

DAY 7 ASSIGNMENT

```
import pandas as pd

import numpy as np

import matplotlib.pyplot as plt

dataset1=pd.read_csv("general_data.csv")

dataset1.head()

print(dataset1.head())

print(dataset1.columns)

print(dataset1.isnull())

print(dataset1.duplicated())

print(dataset1.drop_duplicates())

dataset3=dataset1[['Age','DistanceFromHome','Education','MonthlyIncome','NumCompaniesWorked','PercentSalaryHike','TotalWorkingYears','TrainingTimesLastYear','YearsAtCompany','YearsSinceLastPromotion','YearsWithCurrManager']].describe()

print(dataset3)

dataset3=dataset1[['Age','DistanceFromHome','Education','MonthlyIncome','NumCompaniesWorked','PercentSalaryHike','TotalWorkingYears','TrainingTimesLastYear','YearsAtCompany','YearsSinceLastPromotion','YearsWithCurrManager']].median()

print(dataset3)

dataset3=dataset1[['Age','DistanceFromHome','Education','MonthlyIncome','NumCompaniesWorked','PercentSalaryHike','TotalWorkingYears','TrainingTimesLastYear','YearsAtCompany','YearsSinceLastPromotion','YearsWithCurrManager']].mode()

print(dataset3)

dataset3=dataset1[['Age','DistanceFromHome','Education','MonthlyIncome','NumCompaniesWorked','PercentSalaryHike','TotalWorkingYears','TrainingTimesLastYear','YearsAtCompany','YearsSinceLastPromotion','YearsWithCurrManager']].var()

print(dataset3)

dataset3=dataset1[['Age','DistanceFromHome','Education','MonthlyIncome','NumCompaniesWorked','PercentSalaryHike','TotalWorkingYears','TrainingTimesLastYear','YearsAtCompany','YearsSinceLastPromotion','YearsWithCurrManager']].skew()

print(dataset3)

dataset3=dataset1[['Age','DistanceFromHome','Education','MonthlyIncome','NumCompaniesWorked','PercentSalaryHike','TotalWorkingYears','TrainingTimesLastYear','YearsAtCompany','YearsSinceLastPromotion','YearsWithCurrManager']].kurt()

print(dataset3)

box_plot=dataset1.Age

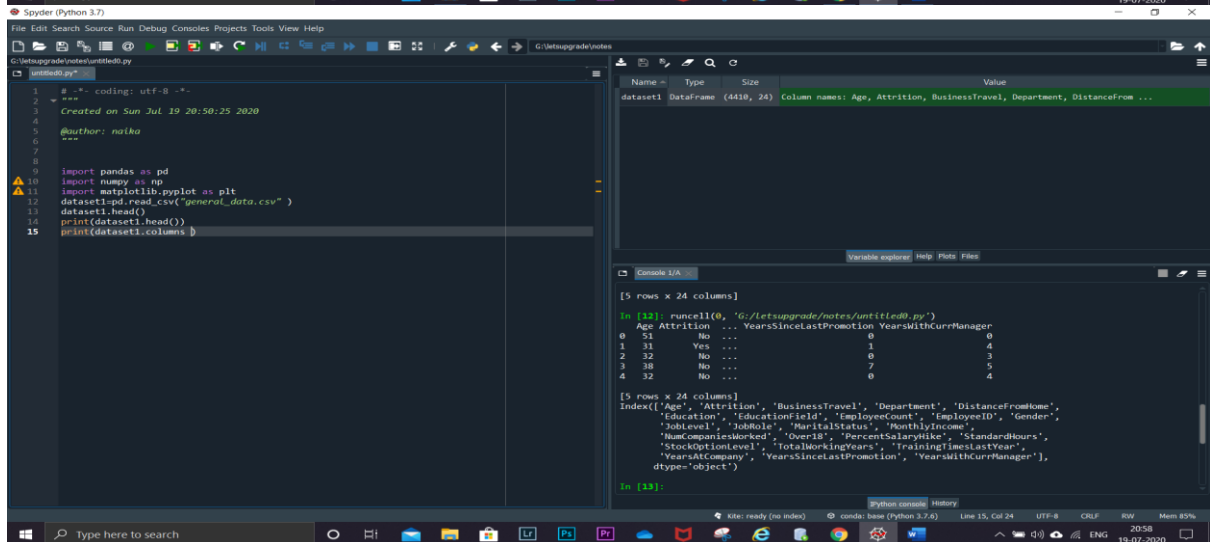
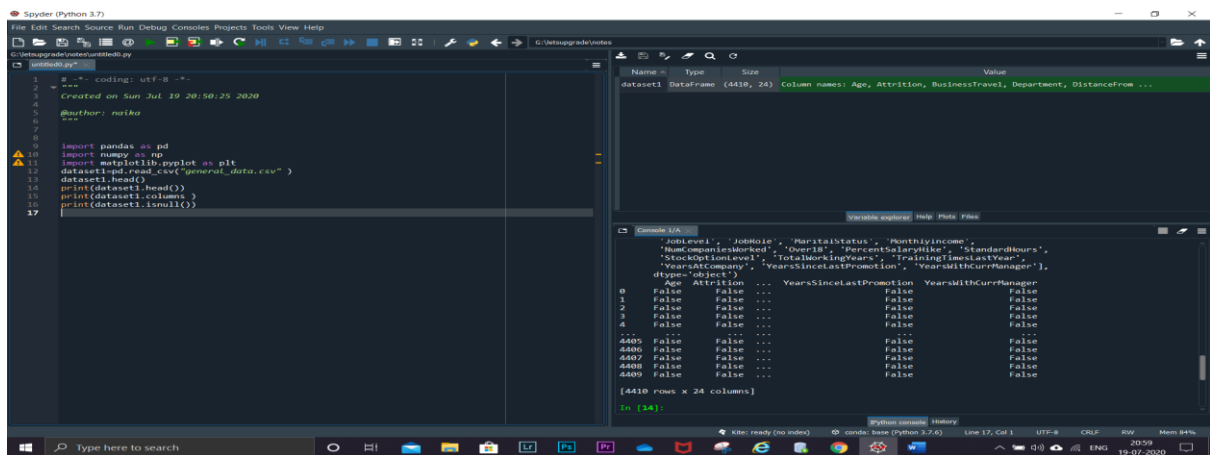
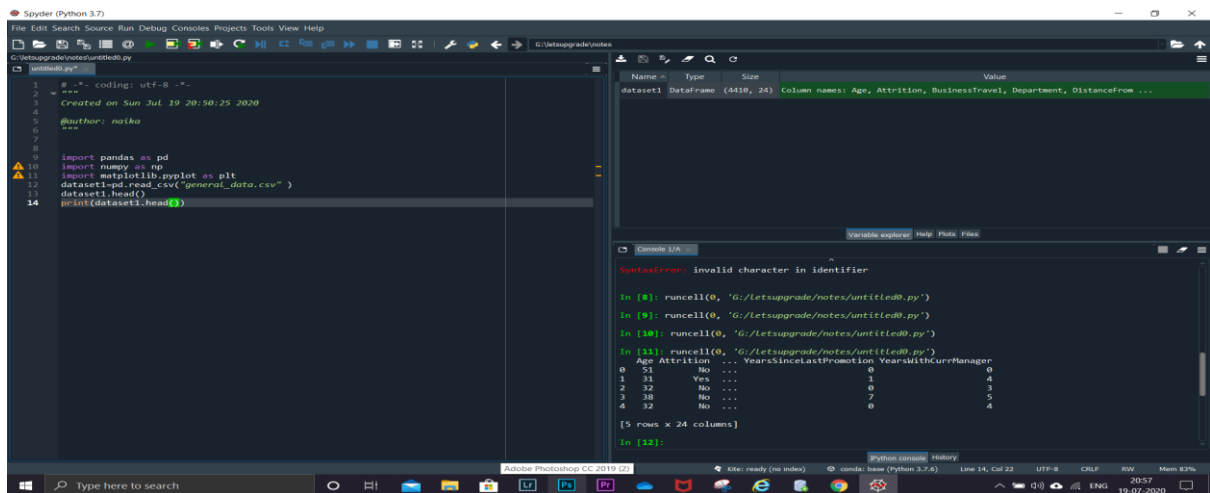
plt.boxplot(box_plot)

box_plot=dataset1.MonthlyIncome

plt.boxplot(box_plot)

box_plot=dataset1.YearsAtCompany

plt.boxplot(box_plot)
```



Spyder (Python 3.7)

```
1 # -*- coding: utf-8 -*-
2 """
3 Created on Sun Jul 19 20:50:25 2020
4
5 @author: naika
6 """
7
8
9
10 import pandas as pd
11 import numpy as np
12 import matplotlib.pyplot as plt
13 dataset1=pd.read_csv("general_data.csv")
14 dataset1.head()
15 print(dataset1.columns)
16 print(dataset1.isnull())
17 print(dataset1.duplicated())
```

Variable explorer

Name	Type	Size	Value
dataset1	DataFrame	(4410, 24)	Column names: Age, Attrition, BusinessTravel, Department, DistanceFrom...

Console I/O

```
4406 False
4407 False
4408 False
4409 False
Length: 4410, dtype: bool
```

Python console

```
In [16]:
```

Spyder (Python 3.7)

```
1 # -*- coding: utf-8 -*-
2 """
3 Created on Sun Jul 19 20:50:25 2020
4
5 @author: naika
6 """
7
8
9
10 import pandas as pd
11 import numpy as np
12 import matplotlib.pyplot as plt
13 dataset1=pd.read_csv("general_data.csv")
14 dataset1.head()
15 print(dataset1.columns)
16 print(dataset1.isnull())
17 print(dataset1.duplicated())
```

Variable explorer

Name	Type	Size	Value
dataset1	DataFrame	(4410, 24)	Column names: Age, Attrition, BusinessTravel, Department, DistanceFrom...

Console I/O

```
4406 False
4407 False
4408 False
4409 False
[4410 rows x 24 columns]
```

Python console

```
In [16]:
```

Spyder (Python 3.7)

```
1 # -*- coding: utf-8 -*-
2 """
3 Created on Sun Jul 19 20:50:25 2020
4
5 @author: naika
6 """
7
8
9
10 import pandas as pd
11 import numpy as np
12 import matplotlib.pyplot as plt
13 dataset1=pd.read_csv("general_data.csv")
14 dataset1.head()
15 print(dataset1.columns)
16 print(dataset1.isnull())
17 print(dataset1.duplicated())
18 dataset1=dataset1.drop_duplicates()
19 dataset3=dataset1[['Age', 'DistanceFromHome', 'Education', 'MonthlyIncome', 'NumCompaniesWorked']]
20 print(dataset3)
21 dataset3=dataset3[['Age', 'DistanceFromHome', 'Education', 'MonthlyIncome', 'NumCompaniesWorked']]
22 print(dataset3)
```

Variable explorer

Name	Type	Size	Value
dataset1	DataFrame	(4410, 24)	Column names: Age, Attrition, BusinessTravel, Department, DistanceFrom...
dataset3	Series	(11,)	Series object of pandas.core.series module

Console I/O

```
4407 25 No ... 1 2
4408 42 No ... 7 8
4409 40 No ... 3 9
[4410 rows x 24 columns]
count 4410.000000 ... 4410.000000
mean 36.923810 ... 4.123129
std 9.133301 ... 3.567127
min 18.000000 ... 0.000000
25% 30.000000 ... 2.000000
50% 36.000000 ... 3.000000
75% 43.000000 ... 7.000000
max 60.000000 ... 17.000000
[8 rows x 11 columns]
Age 36.0
DistanceFromHome 7.0
Education 3.0
MonthlyIncome 49190.0
NumCompaniesWorked 2.0
```

Python console

```
In [17]:
```

Spyder (Python 3.7)

```
1 # -*- coding: utf-8 -*-
2 """
3 Created on Sun Jul 19 20:50:25 2020
4
5 @author: naika
6 """
7
8
9
10 import pandas as pd
11 import numpy as np
12 import matplotlib.pyplot as plt
13 dataset1=pd.read_csv("general_data.csv")
14 dataset1.head()
15 print(dataset1.columns)
16 print(dataset1.isnull())
17 print(dataset1.duplicated())
18 dataset1=dataset1.drop_duplicates()
19 dataset3=dataset1[['Age', 'DistanceFromHome', 'Education', 'MonthlyIncome', 'NumCompaniesWorked']]
20 print(dataset3)
21 dataset3=dataset3[['Age', 'DistanceFromHome', 'Education', 'MonthlyIncome', 'NumCompaniesWorked']]
22 print(dataset3)
```

Variable explorer

Name	Type	Size	Value
dataset1	DataFrame	(4410, 24)	Column names: Age, Attrition, BusinessTravel, Department, DistanceFrom...
dataset3	DataFrame	(8, 11)	Column names: Age, DistanceFromHome, Education, MonthlyIncome, NumComp...

Console I/O

```
4406 42 No ... 0 2
4407 25 No ... 0 2
4408 42 No ... 1 2
4409 40 No ... 7 8
[4410 rows x 24 columns]
count 4410.000000 ... 4410.000000
mean 36.923810 ... 4.123129
std 9.133301 ... 3.567127
min 18.000000 ... 0.000000
25% 30.000000 ... 2.000000
50% 36.000000 ... 3.000000
75% 43.000000 ... 7.000000
max 60.000000 ... 17.000000
[8 rows x 11 columns]
```

Python console

```
In [17]:
```

Spyder (Python 3.7)

File Edit Search Source Run Debug Consoles Projects Tools View Help

G:\datascience\datascience.py

```
1 # -*- coding: utf-8 -*-
2 """
3 Created on Sun Jul 19 20:50:25 2020
4
5 @author: raika
6 """
7
8
9 import pandas as pd
10 import numpy as np
11 import matplotlib.pyplot as plt
12 dataset1=pd.read_csv("general_data.csv")
13 dataset1.head()
14 print(dataset1.head())
15 print(dataset1.columns)
16 print(dataset1.isnull())
17 print(dataset1.duplicated())
18 print(dataset1.drop_duplicates())
19 dataset3=dataset1[['Age', 'DistanceFromHome', 'Education', 'MonthlyIncome', 'NumCompaniesWorked',
20 print(dataset3)
21 dataset3=dataset1[['Age', 'DistanceFromHome', 'Education', 'MonthlyIncome', 'NumCompaniesWorked',
22 print(dataset3)
23 dataset3=dataset1[['Age', 'DistanceFromHome', 'Education', 'MonthlyIncome', 'NumCompaniesWorked',
24 print(dataset3)
```

Variable explorer

Name	Type	Size	Value
dataset1	DataFrame	(4410, 24)	Column names: Age, Attrition, BusinessTravel, Department, DistanceFromHome, ...
dataset3	DataFrame	(1, 11)	Column names: Age, DistanceFromHome, Education, MonthlyIncome, NumComp...

Console I/O

```
max 0.0 UNKNOWN ... 1.0 UNKNOWN
[8 rows x 11 columns]
Age 36.0
DistanceFromHome 7.0
Education 3.0
MonthlyIncome 49100.0
NumCompaniesWorked 2.0
PercentSalaryHike 14.0
TotalWorkingYears 18.0
TrainingTimesLastYear 3.0
YearsAtCompany 5.0
YearsSinceLastPromotion 1.0
YearsWithCurrManager 3.0
dtype: float64
0 35 DistanceFromHome 2 YearsSinceLastPromotion 0 YearsWithCurrManager 2
[1 rows x 11 columns]
In [18]:
```

Spyder (Python 3.7)

File Edit Search Source Run Debug Consoles Projects Tools View Help

G:\datascience\datascience.py

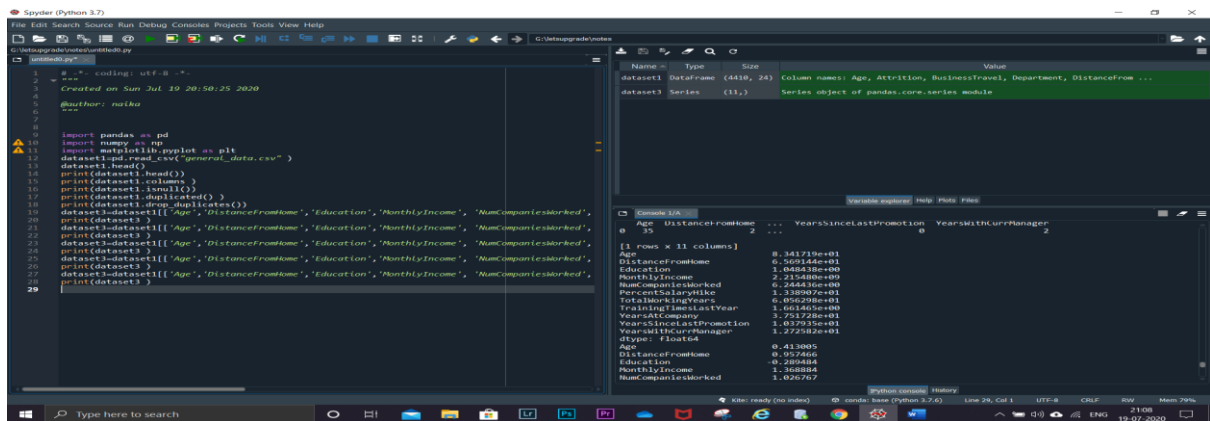
```
1 # -*- coding: utf-8 -*-
2 """
3 Created on Sun Jul 19 20:50:25 2020
4
5 @author: raika
6 """
7
8
9 import pandas as pd
10 import numpy as np
11 import matplotlib.pyplot as plt
12 dataset1=pd.read_csv("general_data.csv")
13 dataset1.head()
14 print(dataset1.head())
15 print(dataset1.columns)
16 print(dataset1.isnull())
17 print(dataset1.duplicated())
18 print(dataset1.drop_duplicates())
19 dataset3=dataset1[['Age', 'DistanceFromHome', 'Education', 'MonthlyIncome', 'NumCompaniesWorked',
20 print(dataset3)
21 dataset3=dataset1[['Age', 'DistanceFromHome', 'Education', 'MonthlyIncome', 'NumCompaniesWorked',
22 print(dataset3)
23 dataset3=dataset1[['Age', 'DistanceFromHome', 'Education', 'MonthlyIncome', 'NumCompaniesWorked',
24 print(dataset3)
25 dataset3=dataset1[['Age', 'DistanceFromHome', 'Education', 'MonthlyIncome', 'NumCompaniesWorked',
26 print(dataset3)
27 print(dataset3)
```

Variable explorer

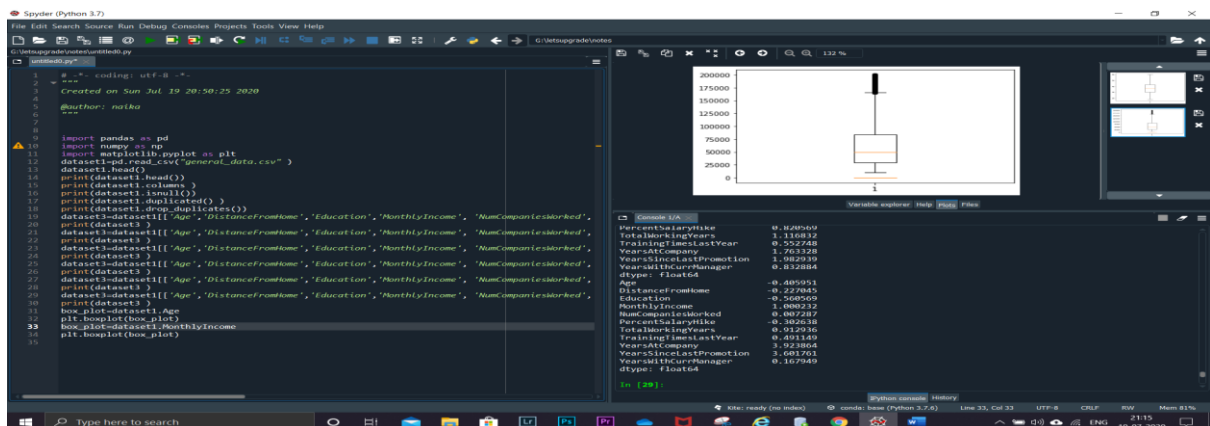
Name	Type	Size	Value
dataset1	DataFrame	(4410, 24)	Column names: Age, Attrition, BusinessTravel, Department, DistanceFromHome, ...
dataset3	Series	(11,)	Series object of pandas.core.series module

Console I/O

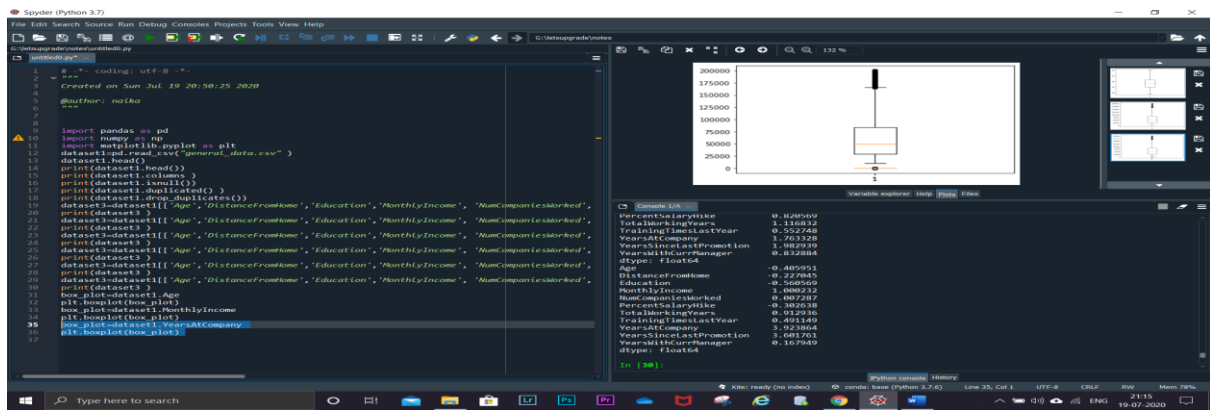
```
YearsSinceLastPromotion 1.0
YearsWithCurrManager 3.0
dtype: float64
0 35 DistanceFromHome 2 YearsSinceLastPromotion 0 YearsWithCurrManager 2
[1 rows x 11 columns]
Age 8.341719e+01
DistanceFromHome 6.565148e+01
Education 1.048438e+00
MonthlyIncome 2.215480e+09
NumCompaniesWorked 6.244430e+00
PercentSalaryHike 1.338007e+01
TotalWorkingYears 6.056298e+01
TrainingTimesLastYear 1.651465e+00
YearsAtCompany 1.751228e+01
YearsSinceLastPromotion 1.407935e+01
YearsWithCurrManager 1.225826e+01
dtype: float64
In [18]:
```



Age is normally distributed without any outliers



Monthly Income is Right skewed with several outliers



Years at company is also Right Skewed with several outliers observed.