

```
In [2]: import pandas as pd
import numpy as np
import seaborn as sns
import matplotlib
from matplotlib import pyplot as plt
import itertools

%matplotlib inline
matplotlib.style.use('ggplot')
# Input data files are available in the "../input/" directory.
# For example, running this (by clicking run or pressing Shift+Enter) will list the files in the input di
rectory

# import os
# print(os.listdir("../input"))

# Any results you write to the current directory are saved as output.
```

```
In [3]: # Reading the dataset
sales_data = pd.read_csv("SalesKaggle3.csv")
```

```
In [4]: sales_data.head()
```

Out[4]:

	Order	File_Type	SKU_number	SoldFlag	SoldCount	MarketingType	ReleaseNumber	New_Release_Flag	StrengthFactor	PriceReg	ReleaseY
0	2	Historical	1737127	0.0	0.0	D	15	1	682743.0	44.99	20
1	3	Historical	3255963	0.0	0.0	D	7	1	1016014.0	24.81	20
2	4	Historical	612701	0.0	0.0	D	0	0	340464.0	46.00	20
3	6	Historical	115883	1.0	1.0	D	4	1	334011.0	100.00	20
4	7	Historical	863939	1.0	1.0	D	2	1	1287938.0	121.95	20

In [5]: *#Statistical description of the dataset*  
 sales\_data.describe()

Out[5]:

	Order	SKU_number	SoldFlag	SoldCount	ReleaseNumber	New_Release_Flag	StrengthFactor	PriceReg	ReleaseY
<b>count</b>	198917.000000	1.989170e+05	75996.000000	75996.000000	198917.000000	198917.000000	1.989170e+05	198917.000000	198917.000000
<b>mean</b>	106483.543242	8.613626e+05	0.171009	0.322306	3.412202	0.642248	1.117115e+06	90.895243	2006.0164
<b>std</b>	60136.716784	8.699794e+05	0.376519	1.168615	3.864243	0.479340	1.522090e+06	86.736367	9.1583
<b>min</b>	2.000000	5.000100e+04	0.000000	0.000000	0.000000	0.000000	6.275000e+00	0.000000	0.000000
<b>25%</b>	55665.000000	2.172520e+05	0.000000	0.000000	1.000000	0.000000	1.614188e+05	42.000000	2003.000000
<b>50%</b>	108569.000000	6.122080e+05	0.000000	0.000000	2.000000	1.000000	5.822240e+05	69.950000	2007.000000
<b>75%</b>	158298.000000	9.047510e+05	0.000000	0.000000	5.000000	1.000000	1.430083e+06	116.000000	2011.000000
<b>max</b>	208027.000000	3.960788e+06	1.000000	73.000000	99.000000	1.000000	1.738445e+07	12671.480000	2018.000000

```
In [6]: # Includes categorical variable
sales_data.describe(include='all')
```

Out[6]:

	Order	File_Type	SKU_number	SoldFlag	SoldCount	MarketingType	ReleaseNumber	New_Release_Flag	StrengthFactor
<b>count</b>	198917.000000	198917	1.989170e+05	75996.000000	75996.000000	198917	198917.000000	198917.000000	1.989170e+05
<b>unique</b>	NaN	2	NaN	NaN	NaN	2	NaN	NaN	NaN
<b>top</b>	NaN	Active	NaN	NaN	NaN	S	NaN	NaN	NaN
<b>freq</b>	NaN	122921	NaN	NaN	NaN	100946	NaN	NaN	NaN
<b>mean</b>	106483.543242	NaN	8.613626e+05	0.171009	0.322306	NaN	3.412202	0.642248	1.117115e+06
<b>std</b>	60136.716784	NaN	8.699794e+05	0.376519	1.168615	NaN	3.864243	0.479340	1.522090e+06
<b>min</b>	2.000000	NaN	5.000100e+04	0.000000	0.000000	NaN	0.000000	0.000000	6.275000e+00
<b>25%</b>	55665.000000	NaN	2.172520e+05	0.000000	0.000000	NaN	1.000000	0.000000	1.614188e+05
<b>50%</b>	108569.000000	NaN	6.122080e+05	0.000000	0.000000	NaN	2.000000	1.000000	5.822240e+05
<b>75%</b>	158298.000000	NaN	9.047510e+05	0.000000	0.000000	NaN	5.000000	1.000000	1.430083e+06
<b>max</b>	208027.000000	NaN	3.960788e+06	1.000000	73.000000	NaN	99.000000	1.000000	1.738445e+07

```
In [7]: sales_data.shape
```

Out[7]: (198917, 14)

```
In [8]: sales_data.nunique()
```

```
Out[8]: Order          198917  
File_Type             2  
SKU_number           133360  
SoldFlag              2  
SoldCount             37  
MarketingType         2  
ReleaseNumber         71  
New_Release_Flag      2  
StrengthFactor        197424  
PriceReg              11627  
ReleaseYear           85  
ItemCount             501  
LowUserPrice          12102  
LowNetPrice           15403  
dtype: int64
```

```
In [9]: sales_data.isnull().values.any()
```

```
Out[9]: True
```

```
In [10]: sales_data.isnull().sum()
```

```
Out[10]: Order          0  
File_Type             0  
SKU_number            0  
SoldFlag             122921  
SoldCount             122921  
MarketingType         0  
ReleaseNumber         0  
New_Release_Flag      0  
StrengthFactor        0  
PriceReg              0  
ReleaseYear           0  
ItemCount             0  
LowUserPrice          0  
LowNetPrice           0  
dtype: int64
```

```
In [11]: sales_data['SoldFlag'].fillna(0, inplace=True)
sales_data['SoldCount'].fillna(0, inplace=True)
```

```
In [12]: sales_data.isnull().sum()
```

```
Out[12]: Order                0
File_Type                0
SKU_number               0
SoldFlag                 0
SoldCount                0
MarketingType            0
ReleaseNumber            0
New_Release_Flag         0
StrengthFactor           0
PriceReg                 0
ReleaseYear              0
ItemCount                0
LowUserPrice             0
LowNetPrice              0
dtype: int64
```

```
In [13]: sales_data[sales_data['File_Type'] == 'Historical']['SKU_number'].count()
```

```
Out[13]: 75996
```

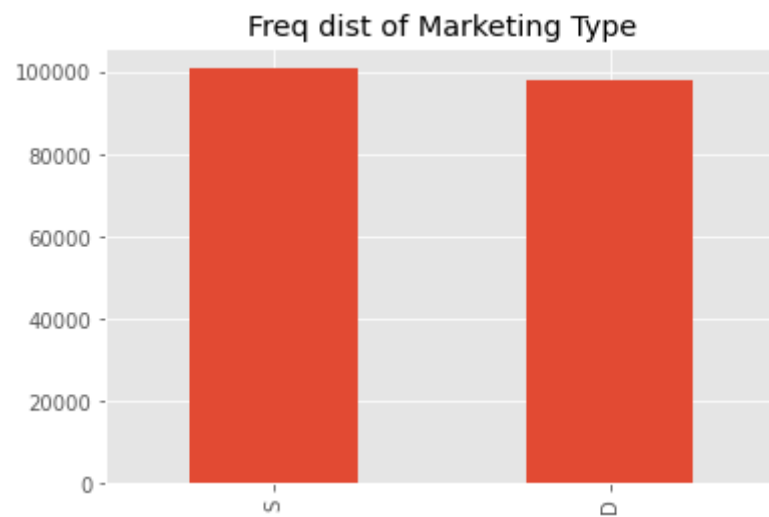
```
In [14]: sales_data[sales_data['File_Type'] == 'Active']['SKU_number'].count()
```

```
Out[14]: 122921
```

```
In [15]: sales_data_hist = sales_data[sales_data['File_Type'] == 'Historical']
sales_data_act = sales_data[sales_data['File_Type'] == 'Active']
```

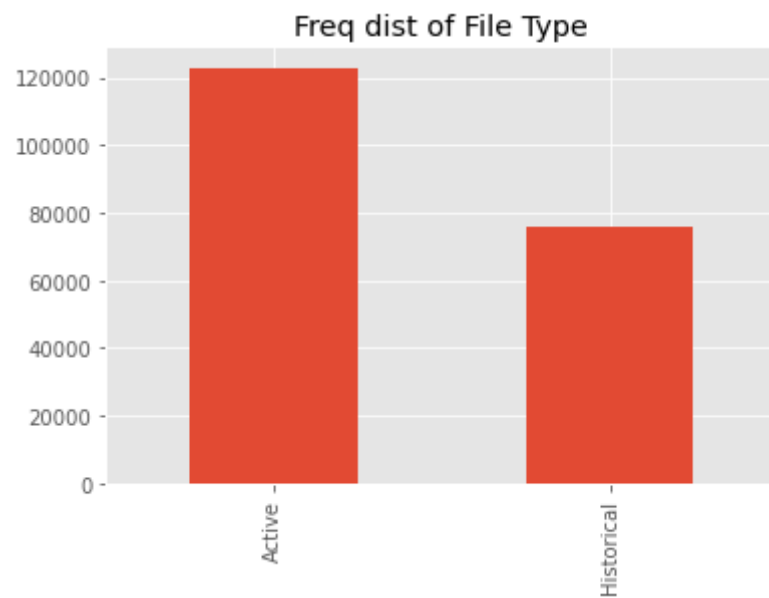
```
In [16]: sales_data['MarketingType'].value_counts().plot.bar(title="Freq dist of Marketing Type")
```

```
Out[16]: <matplotlib.axes._subplots.AxesSubplot at 0x7f9b991995d0>
```



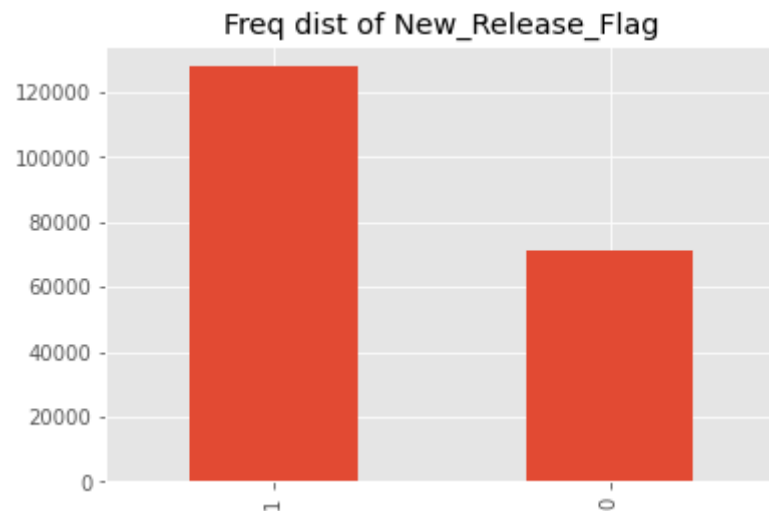
```
In [17]: sales_data['File_Type'].value_counts().plot.bar(title="Freq dist of File Type")
```

```
Out[17]: <matplotlib.axes._subplots.AxesSubplot at 0x7f9b9ed9b450>
```



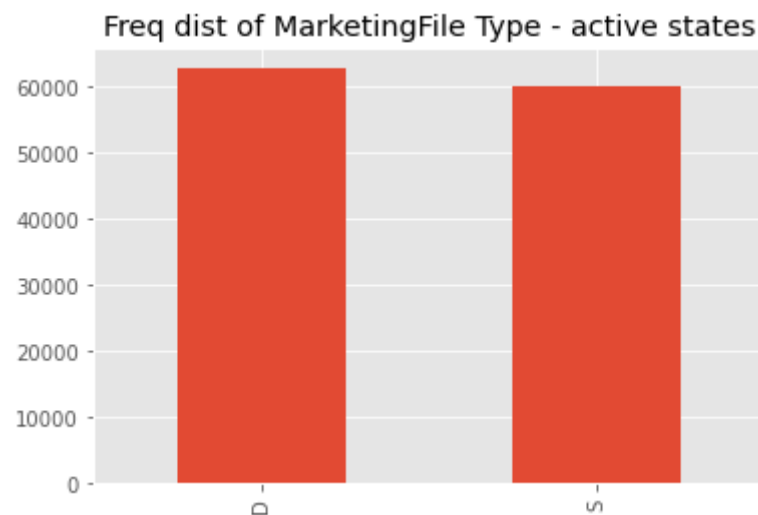
```
In [18]: sales_data['New_Release_Flag'].value_counts().plot.bar(title="Freq dist of New_Release_Flag")
```

```
Out[18]: <matplotlib.axes._subplots.AxesSubplot at 0x7f9b9eb91250>
```



```
In [19]: sales_data_act['MarketingType'].value_counts().plot.bar(title="Freq dist of MarketingFile Type - active states")
```

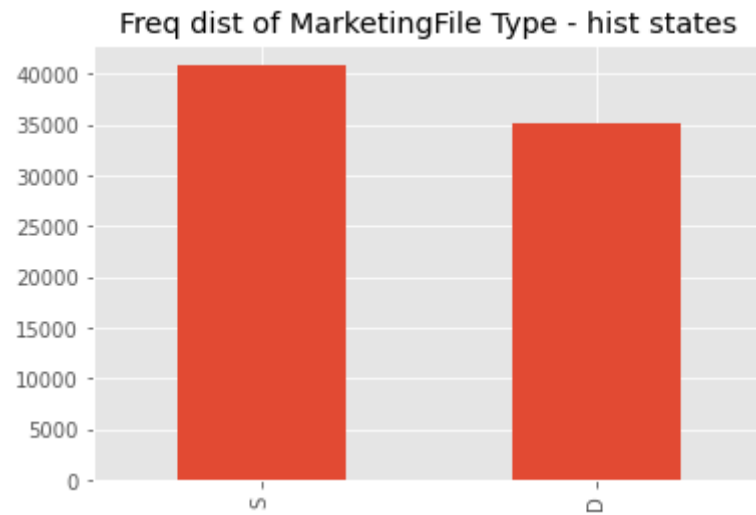
```
Out[19]: <matplotlib.axes._subplots.AxesSubplot at 0x7f9b9eb733d0>
```





```
In [20]: sales_data_hist['MarketingType'].value_counts().plot.bar(title="Freq dist of MarketingFile Type - hist states")
```

```
Out[20]: <matplotlib.axes._subplots.AxesSubplot at 0x7f9b9eb21190>
```



```
In [21]: col_names = ['StrengthFactor', 'PriceReg', 'ReleaseYear', 'ItemCount', 'LowUserPrice', 'LowNetPrice']

fig, ax = plt.subplots(len(col_names), figsize=(16,12))

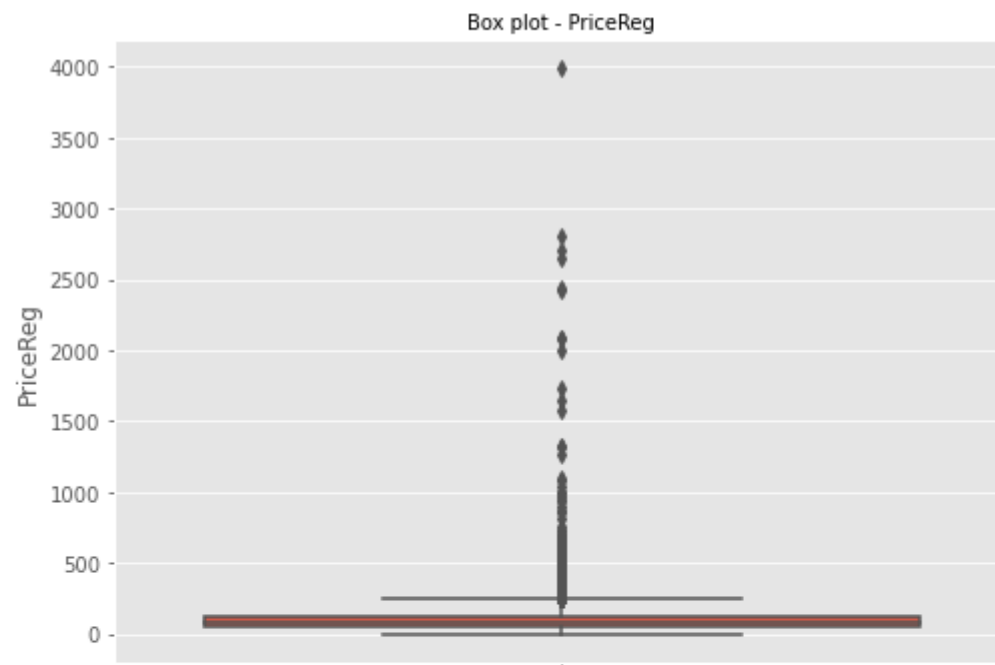
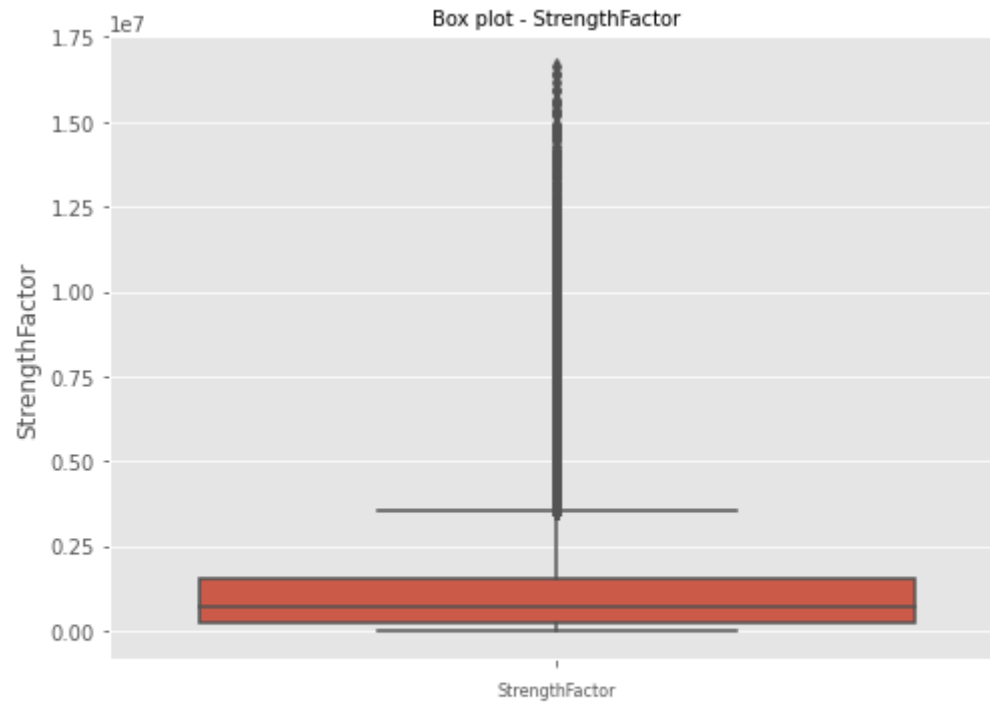
for i, col_val in enumerate(col_names):

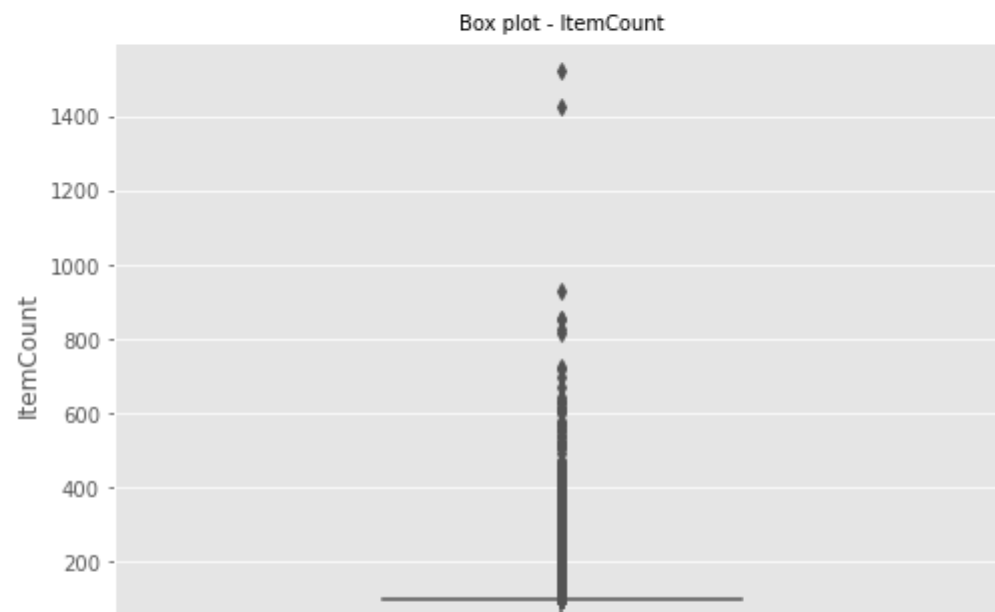
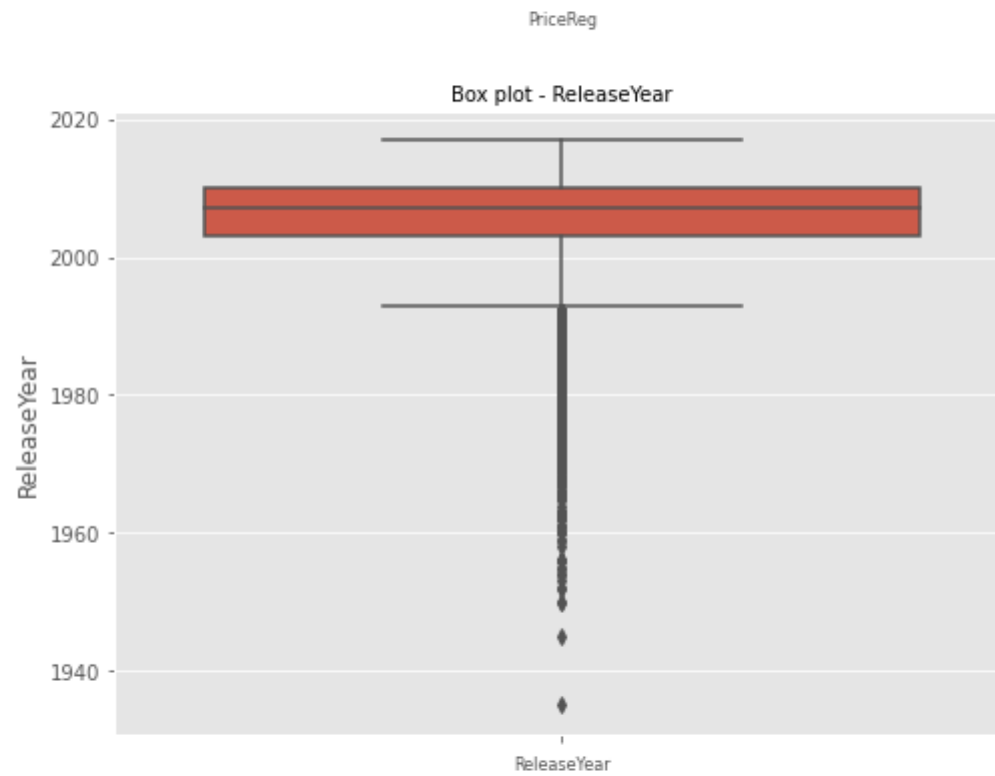
    sns.distplot(sales_data_hist[col_val], hist=True, ax=ax[i])
    ax[i].set_title('Freq dist '+col_val, fontsize=10)
    ax[i].set_xlabel(col_val, fontsize=8)
    ax[i].set_ylabel('Count', fontsize=8)

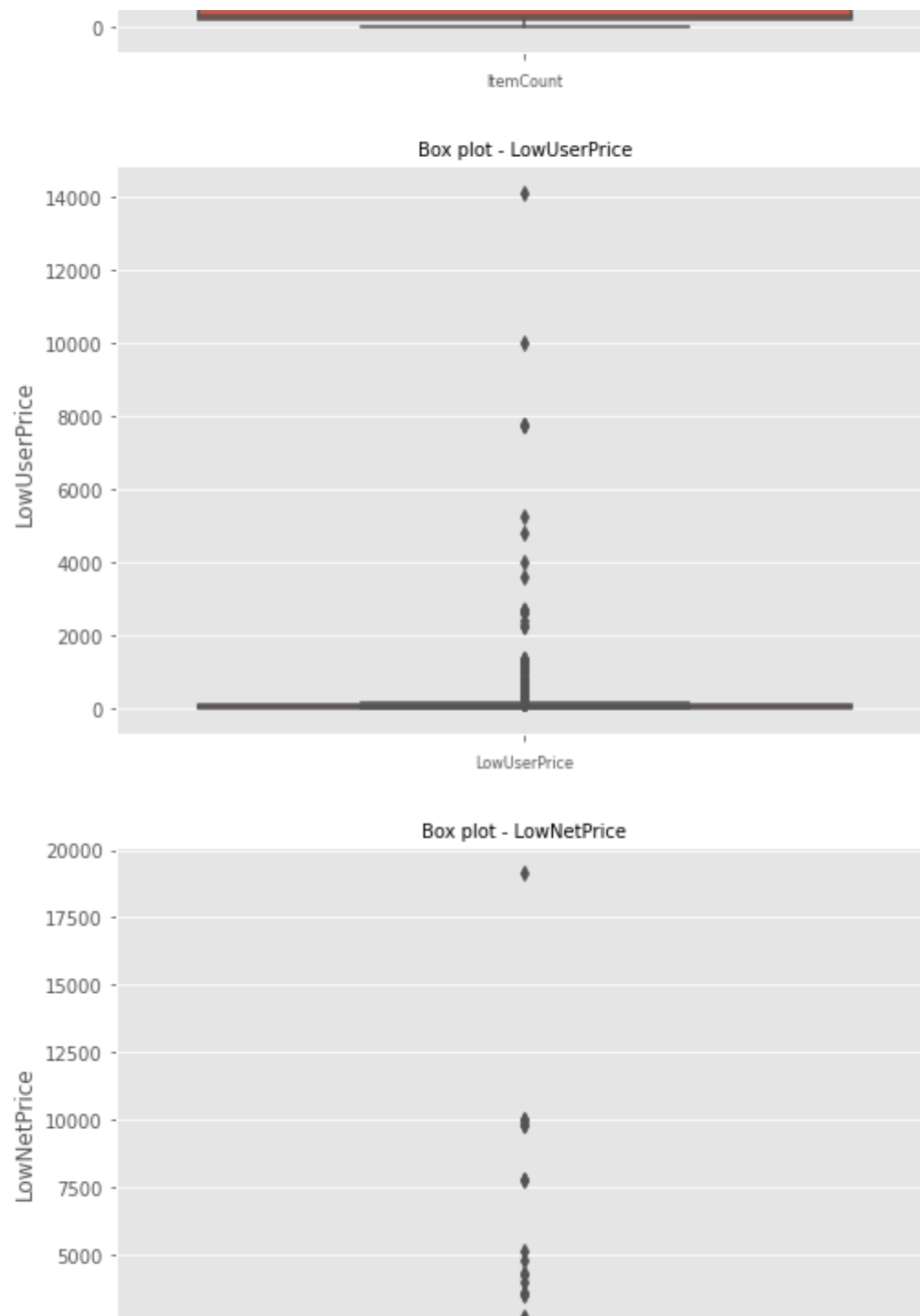
plt.show()
```

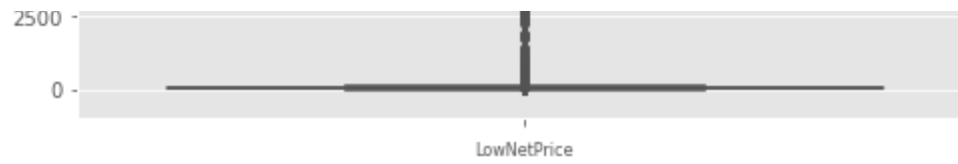


```
In [22]: col_names = ['StrengthFactor', 'PriceReg', 'ReleaseYear', 'ItemCount', 'LowUserPrice', 'LowNetPrice']  
  
fig, ax = plt.subplots(len(col_names), figsize=(8,40))  
  
for i, col_val in enumerate(col_names):  
  
    sns.boxplot(y=sales_data_hist[col_val], ax=ax[i])  
    ax[i].set_title('Box plot - '+col_val, fontsize=10)  
    ax[i].set_xlabel(col_val, fontsize=8)  
  
plt.show()
```







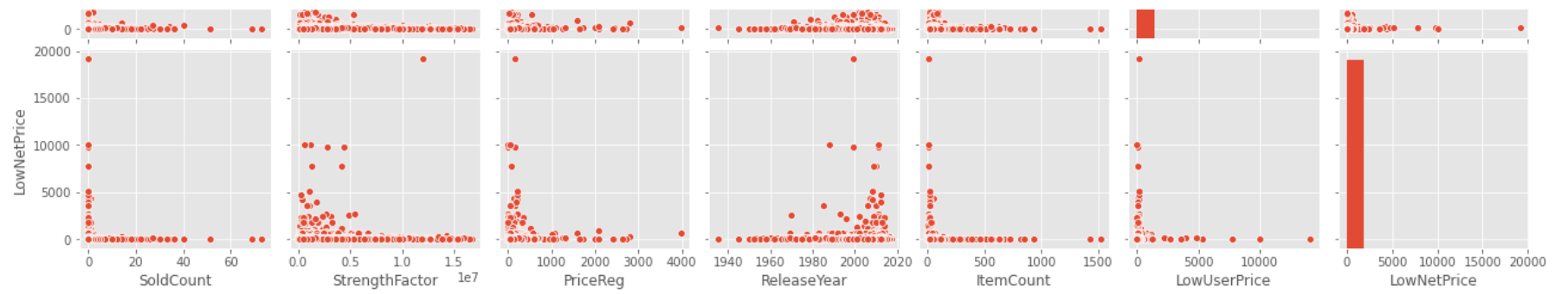




```
In [23]: sales_data_hist = sales_data_hist.drop(['Order', 'File_Type', 'SKU_number', 'SoldFlag', 'MarketingType', 'ReleaseNumber', 'New_Release_Flag'], axis=1)
sns.pairplot(sales_data_hist)
```

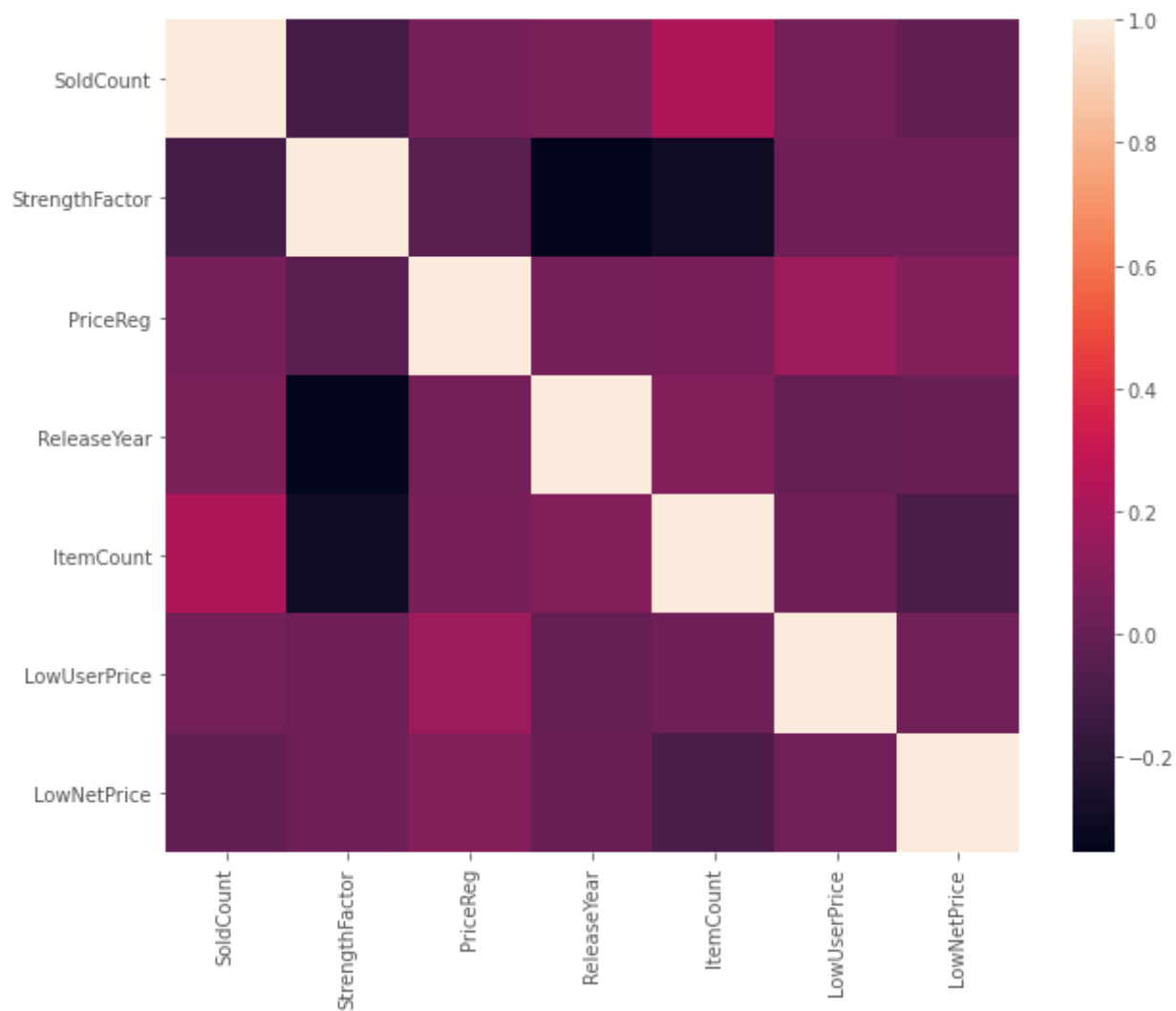
```
Out[23]: <seaborn.axisgrid.PairGrid at 0x7f9b99199850>
```





```
In [24]: f, ax = plt.subplots(figsize=(10, 8))  
corr = sales_data_hist.corr()  
sns.heatmap(corr,  
            xticklabels=corr.columns.values,  
            yticklabels=corr.columns.values)
```

```
Out[24]: <matplotlib.axes._subplots.AxesSubplot at 0x7f9b9a65e310>
```



```
In [25]: # Percentile based outlier removal
def percentile_based_outlier(data, threshold=95):
    diff = (100 - threshold) / 2.0
    minval, maxval = np.percentile(data, [diff, 100 - diff])
    return (data < minval) | (data > maxval)

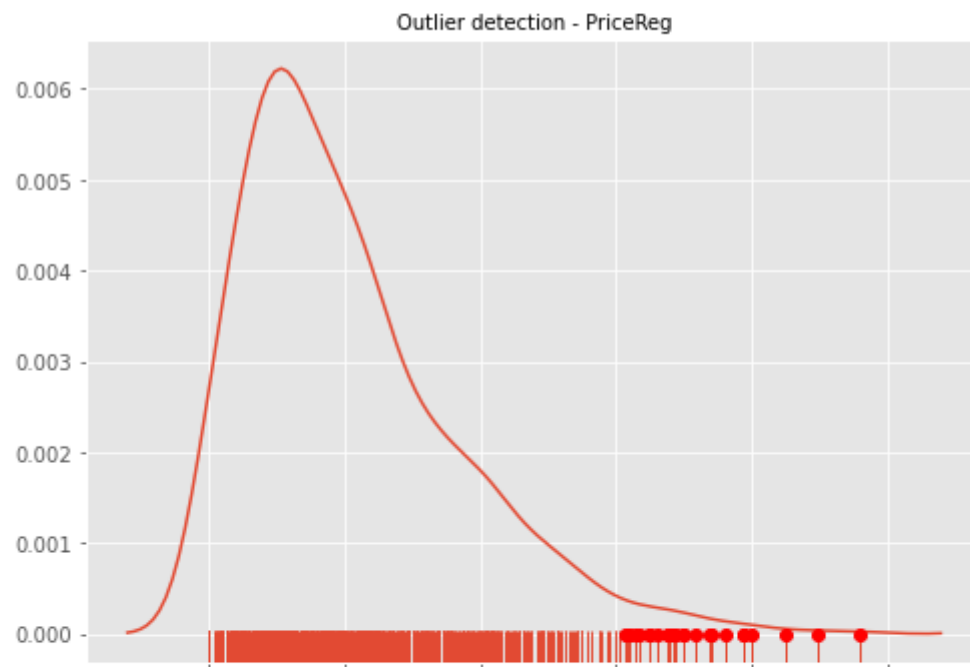
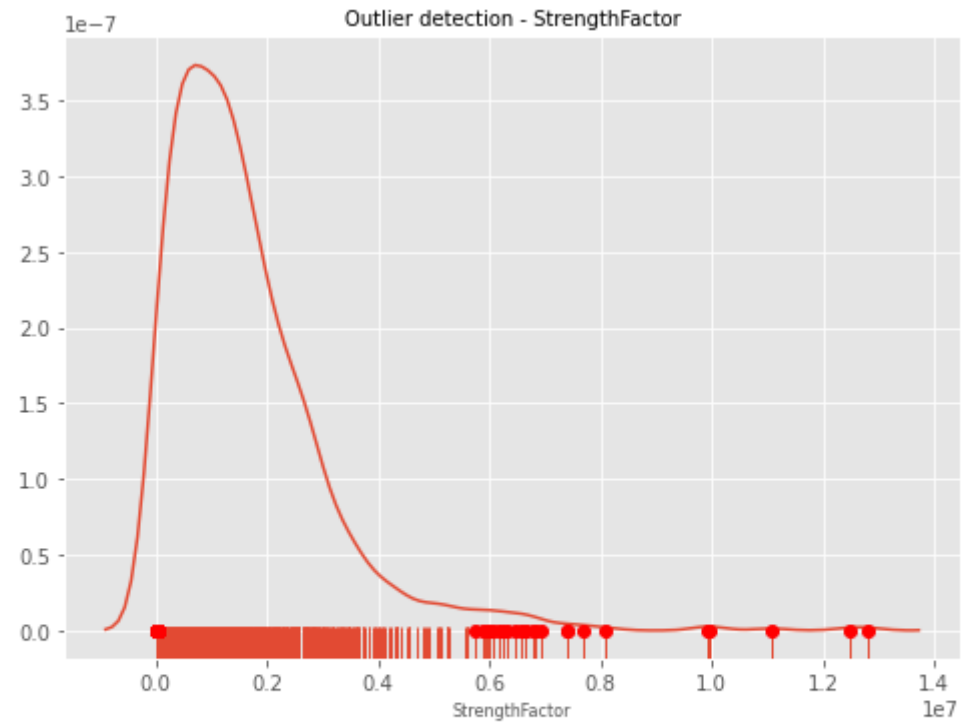
col_names = ['StrengthFactor', 'PriceReg', 'ReleaseYear', 'ItemCount', 'LowUserPrice', 'LowNetPrice']

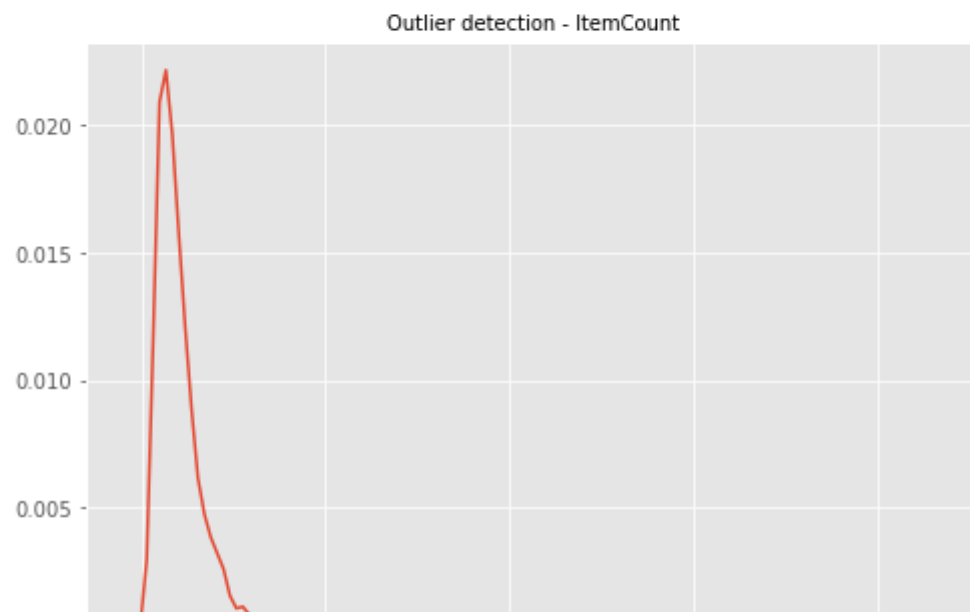
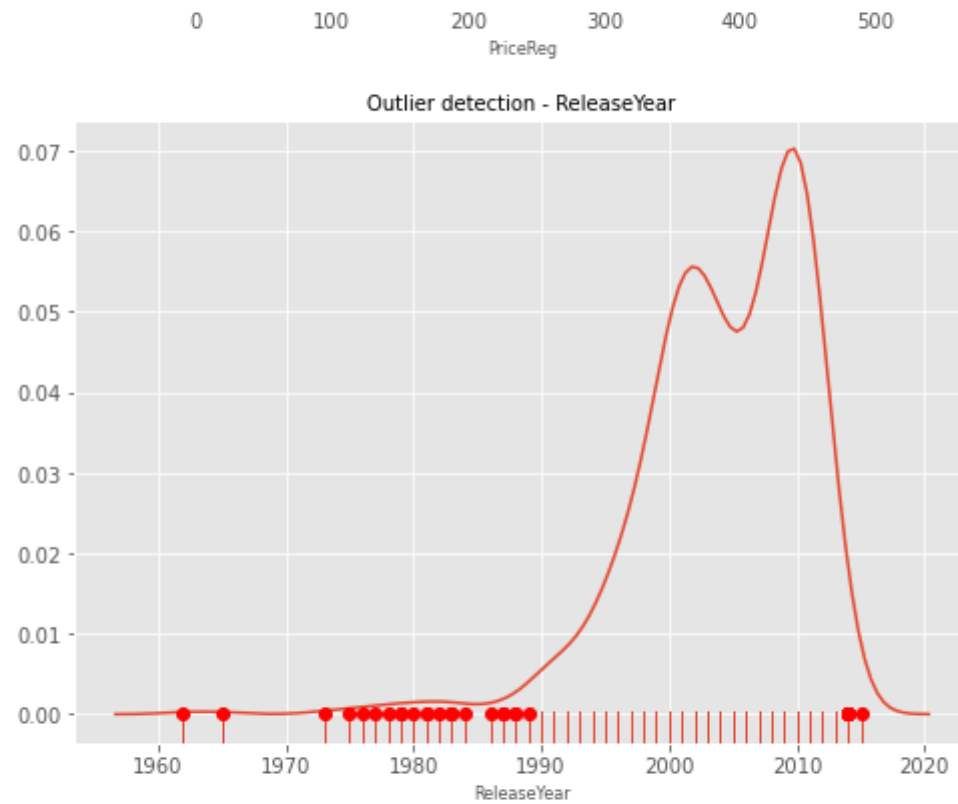
fig, ax = plt.subplots(len(col_names), figsize=(8,40))

for i, col_val in enumerate(col_names):
    x = sales_data_hist[col_val][:1000]
    sns.distplot(x, ax=ax[i], rug=True, hist=False)
    outliers = x[percentile_based_outlier(x)]
    ax[i].plot(outliers, np.zeros_like(outliers), 'ro', clip_on=False)

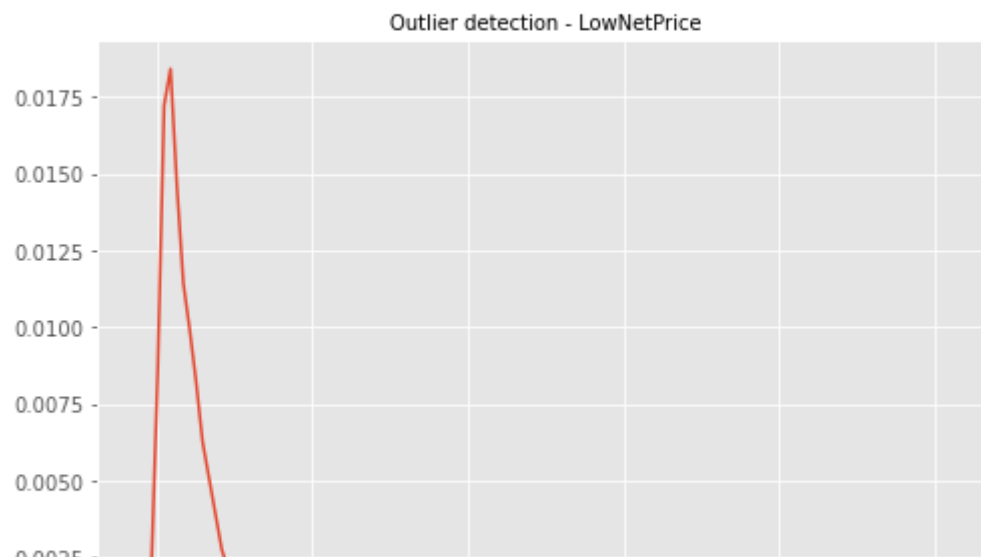
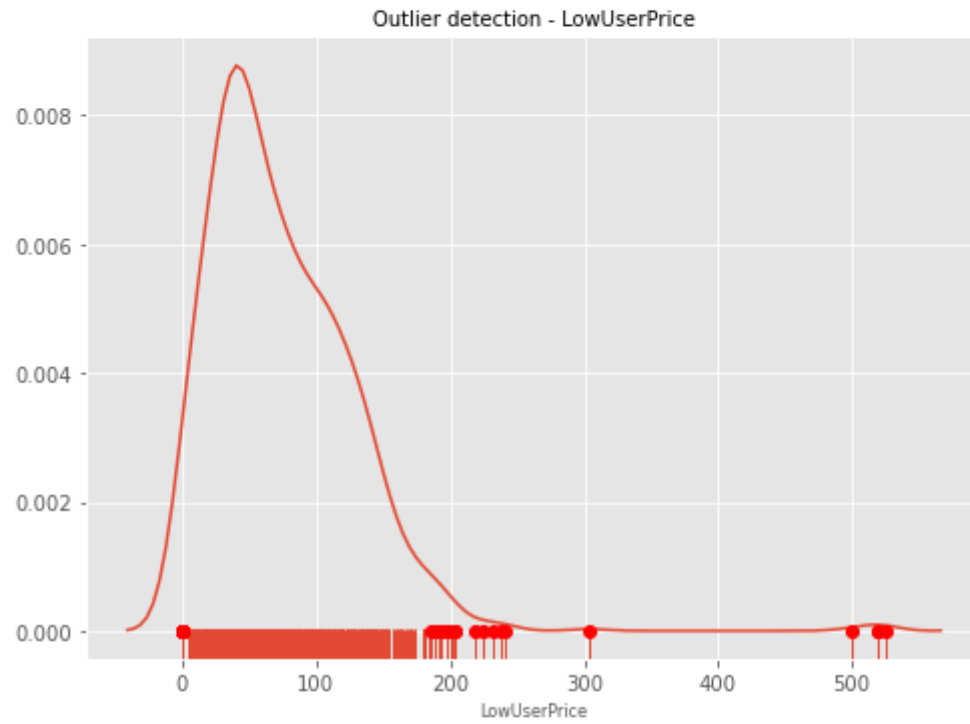
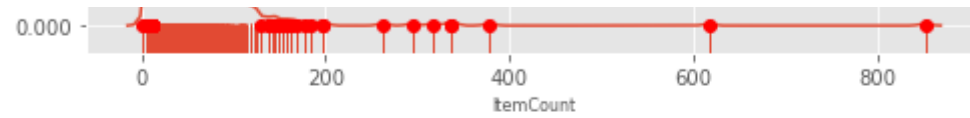
    ax[i].set_title('Outlier detection - '+col_val, fontsize=10)
    ax[i].set_xlabel(col_val, fontsize=8)

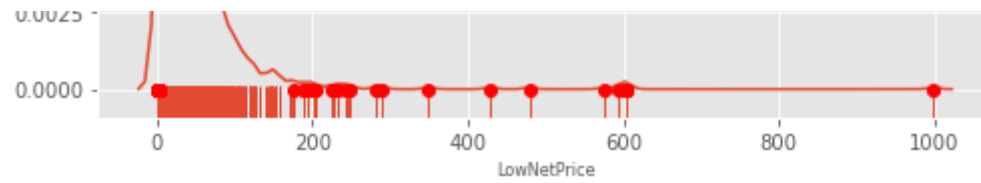
plt.show()
```











In [ ]: