

Welcome to Stock Market Analysis project!

In this portfolio project I have done an analysis on stocks from five different companies in healthcare industry over the stock data for one year time period. I have used pandas to get stock information from Yahoo Finance, visualize different aspects of it, and finally I have analyzed the risk of a stock, based on its previous performance history. I have also used Monte Carlo method to predict future stock price.

In the projects the following questions will be answered through analysis and visualization.

- 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 correlation between different stocks' closing prices?
- 4) What was the correlation between different stocks' daily returns?
- 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?

In [5]: `#Importing Libraries for the Analysis`

```
import pandas as pd
from pandas import Series,DataFrame
import numpy as np

# For visualization
import matplotlib.pyplot as plt
import seaborn as sns
sns.set_style('whitegrid')
%matplotlib inline

# For time stamps
from datetime import datetime

# For division
from _future_ import division
```

In [6]: `# For reading data for Yahoo Finance`

```
import pandas_datareader as dr
```

Creating dataframe for each companies in Health Care industry with the data taken from Yahoo finance within the period '2021-07-10'-2022-07-10' Healthcare industry companies: United Health Group Incorporated (UNH), Eli & Lilly (LLY), Johnson & Johnson (J&J), Pfizer (PFE), Novo Nordisk (NVO)

In [7]: `UNH=dr.get_data_yahoo('UNH', start='2021-07-10', end='2022-07-10')
LLY=dr.get_data_yahoo('LLY', start='2021-07-10', end='2022-07-10')
JNJ=dr.get_data_yahoo('JNJ', start='2021-07-10', end='2022-07-10')
PFE=dr.get_data_yahoo('PFE', start='2021-07-10', end='2022-07-10')
NVO=dr.get_data_yahoo('NVO', start='2021-07-10', end='2022-07-10')`

In [8]: `# Different way : Getting data for all companies with For Loops
health_list=['UNH','LLY','JNJ','PFE','NVO']
for stock in health_list:
 global s[stock]=dr.get_data_yahoo(stock, start='2021-07-10', end='2022-07-10')`

In [10]: `PFE.head()`

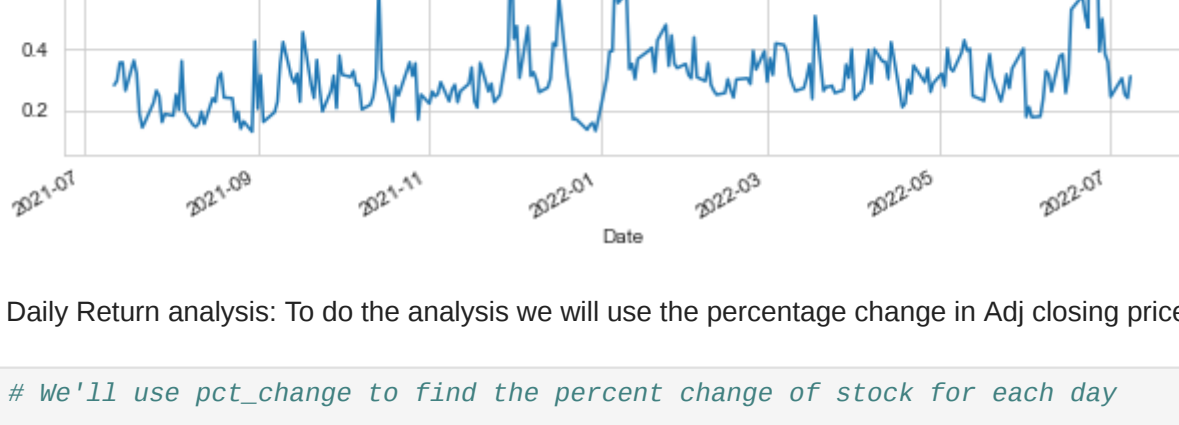
	High	Low	Open	Close	Volume	Adj Close
Date						
2021-07-12	40.250000	39.599998	39.660000	39.759998	24513000.0	38.451939
2021-07-13	39.799999	39.599999	39.770000	39.650002	12222000.0	38.345566
2021-07-14	40.029999	39.660000	39.720001	39.950001	15886600.0	38.635666
2021-07-15	40.189999	39.860001	39.960001	40.090000	20258600.0	38.771080
2021-07-16	40.349998	39.830000	40.139999	40.349998	24270400.0	39.022530

In [11]: `#Summary stats
UNH.describe()`

	High	Low	Open	Close	Volume	Adj Close
count	251.000000	251.000000	251.000000	251.000000	2.510000e+02	251.000000
mean	489.266684	459.537051	464.126215	464.661076	3.076912e+06	461.428663
std	39.791978	38.122297	38.861917	39.129624	1.261039e+06	40.336668
min	393.690002	383.119995	389.339996	387.010010	1.284200e+06	383.246979
25%	425.169998	419.580002	422.049988	423.134995	2.409400e+06	417.556054
50%	473.179993	462.000000	467.190002	468.410004	2.977400e+06	465.368378
75%	499.970001	488.504990	494.050005	495.060003	3.535400e+06	492.562822
max	553.289978	539.000000	545.000000	546.010010	1.589110e+07	544.069763

In stock we see two columns, closing price and adjusted closing price. The closing price refers to the price the stock was traded at last whereas Adjusted closing price is the closing price adjusted for corporate actions such as dividend payouts, stock splits, or the issuance of more shares.

In [12]: `#Playing a look through plot of the Adj closing price of United Health Group
UNH['Adj Close'].plot(legend=True, figsize=(10,4))`



In [13]: `#Volume of stocks traded over the 1 year period
UNH['Volume'].plot(legend=True, figsize=(10,5))`

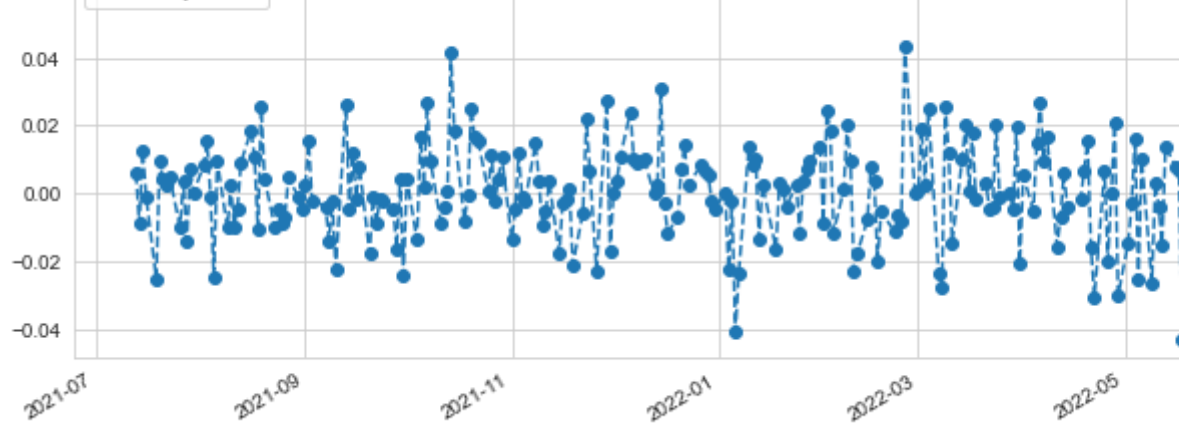
plt.legend(['UNH'])

plt.title('Stocks Volume traded over the period')

(an axes-level function for histograms).

(visibly in July 2022 there is a huge peak in the stock volume being traded for UNH

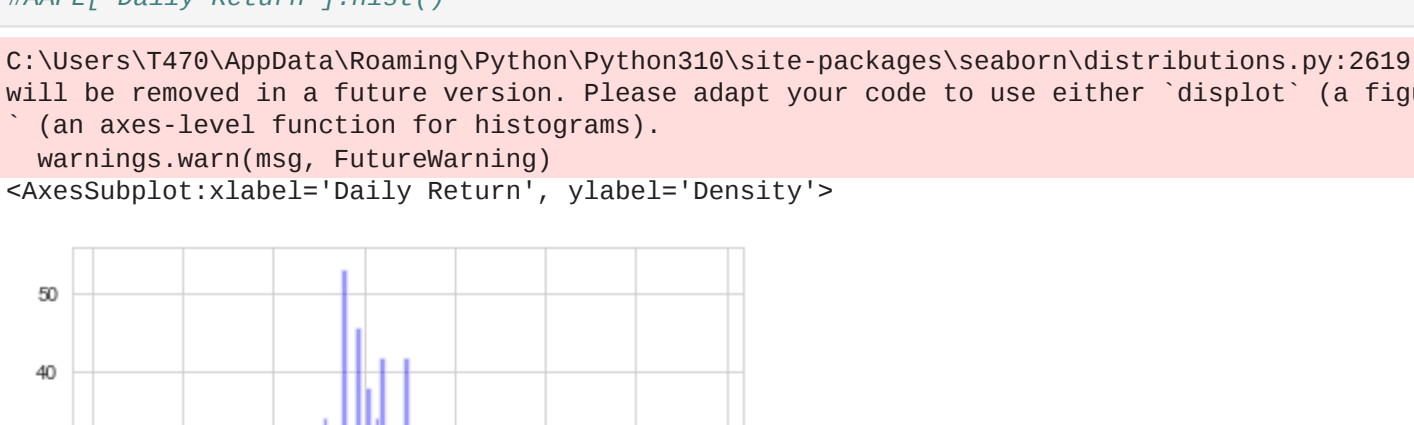
Text(0.5, 1.0, 'Stocks Volume traded over the period')



Daily Return analysis: To do the analysis we will use the percentage change in Adj closing price in each day.

In [14]: `# We'll use pct_change to find the percent change of stock for each day
UNH['Daily Return'] = UNH['Adj Close'].pct_change()`

In [15]: `# Then we'll plot the daily return percentage
UNH['Daily Return'].plot(figsize=(12,4), legend=True, linestyle='--', markers='o')`



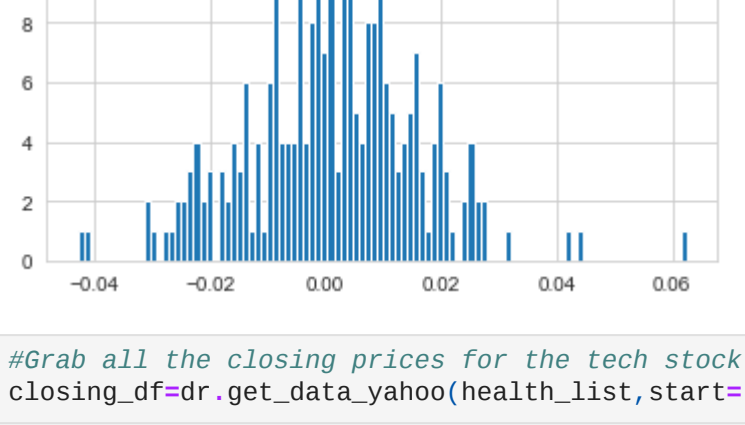
In [16]: `# We use the 'dropna()' here, otherwise the NaN values can't be read by seaborn
sns.displot(UNH['Daily Return'], dropna(), bins=100, color='blue')`

Could have also done: `MAPPA['Daily Return'].hist()`

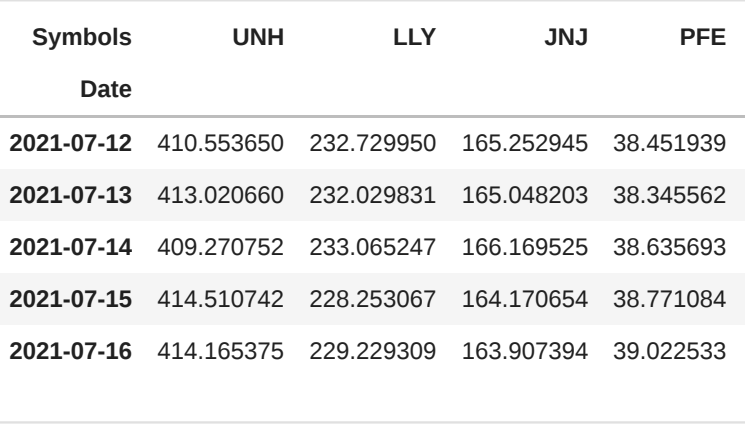
C:\Users\T476\AppData\Local\Programs\Python\Python318\site-packages\seaborn\distributions.py:2619: FutureWarning: 'displot' is a deprecated function and will be removed in a future version. Please adapt your code to use either 'displot' (a figure-level function with similar flexibility) or 'histplot' (an axes-level function for histograms).

warnings.warn(msg, FutureWarning)

<AxesSubplot: xlabel='Daily Return', ylabel='Density'>



In [17]: `UNH['Daily Return'].hist(bins=100)`



In [18]: `#Grab all the closing prices for the tech stock list into one DataFrame
closing_df=dr.get_data_yahoo(health_list, start='2021-07-10', end='2022-07-10')['Adj Close']`

In [19]: `closing_df.head()`

Out [19]:

Date	UNH	LLY	JNJ	PFE	NVO
2021-07-12	410.553660	232.729989	165.252945	38.451939	85.311806
2021-07-13	413.020660	232.028031	165.048203	38.345562	85.370918
2021-07-14	409.270782	231.065247	166.109525	38.635663	85.646751
2021-07-15	414.510742	228.253067	164.170654	38.771084	86.306786
2021-07-16	414.185779	229.229309	163.907394	39.022533	87.341164

In [20]: `#make a new health returns dataframe with percentage change data
health_returns=closing_df.pct_change()`

In [21]: `health_returns.head()`

Out [21]:

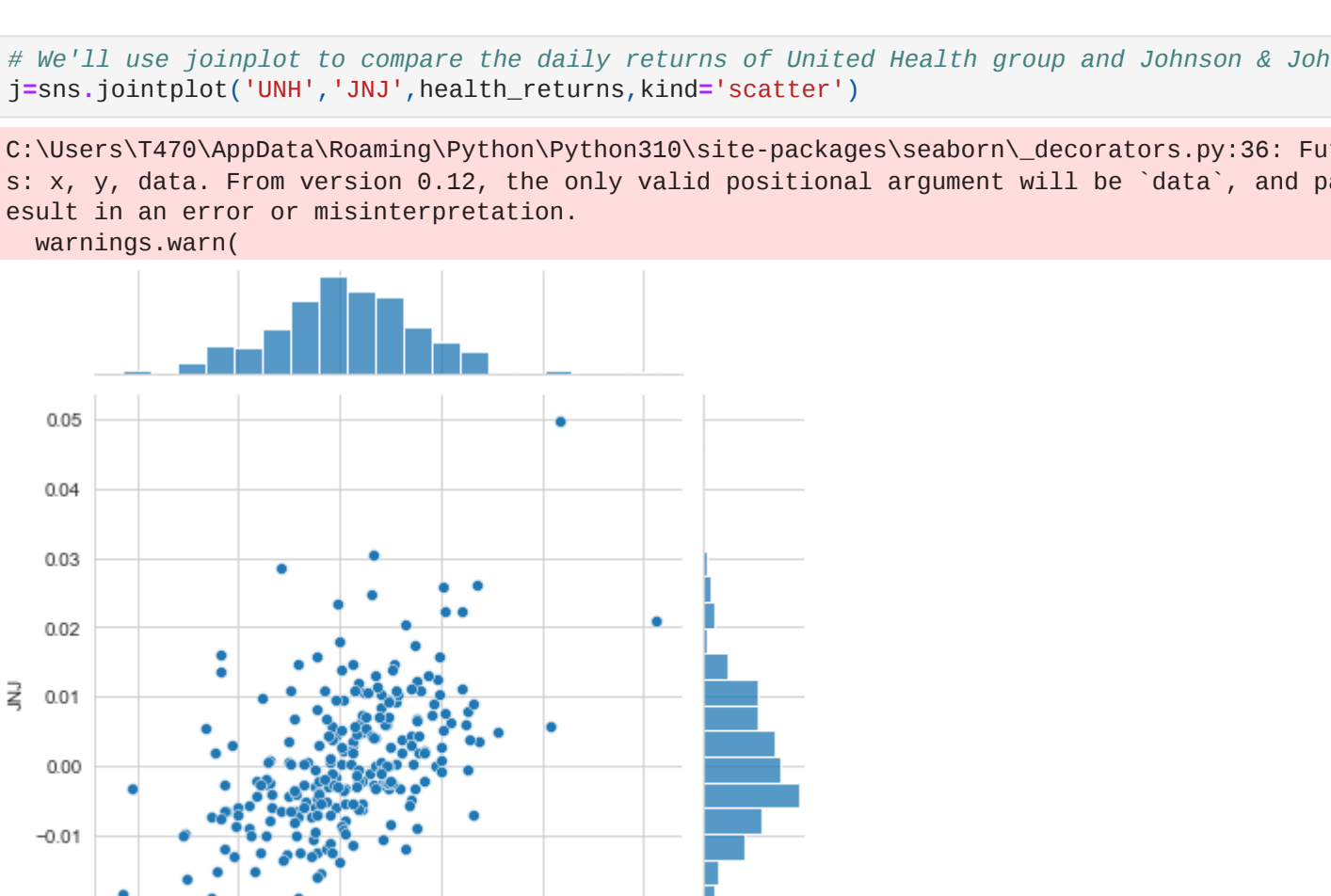
Symbols	UNH	LLY	JNJ	PFE	NVO
2021-07-12	NaN	NaN	NaN	NaN	NaN
2021-07-13	0.000009	-0.003008	-0.001239	-0.002766	0.000693
2021-07-14	-0.009079	0.004462	0.006794	0.007566	0.003221
2021-07-15	0.012803	-0.005647	-0.012029	0.003504	0.007706
2021-07-16	-0.000833	0.004277	-0.001604	0.006485	0.011985

In [22]: `# Comparing UnitedHealth Group Incorporated to itself should show a perfectly linear relationship
sns.jointplot(UNH, UNH, health_returns, kind='scatter', color='seagreen')`

C:\Users\T476\AppData\Local\Programs\Python\Python318\site-packages\seaborn_decorators.py:38: FutureWarning: Pass the following variables as keyword arg s: x, y, data. From version 0.12, the only valid positional argument will be 'data', and passing other arguments without an explicit keyword will r esult in an error or misinterpretation.

warnings.warn(msg, FutureWarning)

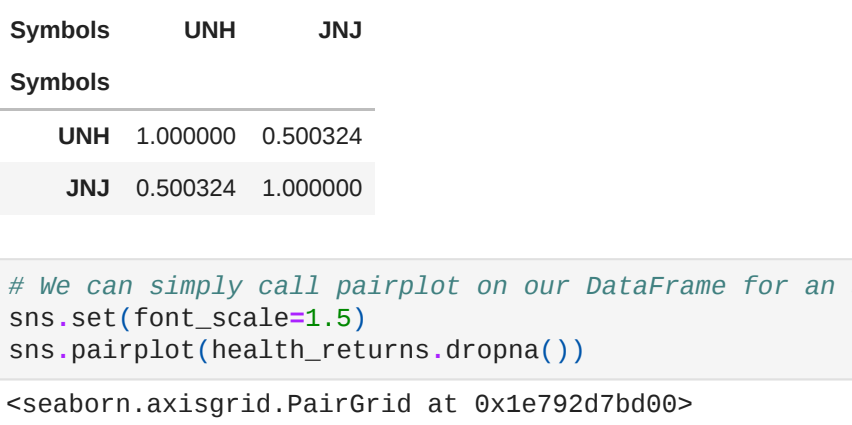
<seaborn.axisgrid.JointGrid at 0x1e792c6787b>



In [23]: `# We'll use jointplot to compare the daily returns of United Health group and Johnson & Johnson
sns.jointplot(UNH, JNJ, health_returns, kind='scatter')`

C:\Users\T476\AppData\Local\Programs\Python\Python318\site-packages\seaborn_decorators.py:38: FutureWarning: Pass the following variables as keyword arg s: x, y, data. From version 0.12, the only valid positional argument will be 'data', and passing other arguments without an explicit keyword will r esult in an error or misinterpretation.

warnings.warn(



In [24]: `import scipy.stats as stats
sns.set(font_scale=1.5)`

In [28]: `health_returns[['UNH', 'JNJ']].corr(method='pearson')`

Out [25]:

Symbols	UNH	JNJ
UNH	1.000000	0.500324
JNJ	0.500324	1.000000

In [26]: `# We can simply call pairplot on our DataFrame for an automatic visual analysis of all the comparisons
sns.set(font_scale=1.5)
sns.pairplot(health_returns.dropna())`

Out [26]: `<seaborn.axisgrid.PairGrid at 0x1e792c67b98>`



In [27]: `health_returns[['LLY', 'NVO']].corr(method='pearson')`

Out [27]:

Symbols	LLY	NVO
LLY	1.000000	0.485204
NVO	0.485204	1.000000

In [28]: `# Set up our figure by naming it returns_fig, call PairPlot on the DataFrame
returns_fig = sns.PairGrid(health_returns.dropna())`

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=1)

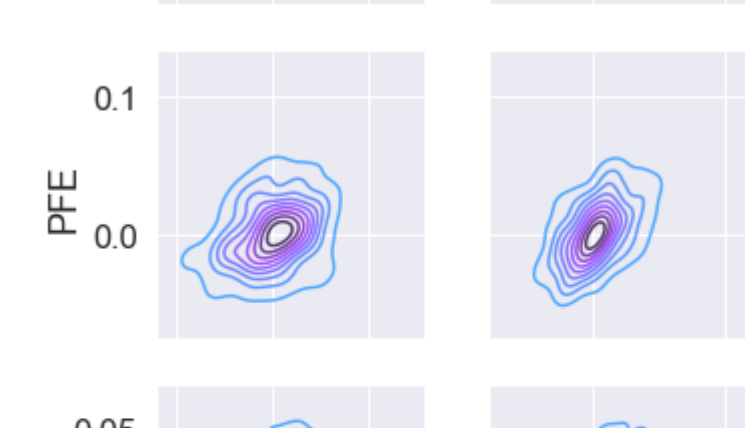
Finally we'll define the diagonal as a series of histogram plots of the daily return

returns_fig.map_diag(plt.hist, bins=30)



In [29]: `# Let's go ahead and use sebrn for a quick correlation plot for the daily returns
sns.heatmap(health_returns.corr(), cmap='YlGnBu', annot=True)`

Out [29]: `<AxesSubplot: xlabel='Symbols', ylabel='Symbols'>`



Just as from the PairPlot we see here numerically and visually that Johnson & Johnson and Eli Lilly & Company had the strongest correlation of daily stock return. It's also interesting to see that all the companies from the industry are positively correlated with JNJ and LLY being the highest correlation and Novo nordisk and Pfizer lowest pearson correlation value.

Now that we've done some daily return analysis, let's go ahead and start looking deeper into actual risk analysis.

In [30]: `sns.set_style('whitegrid')
rets = health_returns[['UNH', 'LLY', 'JNJ', 'PFE', 'NVO']]
rets = rets.dropna()
Let's start by creating a new DataFrame as a cleaned version of the original clean_rets DataFrame
rets = np.percentile(rets, [0.5, 99.5])
area = np.pi*20`

plt.scatter(rets.mean(), rets.std(), alpha = 0.5, s=area)

plt.ylim([0.01, 0.025])

plt.xlim([-0.003, 0.004])

#Set the plot axis titles

plt.xlabel('Expected returns')

plt.ylabel('Risk')

for label, x, y in zip(rets.columns, rets.mean(), rets.std()):

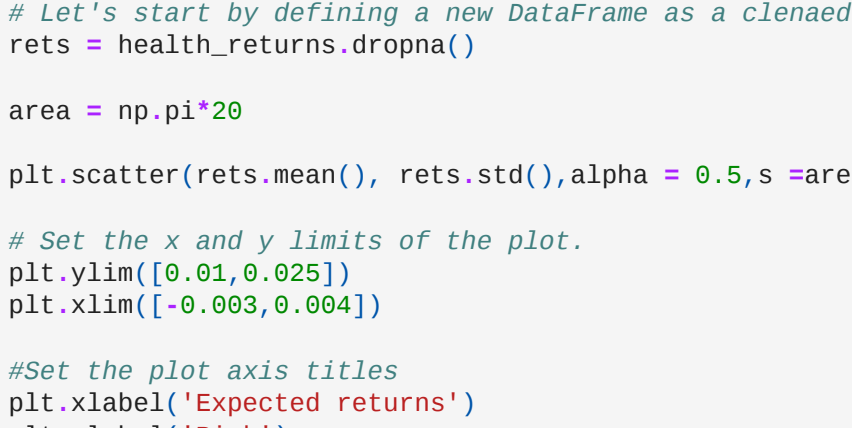
plt.annotate(

label,

xy = (x, y), xytext = (50, 50),

textcoords = 'offset points', ha = 'right', va = 'bottom',

arrowprops = dict(arrowstyle = '-', connectionstyle = 'arc3,rad=0.3'))



From the risk analysis we prefer to invest in stocks that has higher expected returns and less Risk. From the analysis, Johnson & Johnson seems to have a very low risk but with lowest expected returns among all the other companies. Eli Lilly seems to have the highest expected returns among all the companies. United Health group althoug had a little less return but posses lower risk than Eli Lilly. Novo nordisk and pfizer both has higher returns but with high risk. Ideal investment would be to invest in United Health Group.

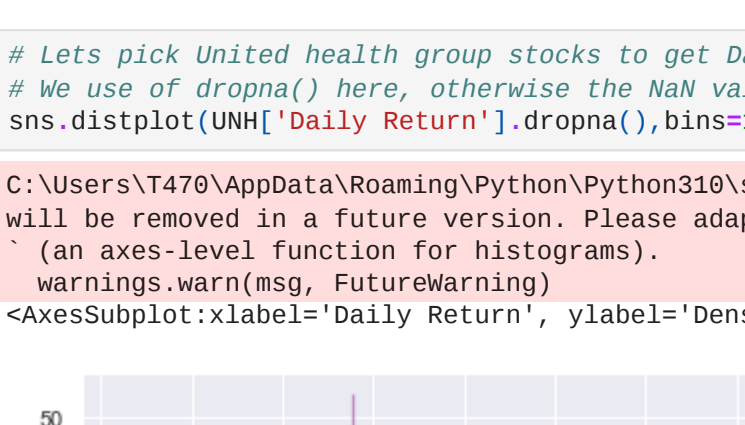
Value at Risk: Let's define a value at risk parameter for our stocks. We can treat value at risk as the amount of money we could expect to lose (aka putting at risk) for a given confidence interval. Here Value at risk is analyzed using the "bootstrap" method. For this method we will calculate the empirical quantiles from a histogram of daily returns

In [31]: `# Lets pick UnitedHealth group stocks to get Daily Return histogram
sns.displot(UNH['Daily Return'], dropna(), bins=100, color='purple')`

C:\Users\T476\AppData\Local\Programs\Python\Python318\site-packages\seaborn\distributions.py:2619: FutureWarning: 'displot' is a deprecated function and will be removed in a future version. Please adapt your code to use either 'displot' (a figure-level function with similar flexibility) or 'histplot' (an axes-level function for histograms).

warnings.warn(msg, FutureWarning)

<AxesSubplot: xlabel='Daily Return', ylabel='Density'>



In [32]: `# To get the risk value of stock we use quantile..
The 0.95 empirical quantile of daily returns
rets['UNH'].quantile(0.95)`

Out [32]: `-0.823516431118474197`

The 0.05 empirical quantile of daily returns is at -0.023. That means that with 95% confidence, our worst daily loss will not exceed 2.3%. If we have a \$1000 investment, our worst daily loss will be \$23. So we with 95% confidence we can say we only going to lose \$23 out of \$1000 investment.

In [33]: `rets['LLY'].quantile(0.95)`

Out [33]: `-0.823596964910883202`

In [34]: `# Set up our time horizon
days = 365`

Now our delta

dt = 1/days

Now let's grab our mu (drift) from the expected return data we got for AAPL

mu = rets.mean()['UNH']

Now let's grab the volatility of the stock from the std() of the average return

sigma = rets.std()['UNH']

In [35]: `def stock_monte_carlo(start_price, days, mu, sigma):
 """ This function takes in starting stock price, days of simulation, mu, sigma, and returns simulated price array"""`

Define a price array

price = np.zeros(days)

price[0] = start_price

Shock and drift

shock = np.zeros(days)

drift = np.zeros(days)

Run price array for number of days

for x in range(1, days):

Calculate Shock

shock[x] = np.random.normal(loc=mu * dt, scale=sigma * np.sqrt(dt))

Calculate drift

drift[x] = mu * dt

Calculate Price

price[x] = price[x-1] + (price[x-1] * (drift[x] + shock[x]))

return price

In [36]: `UNH=np.round(UNH.loc[:,:], decimals=2)`

Out [36]:

Date	High	Low	Open	Close	Volume	Adj Close	Daily Return
2021-07-12	416.399994	410.309998	411.019989	416.040004	28026000.0	410.553619	NaN
2021-07-13	419.980002	415.709991	416.380005	418.540009	29489000.0	413.020660	0.006009
2021-07-14	422.529999	413.350006	420.769989	414.739990	35387000.0	409.270721	-0.009079
2021-07-15	427.120011	407.200012	411.910004	420.049988	35468000.0	414.510742	0.012803
2021-07-16	422.950009	417.640015	421.359985	419.700012	26254000.0	414.185779	-0.000833
2022-07-01	517.330029	502.309998	512.320007	517.400024	2426000.0	517.400024	0.007340
2022-07-05	511.000000	492.250000	507.640015	505.239990	3029600.0	505.239990	-0.023502
2022-07-06	517.409973	504.299988	505.110004	515.289978	2506700.0	515.289978	0.019892
2022-07-07	517.299978	502.229980	515.250000	514.380005	2385400.0	514.380005	-0.001766
2022-07-08	528.369995	511.010011	512.309998	518.630005	3992100.0	518.630005	0.006262

251 rows x 7 columns

In [37]: `# Get start price from UNH.head()
start_price = 411.82`

for run in range(100):

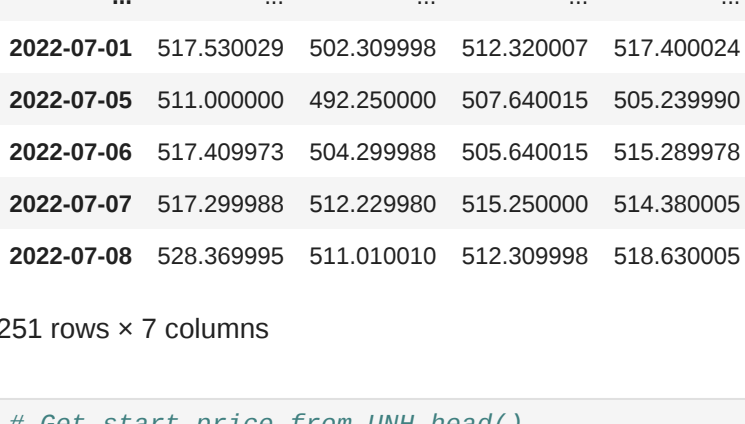
plt.plot(stock_monte_carlo(start_price, days, mu, sigma))

plt.xlabel('Days')

plt.ylabel('Price')

plt.title('Monte Carlo Analysis for United Health Group')

Text(0.5, 1.0, 'Monte Carlo Analysis for United Health Group')



In [38]: `# Set a large number of runs
runs = 10000`

Create an empty matrix to hold the end price data

simulations = np.zeros(runs)

<