

Rohan Nyati

500075940

R177219148

Batch-5 (AI & MI)

Lab File

Experiment 1

Anaconda & Spyder Installation for windows

1. Click on the link below to open the download page

<https://www.anaconda.com/download/#windows>



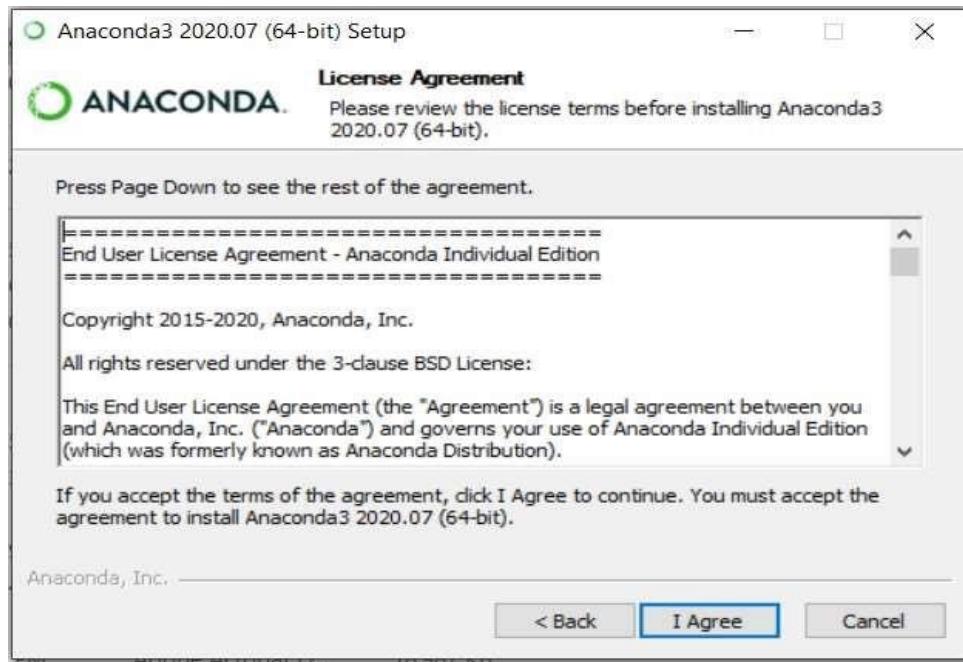
2. Click on the **Download** button and check for the compatibility of your system. Then, it will start downloading.



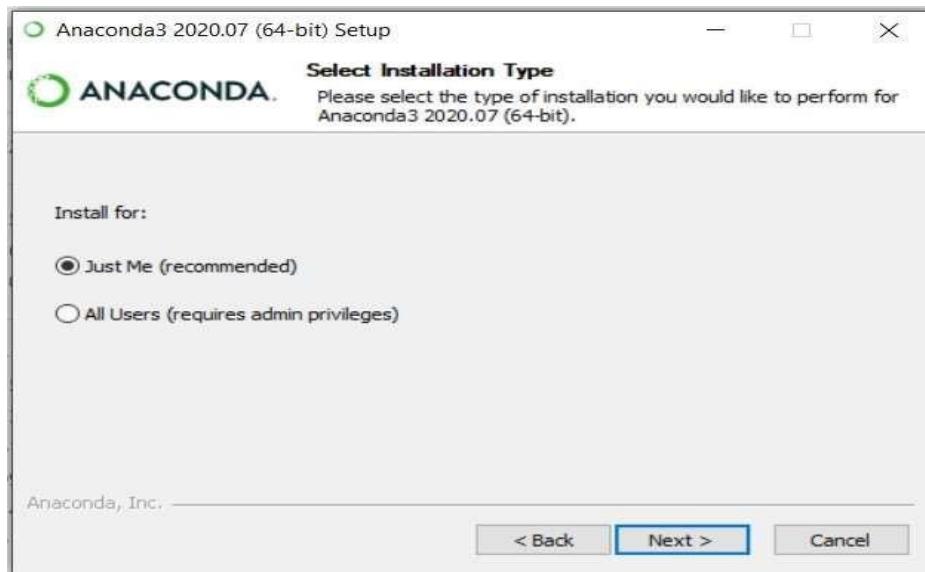
3. Double click the installer to launch.
4. Click on Next.



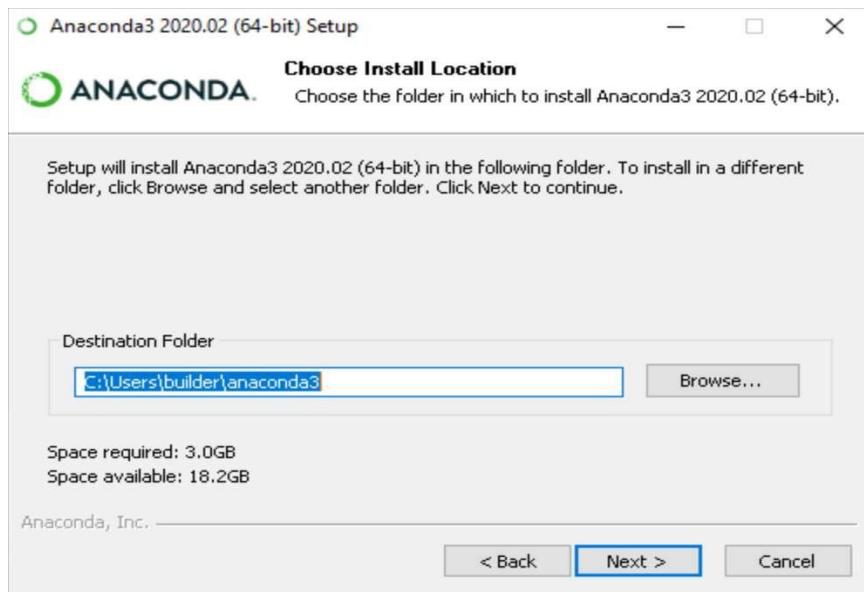
5. Read the license agreement and click on "I Agree".



6. Select installation type “**Just Me**” unless you’re installing it for all users (which require Windows Administrator privileges) and click on **Next**.



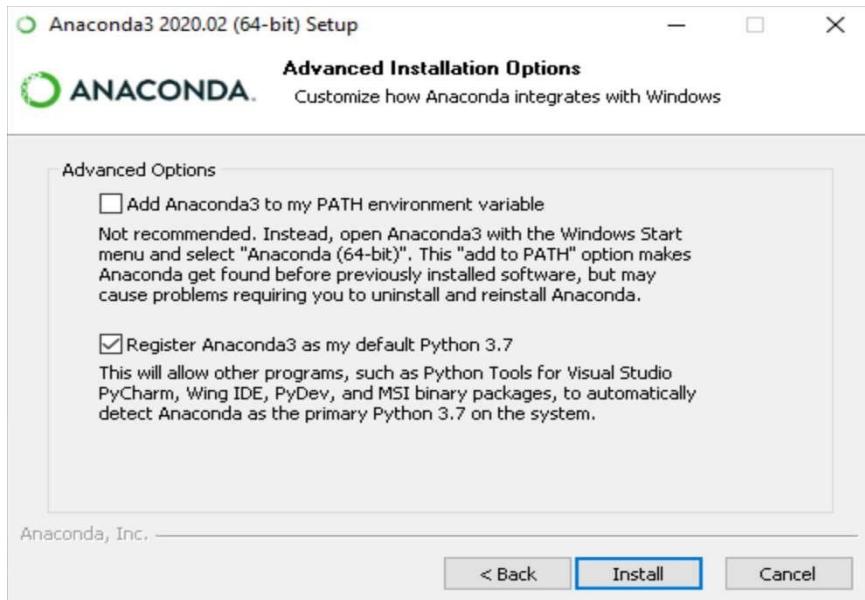
7. Select a destination folder to install Anaconda and click the **Next** button.



8. Choose whether to add Anaconda to your **PATH** environment variable. We recommend NOT adding Anaconda to the **PATH** environment variable, since this can interfere with other softwares. Instead, use Anaconda software by opening [Anaconda Navigator](#) or the [Anaconda Prompt](#) from the Start Menu

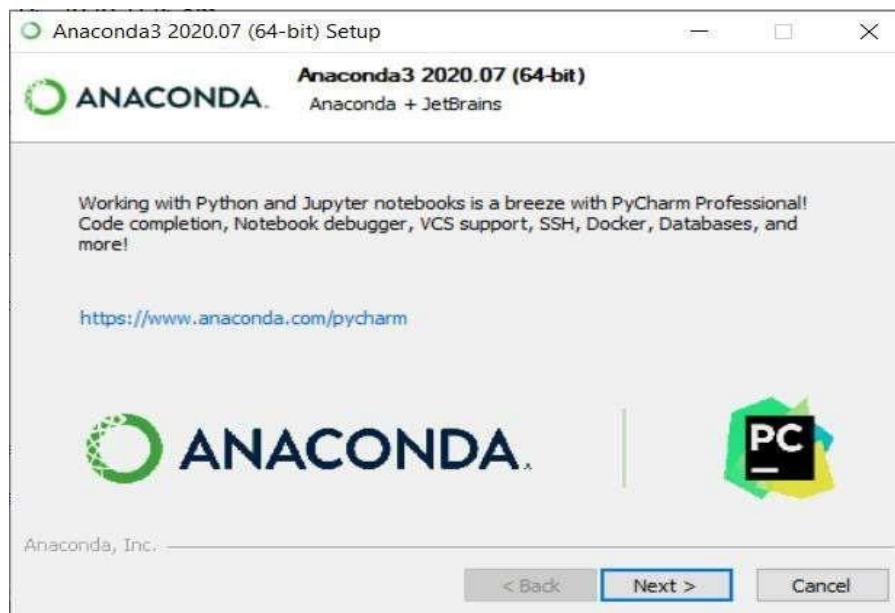
Choose whether to register Anaconda as your default Python. Unless you plan to install and run multiple versions of Anaconda or multiple versions of Python, accept the default version and leave this box checked.

9. Click the **Install** button.

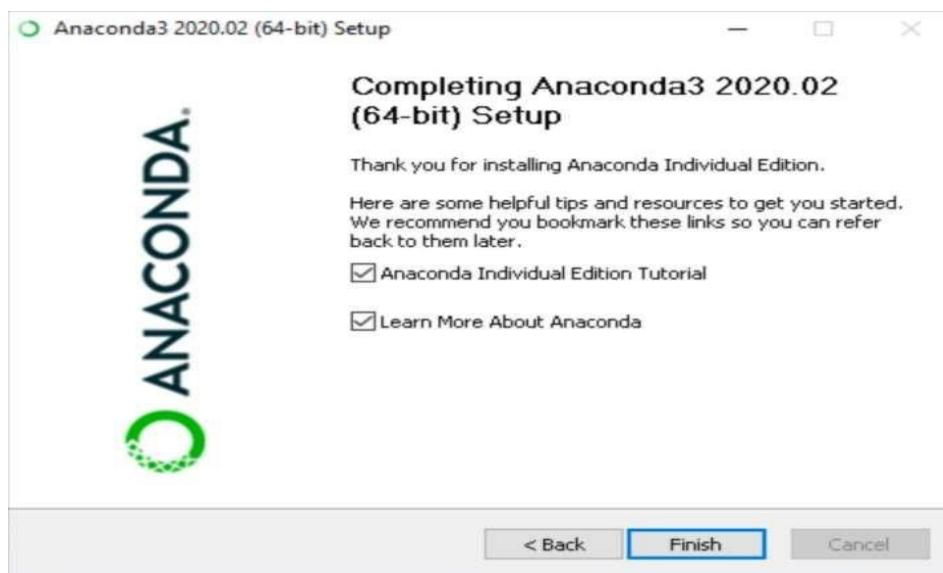


If you want to watch the packages Anaconda is installing, click on **Show Details**.

10. Click on the **Next** button.

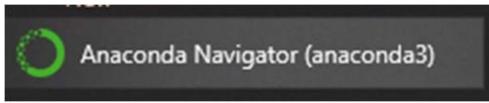


11. And then click the **Finish** button.

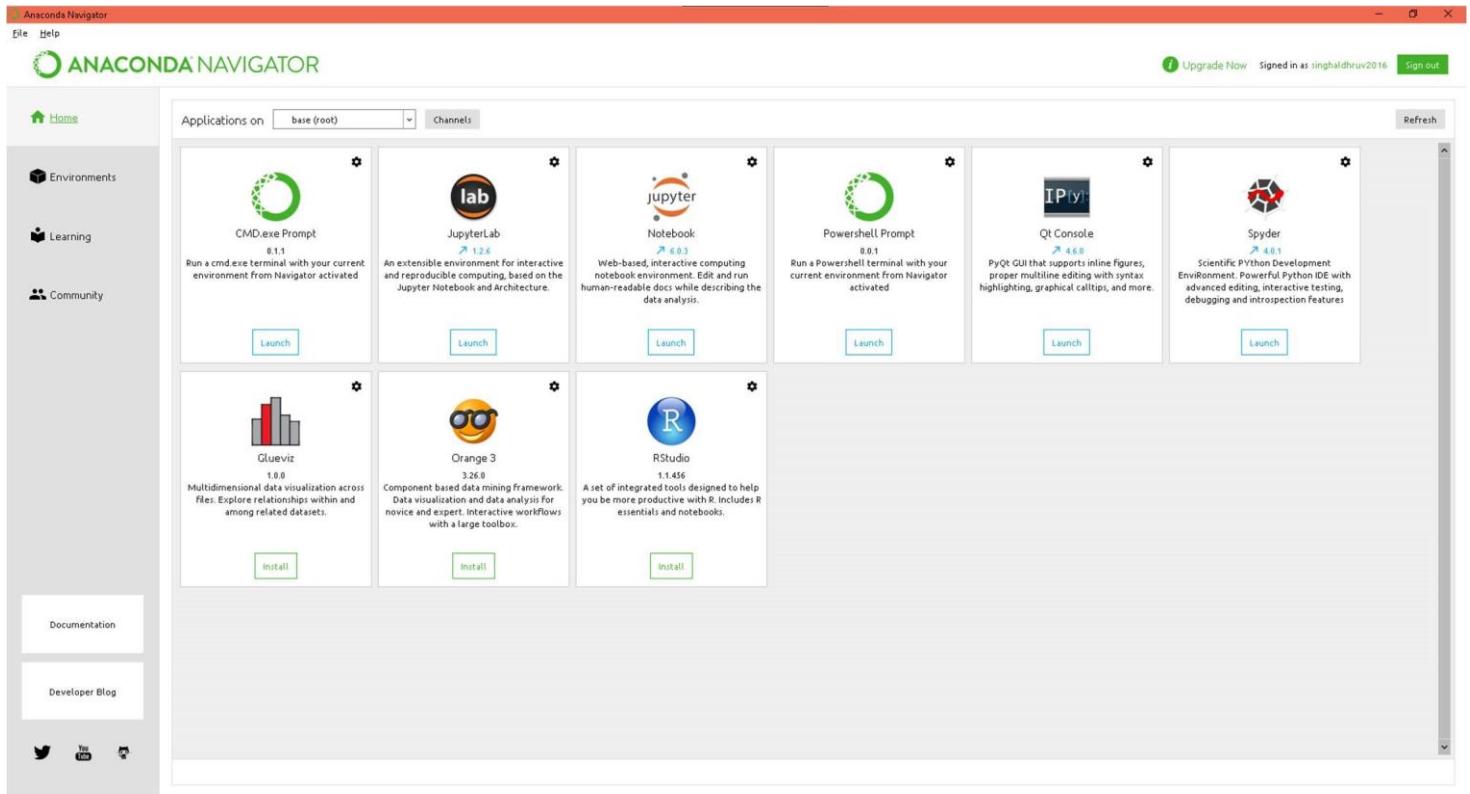


12. After a successful installation you will see the "**Thanks for installing Anaconda**" dialog box.

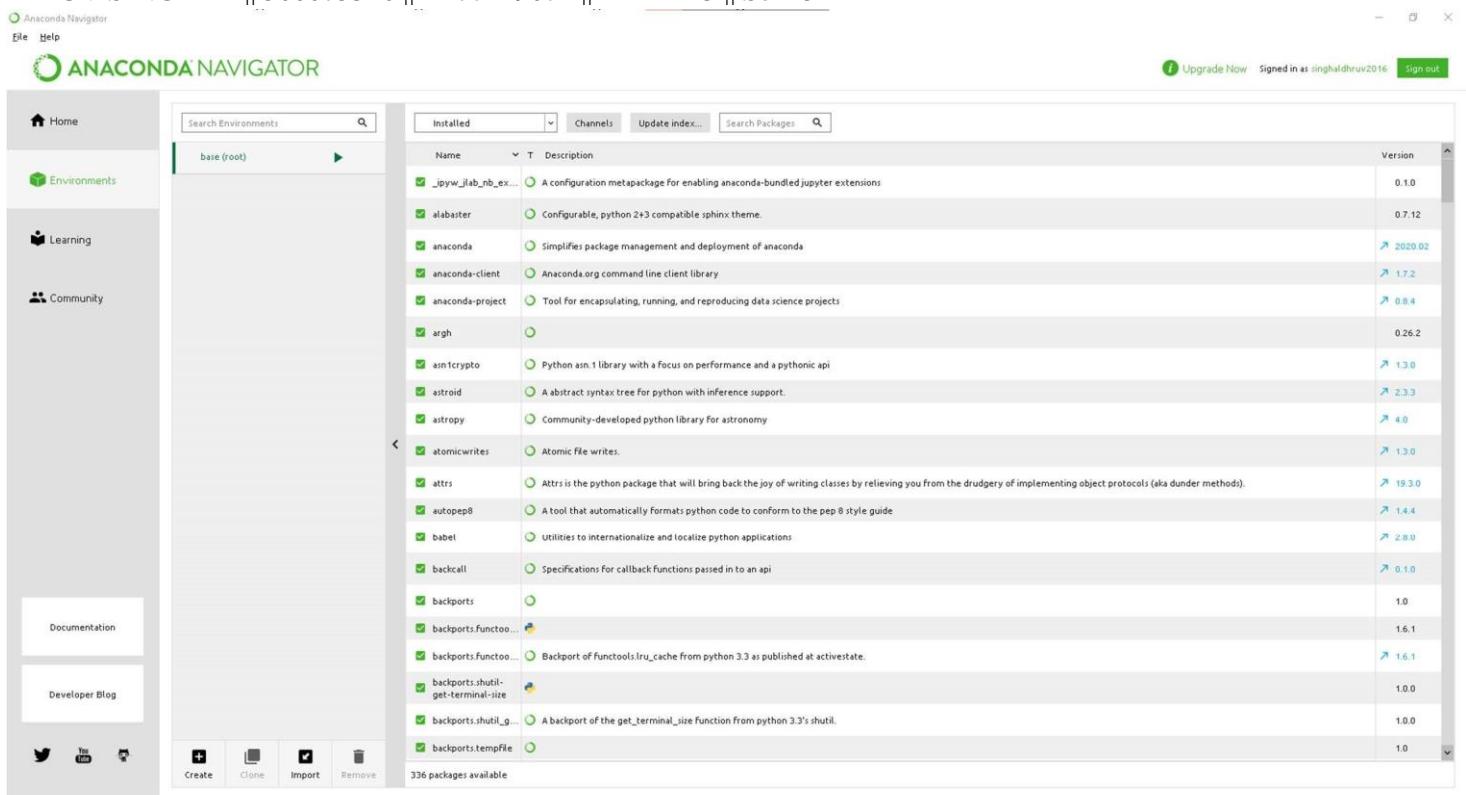
13. Open Anaconda navigator.



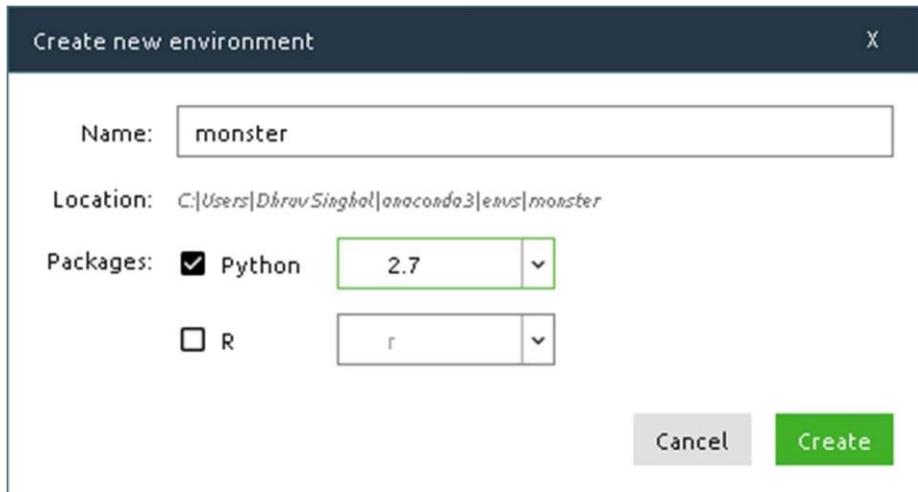
14. This screen will pop up.



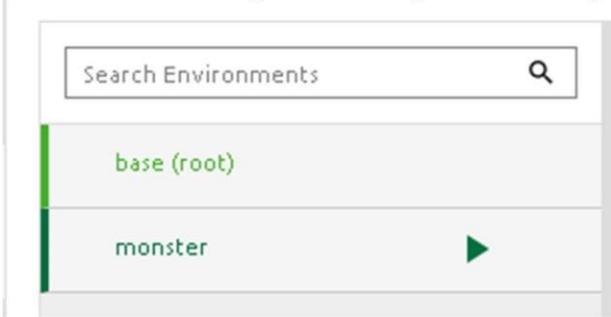
15. Go to “Environments” from Left Side panel.



16. From Bottom Left click “Create” and choose version and name of the environment.



17. Create environment.



18. Activate Environment from Anaconda Prompt and Launch Spyder(After installing From Anaconda Navigator)

Spyder:

Spyder, the Scientific Python Development Environment, is a free integrated development environment (IDE) that is included with Anaconda.

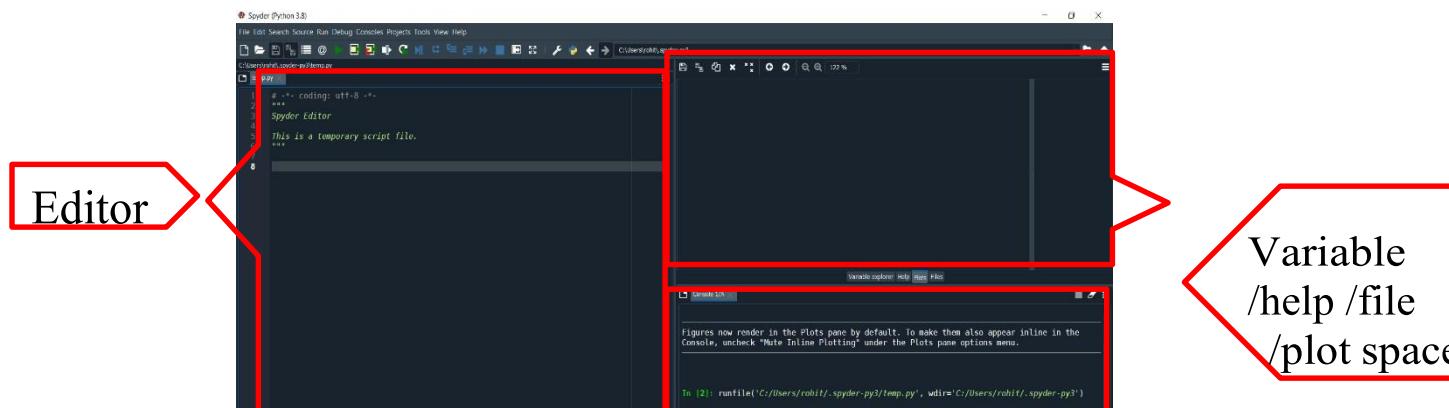
It includes:

- Editing
- Interactive testing
- Debugging
- Introspection features

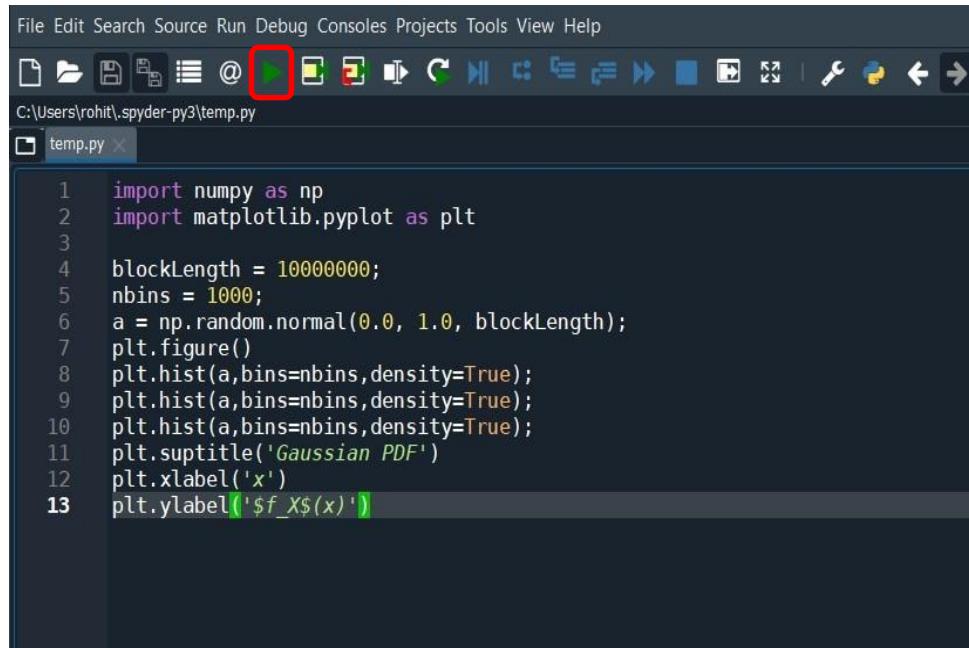
Steps for Spyder setup and run a test code:

1. In Window search box, type **Spyder** and press **Enter**. 2. Spyder IDE opened and you can see a total of 3 area:

- a. Editor
- b. Console
- c. Variable/help/file/plot space



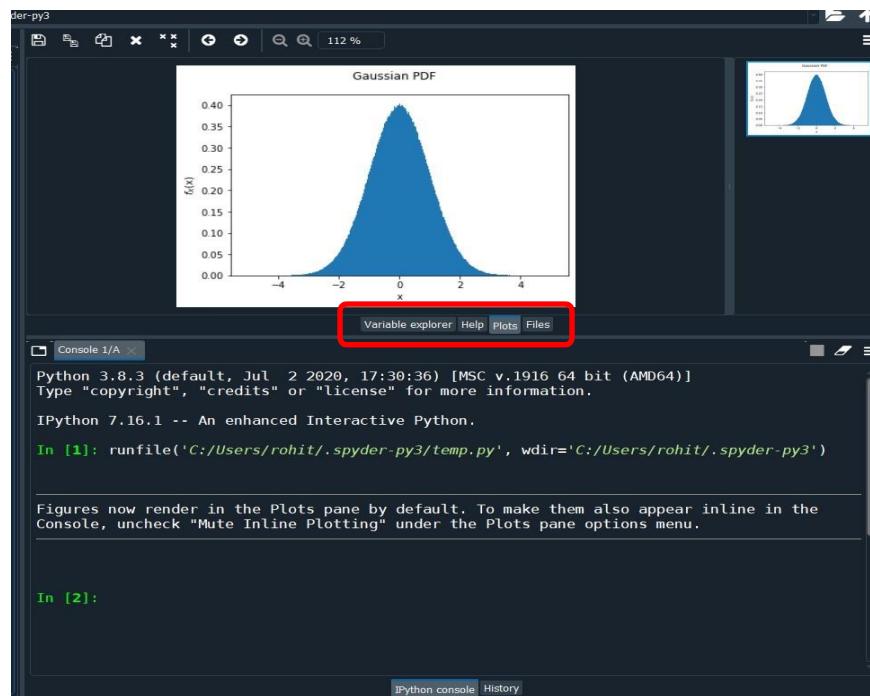
3. Let's write a test code in the Editor and run the code by clicking on Run button.



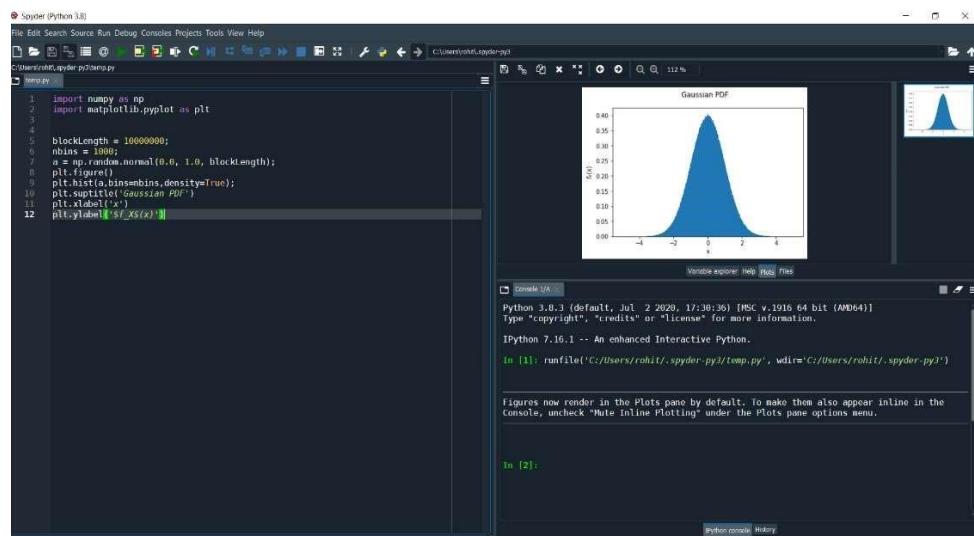
The screenshot shows the Spyder Python IDE interface. The menu bar includes File, Edit, Search, Source, Run, Debug, Consoles, Projects, Tools, View, and Help. The toolbar contains various icons for file operations, including a green play button icon which is circled in red to indicate it's the 'Run' button. The current working directory is C:\Users\rohit\spyder-py3\temp.py. A single file tab labeled 'temp.py' is open. The code in the editor is:

```
1 import numpy as np
2 import matplotlib.pyplot as plt
3
4 blockLength = 10000000;
5 nbins = 1000;
6 a = np.random.normal(0.0, 1.0, blockLength);
7 plt.figure()
8 plt.hist(a,bins=nbins,density=True);
9 plt.hist(a,bins=nbins,density=True);
10 plt.hist(a,bins=nbins,density=True);
11 plt.suptitle('Gaussian PDF')
12 plt.xlabel('x')
13 plt.ylabel('$f_X(x)$')
```

4. You can see the variable, plot, files on right side of IDE by clicking appropriate tabs as highlighted with Red color below.



5. As a whole **Spyder screen** looks like as below.



Experiment 2

Spyder (Python 3.7)

File Edit Search Source Run Debug Consoles Projects Tools View Help

B:\3rd year\5th sem\P&AD lab\Labexp1.py

labexp1.2.py x Labexp1.py x

```
1 import pandas as pd
2 #%% Q1
3 fruits=pd.DataFrame({'Apples':[30], 'Bananas':[21]})
4 print(fruits)
5
6 #%% Q2
7 data={'Apples':[35,41], 'Bananas':[21,34]}
8 fruit_sales=pd.DataFrame(data,index=['2017 Sales','2018 Sales'])
9 print(fruit_sales)
10
11 #%% Q3
12 ingredients={ 'Dinner':pd.Series(['4 cups','1 cup','2 Large','1 can'],
13                                index=['Flour','Milk','Eggs','Spam'])}
14 df=pd.DataFrame(ingredients)
15 print(df['Dinner'])
16
17 #%% Q4
18 data=pd.read_csv('C:/Users/Dhruv Singhania/winemag-data_first150k.csv')
19 df=pd.DataFrame(data)
20 print(df)
21
22 #%% Q5
23 animals = pd.DataFrame({'Cows': [12, 20], 'Goats': [22, 19]}, index=['Year 1', 'Year 2'])
24 print(animals)
25 print(type(animals))
26 animals.to_csv("B:/3rd year/5th sem/P&AD lab/hello.csv")
27
```

Source Console Object

Usage

Help Variable Explorer Plots Files

Console 5/A x

```
Python 3.7.6 (default, Jan  8 2020, 20:23:39) [MSC v.1916 64 bit (AMD64)]
Type "copyright", "credits" or "license" for more information.

IPython 7.26.0 -- An enhanced Interactive Python.

In [1]: runcell(0, 'B:/3rd year/5th sem/P&AD Lab/Labexp1.py')

In [2]: runcell('Q1', 'B:/3rd year/5th sem/P&AD Lab/Labexp1.py')
      Apples Bananas
0       30      21

In [3]: runcell('Q2', 'B:/3rd year/5th sem/P&AD Lab/Labexp1.py')
      Apples Bananas
2017 Sales   35     21
2018 Sales   41     34

In [4]: runcell('Q3', 'B:/3rd year/5th sem/P&AD Lab/Labexp1.py')
Flour    4 cups
Milk     1 cup
Eggs    2 large
Spam     1 can
Name: Dinner, dtype: object

In [5]: runcell('Q4', 'B:/3rd year/5th sem/P&AD Lab/Labexp1.py')
      Unnamed: 0  country ... variety          winery
0            0    US    ... Cabernet Sauvignon  Heitz
1            1    Spain   ... Tinta de Toro  Bodega Carmen Rodriguez
2            2    US    ... Sauvignon Blanc  Macauley
3            3    US    ... Pinot Noir        Ponzi
4            4    France   ... Provence red blend  Domaine de la Béguade
...
150925    150925  Italy    ... White Blend  Feudi di San Gregorio
150926    150926  France   ... Champagne Blend  H.Germain
150927    150927  Italy    ... White Blend  Terredora
150928    150928  France   ... Champagne Blend  Gosset
150929    150929  Italy    ... Pinot Grigio  Alois Lageder
[150930 rows x 11 columns]

In [6]: runcell('Q5', 'B:/3rd year/5th sem/P&AD Lab/Labexp1.py')
      Cows Goats
Year 1    12    22
Year 2    20    19
<class 'pandas.core.frame.DataFrame'>
```

LSP Python: ready conda (Python 3.7.6) Line 27, Col 1 ASCII CRLF RW Mem 44%

Pandas

```

24 print(animals)
25 print(type(animals))
26 animals.to_csv("B:/3rd year/5th sem/P&AD lab/hello.csv")
27
28
29 reviews=pd.read_csv('C:/Users/Dhruv Singhal/winemag-data_first150k.csv')
30 reviews.head()
31
32 #%% Q1: Select the description column from reviews and assign the result to the variable desc
33 desc = reviews.description
34 print(desc)
35
36 #Q2: Select the first value from the description column of reviews , assigning it to the variable
37 first_desc=reviews.description.iloc[0]
38 print(first_desc)
39
40 #Q3: Select the first row of data from reviews , assigning it to the variable
41 first_row=reviews.iloc[0]
42 print(first_row)
43
44
45 #Q4: Select the first 10 values from description column in review, assigning the result to the variable
46 first_desc=reviews['description'].head(10)
47 print(first_desc)
48
49
50 #Q5: Select the records with index labels 1,2,3,5,8 assigning the result to variable
51 Sample_reviews=reviews.iloc[[1,2,3,5,8]]
52 print(Sample_reviews)
53
54
55
56
57
58
59
60
61
62
63
64
65
66
67
68
69

```

```

In [3]: runcell('q1: select the description column from reviews and assign the result to the variable desc', 'B:/3rd year/5th sem/P&AD lab/Labexp1.py')
0 This tremendous 98% varietal wine hails from ...
1 Ripe aromas of fig, blackberry and cassis are ...
2 Mac Watson honors the memory of a wine once ma...
3 This spent 20 months in 38% new French oak, an...
4 This is the top wine from La Béguile, named aft...
...
158925 Many people feel France represents southern Italy...
158926 Off-the-charts ripe figs with a hint of lime and...
158927 This classic red wine comes from a small vineyard...
158928 A perfect salmon shade, with scents of peaches...
158929 More Pinot Grigios should taste like this. A m...
Name: description, Length: 158930, dtype: object

```

```
In [4]: runcell('q1: select the description column from reviews and assign the result to the variable desc', 'B:/3rd year/5th sem/P&AD lab/Labexp1.py')
```

This tremendous 98% varietal wine hails from Oakville and was aged over three years in oak. Juicy red-cherry fruit and a compelling hint of caramel greet the palate, framed by elegant, fine tannins and a subtle minty tone in the background. Balanced and rewarding from start to finish, it has years ahead of it to develop further nuance. Enjoy 2022-2038.

```
In [5]: runcell('q1: select the description column from reviews and assign the result to the variable desc', 'B:/3rd year/5th sem/P&AD lab/Labexp1.py')
```

Unnamed: 0	country	description	designation	points	price	province	region_1	region_2	variety	winery
0	US	This tremendous 98% varietal wine hails from ...	Martha's Vineyard	96	235.0	California	Napa Valley	Napa	Cabernet Sauvignon	Hetz

```
In [6]: runcell('q1: select the description column from reviews and assign the result to the variable desc', 'B:/3rd year/5th sem/P&AD lab/Labexp1.py')
```

Unnamed: 0	country	variety	winery
0	Spain	Tinta de Toro	Bodega Carmen Rodriguez
1	US	Sauvignon Blanc	MacAuley
2	US	Pinot Noir	Ponzi
3	Spain	Tinta de Toro	Numanthia
4	US	Pinot Noir	Bengström

[5 rows x 11 columns]

```

54: #%%%
55: #06: Create a variable df containing the country,provience,region_1 and region_2 column of records
56: cols=['country','provience','region_1','region_2']
57: indices=[0,1,10,100]
58: df=reviews.loc[indices,cols]
59: print(df)
60: #%%%
61: #07: Create a variable containing the country and variety col of first 100 records
62: cols_ine=['country','variety']
63: df=reviews.loc[:100,cols_ine]
64: print(df)
65: #%%%
66: #08: create a dataframe italian_wines containing reviews of wine made3 in italy
67: italian_wine=reviews[reviews['country']=='Italy']
68: print(italian_wine)
69: #%%%
70: #09:create a dataframe containing all reviews with at least 95 points (out of 100) for wines from Australia
71: top_oceania_wines = reviews.loc[reviews.country.isin(['Australia', 'New Zealand']) & (reviews.points >= 95)]
72: print(top_oceania_wines)
73:
74:
75:
76:
77:
78:
79:
80:
81:
82:
83:
84:
85:
86:
87:
88:
89:
90:
91:
92:
93:
94:
95:
96:
97:
98:
99:
100:
101:
102:
103:
104:
105:
106:
107:

```

Name	Type	Size	Value
cols	list	4	['country', 'provience', 'region_1', 'region_2']
cols_ine	list	2	['country', 'variety']

Console 10/A x

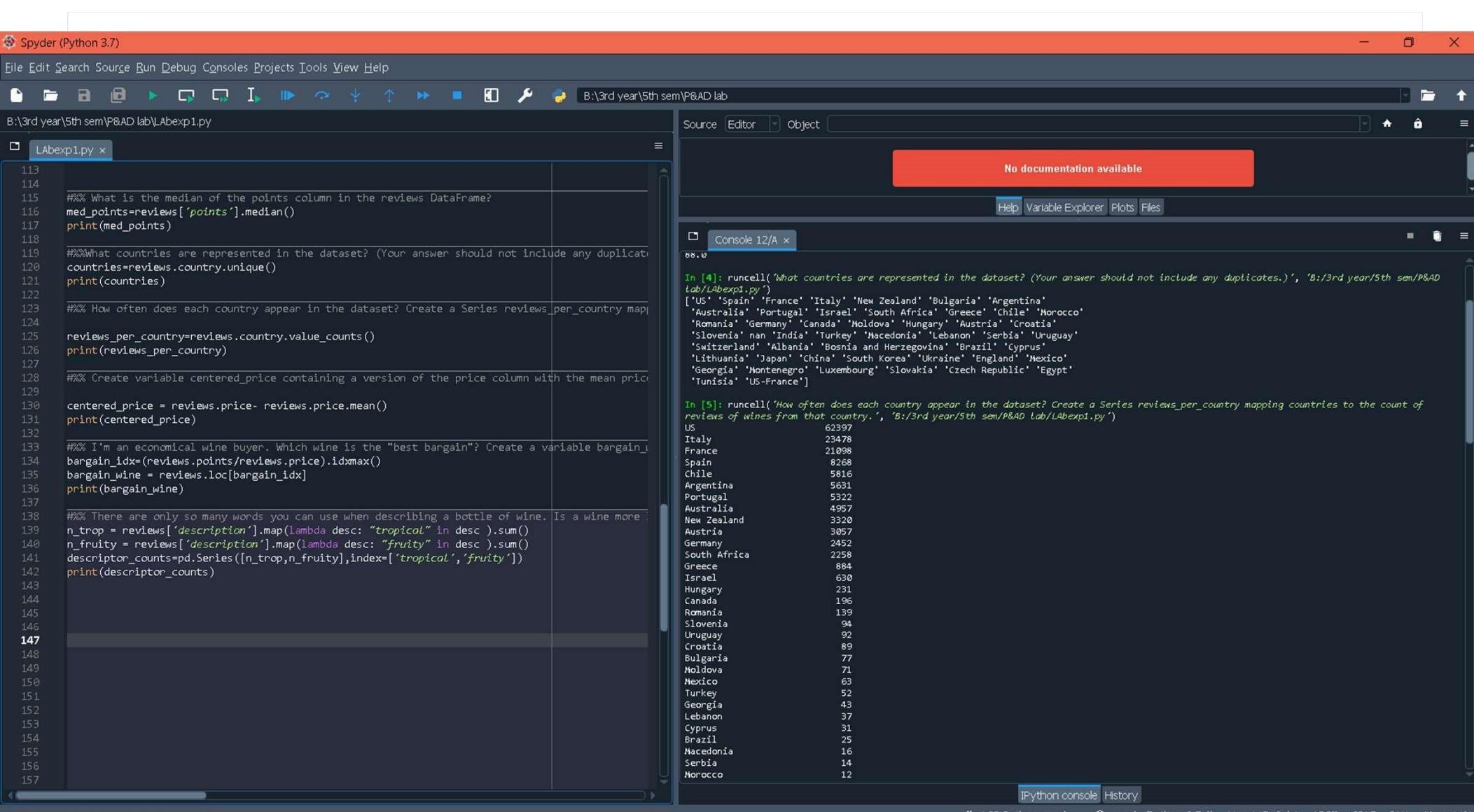
```
In [13]: runcell15 'B:\3rd year\5th sem\P&AD lab\Labexp1.py'
In [13]: runcell15 'B:\3rd year\5th sem\P&AD lab\Labexp1.py'
Unamed: 0 country ... variety winery
2148 2148 Australia ... Shiraz Torbreck
2458 2458 Australia ... Cabernet-Shiraz Hickinbotham
3833 3833 Australia ... Cabernet Sauvignon Penfolds
3844 3844 Australia ... Shiraz Henschke
3847 3847 Australia ... Red Blend Heartland
...
122779 122779 Australia ... Shiraz Standish
127634 127634 Australia ... Shiraz Henschke
137383 137383 Australia ... Syrah Clarendon Hills
159562 159562 Australia ... Shiraz Penfolds
159563 159563 Australia ... Shiraz Penfolds
[82 rows x 11 columns]
```

```
In [14]: runcell15 'B:\3rd year\5th sem\P&AD lab\Labexp1.py'
In [14]: runcell15 'B:\3rd year\5th sem\P&AD lab\Labexp1.py'
country variety
0 US Cabernet Sauvignon
1 Spain Tinta de Toro
2 US Sauvignon Blanc
3 US Pinot noir
4 France Provence red blend
...
96 US Chardonnay
97 US Cabernet Sauvignon
98 France Merlot-Velbec
99 France Ugni Blanc-Colombard
100 US Viognier
[101 rows x 2 columns]
```

```
In [15]: runcell15 'B:\3rd year\5th sem\P&AD lab\Labexp1.py'
In [15]: runcell15 'B:\3rd year\5th sem\P&AD lab\Labexp1.py'
Unamed: 0 country ... variety winery
38 38 Italy ... Friuli ... Borgo del Taglio
32 32 Italy ... Sangiovese Abbadia Ardenga
35 35 Italy ... Sangiovese Carillon
37 37 Italy ... Sangiovese Avignonesi
38 38 Italy ... Sangiovese Casina di Cornia
...
159928 159928 Italy ... Champagne Blend Letrari
159922 159922 Italy ... Tokai Ronchi di Manzano
159925 159925 Italy ... White Blend Feudi di San Gregorio
159927 159927 Italy ... White Blend Terredora
159929 159929 Italy ... Pinot Grigio Alois Lageder
[23478 rows x 11 columns]
```

```
In [16]: runcell15 'B:\3rd year\5th sem\P&AD lab\Labexp1.py'
In [16]: runcell15 'B:\3rd year\5th sem\P&AD lab\Labexp1.py'
Unamed: 0 country ... variety winery
2148 2148 Australia ... Shiraz Torbreck
2458 2458 Australia ... Cabernet-Shiraz Hickinbotham
3833 3833 Australia ... Cabernet Sauvignon Penfolds
3844 3844 Australia ... Shiraz Henschke
3847 3847 Australia ... Red Blend Heartland
...
122779 122779 Australia ... Shiraz Standish
[82 rows x 11 columns]
```

IPython console History



B:\3rd year\5th sem\P&AD lab\Labexp1.py

```

93 #%% What is the best wine I can buy for a given amount of money? Create a Series whose index is wine prices an
94 best_rating_per_price = reviews.groupby('price')['points'].max().sort_index()
95 print(best_rating_per_price)
96
97 #%% What are the minimum and maximum prices for each variety of wine? Create a DataFrame whose index is the v
98 price_extremes = reviews.groupby('variety').price.agg([min, max])
99 print(price_extremes)
100
101 #%% What are the most expensive wine varieties? Create a variable sorted_varieties containing a copy of the
102 sorted_varieties = price_extremes.sort_values(by=['min', 'max'], ascending=False)
103 print(sorted_varieties)
104
105 #%% What combination of countries and varieties are most common? Create a Series whose index is a MultiIndex
106 country_variety_counts = reviews.groupby(['country', 'variety']).size().sort_values(ascending=False)
107 print(country_variety_counts)
108
109
110
111
112
113
114
115
116
117
118
119
120
121
122
123
124
125
126
127
128
129
130
131
132
133
134
135
136
137
138
139
140
141
142
143
144
145
146

```

B:\3rd year\5th sem\P&AD lab\Labexp1.py

Name	Type	Size	Value
best_rating_per_price	Series	(357,)	Series object of pandas.core.series module
cols	List	A	['country', 'variety', 'rating_1', 'rating_2']

Console 10/A x

```

2013.0 91
2000.0 99
Name: points, Length: 357, dtype: int64

In [10]: runcell('What are the minimum and maximum prices for each variety of wine? Create a DataFrame whose index is the variety category from the dataset and whose values are the min and max values thereof.', 'B:\3rd year\5th sem\P&AD lab\Labexp1.py')
          min      max
variety
Aegeanpotiski    8.0   65.0
Aglianico        6.0  130.0
Aidani            27.0   27.0
Airen             8.0   30.0
Albana            8.0   66.0
...
Zierfandler-Rotgipfler  28.0   25.0
Zinfandel         4.0  300.0
Zlahtina          13.0   17.0
 Zweigelt          9.0   70.0
Zilavka           13.0   15.0

[632 rows x 2 columns]

In [11]: runcell('What are the most expensive wine varieties? Create a variable sorted_varieties containing a copy of the dataframe from the previous question where varieties are sorted in descending order based on minimum price, then on maximum price (to break ties).', 'B:\3rd year\5th sem\P&AD lab\Labexp1.py')
          min      max
variety
Cabernet-Shiraz  159.0  159.0
Mazuelo           92.0  285.0
Carignan-Syrah   38.0   38.0
Syrah-Cabernet Franc  60.0  60.0
Nasco              65.0  65.0
...
Rabigato           NaN   NaN
Sacy               NaN   NaN
Sauvignon Blanc-Sauvignon Gris  NaN   NaN
Terret Blanc       NaN   NaN
Zelen              NaN   NaN

[632 rows x 2 columns]

In [12]: runcell('What combination of countries and varieties are most common? Create a Series whose index is a MultiIndex of (country, variety) pairs. For example, a pinot noir produced in the US should map to ("US", "Pinot Noir"). Sort the values in the series in descending order based on wine count.', 'B:\3rd year\5th sem\P&AD lab\Labexp1.py')
country  variety
US        Pinot Noir      38348
                  Cabernet Sauvignon  9178
                  Chardonnay        8127
France    Bordeaux-style Red Blend  4908
US        Syrah            4274
...
France    Pied de Pendix      1
Tunisia   White Blend       1
Rose                 1
Switzerland  White Blend     1
US        Carignan-Grenache     1
Length: 3475, dtype: int64

In [13]:

```

IPython console History

LSP Python: ready conda (Python 3.7.6) Line 122, Col 1 ASCII CRLF RW Mem 45%



LABexp1.py x

```
79
80 %% What is the data type of the points column in the dataset?
81 print(reviews['points'].dtype)
82 %% Create a Series from entries in the points column, but convert the entries to strings. Hint: str
83 point_strings=reviews['points'].astype(str)
84 print(type(point_strings[0]))
85 %% Sometimes the price column is null. How many reviews in the dataset are missing a price?
86 n_missing_prices=reviews['price'].isnull().sum()
87 print(n_missing_prices)
88 %% What are the most common wine-producing regions? Create a Series counting the number of times ea
89 rev_region=reviews['region_1'].fillna('Unknown').value_counts().sort_values(ascending=False)
90 print(rev_region)
91
92
93
94
95
96
97
98
99
100
101
102
103
104
105
106
107
108
109
110
111
112
113
114
115
116
117
118
119
120
121
122
123
```

Name	Type	Size	Value
n_missing_prices	int64	1	13695
point_strings	Series	(1500000)	Series object of pandas.core.series module

Console 10/A x

In [8]: runcell('What is the data type of the points column in the dataset?', 'B:/3rd year/5th sem/P&AD lab/LABexp1.py')
int64

In [9]: runcell('Create a Series from entries in the points column, but convert the entries to strings. Hint: strings are str in native Python.', 'B:/3rd year/5th sem/P&AD Lab/LABexp1.py')
<class 'str'>

In [10]: runcell('Sometimes the price column is null. How many reviews in the dataset are missing a price?', 'B:/3rd year/5th sem/P&AD Lab/LABexp1.py')
13695

In [11]: runcell('What are the most common wine-producing regions? Create a Series counting the number of times each value occurs in the region_1 field. This field is often missing data, so replace missing values with Unknown. Sort in descending order. Your output should look something like this:', 'B:/3rd year/5th sem/P&AD Lab/LABexp1.py')

Unknown	25060
Napa Valley	6209
Columbia Valley (WA)	4975
Mendoza	3586
Russian River Valley	3571
	...
Vin de Pays de Côtes du Tarn	1
Clos de Lambrays	1
California-Oregon	1
Mâcon-Mâconey	1
Coteaux du Tricastin	1
Name: region_1, Length: 1237, dtype: int64	

In [12]:

IPython console History

Spyder (Python 3.7)

File Edit Search Source Run Debug Consoles Projects Tools View Help

B:\3rd year\5th sem\P&AD lab\Labexp1.py

Source Editor Object

arr

Help Variable Explorer Plots Files

Console 10/A ×

Python 3.7.6 (default, Jan 8 2020, 20:23:39) [MSC v.1916 64 bit (AMD64)]
Type "copyright", "credits" or "license" for more information.

IPython 7.26.0 -- An enhanced Interactive Python.

In [1]: runcell(0, 'B:/3rd year/5th sem/P&AD Lab/Labexp1.py')

In [2]: runcell(6, 'B:/3rd year/5th sem/P&AD Lab/Labexp1.py')

In [3]: runcell('region_1 and region_2 are pretty uninformative names for locale columns in the dataset.
Create a copy of reviews with these columns renamed to region and locale, respectively.', 'B:/3rd year/5th
sem/P&AD Lab/Labexp1.py')

	Unnamed: 0	country	...	variety	winery
0	0	US	...	Cabernet Sauvignon	Heitz
1	1	Spain	...	Tinta de Toro	Bodega Carmen Rodriguez
2	2	US	...	Sauvignon Blanc	MacCauley
3	3	US	...	Pinot Noir	Ponzi
4	4	France	...	Provence red blend	Domaine de la Béguade

[5 rows x 11 columns]

In [4]: runcell('Set the index name in the dataset to wines.', 'B:/3rd year/5th sem/P&AD Lab/Labexp1.py')

wine	Unnamed: 0	country	...	variety	winery
0	0	US	...	Cabernet Sauvignon	Heitz
1	1	Spain	...	Tinta de Toro	Bodega Carmen Rodriguez
2	2	US	...	Sauvignon Blanc	MacCauley
3	3	US	...	Pinot Noir	Ponzi
4	4	France	...	Provence red blend	Domaine de la Béguade

[5 rows x 11 columns]

In [5]:

IPython console History

LSP Python: ready conda (Python 3.7.6) Line 97, Col 1 ASCII CRLF RW Mem 45%

```
73 #%% region_1 and region_2 are pretty uninformative names for locale columns in the dataset. Create a
74 renamed=reviews.rename(columns=dict(region_1='region',region_2='Locale'))
75 print(renamed.head())
76 #%% Set the index name in the dataset to wines.
77 reindexed=reviews.rename_axis('wine',axis=0)
78 print(reindexed.head())
79
80
81
82
83
84
85
86
87
88
89
90
91
92
93
94
95
96
97
98
99
100
101
102
103
104
105
106
107
108
109
110
111
112
113
114
115
116
117
```

NumPy

Spyder (Python 3.7)

File Edit Search Run Debug Consoles Projects Tools View Help

B:\3rd year\5th sem\P&AD lab\labexp1.2.py

labexp1.2.py x LABexp1.py x

```
1 import numpy as np
2 #%% Insert the correct method for creating a NumPy array.
3 arr = np.array([1, 2, 3, 4, 5])
4 print(arr)
5 #%% Insert the correct argument for creating a NumPy array with 2 dimensions.
6 arr = np.array([1, 2, 3, 4], ndmin=2)
7 print(arr)
8 #%% Insert the correct syntax for checking the number of dimension of a NumPy array.
9 arr = np.array([1, 2, 3, 4])
10 print(arr.ndim)
11 #%% Insert the correct syntax for printing the first item in the array.
12 arr = np.array([1, 2, 3, 4, 5])
13 print(arr[0])
14 #%% Insert the correct syntax for printing the number 50 from the array.
15 arr = np.array([1, 2, 3, 4, 5])
16 print(arr[4])
17 #%% Insert the correct syntax for printing the number 50 from the array.
18 arr = np.array([[10, 20, 30, 40], [50, 60, 70, 80]])
19 print(arr[1,0])
20 #%% Use negative index to print the last item in the array.
21 arr = np.array([10, 20, 30, 40, 50])
22 print(arr[-1])
23 #%% Insert the correct slicing syntax to print the following selection of the array:
24 #Everything from (including) the second item to (not including) the fifth item.
25 arr = np.array([10, 15, 20, 25, 30, 35, 40])
26 print(arr[1:4])
27 #%% Insert the correct slicing syntax to print the following selection of the array:
28 #Everything from (including) the third item to (not including) the fifth item.
29 arr = np.array([10, 15, 20, 25, 30, 35, 40])
30 print(arr[2:4])
31 #%% Insert the correct slicing syntax to print the following selection of the array:
32 #Every other item from (including) the second item to (not including) the fifth item.
33 arr = np.array([10, 15, 20, 25, 30, 35, 40])
34 print(arr[1:5:2])
35 #%% Insert the correct slicing syntax to print the following selection of the array:
36 #Every other item from the entire array
37 arr = np.array([10, 15, 20, 25, 30, 35, 40])
38 print(arr[::-2])
39
40
41
42
43
44
45
46
```

Source Editor Object

arr

Help Variable Explorer Plots Files

Console 9/A x

```
Python 3.7.6 (default, Jan  8 2020, 20:23:39) [MSC v.1916 64 bit (AMD64)]
Type "copyright", "credits" or "license" for more information.

IPython 7.26.0 -- An enhanced Interactive Python.

In [1]: runcell(0, 'B:/3rd year/5th sem/P&AD Lab/labexp1.2.py')
In [2]: runcell('Insert the correct method for creating a NumPy array.', 'B:/3rd year/5th sem/P&AD Lab/labexp1.2.py')
[1 2 3 4 5]

In [3]: runcell('Insert the correct argument for creating a NumPy array with 2 dimensions.', 'B:/3rd year/5th sem/P&AD Lab/labexp1.2.py')
[[1 2 3 4 5]]

In [4]: runcell('Insert the correct syntax for checking the number of dimension of a NumPy array.', 'B:/3rd year/5th sem/P&AD Lab/labexp1.2.py')
1

In [5]: runcell('Insert the correct syntax for printing the first item in the array.', 'B:/3rd year/5th sem/P&AD Lab/labexp1.2.py')
1

In [6]: runcell('Insert the correct syntax for printing the number 50 from the array., #1', 'B:/3rd year/5th sem/P&AD Lab/labexp1.2.py')
5

In [7]: runcell('Insert the correct syntax for printing the number 50 from the array., #2', 'B:/3rd year/5th sem/P&AD Lab/labexp1.2.py')
50

In [8]: runcell('Use negative index to print the last item in the array.', 'B:/3rd year/5th sem/P&AD Lab/labexp1.2.py')
50

In [9]: runcell('Insert the correct slicing syntax to print the following selection of the array:, #1', 'B:/3rd year/5th sem/P&AD Lab/labexp1.2.py')
[15 20 25]

In [10]: runcell('Insert the correct slicing syntax to print the following selection of the array:, #2', 'B:/3rd year/5th sem/P&AD Lab/labexp1.2.py')
[20 25]

In [11]: runcell('Insert the correct slicing syntax to print the following selection of the array:, #3', 'B:/3rd year/5th sem/P&AD Lab/labexp1.2.py')
[15 25]

In [12]: runcell('Insert the correct slicing syntax to print the following selection of the array:, #4', 'B:/3rd year/5th sem/P&AD Lab/labexp1.2.py')
[10 20 30 40]

In [13]:
```

IPython console History

LSP Python: ready conda (Python 3.7.6) Line 39, Col 1 ASCII CRLF RW Mem 43%

B:\3rd year\5th sem\P&AD lab\labexp1.2.py

labexp1.2.py x Labexp1py x

```
48
49 %% NumPy uses a character to represent each of the following data types, which one?
50 # i = integer
51 # b = boolean
52 # u = unsigned integer
53 # f = float
54 # c = complex float
55 # m = timedelta
56 # M = datetime
57 # O= object
58 # S = string
59 %% Insert the correct NumPy syntax to print the data type of an array.
60 arr = np.array([1, 2, 3, 4])
61 print(arr.dtype)
62 %% Insert the correct argument to specify that the array should be of type STRING.
63 arr = np.array([1, 2, 3, 4], dtype='S')
64 print(arr)
65 %% Insert the correct method to change the data type to integer.
66 arr = np.array([1.1, 2.1, 3.1])
67 newarr = arr.astype('i')
68 print(newarr)
69 %% Use the correct method to make a copy of the array.
70 arr = np.array([1, 2, 3, 4, 5])
71 x = arr.copy()
72 print(x)
73 %% Use the correct method to make a view of the array.
74 arr = np.array([1, 2, 3, 4, 5])
75 x = arr.view()
76 print(x)
77 %% Use the correct NumPy syntax to check the shape of an array.
78 arr = np.array([1, 2, 3, 4, 5])
79 print(arr.shape)
80 %% Use the correct NumPy method to change the shape of an array from 1-D to 2-D.
81 arr = np.array([1, 2, 3, 4, 5, 6, 7, 8, 9, 10, 11, 12])
82 newarr = arr.reshape(4, 3)
83 print(newarr)
84 %% Use a correct NumPy method to change the shape of an array from 2-D to 1-D.
85 arr = np.array([1, 2, 3, 4, 5, 6, 7, 8, 9, 10, 11, 12])
86 newarr = arr.reshape(-1)
87 print(newarr)
88 %% Use a correct NumPy method to join two arrays into a single array.
89 arr1 = np.array([1, 2, 3])
90 arr2 = np.array([4, 5, 6])
91 arr = np.concatenate((arr1, arr2))
92 print(arr)
```

Source Editor Object

arr

Help Variable Explorer Plots Files

Console 9/A x

In [13]: runcell('NumPy uses a character to represent each of the following data types, which one?', 'B:/3rd year/5th sem/P&AD Lab/labexp1.2.py')
Nothing to execute, this cell is empty.

In [14]: runcell('Insert the correct NumPy syntax to print the data type of an array.', 'B:/3rd year/5th sem/P&AD Lab/labexp1.2.py')
int32

In [15]: runcell('Insert the correct argument to specify that the array should be of type STRING.', 'B:/3rd year/5th sem/P&AD Lab/labexp1.2.py')
[b'1' b'2' b'3' b'4']

In [16]: runcell('Insert the correct method to change the data type to integer.', 'B:/3rd year/5th sem/P&AD Lab/labexp1.2.py')
[1 2 3]

In [17]: runcell('Use the correct method to make a copy of the array.', 'B:/3rd year/5th sem/P&AD Lab/labexp1.2.py')
[1 2 3 4 5]

In [18]: runcell('Use the correct NumPy syntax to check the shape of an array.', 'B:/3rd year/5th sem/P&AD Lab/labexp1.2.py')
(5,)

In [19]: runcell('Use the correct NumPy method to change the shape of an array from 1-D to 2-D.', 'B:/3rd year/5th sem/P&AD Lab/labexp1.2.py')
[[1 2 3]
 [4 5 6]
 [7 8 9]
 [10 11 12]]

In [20]: runcell('Use a correct NumPy method to change the shape of an array from 2-D to 1-D.', 'B:/3rd year/5th sem/P&AD Lab/labexp1.2.py')
[1 2 3 4 5 6 7 8 9 10 11 12]

In [21]: runcell('Use a correct NumPy method to join two arrays into a single array.', 'B:/3rd year/5th sem/P&AD Lab/labexp1.2.py')
[1 2 3 4 5 6]

In [22]:

IPython console History

LSP Python: ready conda (Python 3.7.6) Line 94, Col 1 ASCII CRLF RW Mem 44%

B:\3rd year\5th sem\P&AD lab\labexp1.2.py

labexp1.2.py x Labexp1.py x

```
94
95
96
97
98
99
100
101
102
103
104
105
106
107
108
109
110
111
112
113
114
115
116
117
118
119
120
121
122
123
124
125
126 #%% Use the correct NumPy method to find all items with the value 4.
127 arr = np.array([1, 2, 3, 4, 5, 4, 4])
128 x = np.where(arr == 4)
129 print(x)
130 #%% Use the correct NumPy method to return a sorted array.
131 arr = np.array([3, 2, 0, 1])
132 x = np.sort(arr)
133 print(x)
134
135
136
137
138
139
```

Source Editor Object

arr

Console 9/A x

In [23]: runcell('Use the correct NumPy method to find all items with the value 4.', 'B:/3rd year/5th sem/P&AD lab/labexp1.2.py')
(array([3, 5, 6], dtype=int64),)

In [24]: runcell('Use the correct NumPy method to return a sorted array.', 'B:/3rd year/5th sem/P&AD lab/labexp1.2.py')
[0 1 2 3]

In [25]:

IPython console History

LSP Python: ready conda (Python 3.7.6) Line 113, Col 1 ASCII CRLF RW Mem 44%

Experiment 3

AIM:

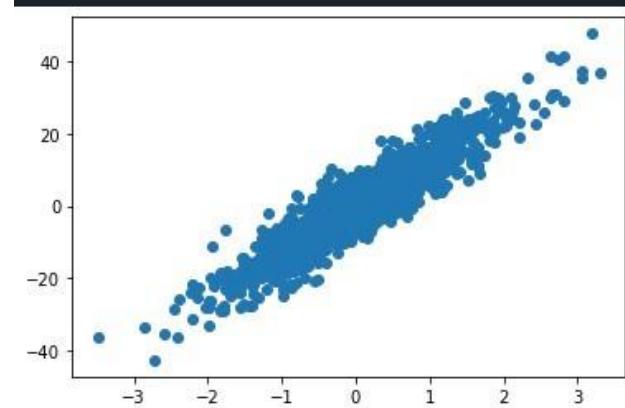
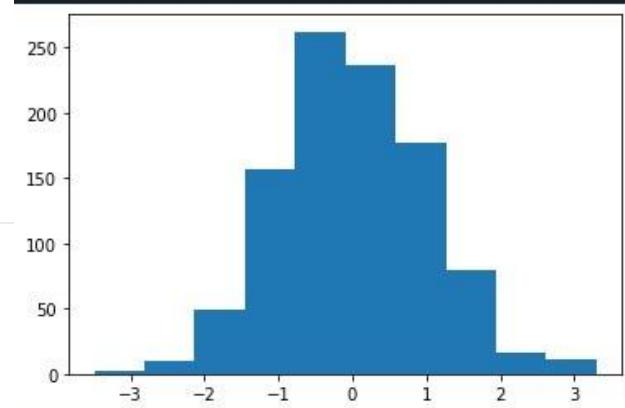
To perform Linear Regression (OLS Model)

CODE:

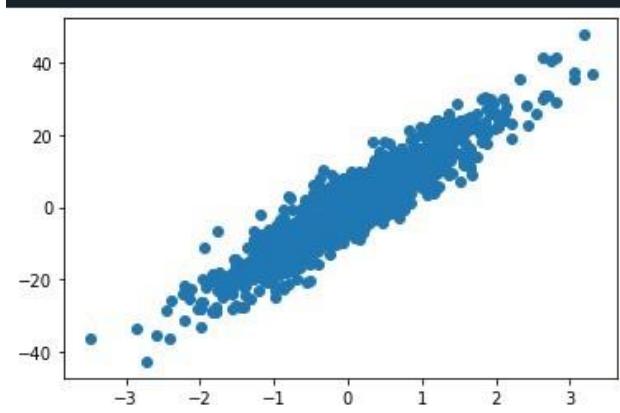
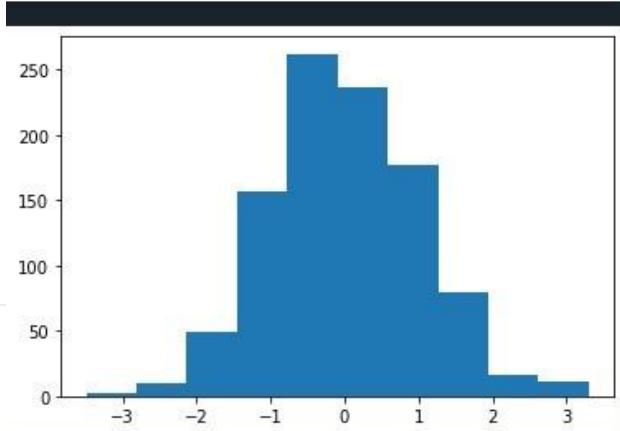
```
1  #%%
2  # import required libraries
3  from sklearn.datasets import make_regression
4  from matplotlib import pyplot as plt
5  from sklearn.model_selection import train_test_split
6  from sklearn.linear_model import LinearRegression
7  from sklearn.metrics import mean_squared_error ,mean_absolute_error
8  import numpy as np
9  import seaborn as sns
10 from math import sqrt
11 import warnings
12 warnings.filterwarnings('ignore')
13 #%%
14 #generate dataset for regression
15 X,y = make_regression(n_samples=1000, n_features = 1, shuffle = True, noise = 5,random_state=86)
16 #%%
17 #Statistics
18 stdX = np.std(X)
19 meanX = np.mean(X)
20
21 stdY = np.std(y)
22 meanY = np.mean(y)
23 print(meanX)
24 print(stdX)
25 print(meanY)
26 print(stdY)
27 #%%
28 # plot generated data
29 plt.hist(X)
30 plt.show()
31 #%%
32 plt.scatter(X,y)
33 plt.show()
34 #%%
35 # As mean of this dataset is already near to 0, so standard scaling is not required
36 #%%
37 # Preprocessing
38 #%%
39 plt.hist(X)
40 plt.show()
41 #%%
42 #Plot for data after preprocessing
43 plt.scatter(X,y)
44 plt.show()
45 #%%
46 # Train Test Split
47 X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.10, random_state=42)
48 #%%
49 sns.distplot(y_train)
50 plt.show()
51 #%%
52 sns.distplot(y_test)
53 plt.show()
54 #%%
55 #Performing Linear Regression || Using Traditional Method
```

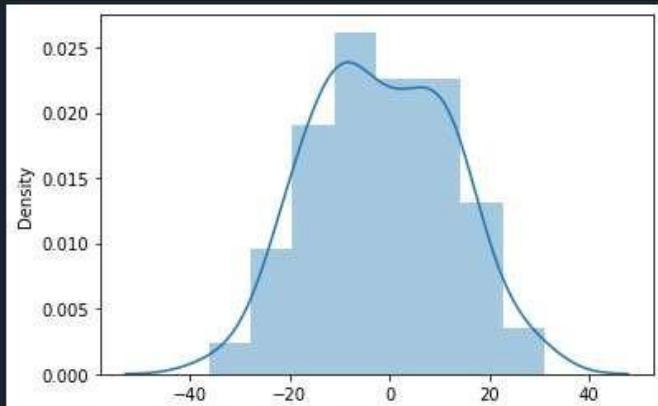
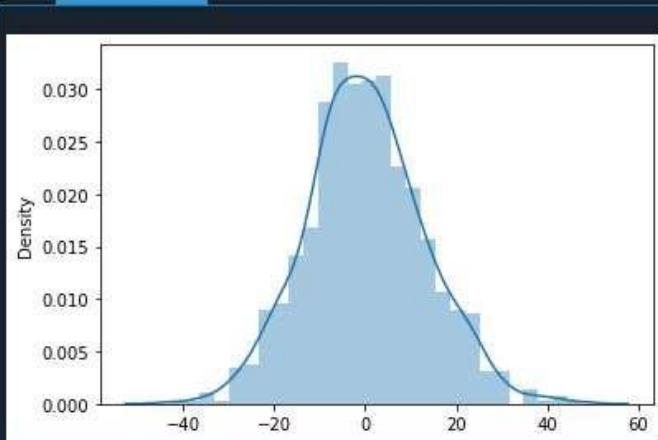
```
56 #Initializing variables
57 n = len(X_train)
58 num = 0
59 denom = 0
60 %%%
61 #Calculating slope & intercept
62 for i in range(n):
63     num += (X_train[i]-meanX) * (y_train[i]-meanY)
64     denom += (X_train[i] - meanX)**2
65 m = num/denom
66 c = meanY - (m*meanX)
67 print(m, ',', c)
68
69 min_X = 1 - np.min(X_train)
70 max_X = 1 - np.max(X_train)
71 print(min_X, ',', max_X)
72 %%%
73 #calculating the values of x & y
74 x = np.linspace(min_X, max_X,1000)
75 y = (m*x)+c
76 %%% Plot Regression Line
77 plt.scatter(X_train,y_train,color='b')
78 plt.plot(x,y,color='r')
79 plt.title('Simple Linear Regression')
80 plt.xlabel('X')
81 plt.ylabel('Y')
82 %%%
83 #calculating error
84 y_pred = []
85 sum_pred = 0
86 sum_actual = 0
87 for i in range(len(X_test)):
88     val = (m*X_test[i]+c)
89     y_pred.append(val)
90 mse = mean_squared_error(y_test, y_pred)
91 print("MSE", mse)
92 r2 = sqrt(mse)
93 print("R2",r2)
94 mae = mean_absolute_error(y_test, y_pred)
95 print("MAE", mae)
96 %%%
97 # Perform Linear Regression || Using sklearn
98 # Build Model
99 LRmodel = LinearRegression()
100 LRmodel.fit(X_train, y_train)
101 %%%
102 # Error Calculation
103 print("-----")
104 y_pred1 = LRmodel.predict(X_test)
105 mse1 = mean_squared_error(y_test, y_pred1)
106 print("MSE with sklearn", mse1)
107 r2new = sqrt(mse1)
108 print("R2 with sklearn",r2new)
109 mael = mean_absolute_error(y_test, y_pred1)
110 print("MAE with sklearn", mael)
```

```
In [1]: runfile('B:/3rd year/5th sem/P&AD/LinearRegression.py', wdir='B:/3rd year/5th sem/P&AD')
0.007610682926160853
0.9942474804457474
-0.05835119770320271
13.08071673806697
```



Output:





[12.21409321] , [-0.15130879]

3.8487665053713127 , -2.303677618501437

MSE 20.304967238809418

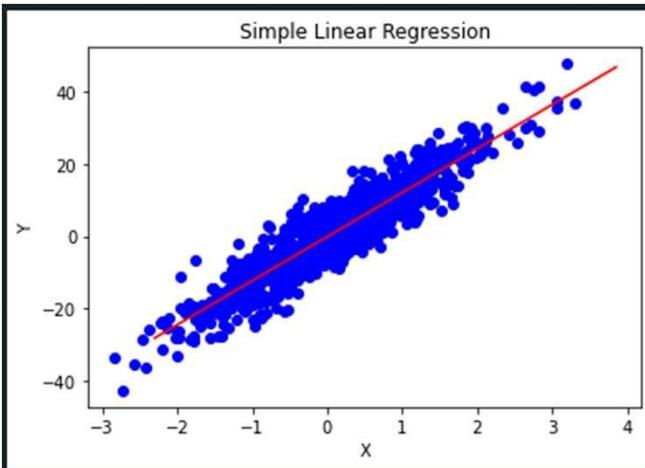
R2 4.506103332016413

MAE 3.611272677585965

MSE with sklearn 20.33856770371379

R2 with sklearn 4.509830119163447

MAE with sklearn 3.610421754936301



In [2]:

Experiment 4

AIM:



To perform Logistic Regression

CODE:

```
1  """
2  @author: Dhruv Singhal
3  """
4
5 #%%
6 from sklearn.datasets import make_classification
7 from sklearn.model_selection import train_test_split
8 from sklearn.linear_model import LogisticRegression
9 import seaborn as sns
10 import pandas as pd
11 from sklearn.metrics import confusion_matrix
12 import matplotlib.pyplot as plt
13 import warnings
14 warnings.filterwarnings('ignore')
15 #%%
16
17 #generate dataset
18 X,y = make_classification(n_samples=1000,n_classes=2)
19 X=pd.DataFrame(X)
20 y=pd.Series(y,name=20)
21 print("X values are:",X.head())
22 print("Y values are:",y.head())
23
24 #%% Visualization
25
26 plt.hist(X)
27 plt.show()
28
29 #%%
30 sns.distplot(y)
31 plt.show()
32
33 #%%
34 sns.distplot(X)
35 plt.show()
36
37 #%% feature extraction preprocessing, correlation matrix
38 _,graph=plt.subplots(figsize=(15,10))
39 sns.heatmap(X.corr(),annot=True,ax=graph,square=True)
40 plt.show()
41
42 #%%
43 df=pd.merge(X,y,right_index=True,left_index=True)
44 print(df.head())
45 print(df.corr()[[20]].abs().sort_values(by=20,ascending=False))
```

```
46 #%% columns with high correlation will be dropped
47 # creating useful data
48 datasetX=df[[1,2,3,4,5,6,7,8,9,10,11,12,13,14,15,16,17,18,19]]
49 datasetY=df[20]
50 print(datasetY.head())
51
52 #%%
53 X_train,X_test,Y_train,Y_test=train_test_split(X,y,test_size=0.15,random_state=42)
54
55 model=LogisticRegression()
56 model.fit(X_train,Y_train)
57 print(model.classes_)
58
59 Y_pred=model.predict(X_test)
60 print(Y_pred)
61
62 print("train Accuracy:",model.score(X_train,Y_train))
63 print("test Accuracy:",model.score(X_test,Y_test))
64
65 print(confusion_matrix(Y_test,Y_pred))
66
67
68
69
70
71
```

OUTPUT:

```
Python 3.7.6 (default, Jan  8 2020, 20:23:39) [MSC v.1916 64 bit (AMD64)]
Type "copyright", "credits" or "license" for more information.

IPython 7.26.0 -- An enhanced Interactive Python.

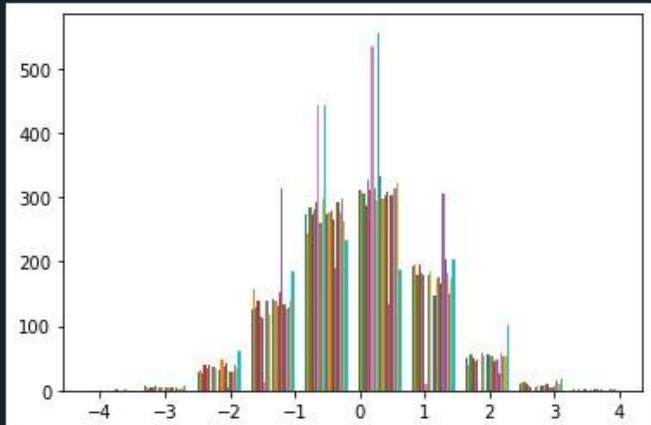
In [1]: runcell(0, 'C:/Users/Dhruv Singhal/untitled0.py')

In [2]: runcell(1, 'C:/Users/Dhruv Singhal/untitled0.py')

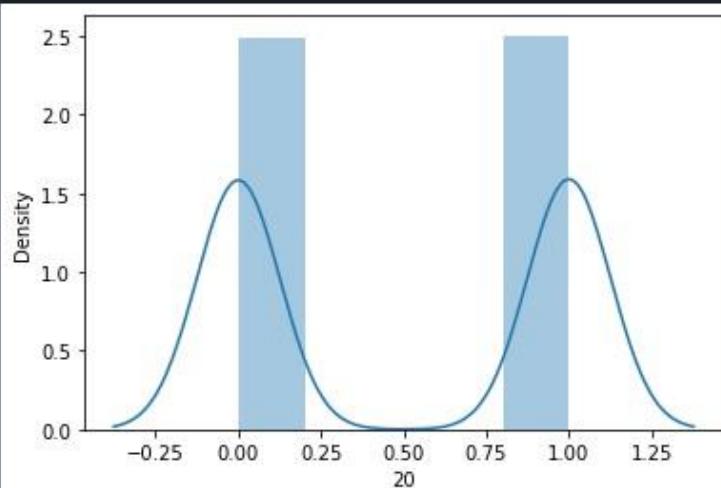
In [3]: runcell(2, 'C:/Users/Dhruv Singhal/untitled0.py')
X values are: 0      1      2     ...    17     18     19
0 -1.259436 -0.372492 0.054760 ... -0.035846 -1.462309 -1.140185
1 -0.800033 0.478840 -0.629627 ... 0.170002 0.445090 -1.299911
2 -0.770212 -0.741631 1.075089 ... 0.075282 -1.240537 -0.914384
3 1.583047 -0.705226 0.710976 ... -1.260828 0.947186 -0.973737
4 0.431166 0.902277 2.205464 ... 0.333836 0.417141 1.348560

[5 rows x 20 columns]
Y values are: 0      0
1      0
2      1
3      0
4      1
Name: 20, dtype: int32

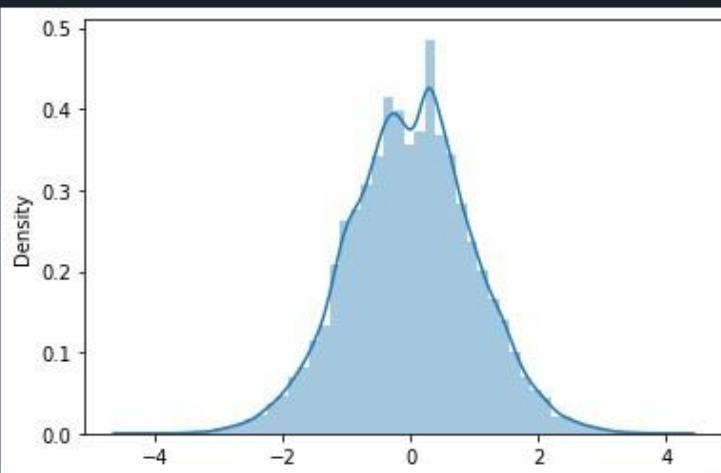
In [4]: runcell('Visualization', 'C:/Users/Dhruv Singhal/untitled0.py')
```



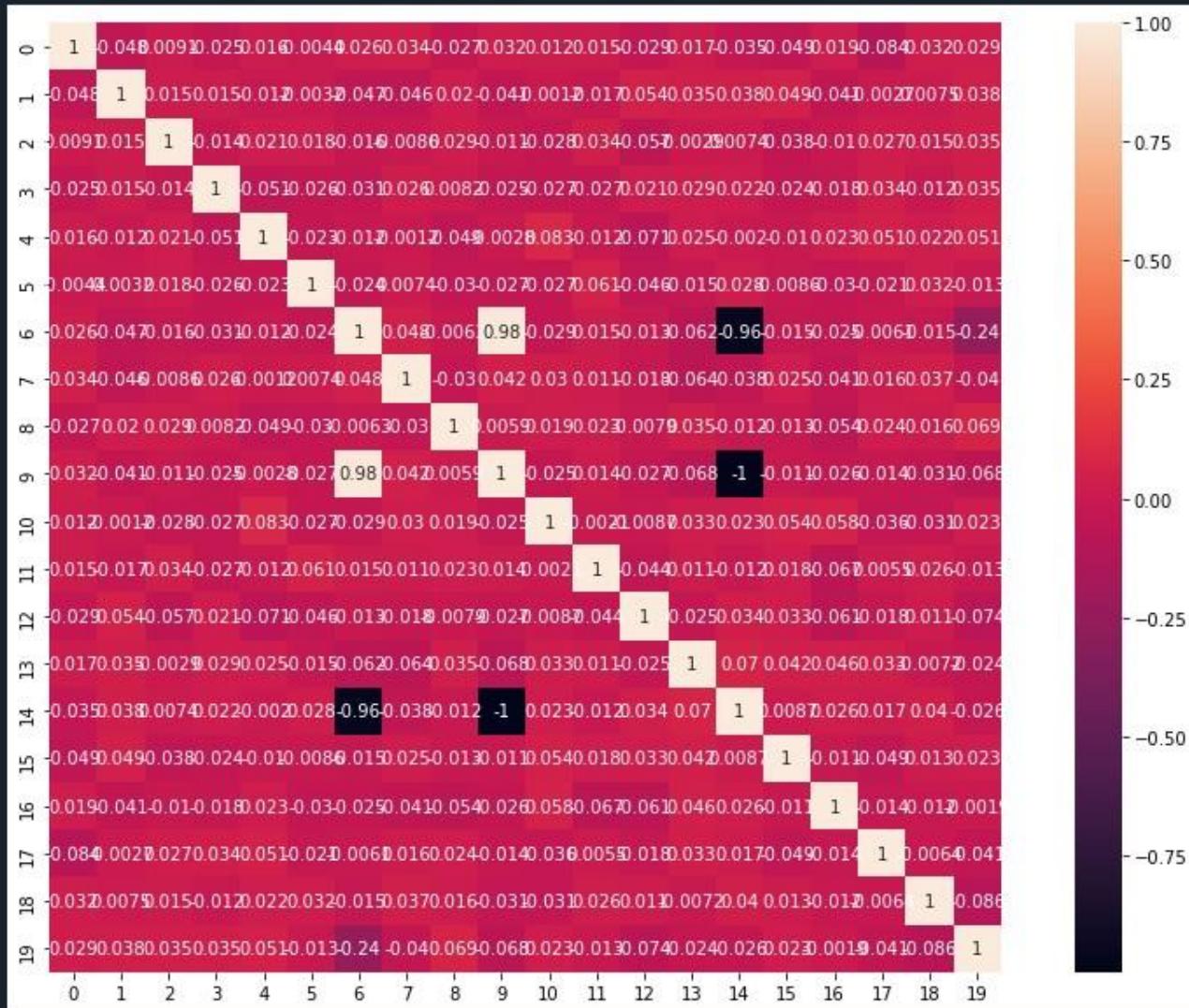
```
In [5]: runcell(4, 'C:/Users/Dhruv Singhal/untitled0.py')
```



```
In [6]: runcell(5, 'C:/Users/Dhruv Singhal/untitled0.py')
```



```
In [7]: runcell('feature extractionpreprocessing, correlation matrix', 'C:/Users/Dhruv Singh/untitled0.py')
```



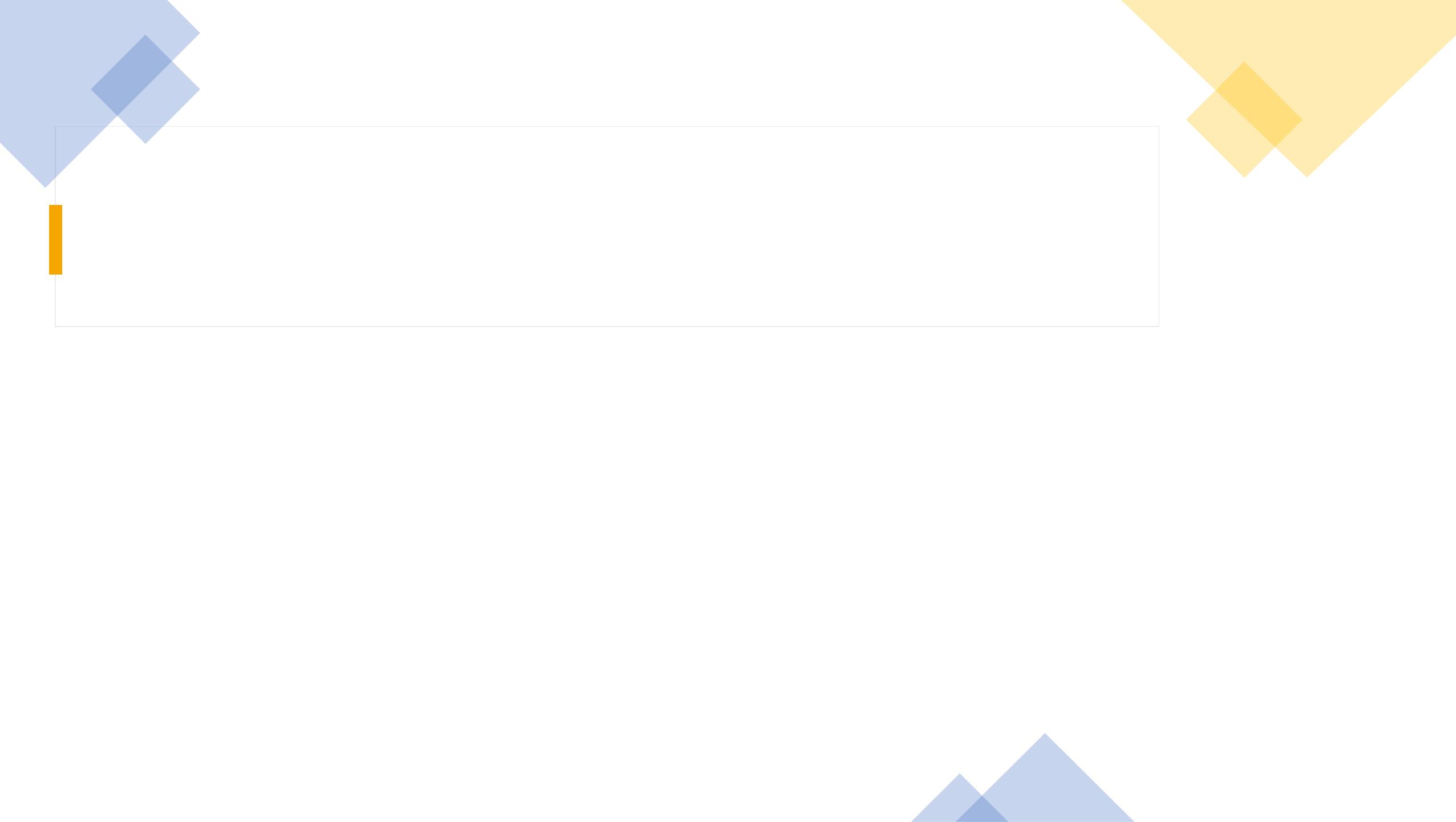
```
In [8]: runcell(7, 'C:/Users/Dhruv Singhal/untitled0.py')
      0         1         2         3       ...        17        18        19        20
0 -1.259436 -0.372492  0.054760 -2.045250   ...  -0.035846 -1.462309 -1.140185  0
1 -0.800033  0.478840 -0.629627 -0.773224   ...  0.170002  0.445090 -1.299911  0
2 -0.770212 -0.741631  1.075089 -1.506821   ...  0.075282 -1.240537 -0.914384  1
3  1.583047 -0.705226  0.710976 -1.241509   ...  -1.260828  0.947186 -0.973737  0
4  0.431166  0.902277  2.205464 -0.400889   ...  0.333836  0.417141  1.348560  1
```

[5 rows x 21 columns]

```
      20
20  1.000000
14  0.907235
9   0.906191
6   0.882873
13  0.067238
12  0.050728
1   0.046343
15  0.036825
3   0.032795
7   0.029817
0   0.029335
11  0.028842
16  0.026296
18  0.018547
5   0.018465
4   0.015409
10  0.011313
19  0.007832
17  0.005624
8   0.001292
2   0.000811
```

```
In [9]: runcell('columns with high correlation will be dropped', 'C:/Users/Dhruv Singhal/untitled0.py')
```

```
0   0
1   0
2   1
3   0
4   1
Name: 20, dtype: int32
```

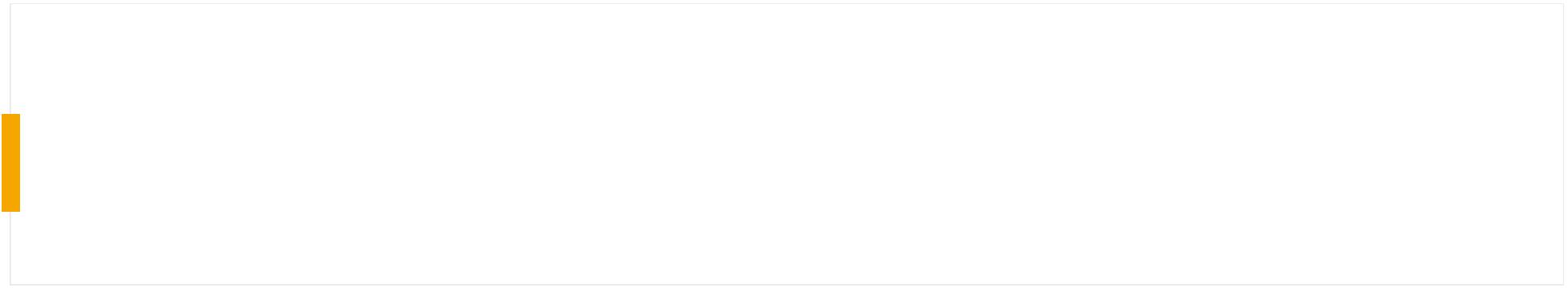


```
In [10]: runcell(9, 'C:/Users/Dhruv Singhal/untitled0.py')
[0 1]
[1 0 1 0 0 1 0 0 1 0 0 1 0 1 0 0 0 1 0 1 1 1 0 0 0 0 1 1 1 1 0
 0 1 1 0 1 1 0 1 0 0 1 1 1 1 0 0 0 0 0 1 1 1 1 1 0 1 0 0 0 1 0 0 0
 0 0 1 0 1 0 0 0 0 1 1 0 0 0 1 0 1 1 1 0 0 0 0 1 1 1 1 0 0 0 1 0 1 0
 0 0 0 0 1 1 0 0 1 0 0 0 1 0 1 1 1 0 1 0 1 0 1 0 1 0 0 0 0 1 0 1 0 0 0 1 1
 0 0]
train Accuracy: 0.9776470588235294
test Accuracy: 0.9733333333333334
[[82  0]
 [ 4 64]]
```

```
In [11]:
```

Experiment 5

```
1 import numpy as np
2 import matplotlib.pyplot as plt
3 from sklearn.model_selection import train_test_split
4 from sklearn.linear_model import LinearRegression
5 from sklearn.preprocessing import PolynomialFeatures
6 from sklearn.datasets import make_regression
7
8 #%%
9 # generate dataset
10 X,y = make_regression(n_samples=100, n_features=5,noise=50)
11
12 #%%
13 #plot the graph
14 plt.scatter(X[:,2],y)
15
16 #%%
17 #Train test Split
18 X_train, X_test, y_train, y_test = train_test_split(X,y,test_size=0.3) #30%test and 70% train dataset.
19
20 #%%
21 #implement Polynomial Regression
22 poly = PolynomialFeatures(degree = 3) # degree of polynomial=3
23 X_poly = poly.fit_transform(X_train)
```



```
57 #%%  
58 #implement through ridgeCV  
59 from sklearn.linear_model import RidgeCV  
60 model01=RidgeCV(alphas=[1e-3,1e-2,1e-1,1])  
61 model01.fit(X_poly,y_train)  
62 #%%  
63 #prediction  
64 y_pred_rid=model01.predict(X_poly_test)  
65 #%%  
66 #calculate rmse ,log error and R^2 score  
67 rms_rid=mean_squared_error(y_test, y_pred_rid)  
68 print("RMSE Error: ",rms_rid)  
69 log_rid=mean_squared_log_error(abs(y_test), abs(y_pred_rid))  
70 print("Log Error: ",log_rid)  
71 r2_sc_ridge=r2_score(y_test,y_pred_rid)  
72 print("R^2 Score: ",r2_sc_ridge)  
73 #%%  
74 #implement cross validation model  
75 from sklearn.model_selection import KFold,cross_val_score  
76 folds = KFold(n_splits = 5, shuffle = True, random_state = 100)  
77 scores = cross_val_score(model, X_train, y_train, scoring='neg_mean_absolute_error', cv=folds)  
78 print(np.average(scores) )  
79  
80  
81  
82  
83
```

Experiment 7

```
1  #-*- coding: utf-8 -*-
2  """
3  Created on Mon Sep 20 15:33:41 2021
4
5  @author: Dhruv Singhal
6  """
7
8  #%%
9  import numpy as np
10 import matplotlib.pyplot as plt
11 #import pandas as pd
12 import seaborn as sns
13 import warnings
14 warnings.filterwarnings('ignore')
15 #%%
16 from sklearn.datasets import make_regression
17 x,y=make_regression(n_samples=10000,n_features=5,noise=30)
18 sns.distplot(y)
19
20 #%%
21 plt.hist(x)
22 #%%
23 from sklearn.model_selection import train_test_split
24 x_train,x_test,y_train,y_test=train_test_split(x,y,train_size=0.15,random_state=42)
25 from sklearn.metrics import r2_score,mean_squared_error
26 #%%
27 from sklearn.linear_model import LinearRegression
28 lin=LinearRegression()
29 lin.fit(x_train,y_train)
30 y_pred=lin.predict(x_test)
31 print("r2 score without tuning",r2_score(y_test,y_pred))
32 print("RMSE without tuning",np.sqrt(mean_squared_error(y_test,y_pred)))
33
34 #%% Hyperparameter Tuning From Here
35 #%%
36 from sklearn.model_selection import GridSearchCV
37 tuned_parameters = [{ 'fit_intercept': [ 'True' ], 'normalize': [ 'True' ] },
38                      { 'fit_intercept': [ 'False' ], 'normalize': [ 'True' ] },
39                      { 'fit_intercept': [ 'True' ], 'normalize': [ 'False' ] },
40                      { 'fit_intercept': [ 'False' ], 'normalize': [ 'False' ] }
41
42                      ]
43 clf=GridSearchCV(LinearRegression(),tuned_parameters,scoring='r2')
44 clf.fit(x_train,y_train)
```

```
45 print("Best parameters set found on development set:")
46 print()
47 print(clf.best_params_)
48 print()
49 print("Best Score:",clf.best_score_)
50 z=clf.cv_results_
51
52 #%% kfold cross validation with hyperparameter tuning
53 from sklearn.model_selection import KFold
54 k = 5
55 kf = KFold(n_splits=k, random_state=None)
56 model = GridSearchCV(LinearRegression(),tuned_parameters,scoring=('r2'))
57
58 for train_index , test_index in kf.split(x):
59     X_train , X_test = x[train_index,:],x[test_index,:]
60     y_train , y_test = y[train_index] , y[test_index]
61
62     model.fit(X_train,y_train)
63     pred_values = model.predict(X_test)
64 print("Best parameters set found on development set:")
65 print()
66 print(model.best_params_)
67 print()
68 print("Best Score:",model.best_score_)
69 z2=model.cv_results_
70
71
```

Experiment 6

CODE:

```
1 # -*- coding: utf-8 -*-
2 """
3
4 @author: Dhruv Singhal
5 """
6
7 """
8 import numpy as np
9 import pandas as pd
10 import matplotlib.pyplot as plt
11 import seaborn as sns
12 from sklearn.datasets import make_regression
13 """
14 X,y=make_regression(n_samples=100,n_features=10,n_informative=4,random_state=0)
15 sns.distplot(y)
16
17 """
18 from sklearn.decomposition import PCA
19 pca=PCA(n_components=5)
20 principalComponents=pca.fit_transform(X)
21 X_new=pd.DataFrame(data=principalComponents,columns=['PC1','PC2','PC3','PC4','PC5'])
22 print(X_new.head())
23 """
24 from sklearn.model_selection import train_test_split
25 from sklearn.linear_model import LinearRegression
26 X_train , X_test , Y_train , Y_test = train_test_split(X_new, y, train_size=0.20,random_state=0)
27 lin=LinearRegression()
28 model=lin.fit(X_train,Y_train)
29 y_pred=model.predict(X_test)
30 """
31 print('Coefficients: ', model.coef_)
32 print('Variance score: {}'.format(model.score(X_test, Y_test)))
33 print("Test Score",model.score(X_train,Y_train))
34 print("Test Score",model.score(X_test,Y_test))
35 """
36 #training
37 plt.scatter(model.predict(X_train),model.predict(X_train) - Y_train,color = "green", s = 10, label = 'Train data')
38 #testing
39 plt.scatter(model.predict(X_test),model.predict(X_test) - Y_test,color = "orange", s = 10, label = 'Testing data')
40
41 """
42 import math
43 from sklearn.metrics import mean_squared_error
44 print("RMSE",math.sqrt(mean_squared_error(Y_test,y_pred)))
45 from sklearn.metrics import r2_score
46 print("R^2",r2_score(Y_test,y_pred))
47
48
49 """
50 print(np.max(Y_train),np.min(Y_train))
51
52 """
53 plt.figure(figsize=(20,10))
54 plt.plot(y_pred)
55 plt.plot(Y_test)
```

```
57
58
59
60
61    #%%%
62    X_train2 , X_test2 , Y_train2 , Y_test2 = train_test_split(X, y, train_size=0.20,random_state=0)
63    lin=LinearRegression()
64    model2=lin.fit(X_train2,Y_train2)
65    y_pred2=model2.predict(X_test2)
66    #%%
67    print('Coefficients: ', model2.coef_)
68    print('Variance score: {}'.format(model2.score(X_test2, Y_test2)))
69    print("Test Score",model2.score(X_train2,Y_train2))
70    print("Test Score",model2.score(X_test2,Y_test2))
71
72    #%%
73    import math
74    from sklearn.metrics import mean_squared_error
75    print("RMSE",math.sqrt(mean_squared_error(Y_test2,y_pred2)))
76    from sklearn.metrics import r2_score
77    print("R^2",r2_score(Y_test2,y_pred2))
78
79
80    #%%
81    print(np.max(Y_train2),np.min(Y_train2))
82
83    #%%
84    plt.figure(figsize=(20,10))
85    plt.plot(y_pred2)
86    plt.plot(Y_test2)
87
88
89
90 |
```

Console 1/A

Python 3.7.6 (default, Jan 8 2020, 20:23:39) [MSC v.1916 64 bit (AMD64)]
Type "copyright", "credits" or "license" for more information.

OUTPUT:

IPython 7.26.0 -- An enhanced Interactive Python.

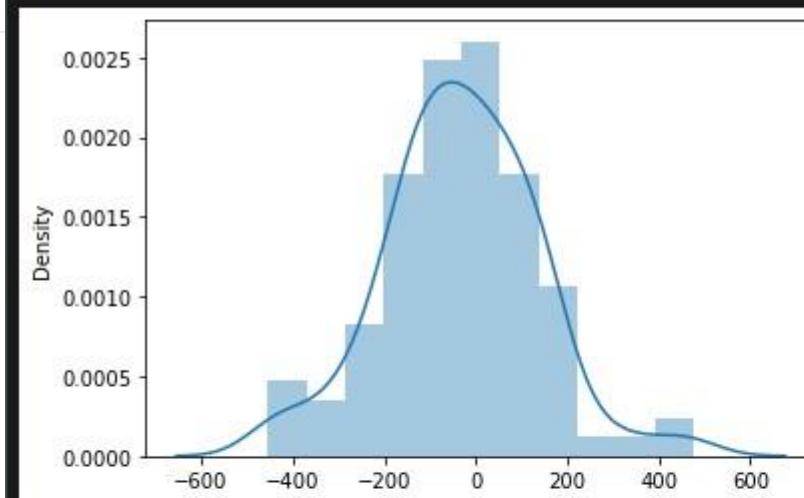
In [1]: runcell(0, 'B:/3rd year/5th sem/P&AD/PCA.py')

In [2]: runcell(1, 'B:/3rd year/5th sem/P&AD/PCA.py')

In [3]: runcell(2, 'B:/3rd year/5th sem/P&AD/PCA.py')

C:\Users\Dhruv Singhal\anaconda3\lib\site-packages\seaborn\distributions.py:2619: FutureWarning: `distplot` 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)



In [4]: runcell(3, 'B:/3rd year/5th sem/P&AD/PCA.py')

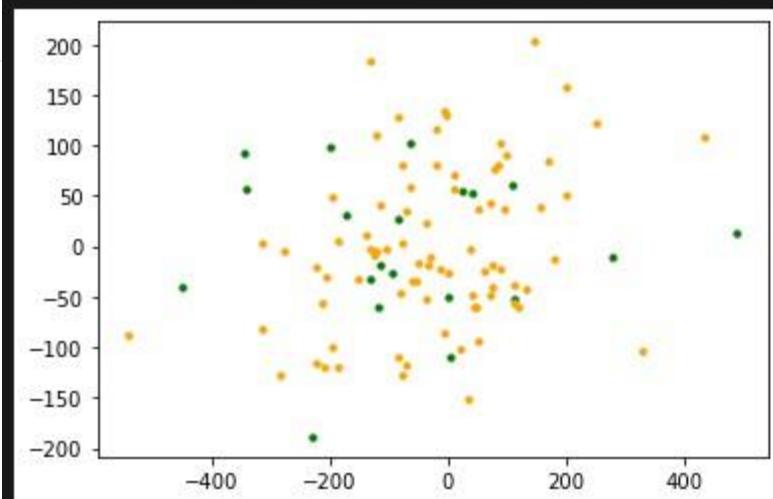
	PC1	PC2	PC3	PC4	PC5
0	2.689187	0.188731	0.553922	0.677884	0.069860
1	-2.718170	-2.513221	0.564428	0.308534	0.148173
2	-0.150486	3.214459	-0.133615	-1.756744	0.593628
3	-1.597885	-2.291277	-0.211204	-0.593296	-0.201198
4	-0.621526	-0.492717	0.588593	0.407506	-0.956159

In [5]: runcell(4, 'B:/3rd year/5th sem/P&AD/PCA.py')

```
In [6]: runcell(5, 'B:/3rd year/5th sem/P&AD/PCA.py')
```

```
Output from spyder call 'get_namespace_view':  
Coefficients: [-14.02798066 29.78262828 -44.53078748 120.20573166 -94.01743173]  
Variance score: 0.715357937146422  
Test Score 0.8948220992728023  
Test Score 0.715357937146422
```

```
In [7]: runcell(6, 'B:/3rd year/5th sem/P&AD/PCA.py')
```



```
In [8]: runcell(7, 'B:/3rd year/5th sem/P&AD/PCA.py')
```

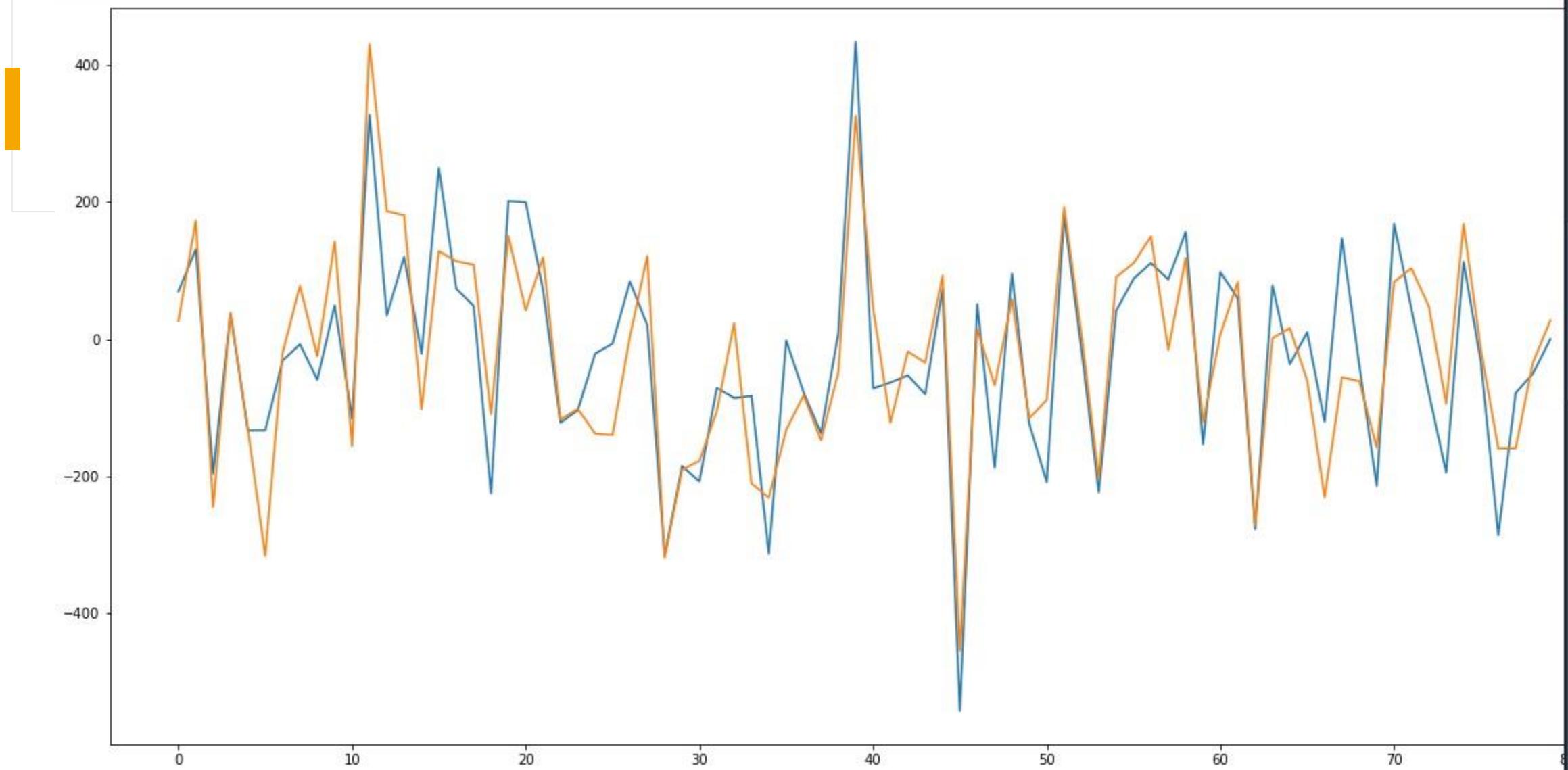
```
RMSE 79.41012724129493  
R^2 0.715357937146422
```

```
In [9]: runcell(8, 'B:/3rd year/5th sem/P&AD/PCA.py')
```

```
476.0496266907876 -436.9209738094973
```

```
In [10]: runcell(9, 'B:/3rd year/5th sem/P&AD/PCA.py')
```

```
In [10]: runcell(9, 'B:/3rd year/5th sem/P&AD/PCA.py')
```



```
In [11]: runcell(10, "B:/3rd year/5th sem/P&AD/PCA.py")
```

```
In [12]: runcell(11, "B:/3rd year/5th sem/P&AD/PCA.py")
```

```
Coefficients: [-1.05459713e-14  9.83472758e-14 -7.33134857e-14  7.00528623e+01  
 8.83077597e+01  9.66575107e+01  8.21903908e+01  3.77754737e-14  
 4.10586745e-14 -4.61401623e-14]
```

```
Variance score: 1.0
```

```
Test Score 1.0
```

```
Test Score 1.0
```

```
In [13]: runcell(12, "B:/3rd year/5th sem/P&AD/PCA.py")
```

```
RMSE 79.41012724129493
```

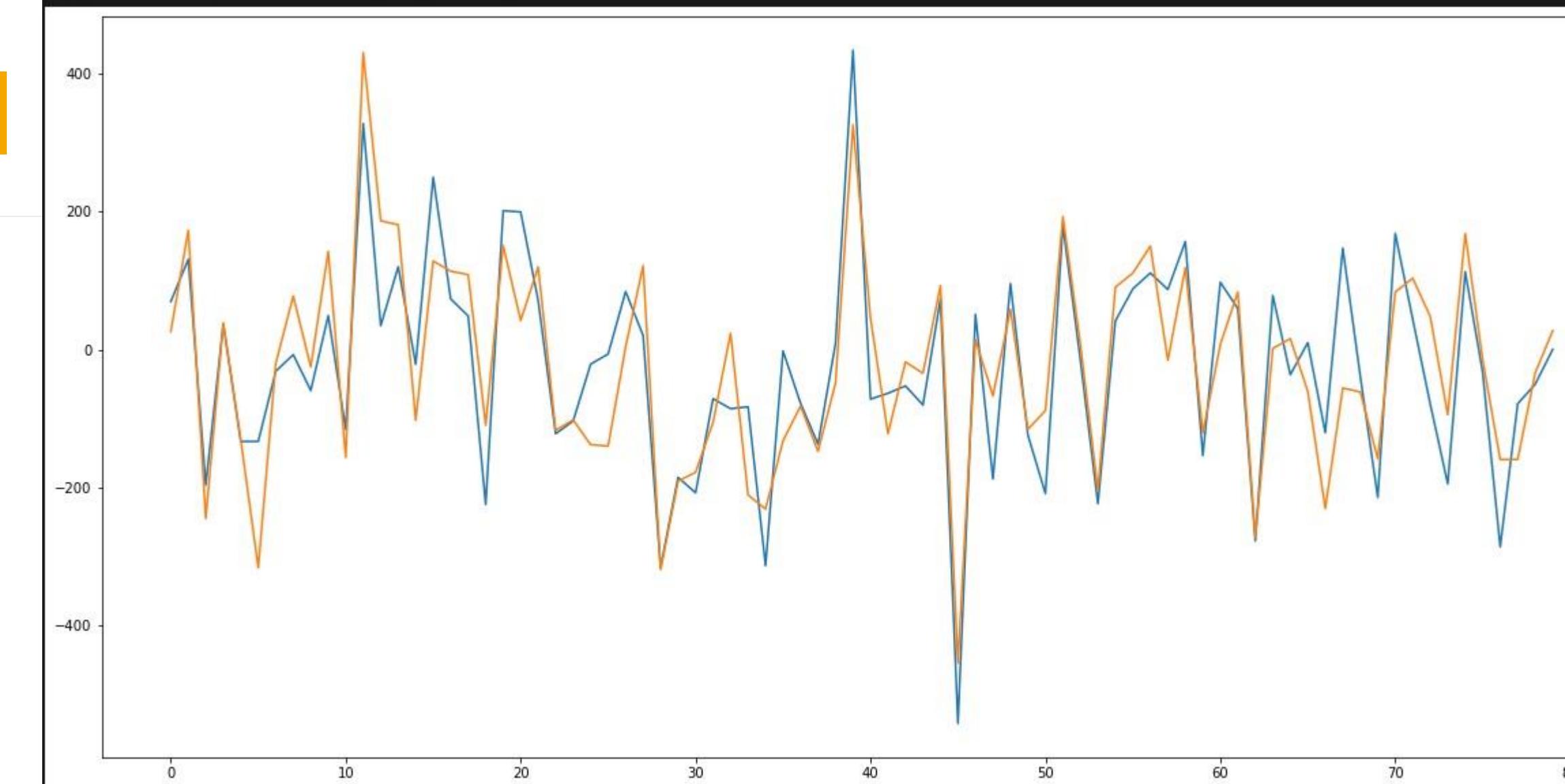
```
R^2 0.715357937146422
```

```
In [14]: runcell(13, "B:/3rd year/5th sem/P&AD/PCA.py")
```

```
476.0496266907876 -436.9209738094973
```

```
In [15]: runcell(14, "B:/3rd year/5th sem/P&AD/PCA.py")
```

```
In [15]: runcell(14, 'B:/3rd year/5th sem/P&AD/PCA.py')
```



Experiment 8

Hyperparameter

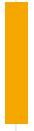
Tuning and Nested Cross validation on



Linear Regression

```
1 import numpy as np
2 import matplotlib.pyplot as plt
3 import seaborn as sns
4 import warnings
5 warnings.filterwarnings('ignore')
6 #%%%
7 from sklearn.datasets import make_regression
8 x,y=make_regression(n_samples=1000,n_features=5,noise=20)
9 sns.distplot(y)
10
11 #%%
12 plt.hist(x)
13 #%%
14 x_train = x
15 y_train = y
16
17 #%%
18 from sklearn.model_selection import GridSearchCV, cross_val_score
19 from sklearn.model_selection import KFold
20 from sklearn.linear_model import Ridge
21 NUM_TRIALS = 30
22
23
24 tuned_parameters = [{"solver": ['svd', 'lsqr'], "fit_intercept": ['True'], "normalize": ['False']},
25                      {"solver": ['sag', 'cholesky'], "fit_intercept": ['False'], "normalize": ['true']}]
26
27 score = 'r2'
28 non_nested_scores = np.zeros(NUM_TRIALS)
29 nested_scores = np.zeros(NUM_TRIALS)
30
31 # Loop for each trial
32 for i in range(NUM_TRIALS):
33
34     # model= GridSearchCV(linear_model.LinearRegression(), tuned_parameters, scoring= score)
35     inner_cv = KFold(n_splits=4, shuffle=True, random_state=i)
36     outer_cv = KFold(n_splits=4, shuffle=True, random_state=i)
37     model= GridSearchCV(estimator = Ridge(), param_grid = tuned_parameters, scoring = score)
38     model.fit(x_train, y_train)
39     non_nested_scores[i] = model.best_score_
40
41
42     # Nested CV with parameter optimization
43     model = GridSearchCV(estimator= Ridge(), param_grid = tuned_parameters, cv=inner_cv, scoring= score)
44     nested_score = cross_val_score(model, X=x_train, y=y_train, cv=outer_cv)
45     nested_scores[i] = nested_score.mean()
46
47
48 score_difference = non_nested_scores - nested_scores
49
50 print("Average difference of {:.6f} with std. dev. of {:.6f}."
51       .format(score_difference.mean(), score_difference.std()))
52
53
54
```

CODE:

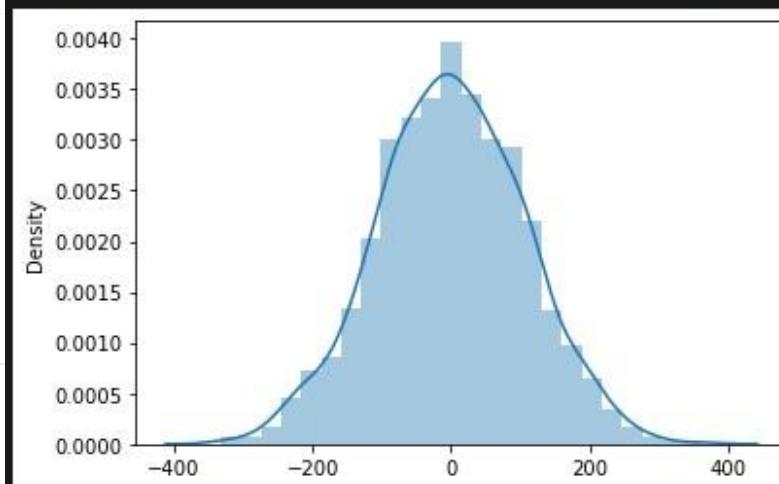


OUTPUT:

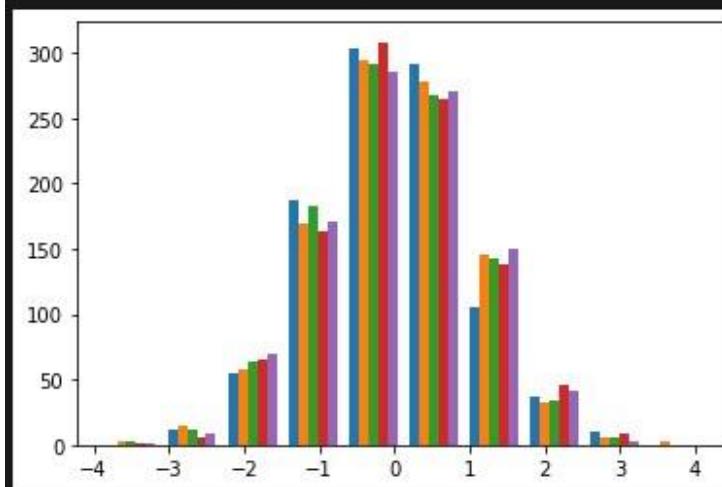
Console 3/A

In [1]: runcell(0, 'B:/3rd year/5th sem/P&AD/exp9.py')

In [2]: runcell(1, 'B:/3rd year/5th sem/P&AD/exp9.py')



In [3]: runcell(2, 'B:/3rd year/5th sem/P&AD/exp9.py')



In [4]: runcell(3, 'B:/3rd year/5th sem/P&AD/exp9.py')

In [5]: runcell(4, 'B:/3rd year/5th sem/P&AD/exp9.py')

Average difference of -0.003284 with std. dev. of 0.001567.

Hyperparameter

Tuning and Nested Cross validation on



Logistic Regression

```
1 import numpy as np
2 from sklearn.datasets import make_classification
3 from sklearn import linear_model
4 #from matplotlib import pyplot as plt
5 from sklearn.model_selection import GridSearchCV
6 from sklearn.model_selection import KFold
7 from sklearn.model_selection import cross_val_score
8 import warnings
9 warnings.filterwarnings('ignore')
10 #%%%
11 # Generating Data
12 X, y = make_classification(n_samples = 1000, n_features = 5, n_classes = 2)
13 x_train = X
14 y_train = y
15 #%%%
16 NUM_TRIALS = 30
17
18 tuned_parameters = [{ 'penalty' : ['l1', 'l2', 'elasticnet', 'none'],
19   'C' : np.logspace(-4, 4, 20),
20   'solver' : ['lbfgs','newton-cg','liblinear','sag','saga']
21   }]
22
23 score = 'accuracy'
24 non_nested_scores = np.zeros(NUM_TRIALS)
25 nested_scores = np.zeros(NUM_TRIALS)
26 # Loop for each trial
27 for i in range(NUM_TRIALS):
28
29     # model= GridSearchCV(linear_model.LogisticRegression(), tuned_parameters, scoring= score)
30     inner_cv = KFold(n_splits=4, shuffle=True, random_state=i)
31     outer_cv = KFold(n_splits=4, shuffle=True, random_state=i)
32     model= GridSearchCV(estimator = linear_model.LogisticRegression(), param_grid = tuned_parameters, scoring = score)
33     model.fit(x_train, y_train)
34     print(model.best_params_)
35     non_nested_scores[i] = model.best_score_
36
37
38     # Nested CV with parameter optimization
39     model = GridSearchCV(estimator= linear_model.LogisticRegression(), param_grid = tuned_parameters, cv=inner_cv, scoring= score)
40     nested_score = cross_val_score(model, X=x_train, y=y_train, cv=outer_cv)
41     nested_scores[i] = nested_score.mean()
42     score_difference = non_nested_scores - nested_scores
43     print("Average difference of {:.6f} with std. dev. of {:.6f}.".format(score_difference.mean(), score_difference.std()))
44
45 %%%
```

CODE:



OUTPUT:

Hyperparameter

Tuning and Nested Cross validation on



Logistic Regression

```
1 import numpy as np
2 #from sklearn.datasets import make_classification
3 from sklearn import linear_model
4 #from matplotlib import pyplot as plt
5 from sklearn.model_selection import GridSearchCV
6 from sklearn.model_selection import KFold
7 from sklearn.model_selection import cross_val_score
8 import warnings
9 warnings.filterwarnings('ignore')
10 #%%
11 # Generating Data
12 from sklearn.datasets import load_breast_cancer
13
14 data = load_breast_cancer()
15
16 X = data.data
17
18 y = data.target
19 x_train = X
20 y_train = y
21 #%%
22 NUM_TRIALS = 30
23
24
25 tuned_parameters = [{ 'penalty' : ['l1', 'l2', 'elasticnet', 'none'],
26 'solver' : ['lbfgs','newton-cg','liblinear','sag','saga']
27 }]
28 score = 'accuracy'
29 non_nested_scores = np.zeros(NUM_TRIALS)
30 nested_scores = np.zeros(NUM_TRIALS)
31 # Loop for each trial
32 for i in range(NUM_TRIALS):
33
34     # model= GridSearchCV(linear_model.LogisticRegression(), tuned_parameters, scoring= score)
35     inner_cv = KFold(n_splits=4, shuffle=True, random_state=i)
36     outer_cv = KFold(n_splits=4, shuffle=True, random_state=i)
37     model= GridSearchCV(estimator = linear_model.LogisticRegression(), param_grid = tuned_parameters, scoring = score)
38     model.fit(x_train, y_train)
39     print(model.best_params_)
40     non_nested_scores[i] = model.best_score_
41
42
43     # Nested CV with parameter optimization
44     model = GridSearchCV(estimator= linear_model.LogisticRegression(), param_grid = tuned_parameters, cv=inner_cv, scoring= score)
45     nested_score = cross_val_score(model, X=x_train, y=y_train, cv=outer_cv)
46     nested_scores[i] = nested_score.mean()
47 score_difference = non_nested_scores - nested_scores
48 print("Average difference of {:.6f} with std. dev. of {:.6f}.".format(score_difference.mean(), score_difference.std()))
49
50
```

(Breast Cancer Dataset)

CODE:



OUTPUT:

Experiment 9



CODE:

```
1 import matplotlib.pyplot as plt
2 import numpy as np
3
4
5 def plot_gpr_samples(gpr_model, n_samples, ax):
6     """Plot samples drawn from the Gaussian process model.
7
8     If the Gaussian process model is not trained then the drawn samples are
9     drawn from the prior distribution. Otherwise, the samples are drawn from
10    the posterior distribution. Be aware that a sample here corresponds to a
11    function. """
12
13    x = np.linspace(0, 5, 100)
14    X = x.reshape(-1, 1)
15
16    y_mean, y_std = gpr_model.predict(X, return_std=True)
17    y_samples = gpr_model.sample_y(X, n_samples)
18
19    for idx, single_prior in enumerate(y_samples.T):
20        ax.plot(
21            x,
22            single_prior,
23            linestyle="--",
24            alpha=0.7,
25            label=f"Sampled function #{idx + 1}",
26        )
27    ax.plot(x, y_mean, color="black", label="Mean")
28    ax.fill_between(
29        x,
30        y_mean - y_std,
31        y_mean + y_std,
32        alpha=0.1,
33        color="black",
34        label=r"$\pm$ 1 std. dev.",
35    )
36    ax.set_xlabel("x")
37    ax.set_ylabel("y")
38    ax.set_ylim([-3, 3])
39
40
41
42    rng = np.random.RandomState(4)
43    X_train = rng.uniform(0, 5, 10).reshape(-1, 1)
44    y_train = np.sin((X_train[:, 0] - 2.5) ** 2)
45    n_samples = 5
46
47
48    ### RBF Kernel
49
50    from sklearn.gaussian_process import GaussianProcessRegressor
51    from sklearn.gaussian_process.kernels import RBF
52
53    kernel = 1.0 * RBF(length_scale=1.0, length_scale_bounds=(1e-1, 10.0))
54    gpr = GaussianProcessRegressor(kernel=kernel, random_state=0)
55
56    fig, axs = plt.subplots(nrows=2, sharex=True, sharey=True, figsize=(10, 8))
```

```
57
58 # plot prior
59 plot_gpr_samples(gpr, n_samples=n_samples, ax=axs[0])
60 axs[0].set_title("Samples from prior distribution")
61
62 # plot posterior
63 gpr.fit(X_train, y_train)
64 plot_gpr_samples(gpr, n_samples=n_samples, ax=axs[1])
65 axs[1].scatter(X_train[:, 0], y_train, color="red", zorder=10, label="Observations")
66 axs[1].legend(bbox_to_anchor=(1.05, 1.5), loc="upper left")
67 axs[1].set_title("Samples from posterior distribution")
68
69 fig.suptitle("Radial Basis Function kernel", fontsize=18)
70 plt.tight_layout()
71
72
73 print(f"Kernel parameters before fit:\n{kernel}")
74 print(
75     f"Kernel parameters after fit: \n{gpr.kernel_}\n"
76     f"Log-likelihood: {gpr.log_marginal_likelihood(gpr.kernel_.theta):.3f}"
77 )
78
79
80
81 #### Rational Quadratic Kernel
82
83 from sklearn.gaussian_process.kernels import RationalQuadratic
84
85 kernel = 1.0 * RationalQuadratic(length_scale=1.0, alpha=0.1, alpha_bounds=(1e-5, 1e15))
86 gpr = GaussianProcessRegressor(kernel=kernel, random_state=0)
87
88 fig, axs = plt.subplots(nrows=2, sharex=True, sharey=True, figsize=(10, 8))
89
90 # plot prior
91 plot_gpr_samples(gpr, n_samples=n_samples, ax=axs[0])
92 axs[0].set_title("Samples from prior distribution")
93
94 # plot posterior
95 gpr.fit(X_train, y_train)
96 plot_gpr_samples(gpr, n_samples=n_samples, ax=axs[1])
97 axs[1].scatter(X_train[:, 0], y_train, color="red", zorder=10, label="Observations")
98 axs[1].legend(bbox_to_anchor=(1.05, 1.5), loc="upper left")
99 axs[1].set_title("Samples from posterior distribution")
100
101 fig.suptitle("Rational Quadratic kernel", fontsize=18)
102 plt.tight_layout()
103
104
105 print(f"Kernel parameters before fit:\n{kernel}")
106 print(
107     f"Kernel parameters after fit: \n{gpr.kernel_}\n"
108     f"Log-likelihood: {gpr.log_marginal_likelihood(gpr.kernel_.theta):.3f}"
109 )
```

```
110 #%%
111 ### ExpSineSquared
112
113
114 from sklearn.gaussian_process.kernels import ExpSineSquared
115
116
117 kernel = 1.0 * ExpSineSquared(length_scale=1,periodicity=1,length_scale_bounds=(1e-5, 1e15),periodicity_bounds=(1e-5, 1e15))
118 gpr = GaussianProcessRegressor(kernel=kernel, random_state=0)
119
120 fig, axs = plt.subplots(nrows=2, sharex=True, sharey=True, figsize=(10, 8))
121
122 # plot prior
123 plot_gpr_samples(gpr, n_samples=n_samples, ax=axs[0])
124 axs[0].set_title("Samples from prior distribution")
125
126 # plot posterior
127 gpr.fit(X_train, y_train)
128 plot_gpr_samples(gpr, n_samples=n_samples, ax=axs[1])
129 axs[1].scatter(X_train[:, 0], y_train, color="red", zorder=10, label="Observations")
130 axs[1].legend(bbox_to_anchor=(1.05, 1.5), loc="upper left")
131 axs[1].set_title("Samples from posterior distribution")
132
133 fig.suptitle("ExpSineSquared", fontsize=18)
134 plt.tight_layout()
135
136 print(f"Kernel parameters before fit:\n{kernel}")
137 print(
138     f"Kernel parameters after fit: \n{gpr.kernel_}\n"
139     f"Log-likelihood: {gpr.log_marginal_likelihood(gpr.kernel_.theta):.3f}"
140 )
141 #%%
142 ### Matern
143
144
145 from sklearn.gaussian_process.kernels import Matern
146
147
148 kernel = 1.0 * Matern(length_scale=1,length_scale_bounds=(1e-5, 1e15),nu=1.5)
149 gpr = GaussianProcessRegressor(kernel=kernel, random_state=0)
150
151 fig, axs = plt.subplots(nrows=2, sharex=True, sharey=True, figsize=(10, 8))
152
153 # plot prior
154 plot_gpr_samples(gpr, n_samples=n_samples, ax=axs[0])
155 axs[0].set_title("Samples from prior distribution")
156
157 # plot posterior
158 gpr.fit(X_train, y_train)
159 plot_gpr_samples(gpr, n_samples=n_samples, ax=axs[1])
160 axs[1].scatter(X_train[:, 0], y_train, color="red", zorder=10, label="Observations")
161 axs[1].legend(bbox_to_anchor=(1.05, 1.5), loc="upper left")
162 axs[1].set_title("Samples from posterior distribution")
163
164 fig.suptitle("Matern", fontsize=18)
165 plt.tight_layout()
```

```
166
167 print(f"Kernel parameters before fit:\n{kernel}")
168 print(
169     f"Kernel parameters after fit: \n{gpr.kernel_} \n"
170     f"Log-likelihood: {gpr.log_marginal_likelihood(gpr.kernel_.theta):.3f}"
171 )
172
173 #%% 
174 ### WhiteKernel
175
176
177 from sklearn.gaussian_process.kernels import WhiteKernel
178
179
180 kernel = 1.0 * WhiteKernel(noise_level=1, noise_level_bounds=(1e-5, 1e15))
181 gpr = GaussianProcessRegressor(kernel=kernel, random_state=0)
182
183 fig, axs = plt.subplots(nrows=2, sharex=True, sharey=True, figsize=(10, 8))
184
185 # plot prior
186 plot_gpr_samples(gpr, n_samples=n_samples, ax=axs[0])
187 axs[0].set_title("Samples from prior distribution")
188
189 # plot posterior
190 gpr.fit(X_train, y_train)
191 plot_gpr_samples(gpr, n_samples=n_samples, ax=axs[1])
192 axs[1].scatter(X_train[:, 0], y_train, color="red", zorder=10, label="Observations")
193 axs[1].legend(bbox_to_anchor=(1.05, 1.5), loc="upper left")
194 axs[1].set_title("Samples from posterior distribution")
195
196 fig.suptitle("WhiteKernel", fontsize=18)
197 plt.tight_layout()
198
199 print(f"Kernel parameters before fit:\n{kernel}")
200 print(
201     f"Kernel parameters after fit: \n{gpr.kernel_} \n"
202     f"Log-likelihood: {gpr.log_marginal_likelihood(gpr.kernel_.theta):.3f}"
203 )
204
205 #%%
206
207 ### Exponentiation
208
209
210 from sklearn.gaussian_process.kernels import Exponentiation
211
212
213 kernel = 1.0 * Exponentiation(RationalQuadratic(), exponent=2)
214 gpr = GaussianProcessRegressor(kernel=kernel, random_state=0)
215
216 fig, axs = plt.subplots(nrows=2, sharex=True, sharey=True, figsize=(10, 8))
217
218 # plot prior
219 plot_gpr_samples(gpr, n_samples=n_samples, ax=axs[0])
220 axs[0].set_title("Samples from prior distribution")
221
222 # plot posterior
```

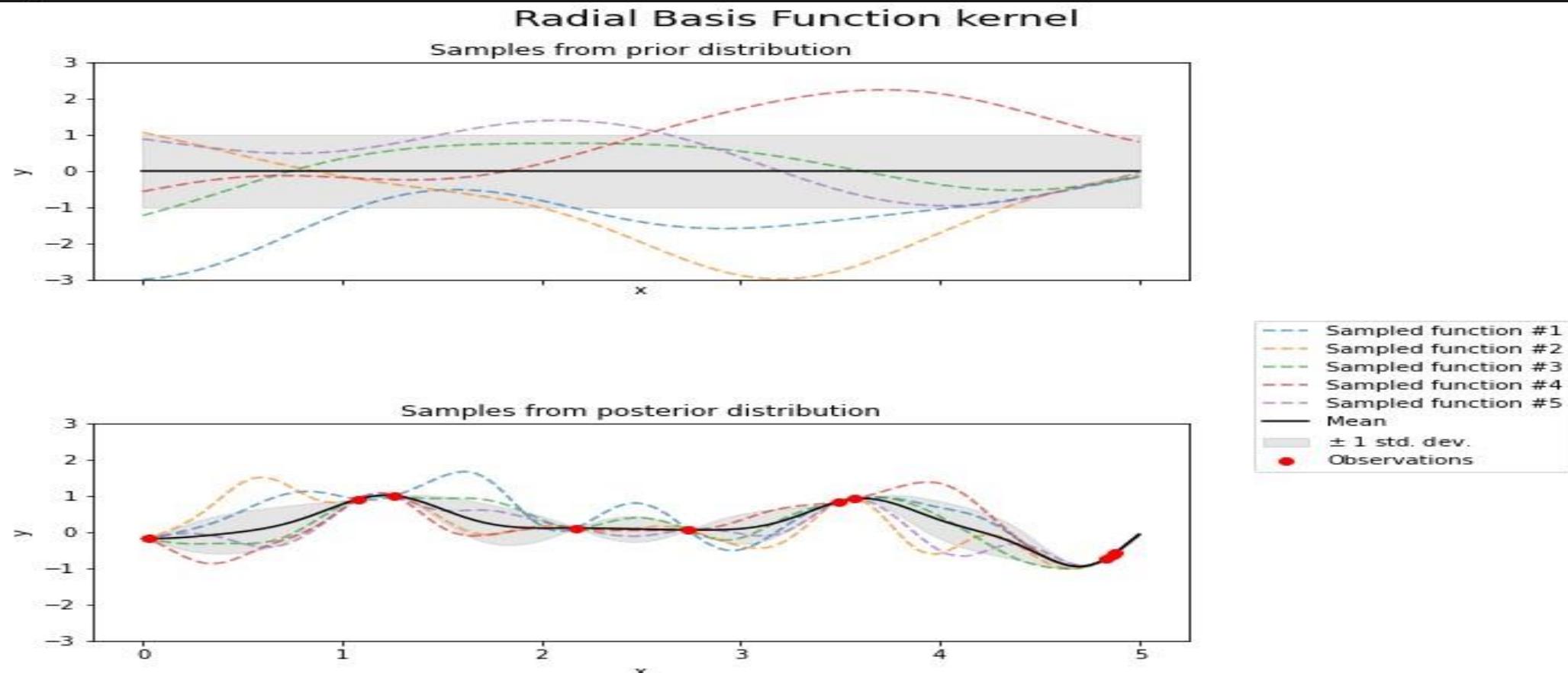
```
221
222 # plot posterior
223 gpr.fit(X_train, y_train)
224 plot_gpr_samples(gpr, n_samples=n_samples, ax=axs[1])
225 axs[1].scatter(X_train[:, 0], y_train, color="red", zorder=10, label="Observations")
226 axs[1].legend(bbox_to_anchor=(1.05, 1.5), loc="upper left")
227 axs[1].set_title("Samples from posterior distribution")
228
229 fig.suptitle("Exponentiation", fontsize=18)
230 plt.tight_layout()
231
232 print(f"Kernel parameters before fit:\n{kernel}")
233 print(
234     f"Kernel parameters after fit: \n{gpr.kernel_} \n"
235     f"Log-likelihood: {gpr.log_marginal_likelihood(gpr.kernel_.theta):.3f}"
236 )
237 #%%
238 ### DotProduct
239
240
241 from sklearn.gaussian_process.kernels import DotProduct
242
243
244 kernel = 1.0 * DotProduct(sigma_0=1,sigma_0_bounds=(1e-5, 1e5))
245 gpr = GaussianProcessRegressor(kernel=kernel, random_state=0)
246
247 fig, axs = plt.subplots(nrows=2, sharex=True, sharey=True, figsize=(10, 8))
248
249 # plot prior
250 plot_gpr_samples(gpr, n_samples=n_samples, ax=axs[0])
251 axs[0].set_title("Samples from prior distribution")
252
253 # plot posterior
254 gpr.fit(X_train, y_train)
255 plot_gpr_samples(gpr, n_samples=n_samples, ax=axs[1])
256 axs[1].scatter(X_train[:, 0], y_train, color="red", zorder=10, label="Observations")
257 axs[1].legend(bbox_to_anchor=(1.05, 1.5), loc="upper left")
258 axs[1].set_title("Samples from posterior distribution")
259
260 fig.suptitle("DotProduct", fontsize=18)
261 plt.tight_layout()
262
263 print(f"Kernel parameters before fit:\n{kernel}")
264 print(
265     f"Kernel parameters after fit: \n{gpr.kernel_} \n"
266     f"Log-likelihood: {gpr.log_marginal_likelihood(gpr.kernel_.theta):.3f}"
267 )
268 #%%
269 ### Sum kernel
270
271
272 from sklearn.gaussian_process.kernels import Sum
273
274
275 kernel = 1.0 * Sum(DotProduct(),WhiteKernel())
276 gpr = GaussianProcessRegressor(kernel=kernel, random_state=0)
```

```
277
278 fig, axs = plt.subplots(nrows=2, sharex=True, sharey=True, figsize=(10, 8))
279
280 # plot prior
281 plot_gpr_samples(gpr, n_samples=n_samples, ax=axs[0])
282 axs[0].set_title("Samples from prior distribution")
283
284 # plot posterior
285 gpr.fit(X_train, y_train)
286 plot_gpr_samples(gpr, n_samples=n_samples, ax=axs[1])
287 axs[1].scatter(X_train[:, 0], y_train, color="red", zorder=10, label="Observations")
288 axs[1].legend(bbox_to_anchor=(1.05, 1.5), loc="upper left")
289 axs[1].set_title("Samples from posterior distribution")
290
291 fig.suptitle("Sum kernel", fontsize=18)
292 plt.tight_layout()
293
294 print(f"Kernel parameters before fit:\n{kernel}")
295 print(
296     f"Kernel parameters after fit: {gpr.kernel_}\n"
297     f"Log-likelihood: {gpr.log_marginal_likelihood(gpr.kernel_.theta):.3f}"
298 )
299 #20%
```

OUTPU

```
In [1]: runcell(0, 'B:/3rd year/5th sem/P&AD/exp9.py')
Kernel parameters before fit:
1**2 * RBF(length_scale=1))+++
Kernel parameters after fit:
0.594**2 * RBF(length_scale=0.279)
Log-likelihood: -0.067
C:\Users\Dhruv Singhal\anaconda3\lib\site-packages\sklearn\gaussian_process\_gpr.py:506: ConvergenceWarning: lbfgs failed to converge (status=2):
ABNORMAL_TERMINATION_IN_LNSRCH.

Increase the number of iterations (max_iter) or scale the data as shown in:
    https://scikit-learn.org/stable/modules/preprocessing.html
    _check_optimize_result("lbfgs", opt_res)
Kernel parameters before fit:
1**2 * RationalQuadratic(alpha=0.1, length_scale=1))
Kernel parameters after fit:
0.594**2 * RationalQuadratic(alpha=1.78e+06, length_scale=0.279)
Log-likelihood: -0.067
```

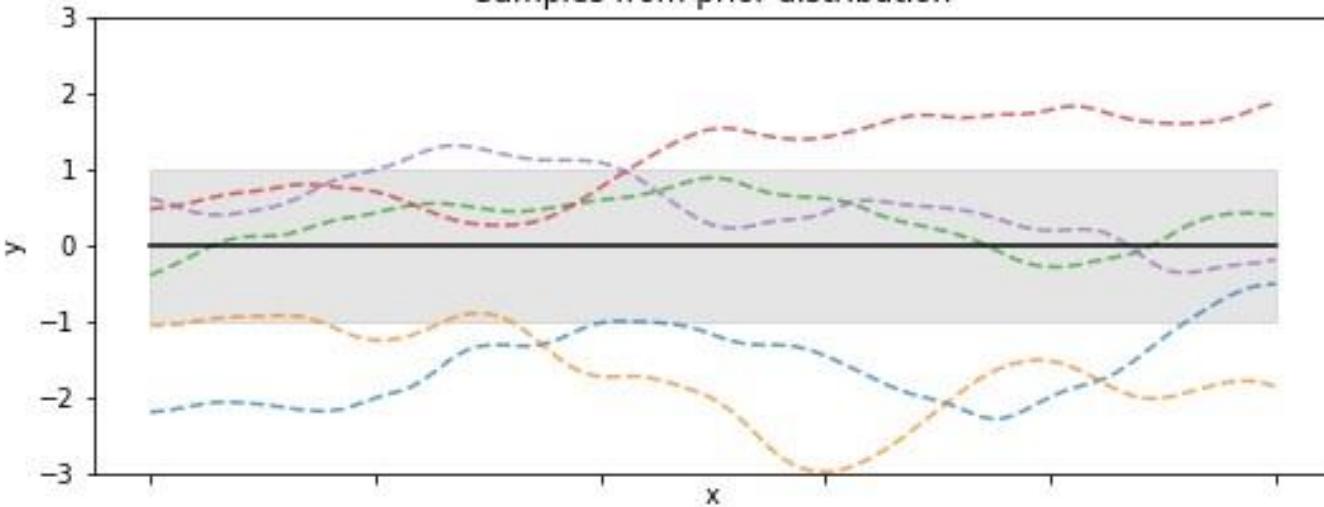


T:

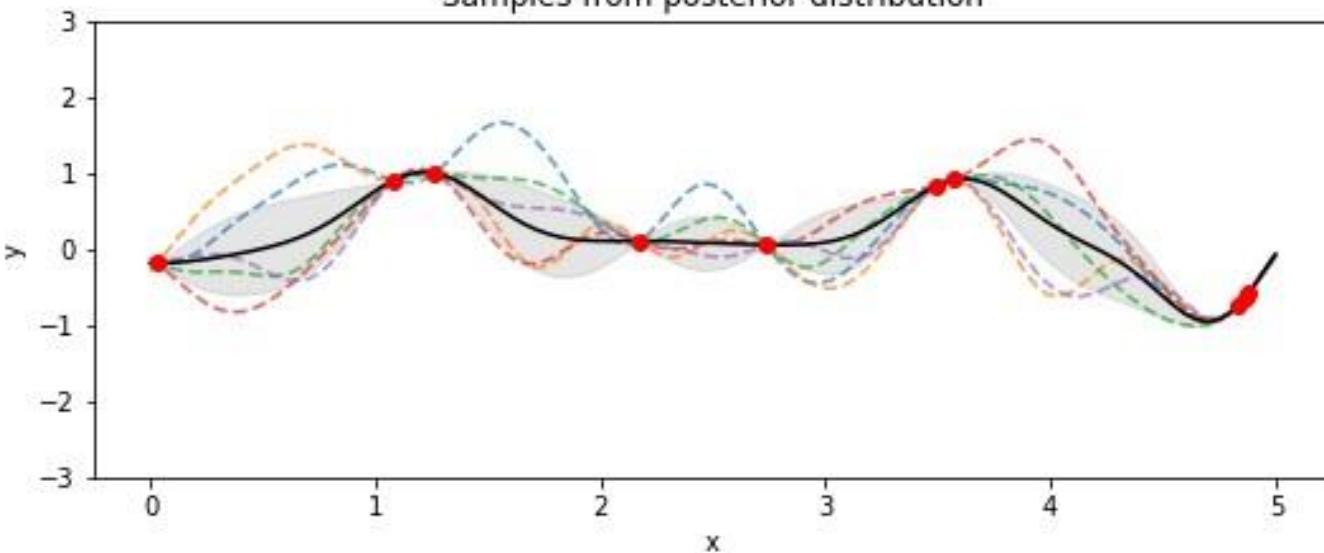


Rational Quadratic kernel

Samples from prior distribution

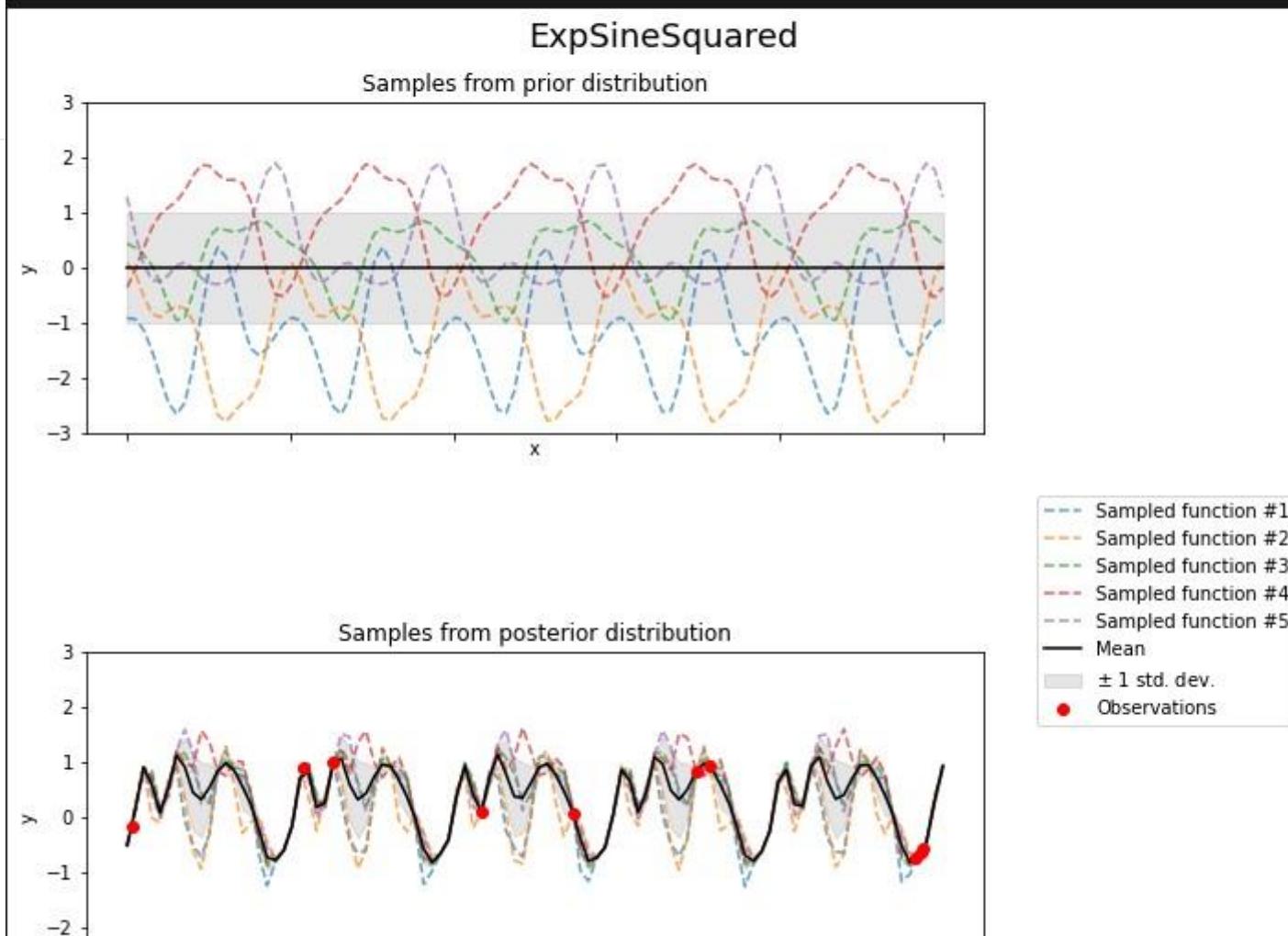


Samples from posterior distribution

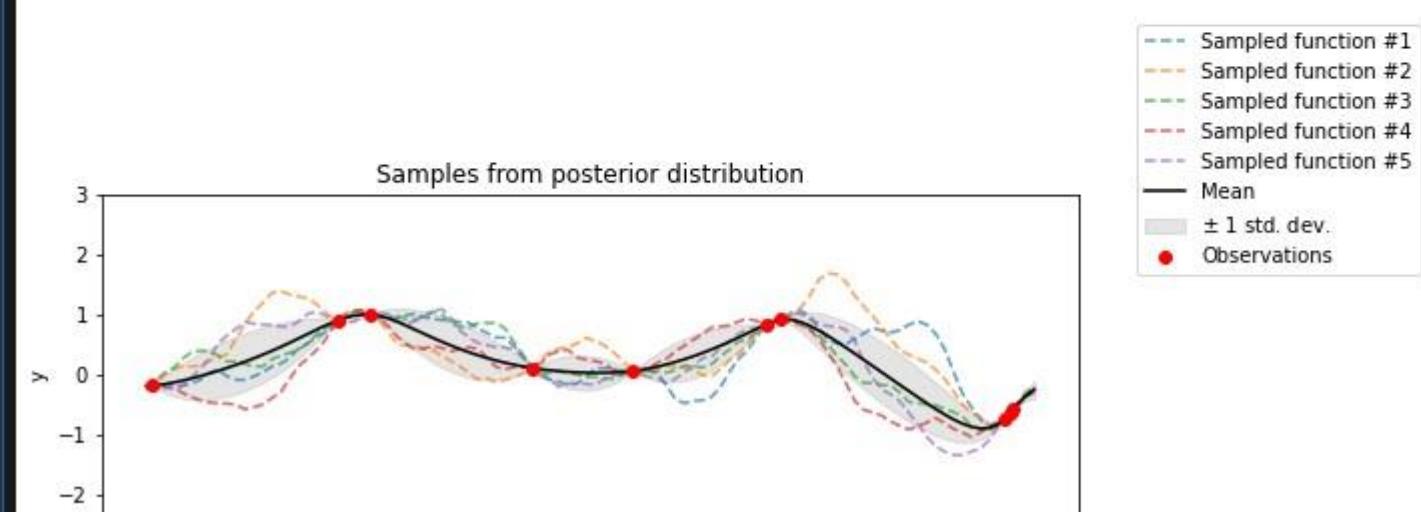
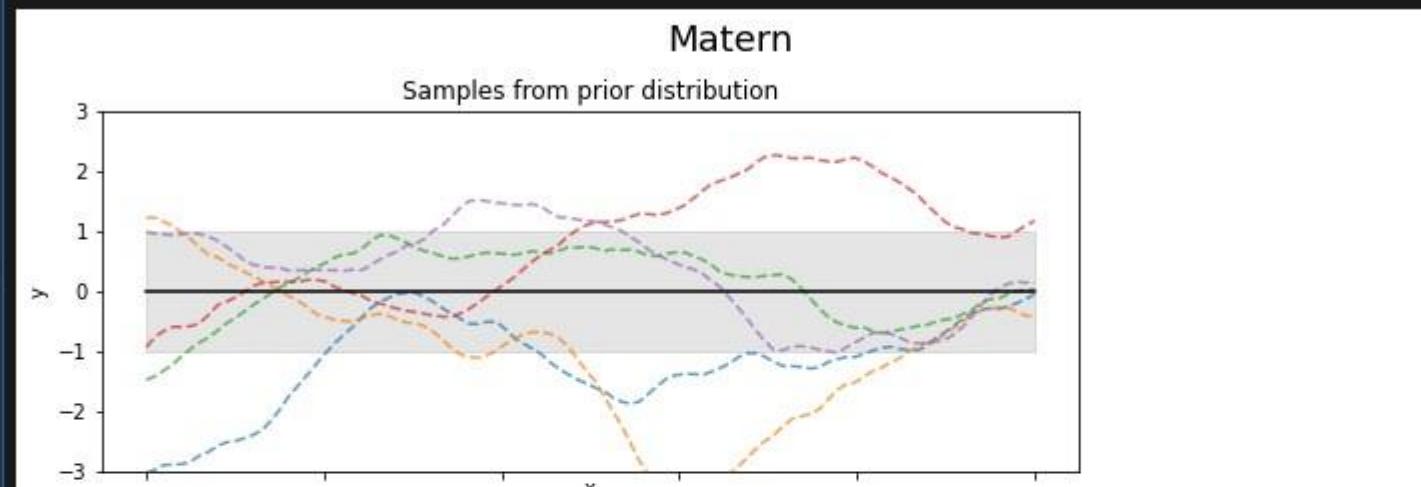


- Sampled function #1
- Sampled function #2
- Sampled function #3
- Sampled function #4
- Sampled function #5
- Mean
- ± 1 std. dev.
- Observations

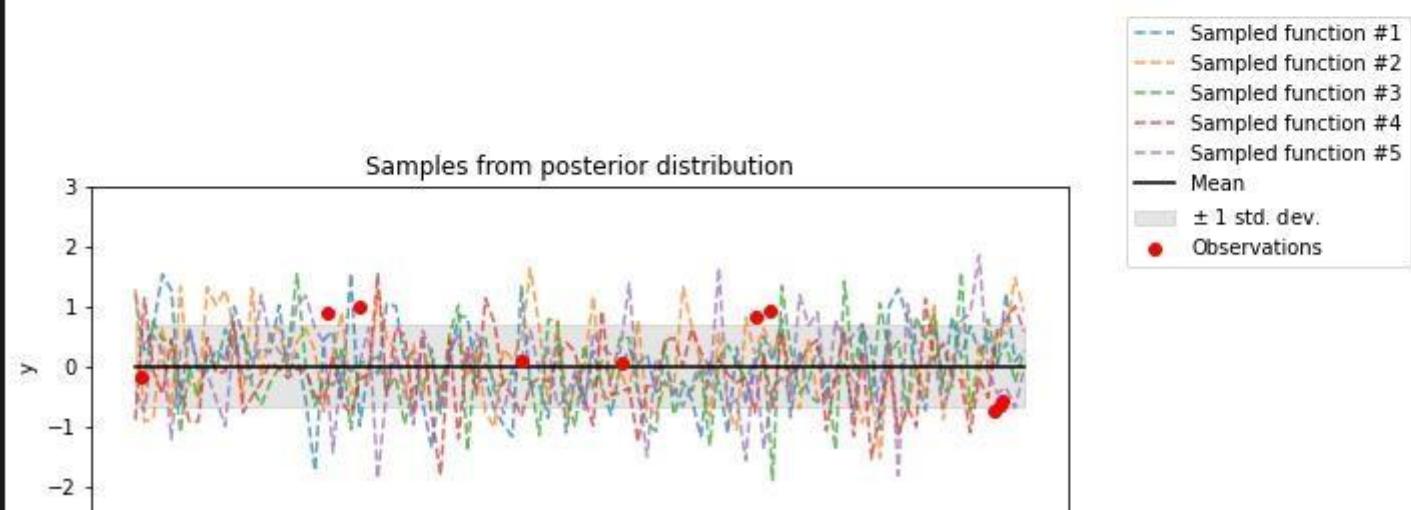
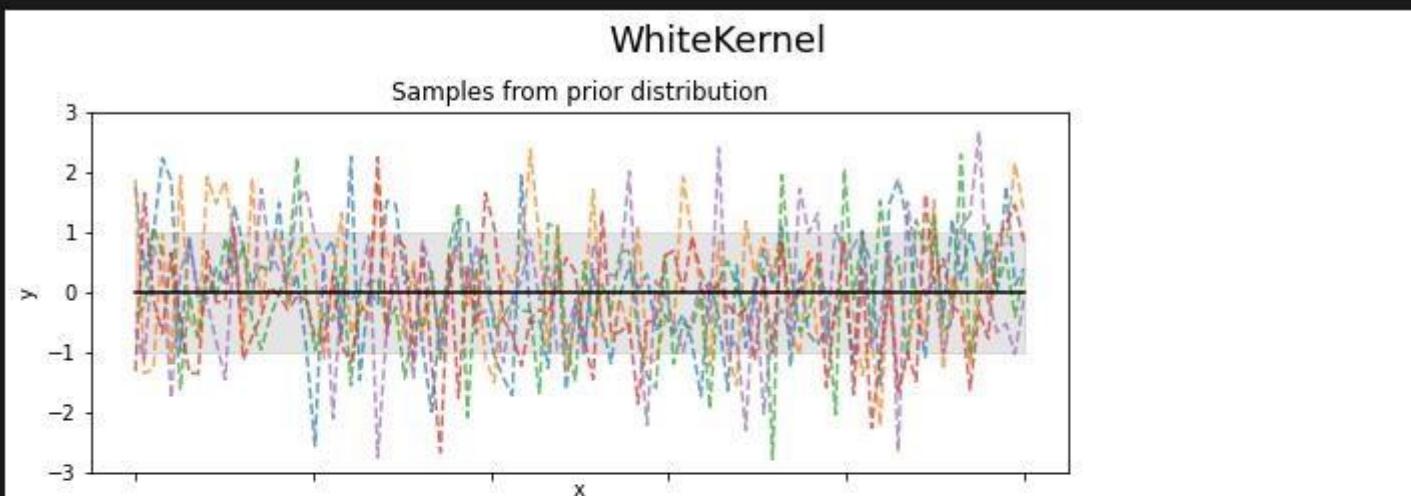
```
In [2]: runcell(1, 'B:/3rd year/5th sem/P&AD/exp9.py')
Kernel parameters before fit:
1**2 * ExpSineSquared(length_scale=1, periodicity=1))
Kernel parameters after fit:
0.736**2 * ExpSineSquared(length_scale=0.438, periodicity=0.978)
Log-likelihood: -5.478
```



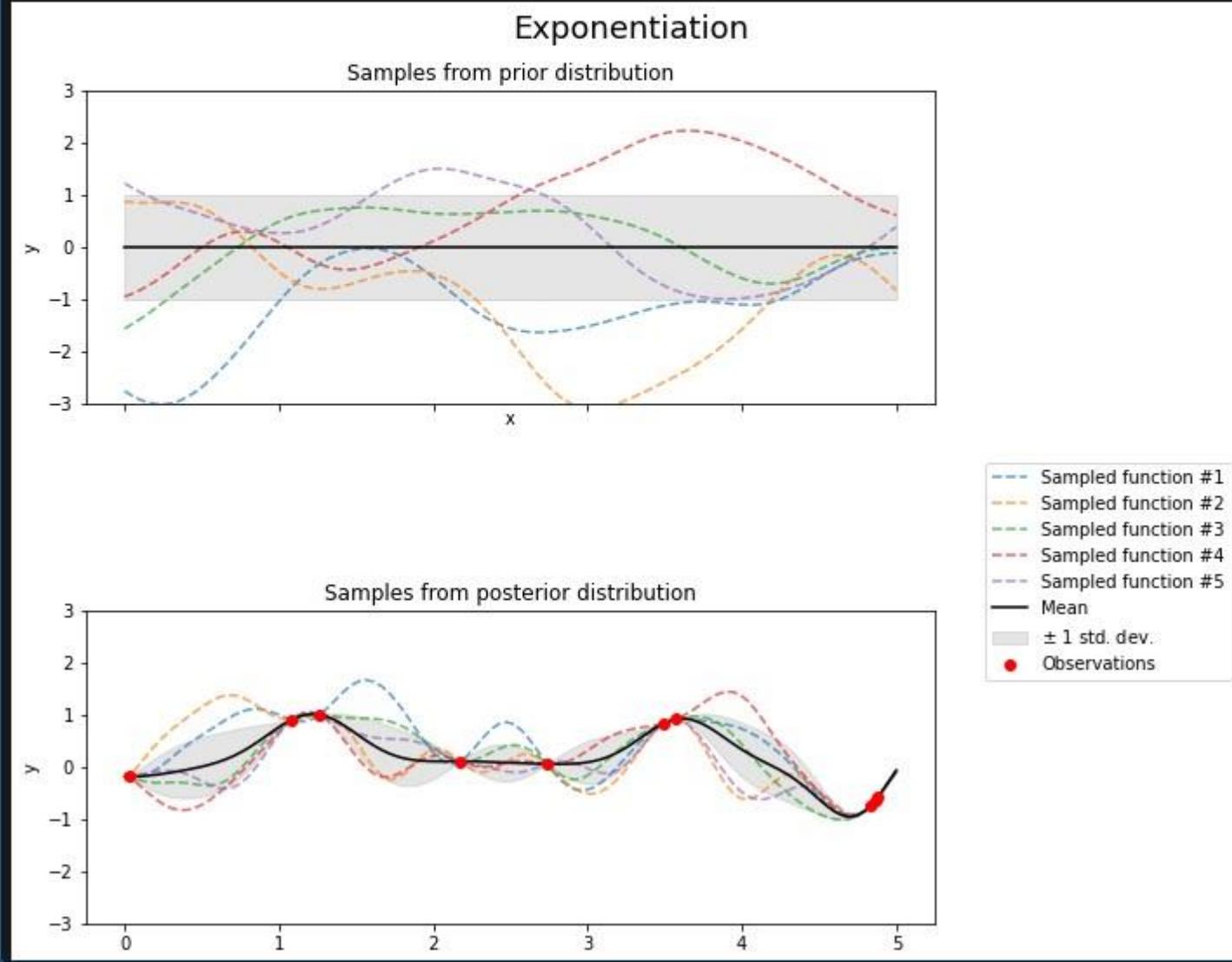
```
In [3]: runcell(2, 'B:/3rd year/5th sem/P&AD/exp9.py')
Kernel parameters before fit:
1**2 * Matern(length_scale=1, nu=1.5))
Kernel parameters after fit:
0.609**2 * Matern(length_scale=0.484, nu=1.5)
Log-likelihood: -1.185
```



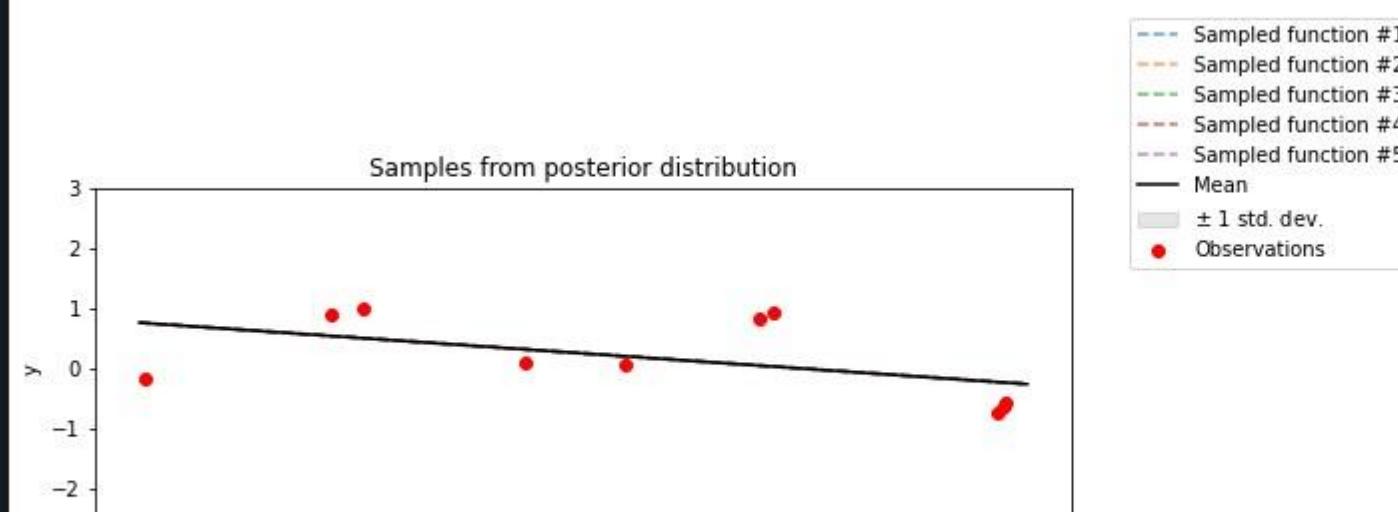
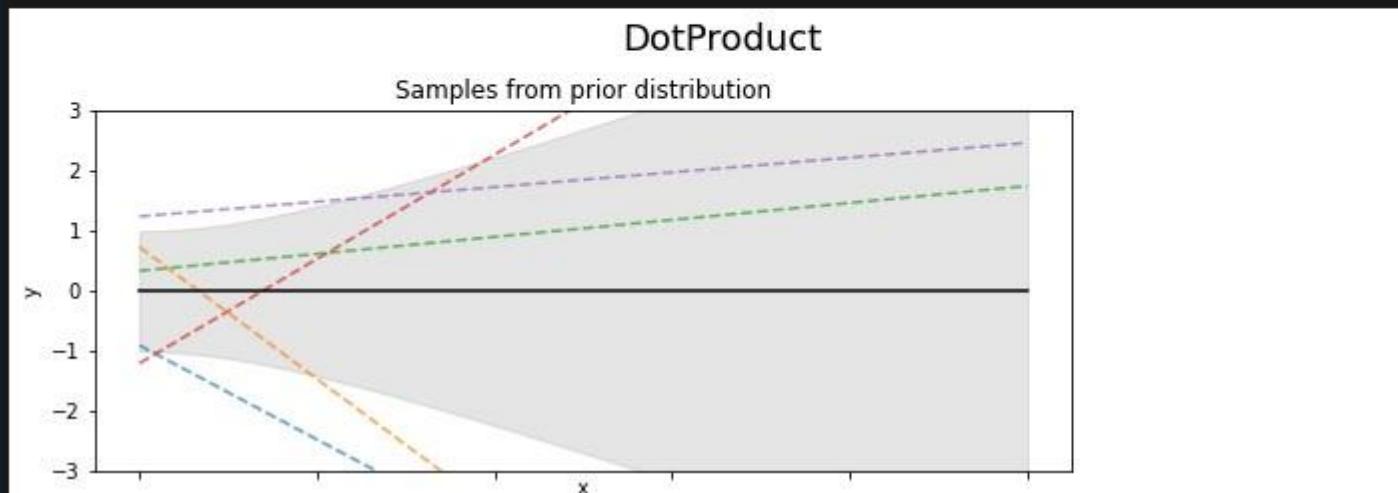
```
In [4]: runcell(3, 'B:/3rd year/5th sem/P&AD/exp9.py')
Kernel parameters before fit:
1**2 * WhiteKernel(noise_level=1))
Kernel parameters after fit:
0.827**2 * WhiteKernel(noise_level=0.684)
Log-likelihood: -10.386
```



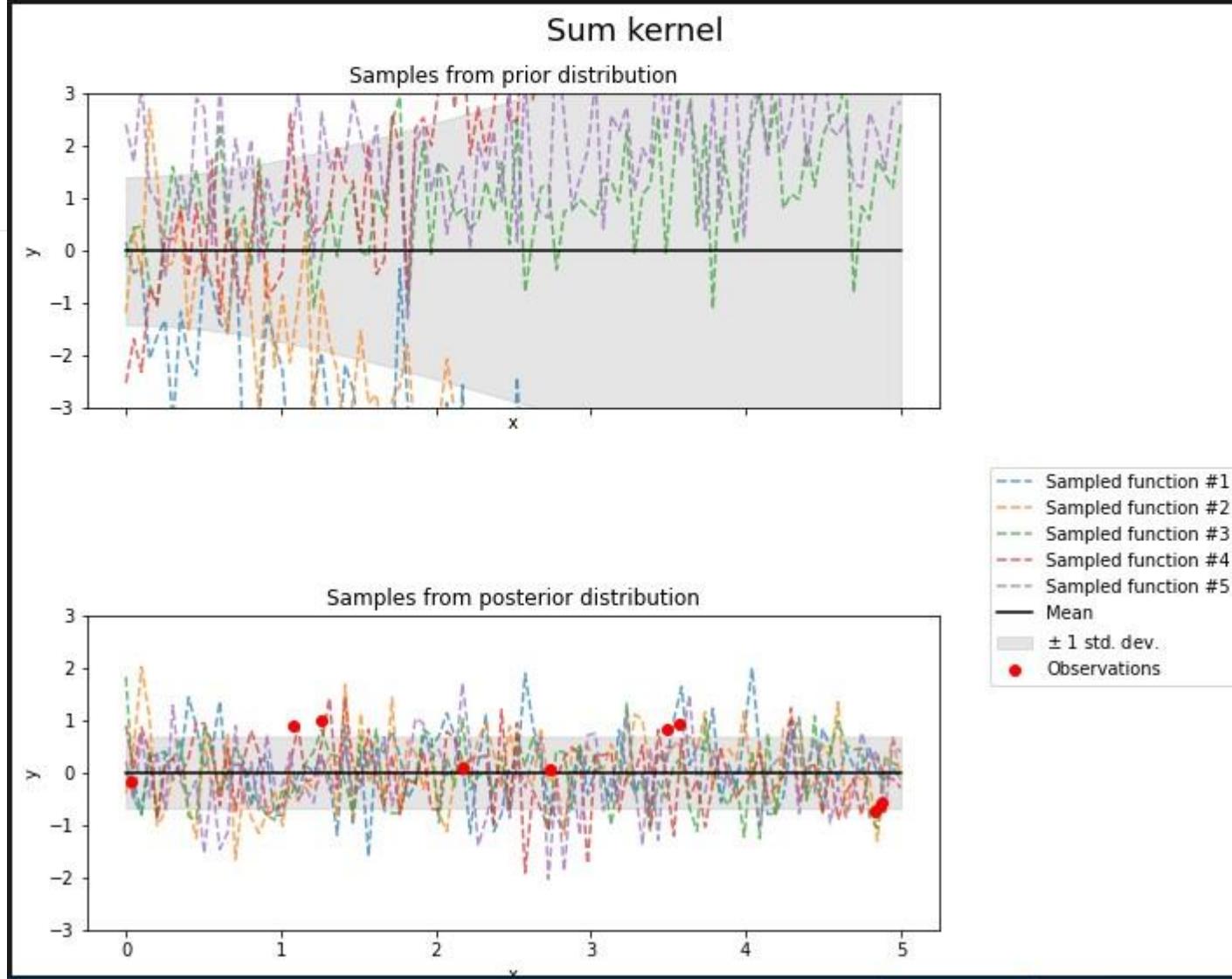
```
In [5]: runcell(4, 'B:/3rd year/5th sem/P&AD/exp9.py')
C:\Users\Dhruv Singhal\anaconda3\lib\site-packages\sklearn\gaussian_process\kernels.py:418: ConvergenceWarning: The optimal value found for dimension 0 of
parameter k2_kernel_alpha is close to the specified upper bound 100000.0. Increasing the bound and calling fit again may find a better value.
  ConvergenceWarning)
Kernel parameters before fit:
1**2 * RationalQuadratic(alpha=1, length_scale=1) ** 2
Kernel parameters after fit:
0.594**2 * RationalQuadratic(alpha=1e+05, length_scale=0.394) ** 2
Log-likelihood: -0.067
```



```
In [6]: runcell(5, 'B:/3rd year/5th sem/P&AD/exp9.py')
Kernel parameters before fit:
1**2 * DotProduct(sigma_0=1))
Kernel parameters after fit:
0.999**2 * DotProduct(sigma_0=0.999)
Log-likelihood: -16319044829.997
```



```
In [7]: runcell(6, 'B:/3rd year/5th sem/P&AD/exp9.py')
C:\Users\Dhruv Singhal\anaconda3\lib\site-packages\sklearn\gaussian_process\kernels.py:409: ConvergenceWarning: The optimal value found for dimension 0 of
parameter k1_constant_value is close to the specified lower bound 1e-05. Decreasing the bound and calling fit again may find a better value.
  ConvergenceWarning)
Kernel parameters before fit:
1**2 * DotProduct(sigma_0=1) + WhiteKernel(noise_level=1)
Kernel parameters after fit:
0.00316**2 * DotProduct(sigma_0=0.339) + WhiteKernel(noise_level=4.67e+04)
Log-likelihood: -10.387
```



Experiment 10



CODE:

```
1 import numpy as np
2
3 x1 = np.random.random( (100, 1))
4 y= 4 + 3*x1 + np.random.randn(100, 1)
5
6 x0 = np.ones( (100, 1))
7 X = np.concatenate( (x0, x1), axis = 1)
8
9
10 %%%
11 temp1 = np.linalg.inv(np.dot (X.T, X))
12 temp2 = np.dot(temp1,X.T)
13 w = np.dot(temp2, y)
14 print("-----")
15 print("Least squares method(Direct) Single Input")
16 print("-----")
17 print("w0",w[0])
18 print("w1",w[1])
19
20 %%
21
22 import numpy as np
23 x1 = np.random.random( (100, 3))
24 X=np.c_[np.ones((100,1)),x1]
25 a=[[4,5,8],
26     [8,5,7],
27     [7,6,3],
28     [1,3,8]]
29 W=np.array(a)
30 y1=np.dot(X,W)
31 temp1 = np.linalg.inv(np.dot (X.T, X))
32 temp2 = np.dot(temp1,X.T)
33 w = np.dot(temp2, y1)
34 print("-----")
35 print("Least squares method(Direct) Multiple Input")
36 print("-----")
37 print("w1's are:\n" ,w)
38
```

```
39 #%%
40 X_ = np.random. random( (100, 3))
41 y1= 4 + 3*X_ + np. random. randn(100, 1)
42 y2= 5 + 2*X_ + np. random. randn(100, 1)
43 y3= 3 + 6*X_ + np. random. randn(100, 1)
44 y4= 7 + 9*X_ + np. random. randn(100, 1)
45 Xwb=np.c_[np.ones((100,1)),X_]
46 W_=Xwb.T.dot(Xwb)
47 tp1 = np.linalg.inv(np.dot (Xwb.T, Xwb))
48 tp2 = np.dot(tp1,Xwb.T)
49 w1 = np. dot(tp2, y1)
50 w2 = np. dot(tp2, y2)
51 w3 = np. dot(tp2, y3)
52 w4 = np. dot(tp2, y4)
53
54
55 print("-----Modified-----")
56 print("Least squares method(Direct) Multiple Input")
57 print("-----")
58 print("W1:\n",w1)
59 print("W2:\n",w2)
60 print("W3:\n",w3)
61 print("W4:\n",w4)
62 #%%
63 print(np.concatenate((w1,w2,w3,w4)))
```

OUTPUT:

```
IPython 7.26.0 -- An enhanced Interactive Python.

In [1]: runcell(0, 'B:/3rd year/5th sem/P&AD/exp10.py')

In [2]: runcell(1, 'B:/3rd year/5th sem/P&AD/exp10.py')
-----
Least squares method(Direct) Single Input
-----
W0 [3.6047836]
W1 [3.50268276]

In [3]: runcell(2, 'B:/3rd year/5th sem/P&AD/exp10.py')
-----
Least squares method(Direct) Multiple Input
-----
W1's are:
[[4. 5. 8.]
 [8. 5. 7.]
 [7. 6. 3.]
 [1. 3. 8.]]
```

```
In [4]: runcell(3, 'B:/3rd year/5th sem/P&AD/exp10.py')
-----Modified-----
Least squares method(Direct) Multiple Input
-----
W1:
[[ 3.87767253  3.87767253  3.87767253]
 [ 3.05152315  0.05152315  0.05152315]
 [ 0.3952086   3.3952086   0.3952086 ]
 [-0.24948184 -0.24948184  2.75051816]]
W2:
[[ 5.67017869  5.67017869  5.67017869]
 [ 1.25767771 -0.74232229 -0.74232229]
 [-0.28164296  1.71835704 -0.28164296]
 [-0.2510575   -0.2510575   1.7489425 ]]
W3:
[[ 3.00627306  3.00627306  3.00627306]
 [ 6.48435087  0.48435087  0.48435087]
 [-0.09223401  5.90776599 -0.09223401]
 [-0.3093984   -0.3093984   5.6906016 ]]
W4:
[[ 6.98119803  6.98119803  6.98119803]
 [ 9.46229368  0.46229368  0.46229368]
 [-0.05479141  8.94520859 -0.05479141]
 [-0.61903883 -0.61903883  8.38096117]]
In [5]: runcell(4, 'B:/3rd year/5th sem/P&AD/exp10.py')
[[ 3.87767253  3.87767253  3.87767253]
 [ 3.05152315  0.05152315  0.05152315]
 [ 0.3952086   3.3952086   0.3952086 ]
 [-0.24948184 -0.24948184  2.75051816]
 [ 5.67017869  5.67017869  5.67017869]
 [ 1.25767771 -0.74232229 -0.74232229]
 [-0.28164296  1.71835704 -0.28164296]
 [-0.2510575   -0.2510575   1.7489425 ]
 [ 3.00627306  3.00627306  3.00627306]
 [ 6.48435087  0.48435087  0.48435087]
 [-0.09223401  5.90776599 -0.09223401]
 [-0.3093984   -0.3093984   5.6906016 ]
 [ 6.98119803  6.98119803  6.98119803]
 [ 9.46229368  0.46229368  0.46229368]
 [-0.05479141  8.94520859 -0.05479141]
 [-0.61903883 -0.61903883  8.38096117]]
```

Experiment 11

Different Anomaly Detection Algorithms on 2D Dataset

Importing required libraries

```
I [3]: importim  
importnumba_n  
importmatplotlib  
importmatplotlibnvploa_nl  
fro sklearimporsv  
fro skleardatasetimpormake_moonmake_blob  
fro sklearcovarianimporellipticEnvelo  
fro sklearensemblimporiisolatlonFore  
fro skleameighboimpordLocalOutlierFact
```

Data

for this experiment we have use make_moons, make_blobs datsets which are already in sklearn.datasets

sklearn.datasets.make_blobs(n_samples=100, n_features=2, *, centers=None, cluster_std=1.0, center_box=(-10.0, 10.0), shuffle=True, random_state=None, return_centers=False) Generate isotropic Gaussian blobs for clustering.

sklearn.datasets.make_moons(n_samples=100, *, shuffle=True, noise=None, random_state=None) Make two interleaving half circles. A simple toy dataset to visualize clustering and classification algorithms.

Preprocessing & Visualization

for this experiment we do not need preprocessing as we need to visualize different anomaly detection algorithms

Model Building

Covariance Matrix

This metric assesses how many standard deviations σ away x_i is from μ , thereby being dimensionless. An extreme observation has a large distance from the center of a distribution. It is useful in predictive maintenance as it enables finding unusual behavior of a system; upcoming down-times and malfunction can be targeted and identified, which is particularly important since finding faults in equipment, machines, and devices early can reduce the extent of the damage.

One Class SVM

One Class SVM mode is trained in only one class, referred to as the normal class. The model learns all the features and patterns of the normal class . When a new observation is introduced to the model, then based on its learning, the One Class SVM detects if the new observation deviates from the normal behaviour then it classifies it as an oulier else the new observation will identified as an inlier. One Class SVM is based on Support Vector Machines where the binary classes are separated by a non-linear hyper plane

Local Outlier Factor

Local outlier factor (LOF) is an algorithm used for Unsupervised outlier detection. It produces an anomaly score that represents data points which are outliers in the data set. It does this

by measuring the local density deviation of a given data point with respect to the data points near it.

Isolation Forest

Isolation Forest is based on the Decision Tree algorithm. It isolates the outliers by randomly selecting a feature from the given set of features and then randomly selecting a split value between the max and min values of that feature. This random partitioning of features will produce shorter paths in trees for the anomalous data points, thus distinguishing them from the rest of the data.

Compile & Train

Altering the hyperparameters to see the effect of the hyperparameter on the working of the algorithm

```
In [4]: matplotlib.rcParams['contour.negative_linestyle'] = 'solid'
n_samples = 300
outliers_fraction = 0.15
n_outliers = int(outliers_fraction * n_samples)
n_inliers = n_samples - n_outliers
anomaly_algorithms = [
    ("Robust covariance", EllipticEnvelope(contamination=outliers_fraction)),
    ("One-Class SVM", svm.OneClassSVM(nu=outliers_fraction, kernel="poly",
                                       gamma=0.1)),
    ("Isolation Forest", IsolationForest(contamination=outliers_fraction,
                                           random_state=42)),
    ("Local Outlier Factor", LocalOutlierFactor(
        n_neighbors=35, contamination=outliers_fraction))]

# Define datasets
blobs_params = dict(random_state=0, n_samples=n_inliers, n_features=2)
datasets = [
    make_blobs(centers=[[0, 0], [0, 0]], cluster_std=0.5,
               **blobs_params)[0],
    make_blobs(centers=[[2, 2], [-2, -2]], cluster_std=[0.5, 0.5],
               **blobs_params)[0],
    make_blobs(centers=[[2, 2], [-2, -2]], cluster_std=[1.5, .3],
               **blobs_params)[0],
    4. * (make_moons(n_samples=n_samples, noise=.05, random_state=0)[0] -
           np.array([0.5, 0.25])),
    14. * (np.random.RandomState(42).rand(n_samples, 2) - 0.5)]
matplotlib.rcParams['contour.negative_linestyle'] = 'solid'
# Compare given classifiers under given settings
xx, yy = np.meshgrid(np.linspace(-7, 7, 150),
                     np.linspace(-7, 7, 150))

plt.figure(figsize=(len(anomaly_algorithms) * 2 + 4, 12.5))
plt.subplots_adjust(left=.02, right=.98, bottom=.001, top=.96, wspace=.05,
                    hspace=.01)

plot_num = 1
rng = np.random.RandomState(42)

for i_dataset, X in enumerate(datasets):
    # Add outliers
    X = np.concatenate([X, rng.uniform(low=-6, high=6, size=(n_outliers, 2))],
                      axis=0)

    for name, algorithm in anomaly_algorithms:
        t0 = time.time()
        algorithm.fit(X)
        t1 = time.time()
        plt.subplot(len(datasets), len(anomaly_algorithms), plot_num)
        if i_dataset == 0:
            plt.title(name, size=18)

        # fit the data and tag outliers
        if name == "Local Outlier Factor":
            y_pred = algorithm.fit_predict(X)
        else:
            y_pred = algorithm.fit(X).predict(X)

        # plot the Levels Lines and the points
        if name != "Local Outlier Factor": # LOF does not implement predict
            Z = algorithm.predict(np.c_[xx.ravel(), yy.ravel()])
            Z = Z.reshape(xx.shape)
            plt.contour(xx, yy, Z, levels=[0], linewidths=2, colors='black')

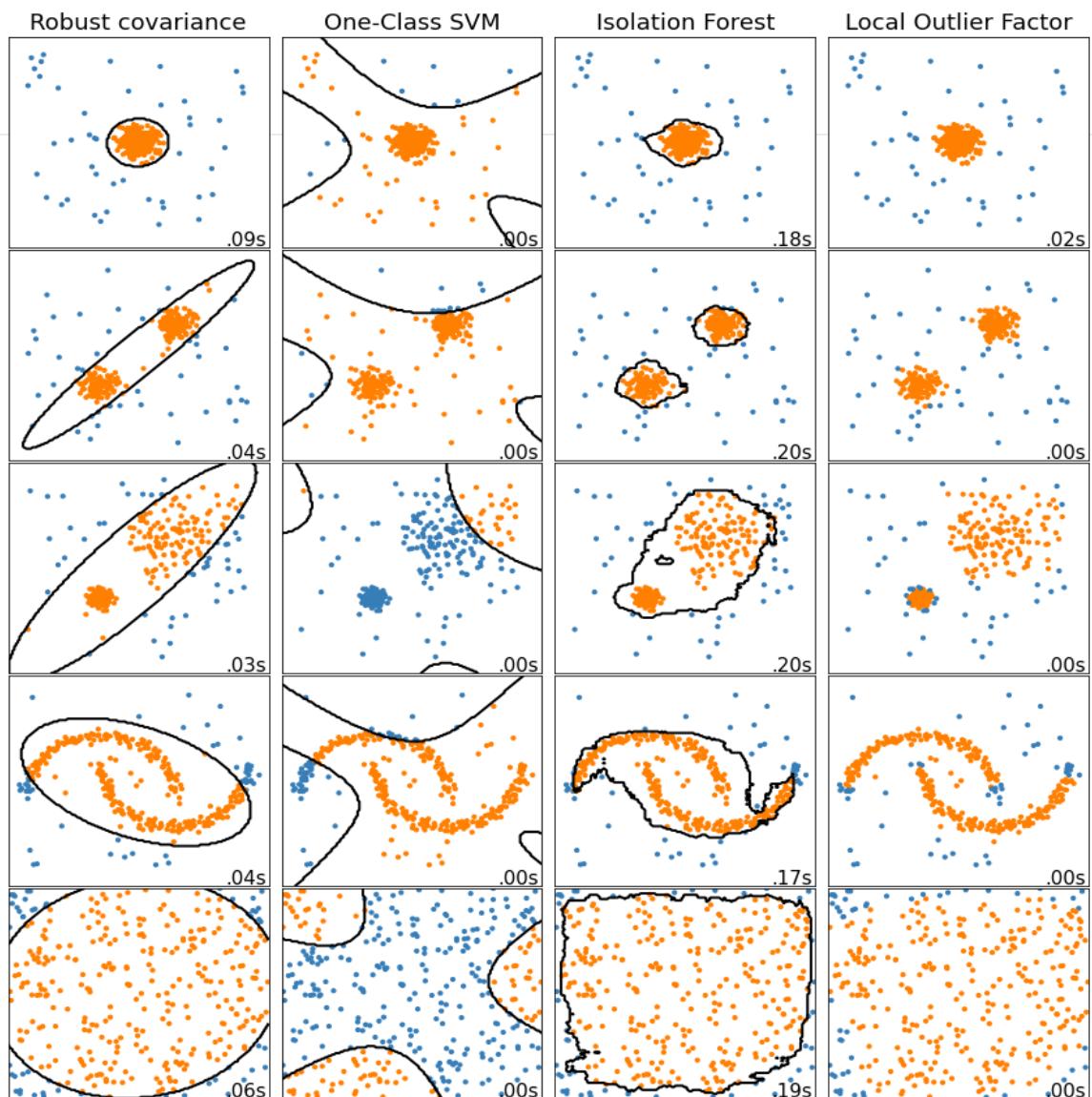
        colors = np.array(['#377eb8', '#ff7f00'])
        plt.scatter(X[:, 0], X[:, 1], s=10, color=colors[(y_pred + 1) // 2])

        plt.xlim(-7, 7)
        plt.ylim(-7, 7)
        plt.xticks(())
        plt.yticks(())
        plt.text(.99, .01, ('%.2fs' % (t1 - t0)).lstrip('0'),
                 transform=plt.gca().transAxes, size=15,
                 horizontalalignment='right')
        plot_num += 1
```

Result

By this experiment i came to know the working of different anomaly detection algorithms and how they differ from each other in working and in time perspective how they find anomaly

In []:



Experiment 12

Support Vector Machines

Importing required libraries

```
In [1]: import numpy as np
import matplotlib.pyplot as plt
from sklearn.model_selection import train_test_split
from sklearn import svm
from sklearn.datasets import make_classification
from sklearn.model_selection import GridSearchCV
from sklearn.model_selection import KFold
from sklearn.model_selection import cross_val_score
```

Creating Data

Creating Dataset with 5000 sample points and 2 features and 2 classes and number of informative features are 2 and number of redundant features are 0 and number of repeated features are 0 and random_state is 50 and difference between 2 class is set as 3.5

```
In [2]: X, y = make_classification(n_samples=5000, n_features=2, n_classes=2
, n_informative=2, n_redundant=0, n_repeated=0, random_state=0, class_sep=2.5)
```

Preprocessing & Visualization

since we have created dataset using make_classification , it by default make dataset's mean tending to 0 and Standard Deviation to 1 which doesn't require any preprocessing.

Model Building

The goal of the SVM algorithm is to create the best line or decision boundary that can segregate n-dimensional space into classes so that we can easily put the new data point in the correct category in the future. This best decision boundary is called a hyperplane.

SVM chooses the extreme points/vectors that help in creating the hyperplane. These extreme cases are called as support vectors, and hence algorithm is termed as Support Vector Machine. In this case we have called the in built function svm.SVC

Compile & Train

Firstly splitting the input dataset into training and test parts following which training set is passed through svm.SVC and is fitted to it. and predicting over the test test to check the model performance.

```
In [3]: X_train,X_test,Y_train,Y_test=train_test_split(X,y,test_size=0.20,random_state=42)
clf = svm.SVC(kernel="linear", C=10)
clf.fit(X_train,Y_train)
Y_pred=clf.predict(X_test)
print("train Accuracy:",clf.score(X_train,Y_train))
print("test Accuracy:",clf.score(X_test,Y_test))
plt.scatter(X[:, 0], X[:, 1], c=y, s=30, cmap=plt.cm.Paired)

# plot the decision function

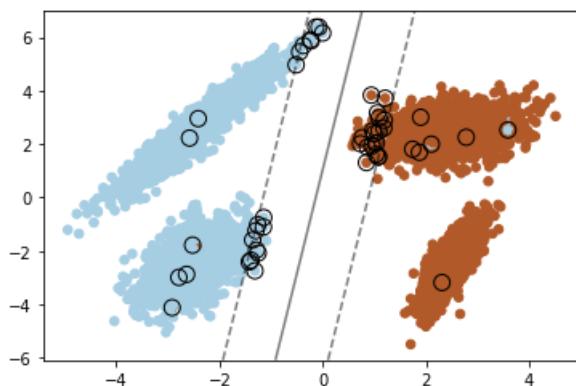
ax = plt.gca()
xlim = ax.get_xlim()
ylim = ax.get_ylim()

# create grid to evaluate model
xx = np.linspace(xlim[0], xlim[1], 30)
yy = np.linspace(ylim[0], ylim[1], 30)
YY, XX = np.meshgrid(yy, xx)
xy = np.vstack([XX.ravel(), YY.ravel()]).T
Z = clf.decision_function(xy).reshape(XX.shape)

# plot decision boundary and margins
ax.contour(
    XX, YY, Z, colors="k", levels=[-1, 0, 1], alpha=0.5, linestyles=["--", "-", "--"])
)
# plot support vectors
ax.scatter(
    clf.support_vectors_[:, 0],
    clf.support_vectors_[:, 1],
    s=100,
    linewidth=1,
    facecolors="none",
    edgecolors="k",
)
)

plt.show()

train Accuracy: 0.99675
test Accuracy: 0.996
```



Result

since training is giving 99.67% accuracy and in testing it is giving 99.60% accuracy so to check if model is overfitting or not or to improve the validation accuracy more and make more generalized model . So to overcome this we are using k fold cross validation and hyperparameter tuning to improve the overall score and generate generalizesd model and knowing which hyperparameter combination gives best overall score

```
In [4]: #Cross Validation without Hyper parameter Tuning
clf = svm.SVC(kernel='linear', C=2)
kf = KFold(n_splits=5)
score = cross_val_score(clf, X, y, cv=kf)
print("Average Cross Validation score ( Accuracy):", np.mean(score))
```

Average Cross Validation score (Accuracy) :0.9965906550413592

```
In [5]: # Hyper parameter tuning on SVM
tuned_parameters = [{"kernel": "linear", "C": [100, 10, 20, 40, 50]}, {"kernel": "rbf", "C": [90, 40, 50, 100, 90]}]

clf = GridSearchCV(SVC()), tuned_parameters, scoring='accuracy')
clf.fit(X, y)
print("Best parameters set found on development set:")
print()
print(clf.best_params_)
print()
print("Best Score:", clf.best_score_)
```

Best parameters set found on development set:

```
{'C': 100, 'kernel': 'linear'}
Best Score: 0.9966000000000002
```

```
In [6]: %% kfold cross validation with hyperparameter tuning
k = 5
kf = KFold(n_splits=k, random_state=None)
model = GridSearchCV(SVC()), tuned_parameters, scoring='accuracy')

for train_index, test_index in kf.split(X):
    X_train, X_test = X[train_index], X[test_index]
    y_train, y_test = y[train_index], y[test_index]

    model.fit(X_train, y_train)
    pred_values = model.predict(X_test)
print("Best parameters set found on development set:")
print()
print(model.best_params_)
print()
print("Best Score:", model.best_score_)
```

Best parameters set found on development set:

```
{'C': 100, 'kernel': 'linear'}
Best Score: 0.9964999999999999
```

Now the overall score
is 99.649% accuracy

In []:

Experiment 13

Random Forest Classifier

Importing required libraries

```
In [1]: from sklearn.datasets import make_classification
from sklearn.model_selection import train_test_split
from sklearn.ensemble import RandomForestClassifier
import seaborn as sns
import pandas as pd
import matplotlib.pyplot as plt
import warnings
warnings.filterwarnings('ignore')
```

Data creation

Creating Dataset using make_classification with 5000 sample points and 5 classes with 20 number of features in which 10 are informative features and 5 are redundant features and shuffling is turned ON

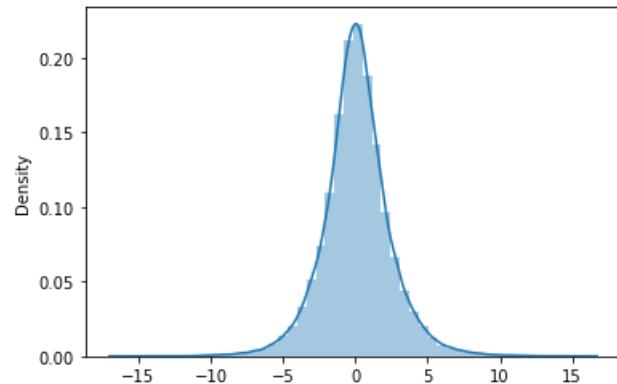
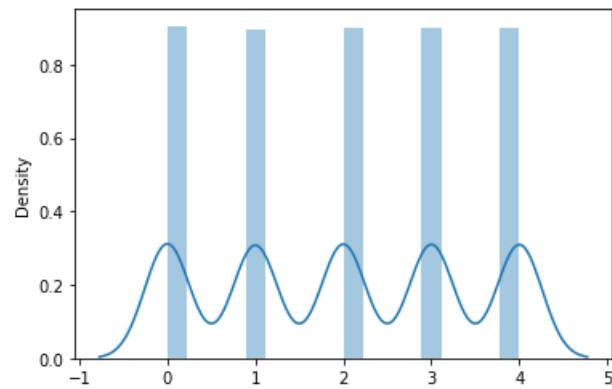
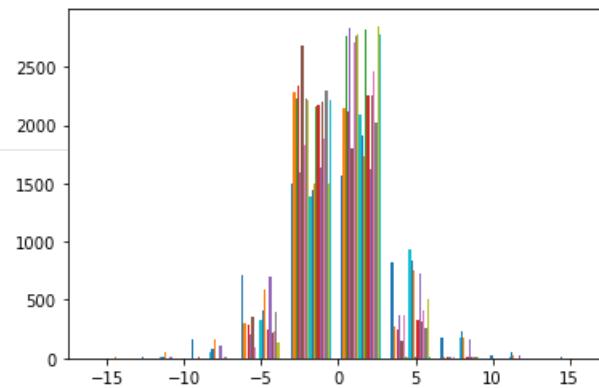
```
In [2]: X,y = make_classification(n_samples=5000,n_classes=5,n_features=20,
n_informative=10,n_redundant=5,shuffle=True,random_state=0)
X=pd.DataFrame(X)
y=pd.DataFrame(y)
```

Preprocessing & Visualization

since we have created dataset using make_classification , it by default make dataset's mean tending to 0 and Standard Deviation to 1 which doesn't require any preprocessing. For Visualization varition of values of X(with 20 features) and y's distplot to observe y's varition of values and distplot of X to check if any pre processing is to be required

In [3]:

```
plt.hist(X)
plt.show()
sns.distplot(y)
plt.show()
sns.distplot(X)
plt.show()
```



feature extraction preprocessing, correlation matrix to check if any pre processing is to be required

```
In [4]: _,graph=plt.subplots(figsize=(15,10))
sns.heatmap(X.corr(),annot=True,ax=graph,square=True)
plt.show()
```



Model Building

Random forest classifier creates a set of decision trees from a randomly selected subset of the training set. It is basically a set of decision trees from a randomly selected subset of the training set and then it collects the votes from different decision trees to decide the final prediction. Assuming dataset has "m" features, the random forest will randomly choose "k" features where $k < m$. Now, the algorithm will calculate the root node among the k features by picking a node that has the highest information gain. After that, the algorithm splits the node into child nodes and repeats this process "n" times. Now we have a forest with n trees. Finally, perform bootstrapping, ie, combine the results of all the decision trees present in the forest.

We have directly used the in built function of sklearn library i.e., from sklearn.ensemble import RandomForestClassifier which directly imports this classifier

Compile & Train

Firstly splitting the input dataset into training and test parts following which training set is passed through RandomForestClassifier and is fitted to it. and predicting over the test test to check the model performance.

```
In [5]: X_train,X_test,Y_train,Y_test=train_test_split (X,y,test_size=0.20,random_state=42)

model=RandomForestClassifier ()
model.fit(X_train,Y_train)
Y_pred=model.predict(X_test)
print("train Accuracy:",model.score(X_train,Y_train))
print("test Accuracy:",model.score(X_test,Y_test))

tr
ai
n
Ac
cu
ra
cy
:
1.
0
te
st
Ac
cu
ra
cy
:
0.
79
```

Result

since training is giving 100% accuracy and in testing it is giving 79% accuracy so it's clearly overfitting so to overcome this we are using k fold cross validation and hyperparameter tuning to improve the overall score and generate generalized model and knowing which hyperparameter combination gives best overall score

```
In [6]: #Cross Validation without Hyper
parameter tuning from
sklearn.model_selection import KFold
from sklearn.model_selection import
cross_val_score clf =
RandomForestClassifier(max_depth=10,
random_state=0) kf=KFold(n_splits=7)
score=cross_val_score(clf, X, y,
scoring='accuracy',cv=kf) print("Cross
Validation Scores are {}".format(score))
```

```

print("Average Cross Validation score
:{}\n".format(score.mean()))

Cross Validation Scores are [0.78461538 0.77062937 0.80672269 0.82633053 0.78431373
0.78011204
0.79411765]
Average Cross Validation score :0.7924059134143168

```

```

In [7]: from sklearn.model_selection import GridSearchCV tuned_parameters =
[{'n_estimators':[10,20,40,100], 'criterion':['gini', 'entropy'],
'max_features':['auto', 'sqrt', 'log2'], 'bootstrap':[True, False]}]
clf=GridSearchCV(RandomForestClassifier(),tuned_parameters,scoring='
accuracy'),verbose=3) clf.fit(X,y)
print("Best parameters set found on
development set:") print()
print(clf.best
_params_)
print()
print("Best
Score:",clf.be
st_score_)

= 0.3s
[CV 1/5] END bootstrap=False, criterion=gini, max_features=auto, n_estimators=40;;
score=0.804 total time
= 0.7s
[CV 2/5] END bootstrap=False, criterion=gini, max_features=auto, n_estimators=40;;
score=0.805 total time
= 0.7s
[CV 3/5] END bootstrap=False, criterion=gini, max_features=auto, n_estimators=40;;
score=0.813 total time
= 0.7s
[CV 4/5] END bootstrap=False, criterion=gini, max_features=auto, n_estimators=40;;
score=0.794 total time
= 0.8s
[CV 5/5] END bootstrap=False, criterion=gini, max_features=auto, n_estimators=40;;
score=0.804 total time
= 0.8s
[CV 1/5] END bootstrap=False, criterion=gini, max_features=auto, n_estimators=100;;
score=0.810 total tim e= 2.0s
[CV 2/5] END bootstrap=False, criterion=gini, max_features=auto, n_estimators=100;;
score=0.827 total tim e= 1.9s
[CV 3/5] END bootstrap=False, criterion=gini, max_features=auto, n_estimators=100;;
score=0.837 total tim e= 2.0s
[CV 4/5] END bootstrap=False, criterion=gini, max_features=auto, n_estimators=100;;
score=0.805 total tim e= 2.0s
[CV 5/5] END bootstrap=False, criterion=gini, max_features=auto, n_estimators=100;;
score=0.821 total tim
2

Best parameters set found on development set:

{'bootstrap': False, 'criterion': 'gini', 'max_features': 'log2', 'n_estimators': 100}

Best Score: 0.8214 Now the overall score is 82.14% accuracy which is above the testing accuracy without
cross validation

```

Experiment 14

DBSCAN Implementation

Importing Required Libraries

```
In [1]: from sklearn.cluster import DBSCAN
import pandas as pd
import matplotlib.pyplot as plt
import warnings
from sklearn.preprocessing import StandardScaler
warnings.filterwarnings('ignore')
```

Data

for this experiment we have use make_blobs datasets which are already in sklearn.datasets

sklearn.datasets.make_blobs(n_samples=100, n_features=2, *, centers=None, cluster_std=1.0, center_box=(-10.0, 10.0), shuffle=True, random_state=None, return_centers=False) Generate isotropic Gaussian blobs for clustering.

```
In [2]: from sklearn.datasets import make_blobs
```

```
In [3]: X,y=make_blobs(n_samples=1000, n_features=2, centers=[[0.5, 2], [-1, -1], [1.5, -1]], cluster_std=0.5, center_box=(-10.0, 10.0), shuffle=True, random_state=42)
```

```
In [4]: X.shape
```

```
Out[4]: (1000, 2)
```

```
In [5]:
```

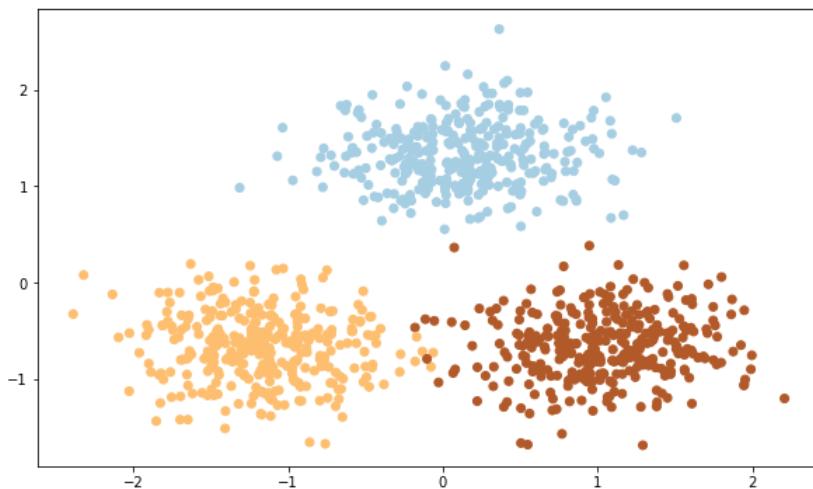
```
y.shape
```

```
Out[5]: (1000,)
```

Data Preprocessing and Visualization

```
In [6]: X = StandardScaler().fit_transform(X)
plt.figure(figsize=(10,6))
plt.scatter(X[:,0], X[:,1], c=y, cmap='Paired')
```

```
Out[6]: <matplotlib.collections.PathCollection at 0x1bc 149f3388>
```



Model Building

```
In [7]: db = DBSCAN(eps=0.45, min_samples=50)
db.fit(X)
y_pred = db.fit_predict(X)
plt.figure(figsize=(10,6))
plt.scatter(X[:,0], X[:,1], c=y_pred, cmap='Paired')
```

```
Out[7]: <matplotlib.collections.PathCollection at 0x1bc16dcf 608>
```

