

# Automate Exploratory Data Analysis

- pandas profiling
- sweetviz
- AutoViz

In [2]:

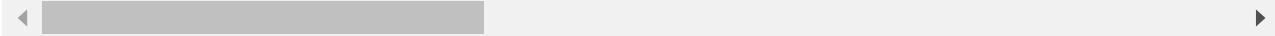
```
import pandas as pd
import numpy as np
```

In [3]:

```
data = pd.read_csv(r"C:\Users\shail\Google Drive\Self\Learning\iNuerons\Project\Dataset\data1.csv")
data.head()
```

Out[3]:

	ProdTaken	Age	TypeofContact	CityTier	DurationOfPitch	Occupation	Gender	NumberOfPersonVisiting
0	1	41.0	Self Enquiry	3	6.0	Salaried	Female	3
1	0	49.0	Company Invited	1	14.0	Salaried	Male	3
2	1	37.0	Self Enquiry	1	8.0	Free Lancer	Male	3
3	0	33.0	Company Invited	1	9.0	Salaried	Female	2
4	0	36.0	Self Enquiry	1	8.0	Small Business	Male	2



In [4]:

```
data.shape
```

Out[4]:

```
(4888, 19)
```

## Pandas Profiling

In [5]:

```
import pandas_profiling
```

In [6]:

```
from pandas_profiling import profile_report
```

In [7]:

```
data.profile_report()
```

# Overview

## Dataset statistics

<b>Number of variables</b>	19
<b>Number of observations</b>	4888
<b>Missing cells</b>	0
<b>Missing cells (%)</b>	0.0%
<b>Duplicate rows</b>	141
<b>Duplicate rows (%)</b>	2.9%
<b>Total size in memory</b>	725.7 KiB
<b>Average record size in memory</b>	152.0 B

## Variable types

<b>Categorical</b>	14
<b>Numeric</b>	5

## Alerts

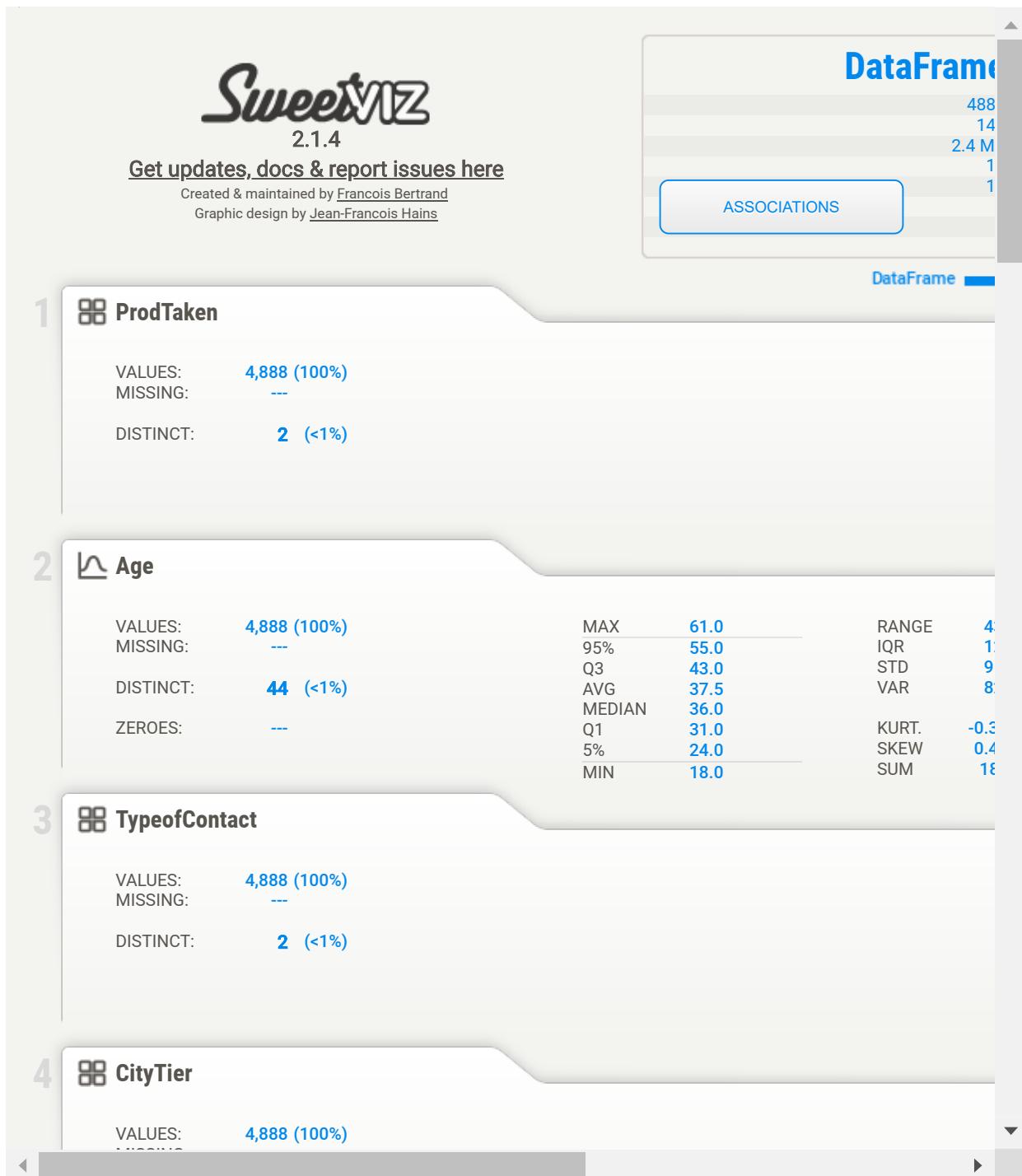
Dataset has 141 (2.9%) duplicate rows	Duplicates
NumberOfPersonVisiting is highly correlated with NumberOfChildrenVisiting	High correlation
NumberOfChildrenVisiting is highly correlated with	High correlation

Out[7]:

sweetviz

In [ ]: pip install sweetviz

In [9]:  
`import sweetviz as sv  
analyze_report = sv.analyze(data)  
analyze_report.show_notebook()`



```
In [ ]: pip install autoviz
```

```
In [11]: from autoviz.AutoViz_Class import AutoViz_Class
AV = AutoViz_Class()
df = AV.AutoViz(r"C:\Users\shail\Google Drive\Self\Learning\iNuerons\Project\Dataset\data1\Tr
```

Imported v0.1.55. After importing, execute '%matplotlib inline' to display charts in Jupyter.

```
AV = AutoViz_Class()
dfte = AV.AutoViz(filename, sep=',', depVar='', dfte=None, header=0, verbose=1, lowess=False,
chart_format='svg', max_rows_analyzed=150000, max_cols_analyzed=30, save_plot_di
```

r=None)  
 Update: verbose=0 displays charts in your local Jupyter notebook.  
 verbose=1 additionally provides EDA data cleaning suggestions. It also displays charts.  
 verbose=2 does not display charts but saves them in AutoViz\_Plots folder in local machine.  
 chart\_format='bokeh' displays charts in your local Jupyter notebook.  
 chart\_format='server' displays charts in your browser: one tab for each chart type  
 chart\_format='html' silently saves interactive HTML files in your local machine

Shape of your Data Set loaded: (4888, 20)

```
#####
##### C L A S S I F Y I N G V A R I A B L E S #####
#####
Classifying variables in data set...
```

	Nuniques	dtype	Nulls	Nullpercent	NuniquePercent	Value counts	Data cleaning improvement suggestions
						Min	
<b>CustomerID</b>	4888	int64	0	0.000000	100.000000	0	possible ID column: drop
<b>MonthlyIncome</b>	2475	float64	233	4.766776	50.634206	0	fill missing values, skewed column: cap or drop possible outliers
<b>Age</b>	44	float64	226	4.623568	0.900164	0	fill missing values
<b>DurationOfPitch</b>	34	float64	251	5.135025	0.695581	0	fill missing values, skewed column: cap or drop possible outliers
<b>NumberOfTrips</b>	12	float64	140	2.864157	0.245499	0	fill missing values, skewed column: cap or drop possible outliers
<b>NumberOfFollowups</b>	6	float64	45	0.920622	0.122750	0	fill missing values
<b>Designation</b>	5	object	0	0.000000	0.102291	230	
<b>PitchSatisfactionScore</b>	5	int64	0	0.000000	0.102291	0	
<b>ProductPitched</b>	5	object	0	0.000000	0.102291	230	
<b>NumberOfPersonVisiting</b>	5	int64	0	0.000000	0.102291	0	
<b>MaritalStatus</b>	4	object	0	0.000000	0.081833	682	
<b>Occupation</b>	4	object	0	0.000000	0.081833	2	
<b>NumberOfChildrenVisiting</b>	4	float64	66	1.350245	0.081833	0	fill missing values
<b>Gender</b>	3	object	0	0.000000	0.061375	155	
<b>CityTier</b>	3	int64	0	0.000000	0.061375	0	

	Nuniques	dtype	Nulls	Nullpercent	NuniquePercent	Value counts	Data cleaning improvement suggestions
						Min	
<b>PreferredPropertyStar</b>	3	float64	26	0.531915	0.061375	0	fill missing values
<b>ProdTaken</b>	2	int64	0	0.000000	0.040917	0	
<b>TypeofContact</b>	2	object	25	0.511457	0.040917	1419	fill missing values, fix mixed data types
<b>Passport</b>	2	int64	0	0.000000	0.040917	0	
<b>OwnCar</b>	2	int64	0	0.000000	0.040917	0	

20 Predictors classified...

1 variables removed since they were ID or low-information variables

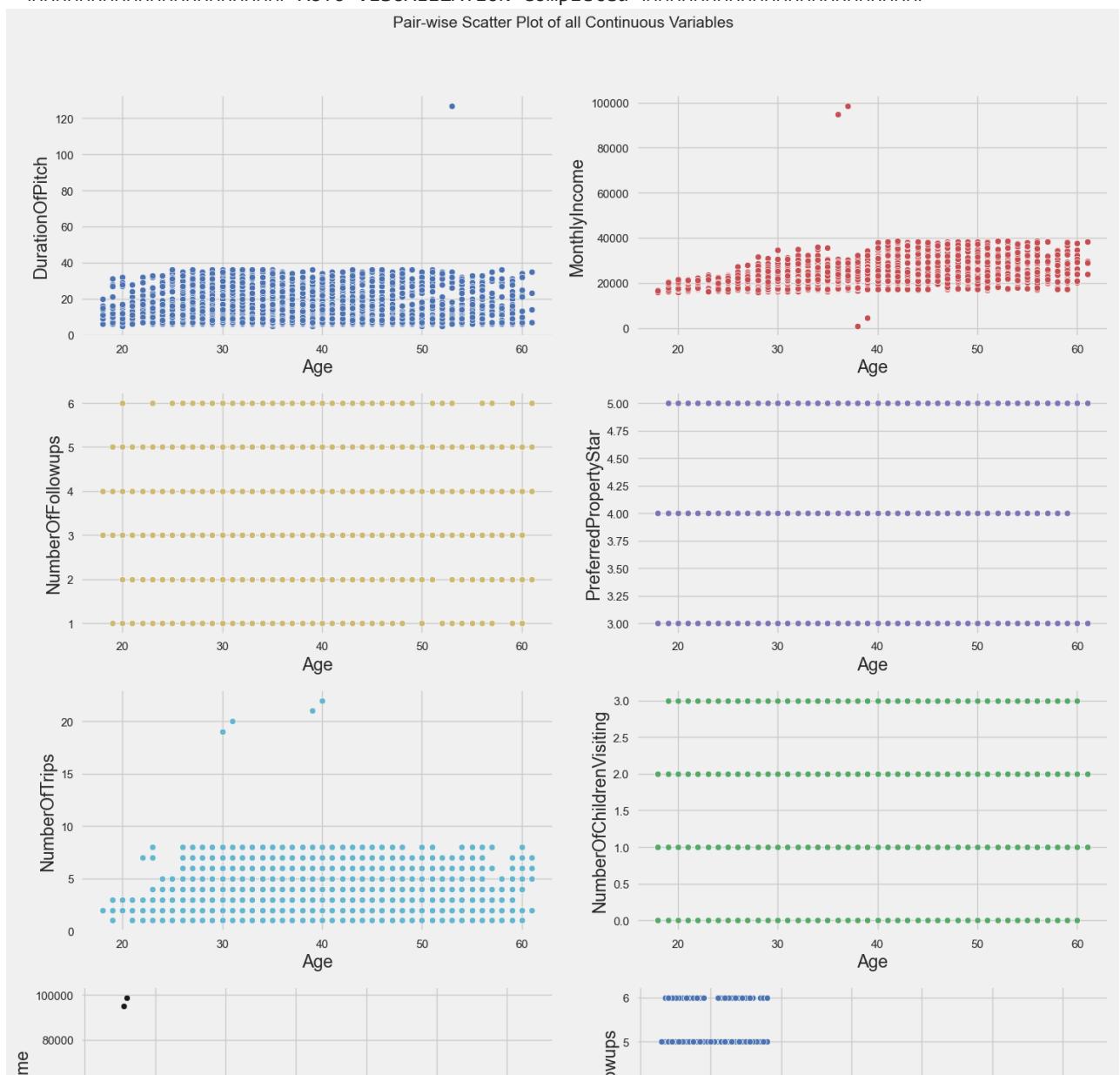
List of variables removed: ['CustomerID']

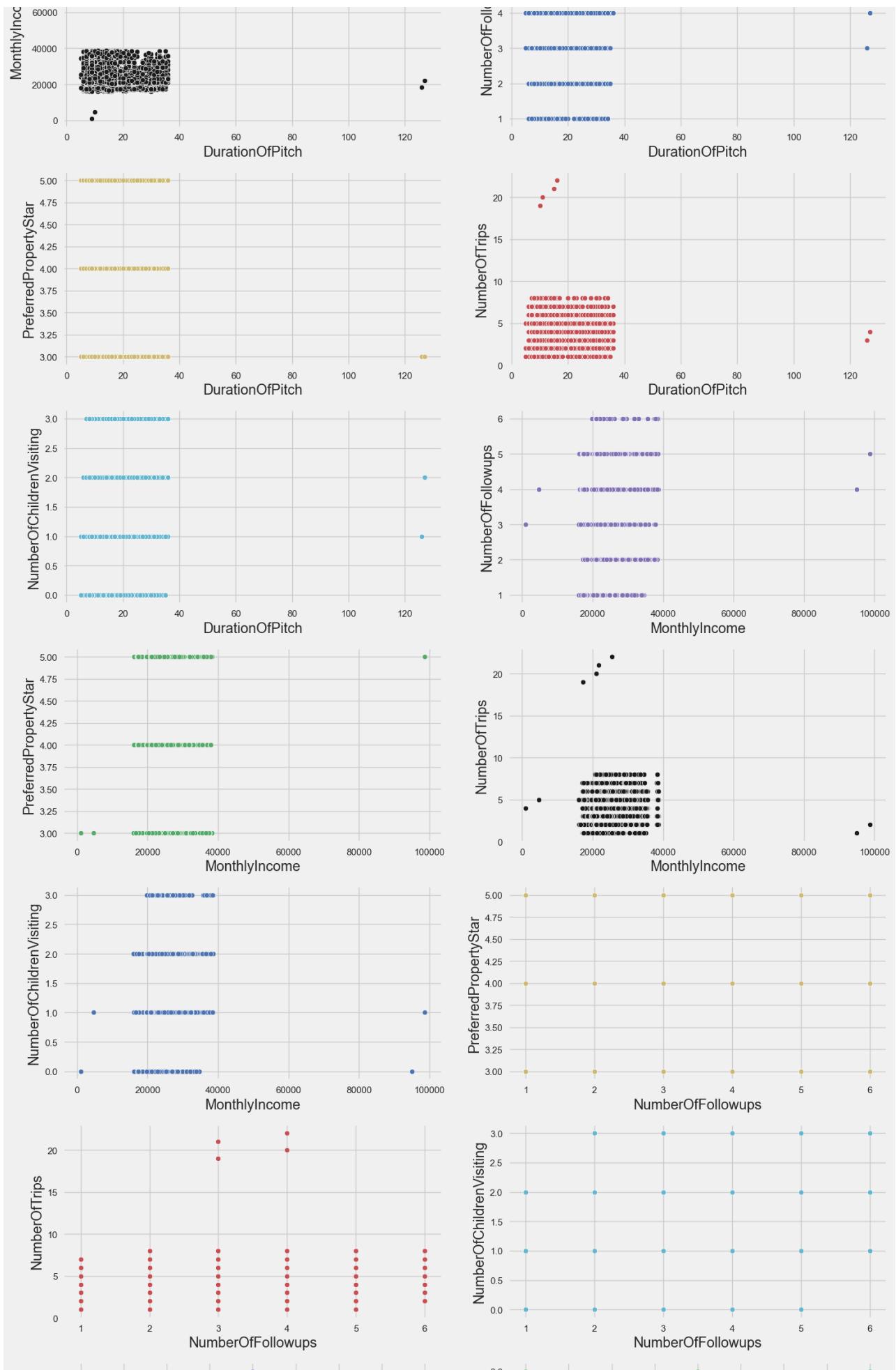
Number of All Scatter Plots = 28

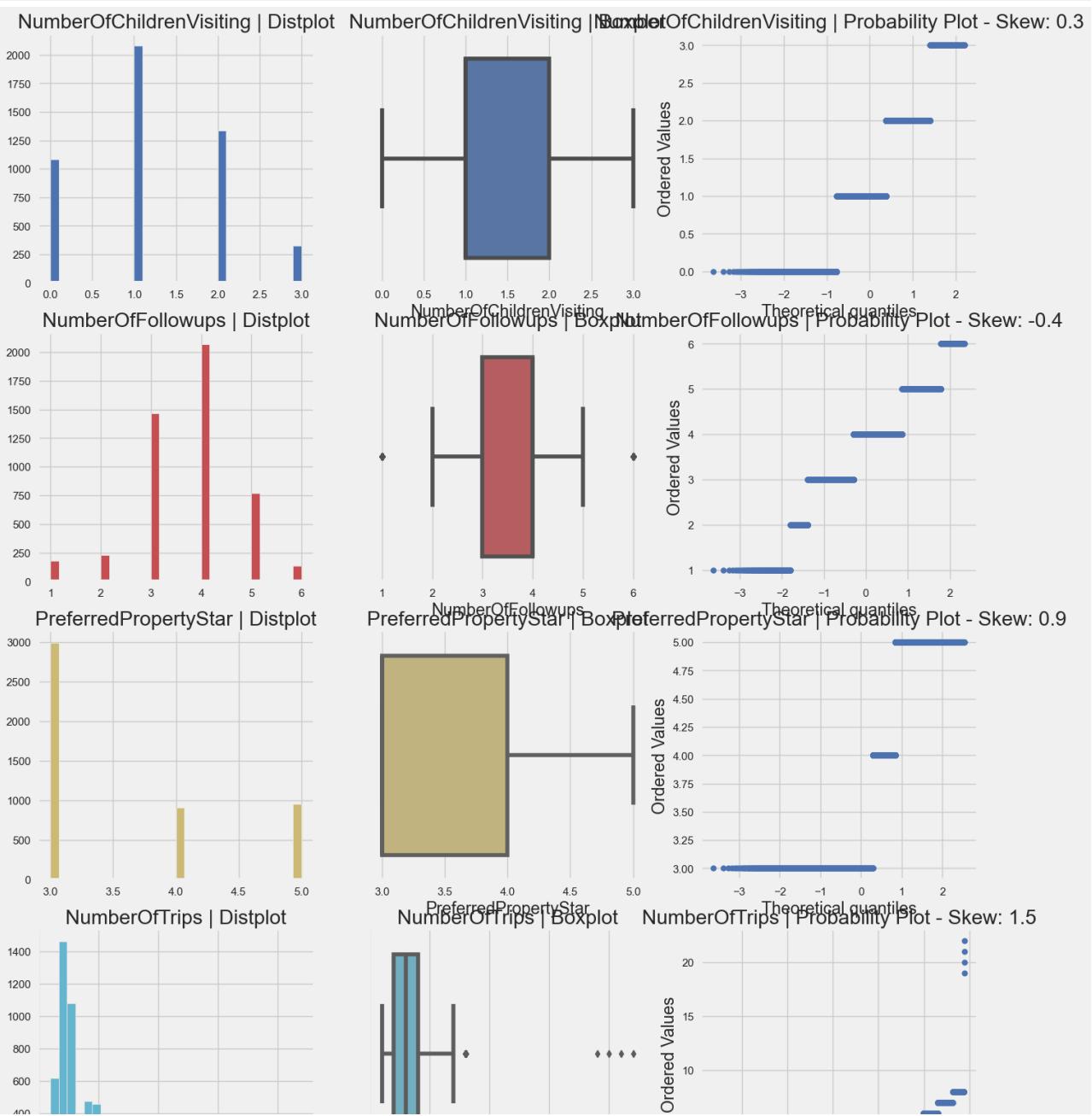
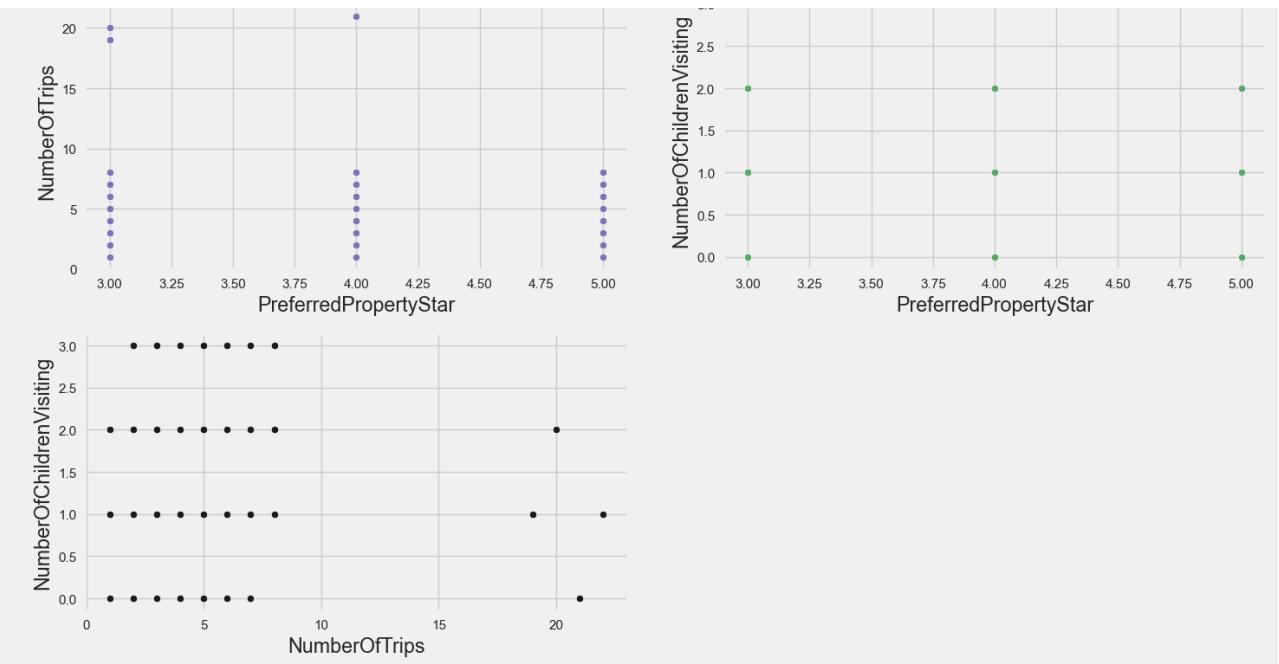
All Plots done

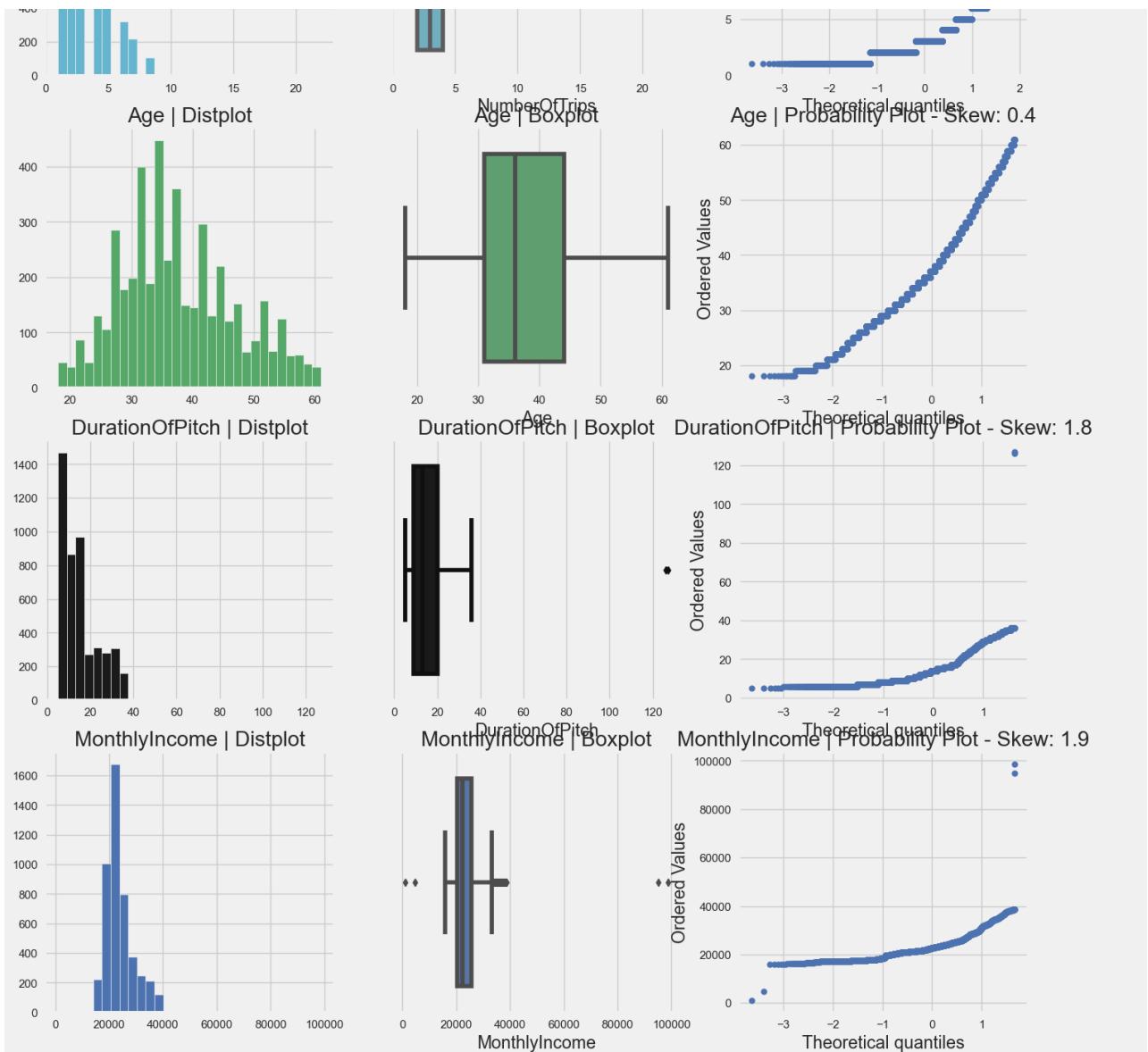
Time to run AutoViz = 6 seconds

##### AUTO VISUALIZATION Completed #####

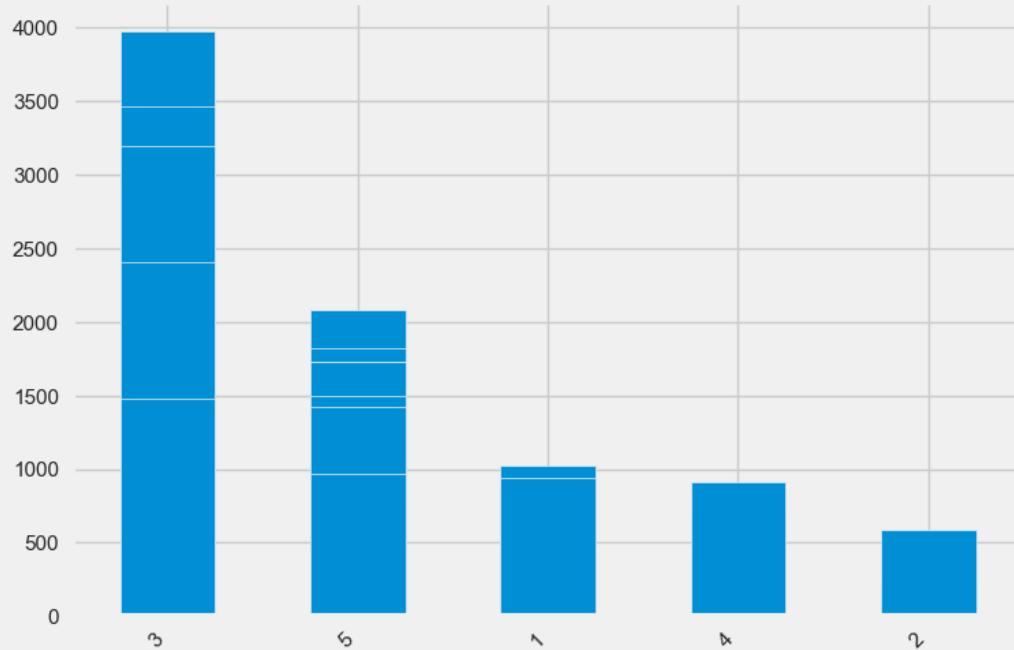




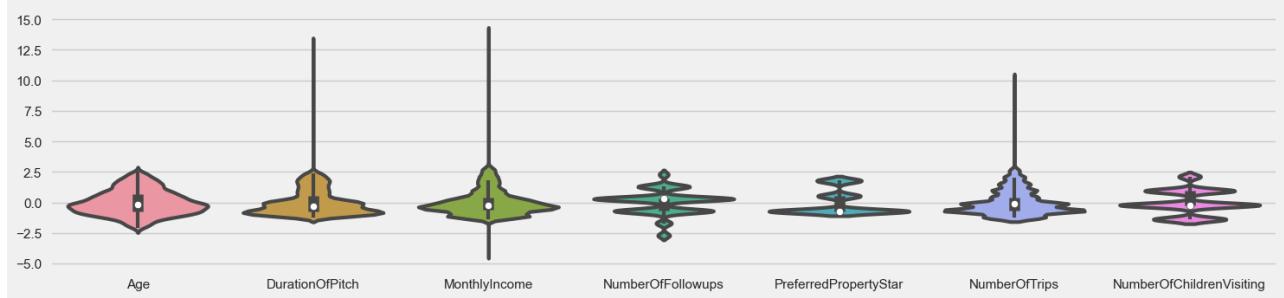


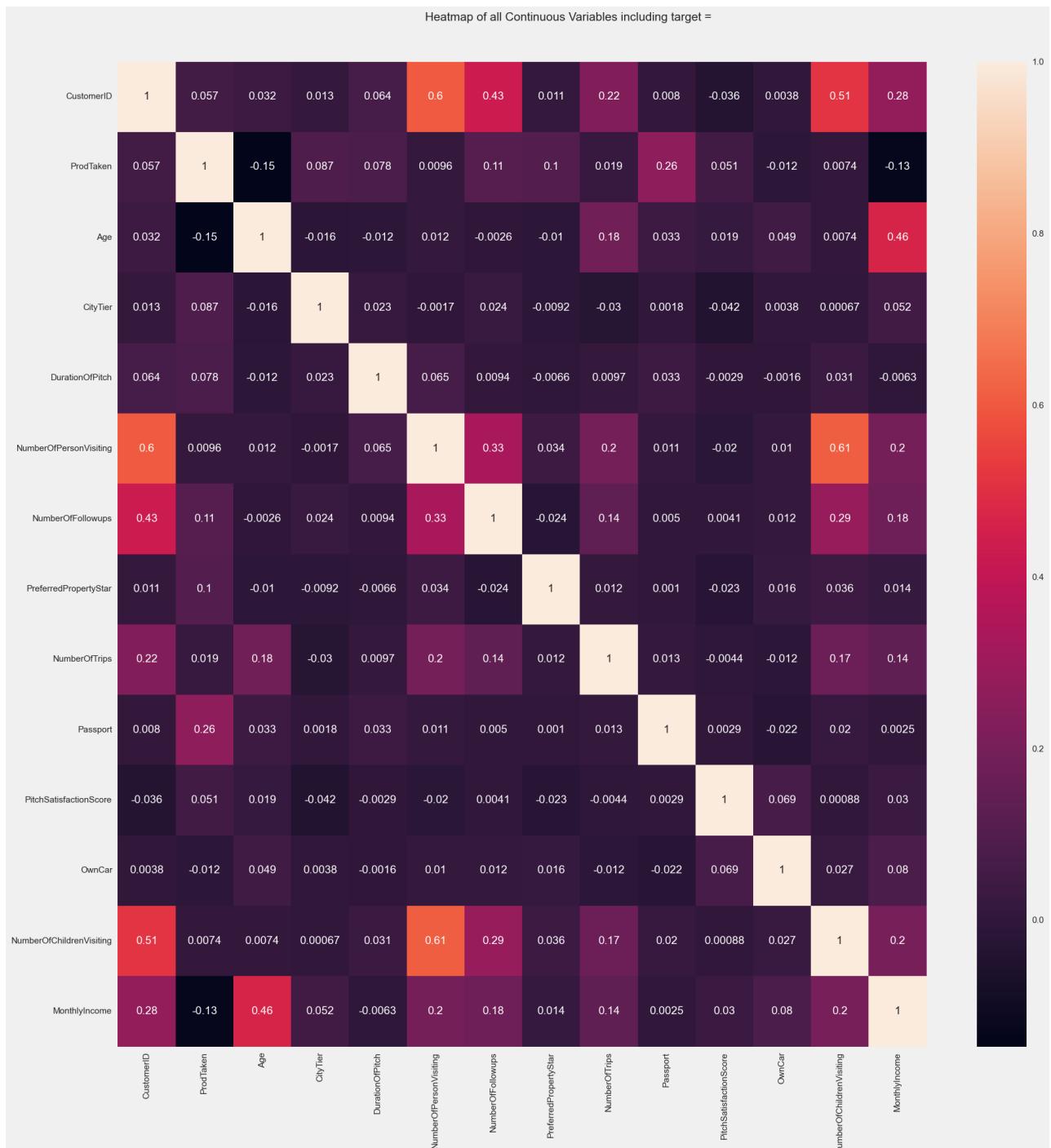


## Distribution of PitchSatisfactionScore (top 15 categories only)



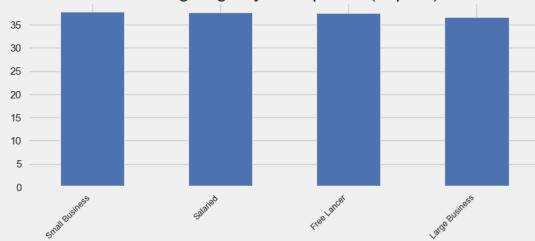
Violin Plot of all Continuous Variables



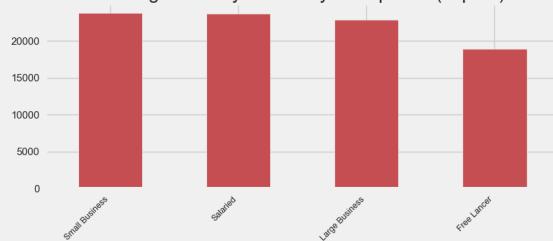


Bar plots for each Continuous by each Categorical variable

Average Age by Occupation (Top 20)



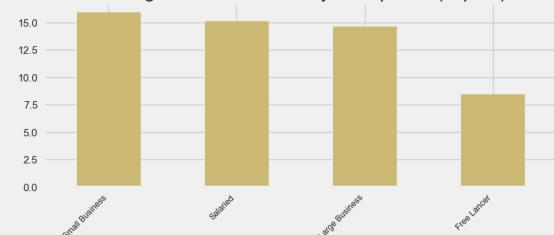
Average MonthlyIncome by Occupation (Top 20)



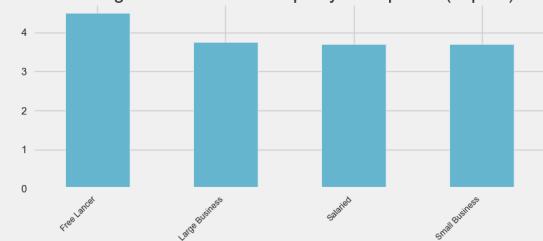
Average PreferredPropertyStar by Occupation (Top 20)



Average DurationOfPitch by Occupation (Top 20)

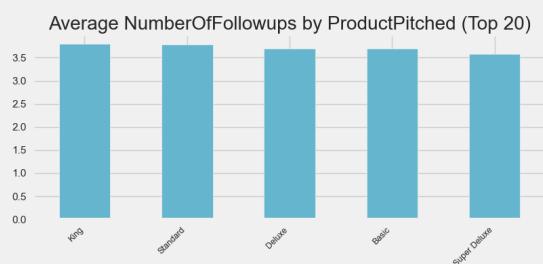
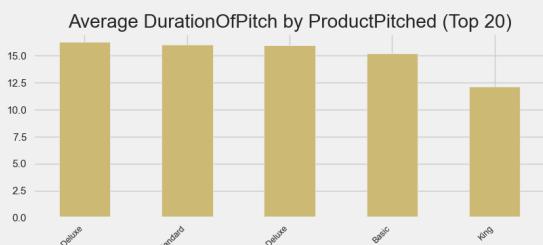
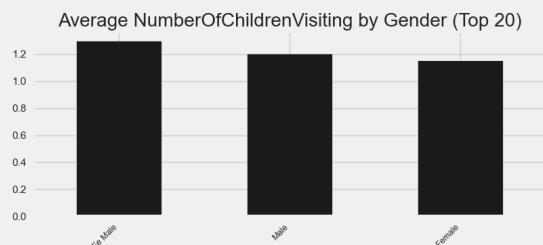
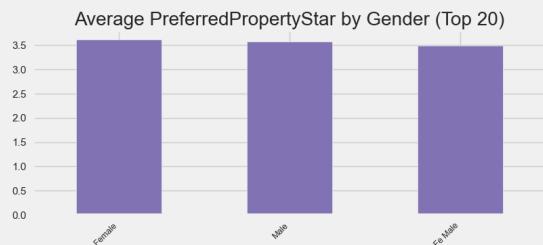
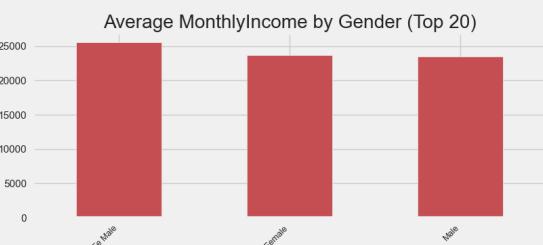
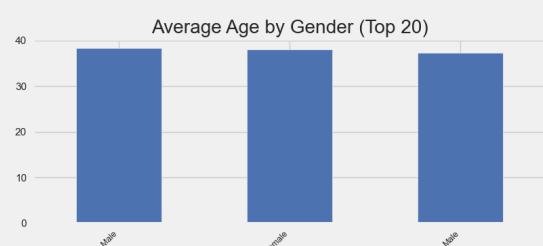
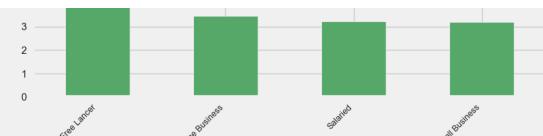
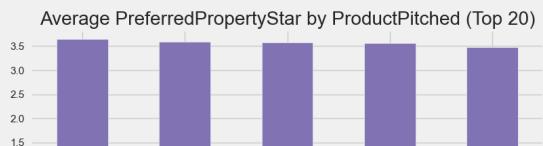
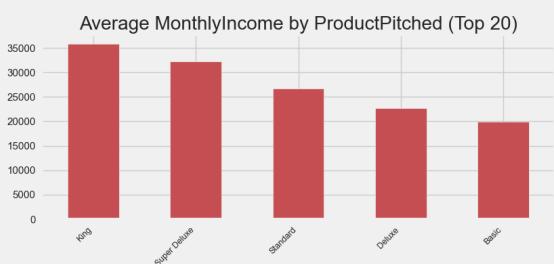
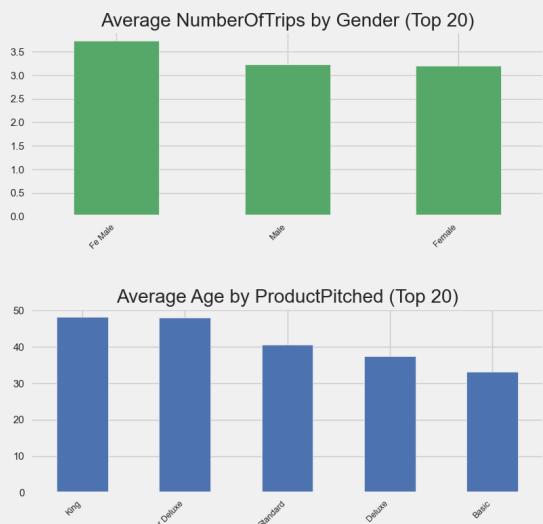
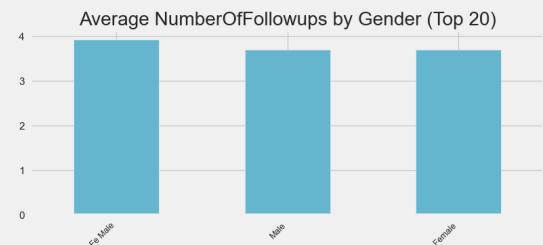
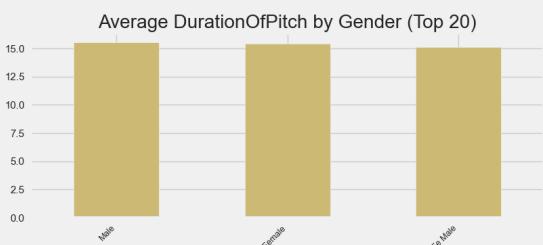
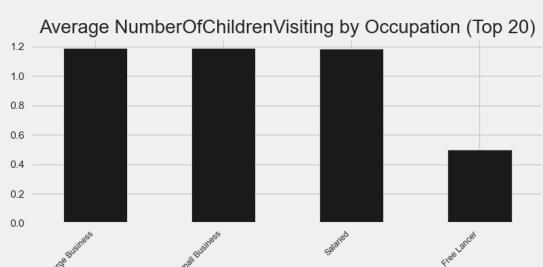
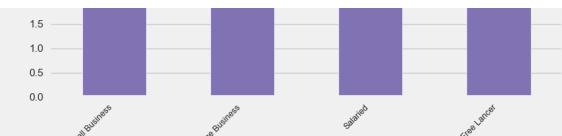


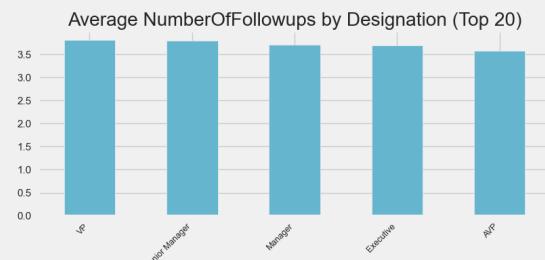
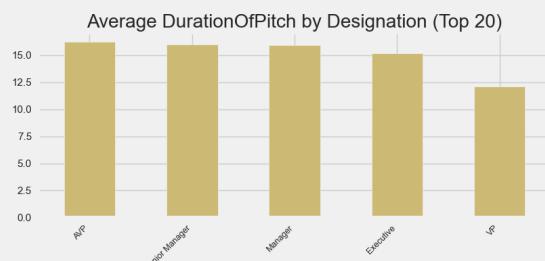
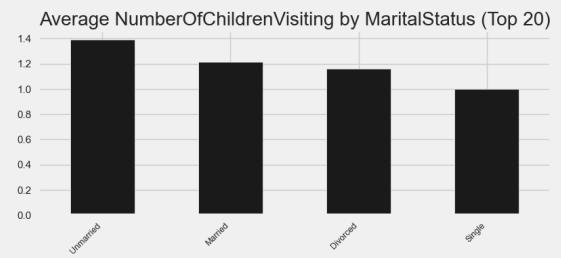
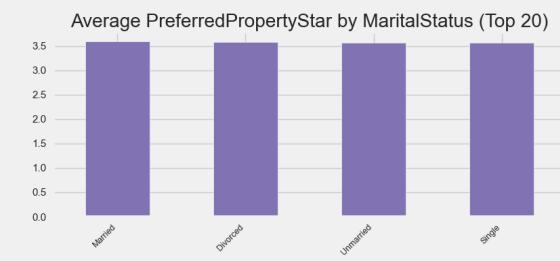
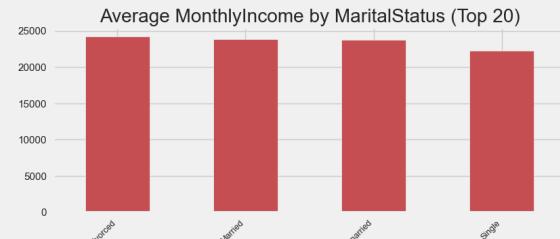
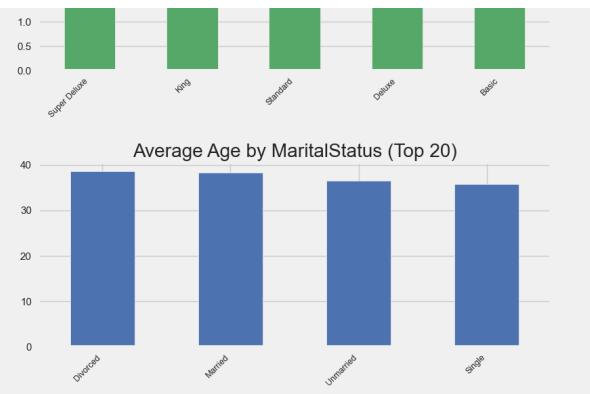
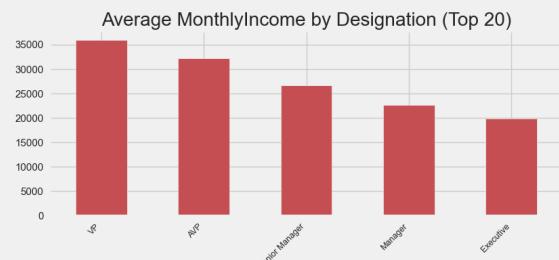
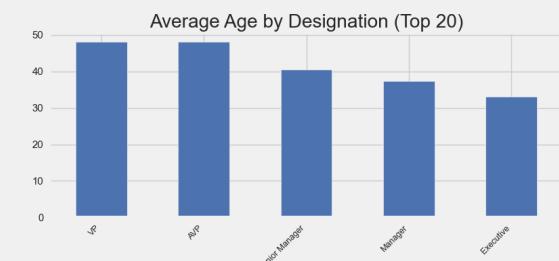
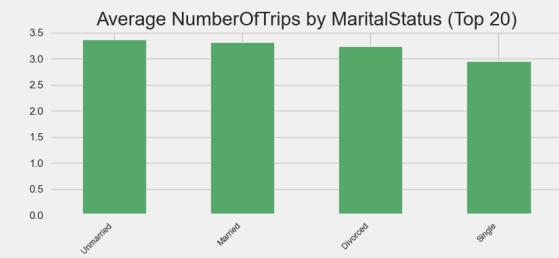
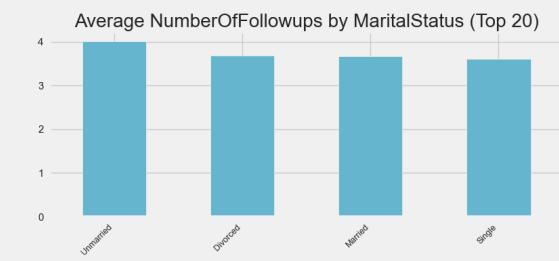
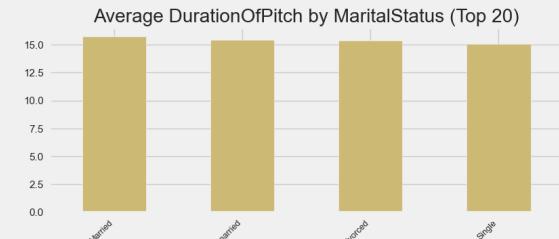
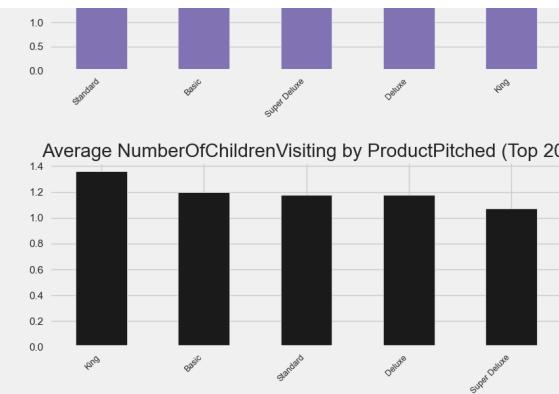
Average NumberOfFollowups by Occupation (Top 20)

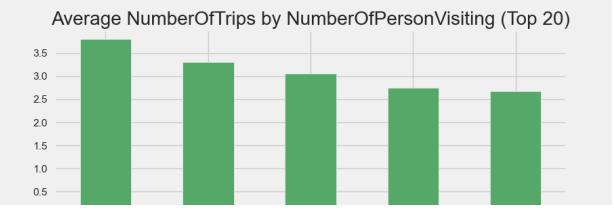
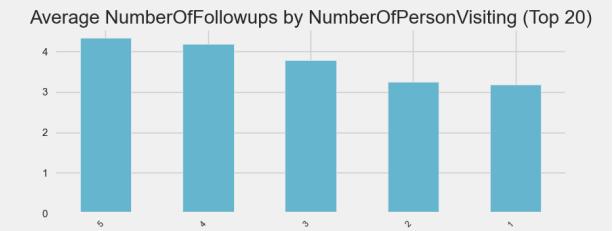
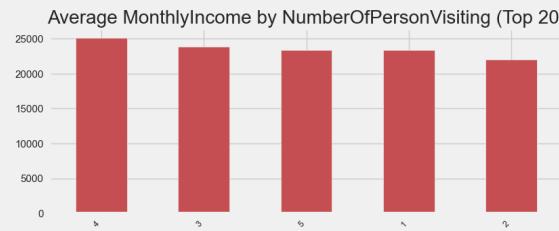
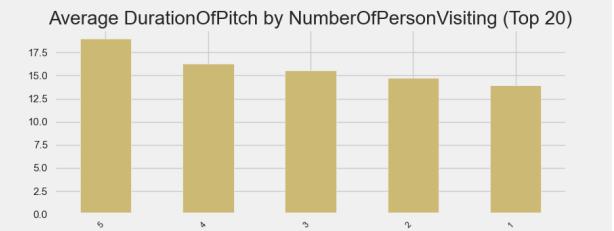
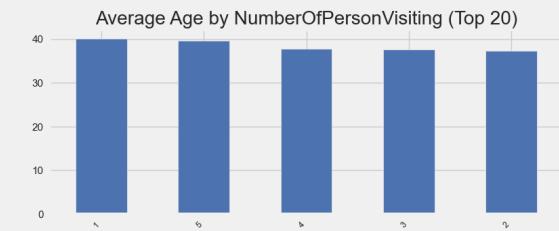
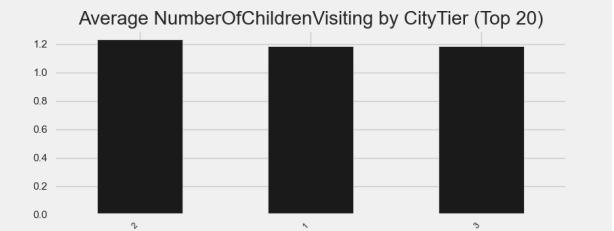
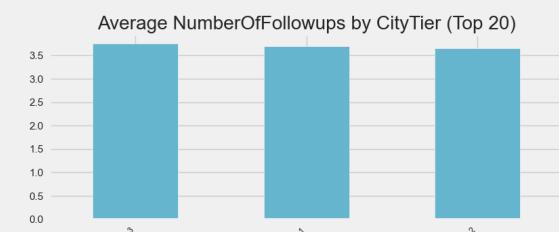
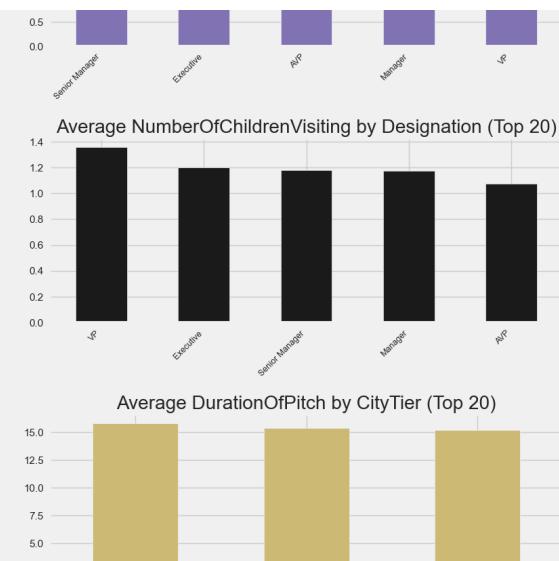


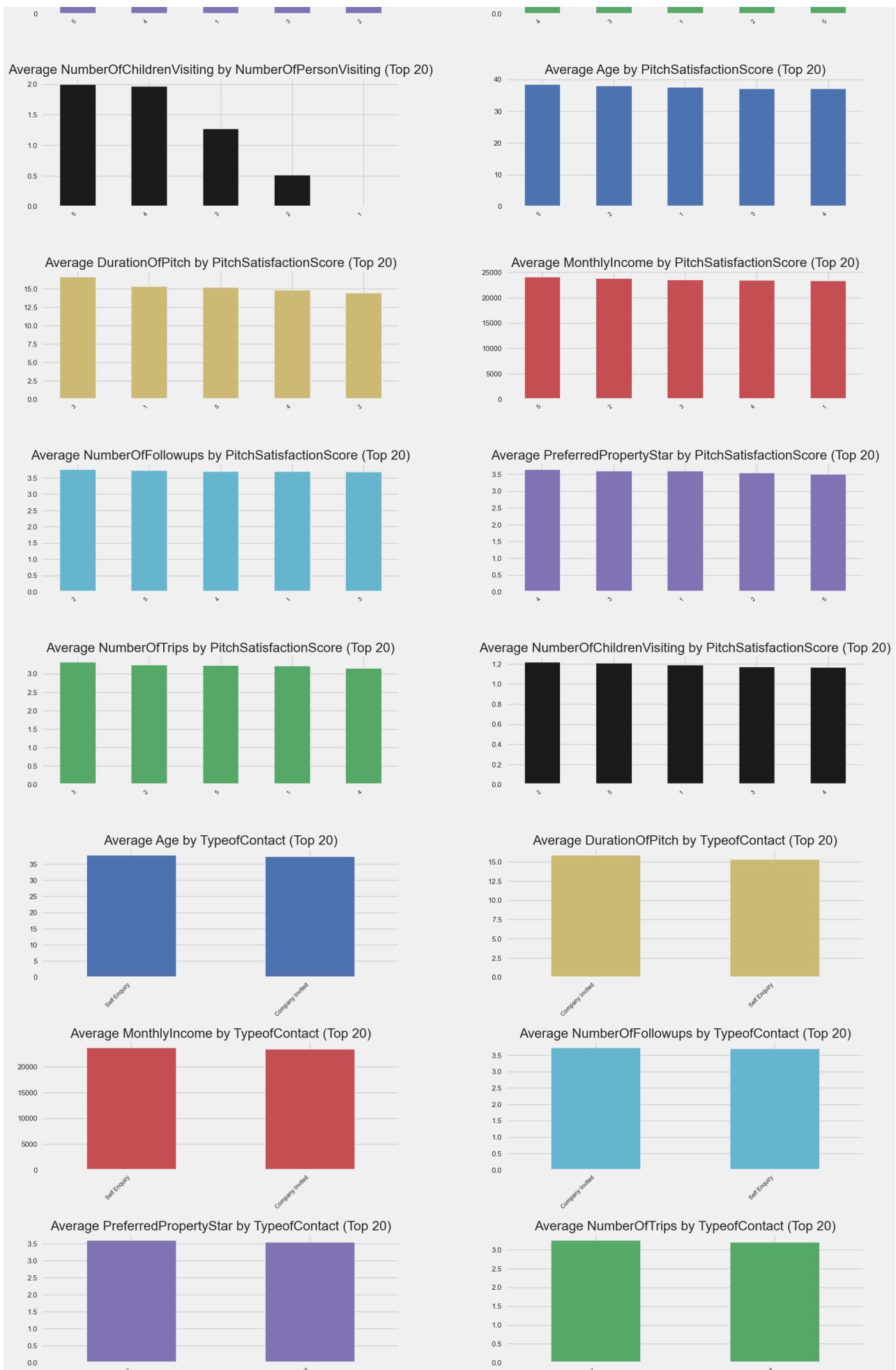
Average NumberOfTrips by Occupation (Top 20)

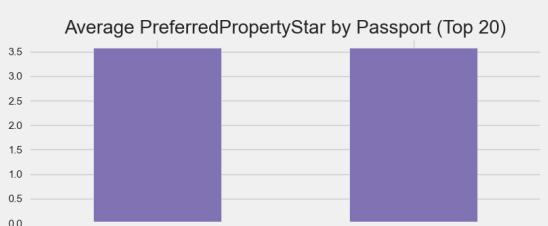
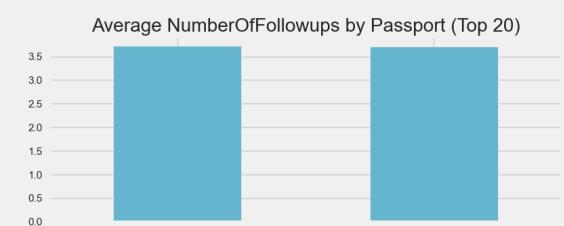
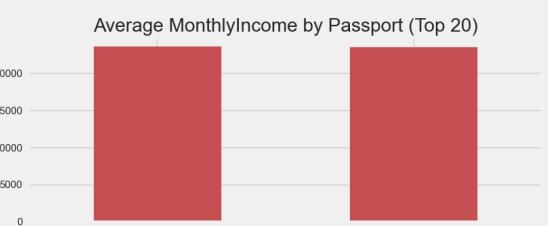
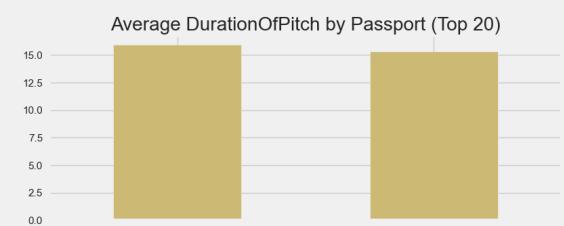
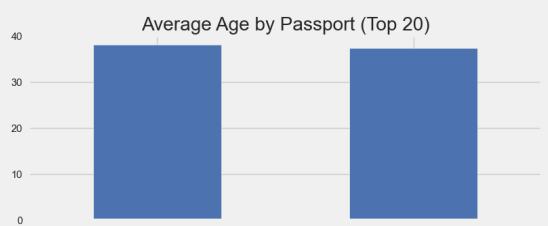
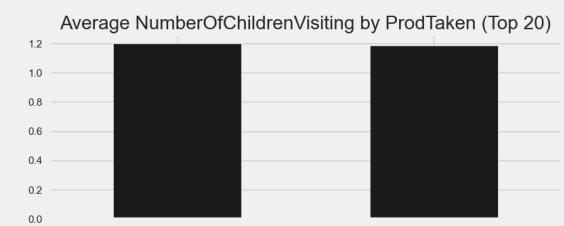
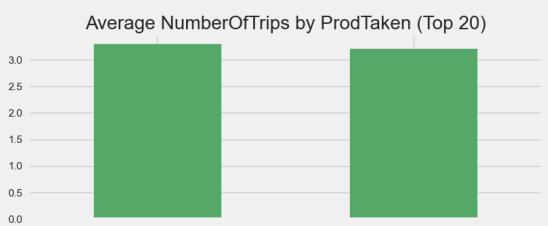
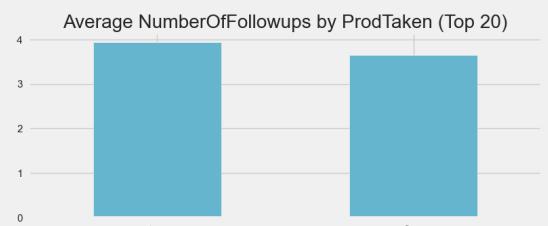
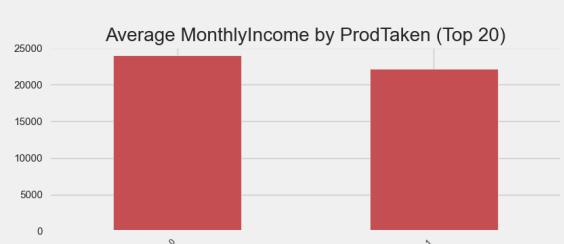
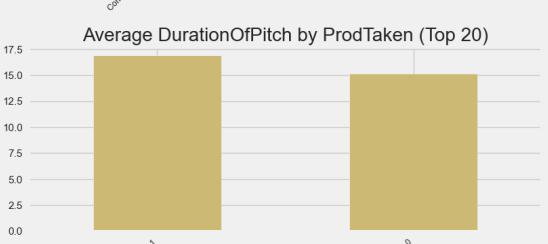
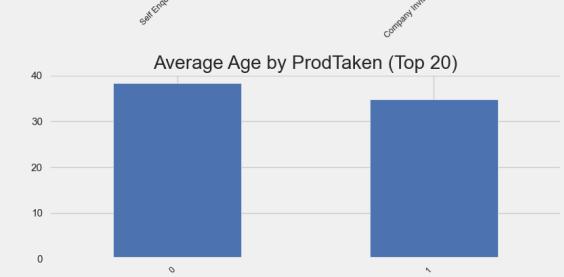
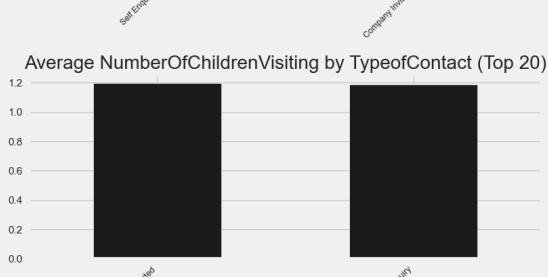


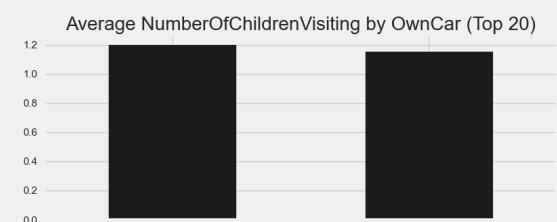
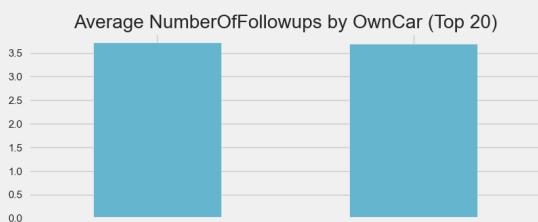
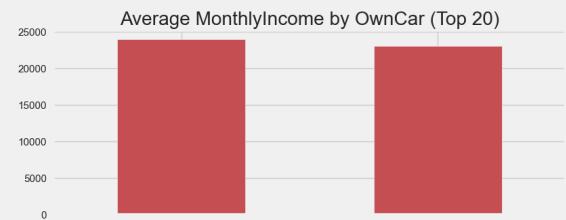
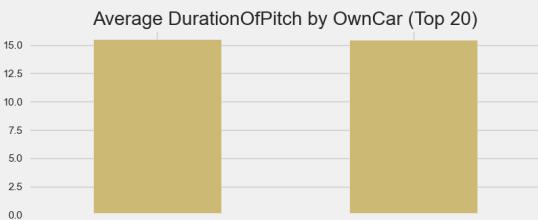
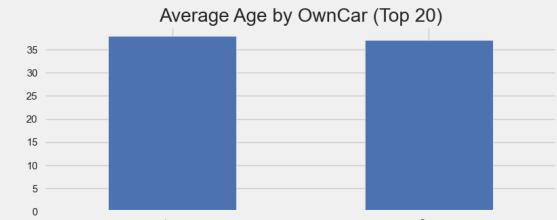












Thank a Lot for being with us!

In [ ]: