```
1  !pip install pandas numpy scikit-learn nltk textblob gensim spacy transformers matplotlib s
2  !python -m spacy download en_core_web_sm
3  !python -m nltk.downloader punkt wordnet omw-1.4 stopwords
```

Show hidden output

```
1
2
3
4
5  import pandas as pd
6  import numpy as np
7  import re
8  import nltk
9  import spacy
10  from textblob import TextBlob
11  from sklearn.feature_extraction.text import CountVectorizer, TfidfVectorizer
12  from sklearn.model_selection import train_test_split
13  from sklearn.metrics import accuracy_score, classification_report
14  import matplotlib.pyplot as plt
15  import seaborn as sns
16
17  # Deep Learning
18  import torch
19  from torch import nn
20  from torch.utils.data import DataLoader, TensorDataset
21  from transformers import BertTokenizer, BertForSequenceClassification, pipeline
22  from gensim.models import Word2Vec, FastText
23  from gensim.utils import simple_preprocess
24  import warnings
25  warnings.filterwarnings("ignore")
26
27  # NLTK & SpaCy
28  nltk.download('punkt')
29  nltk.download('punkt_tab')
30  nltk.download('stopwords')
31  nltk.download('wordnet')
32  nlp = spacy.load("en_core_web_sm")
33
34  # Set style
35  plt.style.use('seaborn-v0_8')
36  sns.set_palette("husl")
```

Show hidden output

## ⌄  1.Text Pre-processing on file.txt

```
1 # a. Import necessary libraries
2 import nltk
3 from nltk.tokenize import word_tokenize, sent_tokenize
4 from nltk.corpus import stopwords
5 from nltk.stem import PorterStemmer, WordNetLemmatizer
6 from textblob import TextBlob
7 import spacy
8
9 # Auto-create file.txt with improper text
10 file_content = """
11 Natural language processing (NLP) is a sub-field of artificial intelligence, concerned with
12
13 NLP is used in many applications such as: machine translation, sentiment analysis, speech
14
15 Some common challenges include: ambiguity in language, sarcasm, slang, and misspellings. Fo
16
17 Here are some sample sentences with errors:
18 - "Ths sentnce has speling erors."
19 - "I luv programing in pyhton!!!"
```

```
20 - "Artifical inteligence is the future..."
21
22 NLP techniques help clean and normalize such texts. Tokenization breaks text into words, s
23
24 Named entity recognition (NER) identifies names of people, organizations, locations. For e;
25
26 Finally, sentence boundary detection splits text into sentences. This is important for summ
27 """
28
29 with open("file.txt", "w", encoding="utf-8") as f:
30     f.write(file_content.strip())
31
32 print("file.txt created successfully!")
33 nltk.download('punkt_tab')
34 nltk.download('averaged_perceptron_tagger_eng')
35 nltk.download('punkt')
36 nltk.download('stopwords')
37 nltk.download('wordnet')
38 nltk.download('averaged_perceptron_tagger')
39
40 nlp = spacy.load("en_core_web_sm")
41
42 # b. Load the text corpus
43 with open("file.txt", "r", encoding="utf-8") as f:
44     text = f.read()
45
46 print("Original text length:", len(text))
47 print(text[:500], "\n")
48
49 # c. Tokenization
50 tokens = word_tokenize(text.lower())
51 print("First 30 tokens:", tokens[:30])
52
53 # d. Spelling correction
54 corrected_tokens = []
55 for token in tokens:
56     corrected = str(TextBlob(token).correct())
57     corrected_tokens.append(corrected)
58
59 corrected_text = " ".join(corrected_tokens)
60 print("\nFirst 10 corrected tokens:", corrected_tokens[:10])
61 print("\nCorrected text (first 200 chars):", corrected_text[:200])
62
63 # e. POS tagging on corrected tokens
64 pos_tags = nltk.pos_tag(corrected_tokens)
65 print("\nPOS tags (first 10):", pos_tags[:10])
66
67 # f. Remove stop words
68 stop_words = set(stopwords.words("english"))
69 filtered_tokens = [t for t in corrected_tokens if t not in stop_words and t.isalpha()]
70 print("\nFirst 20 tokens after stop-word removal:", filtered_tokens[:20])
71
72 # g. Stemming & Lemmatization
73 stemmer = PorterStemmer()
74 lemmatizer = WordNetLemmatizer()
75
76 stemmed = [stemmer.stem(t) for t in filtered_tokens]
77 lemmatized = [lemmatizer.lemmatize(t) for t in filtered_tokens]
78
79 print("\nStemmed (first 20):", stemmed[:20])
80 print("Lemmatized (first 20):", lemmatized[:20])
81
82 # h. Sentence boundaries
83 sentences = sent_tokenize(text)
84 print(f"\nTotal number of sentences: {len(sentences)}")
```

```
file.txt created successfully!
[nltk_data] Downloading package punkt_tab to /root/nltk_data...
[nltk_data]   Package punkt_tab is already up-to-date!
```

```
[nltk_data] Downloading package averaged_perceptron_tagger_eng to
[nltk_data]     /root/nltk_data...
[nltk_data]   Unzipping taggers/averaged_perceptron_tagger_eng.zip.
[nltk_data] Downloading package punkt to /root/nltk_data...
[nltk_data]   Package punkt is already up-to-date!
[nltk_data] Downloading package stopwords to /root/nltk_data...
[nltk_data]   Package stopwords is already up-to-date!
[nltk_data] Downloading package wordnet to /root/nltk_data...
[nltk_data]   Package wordnet is already up-to-date!
[nltk_data] Downloading package averaged_perceptron_tagger to
[nltk_data]     /root/nltk_data...
[nltk_data]   Package averaged_perceptron_tagger is already up-to-
[nltk_data]       date!
Original text length: 1253
Natural language processing (NLP) is a sub-field of artificial intelligence, concerned with the interactions between compute

NLP is used in many applications such as: machine translation, sentiment analysis, speech recognition, chatbots and informat

Some common challenges include: ambiguity in language, sarcasm, slang, and misspellin

First 30 tokens: ['natural', 'language', 'processing', '(', 'nlp', ')', 'is', 'a', 'sub-field', 'of', 'artificial', 'intelli

First 10 corrected tokens: ['natural', 'language', 'processing', '(', 'nap', ')', 'is', 'a', 'sub-field', 'of']

Corrected text (first 200 chars): natural language processing ( nap ) is a sub-field of artificial intelligence , concerned

POS tags (first 10): [('natural', 'JJ'), ('language', 'NN'), ('processing', 'NN'), ('(', '('), ('nap', 'JJ'), (')', ')'), ('

First 20 tokens after stop-word removal: ['natural', 'language', 'processing', 'nap', 'artificial', 'intelligence', 'concerr

Stemmed (first 20): ['natur', 'languag', 'process', 'nap', 'artifici', 'intellig', 'concern', 'interact', 'comput', 'human',
Lemmatized (first 20): ['natural', 'language', 'processing', 'nap', 'artificial', 'intelligence', 'concerned', 'interaction'

Total number of sentences: 14
```

## 2.Feature Extraction on 20newsgroups

```python
1 # a. Import packages
2 from sklearn.datasets import fetch_20newsgroups
3 import pandas as pd
4 import re
5 from sklearn.feature_extraction.text import CountVectorizer, TfidfVectorizer
6
7 # b. Fetch dataset
8 data = fetch_20newsgroups(subset='all', remove=('headers', 'footers', 'quotes'))
9 df = pd.DataFrame({'text': data.data, 'target': data.target})
10
11 # c. Clean data
12 def clean_text(txt):
13     txt = txt.lower()
14     txt = re.sub(r'[^a-z\s]', '', txt)
15     txt = re.sub(r'\s+', ' ', txt).strip()
16     return txt
17
18 df['cleaned'] = df['text'].apply(clean_text)
19
20 # d. BoW model
21 bow_vectorizer = CountVectorizer(max_features=5000, stop_words='english')
22 bow_matrix = bow_vectorizer.fit_transform(df['cleaned'])
23 bow_df = pd.DataFrame(bow_matrix.toarray(), columns=bow_vectorizer.get_feature_names_out())
24
25 # e. TF-IDF model
26 tfidf_vectorizer = TfidfVectorizer(max_features=5000, stop_words='english')
27 tfidf_matrix = tfidf_vectorizer.fit_transform(df['cleaned'])
28 tfidf_df = pd.DataFrame(tfidf_matrix.toarray(), columns=tfidf_vectorizer.get_feature_names_
29
30 # f. Compare top 20 frequent words
31 bow_sums = bow_df.sum().sort_values(ascending=False).head(20)
32 tfidf_sums = tfidf_df.sum().sort_values(ascending=False).head(20)
33
34 print("BoW Top 20 words:\n", bow_sums)
35 print("\nTF-IDF Top 20 words:\n", tfidf_sums)
```

```
BoW Top 20 words:
 dont                    6433
like                     6391
```

```
people                               6294
just                                 6155
know                                 5752
use                                  4989
think                                4983
time                                 4601
does                                 4489
new                                  4034
im                                   4004
good                                 3878
make                                 3334
way                                  3328
maxaxaxaxaxaxaxaxaxaxaxaxaxaxaxax    3317
used                                 3028
say                                  3007
did                                  2989
god                                  2949
right                                2944
dtype: int64

TF-IDF Top 20 words:
 like       324.598339
just        322.250787
know        318.619716
dont        317.050850
think       271.152683
does        267.344562
people      265.967689
im          265.655044
use         235.307722
thanks      232.253079
good        223.114082
time        215.373138
new         200.757875
make        177.541428
need        175.007184
way         174.105098
want        171.583419
did         171.285917
right       169.222452
ive         168.423608
dtype: float64
```

### 3.Amazon Musical Instruments Reviews

```python
1
2 # Import necessary libraries
3 import pandas as pd
4 import numpy as np
5 import re
6 import matplotlib.pyplot as plt
7 import seaborn as sns
8
9 from sklearn.model_selection import train_test_split
10 from sklearn.feature_extraction.text import TfidfVectorizer
11 from sklearn.metrics import (
12     classification_report, confusion_matrix, accuracy_score,
13     roc_curve, auc, mean_squared_error
14 )
15
16 # Classifiers
17 from sklearn.linear_model import LogisticRegression
18 from sklearn.naive_bayes import MultinomialNB
19 from sklearn.neighbors import KNeighborsClassifier
20 from sklearn.tree import DecisionTreeClassifier
21 from sklearn.ensemble import RandomForestClassifier, GradientBoostingClassifier
22 import xgboost as xgb
23
24 # Regressors
25 from sklearn.linear_model import LinearRegression
26 from sklearn.tree import DecisionTreeRegressor
27 from sklearn.ensemble import RandomForestRegressor, GradientBoostingRegressor
28 import xgboost as xgb
29
30 # Set style
31 sns.set_palette("husl")
32
33 # ============================================
```

```python
34 # a. Data Input
35 # ===========================================
36
37 data = {
38     "reviewText": [
39         "Not much to write about here, but it does exac...",
40         "The product does exactly as it should and is q...",
41         "The primary job of this device is to block the...",
42         "Nice windscreen protects my MXL mic and preven...",
43         "This pop filter is great. It looks and perform...",
44         "So good that I bought another one. Love the h...",
45         "I have used monster cables for years, and with...",
46         "I now use this cable to run from the output of...",
47         "Perfect for my Epiphone Sheraton II. Monster ...",
48         "Monster makes the best cables and a lifetime w..."
49     ],
50     "Overall": [5, 4, 3, 5, 4, 4, 5, 5, 3, 4]
51 }
52
53 df = pd.DataFrame(data)
54 print("Original Data:")
55 print(df)
56 print("\n" + "="*60 + "\n")
57
58 # ===========================================
59 # b. Clean reviewText
60 # ===========================================
61
62 def clean_review(text):
63     text = text.lower()
64     text = re.sub(r'[^a-z\s]', '', text)  # Remove punctuation & numbers
65     text = re.sub(r'\s+', ' ', text).strip()  # Normalize whitespace
66     return text
67
68 df['cleaned_review'] = df['reviewText'].apply(clean_review)
69
70 print("Cleaned Reviews (first 3):")
71 for i in range(3):
72     print(f"{i+1}. {df['cleaned_review'].iloc[i]}")
73 print("\n" + "="*60 + "\n")
74
75 # ===========================================
76 # c. TF-IDF Feature Extraction
77 # ===========================================
78
79 vectorizer = TfidfVectorizer(
80     max_features=100,
81     stop_words='english',
82     ngram_range=(1,2),  # Include bigrams for richer features
83     min_df=1
84 )
85
86 X = vectorizer.fit_transform(df['cleaned_review'])
87 feature_names = vectorizer.get_feature_names_out()
88
89 print(f"TF-IDF Matrix Shape: {X.shape}")
90 print(f"Features (sample): {feature_names[:10]}")
91 print("\n" + "="*60 + "\n")
92
93 # Target for classification: 1 = Positive (4–5), 0 = Negative (1–3)
94 y_class = (df['Overall'] >= 4).astype(int)  # Critical Fix: Use 0 and 1
95 y_reg = df['Overall']
96
97 # Train-test split
98 X_train_c, X_test_c, y_train_c, y_test_c = train_test_split(
99     X, y_class, test_size=0.3, random_state=42, stratify=y_class
100 )
101 X_train_r, X_test_r, y_train_r, y_test_r = train_test_split(
```

```
102      X, y_reg, test_size=0.3, random_state=42
103 )
104
105 # ===========================================
106 # d. Models Definition
107 # ===========================================
108
109 classifiers = {
110     "Logistic Regression": LogisticRegression(max_iter=1000),
111     "Naive Bayes": MultinomialNB(),
112     "KNN": KNeighborsClassifier(n_neighbors=3),
113     "Decision Tree": DecisionTreeClassifier(random_state=42),
114     "Random Forest": RandomForestClassifier(n_estimators=100, random_state=42),
115     "GBM": GradientBoostingClassifier(random_state=42),
116     "XGBoost": xgb.XGBClassifier(
117         use_label_encoder=False,
118         eval_metric='logloss',
119         random_state=42
120     )
121 }
122
123 regressors = {
124     "Linear Regression": LinearRegression(),
125     "Decision Tree Reg": DecisionTreeRegressor(random_state=42),
126     "Random Forest Reg": RandomForestRegressor(n_estimators=100, random_state=42),
127     "GBM Reg": GradientBoostingRegressor(random_state=42),
128     "XGB Reg": xgb.XGBRegressor(random_state=42)
129 }
130
131 # ===========================================
132 # e. Evaluation Function (Classification)
133 # ===========================================
134
135 def evaluate_classification(model, X_test, y_test, name):
136     y_pred = model.predict(X_test)
137     y_prob = model.predict_proba(X_test)[:, 1] if hasattr(model, "predict_proba") else Nor
138
139     print(f"\n{'='*20} {name} {'='*20}")
140     print(f"Accuracy: {accuracy_score(y_test, y_pred):.3f}")
141     print("\nClassification Report:")
142     print(classification_report(y_test, y_pred, target_names=['Negative', 'Positive']))
143
144     # Confusion Matrix
145     cm = confusion_matrix(y_test, y_pred)
146     plt.figure(figsize=(5,4))
147     sns.heatmap(cm, annot=True, fmt='d', cmap='Blues',
148                 xticklabels=['Neg', 'Pos'], yticklabels=['Neg', 'Pos'])
149     plt.title(f"{name} - Confusion Matrix")
150     plt.xlabel("Predicted")
151     plt.ylabel("Actual")
152     plt.show()
153
154     # ROC Curve
155     if y_prob is not None:
156         fpr, tpr, _ = roc_curve(y_test, y_prob)
157         roc_auc = auc(fpr, tpr)
158         plt.figure(figsize=(5,4))
159         plt.plot(fpr, tpr, label=f'ROC Curve (AUC = {roc_auc:.3f})')
160         plt.plot([0,1], [0,1], 'k--')
161         plt.xlim([0,1]); plt.ylim([0,1.05])
162         plt.xlabel('False Positive Rate')
163         plt.ylabel('True Positive Rate')
164         plt.title(f"{name} - ROC Curve")
165         plt.legend()
166         plt.show()
167
168 # ===========================================
169 # CLASSIFICATION RESULTS
```

```
170 # ============================================
171
172 print("\nCLASSIFICATION RESULTS\n" + "-"*60)
173 for name, clf in classifiers.items():
174     clf.fit(X_train_c, y_train_c)
175     evaluate_classification(clf, X_test_c, y_test_c, name)
176
177 # ============================================
178 # REGRESSION RESULTS
179 # ============================================
180
181 print("\n\nREGRESSION RESULTS\n" + "-"*60)
182
183 for name, reg in regressors.items():
184     reg.fit(X_train_r, y_train_r)
185     y_pred = reg.predict(X_test_r)
186     mse = mean_squared_error(y_test_r, y_pred)
187     print(f"{name:25} → MSE: {mse:.3f} | RMSE: {np.sqrt(mse):.3f}")
188
189 # Optional: Predict on first review
190 print("\n" + "="*60)
191 print("SAMPLE PREDICTION (First Review):")
192 sample = vectorizer.transform([df['cleaned_review'].iloc[0]])
193 for name, clf in classifiers.items():
194     pred = clf.predict(sample)[0]
195     prob = clf.predict_proba(sample)[0]
```

```
Original Data:
                                   reviewText  Overall
0  Not much to write about here, but it does exac...        5
1  The product does exactly as it should and is q...        4
2  The primary job of this device is to block the...        3
3  Nice windscreen protects my MXL mic and preven...        5
4  This pop filter is great. It looks and perform...        4
5   So good that I bought another one. Love the h...        4
6  I have used monster cables for years, and with...        5
7  I now use this cable to run from the output of...        5
8   Perfect for my Epiphone Sheraton II. Monster ...        3
9  Monster makes the best cables and a lifetime w...        4


=============================================================

Cleaned Reviews (first 3):
1. not much to write about here but it does exac
2. the product does exactly as it should and is q
3. the primary job of this device is to block the


=============================================================

TF-IDF Matrix Shape: (10, 70)
Features (sample): ['best' 'best cables' 'block' 'bought' 'bought love' 'cable' 'cable run'
 'cables' 'cables lifetime' 'cables years']


=============================================================


CLASSIFICATION RESULTS
-------------------------------------------------------------

==================== Logistic Regression ====================
Accuracy: 0.667

Classification Report:
              precision    recall  f1-score   support

    Negative       0.00      0.00      0.00         1
    Positive       0.67      1.00      0.80         2

    accuracy                           0.67         3
   macro avg       0.33      0.50      0.40         3
weighted avg       0.44      0.67      0.53         3
```
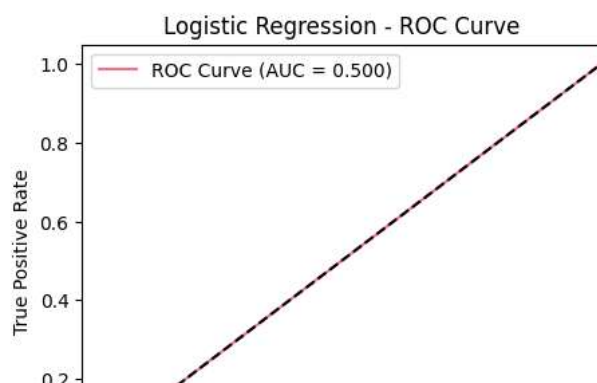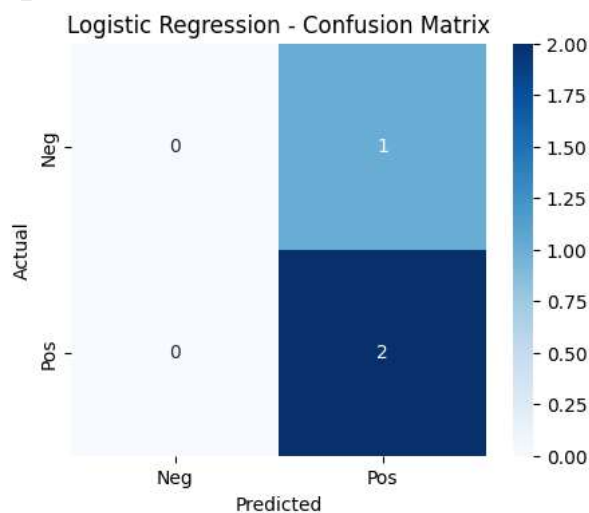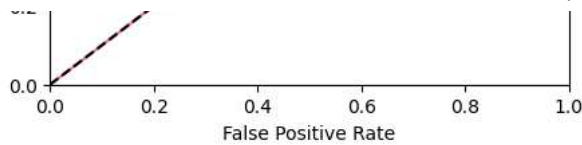
/usr/local/lib/python3.12/dist-packages/sklearn/metrics/_classification.py:1565: UndefinedMetricWarning: Precision is ill-
 _warn_prf(average, modifier, f"{metric.capitalize()} is", len(result))
/usr/local/lib/python3.12/dist-packages/sklearn/metrics/_classification.py:1565: UndefinedMetricWarning: Precision is ill-
 _warn_prf(average, modifier, f"{metric.capitalize()} is", len(result))
/usr/local/lib/python3.12/dist-packages/sklearn/metrics/_classification.py:1565: UndefinedMetricWarning: Precision is ill-
 _warn_prf(average, modifier, f"{metric.capitalize()} is", len(result))



Logistic Regression - Confusion Matrix



Logistic Regression - ROC Curve

```
==================== Naive Bayes ====================
Accuracy: 0.667

Classification Report:
              precision    recall  f1-score   support

    Negative       0.00      0.00      0.00         1
    Positive       0.67      1.00      0.80         2

    accuracy                           0.67         3
   macro avg       0.33      0.50      0.40         3
weighted avg       0.44      0.67      0.53         3
```
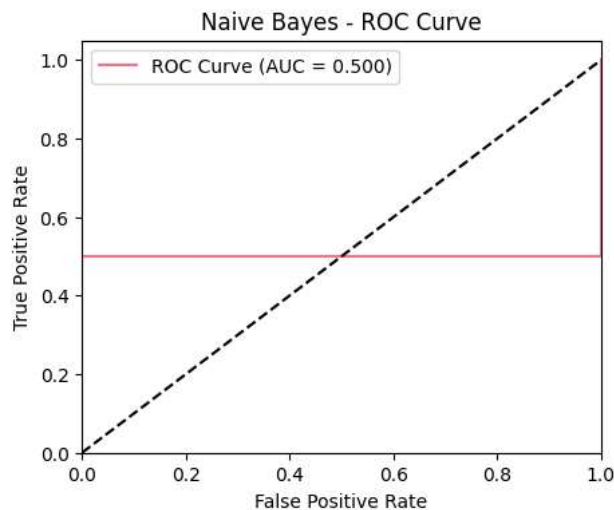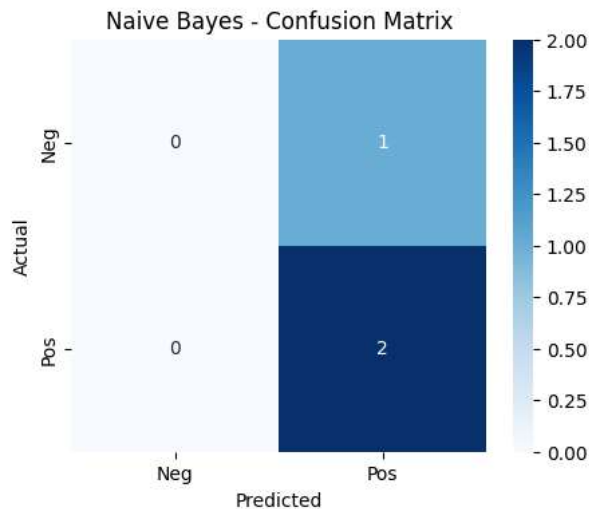
/usr/local/lib/python3.12/dist-packages/sklearn/metrics/_classification.py:1565: UndefinedMetricWarning: Precision is ill-
  _warn_prf(average, modifier, f"{metric.capitalize()} is", len(result))
/usr/local/lib/python3.12/dist-packages/sklearn/metrics/_classification.py:1565: UndefinedMetricWarning: Precision is ill-
  _warn_prf(average, modifier, f"{metric.capitalize()} is", len(result))
/usr/local/lib/python3.12/dist-packages/sklearn/metrics/_classification.py:1565: UndefinedMetricWarning: Precision is ill-
  _warn_prf(average, modifier, f"{metric.capitalize()} is", len(result))





```
==================== KNN ====================
Accuracy: 0.667

Classification Report:
              precision    recall  f1-score   support

    Negative       0.00      0.00      0.00         1
    Positive       0.67      1.00      0.80         2

    accuracy                           0.67         3
   macro avg       0.33      0.50      0.40         3
weighted avg       0.44      0.67      0.53         3
```
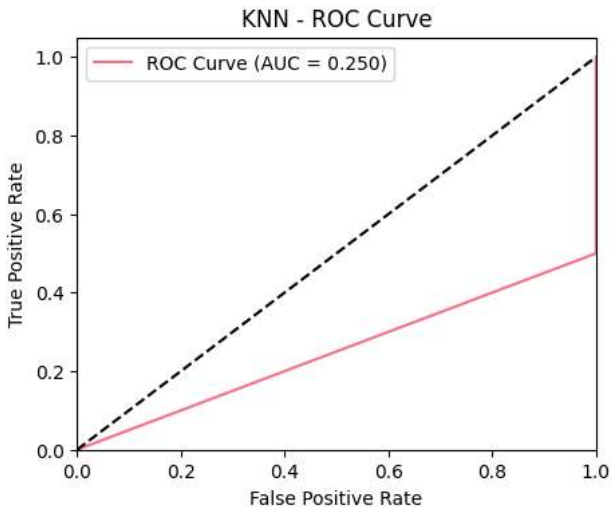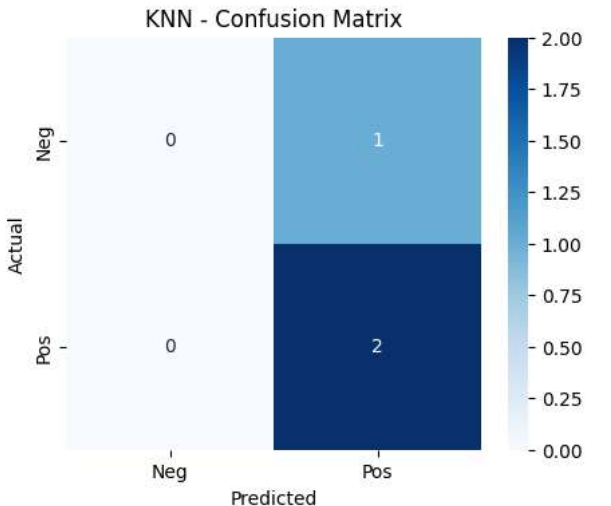
/usr/local/lib/python3.12/dist-packages/sklearn/metrics/_classification.py:1565: UndefinedMetricWarning: Precision is ill-
  _warn_prf(average, modifier, f"{metric.capitalize()} is", len(result))
/usr/local/lib/python3.12/dist-packages/sklearn/metrics/_classification.py:1565: UndefinedMetricWarning: Precision is ill-
  _warn_prf(average, modifier, f"{metric.capitalize()} is", len(result))
/usr/local/lib/python3.12/dist-packages/sklearn/metrics/_classification.py:1565: UndefinedMetricWarning: Precision is ill-

_warn_prf(average, modifier, f"{metric.capitalize()} is", len(result))

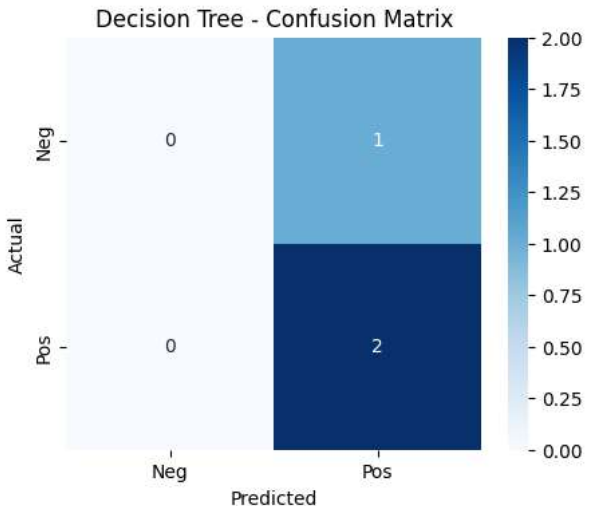### KNN - Confusion Matrix



### KNN - ROC Curve



```
==================== Decision Tree ====================
Accuracy: 0.667

Classification Report:
              precision    recall  f1-score   support

    Negative       0.00      0.00      0.00         1
    Positive       0.67      1.00      0.80         2

    accuracy                           0.67         3
   macro avg       0.33      0.50      0.40         3
weighted avg       0.44      0.67      0.53         3
```
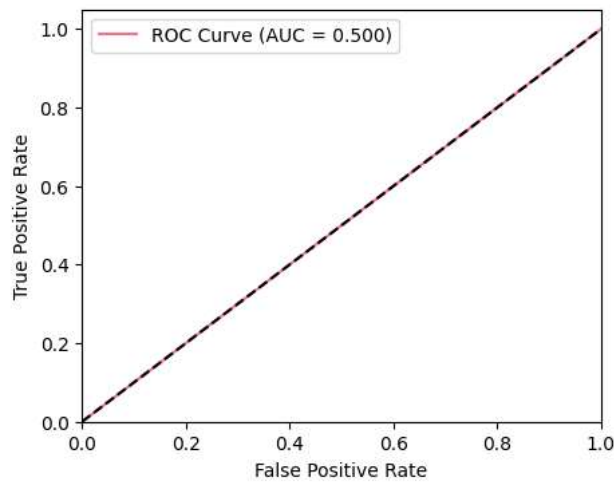
/usr/local/lib/python3.12/dist-packages/sklearn/metrics/_classification.py:1565: UndefinedMetricWarning: Precision is ill-
  _warn_prf(average, modifier, f"{metric.capitalize()} is", len(result))
/usr/local/lib/python3.12/dist-packages/sklearn/metrics/_classification.py:1565: UndefinedMetricWarning: Precision is ill-
  _warn_prf(average, modifier, f"{metric.capitalize()} is", len(result))
/usr/local/lib/python3.12/dist-packages/sklearn/metrics/_classification.py:1565: UndefinedMetricWarning: Precision is ill-
  _warn_prf(average, modifier, f"{metric.capitalize()} is", len(result))

### Decision Tree - Confusion Matrix



### Decision Tree - ROC Curve

```
==================== Random Forest ====================
Accuracy: 0.667

Classification Report:
              precision    recall  f1-score   support

    Negative       0.00      0.00      0.00         1
    Positive       0.67      1.00      0.80         2

    accuracy                           0.67         3
   macro avg       0.33      0.50      0.40         3
weighted avg       0.44      0.67      0.53         3
```
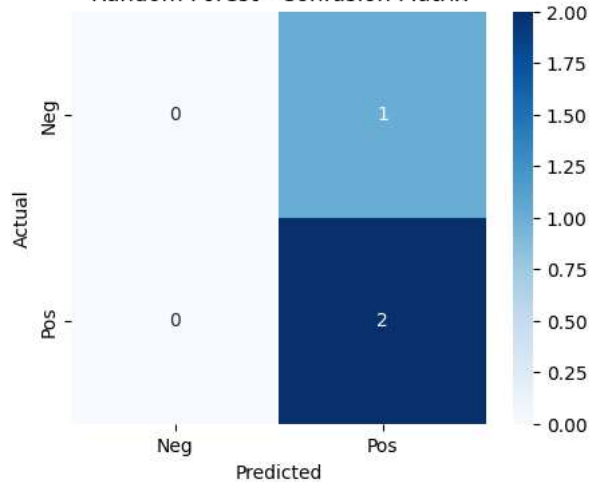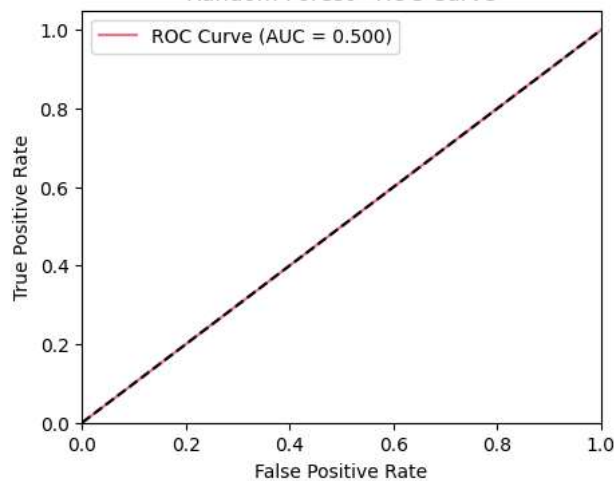
/usr/local/lib/python3.12/dist-packages/sklearn/metrics/_classification.py:1565: UndefinedMetricWarning: Precision is ill-
  _warn_prf(average, modifier, f"{metric.capitalize()} is", len(result))
/usr/local/lib/python3.12/dist-packages/sklearn/metrics/_classification.py:1565: UndefinedMetricWarning: Precision is ill-
  _warn_prf(average, modifier, f"{metric.capitalize()} is", len(result))
/usr/local/lib/python3.12/dist-packages/sklearn/metrics/_classification.py:1565: UndefinedMetricWarning: Precision is ill-
  _warn_prf(average, modifier, f"{metric.capitalize()} is", len(result))





```
==================== GBM ====================
Accuracy: 0.667
```

```
Classification Report:
              precision    recall  f1-score   support

    Negative       0.00      0.00      0.00         1
    Positive       0.67      1.00      0.80         2

    accuracy                           0.67         3
   macro avg       0.33      0.50      0.40         3
weighted avg       0.44      0.67      0.53         3
```
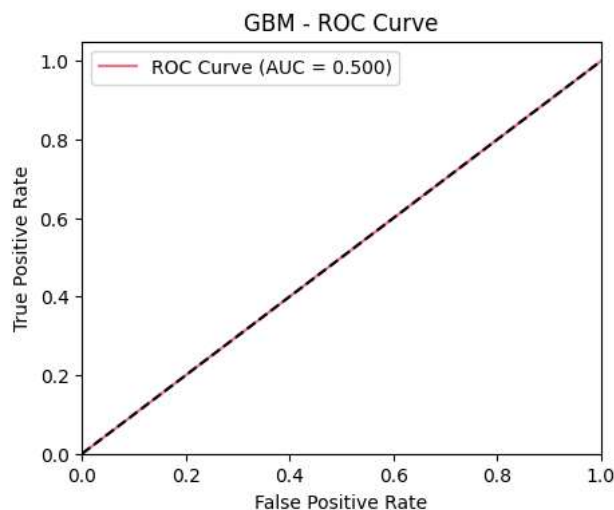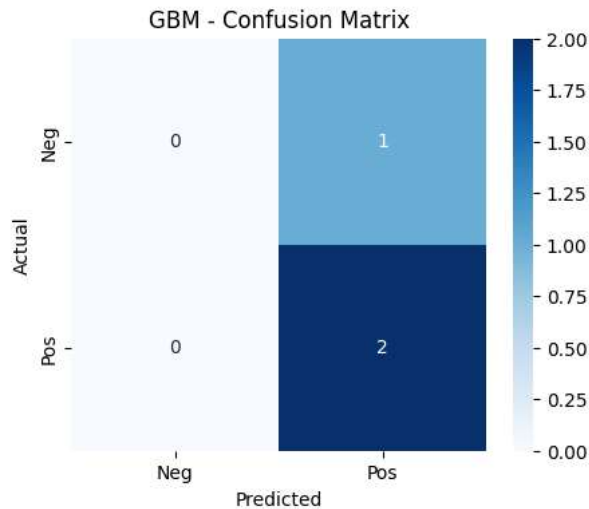
```
/usr/local/lib/python3.12/dist-packages/sklearn/metrics/_classification.py:1565: UndefinedMetricWarning: Precision is ill-
  _warn_prf(average, modifier, f"{metric.capitalize()} is", len(result))
/usr/local/lib/python3.12/dist-packages/sklearn/metrics/_classification.py:1565: UndefinedMetricWarning: Precision is ill-
  _warn_prf(average, modifier, f"{metric.capitalize()} is", len(result))
/usr/local/lib/python3.12/dist-packages/sklearn/metrics/_classification.py:1565: UndefinedMetricWarning: Precision is ill-
  _warn_prf(average, modifier, f"{metric.capitalize()} is", len(result))
```



GBM - Confusion Matrix



GBM - ROC Curve

```
=================== XGBoost ===================
Accuracy: 0.667

Classification Report:
              precision    recall  f1-score   support

    Negative       0.00      0.00      0.00         1
    Positive       0.67      1.00      0.80         2

    accuracy                           0.67         3
   macro avg       0.33      0.50      0.40         3
weighted avg       0.44      0.67      0.53         3
```

```
/usr/local/lib/python3.12/dist-packages/xgboost/training.py:199: UserWarning: [15:17:15] WARNING: /workspace/src/learner.c
Parameters: { "use_label_encoder" } are not used.

  bst.update(dtrain, iteration=i, fobj=obj)
/usr/local/lib/python3.12/dist-packages/sklearn/metrics/_classification.py:1565: UndefinedMetricWarning: Precision is ill-
  _warn_prf(average, modifier, f"{metric.capitalize()} is", len(result))
/usr/local/lib/python3.12/dist-packages/sklearn/metrics/_classification.py:1565: UndefinedMetricWarning: Precision is ill-
  _warn_prf(average, modifier, f"{metric.capitalize()} is", len(result))
/usr/local/lib/python3.12/dist-packages/sklearn/metrics/_classification.py:1565: UndefinedMetricWarning: Precision is ill-
  _warn_prf(average, modifier, f"{metric.capitalize()} is", len(result))
```



XGBoost - Confusion Matrix

## 4.IMDB Sentiment Analysis (Kaggle)

```python
1 from google.colab import files
2 import os
3
4 if not os.path.exists("/root/.kaggle"):
5     os.makedirs("/root/.kaggle", exist_ok=True)
6
7 print("Please upload your kaggle.json file:")
8 uploaded = files.upload()                    # <-- Upload kaggle.json
9 for fn in uploaded.keys():
10     os.rename(fn, "/root/.kaggle/kaggle.json")
11 os.chmod("/root/.kaggle/kaggle.json", 600)
12
13 # ----------------------------------------------------------------
14 # 3. Download the exact dataset you mentioned
15 # ----------------------------------------------------------------
16 !kaggle datasets download -d columbine/imdb-dataset-sentiment-analysis-in-csv-format --un:
17
18 # The dataset contains Train.csv and Test.csv
19 # We will use Test.csv as per your instruction
20 csv_path = "Test.csv"   # <-- This is the file we need
21 print(f"\nDataset ready: {csv_path}")
22
23 # ----------------------------------------------------------------
24 # 4. Import libraries
25 # ----------------------------------------------------------------
26 import pandas as pd
27 import re
28 from textblob import TextBlob
29 import warnings
30 warnings.filterwarnings("ignore")
31
32 # ----------------------------------------------------------------
33 # 5. Load the CSV file
34 # ----------------------------------------------------------------
35 df = pd.read_csv(csv_path)
36 print(f"Loaded {len(df)} rows from {csv_path}")
37 print(df.head())
38
39 # ----------------------------------------------------------------
40 # 6. Fetch the text column
41 # ----------------------------------------------------------------
42 # The column is named 'text' in this dataset
43 texts = df['text'].astype(str).copy()
44
45 # ----------------------------------------------------------------
46 # 7. Extract and remove @handles
47 # ----------------------------------------------------------------
48 def remove_handles(t):
49     return re.sub(r'@\w+', '', t)
50
51 texts_clean = texts.apply(remove_handles)
52
53 # ----------------------------------------------------------------
54 # 8. Perform sentiment analysis using TextBlob
55 # ----------------------------------------------------------------
56 def get_sentiment(txt):
57     blob = TextBlob(txt)
58     polarity = blob.sentiment.polarity
59     if polarity > 0.1:
60         return 'positive'
61     elif polarity < -0.1:
62         return 'negative'
63     else:
64         return 'neutral'
65
```