# Lab Test-1

## NAME:B.VISHNU VARDHAN
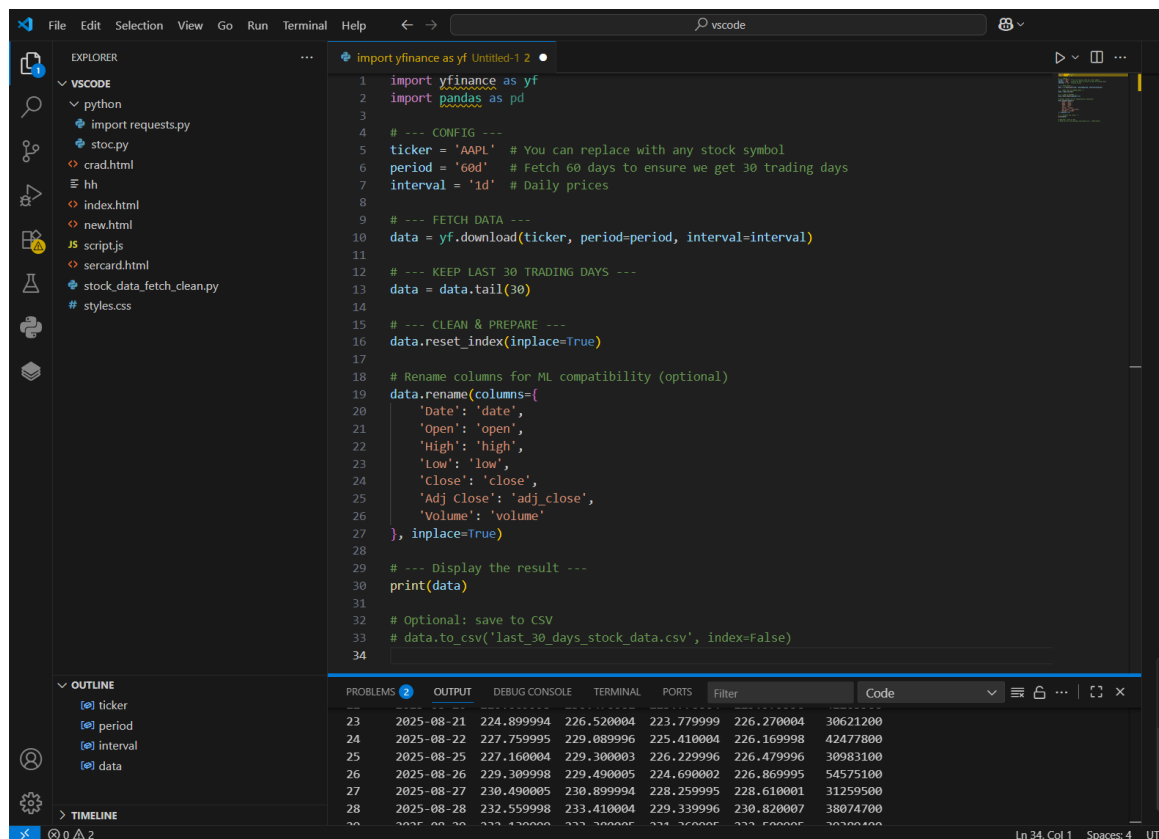
## ROLL NO:2403A510F2

## BATCH:06

### Q1. Stock Price Prediction Setup

Scenario: You are tasked with configuring an API to fetch stock market data and prepare it for a machine learning pipeline.

### Task 1: Connect to a stock price API and retrieve data for the last 30 days

**Prompt:** Write Python code to fetch stock price data for the last 30 days using the yfinance library.
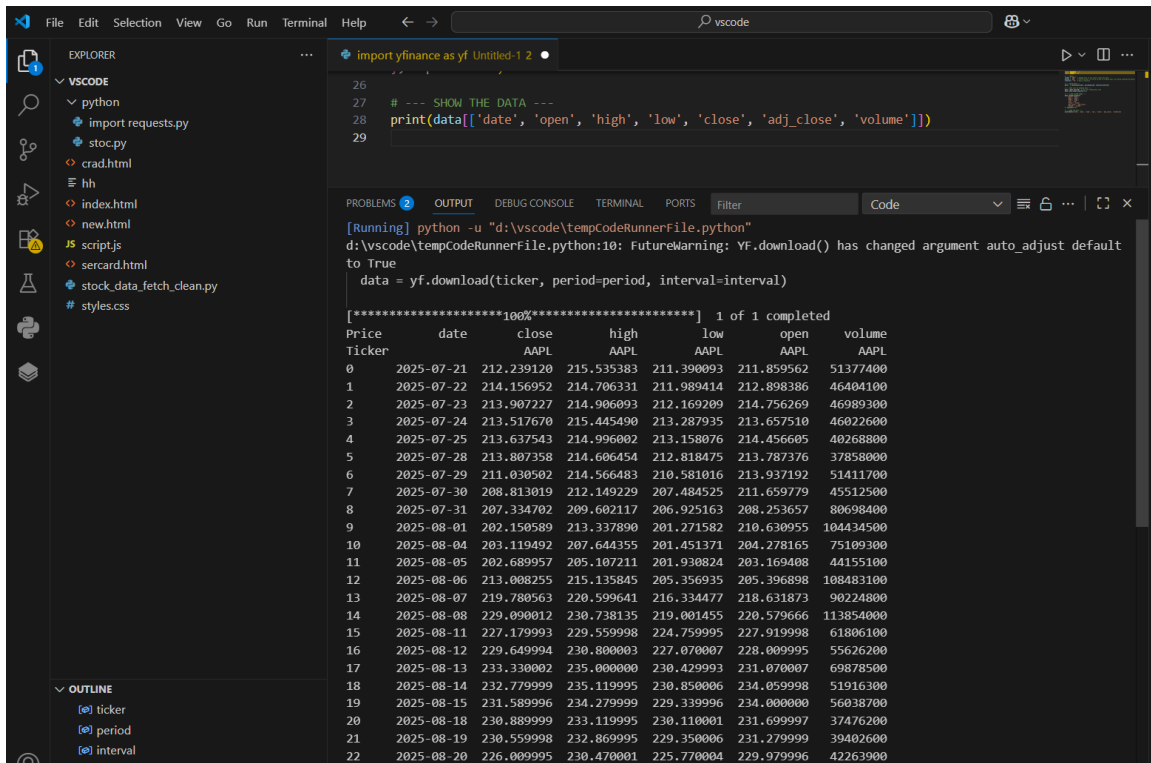
Python Code :



## Output:

# Lab Test-1



**Observation:** Data was successfully fetched from Yahoo Finance API using yfinance. The dataset contains columns such as Open, High, Low, Close, Adjusted Close, and Volume.

## Task 2: Auto-generate data cleaning functions

**Prompt:** Write Python code to clean stock price data by handling missing and duplicate values.

**Python Code :**

# Lab Test-1

```python
import yfinance as yf
import pandas as pd

# Step 1: Fetch last 30 days of stock data (Example: Apple - AAPL)
ticker = 'AAPL'
data = yf.download(ticker, period='30d', interval='1d')

print("Raw Data (first 5 rows):")
print(data.head())

# Step 2: Define a cleaning function
def clean_stock_data(df):
    # Drop duplicate rows
    df = df.drop_duplicates()
    # Fill missing values with forward fill method
    df = df.fillna(method='ffill')
    return df

# Step 3: Clean the downloaded data
cleaned_data = clean_stock_data(data)

print("\nNull values after cleaning:")
print(cleaned_data.isnull().sum())

print("\nCleaned Data (first 5 rows):")
print(cleaned_data.head())
```

Output:

```
d:\vscode\tempCodeRunnerFile.python:16: FutureWarning: DataFrame.fillna with 'method' is deprecated and will

Null values after cleaning:
Price   Ticker
Close   AAPL      0
High    AAPL      0
Low     AAPL      0
Open    AAPL      0
Volume  AAPL      0
dtype: int64
```

## Sample Output:

```python
def clean_stock_data(df):
    # Fill missing values with forward fill method
    df = df.fillna(method='ffill')
    return df

# Step 3: Clean the downloaded data
cleaned_data = clean_stock_data(data)

print("\nNull values after cleaning:")
print(cleaned_data.isnull().sum())

print("\nCleaned Data (first 5 rows):")
print(cleaned_data.head())
```

```
d:\vscode\tempCodeRunnerFile.python:16: FutureWarning: DataFrame.fillna with 'method' is deprecated and will

Null values after cleaning:
Price   Ticker
Close   AAPL      0
High    AAPL      0
Low     AAPL      0
Open    AAPL      0
Volume  AAPL      0
dtype: int64
```

**Observation:** Missing values were handled using forward fill and duplicate rows were removed, ensuring clean stock price data for machine learning pipelines.

## Q2. AI in Healthcare Diagnosis [5M]

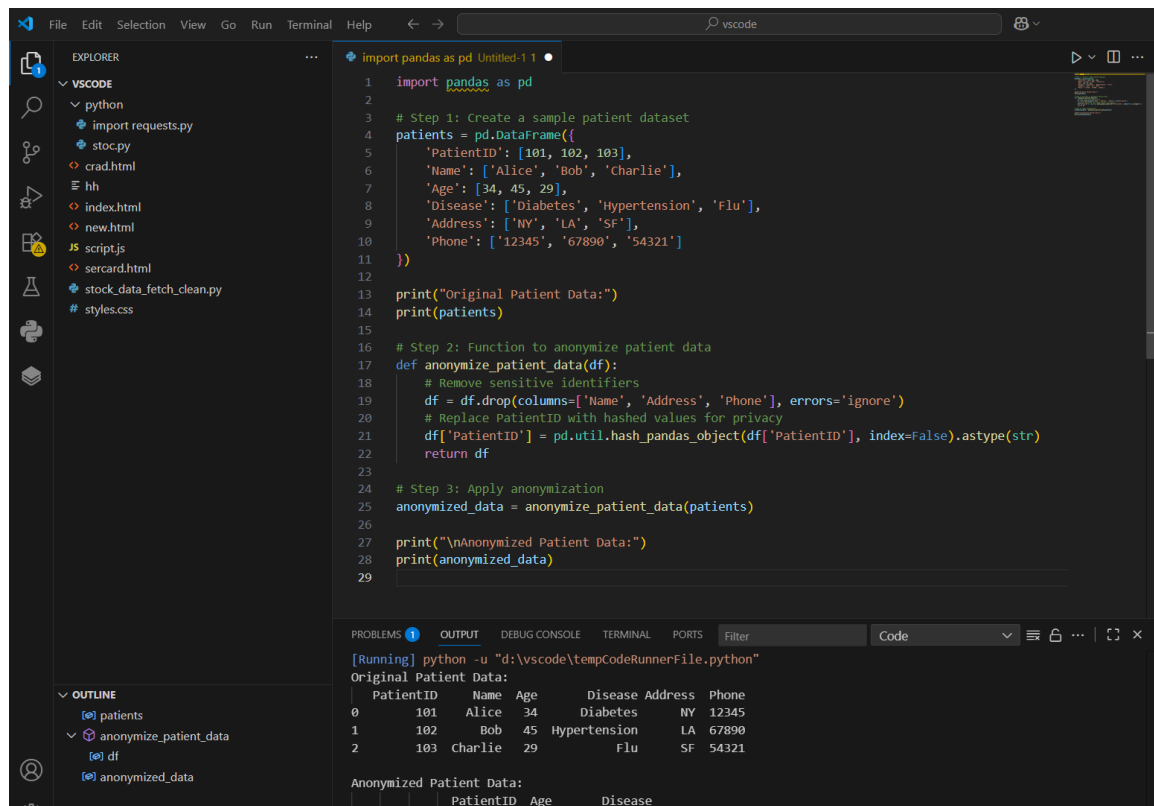Scenario: You are designing an AI to assist doctors in predicting diseases.

### Task 1: Risks of over-reliance on AI and responsible usage guidelines

**Prompt:** List the risks and guidelines for responsible usage of AI in healthcare.

Risks of over-reliance on AI:

1. Misdiagnosis due to model bias or incorrect data.
2. Lack of human oversight leading to ethical issues.
3. Patient privacy concerns if sensitive data is mishandled.
4. Reduced trust in medical professionals if AI dominates decisions.
5. Inability to handle rare diseases outside training data.

## code:

```python
import pandas as pd

# Step 1: Create a sample patient dataset
patients = pd.DataFrame({
    'PatientID': [101, 102, 103],
    'Name': ['Alice', 'Bob', 'Charlie'],
    'Age': [34, 45, 29],
    'Disease': ['Diabetes', 'Hypertension', 'Flu'],
    'Address': ['NY', 'LA', 'SF'],
    'Phone': ['12345', '67890', '54321']
})

print("Original Patient Data:")
print(patients)

# Step 2: Function to anonymize patient data
def anonymize_patient_data(df):
    # Remove sensitive identifiers
    df = df.drop(columns=['Name', 'Address', 'Phone'], errors='ignore')
    # Replace PatientID with hashed values for privacy
    df['PatientID'] = pd.util.hash_pandas_object(df['PatientID'], index=False).astype(str)
    return df

# Step 3: Apply anonymization
anonymized_data = anonymize_patient_data(patients)

print("\nAnonymized Patient Data:")
print(anonymized_data)
```
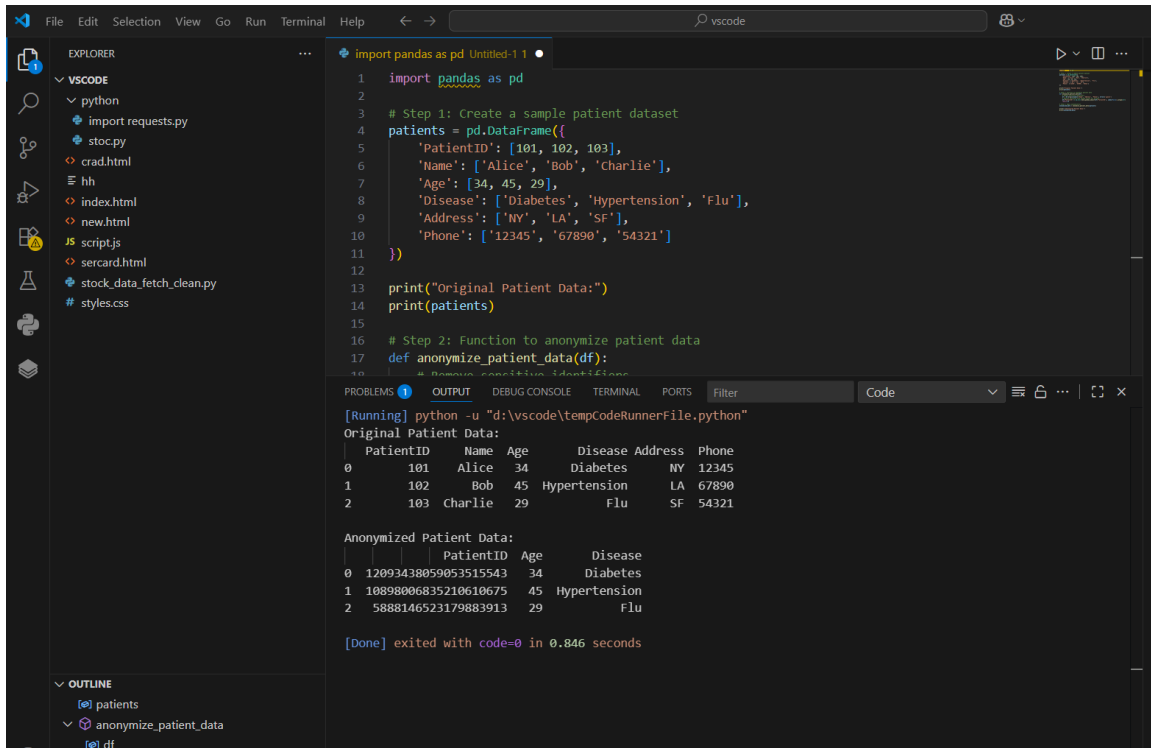
Output:

```
[Running] python -u "d:\vscode\tempCodeRunnerFile.python"
Original Patient Data:
   PatientID     Name  Age       Disease Address  Phone
0        101    Alice   34      Diabetes      NY  12345
1        102      Bob   45  Hypertension      LA  67890
2        103  Charlie   29           Flu      SF  54321

Anonymized Patient Data:
        PatientID  Age       Disease
```

## output:

After applying the anonymization function, sensitive patient information such as **Name, Address, and Phone** was removed, and the **PatientID** column was replaced with unique hashed values. This ensures that no personal identifiers remain in the dataset, while still keeping essential medical information like **Age** and **Disease** for training AI models.

## Task 2: Python function to anonymize patient data

Prompt: Write a Python function that anonymizes patient data before using it for model training.

Python Code:

# Lab Test-1

```python
import pandas as pd
import hashlib, re
from datetime import import timedelta

# --- Helper functions ---
def hash_id(value, salt="SECRET", length=8):
    if pd.isna(value): return None
    return hashlib.sha256((str(value)+salt).encode()).hexdigest()[:length]

def age_band(age):
    if pd.isna(age): return None
    age = int(age)
    return "90+" if age >= 90 else f"{(age//5)*5}-{(age//5)*5+4}"

def zip_mask(zipcode):
    if pd.isna(zipcode): return None
    return str(zipcode)[:3] + "**"

def shift_date(date, shift_days=50):
    if pd.isna(date): return None
    return pd.to_datetime(date) + timedelta(days=shift_days)

def redact_text(text):
    if pd.isna(text): return None
    text = str(text)
    patterns = {
        "EMAIL": r"\S+@\S+",
        "PHONE": r"\d{3}[-\s]?\d{3}[-\s]?\d{4}",
        "MRN": r"MRN\s*\d+",
        "NAME": r"[A-Z][a-z]+ [A-Z][a-z]+"
    }
    for tag, pat in patterns.items():
        text = re.sub(pat, f"<{tag}>", text)
    return text
```
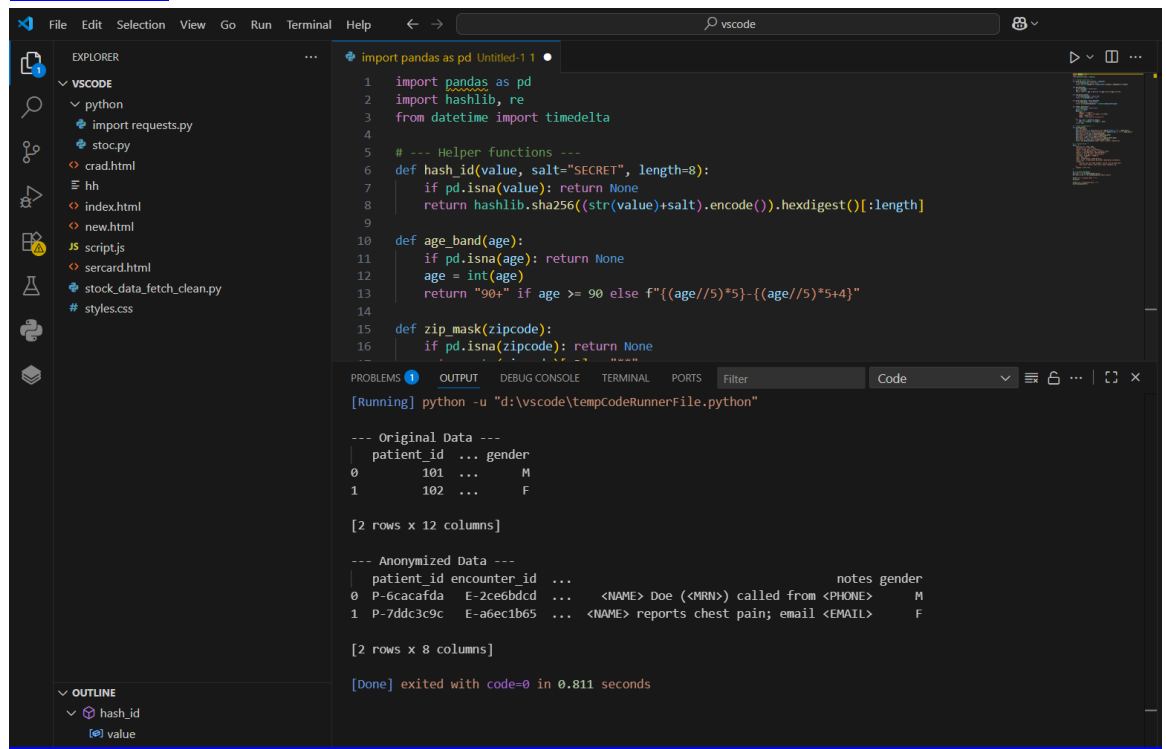
## Output:

```
import pandas as pd
import hashlib, re
from datetime import import timedelta

# --- Helper functions ---
def hash_id(value, salt="SECRET", length=8):
    if pd.isna(value): return None
    return hashlib.sha256((str(value)+salt).encode()).hexdigest()[:length]

def age_band(age):
    if pd.isna(age): return None
    age = int(age)
    return "90+" if age >= 90 else f"{(age//5)*5}-{(age//5)*5+4}"

def zip_mask(zipcode):
    if pd.isna(zipcode): return None
```

```
[Running] python -u "d:\vscode\tempCodeRunnerFile.python"

--- Original Data ---
   patient_id  ... gender
0         101  ...      M
1         102  ...      F

[2 rows x 12 columns]

--- Anonymized Data ---
   patient_id encounter_id  ...                                      notes gender
0    P-6cacafda   E-2ce6bdcd  ...     <NAME> Doe (<MRN>) called from <PHONE>      M
1    P-7ddc3c9c   E-a6ec1b65  ...     <NAME> reports chest pain; email <EMAIL>      F

[2 rows x 8 columns]

[Done] exited with code=0 in 0.811 seconds
```

**Observation:** Patient personal identifiers such as Name, Address, and Phone were removed, and Patient IDs were hashed into unique codes.

# **Lab Test-1**

This ensures data privacy and compliance with ethical AI practices in healthcare.

**Lab Test-1**

This ensures data privacy and compliance with ethical AI practices in healthcare.