

In [1]:

```
1 import pandas as pd
2 import matplotlib.pyplot as plt
3 import numpy as np
4 import seaborn as sns
5 import scipy.stats as stats
```

In [2]:

```
1 dataset=pd.read_csv("PreProcessed_Placement.csv")
2 dataset
```

Out [2]:

	sl_no	gender	ssc_p	ssc_b	hsc_p	hsc_b	hsc_s	degree_p	degree_t	workex	etest_p	specialisation	mba_p	status	salary	
	0	1	M	67.00	Others	91.00	Others	Commerce	58.00	Sci&Tech	No	55.0	Mkt&HR	58.80	Placed	270000.0
	1	2	M	79.33	Central	78.33	Others	Science	77.48	Sci&Tech	Yes	86.5	Mkt&Fin	66.28	Placed	200000.0
	2	3	M	65.00	Central	68.00	Central	Arts	64.00	Comm&Mgmt	No	75.0	Mkt&Fin	57.80	Placed	250000.0
	3	4	M	56.00	Central	52.00	Central	Science	52.00	Sci&Tech	No	66.0	Mkt&HR	59.43	Not Placed	0.0
	4	5	M	85.80	Central	73.60	Central	Commerce	73.30	Comm&Mgmt	No	96.8	Mkt&Fin	55.50	Placed	425000.0

	210	211	M	80.60	Others	82.00	Others	Commerce	77.60	Comm&Mgmt	No	91.0	Mkt&Fin	74.49	Placed	400000.0
	211	212	M	58.00	Others	60.00	Others	Science	72.00	Sci&Tech	No	74.0	Mkt&Fin	53.62	Placed	275000.0
	212	213	M	67.00	Others	67.00	Others	Commerce	73.00	Comm&Mgmt	Yes	59.0	Mkt&Fin	69.72	Placed	295000.0
	213	214	F	74.00	Others	66.00	Others	Commerce	58.00	Comm&Mgmt	No	70.0	Mkt&HR	60.23	Placed	204000.0
	214	215	M	62.00	Central	58.00	Others	Science	53.00	Comm&Mgmt	No	89.0	Mkt&HR	60.22	Not Placed	0.0

215 rows × 15 columns

In [3]:

```
1 dataset.isnull().sum()
```

Out[3]:

```
sl_no      0
gender      0
ssc_p      0
ssc_b      0
hsc_p      0
hsc_b      0
hsc_s      0
degree_p   0
degree_t   0
workex     0
etest_p    0
specialisation  0
mba_p      0
status     0
salary     0
dtype: int64
```

Q.1 - Replace the NaN values with correct value. And justify why you have chosen the same.

```
dataset["salary"].fillna(0,inplace=True) dataset
```

```
dataset.isnull().sum()
```

ANS : Found NaN values in the salary column. Replaced the NaN values with 0 because they wont receive salary since they are not placed. Besides the meaning of original dataset is retained by the replacement with zero.

In []:

```
1
```

Q.2 - How many of them are not placed?

In [4]:

```
1 count = dataset['status'].value_counts().get('Not Placed', 0)
2 count
```

Out[4]:

```
67
```

Ans : 67

In []:

```
1
```

Q.3- Find the reason for non placement from the dataset.

Defining Quan & Qual

In [5]:

1

2

3

4

5

6

7

8

9

10

quan=[]

qual=[]

for columnName in dataset.columns:

#print(columnName)

if(dataset[columnName].dtype=='0'):

#print("qual")

qual.append(columnName)

else:

#print("quan")

quan.append(columnName)

In [6]:

1

quan

Out[6]: ['sl_no', 'ssc_p', 'hsc_p', 'degree_p', 'etest_p', 'mba_p', 'salary']

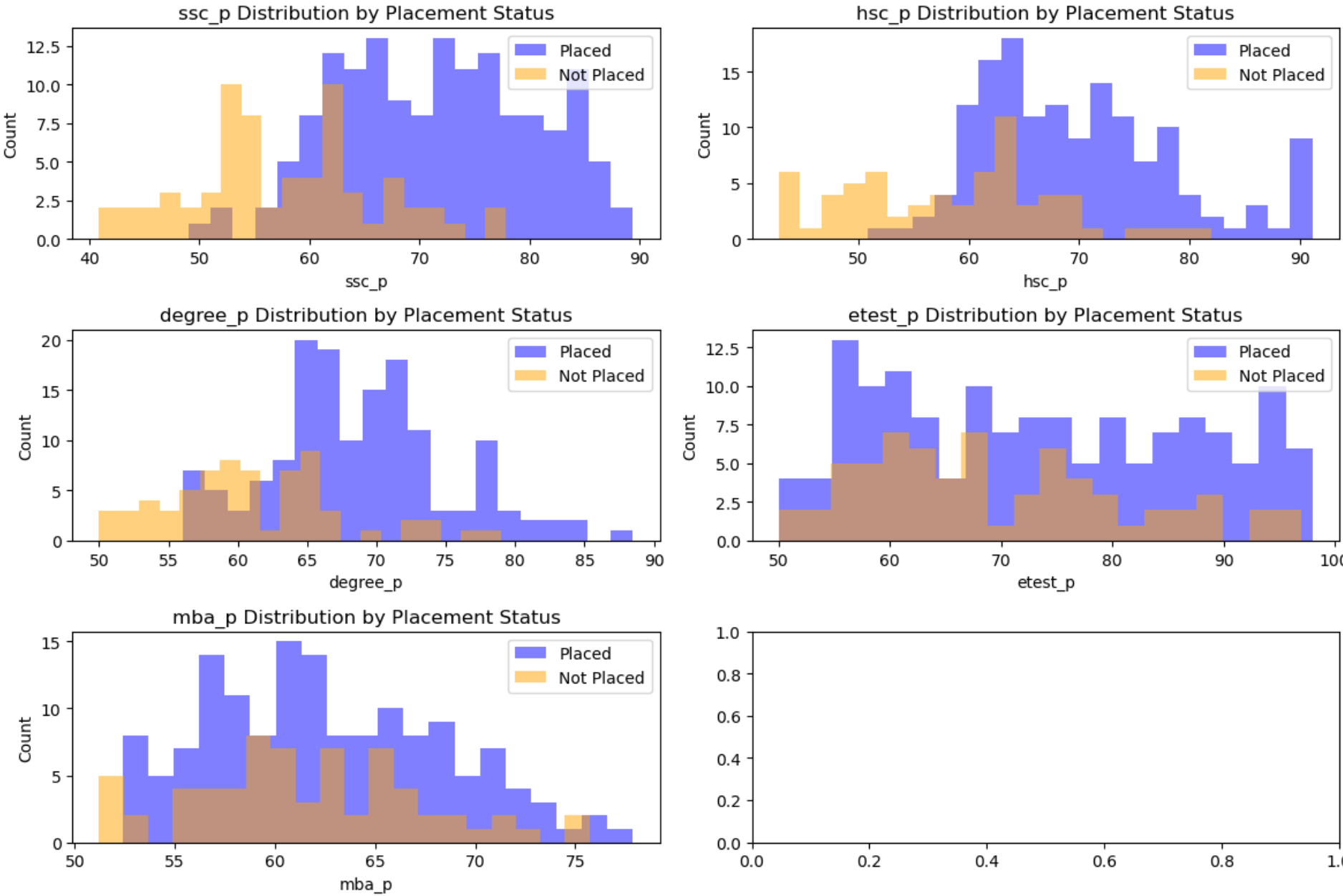
In [7]:

1

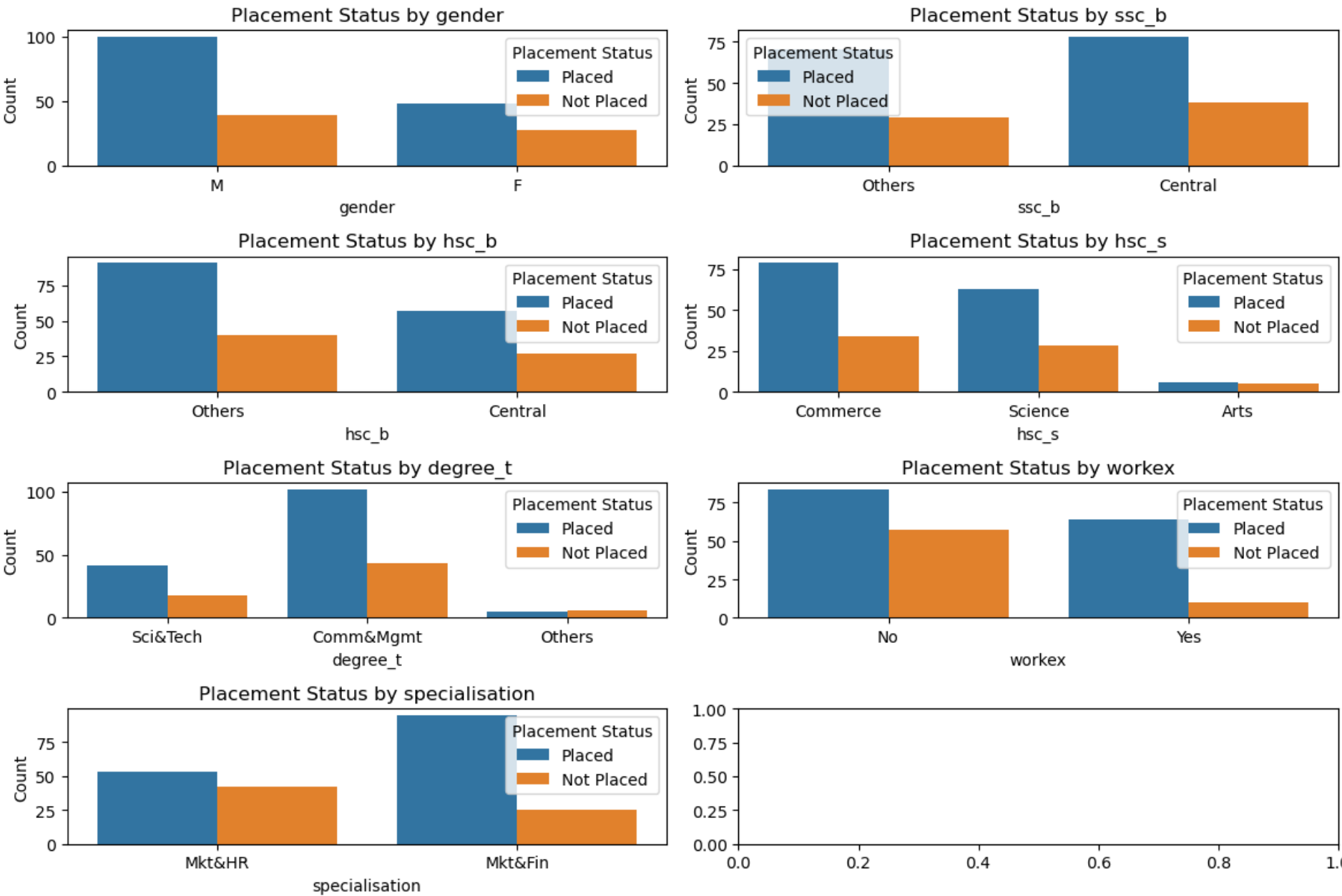
qual

Out[7]: ['gender',
'ssc_b',
'hsc_b',
'hsc_s',
'degree_t',
'workex',
'specialisation',
'status']

```
In [8]: 1 import matplotlib.pyplot as plt
2 import numpy as np
3
4 # Define the list of numerical columns you want to include
5 numerical_columns = quan
6
7 # Define the columns to exclude
8 exclude_columns = ['sl_no', 'salary']
9
10 # Filter out the excluded columns
11 filtered_numerical_columns = [col for col in numerical_columns if col not in exclude_columns]
12
13 # Calculate the number of rows and columns for subplots
14 num_plots = len(filtered_numerical_columns)
15 num_cols = 2 # Two columns for 'Placed' and 'Not Placed'
16 num_rows = (num_plots + 1) // num_cols # Calculate the number of rows needed
17
18 # Create a single figure with subplots
19 fig, axes = plt.subplots(num_rows, num_cols, figsize=(12, 8))
20
21 # Flatten the axes array if there's only one row
22 if num_rows == 1:
23     axes = [axes]
24
25 # Loop through each numerical column and create histograms
26 for idx, column in enumerate(filtered_numerical_columns):
27     row_idx = idx // num_cols
28     col_idx = idx % num_cols
29
30     # Create histograms for the current numerical column, grouped by placement status
31     axes[row_idx][col_idx].hist(dataset[dataset['status'] == 'Placed'][column], bins=20, alpha=0.5, label='Placed')
32     axes[row_idx][col_idx].hist(dataset[dataset['status'] == 'Not Placed'][column], bins=20, alpha=0.5, label='Not Placed')
33
34     # Customize the subplot
35     axes[row_idx][col_idx].set_xlabel(column)
36     axes[row_idx][col_idx].set_ylabel('Count')
37     axes[row_idx][col_idx].set_title(f'{column} Distribution by Placement Status')
38     axes[row_idx][col_idx].legend()
39
40 # Adjust the layout to prevent overlapping labels
41 plt.tight_layout()
42
43 # Show the plots
44 plt.show()
45
```



```
In [9]: 1 import matplotlib.pyplot as plt
2 import seaborn as sns
3
4 # Define the list of categorical columns you want to include
5 categorical_columns = qual
6
7 # Define the column to exclude
8 exclude_column = 'status'
9
10 # Calculate the number of rows and columns for subplots
11 num_plots = len(categorical_columns) - (1 if exclude_column in categorical_columns else 0)
12 num_cols = 2 # Two columns for 'Placed' and 'Not Placed'
13 num_rows = (num_plots + 1) // num_cols # Calculate the number of rows needed
14
15 # Create a single figure with subplots
16 fig, axes = plt.subplots(num_rows, num_cols, figsize=(12, 8))
17
18 # Flatten the axes array if there's only one row
19 if num_rows == 1:
20     axes = [axes]
21
22 # Loop through each categorical column and create count plots
23 plot_idx = 0 # Index for tracking the subplot position
24 for column in categorical_columns:
25     if column != exclude_column:
26         row_idx = plot_idx // num_cols
27         col_idx = plot_idx % num_cols
28
29         # Create a count plot for the current categorical column
30         sns.countplot(data=dataset, x=column, hue='status', ax=axes[row_idx][col_idx])
31
32         # Customize the subplot
33         axes[row_idx][col_idx].set_xlabel(column)
34         axes[row_idx][col_idx].set_ylabel('Count')
35         axes[row_idx][col_idx].set_title(f'Placement Status by {column}')
36         axes[row_idx][col_idx].legend(title='Placement Status')
37
38         plot_idx += 1
39
40 # Adjust the layout to prevent overlapping labels
41 plt.tight_layout()
42
43 # Show the plots
44 plt.show()
```



Ans: By analysing the relationship of categorical column & numerical column with respect to placement status from the above plots, we can clearly see that there is no significant influence of numerical column on the not placed placement status, where as categorical columns like degree_t(Comm & Mmnt), work experience & specialisation(MKT & HR) has greater influence on the not placed placement status.

In []:

1

Q.4- What is the kind of relation is between salary and mba_p?

In [10]:

1 dataset.corr()
2 #numeric_dataset = dataset.select_dtypes(include=['number'])
3 #correlation_matrix = numeric_dataset.corr()
4 #correlation_matrix

/var/folders/07/ykgp85052b11h5kz22ghn8l40000gn/T/ipykernel_89315/3795343161.py:1: FutureWarning: The default value of numeric_only in DataFrame.corr is deprecated. In a future version, it will default to False. Select only valid columns or specify the value of numeric_only to silence this warning.
dataset.corr()

Out[10]:

	sl_no	ssc_p	hsc_p	degree_p	etest_p	mba_p	salary
sl_no	1.000000	-0.078155	-0.090500	-0.091039	0.063636	0.022327	0.001217
ssc_p	-0.078155	1.000000	0.513478	0.538686	0.261993	0.388478	0.558475
hsc_p	-0.090500	0.513478	1.000000	0.434606	0.240775	0.348452	0.459424
degree_p	-0.091039	0.538686	0.434606	1.000000	0.227147	0.402376	0.423762
etest_p	0.063636	0.261993	0.240775	0.227147	1.000000	0.218055	0.186775
mba_p	0.022327	0.388478	0.348452	0.402376	0.218055	1.000000	0.141417
salary	0.001217	0.558475	0.459424	0.423762	0.186775	0.141417	1.000000

Ans: mba_p & salary exhibits a low degree positive correlation.This implies that individuals with higher MBA_p tend to have slightly higher salaries, but the relationship is not very strong.

In []:

1

Q.5- Which specialization is getting minimum salary?

In [11]:

1 min_salary_by_specialization = dataset.groupby('specialisation')['salary'].min()
2 min_salary_specializations = min_salary_by_specialization[min_salary_by_specialization == min_salary_by_
3 min_salary_specializations

Out[11]: specialisation
Mkt&Fin 0.0
Mkt&HR 0.0
Name: salary, dtype: float64

In [12]:

1 # Replace with the value you want to exclude
2 exclude_value = 0.0
3
4 # Filter the DataFrame to exclude the specified value
5 filtered_dataset = dataset[dataset['salary'] != exclude_value]
6
7 # Check if there are any remaining rows after excluding the value
8 if not filtered_dataset.empty:
9
10 # Find the minimum salary in the filtered DataFrame
11 min_salary_by_specialization = filtered_dataset.groupby('specialisation')['salary'].min()
12 min_salary_specializations = min_salary_by_specialization[min_salary_by_specialization == min_salar
13 min_salary_specializations

In [13]:

1 min_salary_specializations

Out[13]: specialisation
Mkt&Fin 200000.0
Mkt&HR 200000.0
Name: salary, dtype: float64

Ans: Mkt&Fin & Mkt&HR are the 2 specialisations which receive a minimum wage of 200000 if we exclude the salary zero since it is the value which has been replaced for the non placed students inorder to retain the meaning of original dataset.

In []:

1

Q.6 - How many of them are getting above 500000 salary?

```
In [14]: 1 threshold_salary = 500000
2 count_above_threshold = len(dataset[dataset['salary'] > threshold_salary])
3 count_above_threshold

Out[14]: 3
```

Ans: 3 members are getting the salary above 500000.

```
In [ ]: 1
```

Q.7 - Test the Analysis of Variance between etest_p and mba_p at signifance level 5%.(Make decision using Hypothesis Testing)

```
In [15]: 1 import scipy.stats as stats
2 stats.f_oneway(dataset['etest_p'],dataset['mba_p'])

Out[15]: F_onewayResult(statistic=98.64487057324706, pvalue=4.672547689133573e-21)
```

Ans : Since p value is < 0.05 we reject null hypothesis(H0) & accept alternate hypothesis(H1) & this implies that there is significant difference between etest_p and mba_p. The F-statistic and the small p-value support this conclusion.

```
In [ ]: 1
```

Q.8 - Test the similarity between the degree_t(Sci&Tech) and specialisation(Mkt&HR) with respect to salary at significance level of 5%.(Make decision using Hypothesis Testing)

```
In [16]: 1 from scipy.stats import ttest_ind
2 dataset=dataset.dropna()
3 Sci_Tech = dataset[dataset['degree_t']=='Sci&Tech']['salary']
4 Mkt_HR = dataset[dataset['specialisation']=='Mkt&HR']['salary']
5 ttest_ind(Sci_Tech, Mkt_HR)

Out[16]: TtestResult(statistic=2.692041243555374, pvalue=0.007897969943471179, df=152.0)
```

```
In [17]: 1 # Define the categorical values you want to compare
2 category_1 = 'Sci&Tech'
3 category_2 = 'Mkt&HR'
4
5 # Filter the DataFrame for each category
6 subset_1 = dataset[dataset['degree_t'] == category_1]
7 subset_2 = dataset[dataset['specialisation'] == category_2]
8
9 # Calculate the mean salary for each subset
10 mean_salary_1 = subset_1['salary'].mean()
11 mean_salary_2 = subset_2['salary'].mean()
12
13 # Print the mean salaries for each category
14 print(f'Mean Salary for {category_1}: {mean_salary_1}')
15 print(f'Mean Salary for {category_2}: {mean_salary_2}')
16
```

Mean Salary for Sci&Tech: 218627.11864406778
Mean Salary for Mkt&HR: 150842.1052631579

Ans : Since p value is < 0.05 we reject null hypothesis(H0) & accept alternate hypothesis(H1) & this implies that there is significant difference between the degree_t(Sci&Tech) and specialisation(Mkt&HR) with respect to salaries mean. Additionally, the positive t-statistic suggests that the mean salary for the 'Science & Technology' group is higher than the mean salary for the 'Marketing & HR' group.

To conclude there is a significant difference in salaries between these two groups, with 'Science & Technology' having a higher mean salary at a 5% significance level.

```
In [ ]: 1
```

Q.9 - Convert the normal distribution to standard normal distribution for salary column.


```
In [18]: 1 import matplotlib.pyplot as plt
2 import seaborn as sns
3
4
5 fig, (ax1, ax2) = plt.subplots(1, 2, figsize=(12, 4))
6
7 sns.distplot(dataset['salary'], ax=ax1)
8 ax1.set_title("Original Normal Distribution")
9
10 # Define the stdNBgraph function
11 def stdNBgraph(dataset, ax):
12     mean = dataset.mean()
13     std = dataset.std()
14     values = [i for i in dataset]
15     z_score = [(j - mean) / std for j in values]
16     sns.distplot(z_score, kde=True, ax=ax)
17     ax.set_title("Standard Normal Distribution")
18
19 # Call the stdNBgraph function on the same dataset and the second subplot
20 stdNBgraph(dataset['salary'], ax2)
21
22 # Show the merged plot
23 plt.show()
24
```

/var/folders/07/ykgp85052b11h5kz22ghn8l40000gn/T/ipykernel_89315/3264537401.py:7: UserWarning:

`distplot` is a deprecated function and will be removed in seaborn v0.14.0.

Please adapt your code to use either `displot` (a figure-level function with similar flexibility) or `histplot` (an axes-level function for histograms).

For a guide to updating your code to use the new functions, please see <https://gist.github.com/mwaskom/de44147ed2974457ad6372750bbe5751> (<https://gist.github.com/mwaskom/de44147ed2974457ad6372750bbe5751>)

sns.distplot(dataset['salary'], ax=ax1)

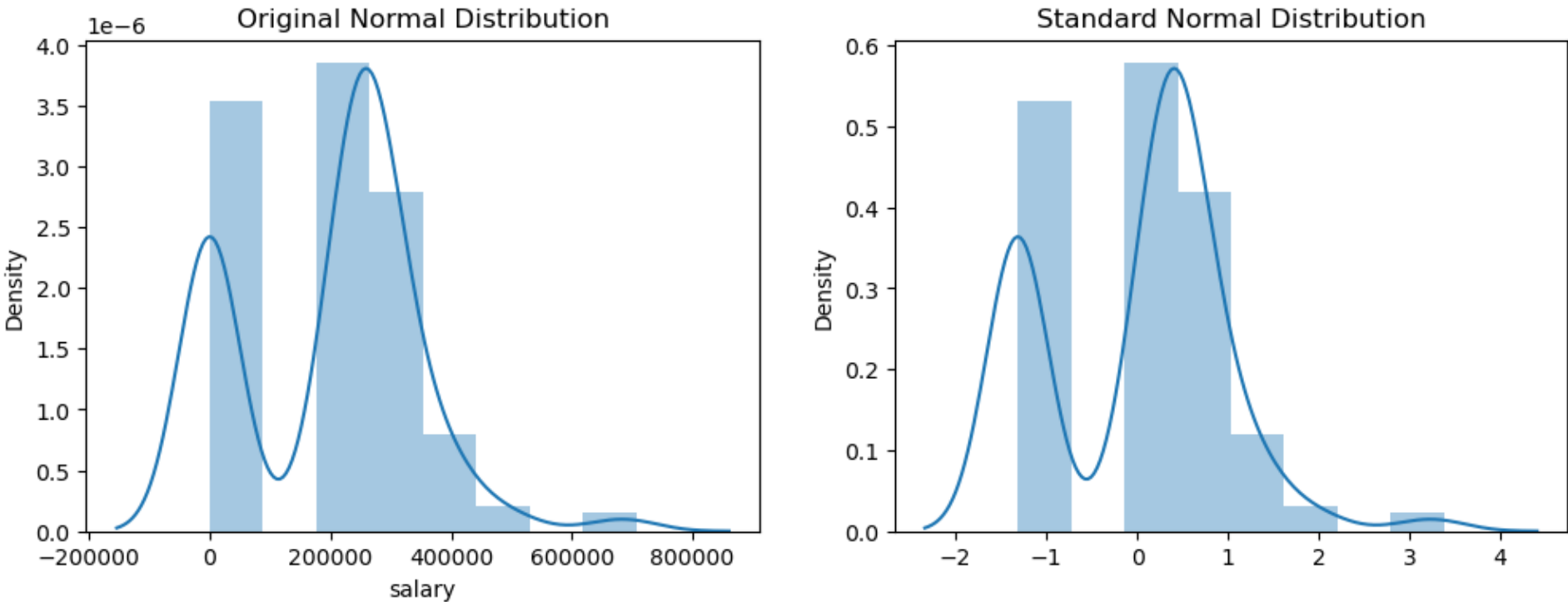
/var/folders/07/ykgp85052b11h5kz22ghn8l40000gn/T/ipykernel_89315/3264537401.py:16: UserWarning:

`distplot` is a deprecated function and will be removed in seaborn v0.14.0.

Please adapt your code to use either `displot` (a figure-level function with similar flexibility) or `histplot` (an axes-level function for histograms).

For a guide to updating your code to use the new functions, please see <https://gist.github.com/mwaskom/de44147ed2974457ad6372750bbe5751> (<https://gist.github.com/mwaskom/de44147ed2974457ad6372750bbe5751>)

sns.distplot(z_score, kde=True, ax=ax)



In []:

1

Q.10 - What is the Probability Density Function of the salary range from 700000 to 900000?

```
In [19]: 1 def get_pdf_probability(dataset,startrange,endrange):
2         from matplotlib import pyplot
3         from scipy.stats import norm
4         import seaborn as sns
5         ax = sns.distplot(dataset,kde=True,kde_kws={'color':'blue'},color='Green')
6         pyplot.axvline(startrange,color='Red')
7         pyplot.axvline(endrange,color='Red')
8         # generate a sample
9         sample = dataset
10        # calculate parameters
11        sample_mean =sample.mean()
12        sample_std = sample.std()
13        print('Mean=%.3f, Standard Deviation=%.3f' % (sample_mean, sample_std))
14        # define the distribution
15        dist = norm(sample_mean, sample_std)
16
17        # sample probabilities for a range of outcomes
18        values = [value for value in range(startrange, endrange)]
19        probabilities = [dist.pdf(value) for value in values]
20        prob=sum(probabilities)
21        print("The area between range({},{}):{}".format(startrange,endrange,sum(probabilities)))
22        return prob
23
```

```
In [20]: 1 get_pdf_probability(dataset["salary"],700000,900000)
```

/var/folders/07/ykgp85052b11h5kz22ghn8l40000gn/T/ipykernel_89315/2842244316.py:5: UserWarning:

`distplot` is a deprecated function and will be removed in seaborn v0.14.0.

Please adapt your code to use either `displot` (a figure-level function with similar flexibility) or `histplot` (an axes-level function for histograms).

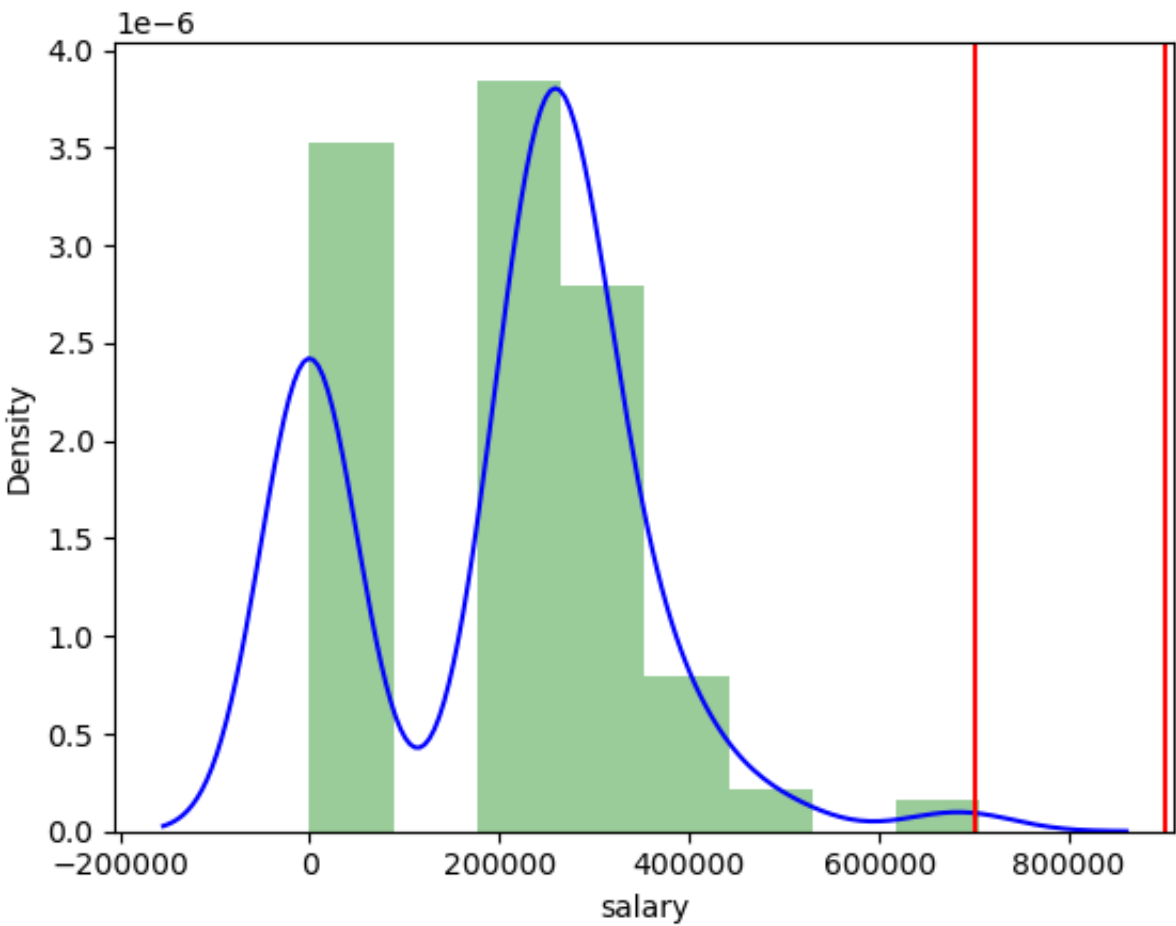
For a guide to updating your code to use the new functions, please see <https://gist.github.com/mwaskom/de44147ed2974457ad6372750bbe5751> (<https://gist.github.com/mwaskom/de44147ed2974457ad6372750bbe5751>)

ax = sns.distplot(dataset,kde=True,kde_kws={'color':'blue'},color='Green')

Mean=197615.116, Standard Deviation=150305.844

The area between range(700000,900000):0.0004137812663830059

Out[20]: 0.0004137812663830059



Ans: The Probability Density Function of the salary range from 7,00,000 to 9,00,000 is 0.000413 which implies an extremely low probability is associated within this range, reflecting that extremely very small number of students receive the salary in this specified range in the distribution.

```
In [ ]: 1
```

Q. 11 - Test the similarity between the degree_t(Sci&Tech)with respect to etest_p and mba_p at significance level of 5%.(Make decision using Hypothesis Testing)

In [21]:

```
1 from scipy.stats import ttest_rel
2 #dataset=dataset.dropna()
3 etest_p = dataset[dataset['degree_t']=='Sci&Tech']['etest_p']
4 mba_p= dataset[dataset['degree_t']=='Sci&Tech']['mba_p']
5 ttest_rel(etest_p, mba_p)
```

Out[21]: TtestResult(statistic=5.0049844583693615, pvalue=5.517920600505392e-06, df=58)

Since p value is < 0.05 we reject null hypothesis(H0) & accept alternate hypothesis(H1) & this implies that there is significant dissimilarity between the etest_p and mba_p wrt degree_t(Sci&Tech).

In []:

1

Q.12 - Which parameter is highly correlated with salary?

In [22]:

```
1 dataset.corr()
```

/var/folders/07/ykgp85052b11h5kz22ghn8l40000gn/T/ipykernel_89315/2191645083.py:1: FutureWarning: The default value of numeric_only in DataFrame.corr is deprecated. In a future version, it will default to False. Select only valid columns or specify the value of numeric_only to silence this warning.

```
dataset.corr()
```

Out[22]:

	sl_no	ssc_p	hsc_p	degree_p	etest_p	mba_p	salary
sl_no	1.000000	-0.078155	-0.090500	-0.091039	0.063636	0.022327	0.001217
ssc_p	-0.078155	1.000000	0.513478	0.538686	0.261993	0.388478	0.558475
hsc_p	-0.090500	0.513478	1.000000	0.434606	0.240775	0.348452	0.459424
degree_p	-0.091039	0.538686	0.434606	1.000000	0.227147	0.402376	0.423762
etest_p	0.063636	0.261993	0.240775	0.227147	1.000000	0.218055	0.186775
mba_p	0.022327	0.388478	0.348452	0.402376	0.218055	1.000000	0.141417
salary	0.001217	0.558475	0.459424	0.423762	0.186775	0.141417	1.000000

Ans: Ssc_p is highly correlated with salary by the value 0.558475. This suggests a moderately positive linear relationship between ssc_p and salary. In other words, on average, the individuals with higher ssc_p tend to have higher salaries.

In []:

1

Q.13 - Plot any useful graph and explain it.

```
In [23]: 1 import seaborn as sns
2 sns.distplot(dataset["mba_p"])
3 mean=dataset['mba_p'].mean()
4 std=dataset['mba_p'].std()
5 skew=dataset["mba_p"].skew()
6 print('Mean=%.3f, Standard Deviation=%.3f, Skew=%.3f' % (mean,std, skew))
```

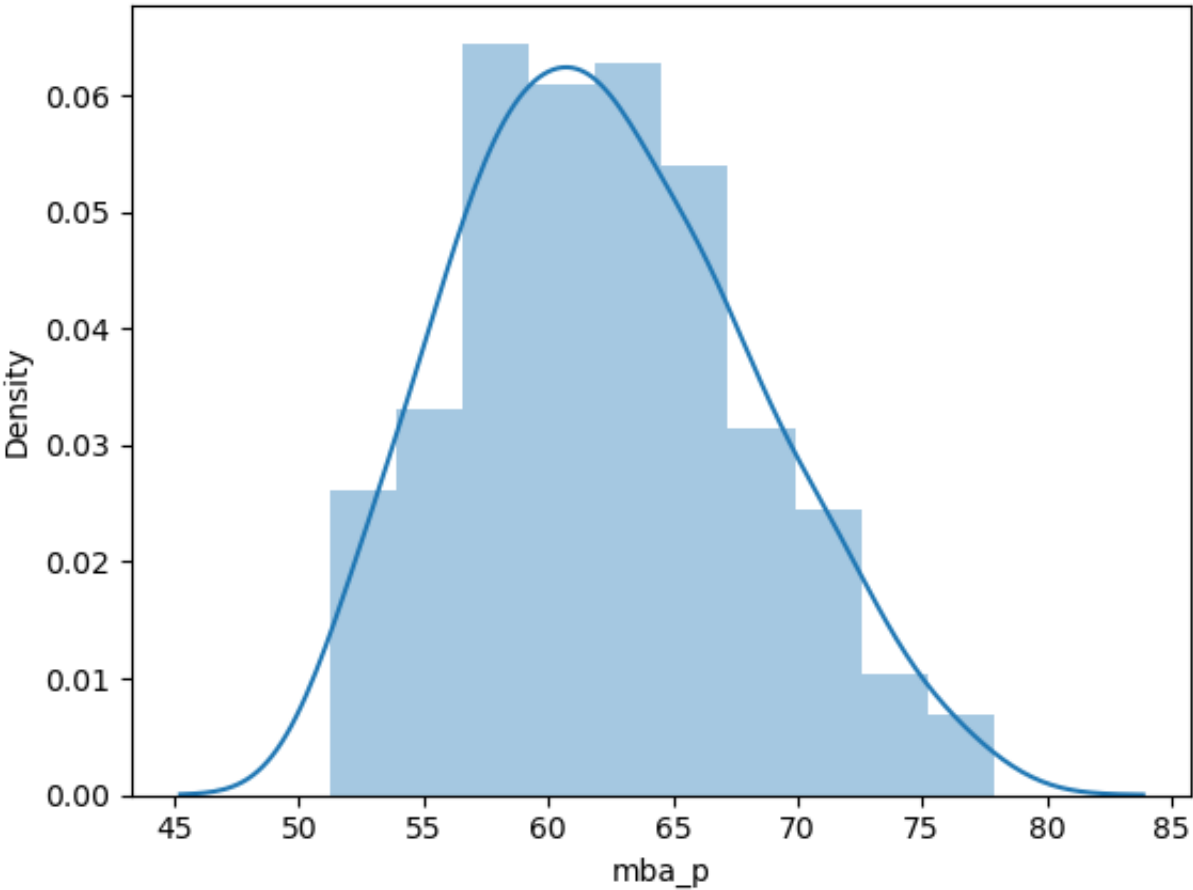
Mean=62.278, Standard Deviation=5.833, Skew=0.314

/var/folders/07/ykgp85052b11h5kz22ghn8l40000gn/T/ipykernel_89315/3102433105.py:2: UserWarning:
`distplot` is a deprecated function and will be removed in seaborn v0.14.0.

Please adapt your code to use either `displot` (a figure-level function with similar flexibility) or `histplot` (an axes-level function for histograms).

For a guide to updating your code to use the new functions, please see <https://gist.github.com/mwaskom/de44147ed2974457ad6372750bbe5751> (<https://gist.github.com/mwaskom/de44147ed2974457ad6372750bbe5751>)

```
sns.distplot(dataset["mba_p"])
```



Ans : The above plot gives the distribution of MBA pass marks.

From the graph , it is obvious that the average score i.e the mean in mba_p is 62.278, and it represents the central tendency of the data.

The standard deviation is 5.833, which indicates the spread or variability of the data around the mean. A larger standard deviation would result in a wider and flatter distribution, while a smaller standard deviation would result in a narrower and taller distribution.

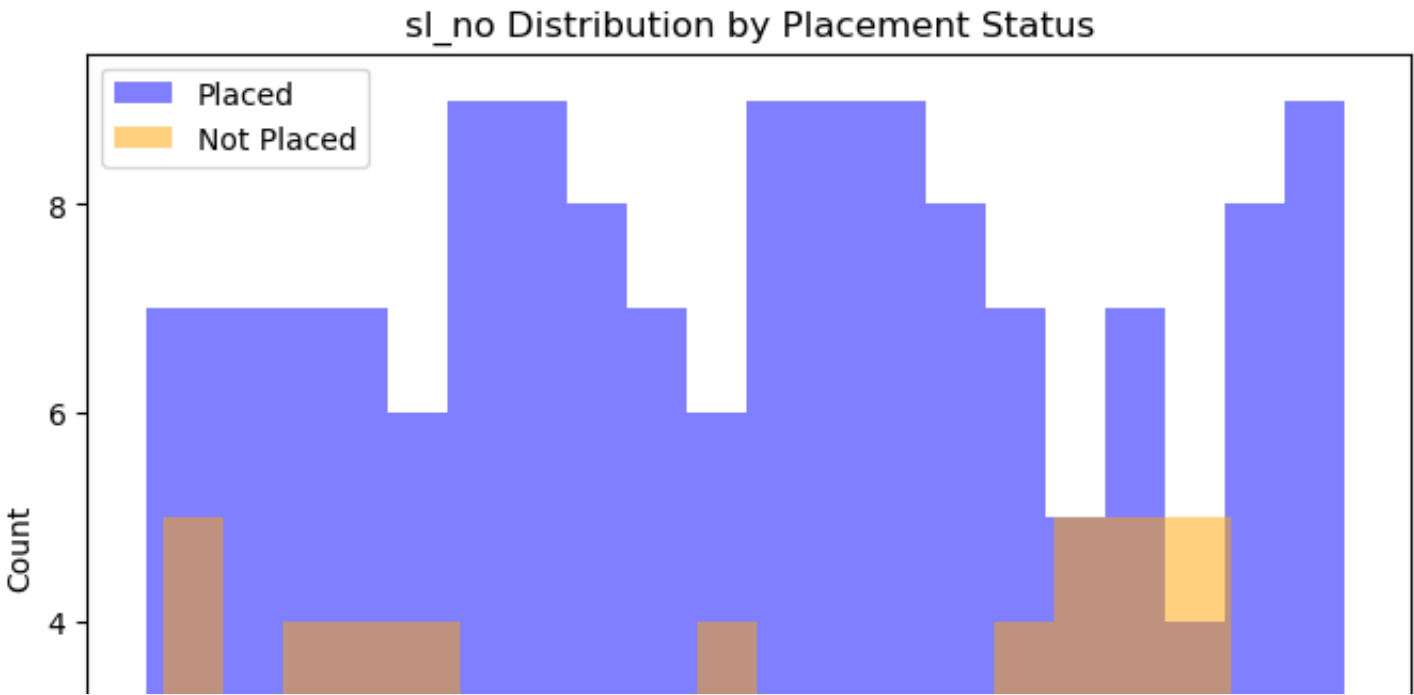
By visualizing the data using this plot, we get the insights about the distribution of MBA pass marks and how they are related to the mean and standard deviation.

Additional Inference: MBA pass marks exhibits a positive skewness value of 0.314 which suggests a slight right-skew in the distribution, indicating the presence of a few high-scoring individuals with exceptionally high scores compared to the majority. However, the skewness is not very pronounced, and the distribution is relatively close to being symmetric.

```
In [ ]: 1
```

Miscellaneous

```
In [24]: 1 import matplotlib.pyplot as plt
2 # Define the list of numerical columns you want to include
3 numerical_columns = quan
4
5 # Define the column to exclude
6 exclude_column = 'sl_no', 'salary'
7
8 for column in numerical_columns:
9     if column != exclude_column:
10         plt.figure(figsize=(8, 6))
11
12         # Create histograms for the current numerical column, grouped by placement status
13         plt.hist(dataset[dataset['status'] == 'Placed'][column], bins=20, alpha=0.5, label='Placed', co
14         plt.hist(dataset[dataset['status'] == 'Not Placed'][column], bins=20, alpha=0.5, label='Not Pla
15
16         # Customize the plot
17         plt.xlabel(column)
18         plt.ylabel('Count')
19         plt.title(f'{column} Distribution by Placement Status')
20         plt.legend()
21
22         plt.show()
23
```



```
In [25]: 1 import seaborn as sns
2
3 # Define the list of categorical columns you want to include
4 categorical_columns = qual
5
6 # Define the column to exclude
7 exclude_column = 'status'
8
9 # Loop through each categorical column and create count plots, excluding the specified column
10 for column in categorical_columns:
11     if column != exclude_column:
12         plt.figure(figsize=(8, 6))
13
14         # Create a count plot for the current categorical column
15         sns.countplot(data=dataset, x=column, hue='status')
16
17         # Customize the plot
18         plt.xlabel(column)
19         plt.ylabel('Count')
20         plt.title(f'Placement Status by {column}')
21         plt.legend(title='Placement Status')
22
23         plt.show()
```

