

기계학습 감성 분석 기말 과제

학번 : 201921030

이름 : 임재석

목차

1. 감성 분석의 기본 이해
 - 1.1. 감성 분석 소개
 - 1.2. 감성 분석의 응용
2. 데이터 전처리
 - 2.1. 데이터 수집
 - 2.2. 데이터 전처리 과정
3. 모델 선택 및 학습
 - 3.1. 사용된 모델 소개
 - 3.2. 모델 학습 방법
 - 3.3. 모델 강점과 약점
4. 평가
 - 4.1. 성능 평가 지표 소개
 - 4.2. 모델 성능 평가 결과
5. 시각화 및 해석
 - 5.1. 결과 시각화
 - 5.2. 결과 해석

1-1

감성분석이란?

감성 분석(Sentiment Analysis)이란 텍스트에 들어있는 의견이나 감성, 평가, 태도 등의 주관적인 정보를 컴퓨터를 통해 분석하는 과정

응용

감성 분석은 마케팅에서 고객 서비스, 임상 의학에 이르기까지 다양한 애플리케이션을 위한 리뷰 및 설문 조사 응답, 온라인 및 소셜 미디어, 의료 자료 등 고객의 소리 자료에 널리 적용

1-2

전처리

모듈 호출 및 세팅

```
%pip install pandas numpy scikit-learn nltk matplotlib seaborn

import pandas as pd
import numpy as np
import re
from sklearn.model_selection import train_test_split, GridSearchCV, cross_val_score
from sklearn.pipeline import Pipeline
from sklearn.feature_extraction.text import TfidfVectorizer
from sklearn.svm import SVC
from sklearn.metrics import classification_report, accuracy_score, confusion_matrix, roc_curve, auc
from nltk.corpus import stopwords
from nltk.stem import WordNetLemmatizer
from nltk.tokenize import word_tokenize
import nltk
import matplotlib.pyplot as plt
import seaborn as sns

# NLTK 리소스 다운로드
nltk.download('punkt')
nltk.download('stopwords')
nltk.download('wordnet')
```

우선 ML에 필요한 모듈 호출과 감성 분석에 필요한 NLTK 소스를 호출

```

# 데이터 불러오기
df = pd.read_csv('./wine_review.csv')

# 필요한 열만 선택
df = df[['reviews.text', 'reviews.rating']]

# 결측값 제거
df.dropna(subset=['reviews.text', 'reviews.rating'], inplace=True)

# 감성 라벨링 (5점 기준 4, 5는 긍정, 1, 2는 부정, 3은 중립으로 간주)
def label_sentiment(rating):
    if rating >= 4:
        return 'positive'
    elif rating <= 2:
        return 'negative'
    else:
        return 'neutral'

df['sentiment'] = df['reviews.rating'].apply(label_sentiment)

```

불러온 데이터를 호출하고 기초적인 결측치 제거와 라벨링을 진행

```

# 텍스트 전처리 함수 정의
def preprocess_text(text):
    # 소문자 변환
    text = text.lower()
    # HTML 태그 제거
    text = re.sub(r'<.*?>', '', text)
    # 특수 문자 및 숫자 제거
    text = re.sub(r'^a-zA-Z\s]', '', text)
    # 토큰화
    words = word_tokenize(text)
    # 불용어 제거
    words = [word for word in words if word not in stopwords.words('english')]
    # 표제어 추출
    lemmatizer = WordNetLemmatizer()
    words = [lemmatizer.lemmatize(word) for word in words]
    return ' '.join(words)

# 텍스트 데이터 전처리 적용
df['cleaned_text'] = df['reviews.text'].apply(preprocess_text)

# 데이터셋 분할
x = df['cleaned_text']
y = df['sentiment']

# 저장할 전처리된 데이터를 새 CSV 파일로 저장
df.to_csv('wine_review_cleaned.csv', index=False)

print("Data preprocessing and CSV file creation completed.")

```

요구사항에 맞춘 토큰화와 각 문자들을 전처리하고 새로운 데이터셋을 csv로 저장

BEFORE Data Set

```
code.ipynb wine_review_cleaned.csv wine_review.csv X
wine_review.csv
1 id,asins,brand,categories,dateAdded,dateUpdated,descriptions,dimension,ean,flavors,keys,manufacturer,manufacturerNumber,name,
  reviews.date,reviews.dateAdded,reviews.dateSeen,reviews.didPurchase,reviews.doRecommend,reviews.id,reviews.numHelpful,reviews.
  rating,reviews.sourceURLs,reviews.text,reviews.title,reviews.userCity,reviews.userProvince,reviews.username,sizes,sourceURLs,upc,
  weight
2 AV13C1KCGV-KLJ3akN68,,Gallo,"Food & Beverage,Beverages,Wine, Beer & Liquor,Wine",2017-07-24T23:59:11Z,2018-01-10T18:06:28Z,,1.0
  in x 1.0 in x 1.0 in,,,492130001994,gallo/13312834",,13312834,Ecco Domani174 Pinot Grigio - 750ml Bottle,2017-08-01T21:13:49.
  000Z,2018-01-09T13:24:04Z,"2017-12-14T19:41:00.000Z,2017-12-19T19:55:00.000Z,2017-12-06T01:40:00.000Z,2017-11-18T20:09:00.000Z,
  2017-10-26T16:13:00.000Z,2017-08-16T11:20:15.746Z,2017-08-18T18:30:15.452Z,2017-09-01T17:04:29.797Z,2017-08-06T14:00:48.156Z,
  2017-08-20T18:01:35.030Z,2017-08-08T11:58:19.262Z,2017-08-04T15:49:50.043Z",,TRUE,,1,5,https://redsky.target.com/
  groot-domain-api/v1/reviews/13312834?sort=helpfulness_desc&limit=200&offset=0,This a fantastic white wine for any occasion!,My
  Favorite White Wine,,Bjh,"http://redsky.target.com/v1/plp/search?kwr=y&count=90&offset=0&category=5xsxvZ56dar,http://redsky.
  target.com/v1/plp/search?kwr=y&count=90&offset=0&category=5xsxv,http://redsky.target.com/v1/plp/search?kwr=y&count=90&offset=0&
  category=5xsxvZ5604a,http://redsky.target.com/v1/plp/search?kwr=y&count=90&offset=0&category=5n5q6,http://redsky.target.com/v1/
  plp/search?kwr=y&count=90&offset=810&category=5xt0r,http://redsky.target.com/v2/pdp/tcin/13312834?excludes=taxonomy,https://
  redsky.target.com/groot-domain-api/v1/reviews/13312834?sort=helpfulness_desc&limit=200&offset=0,http://redsky.target.com/v1/plp/
  search?kwr=y&count=90&offset=900&category=5xt0r,http://redsky.target.com/v1/plp/search?kwr=y&count=90&offset=990&category=5xt0r,
  http://redsky.target.com/v1/plp/search?kwr=y&count=90&offset=0&category=5xsxvZ5zkq6,http://redsky.target.com/v1/plp/search?kwr=y&
  count=90&offset=540&category=5xsxv,http://redsky.target.com/v1/plp/search?kwr=y&count=90&offset=0&category=5xsxvZ5zkq4,http://
  redsky.target.com/v1/plp/search?kwr=y&count=90&offset=0&category=5xsxvZ5zkq5,http://redsky.target.com/v1/plp/search?kwr=y&
  count=90&offset=0&category=5xsxvZ5zkq3",4.9213E+11,1.0 lbs
3 AV13C5vW-jtxr-f38AQ0,,Fresh Craft Co.,,"Food & Beverage,Beverages,Wine, Beer & Liquor,Wine",2017-07-24T23:59:42Z,
  2018-01-10T05:38:33Z,"[{"dateSeen":["2017-12-21T05:43:00.000Z","2017-12-16T19:47:00.000Z","2017-12-07T17:57:00.000Z",
  "2017-11-19T16:52:00.000Z","2017-10-26T04:50:00.000Z","2017-10-26T18:16:00.000Z","2017-09-17T00:29:40.578Z",
  "2017-09-21T05:03:20.394Z","2017-09-21T05:03:19.626Z","2017-09-01T19:14:26.079Z","2017-09-01T21:55:58.715Z",
  "2017-08-23T02:54:07.671Z","2017-08-20T19:32:08.011Z","2017-08-23T02:01:23.533Z","2017-08-22T15:22:26.562Z",
```

전처리 전 데이터는 각 컬럼과 필요없는 날짜, 링크 등 데이터가 존재

AFTER Data Set

```
code.ipynb wine_review_cleaned.csv X wine_review.csv
wine_review_cleaned.csv
1 reviews.text,reviews.rating,sentiment,cleaned_text
2 This a fantastic white wine for any occasion!,5.0,positive,fantastic white wine occasion
3 "Tart, not sweet...very refreshing and delicious!",5.0,positive,tart sweetvery refreshing delicious
4 I was given this wine so it was a delightful surprise to find that it has a flavorful and delicious taste! A new favorite!!!!,5.0,
5 This is a phenomenal wine and my new favorite red.,5.0,positive,phenomenal wine new favorite red
6 4 750ml bottles for the price of two With way less packaging YES PLEASE! I was nervous it was too good to be true and I wouldn't l
7 I LOVE Becks NA. It tastes just like a regular ale. It smells like one too. It tastes great. It's the only NA that I drink. My fri
8 This wine has a wonderful but strong aroma its a bit bitter has a bite to it but still goodit was worth it in the end,3.0,neutral,
9 "I would give one more star if it came clean on the bottle and called itself a 'sweet red'. Instead it get's poetic and grandiose
10 Delicious and very affordable,5.0,positive,delicious affordable
11 "This is a very smooth red with Aromas of cocoa, coffee, tobacco sweet black cherry. Bold but soft. I highly recommend this one.",
12 "Based on positive reviews, I served all of the Fancy Pants varieties at a recent party. My guests and I loved them! Since I am no
13 Nice fruity and sweet taking sparkling wine. Goes great in mimosa. You have to try it to see how you will like it. Even the guys w
14 A rich amber Ale with a subtle but noticeable hop character. Great beer to pair with food!,5.0,positive,rich amber ale subtle noti
15 "This light bodied IPA is dry hopped with cascade for a bold and bright HP aroma, balanced by a refreshing and crisp Ale. Great fo
16 "This light bodied IPA is dry hopped with cascade for a bold and bright HP aroma, balanced by a refreshing and crisp Ale. Great fo
17 Great for a perfect summer day or a get together with the girls.,5.0,positive,great perfect summer day get together girl
18 I'm not sure if this is a bad bottle or if it was intentionally made effervescent. I don't want bubbles when I drink Pinot Grigio.
19 "Obsessed with this new drink---However, it has been sold out at every Target I have been to lately :(",5.0,positive,obsessed new
20 "This is an excellent red wine, not to dry and not real sweet. More of a red blend. Great go with anything anytime. In the right p
21 One of the best Cabernets I had in quite sometime.....!!!!!!!,5.0,positive,one best cabernet quite sometime
22 "This is a great buy for the price. Someone had given me a bottle and I liked the look of it so much that I didn't want to open it
23 Gluten free plus lower on calories plus very refreshing!,5.0,positive,gluten free plus lower calorie plus refreshing
24 Best of the bunch of all of The wine cube sku's. Target Wine Buyer: please bring back!,5.0,positive,best bunch wine cube skus targ
25 "I like the wine and I like the idea of a sealed cube during the week when only drinking a glass or two at a time. But when they r
26 "I have tried my share of Moscato D' Astis before, this has to be the smoothest that I have ever tried. Not too sweet or too fizzy
27 "Highly, highly recommend!! It's sweet but not syrupy- just right! And I love that it's bubbly, I don't have to mix my wine with s
28 My favorite bottle design. The best cheap lager.,5.0,positive,favorite bottle design best cheap lager
29 Delicious. I hate all the other flavors but this one tastes great. Now if I can only find it in the store.,5.0,positive,delicious
```

토큰 제거와 필요한 평가, 감성 텍스트만 남은 데이터셋을 확인가능

3 -1 모델

모델 생성 및 튜닝

```
# 데이터 불러오기 (전처리된 데이터 사용)
df = pd.read_csv('./wine_review_cleaned.csv')

# 결측값 처리
df.dropna(subset=['cleaned_text'], inplace=True)

# 데이터셋 분할
X = df['cleaned_text']
y = df['sentiment']

# 학습 및 테스트 세트로 분할
X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2, random_state=42, stratify=y)

# 파이프라인 설정
pipeline = Pipeline([
    ('tfidf', TfidfVectorizer()),
    ('svc', SVC())
])
```

위 코드를 통해 전처리를 거친 데이터를 호출하고 학습을 위해 분할.

분할은 다음과 같음

X = df['cleaned_text']: 텍스트 데이터를 특징 변수

y = df['sentiment']: 감성 데이터를 목표 변수

train_test_split(X, y, test_size=0.2, random_state=42, stratify=y): 데이터를 학습 세트와 테스트 세트로 80:20 비율로 분할 stratify=y는 원본 데이터의 클래스 비율을 유지하기 위해 사용

Pipeline: 여러 처리 단계를 묶어주는 역할

SVC(): 서포트 벡터 머신 분류기를 사용

TfidfVectorizer(): 텍스트 데이터를 TF-IDF(Term Frequency-Inverse Document Frequency) 벡터로 변환

하이퍼파라미터 튜닝

```
17
# 하이퍼파라미터 튜닝
param_grid = {
    'tfidf__max_df': [0.8, 0.9, 1.0],
    'tfidf__ngram_range': [(1, 1), (1, 2)],
    'svc__C': [0.1, 1, 10],
    'svc__kernel': ['linear', 'rbf']
}

# GridSearchCV를 사용하여 최적의 하이퍼파라미터 탐색
grid_search = GridSearchCV(pipeline, param_grid, cv=5, n_jobs=-1, scoring='accuracy')
grid_search.fit(X_train, y_train)

# 최적의 모델로 예측 수행
best_model = grid_search.best_estimator_
y_pred = best_model.predict(X_test)

# 모델 평가
print(f"Best Parameters: {grid_search.best_params_}")
print(f"Accuracy: {accuracy_score(y_test, y_pred)}")
print(classification_report(y_test, y_pred))

# 모델 평가 (교차 검증)
cross_val_scores = cross_val_score(best_model, X, y, cv=5, scoring='accuracy')
print(f"Cross-validation scores: {cross_val_scores}")
print(f"Mean cross-validation score: {cross_val_scores.mean()}")
```

GridSearchCV에서 사용할 하이퍼파라미터 그리드 사용

GridSearchCV: 교차 검증을 통해 최적의 하이퍼파라미터를 탐색

grid_search.best_params_: 최적의 하이퍼파라미터를 출력

accuracy_score(y_test, y_pred): 테스트 세트의 정확도를 계산

classification_report(y_test, y_pred): 정밀도, 재현율, F1 점수를 포함한 상세한 분류 보고서를 출력

cross_val_score(best_model, X, y, cv=5, scoring='accuracy'): 전체 데이터에 대해 5-폴드 교차 검증을 수행하여 정확도를 평가

cross_val_scores.mean(): 교차 검증의 평균 정확도를 계산

튜닝 결과

```
Best Parameters: {'svc__C': 10, 'svc__kernel': 'rbf', 'tfidf__max_df': 0.8, 'tfidf__ngram_range': (1, 1)}
Accuracy: 0.934560327198364

      precision    recall  f1-score   support

negative      0.75      0.14      0.23         22
neutral       0.00      0.00      0.00          13
positive      0.94      1.00      0.97        454

accuracy              0.93         489
macro avg      0.56      0.38      0.40         489
weighted avg    0.90      0.93      0.91         489

/home/codespace/.local/lib/python3.10/site-packages/sklearn/metrics/_classification.py:1517: UndefinedMetricWarning: Precision-Recall score was calculated with True Positives=0 while there were no samples of this class
_warn_prf(average, modifier, f"{metric.capitalize()} is", len(result))
/home/codespace/.local/lib/python3.10/site-packages/sklearn/metrics/_classification.py:1517: UndefinedMetricWarning: Precision-Recall score was calculated with True Positives=0 while there were no samples of this class
_warn_prf(average, modifier, f"{metric.capitalize()} is", len(result))
/home/codespace/.local/lib/python3.10/site-packages/sklearn/metrics/_classification.py:1517: UndefinedMetricWarning: Precision-Recall score was calculated with True Positives=0 while there were no samples of this class
_warn_prf(average, modifier, f"{metric.capitalize()} is", len(result))
Cross-validation scores: [0.93047035 0.92229039 0.92842536 0.93237705 0.93237705]
Mean cross-validation score: 0.9291880384860371
```

최적의 하이퍼파라미터

svc__C: 10: SVM 분류기의 규제 매개변수 CCC가 10으로 설정 높은 CCC 값은 모델이 데이터에 더 잘 맞게 되지만, 과적합(overfitting)의 위험 가능

- **svc__kernel: 'rbf':** RBF (Radial Basis Function) 커널을 사용하여 비선형 경계를 학습
- **tfidf__max_df: 0.8:** TF-IDF 벡터라이저에서 단어가 전체 문서의 80% 이하에서 나타나면 포함
- **tfidf__ngram_range: (1, 1):** Unigram (단일 단어)를 사용

테스트 세트에서 93.46%의 높은 정확도를 달성

negative 클래스:

- **precision: 0.75:** 모델이 negative 클래스로 예측한 것 중 75%가 실제로 negative
- **recall: 0.14:** 실제 negative 샘플 중 14%만이 정확히 예측
- **f1-score: 0.23:** 정밀도와 재현율의 조화 평균을 나타냄. 낮은 값은 모델이 negative 클래스에서 잘 작동하지 않음을 나타냈음

neutral 클래스:

- **precision: 0.00, recall: 0.00, f1-score: 0.00:** 모델이 neutral 클래스를 전혀 예측하지 못함. 이는 neutral 샘플이 너무 적어서 모델이 이 클래스를 제대로 학습하지 못한 결과를 도출함

positive 클래스:

- **precision: 0.94:** 모델이 positive 클래스로 예측한 것 중 94%가 실제로 positive
- **recall: 1.00:** 실제 positive 샘플 중 100%가 정확히 예측
- **f1-score: 0.97:** 매우 높은 값으로, 모델이 positive 클래스에서 매우 잘 작동하는 것을 알 수있음

전체 정확도: 0.93으로, 모델이 전반적으로 높은 정확도를 달성

macro avg: 각 클래스의 평균 성능을 나타내며, 이 값이 낮은 이유는 neutral 클래스의 성능이 매우 저조함

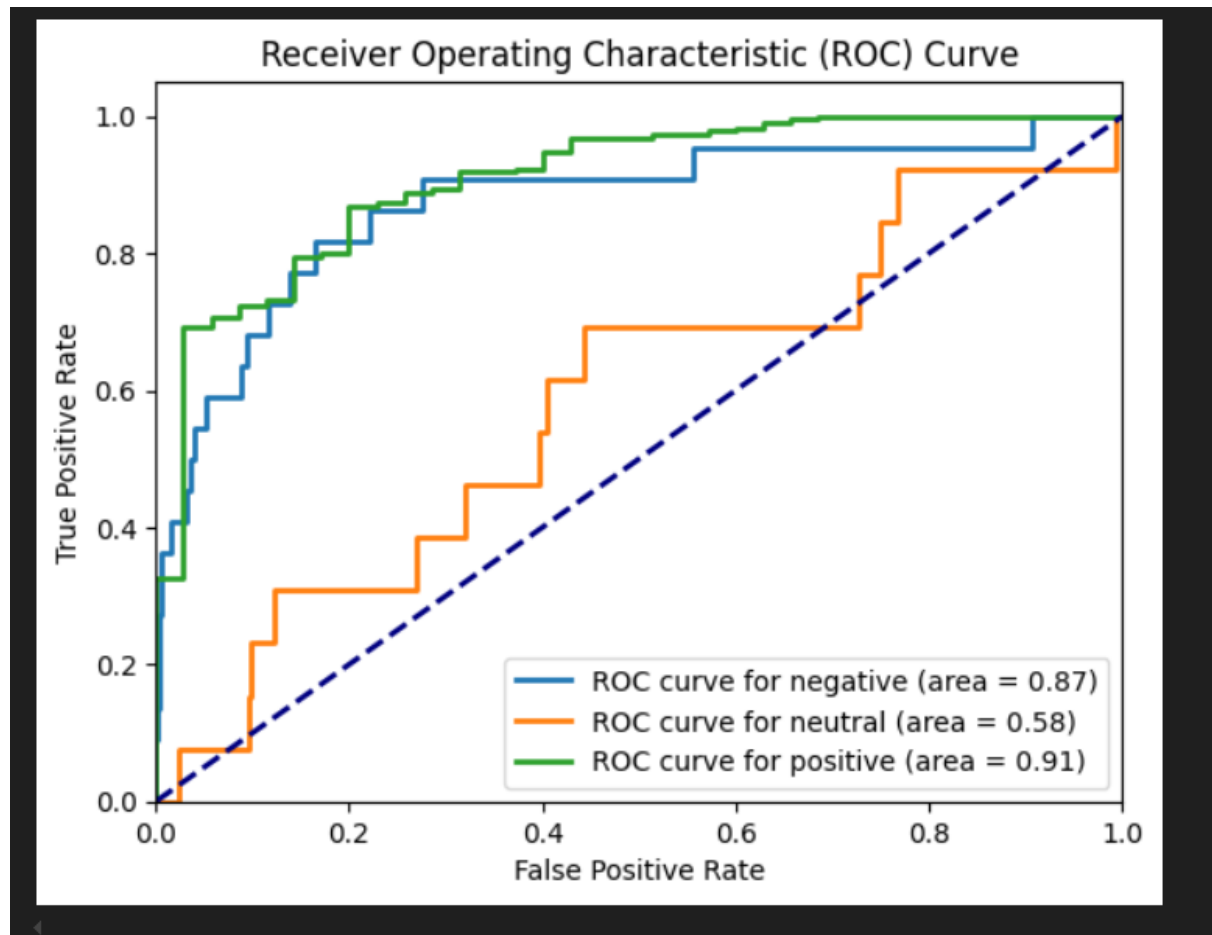
weighted avg: 클래스의 지원 수를 고려한 가중 평균입니다. 전체적인 성능을 더 잘 나타냅니다.

교차 검증 정확도: 각 폴드에서 높은 정확도를 유지하고 있으며, 평균 교차 검증 정확도가 약 92.92%입니다. 이는 모델이 일관되게 잘 작동함을 나타냄

튜닝 결론

- 모델은 positive 클래스에서 매우 잘 작동하며, 전반적인 정확도도 높음
- negative와 neutral 클래스의 성능은 낮으며, 특히 neutral 클래스는 전혀 예측되지 않았음
- 데이터 불균형 문제가 있으며, 이를 해결하기 위해 데이터 증강이나 클래스 가중치를 조정하는 방법을 고려해야함
- 모델의 전반적인 성능은 우수하지만, 모든 클래스에서 균형 잡힌 성능을 얻기 위해 추가적인 조정이 필요하다 판단됨

5 시각화



분석

각 클래스별 ROC 곡선 및 AUC 값

부정 클래스 (Negative)

ROC 곡선: 파란색

AUC: 0.87

분석: AUC 값이 0.87로 높아 부정 클래스에 대한 모델의 성능이 좋음을 나타냄

중립 클래스 (Neutral)

ROC 곡선: 주황색

AUC: 0.58

분석: AUC 값이 0.58로, 중립 클래스에 대한 모델의 성능이 낮음. 이 클래스에 대한 예측 성능이 저조함을 알수있음.

긍정 클래스 (Positive)

ROC 곡선: 초록색

AUC: 0.91

분석: AUC 값이 0.91로 매우 높아 긍정 클래스에 대한 모델의 성능이 매우 좋음을 나타냅니다.

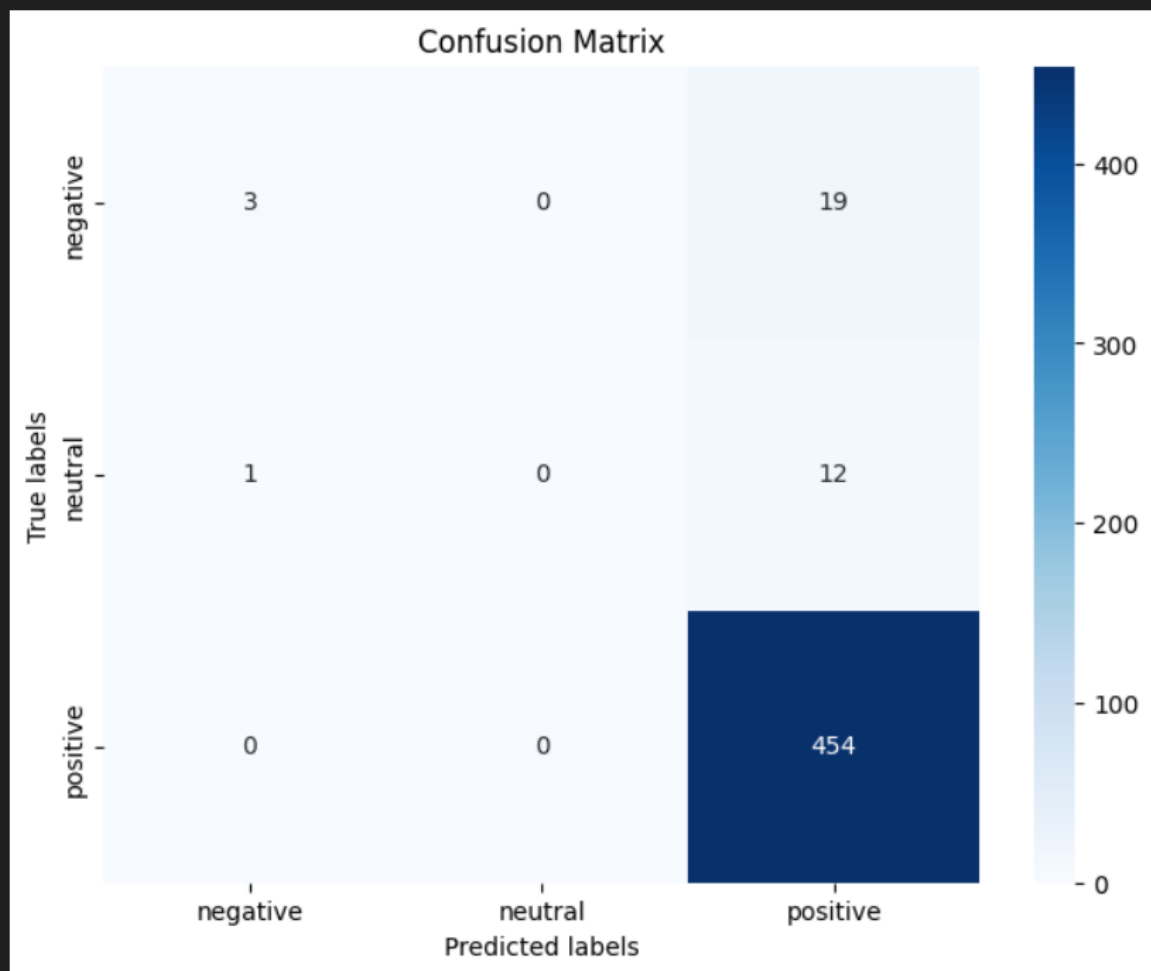
요약 및 결론

긍정 클래스는 AUC 값이 0.91로 가장 높아, 이 클래스에 대한 예측 성능이 매우 뛰어남

부정 클래스도 AUC 값이 0.87로 높아, 모델이 부정 클래스를 비교적 잘 예측하고 있음

중립 클래스는 AUC 값이 0.58로, 모델이 중립 클래스를 잘 예측하지 못지만. 이는 앞서 혼동행렬 분석에서 확인한 바와 같이 중립 클래스를 거의 긍정 클래스로 예측하는 경향과 일치

혼동 행렬 분석



Negative (부정) 클래스:

총 22 개 ($3 + 0 + 19$) 중 실제 부정인 22 개 데이터 중 3 개는 부정으로 정확하게 예측되었고, 19 개는 긍정으로 잘못 예측.
단, 중립으로 예측된 경우는 없음

Neutral (중립) 클래스:

총 13 개 ($1 + 0 + 12$) 중 실제 중립인 13 개 데이터 중 1 개는 부정으로, 12 개는 긍정으로 잘못 예측.
중립으로 예측된 경우는 없음

Positive (긍정) 클래스:

총 454 개 ($0 + 0 + 454$) 중 실제 긍정인 454 개 데이터는 모두 긍정으로 정확하게 예측.