

```
(224) # Initialize Outer Grader
def grader_check():
    grader = otter.Notebook()
```



NAME: Talia Arauzo

STUDENT ID: 3031791660

Data-X Fall 2020: Homework 04

Linear regression and matplotlib.

Author: Ming-Yen Kao.

In this homework, you will do some exercises with regression and plotting.

Only submit your jupyter notebook to Gradescope

Question 1 - Data Pre-processing

Data Source of the S&P 500 Index is Yahoo Finance; Data Source of Gold.csv is Gold.org; Data Source of 30YTBond.csv is Macrotrends.net. SP500 is a stock market index that measures the stock performance of 500 large companies listed on stock exchanges in the United States. The unit of Gold price is USD per oz. 30YTBond stands for US Treasury Bond future price with maturity of 30 years.

1a) Use pandas to load data from 'csv' to DataFrame. Please load 'SP500.csv' to df.SP500, load 'Gold.csv' to df.Gold, and load '30YTBond.csv' to df.30YTB.

```
In [58]: # Load required modules
import numpy as np
import pandas as pd
import matplotlib.pyplot as plt
matplotlib inline
```

```
In [112]: df.SP500 = pd.read_csv("SP500.csv")
df.Gold = pd.read_csv("Gold.csv")
df.30YTB = pd.read_csv("30YTBond.csv")
```

```
In [60]: print(df.SP500.HEAD())
display(df.SP500.head())
print(df.Gold)
display(df.Gold.HEAD())
print(df.30YTBond)
display(df.30YTBond.head())
```

SP500:

	Date	Open	High	Low	Close	Adj Close	Volume
0	1928-12-01	17.660000	17.660000	17.660000	17.660000	17.660000	0
1	1928-02-01	17.760000	17.760000	17.760000	17.760000	17.760000	0
2	1928-02-01	17.630001	17.629999	16.950001	17.260000	17.260000	0
3	1928-03-01	17.299999	19280001	17.299999	19280001	19.280001	0
4	1928-04-01	18.910000	19.750000	18.910000	19.750000	19.750000	0

Gold:

	Date	US dollar	Euro	Pound sterling	Japanese yen
0	1999/12/31	35.20	27.24	14.66	12592.44
1	1970/1/30	34.99	23.20	14.54	12496.35
2	1970/2/27	35.02	23.26	14.55	12529.95
3	1970/3/31	35.30	23.42	14.67	12627.65
4	1970/4/30	35.85	23.79	14.93	12658.36

30YTBond:

	Date	Value
0	1977/8/2	102.375
1	1977/8/23	102.781
2	1977/8/24	102.658
3	1977/8/25	103.094
4	1977/8/26	103.062

```
In [61]: grader.check("q1a")
```

Out [61]: All tests passed!

1b) Only leave 'Date' and 'USDollar price' in df.Gold and only leave 'Date' and 'closing price' in df.SP500. In other word 'Open', 'High', 'Low', 'Adj Close', and 'Volume' have to be dropped from df.SP500. 'Euro', 'Pound sterling', and 'Japanese yen' have to be dropped from df.Gold.

```
In [113]: df.SP500 = df.SP500[["Date", "Close"]]
df.Gold = df.Gold[["Date", "US dollar"]]
```

```
In [63]: print(df.SP500)
display(df.SP500.head())
print(df.Gold)
display(df.Gold.head())
```

SP500:

	Date	Close
0	1927-12-01	17.660000
1	1928-01-01	17.760000
2	1928-02-01	17.260000
3	1928-03-01	19.280001
4	1928-04-01	19.750000

Gold:

	Date	US dollar
0	1999/12/31	35.20
1	1970/1/30	34.99
2	1970/2/27	35.02
3	1970/3/31	35.30
4	1970/4/30	35.85

```
In [64]: grader.check("q1b")
```

Out [64]: All tests passed!

1c) Sample the very first price of the asset in each month to represent that month.

For example, 1999/1/2 price of bananan is 90, 1999/1/6 price of bananan is 85, 1999/1/6 price of bananan is 91. You should use 90 to represent the price of 1999/1. Please convert the Date to pandas DateTime format with day always = 1 ex: 1970-02-01 or 2000-11-01. Replace the original data frames of df.SP500, df.Gold, and df.SP500 with the sampled data frames. For example, we have [1999/1/2: 90, 1999/1/5: 85, 1999/1/6: 91, 1999/2/7: 89], after sampling, we should have [1999/1/1: 90, 1999/2/1: 70].

- Remember to keep the Date in the format of Pandas timestamp!

df.30YTB after sampling should be

```
In [116]: # df.SP500
df.SP500["Date"] = pd.to_datetime(df.SP500["Date"])
df.SP500["year_month"] = pd.DateTimeIndex(df.SP500["Date"]).to_period('M')
df.SP500 = df.SP500[["Date", "Close"]].isin(df.SP500.groupby("year_month")["Date"].min()).drop(columns = ["year_month"])
```

```
In [118]: # df.Gold
df.Gold["Date"] = pd.to_datetime(df.Gold["Date"])
df.Gold["year_month"] = pd.DateTimeIndex(df.Gold["Date"]).to_period('M')
df.Gold = df.Gold[["Date", "US dollar"]].isin(df.Gold.groupby("year_month")["Date"].min()).drop(columns = ["year_month"])
```

```
In [114]: # df.30YTB
df.30YTB["Date"] = pd.to_datetime(df.30YTB["Date"])
df.30YTB["year_month"] = pd.DateTimeIndex(df.30YTB["Date"]).to_period('M')
df.30YTB = df.30YTB[["Date", "Value"]].isin(df.30YTB.groupby("year_month")["Date"].min()).drop(columns = ["year_month"])
```

```
In [119]: grader.check("q1c")
```

Out [119]: All tests passed!

Question 2 - Data Visualization & Linear Regression

```
In [120]: import matplotlib.pyplot as plt
```

2a) Observe the example plot below and try to reproduce it: y-left is df.SP500[Date] in linear scale and y-right is df.SP500[Date] in log scale.

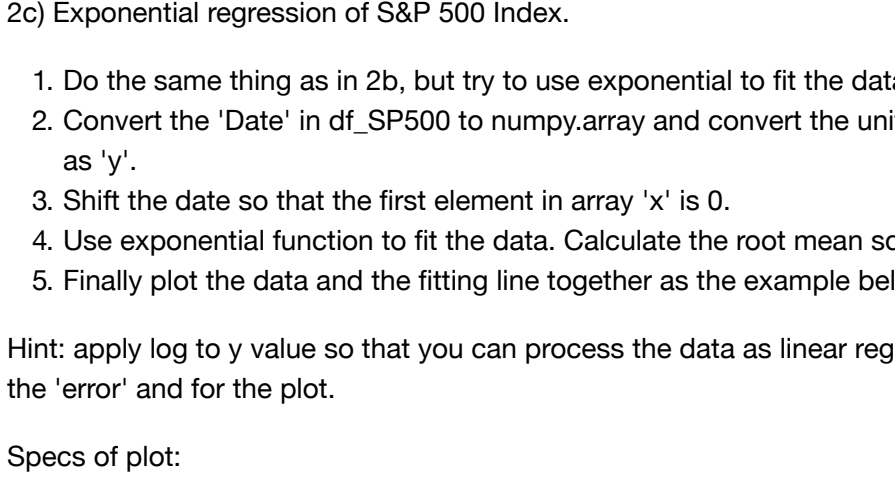
- Specs:
- x-label and x-scale in black
 - y-label and x-scale in black
 - Left label, y-right label, and the corresponding line-plot in blue
 - Right label, y-right label, and the corresponding line-plot in red
 - plot title is 'Price of the S&P 500 Index'
 - y-left label is 'Price (linear in \$ US Dollars)'
 - y-right label is 'Price (log in \$ US Dollars)'
 - x label is 'Date'

Some helpful functions: axis.set_xlabel(), axis.set_ylabel(), axis.tick_params(), and axis.twinx()

```
In [146]: fig, ax1 = plt.subplots()
plt.title("Price of the S&P 500 Index")

# Blue side
ax1.set_xlabel("Date")
ax1.set_ylabel("Price (linear in $ US Dollars)", color = color1)
ax1.plot(df.SP500["Date"], df.SP500["Close"], color = color1)
ax1.tick_params(axis='y', labelcolor = color1)

# Red side
ax2 = ax1.twinx()
color2 = "tab:red"
ax2.set_ylabel("Price (log in $ US Dollars)", color = color2)
ax2.plot(df.SP500["Date"], df.SP500["Close"], color = color2)
ax2.tick_params(axis='y', labelcolor = color2)
ax2.set_yscale("log");
```



```
In [147]: grader.check("q2a")
```

Out [147]: All tests passed!

2b) Linear regression of SP500 Index.

- Convert the 'Date' in df.SP500 to numpy.array and convert the unit to day as 'x' and convert the 'Close' in df.SP500 to numpy.array as 'y'.
- Shift the date so that the first element in array 'x' is 0.
- Find the linear regression procedure in the slides <https://datax.berkeley.edu/wp-content/uploads/2020/09/slides-m100-linear.pdf>.
- Calculate the root mean square error and store it in 'error'.
- Finally plot the data and the fitting line together as the example plot below.

Specs of plot:

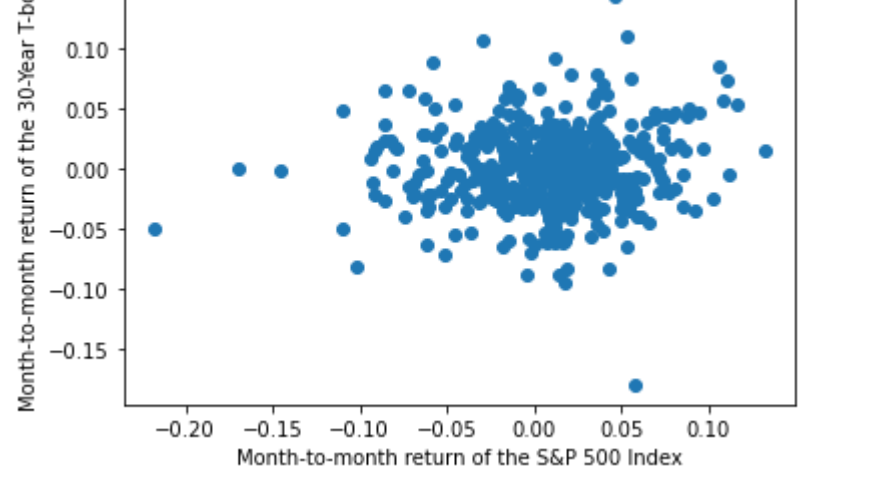
- Color in default (please plot the data first so that the color of data will be blue)
- plot title is 'Price of the S&P 500 Index'
- x-label is 'Date'
- y-label is 'Price (in \$ US Dollars)'

```
In [194]: # Linear reg
x = df.SP500["Date"].to_numpy().astype("int") / 86400
x = x - x[0]
y = df.SP500["Close"].to_numpy()
```

```
In [204]: # Linear reg
E_x = np.mean(x)
E_y = np.mean(y)
cov_xy = np.mean(x * y) - (E_x * E_y)
y_0 = E_y - (cov_xy / np.var(x)) * E_x
m = cov_xy / np.var(x)
y_pred = y_0 + (m * x)
error = np.sqrt(np.mean(np.square(y_pred - y)))
print ("Root Mean Square Error: ", error)
```

Root Mean Square Error: 4.7347997247592337

```
In [205]: # Linear reg plot
plt.plot(df.SP500["Date"], df.SP500["Close"])
plt.plot(df.SP500["Date"], y_pred)
plt.title("Price of the S&P 500 Index")
plt.xlabel("Date")
plt.ylabel("Price (in $ US Dollars)");
```



```
In [206]: grader.check("q2b")
```

Out [206]: All tests passed!

2c) Exponential regression of 2b Plot Index.

- Do the same thing as in 2b, but try to use exponential to fit the data this time.
- Convert the 'Date' in df.SP500 to numpy.array and convert the unit to day as 'x' and convert the 'Close' in df.SP500 to numpy.array as 'y'.
- Shift the date so that the first element in array 'x' is 0.
- Use exponential function to fit the data. Calculate the root mean square error and store it in 'error'.
- Finally plot the data and the fitting line together as the example below.

Hint: apply log to y value so that you can process the data as linear regression in 2(b). Remember to convert it back for the calculation for the 'error' and for the plot.

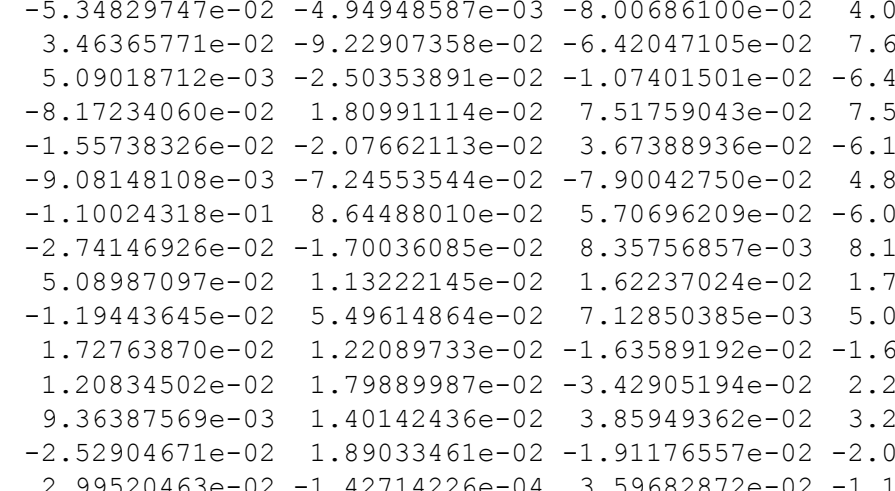
Specs of plot:

- Color in default (please plot the data first so that the color of data will be blue)
- plot title is 'Price of the S&P 500 Index'
- x-label is 'Date'
- y-label is 'Price (in \$ US Dollars)'

```
In [219]: # Exponential reg
y_log = np.log(y)
E_x = np.mean(x)
E_y = np.mean(y_log)
cov_xy = np.mean(x * y_log) - (E_x * E_y)
y_0 = E_y - (cov_xy / np.var(x)) * E_x
m = cov_xy / np.var(x)
y_pred = np.exp(y_0 + (m * x))
error = np.sqrt(np.mean(np.square(y_pred - y)))
print ("Root Mean Square Error: ", error)
```

Root Mean Square Error: 184.58637172480607

```
In [222]: # Linear reg plot
plt.plot(df.SP500["Date"], df.SP500["Close"])
plt.plot(df.SP500["Date"], y_pred)
plt.title("Price of the S&P 500 Index")
plt.xlabel("Date")
plt.ylabel("Price (in $ US Dollars)");
```



```
In [223]: grader.check("q2c")
```

Out [223]: All tests passed!

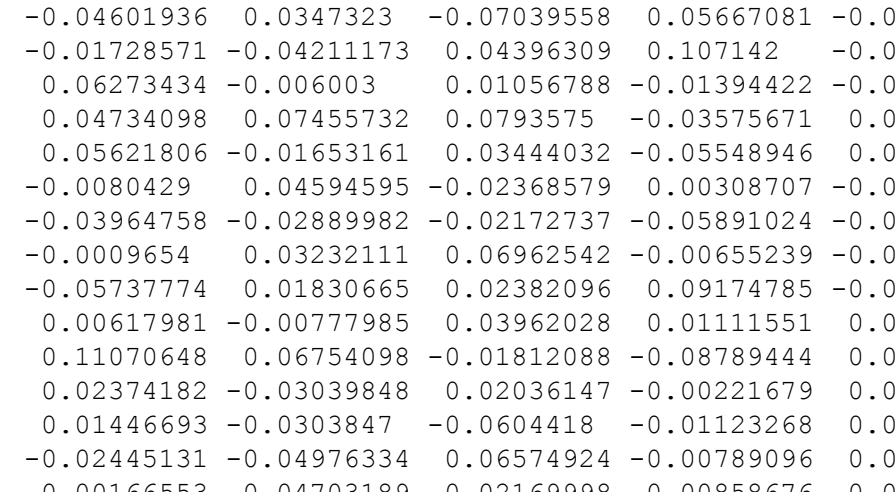
2d) Scatter plot of month-to-month return of S&P 500 Index and 30 Year Treasury Bond.

- Compute the month-to-month returns on S&P 500 and 30 Year Treasury Bond between 1977-08-01 (inclusive) and 2019-12-31 (inclusive). The month-to-month returns should be in the net-increment. For example, price on 2014-1-1 is 100 and price on 2014-1-1 is 105. The month-to-month return from 2014-1-1 to 2014-2-1 should be 0.05. Another example, if the array of the price of S&P 500 is [100, 105, 100, 105], the month-to-month return array should be [0.05, -0.0476190476, 0.05].
- Save the month-to-month returns in the format of numpy array to df.x3' and 'y3', respectively. Note that if the dimension of the array of price is n, the dimension of month-to-month return array should be n-1.
- Finally plot 'x3' and 'y3' as scatter plot. Try to reproduce the plot below.

Specs of plot:

- Color in default
- plot title is 'Scatter Plot of the Returns'
- x-label is 'Month-to-month return of the S&P 500 Index'
- y-label is 'Month-to-month return of the 30-Year T-bond'

```
In [293]: mask_SP500 = (df.SP500["Date"] >= pd.Timestamp(1977, 8, 1)) & (df.SP500["Date"] <= pd.Timestamp(2019, 12, 31))
x3 = df.SP500.sort_values(by = "Date").loc[mask_SP500][["Close"]].values
x3 = np.array([x3[idx + 1] - val / val for idx, val in enumerate(x3) if idx < len(x3) - 1])
mask_30YTB = (df.30YTB["Date"] >= pd.Timestamp(1977, 8, 1)) & (df.30YTB["Date"] <= pd.Timestamp(2019, 12, 31))
y3 = df.30YTB.sort_values(by = "Date").loc[mask_30YTB][["Value"]].values
y3 = np.array([y3[idx + 1] - val / val for idx, val in enumerate(y3) if idx < len(y3) - 1])
plt.scatter(x3, y3)
plt.title("Scatter Plot of the Returns")
plt.xlabel("Month-to-month return of the S&P 500 Index")
plt.ylabel("Month-to-month return of the 30-Year T-bond");
```



```
In [294]: print (x3, y3)
```

[-2.48008688e-03 -4.34062265e-02 2.69656282e-02 2.84715801e-03
-6.15141758e-02 -2.47618936e-02 2.49310429e-02 8.54164677e-02
4.23483533e-03 -1.75853459e-02 5.39097776e-02 2.59237266e-02
7.26110343e-02 -3.91500392e-02 1.66379742e-02 1.48897636e-02
3.87461134e-02 -3.65255779e-02 3.51516105e-02 3.87452215e-02
2.63646775e-02 3.8655512e-02 8.74544714e-03 5.30775585e-02
0.00000000e+00 -6.86059276e-02 4.26242781e-02 1.67671245e-02
5.76246237e-02 -4.37981765e-03 -1.01794895e-01 4.11402210e-02
4.6570741e-02 2.69687168e-02 3.49259365e-02 9.09120855e-02
2.51675280e-02 1.60210586e-01 1.02377053e-01 -2.34542447e-02
-4.57424295e-02 1.320767345e-02 3.60325730e-02 -2.34558971e-02
-1.6561889e-03 -1.040793467e-02 -2.21026587e-03 -6.20989989e-02
5.38517530e-02 4.93478654e-02 1.6359362e-02 1.66320363e-01
-1.7538674e-02 -6.05481801e-02 -0.1071116e-02 4.00143180e-02
-3.91618423e-02 -2.02895608e-02 -0.23906485e-02 1.15937276e-01
7.6131937e-02 1.10446796e-01 3.59706588e-02 1.52313579e-02
2.3132721e-02 0.0098839 0.03737385 -0.0093994 -0.01234916 -0.0396465
-1.24064592e-02 3.23295771e-02 -0.30303036e-02 1.13188732e-02
-0.1582303e-02 -1.51743475e-02 1.74558083e-02 -8.83410494e-03
-1.91526474e-02 -3.88593467e-02 1.34979472e-02 5.46537576e-02
3.93564500e-02 2.50253991e-02 1.6359362e-02 1.66320363e-01
-3.47964379e-02 -6.02648985e-02 -0.15122528e-02 2.3734984e-02
7.40851449e-02 6.62878114e-03 -0.28701336e-02 -0.59427644e-02
5.40510532e-02 1.21340172e-02 -4.84757689e-03 -1.19945162e-02
-3.47240780e-02 4.25088143e-02 6.50615975e-02 4.50610526e-02
2.36652784e-02 3.71981603e-02 3.49259365e-02 9.09120855e-02
5.0229875e-02 1.41095206e-02 -5.8682836e-02 7.11926070e-02
-8.54386059e-02 5.47293300e-02 2.14771911e-02 -2.82882713e-02
1.31766896e-02 3.69236189e-02 2.63986575e-02 -1.45021212e-02
0.03419715e-02 4.79145906e-02 1.6359362e-02 1.66320363e-01
-2.41661652e-02 -2.17630416e-01 -8.5396545e-02 7.28614797e-02
4.04322686e-02 4.18174027e-02 -0.33432595e-02 1.55167164e-02
3.17621631e-02 4.32560109e-02 -5.41137477e-02 -3.86010454e-02
2.16083420e-02 0.0662535 0.0293563 0.0293563 0.0293563 0.0293563
7.11147916e-02 -2.89441489e-02 2.08059524e-02 5.00893988e-02
3.5174935e-02 -7.92454512e-03 8.83702359e-02 1.85167164e-02
-6.5443747e-02 -2.51745226e-02 1.65413274e-02 2.24168136e-02
8.725253e-02 8.5396252e-02 2.42549840e-02 1.68971878e-02
1.98891902e-02 -8.88635767e-03 -5.22315808e-03 4.4814383e-02
-5.11843071e-02 -6.69821297e-02 5.99342138e-02 -0.0448952 -0.0562928
4.15177571e-02 6.72811749e-02 2.22028238e-02 3.19799050e-02
3.86049746e-02 4.78926292e-02 4.48593432e-02 1.94887874e-02
-9.1437375e-02 1.18342370e-02 3.39037087e-02 1.1538755e-02
-1.99237505e-02 9.58954208e-03 -2.18318627e-02 2.78927146e-02
6.39597075e-04 1.74588321e-02 3.39373641e-02 -2.39973500e-02
1.31766896e-02 3.69236189e-02 2.63986575e-02 -1.45021212e-02
0.03419715e-02 4.79145906e-02 1.6359362e-02 1.66320363e-01
-2.41661652e-02 -2.17630416e-01 -8.5396545e-02 7.28614797e-02
4.04322686e-02 4.18174027e-02 -0.33432595e-02 1.55167164e-02
3.17621631e-02 4.32560109e-02 -5.41137477e-02 -3.86010454e-02
2.16083420e-02 0.0662535 0.0293563 0.0293563 0.0293563 0.0293563
7.11147916e-02 -2.89441489e-02 2.08059524e-02 5.00893988e-02
3.5174935e-02 -7.92454512e-03 8.83702359e-02 1.85167164e-02
-6.5443747e-02 -2.51745226e-02 1.65413274e-02 2.24168136e-02
8.725253e-02 8.5396252e-02 2.42549840e-02 1.68971878e-02
1.98891902e-02 -8.88635767e-03 -5.22315808e-03 4.4814383e-02
-5.11843071e-02 -6.69821297e-02 5.99342138e-02 -0.0448952 -0.0562928
4.15177571e-02 6.72811749e-02 2.22028238e-02 3.19799050e-02
3.86049746e-02 4.78926292e-02 4.48593432e-02 1.94887874e-02
-9.1437375e-02 1.18342370e-02 3.39037087e-02 1.1538755e-02
-1.99237505e-02 9.58954208e-03 -2.18318627e-02 2.78927146e-02
6.39597075e-04 1.74588321e-02 3.39373641e-02 -2.39973500e-02
1.31766896e-02 3.69236189e-02 2.63986575e-02 -1.45021212e-02
0.03419715e-02 4.79145906e-02 1.6359362e-02 1.66320363e-01
-2.41661652e-02 -2.17630416e-01 -8.5396545e-02 7.28614797e-02
4.04322686e-02 4.18174027e-02 -0.33432595e-02 1.55167164e-02
3.17621631e-02 4.32560109e-02 -5.41137477e-02 -3.86010454e-02
2.16083420e-02 0.0662535 0.0293563 0.0293563 0.0293563 0.0293563
7.11147916e-02 -2.89441489e-02 2.08059524e-02 5.00893988e-02
3.5174935e-02 -7.92454512e-03 8.83702359e-02 1.85167164e-02
-6.5443747e-02 -2.51745226e-02 1.65413274e-02 2.24168136e-02
8.725253e-02 8.5396252e-02 2.42549840e-02 1.68971878e-02
1.98891902e-02 -8.88635767e-03 -5.22315808e-03 4.4814383e-02
-5.11843071e-02 -6.69821297e-02 5.99342138e-02 -0.0448952 -0.0562928
4.1