

Data Visualization for Bank Churn Data

```
In [49]: 1 import pandas as pd
2 import numpy as np
3
4 import matplotlib.pyplot as plt
5 import seaborn as sb
6
7 from scipy import stats
```

```
In [2]: 1 df=pd.read_csv("train.csv")
```

```
In [3]: 1 df.head()
```

Out[3]:

	id	CustomerId	Surname	CreditScore	Geography	Gender	Age	Tenure	Balance	NumOff
0	15674932	Okwudilichukwu		668	France	Male	33.0	3	0.00	
1	15749177	Okwudilolisa		627	France	Male	33.0	1	0.00	
2	15694510	Hsueh		678	France	Male	40.0	10	0.00	
3	15741417	Kao		581	France	Male	34.0	2	148882.54	
4	15766172	Chiemenam		716	Spain	Male	33.0	5	0.00	

```
In [4]: 1 df.info()
```

```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 165034 entries, 0 to 165033
Data columns (total 14 columns):
 #  Column            Non-Null Count  Dtype  
--- 
 0  id                165034 non-null   int64  
 1  CustomerId        165034 non-null   int64  
 2  Surname           165034 non-null   object  
 3  CreditScore       165034 non-null   int64  
 4  Geography          165034 non-null   object  
 5  Gender             165034 non-null   object  
 6  Age                165034 non-null   float64 
 7  Tenure             165034 non-null   int64  
 8  Balance            165034 non-null   float64 
 9  NumOfProducts      165034 non-null   int64  
 10  HasCrCard         165034 non-null   float64 
 11  IsActiveMember    165034 non-null   float64 
 12  EstimatedSalary   165034 non-null   float64 
 13  Exited            165034 non-null   int64  
dtypes: float64(5), int64(6), object(3)
memory usage: 17.6+ MB
```

In [5]: 1 df.describe()

Out[5]:

	id	CustomerId	CreditScore	Age	Tenure	Balance	IsChurn
count	165034.0000	1.650340e+05	165034.000000	165034.000000	165034.000000	165034.000000	165034.000000
mean	82516.5000	1.569201e+07	656.454373	38.125888	5.020353	55478.086689	0.000000
std	47641.3565	7.139782e+04	80.103340	8.867205	2.806159	62817.663278	0.000000
min	0.0000	1.556570e+07	350.000000	18.000000	0.000000	0.000000	0.000000
25%	41258.2500	1.563314e+07	597.000000	32.000000	3.000000	0.000000	0.000000
50%	82516.5000	1.569017e+07	659.000000	37.000000	5.000000	0.000000	0.000000
75%	123774.7500	1.575682e+07	710.000000	42.000000	7.000000	119939.517500	0.000000
max	165033.0000	1.581569e+07	850.000000	92.000000	10.000000	250898.090000	0.000000

In [6]: 1 df_num = df.copy()

In [7]: 1 cat_var = ['Surname', 'Geography', 'Gender']
2 num_var = df.columns.difference(cat_var)

In [8]: 1 from sklearn.preprocessing import LabelEncoder

In [9]: 1 for features in cat_var:
2 le=LabelEncoder()
3 df_num[features] = le.fit_transform(df[features])

In [10]:

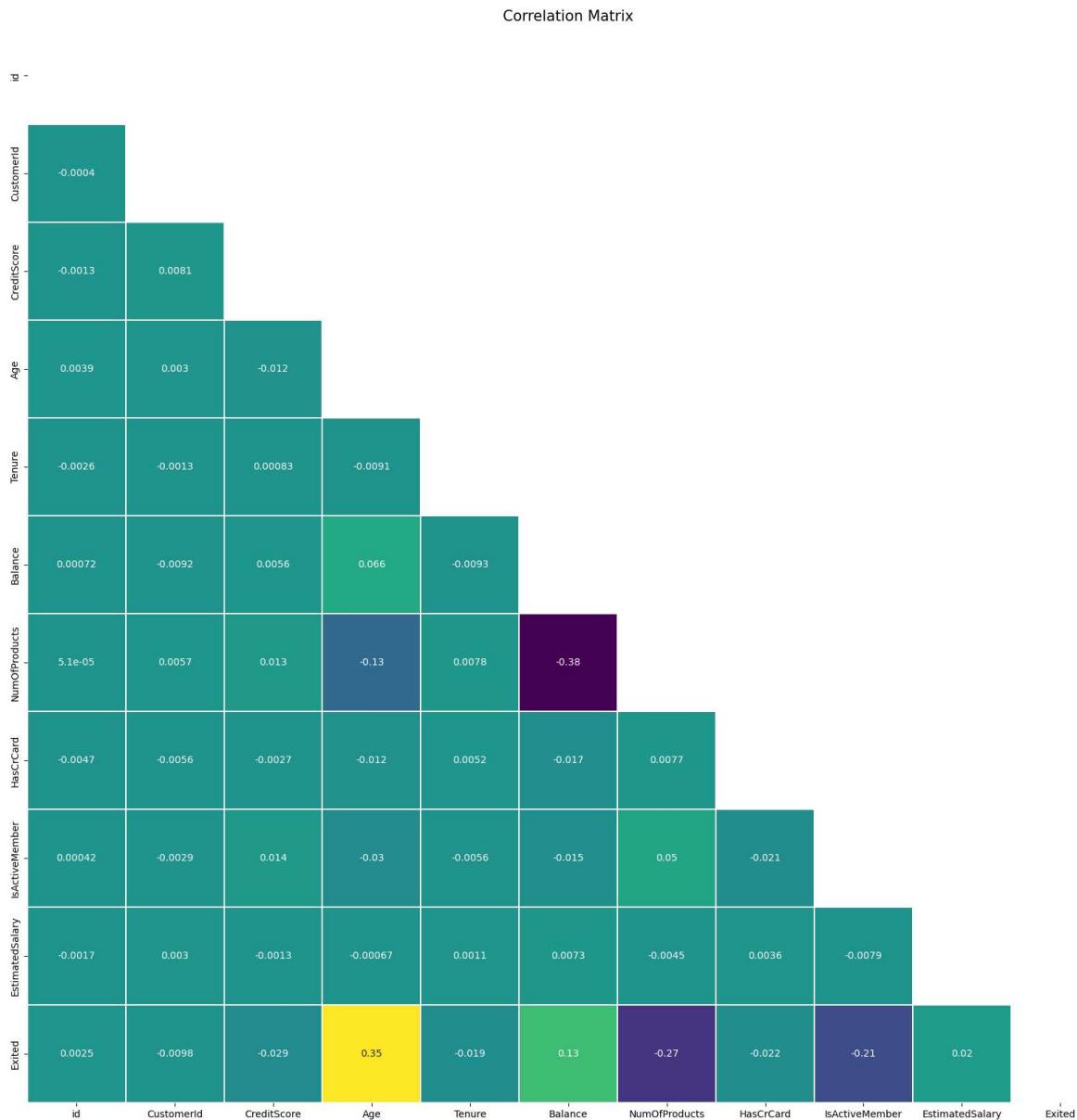
```

1 fig,ax= plt.subplots(figsize=(20,20))
2 corr_matrix = df.corr(method='spearman')
3 sb.heatmap(corr_matrix,annot=True,linewidths=0.1,cbar=False, ax=ax, mask=numpy.triu(corr_matrix))
4
5 ax.set_title("Correlation Matrix",fontsize=15)
6 plt.show()

```

C:\Users\deepa\AppData\Local\Temp\ipykernel_25308\1309991678.py:2: 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.

```
corr_matrix = df.corr(method='spearman')
```



The Correlation Matrix Shows:

- a strong correlation between Age and Exited
- a strong correlation between Balance and NumOfProducts
- a strong correlation between NumOfProducts and Exited

- a strong correlation between IsActiveMemeber and Exited
- a small correlation between NumofProducts and Age
- a tiny correlation between IsActiveMember and NumOfProducts

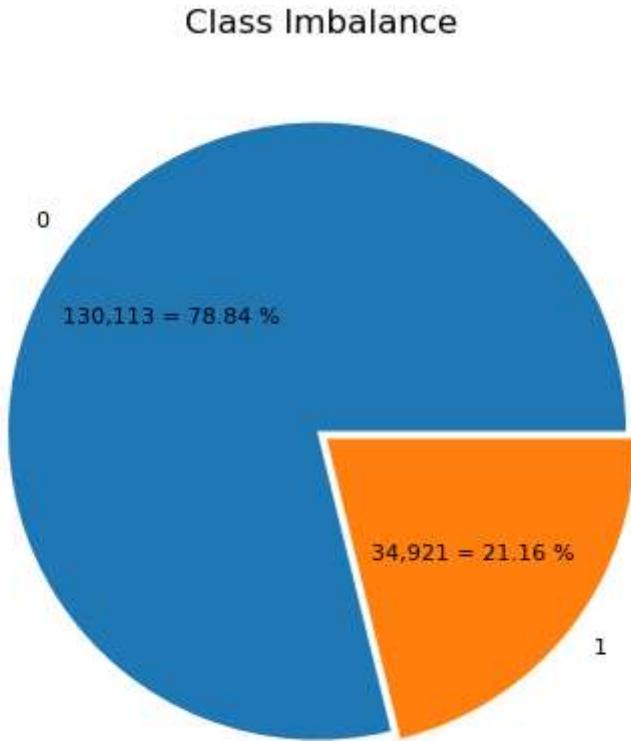
Let's explore these relations

Class Imbalance

Here we check the imbalance in data for variabel "Exited" using a pie chart

```
In [11]: 1 targets=df["Exited"].unique()
2 classes=[df[(df["Exited"] == target)][ "id"].count() for target in targets]
```

```
In [12]: 1 fix,ax = plt.subplots(figsize=(5,5))
2 plt.pie(classes, labels=targets, explode=[0.02]*2, autopct= lambda x: "{:,0.0f} \n" + str(x) + "%")
3 ax.set_title('Class Imbalance')
4 plt.show()
5
```



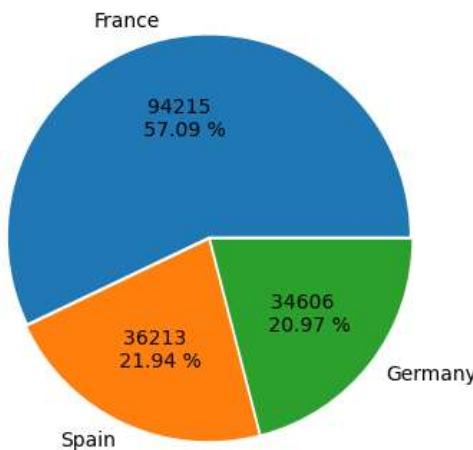
Here we can see that the 78.84% remains a customer of the bank where rest 21.16% exits from the bank. This shows that data is highly imbalanced.

Churn by Geography

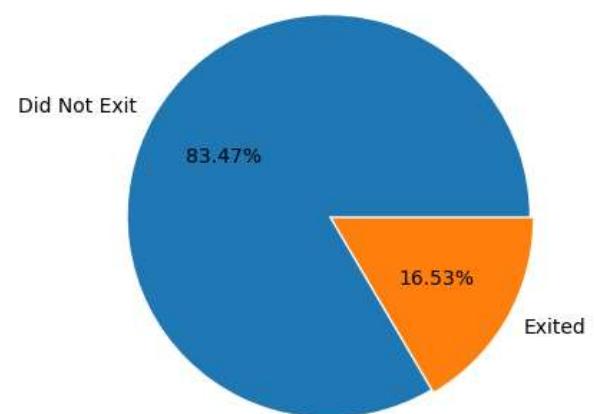
In [13]:

```
1 geo_target = df['Geography'].unique()
2 fig, ax = plt.subplots(nrows=2, ncols=2, figsize=(10,10))
3
4 geo_count = [df[(df['Geography'] == target)]['id'].count() for target in g
5
6
7 ax = ax.flatten()
8
9 _ = ax[0].pie(geo_count, labels=geo_target, autopct = lambda x: "{:.0f} \n"
10 _ = ax[0].set_title("Population by Geography")
11 for i, target in enumerate (geo_target):
12
13     classes = [df[(df['Geography']==target) & (df['Exited']==0)]['id'].cou
14
15     label = ['Did Not Exit', 'Exited']
16
17     _ = ax[i+1].pie(classes, labels = label, autopct = "%2f%%", explode= [0
18     _ = ax[i+1].set_title(target)
```

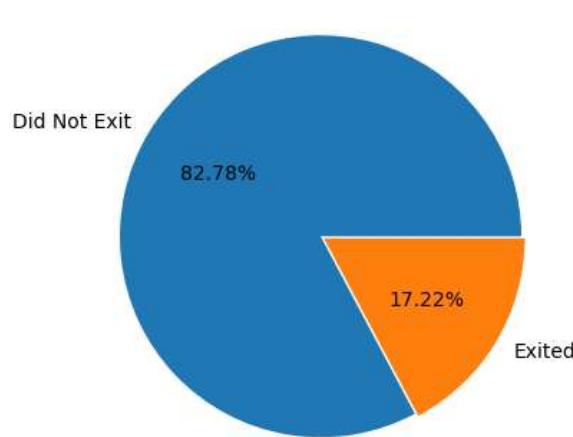
Population by Geography



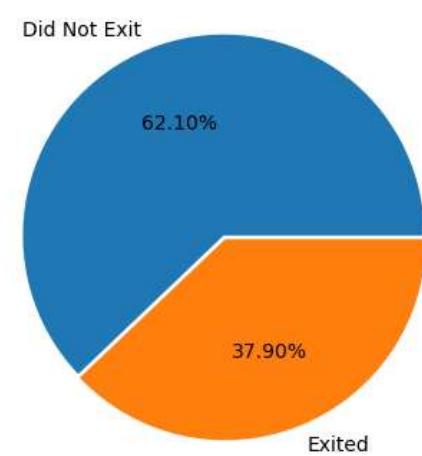
France



Spain



Germany



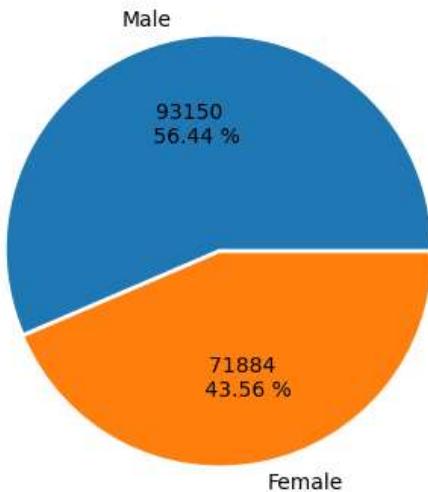
We see that Most of the people in germany exits, this can be an important feature for predictions

Gender distribution across different countries

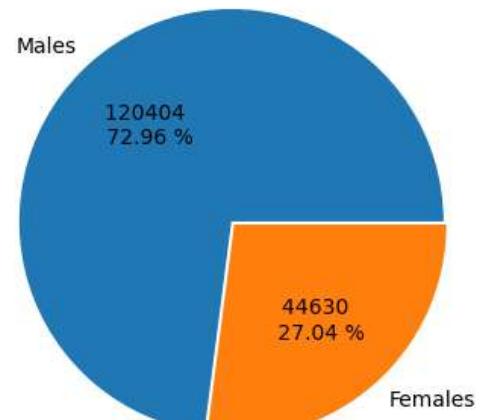
In [14]:

```
1 gen_target = df['Gender'].unique()
2 fig, ax = plt.subplots(nrows=2, ncols=2, figsize=(10,10))
3
4 gen_count = [df[(df['Gender'] == target)]['id'].count() for target in gen_
5
6
7 ax = ax.flatten()
8
9 _ = ax[0].pie(gen_count, labels=gen_target, autopct = lambda x:"{:.0f} \n"
10 _ = ax[0].set_title("Population by Geography")
11
12 for i, target in enumerate (geo_target):
13
14     classes = [df[(df['Geography']==target) & (df['Gender']=='Male')]['id'
15
16     label = ['Males','Females']
17
18     _ = ax[i+1].pie(classes,labels = label, autopct = lambda x:"{:.0f} \n"
19     _ = ax[i+1].set_title(target)
```

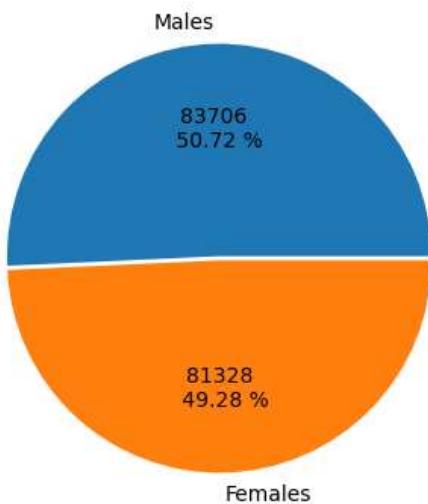
Population by Geography



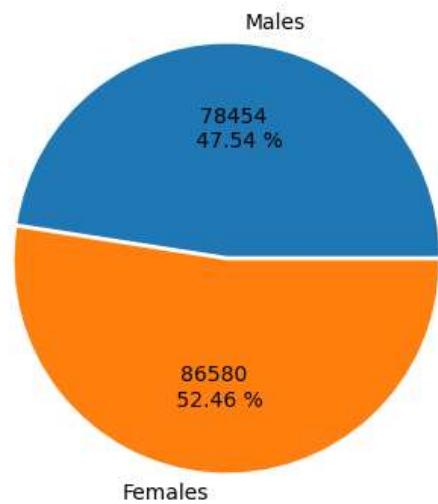
France



Spain

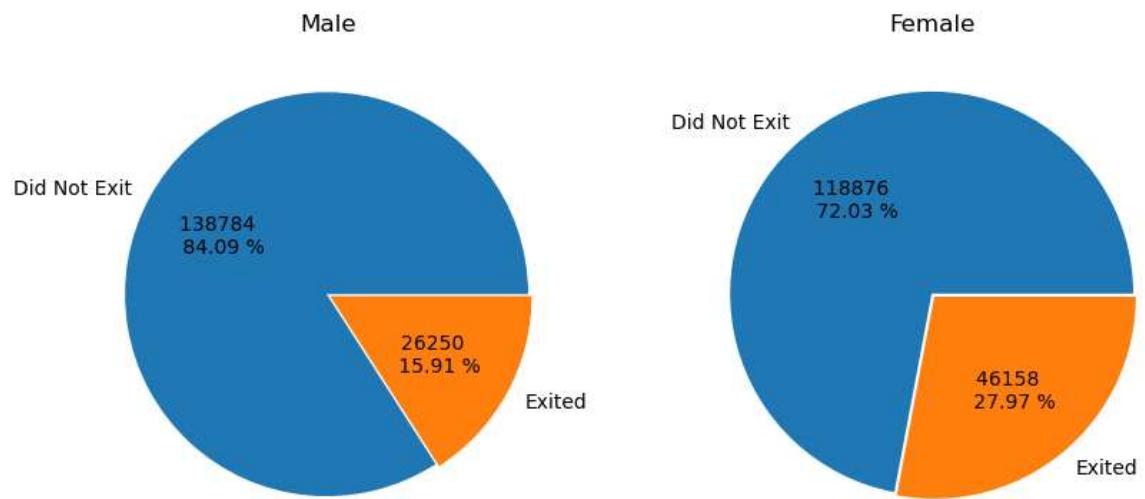


Germany



- The percentage of Male population (56%) is slightly larger than the Female population (44%)
- For France we see a high percentage (73%) of Males compares to only 27% of females. Whereas for Spain and Germany, the gender distribution is more or less equal.

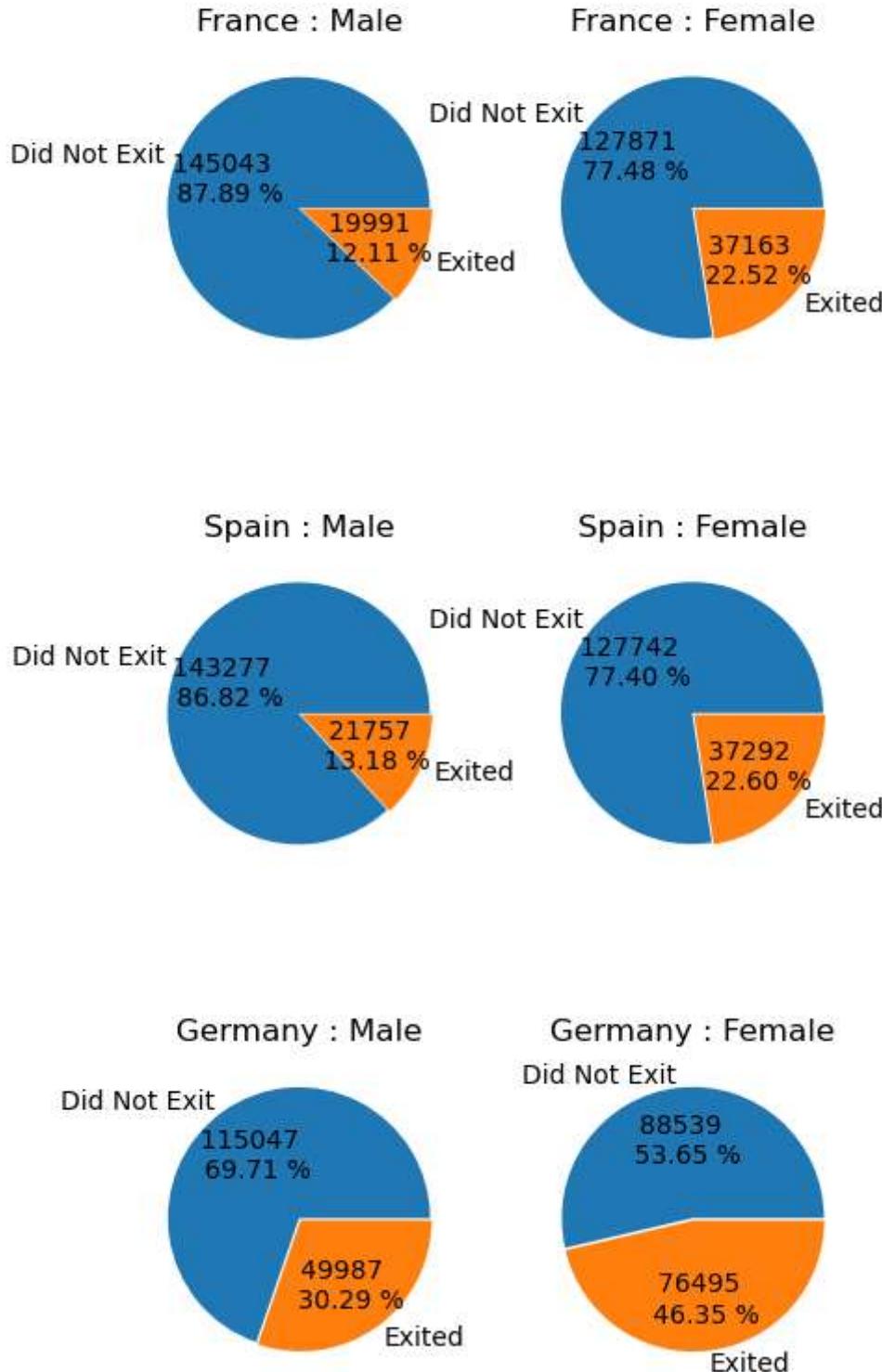
```
In [15]: 1 gen_target = df['Gender'].unique()
2
3 fix, ax = plt.subplots(nrows=1, ncols=2, figsize=(10,6))
4
5 gen_count = [df[(df["Gender"]==target)]['id'].count() for target in gen_ta
6
7 ax=ax.flatten()
8
9
10 for i, target in enumerate (gen_target):
11
12     classes = [df[(df['Gender']==target) & (df['Exited']==0)]['id'].count()
13
14     label = ['Did Not Exit', 'Exited']
15
16     _= ax[i].pie(classes, labels = label, autopct = lambda x: "{:.0f} \n {:
17     _= ax[i].set_title(target)
```



The above pie chart shows that larger number of females tend to terminate bank's services than Males. Approx 28% of the female population exits whereas only 15.9% of the males exits the bank services

```
In [16]: 1 gen_target = df['Gender'].unique()
2 geo_target = df["Geography"].unique()
3
4 fix, ax = plt.subplots(nrows=3, ncols=2, figsize=(5,10))
5
6 gen_count = [df[(df["Gender"]==target)]['id'].count() for target in gen_ta
7
8 ax=ax.flatten()
9
10 #_ = ax[0].pie(gen_count, labels=gen_target, autopct = lambda x:"{:.0f} \n
11 #_ = ax[0].set_title("Data by Gender")
12
13 temp=0
14 for j in range(3):
15     for i, target in enumerate (gen_target):
16
17         classes = [df[(df['Gender']==target) & (df['Exited']==0) & (df['Ge
18
19         label = ['Did Not Exit','Exited']
20
21         _ = ax[temp].pie(classes,labels = label, autopct = lambda x:"{:.0f
22         _ = ax[temp].set_title("{} : {}".format(geo_target[j],target))
23         temp+=1
24
25 plt.suptitle("Churn by Geogaphy and Gender")
```

Churn by Geography and Gender



- We see the same trend of Higher number of exits in Germany as well as higher number of exits by females compared to males.
- One thing to notice here is percentage of females exiting in Germany is only about twice as that of percentage of females exiting in any other countries, percentage of males exiting in

germany is 2.3 times higher than percentage of males exiting in other countries

In [50]:

```

1  def dens_hist_exit(var,bins=100):
2      fig, ax = plt.subplots(nrows=2,ncols=1,figsize=(8,6),gridspec_kw={'hspace':0.5})
3      ax=ax.flatten()
4      sb.set_style("whitegrid")
5      _= sb.kdeplot(df[df["Exited"]==0][var], fill = True, label = "Not Exited",)
6      _= sb.kdeplot(df[df["Exited"]==1][var], fill = True, label = "Exited",)
7      _= ax[0].set_title("{} Distribution".format(var))
8      _= ax[0].set_xlabel("")
9      _= ax[0].legend()
10
11     _= ax[1].hist(df[df["Exited"]==0][var], fill = True, label = "Not Exited",)
12     _= ax[1].hist(df[df["Exited"]==1][var], fill = True, label = "Exited",)
13     _= ax[1].legend()
14     #_= ax[1].set_title("{} Distribution".format(var))
15     _= ax[1].set_xlabel(var)
16     _= ax[1].set_ylabel("Frequency")
17
18 def dens_hist_gen(var,var2,bins=100):
19
20     cond=df[var2].unique()
21
22
23     fig, ax = plt.subplots(nrows=2,ncols=1,figsize=(8,6),gridspec_kw={'hspace':0.5})
24     ax=ax.flatten()
25
26     for i in range(len(cond)):
27         #print(cond[i])
28         _= sb.kdeplot(df[df[var2]==cond[i]][var], fill = True, label = cond[i],)
29         _= ax[1].hist(df[df[var2]==cond[i]][var], fill = True, label = cond[i],)
30
31     #_= sb.kdeplot(df[df[var2]==cond[i]][var], fill = True, label = var2,)
32     _= ax[0].set_title("{} Distribution by {}".format(var,var2))
33     _= ax[0].set_xlabel("")
34     _= ax[0].legend()
35
36
37     #_= ax[1].hist(df[df[var2]==cond[i]][var], fill = True, label = var2,)
38     _= ax[1].legend()
39     #_= ax[1].set_title("{} Distribution by {}".format(var, var2))
40     _= ax[1].set_xlabel(var)
41     _= ax[1].set_ylabel("Frequency")
42

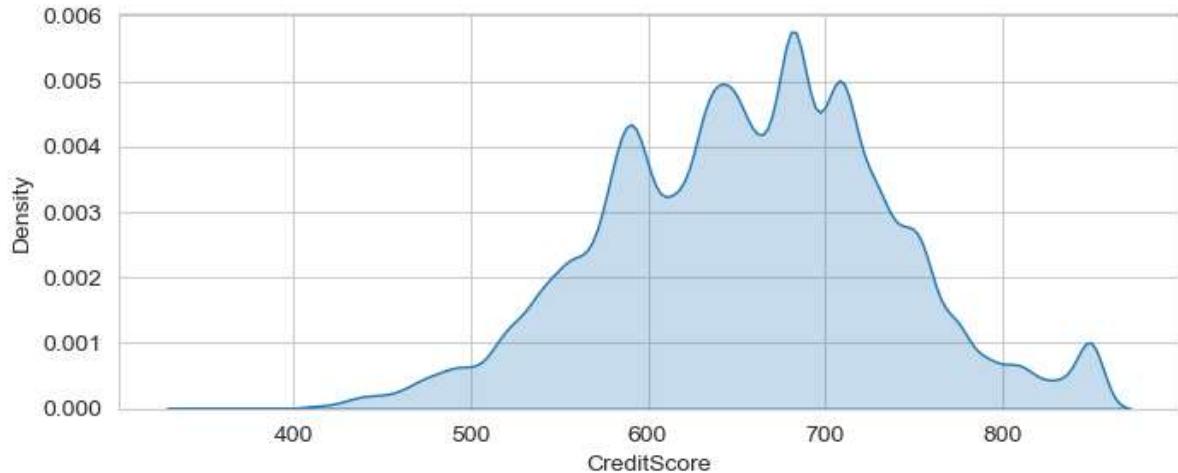
```

In []:

1

Credit Score

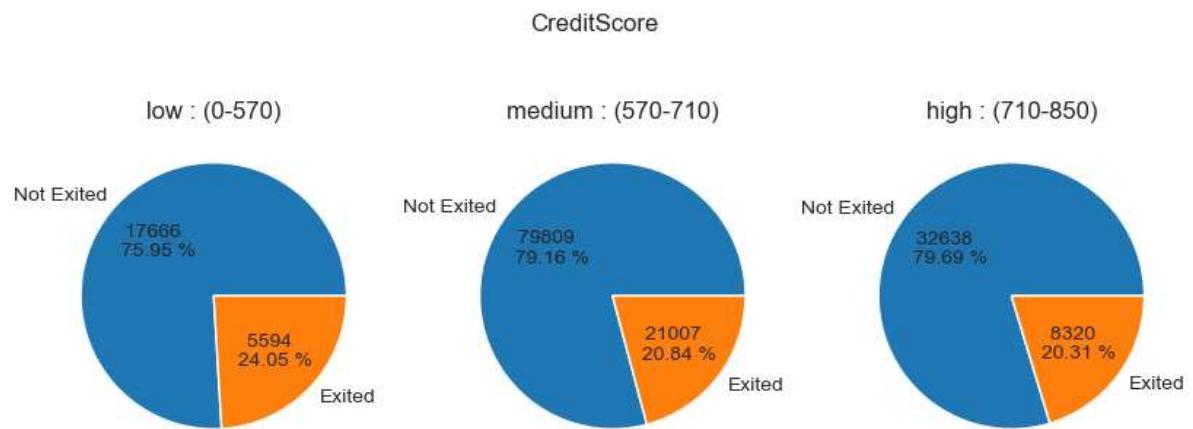
```
In [51]: 1 fig, ax = plt.subplots(figsize=(8,3))
2 sb.set_style("darkgrid")
3 _ = sb.kdeplot(df["CreditScore"], fill=True)
```



We see that very small population has creditscore in the range 0-580, a good portion of population lies within the creditscore of 580-700 and very few people lies in the very high credit score region. Based on this information we can categorize people based on low, medium and high Creditscore

```
In [52]: 1 def credit_score(score):
2     if score<=570:
3         return 'low'
4     elif score<=710:
5         return 'medium'
6     else:
7         return 'high'
8
9 df["CrScoreCat"] = df["CreditScore"].apply(credit_score)
10
```

```
In [53]: 1 fig, ax = plt.subplots(nrows=1,ncols=3,figsize=(10,4))
2 ax=ax.flatten()
3 Cr_cat = ['low','medium','high']
4 Cr_score = ['0-570','570-710','710-850']
5
6 for i in range(3):
7     target = [df[(df["CrScoreCat"]==Cr_cat[i]) & (df["Exited"]==0)]["id"]].
8     label = ['Not Exited','Exited']
9     _=ax[i].pie(target,labels=label,autopct=lambda x:"{:.0f} \n {:.2f} %"
10     _=ax[i].set_title("{} : ({})".format(Cr_cat[i],Cr_score[i]))
11     _=plt.suptitle("CreditScore")
12
```

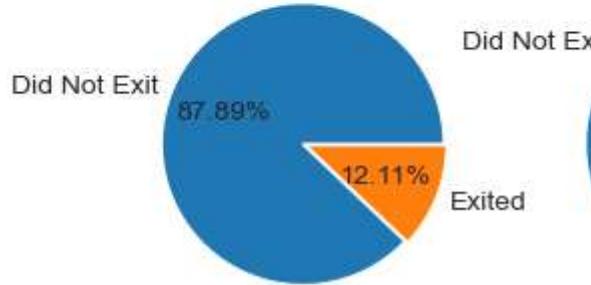


The percentage of people exited is more or less same throughout different credit score groups with percentage of exited population in low credit-score being slightly higher than other credit score groups

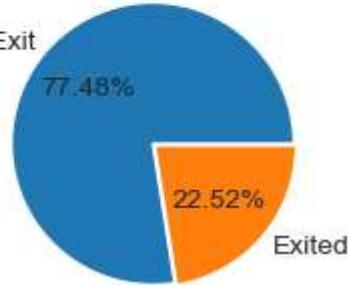
```
In [54]: 1 gen_target = df['Gender'].unique()
2 geo_target = df["Geography"].unique()
3
4 fix, ax = plt.subplots(nrows=3, ncols=2, figsize=(5,10))
5
6 gen_count = [df[(df["Gender"]==target)]['id'].count() for target in gen_ta
7
8 ax=ax.flatten()
9
10 #_ = ax[0].pie(gen_count, labels=gen_target, autopct = lambda x:"{:.0f} \n
11 #_ = ax[0].set_title("Data by Gender")
12
13 temp=0
14 for j in range(3):
15     for i, target in enumerate (gen_target):
16
17         classes = [df[(df['Gender']==target) & (df['Exited']==0) & (df['Ge
18
19         label = ['Did Not Exit','Exited']
20
21         _ = ax[temp].pie(classes,labels = label, autopct = "%.2f%%", explode
22         _ = ax[temp].set_title("{} : {}".format(geo_target[j],target))
23         temp+=1
24
25 plt.suptitle("Churn by Geogarphy and Gender")
```

Churn by Geography and Gender

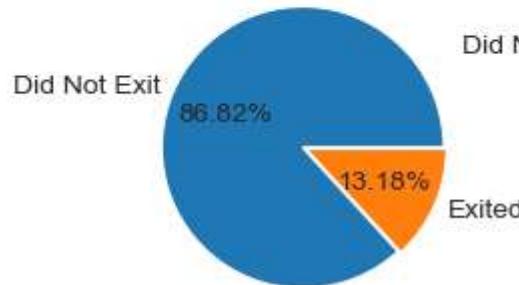
France : Male



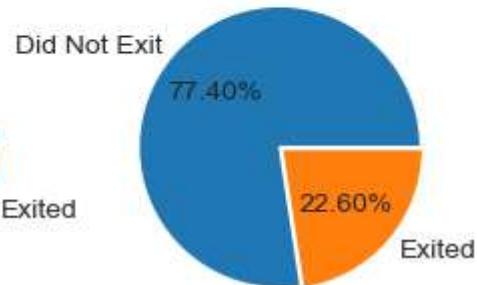
France : Female



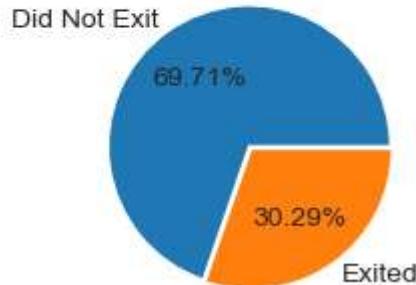
Spain : Male



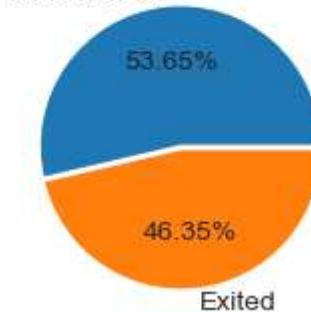
Spain : Female



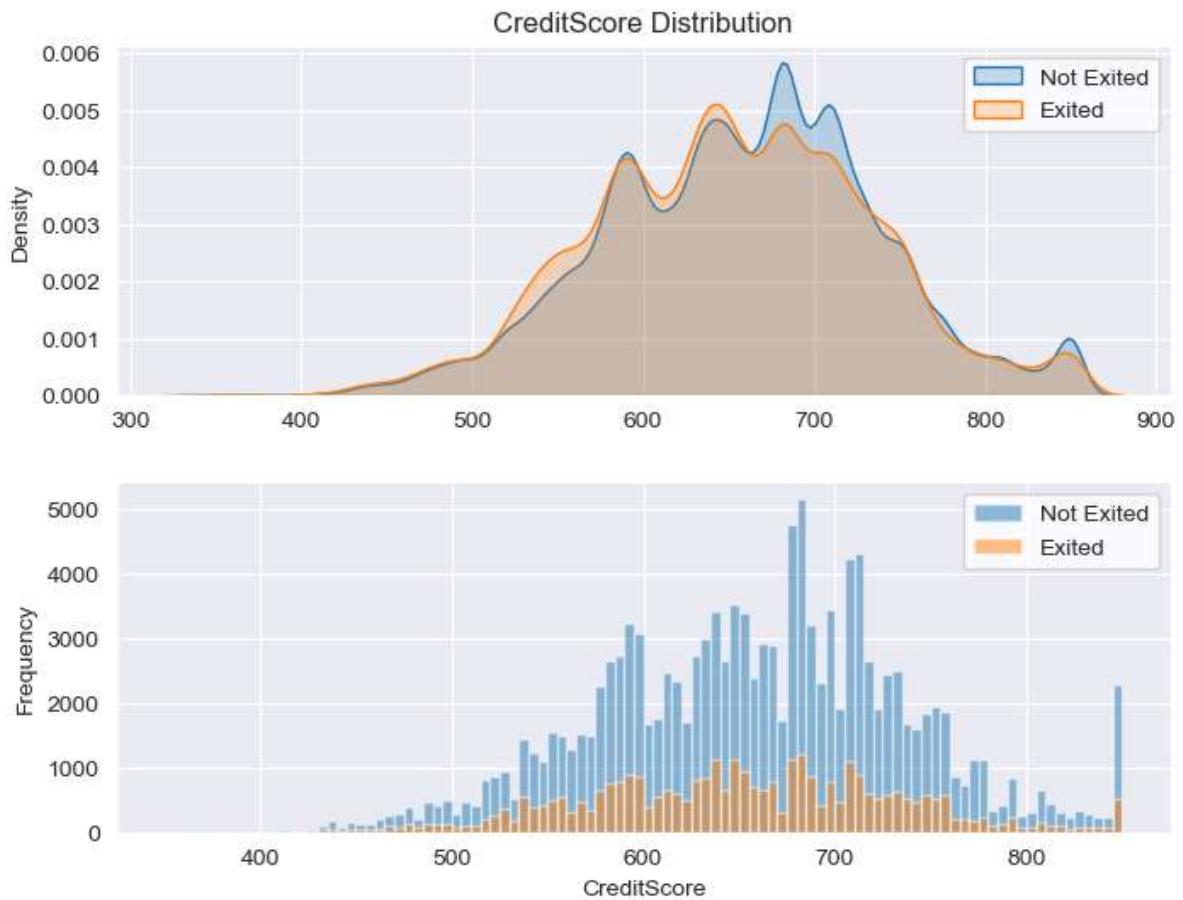
Germany : Male



Germany : Female



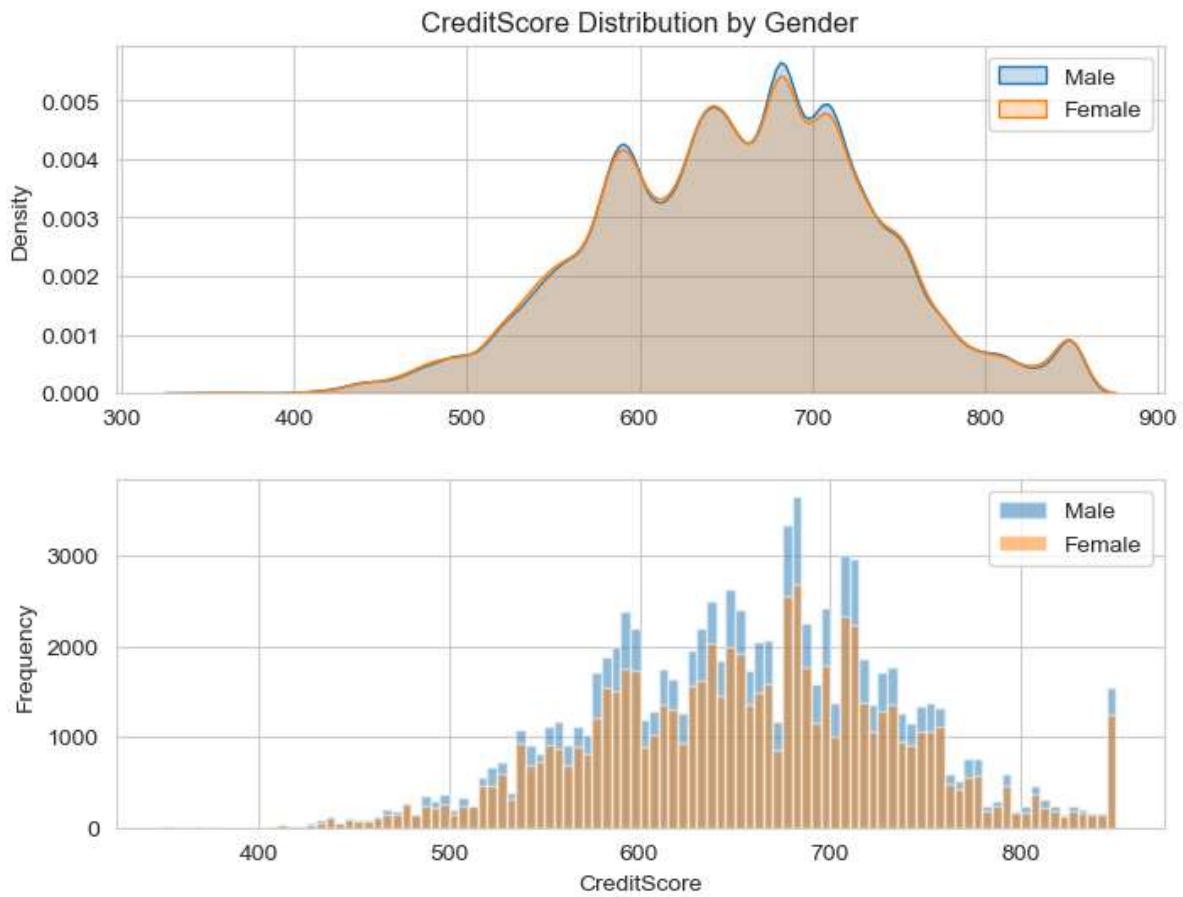
```
In [55]: 1 dens_hist_exit("CreditScore")
```



People having credit score in the range 660 - 720 are more likely to not exit

Credit Score By Gender

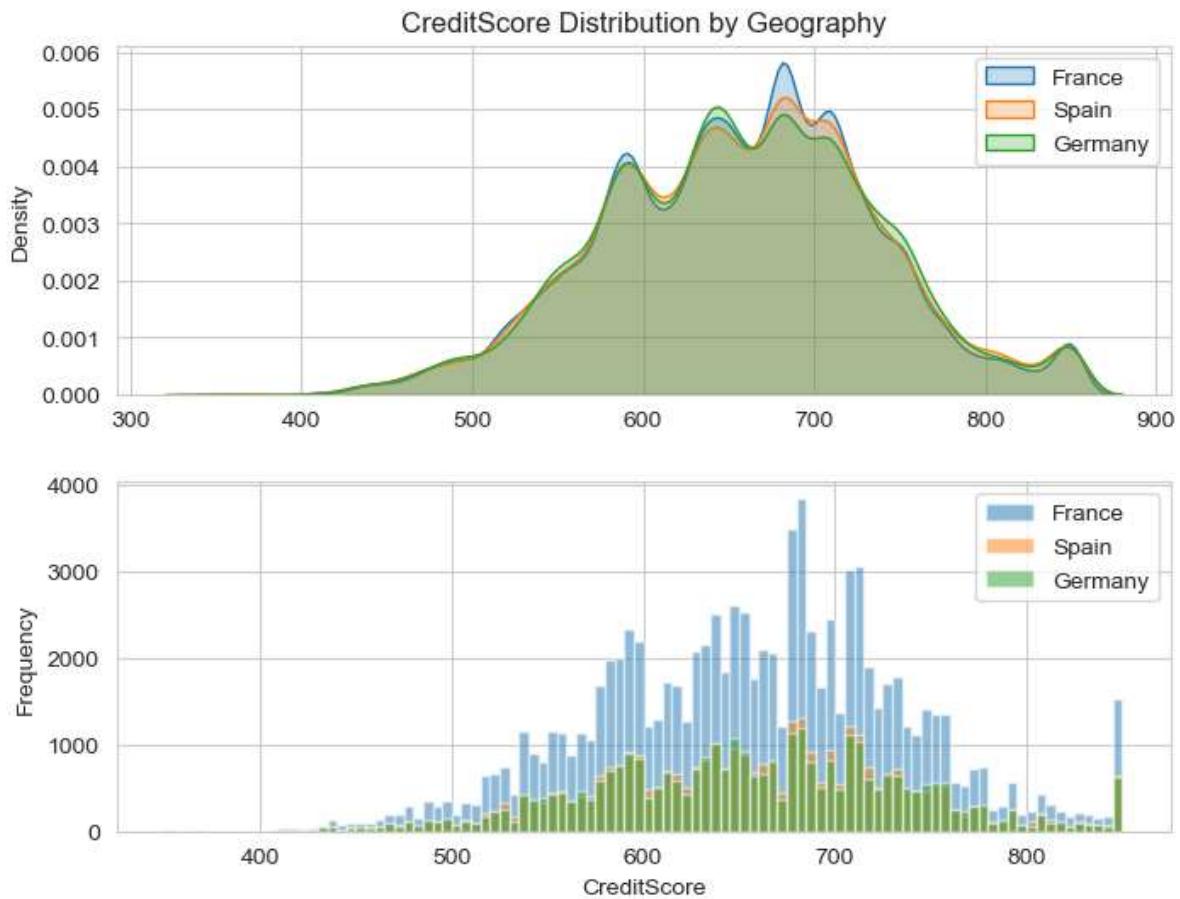
In [56]: 1 dens_hist_gen("CreditScore", "Gender")



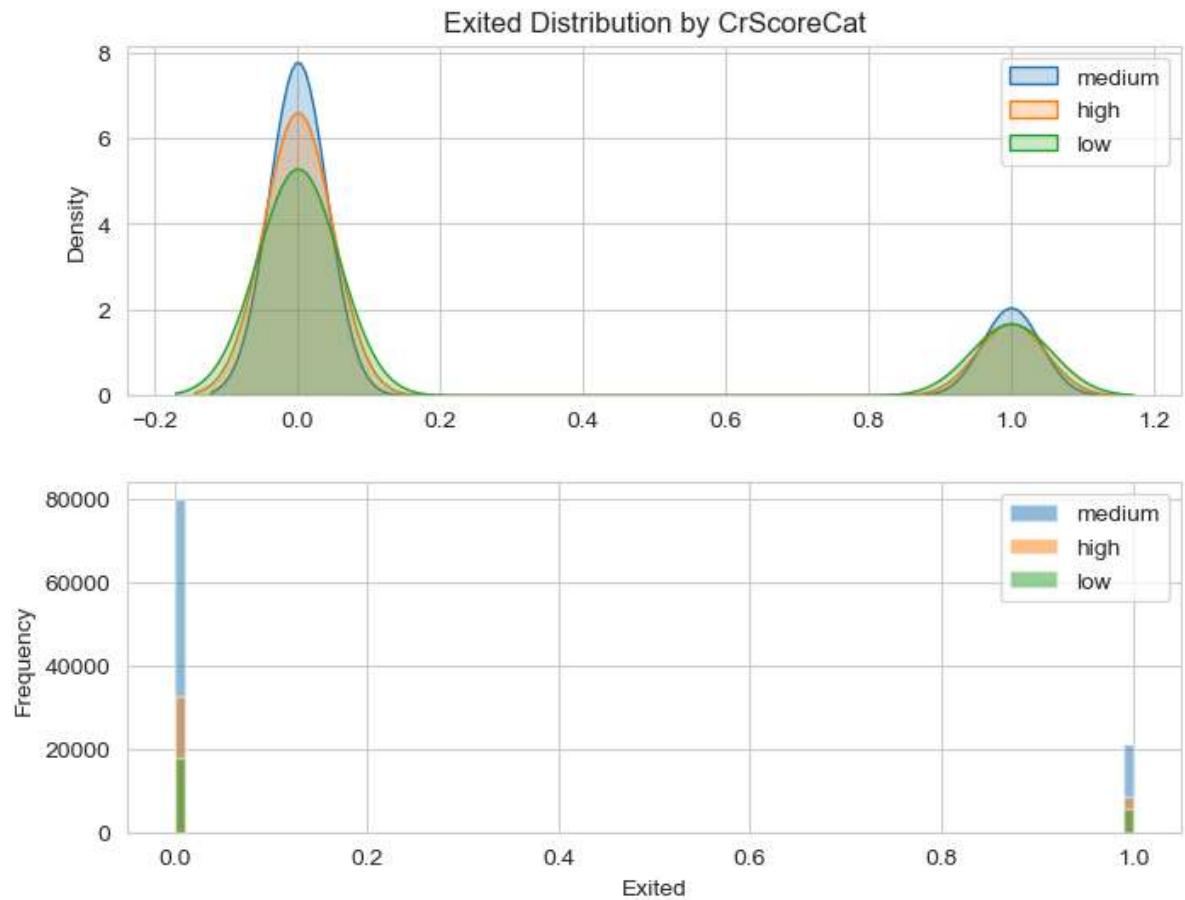
There is no clear distinction for Credit score distribution for different genders

Credit Score By Geography

```
In [57]: 1 dens_hist_gen('CreditScore', 'Geography')
```



In [58]: 1 dens_hist_gen("Exited", "CrScoreCat")



This shows that people having medium credit score is more likely to stay with the bank whereas those having low credit score terminates their services with the bank

In [59]: 1 df

Out[59]:

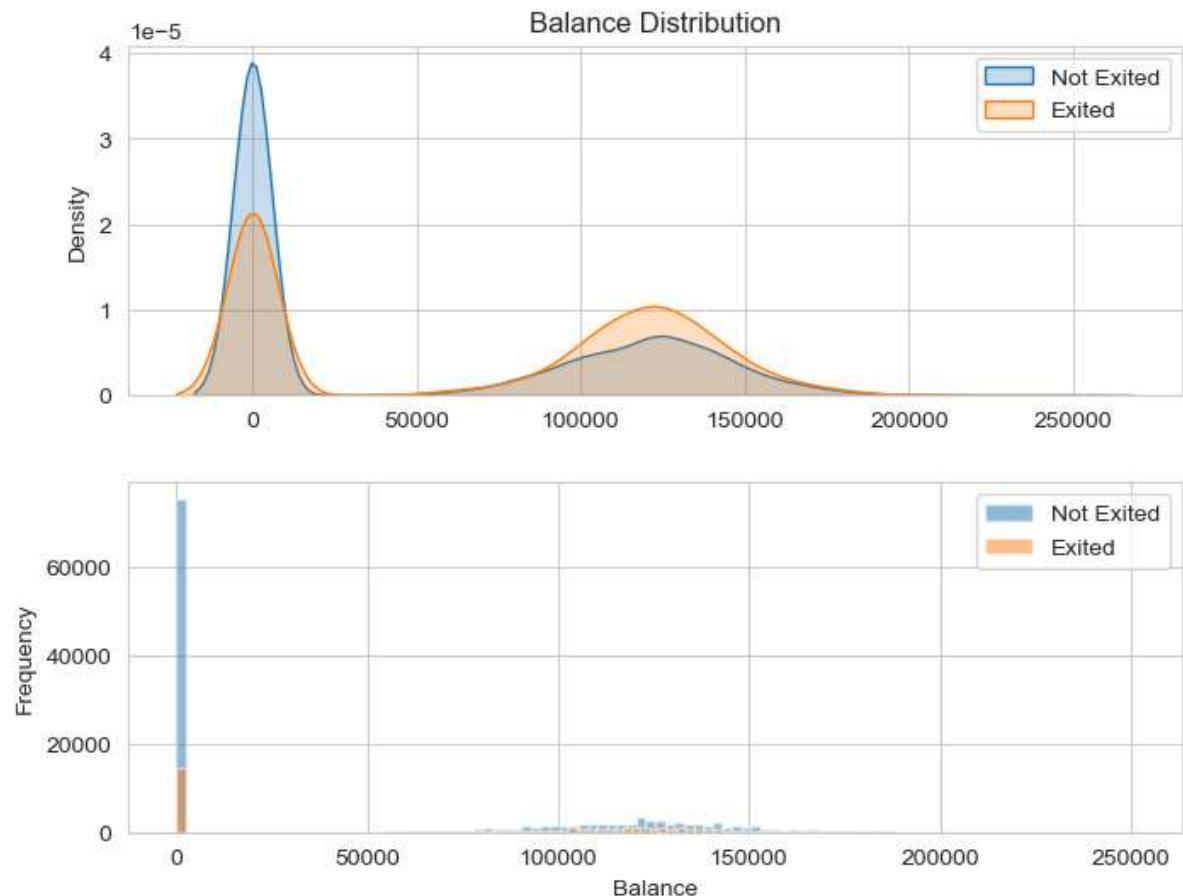
	id	CustomerId	Surname	CreditScore	Geography	Gender	Age	Tenure	Bal
0	0	15674932	Okwudilichukwu	668	France	Male	33.0	3	
1	1	15749177	Okwudilolisa	627	France	Male	33.0	1	
2	2	15694510	Hsueh	678	France	Male	40.0	10	
3	3	15741417	Kao	581	France	Male	34.0	2	1488
4	4	15766172	Chiemenam	716	Spain	Male	33.0	5	
...
165029	165029	15667085	Meng	667	Spain	Female	33.0	2	
165030	165030	15665521	Okechukwu	792	France	Male	35.0	3	
165031	165031	15664752	Hsia	565	France	Male	31.0	5	
165032	165032	15689614	Hsiung	554	Spain	Female	30.0	7	1615
165033	165033	15732798	Ulyanov	850	France	Male	31.0	1	

165034 rows × 16 columns

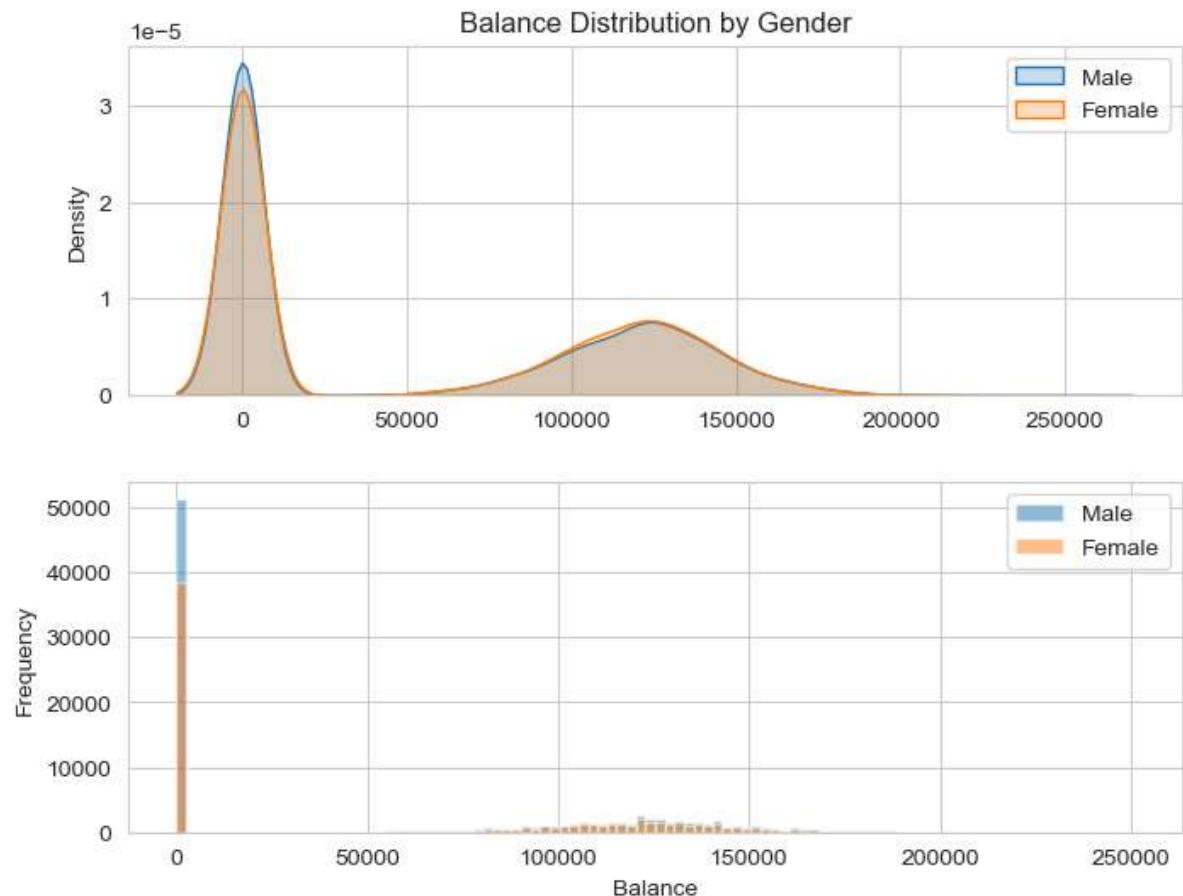


Balance

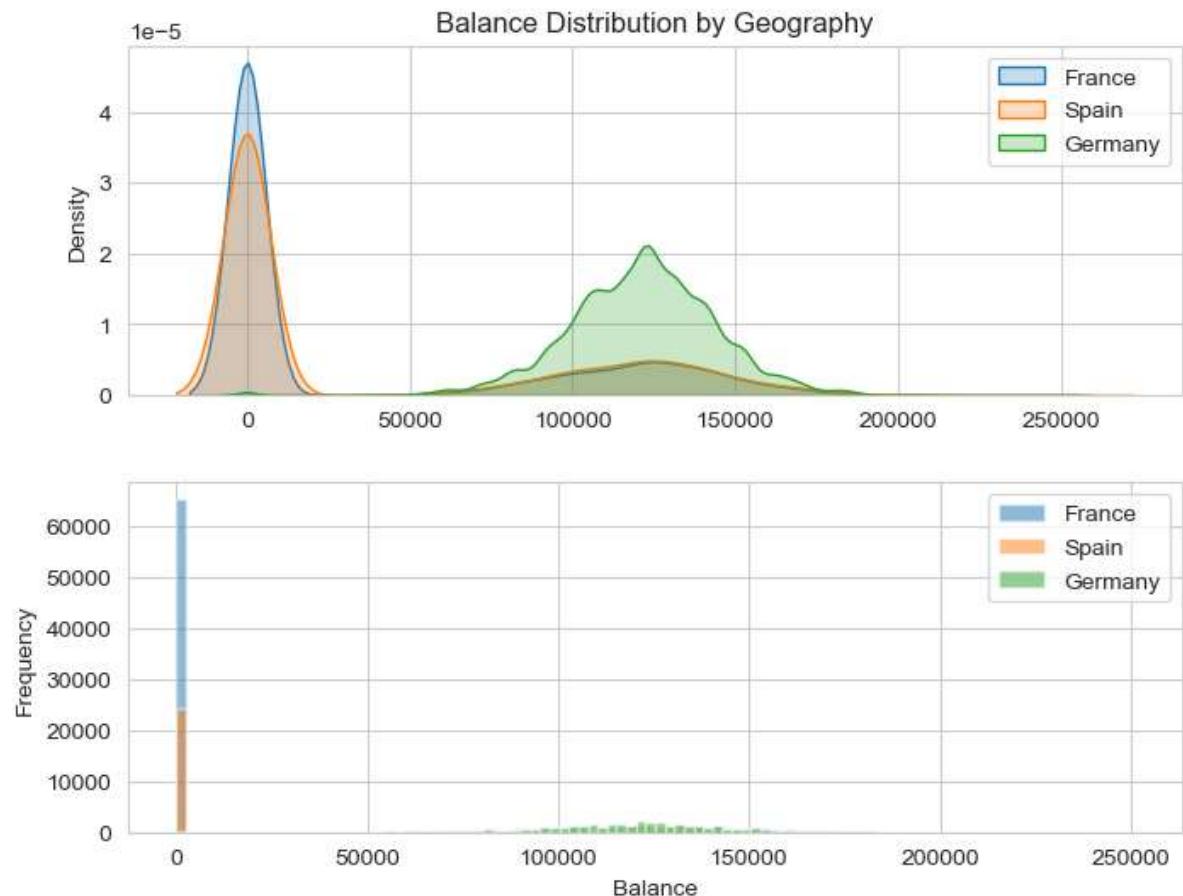
In [60]: 1 dens_hist_exit('Balance')



In [61]: 1 dens_hist_gen("Balance", "Gender")

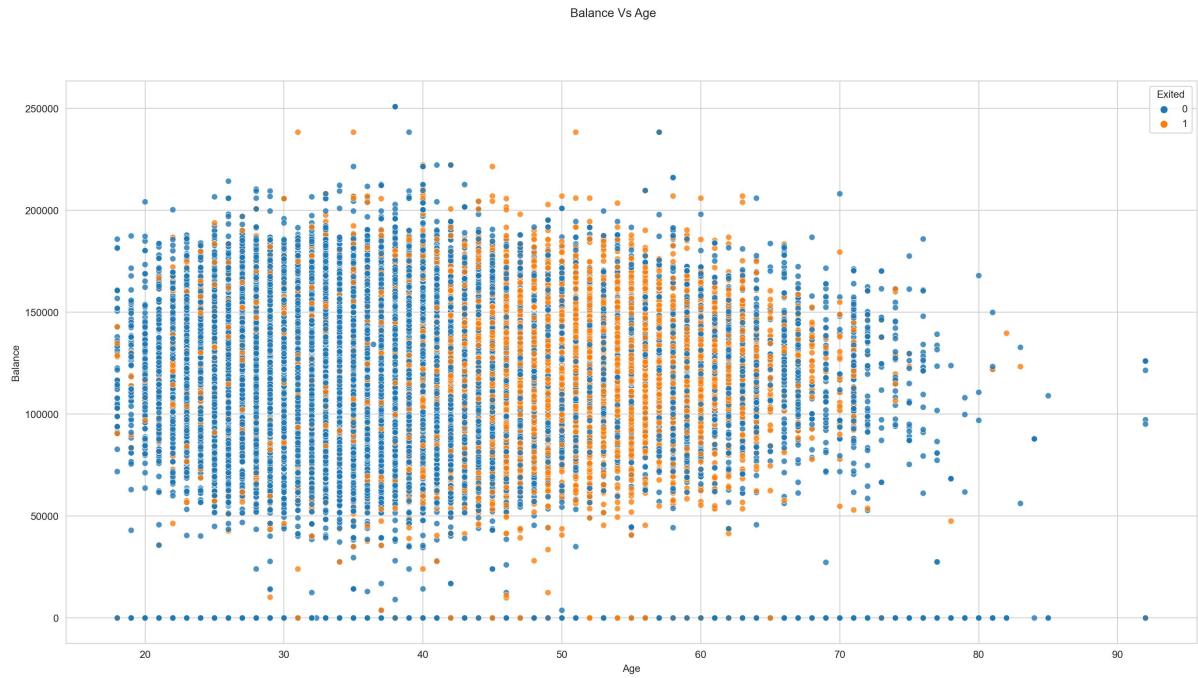


In [62]: 1 dens_hist_gen("Balance", "Geography")



We see that Germany has more number of people having higher bank balance compared to France and Spain

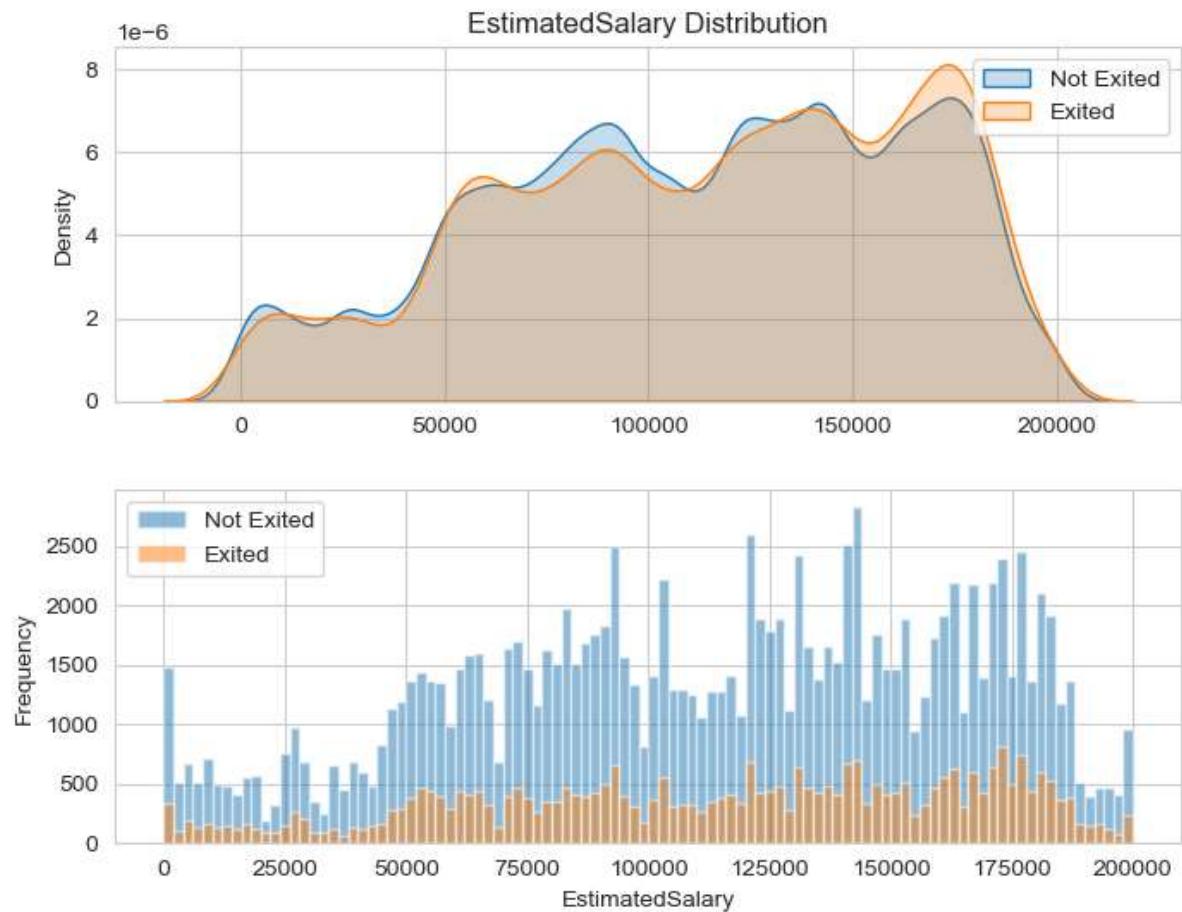
```
In [63]: 1 fix, ax = plt.subplots(figsize=(20,10),dpi=200)
2 _=sb.scatterplot(df,x='Age',y='Balance',hue='Exited',alpha=0.8)
3 _=plt.suptitle("Balance Vs Age")
4 #plt.plot(df["Age"],df["Balance"])
5
```



We see that people older than age 45 and having higher balance (more than 50000) are more likely to terminate their services with the bank compared to young people having the same bank balance

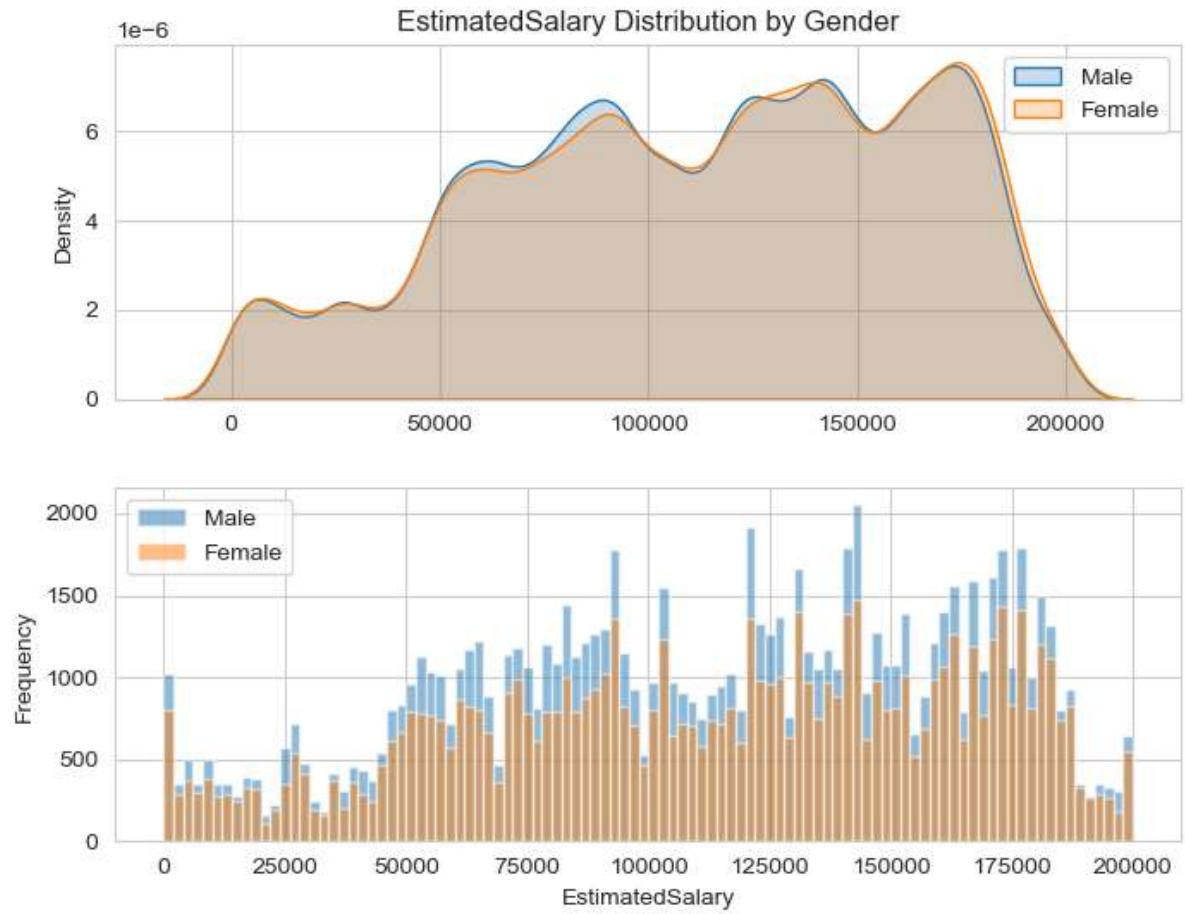
Estimated Salary

In [64]: 1 dens_hist_exit("EstimatedSalary")



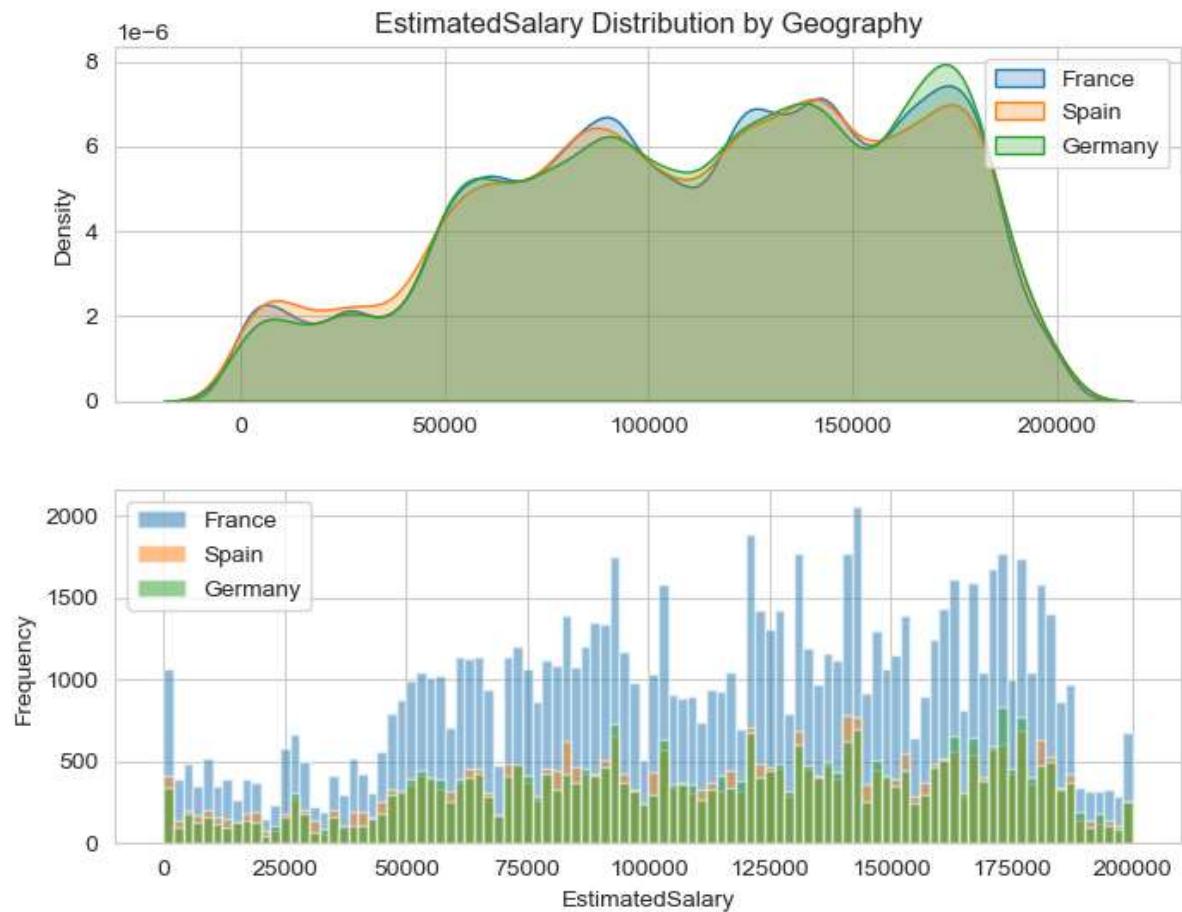
Estimated Salary By Gender

```
In [65]: 1 dens_hist_gen("EstimatedSalary", "Gender")
```



Estimated Salary by Geography

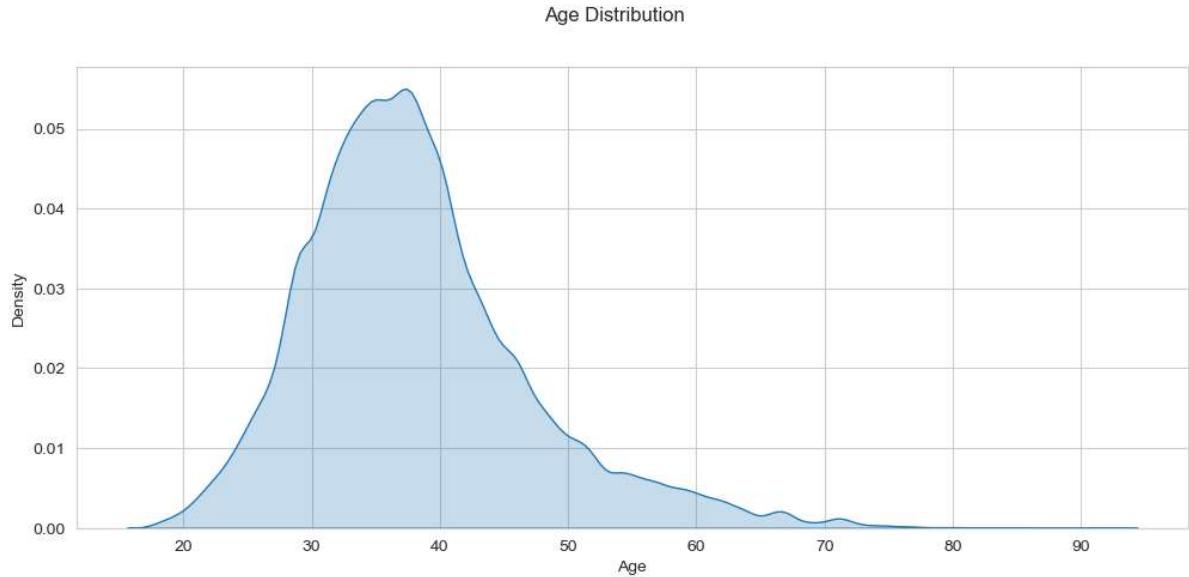
In [66]: 1 dens_hist_gen("EstimatedSalary", "Geography")



There is no clear distinction for estimated salary based on geography or gender

AGE

```
In [67]: 1 fig, ax = plt.subplots(figsize=(12,5))
2
3 sb.set_style("whitegrid")
4
5 _ = sb.kdeplot(df["Age"], fill=True)
6 _ = plt.suptitle('Age Distribution')
```

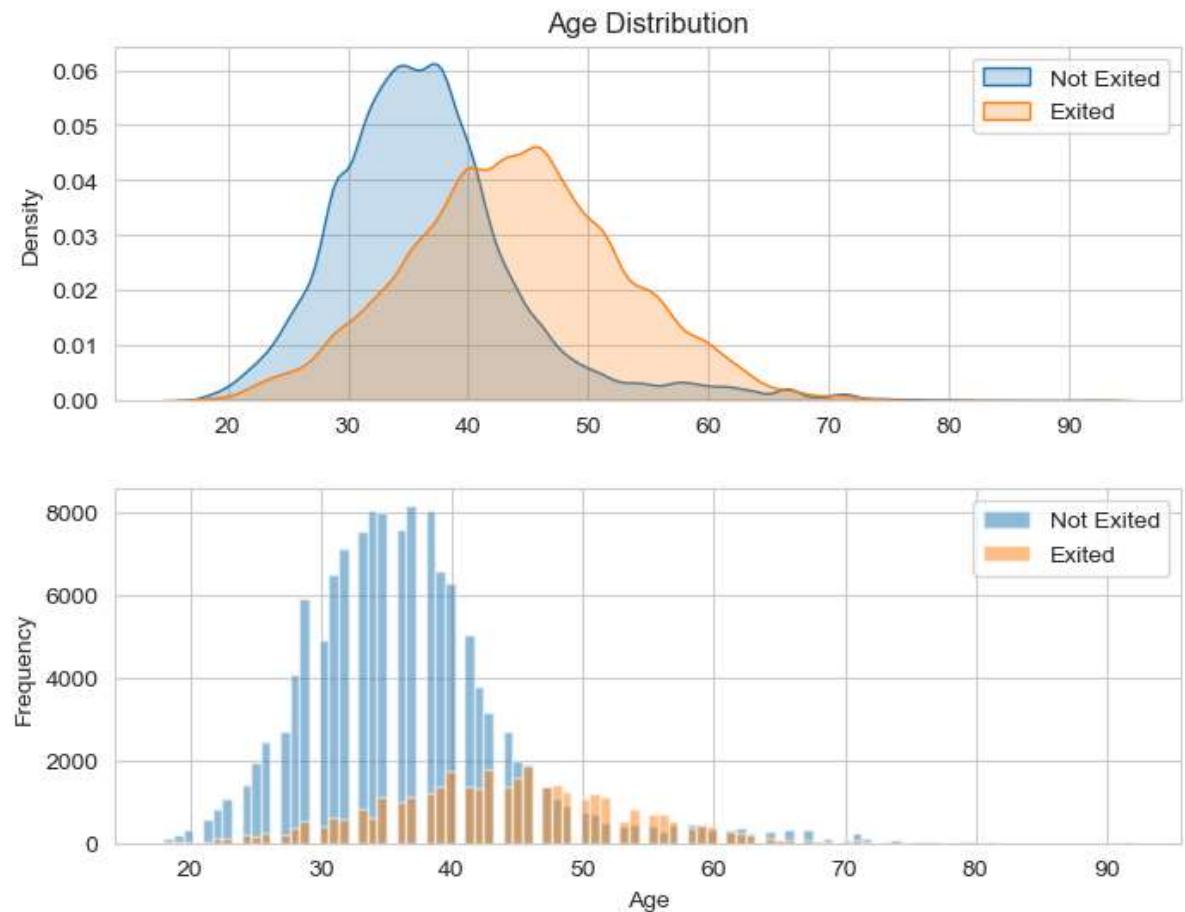


We see that most of our customer data lies between the age of 25-45

```
In [68]: 1 # Creating age groups
2
3 def age_group(age):
4     if age<=35:
5         return "Young"
6     elif age<=45:
7         return "Middle_Aged"
8     elif age<=60:
9         return "Senior"
10    else:
11        return "Old"
```

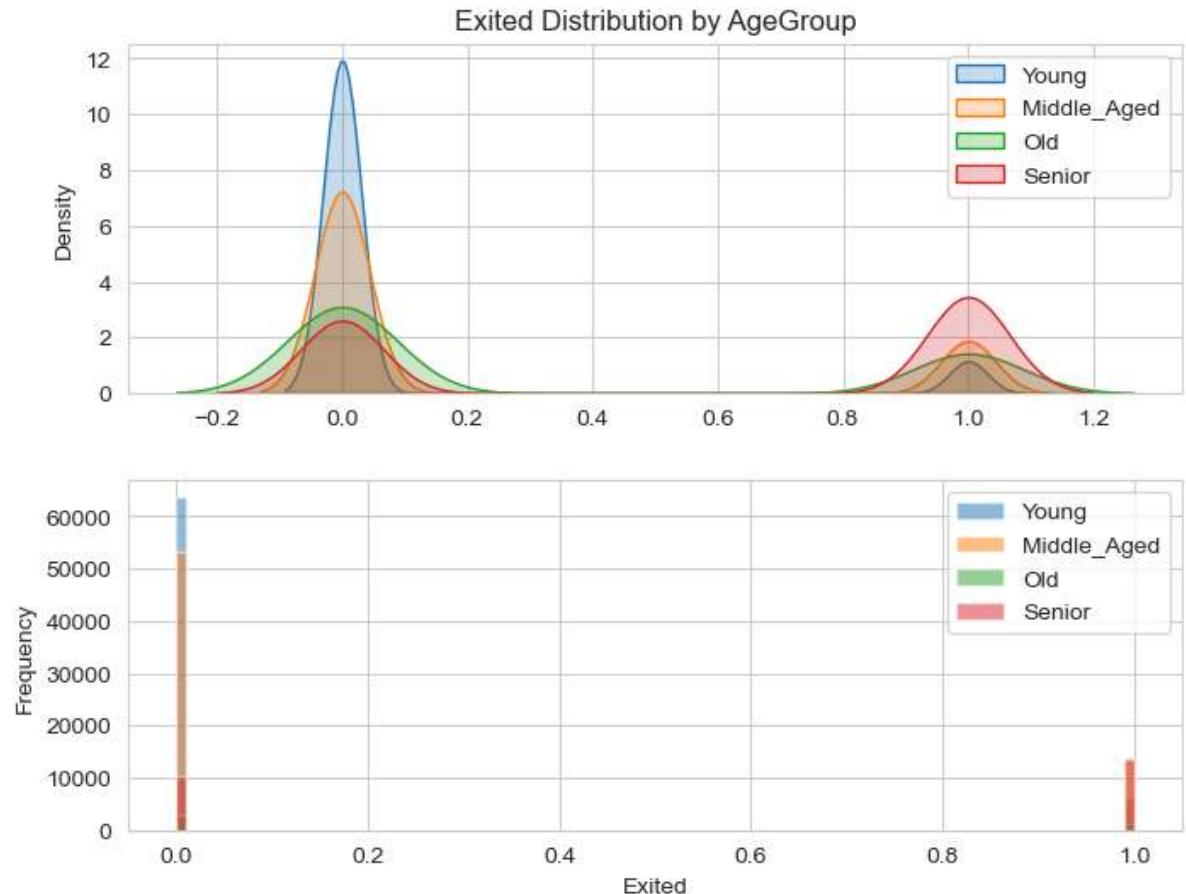
```
In [69]: 1 df["AgeGroup"] = df["Age"].apply(age_group)
```

In [70]: 1 dens_hist_exit("Age")



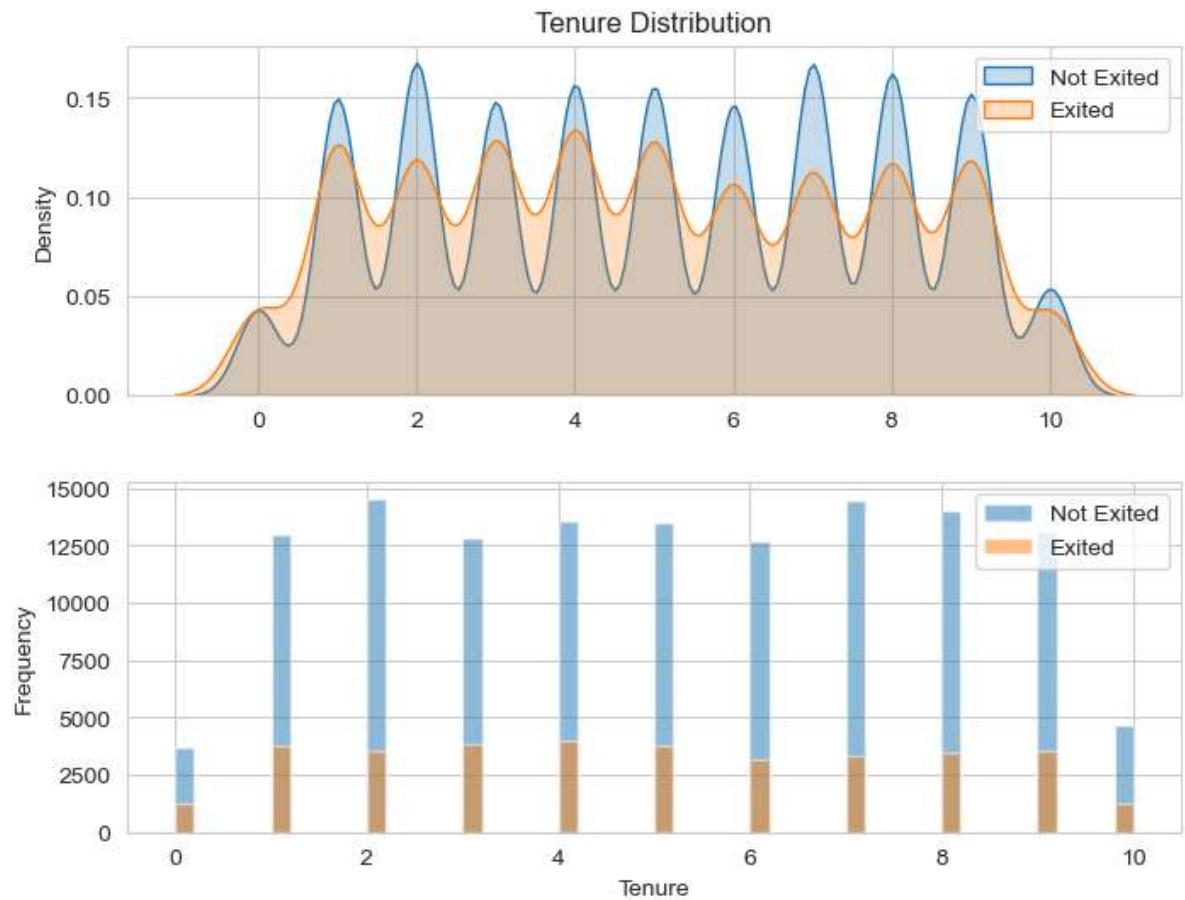
There is a clear distinction between the age groups who terminate their services Vs those who stay with the bank. This is an important feature for Churn determination which was also showed by the correlation matrix

In [71]: 1 dens_hist_gen("Exited", "AgeGroup")



We see that based on our defined age groups, Young and middle-Aged people are more likely to stay with the bank whereas senior people are more likely to exit

In [72]: 1 dens_hist_exit("Tenure", 50)

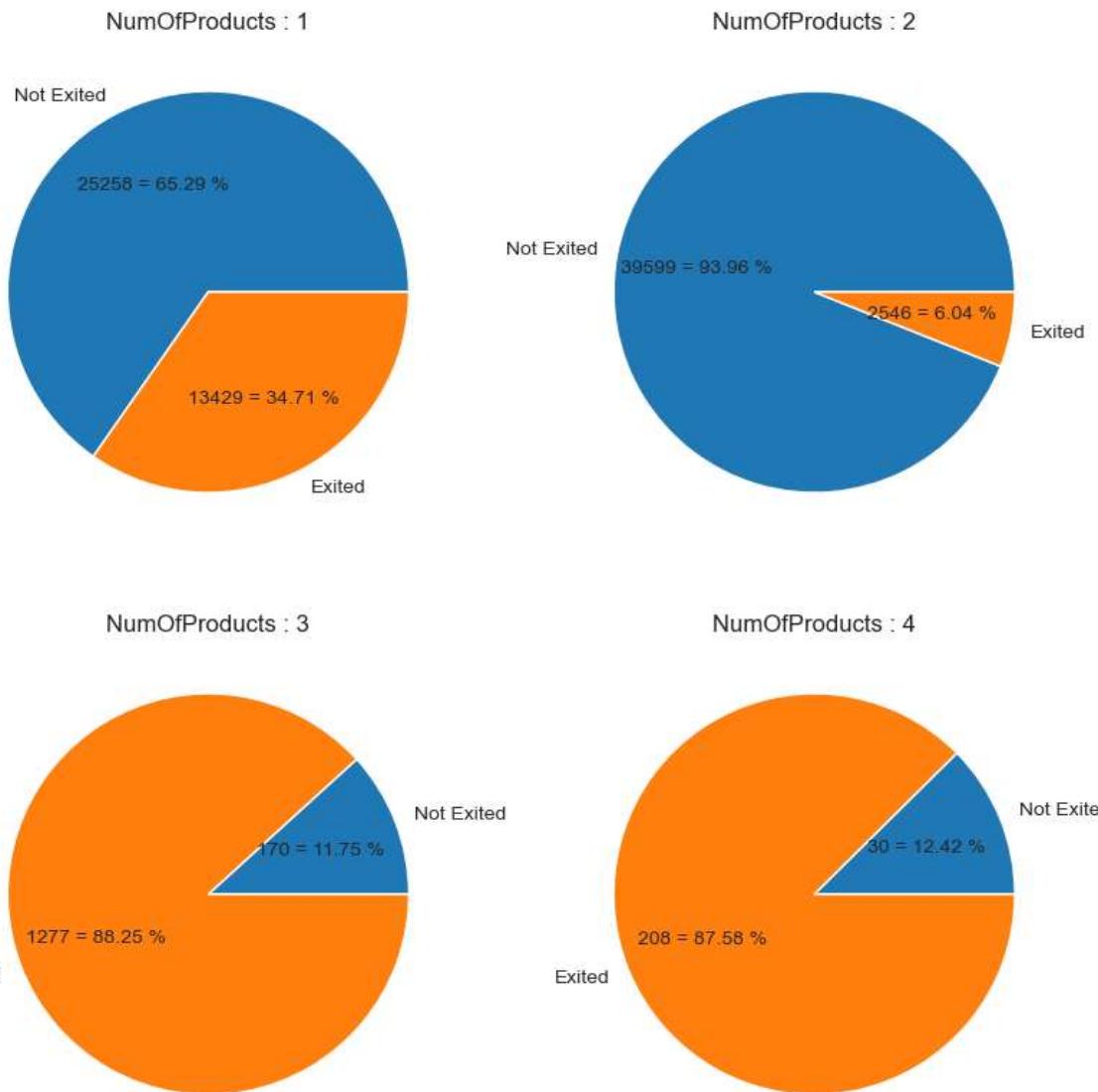


This shows that people staying with bank for 1 to 9 years are more likely to stay with the bank

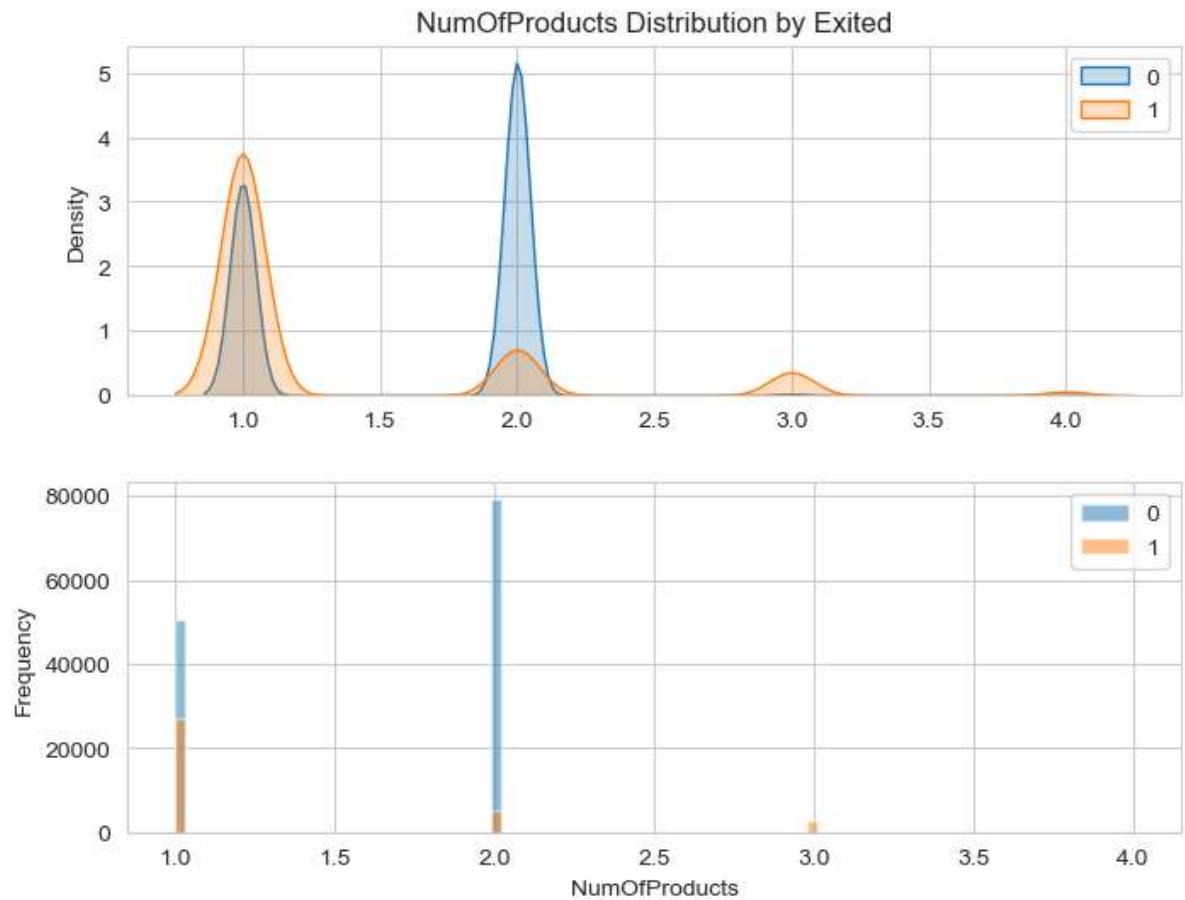
Number of products

In [73]:

```
1 fig,ax = plt.subplots(nrows=2,ncols=2,figsize=(10,10))
2
3 ax=ax.flatten()
4
5 for i in range(4):
6
7     NOP_count=[df[(df["NumOfProducts"]== i+1) & (df["Exited"]==0)]["id"].c
8     NOP_label=["Not Exited","Exited"]
9     _=ax[i].pie(NOP_count,labels=NOP_label,autopct=lambda x: "{:.0f} = {:.0f}%".format(x,NOP_count[x]))
10   _=ax[i].set_title("NumOfProducts : {}".format(i+1))
```

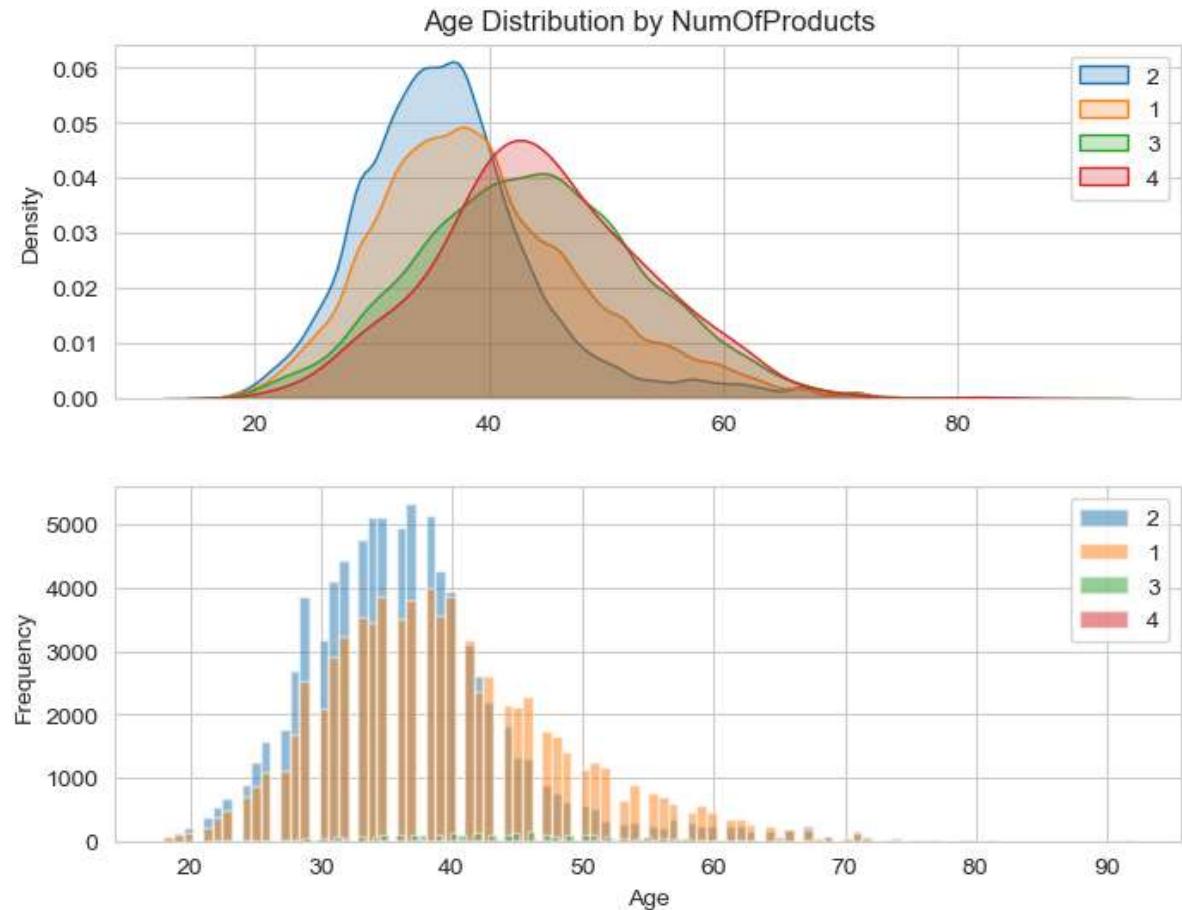


In [74]: 1 dens_hist_gen("NumOfProducts", "Exited")



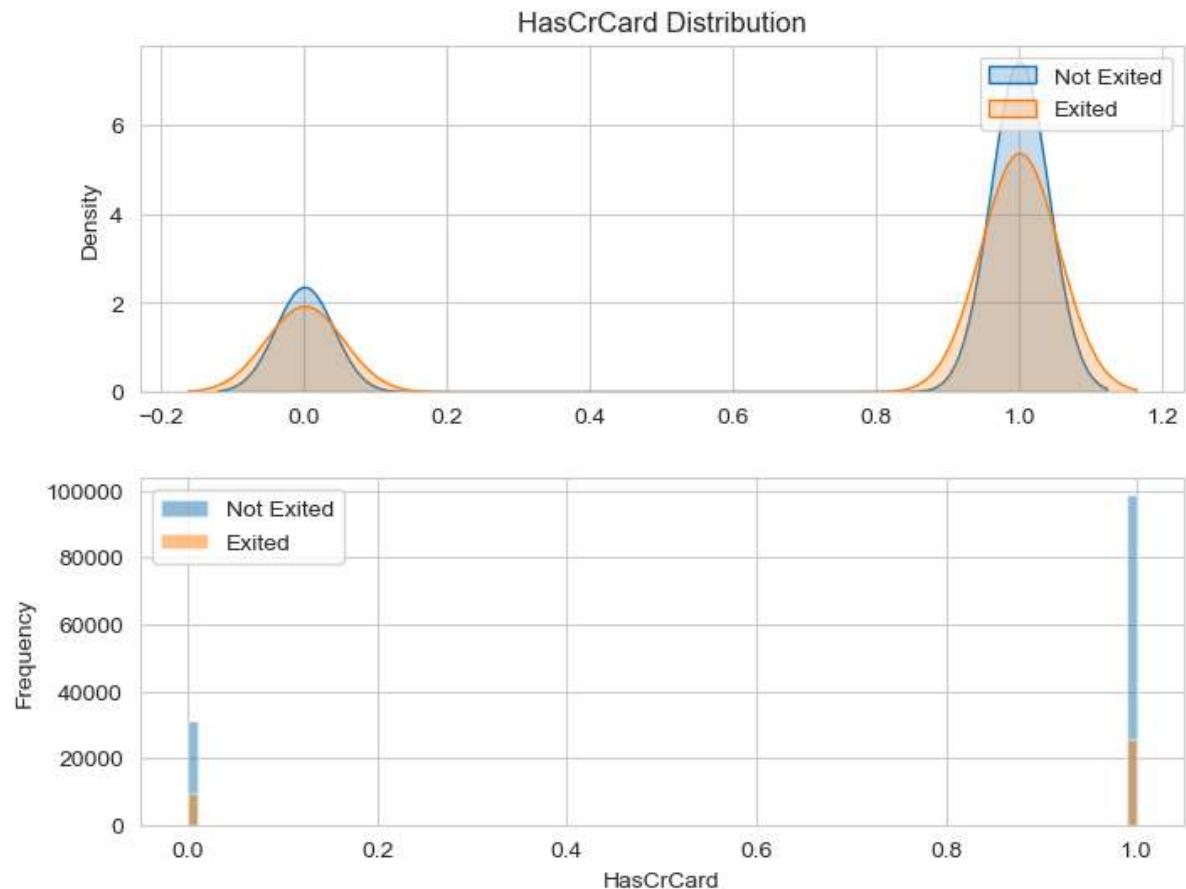
Larger part of the population who has 2 products is more likely to stay with bank whereas those having more than 2 products are more likely to exit

```
In [75]: 1 dens_hist_gen("Age", "NumOfProducts")
```

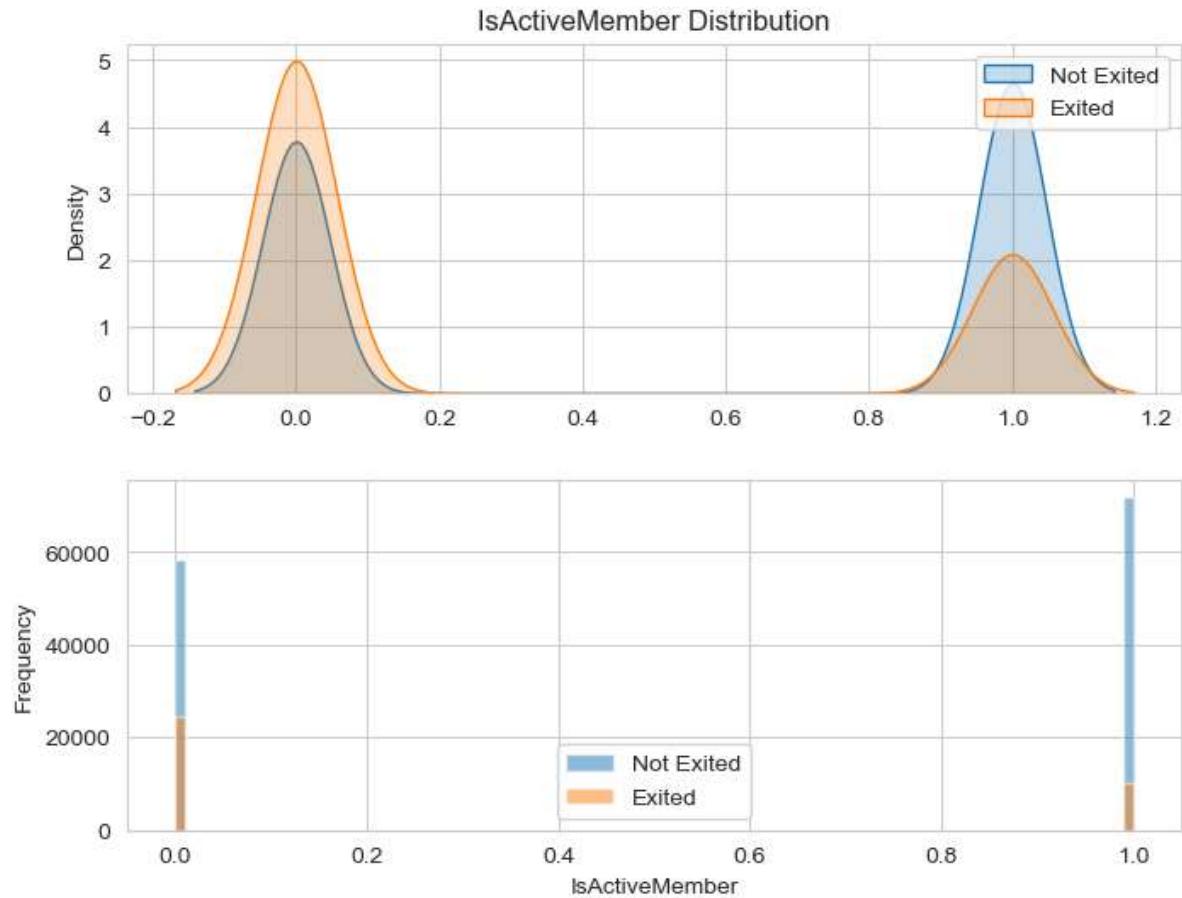


- The distribution of age for people who own 2 products (blue) skews younger, with a peak around the late 20s to early 30s.
- Those with 1 product (orange) have a broader age distribution, peaking in the 30s and then slowly declining.
- The distribution for 3 products (green) and 4 products (red) are more similar to each other, both peaking around the mid-30s to 40s but with the 4 products curve having a slightly broader spread.

```
In [76]: 1 dens_hist_exit("HasCrCard")
```

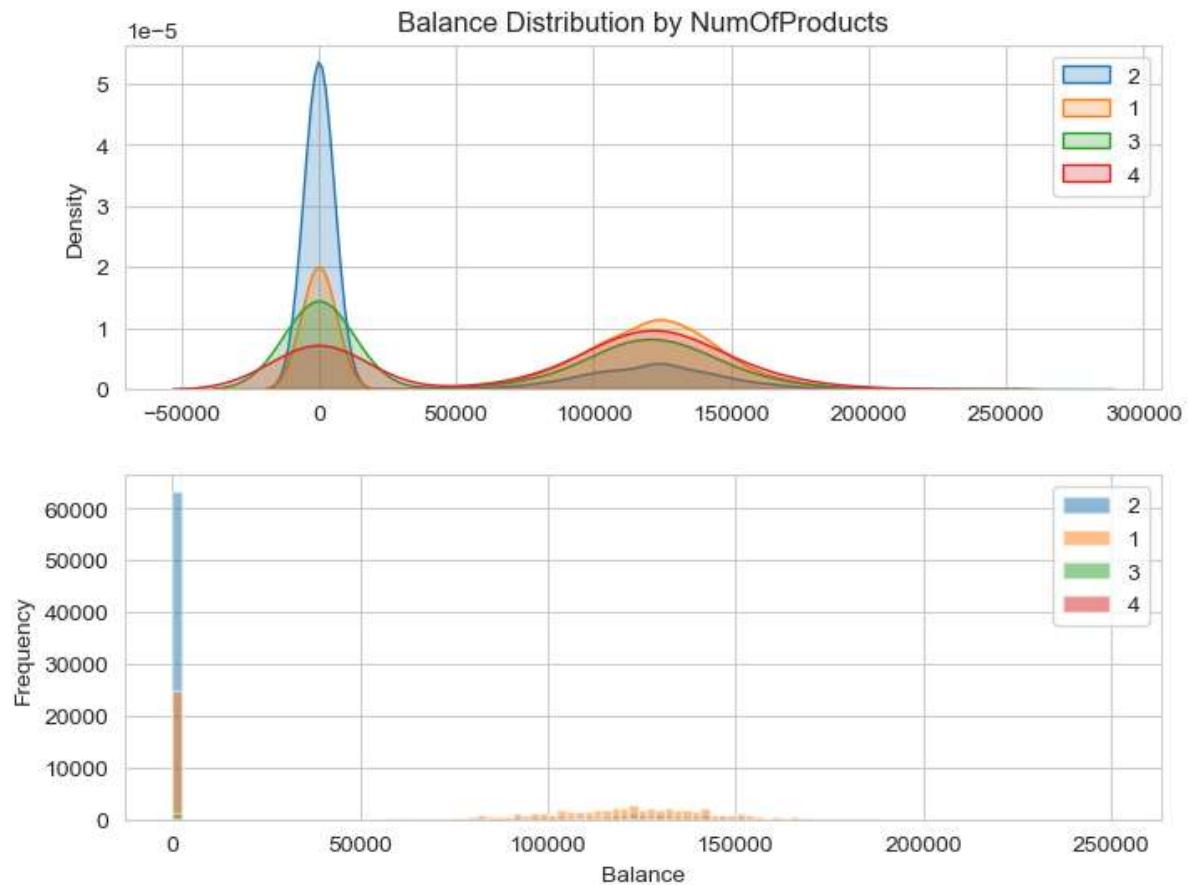


In [77]: 1 dens_hist_exit("IsActiveMember")



- From the bottom histogram, we can see that a larger number of active members have not exited (indicated by the blue bar at 1 on the x-axis). The number of non-active members (0 on the x-axis) who have not exited is also significant, but smaller than that of active members.
- Conversely, a smaller number of active members have exited (orange bar at 1 on the x-axis), while the number of non-active members who have exited is relatively smaller.
- If we consider the density plot, it suggests that the proportion of members who have exited is greater in the non-active members than in the active members, which could indicate that active membership is potentially a factor in retaining members.

In [78]: 1 dens_hist_gen("Balance", "NumOfProducts")



- Individuals owning 2 products (blue) have a sharp peak around a balance of 0, suggesting that a large number of these individuals either maintain a low balance or the data might be reflecting an account status (like zero-balance accounts).
- The distributions for 1 product (orange), 3 products (green), and 4 products (red) are broader with peaks at higher balances compared to those with 2 products. This indicates that individuals with 1, 3, or 4 products tend to maintain higher account balances on average.
- Individuals with 1 product have both high frequency and high density at lower balances, indicating that they are the largest group within the population with a balance closer to 0.