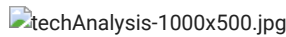


## ✓ Data Project - Stock Market Analysis



Time Series data is a series of data points indexed in time order. Time series data is everywhere, so manipulating them is important for any data analyst or data scientist.

In this notebook, we will discover and explore data from the stock market, particularly some technology stocks (Apple, Amazon, Google, and Microsoft). We will learn how to use yfinance to get stock information, and visualize different aspects of it using Seaborn and Matplotlib. We will look at a few ways of analyzing the risk of a stock, based on its previous performance history. We will also be predicting future stock prices through a Long Short Term Memory (LSTM) method!

We'll be answering the following questions along the way:

- 1.) What was the change in price of the stock over time?
- 2.) What was the daily return of the stock on average?
- 3.) What was the moving average of the various stocks?
- 4.) What was the correlation between different stocks'?
- 5.) How much value do we put at risk by investing in a particular stock?
- 6.) How can we attempt to predict future stock behavior? (Predicting the closing price stock price of APPLE inc using LSTM)

---

### Getting the Data

The first step is to get the data and load it to memory. We will get our stock data from the Yahoo Finance website. Yahoo Finance is a rich resource of financial market data and tools to find compelling investments. To get the data from Yahoo Finance, we will be using yfinance library which offers a threaded and Pythonic way to download market data from Yahoo. Check this article to learn more about yfinance: [Reliably download historical market data from with Python](#)

## ✓ 1. What was the change in price of the stock overtime?

In this section we'll go over how to handle requesting stock information with pandas, and how to analyze basic attributes of a stock.

```
!pip install -q yfinance
```

```
import pandas as pd
import numpy as np

import matplotlib.pyplot as plt
import seaborn as sns
sns.set_style('whitegrid')
plt.style.use("fivethirtyeight")
%matplotlib inline

# For reading stock data from yahoo
from pandas_datareader.data import DataReader
import yfinance as yf
from pandas_datareader import data as pdr

yf.pdr_override()

# For time stamps
from datetime import datetime

# The tech stocks we'll use for this analysis
tech_list = ['AAPL', 'GOOG', 'MSFT', 'AMZN']

# Set up End and Start times for data grab
tech_list = ['AAPL', 'GOOG', 'MSFT', 'AMZN']

end = datetime.now()
start = datetime(end.year - 1, end.month, end.day)

for stock in tech_list:
    globals()[stock] = yf.download(stock, start, end)

company_list = [AAPL, GOOG, MSFT, AMZN]
company_name = ["APPLE", "GOOGLE", "MICROSOFT", "AMAZON"]

for company, com_name in zip(company_list, company_name):
    company["company_name"] = com_name

df = pd.concat(company_list, axis=0)
df.tail(10)
```

/usr/local/lib/python3.10/dist-packages/yfinance/base.py:48: FutureWarning: The default dtype for empty Series will be 'object' inst

\_empty\_series = pd.Series()

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

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

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

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

	Open	High	Low	Close	Adj Close	Volume	company_name
Date							
2024-01-22	156.889999	157.050003	153.899994	154.779999	154.779999	43687500	AMAZON
2024-01-23	154.850006	156.210007	153.929993	156.020004	156.020004	37986000	AMAZON
2024-01-24	157.800003	158.509995	156.479996	156.869995	156.869995	48547300	AMAZON
2024-01-25	156.949997	158.509995	154.550003	157.750000	157.750000	43638600	AMAZON
2024-01-26	158.419998	160.720001	157.910004	159.119995	159.119995	51047400	AMAZON
2024-01-29	159.339996	161.289993	158.899994	161.259995	161.259995	45270400	AMAZON
2024-01-30	160.699997	161.729996	158.490005	159.000000	159.000000	45207400	AMAZON
2024-01-31	157.000000	159.009995	154.809998	155.199997	155.199997	50284400	AMAZON
2024-02-01	155.869995	159.759995	155.619995	159.279999	159.279999	76542400	AMAZON
2024-02-02	169.190002	172.500000	167.330002	171.809998	171.809998	117154900	AMAZON

Reviewing the content of our data, we can see that the data is numeric and the date is the index of the data. Notice also that weekends are missing from the records.

**Quick note:** Using `globals()` is a sloppy way of setting the `DataFrame` names, but it's simple. Now we have our data, let's perform some basic data analysis and check our data.

✓ Descriptive Statistics about the Data

`.describe()` generates descriptive statistics. Descriptive statistics include those that summarize the central tendency, dispersion, and shape of a dataset's distribution, excluding NaN values.

Analyzes both numeric and object series, as well as `DataFrame` column sets of mixed data types. The output will vary depending on what is provided. Refer to the notes below for more detail.

```
# Summary Stats
AAPL.describe()
```

	Open	High	Low	Close	Adj Close	Volume
count	250.000000	250.000000	250.000000	250.000000	250.000000	2.500000e+02
mean	176.896561	178.428160	175.618280	177.143360	176.785870	5.745947e+07
std	13.677455	13.528870	13.621885	13.551779	13.714590	1.591402e+07
min	144.380005	146.710007	143.899994	145.309998	144.722931	2.404830e+07
25%	169.342499	170.904995	167.684998	169.612495	168.927254	4.735310e+07
50%	178.035004	179.709999	176.969994	178.290001	178.050423	5.381205e+07
75%	189.140003	189.964996	187.570004	189.095001	188.717094	6.405040e+07
max	198.020004	199.619995	197.000000	198.110001	198.110001	1.282567e+08

We have only 255 records in one year because weekends are not included in the data.

### Information About the Data

`.info()` method prints information about a `DataFrame` including the index `dtype` and columns, non-null values, and memory usage.

```
# General info
AAPL.info()

<class 'pandas.core.frame.DataFrame'>
DatetimeIndex: 250 entries, 2023-02-06 to 2024-02-02
Data columns (total 7 columns):
#   Column          Non-Null Count  Dtype
---  ---
0   Open            250 non-null    float64
1   High            250 non-null    float64
2   Low             250 non-null    float64
3   Close           250 non-null    float64
4   Adj Close       250 non-null    float64
5   Volume          250 non-null    int64
6   company_name    250 non-null    object
dtypes: float64(5), int64(1), object(1)
memory usage: 15.6+ KB
```

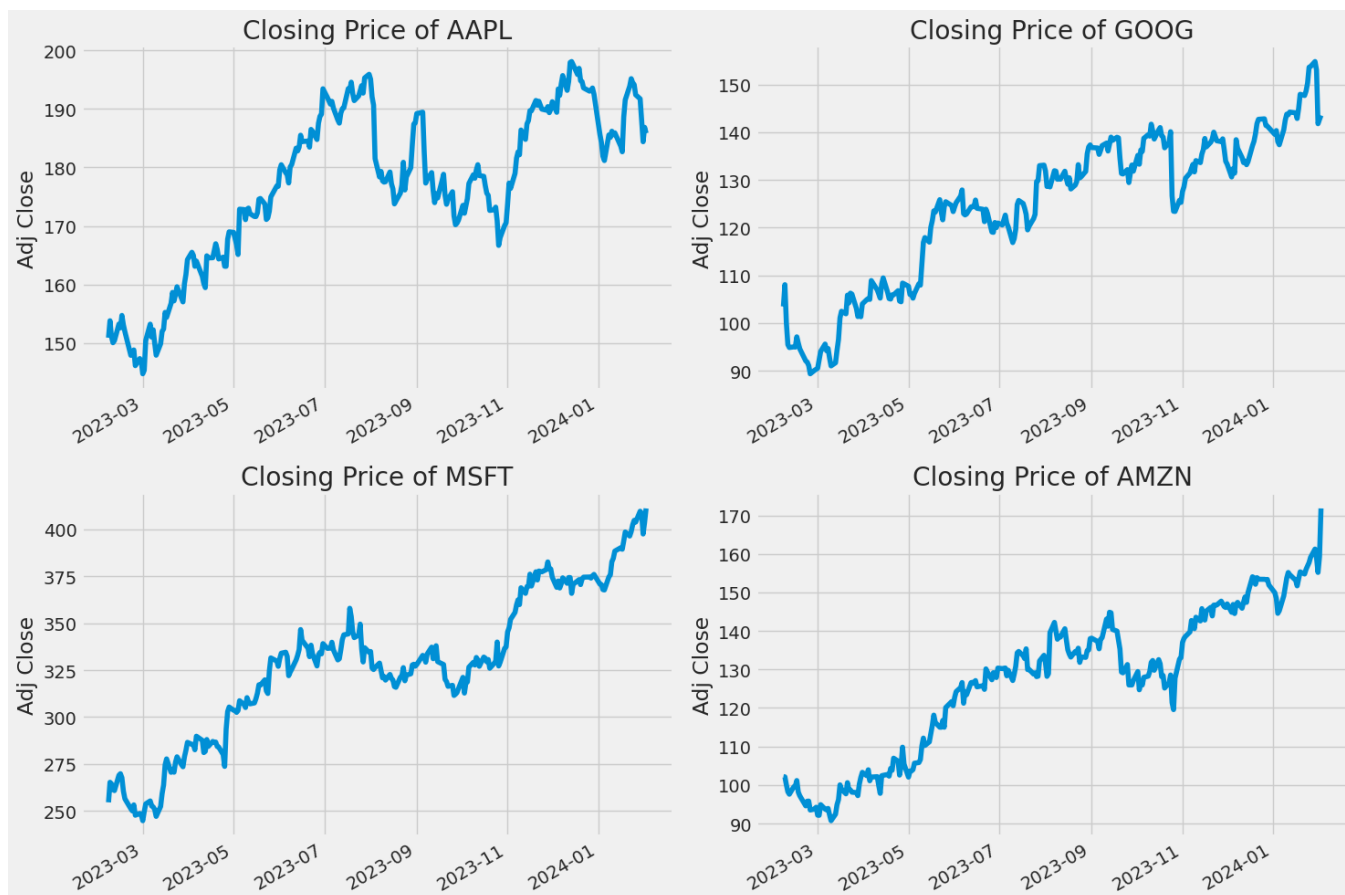
### Closing Price

The closing price is the last price at which the stock is traded during the regular trading day. A stock's closing price is the standard benchmark used by investors to track its performance over time.

```
# Let's see a historical view of the closing price
plt.figure(figsize=(15, 10))
plt.subplots_adjust(top=1.25, bottom=1.2)

for i, company in enumerate(company_list, 1):
    plt.subplot(2, 2, i)
    company['Adj Close'].plot()
    plt.ylabel('Adj Close')
    plt.xlabel(None)
    plt.title(f"Closing Price of {tech_list[i - 1]}")

plt.tight_layout()
```



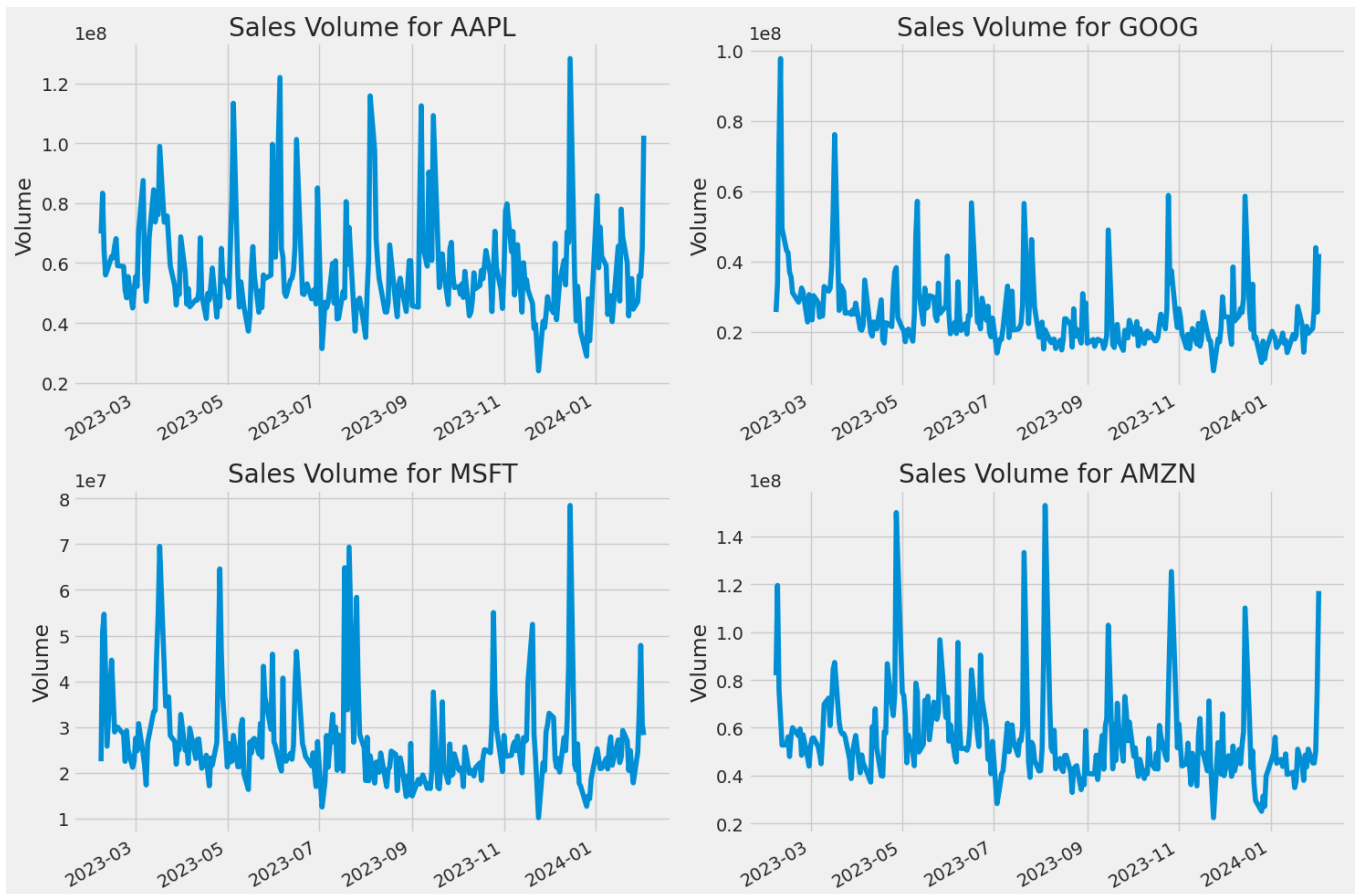
## Volume of Sales

Volume is the amount of an asset or security that changes hands over some period of time, often over the course of a day. For instance, the stock trading volume would refer to the number of shares of security traded between its daily open and close. Trading volume, and changes to volume over the course of time, are important inputs for technical traders.

```
# Now let's plot the total volume of stock being traded each day
plt.figure(figsize=(15, 10))
plt.subplots_adjust(top=1.25, bottom=1.2)

for i, company in enumerate(company_list, 1):
    plt.subplot(2, 2, i)
    company['Volume'].plot()
    plt.ylabel('Volume')
    plt.xlabel(None)
    plt.title(f"Sales Volume for {tech_list[i - 1]}")

plt.tight_layout()
```



Now that we've seen the visualizations for the closing price and the volume traded each day, let's go ahead and calculate the moving average for the stock.

## ✓ 2. What was the moving average of the various stocks?

The moving average (MA) is a simple technical analysis tool that smooths out price data by creating a constantly updated average price. The average is taken over a specific period of time, like 10 days, 20 minutes, 30 weeks, or any time period the trader chooses.

```
ma_day = [10, 20, 50]

for ma in ma_day:
    for company in company_list:
        column_name = f"MA for {ma} days"
        company[column_name] = company['Adj Close'].rolling(ma).mean()

fig, axes = plt.subplots(nrows=2, ncols=2)
fig.set_figheight(10)
fig.set_figwidth(15)

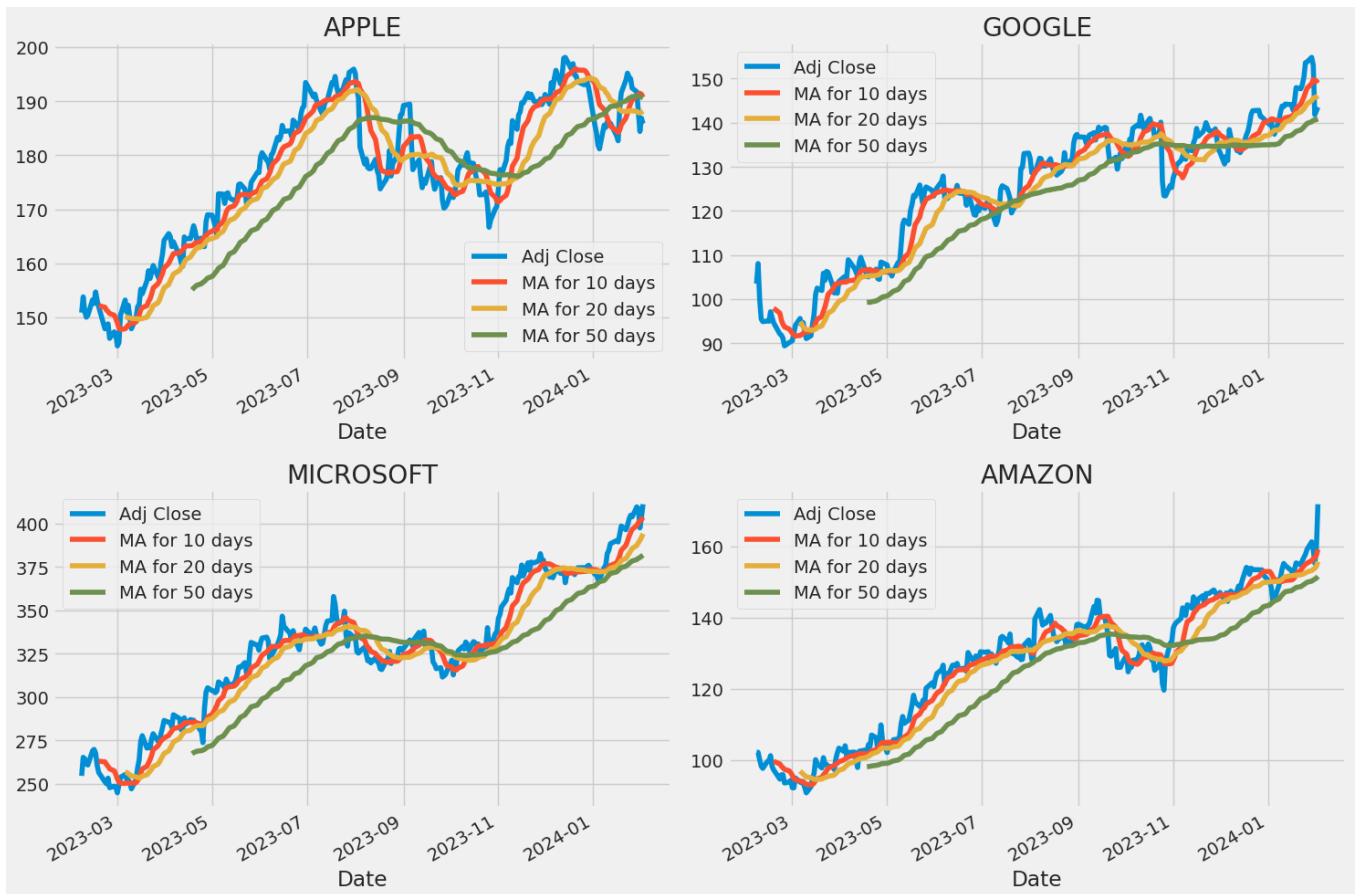
AAPL[['Adj Close', 'MA for 10 days', 'MA for 20 days', 'MA for 50 days']].plot(ax=axes[0,0])
axes[0,0].set_title('APPLE')

GOOG[['Adj Close', 'MA for 10 days', 'MA for 20 days', 'MA for 50 days']].plot(ax=axes[0,1])
axes[0,1].set_title('GOOGLE')

MSFT[['Adj Close', 'MA for 10 days', 'MA for 20 days', 'MA for 50 days']].plot(ax=axes[1,0])
axes[1,0].set_title('MICROSOFT')

AMZN[['Adj Close', 'MA for 10 days', 'MA for 20 days', 'MA for 50 days']].plot(ax=axes[1,1])
axes[1,1].set_title('AMAZON')

fig.tight_layout()
```



We see in the graph that the best values to measure the moving average are 10 and 20 days because we still capture trends in the data without noise.

### ✓ 3. What was the daily return of the stock on average?

Now that we've done some baseline analysis, let's go ahead and dive a little deeper. We're now going to analyze the risk of the stock. In order to do so we'll need to take a closer look at the daily changes of the stock, and not just its absolute value. Let's go ahead and use pandas to retrieve the daily returns for the Apple stock.

```
# We'll use pct_change to find the percent change for each day
for company in company_list:
    company['Daily Return'] = company['Adj Close'].pct_change()

# Then we'll plot the daily return percentage
fig, axes = plt.subplots(nrows=2, ncols=2)
fig.set_figheight(10)
fig.set_figwidth(15)

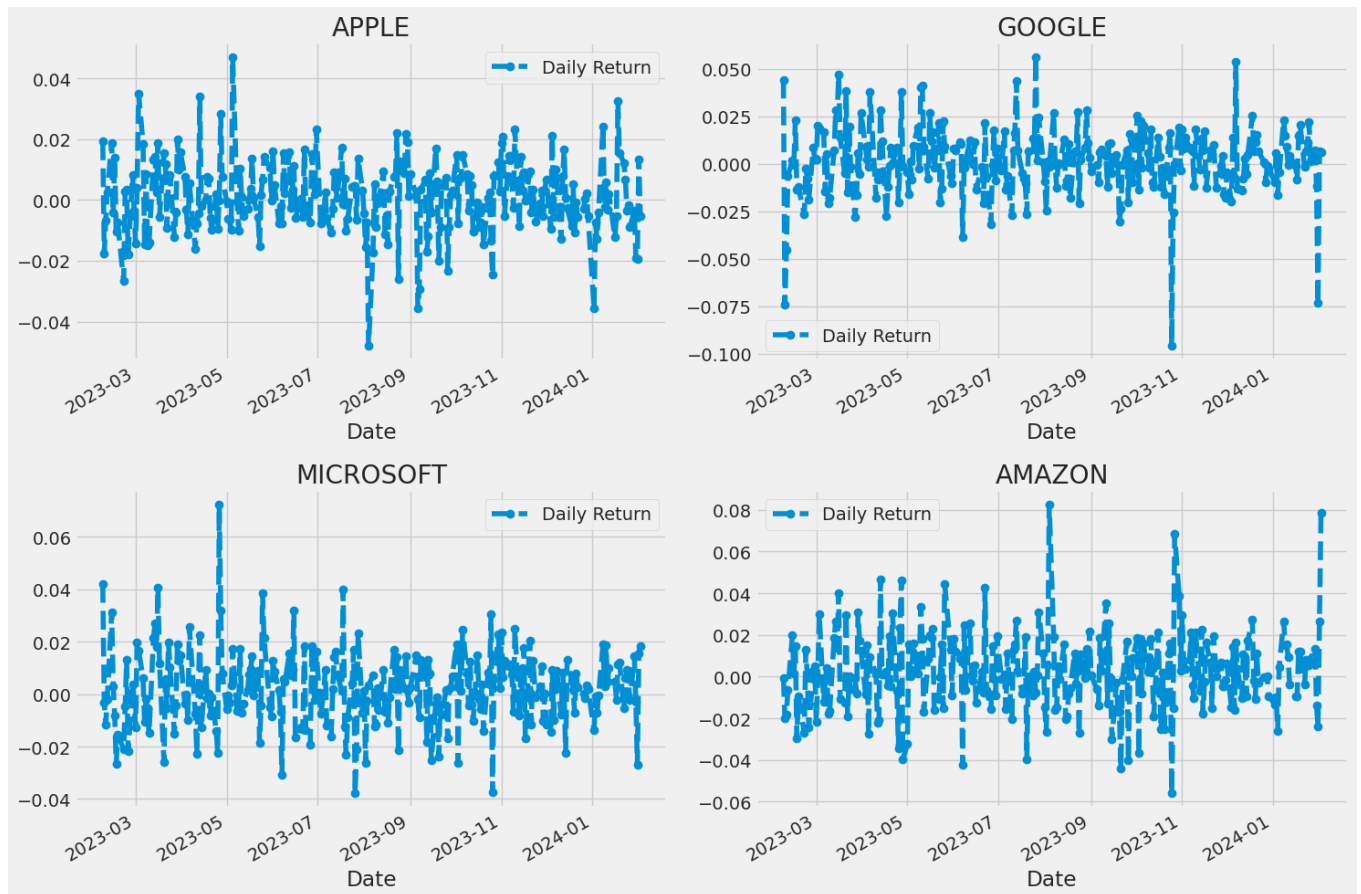
AAPL['Daily Return'].plot(ax=axes[0,0], legend=True, linestyle='--', marker='o')
axes[0,0].set_title('APPLE')

GOOG['Daily Return'].plot(ax=axes[0,1], legend=True, linestyle='--', marker='o')
axes[0,1].set_title('GOOGLE')

MSFT['Daily Return'].plot(ax=axes[1,0], legend=True, linestyle='--', marker='o')
axes[1,0].set_title('MICROSOFT')

AMZN['Daily Return'].plot(ax=axes[1,1], legend=True, linestyle='--', marker='o')
axes[1,1].set_title('AMAZON')

fig.tight_layout()
```

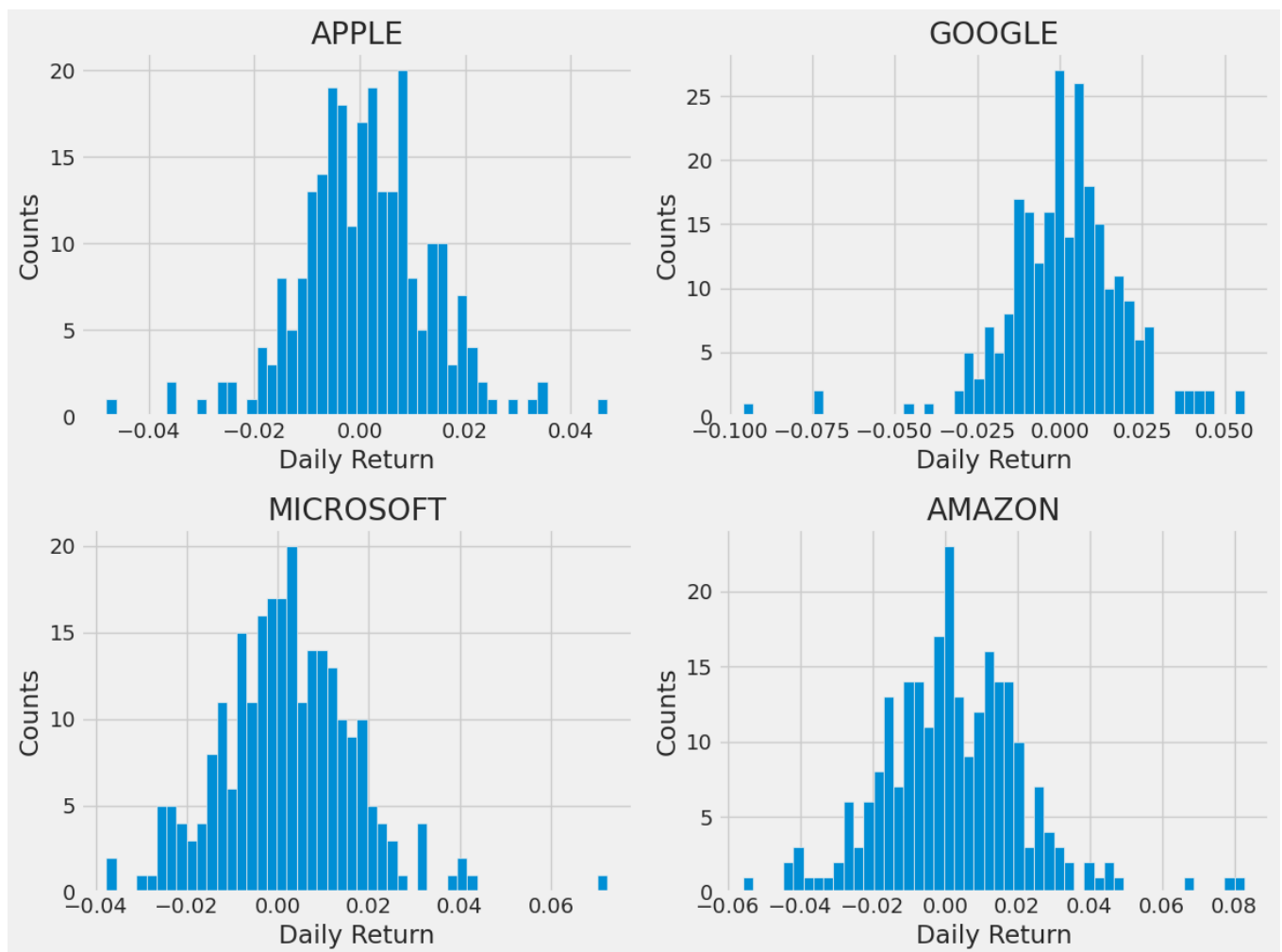


Great, now let's get an overall look at the average daily return using a histogram. We'll use seaborn to create both a histogram and kde plot on the same figure.

```
plt.figure(figsize=(12, 9))

for i, company in enumerate(company_list, 1):
    plt.subplot(2, 2, i)
    company['Daily Return'].hist(bins=50)
    plt.xlabel('Daily Return')
    plt.ylabel('Counts')
    plt.title(f'{company_name[i - 1]}')

plt.tight_layout()
```



#### 4. What was the correlation between different stocks closing prices?

Correlation is a statistic that measures the degree to which two variables move in relation to each other which has a value that must fall between -1.0 and +1.0. Correlation measures association, but doesn't show if x causes y or vice versa – or if the association is caused by a third factor[1].

Now what if we wanted to analyze the returns of all the stocks in our list? Let's go ahead and build a DataFrame with all the ['Close'] columns for each of the stocks dataframes.

```
# Grab all the closing prices for the tech stock list into one DataFrame

closing_df = pdr.get_data_yahoo(tech_list, start=start, end=end)['Adj Close']

# Make a new tech returns DataFrame
tech_rets = closing_df.pct_change()
tech_rets.head()
```

[\*\*\*\*\*100%\*\*\*\*\*] 4 of 4 completed

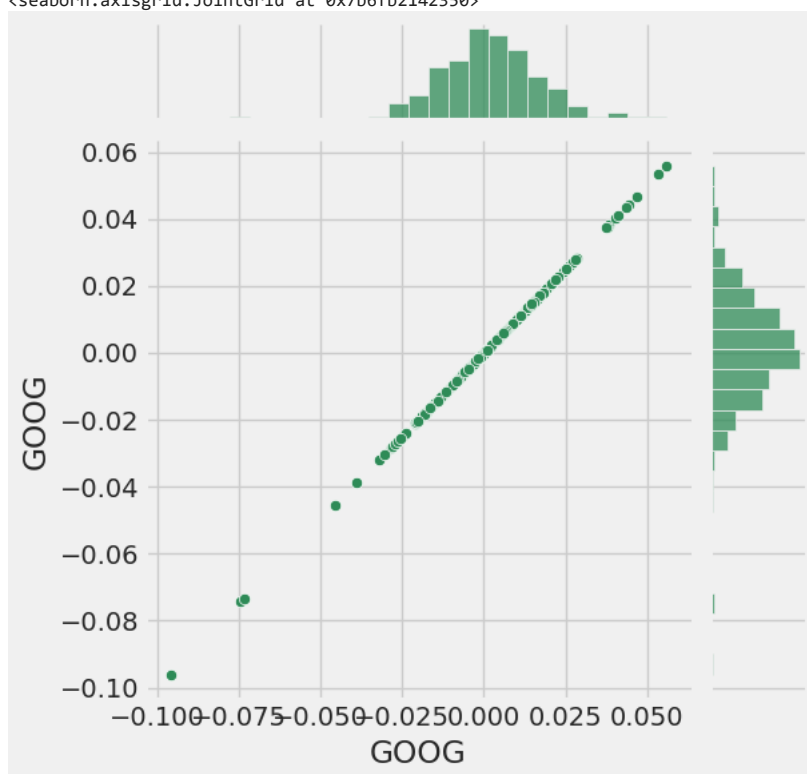
Ticker	AAPL	AMZN	GOOG	MSFT
Date				
2023-02-06	NaN	NaN	NaN	NaN
2023-02-07	0.019245	-0.000685	0.044167	0.042022
2023-02-08	-0.017653	-0.020174	-0.074417	-0.003102
2023-02-09	-0.006912	-0.018091	-0.045400	-0.011660
2023-02-10	0.002456	-0.006413	-0.006285	-0.001972



Now we can compare the daily percentage return of two stocks to check how correlated. First let's see a sotck compared to itself.

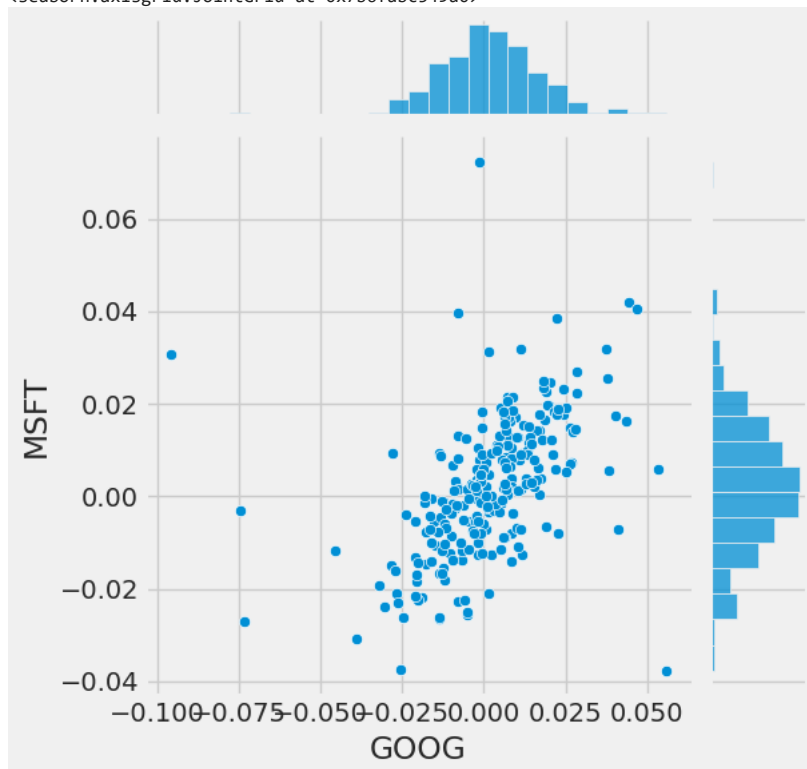
```
# Comparing Google to itself should show a perfectly linear relationship
sns.jointplot(x='GOOG', y='GOOG', data=tech_rets, kind='scatter', color='seagreen')
```

<seaborn.axisgrid.JointGrid at 0x7b6fb2142350>



```
# We'll use jointplot to compare the daily returns of Google and Microsoft
sns.jointplot(x='GOOG', y='MSFT', data=tech_rets, kind='scatter')
```

<seaborn.axisgrid.JointGrid at 0x7b6fab549a0>



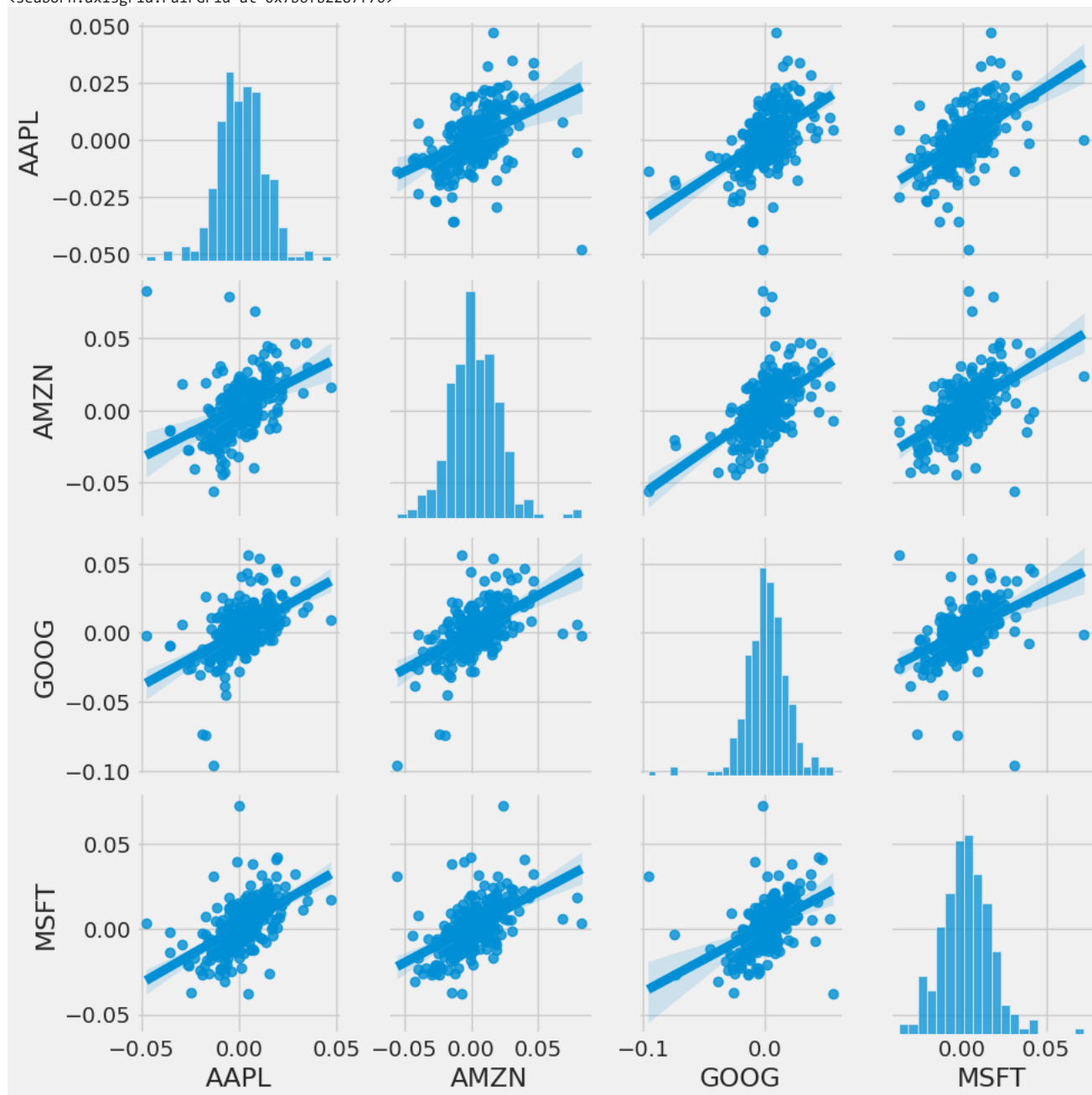
So now we can see that if two stocks are perfectly (and positivley) correlated with each other a linear relationship bewteen its daily return values should occur.

Seaborn and pandas make it very easy to repeat this comparison analysis for every possible combination of stocks in our technology stock ticker list. We can use `sns.pairplot()` to automatically create this plot

```
# We can simply call pairplot on our DataFrame for an automatic visual analysis
# of all the comparisons
```

```
sns.pairplot(tech_rets, kind='reg')
```

```
<seaborn.axisgrid.PairGrid at 0x7b6fb2287f70>
```



Above we can see all the relationships on daily returns between all the stocks. A quick glance shows an interesting correlation between Google and Amazon daily returns. It might be interesting to investigate that individual comparison.

While the simplicity of just calling `sns.pairplot()` is fantastic we can also use `sns.PairGrid()` for full control of the figure, including what kind of plots go in the diagonal, the upper triangle, and the lower triangle. Below is an example of utilizing the full power of seaborn to achieve this result.

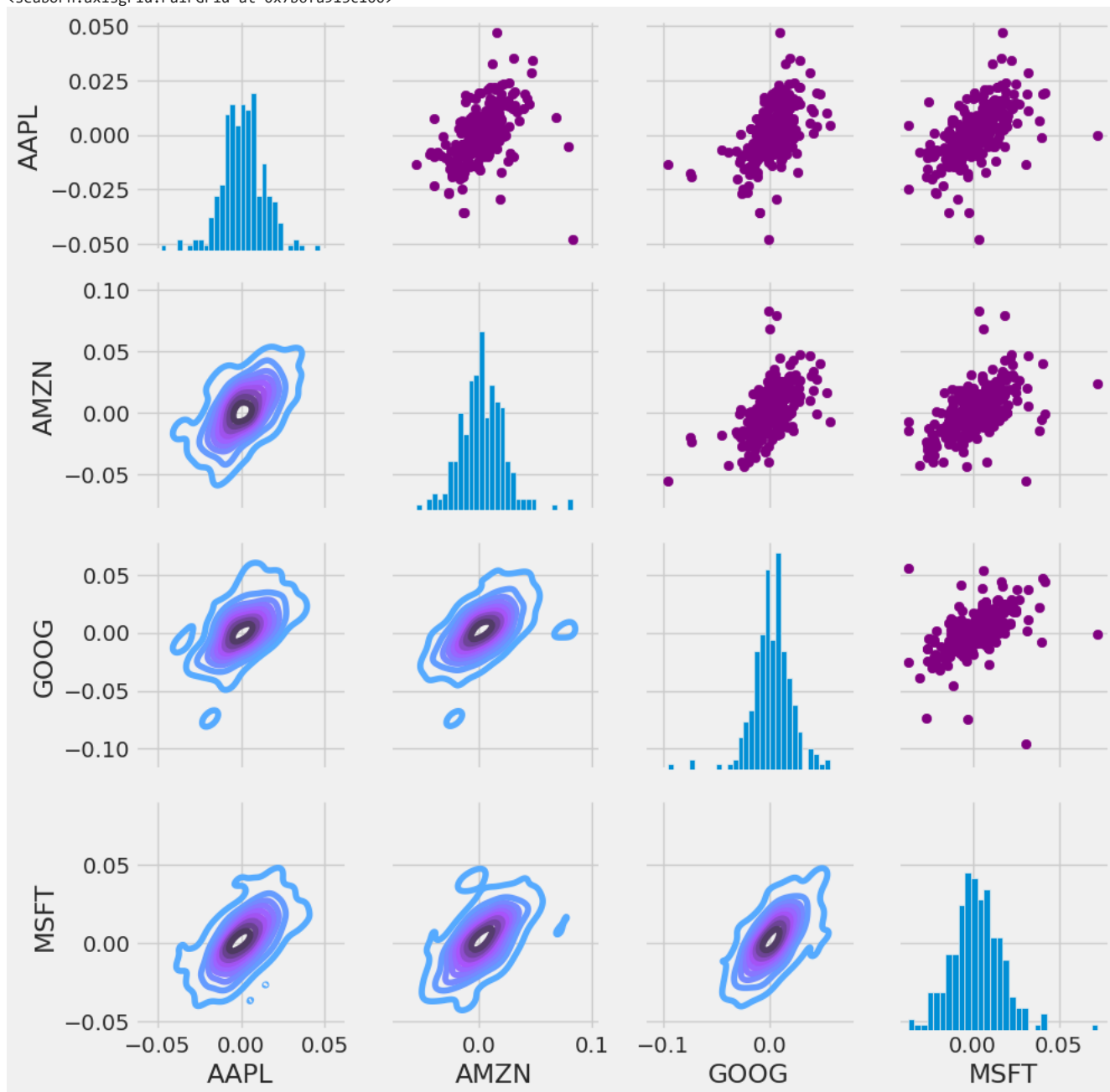
```
# Set up our figure by naming it returns_fig, call PairPlot on the DataFrame
return_fig = sns.PairGrid(tech_rets.dropna())
```

```
# Using map_upper we can specify what the upper triangle will look like.
return_fig.map_upper(plt.scatter, color='purple')
```

```
# We can also define the lower triangle in the figure, including the plot type (kde)
# or the color map (BluePurple)
return_fig.map_lower(sns.kdeplot, cmap='cool_d')
```

```
# Finally we'll define the diagonal as a series of histogram plots of the daily return
return_fig.map_diag(plt.hist, bins=30)
```

&lt;seaborn.axisgrid.PairGrid at 0x7b6fa913c100&gt;



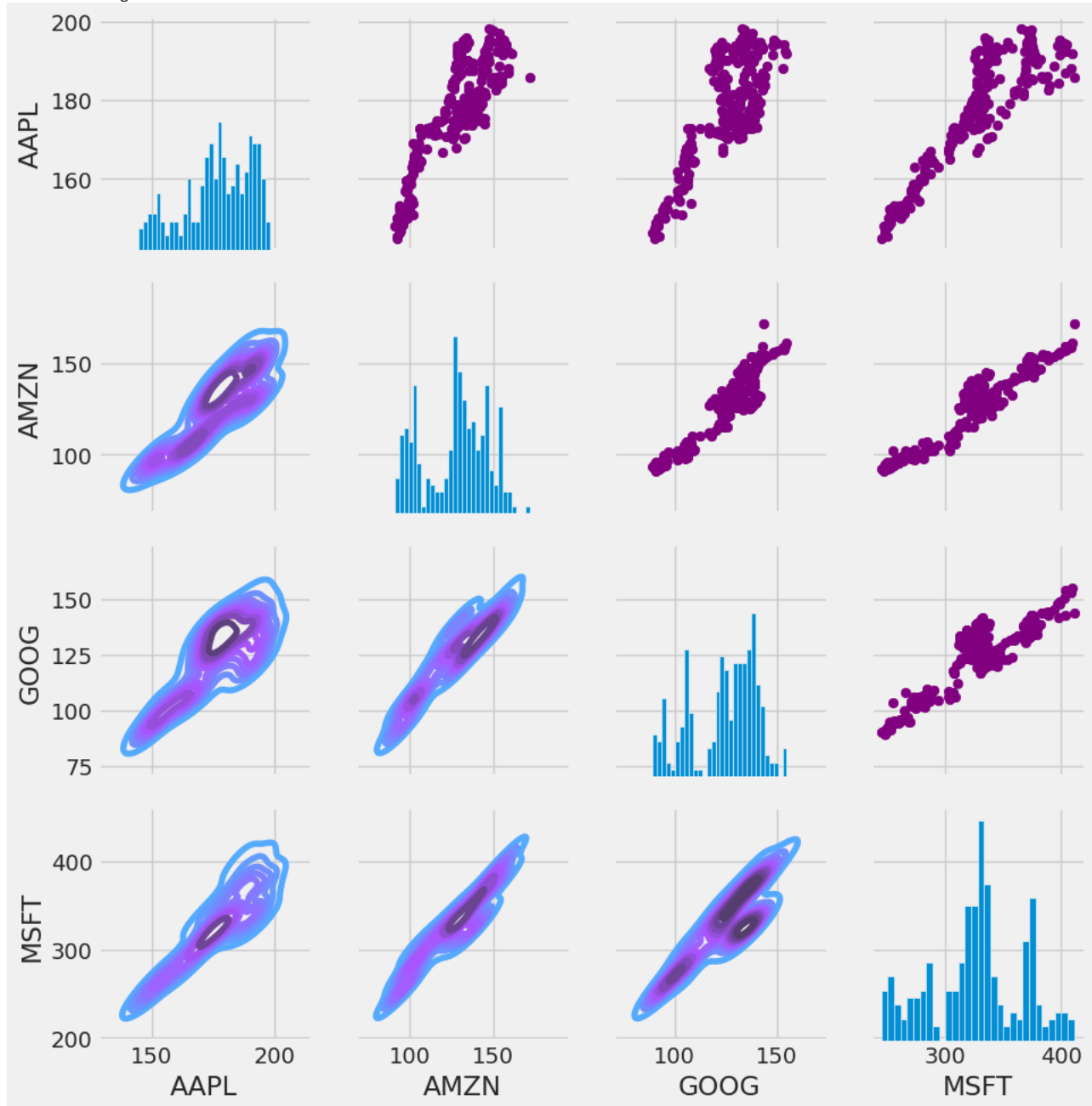
```
# Set up our figure by naming it returns_fig, call PairPlot on the DataFrame
returns_fig = sns.PairGrid(closing_df)
```

```
# Using map_upper we can specify what the upper triangle will look like.
returns_fig.map_upper(plt.scatter,color='purple')
```

```
# We can also define the lower triangle in the figure, including the plot type (kde) or the color map (BluePurple)
returns_fig.map_lower(sns.kdeplot,cmap='cool_d')
```

```
# Finally we'll define the diagonal as a series of histogram plots of the daily return
returns_fig.map_diag(plt.hist,bins=30)
```

&lt;seaborn.axisgrid.PairGrid at 0x7b6fa87baa10&gt;

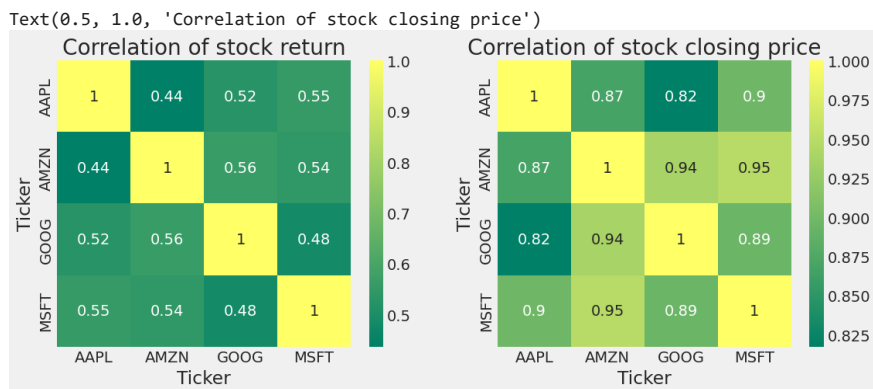


Finally, we could also do a correlation plot, to get actual numerical values for the correlation between the stocks' daily return values. By comparing the closing prices, we see an interesting relationship between Microsoft and Apple.

```
plt.figure(figsize=(12, 10))

plt.subplot(2, 2, 1)
sns.heatmap(tech_rets.corr(), annot=True, cmap='summer')
plt.title('Correlation of stock return')

plt.subplot(2, 2, 2)
sns.heatmap(closing_df.corr(), annot=True, cmap='summer')
plt.title('Correlation of stock closing price')
```



Just like we suspected in our `PairPlot` we see here numerically and visually that Microsoft and Amazon had the strongest correlation of daily stock return. It's also interesting to see that all the technology companies are positively correlated.

## ✓ 5. How much value do we put at risk by investing in a particular stock?

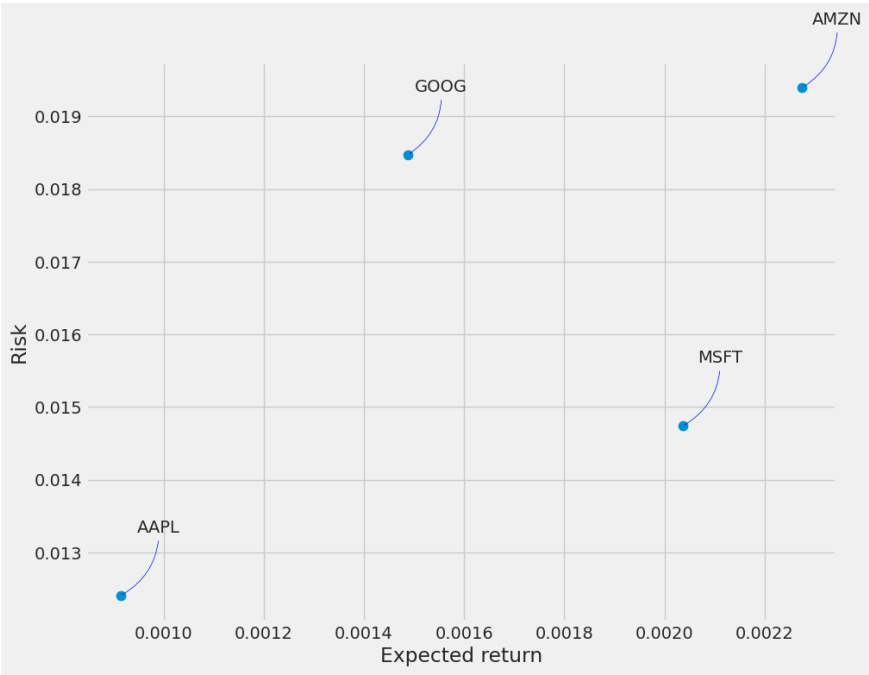
There are many ways we can quantify risk, one of the most basic ways using the information we've gathered on daily percentage returns is by comparing the expected return with the standard deviation of the daily returns.

```
rets = tech_rets.dropna()

area = np.pi * 20

plt.figure(figsize=(10, 8))
plt.scatter(rets.mean(), rets.std(), s=area)
plt.xlabel('Expected return')
plt.ylabel('Risk')

for label, x, y in zip(rets.columns, rets.mean(), rets.std()):
    plt.annotate(label, xy=(x, y), xytext=(50, 50), textcoords='offset points', ha='right', va='bottom',
                 arrowprops=dict(arrowstyle='-', color='blue', connectionstyle='arc3,rad=-0.3'))
```



6. Predicting the closing price stock price of APPLE inc:

```
# Get the stock quote
df = pdr.get_data_yahoo('AAPL', start='2012-01-01', end=datetime.now())
# Show teh data
df
```

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

	Open	High	Low	Close	Adj Close	Volume	
Date							
2012-01-03	14.621429	14.732143	14.607143	14.686786	12.449689	302220800	
2012-01-04	14.642857	14.810000	14.617143	14.765714	12.516596	260022000	
2012-01-05	14.819643	14.948214	14.738214	14.929643	12.655560	271269600	
2012-01-06	14.991786	15.098214	14.972143	15.085714	12.787855	318292800	
2012-01-09	15.196429	15.276786	15.048214	15.061786	12.767571	394024400	
...	...	...	...	...	...	...	
2024-01-29	192.009995	192.199997	189.580002	191.729996	191.729996	47145600	
2024-01-30	190.940002	191.800003	187.470001	188.039993	188.039993	55859400	
2024-01-31	188.010000	188.160000	187.110000	187.160000	187.160000	55859400	

```
plt.figure(figsize=(16,6))
plt.title('Close Price History')
plt.plot(df['Close'])
plt.xlabel('Date', fontsize=18)
plt.ylabel('Close Price USD ($)', fontsize=18)
plt.show()
```



```
# Create a new dataframe with only the 'Close column
data = df.filter(['Close'])
# Convert the dataframe to a numpy array
dataset = data.values
# Get the number of rows to train the model on
training_data_len = int(np.ceil( len(dataset) * .95 ))
```

```
training_data_len
```

```
2889
```

```
# Scale the data
from sklearn.preprocessing import MinMaxScaler
```

```
scaler = MinMaxScaler(feature_range=(0,1))
scaled_data = scaler.fit_transform(dataset)
```

```
scaled_data
```

```
array([[0.00401431],
       [0.00444289],
       [0.00533302],
       ...,
       [0.92555484],
       [0.93891265],
       [0.93342839]])
```

```
# Create the training data set
# Create the scaled training data set
train_data = scaled_data[0:int(training_data_len), :]
# Split the data into x_train and y_train data sets
x_train = []
y_train = []
```

```
for i in range(60, len(train_data)):
    x_train.append(train_data[i-60:i, 0])
    y_train.append(train_data[i, 0])
    if i<= 61:
        print(x_train)
        print(y_train)
        print()
```

```
# Convert the x_train and y_train to numpy arrays
x_train, y_train = np.array(x_train), np.array(y_train)
```

```
# Reshape the data
x_train = np.reshape(x_train, (x_train.shape[0], x_train.shape[1], 1))
# x_train.shape
```

```
[array([0.00401431, 0.00444289, 0.00533302, 0.00618049, 0.00605056,
        0.00634339, 0.00620958, 0.00598462, 0.00567821, 0.00662652,
        0.00748175, 0.007218 , 0.00577323, 0.00715207, 0.00579457,
        0.01088518, 0.01049151, 0.01100542, 0.01211663, 0.01278955,
        0.01273332, 0.01252582, 0.01341013, 0.01424207, 0.01518457,
        0.01670691, 0.01990478, 0.01995326, 0.02173353, 0.02306387,
        0.02077746, 0.02165789, 0.02164044, 0.02410915, 0.02375813,
        0.02440779, 0.02557523, 0.0262249 , 0.02809631, 0.02945961,
        0.02985329, 0.02999098, 0.02765997, 0.02709757, 0.02718096,
        0.02937236, 0.02998905, 0.03131358, 0.03443581, 0.03860139,
        0.0378218 , 0.03782373, 0.04083544, 0.04177794, 0.04110694,
        0.04049413, 0.03985611, 0.04197573, 0.0434302 , 0.04403914]))]
[0.042534249860459186]

[array([0.00401431, 0.00444289, 0.00533302, 0.00618049, 0.00605056,
        0.00634339, 0.00620958, 0.00598462, 0.00567821, 0.00662652,
        0.00748175, 0.007218 , 0.00577323, 0.00715207, 0.00579457,
        0.01088518, 0.01049151, 0.01100542, 0.01211663, 0.01278955,
        0.01273332, 0.01252582, 0.01341013, 0.01424207, 0.01518457,
        0.01670691, 0.01990478, 0.01995326, 0.02173353, 0.02306387,
        0.02077746, 0.02165789, 0.02164044, 0.02410915, 0.02375813,
        0.02440779, 0.02557523, 0.0262249 , 0.02809631, 0.02945961,
        0.02985329, 0.02999098, 0.02765997, 0.02709757, 0.02718096,
        0.02937236, 0.02998905, 0.03131358, 0.03443581, 0.03860139,
        0.0378218 , 0.03782373, 0.04083544, 0.04177794, 0.04110694,
        0.04049413, 0.03985611, 0.04197573, 0.0434302 , 0.04403914]), array([0.00444289, 0.00533302, 0.00618049, 0.00605056, 0.00634339,
        0.00620958, 0.00598462, 0.00567821, 0.00662652, 0.00748175,
        0.007218 , 0.00577323, 0.00715207, 0.00579457, 0.01088518,
        0.01049151, 0.01100542, 0.01211663, 0.01278955, 0.01273332,
        0.01252582, 0.01341013, 0.01424207, 0.01518457, 0.01670691,
        0.01990478, 0.01995326, 0.02173353, 0.02306387, 0.02077746,
        0.02165789, 0.02164044, 0.02410915, 0.02375813, 0.02440779,
        0.02557523, 0.0262249 , 0.02809631, 0.02945961, 0.02985329,
        0.02999098, 0.02765997, 0.02709757, 0.02718096, 0.02937236,
        0.02998905, 0.03131358, 0.03443581, 0.03860139, 0.0378218 ,
        0.03782373, 0.04083544, 0.04177794, 0.04110694, 0.04049413,
        0.03985611, 0.04197573, 0.0434302 , 0.04403914, 0.04253425]))]
[0.042534249860459186 0.040494130474309751]
```