

# Математичні моделі в продуктовому маркетингу

## Сегментація і кластерний аналіз

### Практична робота 2 (частина 1)

1. CRM Analytics: Customer Segmentation with RFM:

- 1.1. Завантажте датасет: <https://medium.com/@zbeyza/crm-analytics-customer-segmentation-with-rfm-208ddc10c623>
- 1.2. Порахуйте RFM Score.
- 1.3. Визначте сегменти за допомогою REGEX та їх статистичні показники і частки.

Студентка Пороскун Олена. Група ПМ-21

In [1]:  
Requirement already satisfied: squarify in c:\users\admin\anaconda3\lib\site-packages (0.4.3)

In [2]:  
import pandas as pd  
import numpy as np  
import matplotlib.pyplot as plt  
import seaborn as sns  
sns.set()  
import squarify  
  
import datetime as dt  
  
import warnings  
warnings.filterwarnings('ignore')  
  
excel\_data\_df = pandas.read\_excel('records.xlsx', sheet\_name='Employees')

In [3]:  
df = pd.read\_excel('online\_retail\_II.xlsx')  
df = df\_.copy()  
df.head()

Out[3]:

	Invoice	StockCode	Description	Quantity		InvoiceDate	Price	Customer ID	Country
0	489434	85048	15CM CHRISTMAS GLASS BALL 20 LIGHTS	12	2009-12-01 07:45:00	6.95	13085.0	United Kingdom	83.4
1	489434	79323P	PINK CHERRY LIGHTS	12	2009-12-01 07:45:00	6.75	13085.0	United Kingdom	81.0
2	489434	79323W	WHITE CHERRY LIGHTS	12	2009-12-01 07:45:00	6.75	13085.0	United Kingdom	100.8
3	489434	22041	RECORD FRAME 7" SINGLE SIZE	48	2009-12-01 07:45:00	2.10	13085.0	United Kingdom	30.0
4	489434	21232	STRAWBERRY CERAMIC TRINKET BOX	24	2009-12-01 07:45:00	1.25	13085.0	United Kingdom	

In [4]:  
# перевірка відсутнього значення  
df.isnull().sum  
  
# Клієнти з відсутнім ідентифікатором клієнта не можуть бути включені в процес сегментації, тому їх необхідно видалити з набору даних.  
df.dropna(inplace = True)

In [5]:  
# додавання загальної ціни на основі продуктів до набору даних як змінної  
df["TotalPrice"] = df[["Quantity"] \* df["Price"]  
df.head()

Out[5]:

	Invoice	StockCode	Description	Quantity		InvoiceDate	Price	Customer ID	Country	TotalPrice
	489434	85048	15CM CHRISTMAS GLASS BALL 20 LIGHTS	12	2009-12-01 07:45:00	6.95	13085.0	United Kingdom		83.4
1	489434	79323P	PINK CHERRY LIGHTS	12	2009-12-01 07:45:00	6.75	13085.0	United Kingdom		81.0
2	489434	79323W	WHITE CHERRY LIGHTS	12	2009-12-01 07:45:00	6.75	13085.0	United Kingdom		100.8
3	489434	22041	RECORD FRAME 7" SINGLE SIZE	48	2009-12-01 07:45:00	2.10	13085.0	United Kingdom		100.8
4	489434	21232	STRAWBERRY CERAMIC TRINKET BOX	24	2009-12-01 07:45:00	1.25	13085.0	United Kingdom		30.0

In [6]:  
# описова статистика набору даних  
df.describe().T  
  
# тут є від'ємні значення, через повернення, потрібно їх позбутися.

Out[6]:

	count	mean	std	min	25%	50%	75%	max
Quantity	417534.0	12.758815	101.220424	-9360.00	2.00	4.00	12.00	19152.00
Price	417534.0	3.887547	71.131797	0.00	1.25	1.95	3.75	25111.09
Customer ID	417534.0	15360.645478	1680.811316	12346.00	13983.00	15311.00	16799.00	18287.00
TotalPrice	417534.0	19.994081	99.915863	-25111.09	4.25	11.25	19.35	15818.40

In [7]:  
# в атрибуті Invoice значення, що починаються з "C", вказать подарункі товари.  
df = df[~df["Invoice"].str.contains("C", na=False)]

In [8]:  
df.describe().T

Out[8]:

	count	mean	std	min	25%	50%	75%	max
Quantity	407695.0	13.586686	96.842229	1.0	2.00	5.00	12.00	19152.0
Price	407695.0	3.294188	34.756655	0.0	1.25	1.95	3.75	10933.5
Customer ID	407695.0	15368.504107	1679.795700	12346.0	13997.00	15321.00	16812.00	18287.0
TotalPrice	407695.0	21.663261	77.147356	0.0	4.95	11.90	19.50	15818.4

Calculating RFM Metrics

In [9]:  
df["InvoiceDate"].max()  
analysis\_date = dt.datetime(2023,6,5)

In [10]:  
rfm = df.groupby("Customer ID").agg({"InvoiceDate": lambda invoice\_date: (analysis\_date - invoice\_date.max()).days,  
"Invoice": lambda invoice: invoice.nunique(),  
"TotalPrice": lambda total\_price: total\_price.sum()})  
rfm.columns = ["recency", "frequency", "monetary"]  
rfm.head()

Out[10]:

	recency	frequency	monetary
Customer ID			
12346.0	4724	11	372.86
12347.0	4562	2	1323.32
12348.0	4633	1	222.16
12349.0	4602	3	2671.14
12351.0	4570	1	300.93

In [11]:  
# Це небезпечна ситуація, коли монетарні показники дорівнюють 0, тому їх слід усунути.  
rfm = rfm[rfm["monetary"] > 0]  
rfm.describe().T

Out[11]:

	count	mean	std	min	25%	50%	75%	max
recency	4312.0	4650.172542	96.861457	4560.00	4577.0000	4612.00	4695.0000	4933.00
frequency	4312.0	4.455705	8.170213	1.00	1.0000	2.00	5.0000	205.00
monetary	4312.0	2048.238236	8914.481280	2.95	307.9875	706.02	1723.1425	349164.35

Calculating RFM Scores

In [12]:  
rfm["recency\_score"] = pd.qcut(rfm["recency"], 5, labels=[5, 4, 3, 2, 1])  
rfm["frequency\_score"] = pd.qcut(rfm["frequency"],rank(method="first"), 5, labels=[1, 2, 3, 4, 5])  
rfm["monetary\_score"] = pd.qcut(rfm["monetary"], 5, labels=[1, 2, 3, 4, 5])  
  
rfm["RF\_SCORE"] = (rfm["recency\_score"].astype(str) + rfm["frequency\_score"].astype(str))  
rfm.head()

Out[12]:

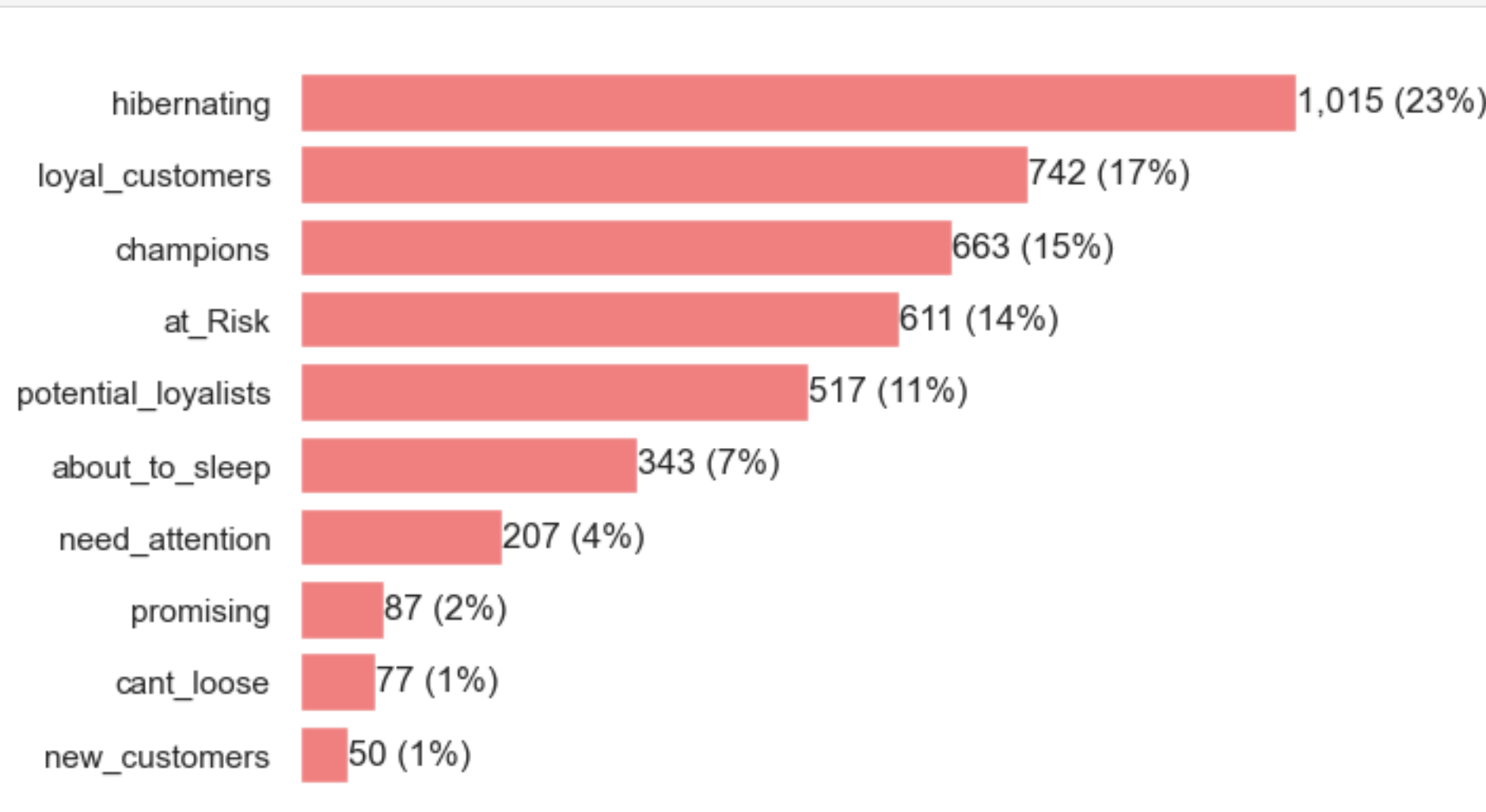
	recency	frequency	monetary	recency_score	frequency_score	monetary_score	RF_SCORE
Customer ID							
12346.0	4724	11	372.86	2	5	2	25
12347.0	4562	2	1323.32	5	2	4	52
12348.0	4633	1	222.16	2	1	1	21
12349.0	4602	3	2671.14	3	3	5	33
12351.0	4570	1	300.93	5	1	2	51

In [13]:  
# перетворення RF\_SCORE в сегмент за допомогою Regex  
seg\_map = {  
r"[1-2][1-2]": 'hibernating',  
r"[1-2][3-4]": 'at\_Risk',  
r"[1-2]5": 'cant\_loose',  
r"[3-5]2": 'about\_to\_sleep',  
r"33": 'need\_attention',  
r"[3-4][4-5]": 'loyal\_customers',  
r"41": 'promising',  
r"51": 'new\_customers',  
r"[4-5][2-3]": 'potential\_loyalists',  
r"5[4-5]": 'champions'  
}  
  
rfm['segment'] = rfm['RF\_SCORE'].replace(seg\_map, regex=True)  
rfm.head()

Out[13]:

	recency	frequency	monetary	recency_score	frequency_score	monetary_score	RF_SCORE	segment
Customer ID								
12346.0	4724	11	372.86	2	5	2	25	cant_loose
12347.0	4562	2	1323.32	5	2	4	52	potential_loyalists
12348.0	4633	1	222.16	2	1	1	21	hibernating
12349.0	4602	3	2671.14	3	3	5	33	need_attention
12351.0	4570	1	300.93	5	1	2	51	new_customers

In [14]:  
# скільки сегментів і який їхній відсоток  
segments\_counts = rfm["segment"].value\_counts().sort\_values(ascending=True)  
  
fig, ax = plt.subplots()  
  
bars = ax.barh(range(len(segments\_counts)),  
segments\_counts,  
color='lightcoral')  
ax.set\_frame\_on(False)  
ax.tick\_params(left=False,  
bottom=False,  
labelbottom=False)  
ax.set\_yticks(range(len(segments\_counts)))  
ax.set\_yticklabels(segments\_counts.index)  
  
for i, bar in enumerate(bars):  
value = bar.get\_width()  
if segments\_counts.index[i] in ['Can\'t loose']:  
bar.set\_color('firebrick')  
bar.set\_color('firebrick')  
ax.text(value,  
bar.get\_y() + bar.get\_height()/2,  
{:,} ({:}%".format(int(value),  
int(value\*100/segments\_counts.sum()))),  
va='center',  
ha='left')  
)  
  
plt.show(block=True)



In [15]:  
segments\_counts.index

Out[15]:  
Index(['new\_customers', 'cant\_loose', 'promising', 'need\_attention',  
'about\_to\_sleep', 'potential\_loyalists', 'at\_Risk', 'champions',  
'loyal\_customers', 'hibernating'],  
dtype='object')

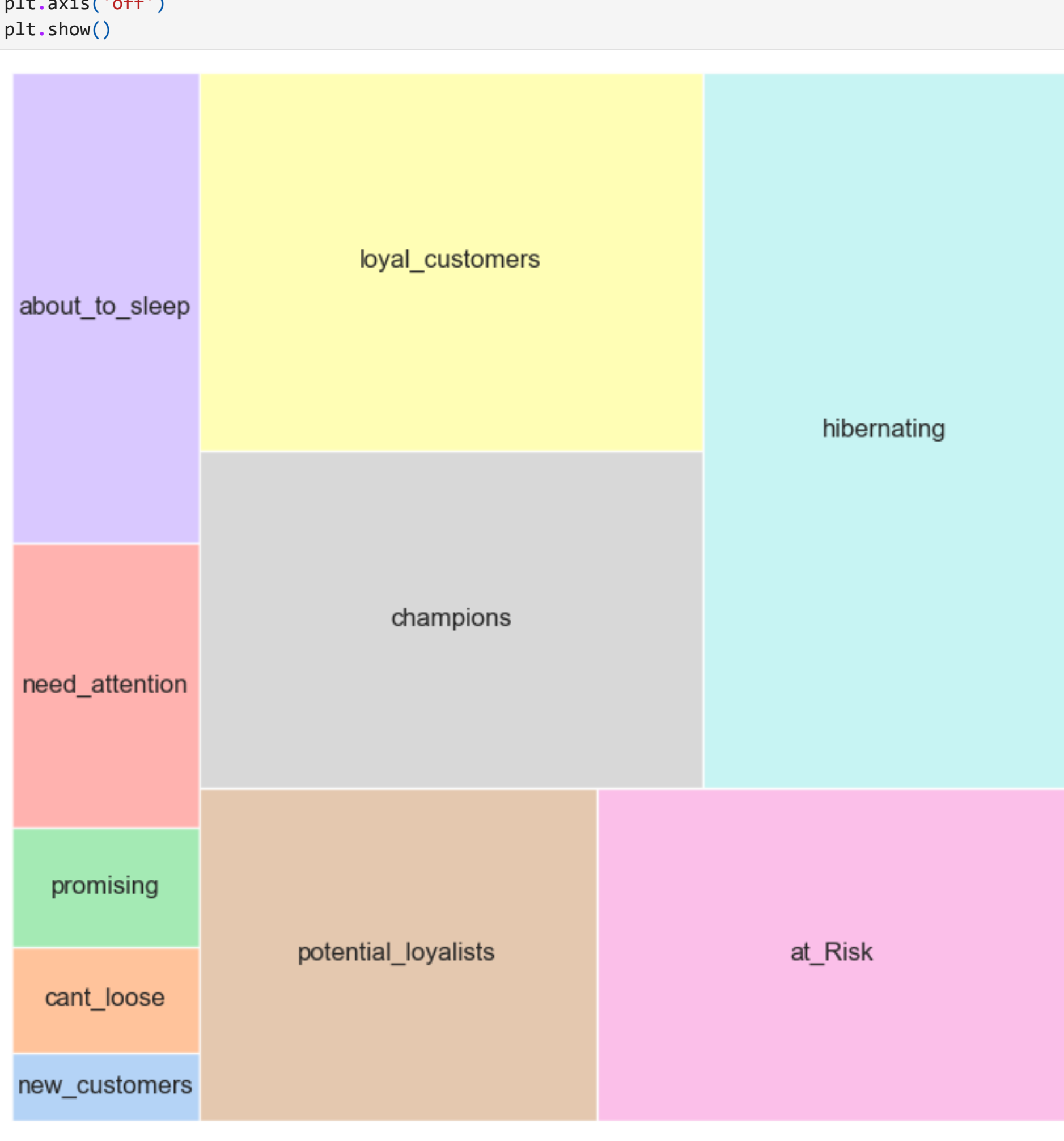
In [16]:  
segments\_counts.values

Out[16]:  
array([ 50, 77, 87, 207, 343, 517, 611, 663, 742, 1015],  
dtype=int64)

In [17]:  
segments\_counts.describe()

Out[17]:  
count 10.000000  
mean 431.298000  
std 329.668116  
min 50.000000  
25% 117.000000  
50% 430.000000  
75% 658.000000  
max 1015.000000  
Name: segment, dtype: float64

In [18]:  
fig, ax = plt.subplots(1, figsize = (9,9))  
colors = sns.color\_palette('pastel')[0:11]  
squarify.plot(sizes=segments\_counts.values,  
label=segments\_counts.index,  
color = colors,  
alpha=.8)  
  
plt.axis('off')  
plt.show()



In [19]:  
plt.figure(figsize = (10,10))  
colors = sns.color\_palette('pastel')[0:11]  
plt.pie(segments\_counts.values,  
labels=segments\_counts.index,  
color = colors,  
autopct='%0.1f%%')  
plt.show()

