Clustering Course
Lab 1: Analysis of financial market data

Contents

Analysis of financial market data
Summary

Hello and welcome to first lab!

Let me give you an overview of what we are going to practice in this lab.

First, we are going to download stock market data for S&P 500 Index.
We will then see summary statistics to have an overall idea about our data set.

We will examine return distribution of the index and we will finish the lab with discussing about volatility and have few words about crisis periods.
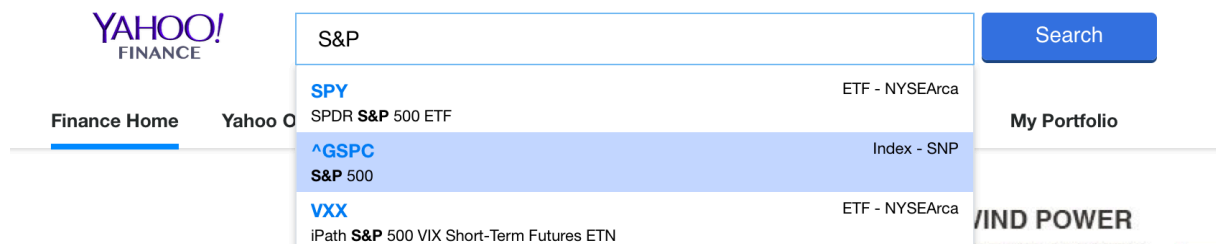
By the end of this lab, you will learn how to use SPSS Modeler for exploratory data analysis.
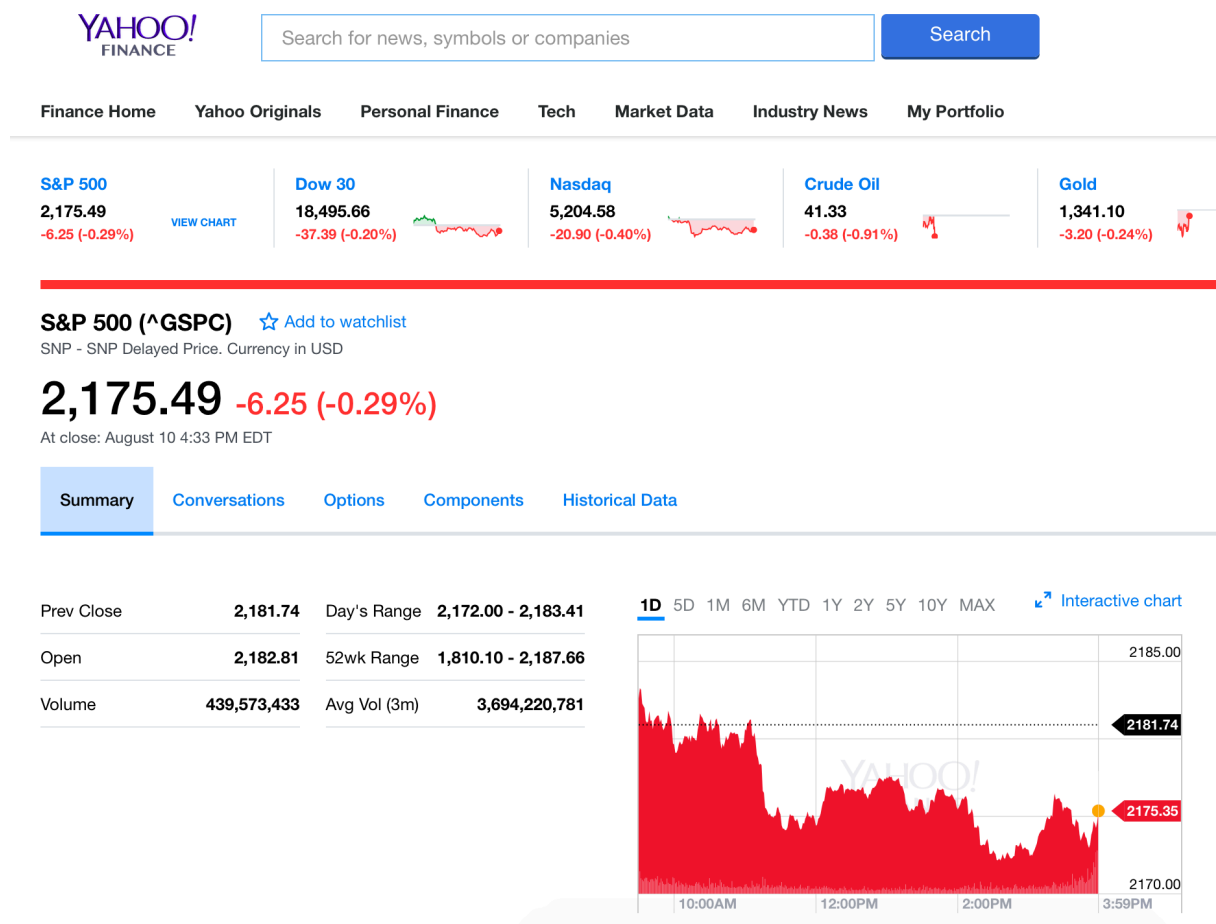
Let's get started.

1. Downloading daily data for S&P 500 Companies

S&P 500 stands for Standard and Poor's 500, and it's stock market index which tracks largest 500 companies based on their market capitalization, listed on the NYSE and NASDAQ which are American stock exchanges.

I will open my browser and navigate to finance.yahoo.com. I will type "S&P" to search box and will choose ^GSPC from available symbols.



I will be presented with following page which gives us S&P 500 Index specific information.



I will click on "Historical Data", in this page, I will adjust time range of dataset by selecting "MAX" and I will click "Apply" and then "Download Data"

| Time Period: | Jan 03, 1950 - Aug 11, 2016 ⌄ | | | | Show: Historical Prices ⌄ | | | Frequency: Daily ⌄ | | | Apply |

| 1D | 5D | 3M | 6M |
| YTD | 1Y | 5Y | MAX |

| Currency in | | | | | | | | | ⬇ Download Data |
| Date | | | | High | Low | Adj Close* | | | Volume |

| Start Date | End Date |
| 1/3/1950 | 8/11/2016 |

Aug 10, 2| | | | 2,183.41 | 2,172.00 | 2,175.49 | | | 3,254,950,000 |

What I will have is comma separated file with S&P 500 Index data. I will rename file as "sp500.csv"

I will open SPSS Modeler and save stream as "Lab_1" into my lab folder.

From "Sources" palette, I will add "Var. File" node to import sp500.csv file.

My import settings looks like following and these are default settings when you add "Var. File" node into your stream.

You can click "Preview" to see first 10 records of dataset.

You can also see in "Data" tab, field types of the data set, and we can see that field types are correctly identified for all fields.
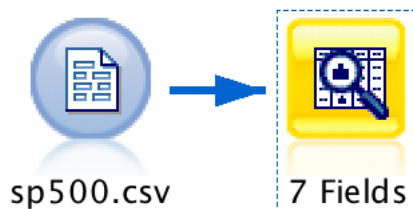


2. Summary Statistics

Before we can continu with summary statistics we can do a little "Data Audit" to have general feel about our dataset.

From "Output" palette, add a "Data Audit" node, double click to open it, check "Advanced Statistics" box and "Calculate Median and Mode" as well and and click "Run"

Settings | Quality | Output | Annotations

◉ Default          ◎ Use custom fields

Fields:

Overlay:

**Display**

☑ Graphs   ☑ Basic statistics   ☑ Advanced statistics

☑ Calculate median and mode (may slow performance on large datasets)

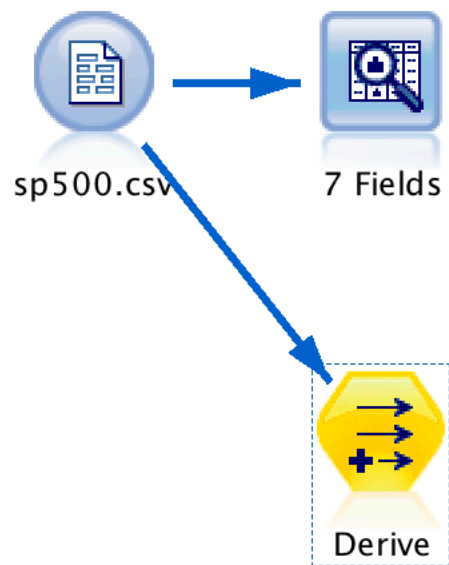OK | ▶ Run | Cancel | Apply | Reset

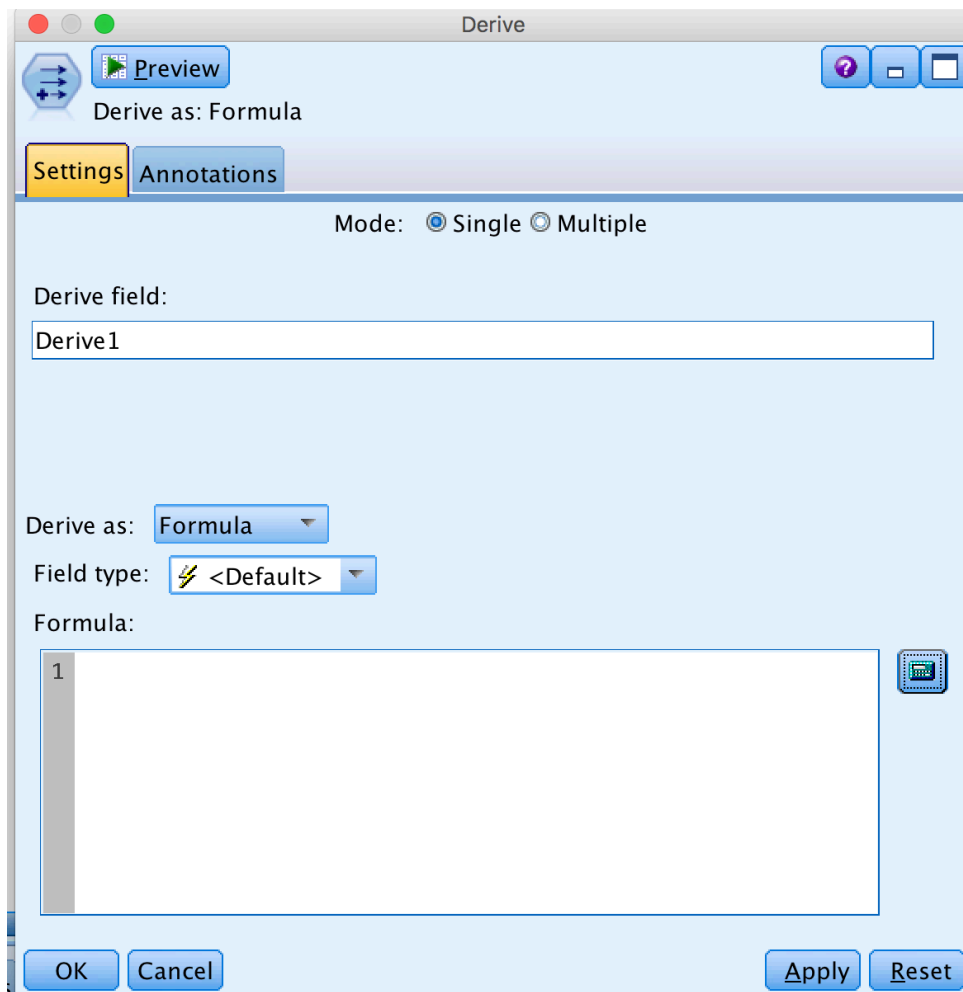| Field | Sample Graph | Measurement | Min | Max | Sum | Range | Mean | Mean Std. Err. | Std. Dev | Variance | Skewness | Skewness Std. Err. |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| Date | | ✎ Continuous | 1950-01-03 | 2016-08-10 | -- | 2101766400... | -- | -- | -- | -- | -- | -- |
| Open | | ✎ Continuous | 16.660 | 2183.760 | 8352737.578 | 2167.100 | 498.373 | 4.432 | 573.779 | 329222.772 | 1.131 | 0.019 |
| High | | ✎ Continuous | 16.660 | 2187.660 | 8405159.744 | 2171.000 | 501.501 | 4.458 | 577.131 | 333080.325 | 1.128 | 0.019 |
| Low | | ✎ Continuous | 16.660 | 2178.610 | 8297530.715 | 2161.950 | 495.079 | 4.404 | 570.172 | 325096.620 | 1.134 | 0.019 |
| Close | | ✎ Continuous | 16.660 | 2182.870 | 8354845.730 | 2166.210 | 498.499 | 4.433 | 573.906 | 329367.754 | 1.131 | 0.019 |
| Volume | | ✎ Continuous | 680000 | 11456230400 | 1385896683... | 11455550400 | 826907329... | 11489640.8... | 1487453720... | 2212518570961... | 2.131 | 0.019 |
| Adj Cl... | | ✎ Continuous | 16.660 | 2182.870 | 8354845.730 | 2166.210 | 498.499 | 4.433 | 573.906 | 329367.754 | 1.131 | 0.019 |

We can see main properties of dataset such as minimum and maximum values, mean, standard deviation, variance, skewness and kurtosis which will be important when we are examining return series.
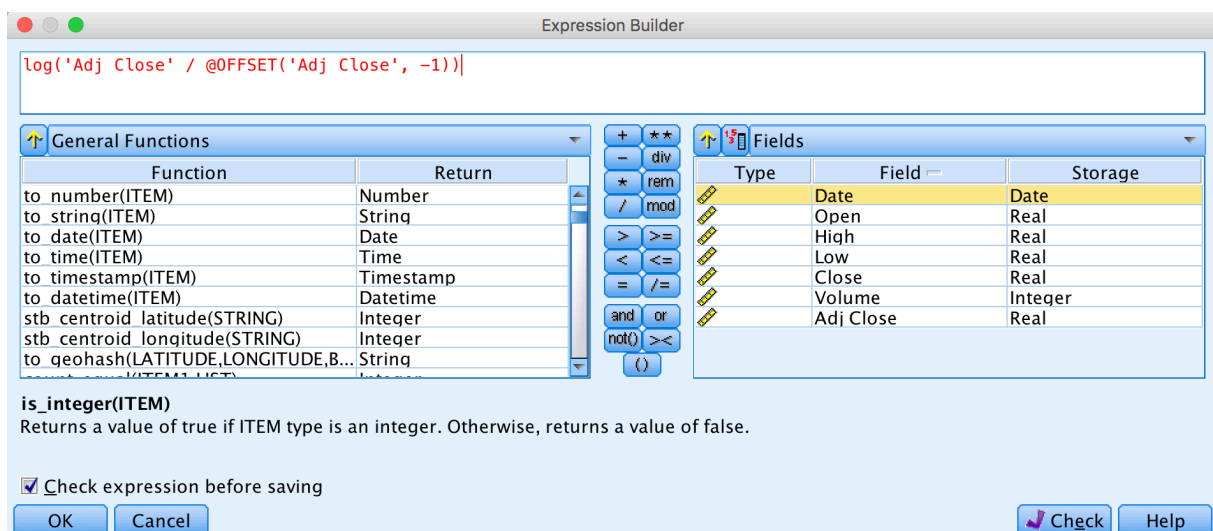
3. Return distributions

In order to see return distributions, we can calculate daily return based on "Adj Close" field by using "Derive" node.



"Derive" node allows us to use CLEM expressions which are set of ready-to-use libraries which you can use to manipulate your data set.

You can click on this button to Launch Expression Builder and see available function also check syntax of your expressions before running them.
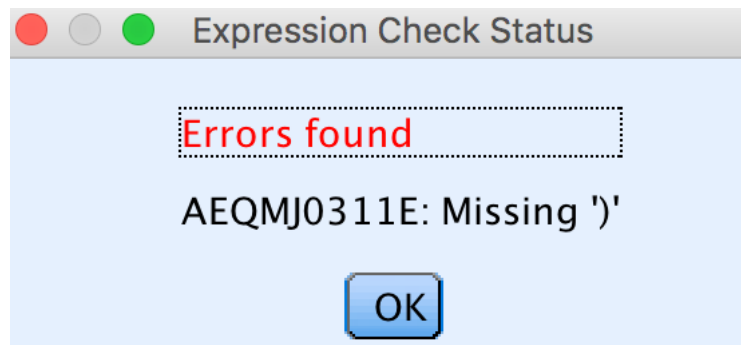


In this case I will use "log" and "OFFSET" function to calculate log returns based on "Adj Close" field.

In text box above, I write my expression and click "Check", if there are no errors, expression color will turn to black

```
log('Adj Close' / @OFFSET('Adj Close', -1))
```
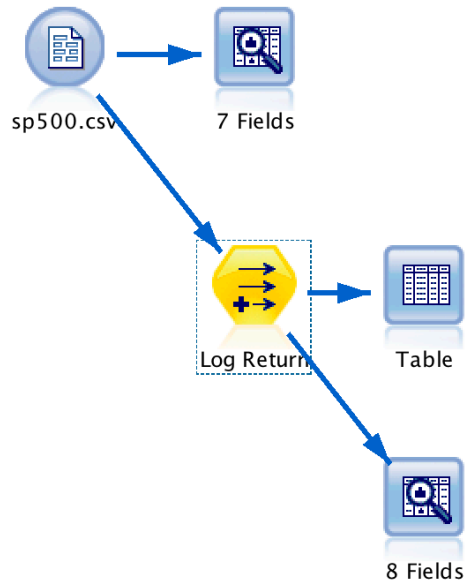
if I delete last paranthesis, then click "Check" again, it will give me an error.

**Expression Check Status**

**Errors found**

AEQMJ0311E: Missing ')'

OK

I can click "OK", and set name for this field in "Derive field:" text box and click "OK"

**Log Return**

Preview

Derive as: Formula

**Settings** Annotations

Mode: ◉ Single ◯ Multiple

Derive field:

Log Return

Derive as: Formula ▾

Field type: ⚡ <Default> ▾

Formula:

```
1 log('Adj Close' / @OFFSET('Adj Close', -1))
```

I can add "Table" node from "Output" palette to see my derived field "Log Return"
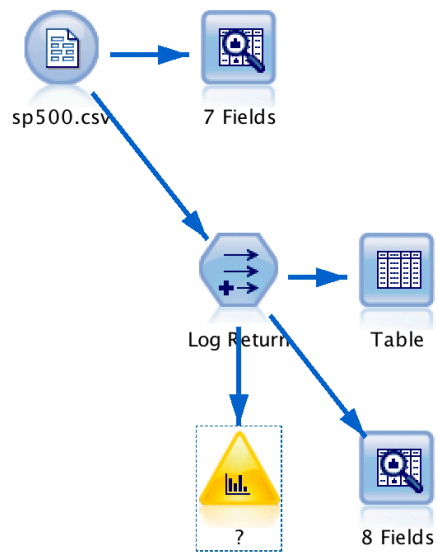


I can add "Data Audit" Node here again to see main properties for "Log Return" field.
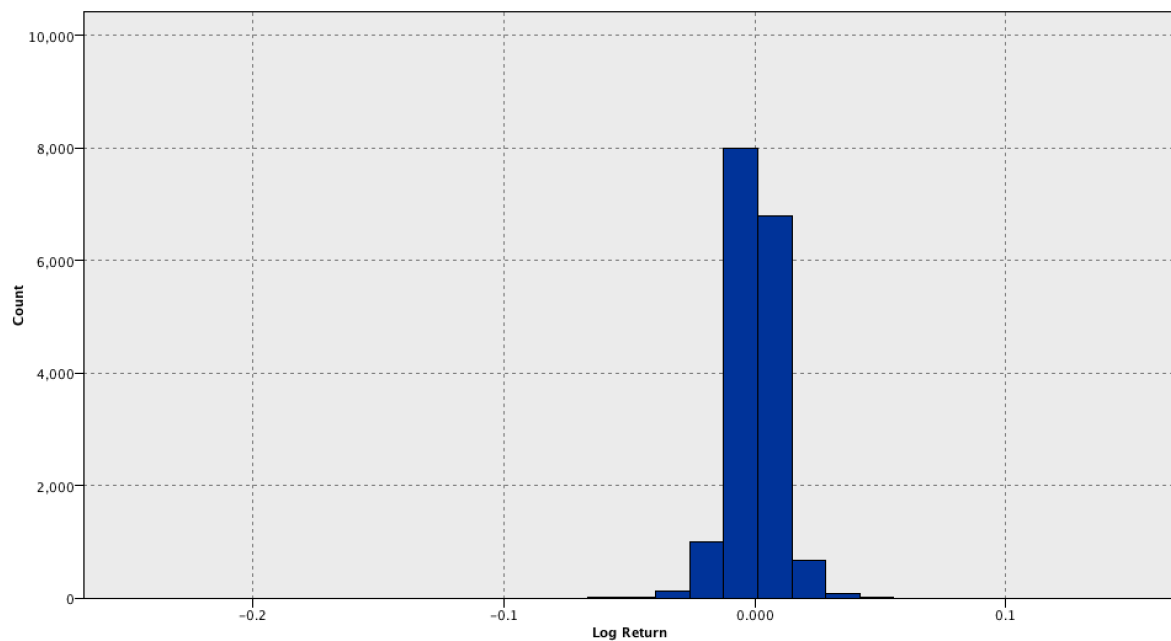


We can see here minimum return and maximum return. Notice how index lost almost 23% of it's value in just one day.

Now I can see return distribution of S&P 500 Index by adding "Histogram" node from "Graphs" palette.
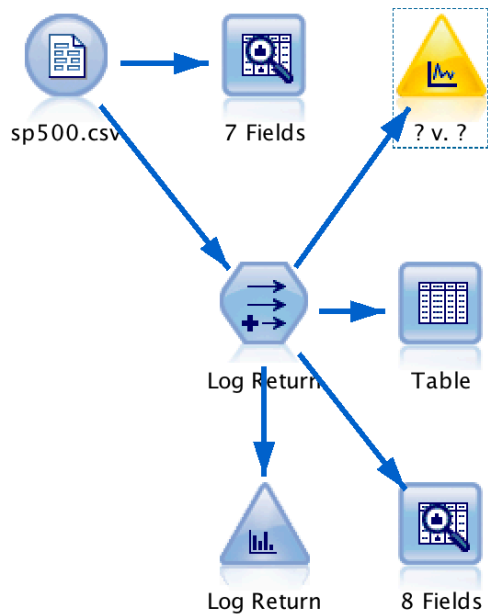


Double click on "Histogram" to open it and only thing you need to do here is selecting "Log Return" field.
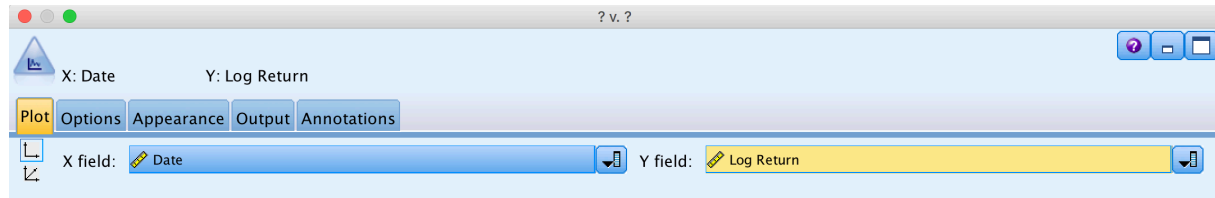


Histogram shows us distribution of daily log returns and notice extreme tails of the distribution.
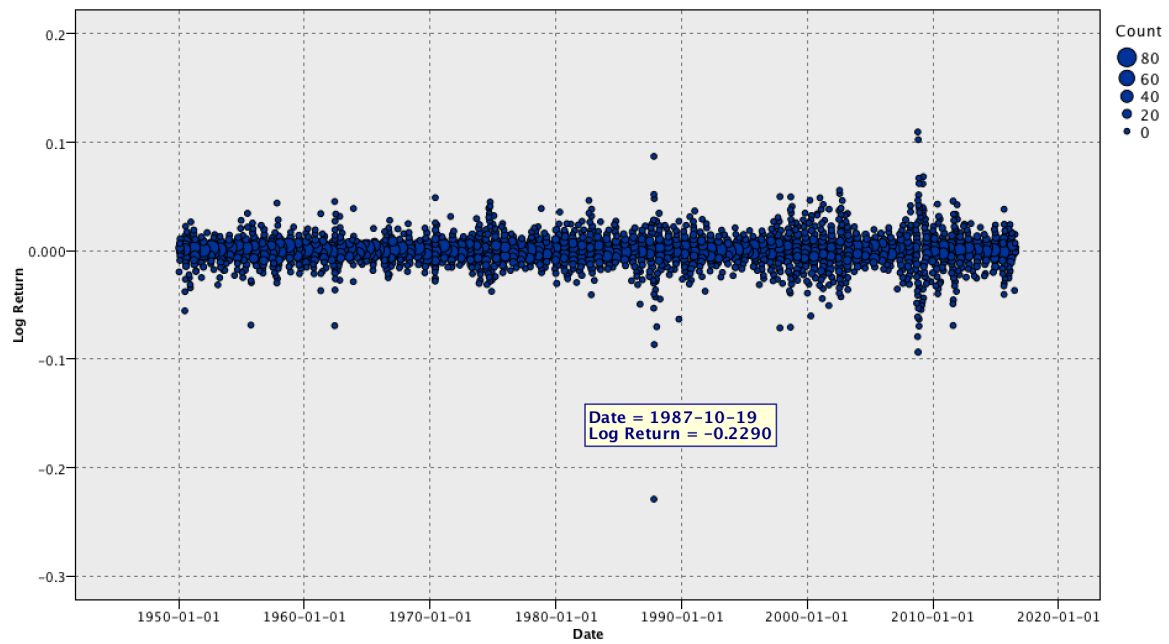
## 4. Plotting return series

If we would like to see dates for those extreme events we can use "Plot" node.



Set field selection as I do and click run



We can see that 23% daily loss happened on 19[th] October, 1987. This event is referred as "Black Monday" where markets recorded huge losses.

Summary

In this lab, you have learned basic functionalities of SPSS Modeler to import and work with your dataset.

Thank you and hope to see you next time.