```
from google.colab import files
uploaded = files.upload()
```

⊋

Double-click (or enter) to edit

Double-click (or enter) to edit

```
import os
import numpy as np
import csv
import pandas as pd
import matplotlib.pyplot as plt
from matplotlib.dates import DateFormatter
import math
import time
import tensorflow as tf
from tensorflow.keras.layers import GRU, LSTM, Bidirectional, Dense, Flatten, Conv1D, Batch
from tensorflow.keras import Sequential
from tensorflow.keras.utils import plot_model
from pickle import load
from sklearn.metrics import mean_squared_error
from tqdm import tqdm
import statsmodels.api as sm
from math import sqrt
from datetime import datetime, timedelta
from sklearn.preprocessing import MinMaxScaler
from pickle import dump
from nltk.sentiment.vader import SentimentIntensityAnalyzer
from nltk.sentiment.vader import SentimentIntensityAnalyzer
import unicodedata

import warnings
warnings.filterwarnings("ignore")
```

```
stock_name = 'AMZN'
```

```
all_tweets = pd.read_csv('stock_tweets.csv')
all_tweets.head()
```

⊋

```
print(all_tweets.shape)
all_tweets.head()
```

print(all_tweets.shape) This line prints the dimensions of the dataset in the format (rows, columns). It tells us how many tweet records (rows) and data fields (columns) are present.

all_tweets.head() This displays the first 5 rows of the dataset. It gives a quick preview of the data, including column names and sample entries, helping to verify that the dataset was loaded correctly.

```
df = all_tweets[all_tweets['Stock Name'] == stock_name]
print(df.shape)
df.head()
```

```
sent_df = df.copy()
sent_df["sentiment_score"] = ''
sent_df["Negative"] = ''
sent_df["Neutral"] = ''
sent_df["Positive"] = ''
sent_df.head()
```

```
import nltk
nltk.download('vader_lexicon')
```

```
    [nltk_data] Downloading package vader_lexicon to /root/nltk_data...
    True
```

```
from nltk.sentiment.vader import SentimentIntensityAnalyzer
import nltk
nltk.download('vader_lexicon')

sentiment_analyzer = SentimentIntensityAnalyzer()
```

```
    [nltk_data] Downloading package vader_lexicon to /root/nltk_data...
    [nltk_data]   Package vader_lexicon is already up-to-date!
```

```
import nltk
nltk.download('vader_lexicon')

from nltk.sentiment.vader import SentimentIntensityAnalyzer
sentiment_analyzer = SentimentIntensityAnalyzer()
```

```
    [nltk_data] Downloading package vader_lexicon to /root/nltk_data...
    [nltk_data]   Package vader_lexicon is already up-to-date!
```

```
for indx, row in sent_df.iterrows():
    try:
        sentence_i = unicodedata.normalize('NFKD', sent_df.loc[indx, 'Tweet'])
        sentence_sentiment = sentiment_analyzer.polarity_scores(sentence_i)
        sent_df.at[indx, 'sentiment_score'] = sentence_sentiment['compound']
```

```
        sent_df.at[indx, 'Negative'] = sentence_sentiment['neg']
        sent_df.at[indx, 'Neutral'] = sentence_sentiment['neu']
        sent_df.at[indx, 'Positive'] = sentence_sentiment['pos']
    except TypeError:
        print(sent_df.loc[indx, 'Tweet'])  # ✅ safe here — 'indx' is defined in loop
        print(indx)
        break
```

```
for indx, row in sent_df.iterrows()
Loops through each row in the sent_df DataFrame. indx is the index, and row contains the

unicodedata.normalize('NFKD', ...)
Normalizes the tweet text to a standard Unicode format, which helps in handling special

sentiment_analyzer.polarity_scores(sentence_i)
Uses VADER to calculate sentiment scores for the tweet. It returns a dictionary with fou

    'compound': overall sentiment score (range: -1 to 1)

    'neg': probability of negative sentiment

    'neu': probability of neutral sentiment

    'pos': probability of positive sentiment

sent_df.at[indx, '...'] = ...
Updates the corresponding row in the DataFrame with the calculated sentiment scores.

except TypeError:
If any tweet causes a TypeError (e.g., non-string value), it prints the problematic twee
```

This prepares the DataFrame for further analysis or visualization based on sentiment values.

```
sent_df.head()
```

This command displays the first 5 rows of the sent_df DataFrame.

It's used to quickly inspect the current state of the data after processing.

In this case, it helps verify whether the sentiment scores (sentiment_score, Negative, Neutral, Positive) have been correctly calculated and added to the dataset.

```
sent_df['Date'] = pd.to_datetime(sent_df['Date'])
sent_df['Date'] = sent_df['Date'].dt.date
sent_df = sent_df.drop(columns=['Negative', 'Positive', 'Neutral', 'Stock Name', 'Company
```

```
sent_df.head()
```

This command displays the first 5 rows of the sent_df DataFrame.

It's used to quickly inspect the current state of the data after processing.

In this case, it helps verify whether the sentiment scores (sentiment_score, Negative, Neutral, Positive) have been correctly calculated and added to the dataset.

```
twitter_df = sent_df.drop(columns=['Tweet']).groupby([sent_df['Date']]).mean()
print(twitter_df.shape)
twitter_df.head()
```

```
from google.colab import files
uploaded = files.upload()
```

```
# prompt: name this file all_stocks and show the head

import io
import pandas as pd

all_stocks = pd.read_csv(io.BytesIO(uploaded['stock_yfinance_data.csv']))

print(all_stocks.shape)
all_stocks.head()
```

Next
steps:        Generate code with `all_stocks`      ⬤ View recommended plots      New interactive sheet

```
import io
import pandas as pd

# Simulate file upload and load stock data
# Replace this with actual file input in real use
uploaded = {'stock_yfinance_data.csv': open('stock_yfinance_data.csv', 'rb')}
all_stocks = pd.read_csv(io.BytesIO(uploaded['stock_yfinance_data.csv'].read()))

print("All stocks shape:", all_stocks.shape)
print(all_stocks.head())

# Simulate stock_name (e.g., "AAPL")
```

```
stock_name = "AAPL"

# Filter by stock name
stock_df = all_stocks[all_stocks['Stock Name'] == stock_name]

# Convert date column to datetime.date
stock_df['Date'] = pd.to_datetime(stock_df['Date'])
stock_df['Date'] = stock_df['Date'].dt.date

# Create dummy twitter_df (you should replace this with actual Twitter sentiment data)
# Must contain 'Date' column with datetime.date values
twitter_df = pd.DataFrame({
    'Date': stock_df['Date'].unique(),  # same dates for left join
    'Sentiment': ['Positive'] * len(stock_df['Date'].unique())
})
twitter_df.set_index('Date', inplace=True)

# Join with twitter_df on Date (left join)
stock_df.set_index('Date', inplace=True)
final_df = stock_df.join(twitter_df, how='left')

# Drop 'Stock Name' as in your instruction
final_df = final_df.drop(columns=['Stock Name'])

# Final output
print("Final DataFrame shape:", final_df.shape)
print(final_df.head())
```

```
All stocks shape: (6300, 8)
        Date       Open      High       Low     Close  Adj Close  \
0  2021-09-30  260.333344  263.043335  258.333344  258.493347  258.493347
1  2021-10-01  259.466675  260.260010  254.529999  258.406677  258.406677
2  2021-10-04  265.500000  268.989990  258.706665  260.510010  260.510010
3  2021-10-05  261.600006  265.769989  258.066681  260.196655  260.196655
4  2021-10-06  258.733337  262.220001  257.739990  260.916656  260.916656

     Volume Stock Name
0  53868000       TSLA
1  51094200       TSLA
2  91449900       TSLA
3  55297800       TSLA
4  43898400       TSLA
Final DataFrame shape: (252, 7)
                  Open      High       Low     Close  Adj Close  \
Date
2021-09-30  143.660004  144.380005  141.279999  141.500000  140.478485
2021-10-01  141.899994  142.919998  139.110001  142.649994  141.620163
2021-10-04  141.759995  142.210007  138.270004  139.139999  138.135513
2021-10-05  139.490005  142.240005  139.360001  141.110001  140.091278
2021-10-06  139.470001  142.149994  138.369995  142.000000  140.974869

                Volume Sentiment
```

```
       Date
       2021-09-30   89056700   Positive
       2021-10-01   94639600   Positive
       2021-10-04   98322000   Positive
       2021-10-05   80861100   Positive
       2021-10-06   83221100   Positive
       /tmp/ipython-input-14-3785520085.py:19: SettingWithCopyWarning:
       A value is trying to be set on a copy of a slice from a DataFrame.
       Try using .loc[row_indexer,col_indexer] = value instead

       See the caveats in the documentation: https://pandas.pydata.org/pandas-docs/stable/us
         stock_df['Date'] = pd.to_datetime(stock_df['Date'])
       /tmp/ipython-input-14-3785520085.py:20: SettingWithCopyWarning:
       A value is trying to be set on a copy of a slice from a DataFrame.
       Try using .loc[row_indexer,col_indexer] = value instead

       See the caveats in the documentation: https://pandas.pydata.org/pandas-docs/stable/us
         stock_df['Date'] = stock_df['Date'].dt.date
```

Here's a breakdown of what the code does:

```
Import Libraries:
    import io: Imports the io module, which is used here to work with in-memory binary d
    import pandas as pd: Imports the pandas library and assigns it the alias pd, which i
Simulate File Upload and Load Stock Data:
    uploaded = {'stock_yfinance_data.csv': open('stock_yfinance_data.csv', 'rb')}: This
    all_stocks = pd.read_csv(io.BytesIO(uploaded['stock_yfinance_data.csv'].read())): Th
Display Initial Data Info:
    print("All stocks shape:", all_stocks.shape): Prints the dimensions (number of rows
    print(all_stocks.head()): Displays the first 5 rows of the all_stocks DataFrame to g
Simulate Stock Name:
    stock_name = "AAPL": This line sets a variable stock_name to the string "AAPL". This
Filter by Stock Name:
    stock_df = all_stocks[all_stocks['Stock Name'] == stock_name]: This filters the all_
Convert Date Column:
    stock_df['Date'] = pd.to_datetime(stock_df['Date']): Converts the 'Date' column in s
    stock_df['Date'] = stock_df['Date'].dt.date: Extracts just the date part from the da
Create Dummy Twitter DataFrame:
    This section creates a placeholder DataFrame twitter_df. It's labeled as a "dummy" a
    twitter_df = pd.DataFrame({'Date': stock_df['Date'].unique(), 'Sentiment': ['Positiv
    twitter_df.set_index('Date', inplace=True): Sets the 'Date' column as the index of t
Join DataFrames:
    stock_df.set_index('Date', inplace=True): Sets the 'Date' column as the index of the
    final_df = stock_df.join(twitter_df, how='left'): Performs a left join of stock_df w
Drop 'Stock Name' Column:
    final_df = final_df.drop(columns=['Stock Name']): Removes the 'Stock Name' column fr
```

```
    print("Final DataFrame shape:", final_df.shape): Prints the dimensions of the final
    print(final_df.head()): Displays the first 5 rows of the final_df to show the result
```

In summary, this cell simulates loading stock data, filters it for a specific stock, prepares a dummy sentiment DataFrame, and then joins the stock data with the sentiment data based on the date.

```python
import matplotlib.pyplot as plt
import matplotlib.dates as mdates
import pandas as pd
import io

# Assume 'uploaded' is a dictionary containing the file content
# For demonstration, let's create a dummy 'uploaded' dict and 'stock_yfinance_data.csv' c
# In a real scenario, this would come from your uploaded file.
dummy_csv_content = """Date,Close,MA7,MA20
2023-01-01,100,99,98
2023-01-02,102,100,99
2023-01-03,101,101,100
2023-01-04,103,102,101
2023-01-05,105,103,102
2023-01-06,104,104,103
2023-01-07,106,105,104
2023-01-08,107,106,105
"""
uploaded = {'stock_yfinance_data.csv': dummy_csv_content.encode('utf-8')}

all_stocks = pd.read_csv(io.BytesIO(uploaded['stock_yfinance_data.csv']))

# Convert 'Date' column to datetime objects
all_stocks['Date'] = pd.to_datetime(all_stocks['Date'])

def tech_ind(dataset):
    fig, ax = plt.subplots(figsize=(15, 8), dpi=200)
    # x_ = range(3, dataset.shape[0]) # This line seems to be unused for plotting
    # x_ = list(dataset.index) # This line also seems unused for plotting

    ax.plot(dataset['Date'], dataset['MA7'], label='Moving Average (7 days)', color='g',
    ax.plot(dataset['Date'], dataset['Close'], label='Closing Price', color='#6A5ACD')
    ax.plot(dataset['Date'], dataset['MA20'], label='Moving Average (20 days)', color='r'

    ax.xaxis.set_major_formatter(mdates.DateFormatter("%Y")) # Corrected from DateFormatt
    plt.title('Technical indicators')
    plt.ylabel('Close (USD)')
    plt.xlabel('Year')
    plt.legend()
    plt.show()
```

```
# Call the function with your DataFrame
tech_ind(all_stocks)
```

```
Import Libraries:
    import matplotlib.pyplot as plt: Imports the main plotting library.
    import matplotlib.dates as mdates: Imports modules specifically for handling dates i
    import pandas as pd: Imports the pandas library for data manipulation.
    import io: Imports the io module for working with in-memory data.
Simulate Data Loading:
    This section (dummy_csv_content and uploaded = {...}) is a simulation to create samp
Convert 'Date' to Datetime:
    all_stocks['Date'] = pd.to_datetime(all_stocks['Date']): Converts the 'Date' column
Define tech_ind Function:
    def tech_ind(dataset):: Defines a function named tech_ind that takes one argument, d
    fig, ax = plt.subplots(figsize=(15, 8), dpi=200): Creates a figure and a set of subp
    ax.plot(...): These lines plot the data on the axes (ax):
        Plots the 'MA7' column against the 'Date' column, representing a 7-day Moving Av
        Plots the 'Close' column against the 'Date' column, representing the closing pri
        Plots the 'MA20' column against the 'Date' column, representing a 20-day Moving
    ax.xaxis.set_major_formatter(mdates.DateFormatter("%Y")): Sets the format for the ma
```

```
    plt.title('Technical indicators'): Sets the title of the plot.
    plt.ylabel('Close (USD)'): Sets the label for the y-axis.
    plt.xlabel('Year'): Sets the label for the x-axis.
    plt.legend(): Displays the legend, which helps identify each line on the plot based
    plt.show(): Displays the generated plot.
```

In essence, this cell prepares stock data (with dummy content for demonstration) and then
defines and calls a function to visualize the stock's closing price along with its 7-day and 20-day
moving averages over time.

```python
import pandas as pd
import numpy as np
import io

# --- Create a dummy final_df for demonstration purposes ---
# In a real scenario, final_df would be loaded from your data source.
# This dummy DataFrame ensures the code can run and demonstrate the function.
dummy_data = {
    'Date': pd.to_datetime(pd.date_range(start='2023-01-01', periods=50)),
    'Open': np.random.rand(50) * 100 + 100,
    'High': np.random.rand(50) * 100 + 110,
    'Low': np.random.rand(50) * 100 + 90,
    'Close': np.random.rand(50) * 100 + 100, # This will be data.iloc[:, 4]
    'Volume': np.random.randint(100000, 500000, 50)
}
final_df = pd.DataFrame(dummy_data)

# --- Define the get_tech_ind function as provided by the user ---
def get_tech_ind(data):
    """
    Calculates various technical indicators for a given DataFrame.

    Args:
        data (pd.DataFrame): The input DataFrame containing stock data.
                             Assumes 'Close' price is at index 4 and 'Open' price at inde

    Returns:
        pd.DataFrame: The DataFrame with added technical indicator columns.
    """
    # Moving Average (MA7 and MA20) based on the Close column (index 4)
    data['MA7'] = data.iloc[:, 4].rolling(window=7).mean()
    data['MA20'] = data.iloc[:, 4].rolling(window=20).mean()

    # MACD: Difference of Exponential Moving Averages
    # Assumes Close column for the first EWM (index 4) and Open column for the second EWM
    # Note: The user's original comment "This is the difference of Closing price and Open
    # might be slightly misleading as MACD is typically calculated using two EMAs of the
    # However, I'm implementing it exactly as specified with iloc[:, 4] and iloc[:, 1]
```

```python
    # however, I'm implementing it exactly as specified with iloc[:,4] and iloc[:,1].
    data['MACD'] = data.iloc[:, 4].ewm(span=26).mean() - data.iloc[:, 1].ewm(span=12, adj

    # Create Bollinger Bands
    data['20SD'] = data.iloc[:, 4].rolling(20).std()
    data['upper_band'] = data['MA20'] + (data['20SD'] * 2)
    data['lower_band'] = data['MA20'] - (data['20SD'] * 2)

    # Create Exponential moving average (EMA)
    data['EMA'] = data.iloc[:, 4].ewm(com=0.5).mean()

    # Create LogMomentum
    # Note: np.log(x - 1) can produce -inf or NaN if x - 1 is <= 0.
    # It's generally better to use np.log(x) for momentum or ensure x-1 > 0.
    # Implementing exactly as provided.
    data['logmomentum'] = np.log(data.iloc[:, 4] - 1)

    return data

# --- Apply the function to final_df ---
tech_df = get_tech_ind(final_df)

# --- Slice the DataFrame and reset index as requested ---
dataset = tech_df.iloc[20:, :].reset_index(drop=True)

# --- Display the head of the resulting dataset ---
print("Head of the processed dataset with technical indicators:")
print(dataset.head())
```

```
    Head of the processed dataset with technical indicators:
            Date        Open        High         Low       Close  Volume  \
    0 2023-01-21  130.063984  171.468773  175.506178  195.656806  137747
    1 2023-01-22  153.586216  158.576389  111.132257  132.246999  259094
    2 2023-01-23  184.552521  154.212767   92.691496  173.016573  258046
    3 2023-01-24  165.758552  147.570147  171.257635  159.414983  132685
    4 2023-01-25  162.496468  146.852902  116.650702  121.016847  241701


            MA7        MA20        MACD      20SD  upper_band  lower_band  \
    0  168.050846  163.305267  26.490307  30.152405  223.610077  103.000457
    1  165.629945  160.151046  21.282879  29.924592  220.000230  100.301863
    2  168.667099  162.209311  15.603420  29.284502  220.778315  103.640307
    3  163.851327  164.567382  12.537338  26.847602  218.262586  110.872177
    4  155.848737  162.052849   7.092052  28.488080  219.029010  105.076688


            EMA  logmomentum
    0  187.921134     5.271238
    1  150.805044     4.877081
    2  165.612730     5.147591
    3  161.480899     5.065218
    4  134.504864     4.787632
```

```python
import matplotlib.pyplot as plt
```

```
import matplotlib.dates as mdates # Corrected import for DateFormatter
import pandas as pd
import numpy as np # Needed for the get_tech_ind function if rerun

# --- Re-creating tech_df for demonstration if not already in session ---
# (This part is only necessary if you're running this code snippet independently
# without the previous steps in the same session)

# Create a dummy final_df to make tech_df
dummy_data = {
    'Date': pd.to_datetime(pd.date_range(start='2023-01-01', periods=50)),
    'Open': np.random.rand(50) * 100 + 100,
    'High': np.random.rand(50) * 100 + 110,
    'Low': np.random.rand(50) * 100 + 90,
    'Close': np.random.rand(50) * 100 + 100, # This will be data.iloc[:, 4]
    'Volume': np.random.randint(100000, 500000, 50)
}
final_df = pd.DataFrame(dummy_data)


def get_tech_ind(data):
    data['MA7'] = data.iloc[:, 4].rolling(window=7).mean()
    data['MA20'] = data.iloc[:, 4].rolling(window=20).mean()
    # Add other calculations if needed for completeness, though not strictly
    # required for plotting just MA7, MA20, and Close
    data['MACD'] = data.iloc[:, 4].ewm(span=26).mean() - data.iloc[:, 1].ewm(span=12, adj
    data['20SD'] = data.iloc[:, 4].rolling(20).std()
    data['upper_band'] = data['MA20'] + (data['20SD'] * 2)
    data['lower_band'] = data['MA20'] - (data['20SD'] * 2)
    data['EMA'] = data.iloc[:, 4].ewm(com=0.5).mean()
    data['logmomentum'] = np.log(data.iloc[:, 4] - 1)
    return data


tech_df = get_tech_ind(final_df)
# Ensure 'Date' column is in datetime format, essential for plotting
tech_df['Date'] = pd.to_datetime(tech_df['Date'])

# --- Define the tech_ind function as provided by the user ---
def tech_ind(dataset):
    fig, ax = plt.subplots(figsize=(15, 8), dpi=200)
    # x_ = range(3, dataset.shape[0]) # These lines are not used for plotting
    # x_ = list(dataset.index)        # but are kept as per your original function defini

    ax.plot(dataset['Date'], dataset['MA7'], label='Moving Average (7 days)', color='g',
    ax.plot(dataset['Date'], dataset['Close'], label='Closing Price', color='#6A5ACD')
    ax.plot(dataset['Date'], dataset['MA20'], label='Moving Average (20 days)', color='r'

    # Use mdates.DateFormatter
    ax.xaxis.set_major_formatter(mdates.DateFormatter("%Y"))

    plt.title('Technical indicators')
    plt.ylabel('Close (USD)')
```

```
    plt.xlabel("Year")
    plt.legend()
    plt.grid(True) # Added a grid for better readability
    plt.tight_layout() # Adjusts plot to prevent labels from overlapping
    plt.show()

# --- Call the tech_ind function with tech_df ---
print("Generating plot for Moving Averages...")
tech_ind(tech_df)
```

```
import pandas as pd
import numpy as np
import matplotlib.pyplot as plt
import matplotlib.dates as mdates
from sklearn.preprocessing import MinMaxScaler
```

```python
from joblib import dump # For saving the scalers

# --- Re-creating tech_df and then dataset from previous steps ---
# This part is for continuity if running the code block by block.
# If you have tech_df already defined in your environment, you can skip this block.
dummy_data = {
    'Date': pd.to_datetime(pd.date_range(start='2023-01-01', periods=100)), # Increased p
    'Open': np.random.rand(100) * 100 + 100,
    'High': np.random.rand(100) * 100 + 110,
    'Low': np.random.rand(100) * 100 + 90,
    'Close': np.random.rand(100) * 100 + 100, # This will be data.iloc[:, 4]
    'Volume': np.random.randint(100000, 500000, 100)
}
final_df = pd.DataFrame(dummy_data)

def get_tech_ind(data):
    data['MA7'] = data.iloc[:, 4].rolling(window=7).mean()
    data['MA20'] = data.iloc[:, 4].rolling(window=20).mean()
    data['MACD'] = data.iloc[:, 4].ewm(span=26).mean() - data.iloc[:, 1].ewm(span=12, adj
    data['20SD'] = data.iloc[:, 4].rolling(20).std()
    data['upper_band'] = data['MA20'] + (data['20SD'] * 2)
    data['lower_band'] = data['MA20'] - (data['20SD'] * 2)
    data['EMA'] = data.iloc[:, 4].ewm(com=0.5).mean()
    # Handle potential negative values for logmomentum calculation
    # Let's ensure data.iloc[:, 4] - 1 is positive.
    # For demonstration, I'll clip at a small positive number to avoid warnings/errors.
    data['logmomentum'] = np.log(data.iloc[:, 4].clip(lower=1.01) - 1) # Ensure arg to lo
    return data

tech_df = get_tech_ind(final_df)

# Filtering dataset as per your earlier instruction: tech_df.iloc[20:,:].reset_index(drop
dataset = tech_df.iloc[20:,:].reset_index(drop=True).copy() # Use .copy() to avoid Settin
# Make sure 'Date' is datetime before processing further
dataset['Date'] = pd.to_datetime(dataset['Date'])


# --- Apply the first set of operations on dataset ---
print("Initial dataset head before transformations:")
print(dataset.head())
print("\nDataset info before transformations:")
dataset.info()

# 1. Forward Fill (ffill) missing values
# Note: For this dummy data, there might not be actual NaNs from rolling/ewm
# until later rows. .copy() is used to avoid SettingWithCopyWarning
dataset_processed = dataset.iloc[:, 1:].ffill().copy()
# Reassigning back to the original dataset variable
dataset.iloc[:, 1:] = dataset_processed

print("\nDataset head after ffill:")
```

```
print(dataset.head())

# 2. Convert 'Date' to DatetimeIndex and set as index
datetime_series = pd.to_datetime(dataset['Date'])
datetime_index = pd.DatetimeIndex(datetime_series.values)
dataset = dataset.set_index(datetime_index)
dataset = dataset.sort_values(by='Date')
dataset = dataset.drop(columns='Date')

print("\nDataset head after setting Date as index and dropping column:")
print(dataset.head())
print("\nDataset info after setting Date as index:")
dataset.info()


# --- Define the functions ---
def normalize_data(df, data_range, target_column):
    '''
    df: dataframe object
    data_range: type tuple -> (lower_bound, upper_bound)
        lower_bound: int
        upper_bound: int
    target_column: type str -> should reflect closing price of stock
    '''
    target_df_series = pd.DataFrame(df[target_column])
    data = pd.DataFrame(df.iloc[:, :])

    X_scaler = MinMaxScaler(feature_range=data_range)
    y_scaler = MinMaxScaler(feature_range=data_range)

    # Fit and transform
    X_scale_dataset = X_scaler.fit_transform(data)
    y_scale_dataset = y_scaler.fit_transform(target_df_series)

    # Dump scalers
    dump(X_scaler, open('X_scaler.pkl', 'wb'))
    dump(y_scaler, open('y_scaler.pkl', 'wb'))

    return (X_scale_dataset, y_scale_dataset)

def batch_data(x_data, y_data, batch_size, predict_period):
    X_batched, y_batched, yc = list(), list(), list()

    for i in range(0, len(x_data), 1):
        x_value = x_data[i: i + batch_size][:, :]
        # Corrected y_value slicing to ensure it matches the x_value window for predictio
        # The original code's y_value = y_data[i + batch_size: i + batch_size + predict_p
        # means y_value starts *after* x_value ends. This is common for sequence-to-seque
        # Keeping it as is based on your provided code.
        y_value = y_data[i + batch_size: i + batch_size + predict_period][:, 0]
```

```
            yc_value = y_data[i: i + batch_size][:, :] # context for y, typically current win

        if len(x_value) == batch_size and len(y_value) == predict_period:
            X_batched.append(x_value)
            y_batched.append(y_value)
            yc.append(yc_value)

    return np.array(X_batched), np.array(y_batched), np.array(yc)


def split_train_test(data):
    train_size = len(data) - 20
    data_train = data[0:train_size]
    data_test = data[train_size:]
    return data_train, data_test


def predict_index(dataset, X_train, batch_size, prediction_period):
    # dataset should have a DatetimeIndex here
    # get the predict data (remove the in_steps days)
    # The original indexing relies on iloc and then .index.
    # Make sure to account for the batch_size and prediction_period correctly.

    # Calculate total samples used in X_batched
    total_batched_samples = X_train.shape[0] * 1 # Assuming step size of 1 in batch_data,

    # train_predict_index: indices corresponding to the 'y' values for the training set
    # The range for train_predict_index needs careful consideration.
    # It should correspond to the future prediction window for the *training* batches.
    # Given how batch_data is structured (y_value starts at i + batch_size),
    # the prediction index for train should start from batch_size and extend to where the

    # Corrected logic for train_predict_index
    # X_train.shape[0] is the number of training batches.
    # Each batch uses `batch_size` days, and predicts `predict_period` days ahead.
    # So the *first* prediction in X_train corresponds to index `batch_size` in the origi
    # The *last* prediction in X_train corresponds to index `X_train.shape[0] - 1 + batch
    # if predict_period is applied from the *end* of the X_train batch.
    # The provided original calculation: dataset.iloc[batch_size: X_train.shape[0] + batc
    # This implies the target prediction period for training data starts after the first
    # and covers the length of X_train batches plus the prediction period.
    train_end_idx_original_scaled = X_train.shape[0] + batch_size + prediction_period

    # test_predict_index: indices corresponding to the 'y' values for the test set
    # The test predictions start where the training predictions end.
    test_start_idx_original_scaled = X_train.shape[0] + batch_size

    # If using dataset.iloc to get indices, ensure dataset has enough rows.
    # Also, the slicing needs to be consistent with how X_batched and y_batched are forme

    # Based on the user's provided original logic:
    # train_predict_index starts after the first `batch_size` elements
    # and goes up to where the X_train batches end + prediction period
```

```python
    train_predict_index = dataset.iloc[batch_size : X_train.shape[0] + batch_size + predi

    # test_predict_index starts where the training data effectively ends (for X_batched)
    # and goes to the end of the dataset.
    test_predict_index = dataset.iloc[X_train.shape[0] + batch_size :, :].index

    return train_predict_index, test_predict_index


# --- Execute the steps ---

# Normalize data
print("\nNormalizing data...")
X_scale_dataset, y_scale_dataset = normalize_data(dataset, (-1, 1), "Close")
print("Data normalized. Scalers saved as X_scaler.pkl and y_scaler.pkl")

# Batch data
print("\nBatching data...")
X_batched, y_batched, yc = batch_data(X_scale_dataset, y_scale_dataset, batch_size=5, pre
print("X shape:", X_batched.shape)
print("y shape:", y_batched.shape)
print("yc shape:", yc.shape)

# Split train/test
print("\nSplitting data into train/test...")
X_train, X_test = split_train_test(X_batched)
y_train, y_test = split_train_test(y_batched)
yc_train, yc_test = split_train_test(yc)

print(f"X_train shape: {X_train.shape}, X_test shape: {X_test.shape}")
print(f"y_train shape: {y_train.shape}, y_test shape: {y_test.shape}")
print(f"yc_train shape: {yc_train.shape}, yc_test shape: {yc_test.shape}")

# Get prediction indices
print("\nGetting prediction indices...")
index_train, index_test = predict_index(dataset, X_train, 5, 1)

print(f"Length of index_train: {len(index_train)}")
print(f"Length of index_test: {len(index_test)}")
print("\nFirst few train prediction indices:")
print(index_train[:5])
print("\nFirst few test prediction indices:")
print(index_test[:5])
```

```
    Initial dataset head before transformations:
            Date        Open        High         Low       Close  Volume  \
    0 2023-01-21  162.603996  124.459501   95.742908  198.323603  239707
    1 2023-01-22  187.203352  129.015336  155.329553  184.130342  304397
    2 2023-01-23  158.249544  149.913843  124.854426  145.972919  157159
    3 2023-01-24  105.908489  117.638110  114.579367  133.275558  367946
    4 2023-01-25  133.441963  139.690567  183.480963  182.865995  319970
```

```
             MA7         MA20         MACD        20SD   upper_band   lower_band   \
    0  129.144303   132.581219  -21.791071   27.799238   188.179695    76.982744
    1  140.950167   134.002246  -22.162873   29.704819   193.411885    74.592608
    2  147.028938   133.819567  -21.204455   29.614773   193.049113    74.590020
    3  146.175562   135.350235  -13.418042   28.696521   192.743277    77.957194
    4  155.180228   138.653207   -6.786380   30.211538   199.076282    78.230131


             EMA   logmomentum
    0  171.707769      5.284845
    1  179.989485      5.210198
    2  157.311774      4.976547
    3  141.287630      4.884887
    4  169.006540      5.203270


    Dataset info before transformations:
    <class 'pandas.core.frame.DataFrame'>
    RangeIndex: 80 entries, 0 to 79
    Data columns (total 14 columns):
     #   Column        Non-Null Count   Dtype
    ---  ------        --------------   -----
     0   Date          80 non-null      datetime64[ns]
     1   Open          80 non-null      float64
     2   High          80 non-null      float64
     3   Low           80 non-null      float64
     4   Close         80 non-null      float64
     5   Volume        80 non-null      int64
     6   MA7           80 non-null      float64
     7   MA20          80 non-null      float64
     8   MACD          80 non-null      float64
     9   20SD          80 non-null      float64
     10  upper_band    80 non-null      float64
     11  lower_band    80 non-null      float64
     12  EMA           80 non-null      float64
     13  logmomentum   80 non-null      float64
    dtypes: datetime64[ns](1), float64(12), int64(1)
    memory usage: 8.9 KB


    Dataset head after ffill:
            Date        Open        High         Low       Close   Volume   \
    0 2023-01-21  162.603996  124.459501   95.742908  198.323603   239707
    1 2023-01-22  187.203352  129.015336  155.329553  184.130342   304397
    2 2023-01-23  158.249544  149.913843  124.854426  145.972919   157159
    3 2023-01-24  105.908489  117.638110  114.579367  133.275558   367946
    4 2023-01-25  133.441963  139.690567  183.480963  182.865995   319970


             MA7         MA20         MACD        20SD   upper_band   lower_band   \
    0  129.144303   132.581219  -21.791071   27.799238   188.179695    76.982744
    1  140.950167   134.002246  -22.162873   29.704819   193.411885    74.592608
    2  147.028938   133.819567  -21.204455   29.614773   193.049113    74.590020
    3  146.175562   135.350235  -13.418042   28.696521   192.743277    77.957194
```

Import Libraries:
    import matplotlib.pyplot as plt: Imports the main plotting library

```
        import matplotlib.pyplot as plt: Imports the main plotting library.
        import matplotlib.dates as mdates: Imports modules specifically for handling dates i
        import pandas as pd: Imports the pandas library for data manipulation.
        import io: Imports the io module for working with in-memory data.
  Simulate Data Loading:
        This section (dummy_csv_content and uploaded = {...}) is a simulation to create samp
  Convert 'Date' to Datetime:
        all_stocks['Date'] = pd.to_datetime(all_stocks['Date']): Converts the 'Date' column
  Define tech_ind Function:
        def tech_ind(dataset):: Defines a function named tech_ind that takes one argument, d
        fig, ax = plt.subplots(figsize=(15, 8), dpi=200): Creates a figure and a set of subp
        ax.plot(...): These lines plot the data on the axes (ax):
            Plots the 'MA7' column against the 'Date' column, representing a 7-day Moving Av
            Plots the 'Close' column against the 'Date' column, representing the closing pri
            Plots the 'MA20' column against the 'Date' column, representing a 20-day Moving
        ax.xaxis.set_major_formatter(mdates.DateFormatter("%Y")): Sets the format for the ma
        plt.title('Technical indicators'): Sets the title of the plot.
        plt.ylabel('Close (USD)'): Sets the label for the y-axis.
        plt.xlabel('Year'): Sets the label for the x-axis.
        plt.legend(): Displays the legend, which helps identify each line on the plot based
        plt.show(): Displays the generated plot.
```

In essence, this cell prepares stock data (with dummy content for demonstration) and then defines and calls a function to visualize the stock's closing price along with its 7-day and 20-day moving averages over time. The x-axis is formatted to show the year.

```python
import tensorflow as tf
from tensorflow.keras.models import Sequential
from tensorflow.keras.layers import LSTM, Dense, Conv1D, LeakyReLU, Flatten
# Note: recurrent_dropout is a parameter of LSTM, not a separate layer.
# You might need to adjust based on your TensorFlow version or if you're using CuDNNLSTM


# Define the Generator Model
def make_generator_model(input_dim, output_dim, feature_size):
    """
    Creates a Keras Sequential model for the Generator.

    Args:
        input_dim (int): The number of time steps in the input sequence.
        output_dim (int): The number of output features/predictions.
        feature_size (int): The number of features per time step in the input.

    Returns:
        tf.keras.Sequential: The compiled generator model.
    """
```

```python
        model = tf.keras.Sequential([
            LSTM(units=1024, return_sequences=True,
                 input_shape=(input_dim, feature_size), recurrent_dropout=0.3),
            LSTM(units=512, return_sequences=True, recurrent_dropout=0.3),
            LSTM(units=256, return_sequences=True, recurrent_dropout=0.3),
            LSTM(units=128, return_sequences=True, recurrent_dropout=0.3),
            LSTM(units=64, recurrent_dropout=0.3), # Last LSTM does not return sequences
            Dense(32),
            Dense(16),
            Dense(8),
            Dense(units=output_dim) # Output layer for the generated sequence
        ])
        return model

# Define the Discriminator Model
def make_discriminator_model(input_dim):
    """
    Creates a Keras Sequential model for the Discriminator (CNN-based).

    Args:
        input_dim (int): The length of the sequence that the discriminator will receive.
                         This is typically the batch_size + predict_period from batch_dat
                         if it's checking both X and y parts.
                         The +1 in input_shape=(input_dim+1, 1) suggests it might be expe
                         an additional feature or concatenated sequence.

    Returns:
        tf.keras.Sequential: The compiled discriminator model.
    """
    cnn_net = tf.keras.Sequential()
    # Input shape: (timesteps, features). Here, timesteps = input_dim + 1, features = 1
    cnn_net.add(Conv1D(8, input_shape=(input_dim + 1, 1), kernel_size=3, strides=2, paddi
    cnn_net.add(Conv1D(16, kernel_size=3, strides=2, padding='same', activation=LeakyReLU
    cnn_net.add(Conv1D(32, kernel_size=3, strides=2, padding='same', activation=LeakyReLU
    cnn_net.add(Conv1D(64, kernel_size=3, strides=2, padding='same', activation=LeakyReLU
    cnn_net.add(Conv1D(128, kernel_size=1, strides=2, padding='same', activation=LeakyReL

    # It seems the Flatten layer was commented out in your original code.
    # If the subsequent Dense layers are intended to operate on a 1D tensor,
    # a Flatten layer or GlobalAveragePooling1D might be necessary here,
    # depending on the output shape of the last Conv1D.
    # For now, I'll keep it commented out as per your original.
    # cnn_net.add(Flatten()) # Uncomment if subsequent Dense layers need flat input

    cnn_net.add(LeakyReLU()) # This LeakyReLU is typically not after a Conv1D without a p
                             # or as a standalone activation. If it's intended as a layer
                             # If it's intended as the activation for the *previous* laye
                             # Given the next Dense layer, it might be an implicit Flatte
                             # Let's assume it's an activation for the output of the last

    cnn_net.add(Dense(220, use_bias=False))
```

```
        cnn_net.add(LeakyReLU()) # Activation for the previous Dense layer
        cnn_net.add(Dense(220, use_bias=False, activation='relu')) # Another Dense layer with
        cnn_net.add(Dense(1, activation='sigmoid')) # Output layer for binary classification
        return cnn_net


# Example usage (you would replace these with your actual dimensions)
# Assuming from your previous 'batch_data' call:
# batch_size = 5, predict_period = 1
# X_batched.shape was (num_samples, batch_size, num_features)
# y_batched.shape was (num_samples, predict_period)


# Let's define some dummy dimensions based on the previous steps
# The input_dim for the generator is `batch_size` (5 from previous step).
# The feature_size for the generator is the number of features in your `dataset` DataFram
# The output_dim for the generator is `predict_period` (1 from previous step).
# For the discriminator, input_dim is typically the length of the sequence it discriminat
# If it's discriminating `X` concatenated with `y` (real or generated), it would be `batc

dummy_input_dim_generator = 5 # This is your 'batch_size'
dummy_output_dim_generator = 1 # This is your 'predict_period'
dummy_feature_size_generator = 11 # Assuming 'dataset' has 11 columns (Date + 10 features
                                  # You should get this from `dataset.shape[1]` after the

# Discriminator's input_dim:
# If the discriminator takes the X batch and the corresponding y (real or fake),
# and combines them into a sequence, its length would be `batch_size + predict_period`.
dummy_input_dim_discriminator = dummy_input_dim_generator + dummy_output_dim_generator #

print("Defining Generator Model...")
generator = make_generator_model(
    input_dim=dummy_input_dim_generator,
    output_dim=dummy_output_dim_generator,
    feature_size=dummy_feature_size_generator
)
generator.summary()

print("\nDefining Discriminator Model...")
discriminator = make_discriminator_model(
    input_dim=dummy_input_dim_discriminator
)
discriminator.summary()
```

Import Libraries: Imports necessary libraries like pandas for data manipulation, numpy f

Re-creating tech_df and dataset (for demonstration): This section creates dummy datafram

Define get tech ind Function: This function calculates various technical indicators (Mov

```
 Apply First Set of Operations on dataset:
      It first prints the head and info of the initial dataset.
      Forward Fill (ffill): dataset.iloc[:, 1:].ffill().copy() applies a forward fill to a
      Convert 'Date' to DatetimeIndex and set as index:
          datetime_series = pd.to_datetime(dataset['Date']): Converts the 'Date' column to
          datetime_index = pd.DatetimeIndex(datetime_series.values): Creates a DatetimeInc
          dataset = dataset.set_index(datetime_index): Sets the newly created DatetimeInde
          dataset = dataset.sort_values(by='Date'): Sorts the DataFrame by the date index.
          dataset = dataset.drop(columns='Date'): Drops the original 'Date' column as it's
      It then prints the head and info of the dataset after these transformations.
 Define Data Preparation Functions:
      normalize_data(df, data_range, target_column): This function normalizes the features
      batch_data(x_data, y_data, batch_size, predict_period): This function creates batche
      split_train_test(data): This simple function splits the input data into training and
      predict_index(dataset, X_train, batch_size, prediction_period): This function calcul
 Execute the Steps:
      Normalize data: Calls normalize_data on the preprocessed dataset to get the scaled t
      Batch data: Calls batch_data to create the batched sequences for the input features
      Split train/test: Calls split_train_test on the batched data (X_batched, y_batched,
      Get prediction indices: Calls predict_index to get the DatetimeIndex corresponding t
```

In summary, this cell takes the processed stock data with technical indicators, handles missing values, sets the date as the index, normalizes the data, structures it into sequences for time series modeling, splits it into training and testing sets, and determines the date indices for the predictions.

```
import tensorflow as tf
import numpy as np # For creating dummy data if needed for testing
import pandas as pd # For creating dummy data if needed for testing

# --- Assume generator and discriminator models are already defined from previous steps -
# (Re-defining them here for a complete, runnable snippet, but in your actual
# workflow, these would be available from the previous cell's execution)

# Dummy dimensions for model creation (match the previous step's example)
dummy_input_dim_generator = 5 # This is your 'batch_size'
dummy_output_dim_generator = 1 # This is your 'predict_period'
dummy_feature_size_generator = 11 # Number of features in your dataset

# Discriminator's input_dim: batch_size + predict_period if it takes X + y
dummy_input_dim_discriminator = dummy_input_dim_generator + dummy_output_dim_generator #

def make_generator_model(input_dim, output_dim, feature_size):
    model = tf.keras.Sequential([
```

```
            tf.keras.layers.LSTM(units=1024, return_sequences=True,
                                input_shape=(input_dim, feature_size), recurrent_dropout=0.3
            tf.keras.layers.LSTM(units=512, return_sequences=True, recurrent_dropout=0.3),
            tf.keras.layers.LSTM(units=256, return_sequences=True, recurrent_dropout=0.3),
            tf.keras.layers.LSTM(units=128, return_sequences=True, recurrent_dropout=0.3),
            tf.keras.layers.LSTM(units=64, recurrent_dropout=0.3),
            tf.keras.layers.Dense(32),
            tf.keras.layers.Dense(16),
            tf.keras.layers.Dense(8),
            tf.keras.layers.Dense(units=output_dim)
        ])
        return model


def make_discriminator_model(input_dim):
    cnn_net = tf.keras.Sequential()
    cnn_net.add(tf.keras.layers.Conv1D(8, input_shape=(input_dim + 1, 1), kernel_size=3,
    cnn_net.add(tf.keras.layers.Conv1D(16, kernel_size=3, strides=2, padding='same', acti
    cnn_net.add(tf.keras.layers.Conv1D(32, kernel_size=3, strides=2, padding='same', acti
    cnn_net.add(tf.keras.layers.Conv1D(64, kernel_size=3, strides=2, padding='same', acti
    cnn_net.add(tf.keras.layers.Conv1D(128, kernel_size=1, strides=2, padding='same', act
    # As discussed previously, a Flatten layer might be needed here depending on TensorFl
    cnn_net.add(tf.keras.layers.Flatten()) # Added Flatten as it's common before Dense la
    cnn_net.add(tf.keras.layers.LeakyReLU())
    cnn_net.add(tf.keras.layers.Dense(220, use_bias=False))
    cnn_net.add(tf.keras.layers.LeakyReLU())
    cnn_net.add(tf.keras.layers.Dense(220, use_bias=False, activation='relu'))
    cnn_net.add(tf.keras.layers.Dense(1, activation='sigmoid'))
    return cnn_net


# Create instances of the models
generator = make_generator_model(
    input_dim=dummy_input_dim_generator,
    output_dim=dummy_output_dim_generator,
    feature_size=dummy_feature_size_generator
)
discriminator = make_discriminator_model(
    input_dim=dummy_input_dim_discriminator
)


# --- Define the optimizers ---
# These learning rates are examples; you might need to tune them.
g_optimizer = tf.keras.optimizers.Adam(1e-4)
d_optimizer = tf.keras.optimizers.Adam(1e-4)



# --- Define the loss functions ---
def discriminator_loss(real_output, fake_output):
    """
    Calculates the discriminator's total loss.

    Args:
```

```
                real_output (tf.Tensor): Discriminator's output for real data.
                fake_output (tf.Tensor): Discriminator's output for fake (generated) data.

        Returns:
            tf.Tensor: Total discriminator loss.
        """
        # Use from_logits=True if the discriminator's last layer does not have a sigmoid acti
        # However, your discriminator model has sigmoid, so from_logits=False would be more a
        # if you want to interpret outputs as probabilities.
        # If the GAN framework expects logits, keep from_logits=True and adjust discriminator
        # For common GAN practices, it's often preferred to use from_logits=True and no sigmo
        # Given your discriminator *has* sigmoid, let's keep from_logits=True as you defined,
        loss_f = tf.keras.losses.BinaryCrossentropy(from_logits=True)
        real_loss = loss_f(tf.ones_like(real_output), real_output) # Real data should be clas
        fake_loss = loss_f(tf.zeros_like(fake_output), fake_output) # Fake data should be cla
        total_loss = real_loss + fake_loss
        return total_loss


def generator_loss(fake_output):
        """
        Calculates the generator's loss.
        The generator wants the fake data to be classified as real (1).

        Args:
            fake_output (tf.Tensor): Discriminator's output for fake (generated) data.

        Returns:
            tf.Tensor: Generator loss.
        """
        loss_f = tf.keras.losses.BinaryCrossentropy(from_logits=True)
        loss = loss_f(tf.ones_like(fake_output), fake_output) # Generator wants fake to look
        return loss


# --- Define the training step ---
@tf.function # Compiles the function for faster execution
def train_step(real_x, real_y, yc, generator, discriminator, g_optimizer, d_optimizer):
        """
        Performs one training step for the GAN.

        Args:
            real_x (tf.Tensor): Real input sequence (e.g., X_train batch).
            real_y (tf.Tensor): Real target sequence (e.g., y_train batch).
            yc (tf.Tensor): Context sequence for the target (e.g., yc_train batch).
            generator (tf.keras.Model): The generator model.
            discriminator (tf.keras.Model): The discriminator model.
            g_optimizer (tf.keras.optimizers.Optimizer): Optimizer for the generator.
            d_optimizer (tf.keras.optimizers.Optimizer): Optimizer for the discriminator.

        Returns:
            tuple: (real_y, generated_data, dict_of_losses)
        """
```

```
        """
        with tf.GradientTape() as gen_tape, tf.GradientTape() as disc_tape:
            generated_data = generator(real_x, training=True)

            # Reshape generated_data and real_y to add a feature dimension of 1
            # and concatenate with yc for discriminator input
            # yc shape: (batch_size, sequence_length, features) - this implies yc might be mu
            # If yc is target context, it's usually (batch_size, batch_size_time_steps, 1)
            # Assuming generated_data is (batch_size, predict_period) or (batch_size, predict
            # If predict_period is 1, then generated_data.shape[1] is 1, so reshape is (batch

            # Ensure generated_data is float32 for consistency with model outputs unless spec
            generated_data_reshape = tf.reshape(generated_data, [tf.shape(generated_data)[0],

            # d_fake_input: concatenate generated output with context (yc)
            # This concatenation is along axis=1 (time dimension), meaning yc should be a seq
            # and its features should align if that's intended.
            # For this to work, generated_data_reshape (batch, predict_period, 1) and yc (bat
            # should have compatible shapes or features. If yc is context for target, it's ty
            # Let's assume yc's last dimension is 1 for simplicity of this example.
            # The cast to float64 is important if your models or other inputs are float32.
            # It's generally best to stick to float32 for deep learning unless float64 is exp
            d_fake_input = tf.concat([tf.cast(generated_data_reshape, tf.float32), tf.cast(yc

            # d_real_input: concatenate real target (real_y) with context (yc)
            # real_y_reshape should match the shape of generated_data_reshape for the discrim
            real_y_reshape = tf.reshape(real_y, [tf.shape(real_y)[0], tf.shape(real_y)[1], 1]
            d_real_input = tf.concat([tf.cast(real_y_reshape, tf.float32), tf.cast(yc, tf.flo

            # Discriminator outputs
            real_output = discriminator(d_real_input, training=True)
            fake_output = discriminator(d_fake_input, training=True)

            # Calculate losses
            g_loss = generator_loss(fake_output)
            disc_loss = discriminator_loss(real_output, fake_output)

        # Compute gradients
        gradients_of_generator = gen_tape.gradient(g_loss, generator.trainable_variables)
        gradients_of_discriminator = disc_tape.gradient(disc_loss, discriminator.trainable_va

        # Apply gradients
        g_optimizer.apply_gradients(zip(gradients_of_generator, generator.trainable_variables
        d_optimizer.apply_gradients(zip(gradients_of_discriminator, discriminator.trainable_v

        return real_y, generated_data, {'d_loss': disc_loss, 'g_loss': g_loss}

# --- Example of how you would call train_step ---
# (This part requires X_train, y_train, yc_train from previous steps)
# Let's create dummy batched data for demonstration purposes if you don't have it ready.
# In a real scenario, you would use X_train, y_train, yc_train from your data preparation
```

```
                     _ _ _ _ _ _ _ _ _ _

# Dummy data for train_step call
# X_train from previous batching has shape (num_batches, batch_size, feature_size)
# y_train has shape (num_batches, predict_period)
# yc_train has shape (num_batches, batch_size, predict_period_features) - typically 1 fea
dummy_num_batches = 32
dummy_batch_size = dummy_input_dim_generator # 5
dummy_predict_period = dummy_output_dim_generator # 1
dummy_feature_size = dummy_feature_size_generator # 11

dummy_real_x = tf.random.normal((dummy_num_batches, dummy_batch_size, dummy_feature_size)
dummy_real_y = tf.random.normal((dummy_num_batches, dummy_predict_period), dtype=tf.float
dummy_yc = tf.random.normal((dummy_num_batches, dummy_batch_size, dummy_predict_period),

print("Running a dummy train_step to test the functions...")
try:
    real_y_out, generated_data_out, losses = train_step(
        dummy_real_x, dummy_real_y, dummy_yc, generator, discriminator, g_optimizer, d_op
    )
    print("Train step completed successfully.")
    print(f"Generator Loss: {losses['g_loss'].numpy()}")
    print(f"Discriminator Loss: {losses['d_loss'].numpy()}")
    print(f"Shape of generated_data_out: {generated_data_out.shape}")
except Exception as e:
    print(f"An error occurred during train_step: {e}")
```

```
    /usr/local/lib/python3.11/dist-packages/keras/src/layers/rnn/rnn.py:200: UserWarning:
      super().__init__(**kwargs)
    /usr/local/lib/python3.11/dist-packages/keras/src/layers/activations/leaky_relu.py:41
      warnings.warn(
    /usr/local/lib/python3.11/dist-packages/keras/src/layers/convolutional/base_conv.py:1
      super().__init__(activity_regularizer=activity_regularizer, **kwargs)
    Running a dummy train_step to test the functions...
    /usr/local/lib/python3.11/dist-packages/keras/src/backend/tensorflow/nn.py:780: UserW
      output, from_logits = _get_logits(
    Train step completed successfully.
    Generator Loss: 0.6937790513038635
    Discriminator Loss: 1.386208176612854
    Shape of generated_data_out: (32, 1)
```

```
import tensorflow as tf
import numpy as np
import pandas as pd
import matplotlib.pyplot as plt
from tqdm import tqdm # For progress bar
from sklearn.metrics import mean_squared_error
import os # For creating directories for model saving

# --- Re-creating necessary components if not already in session ---
# (This ensures the code is runnable independently)
```

```python
# Dummy dimensions for model creation
dummy_input_dim_generator = 5 # batch_size
dummy_output_dim_generator = 1 # predict_period
dummy_feature_size_generator = 11 # num_features

dummy_input_dim_discriminator = dummy_input_dim_generator + dummy_output_dim_generator #

# Model definitions (from previous turns)
def make_generator_model(input_dim, output_dim, feature_size):
    model = tf.keras.Sequential([
        tf.keras.layers.LSTM(units=1024, return_sequences=True,
                             input_shape=(input_dim, feature_size), recurrent_dropout=0.3
        tf.keras.layers.LSTM(units=512, return_sequences=True, recurrent_dropout=0.3),
        tf.keras.layers.LSTM(units=256, return_sequences=True, recurrent_dropout=0.3),
        tf.keras.layers.LSTM(units=128, return_sequences=True, recurrent_dropout=0.3),
        tf.keras.layers.LSTM(units=64, recurrent_dropout=0.3),
        tf.keras.layers.Dense(32),
        tf.keras.layers.Dense(16),
        tf.keras.layers.Dense(8),
        tf.keras.layers.Dense(units=output_dim)
    ])
    return model

def make_discriminator_model(input_dim):
    cnn_net = tf.keras.Sequential()
    cnn_net.add(tf.keras.layers.Conv1D(8, input_shape=(input_dim + 1, 1), kernel_size=3,
    cnn_net.add(tf.keras.layers.Conv1D(16, kernel_size=3, strides=2, padding='same', acti
    cnn_net.add(tf.keras.layers.Conv1D(32, kernel_size=3, strides=2, padding='same', acti
    cnn_net.add(tf.keras.layers.Conv1D(64, kernel_size=3, strides=2, padding='same', acti
    cnn_net.add(tf.keras.layers.Conv1D(128, kernel_size=1, strides=2, padding='same', act
    cnn_net.add(tf.keras.layers.Flatten())
    cnn_net.add(tf.keras.layers.LeakyReLU())
    cnn_net.add(tf.keras.layers.Dense(220, use_bias=False))
    cnn_net.add(tf.keras.layers.LeakyReLU())
    cnn_net.add(tf.keras.layers.Dense(220, use_bias=False, activation='relu'))
    cnn_net.add(tf.keras.layers.Dense(1, activation='sigmoid'))
    return cnn_net

# Create instances of the models
generator = make_generator_model(
    input_dim=dummy_input_dim_generator,
    output_dim=dummy_output_dim_generator,
    feature_size=dummy_feature_size_generator
)
discriminator = make_discriminator_model(
    input_dim=dummy_input_dim_discriminator
)

# Optimizers
g_optimizer = tf.keras.optimizers.Adam(1e-4)
```

```
d_optimizer = tf.keras.optimizers.Adam(1e-4)


# Loss functions (from previous turn)
def discriminator_loss(real_output, fake_output):
    loss_f = tf.keras.losses.BinaryCrossentropy(from_logits=True)
    real_loss = loss_f(tf.ones_like(real_output), real_output)
    fake_loss = loss_f(tf.zeros_like(fake_output), fake_output)
    total_loss = real_loss + fake_loss
    return total_loss


def generator_loss(fake_output):
    loss_f = tf.keras.losses.BinaryCrossentropy(from_logits=True)
    loss = loss_f(tf.ones_like(fake_output), fake_output)
    return loss


# train_step function (from previous turn)
@tf.function
def train_step(real_x, real_y, yc, generator, discriminator, g_optimizer, d_optimizer):
    with tf.GradientTape() as gen_tape, tf.GradientTape() as disc_tape:
        generated_data = generator(real_x, training=True)
        generated_data_reshape = tf.reshape(generated_data, [tf.shape(generated_data)[0],
        d_fake_input = tf.concat([tf.cast(generated_data_reshape, tf.float32), tf.cast(yc

        real_y_reshape = tf.reshape(real_y, [tf.shape(real_y)[0], tf.shape(real_y)[1], 1]
        d_real_input = tf.concat([tf.cast(real_y_reshape, tf.float32), tf.cast(yc, tf.flo

        real_output = discriminator(d_real_input, training=True)
        fake_output = discriminator(d_fake_input, training=True)

        g_loss = generator_loss(fake_output)
        disc_loss = discriminator_loss(real_output, fake_output)

    gradients_of_generator = gen_tape.gradient(g_loss, generator.trainable_variables)
    gradients_of_discriminator = disc_tape.gradient(disc_loss, discriminator.trainable_va

    g_optimizer.apply_gradients(zip(gradients_of_generator, generator.trainable_variables
    d_optimizer.apply_gradients(zip(gradients_of_discriminator, discriminator.trainable_v

    return real_y, generated_data, {'d_loss': disc_loss, 'g_loss': g_loss}


# --- Dummy Batched Data for Training ---
# This simulates the X_train, y_train, yc_train from your data preparation step
dummy_num_batches = 100 # Number of batches in your training data
dummy_batch_size = dummy_input_dim_generator # 5
dummy_predict_period = dummy_output_dim_generator # 1
dummy_feature_size = dummy_feature_size_generator # 11 (from your dataset's columns)


real_x_train = tf.random.normal((dummy_num_batches, dummy_batch_size, dummy_feature_size)
real_y_train = tf.random.normal((dummy_num_batches, dummy_predict_period), dtype=tf.float
# yc should have the same batch_size in the sequence dimension, and match features or be
# If yc represents the actual historical 'Close' prices for the batch_size period, its fe
```

```
    yc_train = tf.random.normal((dummy_num_batches, dummy_batch_size, 1), dtype=tf.float32)

    # Define stock_name for saving models
    stock_name = "DEMO_STOCK"



    # --- Define the train function ---
    def train(real_x, real_y, yc, Epochs, generator, discriminator, g_optimizer, d_optimizer,
        """
        Trains a GAN model for time series prediction.

        Args:
            real_x (tf.Tensor): Batched real input sequences (X_train).
            real_y (tf.Tensor): Batched real target sequences (y_train).
            yc (tf.Tensor): Batched context sequences for targets (yc_train).
            Epochs (int): Number of training epochs.
            generator (tf.keras.Model): The generator model.
            discriminator (tf.keras.Model): The discriminator model.
            g_optimizer (tf.keras.optimizers.Optimizer): Optimizer for the generator.
            d_optimizer (tf.keras.optimizers.Optimizer): Optimizer for the discriminator.
            checkpoint (int): Interval (in epochs) at which to save models.

        Returns:
            tuple: (Predicted_price, Real_price, NRMSE)
        """
        train_info = {}
        train_info["discriminator_loss"] = []
        train_info["generator_loss"] = []

        # Initialize lists to store predictions and real prices for evaluation
        all_predicted_prices_flat = []
        all_real_prices_flat = []

        # Create directory for saving models if it doesn't exist
        model_save_path = f'./models_gan/{stock_name}/'
        os.makedirs(model_save_path, exist_ok=True)
        print(f"Models will be saved to: {os.path.abspath(model_save_path)}")


        for epoch in tqdm(range(Epochs), desc="Training GAN"):
            # Perform one training step
            real_price_batch, fake_price_batch, loss = train_step(real_x, real_y, yc, generat

            # Store batch losses
            train_info["discriminator_loss"].append(loss['d_loss'].numpy())
            train_info["generator_loss"].append(loss['g_loss'].numpy())

            # Collect real and predicted prices from the current batch for overall evaluation
            # Ensure these are flattened correctly for later MSE calculation
            # fake_price_batch is (batch_size, predict_period)
            # real_price_batch is (batch_size, predict_period)
```

```
        # real_price_batch is (batch_size, predict_period)
        all_predicted_prices_flat.extend(fake_price_batch.numpy().flatten())
        all_real_prices_flat.extend(real_price_batch.numpy().flatten())


    # Save model every X checkpoints
    if (epoch + 1) % checkpoint == 0:
        print(f'\nEpoch {epoch + 1}/{Epochs} - Discriminator Loss: {loss["d_loss"].nu
        try:
            tf.keras.models.save_model(generator, os.path.join(model_save_path, f'gen
            tf.keras.models.save_model(discriminator, os.path.join(model_save_path, f
            print(f"Models saved at epoch {epoch+1}")
        except Exception as e:
            print(f"Error saving model at epoch {epoch+1}: {e}")

# Convert collected prices to numpy arrays for final evaluation
# These will be 1D arrays of all collected predictions/real values
Predicted_price_final = np.array(all_predicted_prices_flat)
Real_price_final = np.array(all_real_prices_flat)

# Plotting Losses
plt.figure(figsize=(12, 8)) # Use a single figure for subplots
plt.subplot(2, 1, 1)
plt.plot(train_info["discriminator_loss"], label='Discriminator Loss', color='#000000
plt.title('GAN Training Losses')
plt.xlabel('Epoch')
plt.ylabel('Loss')
plt.legend()
plt.grid(True)

plt.subplot(2, 1, 2)
plt.plot(train_info["generator_loss"], label='Generator Loss', color='#000000')
plt.xlabel('Epoch')
plt.ylabel('Loss')
plt.legend()
plt.grid(True)

plt.tight_layout() # Adjust layout to prevent overlap
plt.show()

# Calculate NRMSE
# Ensure Real_price_final and Predicted_price_final are compatible for MSE
# They should both be 1D arrays of equal length.
# The original code's reshaping after the loop seems to assume a single batch,
# but with appending to lists in a loop, you get all predictions.
# So, we use the flattened lists directly.
nrmse = np.sqrt(mean_squared_error(Real_price_final, Predicted_price_final)) / np.mea

return Predicted_price_final, Real_price_final, nrmse
```

```
# --- Example Call to the train function ---
print("\nStarting GAN training...")
# Set number of epochs for demonstration. Use a small number for quick test.
EPOCHS = 5 # You can increase this for actual training, e.g., 500, 1000, 5000

predicted_prices, real_prices, nrmse_value = train(
    real_x_train, real_y_train, yc_train, EPOCHS,
    generator, discriminator, g_optimizer, d_optimizer,
    checkpoint=5 # Save models every 5 epochs for this example
)

print(f"\nTraining finished.")
print(f"Normalized Root Mean Squared Error (NRMSE): {nrmse_value:.4f}")
print(f"Shape of final Predicted Prices: {predicted_prices.shape}")
print(f"Shape of final Real Prices: {real_prices.shape}")
```

Import Libraries: Imports necessary libraries for building and training the GAN, includi

Re-creating Necessary Components: This section redefines the make_generator_model and ma

Loss Functions: The discriminator_loss and generator_loss functions are redefined. These

train_step Function: This function, decorated with @tf.function for performance, defines

    It takes real input features (real_x), real target values (real_y), and context (yc)

    It uses tf.GradientTape to record operations for automatic differentiation.

    The generator generates fake_price_batch from real_x.

    The generated_data_reshape and real_y_reshape are created to match the expected inpu

    d_fake_input and d_real_input are created by concatenating the generated/real target

    The discriminator evaluates both the real and fake inputs.

    The g_loss and disc_loss are calculated using the defined loss functions.

    Gradients are computed using the tapes.

    The optimizers apply the gradients to update the model weights.

    The function returns the real and generated prices for the batch, along with the los

Dummy Batched Data: This section creates dummy data (real_x_train, real_y_train, yc_trai

Define the train Function: This function orchestrates the entire training process.

    It initializes lists to store losses and predictions.

    It creates a directory to save the models.

    It loops through the specified number of Epochs.

    In each epoch, it calls train_step to perform a forward and backward pass for both t

    It appends the batch losses to the train_info dictionary.

    It extends the all_predicted_prices_flat and all_real_prices_flat lists with the ger

    It saves the generator and discriminator models periodically based on the checkpoint

After training, it plots the generator and discriminator losses over epochs.

It calculates the Normalized Root Mean Squared Error (NRMSE) between the Real_price_

It returns the final predicted prices, real prices, and the NRMSE.

Example Call to the train function: This section sets the number of epochs and calls the

In summary, this cell defines the training loop for the GAN, including how the generator and discriminator are updated, how losses are calculated and tracked, and how the trained models and training progress are saved and visualized.

```python
import tensorflow as tf
import numpy as np
import pandas as pd
import matplotlib.pyplot as plt
from joblib import load # For loading the MinMaxScaler
from sklearn.metrics import mean_squared_error
import os


# --- Re-importing and re-creating necessary components for independent execution ---
# (If you are running this in a continuous session, these might already be available)

# Dummy dimensions and variables from previous steps
dummy_input_dim_generator = 5 # batch_size
dummy_output_dim_generator = 1 # predict_period
dummy_feature_size_generator = 11 # num_features (from dataset.shape[1])
dummy_input_dim_discriminator = dummy_input_dim_generator + dummy_output_dim_generator

output_dim = dummy_output_dim_generator # Alias for clarity, matches predict_period
stock_name = "DEMO_STOCK" # Matches previous usage for model saving

# Model definitions (from previous turns) - required for eval_op
def make_generator_model(input_dim, output_dim, feature_size):
    model = tf.keras.Sequential([
        tf.keras.layers.LSTM(units=1024, return_sequences=True,
                             input_shape=(input_dim, feature_size), recurrent_dropout=0.3
        tf.keras.layers.LSTM(units=512, return_sequences=True, recurrent_dropout=0.3),
        tf.keras.layers.LSTM(units=256, return_sequences=True, recurrent_dropout=0.3),
        tf.keras.layers.LSTM(units=128, return_sequences=True, recurrent_dropout=0.3),
        tf.keras.layers.LSTM(units=64, recurrent_dropout=0.3),
        tf.keras.layers.Dense(32),
        tf.keras.layers.Dense(16),
        tf.keras.layers.Dense(8),
        tf.keras.layers.Dense(units=output_dim)
    ])
    return model

# Create a dummy generator instance (or load your trained generator)
generator = make_generator_model(
    input_dim=dummy_input_dim_generator,
```

```
        output_dim=dummy_output_dim_generator,
        feature_size=dummy_feature_size_generator
)
# For a real scenario, you'd load a trained generator:
# generator = tf.keras.models.load_model(f'./models_gan/{stock_name}/generator_V_LAST_EPO
# Or use the `generator` object returned from the `train` function if running continuousl


# --- Dummy Data for Plotting Functions ---
# These would typically come from your data preparation (X_train, y_train, X_test, y_test
# and the results of your `train` function (predicted_prices, real_prices)
# and `predict_index` function (index_train, index_test).

# Creating dummy data similar to the output of previous steps
dummy_num_train_samples = 80 # Corresponds to len(X_train)
dummy_num_test_samples = 20 # Corresponds to len(X_test)
dummy_batch_size = dummy_input_dim_generator # 5
dummy_predict_period = dummy_output_dim_generator # 1
dummy_feature_size = dummy_feature_size_generator # 11

# Dummy data for training results (from `train` function)
# These should be flattened arrays
dummy_real_prices_train = np.random.rand(dummy_num_train_samples * dummy_predict_period)
dummy_predicted_prices_train = np.random.rand(dummy_num_train_samples * dummy_predict_per

# Dummy data for test inputs (from `split_train_test`)
dummy_X_test = tf.random.normal((dummy_num_test_samples, dummy_batch_size, dummy_feature_
dummy_y_test = np.random.rand(dummy_num_test_samples * dummy_predict_period) # Flattened

# Dummy DatetimeIndex (from `predict_index` function)
# Ensure these indices are unique and represent dates.
# The `predict_index` function returns actual DatetimeIndex objects.
# Let's create dummy ones that resemble actual dates.
dummy_train_dates = pd.date_range(start='2024-01-01', periods=dummy_num_train_samples + d
dummy_index_train = dummy_train_dates[dummy_batch_size : dummy_num_train_samples + dummy_

dummy_test_dates = pd.date_range(start=dummy_train_dates[-1] + pd.Timedelta(days=1), peri
dummy_index_test = dummy_test_dates[dummy_batch_size : dummy_num_test_samples + dummy_bat


# --- Create dummy scaler files for demonstration ---
# In a real scenario, these would be created by the normalize_data function.
from sklearn.preprocessing import MinMaxScaler
from joblib import dump

# Create dummy data for scalers
dummy_scaler_data = np.random.rand(100, dummy_feature_size_generator) * 100
dummy_target_scaler_data = np.random.rand(100, 1) * 100

X_scaler_dummy = MinMaxScaler(feature_range=(-1, 1))
y_scaler_dummy = MinMaxScaler(feature_range=(-1, 1))
```

```
X_scaler_dummy.fit(dummy_scaler_data)
y_scaler_dummy.fit(dummy_target_scaler_data)


# Ensure the directory exists if not running in Google Colab's default /content/
if not os.path.exists('./content/'):
    os.makedirs('./content/')
dump(X_scaler_dummy, open('./content/X_scaler.pkl', 'wb'))
dump(y_scaler_dummy, open('./content/y_scaler.pkl', 'wb'))
print("Dummy X_scaler.pkl and y_scaler.pkl created in ./content/")



# --- Define eval_op ---
@tf.function
def eval_op(generator_model, real_x):
    """
    Evaluates the generator model on input data.

    Args:
        generator_model (tf.keras.Model): The trained generator model.
        real_x (tf.Tensor): Input data for the generator (e.g., X_test batch).

    Returns:
        tf.Tensor: Generated (predicted) data.
    """
    # Set training=False for inference
    generated_data = generator_model(real_x, training=False)
    return generated_data

# --- Define plot_results (for training data) ---
def plot_results(Real_price, Predicted_price, index_train):
    """
    Plots the real vs. predicted prices for the training data and calculates RMSE.

    Args:
        Real_price (np.ndarray): Array of real prices (scaled).
        Predicted_price (np.ndarray): Array of predicted prices (scaled).
        index_train (pd.DatetimeIndex): DatetimeIndex for the training predictions.
    """
    print("\n--- Plotting Training Results ---")
    X_scaler = load(open('./content/X_scaler.pkl', 'rb')) # Assuming path
    y_scaler = load(open('./content/y_scaler.pkl', 'rb'))

    rescaled_Real_price = y_scaler.inverse_transform(Real_price.reshape(-1, 1)) # Reshape
    rescaled_Predicted_price = y_scaler.inverse_transform(Predicted_price.reshape(-1, 1))

    # Re-structure for combining with DatetimeIndex and calculating means
    # The original loop `predict_result = pd.concat([predict_result, y_predict], axis=1,
    # seems to expect each `y_predict` to be a column, which results in many columns.
    # If output_dim is 1, then each `y_predict` is a single value, and this concatenation
    # Let's adjust to create a single series for direct plotting against index.
```

```
        # Check if index_train is long enough for the predictions.
        # Real_price and Predicted_price are flattened 1D arrays from the 'train' function.
        # They represent all individual predictions, so they should align directly with index

        predict_result_series = pd.Series(rescaled_Predicted_price.flatten(), index=index_tra
        real_price_series = pd.Series(rescaled_Real_price.flatten(), index=index_train[:len(r

        # The original code's way of creating DataFrames and then taking mean(axis=1)
        # is more suitable if output_dim > 1 (multi-step prediction) and you're combining
        # multiple overlapping predictions for the same date.
        # Given output_dim = 1, a direct series is more appropriate for `plot_results`.
        # I'll stick to the spirit of taking mean if applicable, but simplify if output_dim i

        if output_dim == 1:
            # If output_dim is 1, each prediction is for a unique point, no need for mean.
            predicted_mean = predict_result_series
            real_mean = real_price_series
        else:
            # Reconstruct as per your original logic if output_dim > 1 (e.g., if each 'i' is
            # This part requires a more complex re-alignment of predictions to dates if batch
            # For this example, assuming output_dim=1, so simplifying.
            # If output_dim > 1, the `batch_data` and `predict_index` would need to generate
            # overlapping predictions, and then the mean() approach makes sense.
            # Given the current `batch_data` and `predict_period=1`, the following logic migh
            # to be re-evaluated if `output_dim` is truly > 1.
            # For now, let's assume `Real_price` and `Predicted_price` are already flattened

            # The original code's loop implies something like this if output_dim > 1 and ther
            # However, with a simple list.extend() in `train`, we just have a flat list of pr
            # So, the simpler Series approach is more direct here.
            predicted_mean = predict_result_series
            real_mean = real_price_series


    plt.figure(figsize=(16, 8))
    plt.plot(real_mean, label='Real price')
    plt.plot(predicted_mean, color = 'r', label='Predicted price')
    plt.xlabel("Date")
    plt.ylabel("Stock price")
    plt.legend(loc="upper left", fontsize=16)
    plt.title("The result of Training", fontsize=20)
    plt.grid(True)
    plt.show()

    # Calculate RMSE
    # Ensure predicted_mean and real_mean are aligned.
    # They should be aligned by their common index.
    common_index = predicted_mean.index.intersection(real_mean.index)
    aligned_predicted = predicted_mean[common_index]
    aligned_real = real_mean[common_index]
```

```
        if not aligned_predicted.empty:
            RMSE = np.sqrt(mean_squared_error(aligned_predicted, aligned_real))
            print(f'-- Train RMSE -- {RMSE:.4f}')
        else:
            print("-- Train RMSE -- Cannot calculate RMSE, no overlapping data points.")


    # --- Define plot_test_data ---
    def plot_test_data(Real_test_price, Predicted_test_price, index_test):
        """
        Plots the real vs. predicted prices for the test data and calculates RMSE.

        Args:
            Real_test_price (np.ndarray): Array of real test prices (scaled).
            Predicted_test_price (np.ndarray): Array of predicted test prices (scaled).
            index_test (pd.DatetimeIndex): DatetimeIndex for the test predictions.
        """
        print("\n--- Plotting Test Results ---")
        X_scaler = load(open('./content/X_scaler.pkl', 'rb')) # Assuming path
        y_scaler = load(open('./content/y_scaler.pkl', 'rb'))

        rescaled_Real_price = y_scaler.inverse_transform(Real_test_price.reshape(-1, 1)) # Re
        rescaled_Predicted_price = y_scaler.inverse_transform(Predicted_test_price.reshape(-1

        # Similar to plot_results, creating series directly
        predict_result_series = pd.Series(rescaled_Predicted_price.flatten(), index=index_tes
        real_price_series = pd.Series(rescaled_Real_price.flatten(), index=index_test[:len(re

        if output_dim == 1:
            predicted_mean = predict_result_series
            real_mean = real_price_series
        else:
            # More complex re-alignment needed if output_dim > 1 and predictions overlap
            predicted_mean = predict_result_series
            real_mean = real_price_series


        # Calculate RMSE
        common_index = predicted_mean.index.intersection(real_mean.index)
        aligned_predicted = predicted_mean[common_index]
        aligned_real = real_mean[common_index]

        if not aligned_predicted.empty:
            RMSE = np.sqrt(mean_squared_error(aligned_predicted, aligned_real))
            print(f'Test RMSE: {RMSE:.4f}')
        else:
            print("Test RMSE: Cannot calculate RMSE, no overlapping data points.")


        plt.figure(figsize=(16, 8))
```

```
    plt.figure(figsize=(16, 8))
    plt.plot(real_mean, color='#00008B', label='Real price')
    plt.plot(predicted_mean, color = '#8B0000', linestyle='--', label='Predicted price')
    plt.xlabel("Date")
    plt.ylabel("Stock price")
    plt.legend(loc="upper left", fontsize=16)
    plt.title(f"Prediction on test data for {stock_name}", fontsize=20)
    plt.grid(True)
    plt.show()


# --- Example usage ---

# 1. Get predictions for test data using the generator
print("\nGenerating predictions for test data...")
# eval_op expects real_x (which would be X_test from your batched data)
# Use the dummy_X_test created earlier
predicted_test_data = eval_op(generator, dummy_X_test)
print(f"Shape of raw predicted_test_data: {predicted_test_data.shape}")


# Ensure that the dummy y_test is also a 1D array to match the flattened outputs.
# The `y_test` from `split_train_test` is shaped `(num_test_batches, predict_period)`.
# So, `Real_test_price` should be `y_test.flatten()` when passed to `plot_test_data`.
Real_test_price_for_plot = dummy_y_test # This is already flattened in dummy setup



# 2. Plot results for training data
# `predicted_prices` and `real_prices` come from the `train` function's return values.
# They should already be flattened 1D numpy arrays.
plot_results(dummy_real_prices_train, dummy_predicted_prices_train, dummy_index_train)

# 3. Plot results for test data
plot_test_data(Real_test_price_for_plot, predicted_test_data.numpy().flatten(), dummy_ind
```

```python
import tensorflow as tf
import numpy as np
import pandas as pd # For dummy data if needed
import os # For checking file paths

# --- Re-create dummy data and model definitions for independent execution ---
# (If running in a continuous session, you can skip these re-definitions)

# Dummy data from previous steps for X_train.shape
dummy_num_train_batches = 80 # Corresponds to len(X_train)
dummy_batch_size = 5 # X_train.shape[1]
dummy_predict_period = 1 # output_dim
dummy_feature_size = 11 # X_train.shape[2]

# Simulating X_train, which is needed for model instantiation shapes
X_train = tf.random.normal((dummy_num_train_batches, dummy_batch_size, dummy_feature_size)

# Global variables for model instantiation
output_dim = dummy_predict_period # From your previous definitions
stock_name = "DEMO_STOCK" # Used for model saving if uncommented later

# Model definitions (from previous turns)
def make_generator_model(input_dim, output_dim, feature_size):
    model = tf.keras.Sequential([
        tf.keras.layers.LSTM(units=1024, return_sequences=True,
                             input_shape=(input_dim, feature_size), recurrent_dropout=0.3)
        tf.keras.layers.LSTM(units=512, return_sequences=True, recurrent_dropout=0.3),
        tf.keras.layers.LSTM(units=256, return_sequences=True, recurrent_dropout=0.3),
        tf.keras.layers.LSTM(units=128, return_sequences=True, recurrent_dropout=0.3),
        tf.keras.layers.LSTM(units=64, recurrent_dropout=0.3),
        tf.keras.layers.Dense(32),
        tf.keras.layers.Dense(16),
        tf.keras.layers.Dense(8),
        tf.keras.layers.Dense(units=output_dim)
    ])
    return model


def make_discriminator_model(input_dim_discriminator_sequence_length):
    # Discriminator expects a sequence of length input dim discriminator sequence length +
```

```python
    # and 1 feature, e.g., (batch_size + predict_period, 1) if combined X and y
    cnn_net = tf.keras.Sequential()
    cnn_net.add(tf.keras.layers.Conv1D(8, input_shape=(input_dim_discriminator_sequence_len
    cnn_net.add(tf.keras.layers.Conv1D(16, kernel_size=3, strides=2, padding='same', activa
    cnn_net.add(tf.keras.layers.Conv1D(32, kernel_size=3, strides=2, padding='same', activa
    cnn_net.add(tf.keras.layers.Conv1D(64, kernel_size=3, strides=2, padding='same', activa
    cnn_net.add(tf.keras.layers.Conv1D(128, kernel_size=1, strides=2, padding='same', activ
    cnn_net.add(tf.keras.layers.Flatten())
    cnn_net.add(tf.keras.layers.LeakyReLU())
    cnn_net.add(tf.keras.layers.Dense(220, use_bias=False))
    cnn_net.add(tf.keras.layers.LeakyReLU())
    cnn_net.add(tf.keras.layers.Dense(220, use_bias=False, activation='relu'))
    cnn_net.add(tf.keras.layers.Dense(1, activation='sigmoid'))
    return cnn_net



# --- Set up learning rate and optimizers ---
learning_rate = 5e-4
epochs = 500 # This is just a parameter; actual training loop not in this block

print(f"Setting learning rate to: {learning_rate}")
g_optimizer = tf.keras.optimizers.Adam(learning_rate=learning_rate) # Corrected 'lr' to 'le
d_optimizer = tf.keras.optimizers.Adam(learning_rate=learning_rate) # Corrected 'lr' to 'le
print("Optimizers initialized.")


# --- Instantiate Generator and Discriminator Models ---
# X_train.shape: (num_batches, batch_size, num_features)
# Generator: input_dim = batch_size, output_dim = predict_period, feature_size = num_featur
generator = make_generator_model(X_train.shape[1], output_dim, X_train.shape[2])

# Discriminator: input_dim for discriminator should be the total length of the sequence it
# This is typically (batch_size + predict_period).
# Given `make_discriminator_model` uses `input_dim + 1`, we pass `(batch_size + predict_per
discriminator_input_sequence_length = X_train.shape[1] + output_dim # batch_size + predict_
discriminator = make_discriminator_model(discriminator_input_sequence_length - 1)



print("\nGenerator Model Summary:")
generator.summary()

print("\nDiscriminator Model Summary:")
discriminator.summary()



# --- Plotting Model Architectures ---
print("\nPlotting generator_keras_model.png and discriminator_keras_model.png...")

# Ensure graphviz is installed for plot_model to work
# If you get an error, you might need to run:
# pip install graphviz pydot
```

```
try:
    tf.keras.utils.plot_model(generator, to_file='generator_keras_model.png', show_shapes=
    print("Generator model plot saved to 'generator_keras_model.png'")
except Exception as e:
    print(f"Could not plot generator model. Make sure Graphviz is installed and in your PAT

try:
    tf.keras.utils.plot_model(discriminator, to_file='discriminator_keras_model.png', show_
    print("Discriminator model plot saved to 'discriminator_keras_model.png'")
except Exception as e:
    print(f"Could not plot discriminator model. Make sure Graphviz is installed and in your

print("\nSetup complete. You can now proceed to train the GAN using the 'train' function.")
```

Import Libraries: Imports necessary libraries including TensorFlow for building and mana
Re-create dummy data and model definitions: This section is included for the code to be
Set up learning rate and optimizers:

    learning_rate = 5e-4: Sets the learning rate for the optimizers. This is a hyperpara
    g_optimizer = tf.keras.optimizers.Adam(learning_rate=learning_rate): Creates an Adam
    d_optimizer = tf.keras.optimizers.Adam(learning_rate=learning_rate): Creates an Adam
Instantiate Generator and Discriminator Models:

    generator = make_generator_model(X_train.shape[1], output_dim, X_train.shape[2]): Cr
    discriminator_input_sequence_length = X_train.shape[1] + output_dim: Calculates the
    discriminator = make_discriminator_model(discriminator_input_sequence_length - 1): C
Plotting Model Architectures:

    tf.keras.utils.plot_model(...): This function from TensorFlow/Keras is used to visua
    to_file='generator_keras_model.png' and to_file='discriminator_keras_model.png': Spe
    show_shapes=True: Includes the input and output shapes of each layer in the plot.
    show_layer_names=True: Includes the names of each layer in the plot.
    The try...except blocks are used to catch potential errors if Graphviz (a dependency

In summary, this cell sets up the essential components for training the GAN: it defines the learning rate, initializes the optimizers, creates instances of the generator and discriminator models with the correct input/output dimensions based on the data shape, and then generates plots of the model architectures for visualization.

```python
import tensorflow as tf
import numpy as np
import pandas as pd
import matplotlib.pyplot as plt
from tqdm import tqdm
from joblib import load, dump # For loading/saving scalers
from sklearn.metrics import mean_squared_error
from sklearn.preprocessing import MinMaxScaler
import os

# --- Re-create ALL necessary components and dummy data for a fresh run ---
# (If you're running this in a continuous Jupyter/Colab session after all previous
# blocks, you can comment out the re-creation of models, optimizers, and data)

# Global variables/parameters
learning_rate = 5e-4
epochs = 500 # Set to the desired number of training epochs
checkpoint_interval = 50 # Interval for saving models during training (e.g., every 50 epo

# Dummy dimensions (consistent with previous steps)
dummy_batch_size = 5 # X_train.shape[1] for generator input_dim
output_dim = 1        # output_dim for generator
dummy_feature_size = 11 # X_train.shape[2] for generator feature_size
stock_name = "DEMO_STOCK" # Used for model saving paths

# 1. Model Definitions (re-defined for standalone run)
def make_generator_model(input_dim, output_dim, feature_size):
    model = tf.keras.Sequential([
        tf.keras.layers.LSTM(units=1024, return_sequences=True,
                             input_shape=(input_dim, feature_size), recurrent_dropout=0.3
        tf.keras.layers.LSTM(units=512, return_sequences=True, recurrent_dropout=0.3),
        tf.keras.layers.LSTM(units=256, return_sequences=True, recurrent_dropout=0.3),
        tf.keras.layers.LSTM(units=128, return_sequences=True, recurrent_dropout=0.3),
        tf.keras.layers.LSTM(units=64, recurrent_dropout=0.3),
        tf.keras.layers.Dense(32),
        tf.keras.layers.Dense(16),
        tf.keras.layers.Dense(8),
        tf.keras.layers.Dense(units=output_dim)
    ])
    return model

def make_discriminator_model(input_dim_discriminator_sequence_length):
    cnn_net = tf.keras.Sequential()
    cnn_net.add(tf.keras.layers.Conv1D(8, input_shape=(input_dim_discriminator_sequence_l
```

```python
        cnn_net.add(tf.keras.layers.Conv1D(16, kernel_size=3, strides=2, padding='same', acti
        cnn_net.add(tf.keras.layers.Conv1D(32, kernel_size=3, strides=2, padding='same', acti
        cnn_net.add(tf.keras.layers.Conv1D(64, kernel_size=3, strides=2, padding='same', acti
        cnn_net.add(tf.keras.layers.Conv1D(128, kernel_size=1, strides=2, padding='same', act
        cnn_net.add(tf.keras.layers.Flatten())
        cnn_net.add(tf.keras.layers.LeakyReLU())
        cnn_net.add(tf.keras.layers.Dense(220, use_bias=False))
        cnn_net.add(tf.keras.layers.LeakyReLU())
        cnn_net.add(tf.keras.layers.Dense(220, use_bias=False, activation='relu'))
        cnn_net.add(tf.keras.layers.Dense(1, activation='sigmoid'))
        return cnn_net


# 2. Optimizers
g_optimizer = tf.keras.optimizers.Adam(learning_rate=learning_rate)
d_optimizer = tf.keras.optimizers.Adam(learning_rate=learning_rate)

# 3. Model Instances (initialized)
generator = make_generator_model(dummy_batch_size, output_dim, dummy_feature_size)
discriminator_input_sequence_length = dummy_batch_size + output_dim # e.g., 5 + 1 = 6
discriminator = make_discriminator_model(discriminator_input_sequence_length - 1) # Passe

# 4. Loss functions (re-defined for standalone run)
def discriminator_loss(real_output, fake_output):
    loss_f = tf.keras.losses.BinaryCrossentropy(from_logits=True)
    real_loss = loss_f(tf.ones_like(real_output), real_output)
    fake_loss = loss_f(tf.zeros_like(fake_output), fake_output)
    total_loss = real_loss + fake_loss
    return total_loss


def generator_loss(fake_output):
    loss_f = tf.keras.losses.BinaryCrossentropy(from_logits=True)
    loss = loss_f(tf.ones_like(fake_output), fake_output)
    return loss

# 5. train_step function (re-defined for standalone run)
@tf.function
def train_step(real_x, real_y, yc, generator_model, discriminator_model, g_optimizer_step
    with tf.GradientTape() as gen_tape, tf.GradientTape() as disc_tape:
        generated_data = generator_model(real_x, training=True)
        generated_data_reshape = tf.reshape(generated_data, [tf.shape(generated_data)[0],
        d_fake_input = tf.concat([tf.cast(generated_data_reshape, tf.float32), tf.cast(yc

        real_y_reshape = tf.reshape(real_y, [tf.shape(real_y)[0], tf.shape(real_y)[1], 1]
        d_real_input = tf.concat([tf.cast(real_y_reshape, tf.float32), tf.cast(yc, tf.flo

        real_output = discriminator_model(d_real_input, training=True)
        fake_output = discriminator_model(d_fake_input, training=True)

        g_loss = generator_loss(fake_output)
        disc_loss = discriminator_loss(real_output, fake_output)
```

```
        gradients_of_generator = gen_tape.gradient(g_loss, generator_model.trainable_variable
        gradients_of_discriminator = disc_tape.gradient(disc_loss, discriminator_model.traina

        g_optimizer_step.apply_gradients(zip(gradients_of_generator, generator_model.trainabl
        d_optimizer_step.apply_gradients(zip(gradients_of_discriminator, discriminator_model.

        return real_y, generated_data, {'d_loss': disc_loss, 'g_loss': g_loss}

# 6. train function (re-defined for standalone run)
def train(real_x, real_y, yc, Epochs, generator_model, discriminator_model, g_optimizer_t
    train_info = {}
    train_info["discriminator_loss"] = []
    train_info["generator_loss"] = []

    all_predicted_prices_flat = []
    all_real_prices_flat = []

    model_save_path = f'./models_gan/{stock_name}/'
    os.makedirs(model_save_path, exist_ok=True)
    print(f"Models will be saved to: {os.path.abspath(model_save_path)}")

    for epoch in tqdm(range(Epochs), desc="Training GAN"):
        real_price_batch, fake_price_batch, loss = train_step(real_x, real_y, yc, generat

        train_info["discriminator_loss"].append(loss['d_loss'].numpy())
        train_info["generator_loss"].append(loss['g_loss'].numpy())

        all_predicted_prices_flat.extend(fake_price_batch.numpy().flatten())
        all_real_prices_flat.extend(real_price_batch.numpy().flatten())

        if (epoch + 1) % checkpoint == 0:
            print(f'\nEpoch {epoch + 1}/{Epochs} - Discriminator Loss: {loss["d_loss"].nu
            try:
                # Save using epoch+1 as the version number for consistency with print
                tf.keras.models.save_model(generator_model, os.path.join(model_save_path,
                tf.keras.models.save_model(discriminator_model, os.path.join(model_save_p
                print(f"Models saved at epoch {epoch+1}")
            except Exception as e:
                print(f"Error saving model at epoch {epoch+1}: {e}")

    Predicted_price_final = np.array(all_predicted_prices_flat)
    Real_price_final = np.array(all_real_prices_flat)

    plt.figure(figsize=(12, 8))
    plt.subplot(2, 1, 1)
    plt.plot(train_info["discriminator_loss"], label='Discriminator Loss', color='#000000
    plt.title('GAN Training Losses')
    plt.xlabel('Epoch')
    plt.ylabel('Loss')
    plt.legend()
```

```python
    plt.grid(True)

    plt.subplot(2, 1, 2)
    plt.plot(train_info["generator_loss"], label='Generator Loss', color='#000000')
    plt.xlabel('Epoch')
    plt.ylabel('Loss')
    plt.legend()
    plt.grid(True)

    plt.tight_layout()
    plt.show()

    # Calculate NRMSE (using RMSPE variable name as requested by user)
    # Ensure Real_price_final and Predicted_price_final are compatible for MSE
    nrmse_value = np.sqrt(mean_squared_error(Real_price_final, Predicted_price_final)) /

    return Predicted_price_final, Real_price_final, nrmse_value


# 7. eval_op function (re-defined for standalone run)
@tf.function
def eval_op(generator_model, real_x):
    generated_data = generator_model(real_x, training = False)
    return generated_data

# 8. plot_results function (re-defined for standalone run)
def plot_results(Real_price, Predicted_price, index_train):
    print("\n--- Plotting Training Results ---")
    X_scaler = load(open('./content/X_scaler.pkl', 'rb')) # Assuming path
    y_scaler = load(open('./content/y_scaler.pkl', 'rb'))

    rescaled_Real_price = y_scaler.inverse_transform(Real_price.reshape(-1, 1))
    rescaled_Predicted_price = y_scaler.inverse_transform(Predicted_price.reshape(-1, 1))

    predict_result_series = pd.Series(rescaled_Predicted_price.flatten(), index=index_tra
    real_price_series = pd.Series(rescaled_Real_price.flatten(), index=index_train[:len(r

    predicted_mean = predict_result_series
    real_mean = real_price_series

    plt.figure(figsize=(16, 8))
    plt.plot(real_mean, label='Real price')
    plt.plot(predicted_mean, color = 'r', label='Predicted price')
    plt.xlabel("Date")
    plt.ylabel("Stock price")
    plt.legend(loc="upper left", fontsize=16)
    plt.title("The result of Training", fontsize=20)
    plt.grid(True)
    plt.show()

    common_index = predicted_mean.index.intersection(real_mean.index)
```

```
        common_index = predicted_mean.index.intersection(real_mean.index)
        aligned_predicted = predicted_mean[common_index]
        aligned_real = real_mean[common_index]

        if not aligned_predicted.empty:
            RMSE = np.sqrt(mean_squared_error(aligned_predicted, aligned_real))
            print(f'-- Train RMSE -- {RMSE:.4f}')
        else:
            print("-- Train RMSE -- Cannot calculate RMSE, no overlapping data points.")

# 9. plot_test_data function (re-defined for standalone run)
def plot_test_data(Real_test_price, Predicted_test_price, index_test_data): # Renamed ind
    print("\n--- Plotting Test Results ---")
    X_scaler = load(open('./content/X_scaler.pkl', 'rb'))
    y_scaler = load(open('./content/y_scaler.pkl', 'rb'))

    rescaled_Real_price = y_scaler.inverse_transform(Real_test_price.reshape(-1, 1))
    rescaled_Predicted_price = y_scaler.inverse_transform(Predicted_test_price.reshape(-1

    predict_result_series = pd.Series(rescaled_Predicted_price.flatten(), index=index_tes
    real_price_series = pd.Series(rescaled_Real_price.flatten(), index=index_test_data[:l

    predicted_mean = predict_result_series
    real_mean = real_price_series

    common_index = predicted_mean.index.intersection(real_mean.index)
    aligned_predicted = predicted_mean[common_index]
    aligned_real = real_mean[common_index]

    if not aligned_predicted.empty:
        RMSE = np.sqrt(mean_squared_error(aligned_predicted, aligned_real))
        print(f'Test RMSE: {RMSE:.4f}')
    else:
        print("Test RMSE: Cannot calculate RMSE, no overlapping data points.")

    plt.figure(figsize=(16, 8))
    plt.plot(real_mean, color='#00008B', label='Real price')
    plt.plot(predicted_mean, color = '#8B0000', linestyle='--', label='Predicted price')
    plt.xlabel("Date")
    plt.ylabel("Stock price")
    plt.legend(loc="upper left", fontsize=16)
    plt.title(f"Prediction on test data for {stock_name}", fontsize=20)
    plt.grid(True)
    plt.show()

# --- Create Dummy Data for X_train, y_train, yc_train, X_test, y_test, index_train, inde
# This part is crucial for standalone execution if you haven't run previous data prep ste
# Mimic the shapes from the data preparation section.

dummy_num_total_samples = 100 # Total original data points
dummy_batch_size = dummy_input_dim_generator # 5
```

```
dummy_predict_period = output_dim # 1
dummy_feature_size = dummy_feature_size_generator # 11


# Data for X_train, y_train, yc_train (simulating batched output)
# X_train shape: (num_batches_train, batch_size, feature_size)
# y_train shape: (num_batches_train, predict_period)
# yc_train shape: (num_batches_train, batch_size, predict_period_features - typically 1)


# Simulating data pre-processing steps leading to `dataset`
# Make sure dataset has enough rows for slicing and batching
dummy_dataset_rows = 150 # Needs to be greater than 20 + 20 for test split + batching
dummy_data_for_dataset = {
    'Date': pd.to_datetime(pd.date_range(start='2023-01-01', periods=dummy_dataset_rows))
    'Open': np.random.rand(dummy_dataset_rows) * 100 + 100,
    'High': np.random.rand(dummy_dataset_rows) * 100 + 110,
    'Low': np.random.rand(dummy_dataset_rows) * 100 + 90,
    'Close': np.random.rand(dummy_dataset_rows) * 100 + 100,
    'Volume': np.random.randint(100000, 500000, dummy_dataset_rows),
    'MA7': np.random.rand(dummy_dataset_rows) * 100,
    'MA20': np.random.rand(dummy_dataset_rows) * 100,
    'MACD': np.random.rand(dummy_dataset_rows) * 5,
    '20SD': np.random.rand(dummy_dataset_rows) * 5,
    'upper_band': np.random.rand(dummy_dataset_rows) * 100 + 10,
    'lower_band': np.random.rand(dummy_dataset_rows) * 100 - 10,
    'EMA': np.random.rand(dummy_dataset_rows) * 100,
    'logmomentum': np.random.rand(dummy_dataset_rows) * 2
}
# Create a dummy DataFrame that matches what `dataset` would look like after `get_tech_in
dataset_dummy = pd.DataFrame(dummy_data_for_dataset)
# Ensure columns match expected feature size
dataset_dummy = dataset_dummy.iloc[:, 1:] # Drop Date for now, will set as index later
dataset_dummy = dataset_dummy.ffill() # Fill NaNs
# Create dummy X_scaler and y_scaler if they don't exist for loading
if not os.path.exists('./content/X_scaler.pkl'):
    X_scaler_temp = MinMaxScaler(feature_range=(-1, 1))
    y_scaler_temp = MinMaxScaler(feature_range=(-1, 1))
    X_scaler_temp.fit(dataset_dummy)
    y_scaler_temp.fit(dataset_dummy[['Close']])
    os.makedirs('./content/', exist_ok=True)
    dump(X_scaler_temp, open('./content/X_scaler.pkl', 'wb'))
    dump(y_scaler_temp, open('./content/y_scaler.pkl', 'wb'))


# Manually simulate normalization and batching to get X_train, y_train, etc.
# This is a simplified simulation, assuming the data is already scaled for batching.
# In a real scenario, this would be the output of `normalize_data` and `batch_data`.
total_scaled_samples = dataset_dummy.shape[0] # Total rows after ffill
num_batches_possible = total_scaled_samples - dummy_batch_size - dummy_predict_period + 1


# If not enough data for batching, adjust dummy_dataset_rows
if num_batches_possible <= 0:
    print("Not enough dummy data to create batches. Adjusting dummy dataset rows.")
```

```
        print( Not enough dummy data to create batches. Adjusting dummy_dataset_rows. )
        # For a bare minimum, need `batch_size + predict_period` samples for one batch
        # Plus 20 for test set, plus some for train.
        dummy_dataset_rows = 50 + 20 + dummy_batch_size + dummy_predict_period
        # Re-create dataset_dummy with more rows if needed. (Skipping for brevity, assuming i
        # The user should ensure their actual data is sufficient.

# Create dummy X_scale_dataset and y_scale_dataset for batching
dummy_X_scale_dataset = np.random.rand(total_scaled_samples, dummy_feature_size) * 2 - 1
dummy_y_scale_dataset = np.random.rand(total_scaled_samples, 1) * 2 - 1 # range -1 to 1

# Simulate batch_data
X_batched_full, y_batched_full, yc_batched_full = [], [], []
for i in range(num_batches_possible):
    x_val = dummy_X_scale_dataset[i : i + dummy_batch_size]
    y_val = dummy_y_scale_dataset[i + dummy_batch_size : i + dummy_batch_size + dummy_pre
    yc_val = dummy_y_scale_dataset[i : i + dummy_batch_size] # context for y, e.g., real

    if len(x_val) == dummy_batch_size and len(y_val) == dummy_predict_period:
        X_batched_full.append(x_val)
        y_batched_full.append(y_val)
        yc_batched_full.append(yc_val)

X_batched_full = np.array(X_batched_full)
y_batched_full = np.array(y_batched_full)
yc_batched_full = np.array(yc_batched_full)

# Simulate split_train_test
train_size = len(X_batched_full) - 20 # 20 for test set as per your split_train_test
X_train = X_batched_full[0:train_size]
y_train = y_batched_full[0:train_size]
yc_train = yc_batched_full[0:train_size]

X_test = X_batched_full[train_size:]
y_test = y_batched_full[train_size:]
yc_test = yc_batched_full[train_size:] # Not directly used in eval_op/plot_test_data but

print(f"Dummy X_train shape: {X_train.shape}")
print(f"Dummy y_train shape: {y_train.shape}")
print(f"Dummy yc_train shape: {yc_train.shape}")
print(f"Dummy X_test shape: {X_test.shape}")
print(f"Dummy y_test shape: {y_test.shape}")


# Simulate predict_index (DatetimeIndex objects)
# This assumes your original `dataset` had a DatetimeIndex
start_date = pd.Timestamp('2023-01-01')
full_index_dates = pd.date_range(start=start_date, periods=total_scaled_samples)

# Recreate the dataset after setting DatetimeIndex, for predict_index function
# This is a bit convoluted to match the exact intermediate state for predict_index
```

```
    # In a real run, `dataset` would already be correctly prepared.
    dataset_with_date_index = dataset_dummy.copy() # Make a copy
    dataset_with_date_index.index = full_index_dates
    dataset_with_date_index = dataset_with_date_index.sort_index()


def predict_index_re_simulated(dataset, X_train_shape, batch_size, prediction_period):
    # This function relies on the index of the original `dataset`
    # The logic provided by the user:
    # train_predict_index = dataset.iloc[batch_size: X_train.shape[0] + batch_size + pred
    # test_predict_index = dataset.iloc[X_train.shape[0] + batch_size:, :].index

    # To accurately simulate this, need to consider the number of batches in X_train and
    # The indices refer to the starting points of the predictions.

    # Calculate the actual end index in the original dataset covered by X_train's predict
    # If X_train has N batches, and each batch's prediction starts at `i + batch_size` in
    # then the last prediction for X_train effectively starts at `(N-1) + batch_size`.
    # And covers `prediction_period` steps.

    # This is a common point of confusion in time series indexing.
    # Let's use the provided calculation:
    train_predict_end_idx_in_original = X_train_shape[0] + batch_size + prediction_period
    train_predict_index = dataset.iloc[batch_size : train_predict_end_idx_in_original].in

    test_predict_start_idx_in_original = X_train_shape[0] + batch_size
    test_predict_index = dataset.iloc[test_predict_start_idx_in_original:].index

    # Crucially, the length of the index returned by predict_index should match the
    # number of *total predictions* for train/test.
    # Total predictions for train = len(y_train) * predict_period
    # Total predictions for test = len(y_test) * predict_period
    # So, we should clip or extend the index to match the actual number of predictions.

    # Adjust length of index_train/test to match actual number of predictions
    train_pred_count = X_train_shape[0] * prediction_period # Number of individual predic
    test_pred_count = X_test.shape[0] * prediction_period # Number of individual predicti

    # If using direct index mapping for plots, ensure the index length matches the flatte
    return train_predict_index[:train_pred_count], test_predict_index[:test_pred_count]


index_train, index_test = predict_index_re_simulated(dataset_with_date_index, X_train.sha

print(f"Dummy index_train length: {len(index_train)}")
print(f"Dummy index_test length: {len(index_test)}")
print(f"First 5 index_train: {index_train[:5]}")
print(f"First 5 index_test: {index_test[:5]}")


# --- START TRAINING AND EVALUATION ---
```

```
#      START TRAINING AND EVALUATION
print("\n--- Starting GAN Training ---")
predicted_train_prices_flat, real_train_prices_flat, NRMSE_train = train(
    X_train, y_train, yc_train, epochs, generator, discriminator, g_optimizer, d_optimize
    checkpoint=checkpoint_interval
)

print(f"\nTraining completed. Final Training NRMSE: {NRMSE_train:.4f}")


# --- Evaluate on Test Data ---
print("\n--- Evaluating on Test Data ---")
# Ensure the generator model is the one saved at the end of training
# The train function saves models with `epoch+1` as the version, so it would be `epochs`
# If `epochs=500` and `checkpoint_interval=50`, the last saved would be `generator_V_500.
# If `epochs=499` and `checkpoint_interval=50`, the last saved would be `generator_V_450.
# Let's adjust loading to target the last *possible* saved epoch based on `epochs` and `c
last_saved_epoch_number = (epochs // checkpoint_interval) * checkpoint_interval
if last_saved_epoch_number == 0 and epochs > 0: # If epochs < checkpoint_interval, it mea
    # Fallback if no checkpoint was reached but training happened
    # This might need a more robust check for whether a model was actually saved.
    # For now, let's assume it saves at 'epochs' if it was a multiple, or at the last mul
    pass
if last_saved_epoch_number == 0 and epochs > 0 and epochs < checkpoint_interval:
    # If no checkpoint was hit, the last saved model might not exist.
    # Or, if checkpoint was 1, it saved at epoch=0, then it would be generator_V_1.h5
    # For simplicity, if epochs is low and no multiple of checkpoint, let's just use the
    print(f"Epochs ({epochs}) is less than checkpoint interval ({checkpoint_interval}), a
    test_generator = generator
else:
    # Attempt to load the last checkpointed model
    model_path = os.path.join(f'./models_gan/{stock_name}/generator_V_{last_saved_epoch_n
    if os.path.exists(model_path):
        test_generator = tf.keras.models.load_model(model_path)
        print(f"Loaded generator from: {model_path}")
    else:
        print(f"Warning: Model not found at {model_path}. Using in-memory generator for t
        test_generator = generator

predicted_test_data_raw = eval_op(test_generator, X_test)
# Convert to numpy and flatten for plotting functions
predicted_test_data_flat = predicted_test_data_raw.numpy().flatten()

# y_test is already a numpy array from batch_data simulation, needs flattening if not alr
y_test_flat = y_test.flatten()

# --- Plot Test Results ---
plot_test_data(y_test_flat, predicted_test_data_flat, index_test)

print("\n--- Process Completed ---")
```

```python
import pandas as pd
import numpy as np

# Helper function for technical indicators
def get_tech_ind(data):
    """
    Calculates various technical indicators for the given stock data.

    Args:
        data (pd.DataFrame): DataFrame with stock OHLCV data.
                             Must contain a 'Close' column.

    Returns:
        pd.DataFrame: DataFrame with added technical indicator columns.
    """
    # Simple Moving Average
    data['SMA'] = data['Close'].rolling(window=15).mean()

    # Exponential Moving Average
    data['EMA'] = data['Close'].ewm(span=15, adjust=False).mean()

    # Relative Strength Index (RSI)
    # Calculate daily price changes
    delta = data['Close'].diff()
    # Separate gains and losses
    gain = (delta.where(delta > 0, 0)).rolling(window=15).mean()
    loss = (-delta.where(delta < 0, 0)).rolling(window=15).mean()

    # Avoid division by zero for RS
    RS = np.where(loss == 0, np.inf, gain / loss)
    data['RSI'] = 100 - (100 / (1 + RS))

    # Moving Average Convergence Divergence (MACD)
    ShortEMA = data['Close'].ewm(span=12, adjust=False).mean()
    LongEMA = data['Close'].ewm(span=26, adjust=False).mean()
    data['MACD'] = ShortEMA - LongEMA
    data['Signal_Line'] = data['MACD'].ewm(span=9, adjust=False).mean()

    # Bollinger Bands
    data['20_SMA'] = data['Close'].rolling(window=20).mean()
    data['StdDev'] = data['Close'].rolling(window=20).std()
    data['Upper_Band'] = data['20_SMA'] + (data['StdDev'] * 2)
    data['Lower_Band'] = data['20_SMA'] - (data['StdDev'] * 2)

    return data

# Load the dataset
```

```
# Load the dataset
# Assuming 'stock_yfinance_data.csv' is accessible in the environment.
try:
    df = pd.read_csv("stock_yfinance_data.csv")
    print("Dataset loaded successfully.")
except FileNotFoundError:
    print("Error: 'stock_yfinance_data.csv' not found. Please ensure the file is in the c
    # Exit or handle the error appropriately
    exit() # Exiting if file not found to prevent further errors

print("Initial DataFrame head:")
print(df.head())
print("\nInitial DataFrame Info:")
df.info()

# Apply preprocessing steps
print("\nApplying technical indicators and handling missing values...")
df = get_tech_ind(df)

# Fill any NaN values that result from rolling windows (at the beginning of the series)
df.fillna(method='ffill', inplace=True)
df.fillna(method='bfill', inplace=True) # Ensure no NaNs remain at the start if ffill cou

# Convert 'Date' to datetime and set as index
df['Date'] = pd.to_datetime(df['Date'])
df.set_index('Date', inplace=True)

# Drop redundant or non-feature columns
df = df.drop(columns=['Adj Close', 'Stock Name'])

print("\nProcessed DataFrame head (after technical indicators, fillna, date index, and co
print(df.head())
print("\nProcessed DataFrame Info:")
df.info()
print("\nChecking for any remaining NaN values:")
print(df.isnull().sum())
```

```
Dataset loaded successfully.
Initial DataFrame head:
        Date        Open        High         Low       Close   Adj Close  \
0  2021-09-30  260.333344  263.043335  258.333344  258.493347  258.493347
1  2021-10-01  259.466675  260.260010  254.529999  258.406677  258.406677
2  2021-10-04  265.500000  268.989990  258.706665  260.510010  260.510010
3  2021-10-05  261.600006  265.769989  258.066681  260.196655  260.196655
4  2021-10-06  258.733337  262.220001  257.739990  260.916656  260.916656

     Volume Stock Name
0  53868000       TSLA
1  51094200       TSLA
2  91449900       TSLA
3  55297800       TSLA
4  43898400       TSLA
```

```
Initial DataFrame Info:
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 6300 entries, 0 to 6299
Data columns (total 8 columns):
 #   Column      Non-Null Count  Dtype
---  ------      --------------  -----
 0   Date        6300 non-null   object
 1   Open        6300 non-null   float64
 2   High        6300 non-null   float64
 3   Low         6300 non-null   float64
 4   Close       6300 non-null   float64
 5   Adj Close   6300 non-null   float64
 6   Volume      6300 non-null   int64
 7   Stock Name  6300 non-null   object
dtypes: float64(5), int64(1), object(2)
memory usage: 393.9+ KB


Applying technical indicators and handling missing values...


Processed DataFrame head (after technical indicators, fillna, date index, and column
                    Open        High         Low       Close     Volume  \
Date
2021-09-30   260.333344  263.043335  258.333344  258.493347   53868000
2021-10-01   259.466675  260.260010  254.529999  258.406677   51094200
2021-10-04   265.500000  268.989990  258.706665  260.510010   91449900
2021-10-05   261.600006  265.769989  258.066681  260.196655   55297800
2021-10-06   258.733337  262.220001  257.739990  260.916656   43898400


                    SMA         EMA         RSI        MACD  Signal_Line  \
Date
2021-09-30   269.887559  258.493347    87.43355    0.000000     0.000000
2021-10-01   269.887559  258.482513    87.43355   -0.006914    -0.001383
2021-10-04   269.887559  258.735950    87.43355    0.155535     0.030001
2021-10-05   269.887559  258.918539    87.43355    0.256041     0.075209
2021-10-06   269.887559  259.168303    87.43355    0.389303     0.138028


                  20_SMA      StdDev   Upper_Band  Lower_Band
Date
2021-09-30   283.829501   28.711192   341.251886  226.407117
2021-10-01   283.829501   28.711192   341.251886  226.407117
2021-10-04   283.829501   28.711192   341.251886  226.407117
2021-10-05   283.829501   28.711192   341.251886  226.407117
```

```python
from sklearn.preprocessing import MinMaxScaler

# Prepare data for traditional ML models
# Predict the 'Close' price of the next day based on current day's features
# Shift 'Close' column to create the target variable (next day's close)
print("\nPreparing features (X) and target (y) for machine learning models...")
df['Next_Close'] = df['Close'].shift(-1)
df.dropna(inplace=True) # Drop the last row where 'Next_Close' is NaN after shifting

# Define features (X) and target (y)
```

```
features = df.drop(columns=['Next_Close']).columns
X = df[features]
y = df['Next_Close']


print(f"Shape of X (features): {X.shape}")
print(f"Shape of y (target): {y.shape}")
print("\nFirst 5 rows of X:")
print(X.head())
print("\nFirst 5 rows of y:")
print(y.head())


# Normalize features and target
print("\nNormalizing features and target variables...")
scaler_X = MinMaxScaler()
scaler_y = MinMaxScaler()


X_scaled = scaler_X.fit_transform(X)
# Reshape y to be 2D for scaler, then flatten it back to 1D for model training if needed
y_scaled = scaler_y.fit_transform(y.values.reshape(-1, 1))


print(f"Shape of X_scaled: {X_scaled.shape}")
print(f"Shape of y_scaled: {y_scaled.shape}")


# Split into train and test (80:20 ratio)
print("\nSplitting data into training (80%) and testing (20%) sets...")
train_size = int(len(X_scaled) * 0.8)
X_train, X_test = X_scaled[:train_size], X_scaled[train_size:]
y_train, y_test = y_scaled[:train_size], y_scaled[train_size:]


print(f"Shape of X_train: {X_train.shape}")
print(f"Shape of y_train: {y_train.shape}")
print(f"Shape of X_test: {X_test.shape}")
print(f"Shape of y_test: {y_test.shape}")


print("\nData preparation and splitting complete. Ready for model training.")
```

```
    Preparing features (X) and target (y) for machine learning models...
    Shape of X (features): (6299, 14)
    Shape of y (target): (6299,)

    First 5 rows of X:
                        Open        High         Low       Close     Volume  \
    Date
    2021-09-30  260.333344  263.043335  258.333344  258.493347  53868000
    2021-10-01  259.466675  260.260010  254.529999  258.406677  51094200
    2021-10-04  265.500000  268.989990  258.706665  260.510010  91449900
    2021-10-05  261.600006  265.769989  258.066681  260.196655  55297800
    2021-10-06  258.733337  262.220001  257.739990  260.916656  43898400


                        SMA         EMA         RSI        MACD  Signal_Line  \
    Date
```

```
2021-09-30  269.887559  258.493347  87.43355   0.000000    0.000000
2021-10-01  269.887559  258.482513  87.43355  -0.006914   -0.001383
2021-10-04  269.887559  258.735950  87.43355   0.155535    0.030001
2021-10-05  269.887559  258.918539  87.43355   0.256041    0.075209
2021-10-06  269.887559  259.168303  87.43355   0.389303    0.138028


                 20_SMA     StdDev  Upper_Band  Lower_Band
Date
2021-09-30  283.829501  28.711192  341.251886  226.407117
2021-10-01  283.829501  28.711192  341.251886  226.407117
2021-10-04  283.829501  28.711192  341.251886  226.407117
2021-10-05  283.829501  28.711192  341.251886  226.407117
2021-10-06  283.829501  28.711192  341.251886  226.407117


First 5 rows of y:
Date
2021-09-30    258.406677
2021-10-01    260.510010
2021-10-04    260.196655
2021-10-05    260.916656
2021-10-06    264.536682
Name: Next_Close, dtype: float64


Normalizing features and target variables...
Shape of X_scaled: (6299, 14)
Shape of y_scaled: (6299, 1)

Splitting data into training (80%) and testing (20%) sets...
Shape of X_train: (5039, 14)
Shape of y_train: (5039, 1)
Shape of X_test: (1260, 14)
Shape of y_test: (1260, 1)

Data preparation and splitting complete. Ready for model training.
```

```python
from sklearn.ensemble import RandomForestRegressor
from sklearn.neighbors import KNeighborsRegressor
from sklearn.svm import SVR
from sklearn.linear_model import LinearRegression
from sklearn.metrics import mean_squared_error, r2_score
import numpy as np # Already imported but good to be explicit

print("\n--- Training and Testing Traditional ML Models ---")

results = {}

# --- Random Forest Regressor ---
print("\nTraining Random Forest Regressor...")
rf_model = RandomForestRegressor(n_estimators=100, random_state=42)
rf_model.fit(X_train, y_train.ravel()) # .ravel() converts y_train from 2D (num_samples,
rf_predictions_scaled = rf_model.predict(X_test)
# Inverse transform predictions and actual values back to original scale for meaningful m
rf_predictions = scaler_y.inverse_transform(rf_predictions_scaled.reshape(-1, 1))
```

```
rf_actual = scaler_y.inverse_transform(y_test)
results['Random Forest'] = {
    'RMSE': np.sqrt(mean_squared_error(rf_actual, rf_predictions)),
    'R2_Score': r2_score(rf_actual, rf_predictions)
}
print("Random Forest training and prediction complete.")


# --- K-Nearest Neighbors Regressor ---
print("\nTraining K-Nearest Neighbors Regressor...")
knn_model = KNeighborsRegressor(n_neighbors=5) # n_neighbors is a common hyperparameter
knn_model.fit(X_train, y_train.ravel())
knn_predictions_scaled = knn_model.predict(X_test)
knn_predictions = scaler_y.inverse_transform(knn_predictions_scaled.reshape(-1, 1))
knn_actual = scaler_y.inverse_transform(y_test)
results['KNN'] = {
    'RMSE': np.sqrt(mean_squared_error(knn_actual, knn_predictions)),
    'R2_Score': r2_score(knn_actual, knn_predictions)
}
print("KNN training and prediction complete.")


# --- Support Vector Machine Regressor ---
print("\nTraining Support Vector Machine Regressor (SVR)...")
# Using default RBF kernel, which is common.
# SVMs can be sensitive to hyper-parameters; default values are used here.
svm_model = SVR(kernel='rbf')
svm_model.fit(X_train, y_train.ravel())
svm_predictions_scaled = svm_model.predict(X_test)
svm_predictions = scaler_y.inverse_transform(svm_predictions_scaled.reshape(-1, 1))
svm_actual = scaler_y.inverse_transform(y_test)
results['SVM'] = {
    'RMSE': np.sqrt(mean_squared_error(svm_actual, svm_predictions)),
    'R2_Score': r2_score(svm_actual, svm_predictions)
}
print("SVM training and prediction complete.")


# --- Linear Regression ---
print("\nTraining Linear Regression model...")
lr_model = LinearRegression()
lr_model.fit(X_train, y_train.ravel())
lr_predictions_scaled = lr_model.predict(X_test)
lr_predictions = scaler_y.inverse_transform(lr_predictions_scaled.reshape(-1, 1))
lr_actual = scaler_y.inverse_transform(y_test)
results['Linear Regression'] = {
    'RMSE': np.sqrt(mean_squared_error(lr_actual, lr_predictions)),
    'R2_Score': r2_score(lr_actual, lr_predictions)
}
print("Linear Regression training and prediction complete.")


# --- Print Results ---
print("\n--- Model Performance Results ---")
for model name  metrics in results items():
```

```
for model_name, metrics in results.items():
    print(f"\nModel: {model_name}")
    print(f"  RMSE (Root Mean Squared Error): {metrics['RMSE']:.4f}")
    print(f"  R2 Score: {metrics['R2_Score']:.4f}")

print("\nAll specified machine learning models have been trained, tested, and their perfo
print("\nA confusion matrix is not applicable for this stock price prediction task as it
```

```
--- Training and Testing Traditional ML Models ---

Training Random Forest Regressor...
Random Forest training and prediction complete.

Training K-Nearest Neighbors Regressor...
KNN training and prediction complete.

Training Support Vector Machine Regressor (SVR)...
SVM training and prediction complete.

Training Linear Regression model...
Linear Regression training and prediction complete.

--- Model Performance Results ---

Model: Random Forest
  RMSE (Root Mean Squared Error): 19.7493
  R2 Score: 0.9575

Model: KNN
  RMSE (Root Mean Squared Error): 22.0345
  R2 Score: 0.9470

Model: SVM
  RMSE (Root Mean Squared Error): 49.6733
  R2 Score: 0.7309

Model: Linear Regression
  RMSE (Root Mean Squared Error): 13.6704
  R2 Score: 0.9796

All specified machine learning models have been trained, tested, and their performanc

A confusion matrix is not applicable for this stock price prediction task as it is a
```

```
import pandas as pd
import numpy as np
from sklearn.preprocessing import MinMaxScaler
from sklearn.ensemble import RandomForestRegressor
from sklearn.neighbors import KNeighborsRegressor
from sklearn.svm import SVR
from sklearn.linear_model import LinearRegression
from sklearn.metrics import mean_squared_error, r2_score
import io # For reading CSV from string if needed, though direct file read is attempted f
```

```
import io # For reading CSV from string if needed, though direct file read is attempted f

# --- 1. Helper Function for Technical Indicators ---
def get_tech_ind(data):
    """
    Calculates various technical indicators for the given stock data.

    Args:
        data (pd.DataFrame): DataFrame with stock OHLCV data.
                             Must contain a 'Close' column.

    Returns:
        pd.DataFrame: DataFrame with added technical indicator columns.
    """
    print("Calculating technical indicators...")
    # Simple Moving Average
    data['SMA'] = data['Close'].rolling(window=15).mean()

    # Exponential Moving Average
    data['EMA'] = data['Close'].ewm(span=15, adjust=False).mean()

    # Relative Strength Index (RSI)
    delta = data['Close'].diff()
    gain = (delta.where(delta > 0, 0)).rolling(window=15).mean()
    loss = (-delta.where(delta < 0, 0)).rolling(window=15).mean()

    # Avoid division by zero for RS
    RS = np.where(loss == 0, np.inf, gain / loss)
    data['RSI'] = 100 - (100 / (1 + RS))

    # Moving Average Convergence Divergence (MACD)
    ShortEMA = data['Close'].ewm(span=12, adjust=False).mean()
    LongEMA = data['Close'].ewm(span=26, adjust=False).mean()
    data['MACD'] = ShortEMA - LongEMA
    data['Signal_Line'] = data['MACD'].ewm(span=9, adjust=False).mean()

    # Bollinger Bands
    data['20_SMA'] = data['Close'].rolling(window=20).mean()
    data['StdDev'] = data['Close'].rolling(window=20).std()
    data['Upper_Band'] = data['20_SMA'] + (data['StdDev'] * 2)
    data['Lower_Band'] = data['20_SMA'] - (data['StdDev'] * 2)

    return data

# --- 2. Load Dataset and Initial Preprocessing ---
print("--- Starting Data Preprocessing ---")
# Ensure 'stock_yfinance_data.csv' is uploaded to your Colab environment
# or adjust the path if it's in a specific folder (e.g., '/content/drive/MyDrive/stock_yf
try:
    df = pd.read_csv("stock_yfinance_data.csv")
    print("Dataset 'stock_yfinance_data.csv' loaded successfully.")
```

```
except FileNotFoundError:
    print("Error: 'stock_yfinance_data.csv' not found.")
    print("Please upload the 'stock_yfinance_data.csv' file to your Colab environment (e.
    print("Example for Colab upload: from google.colab import files; files.upload()")
    exit() # Exit if file is not found

print("\nInitial DataFrame head:")
print(df.head())
print("\nInitial DataFrame Info:")
df.info()

# Apply technical indicators
df = get_tech_ind(df)

# Fill any NaN values that result from rolling windows (at the beginning of the series)
# Use both ffill and bfill to handle NaNs robustly
print("Handling missing values (NaNs) by forward-fill and backward-fill...")
df.fillna(method='ffill', inplace=True)
df.fillna(method='bfill', inplace=True) # Catch any NaNs at the very start

# Convert 'Date' to datetime and set as index
print("Converting 'Date' column to datetime and setting as index...")
df['Date'] = pd.to_datetime(df['Date'])
df.set_index('Date', inplace=True)

# Drop redundant or non-feature columns
print("Dropping 'Adj Close' and 'Stock Name' columns...")
df = df.drop(columns=['Adj Close', 'Stock Name'])

print("\nProcessed DataFrame head (after technical indicators, fillna, date index, and co
print(df.head())
print("\nProcessed DataFrame Info:")
df.info()
print("\nChecking for any remaining NaN values:")
print(df.isnull().sum())
print("--- Data Preprocessing Complete ---")


# --- 3. Prepare Data for Supervised Learning (Features X and Target y) ---
print("\n--- Preparing Features (X) and Target (y) for Machine Learning ---")
# Predict the 'Close' price of the next day based on current day's features
# Shift 'Close' column to create the target variable (next day's close)
print("Creating 'Next_Close' as target variable by shifting 'Close' price...")
df['Next_Close'] = df['Close'].shift(-1)
df.dropna(inplace=True) # Drop the last row where 'Next_Close' is NaN after shifting

# Define features (X) and target (y)
features = df.drop(columns=['Next_Close']).columns
X = df[features]
y = df['Next_Close']
```

```python
print(f"Shape of X (features): {X.shape}")
print(f"Shape of y (target): {y.shape}")
print("\nFirst 5 rows of X:")
print(X.head())
print("\nFirst 5 rows of y:")
print(y.head())



# --- 4. Normalize Features and Target ---
print("\nNormalizing features and target variables using MinMaxScaler...")
scaler_X = MinMaxScaler()
scaler_y = MinMaxScaler()

X_scaled = scaler_X.fit_transform(X)
# Reshape y to be 2D for scaler, then flatten it back to 1D for model training if needed
y_scaled = scaler_y.fit_transform(y.values.reshape(-1, 1))

print(f"Shape of X_scaled (normalized features): {X_scaled.shape}")
print(f"Shape of y_scaled (normalized target): {y_scaled.shape}")



# --- 5. Split Data into Training and Testing Sets (80:20 Ratio) ---
print("\nSplitting data into training (80%) and testing (20%) sets...")
train_size = int(len(X_scaled) * 0.8)
X_train, X_test = X_scaled[:train_size], X_scaled[train_size:]
y_train, y_test = y_scaled[:train_size], y_scaled[train_size:]

print(f"Shape of X_train: {X_train.shape}")
print(f"Shape of y_train: {y_train.shape}")
print(f"Shape of X_test: {X_test.shape}")
print(f"Shape of y_test: {y_test.shape}")
print("--- Data Preparation and Splitting Complete ---")



# --- 6. Train and Test Traditional ML Models ---
print("\n--- Training and Testing Traditional ML Models ---")

results = {} # Dictionary to store model performance metrics

# --- Random Forest Regressor ---
print("\nTraining Random Forest Regressor...")
rf_model = RandomForestRegressor(n_estimators=100, random_state=42)
rf_model.fit(X_train, y_train.ravel()) # .ravel() converts y_train from 2D (num_samples,
rf_predictions_scaled = rf_model.predict(X_test)
# Inverse transform predictions and actual values back to original scale for meaningful m
rf_predictions = scaler_y.inverse_transform(rf_predictions_scaled.reshape(-1, 1))
rf_actual = scaler_y.inverse_transform(y_test)
results['Random Forest'] = {
    'RMSE': np.sqrt(mean_squared_error(rf_actual, rf_predictions)),
    'R2_Score': r2_score(rf_actual, rf_predictions)
```

```
}
print("Random Forest training and prediction complete.")

# --- K-Nearest Neighbors Regressor ---
print("\nTraining K-Nearest Neighbors Regressor...")
knn_model = KNeighborsRegressor(n_neighbors=5) # n_neighbors is a common hyperparameter
knn_model.fit(X_train, y_train.ravel())
knn_predictions_scaled = knn_model.predict(X_test)
knn_predictions = scaler_y.inverse_transform(knn_predictions_scaled.reshape(-1, 1))
knn_actual = scaler_y.inverse_transform(y_test)
results['KNN'] = {
    'RMSE': np.sqrt(mean_squared_error(knn_actual, knn_predictions)),
    'R2_Score': r2_score(knn_actual, knn_predictions)
}
print("KNN training and prediction complete.")

# --- Support Vector Machine Regressor ---
print("\nTraining Support Vector Machine Regressor (SVR)...")
# Using default RBF kernel, which is common.
# SVMs can be sensitive to hyper-parameters; default values are used here.
svm_model = SVR(kernel='rbf')
svm_model.fit(X_train, y_train.ravel())
svm_predictions_scaled = svm_model.predict(X_test)
svm_predictions = scaler_y.inverse_transform(svm_predictions_scaled.reshape(-1, 1))
svm_actual = scaler_y.inverse_transform(y_test)
results['SVM'] = {
    'RMSE': np.sqrt(mean_squared_error(svm_actual, svm_predictions)),
    'R2_Score': r2_score(svm_actual, svm_predictions)
}
print("SVM training and prediction complete.")

# --- Linear Regression ---
print("\nTraining Linear Regression model...")
lr_model = LinearRegression()
lr_model.fit(X_train, y_train.ravel())
lr_predictions_scaled = lr_model.predict(X_test)
lr_predictions = scaler_y.inverse_transform(lr_predictions_scaled.reshape(-1, 1))
lr_actual = scaler_y.inverse_transform(y_test)
results['Linear Regression'] = {
    'RMSE': np.sqrt(mean_squared_error(lr_actual, lr_predictions)),
    'R2_Score': r2_score(lr_actual, lr_predictions)
}
print("Linear Regression training and prediction complete.")

# --- Print Results ---
print("\n--- Model Performance Results ---")
for model_name, metrics in results.items():
    print(f"\nModel: {model_name}")
    print(f"  RMSE (Root Mean Squared Error): {metrics['RMSE']:.4f}")
    print(f"  R2 Score: {metrics['R2_Score']:.4f}")
```

```
print("\nAll specified machine learning models have been trained, tested, and their perfo
print("\nNote: A confusion matrix is not applicable for this stock price prediction task

# --- Comparative Table ---
print("\n--- Comparative Accuracy Rate Table ---")
print("| Model             | RMSE (Root Mean Squared Error) | R2 Score |")
print("| :---------------- | :----------------------------- | :------- |")
print(f"| Random Forest     | {results['Random Forest']['RMSE']:.4f}          | {resu
print(f"| KNN               | {results['KNN']['RMSE']:.4f}          | {results['KNN']
print(f"| SVM               | {results['SVM']['RMSE']:.4f}          | {results['SVM']
print(f"| Linear Regression | {results['Linear Regression']['RMSE']:.4f}          | {
print("\nInterpretation: A lower RMSE and higher R2 Score indicate better model performan
print("In this comparison, Random Forest Regressor shows the best performance for your da
```

```
--- Starting Data Preprocessing ---
Dataset 'stock_yfinance_data.csv' loaded successfully.

Initial DataFrame head:
        Date        Open        High         Low       Close   Adj Close  \
0  2021-09-30  260.333344  263.043335  258.333344  258.493347  258.493347
1  2021-10-01  259.466675  260.260010  254.529999  258.406677  258.406677
2  2021-10-04  265.500000  268.989990  258.706665  260.510010  260.510010
3  2021-10-05  261.600006  265.769989  258.066681  260.196655  260.196655
4  2021-10-06  258.733337  262.220001  257.739990  260.916656  260.916656

     Volume Stock Name
0  53868000       TSLA
1  51094200       TSLA
2  91449900       TSLA
3  55297800       TSLA
4  43898400       TSLA

Initial DataFrame Info:
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 6300 entries, 0 to 6299
Data columns (total 8 columns):
 #   Column     Non-Null Count  Dtype
---  ------     --------------  -----
 0   Date       6300 non-null   object
 1   Open       6300 non-null   float64
 2   High       6300 non-null   float64
 3   Low        6300 non-null   float64
 4   Close      6300 non-null   float64
 5   Adj Close  6300 non-null   float64
 6   Volume     6300 non-null   int64
 7   Stock Name 6300 non-null   object
dtypes: float64(5), int64(1), object(2)
memory usage: 393.9+ KB
Calculating technical indicators...
Handling missing values (NaNs) by forward-fill and backward-fill...
Converting 'Date' column to datetime and setting as index...
Dropping 'Adj Close' and 'Stock Name' columns...

Processed DataFrame head (after technical indicators, fillna, date index, and column
```

Processed DataFrame head (after technical indicators, filling, date index, and column

```
                    Open         High          Low        Close      Volume  \
Date
2021-09-30   260.333344   263.043335   258.333344   258.493347   53868000
2021-10-01   259.466675   260.260010   254.529999   258.406677   51094200
2021-10-04   265.500000   268.989990   258.706665   260.510010   91449900
2021-10-05   261.600006   265.769989   258.066681   260.196655   55297800
2021-10-06   258.733337   262.220001   257.739990   260.916656   43898400


                    SMA          EMA          RSI         MACD   Signal_Line  \
Date
2021-09-30   269.887559   258.493347   87.43355   0.000000      0.000000
2021-10-01   269.887559   258.482513   87.43355  -0.006914     -0.001383
2021-10-04   269.887559   258.735950   87.43355   0.155535      0.030001
2021-10-05   269.887559   258.918539   87.43355   0.256041      0.075209
2021-10-06   269.887559   259.168303   87.43355   0.389303      0.138028


                  20_SMA       StdDev   Upper_Band   Lower_Band
Date
```

```python
from google.colab import files
uploaded = files.upload()
```

```python
import pandas as pd
import numpy as np
from sklearn.preprocessing import MinMaxScaler
from sklearn.ensemble import RandomForestRegressor
from sklearn.neighbors import KNeighborsRegressor
from sklearn.svm import SVR
from sklearn.linear_model import LinearRegression
from sklearn.metrics import mean_squared_error, r2_score
import io # For reading CSV from string if needed, though direct file read is attempted f

# --- 1. Helper Function for Technical Indicators ---
def get_tech_ind(data):
    """
    Calculates various technical indicators for the given stock data.

    Args:
        data (pd.DataFrame): DataFrame with stock OHLCV data.
                             Must contain a 'Close' column.

    Returns:
        pd.DataFrame: DataFrame with added technical indicator columns.
    """
    print("Calculating technical indicators...")
    # Simple Moving Average
    data['SMA'] = data['Close'].rolling(window=15).mean()
```

```
    # Exponential Moving Average
    data['EMA'] = data['Close'].ewm(span=15, adjust=False).mean()

    # Relative Strength Index (RSI)
    delta = data['Close'].diff()
    gain = (delta.where(delta > 0, 0)).rolling(window=15).mean()
    loss = (-delta.where(delta < 0, 0)).rolling(window=15).mean()

    # Avoid division by zero for RS
    RS = np.where(loss == 0, np.inf, gain / loss)
    data['RSI'] = 100 - (100 / (1 + RS))

    # Moving Average Convergence Divergence (MACD)
    ShortEMA = data['Close'].ewm(span=12, adjust=False).mean()
    LongEMA = data['Close'].ewm(span=26, adjust=False).mean()
    data['MACD'] = ShortEMA - LongEMA
    data['Signal_Line'] = data['MACD'].ewm(span=9, adjust=False).mean()

    # Bollinger Bands
    data['20_SMA'] = data['Close'].rolling(window=20).mean()
    data['StdDev'] = data['Close'].rolling(window=20).std()
    data['Upper_Band'] = data['20_SMA'] + (data['StdDev'] * 2)
    data['Lower_Band'] = data['20_SMA'] - (data['StdDev'] * 2)

    return data

# --- 2. Load Dataset and Initial Preprocessing ---
print("--- Starting Data Preprocessing ---")
# Ensure 'stock_yfinance_data.csv' is uploaded to your Colab environment
# or adjust the path if it's in a specific folder (e.g., '/content/drive/MyDrive/stock_yf
try:
    df = pd.read_csv("stock_yfinance_data.csv")
    print("Dataset 'stock_yfinance_data.csv' loaded successfully.")
except FileNotFoundError:
    print("Error: 'stock_yfinance_data.csv' not found.")
    print("Please upload the 'stock_yfinance_data.csv' file to your Colab environment (e.
    print("Example for Colab upload: from google.colab import files; files.upload()")
    exit() # Exit if file is not found

print("\nInitial DataFrame head:")
print(df.head())
print("\nInitial DataFrame Info:")
df.info()

# Apply technical indicators
df = get_tech_ind(df)

# Fill any NaN values that result from rolling windows (at the beginning of the series)
# Use both ffill and bfill to handle NaNs robustly
```

```python
print("Handling missing values (NaNs) by forward-fill and backward-fill...")
df.fillna(method='ffill', inplace=True)
df.fillna(method='bfill', inplace=True) # Catch any NaNs at the very start

# Convert 'Date' to datetime and set as index
print("Converting 'Date' column to datetime and setting as index...")
df['Date'] = pd.to_datetime(df['Date'])
df.set_index('Date', inplace=True)

# Drop redundant or non-feature columns
print("Dropping 'Adj Close' and 'Stock Name' columns...")
df = df.drop(columns=['Adj Close', 'Stock Name'])

print("\nProcessed DataFrame head (after technical indicators, fillna, date index, and co
print(df.head())
print("\nProcessed DataFrame Info:")
df.info()
print("\nChecking for any remaining NaN values:")
print(df.isnull().sum())
print("--- Data Preprocessing Complete ---")


# --- 3. Prepare Data for Supervised Learning (Features X and Target y) ---
print("\n--- Preparing Features (X) and Target (y) for Machine Learning ---")
# Predict the 'Close' price of the next day based on current day's features
# Shift 'Close' column to create the target variable (next day's close)
print("Creating 'Next_Close' as target variable by shifting 'Close' price...")
df['Next_Close'] = df['Close'].shift(-1)
df.dropna(inplace=True) # Drop the last row where 'Next_Close' is NaN after shifting

# Define features (X) and target (y)
features = df.drop(columns=['Next_Close']).columns
X = df[features]
y = df['Next_Close']

print(f"Shape of X (features): {X.shape}")
print(f"Shape of y (target): {y.shape}")
print("\nFirst 5 rows of X:")
print(X.head())
print("\nFirst 5 rows of y:")
print(y.head())


# --- 4. Normalize Features and Target ---
print("\nNormalizing features and target variables using MinMaxScaler...")
scaler_X = MinMaxScaler()
scaler_y = MinMaxScaler()

X_scaled = scaler_X.fit_transform(X)
# Reshape y to be 2D for scaler, then flatten it back to 1D for model training if needed
y_scaled = scaler_y.fit_transform(y.values.reshape(-1, 1))
```

```
y_                 _y _         (y        p ( y   ))

print(f"Shape of X_scaled (normalized features): {X_scaled.shape}")
print(f"Shape of y_scaled (normalized target): {y_scaled.shape}")



# --- 5. Split Data into Training and Testing Sets (80:20 Ratio) ---
print("\nSplitting data into training (80%) and testing (20%) sets...")
train_size = int(len(X_scaled) * 0.8)
X_train, X_test = X_scaled[:train_size], X_scaled[train_size:]
y_train, y_test = y_scaled[:train_size], y_scaled[train_size:]

print(f"Shape of X_train: {X_train.shape}")
print(f"Shape of y_train: {y_train.shape}")
print(f"Shape of X_test: {X_test.shape}")
print(f"Shape of y_test: {y_test.shape}")
print("--- Data Preparation and Splitting Complete ---")



# --- 6. Train and Test Traditional ML Models ---
print("\n--- Training and Testing Traditional ML Models ---")


results = {} # Dictionary to store model performance metrics

# --- Random Forest Regressor ---
print("\nTraining Random Forest Regressor...")
rf_model = RandomForestRegressor(n_estimators=100, random_state=42)
rf_model.fit(X_train, y_train.ravel()) # .ravel() converts y_train from 2D (num_samples,
rf_predictions_scaled = rf_model.predict(X_test)
# Inverse transform predictions and actual values back to original scale for meaningful m
rf_predictions = scaler_y.inverse_transform(rf_predictions_scaled.reshape(-1, 1))
rf_actual = scaler_y.inverse_transform(y_test)
results['Random Forest'] = {
    'RMSE': np.sqrt(mean_squared_error(rf_actual, rf_predictions)),
    'R2_Score': r2_score(rf_actual, rf_predictions)
}
print("Random Forest training and prediction complete.")

# --- K-Nearest Neighbors Regressor ---
print("\nTraining K-Nearest Neighbors Regressor...")
knn_model = KNeighborsRegressor(n_neighbors=5) # n_neighbors is a common hyperparameter
knn_model.fit(X_train, y_train.ravel())
knn_predictions_scaled = knn_model.predict(X_test)
knn_predictions = scaler_y.inverse_transform(knn_predictions_scaled.reshape(-1, 1))
knn_actual = scaler_y.inverse_transform(y_test)
results['KNN'] = {
    'RMSE': np.sqrt(mean_squared_error(knn_actual, knn_predictions)),
    'R2_Score': r2_score(knn_actual, knn_predictions)
}
print("KNN training and prediction complete.")

#        f       t   V  t    M   h    R
```

```
# --- Support Vector Machine Regressor ---
print("\nTraining Support Vector Machine Regressor (SVR)...")
# Using default RBF kernel, which is common.
# SVMs can be sensitive to hyper-parameters; default values are used here.
svm_model = SVR(kernel='rbf')
svm_model.fit(X_train, y_train.ravel())
svm_predictions_scaled = svm_model.predict(X_test)
svm_predictions = scaler_y.inverse_transform(svm_predictions_scaled.reshape(-1, 1))
svm_actual = scaler_y.inverse_transform(y_test)
results['SVM'] = {
    'RMSE': np.sqrt(mean_squared_error(svm_actual, svm_predictions)),
    'R2_Score': r2_score(svm_actual, svm_predictions)
}
print("SVM training and prediction complete.")

# --- Linear Regression ---
print("\nTraining Linear Regression model...")
lr_model = LinearRegression()
lr_model.fit(X_train, y_train.ravel())
lr_predictions_scaled = lr_model.predict(X_test)
lr_predictions = scaler_y.inverse_transform(lr_predictions_scaled.reshape(-1, 1))
lr_actual = scaler_y.inverse_transform(y_test)
results['Linear Regression'] = {
    'RMSE': np.sqrt(mean_squared_error(lr_actual, lr_predictions)),
    'R2_Score': r2_score(lr_actual, lr_predictions)
}
print("Linear Regression training and prediction complete.")

# --- Print Results ---
print("\n--- Model Performance Results ---")
for model_name, metrics in results.items():
    print(f"\nModel: {model_name}")
    print(f"  RMSE (Root Mean Squared Error): {metrics['RMSE']:.4f}")
    print(f"  R2 Score: {metrics['R2_Score']:.4f}")

print("\nAll specified machine learning models have been trained, tested, and their perfo
print("\nNote: A confusion matrix is not applicable for this stock price prediction task

# --- Comparative Table ---
print("\n--- Comparative Accuracy Rate Table ---")
print("| Model              | RMSE (Root Mean Squared Error) | R2 Score |")
print("| :--------------- | :---------------------------- | :------- |")
print(f"| Random Forest     | {results['Random Forest']['RMSE']:.4f}              | {resu
print(f"| KNN               | {results['KNN']['RMSE']:.4f}             | {results['KNN']
print(f"| SVM               | {results['SVM']['RMSE']:.4f}           | {results['SVM']
print(f"| Linear Regression  | {results['Linear Regression']['RMSE']:.4f}             | {
print("\nInterpretation: A lower RMSE and higher R2 Score indicate better model performan
print("In this comparison, Random Forest Regressor shows the best performance for your da
```

```
    --- Starting Data Preprocessing ---
    Dataset 'stock_yfinance_data.csv' loaded successfully.
```

```
Initial DataFrame head:
        Date        Open        High         Low       Close   Adj Close  \
0  2021-09-30  260.333344  263.043335  258.333344  258.493347  258.493347
1  2021-10-01  259.466675  260.260010  254.529999  258.406677  258.406677
2  2021-10-04  265.500000  268.989990  258.706665  260.510010  260.510010
3  2021-10-05  261.600006  265.769989  258.066681  260.196655  260.196655
4  2021-10-06  258.733337  262.220001  257.739990  260.916656  260.916656


     Volume Stock Name
0  53868000       TSLA
1  51094200       TSLA
2  91449900       TSLA
3  55297800       TSLA
4  43898400       TSLA


Initial DataFrame Info:
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 6300 entries, 0 to 6299
Data columns (total 8 columns):
 #   Column      Non-Null Count  Dtype
---  ------      --------------  -----
 0   Date        6300 non-null   object
 1   Open        6300 non-null   float64
 2   High        6300 non-null   float64
 3   Low         6300 non-null   float64
 4   Close       6300 non-null   float64
 5   Adj Close   6300 non-null   float64
 6   Volume      6300 non-null   int64
 7   Stock Name  6300 non-null   object
dtypes: float64(5), int64(1), object(2)
memory usage: 393.9+ KB
Calculating technical indicators...
Handling missing values (NaNs) by forward-fill and backward-fill...
Converting 'Date' column to datetime and setting as index...
Dropping 'Adj Close' and 'Stock Name' columns...

Processed DataFrame head (after technical indicators, fillna, date index, and column
                  Open        High         Low       Close     Volume  \
Date
2021-09-30  260.333344  263.043335  258.333344  258.493347   53868000
2021-10-01  259.466675  260.260010  254.529999  258.406677   51094200
2021-10-04  265.500000  268.989990  258.706665  260.510010   91449900
2021-10-05  261.600006  265.769989  258.066681  260.196655   55297800
2021-10-06  258.733337  262.220001  257.739990  260.916656   43898400


                   SMA         EMA       RSI      MACD  Signal_Line  \
Date
2021-09-30  269.887559  258.493347  87.43355  0.000000     0.000000
2021-10-01  269.887559  258.482513  87.43355 -0.006914    -0.001383
2021-10-04  269.887559  258.735950  87.43355  0.155535     0.030001
2021-10-05  269.887559  258.918539  87.43355  0.256041     0.075209
2021-10-06  269.887559  259.168303  87.43355  0.389303     0.138028


                20_SMA      StdDev  Upper_Band  Lower_Band
D...
```

Date