking answer terrific show yanked 4... | 1 | network thinking answer terrific show yanked 4... |
| **39431** | shave cream product great follow close shave v... | 1 | shave cream product great follow close shave v... |
| **39432** | zzzzzzzzzzzzzzzzzzzzz thought going fall asleep... | 0 | zzzzzzzzzzzzzzzzzzzzz thought going fall asleep... |
| **39433** | put together blind man received book quickly h... | 0 | put together blind man received book quickly h... |

39434 rows × 3 columns

Train dataframe sau khi loại bỏ các mẫu nhãn sai

# EDA



Distribution of Positive and Negative Instances (Train Dataset)

# EDA



Distribution of Positive and Negative Instances (Test Dataset)

# EDA



Distribution of Sequence Lengths (Training Set)

# EDA



Distribution of Sequence Lengths (Test Set)

```python
1 # Sử dụng cột 'preprocess_sentence' làm đặc trưng và 'label' làm mục tiêu
2 X_train = train_df['preprocess_sentence']
3 y_train = train_df['label']
4 X_test = test_df['text']
5 y_test = test_df['label']
```

# NAIVE BAYES

```python
1 # Cài đặt Pipeline với tham số chọn từ GridSearch
2 model_NB = Pipeline([
3     ('tfidf', TfidfVectorizer(min_df=1, ngram_range=(1, 2))),
4     ('nb', MultinomialNB())
5 ])
6
7 # Huấn luyện mô hình trên tập huấn luyện
8 model_NB.fit(X_train, y_train)
9
10 # Dự đoán nhãn trên tập kiểm thử
11 y_pred = model_NB.predict(X_test)
12 predicted_probabilities = model_NB.predict_proba(X_test)
13
14 # Đánh giá mô hình
15 accuracy = accuracy_score(y_test, y_pred)
16 report = classification_report(y_test, y_pred)
17
18 print(f"Accuracy: {accuracy}")
19 print(report)
```

# LOGISTIC REGRESSION

```python
1  # Các siêu tham số
2  best_parameters = {'logisticregression__C': 100, 'tfidfvectorizer__ngram_range': (1, 2)}
3
4  # Cài đặt Pipeline với tham số chọn từ GridSearch
5
6  model_LR = Pipeline([
7      ('tfidfvectorizer', TfidfVectorizer(ngram_range=best_parameters['tfidfvectorizer__ngram_range'])),
8      ('logisticregression', LogisticRegression(C=best_parameters['logisticregression__C']))
9  ])
10
11 # Huấn luyện mô hình trên tập huấn luyện
12 model_LR.fit(X_train, y_train)
13
14 # Dự đoán nhãn trên tập kiểm thử
15 y_pred = model_LR.predict(X_test)
16
17 # Đánh giá mô hình
18 accuracy = accuracy_score(y_test, y_pred)
19 report = classification_report(y_test, y_pred)
20
21 print(f"Accuracy: {accuracy}")
22 print(report)
23
```

# TRANSFORMER (ENCODER)

```
1 model_TC
```

```
TransformerEncoderCls(
  (embd_layer): TokenAndPositionEmbedding(
    (word_emb): Embedding(10000, 200)
    (pos_emb): Embedding(100, 200)
  )
  (transformer_layer): TransformerEncoder(
    (attn): MultiheadAttention(
      (out_proj): NonDynamicallyQuantizableLinear(in_features=200, out_features=200, bias=True)
    )
    (ffn): Sequential(
      (0): Linear(in_features=200, out_features=128, bias=True)
      (1): ReLU()
      (2): Linear(in_features=128, out_features=200, bias=True)
    )
    (layernorm_1): LayerNorm((200,), eps=1e-06, elementwise_affine=True)
    (layernorm_2): LayerNorm((200,), eps=1e-06, elementwise_affine=True)
    (dropout_1): Dropout(p=0.1, inplace=False)
    (dropout_2): Dropout(p=0.1, inplace=False)
  )
  (pooling): AvgPool1d(kernel_size=(100,), stride=(100,), padding=(0,))
  (fc1): Linear(in_features=200, out_features=20, bias=True)
  (fc2): Linear(in_features=20, out_features=2, bias=True)
  (dropout): Dropout(p=0.1, inplace=False)
  (relu): ReLU()
)
```
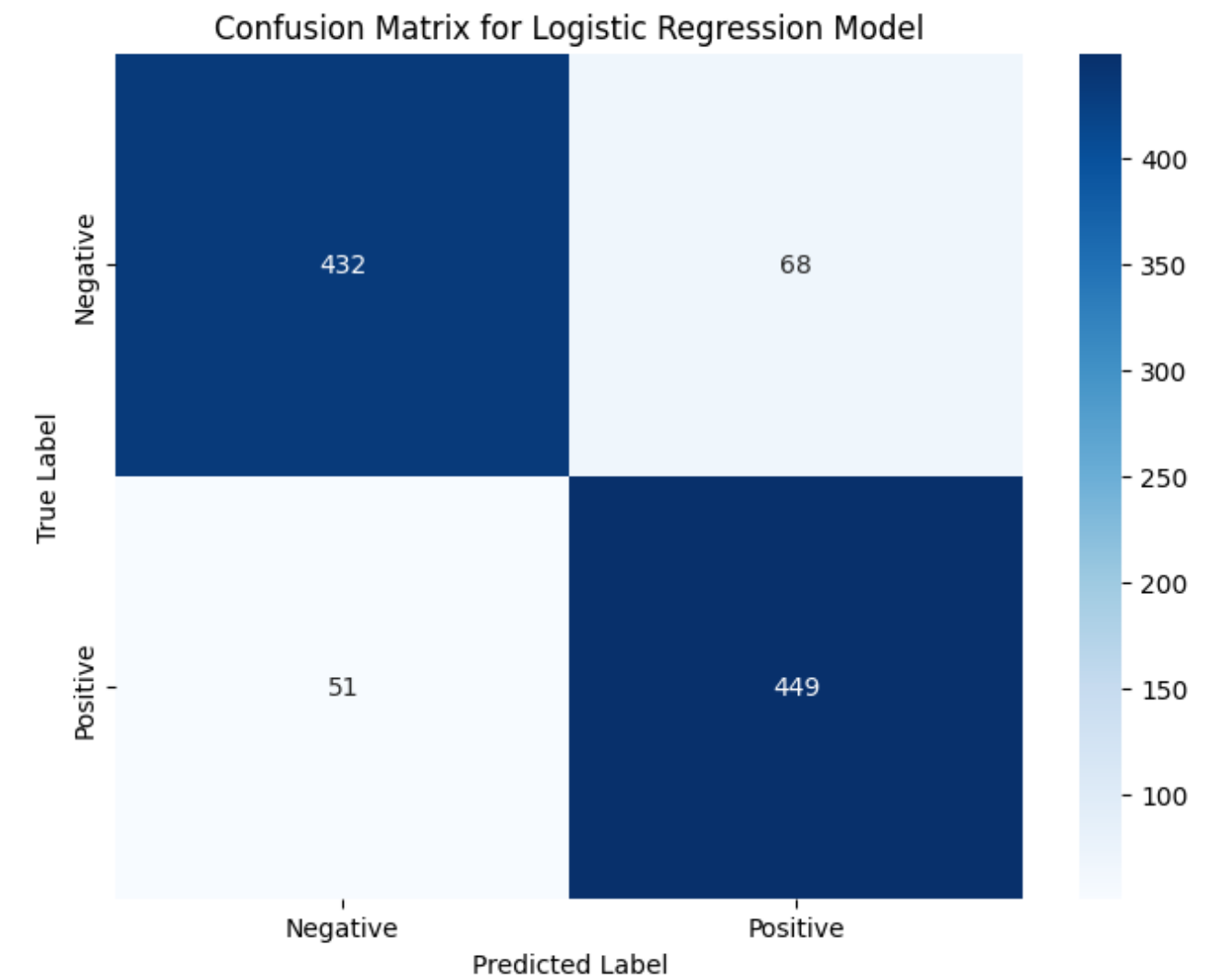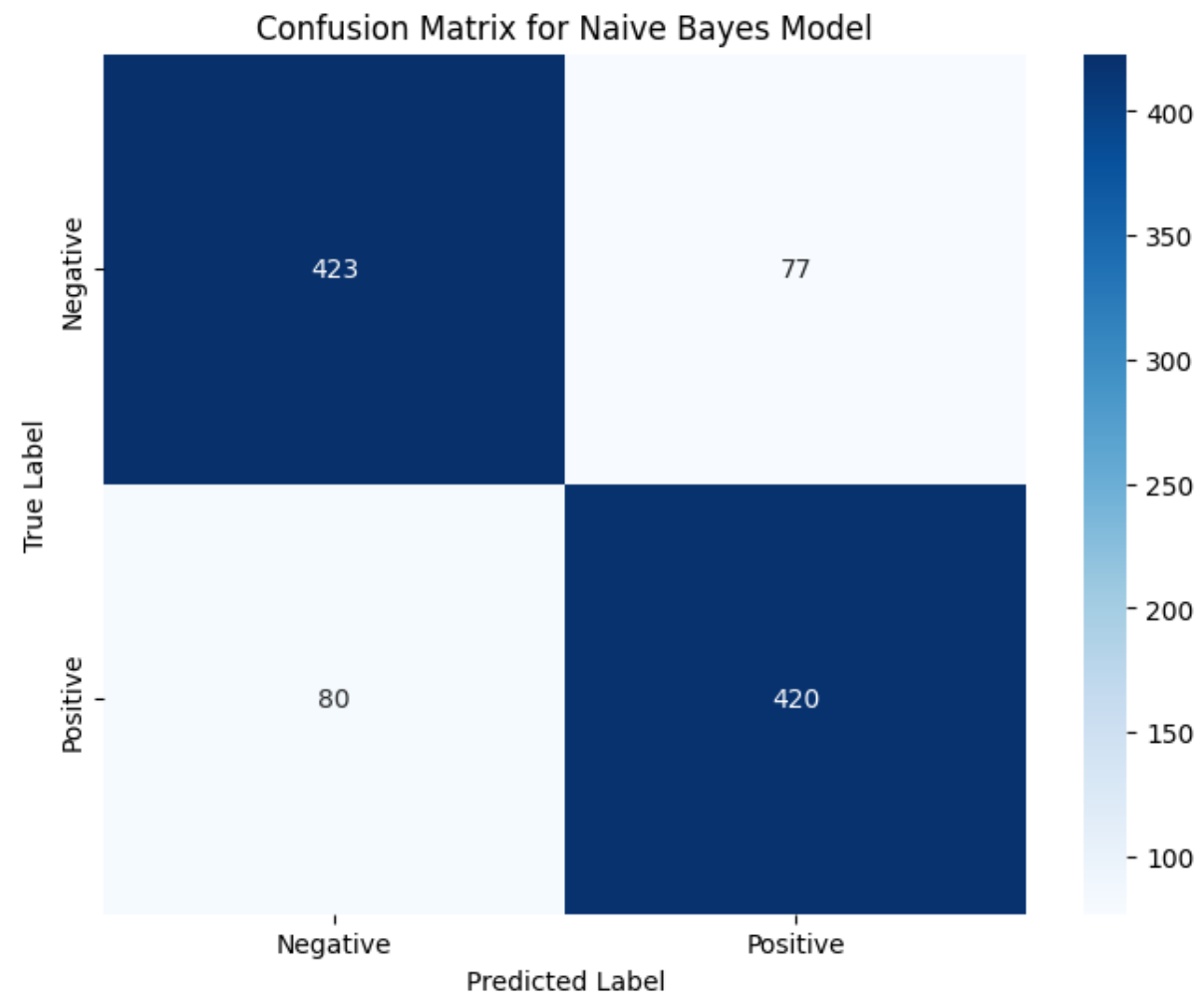
# DistilBERT

```
DistilBERTClass(
  (l1): DistilBertModel(
    (embeddings): Embeddings(
      (word_embeddings): Embedding(30522, 768, padding_idx=0)
      (position_embeddings): Embedding(512, 768)
      (LayerNorm): LayerNorm((768,), eps=1e-12, elementwise_affine=True)
      (dropout): Dropout(p=0.1, inplace=False)
    )
    (transformer): Transformer(
      (layer): ModuleList(
        (0-5): 6 x TransformerBlock(
          (attention): MultiHeadSelfAttention(
            (dropout): Dropout(p=0.1, inplace=False)
            (q_lin): Linear(in_features=768, out_features=768, bias=True)
            (k_lin): Linear(in_features=768, out_features=768, bias=True)
            (v_lin): Linear(in_features=768, out_features=768, bias=True)
            (out_lin): Linear(in_features=768, out_features=768, bias=True)
          )
          (sa_layer_norm): LayerNorm((768,), eps=1e-12, elementwise_affine=True)
          (ffn): FFN(
            (dropout): Dropout(p=0.1, inplace=False)
            (lin1): Linear(in_features=768, out_features=3072, bias=True)
            (lin2): Linear(in_features=3072, out_features=768, bias=True)
            (activation): GELUActivation()
          )
          (output_layer_norm): LayerNorm((768,), eps=1e-12, elementwise_affine=True)
        )
      )
    )
  )
  (pre_classifier): Linear(in_features=768, out_features=768, bias=True)
  (dropout): Dropout(p=0.1, inplace=False)
  (classifier): Linear(in_features=768, out_features=1, bias=True)
)
```
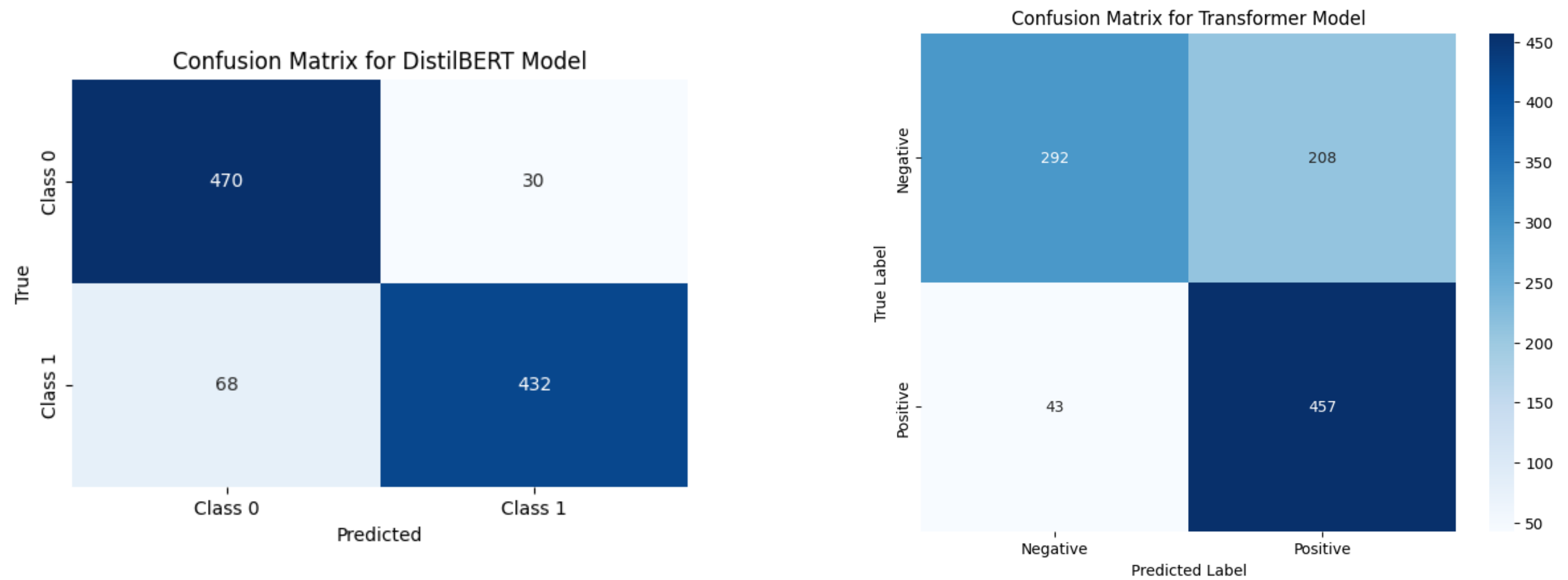
RESULT

# CONFUSION MATRIX

Confusion Matrix for Naive Bayes Model

|  | Negative | Positive |
|---|---|---|
| **Negative** | 423 | 77 |
| **Positive** | 80 | 420 |

True Label / Predicted Label

Confusion Matrix for Logistic Regression Model

|  | Negative | Positive |
|---|---|---|
| **Negative** | 432 | 68 |
| **Positive** | 51 | 449 |

True Label / Predicted Label

# CONFUSION MATRIX

# OTHER METRICS

| Model | Accuracy | Precision | Recall | F1-Score |
|---|---|---|---|---|
| **Model DB** | 0.90 | 0.94 | 0.86 | 0.90 |
| **Model LR** | 0.88 | 0.87 | 0.90 | 0.88 |
| **Model NB** | 0.84 | 0.85 | 0.84 | 0.84 |
| **Model TC** | 0.75 | 0.69 | 0.91 | 0.78 |

# CONCLUSIONS

Sau khi xem xét cả bốn chỉ số đánh giá để có cái nhìn toàn diện về hiệu suất của mỗi mô hình. Trong trường hợp này, mô hình finetune DistilBERT có vẻ là mô hình có hiệu suất tốt nhất dựa trên các chỉ số này.

# DEMO (Anvil)

# DEMO (HuggingFace)

# REFERENCES

1. https://huggingface.co/spaces/perman2011/UEH_SentimentAnalysis
2. https://understated-downright-contest.anvil.app/
3. https://colab.research.google.com/drive/1yRTOP5clrG9OmAglrLpvysv5w4C2gXE0#scrollTo=M0c7nrbF3_rR

UEH
UNIVERSITY

College of
Technology and Design

SCHOOL OF BUSINESS INFORMATION TECHNOLOGY

CHEERS