

In [1]:

```
import pandas as pd
import numpy as np
import seaborn as sns
import matplotlib.pyplot as plt
%matplotlib inline
from statsmodels.formula.api import ols
from statsmodels.stats.anova import _get_covariance,anova_lm
import os
```

Problem 1

Salaries of 40 individual, ANOVA test for Salary is hypothesized to depend on educational qualification and occupation as given data.

In [2]:

```
empdata=pd.read_csv("SalaryData.csv")
```

In [3]:

```
empdata.info()
```

```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 40 entries, 0 to 39
Data columns (total 3 columns):
 #   Column      Non-Null Count  Dtype  
 ---  --          --          --      
 0   Education    40 non-null    object  
 1   Occupation   40 non-null    object  
 2   Salary       40 non-null    int64  
dtypes: int64(1), object(2)
memory usage: 1.1+ KB
```

In [4]:

```
empdata.isnull().sum()
```

Out[4]:

```
Education      0
Occupation     0
Salary         0
dtype: int64
```

In [5]:

```
empdata.head(8)
```

Out[5]:

	Education	Occupation	Salary
0	Doctorate	Adm-clerical	153197
1	Doctorate	Adm-clerical	115945
2	Doctorate	Adm-clerical	175935
3	Doctorate	Adm-clerical	220754
4	Doctorate	Sales	170769
5	Doctorate	Sales	219420
6	Doctorate	Sales	237920
7	Doctorate	Sales	160540

In [6]:

```
pd.DataFrame(empdata.groupby('Education').Occupation.value_counts())
```

Out[6]:

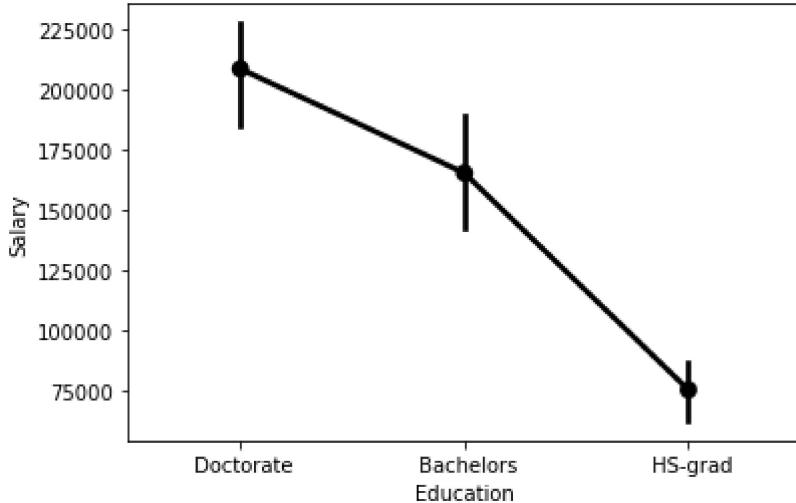
Education	Occupation	Occupation		
		Bachelors	Doctorate	HS-grad
Bachelors	Exec-managerial	4		
	Prof-specialty	4		
	Sales	4		
	Adm-clerical	3		
Doctorate	Prof-specialty	6		
	Sales	5		
	Adm-clerical	4		
	Exec-managerial	1		
HS-grad	Adm-clerical	3		
	Prof-specialty	3		
	Sales	3		

In [9]:

```
sns.pointplot(x='Education', y='Salary', data=empdata,color="black")
```

Out[9]:

```
<AxesSubplot:xlabel='Education', ylabel='Salary'>
```

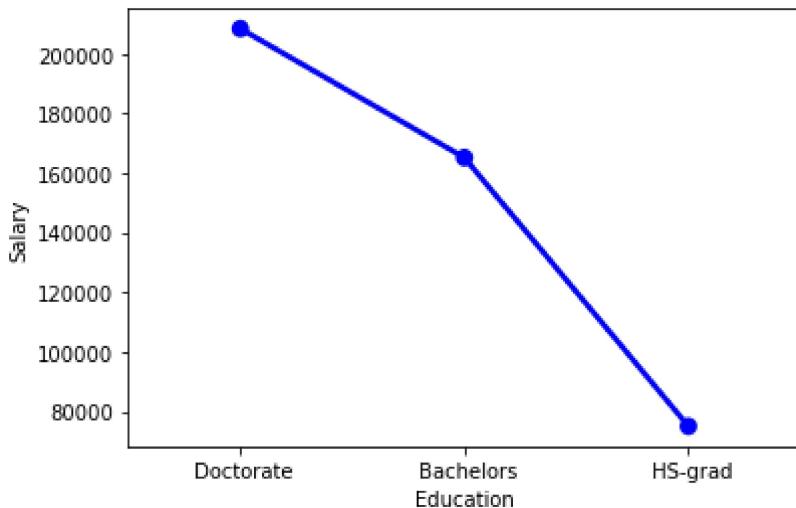


In [10]:

```
sns.pointplot(x='Education', y='Salary', data=empdata, ci=None,color="blue")
```

Out[10]:

```
<AxesSubplot:xlabel='Education', ylabel='Salary'>
```



Q1.3 Perform one-way ANOVA for variable Occupation with respect to the variable ‘Salary’. State whether the

null hypothesis is accepted or rejected based on the ANOVA results.

In [11]:

```
formula = 'Salary ~ C(Occupation)'  
model = ols(formula, empdata).fit()  
aov_table = anova_lm(model)  
print(aov_table)
```

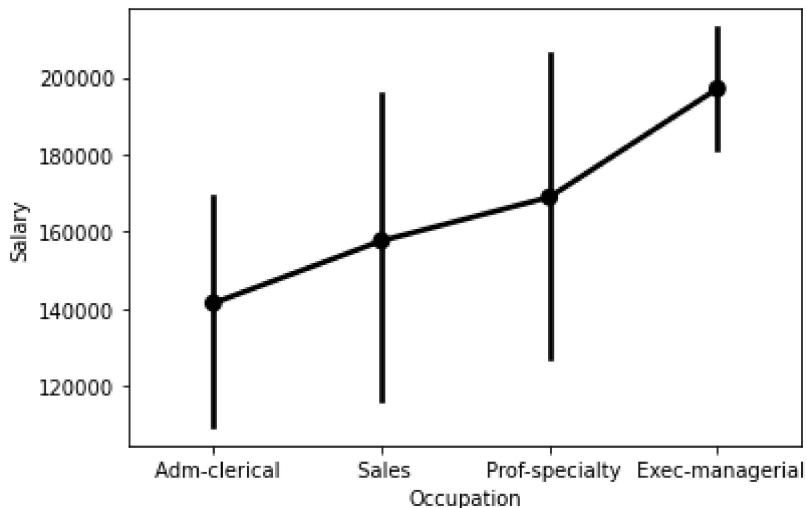
	df	sum_sq	mean_sq	F	PR(>F)
C(Occupation)	3.0	1.125878e+10	3.752928e+09	0.884144	0.458508
Residual	36.0	1.528092e+11	4.244701e+09	NaN	NaN

In [12]:

```
sns.pointplot(x='Occupation', y='Salary', data=empdata,color="black")
```

Out[12]:

```
<AxesSubplot:xlabel='Occupation', ylabel='Salary'>
```

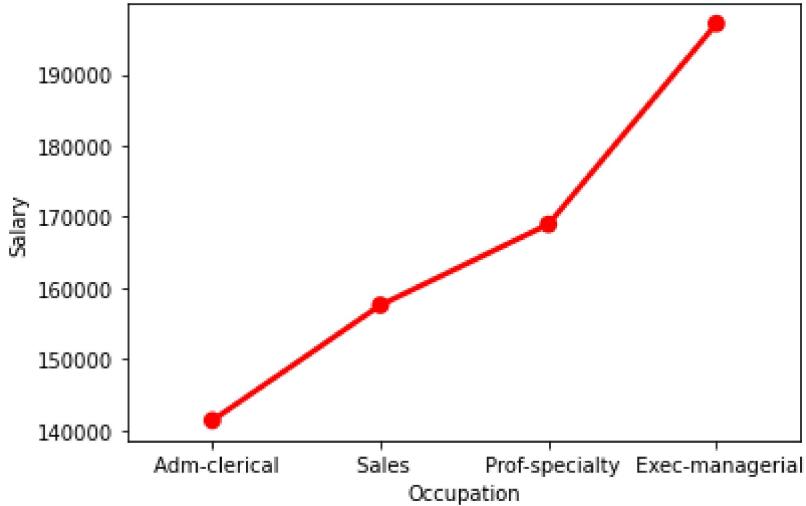


In [13]:

```
sns.pointplot(x='Occupation', y='Salary', data=empdata, ci=None,color="red")
```

Out[13]:

```
<AxesSubplot:xlabel='Occupation', ylabel='Salary'>
```



In []:

Q 1.4 If the null hypothesis is rejected in either (1.2) or in (1.3), find out which class means are significantly different. Interpret the result.

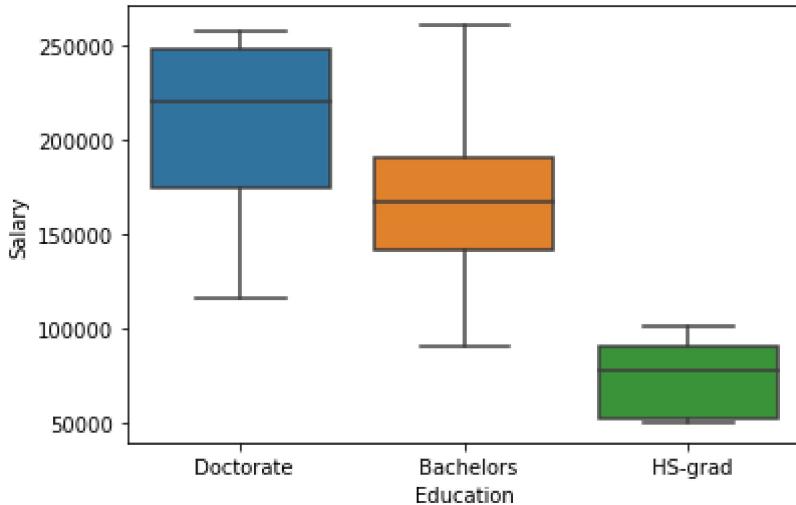
Ans:- 3 mean are differnet below mentioned boxplot

In [14]:

```
sns.boxplot(x='Education',y='Salary',data=empdata)
```

Out[14]:

```
<AxesSubplot:xlabel='Education', ylabel='Salary'>
```



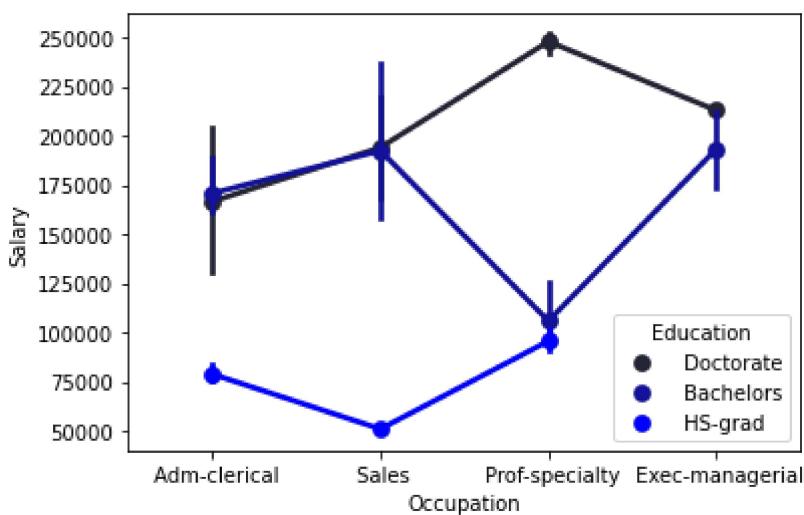
Q 1.5 What is the interaction between the two treatments? Analyze the effects of one variable on the other (Education and Occupation) with the help of an interaction plot.

In [15]:

```
sns.pointplot(x='Occupation', y='Salary',hue='Education',data=empdata,color="blue")
```

Out[15]:

```
<AxesSubplot:xlabel='Occupation', ylabel='Salary'>
```

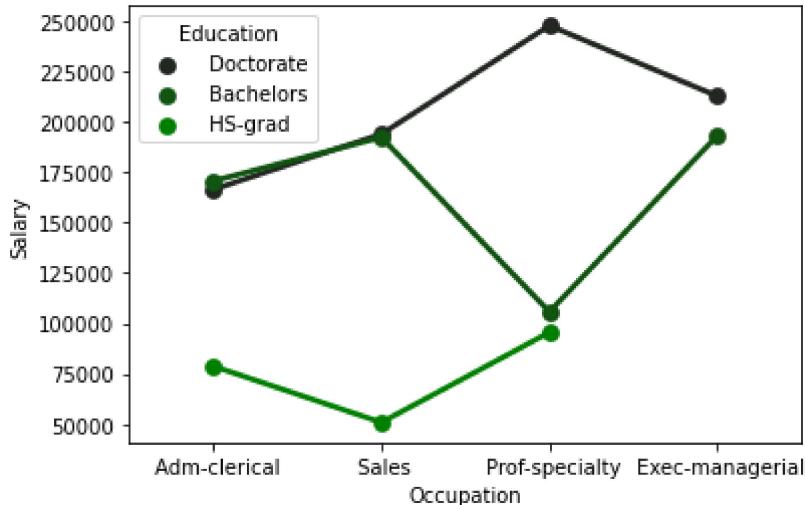


In [16]:

```
sns.pointplot(x='Occupation', y='Salary', hue='Education', data=empdata, ci=None, color="green")
```

Out[16]:

```
<AxesSubplot:xlabel='Occupation', ylabel='Salary'>
```



In []:

Q 1.6 Perform a two-way ANOVA based on the Education and Occupation (along with their interaction Education*Occupation) with the variable ‘Salary’. State the null and alternative hypotheses and state your results. How will you interpret this result?

In [17]:

```
formula = 'Salary ~ C(Education) + C(Occupation) + C(Education):C(Occupation)'
model = ols(formula, empdata).fit()
aov_table = anova_lm(model)
(aov_table)
```

Out[17]:

	df	sum_sq	mean_sq	F	PR(>F)
C(Education)	2.0	1.026955e+11	5.134773e+10	72.211958	5.466264e-12
C(Occupation)	3.0	5.519946e+09	1.839982e+09	2.587626	7.211580e-02
C(Education):C(Occupation)	6.0	3.634909e+10	6.058182e+09	8.519815	2.232500e-05
Residual	29.0	2.062102e+10	7.110697e+08	NaN	NaN

Q 1.7 Explain the business implications of performing ANOVA for this particular case study.

Ans:- mention in report

..... Problem 2.....

Education - Post 12th Standard, expected to do a Principal Component Analysis for this case study according to the instructions

In [18]:

```
edu_data=pd.read_csv("Education---Post+12th+Standard.csv")
```

Q2.1 Perform Exploratory Data Analysis [both univariate and multivariate analysis to be performed]. What insight do you draw from the EDA?

In [19]:

```
edu_data.info()
```

```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 777 entries, 0 to 776
Data columns (total 18 columns):
 #   Column      Non-Null Count  Dtype  
--- 
 0   Names        777 non-null    object  
 1   Apps         777 non-null    int64  
 2   Accept       777 non-null    int64  
 3   Enroll       777 non-null    int64  
 4   Top10perc    777 non-null    int64  
 5   Top25perc    777 non-null    int64  
 6   F.Undergrad  777 non-null    int64  
 7   P.Undergrad  777 non-null    int64  
 8   Outstate     777 non-null    int64  
 9   Room.Board   777 non-null    int64  
 10  Books        777 non-null    int64  
 11  Personal     777 non-null    int64  
 12  PhD          777 non-null    int64  
 13  Terminal     777 non-null    int64  
 14  S.F.Ratio    777 non-null    float64 
 15  perc.alumni  777 non-null    int64  
 16  Expend       777 non-null    int64  
 17  Grad.Rate    777 non-null    int64  
dtypes: float64(1), int64(16), object(1)
memory usage: 109.4+ KB
```

In [20]:

```
edu_data.shape
```

Out[20]:

```
(777, 18)
```

Check for missing value

In [21]:

```
edu_data.isnull().sum()
```

Out[21]:

```
Names          0
Apps          0
Accept         0
Enroll         0
Top10perc     0
Top25perc     0
F.Undergrad    0
P.Undergrad    0
Outstate       0
Room.Board     0
Books          0
Personal        0
PhD             0
Terminal        0
S.F.Ratio       0
perc.alumni     0
Expend          0
Grad.Rate       0
dtype: int64
```

In [22]:

```
edu_data.head()
```

Out[22]:

	Names	Apps	Accept	Enroll	Top10perc	Top25perc	F.Undergrad	P.Undergrad	Outstate
0	Abilene Christian University	1660	1232	721	23	52	2885	537	7440
1	Adelphi University	2186	1924	512	16	29	2683	1227	12280
2	Adrian College	1428	1097	336	22	50	1036	99	11250
3	Agnes Scott College	417	349	137	60	89	510	63	12960
4	Alaska Pacific University	193	146	55	16	44	249	869	7560

In [25]:

```
dups = edu_data.duplicated()  
print('Number of duplicate rows = %d' % (dups.sum()))
```

Number of duplicate rows = 0

Data visulitiation

Univariate – Analysis

In [26]:

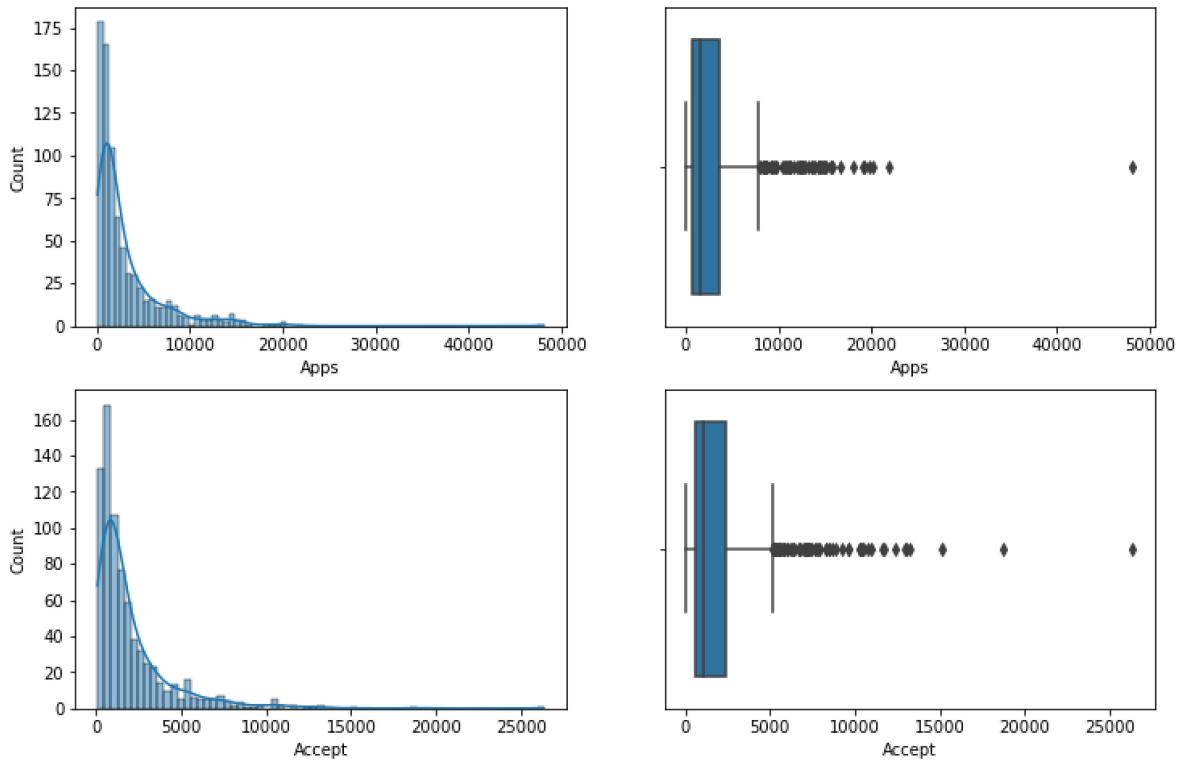
```
edu_data[['Apps', 'Accept']].describe()
```

Out[26]:

	Apps	Accept
count	777.000000	777.000000
mean	3001.638353	2018.804376
std	3870.201484	2451.113971
min	81.000000	72.000000
25%	776.000000	604.000000
50%	1558.000000	1110.000000
75%	3624.000000	2424.000000
max	48094.000000	26330.000000

In [27]:

```
fig,axes=plt.subplots(nrows=2,ncols=2)
fig.set_size_inches(12,8)
sns.histplot(edu_data['Apps'],kde=True,ax=axes[0][0])
sns.boxplot(x='Apps',data=edu_data,ax=axes[0][1])
sns.histplot(edu_data['Accept'],kde=True,ax=axes[1][0])
sns.boxplot(x='Accept',data=edu_data,ax=axes[1][1])
plt.show()
```



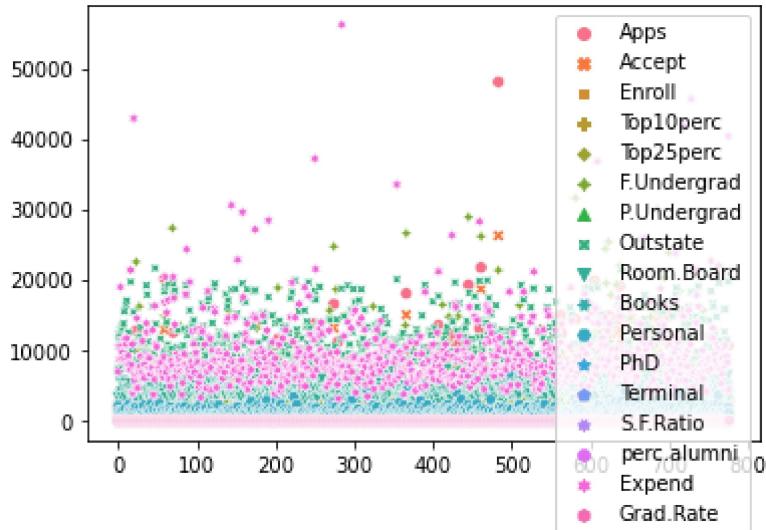
multivariate – analysis

In [28]:

```
sns.scatterplot(data=edu_data)
```

Out[28]:

<AxesSubplot:>

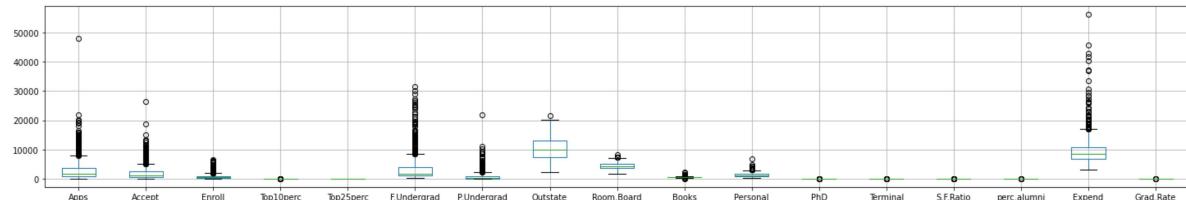


In [29]:

```
edu_data.boxplot(figsize=(25,4))
```

Out[29]:

<AxesSubplot:>

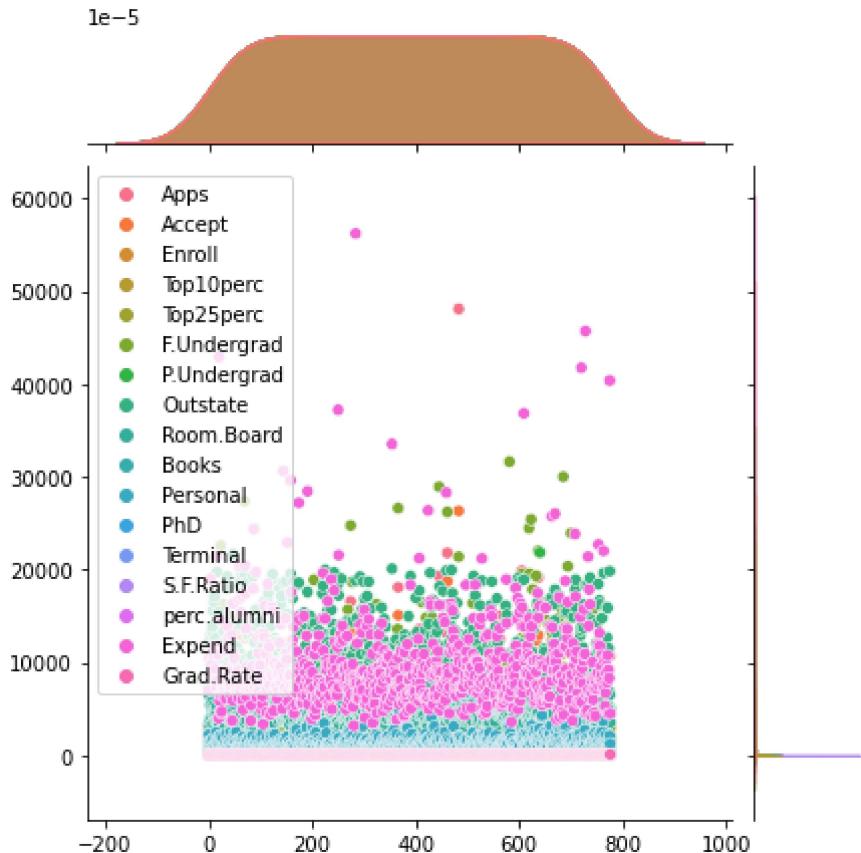


In [30]:

```
sns.jointplot(data=edu_data)
```

Out[30]:

```
<seaborn.axisgrid.JointGrid at 0x12fc94580>
```

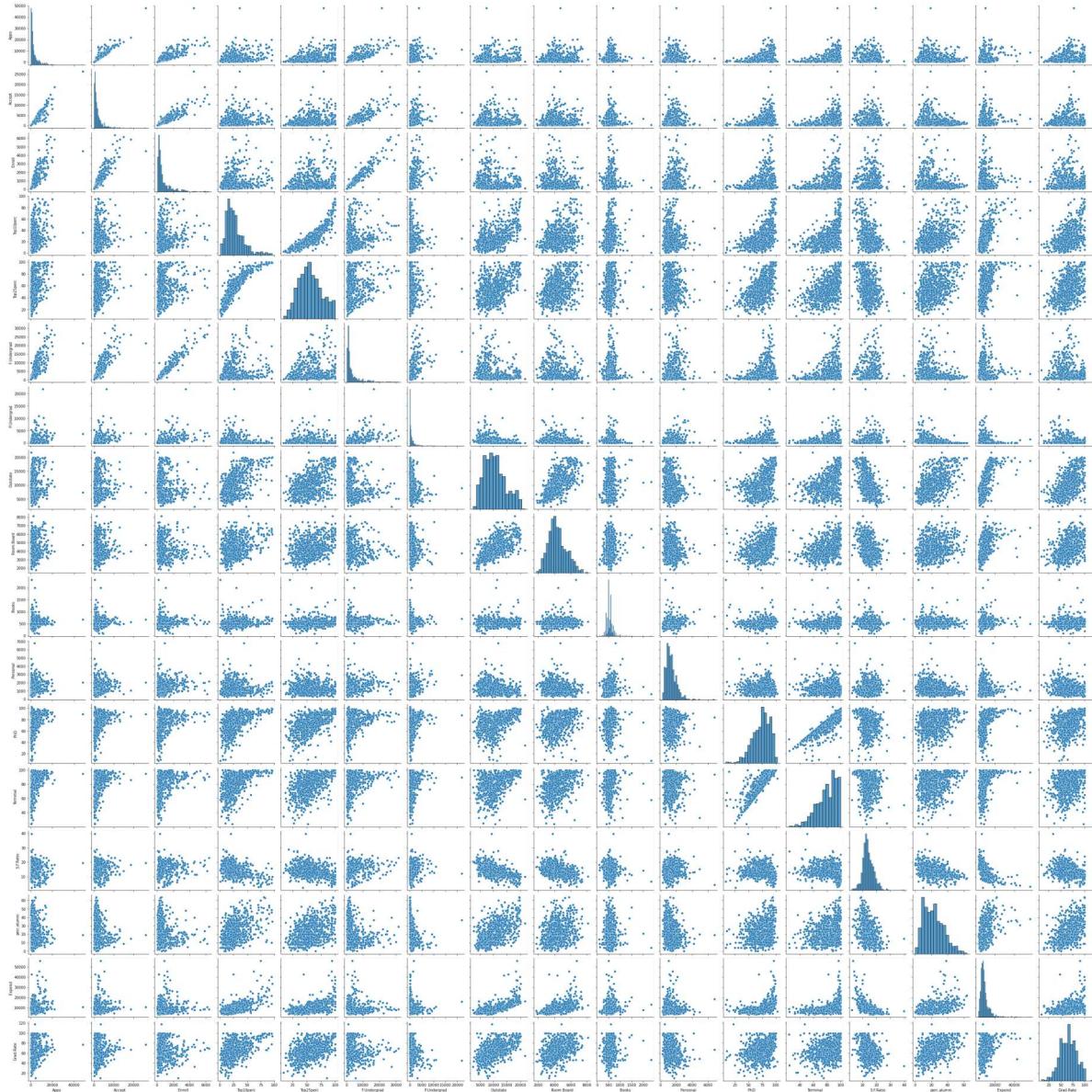


In [31]:

```
sns.pairplot(edu_data)
```

Out[31]:

```
<seaborn.axisgrid.PairGrid at 0x12fcc073b50>
```



In [32]:

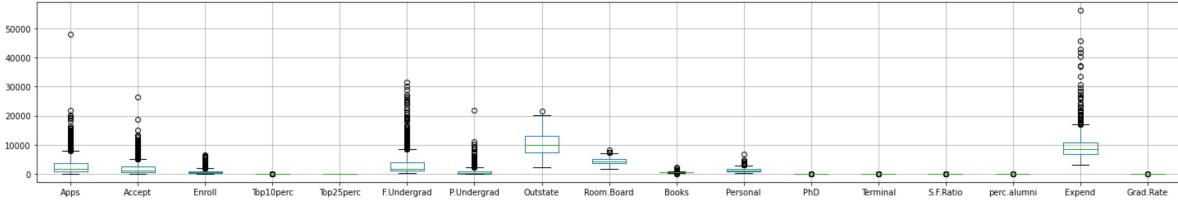
```
corr=edu_data.corr()
```


In [39]:

```
edu_data.boxplot(figsize=(25,4))
```

Out[39]:

<AxesSubplot:>



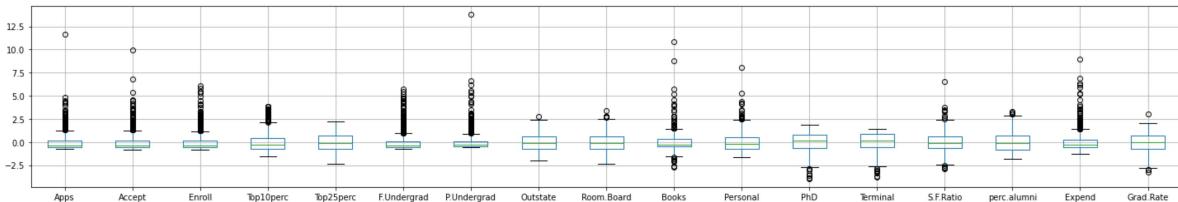
After Scaling

In [40]:

```
data_scaled.boxplot(figsize=(25,4))
```

Out[40]:

<AxesSubplot:>



pramer is reduest and all mean value is near by zero ...

Q 2.5 Extract the eigenvalues and eigenvectors.[print both]

In [42]:

```
print('\n Eigen Values \n%s', pd.DataFrame(eig_vals))
```

```
Eigen Values
%s          0
0    2.250938e+02+0.000000e+00j
1    8.626531e+01+0.000000e+00j
2    5.341962e+01+0.000000e+00j
3    4.858358e+01+0.000000e+00j
4    4.525109e+01+0.000000e+00j
...
772   8.482920e-18-5.859576e-17j
773   3.836886e-17+0.000000e+00j
774  -2.116932e-17+2.654029e-17j
775  -2.116932e-17-2.654029e-17j
776  -2.688490e-17+0.000000e+00j
```

[777 rows x 1 columns]

Q 2.6 Perform PCA and export the data of the Principal Component (eigenvectors) into a data frame with the original features

In [46]:

```
tot = sum(eig_vals)
var_exp = [( i /tot ) * 100 for i in sorted(eig_vals, reverse=True)]
cum_var_exp = np.cumsum(var_exp)
print("Cumulative Variance Explained",pd.DataFrame(cum_var_exp))
```

Cumulative Variance Explained

```
0    36.165187+0.000000j
1    50.025194+0.000000j
2    58.607976+0.000000j
3    66.413765+0.000000j
4    73.684131+0.000000j
..
772  100.000000+0.000000j
773  100.000000+0.000000j
774  100.000000+0.000000j
775  100.000000+0.000000j
776  100.000000+0.000000j
```

[777 rows x 1 columns]

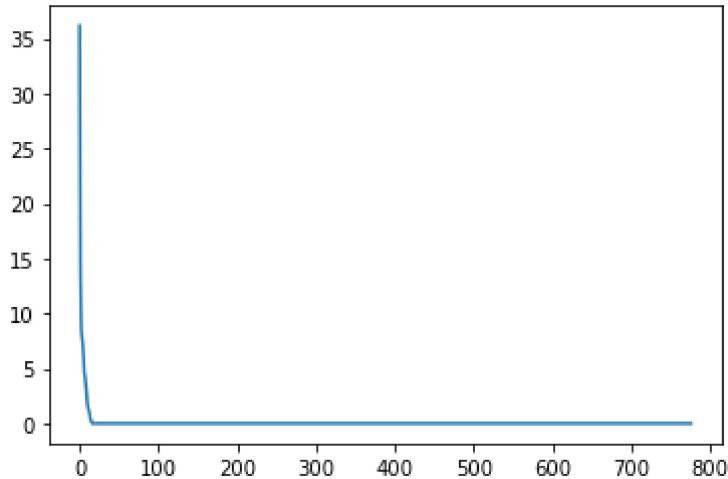
In [47]:

```
plt.plot(var_exp)
```

```
C:\Users\rahul\anaconda3\lib\site-packages\numpy\core\_asarray.py:83: ComplexWarning: Casting complex values to real discards the imaginary part
  return array(a, dtype, copy=False, order=order)
```

Out[47]:

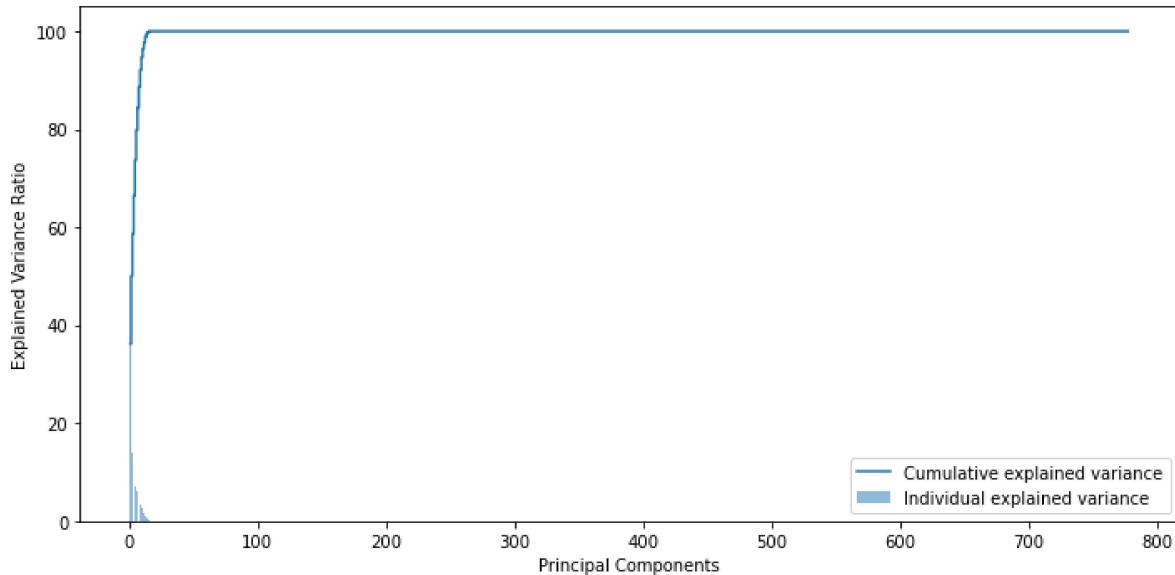
[<matplotlib.lines.Line2D at 0x12fdea57040>]



In [48]:

```
# Ploting
plt.figure(figsize=(10 , 5))
plt.bar(range(1, eig_vals.size + 1), var_exp, alpha = 0.5, align = 'center', label = 'Individual explained variance')
plt.step(range(1, eig_vals.size + 1), cum_var_exp, where='mid', label = 'Cumulative explained variance')
plt.ylabel('Explained Variance Ratio')
plt.xlabel('Principal Components')
plt.legend(loc = 'best')
plt.tight_layout()
plt.show()
```

C:\Users\rahul\anaconda3\lib\site-packages\numpy\core_asarray.py:83: ComplexWarning: Casting complex values to real discards the imaginary part
return array(a, dtype, copy=False, order=order)



Q 2.9 Explain the business implication of using the Principal Component Analysis for this case study. How may PCs help in the further analysis? [Hint: Write Interpretations of the Principal Components Obtained]

ANS :-Mention in reportthanks

In []: