



**Institute for Advanced Studies
in Basic Sciences
GavaZang, Zanjan, Iran**

Sentiment Analysis Model Loader

Prof. Majid Ramezani

Report for a Machine Learning Project <1>

Faeze Ahmadi

14044141

November 2025

تحلیل احساسات متن (Sentiment Analysis)

کد کامل این پروژه و فایل‌های مربوطه در ریپازیتوری گیت‌هاب قابل دسترسی هست:

<https://github.com/Faeze-Ahmadi/sentiment-analysis-project>

این پروژه یک سیستم ساده تحلیل احساسات (Sentiment Analysis) برای تشخیص مثبت یا منفی بودن متنی هست که از کاربر می‌گیریم. ایده اصلی این پروژه اینه که به کامپیوتر یاد بدیم محتوای یک متن رو بخونه و تشخیص بده که متن حس مثبت داره یا منفی؛ مثلاً این متن می‌تونه نظری در مورد فیلم باشه. البته که دیتاست ما هم کامنت‌های مربوط به فیلم و استوری هست. یعنی ما با داشتن حجم زیادی از نظرهای مثبت و منفی و نشون دادن اون‌ها به مدل‌مون آموزش می‌دیم که چی مثبت و چی منفیه، تا وقتی کاربر بهش یه نظر داد، مدل ما خودش پیش‌بینی کنه که اون نظر مثبت یا منفی.

این پروژه دو بخش اصلی داره:

1. تبدیل متن خام به اعداد قابل فهم برای ماشین (به کمک TF-IDF): چون کامپیوتر نمی‌تونه متن رو بخونه و بخاطر همین ما باید متن‌ها رو به عدد تبدیل کنیم و اینطوری مدل‌مون رو آموزش بدیم. TF-IDF مخفف عبارت (Term Frequency – Inverse Document Frequency) هست که به عنوان یک روش معروف برای تبدیل متن به عدد شناخته می‌شه. Term Frequency (فراوانی اصطلاح) یعنی یه متن چندبار داخل متن تکرار شده و Inverse Document Frequency (معکوس فراوانی متن) یعنی این کلمه بین همه متن‌ها چقدر خاص یا کمیاب هست.

2. آموزش یک مدل Machine Learning برای تشخیص احساسات متن‌ها: بحث آموزش دادن هم بر اساس برجسب زدن به متن‌هاست. ما از مدل Logistic Regression برای این کار استفاده می‌کنیم. این مدل، برای classification هست. در واقع ما نیاز داریم که اینجا 0 و 1 ای فکر کنیم، یعنی اگر مدل مثبته پس 1 هست و اگر مدل منفیه پس 0 هست. اما بحث اینجاست که توی متن ما به تحلیل نیاز داریم، پس مدل رگرسیون لجستیک میاد می‌گه متن به دسته 1 تعلق داره یا 0؟.. اینجاست که ما متن یه عدد رو نسبت می‌دیم، مثلاً می‌گیم متن به اندازه 0.87 مثبته؛ پس چون به متن $0.5 <$ هست، مثبته. حالا در ادامه بیشتر در مورد این صحبت می‌کنیم. بعد از آموزش مدل، خود مدل و ابزار تبدیل متن رو ذخیره می‌کنیم تا بتونیم بعدها بدون نیاز به آموزش دوباره، متن‌های جدید رو پیش‌بینی کنیم. یعنی ما یک‌بار دیتاست رو به مدل می‌دیم، اون رو آموزش می‌دیم و بعداً دیگه فقط از مدل استفاده می‌کنیم. چون حجم دیتاها زیاده و آموزش هرباره، هم به زمان و هم به فضا نیاز داره و اصلاً هم کار حرفه‌ای نیست.

ساختار پروژه

فولدر /data

این بخش شامل داده‌های اصلی پروژه هست. دو فولدر جداگونه برای نظرات مثبت و منفی وجود دارن. جمعا 25 هزار متن کوتاه درباره فیلم‌ها در این قسمت قرار داده شده (12,500 متن مثبت و 12,500 متن منفی) و مدل با استفاده از همین داده‌ها آموزش می‌بینه.

/neg ← شامل متن‌های منفی

/pos ← شامل متن‌های مثبت

فایل model.pkl

توی این فایل، مدل آموزش داده شده ذخیره شده.

این مدل همون چیزیه که بعد از آموزش می‌تونه تشخیص بده که متن مثبت یا منفی.

فایل vectorizer.pkl

توی این فایل، TF-IDF vectorizer ذخیره شده.

این ابزار وظیفه داره متن خام رو به یه بردار عددی تبدیل کنه تا مدل اون رو بفهمه. بدون این فایل، مدل نمی‌تونه متن جدید رو تحلیل کنه.

فولدر /src

توی این فولدر تمام فایل‌های مربوط به منطق برنامه وجود دارن. هر فایل به وظیفه مشخص داره.

- فایل load_data.py

وظیفه این فایل خوندن داده‌ها از فولدر `data` و برچسب‌گذاری آنهاست.

به زبان ساده، این فایل متن‌ها رو می‌خونه و مشخص می‌کنه که کدومشون مثبت و کدومشون منفی هستن.

- فایل vectorize.py

توی این فایل متون خام به مقادیر عددی تبدیل می‌شن.

این مرحله شامل ساخت TF-IDF vectorizer و تقسیم داده‌ها به train و test هست

یعنی این فایل متن رو به چیزی که مدل بتونه پردازش کنه تبدیل می‌کنه.

- فایل train_model.py

این فایل مدل ماشین لرنینگ رو آموزش می‌ده.

مدلی که استفاده شده Logistic Regression هست.

اینجا مدل با استفاده از داده‌ها یاد می‌گیره که چطور احساس یک متن رو تشخیص بده.

- فایل save_model.py

توی این فایل مدل آموزش داده شده و vectorizer توی دو فایل جداگانه ذخیره می‌شن ('model.pkl' و 'vectorizer.pkl').

این کار باعث می‌شه که دفعه بعد نیاز نباشه مدل دوباره آموزش داده بشه.

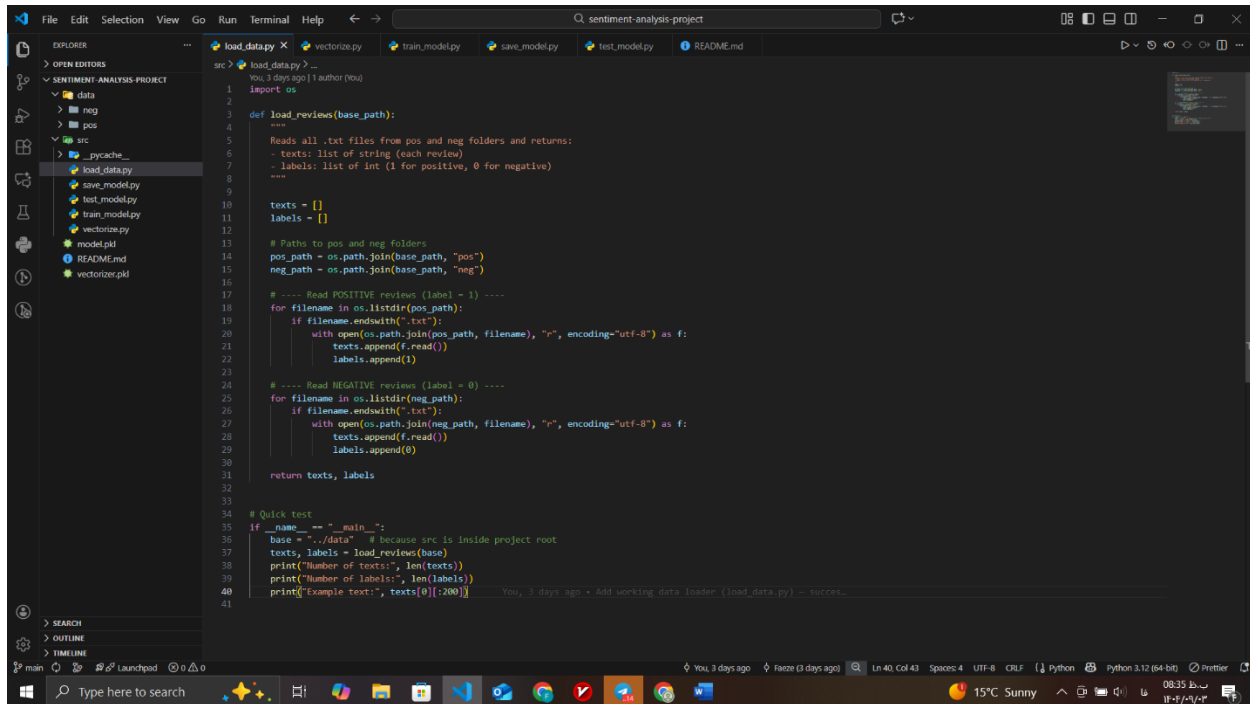
- فایل test_model.py

این فایل یه برنامه تعاملی ساده‌ست؛ برای اینکه ما مدلمون رو تست کنیم که درست کار می‌کنه یا نه. یعنی این فایل نتیجه نهایی

پروژه رو به ما نشان می‌ده.

به این صورت که کاربر یه جمله وارد می‌کنه و مدل پیش‌بینی می‌کنه که جمله مثبت هست یا منفی.

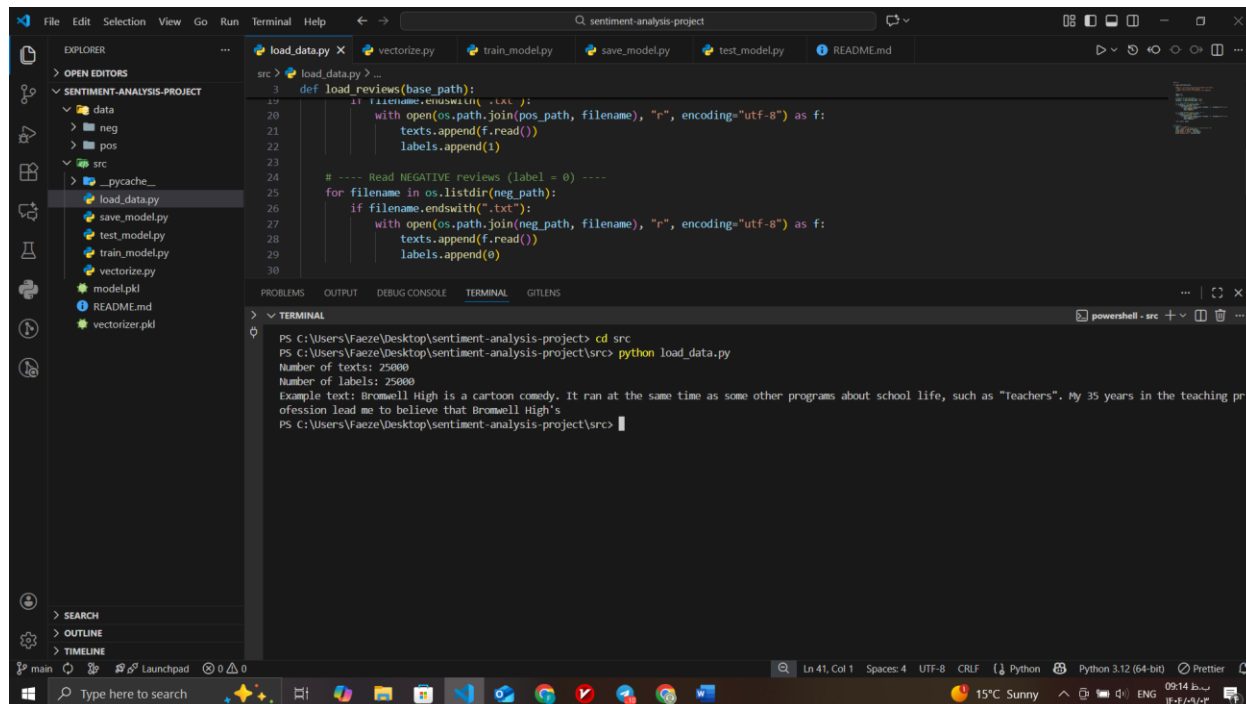
مرحله اول: فایل load_data.py



```
1  src > load_data.py > ...
2  You, 3 days ago | 1 author (You)
3  import os
4
5  def load_reviews(base_path):
6      """
7      Reads all .txt files from pos and neg folders and returns:
8      - texts: list of string (each review)
9      - labels: list of int (1 for positive, 0 for negative)
10     """
11
12     texts = []
13     labels = []
14
15     # Paths to pos and neg folders
16     pos_path = os.path.join(base_path, "pos")
17     neg_path = os.path.join(base_path, "neg")
18
19     # ---- Read POSITIVE reviews (label = 1) ----
20     for filename in os.listdir(pos_path):
21         if filename.endswith(".txt"):
22             with open(os.path.join(pos_path, filename), "r", encoding="utf-8") as f:
23                 texts.append(f.read())
24                 labels.append(1)
25
26     # ---- Read NEGATIVE reviews (label = 0) ----
27     for filename in os.listdir(neg_path):
28         if filename.endswith(".txt"):
29             with open(os.path.join(neg_path, filename), "r", encoding="utf-8") as f:
30                 texts.append(f.read())
31                 labels.append(0)
32
33     return texts, labels
34
35 # Quick test
36 if __name__ == "__main__":
37     base = "../data" # because src is inside project root
38     texts, labels = load_reviews(base)
39     print("Number of texts:", len(texts))
40     print("Number of labels:", len(labels))
41     print("Example text:", texts[0][:200])
```

توی این کد، ابتدا کتابخانه استاندارد پایتون به نام `os` وارد شده که برای کار با مسیرها و فایل‌ها کاربرد داره. سپس تابعی به نام `load_reviews` تعریف شده که وظیفه‌اش خوندن همه فایل‌های متنی با پسوند `.txt` از دو فولدر `pos` و `neg` هست؛ این دو فولدر شامل نظرات مثبت و منفی هستن. ورودی این تابع مسیر پایه‌ای هست که این دو فولدر داخل اون قرار دارن. درون تابع، دو لیست خالی به نام‌های `texts` و `labels` ساخته شدن تا متن‌ها و برچسب‌هاشون (مثبت یا منفی بودن) توی اون‌ها ذخیره شدن. برای اینکه مسیر دقیق به فولدرهای `pos` و `neg` ساخته بشه، از تابع `os.path.join` استفاده شده تا کد مستقل از سیستم عامل و قابل اجرا روی ویندوز یا لینوکس باشه. بعد، با کمک تابع `os.listdir` لیست همه فایل‌های داخل پوشه مثبت خونده می‌شن و هر فایل با پسوند `.txt` باز می‌شه؛ متن داخل فایل با انکودینگ UTF-8 خونده شده و به لیست `texts` اضافه می‌شه و عدد 1 به عنوان برچسب مثبت به لیست `labels` افزوده می‌شه. همین روند برای پوشه نظرات منفی هم تکرار شده با این تفاوت که برچسب صفر (نمایانگر نظر منفی) به لیست برچسب‌ها اضافه می‌شه. در نهایت، دو لیست `texts` و `labels` که به ترتیب شامل تمام متن‌ها و برچسب‌های اون‌ها هستن، برگردونده می‌شن. در انتهای کد یک بخش تست سریع نوشته شده که در صورت اجرای مستقیم فایل، مسیر داده‌ها به صورت نسبی تنظیم می‌شه، سپس تابع `load_reviews` فراخوانی شده و تعداد کل متون و برچسب‌ها چاپ می‌شه، همچنین نمونه‌ای از متن اولین نظر هم برای اطمینان از صحت خوندن داده‌ها نمایش داده می‌شه. به طور کلی این کد وظیفه آماده‌سازی داده‌های اولیه پروژه رو بر عهده داره تا بعداً بتونیم اون‌ها رو برای آموزش مدل استفاده کنیم.

خروجی فایل load_data.py



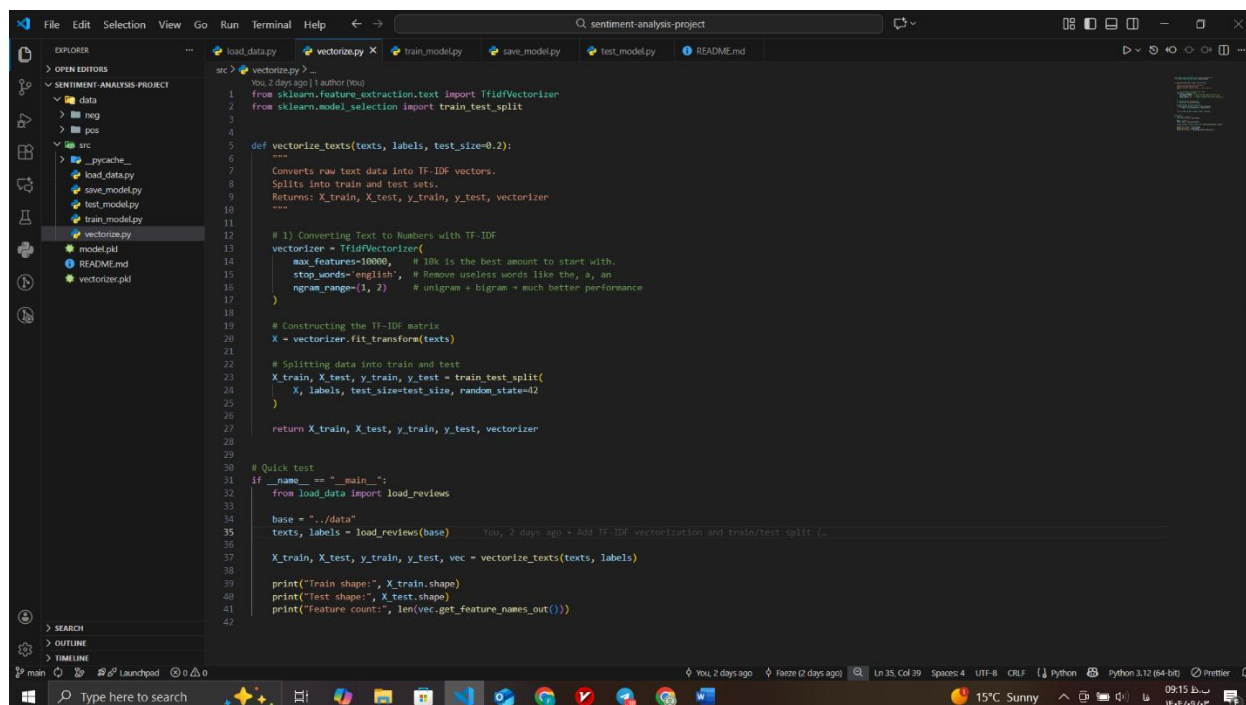
The screenshot shows a Visual Studio Code editor window with the file `load_data.py` open. The file is located in the `src` directory of a project named `sentiment-analysis-project`. The code in `load_data.py` defines a function `load_reviews` that reads reviews from a directory. The terminal output shows the execution of the script, which reads 25000 texts and appends them to a list. The output also shows an example text: "Bromwell High is a cartoon comedy. It ran at the same time as some other programs about school life, such as 'Teachers'. My 35 years in the teaching profession lead me to believe that Bromwell High's".

```
src> load_data.py
3 def load_reviews(base_path):
4     # Read POSITIVE reviews (label = 1) ---
5     for filename in os.listdir(pos_path):
6         if filename.endswith('.txt'):
7             with open(os.path.join(pos_path, filename), "r", encoding="utf-8") as f:
8                 texts.append(f.read())
9                 labels.append(1)
10
11 # --- Read NEGATIVE reviews (label = 0) ---
12 for filename in os.listdir(neg_path):
13     if filename.endswith('.txt'):
14         with open(os.path.join(neg_path, filename), "r", encoding="utf-8") as f:
15             texts.append(f.read())
16             labels.append(0)
17
18 # Print the number of texts and labels
19 print("Number of texts: %d" % len(texts))
20 print("Number of labels: %d" % len(labels))
21
22 # Print an example text
23 print("Example text: %s" % texts[0])
24
25 # Save the data to a file
26 save_data(texts, labels)
```

```
PS C:\Users\Faeze\Desktop\sentiment-analysis-project> cd src
PS C:\Users\Faeze\Desktop\sentiment-analysis-project\src> python load_data.py
Number of texts: 25000
Number of labels: 25000
Example text: Bromwell High is a cartoon comedy. It ran at the same time as some other programs about school life, such as "Teachers". My 35 years in the teaching pr
ofession lead me to believe that Bromwell High's
PS C:\Users\Faeze\Desktop\sentiment-analysis-project\src>
```

خروجی به ما نشون می‌ده که تعداد متن‌ها و برچسب‌ها چقدره و این نتیجه لود شدن کامل دیتاست ماست. یه نمونه تکست هم به عنوان مثال نشون می‌ده که بارگذاری کامل رو تایید کنه.

مرحله دوم: vectorize.py



```
1 from sklearn.feature_extraction.text import TfidfVectorizer
2 from sklearn.model_selection import train_test_split
3
4
5
6
7 def vectorize_texts(texts, labels, test_size=0.2):
8     """
9     Converts raw text data into TF-IDF vectors.
10    Splits into train and test sets.
11    Returns: X_train, X_test, y_train, y_test, vectorizer
12    """
13
14    # 1) Converting Text to Numbers with TF-IDF
15    vectorizer = TfidfVectorizer(
16        max_features=10000, # 10k is the best amount to start with.
17        stop_words='english', # Remove useless words like the, a, an
18        ngram_range=(1, 2) # unigram + bigram = much better performance
19    )
20
21    # Constructing the TF-IDF matrix
22    X = vectorizer.fit_transform(texts)
23
24    # Splitting data into train and test
25    X_train, X_test, y_train, y_test = train_test_split(
26        X, labels, test_size=test_size, random_state=42
27    )
28
29    return X_train, X_test, y_train, y_test, vectorizer
30
31 # Quick test
32 if __name__ == "__main__":
33     from load_data import load_reviews
34
35     base = "../data"
36     texts, labels = load_reviews(base)
37     X_train, X_test, y_train, y_test, vec = vectorize_texts(texts, labels)
38
39     print("Train shape:", X_train.shape)
40     print("Test shape:", X_test.shape)
41     print("Feature count:", len(vec.get_feature_names_out()))
```

این کد وظیفه داره که متن‌های خام رو به شکل بردارهای عددی TF-IDF تبدیل کنه تا مدل‌های یادگیری ماشین بتونن اون‌ها رو پردازش کنن. برای این کار از کلاس `TfidfVectorizer` استفاده شده که ویژگی‌هایی مثل حذف کلمات اضافه مثل a و the ... (stop words)، استفاده از unigram و bigram، و محدود کردن تعداد ویژگی‌ها به 10 هزار تا رو داره تا کیفیت و سرعت کار بهینه باشه. بعد از ساخت ماتریس TF-IDF، داده‌ها به دو بخش آموزش و تست تقسیم میشن تا مدل بتونه روی آموزش یاد بگیره و روی تست ارزیابی بشه.

Unigram چیست؟

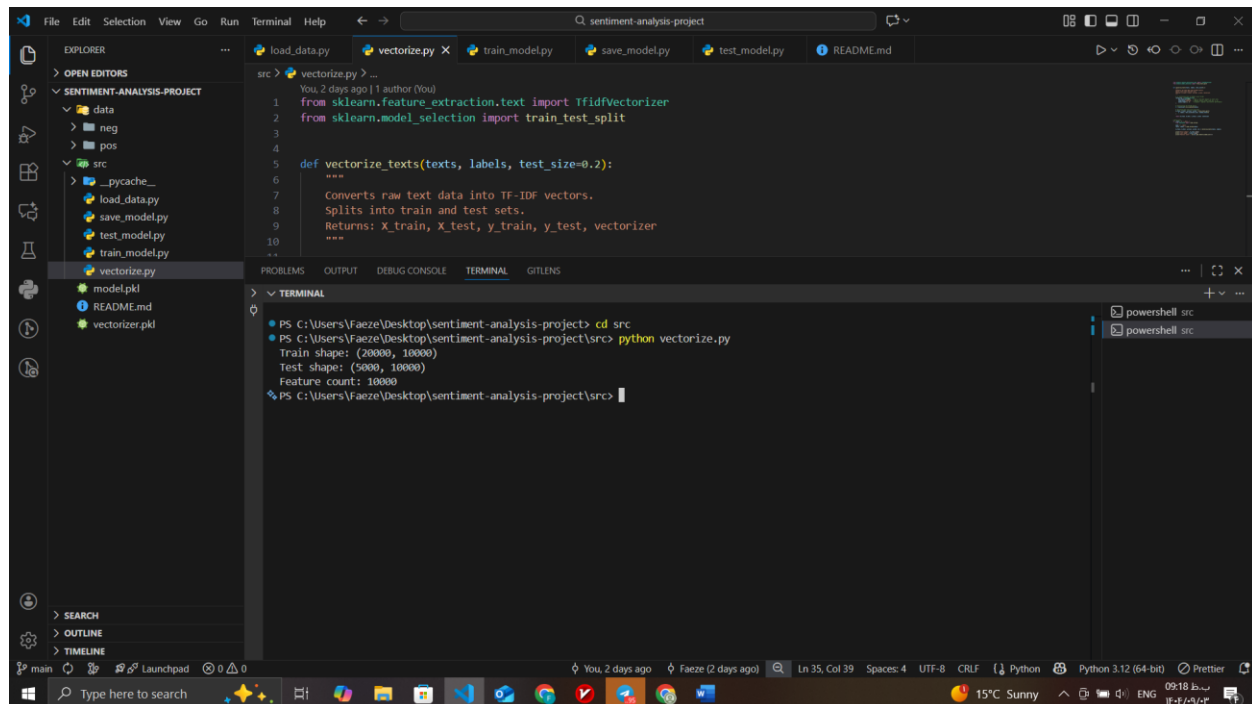
Unigram یعنی وقتی متن رو به کلمات جداگونه تقسیم می‌کنیم، هر کلمه به تنهایی یک واحد (Token) محسوب می‌شه. مثلاً جمله "I love movies" رو وقتی با Unigram نگاه کنیم، سه کلمه داریم: "I"، "love" و "movies". پس هر کلمه به صورت جداگونه بررسی می‌شه و ویژگی‌ها براساس همین تک کلمه‌ها ساخته می‌شن.

Bigram چیست؟

Bigram یعنی به جای نگاه کردن فقط به تک کلمه‌ها، به جفت‌های متوالی کلمات نگاه می‌کنیم. مثلاً همون جمله "I love movies" رو با Bigram به این صورت تقسیم می‌کنیم: "I love" و "love movies". در واقع، ترکیب دو کلمه پشت سر هم رو به عنوان یه واحد در نظر می‌گیریم.

حالا چرا از Bigram استفاده می‌شه؟ چون گاهی معنای یک کلمه به تنهایی کافی نیست و ترکیب دو کلمه مفهوم دقیق‌تری به ما میدن. مثلاً توی جمله "not good"، اگه فقط به Unigram نگاه کنیم، کلمه "good" بار مثبت دارد. ولی "not good" ترکیبیه که معنی منفی می‌ده. Bigram کمک می‌کنه مدل این ترکیب‌ها رو بشناسه و نتیجه بهتری به ما بده.

خروجی فایل vectorize.py



```
src > vectorize.py > ...
1 from sklearn.feature_extraction.text import TfidfVectorizer
2 from sklearn.model_selection import train_test_split
3
4
5 def vectorize_texts(texts, labels, test_size=0.2):
6     """
7     Converts raw text data into TF-IDF vectors.
8     Splits into train and test sets.
9     Returns: X_train, X_test, y_train, y_test, vectorizer
10    """
11    ...
```

```
PS C:\Users\Faeze\Desktop\sentiment-analysis-project> cd src
PS C:\Users\Faeze\Desktop\sentiment-analysis-project\src> python vectorize.py
train shape: (20000, 10000)
Test shape: (5000, 10000)
Feature count: 10000
PS C:\Users\Faeze\Desktop\sentiment-analysis-project\src> |
```

خروجی تابع پنج تا متغیر هست: بردارهای آموزش و تست (X_{train} , X_{test})، برچسب‌های مربوط به هر کدوم (y_{train} , y_{test}) و خود شیء `vectorizer` که برای تبدیل متن‌های جدید به بردارهای TF-IDF لازم هست. به طور خلاصه، این کد متون رو آماده و تقسیم‌بندی می‌کنه تا برای آموزش و ارزیابی مدل آماده باشن.

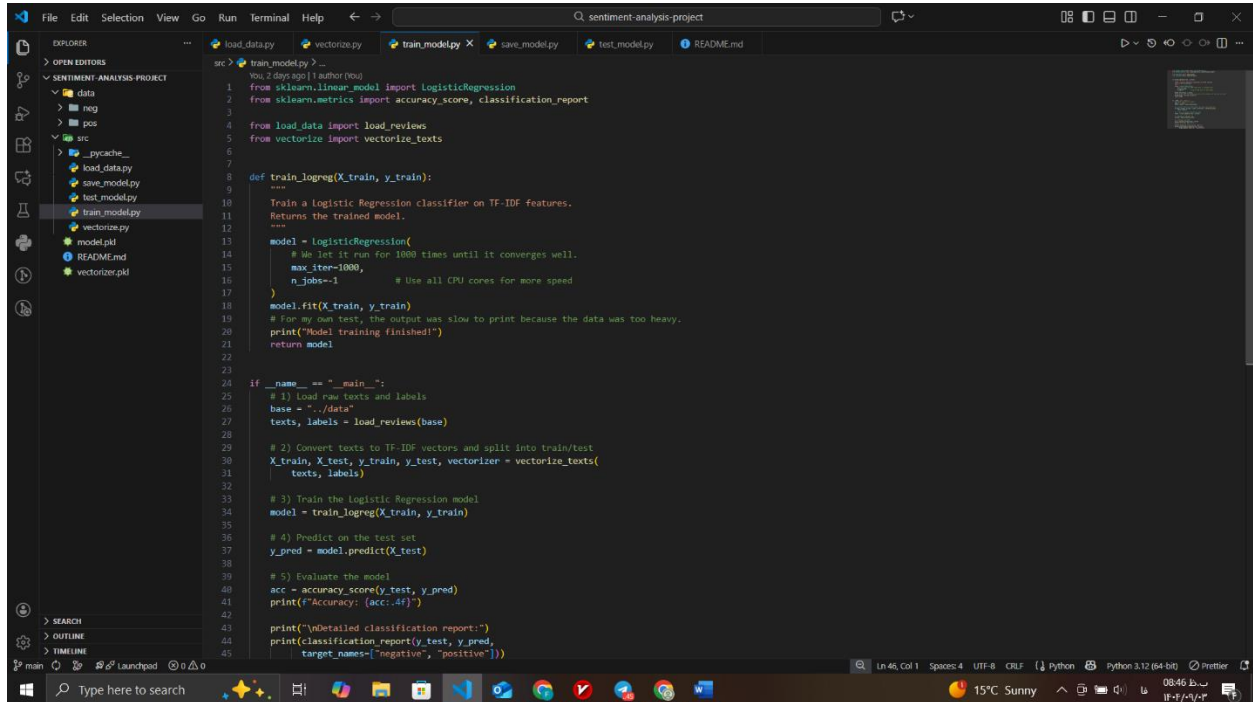
X_{train} : داده‌های آموزش.

X_{test} : داده‌های تست.

y_{train} : برچسب‌های آموزش.

y_{test} : برچسب‌های تست.

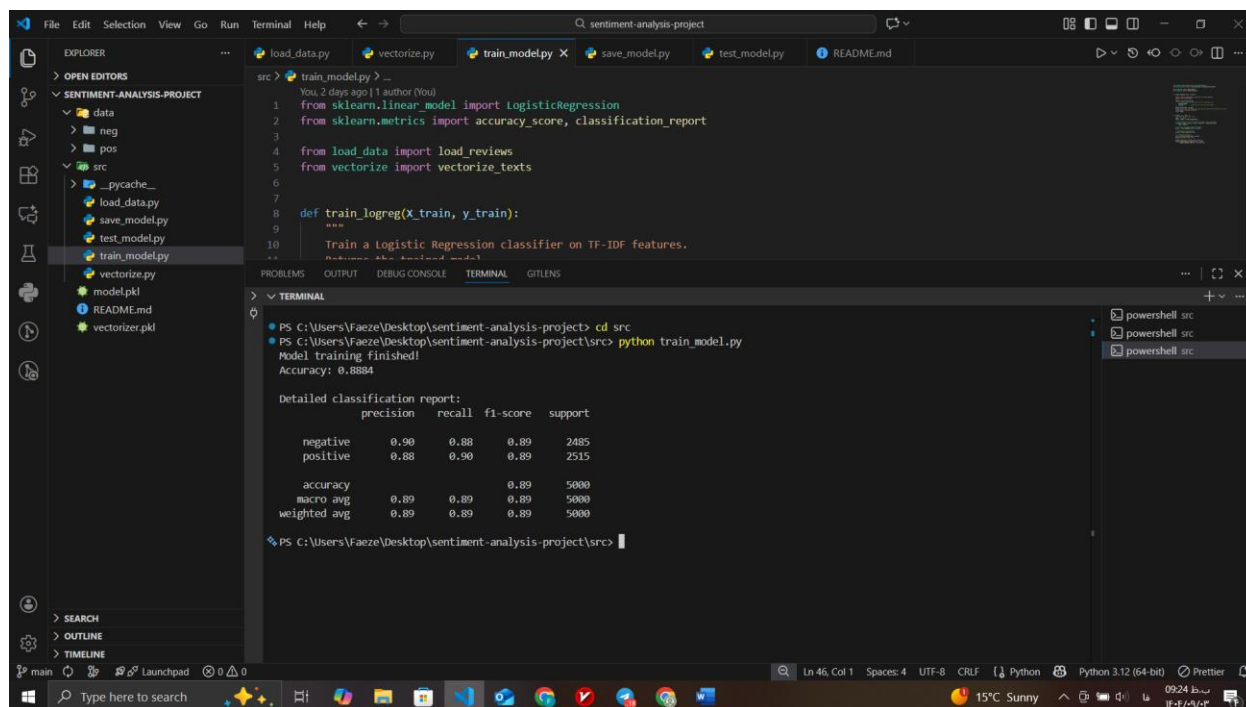
مرحله سوم: فایل train_model.py



```
1 from sklearn.linear_model import LogisticRegression
2 from sklearn.metrics import accuracy_score, classification_report
3
4 from load_data import load_reviews
5 from vectorize import vectorize_texts
6
7
8 def train_logreg(X_train, y_train):
9     """
10     Train a logistic Regression classifier on TF-IDF features.
11     Returns the trained model.
12     """
13     model = LogisticRegression()
14     # We let it run for 1000 times until it converges well.
15     max_iter=1000,
16     n_jobs=-1 # Use all CPU cores for more speed
17 )
18 model.fit(X_train, y_train)
19 # For my own test, the output was slow to print because the data was too heavy.
20 print("Model training finished!")
21 return model
22
23
24 if __name__ == "__main__":
25     # 1) Load raw texts and labels
26     base = "../data"
27     texts, labels = load_reviews(base)
28
29     # 2) Convert texts to TF-IDF vectors and split into train/test
30     X_train, X_test, y_train, y_test, vectorizer = vectorize_texts(
31         texts, labels)
32
33     # 3) Train the logistic Regression model
34     model = train_logreg(X_train, y_train)
35
36     # 4) Predict on the test set
37     y_pred = model.predict(X_test)
38
39     # 5) Evaluate the model
40     acc = accuracy_score(y_test, y_pred)
41     print(f"Accuracy: {acc:.4f}")
42
43     print("\nDetailed classification report:")
44     print(classification_report(y_test, y_pred,
45                               target_names=["negative", "positive"]))
```

این کد برای آموزش یک مدل Logistic Regression استفاده می‌شود که روی ویژگی‌های TF-IDF کار می‌کند. اول با استفاده از تابع `load_reviews` داده‌های متنی و برچسب‌ها رو از پوشه داده‌ها می‌خونیم. بعدش متن‌ها رو با تابع `vectorize_texts` تبدیل به بردارهای عددی TF-IDF می‌کنیم و داده‌ها رو به مجموعه آموزش و تست تقسیم می‌کنیم. توی تابع `train_logreg` یک مدل Logistic Regression ساخته شده که اجازه می‌ده تا 1000 بار روی داده‌ها تکرار کنه تا به بهترین حالت برسه و از تمام هسته‌های پردازنده برای سرعت بیشتر استفاده می‌کند. مدل روی داده‌های آموزشی آموزش داده می‌شه و پس از اتمام، پیام "Model training finished" نمایش داده می‌شه. سپس مدل روی داده‌های تست پیش‌بینی انجام می‌ده و دقت (accuracy) مدل محاسبه و چاپ می‌شه. در نهایت گزارش جزئی‌تری از عملکرد مدل با استفاده از `classification_report` نمایش داده می‌شه که نشون میده مدل چقدر در تشخیص درست نظرات مثبت و منفی موفق بوده. این کد به طور کلی مراحل اصلی آموزش، تست و ارزیابی مدل رو به صورت ساده و شفاف انجام می‌ده.

خروجی فایل train_model.py



```
src > train_model.py > _
You, 2 days ago | 1 author (You)
1 from sklearn.linear_model import LogisticRegression
2 from sklearn.metrics import accuracy_score, classification_report
3
4 from load_data import load_reviews
5 from vectorize import vectorize_texts
6
7
8 def train_logreg(X_train, y_train):
9     """
10    Train a Logistic Regression classifier on TF-IDF features.
11    """
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60
61
62
63
64
65
66
67
68
69
70
71
72
73
74
75
76
77
78
79
80
81
82
83
84
85
86
87
88
89
90
91
92
93
94
95
96
97
98
99
100
```

Model training finished!
Accuracy: 0.8884

Detailed classification report:

	precision	recall	f1-score	support
negative	0.90	0.88	0.89	2485
positive	0.88	0.90	0.89	2515
accuracy			0.89	5000
macro avg	0.89	0.89	0.89	5000
weighted avg	0.89	0.89	0.89	5000

خروجی فایل 'train_model.py' در اصل مدل آموزش داده شده Logistic Regression هست که آماده استفاده برای پیش‌بینی احساسات متون جدید. علاوه بر این، هنگام اجرای فایل، دقت مدل روی داده‌های تست چاپ می‌شود و گزارشی کامل از عملکرد مدل شامل معیارهایی مثل دقت، بازیابی و F1-score برای هر کلاس (مثبت و منفی) نمایش داده می‌شود. این خروجی‌ها به ما کمک می‌کنند که بفهمیم مدل چقدر خوب کار می‌کند و چقدر قابل اعتماد است.

Precision (دقت):

می‌گه از بین تمام نمونه‌هایی که مدل گفته «مثبت» یا «یک» هستند، چند درصد واقعا درست پیش‌بینی شدن. یعنی وقتی مدل میگه «مثبت»، چقدر مطمئنیم که واقعا مثبت باشه.

Recall (بازخوانی یا حساسیت):

می‌گه از بین همه نمونه‌های واقعی مثبت، مدل چقدر توانسته درست تشخیص بده. یعنی مدل چقدر خوب نمونه‌های مثبت واقعی رو پیدا کرده.

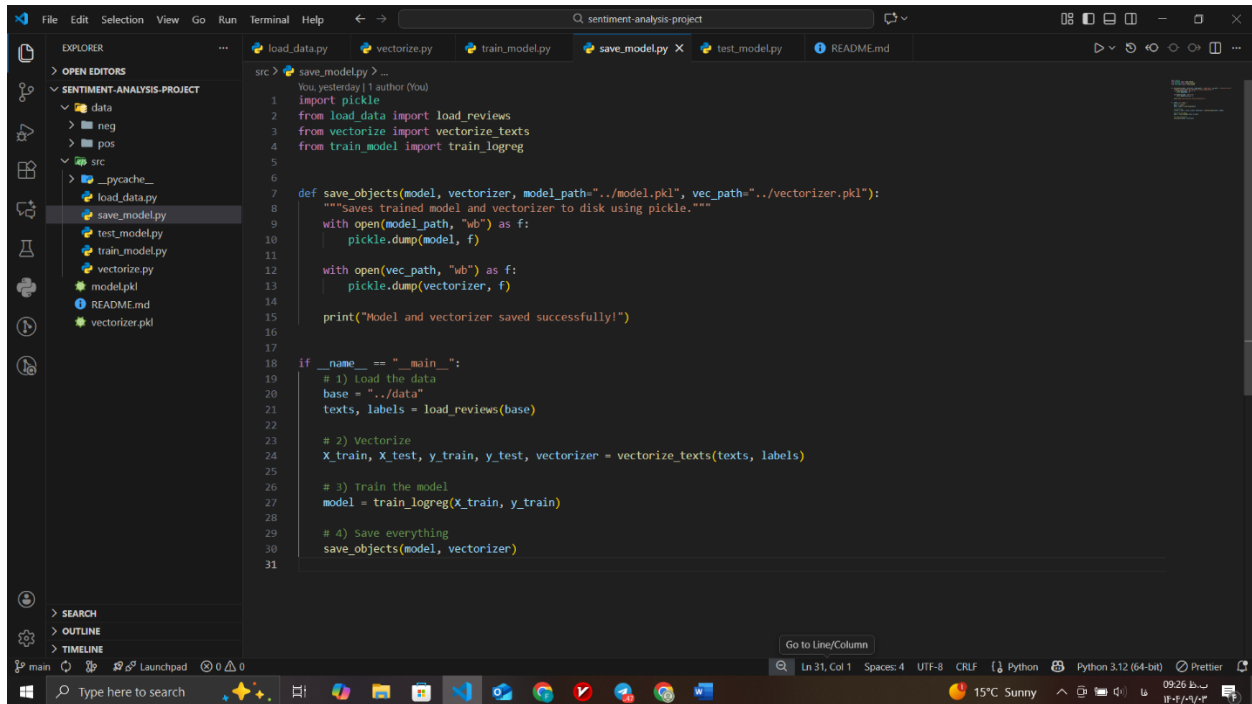
F1-score (امتیاز F1):

میانگینی بین Precision و Recall هست که تعادل بین این دو رو نشون می‌ده. وقتی هم دقت بالا باشه و هم بازخوانی بالا، F1 هم بالا میره.

Support (حمایت):

تعداد واقعی نمونه‌هایی که هر کلاس داره. یعنی مثلا چند داده واقعی مثبت داریم و چند داده واقعی منفی. این عدد فقط برای اطلاع هست و معیار عملکرد نیست.

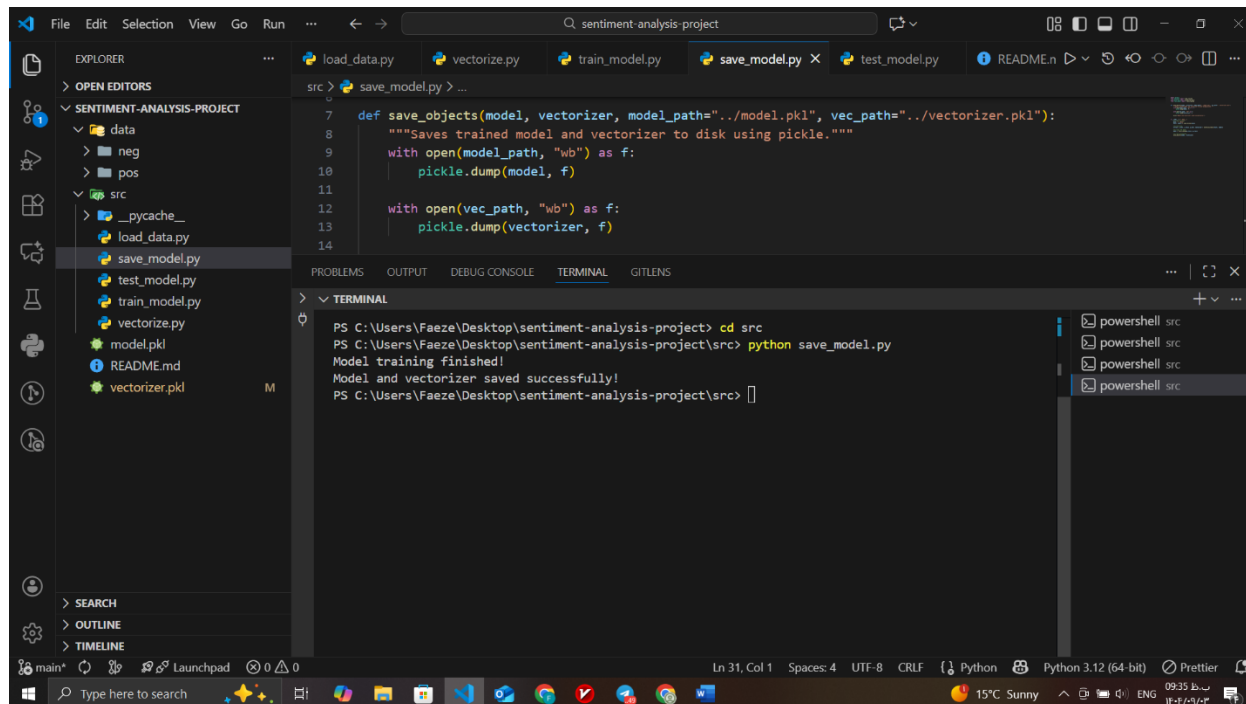
مرحله چهارم: فایل save_model.py



```
src> save_model.py> ...
1 You, yesterday | 1 author (You)
2 import pickle
3 from load_data import load_reviews
4 from vectorize import vectorize_texts
5 from train_model import train_logreg
6
7 def save_objects(model, vectorizer, model_path='../model.pkl', vec_path='../vectorizer.pkl'):
8     """Saves trained model and vectorizer to disk using pickle."""
9     with open(model_path, "wb") as f:
10         pickle.dump(model, f)
11
12     with open(vec_path, "wb") as f:
13         pickle.dump(vectorizer, f)
14
15     print("Model and vectorizer saved successfully!")
16
17
18 if __name__ == "__main__":
19     # 1) Load the data
20     base = "../data"
21     texts, labels = load_reviews(base)
22
23     # 2) Vectorize
24     X_train, X_test, y_train, y_test, vectorizer = vectorize_texts(texts, labels)
25
26     # 3) Train the model
27     model = train_logreg(X_train, y_train)
28
29     # 4) Save everything
30     save_objects(model, vectorizer)
31
```

این کد به پروژه ساده برای آموزش مدل تشخیص احساسات متون است که مراحل اصلی رو از بارگذاری داده‌ها تا ذخیره مدل نهایی پوشش می‌ده. ابتدا با استفاده از تابع 'load_reviews'، متن‌های خام و برچسب‌های مربوط به اون‌ها (مثبت یا منفی) از پوشه داده‌ها خوانده میشن. سپس این متن‌ها به کمک تابع 'vectorize_texts' به بردارهای عددی TF-IDF تبدیل می‌شن و داده‌ها به دو بخش آموزش و تست تقسیم می‌شن تا مدل بتونه روی داده‌های آموزشی یاد بگیره و روی داده‌های تست ارزیابی بشه. در مرحله بعد، مدل Logistic Regression با تابع 'train_logreg' روی داده‌های آموزش آموزش داده می‌شه تا بتونه الگوهای موجود در داده‌ها رو یاد بگیره. در نهایت، برای اینکه نیاز نباشه هر بار برنامه اجرا بشه و مدل دوباره آموزش ببینه، مدل آموزش‌دیده و همچنین بردارساز TF-IDF به کمک کتابخانه 'pickle' به صورت باینری در فایل‌های جداگانه ذخیره می‌شن ('model.pkl' برای مدل و 'vectorizer.pkl' برای بردارساز). این کار باعث می‌شه که در دفعات بعدی فقط با بارگذاری این فایل‌ها، بتوان مدل رو بدون آموزش مجدد به کار گرفت.

خروجی فایل save_model.py



The screenshot shows the Visual Studio Code interface for a project named 'sentiment-analysis-project'. The Explorer sidebar on the left shows the project structure with folders 'data' (containing 'neg' and 'pos') and 'src' (containing 'load_data.py', 'save_model.py', 'test_model.py', 'train_model.py', 'vectorize.py', 'model.pkl', 'README.md', and 'vectorizer.pkl'). The main editor window displays the 'save_model.py' file with the following code:

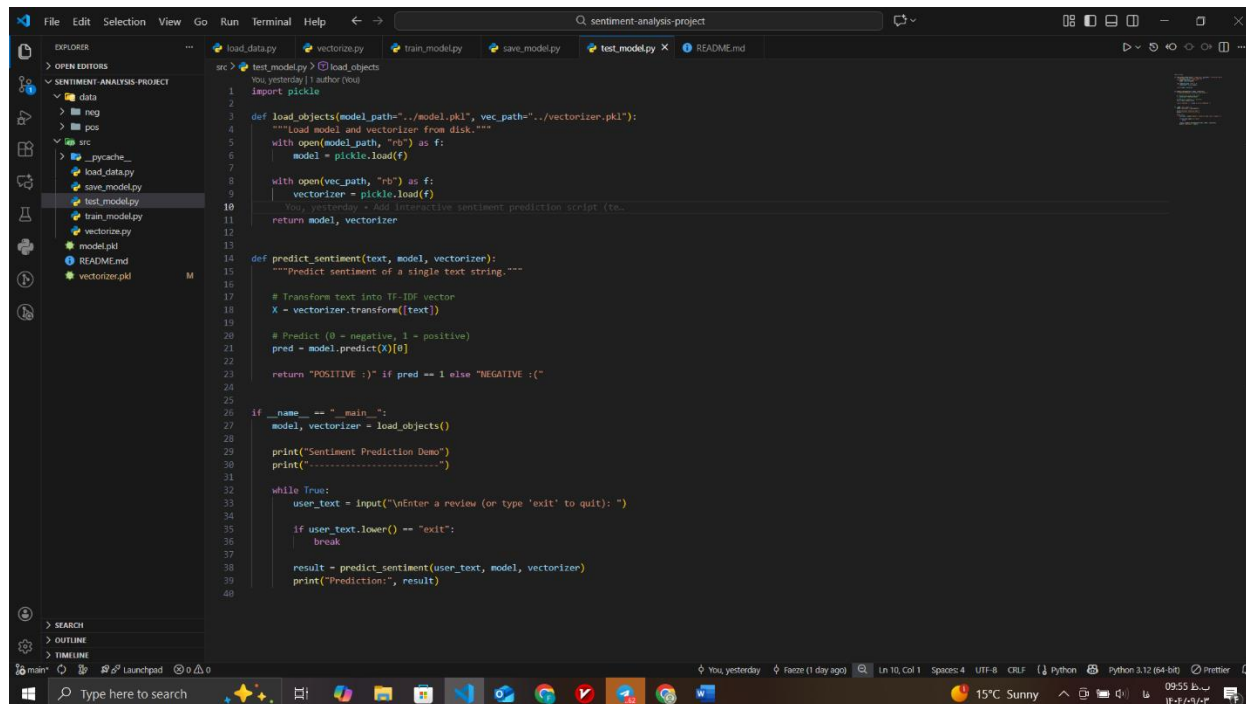
```
7 def save_objects(model, vectorizer, model_path='../model.pkl', vec_path='../vectorizer.pkl'):  
8     """Saves trained model and vectorizer to disk using pickle."""  
9     with open(model_path, "wb") as f:  
10         pickle.dump(model, f)  
11  
12     with open(vec_path, "wb") as f:  
13         pickle.dump(vectorizer, f)  
14
```

The TERMINAL panel at the bottom shows the execution of the script:

```
PS C:\Users\Faeze\Desktop\sentiment-analysis-project> cd src  
PS C:\Users\Faeze\Desktop\sentiment-analysis-project\src> python save_model.py  
Model training finished!  
Model and vectorizer saved successfully!  
PS C:\Users\Faeze\Desktop\sentiment-analysis-project\src>
```

خروجی نهایی این برنامه دو فایل ذخیره شده روی دیسک هست که شامل مدل آموزش دیده و ابزار تبدیل متن به بردار عددی هست و پایه و اساس مراحل بعدی پیش بینی احساسات متن های جدید هست. به طور خلاصه، این کد از داده های خام شروع می کنه، مدل یادگیری ماشین رو آموزش می ده و اون رو برای استفاده در آینده ذخیره می کنه.

مرحله پنجم: فایل test-model.py



```
1 # test_model.py
2 You yesterday | I author (rou)
3 import pickle
4
5 def load_objects(model_path="model.pkl", vec_path="vectorizer.pkl"):
6     """Load model and vectorizer from disk."""
7     with open(model_path, "rb") as f:
8         model = pickle.load(f)
9
10    with open(vec_path, "rb") as f:
11        vectorizer = pickle.load(f)
12    return model, vectorizer
13
14 def predict_sentiment(text, model, vectorizer):
15     """Predict sentiment of a single text string."""
16
17     # Transform text into TF-IDF vector
18     X = vectorizer.transform([text])
19
20     # Predict (0 = negative, 1 = positive)
21     pred = model.predict(X)[0]
22
23     return "POSITIVE :)" if pred == 1 else "NEGATIVE :("
24
25 if __name__ == "__main__":
26     model, vectorizer = load_objects()
27
28     print("Sentiment Prediction Demo")
29     print("-----")
30
31     while True:
32         user_text = input("Enter a review (or type 'exit' to quit): ")
33
34         if user_text.lower() == "exit":
35             break
36
37         result = predict_sentiment(user_text, model, vectorizer)
38         print("Prediction:", result)
```

توی این فایل عملاً مرحله «استفاده از مدل» انجام می‌شه؛ یعنی جایی که دیگر قرار نیست مدل رو آموزش بدیم، بلکه می‌خواهیم مدل ذخیره‌شده رو بالا ببریم و روی متن‌های جدید پیش‌بینی انجام بدیم. همه‌چیز از دو فایل 'model.pkl' و 'vectorizer.pkl' شروع می‌شه؛ این دو فایل همون مدل آموزش‌داده‌شده و بردارساز TF-IDF هستن که قبلاً ذخیره‌شون کرده بودیم.

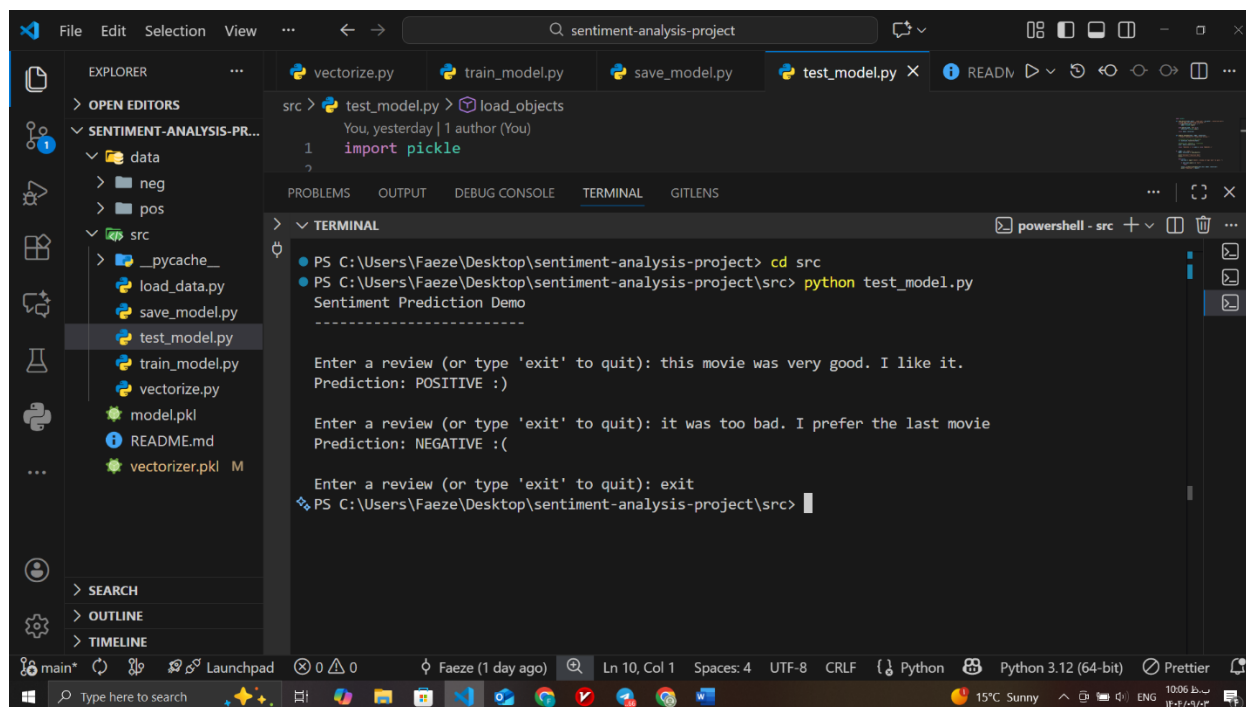
در ابتدای کد، تابع 'load_objects' با استفاده از کتابخانه pickle کارش اینه که این دو فایل رو از دیسک باز کنه و دوباره تبدیلشون کنه به دو شیء قابل استفاده توی پایتون؛ یکی مدل و یکی vectorizer. بدون این مرحله، هر بار باید مدل رو از صفر آموزش می‌دادیم، اما با لود کردن، مستقیم می‌تونیم از اون استفاده کنیم.

بعد از این، تابع 'predict_sentiment' نوشته شده که وظیفه‌اش گرفتن یه متن خام و تبدیل اون به یه بردار TF-IDF (با استفاده از همان vectorizer ذخیره‌شده) هست. سپس مدل روی این بردار پیش‌بینی انجام می‌ده و مشخص می‌کنه که این متن مثبت یا منفی. خروجی تابع به صورت یه رشته ساده نمایش داده می‌شه:

"POSITIVE :)" یا "NEGATIVE :("

توی بخش اصلی برنامه ('__if __name__ == "__main__"')، ابتدا مدل و بردارساز با 'load_objects' لود می‌شن. سپس یک حلقه تعاملی شروع می‌شه که کاربر می‌تونه هر جمله‌ای تایپ کنه و برنامه فوراً احساس متن رو پیش‌بینی و چاپ می‌کنه. اگر کاربر کلمه 'exit' رو وارد کنه، برنامه متوقف می‌شه.

خروجی فایل test_model.py



The screenshot shows a Visual Studio Code editor window with the file explorer on the left displaying a project structure for 'SENTIMENT-ANALYSIS-PR...'. The file explorer shows folders 'data' (with 'neg' and 'pos' subfolders) and 'src' (containing '_pycache_', 'load_data.py', 'save_model.py', 'test_model.py', 'train_model.py', 'vectorize.py', 'model.pkl', 'README.md', and 'vectorizer.pkl'). The main editor area shows the 'test_model.py' file with the following code:

```
src > test_model.py > load_objects
You, yesterday | 1 author (You)
1 import pickle
?
```

The terminal window at the bottom shows the execution of the script:

```
PS C:\Users\Faeze\Desktop\sentiment-analysis-project> cd src
PS C:\Users\Faeze\Desktop\sentiment-analysis-project\src> python test_model.py
Sentiment Prediction Demo
-----

Enter a review (or type 'exit' to quit): this movie was very good. I like it.
Prediction: POSITIVE :)

Enter a review (or type 'exit' to quit): it was too bad. I prefer the last movie
Prediction: NEGATIVE :(

Enter a review (or type 'exit' to quit): exit
PS C:\Users\Faeze\Desktop\sentiment-analysis-project\src>
```

خروجی اصلی این فایل برای متنی هست که کاربر وارد می‌کند. یعنی هر بار کاربر به جمله تایپ میکند و برنامه بهش میگه مثبت یا منفی:

...

Prediction: POSITIVE):

...

یا

...

Prediction: NEGATIVE):

...

یعنی خروجی نهایی این فایل، یک سیستم تعاملی تشخیص احساسات هست که از مدل آموزش داده شده استفاده می‌کند و بدون نیاز به آموزش مجدد، روی هر متن دلخواه پیش‌بینی انجام می‌دهد.

توضیحات کلی من راجع به این پروژه تموم شد، برای اطلاع بیشتر از روند پروژه و کدها، می‌تونین فایل Readme داخل ریپازیتوری گیت هابم رو بخونین.

من برای این پروژه از ویدئوهای مربوط به ماشین لرنینگ و پردازش زبان طبیعی یوتیوب، سایت w3schools و مدل‌های زبانی کمک گرفتم. شروع پروژه از روز جمعه 30 آبان ساعت 3 و نیم آغاز شد و الان که دارم گزارشم رو تکمیل می‌کنم، ساعت 10 و نیم شب روز 3 آذر هست. این فقط فرآیند برنامه نویسی بود و من از قبل این زمان هم برای درک دقیق پروژه، چرایی انجام اون و مواردی که لازم بود آموزش ببینم هم کلی تحقیق کردم. واقعیتش درک این پروژه و یاد گرفتن مسیر انجام اون، خیلی لذت بخش بود، چون حالا دیگه می‌دونم مبنای ماشین لرنینگ چیه و اصلاً چجوری می‌شه یه ماشین رو آموزش داد.

بدون شک نتیجه نهایی قابل گسترش هست و من تمام سعیم رو به کار می‌گیرم که به محض یادگیری بیشتر، این پروژه رو هم توسعه بدم.

