

## ✓ Homework 2

Group Members: Yash Shah, Zheng Wan, Linhui Sang

### ✓ Import Libraries

```
import numpy as np
import yfinance as yf
import zipfile
import pandas as pd
import requests
from io import BytesIO
from scipy.interpolate import interp1d
from scipy import integrate
import scipy.stats as ss
from scipy.stats import norm
import matplotlib.pyplot as plt
from google.colab import drive
import torch
import torch.nn as nn
import torch.optim as optim
from torch.utils.data import TensorDataset, DataLoader
```

### ✓ Download Data

```

import yfinance as yf

# Define the ticker symbol for the VIX (CBOE Volatility Index)
ticker_symbol = '^VIX'

# Define the start and end dates
start_date = '2016-01-01'
end_date = '2016-04-30'

# Download the data from Yahoo Finance
vix_data = yf.download(ticker_symbol, start=start_date, end=end_date)

# Display the first few rows of the downloaded data
print(vix_data.head())

```

[\*\*\*\*\*100%\*\*\*\*\*] 1 of 1 completed

Date

2016-01-04	22.480000	23.360001	20.670000	20.700001	20.700001	0
2016-01-05	20.750000	21.059999	19.250000	19.340000	19.340000	0
2016-01-06	21.670000	21.860001	19.799999	20.590000	20.590000	0
2016-01-07	23.219999	25.860001	22.400000	24.990000	24.990000	0
2016-01-08	22.959999	27.080000	22.480000	27.010000	27.010000	0

```

import zipfile
import pandas as pd
import requests
from io import BytesIO

# URL of the ZIP file
zip_file_url = 'https://github.com/larrysangfake/financial_econometrics/blob/main/SPX_2016_options.zip'

# Download the ZIP file
response = requests.get(zip_file_url)
zip_data = BytesIO(response.content)

# Path to the CSV file inside the ZIP file
csv_file_name = 'SPX_2016_options.csv'

# Open the ZIP file
with zipfile.ZipFile(zip_data, 'r') as zip_ref:
    # Extract the CSV file
    zip_ref.extract(csv_file_name, path='.')

# Read the CSV file into a DataFrame
df = pd.read_csv(csv_file_name, header=None)

# Now you can work with the DataFrame (df)
print(df.head())

```

	0	1	2	3	4	5	6	7	8	9	\
0	108105	20160104	736333	736344	11	1	1000	1007.4	1010.8	2000	
1	108105	20160104	736333	736344	11	1	1025	982.3	985.8	0	
2	108105	20160104	736333	736344	11	1	1050	957.3	960.9	0	
3	108105	20160104	736333	736344	11	1	1075	932.3	935.9	0	
4	108105	20160104	736333	736344	11	1	1100	907.4	911.1	0	

	10	11	12	13	14	15	16
0	42291	NaN	2012.66	-0.015304	1	1	0.007121
1	0	NaN	2012.66	-0.015304	1	1	0.007121
2	0	NaN	2012.66	-0.015304	1	1	0.007121
3	0	NaN	2012.66	-0.015304	1	1	0.007121
4	20	NaN	2012.66	-0.015304	1	1	0.007121

```
df.columns = [
    'ID', 'Date', 'Julian Date', 'Julian Maturity Date', 'Time Difference in Days',
    'Strike Price', 'bid', 'ask', 'Unnamed:9', 'Unnamed:10', 'volatility',
    'stock price', 'Unnamed:13', 'Unnamed: 14', 'Unnamed: 15', 'interest rate'
]

# Print the DataFrame to check the new column names
print(df.head())
```

	ID	Date	Julian Date	Julian Maturity Date	\
0	108105	20160104	736333	736344	
1	108105	20160104	736333	736344	
2	108105	20160104	736333	736344	
3	108105	20160104	736333	736344	
4	108105	20160104	736333	736344	

	Time Difference in Days	Call/Put	Strike Price	bid	ask	Unnamed: 9
0	11	1	1000	1007.4	1010.8	20
1	11	1	1025	982.3	985.8	
2	11	1	1050	957.3	960.9	
3	11	1	1075	932.3	935.9	
4	11	1	1100	907.4	911.1	

	Unnamed:10	volatility	stock price	Unnamed:13	Unnamed: 14	Unnamed: 15
0	42291	NaN	2012.66	-0.015304	1	
1	0	NaN	2012.66	-0.015304	1	
2	0	NaN	2012.66	-0.015304	1	
3	0	NaN	2012.66	-0.015304	1	
4	20	NaN	2012.66	-0.015304	1	

	interest rate
0	0.007121
1	0.007121
2	0.007121
3	0.007121
4	0.007121

## ✓ 1. Data Cleaning

```
cleaned_data = df.dropna(subset=['volatility'])
print(cleaned_data.head())
```

	ID	Date	Julian Date	Julian Maturity Date	\
128	108105	20160104	736333	736344	
129	108105	20160104	736333	736344	
130	108105	20160104	736333	736344	
131	108105	20160104	736333	736344	
132	108105	20160104	736333	736344	

	Time Difference in Days	Call/Put	Strike Price	bid	ask	Unnamed
128	11	1	1875	135.3	138.2	5
129	11	1	1880	130.1	133.9	
130	11	1	1885	125.4	129.1	
131	11	1	1890	120.7	124.4	
132	11	1	1895	116.5	119.3	

	Unnamed:10	volatility	stock price	Unnamed:13	Unnamed: 14	\
128	311	0.149708	2012.66	-0.015304	1	
129	29	0.184996	2012.66	-0.015304	1	
130	7	0.196155	2012.66	-0.015304	1	
131	30	0.203980	2012.66	-0.015304	1	
132	5	0.209911	2012.66	-0.015304	1	

	Unnamed: 15	interest rate
128	1	0.007121
129	1	0.007121
130	1	0.007121
131	1	0.007121
132	1	0.007121

## ✓ 2. Create a new column

```
cleaned_data['Average Price'] = (cleaned_data['bid'] + cleaned_data['ask']) / 2
```

```
# To see the updated DataFrame with the new 'Average Price' column
print(cleaned_data.head())
```

	ID	Date	Julian Date	Julian Maturity Date	\
128	108105	20160104	736333	736344	
129	108105	20160104	736333	736344	
130	108105	20160104	736333	736344	
131	108105	20160104	736333	736344	
132	108105	20160104	736333	736344	

	Time Difference in Days	Call/Put	Strike Price	bid	ask	Unnamed
128	11	1	1875	135.3	138.2	5
129	11	1	1880	130.1	133.9	
130	11	1	1885	125.4	129.1	
131	11	1	1890	120.7	124.4	
132	11	1	1895	116.5	119.3	

	Unnamed:10	volatility	stock price	Unnamed:13	Unnamed: 14	\
128	311	0.149708	2012.66	-0.015304	1	
129	29	0.184996	2012.66	-0.015304	1	
130	7	0.196155	2012.66	-0.015304	1	
131	30	0.203980	2012.66	-0.015304	1	
132	5	0.209911	2012.66	-0.015304	1	

	Unnamed: 15	interest rate	Average Price
128	1	0.007121	136.75
129	1	0.007121	132.00
130	1	0.007121	127.25
131	1	0.007121	122.55
132	1	0.007121	117.90

```
<ipython-input-6-d0d5f6261de8>:1: 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>

```
cleaned_data['Average Price'] = (cleaned_data['bid'] + cleaned_data['ask']
```

### ✓ 3. Filter the columns

```
filtered_df = cleaned_data[(cleaned_data['Average Price'] > 0.05)]
```

```
# Check the filtered data
print(filtered_df.head())
```

	ID	Date	Julian Date	Julian Maturity Date	\
128	108105	20160104	736333	736344	
129	108105	20160104	736333	736344	
130	108105	20160104	736333	736344	
131	108105	20160104	736333	736344	
132	108105	20160104	736333	736344	

	Time Difference in Days	Call/Put	Strike Price	bid	ask	Unnamed
128	11	1	1875	135.3	138.2	5
129	11	1	1880	130.1	133.9	
130	11	1	1885	125.4	129.1	
131	11	1	1890	120.7	124.4	
132	11	1	1895	116.5	119.3	

	Unnamed:10	volatility	stock price	Unnamed:13	Unnamed: 14	\
128	311	0.149708	2012.66	-0.015304	1	
129	29	0.184996	2012.66	-0.015304	1	
130	7	0.196155	2012.66	-0.015304	1	
131	30	0.203980	2012.66	-0.015304	1	
132	5	0.209911	2012.66	-0.015304	1	

	Unnamed: 15	interest rate	Average Price
128	1	0.007121	136.75
129	1	0.007121	132.00
130	1	0.007121	127.25
131	1	0.007121	122.55
132	1	0.007121	117.90

## ✓ 4. Out the Money

```
# Filter Out of the Money Call Options (S < K)
otm_calls = filtered_df[(filtered_df['Call/Put'] == 1) & (filtered_df['stock pr

# Filter Out of the Money Put Options (S > K)
otm_puts = filtered_df[(filtered_df['Call/Put'] == -1) & (filtered_df['stock pr

# Combine OTM calls and puts into a single DataFrame
otm_options = pd.concat([otm_calls, otm_puts])

# Check the resulting DataFrame
print(otm_options.head())
```

	ID	Date	Julian Date	Julian Maturity Date	\
156	108105	20160104	736333	736344	
157	108105	20160104	736333	736344	
158	108105	20160104	736333	736344	
159	108105	20160104	736333	736344	
160	108105	20160104	736333	736344	

	Time Difference in Days	Call/Put	Strike Price	bid	ask	Unnamed:9
156	11	1	2015	23.4	25.5	1852
157	11	1	2020	20.6	22.4	1105
158	11	1	2025	18.1	20.0	4957
159	11	1	2030	15.9	17.7	187
160	11	1	2035	13.8	15.6	180

	Unnamed:10	volatility	stock price	Unnamed:13	Unnamed: 14	\
156	4389	0.195975	2012.66	-0.015304	1	
157	2124	0.190988	2012.66	-0.015304	1	
158	29268	0.188562	2012.66	-0.015304	1	
159	1741	0.186419	2012.66	-0.015304	1	
160	628	0.184163	2012.66	-0.015304	1	

	Unnamed: 15	interest rate	Average Price
156	1	0.007121	24.45
157	1	0.007121	21.50
158	1	0.007121	19.05
159	1	0.007121	16.80
160	1	0.007121	14.70

## ✓ 5.The implied volatility



#Black and Scholes

```
def BlackScholes(CallPutFlag,S,X,v,r,T):
```

```
    d1 = (np.log(S/X)+(r+v*v/2.)*T)/(v*np.sqrt(T))
```

```
    d2 = d1-v*np.sqrt(T)
```

```
    if CallPutFlag=="C":
```

```
        P = S*norm.cdf(d1) - X*np.exp(-r*T)*norm.cdf(d2)
```

```
    else:
```

```
        P = -S*norm.cdf(-d1) + X*np.exp(-r*T)*norm.cdf(-d2)
```

```
    return P
```

```
def ivol(K,IV,Kall):
```

```
    if Kall>=K[len(K)-1]:
```

```
        Kall=K[len(K)-1]
```

```
    if Kall<=K[0]:
```

```
        Kall=K[0]
```

```
    funy = interp1d(K,IV, kind='cubic', fill_value="extrapolate")
```

```
    y=funy(Kall)
```

```
    if (np.sum(y<0)>0):
```

```
        if Kall>=K[len(K)-1]:
```

```
            Kall=K[len(K)-1]
```

```
        if Kall<=K[0]:
```

```
            Kall=K[0]
```

```
        funy = interp1d(K,IV, kind='linear', fill_value="extrapolate")
```

```
        y=funy(Kall)
```

```
    return(y)
```

```

def RiskNeutralVolatilitySkewKurt_JVCR(Kvector, IVvector, S0, T, r):
    kmin=.1*S0;
    kmax=3.5*S0;

    def V1(K):
        V1=2*(1-np.log(K/S0))*BlackScholes("C",S0, K, ivol(Kvector,IVvector,K),
        return(V1)
    def V2(K):
        V2=2*(1+np.log(S0/K))*BlackScholes("P",S0, K, ivol(Kvector,IVvector,K),
        return(V2)

    def W1(K):
        W1=(6*np.log(K/S0)-3*np.power(np.log(K/S0),2))*BlackScholes("C",S0, K,
        return(W1)

    def W2(K):
        W2=(6*np.log(S0/K)+3*np.power(np.log(S0/K),2))*BlackScholes("P",S0, K,
        return(W2)

    def X1(K):
        X1=((12*np.power(np.log(K/S0),2) - 4*np.power(np.log(K/S0),3))*(Black
        return(X1)

    def X2(K):
        X2=((12*np.power(np.log(S0/K),2) + 4*np.power(np.log(S0/K),3))*(Black
        return(X2)
    V=integrate.quad(V1,S0,kmax)[0]+integrate.quad(V2,kmin,S0)[0]
    W=integrate.quad(W1,S0,kmax)[0]-integrate.quad(W2,kmin,S0)[0]
    X=integrate.quad(X1,S0,kmax)[0]+integrate.quad(X2,kmin,S0)[0]
    mu=np.exp(r*T)-1-np.exp(r*T)*V/2-np.exp(r*T)*W/6-np.exp(r*T)*X/24;
    #print(V,W,X,mu)

    vol=np.sqrt(1/T * V);
    skew=( np.exp(r*T)*W - 3*mu*np.exp(r*T)*V + 2*np.power(mu,3)) / np.power(np
    kurt=( np.exp(r*T)*X - 4*mu*np.exp(r*T)*W + 6*np.exp(r*T)*np.power(mu,2)*V

    return([vol,skew,kurt]);

```

```

# Convert dates and calculate time to maturity in years
otm_options['Date'] = pd.to_datetime(otm_options['Date'], format='%Y%m%d')
otm_options['Time to Maturity'] = otm_options['Time Difference in Days'] / 365.

# Define a function to apply risk-neutral volatility, skewness, and kurtosis ca
def apply_risk_neutral_vol_skew_kurt(group):
    # Extract group-level constants
    S0 = group['stock price'].iloc[0]
    r = group['interest rate'].iloc[0]
    T = group['Time to Maturity'].iloc[0]

    # Calculate implied volatility, skewness, and kurtosis
    if len(group) < 2:
        return pd.Series([np.nan] * len(group), index=group.index) # Return Na
    else:
        vol, skew, kurt = RiskNeutralVolatilitySkewKurt_JVCR(group['Strike Pric
        return pd.Series([vol] * len(group), index=group.index) # Repeat the c

# Apply the function to each group and create a new DataFrame with the results
volatility_results = otm_options.groupby(['Date', 'Julian Maturity Date']).appl

# Merge the results back to the original DataFrame
otm_options['Implied Volatility'] = volatility_results.reset_index(level=[0, 1]

# Display the results
print(otm_options[['Date', 'Julian Maturity Date', 'Implied Volatility']].head(


```

	Date	Julian Maturity Date	Implied Volatility
156	2016-01-04	736344	0.195965
157	2016-01-04	736344	0.195965
158	2016-01-04	736344	0.195965
159	2016-01-04	736344	0.195965
160	2016-01-04	736344	0.195965

```

# Merge the results back to the original DataFrame
otm_options = pd.merge(otm_options, volatility_results, on=['Date', 'Julian Mat

# Display the results
print(otm_options[['Date', 'Julian Maturity Date', 'Implied Volatility']].head(

import pandas as pd
import numpy as np
from scipy.interpolate import interp1d

# Assuming otm_options is predefined
unique_dates = otm_options['Julian Date'].unique()
unique_maturities = otm_options['Julian Maturity Date'].unique()

results = []

```

```

for date in unique_dates:
    for maturity in unique_maturities:
        subset = otm_options[(otm_options['Julian Date'] == date) & (otm_option

if len(subset) > 1:
    K = subset['Strike Price'].values
    IV = subset['volatility'].values # Ensure this column is correctly

# Safely apply interpolation
try:
    # Assuming ivol is already defined and handles extrapolation or
    interp_func = interp1d(K, IV, kind='linear', fill_value='extrap
    subset = subset.copy() # Work on a copy to avoid SettingWithCo
    subset['implied_volatility'] = subset['Strike Price'].apply(lam
    results.append(subset)
except Exception as e:
    print(f"Error in interpolation for date {date} and maturity {ma

# Combine all results into a single DataFrame
result_df = pd.concat(results, ignore_index=True)

print(result_df.head())

```

	ID	Date	Julian Date	Julian Maturity Date	\
0	108105	20160104	736333	736344	
1	108105	20160104	736333	736344	
2	108105	20160104	736333	736344	
3	108105	20160104	736333	736344	
4	108105	20160104	736333	736344	

	Time Difference in Days	Call/Put	Strike Price	bid	ask	Unnamed:9
0	11	1	2015	23.4	25.5	1852
1	11	1	2020	20.6	22.4	1105
2	11	1	2025	18.1	20.0	4957
3	11	1	2030	15.9	17.7	187
4	11	1	2035	13.8	15.6	180

	Unnamed:10	volatility	stock price	Unnamed:13	Unnamed: 14	Unnamed: 1
0	4389	0.195975	2012.66	-0.015304	1	
1	2124	0.190988	2012.66	-0.015304	1	
2	29268	0.188562	2012.66	-0.015304	1	
3	1741	0.186419	2012.66	-0.015304	1	
4	628	0.184163	2012.66	-0.015304	1	

	interest rate	Average Price	implied_volatility
0	0.007121	24.45	0.195975
1	0.007121	21.50	0.190988
2	0.007121	19.05	0.188562
3	0.007121	16.80	0.186419
4	0.007121	14.70	0.184163

## ✓ 6. 30-day volatility

```
# Check how often we have exactly 30 days to maturity
days_30 = otm_options[otm_options['Time Difference in Days'] == 30]

# See how many such entries exist
print(f"Entries with exactly 30 days to maturity: {days_30.shape[0]}")

# Optional: View some of these entries to verify
print(days_30.head())
```

Entries with exactly 30 days to maturity: 3088

	ID	Date	Julian Date	Julian Maturity Date	\
18638	108105	2016-01-06	736335	736365	
18639	108105	2016-01-06	736335	736365	
18640	108105	2016-01-06	736335	736365	
18641	108105	2016-01-06	736335	736365	
18642	108105	2016-01-06	736335	736365	

	Time Difference in Days	Call/Put	Strike Price	bid	ask	Unnamed
18638	30	1	1995	37.4	38.1	
18639	30	1	2000	34.6	35.3	
18640	30	1	2005	32.0	32.6	
18641	30	1	2010	29.4	30.0	
18642	30	1	2015	26.9	27.5	10

	Unnamed:10	volatility	stock price	Unnamed:13	Unnamed: 14	\
18638	62	0.184500	1990.26	-0.013115	1	
18639	1606	0.182070	1990.26	-0.013115	1	
18640	49	0.179860	1990.26	-0.013115	1	
18641	51	0.177407	1990.26	-0.013115	1	
18642	1044	0.174921	1990.26	-0.013115	1	

	Unnamed: 15	interest rate	Average Price	Time to Maturity	\
18638	1	0.006932	37.75	0.082136	
18639	1	0.006932	34.95	0.082136	
18640	1	0.006932	32.30	0.082136	
18641	1	0.006932	29.70	0.082136	
18642	1	0.006932	27.20	0.082136	

	Implied Volatility
18638	0.18447
18639	0.18447
18640	0.18447
18641	0.18447
18642	0.18447

```
print(otm_options.groupby('Date')['Time Difference in Days'].agg([min, max]))
```

Date	min	max
2016-01-04	4	1082
2016-01-05	3	1081
2016-01-06	2	1080
2016-01-07	8	1079
2016-01-08	7	1078
...	...	...
2016-04-25	2	970
2016-04-26	3	969
2016-04-27	2	968
2016-04-28	6	967
2016-04-29	5	966

[82 rows x 2 columns]

```
# Check the distribution of 'Time Difference in Days' across all data  
print(otm_options['Time Difference in Days'].describe())
```

```
count    272078.000000  
mean      137.898290  
std       203.938153  
min        2.000000  
25%       32.000000  
50%       60.000000  
75%      130.000000  
max      1082.000000  
Name: Time Difference in Days, dtype: float64
```

```

import pandas as pd
from scipy.interpolate import interp1d
import numpy as np

# Sample data setup
data = {
    'Date': ['2024-01-01', '2024-01-01', '2024-01-01', '2024-01-02', '2024-01-02'],
    'Time Difference in Days': [25, 28, 35, 27, 30, 33],
    'Implied Volatility': [0.20, 0.19, 0.18, 0.22, 0.21, 0.20]
}
sample_df = pd.DataFrame(data)
sample_df['Date'] = pd.to_datetime(sample_df['Date'])

def interpolate_iv_at_30(group):
    days = group['Time Difference in Days'].values
    ivs = group['Implied Volatility'].values

    if 30 in days:
        return ivs[np.where(days == 30)[0][0]]
    else:
        interp_function = interp1d(days, ivs, kind='linear', bounds_error=False,
                                     interpolated_value = float(interp_function(30)))

        if min(days) < 30 < max(days):
            return interpolated_value
        else:
            closest_idx = (np.abs(days - 30)).argmin()
            return ivs[closest_idx]

# Apply the function and store the result
result_sample = sample_df.groupby('Date').apply(interpolate_iv_at_30)

# Merge results back into the original DataFrame
sample_df['IV at 30 Days'] = sample_df['Date'].map(result_sample)

# Print the result
print(sample_df[['Date', 'IV at 30 Days']].drop_duplicates())

```

	Date	IV at 30 Days
0	2024-01-01	0.187143
3	2024-01-02	0.210000

```

result = otm_options.groupby('Date').apply(interpolate_iv_at_30)

# Merge results back into the original DataFrame
otm_options['IV at 30 Days'] = otm_options['Date'].map(result)

# Print the result
print(otm_options[['Date', 'IV at 30 Days']].drop_duplicates())

```

	Date	IV at 30 Days
156	2016-01-04	0.178551
7325	2016-01-05	0.172282
14661	2016-01-06	0.184470
21993	2016-01-07	0.225709
29751	2016-01-08	0.226559
...	...	...
582448	2016-04-25	0.116834
589881	2016-04-26	0.116640
597494	2016-04-27	0.115469
605102	2016-04-28	0.124032
612548	2016-04-29	0.130302

```
[82 rows x 2 columns]
```

## 7. Plot

```

import matplotlib.pyplot as plt
import pandas as pd

otm_options['IV at 30 Days'] *= 100 # Convert to percentage

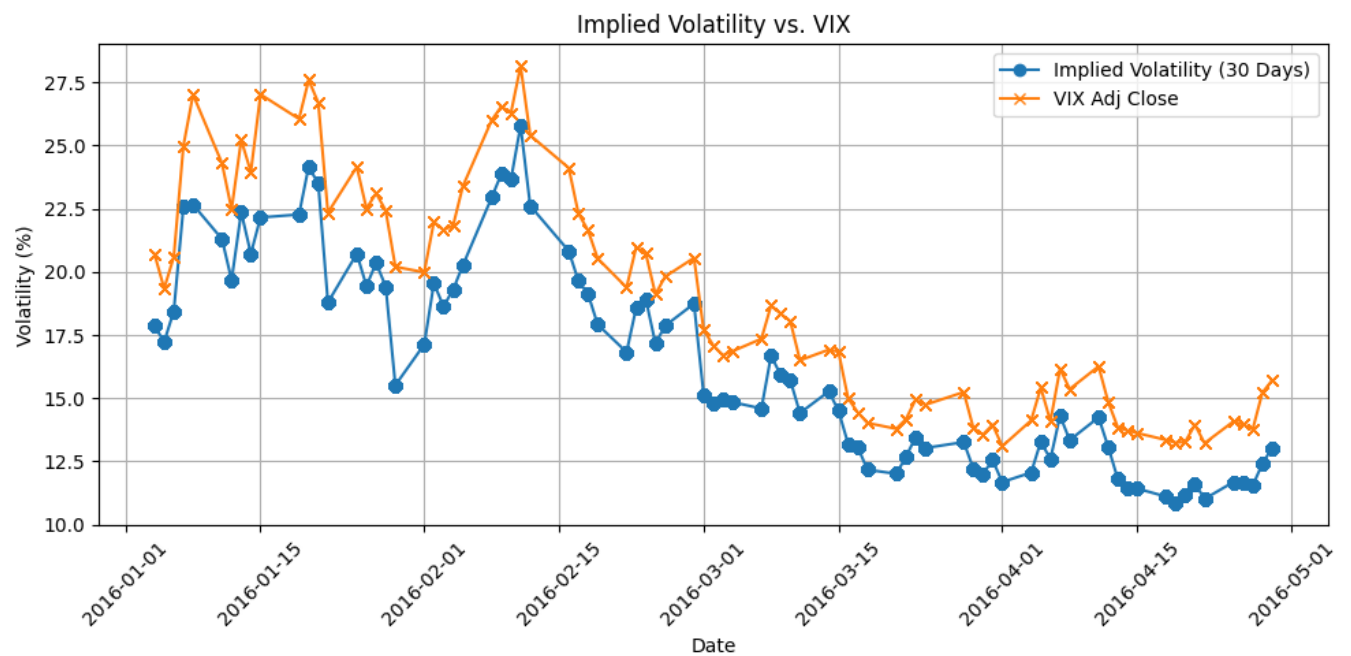
# Merge VIX data with your options data
combined_df = pd.merge(otm_options, vix_data, on='Date', how='inner')

# Plotting both implied volatility and VIX
plt.figure(figsize=(10, 5))
plt.plot(combined_df['Date'], combined_df['IV at 30 Days'], label='Implied Volatility')
plt.plot(combined_df['Date'], combined_df['Adj Close'], label='VIX Adj Close',
plt.title('Implied Volatility vs. VIX')
plt.xlabel('Date')
plt.ylabel('Volatility (%)')
plt.legend()
plt.grid(True)
plt.xticks(rotation=45)
plt.tight_layout()
plt.show()

# Calculate the correlation
correlation = combined_df['IV at 30 Days'].corr(combined_df['Close'])
print("Correlation between IV at 30 Days and VIX Close:", correlation)

```





Correlation between IV at 30 Days and VIX Close: 0.9925817967753203

## ✓ Paper Code

```
options_df =pd.read_csv('./options_combined.csv')
```

```
pd.set_option('display.max_columns', None)
options_df.head()
```

	date	exdate	cp_flag	optiondate	expirationdate	week_day	strike_p
0	20160104	20160115	1	736333	736344	2	
1	20160104	20160115	1	736333	736344	2	
2	20160104	20160115	1	736333	736344	2	
3	20160104	20160115	1	736333	736344	2	
4	20160104	20160115	1	736333	736344	2	

```
options_df_filtered = options_df[(options_df['volume'] != 0) | (options_df['best_bid'] > 0)]
```

```
options_df_filtered.head()
```

	date	exdate	cp_flag	optiondate	expirationdate	week_day	strike_p
0	20160104	20160115	1	736333	736344	2	
1	20160104	20160115	1	736333	736344	2	
2	20160104	20160115	1	736333	736344	2	
3	20160104	20160115	1	736333	736344	2	
4	20160104	20160115	1	736333	736344	2	

```
options_df_filtered[options_df_filtered['callprice'] - options_df_filtered['putprice'] > 0]
```

	date	exdate	cp_flag	optiondate	expirationdate	week_day	strike_price
--	------	--------	---------	------------	----------------	----------	--------------

```
options_df_filtered.columns
```

```
Index(['date', 'exdate', 'cp_flag', 'optiondate', 'expirationdate',
      'week_day',
      'strike_price', 'maturity', 'moneyness', 'index_level',
      'index_return',
      'div_yield', 'riskfree', 'riskfree_option', 'best_bid',
      'best_offer',
      'callprice', 'putprice', 'midquote', 'volume', 'open_interest',
      'impl_volatility', 'vega'],
      dtype='object')
```

```
final_df = options_df_filtered[['optiondate', 'maturity', 'strike_price', 'moneyness', 'impl_volatility']]
final_df.head()
```

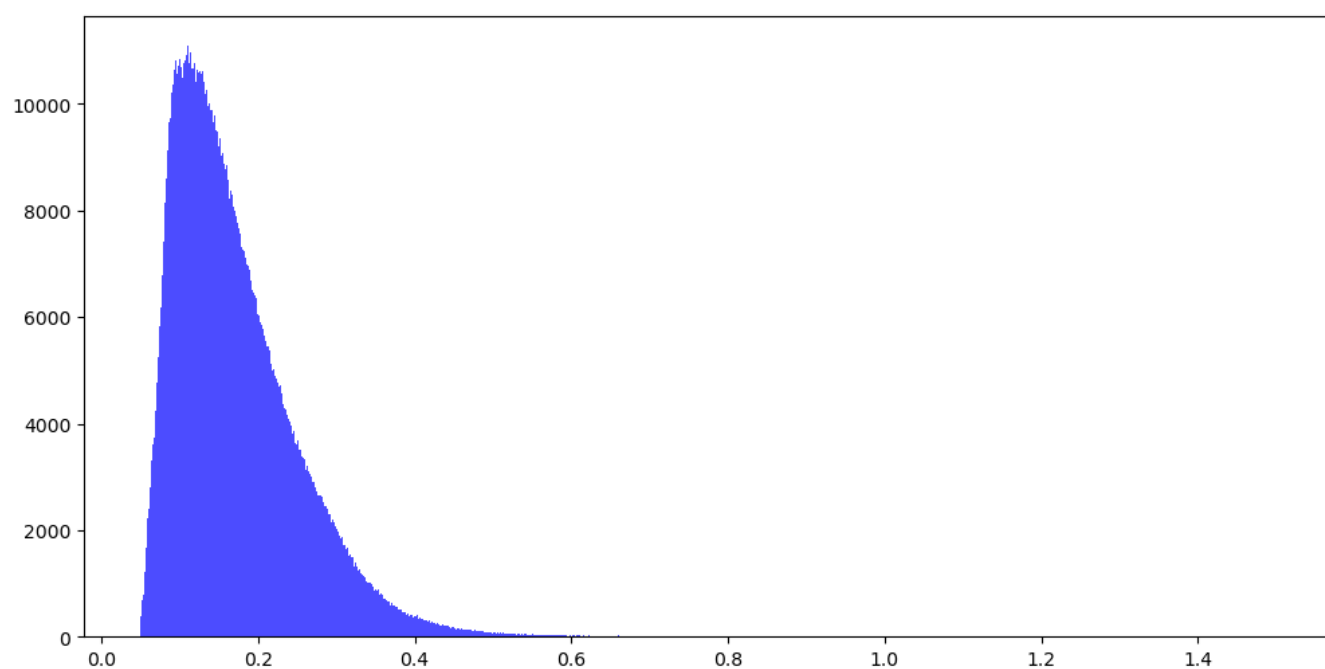
	optiondate	maturity	strike_price	moneyness	impl_volatility
0	736333	11	1875	0.931603	0.149723
1	736333	11	1880	0.934087	0.184999
2	736333	11	1890	0.939056	0.203981
3	736333	11	1900	0.944024	0.211677
4	736333	11	1910	0.948993	0.217900

```
final_df['impl_volatility'].describe()
```

```
count    1.772317e+06  
mean      1.717228e-01  
std       8.268354e-02  
min       5.000100e-02  
25%       1.113900e-01  
50%       1.536130e-01  
75%       2.132350e-01  
max       1.498417e+00  
Name: impl_volatility, dtype: float64
```

```
import matplotlib.pyplot as plt  
import seaborn as sns
```

```
plt.figure(figsize=(12, 6))  
plt.hist(final_df['impl_volatility'], bins=1500, color='blue', alpha=0.7)  
plt.show()
```

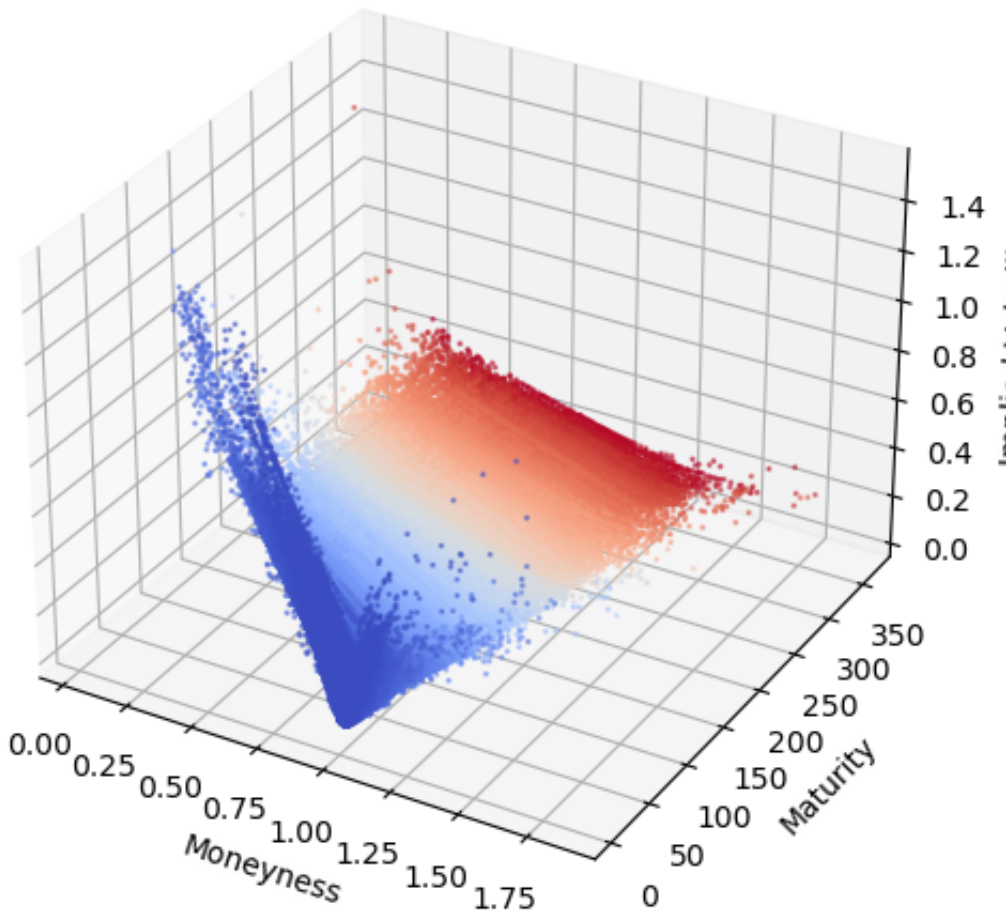


```
final_df.sort_values(['optiondate', 'maturity', 'strike_price', 'moneyness'], i
final_df.head()
```

	optiondate	maturity	strike_price	moneyness	impl_volatility
40	736333	11	1710	0.849622	0.439677
41	736333	11	1715	0.852106	0.432749
42	736333	11	1725	0.857075	0.421268
43	736333	11	1730	0.859559	0.414337
44	736333	11	1745	0.867012	0.397865

```
fig = plt.figure(figsize=(12, 6))
ax = fig.add_subplot(111, projection='3d')
ax.scatter(final_df['moneyness'], final_df['maturity'], final_df['impl_volatili
ax.set_xlabel('Moneyness')
ax.set_ylabel('Maturity')
ax.set_zlabel('Implied Volatility')

plt.show()
```



```
from google.colab import drive
import pandas as pd
import numpy as np

drive.mount('/content/drive')

final_df.to_csv('/content/drive/My Drive/Term 3/Financial Econ/final_df.csv')
```

## ✓ Lets start training the NN

```
final_df = pd.read_csv('/content/drive/My Drive/Term 3/Financial Econ/final_df.csv')

Mounted at /content/drive
```

## ✓ 1 day lag

```
import torch
import torch.nn as nn
import torch.optim as optim
from torch.utils.data import TensorDataset, DataLoader

torch.manual_seed(42)
if torch.cuda.is_available():
    torch.cuda.manual_seed_all(42)

class MyModel(nn.Module):
    def __init__(self, input_size, layer_sizes):
        super(MyModel, self).__init__()
        layers = []
        layers.append(nn.Linear(input_size, layer_sizes[0]))
        layers.append(nn.ReLU())
        for i in range(len(layer_sizes) - 1):
            layers.append(nn.Linear(layer_sizes[i], layer_sizes[i+1]))
            layers.append(nn.ReLU())
        self.model = nn.Sequential(*layers)
        self.output_layer = nn.Linear(layer_sizes[-1], 1)

    def forward(self, x):
        x = self.model(x)
        x = self.output_layer(x)
        return x.squeeze(1)
```

```
final_df.head()
```

	optiondate	maturity	strike_price	moneyness	impl_volatility	impl_volat
0	736333	11	1710	0.849622	0.439677	
1	736333	11	1715	0.852106	0.432749	
2	736333	11	1725	0.857075	0.421268	
3	736333	11	1730	0.859559	0.414337	
4	736333	11	1745	0.867012	0.397865	

```
inplace_final_df = final_df[['optiondate', 'strike_price', 'moneyness', 'maturity', 'impl_volatility']]
date_group = inplace_final_df.groupby('optiondate')
inplace_final_df['bs_vol_daily'] = date_group['impl_volatility'].transform('mean')
inplace_final_df['error_vol'] = inplace_final_df['impl_volatility']-inplace_final_df['bs_vol_daily']
inplace_final_df.head()
```

	optiondate	strike_price	moneyness	maturity	impl_volatility	bs_vol_daily
0	736333	1710	0.849622	11	0.439677	0.225
1	736333	1715	0.852106	11	0.432749	0.225
2	736333	1725	0.857075	11	0.421268	0.225
3	736333	1730	0.859559	11	0.414337	0.225
4	736333	1745	0.867012	11	0.397865	0.225

```
device = torch.device("cuda" if torch.cuda.is_available() else "cpu")
device

model = nn.Sequential(
    nn.Linear(2, 10),
    nn.ReLU(),
    nn.Linear(10, 1)
).to(device(type='cpu'))

target = 'error_vol'
features = ['moneyness', 'maturity']
```

✓ No NN

0 day

```

date_groups = inplace_final_df.groupby('optiondate')
date_keys = list(date_groups.groups.keys())
abs_error = []
preds = []
total_obs = 0
for i in range(len(date_keys)):
    group_train = inplace_final_df.iloc[date_groups.groups[date_keys[i]]]
    avg_y = group_train[target].mean()
    preds = [avg_y]*group_train.shape[0]
    obs = group_train[target].tolist()
    abs_error += np.abs(np.array(preds) - np.array(obs)).tolist()
    total_obs += len(obs)

print(f'RMSE of No NN: {np.sqrt(np.mean(np.square(abs_error)))*100}')

```

RMSE of No NN: 7.3736693813654375

1 day

```

date_groups = inplace_final_df.groupby('optiondate')
date_keys = list(date_groups.groups.keys())
abs_error = []
preds = []
total_obs = 0
for i in range(len(date_keys)-1):
    group_train = inplace_final_df.iloc[date_groups.groups[date_keys[i]]]
    group_test = inplace_final_df.iloc[date_groups.groups[date_keys[i+1]]]
    avg_y = group_train[target].mean()
    preds = [avg_y]*group_test.shape[0]
    obs = group_test[target].tolist()
    abs_error += np.abs(np.array(preds) - np.array(obs)).tolist()
    total_obs += len(obs)

print(f'RMSE of No NN: {np.sqrt(np.mean(np.square(abs_error)))*100}')

```

RMSE of No NN: 7.373097588564792

## ✓ Creating a generic method that will take model as input

```

def train_test(dataframe, feature_list, target_column, model_layer_sizes, epoch
    model = MyModel(input_size=len(features), layer_sizes=model_layer_sizes)
    model = model.to(device)

    date_groups = dataframe.groupby('optiondate')
    date_keys = list(date_groups.groups.keys())
    abs_error = []
    total_obs = 0

```

```

criterion = nn.MSELoss()
optimizer = optim.Adam(model.parameters())

for i in range(len(date_keys)-1):
    group_train = dataframe.iloc[date_groups.groups[date_keys[i]]]
    group_test = dataframe.iloc[date_groups.groups[date_keys[i+1]]]
    X_train = group_train[feature_list]
    y_train = group_train[target_column]
    X_test = group_test[feature_list]
    y_test = group_test[target_column]
    X_train_tensor = torch.tensor(X_train.values, dtype=torch.float32)
    y_train_tensor = torch.tensor(y_train.values, dtype=torch.float32)
    X_test_tensor = torch.tensor(X_test.values, dtype=torch.float32)
    y_test_tensor = torch.tensor(y_test.values, dtype=torch.float32)

    train_dataset = TensorDataset(X_train_tensor, y_train_tensor)
    train_loader = DataLoader(train_dataset, batch_size=32)
    test_dataset = TensorDataset(X_test_tensor, y_test_tensor)
    test_loader = DataLoader(test_dataset, batch_size=32)

    for epoch in range(epochs):
        model.train()
        running_loss = 0.0
        for inputs, labels in train_loader:
            inputs, labels = inputs.to(device), labels.to(device)
            optimizer.zero_grad()
            outputs = model(inputs)
            loss = criterion(outputs, labels)
            loss.backward()
            optimizer.step()
            running_loss += loss.item() * inputs.size(0)
        # print(f"{date_keys[i+1]}: ----> Epoch {epoch+1}/{epochs}, Loss: {running_loss}")

    model.eval()
    preds = []
    obs = []
    for inputs, labels in test_loader:
        inputs, labels = inputs.to(device), labels.to(device)
        outputs = model(inputs)
        preds += outputs.tolist()
        obs += labels.tolist()

    abs_error += np.abs(np.array(preds) - np.array(obs)).tolist()
rmse = np.sqrt(np.mean(np.square(abs_error)))*100
print(f'RMSE of NN{len(model_layer_sizes)}: {rmse}')
return rmse

```



```
train_test(inplace_final_df, features, target, [32], 2)
```

```
RMSE of NN1: 13.484864229159394  
13.484864229159394
```

```
train_test(inplace_final_df, features, target, [32, 16], 10)  
train_test(inplace_final_df, features, target, [32, 16, 8], 10)  
train_test(inplace_final_df, features, target, [32, 16, 8, 4], 10)  
train_test(inplace_final_df, features, target, [32, 16, 8, 4, 2], 10)
```

```
RMSE of NN2: 4.227207829378143  
RMSE of NN3: 3.642886578882718  
RMSE of NN4: 7.378633357171569  
RMSE of NN5: 7.375819690945528  
7.375819690945528
```