

# Project 2.1: Data Cleanup

By: Fairoza Amira Binti Hamzah

## Step 1: Business and Data Understanding

Key Decisions:

1. What decisions needs to be made?

We need to predict which city is the best to open the 14<sup>th</sup> store based on the previous sales data of each city.

2. What data is needed to inform those decisions?

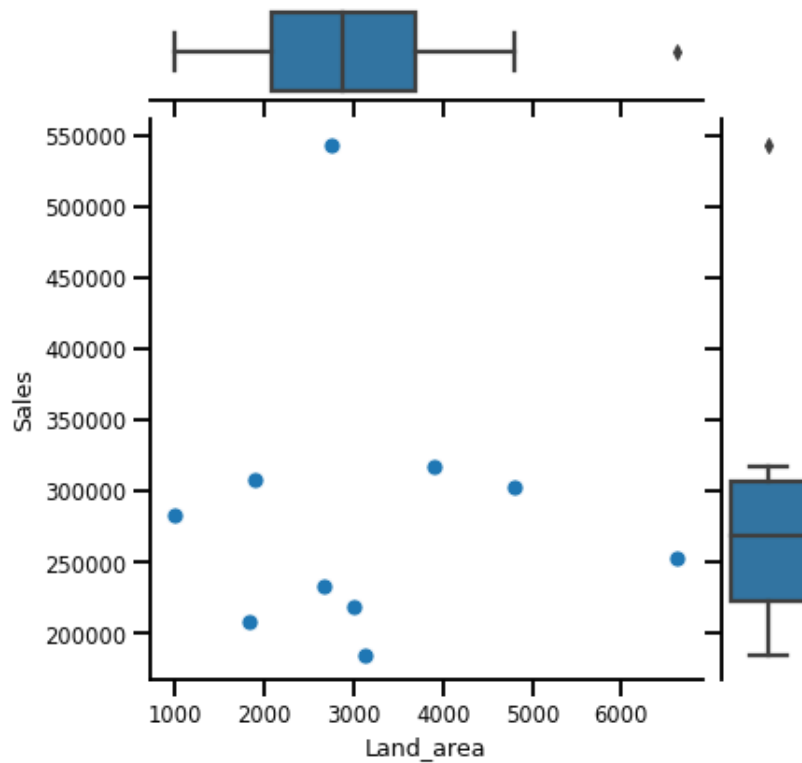
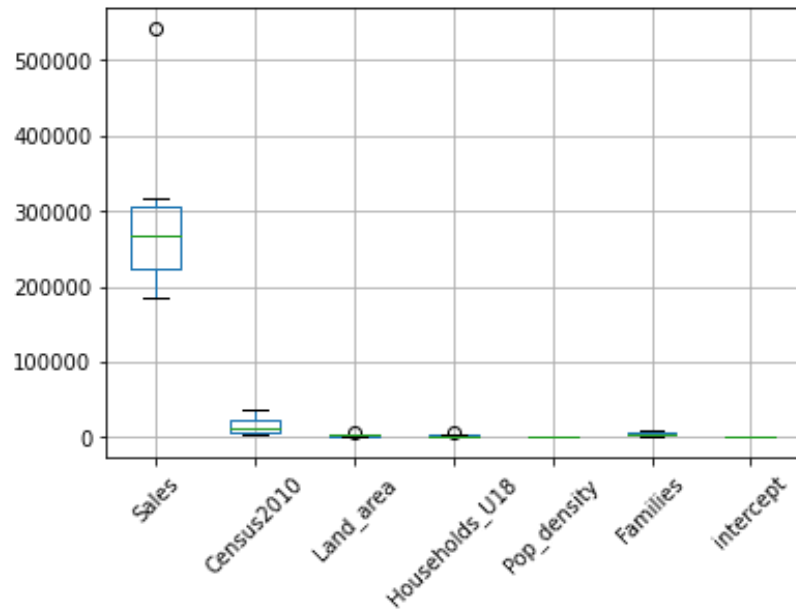
City  
2010 census population  
Pawdacity sales in other stores  
competitor sales  
household with under 18  
land area  
population density  
total families

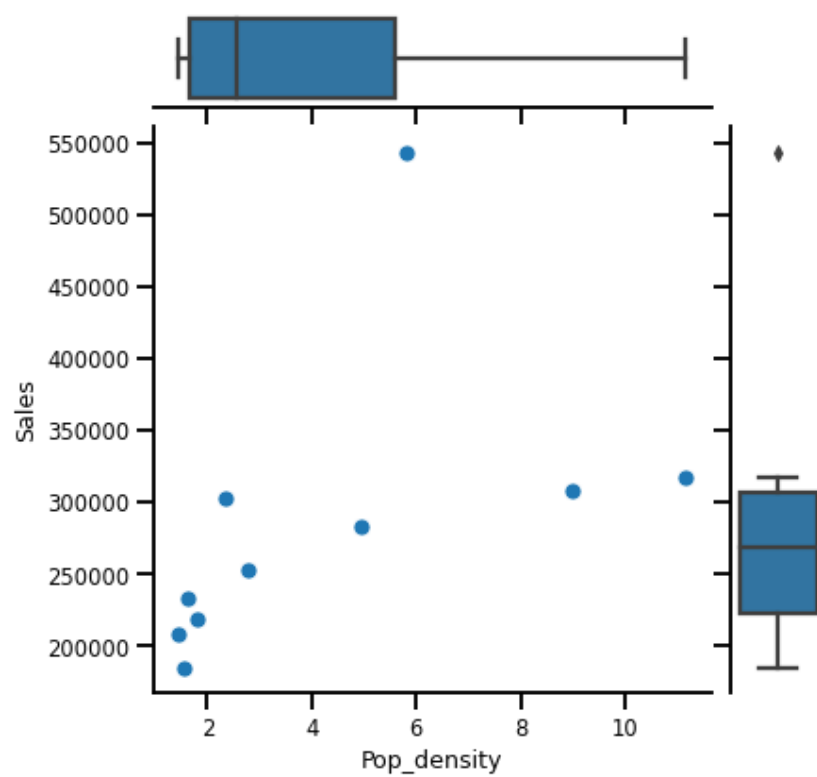
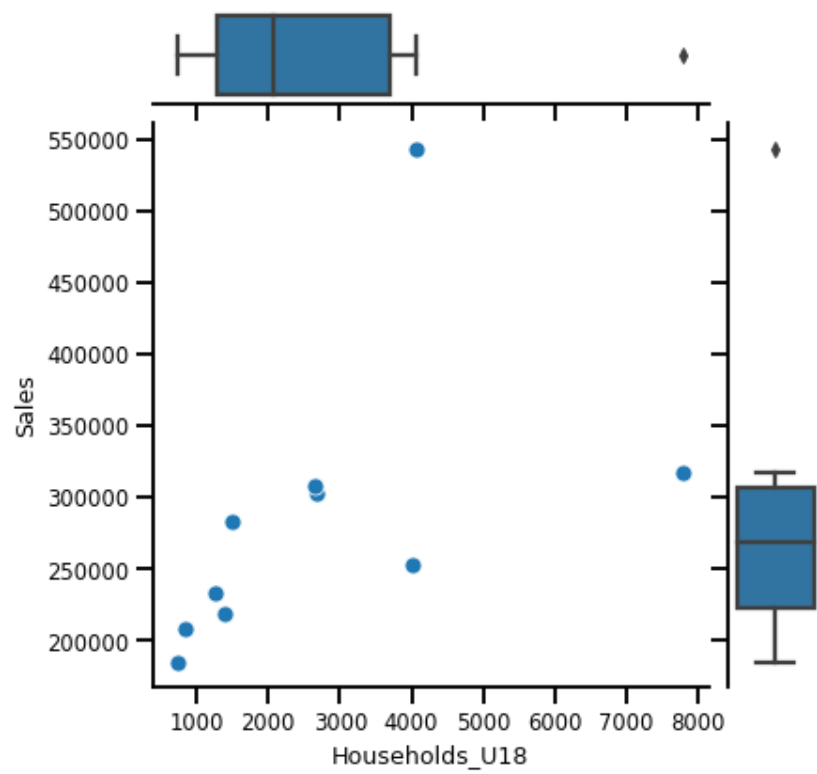
## Step 2: Building the Training Set

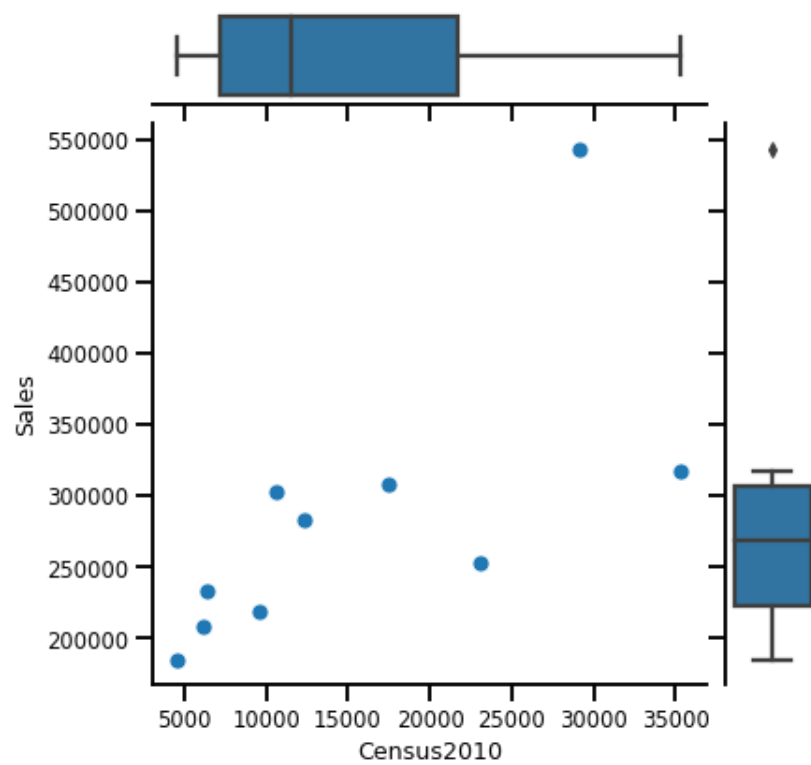
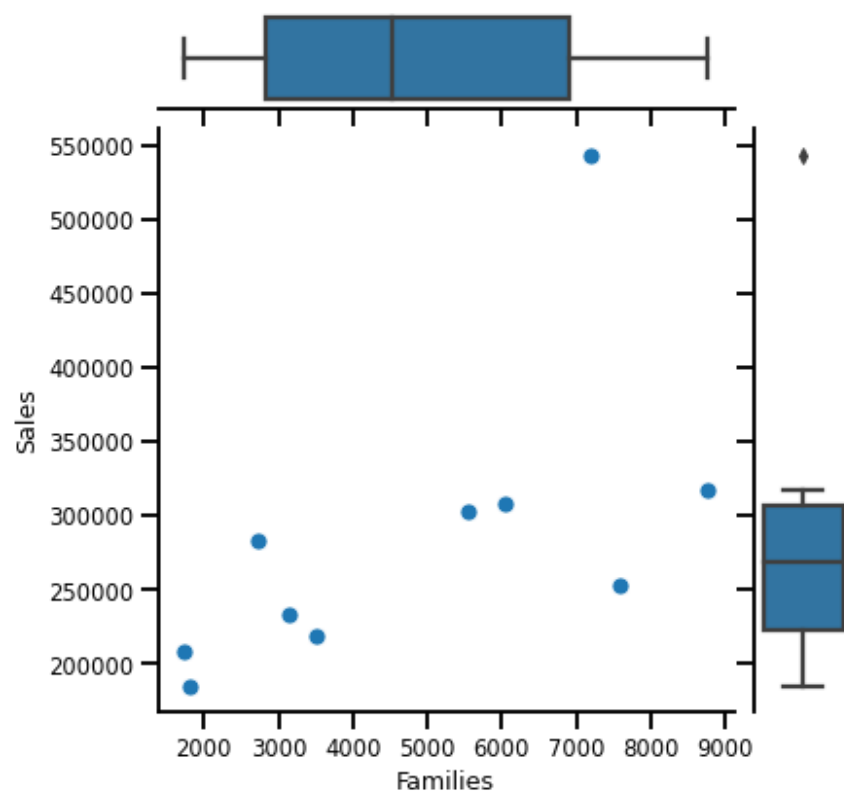
Column	Sum	Average
<i>Census Population</i>	213,862	19442
<i>Total Pawdacity Sales</i>	3,773,304	343027.64
<i>Households with Under 18</i>	34,064	3096.73
<i>Land Area</i>	33,071	3006.49
<i>Population Density</i>	63	5.71
<i>Total Families</i>	62,653	5695.71

## Step 3: Dealing with Outliers

Are there any cities that are outliers in the training set? Which outlier have you chosen to remove or impute? Because this dataset is a small data set (11 cities), **you should only remove or impute one outlier**. Please explain your reasoning.







There are 1 outlier in all the data plotted by using scatter plot and box plot, which is the Gillette city. Thus, Gillette would be the outlier in this case when compared against all other cities due to its greatest distance from the linear trend. Since the relationships between Gillette's population related variables and total sales are still correlated, Gillette should be kept for prediction and analysis.