

[1]: import numpy as np
import pandas as pd
import matplotlib.pyplot as plt
import seaborn as sns

In [2]: df = pd.read_csv('Job_Placement_Data.csv')

Out[2]:

	gender	ssc_percentage	ssc_board	hsc_percentage	hsc_board	hsc_subject	degree_percentage	undergrad_degree	work_experience	emp_test_percentage	specialisation	mba_percent	status
0	M	67.00	Others	91.00	Others	Commerce	58.00	Sci&Tech	No	55.0	Mkt&HR	58.80	Placed
1	M	79.33	Central	78.33	Others	Science	77.48	Sci&Tech	Yes	86.5	Mkt&Fin	66.28	Placed
2	M	65.00	Central	68.00	Central	Arts	64.00	Comm&Mgmt	No	75.0	Mkt&Fin	57.80	Placed
3	M	85.80	Central	73.60	Central	Commerce	52.00	Sci&Tech	No	66.0	Mkt&HR	59.43	Not Placed
4	M	56.00	Central	52.00	Central	Commerce	73.30	Comm&Mgmt	No	96.8	Mkt&Fin	55.50	Placed
...
210	M	80.60	Others	82.00	Others	Commerce	77.60	Comm&Mgmt	No	91.0	Mkt&Fin	74.49	Placed
211	M	58.00	Others	60.00	Others	Science	72.00	Sci&Tech	No	74.0	Mkt&Fin	53.62	Placed
212	M	67.00	Others	67.00	Others	Commerce	73.00	Comm&Mgmt	Yes	59.0	Mkt&Fin	69.72	Placed
213	F	74.00	Others	66.00	Others	Commerce	58.00	Comm&Mgmt	No	70.0	Mkt&HR	60.23	Placed
214	M	62.00	Central	58.00	Others	Science	53.00	Comm&Mgmt	No	89.0	Mkt&HR	60.22	Not Placed

215 rows × 13 columns

In [3]: df.info()

```
<class 'pandas.core.frame.DataFrame'>  
RangeIndex: 215 entries, 0 to 214  
Data columns (total 13 columns):  
#   column                Non-Null Count  Dtype    
--  --  --  --  --  --  --  --  --  --  --  --  --  --    
0   gender                215 non-null    object   
1   ssc_percentage        215 non-null    float64  
2   ssc_board             215 non-null    object   
3   hsc_percentage        215 non-null    float64  
4   hsc_board             215 non-null    object   
5   hsc_subject           215 non-null    object   
6   degree_percentage     215 non-null    float64  
7   undergrad_degree      215 non-null    object   
8   work_experience        215 non-null    object   
9   emp_test_percentage   215 non-null    float64  
10  specialisation         215 non-null    object   
11  mba_percent           215 non-null    float64  
12  status                215 non-null    object   
dtypes: float64(8), object(8)  
memory usage: 22.0+ KB
```

In [4]: df.describe()

	ssc_percentage	hsc_percentage	degree_percentage	emp_test_percentage	mba_percent
count	215.000000	215.000000	215.000000	215.000000	215.000000
mean	67.303395	66.331163	66.370186	72.100558	62.278186
std	10.827205	10.897509	7.358743	13.275956	5.833385
min	40.800000	37.000000	50.000000	50.000000	51.210000
25%	60.600000	60.900000	61.000000	60.000000	57.945000
50%	67.000000	65.000000	66.000000	71.000000	62.000000
75%	75.700000	73.000000	72.000000	83.500000	66.250000
max	89.400000	97.700000	91.000000	98.000000	77.800000

In [5]: # calculate rows and columns
df.shape

Out[5]: (215, 13)

In [6]: # to calculate top 5 entries
df.head()

Out[6]:

	gender	ssc_percentage	ssc_board	hsc_percentage	hsc_board	hsc_subject	degree_percentage	undergrad_degree	work_experience	emp_test_percentage	specialisation	mba_percent	status
0	M	67.00	Others	91.00	Others	Commerce	58.00	Sci&Tech	No	55.0	Mkt&HR	58.80	Placed
1	M	79.33	Central	78.33	Others	Science	77.48	Sci&Tech	Yes	86.5	Mkt&Fin	66.28	Placed
2	M	65.00	Central	68.00	Central	Arts	64.00	Comm&Mgmt	No	75.0	Mkt&Fin	57.80	Placed
3	M	85.80	Central	52.00	Central	Science	52.00	Sci&Tech	No	66.0	Mkt&HR	59.43	Not Placed
4	M	56.00	Central	73.60	Central	Commerce	73.30	Comm&Mgmt	No	96.8	Mkt&Fin	55.50	Placed

In [7]: # to calculate bottom 5 entries
df.tail()

Out[7]:

	gender	ssc_percentage	ssc_board	hsc_percentage	hsc_board	hsc_subject	degree_percentage	undergrad_degree	work_experience	emp_test_percentage	specialisation	mba_percent	status
210	M	80.60	Others	82.00	Others	Commerce	77.6	Comm&Mgmt	No	91.0	Mkt&Fin	74.49	Placed
211	M	58.0	Others	60.0	Others	Science	72.0	Sci&Tech	No	74.0	Mkt&Fin	53.62	Placed
212	M	67.00	Others	67.0	Others	Commerce	73.0	Comm&Mgmt	Yes	59.0	Mkt&Fin	69.72	Placed
213	F	74.0	Others	66.0	Others	Commerce	58.0	Comm&Mgmt	No	70.0	Mkt&HR	60.23	Placed
214	M	62.0	Central	58.0	Others	Science	53.0	Comm&Mgmt	No	89.0	Mkt&HR	60.22	Not Placed

In above data we have 215 candidates job placement details

In [8]: # to find missing data from dataset
df.isna().sum()

Out[8]:

gender	0
ssc_percentage	0
ssc_board	0
hsc_percentage	0
hsc_board	0
hsc_subject	0
degree_percentage	0
undergrad_degree	0
work_experience	0
emp_test_percentage	0
specialisation	0
mba_percent	0
status	0
dtype: int64	

In [9]: # how many Male and Female candidate
df['gender'].value_counts()

Out[9]:

```
M    139  
F     76  
Name: gender, dtype: int64
```

Above data have 139 male candidate and 76 female candidate

In [10]: # columns details
df.columns

Out[10]: Index(['gender', 'ssc_percentage', 'ssc_board', 'hsc_percentage', 'hsc_board', 'hsc_subject', 'degree_percentage', 'undergrad_degree', 'work_experience', 'emp_test_percentage', 'specialisation', 'mba_percent', 'status'], dtype='object')

In [11]: df.dtypes

Out[11]:

gender	object
ssc_percentage	float64
ssc_board	object
hsc_percentage	float64
hsc_board	object
hsc_subject	object
degree_percentage	float64
undergrad_degree	object
work_experience	object
emp_test_percentage	float64
specialisation	object
mba_percent	float64
status	object
dtype: object	

In [12]: # find outlier in data
df.boxplot(figsize=(10,6))
plt.show()

100

90

80

70

60

50

40

ssc_percentage hsc_percentage degree_percentage emp_test_percentage mba_percent

Above boxplot shows outlier in hsc_percentage and degree percentage

In [13]: a = df[df.hsc_percentage >85.0]

Out[13]:

	gender	ssc_percentage	ssc_board	hsc_percentage	hsc_board	hsc_subject	degree_percentage	undergrad_degree	work_experience	emp_test_percentage	specialisation	mba_percent	status
0	M	67.00	Others	91.00	Others	Commerce	58.00	Sci&Tech	No	55.00	Mkt&HR	58.80	Placed
13	F	77.00	Central	87.00	Central	Commerce	59.00	Comm&Mgmt	No	86.00	Mkt&Fin	68.63	Placed
24	M	76.50	Others	97.70	Others	Science	79.86	Sci&Tech	No	97.40	Mkt&Fin	66.28	Placed
42	M	87.00	Others	87.00	Others	Commerce	68.00	Sci&Tech	No	95.90	Mkt&HR	62.90	Placed
78	M	84.00	Others	90.90	Others	Science	64.50	Sci&Tech	No	86.04	Mkt&Fin	59.42	Placed
85	F	83.84	Others	89.83	Others	Commerce	72.20	Sci&Tech	Yes	78.74	Mkt&Fin	76.18	Placed
90	F	85.00	Others	90.00	Others	Commerce	82.00	Comm&Mgmt	No	92.00	Mkt&Fin	68.03	Placed
107	M	82.00	Others	90.00	Others	Commerce	83.00	Comm&Mgmt	No	80.00	Mkt&HR	73.52	Placed
129	M	76.70	Central	89.70	Others	Commerce	66.00	Comm&Mgmt	Yes	94.00	Mkt&Fin	68.55	Placed
134	F	77.44	Central	92.00	Others	Commerce	72.00	Comm&Mgmt	Yes	90.00	Mkt&Fin	67.13	Placed
148	F	77.00	Central	86.00	Central	Arts	56.00	Others	No	57.00	Mkt&Fin	55.47	Placed
162	M	74.20	Central	87.60	Others	Commerce	77.25	Comm&Mgmt	Yes	75.20	Mkt&Fin	66.06	Placed
177	F	73.00	Central	97.00	Others	Commerce	79.00	Comm&Mgmt	Yes	89.00	Mkt&Fin	70.81	Placed

In [14]: # change value of outlier
medianoff['hsc_percentage'].median()

Out[14]: 65.0

In [15]: df.loc[df.hsc_percentage>90, 'hsc_percentage'] = np.nan

In [16]: df.fillna(median,inplace=True)

outlier value is replaced by median

In [17]: a = df[df['degree_percentage']>90]

Out[17]:

	gender	ssc_percentage	ssc_board	hsc_percentage	hsc_board	hsc_subject	degree_percentage	undergrad_degree	work_experience	emp_test_percentage	specialisation	mba_percent	status
197	F	83.86	Others	53.0	Others	Science	91.0	Sci&Tech	No	59.32	Mkt&HR	69.71	Placed

In [18]: df.boxplot(figsize=(10,6))
plt.show()

100

90

80

70

60

50

40

ssc_percentage hsc_percentage degree_percentage emp_test_percentage mba_percent

NO more outlier present in data

In [19]: # how many male candidate placed
a = df.groupby(['gender', 'status']).get_group(('M','Placed'))
a

Out[21]:

	gender	ssc_percentage	ssc_board	hsc_percentage	hsc_board	hsc_subject	degree_percentage	undergrad_degree	work_experience	emp_test_percentage	specialisation	mba_percent	status
0	M	67.00	Others	65.00	Others	Commerce	58.00	Sci&Tech	No	55.00	Mkt&HR	58.80	Placed
1	M	79.33	Central	78.33	Others	Science	77.48	Sci&Tech	Yes	86.50	Mkt&Fin	66.28	Placed
2	M	65.00	Central	68.00	Central	Arts	64.00	Comm&Mgmt	No	75.00	Mkt&Fin	57.80	Placed
4	M	85.80	Central	73.60	Central	Commerce	73.30	Comm&Mgmt	No	96.80	Mkt&Fin	55.50	Placed
7	M	82.00	Central	64.00	Central	Science	66.00	Sci&Tech	Yes	67.00	Mkt&Fin	62.14	Placed
...
207	M	83.33	Central	78.00	Others	Commerce	61.00	Comm&Mgmt	Yes	88.56	Mkt&Fin	71.55	Placed
209	M	62.00	Central	72.00	Central	Commerce	66.00	Comm&Mgmt	No	67.00	Mkt&Fin	56.49	Placed
210	M	80.60	Others	82.00	Others	Commerce	77.60	Comm&Mgmt	No	91.00	Mkt&Fin	74.49	Placed
211	M	58.00	Others	60.00	Others	Science	72.00	Sci&Tech	No	74.00	Mkt&Fin	53.62	Placed
212	M	67.00	Others	67.00	Others	Commerce	73.00	Comm&Mgmt	Yes	59.00	Mkt&Fin	69.72	Placed

100 rows × 13 columns

In [22]: a.shape

Out[22]: (180, 13)

100 Male candidates are placed out of 139

In [23]: df.nunique()

Out[23]:

gender	2
ssc_percentage	183
ssc_board	2
hsc_percentage	92
hsc_board	2
hsc_subject	3
degree_percentage	89
undergrad_degree	3
work_experience	2
emp_test_percentage	190
specialisation	2
mba_percent	285
status	2
dtype: int64	

In [24]: # how many Female candidate placed
b = df.groupby(['gender', 'status']).get_group(('F','Placed'))
b

Out[24]:

	gender	ssc_percentage	ssc_board	hsc_percentage	hsc_board	hsc_subject	degree_percentage	undergrad_degree	work_experience	emp_test_percentage	specialisation	mba_percent	status
13	F	77.00	Central	87.00	Central	Commerce	59.00	Comm&Mgmt	No	86.00	Mkt&Fin	68.63	Placed
15	F	65.00	Central	75.00	Central	Commerce	69.00	Comm&Mgmt	Yes	72.00	Mkt&Fin	51.58	Not Placed
21	F	79.00	Others	76.00	Others	Commerce	85.00	Sci&Tech	No	95.00	Mkt&Fin	69.06	Placed
22	F	69.80	Others	60.80	Others	Science	72.23	Sci&Tech	No	55.53	Mkt&HR	68.81	Placed
23	F	77.40	Others	60.00	Others	Science	64.74	Sci&Tech	Yes	82.00	Mkt&Fin	63.62	Placed
30	F	64.00	Central	73.50	Central	Commerce	73.00	Comm&Mgmt	No	52.00	Mkt&HR	56.70	Placed
32	F	61.00	Central	81.00	Central	Commerce	66.40	Comm&Mgmt	No	50.89	Mkt&HR	62.21	Placed
33	F	87.00	Others	65.00	Others	Science	81.00	Comm&Mgmt	Yes	88.00	Mkt&Fin	72.78	Placed
35	F	69.00	Central	78.00	Central	Commerce	72.00	Comm&Mgmt	No	71.00	Mkt&HR	62.74	Placed
37	F	79.00	Central	76.00	Central	Science	65.60	Sci&Tech	No	58.00	Mkt&Fin	66.68	Placed
38	F	73.00	Others	58.00	Others	Science	66.00	Comm&Mgmt	No	53.70	Mkt&HR	56.86	Placed
40	F	78.00	Central	77.00	Others	Commerce	80.00	Comm&Mgmt	No	60.00	Mkt&Fin	66.72	Placed
44	F	77.00	Others	73.00	Others	Commerce	81.00	Comm&Mgmt	Yes	89.00	Mkt&Fin	69.70	Placed
50	F	75.20	Central	73.20	Central	Science	68.40	Comm&Mgmt	No	65.00	Mkt&HR	62.98	Placed
54	F	74.00	Central	60.00	Others	Science	69.00	Comm&Mgmt	No	79.00	Mkt&HR	65.56	Placed
62	F	86.50	Others	64.20	Others	Science	67.40	Sci&Tech	No	79.00	Mkt&Fin	59.69	Placed
76	F	66.50	Others	70.40	Central	Arts	71.93	Sci&Tech	No	61.00	Mkt&Fin	64.27	Placed
80	F	69.00	Others	62.00	Others	Commerce	69.00	Comm&Mgmt	Yes	67.00	Mkt&HR	62.35	Placed
86	F	83.84	Others	89.83	Others	Commerce	77.20	Comm&Mgmt	Yes	78.74	Mkt&Fin	71.77	Placed
88	F	66.00	Central	62.00	Central	Commerce	73.00	Comm&Mgmt	No	56.00	Mkt&HR	64.36	Placed
89	F	84.00	Others	75.00	Others	Science	69.00	Sci&Tech	Yes	62.00	Mkt&HR	62.36	Placed
92	F	85.00	Others	90.00	Others	Commerce	82.00	Comm&Mgmt	No	92.00	Mkt&Fin	68.03	Placed
92	F	60.23	Central	69.00	Central	Science	66.00	Comm&Mgmt	No	72.00	Mkt&Fin	59.47	Placed
96	F	76.00	Central	70.00	Central	Science	75.00	Comm&Mgmt	Yes	66.00	Mkt&Fin	64.44	Placed
98	F	69.00	Central	73.00	Central	Commerce	65.00	Comm&Mgmt	No	70.00	Mkt&Fin	57.31	Placed
102	F	77.00	Others	61.00	Others	Commerce	68.00	Comm&Mgmt	Yes	57.50	Mkt&Fin	61.31	Placed
110	F	69.50	Central	70.00	Central	Science	72.00	Sci&Tech	No	57.20	Mkt&HR	54.80	Placed
113	F	73.96	Others	79.00	Others	Commerce	67.00	Comm&Mgmt	No	72.15	Mkt&Fin	63.08	Placed
115	F	73.00	Others	63.00	Others	Science	66.00	Comm&Mgmt	No	89.00	Mkt&Fin	60.50	Placed
121	F	64.00	Central	67.00	Others	Science	69.60	Sci&Tech	Yes	55.67	Mkt&HR	77.89	Placed
122	F	65.00	Central	66.80	Central	Arts	69.30	Comm&Mgmt	Yes	80.40	Mkt&Fin	71.00	Placed
125	F												