ASSIGNMENT 1

Name: Aakanksha Bhondve

Rollno: 322004

Grno: 21810939

Comp B1

Aim -

Perform the following operations using Python on suitable data sets. Read a data from different formats(like csv,xls),indexing and selecting data, sort data, describe attributes of data, checking data types of each column, counting unique values of data, format of each column, converting variable data type (e.g. from long to short, vise versa), identifying missing values and fill in the missing values.

Theory -

Numpy: Numpy is Python library that provides mathematical function to handle large dimension array. It provides various method/function for Array, Metrics, and linear algebra.

Pandas: Pandas is one of the most popular Python library for data manipulation and analysis. Pandas provide useful functions to manipulate large amount of structured data. Pandas provide easiest method to perform analysis. It provide large data structures and manipulating numerical tables and time series data. Pandas is a perfect tool for data wrangling. Pandas is designed for quick and easy data manipulation, aggregation, and visualization. There two data structures in Pandas –

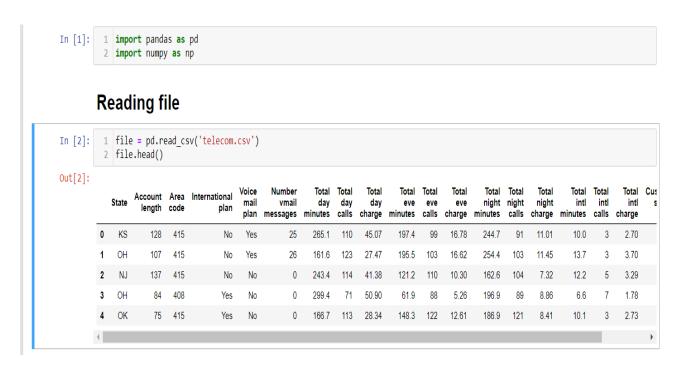
Series – It Handle and store data in one-dimensional data.

DataFrame – It Handle and store Two dimensional data.

Matplotlib: Matplolib is another useful Python library for Data Visualization. Descriptive analysis and visualizing data is very important for any organization. Matplotlib provides various method to Visualize data in more effective way. Matplotlib allows to quickly make line graphs,

pie charts, histograms, and other professional grade figures. Using Matplotlib, one can customize every aspect of a figure. Matplotlib has interactive features like zooming and planning and saving the Graph in graphics format.

Code and Output -



Indexing and Selecting Data



```
In [3]: 1 file.tail(5)
Out[3]:
                                                                                                                                                             Total
intl
                                                    Voice
                                                                         Total Total
                                                                                       Total
                                                                                                Total Total
                                                                                                              Total
                                                                                                                       Total Total
                                                                                                                                      Total
                                                                                                                                               Total Total
                                       International
                                Area
                 State
                                                                     day
minutes
                                                                                                                    night
minutes
                                                                                                                                    night
charge
                                                     mail
                                                               vmail
                                                                                        day
                                                                                                 eve
                                                                                                               eve
                                                                                                                                                intl
                         length code
                                              plan
                                                                                             minutes
                                                                               calls
                                                                                                             charge
                                                                                                                                                           charge
           3328
                                                                                               215.5
                                                                  36
                                                                         156.2
                                                                                                                       279.1
                                                                                                                                                 9.9
                                                                                                                                                              2.67
                            192
                                  415
                                                No
                                                                                       26.55
                                                                                                        126
                                                                                                              18.32
                                                                                                                                      12.56
           3329
                             68
                                  415
                                                No
                                                                   0
                                                                        231.1
                                                                                       39.29
                                                                                               153.4
                                                                                                              13.04
                                                                                                                       191.3
                                                                                                                               123
                                                                                                                                      8.61
                                                                                                                                                 9.6
                                                                                                                                                              2.59
                   RI
                             28
                                  510
                                                      No
                                                                   0
                                                                                       30.74
                                                                                               288.8
                                                                                                              24.55
                                                                                                                       191.9
                                                                                                                                                              3.81
           3330
                                                No
                                                                        180.8
                                                                                109
                                                                                                        58
                                                                                                                                91
                                                                                                                                      8.64
                                                                                                                                                14.1
           3331
                   CT
                            184
                                  510
                                                                   0
                                                                        213.8
                                                                                105
                                                                                       36.35
                                                                                                159.6
                                                                                                              13.57
                                                                                                                       139.2
                                                                                                                               137
                                                                                                                                      6.26
                                                                                                                                                 5.0
                                                                                                                                                              1.35
                                                                                                              22.60
                                                                                                                                                13.7
           3332
                   TN
                             74
                                  415
                                                No
                                                      Yes
                                                                  25
                                                                        234.4
                                                                                113
                                                                                      39.85
                                                                                               265.9
                                                                                                        82
                                                                                                                       241.4
                                                                                                                                77
                                                                                                                                      10.86
                                                                                                                                                              3.70
In [4]:
           1 file.iloc[1:5,0:2]
Out[4]:
              State Account length
           2
               N.J
                               137
               ОН
               OK
                                75
         Sorting Values
In [6]: 1 file['Total day minutes'].sort_values(ascending = False)
Out[6]: 365
                   350.8
          985
                   346.8
         2594
                   345.3
                   337.4
         156
         605
                   335.5
          1986
          2753
          2736
         1397
                     0.0
          1345
         Name: Total day minutes, Length: 3333, dtype: float64
In [8]: 1 file.sort_values(['Total day calls','Total day charge']).head()
Out[8]:
                                                   Voice
mail
                                                            Number
                                                                        Total Total
                                                                                      Total
                                                                                               Total Total
                                                                                                             Total
                                                                                                                     Total Total
                                                                                                                                    Total
                                                                                                                                             Total Total
                                                                                                                                                           Total
                       Account Area length code
                                      International
                                                                         day
                                                                             day
                                                                                                                                    night
                                                                                                    eve
calls
                                                                                                                     night
                                              plan
                                                                                    charge
                                                                                                           charge
                                                                                                                                                         charge
                                                          messages
                                                                    minutes
                                                                                                                                  charge
                                                                                                                                          minutes
                                                                                                                                                   calls
                                                    plan
                                                                                            minutes
          1345
                                                                                              159.6
                                                                                                                                                            1.84
                  SD
                            98
                                 415
                                               No
                                                      No
                                                                  0
                                                                         0.0
                                                                                 0
                                                                                      0.00
                                                                                                      130
                                                                                                            13.57
                                                                                                                     167.1
                                                                                                                                    7.52
                                                                                                                                              6.8
           1397
                           101
                                 510
                                                                         0.0
                                                                                 0
                                                                                      0.00
                                                                                              192.1
                                                                                                                     168.8
                                                                                                                                    7.60
                                                                                                                                              7.2
                                                                                                                                                            1.94
                           155
                                 408
                                               No
                                                      No
                                                                       216.7
                                                                                                                     135.3
                                                                                                                                    6.09
                                                                                                                                              10.8
                                                                                                                                                           2.92
```

144.5

185.8

35 24.57

36 31.59

30

262.3

276.5

101 22.30

134

23.50

82

10.19

8.64

226.5

192.1 104

12.0

5.7

3.24

1.54

МТ

NE

1989

692

415

408

No Yes

No

124

82

```
In [9]: 1 file.sort_values(['Total day calls','Total eve calls'], ascending=[True,False]).head()
Out[9]:
                                                                                                                                                            Total
intl
                                                                                                                             Total
                                                                                                                                     Total
                                      International
                                Area
code
                       Account
                 State
                                                                         day
                                                                              day
calls
                                                                                       day
                                                                                                eve
                                                                                                     eve
calls
                                                                                                                      night
                                                                                                                            night
calls
                                                                                                                                               intl
                         length
                                              plan
                                                                                                                                  charge
                                                    plan
                                                          messages
                                                                     minutes
                                                                                    charge
                                                                                            minutes
                                                                                                           charge
                                                                                                                   minutes
                                                                                                                                           minutes
                                                                                                                                                    calls
                                                                                                                                                          charge
          1345
                  SD
                            98
                                  415
                                               No
                                                      No
                                                                  0
                                                                          0.0
                                                                                 0
                                                                                      0.00
                                                                                               159.6
                                                                                                      130
                                                                                                             13.57
                                                                                                                      167.1
                                                                                                                               88
                                                                                                                                     7.52
                                                                                                                                               6.8
                                                                                                                                                             1.84
           1397
                  VT
                            101
                                  510
                                               No
                                                      No
                                                                  0
                                                                          0.0
                                                                                 0
                                                                                      0.00
                                                                                               192.1
                                                                                                      119
                                                                                                             16.33
                                                                                                                      168.8
                                                                                                                               95
                                                                                                                                     7.60
                                                                                                                                               7.2
                                                                                                                                                             1.94
           1144
                  NH
                           155
                                  408
                                               No
                                                                  0
                                                                        216.7
                                                                                30
                                                                                      36.84
                                                                                               144.3
                                                                                                      125
                                                                                                             12.27
                                                                                                                      135.3
                                                                                                                              106
                                                                                                                                     6.09
                                                                                                                                               10.8
                                                                                                                                                             2.92
           1989
                  МТ
                           124
                                415
                                               No
                                                     Yes
                                                                 30
                                                                        144.5
                                                                                35
                                                                                      24 57
                                                                                               262.3
                                                                                                      101
                                                                                                             22.30
                                                                                                                      226.5
                                                                                                                               82
                                                                                                                                    10.19
                                                                                                                                               12.0
                                                                                                                                                             3.24
           692
                            82
                                408
                                               No
                                                                  0
                                                                        185.8
                                                                                36
                                                                                      31.59
                                                                                              276.5
                                                                                                     134
                                                                                                            23.50
                                                                                                                      192.1
                                                                                                                              104
                                                                                                                                     8.64
                                                                                                                                               5.7
```

Describing the attributes of data

In [10]: 1 file.describe()

Out[10]:

	Account length	Area code	Number vmail messages	Total day minutes	Total day calls	Total day charge	Total eve minutes	Total eve calls	Total eve charge	Total night minutes	Total night calls	Tota
count	3333.000000	3333.000000	3333.000000	3333.000000	3333.000000	3333.000000	3333.000000	3333.000000	3333.000000	3333.000000	3333.000000	3333.0
mean	101.064806	437.182418	8.099010	179.775098	100.435644	30.562307	200.980348	100.114311	17.083540	200.872037	100.107711	9.0
std	39.822106	42.371290	13.688365	54.467389	20.069084	9.259435	50.713844	19.922625	4.310668	50.573847	19.568609	2.2
min	1.000000	408.000000	0.000000	0.000000	0.000000	0.000000	0.000000	0.000000	0.000000	23.200000	33.000000	1.0
25%	74.000000	408.000000	0.000000	143.700000	87.000000	24.430000	166.600000	87.000000	14.160000	167.000000	87.000000	7.
50%	101.000000	415.000000	0.000000	179.400000	101.000000	30.500000	201.400000	100.000000	17.120000	201.200000	100.000000	9.0
75%	127.000000	510.000000	20.000000	216.400000	114.000000	36.790000	235.300000	114.000000	20.000000	235.300000	113.000000	10.
max	243.000000	510.000000	51.000000	350.800000	165.000000	59.640000	363.700000	170.000000	30.910000	395.000000	175.000000	17.

In [13]: 1 file.shape
Out[13]: (3333, 20)

Checking datatypes of each column

RangeIndex: 3333 entries, 0 to 3332 Data columns (total 20 columns):

```
Column
                                   Non-Null Count Dtype
                                                       object
int64
     State
                                    3333 non-null
     Account length
                                    3333 non-null
     Area code
International plan
                                    3333 non-null
                                                       int64
object
                                    3333 non-null
     Voice mail plan
Number vmail messages
                                    3333 non-null
3333 non-null
                                                       object
int64
      Total day minutes
                                    3333 non-null
                                                        float64
      Total day calls
                                    3333 non-null
                                                        int64
      Total day charge
                                    3333 non-null
                                                        float64
      Total eve minutes
                                    3333 non-null
                                                        float64
     Total eve calls
Total eve charge
Total night minutes
                                    3333 non-null
                                                        int64
float64
 11
                                    3333 non-null
                                    3333 non-null
                                                        float64
     Total night calls
Total night charge
                                    3333 non-null
                                                        int64
                                    3333 non-null
                                                        float64
 15
     Total intl minutes
Total intl calls
                                    3333 non-null
                                                        float64
                                    3333 non-null
     Total intl charge
Customer service calls
                                                        float64
                                    3333 non-null
                                    3333 non-null
                                                        int64
 19 Churn
                                    3333 non-null
                                                       bool
dtypes: bool(1), float64(8), int64(8), object(3)
memory usage: 498.1+ KB
```

```
In [16]: 1 file.columns
dtype='object')
In [17]: 1 file.dtypes
Out[17]: State
                                      object
                                       int64
int64
          Account length
          Area code
          International plan
                                      object
          Voice mail plan
Number vmail messages
                                      object
                                       int64
          Total day minutes
                                     float64
          Total day calls
Total day charge
Total eve minutes
                                       int64
                                     float64
                                     float64
          Total eve calls
Total eve charge
                                       int64
                                     float64
          Total night minutes
                                     float64
          Total night calls
Total night charge
                                       int64
                                     float64
          Total intl minutes
                                     float64
          Total intl calls
Total intl charge
                                       int64
                                     float64
          Customer service calls
                                       int64
          Churn
                                        bool
          dtype: object
```

Checking and counting the unique values

```
In [18]: 1 file['State'].unique()
Out[18]: array(['KS', 'OH', 'NJ', 'OK', 'AL', 'MA', 'MO', 'LA', 'WV', 'IN', 'RI', 'IA', 'MT', 'NY', 'ID', 'VT', 'VA', 'TX', 'FL', 'CO', 'AZ', 'SC', 'NE', 'WY', 'HI', 'IL', 'NH', 'GA', 'AK', 'MD', 'AR', 'WI', 'OR', 'MI', 'DE', 'UT', 'CA', 'MN', 'SD', 'NC', 'WA', 'NM', 'NV', 'DC', 'KY', 'ME', 'MS', 'TN', 'PA', 'CT', 'ND'], dtype=object)
In [19]: 1 file['State'].nunique()
Out[19]: 51
In [20]: 1 file.nunique()
Out[20]: State
                                                    51
              Account length
                                                   212
             Area code
                                                     3
             International plan
                                                     2
             Voice mail plan
             Number vmail messages
                                                    46
             Total day minutes
                                                  1667
              Total day calls
                                                  119
             Total day charge
                                                  1667
              Total eve minutes
                                                 1611
              Total eve calls
                                                  123
              Total eve charge
                                                  1440
              Total night minutes
                                                  1591
             Total night calls
                                                   120
              Total night charge
                                                   933
             Total intl minutes
                                                   162
              Total intl calls
                                                    21
             Total intl charge
                                                   162
             Customer service calls
                                                    10
             Churn
                                                     2
             dtype: int64
```

Converting variable datatypes

```
In [20]: 1 file.convert_dtypes().dtypes
Out[20]: State
         Account length
                                      Int64
         Area code
                                      Int64
                                     string
         International plan
                                     string
         Voice mail plan
         Number vmail messages
                                      Int64
         Total day minutes
                                    float64
         Total day calls
                                      Int64
                                    float64
         Total day charge
         Total eve minutes
                                    float64
                                      Tnt64
         Total eve calls
                                    float64
         Total eve charge
         Total night minutes
                                    float64
         Total night calls
                                      Int64
         Total night charge
Total intl minutes
                                    float64
                                    float64
         Total intl calls
                                      Int64
         Total intl charge
                                    float64
         Customer service calls
                                      Int64
         Churn
                                    boolean
         dtype: object
In [21]: 1 file['Churn'] = file['Churn'].astype(int)
```

Checking null values and filling them

orders = pd.read_excel('NewOrders.xlsx')
orders.head() In [27]: Out[27]: Produc Cat Order Date Order Priority Order Quantity Unit Price Customer Name Customer Segment Product Category Sales Discount Province Region 3 2010-10-13 Muhammed Small Office Sto 6 261.5400 0.04 -213.250 38.94 Nunavut Nunavut Business Supplies Organi Sci Office Regular 6 2012- Not 02-20 Specified Ruben Dartt 6.9300 0.01 -4.640 2.08 West Corporate Rule Trir Alberta Supplies Regular 26 2808.0800 1054.820 107.53 Liz Pelletier High West Corporate Furniture Alberta Furni 24 1761.4000 0.09 -1748.560 70.89 Liz Pelletier Alberta West Corporate Furniture 32 2011-07-15 Regular Air 23 160.2335 -85.129 7.99 Liz Pelletier West Corporate Technology High Alberta Communi 5 rows × 21 columns

```
orders['Product Base Margin'].fillna( method ='ffill', inplace = True)
              2 orders.head()
Out[35]:
                Row Order Order
ID ID Date
                                       Order Order
Priority Quantity
                                                                                                        Unit
Price ...
                                                                                    Ship
Mode
                                                                                                                    Customer
Name
                                                                                                                                                    Customer
Segment
                                                                                                                                                                  Product 
Category
                                                                                                                                                                              Produc
Cat
                                                               Sales Discount
                                                                                               Profit
                                                                                                                               Province Region
                           3 2010-
10-13
                                                                                  Regular
Air
                                                                                                                   Muhammed
                                                                                                                                                         Small
                                                                                                                                                                     Office
                                                                                                                                                                                  Stor
                                                                            0.04
                                                                                             -213.250 38.94
                                                        6 261.5400
                                          Low
                                                                                                                                 Nunavut Nunavut
                                                                                                                                                      Business
                                                                                                                     MacIntyre
                                                                                                                                                                   Supplies
                                                                                                                                                                                Organi
                                                                                                                                                                                 Sci
Rule
Trir
                           6 2012- Not 
02-20 Specified
                                                                                  Regular
Air
                                                                                                                        Ruben
Dartt
                                                              6.9300
                                                                            0.01
                                                                                               -4.640
                                                                                                       2.08
                                                                                                                                  Alberta
                                                                                                                                              West Corporate
                                                                                                                                                                  Supplies
                          32 2011-
07-15
                                                                                  Regular
Air
                                                       26 2808.0800
                                                                                            1054.820 107.53 ... Liz Pelletier
                                          High
                                                                            0.07
                                                                                                                                  Alberta
                                                                                                                                              West Corporate
                                                                                                                                                                  Furniture
                                                                                                                                                                                Furnis
                          32 2011-
07-15
                                                                                  Delivery
Truck
                                                                                           -1748.560 70.89 ... Liz Pelletier
                                          High
                                                       24 1761.4000
                                                                            0.09
                                                                                                                                  Alberta
                                                                                                                                              West Corporate
                                                                                                                                                                  Furniture
                                                                                                                                                                                Telep
                          32 2011-
07-15
                                                                            0.04 Regular
Air
                                                       23 160.2335
                                          High
                                                                                              -85.129
                                                                                                       7.99 ... Liz Pelletier
                                                                                                                                  Alberta
                                                                                                                                              West Corporate Technology
                                                                                                                                                                             Communi
            5 rows × 21 columns
```

```
In [37]: 1 orders['Product Base Margin'].isnull()
Out[37]: 0
                 False
                 False
         1
         2
                 False
                 False
         3
         4
                 False
         8394
                 False
         8395
                 False
         8396
                 False
         8397
                 False
         8398
                False
         Name: Product Base Margin, Length: 8399, dtype: bool
```