# Project 2.1: Data Cleanup

## The Business Problem

Pawdacity is a leading pet store chain in Wyoming with 13 stores throughout the state. This year, Pawdacity would like to expand and open a 14th store. Your manager has asked you to perform an analysis to recommend the city for Pawdacity's newest store, based on predicted yearly sales.

## Step 1: Business and Data Understanding

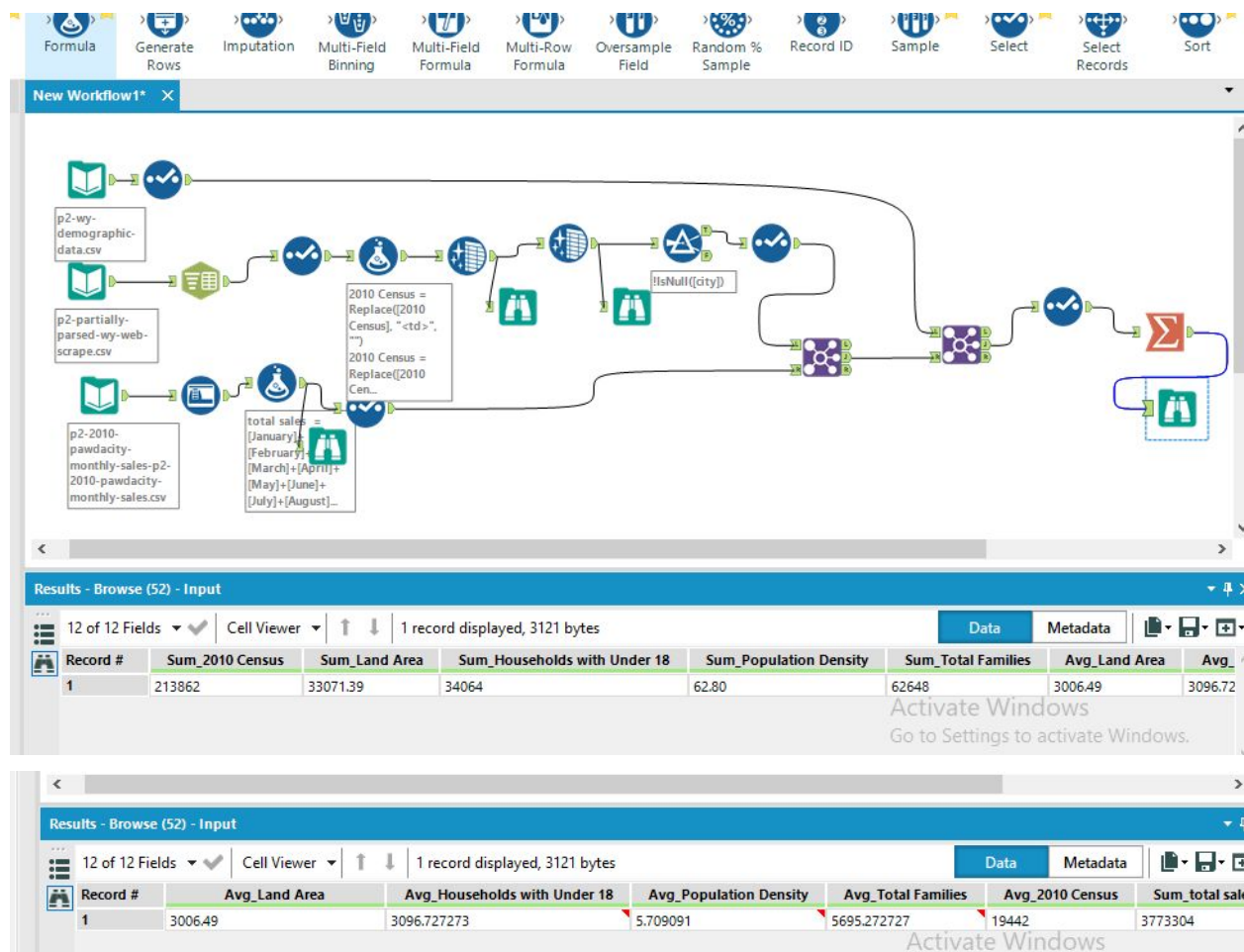### Key Decisions:

*Answer these questions*

1. What decisions needs to be made?

The decision needs to be made where to open the 14th Pawdacity pet store.

2. What data is needed to inform those decisions?

Some of the data required in order to inform this decision are city, population, Pawdacity sales in other stores, competitor sales, household with under 18, land area, total families.
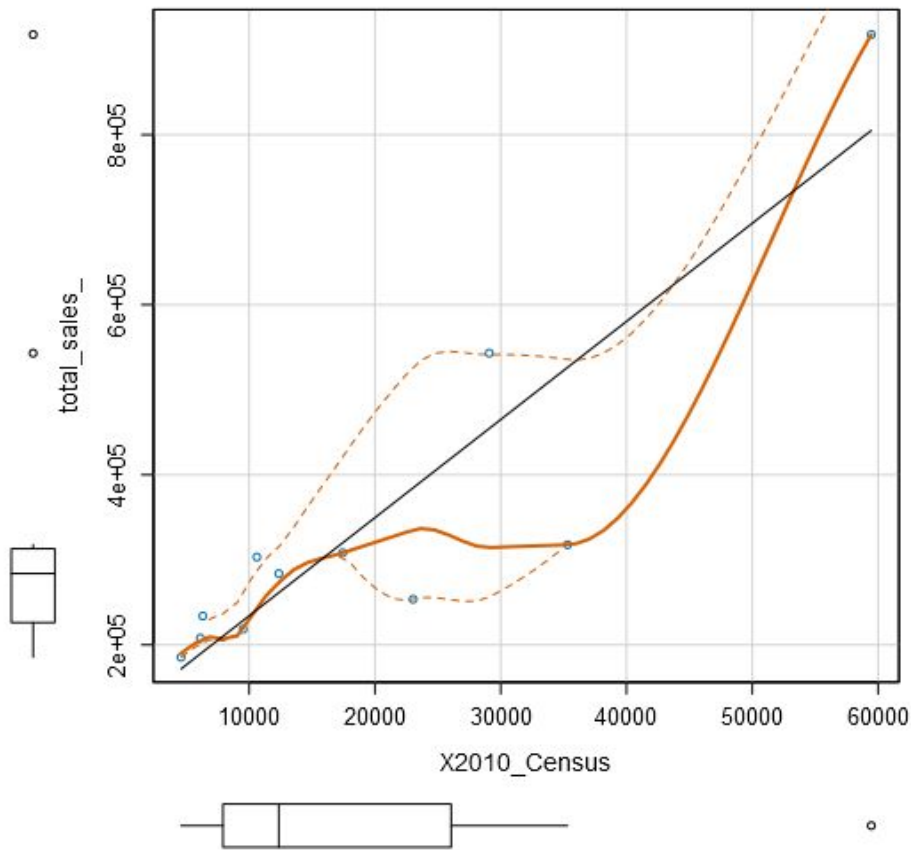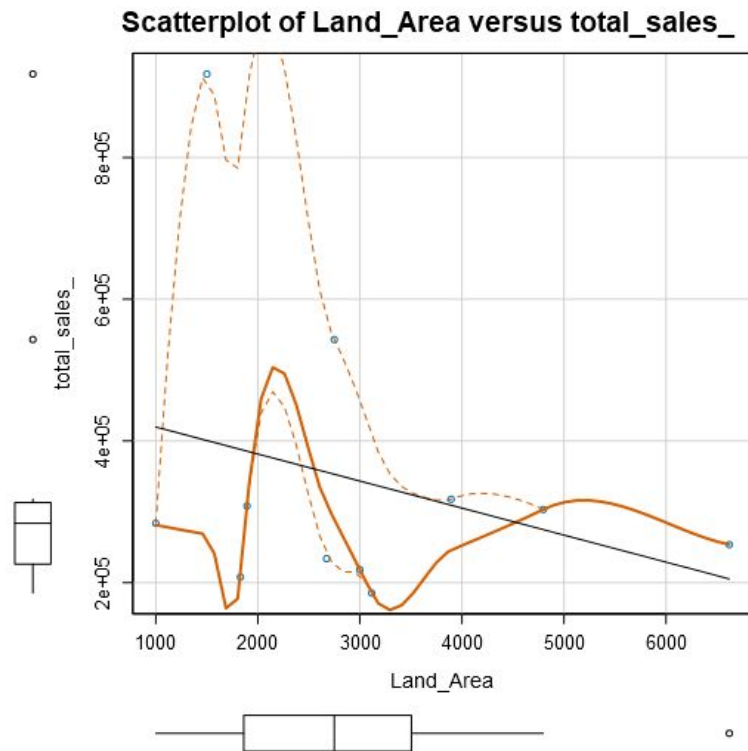
## Step 2: Building the Training Set

| Column | Sum | Average |
|---|---|---|
| *Census Population* | 213,862 | 19442 |
| *Total Pawdacity Sales* | 3,773,304 | 343027.63 |
| *Households with Under 18* | 34,064 | 3096.72 |
| *Land Area* | 33,071 | 3006.49 |
| *Population Density* | 63 | 5.70 |
| *Total Families* | 62,653 | 5695.27 |

# Step 3: Dealing with Outliers

**Scatterplot of X2010_Census versus total_sales_**

## Scatterplot of Land_Area versus total_sales_

total_sales_

8e+05

6e+05

4e+05

2e+05

Land_Area

1000   2000   3000   4000   5000   6000

## Scatterplot of Households_with_Under_18 versus total_sa

total_sales_

8e+05

6e+05

4e+05

2e+05

Households_with_Under_18

1000   2000   3000   4000   5000   6000   7000   8000

**Scatterplot of Population_Density versus total_sales_**



**Scatterplot of Total_Families versus total_sales_**



Based on the scatterplots above, the city of **Gillette** and **Cheyenne** seem to be the outliers as their sales data are higher than expected.

Since the relationships between Gillette's population related variables and total sales are still correlated, Gillette should be kept for analysis.

My suggestion would be to remove the city of Cheyenne and keep the city of Gillette for further analysis.