# Power BI – Extract, Transform and Load

## ETL 1 – Examine Data with Power Query Editor
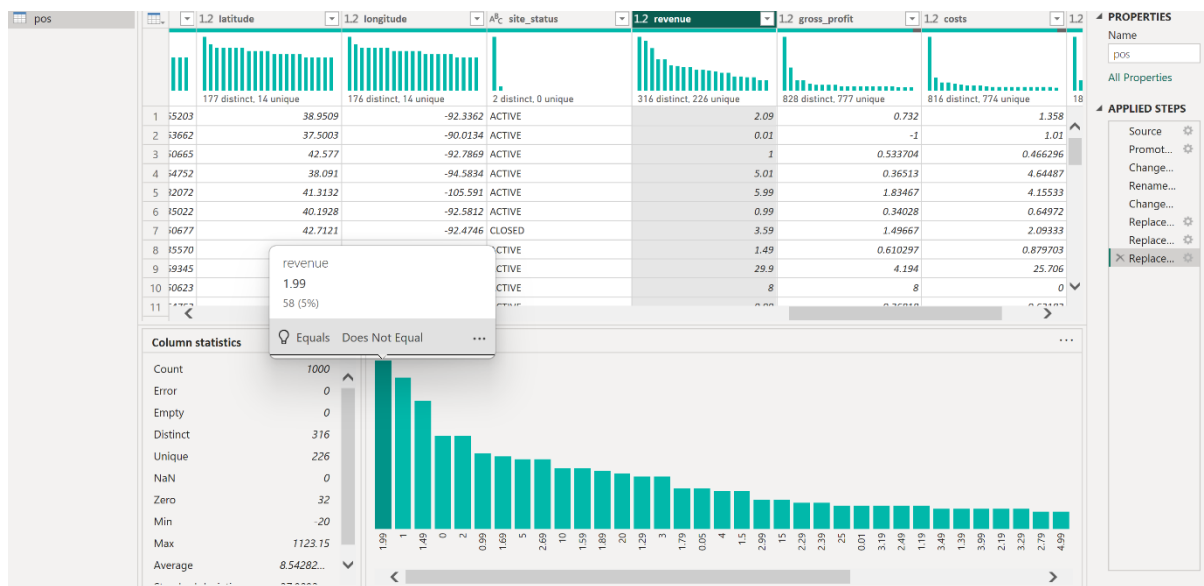
Load the data -> pos for (Point of Sale). The data is from a fictitious convenience store/gas station.
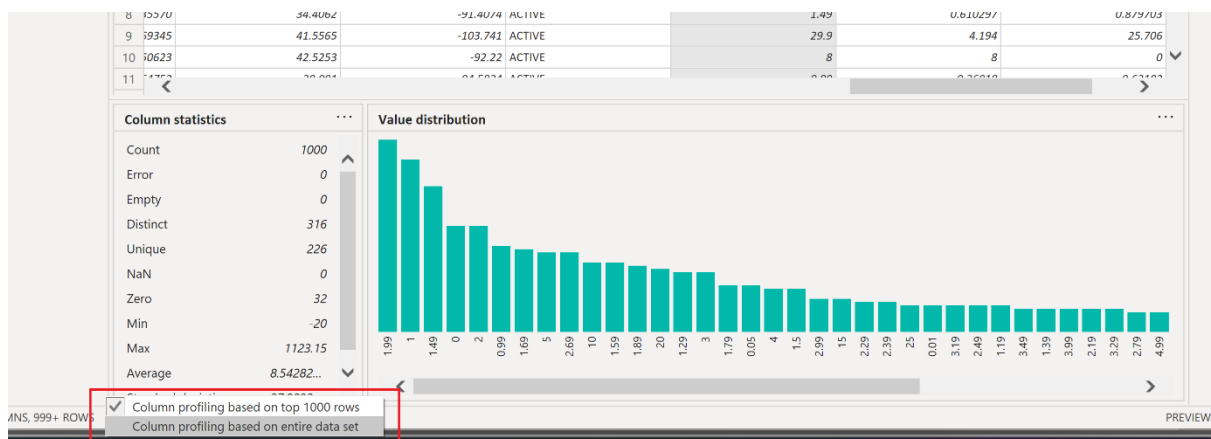


View a column distribution by clicking on a particular column, clicking on 'View' and selecting 'Column distribution'. Transformation steps are shown on right-hand side under 'Applied Steps'.

Clicking on 'View', then checking 'Column profile' gives the following view



We can base the column profiling on the entire dataset



Now tick the 'Column quality' checkbox. This gives the percentage of Valid, Error and Empty entries

Change number columns to correct type using Data Type drop-down.



This can also be changed using the type symbol on the column header



Replace errors with null

## ETL 2 – Dates and Calculated Columns

Add new column from selection



Add new column for Day of Week

Add new Custom Column for gross profit margin



Use custom formula to calculate value for column

New column added



## ETL 3 – Checking for and Eliminating Outliers

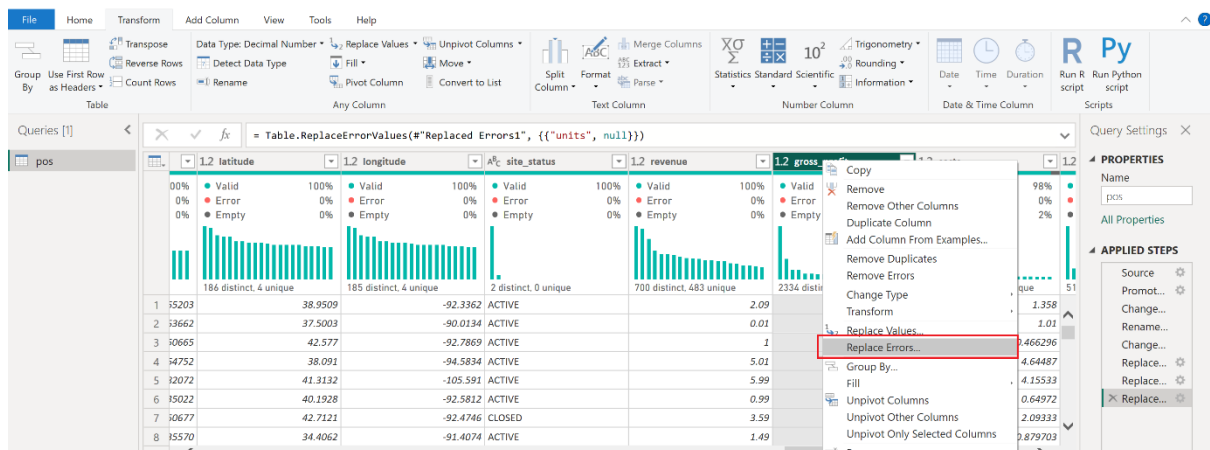The highlighted rows are 3 standard deviations more than the mean (>87), therefore outliers.

These outliers are related to fuel



We can remove these outliers by manually deselecting them from the column



Another way to do this is to go to 'Reduce Rows' in the Home tool and select 'Remove Top Rows' and specify the number of rows to remove. Note this will only work if you've sorted the data correctly.

The column parent_name has high cardinality, meaning that there are 60 distinct types.



We can create a new column to group some of these items

Now we have just 5 different items



We can create a new hierarchy for our new column

## ETL 4 : Data Models and Joins

Load up dataset states.csv that contains a list of states and join with the pos.csv dataset.

Note that the states dataset does not have a zip code column. It is instead called postal_code



However, some of the values contain a 9-digit postal_code



Go to the Model View and create a connection between the two datasets

We can see that there is a One-To-Many relationship between the states and pos table.

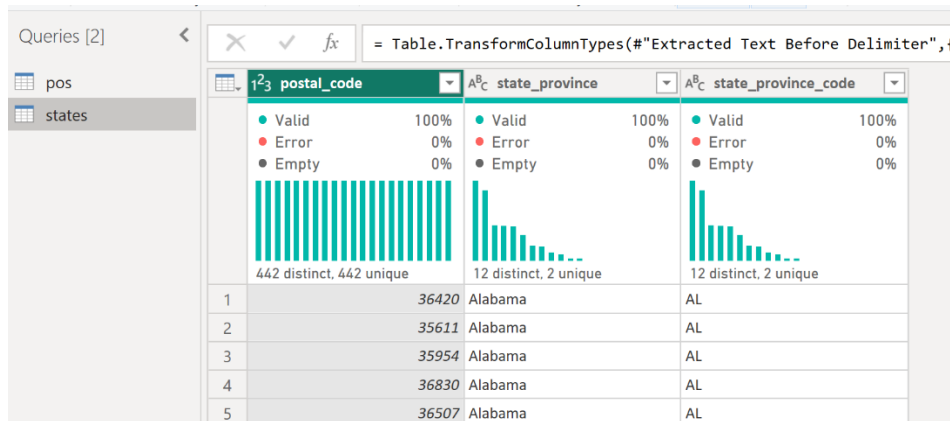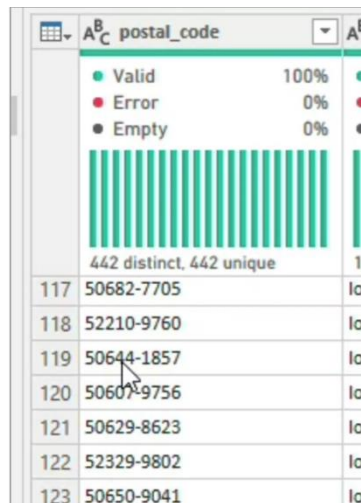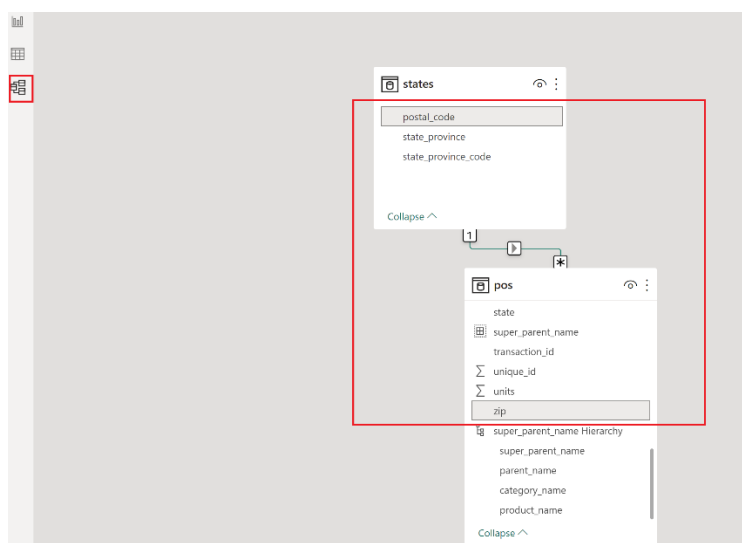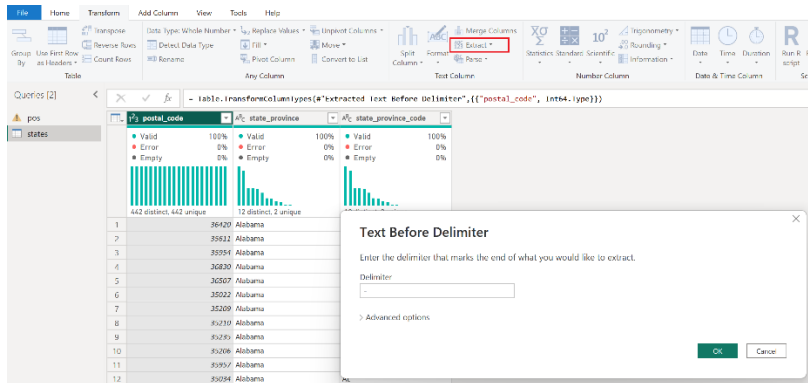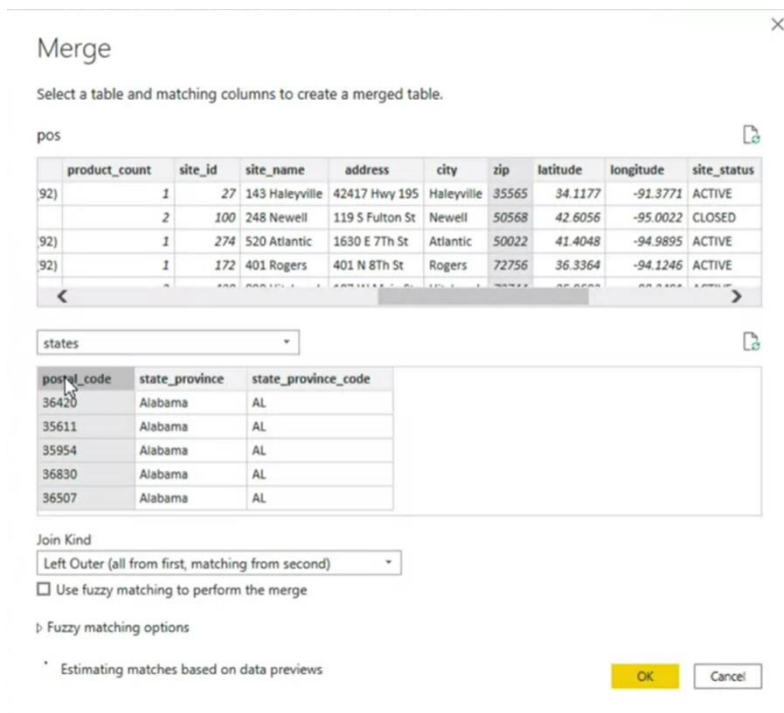Each transaction in pos matches up with only one observation in the states table. But one entry in the states table could match up with many entries in the pos table.

To remove the last 4 digits of the post codes with 9 digits. Highlight the column, click on 'Transform', then 'Extract', 'Text Before Delimiter' and specify – (hyphen)
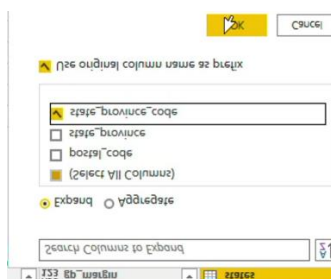


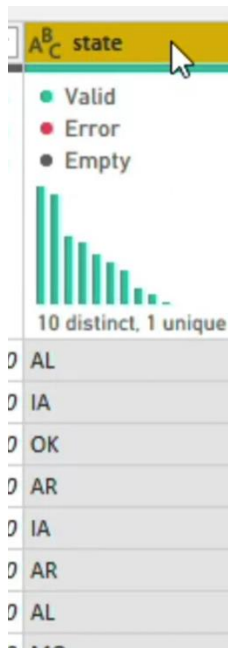Use Merge Queries to link the zip from the pos dataset to the postal_code of the states dataset.

Note that the columns must be the same type.



We only want to keep the state_province_code column

The column name is a bit long, so rename it to state



Use star schema, one main dataset, and any other dataset is joined to that key dataset to avoid querying lots of hierarchical relationships.