# Project 2.1: Data Cleanup

Make a copy of this document. Complete each section. When you are ready, save your file as a PDF document and submit it here:

## Step 1: Business and Data Understanding

*Provide an explanation of the key decisions that need to be made. (250 word limit)*

### Key Decisions:

*Answer these questions*

1. What decisions needs to be made?
   ➔ We need to predict the yearly sales of Pawdacity for different cities.
   ➔ In which city they should open their new store.

2. What data is needed to inform those decisions?
   ➔ First of all ,we need the monthly sales for all the cities so that we can calculate the yearly sales of each city.The data should be of the most recent year.
   ➔ We need population records of the cities.
   ➔ We need demographic data which encompass population numbers and some more specific data (e.g. land area, population density and so on) for different counties and cities in the state.
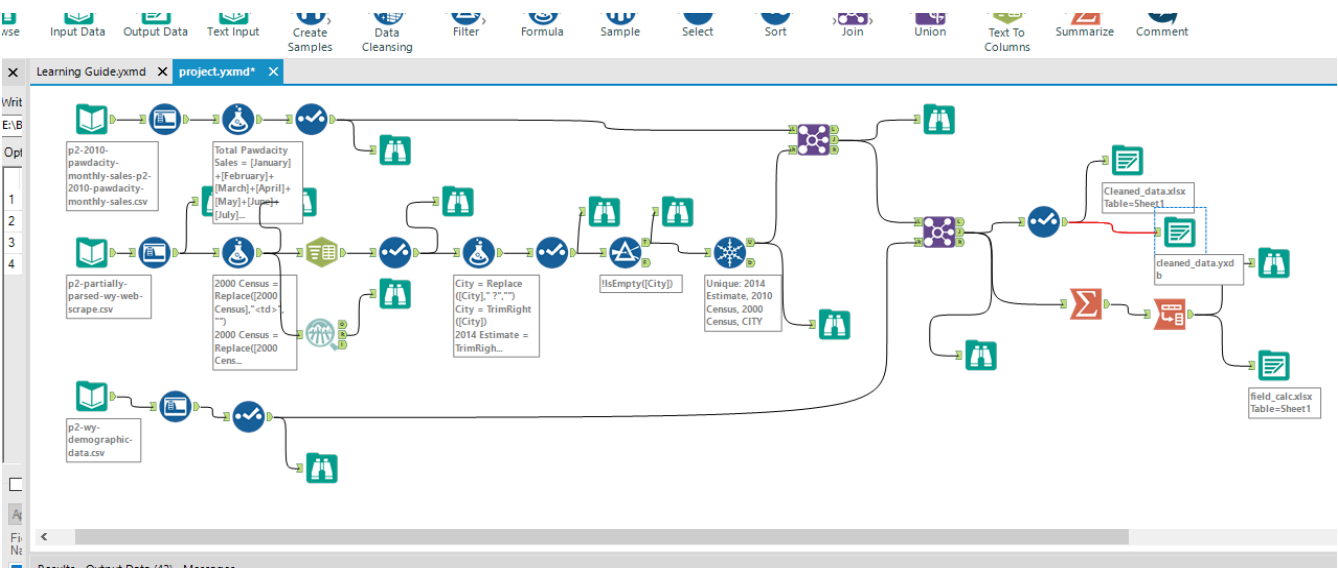
## Step 2: Building the Training Set

*Build your training set given the data provided to you. Your column sums of your dataset should match the sums in the table below.*

*In addition provide the averages on your data set here to help reviewers check your work. You should round up to two decimal places, ex: 1.24*

| Column | Sum | Average |
|---|---|---|
| *Census Population* | *213,862* | 19,442 |
| *Total Pawdacity Sales* | *3,773,304* | 343,027.63 |
| *Households with Under 18* | *34,064* | 3,096.73 |
| *Land Area* | *33,071* | 3,006.49 |
| *Population Density* | *63* | 5.70 |
| *Total Families* | *62,653* | 5,695.70 |

**Justification:-** I have drawn the following workflow for the data wrangling process.



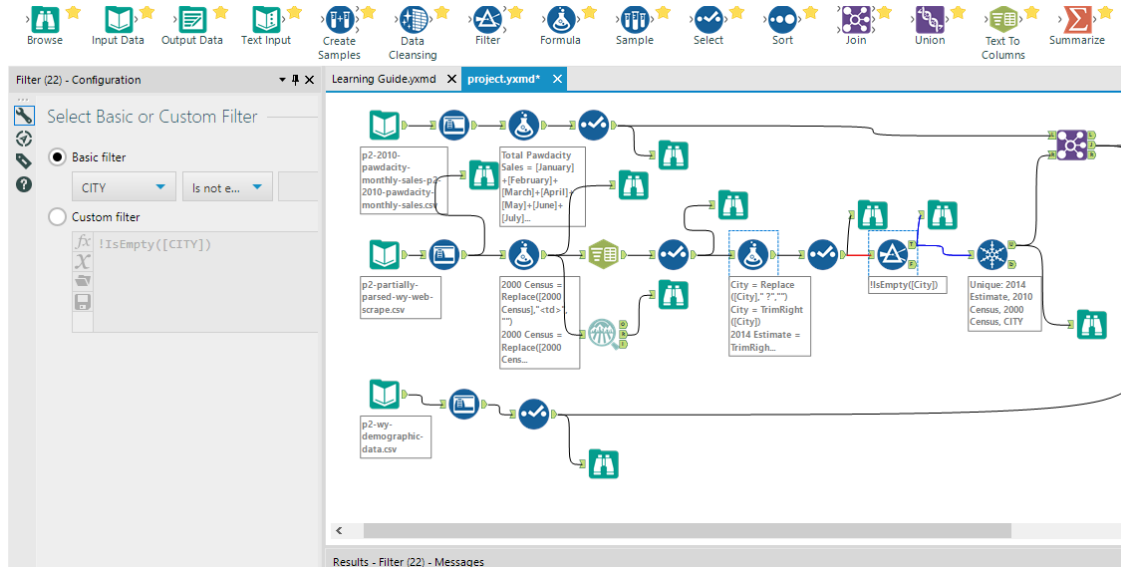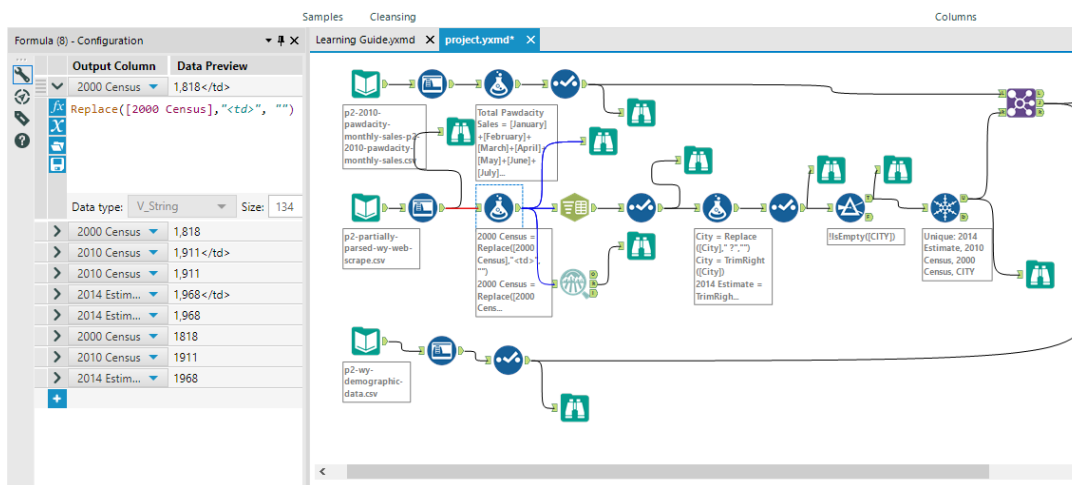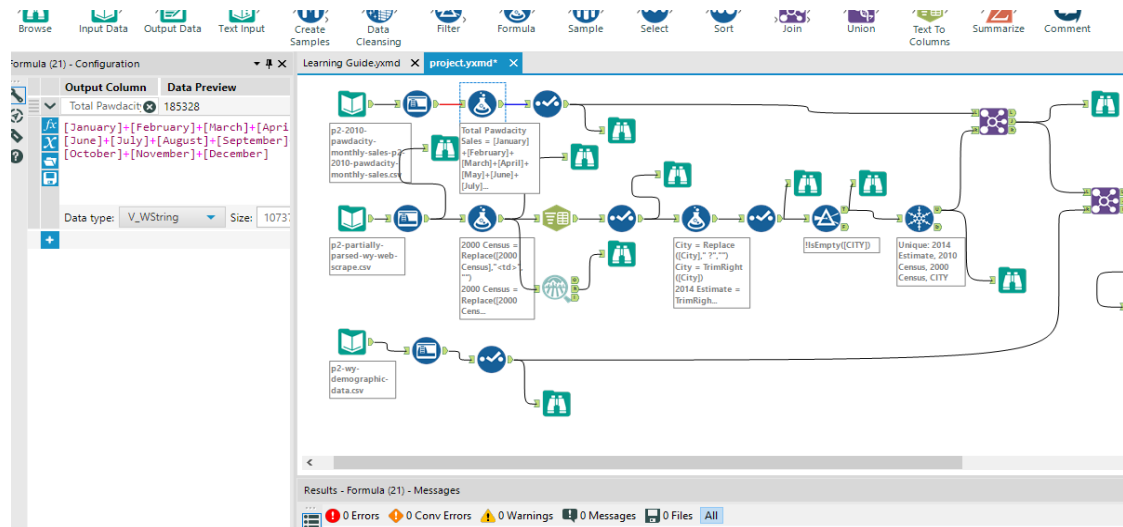**Data Cleaning :-** Field Summary showing % of null values in the given data set.

Record Report

1

## String/Character Fields

| Name | % Missing | Unique Values | Shortest Value | Longest Value | Min Value Count | Max Value Count | Remarks |
|------|-----------|---------------|----------------|---------------|-----------------|-----------------|---------|
| 2000 Census | 0.0% | 96 | - | `<td class="navbox-list navbox-odd hlist" style="text-align:left;border-left-width:2px;border-left-style:solid;width:100%;padding:0px">` | 1 | 3 | Some values of this field have a small number of value counts. If Appropriate, consider combining some value levels together. |
| 2010 Census | 0.0% | 102 | 4 | `<td class="navbox-list navbox-even hlist" style="text-align:left;border-left-width:2px;border-left-style:solid;width:100%;padding:0px;background:transparent;">` | 1 | 2 | Some values of this field have a small number of value counts. If Appropriate, consider combining some value levels together. |
| 2014 Estimate | 0.0% | 98 | 4 | `<td class="navbox-list navbox-even hlist" style="text-align:left;border-left-width:2px;border-left-style:solid;width:100%;padding:0px">` | 1 | 3 | Some values of this field have a small number of value counts. If Appropriate, consider combining some value levels together. |
| City|County | 3.9% | 100 | Cody ?|Park | East Thermopolis|Hot Springs | 1 | 4 | Some values of this field have a small number of value counts. If Appropriate, consider combining some value levels together. |

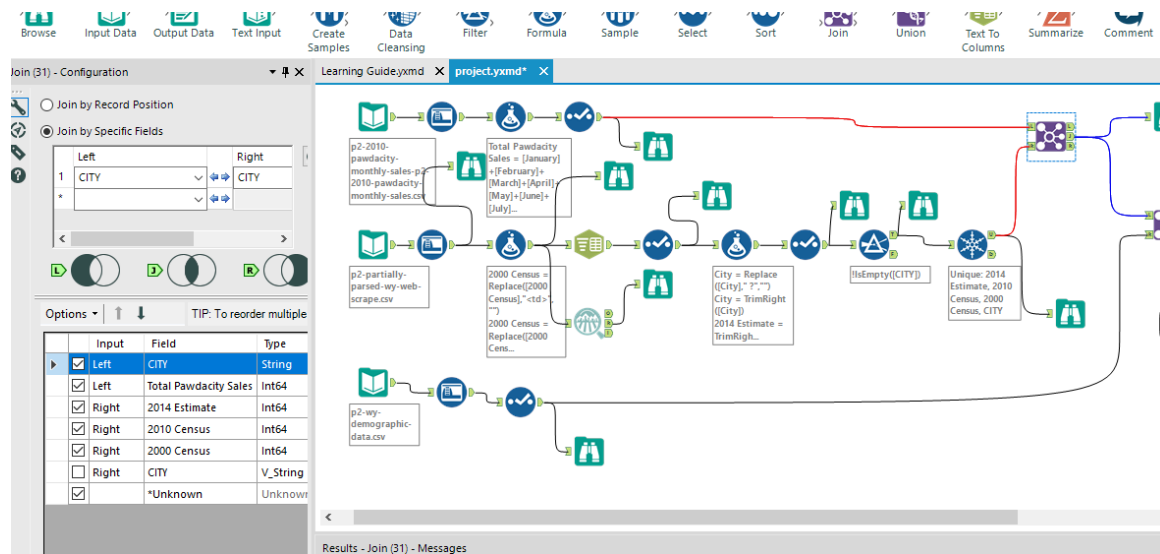Used the not null function to drop the rows with null values.



**Data Formatting :-** Used the string functions to format the data in the ill-formatted rows.

**Data Blending :-** Used the join function to blend the data from three files into single dataset.

**Output :-** Formatted and cleaned data was obtained in an excel file and the characteristics need to be calculated were outputted to another excel file.



| | A | B | C | D | E | F | G | H |
|---|---|---|---|---|---|---|---|---|
| 1 | City | Total Pawdacity Sales | 2010 Census Population | Land Area | Households with Under 18 | Population Density | Total Families | |
| 2 | Buffalo | 185328 | 4585 | 3115.5075 | 746 | 1.55 | 1819.5 | |
| 3 | Casper | 317736 | 35316 | 3894.3091 | 7788 | 11.16 | 8756.32 | |
| 4 | Cheyenne | 917892 | 59466 | 1500.1784 | 7158 | 20.34 | 14612.64 | |
| 5 | Cody | 218376 | 9520 | 2998.95696 | 1403 | 1.82 | 3515.62 | |
| 6 | Douglas | 208008 | 6120 | 1829.4651 | 832 | 1.46 | 1744.08 | |
| 7 | Evanston | 283824 | 12359 | 999.4971 | 1486 | 4.95 | 2712.64 | |
| 8 | Gillette | 543132 | 29087 | 2748.8529 | 4052 | 5.8 | 7189.43 | |
| 9 | Powell | 233928 | 6314 | 2673.57455 | 1251 | 1.62 | 3134.18 | |
| 10 | Riverton | 303264 | 10615 | 4796.859815 | 2680 | 2.34 | 5556.49 | |
| 11 | Rock Sprir | 253584 | 23036 | 6620.201916 | 4022 | 2.78 | 7572.18 | |
| 12 | Sheridan | 308232 | 17444 | 1893.977048 | 2646 | 8.98 | 6039.71 | |
| 13 | | | | | | | | |
| 14 | | | | | | | | |
| 15 | | | | | | | | |
| 16 | | | | | | | | |
| 17 | | | | | | | | |
| 18 | | | | | | | | |

A1          fx     Name

| | A | B | C | D |
|---|---|---|---|---|
| 1 | Name | Value | | |
| 2 | Sum_Total Pawdacity Sales | 3773304 | | |
| 3 | Avg_Total Pawdacity Sales | 343027.6364 | | |
| 4 | Sum_2010 Census | 213862 | | |
| 5 | Avg_2010 Census | 19442 | | |
| 6 | Sum_Land Area | 33071.38039 | | |
| 7 | Avg_Land Area | 3006.489126 | | |
| 8 | Sum_Households with Under 18 | 34064 | | |
| 9 | Avg_Households with Under 18 | 3096.727273 | | |
| 10 | Sum_Population Density | 62.8 | | |
| 11 | Avg_Population Density | 5.709090909 | | |
| 12 | Sum_Total Families | 62652.79 | | |
| 13 | Avg_Total Families | 5695.708182 | | |
| 14 | | | | |
| 15 | | | | |
| 16 | | | | |
| 17 | | | | |
| 18 | | | | |
| 19 | | | | |
| 20 | | | | |

And the answers to question no.2 were obtained.

# Step 3: Dealing with Outliers

*Answer these questions*

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.

## Calculating the interquartile ranges:-

M7          fx     =IF(AND($E7>=$E$19,$E7<=$E$18),"","Y")

| | A | B | C | D | E | F | G | H |
|---|---|---|---|---|---|---|---|---|
| 1 | City | Total Pawdacity Sales | 2010 Census Population | Land Area | Households with Under 18 | Population Density | Total Families | |
| 2 | Buffalo | 185328 | 4585 | 3115.5075 | 746 | 1.55 | 1819.5 | |
| 3 | Casper | 317736 | 35316 | 3894.3091 | 7788 | 11.16 | 8756.32 | |
| 4 | Cheyenne | 917892 | 59466 | 1500.1784 | 7158 | 20.34 | 14612.64 | |
| 5 | Cody | 218376 | 9520 | 2998.95696 | 1403 | 1.82 | 3515.62 | |
| 6 | Douglas | 208008 | 6120 | 1829.4651 | 832 | 1.46 | 1744.08 | |
| 7 | Evanston | 283824 | 12359 | 999.4971 | 1486 | 4.95 | 2712.64 | |
| 8 | Gillette | 543132 | 29087 | 2748.8529 | 4052 | 5.8 | 7189.43 | |
| 9 | Powell | 233928 | 6314 | 2673.57455 | 1251 | 1.62 | 3134.18 | |
| 10 | Riverton | 303264 | 10615 | 4796.859815 | 2680 | 2.34 | 5556.49 | |
| 11 | Rock Springs | 253584 | 23036 | 6620.201916 | 4022 | 2.78 | 7572.18 | |
| 12 | Sheridan | 308232 | 17444 | 1893.977048 | 2646 | 8.98 | 6039.71 | |
| 13 | | | | | | | | |
| 14 | q1 | 226152 | 7917 | 1861.721074 | 1327 | 1.72 | 2923.41 | |
| 15 | q3 | 312984 | 26061.5 | 3504.9083 | 4037 | 7.39 | 7380.805 | |
| 16 | iqr | 86832 | 18144.5 | 1643.187226 | 2710 | 5.67 | 4457.395 | |
| 17 | 1.5*iqr | 130248 | 27216.75 | 2464.780839 | 4065 | 8.505 | 6686.0925 | |
| 18 | upper | 443232 | 53278.25 | 5969.689139 | 8102 | 15.895 | 14066.8975 | |
| 19 | lower | 95904 | -19299.75 | -603.059765 | -2738 | -6.785 | -3762.6825 | |
| 20 | | | | | | | | |
| 21 | | | | | | | | |

## Outliers found :-

| City | Total Pawdacity Sales | 2010 Census Population | Land Area | Households with Under 18 | Population Density | Total Families |
|---|---|---|---|---|---|---|
| Buffalo | | | | | | |
| Casper | | | | | | |
| Cheyenne | Y | Y | | | Y | Y |
| Cody | | | | | | |
| Douglas | | | | | | |
| Evanston | | | | | | |
| Gillette | Y | | | | | |
| Powell | | | | | | |
| Riverton | | | | | | |
| Rock Springs | | | Y | | | |
| Sheridan | | | | | | |

➔ There are 3 cities which have outliers :-
1) Cheyenne
2) Gillette
3) Rock Springs

➔ Cheyenne has 4 outliers and it seems to be a big city ,so I will not remove this city .

➔ Rock –springs have land area as an outlier ,but its all other parameters are in the interquartile range. So it might be accurate date ,because of having large land area,and it won't be removed.

➔ Gillette has all its parameters in interquartile range but its total pawdacity sales is out of interquartile range,which I don't think is possible as compared to the other cities. So,this city should be removed.