# Milestone 3 Processing of Data

WQD7005 – DATA MINING

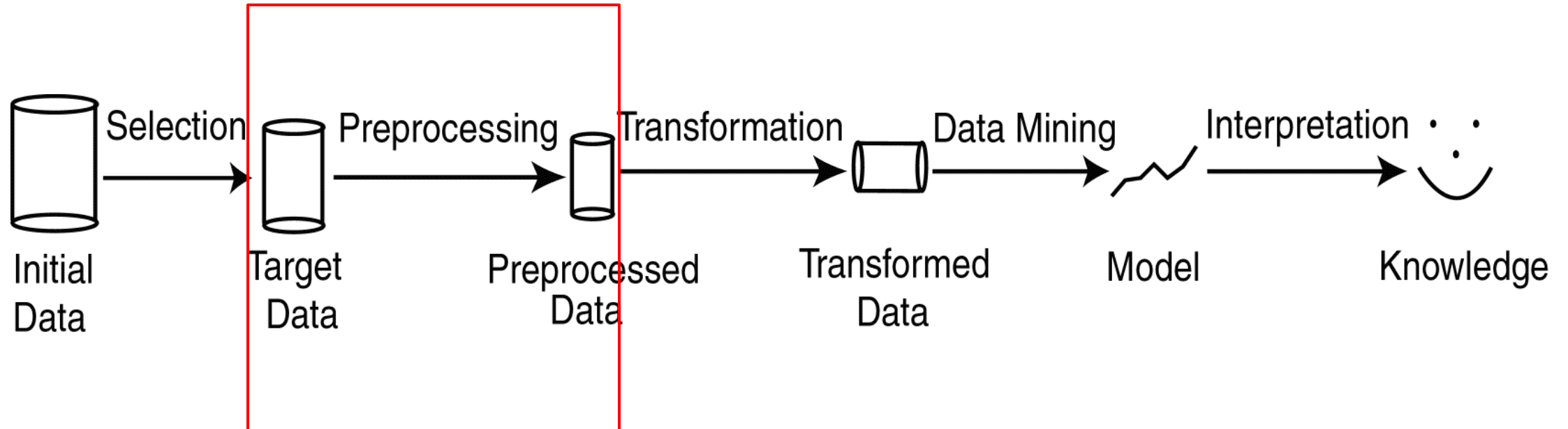ZULKANAIN BIN HASAN

WQD180031

NUR ATHIRAH BT. RAZAK

WQD170072

Knowledge Discovery in Databases (KDD): process of finding useful information and patterns in data:
1. Selection ( Pre-Mining 1): Obtain data from various sources
2. Preprocessing (Pre-Mining 2):  Cleanse data

**Milestone 3**



Selection → Preprocessing → Transformation → Data Mining → Interpretation

Initial Data → Target Data → Preprocessed Data → Transformed Data → Model → Knowledge

# Analysis Goal:
Findings a good stock counter for short term and long-term investment. Using historical data from KLSE Stock Market to predict future trend.
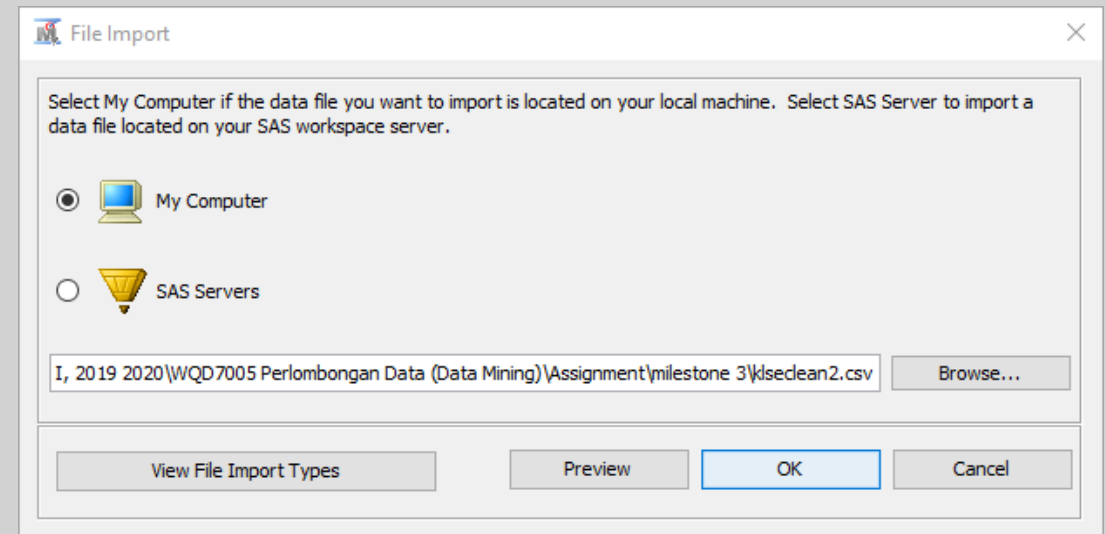
---

Analysis Data:

- Extract data from historical KLSE Stock Market
- Sample target based on adj-close data trend
- Actual sample target approximately 10%

# Create SAS data file

Step 1: Import .csv file from local storage to project diagram by using File Import.

- Data have been pre-cleaning first before import to project diagram.

- Only target data have been used in this step.

Step 2: Create linkages between File Import to Save Data in order to convert .csv file to SAS data file.
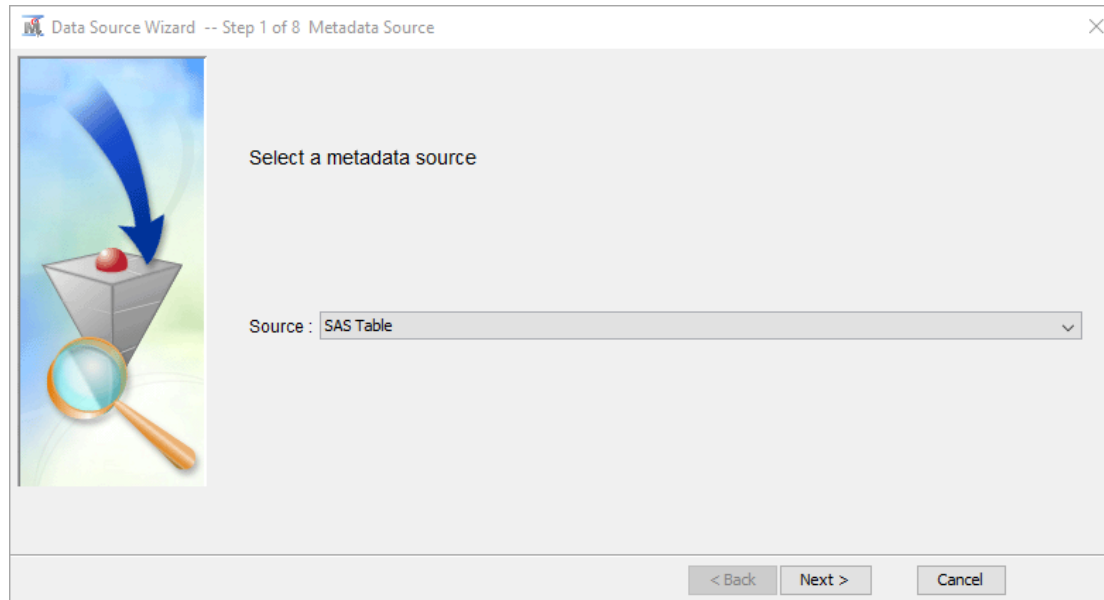- SAS Library Name have been renamed to Stock at Save Data property menu.

Step 3: Select a metadata source that have been create by previous step.
- For this case, data stored in SAS Library – Stock as per previous setup.

- Screening Column Metadata by inspecting all variables for their Role and Level.
- Target have been identified for this data and can be proceed with setting the Metadata Advisor Options (Basic or Advance).

In Advance Advisor Options, 2 property have been changed:

i. Class Levels Count Threshold from 2 to 20

ii. Reject Levels Count Threshold from 20 to 100

1

**Data Source Wizard -- Step 6 of 8  Create Sample**

Do you wish to create a sample data set?
◉ No    ○ Yes

**Table Info**
Columns    11
Rows        16338

**Sample Size**

Type      Percent

Percent   20

Rows

< Back    Next >    Cancel

2

**Data Source Wizard -- Step 7 of 8  Data Source Attributes**

You may change the name and the role, and can specify a population segment identifier for source to be created.

Name :    EM_SAVE_TRAIN
Role :     Raw
Segment :
Notes :

< Back   Next >   Cancel

Final

**Data Source Wizard -- Step 8 of 8  Summary**

Metadata Completed.

**Library:**       STOCK
**Data Source:**   EM_SAVE_TRAIN
**Role:**          Raw

| Role | Level | Count |
| --- | --- | --- |
| Input | Interval | 7 |
| Rejected | Nominal | 2 |
| Target | Interval | 1 |
| Time ID | Interval | 1 |

< Back   Finish   Cancel

Final step for this part shows:
i.    7 role as Input
ii.   2 role as Rejected for Nominal Level
iii.  1 role as Target

# Exploring source data

In this part, 2 property have been modified:

1. Sample Method to random
2. Fetch Size to max

In this explore menu it shows total 11 columns and as variables available in this source data. No. of rows available is 16338.
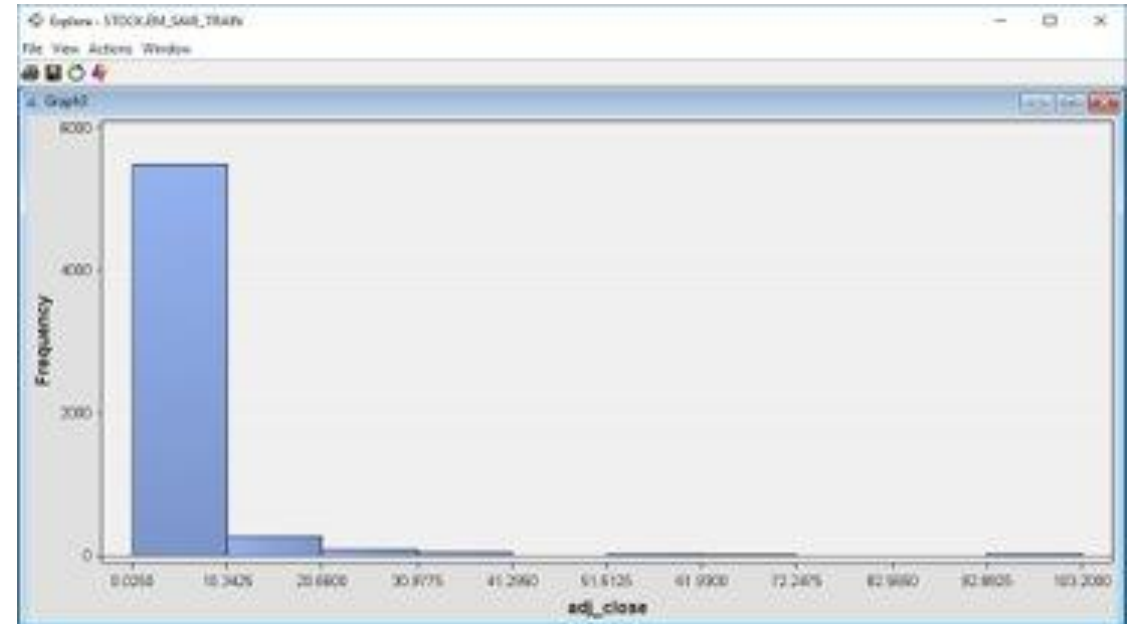
Creating a Histogram for a Single Variable – adj_close
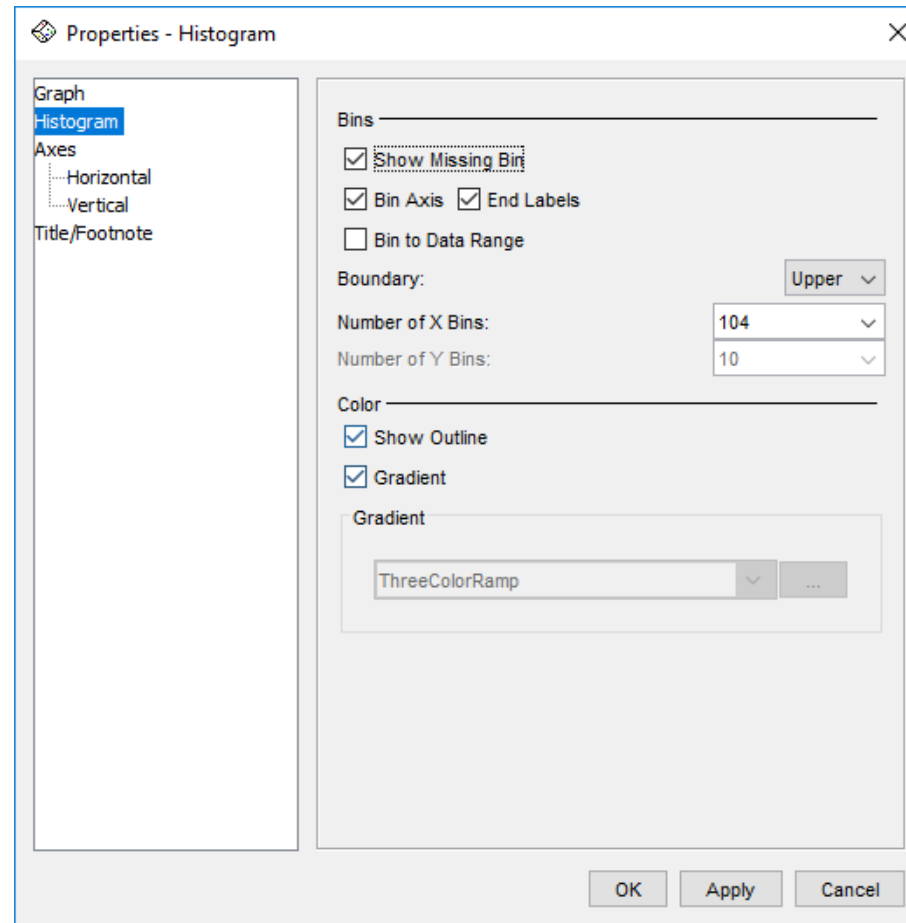
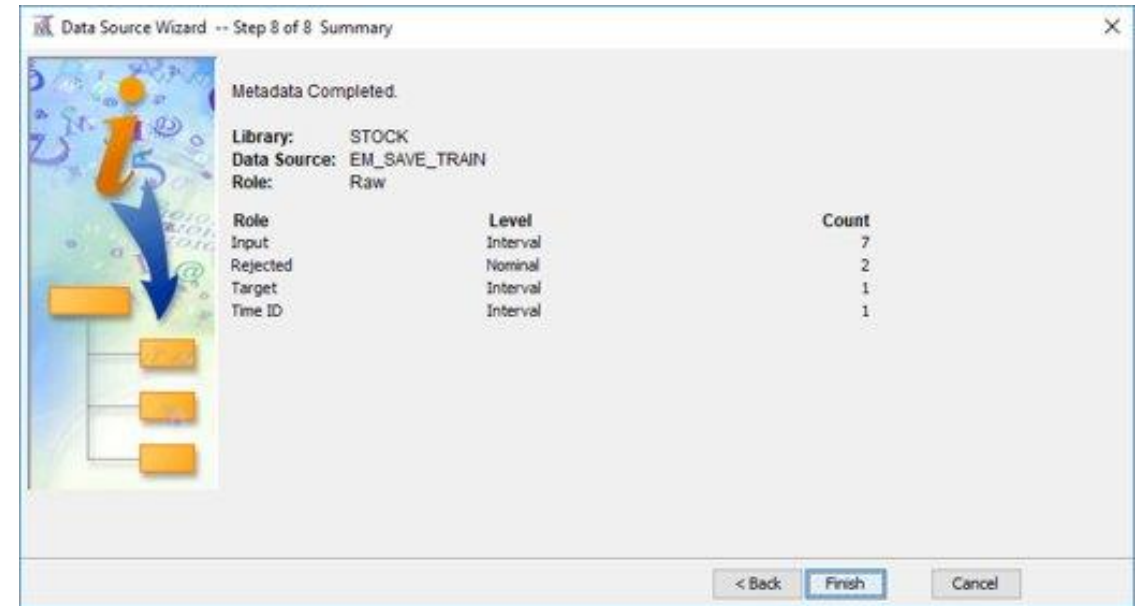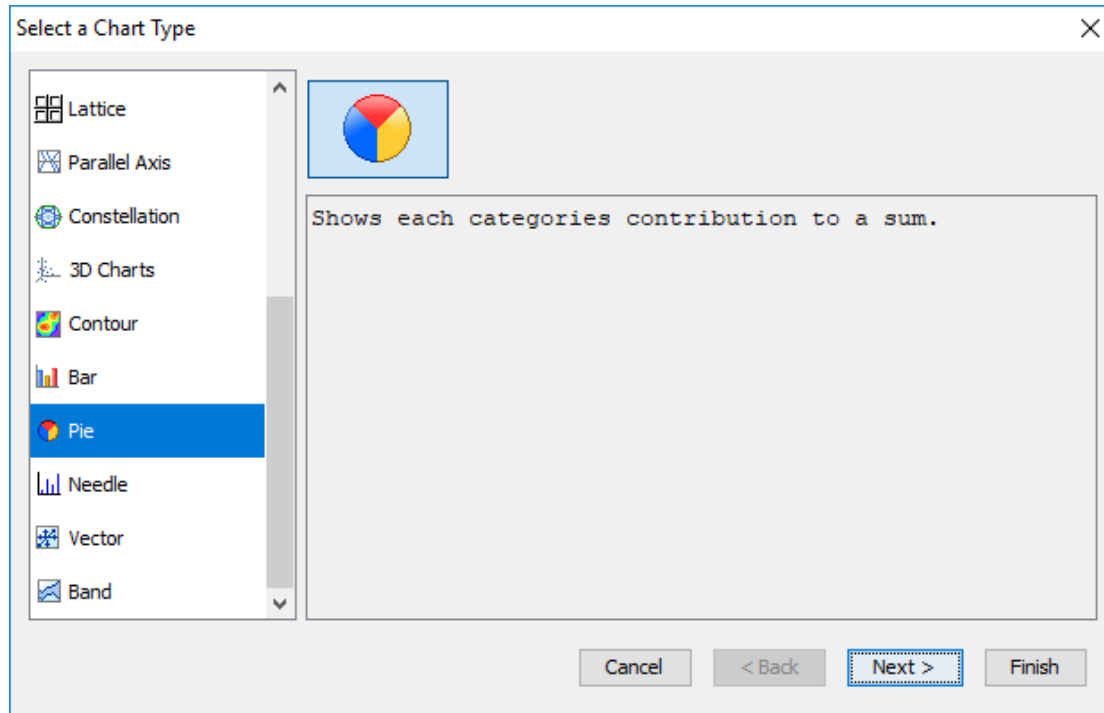- Set adj_close as Role X

Histogram shows:
- The minimum value is 0.025
- The maximum value is 103.2



Left histogram is a default view while right histogram is an individual value after setting to histogram property
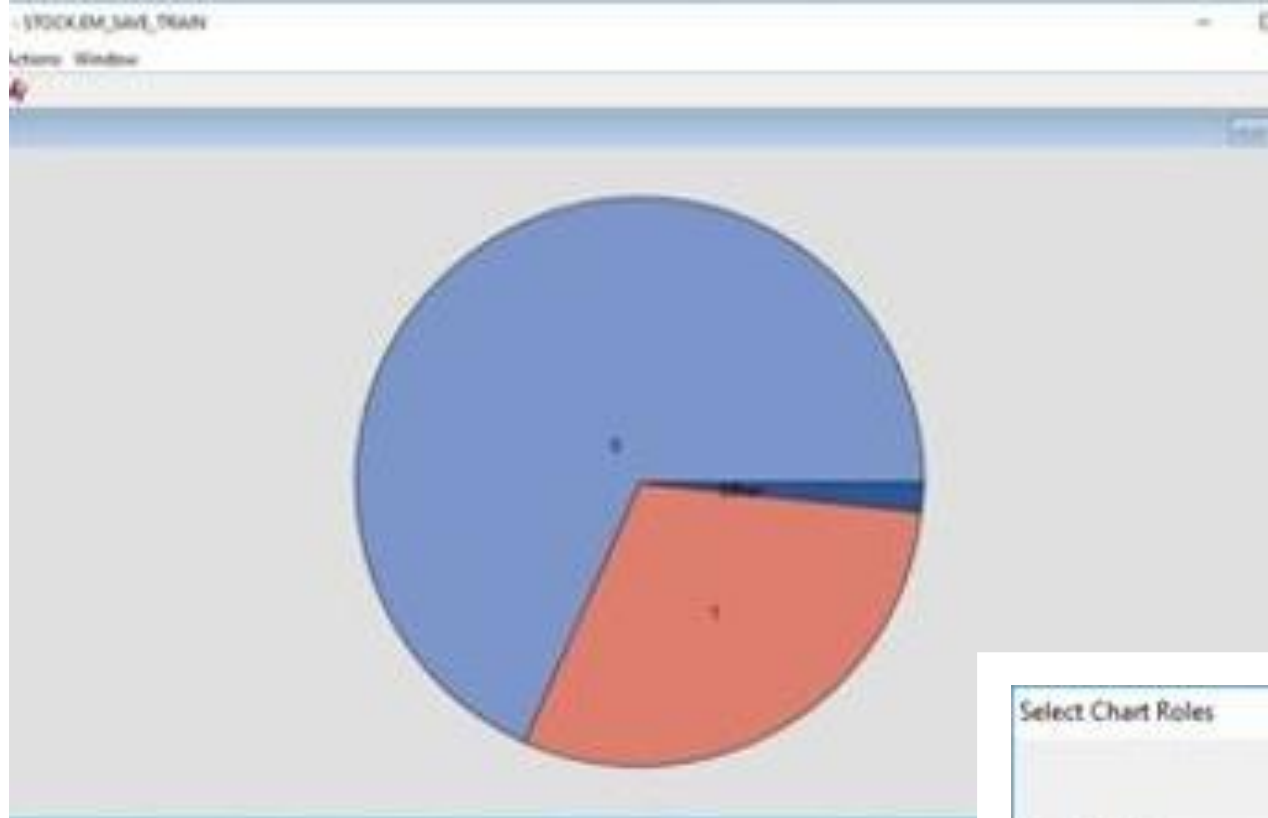
# Adding a "Missing" Bin to a Histogram

Adding Plots to the Explore Window
- Create Pie chart

This step require to set Target Role as Category

As a results, Pie chart shows (60:40):
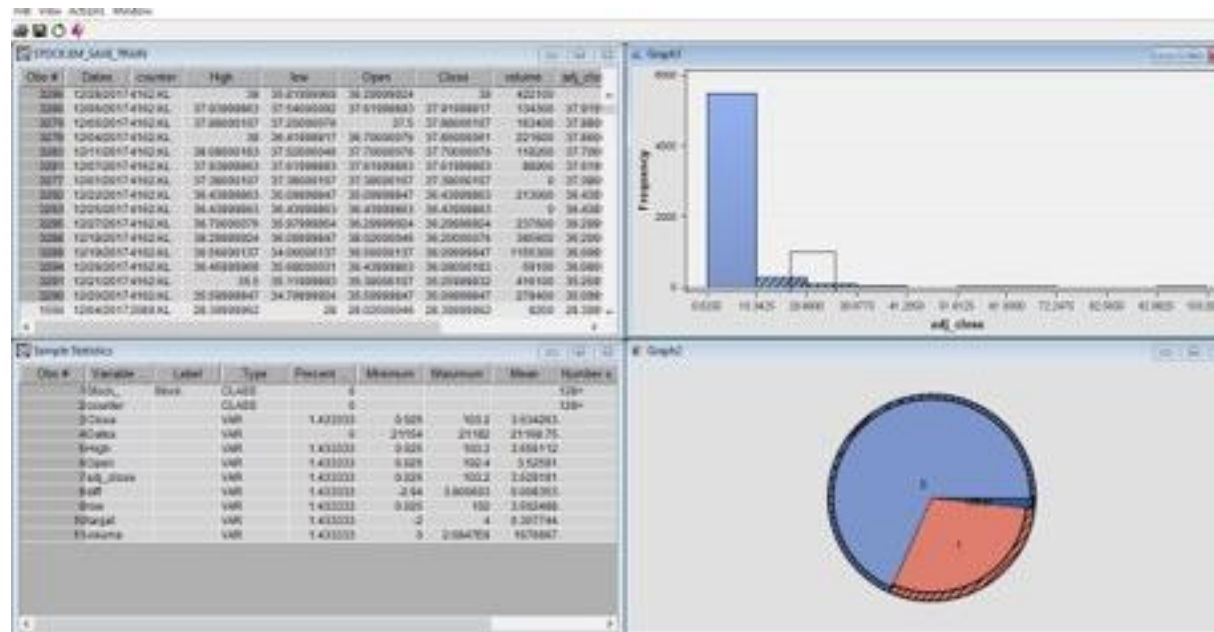Target 0 = 4098 frequency
Target 1 = 1799 frequency

**Select Chart Roles**

| ▲ Variable | Role | Type | Description | Format |
|---|---|---|---|---|
| adj_close | | Numeric | adj_close | BEST12. |
| Close | | Numeric | Close | BEST12. |
| counter | | Character | counter | $7. |
| Dates | | Numeric | Dates | MMDDYY10. |
| diff | | Numeric | diff | BEST12. |
| High | | Numeric | High | BEST12. |
| low | | Numeric | low | BEST12. |
| Open | | Numeric | Open | BEST12. |
| Stock_ | | Character | Stock | |
| target | Category | Numeric | target | BEST12. |
| volume | | Numeric | volume | BEST12. |

Use default assignments

☐ Allow multiple role assignments

Exploring Variable Associations

- This shows all view about this data and how the variable associates. As an example, when point selected to histogram, pie chart and table shows same proportion to its related point.