

```
[In [1]] #to add a new row to existing array
import numpy as np
arr = np.array([1,2,3],[4,5,7],[9,6,10]])
print(arr)

Out[1]:
[[ 1  2  3]
 [ 4  5  7]
 [ 9  6 10]]

[In [2]] #frequency of all elements in series
import pandas as pd
s = pd.Series([1,1,1,2,2,2,2,3,3,3,4,4,5])
np.unique(s, return_counts=True)

Out[2]:
(array([1, 2, 3, 4, 5], dtype=int64), array([3, 4, 3, 2, 1], dtype=int64))

[In [3]] df=pd.DataFrame(np.random.randn(5,4), ['A','B','C','D'], ['w','x','y','z'])
print(df) #df['w'] #>-> to access specific column
df.loc[0,df['w']] #>-> to access specific row

Out[3]:
      w         x         y         z
A  1.467851  1.358444  0.284435  0.476161
B  0.426812  0.165824  0.262149  -1.478847
C -1.092043  0.492741  1.657188  0.222142
D -1.628623  1.435599  0.339862  0.988139
E  0.155973  0.807121  0.760574  0.281043

Name: w, dtype: float64

w      1.467851
x      1.358444
y      0.284435
z      0.476161
dtype: float64

[In [4]] #to access multiple columns
df.loc[['A','C']]

Out[4]:
      w         x         y         z
A  1.467851  1.358444  0.284435  0.476161
C -1.092043  0.492741  1.657188  0.222142

[In [5]] #to access rows based on index position
df.iloc[3] #>-> 3 is index value

Out[5]:
w      -1.628623
x      1.435599
y      0.339862
z      0.988139
dtype: float64

[In [6]] #to get specific value of co-ordinate in dataframe
df.loc['B','y'] #Syntax -> df.loc[row,column]

Out[6]:
0.26213988873531284

[In [7]] #to get multiple co-ordinates
df.loc[['A','B'],['w','z']] #df.loc[[rows],[columns]]

Out[7]:
      w         z
A  1.467851  0.476161
B  0.426812 -1.478847

[In [8]] #retrieving the data based on condition
print(df[df['w']>0]) #then for false
print(df[df['w']<0])

Out[8]:
      w         x         y         z
A  1.467851  1.358444  0.284435  0.476161
B  0.426812  0.165824  0.262149   NaN
C   NaN         NaN         NaN  0.222142
D   NaN         NaN  0.339862  0.988139
E  0.155973   NaN  0.760574  0.281043

      w         x         y         z
A  1.467851  1.358444  0.284435  0.476161
B  0.426812  0.165824  0.262149  -1.478847
E  0.155973  0.807121  0.760574  0.281043

[In [9]] #to convert dictionary to dataframe
df2 = {'id': [1,2,3,4,5], "id": [5,np.nan,np.nan], "C": [1,2,3], "D": [np.nan,np.nan,np.nan]}
df1 = pd.DataFrame(df2)

Out[9]:
   id  A  B  C  D
0   1  0  5  1  NaN
1   2  0  NaN  2  NaN
2   3  NaN  NaN  3  NaN

[In [10]] #to drop values with Nan
df.dropna(thresh,axis,inplace)
print(df.dropna())
#show'all' or 'any'
print("\n",df.dropna(how='all')) #it will be checking the row with all
print("\n",df.dropna(how='any')) #it will delete the rows which have single Nan value

Empty DataFrame
Columns: [A, B, C, D]
Index: []

      A  B  C  D
0  1.0  5.0  1  NaN
1  2.0  NaN  2  NaN
2  NaN  NaN  3  NaN

Empty DataFrame
Columns: [A, B, C, D]
Index: []

[In [11]] df1.dropna(how='all',axis=1)

Out[11]:
   A  B  C
0  1.0  5.0  1
1  2.0  NaN  2
2  NaN  NaN  3

[In [12]] #to check the sum of Nan values in each column
df1.isna().sum()

Out[12]:
A    1
B    2
C    0
D    0
dtype: int64

[In [13]] #thresh
print(df1.dropna(thresh=2,axis=1))

Out[13]:
   A  C
0  1.0  1.0
1  2.0  2
2  NaN  3

[In [14]] #fillna()
df1.fillna(value=0.0)

Out[14]:
   A  B  C  D
0  1.0  5.0  1  0.0
1  2.0  0.0  2  0.0
2  0.0  0.0  3  0.0

[In [15]] df1['A'].fillna(value=df1['A'].mean())

Out[15]:
0    1.0
1    2.6
2    1.5
dtype: float64

[In [16]] # Group by
#to group based on column and perform aggregate functions
data = {'company': ['Google','Google','Meta','Meta','Fb','Fb'],
        'person': ['sam','mani','sai','prasanth','mohan','mounika'],
        'sales': [500, 580, 386, 298, 600, 20]}
df2 = pd.DataFrame(data)

Out[16]:
   company person sales
0  Google    sam    300
1  Google    mani    500
2  Meta     sai    120
3  Meta  prasanth    200
4  Fb      mohan    600
5  Fb     mounika    20

[In [17]] bycomp = df2.groupby('company')
bycomp
bycomp.sum()[['sales']]
bycomp.value_counts()
bycomp['company'].value_counts()

Out[17]:
company
Fb      2
Google  2
Meta   2
Name: count, dtype: int64

[In [18]] #to get the maximum salary person from each company
df2.loc[bycomp['sales'].idxmax()][['person','sales']]

Out[18]:
   person sales
0  mohan    600
1  prasnth   200

[In [19]] bycomp.describe()

Out[19]:
              count  mean      std  min    25%    50%    75%    max
company
Fb           2.0    310.0  410.12193    20.0  165.0  310.0  450.0  600.0
Google       2.0  400.0  141.42186   300.0  350.0  400.0  450.0  500.0
Meta         2.0  150.0   70.71078   100.0  125.0  150.0  175.0  200.0

[In [20]] #to access a csv file
df = pd.read_csv("samplecsv.csv")
df.head() #to get first 5 rows

Out[20]:
   Year  Industry aggregation  NZSIOC  Industry code  NZSIOC  Industry name  NZSIOC  Units  Variable code  Variable name  Variable category  Value  Industry code  ANZSIC06
0  2020  Government           Level1  99999          All industries  Dollars (millions)  HQ1  Total income  Financial performance  733.258  ANZSIC06 divisions A-5 (excluding classes K633...
1  2020  Government           Level1  99999          All industries  Dollars (millions)  HQ4  Sales, government funding, grants and subsidies  Financial performance  660.630  ANZSIC06 divisions A-5 (excluding classes K633...
2  2020  Government           Level1  99999          All industries  Dollars (millions)  HQ5  Interest, dividends and donations  Financial performance  64.342  ANZSIC06 divisions A-5 (excluding classes K633...
3  2020  Government           Level1  99999          All industries  Dollars (millions)  HQ7  Non-operating income  Financial performance  1.285  ANZSIC06 divisions A-5 (excluding classes K633...
4  2020  Government           Level1  99999          All industries  Dollars (millions)  HQ8  Total expenditure  Financial performance  644.872  ANZSIC06 divisions A-5 (excluding classes K633...

[In [21]] #to access a xlsx file
dfl = pd.ExcelFile("sample.xlsx")
dfl.sheet_names()
dfl.parse(0) #to get last five rows

Out[21]:
   Segment  Country  Product  Discount  Band  Units Sold  Manufacturing Price  Sale Price  Gross Sales  Discounts  Sales  COGS  Profit  Date  Month Number  Month Name  Year
695  Small Business  France  Amarilla  High  2475.0  260  300  742500.0  11375.00  631125.00  618750.0  12375.00  2014-03-01  3  March  2014
696  Small Business  Mexico  Amarilla  High  545.0  260  300  163500.0  24570.00  139230.00  136500.0  2730.00  2014-02-01  10  October  2014
697  Government  Mexico  Montana  High  1368.0  5  7  9576.0  1436.40  8139.60  6840.0  1296.00  2014-05-01  2  February  2014
698  Government  Canada  Paseo  High  723.0  10  7  5061.0  759.15  4301.85  3615.0  686.85  2014-04-01  4  April  2014
699  Channel Partners  United States of America  VIT  High  1806.0  250  12  21672.0  3250.80  18421.20  5418.0  13003.20  2014-05-01  5  May  2014

[In [22]] #to access a csv file
df = pd.read_csv("Bellintercs.csv", sep=';', names=['sid','ages','places'])
df.drop(0) #to delete a row

Out[22]:
   sid  ages  places
1  Alice   25  New York
2  Bob    30  San Francisco
3  Charlie  22  Los Angeles

[In [23]] #to access a csv file
df_titanic = pd.read_csv("titanic_train.csv")
df_titanic.columns
df_titanic.head()
df_titanic['Pclass'].unique()
df_titanic['Pclass'].value_counts()

Out[23]:
Pclass
3    491
1    215
2    184
Name: count, dtype: int64

[In [24]] df_titanic.info()

Out[24]:
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 891 entries, 0 to 890
Data columns (total 12 columns):
 #   Column                Non-Null Count  Dtype
---  --
 0   PassengerId           891 non-null    int64
 1   Survived              891 non-null    int64
 2   Pclass                891 non-null    int64
 3   Name                  891 non-null    object
 4   Sex                   891 non-null    object
 5   Age                   714 non-null    float64
 6   SibSp                 891 non-null    int64
 7   Parch                891 non-null    int64
 8   Ticket                891 non-null    object
 9   Fare                  891 non-null    float64
10   Cabin                204 non-null    object
11   Embarked              889 non-null    object
dtypes: float64(2), int64(5), object(5)
memory usage: 81.7+ KB

[In [25]] df_titanic.describe()

Out[25]:
   PassengerId  Survived  Pclass  Age  SibSp  Parch  Fare
mean    166.009258    0.810000    0.910000    29.699118    0.520000    0.381994    32.206268
std     295.353842    0.406982    0.836071    14.570487    1.102743    0.809077    49.693428
min         0.000000    0.000000    1.000000    0.430000    0.000000    0.000000    0.000000
25%    223.500000    0.000000    2.000000    20.125000    0.000000    0.000000    7.914040
50%    446.000000    0.000000    3.000000    28.000000    0.000000    0.000000    14.454200
75%    668.500000    1.000000    3.000000    38.000000    1.000000    0.000000    31.000000
max    851.000000    1.000000    3.000000    80.000000    8.000000    6.000000    512.329200

[In [26]] #sum of Nan's
df_titanic.isna().sum()

Out[26]:
PassengerId    0
Survived        0
Pclass          0
Name            0
Sex             0
Age            177
SibSp           0
Parch           0
Ticket          0
Fare            0
Cabin          687
Embarked        2
dtype: int64

[In [27]] df_titanic.isna().sum()['109:len(df_titanic)']

Out[27]:
PassengerId    0.000000
Survived        0.000000
Pclass          0.000000
Name            0.000000
Sex             0.000000
Age            19.663209
SibSp           0.000000
Parch           0.000000
Ticket          0.000000
Fare            0.000000
Cabin          77.184377
Embarked        0.224467
dtype: float64

[In [28]] #to drop a column
df_titanic.drop(['Cabin'],axis=1)

Out[28]:
   PassengerId  Survived  Pclass  Name  Sex  Age  SibSp  Parch  Ticket  Fare  Embarked
0            0         1         3  Braund, Mr. Owen Harris  male  22.0  1    0  A/5 21171  7.2500  S
1            1         0         3  Cumings, Mrs. John Bradley (Florence Briggs Th... female  38.0  1    0  PC 17599  71.2833  C
2            2         1         3  Heikinen, Miss. Laina  female  26.0  0    0  STON/O2 3101282  7.9250  S
3            3         1         1  Futrelle, Mrs. Jacques Heath (Lily May Peel) female  35.0  1    0  113803  53.1000  S
4            4         0         3  Allen, Mr. William Henry  male  35.0  0    0  373450  8.0500  S
...         ...         ...         ...  ...  ...  ...  ...  ...  ...  ...  ...
886         887         0         2  Monrville, Rev. Juzas  male  27.0  0    0  211536  13.0000  S
887         888         1         1  Graham, Miss. Margaret Edith  female  13.0  0    0  112053  30.0000  S
888         889         0         3  Johnston, Miss. Catherine Helen "Cater" female  18.0  1    2  W/C 667  23.4000  S
889         890         1         1  Behr, Mr. Karl Howell  male  26.0  0    0  111369  30.0000  C
890         891         0         3  Dooley, Mr. Patrick  male  32.0  0    0  270376  7.7500  Q
891 rows x 11 columns

[In [29]] #to find number of people dead
len(df_titanic[df_titanic['Survived']==0].sum())

Out[29]:
549

[In [30]] #matplotlib
#it is mostly used for data visualization library, whi was inspired by MATLAB
#it is used for data visualization in the form of various plots
import matplotlib.pyplot as plt

[In [31]] x=[1,2,3,4]
y=[1,4,2,5]
plt.figure()
plt.plot(x,y)

Out[31]:
[<matplotlib.lines.Line2D at 0x1d38cf76d390>]

5.0
4.5
4.0
3.5
3.0
2.5
2.0
1.5
1.0
1.0 1.5 2.0 2.5 3.0 3.5 4.0

[In [32]] plt.subplot(1,2,1)
plt.plot(x,y)
plt.subplot(1,2,2)
plt.plot(y,x,'r')

Out[32]:
[<matplotlib.lines.Line2D at 0x1d38cf76d990>]

5.0
4.5
4.0
3.5
3.0
2.5
2.0
1.5
1.0
1.0 1.5 2.0 2.5 3.0 3.5 4.0

5.0
4.5
4.0
3.5
3.0
2.5
2.0
1.5
1.0
1.0 1.5 2.0 2.5 3.0 3.5 4.0

[In [33]] plt.subplot(2,1,1)
plt.plot(x,y)
plt.subplot(2,1,2)
plt.plot(y,x,'r')

Out[33]:
[<matplotlib.lines.Line2D at 0x1d38cf76d990>]

5
4
3
2
1
1.0 1.5 2.0 2.5 3.0 3.5 4.0

4
3
2
1
1.0 1.5 2.0 2.5 3.0 3.5 4.0 4.5 5.0

[In [34]] plt.plot(x,y,'v')

Out[34]:
[<matplotlib.lines.Line2D at 0x1d38cf77f9d0>]

5.0
4.5
4.0
3.5
3.0
2.5
2.0
1.5
1.0
1.0 1.5 2.0 2.5 3.0 3.5 4.0

[In [35]] years=np.arange(2000,2005,dtype=int)
y_innova = np.array([20,30,45,19,23])
years,y_innova
plt.plot(years,y_innova,'-o')

Out[35]:
[<matplotlib.lines.Line2D at 0x1d38cf78f9d0>]

30
28
26
24
22
20
18
16
14
2000.0 2000.5 2001.0 2001.5 2002.0 2002.5 2003.0 2003.5 2004.0

[In [36]] years=np.arange(2000,2005,dtype=int)
y_innova = np.array([20,30,45,19,23])
years,y_innova
plt.plot(years,y_innova,'-g')
plt.plot(years,honda,'-r')
plt.title('Car statistics')
plt.xlabel('years')
plt.ylabel('Sales')

Out[36]:
<matplotlib.legend.Legend at 0x1d38cf79dca0>

Car statistics
years
Innova
City
2000.0 2000.5 2001.0 2001.5 2002.0 2002.5 2003.0 2003.5 2004.0

[In [37]] products=['TV','Laptop','Smart phone','Sid shoes','Shirt']
sales=np.random.randint(10,200,size=len(products))
plt.bar(products,sales,color=['violet','blue','green','purple','orange'],edgecolor='yellow')
plt.title('Harika statistics')
plt.xlabel('Products')
plt.ylabel('Sales')

Out[37]:
TV Laptop Smart phone Sid shoes Shirt

[In [40]] products=['TV','Laptop','Smart phone','Sid shoes','Shirt']
sales=np.random.randint(10,200,size=len(products))
plt.bar(products,sales,color=['violet','blue','green','purple','orange'],edgecolor='black',height=0.5)
plt.title('Harika statistics')
plt.xlabel('Products')
plt.ylabel('Sales')

Out[40]:
TV Laptop Smart phone Sid shoes Shirt

[In [41]] names=['Ram','Shyam','Ravi','Satya','Radhey']
slices=[0.85,0.05,0.05,0.05,0.05]
scores=[10,25,40,55,70]
plt.bar(scores,slices,labels=names,startangle=90,shadow=True)

Out[41]:
Ram
Shyam
Ravi
Satya
Radhey

[In [42]] genders=['Male','Female']
male = (df_titanic['Sex']=='male').sum()
female = (df_titanic['Sex']=='female').sum()
count=[male,female]
plt.bar(gender,count,color='blue',width=0.3)

Out[42]:
<BarContainer object of 2 artists>

600
500
400
300
200
100
0
Male Female

[In [43]] genders=['Male','Female']
male = (df_titanic['Sex']=='male') & df_titanic['Survived']
female = (df_titanic['Sex']=='female') & df_titanic['Survived']
count=[male,female]
plt.bar(gender,count,color='blue',width=0.3)

Out[43]:
<BarContainer object of 2 artists>

200
150
100
50
0
Male Female

[In [44]] np.random.seed(0)
arr=np.random.randn(168,15,321)
arr2=np.random.randn(208,25,560)
color=plt.cm.rpy
fig=plt.figure()
axes=fig.add_axes([0,0,1,1])
vmaxes=violinplot(c0)

Out[44]:
275
250
225
200
175
150
125
0.8 1.0 1.2 1.4 1.6 1.8 2.0 2.2
```