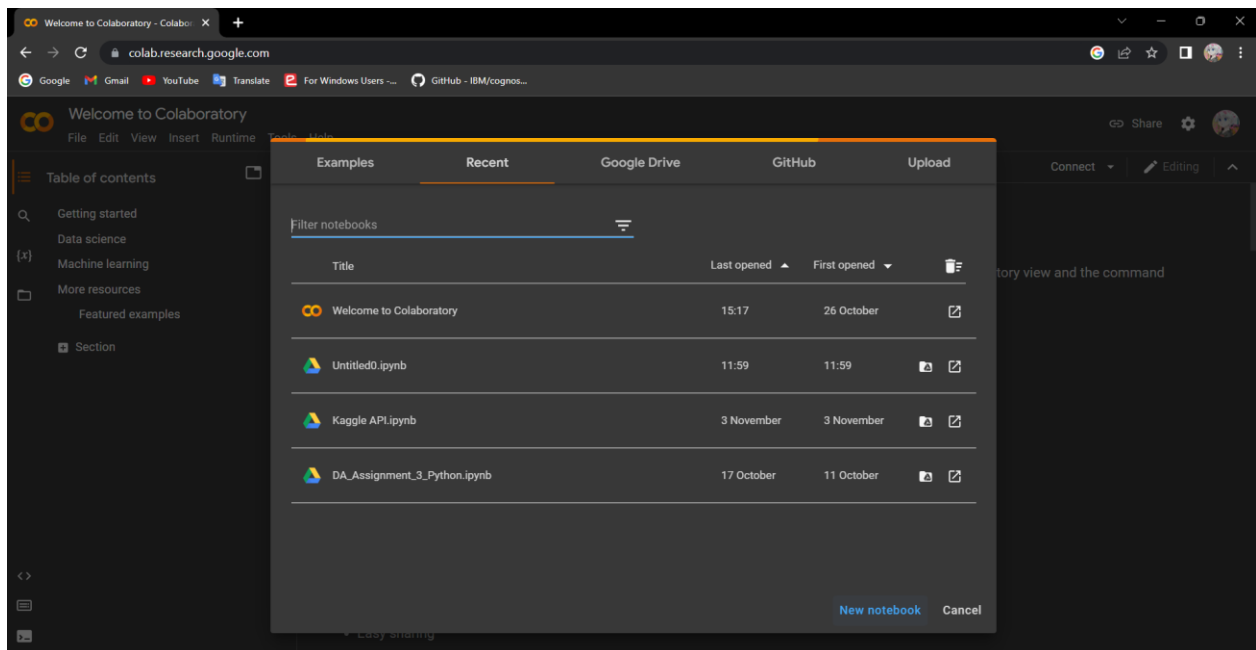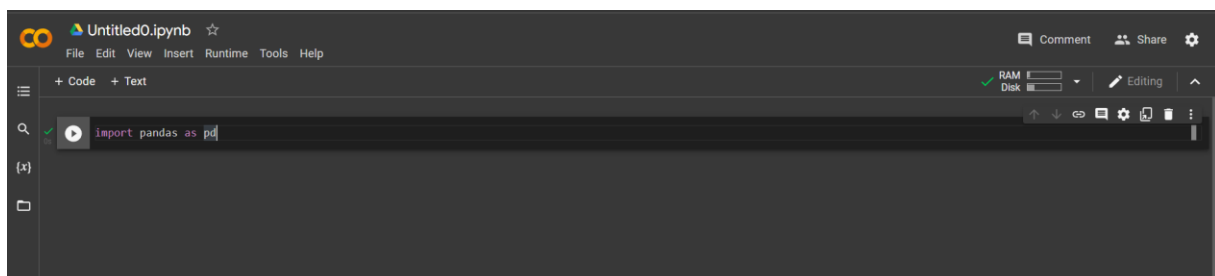| Date | 03 November 2022 |
|---|---|
| Team ID | PNT2022TMID30128 |
| Project Name | Project - Global Sales Data Analytics |

# Sprint 3 - Dataset exploration

**Dataset exploration:** It is the process of finding what are the things that are present in the dataset.

We have used Google collab to explore the dataset.

**Step 1:** Creating a new notebook



**Step 2:** Loading the data into jupyter notebook by using a python library called Pandas (for Python Data Analysis Library).

**Step 3:** Exploring the data by understanding the following things:



1. Finding out the total number of rows and columns in the dataset



**Rows =** 51290 **Columns =** 24

2. Finding if there is any duplicate value

3. How the dataset looks like (Finding out how the variables are organized)



First five rows of the dataset



Last five rows of the dataset

4. What kind of data types are present in the dataset

```
saledata.dtypes
```

```
Row ID             int64
Order ID          object
Order Date        object
Ship Date         object
Ship Mode         object
Customer ID       object
Customer Name     object
Segment           object
City              object
State             object
Country           object
Postal Code      float64
Market            object
Region            object
Product ID        object
Category          object
Sub-Category      object
Product Name      object
Sales            float64
Quantity           int64
Discount         float64
Profit           float64
Shipping Cost    float64
Order Priority    object
dtype: object
```

5. Finding out the overview information about the dataset

```
[14] saledata.info()
```

```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 51290 entries, 0 to 51289
Data columns (total 24 columns):
 #   Column          Non-Null Count  Dtype
---  ------          --------------  -----
 0   Row ID          51290 non-null  int64
 1   Order ID        51290 non-null  object
 2   Order Date      51290 non-null  object
 3   Ship Date       51290 non-null  object
 4   Ship Mode       51290 non-null  object
 5   Customer ID     51290 non-null  object
 6   Customer Name   51290 non-null  object
 7   Segment         51290 non-null  object
 8   City            51290 non-null  object
 9   State           51290 non-null  object
 10  Country         51290 non-null  object
 11  Postal Code     9994 non-null   float64
 12  Market          51290 non-null  object
 13  Region          51290 non-null  object
 14  Product ID      51290 non-null  object
 15  Category        51290 non-null  object
 16  Sub-Category    51290 non-null  object
 17  Product Name    51290 non-null  object
 18  Sales           51290 non-null  float64
 19  Quantity        51290 non-null  int64
 20  Discount        51290 non-null  float64
 21  Profit          51290 non-null  float64
 22  Shipping Cost   51290 non-null  float64
 23  Order Priority  51290 non-null  object
```
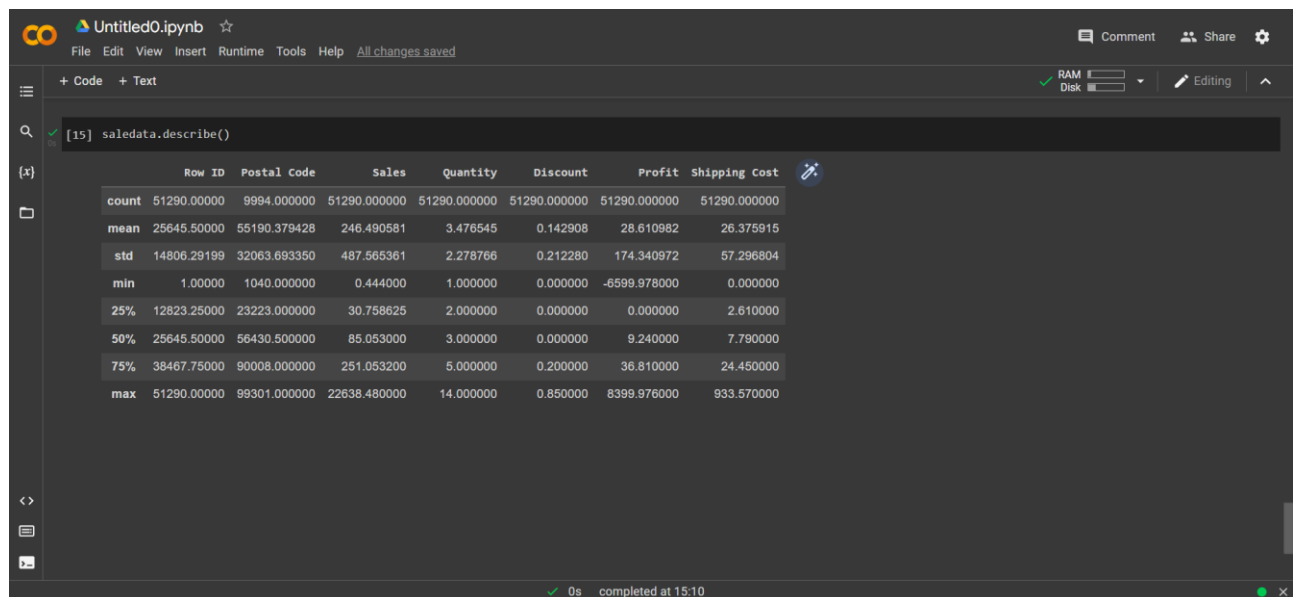
✓ 0s   completed at 15:10

6. Finding out minimum, maximum and mean value of all numerical variables in the dataset.



Each column in the dataset represents different information about the sale

- **Row ID -** It is used to uniquely identify a row in a table.
- **Order ID -** This ID is generated when the order is placed.
- **Order Date -** It shoe the date on when the order is Placed.
- **Ship Date -** The Shipment date of the product.
- **Ship Mode -** In which mode the shipment process is carried out.
- **Customer ID -** ID is generated when the customer Places the first order.
- **Customer Name -** It shows the name of the customer.
- **Segment -** It shows the segment of the customer.
- **City -** The city in which the customer lives.
- **State -** The state in which the customer resides.
- **Country -** It gives the country of the customer.
- **Market -** Market where the order is placed.
- **Region -** It gives the region of the market.
- **Product ID -** This is a unique ID generated for each Product.
- **Category -** It shows which category of the product.
- **Sub-Category -** Sub-Category to which the product Belongs.
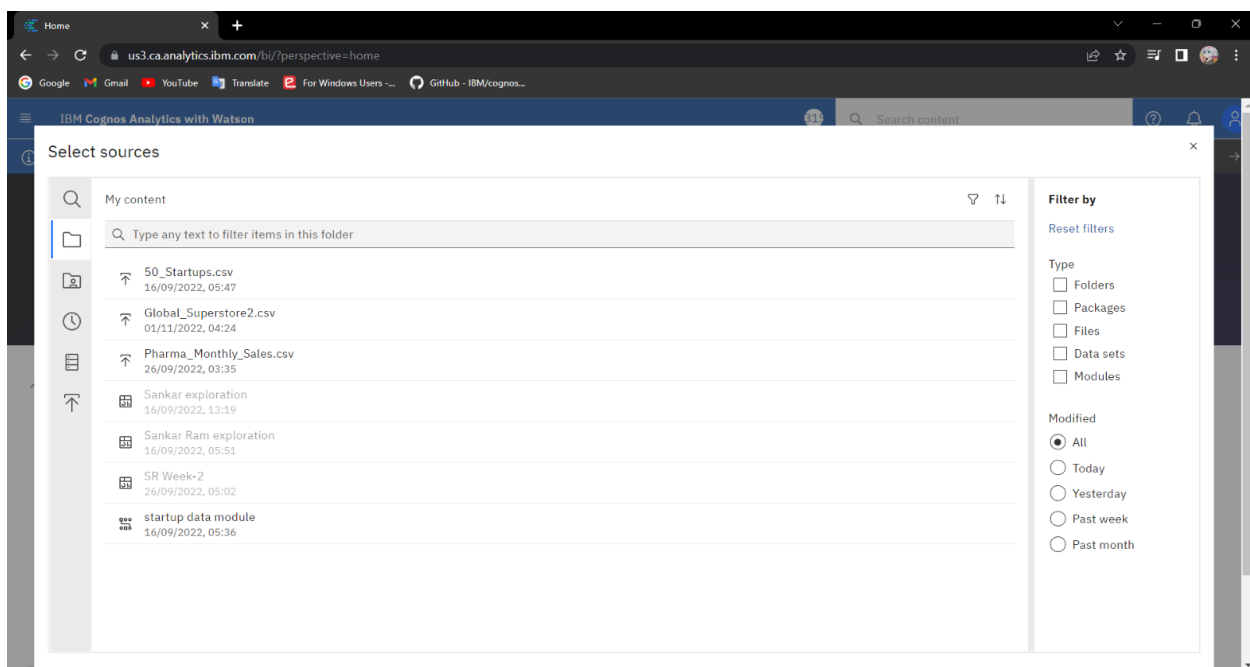- **Product Name -** The name of the product is mentioned.

- **Sales** - It shows the sales details of the product.
- **Quantity** - It shows the number of products ordered.
- **Discount** - How much discount is provided.
- **Profit** - The profit earned by the retailer.
- **Shipping cost** - The cost of shipping.
- **Order Priority** - It shows the priority level of the order.

**Preparing the dataset:** It's the process of cleaning and transforming raw data prior to processing and analysis.
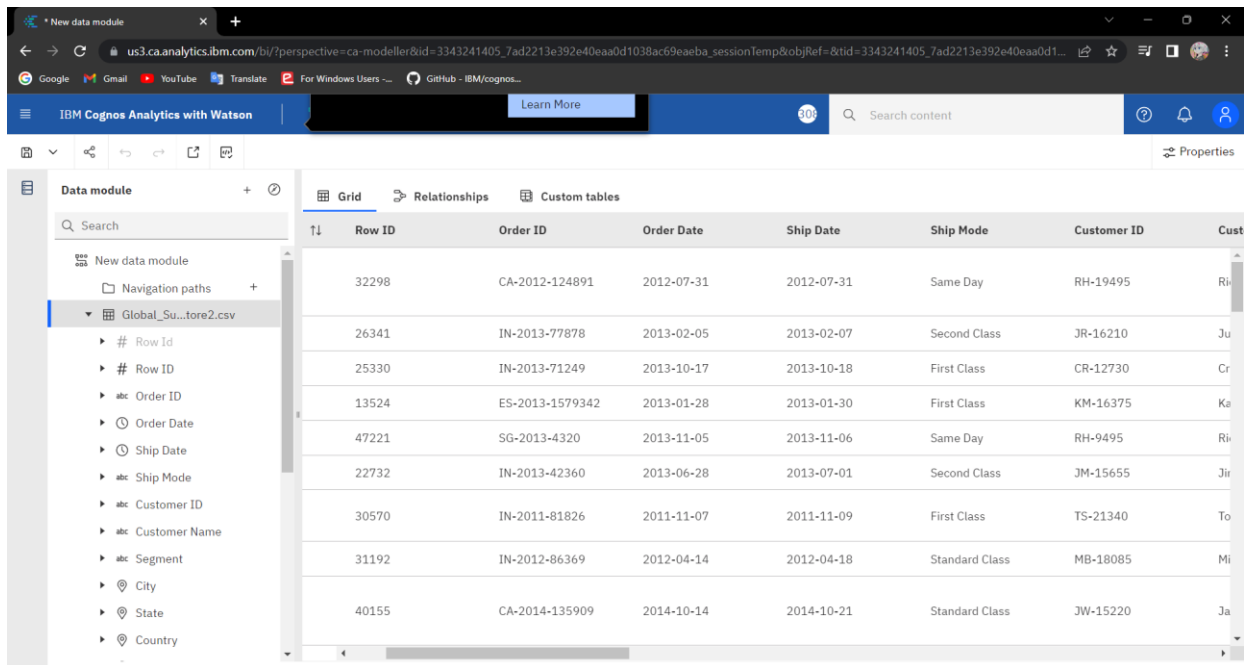
Once we load the data into IBM Cognos, we need to prepare the following:

- Prepare Calculations of Year, Month, Day fields and the related Navigation path.

Select the dataset from the prepare data section.

Once the dataset is loaded make the necessary calculations and navigation path
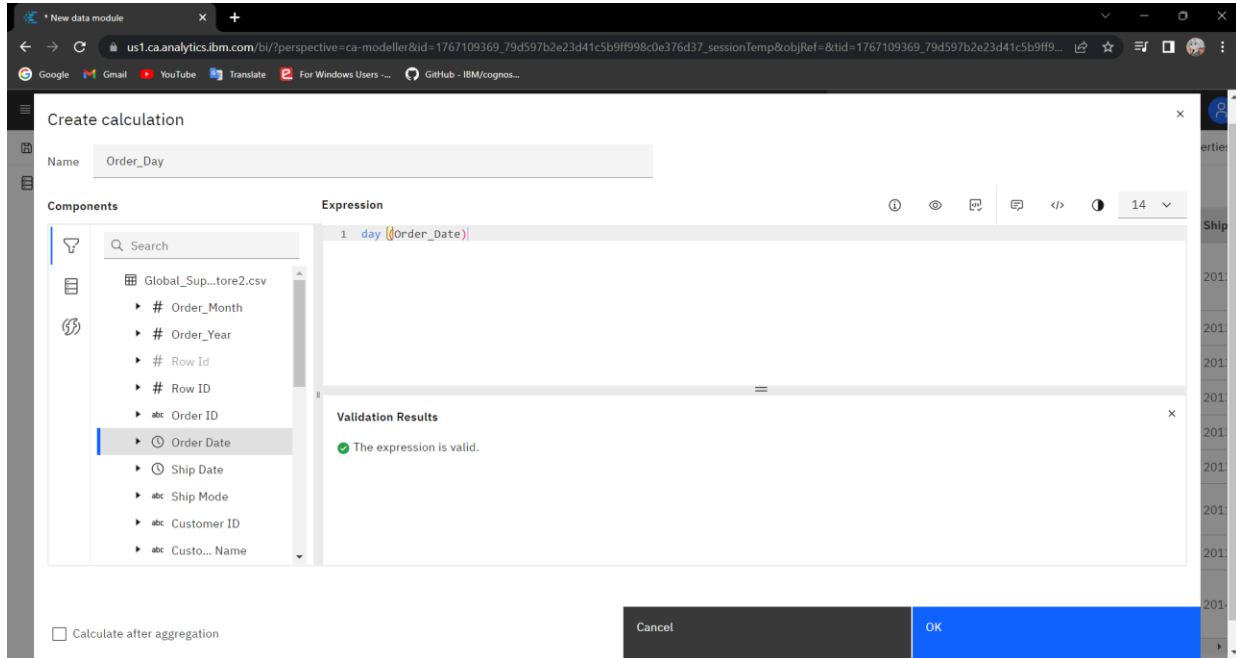


Creating new calculation

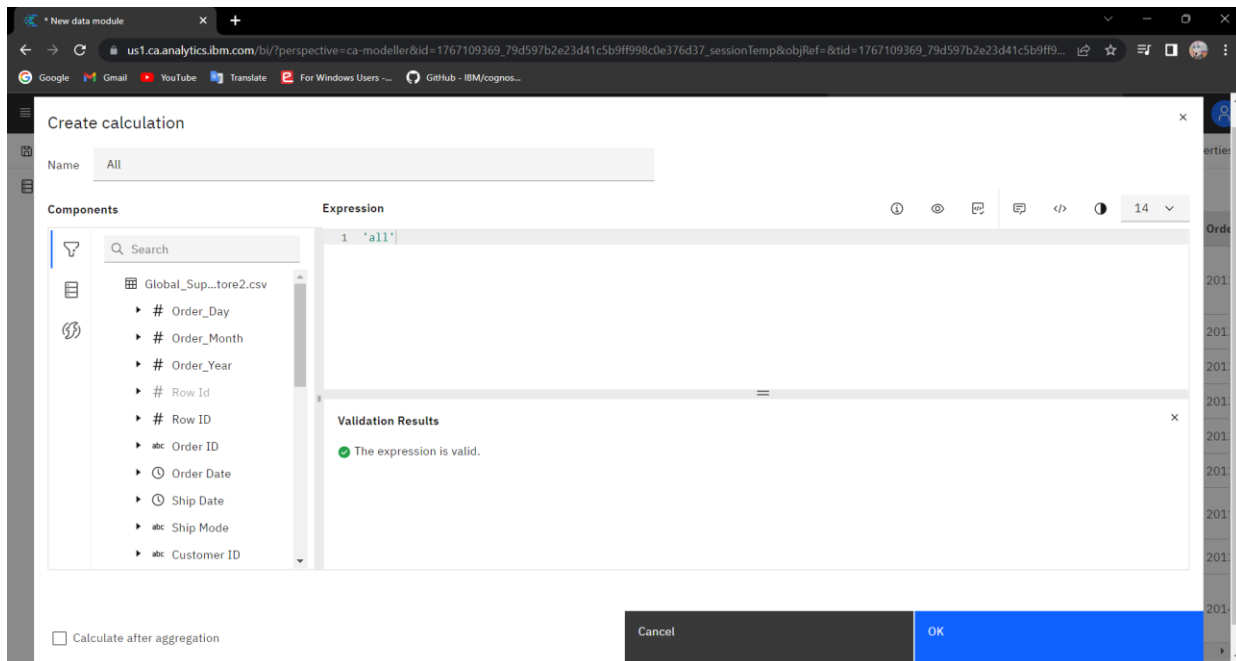# Creating "Order_Year" calculation
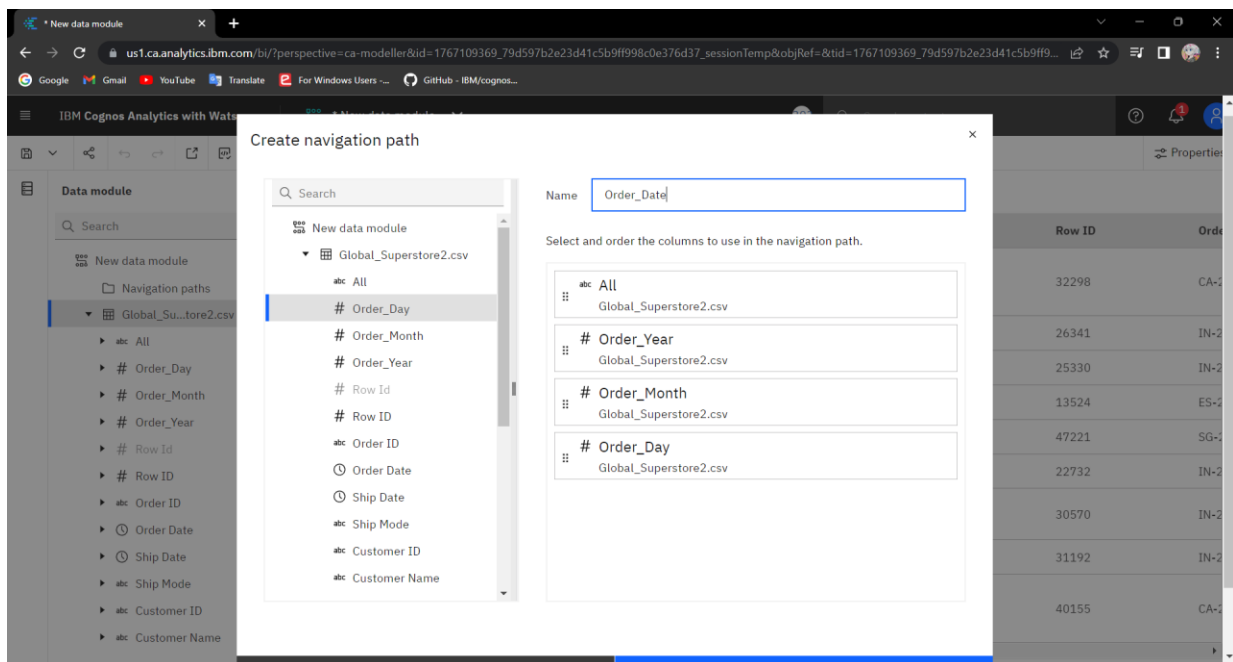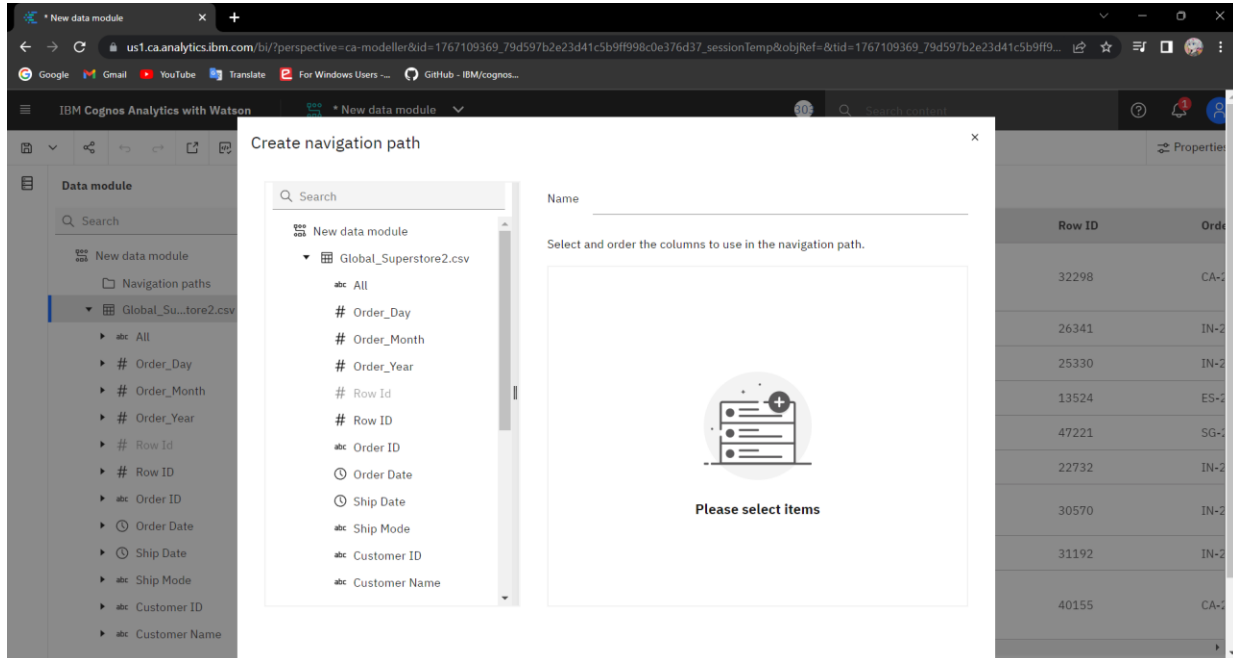


# Creating "Order_Month" calculation

# Creating "Order_Day" calculation



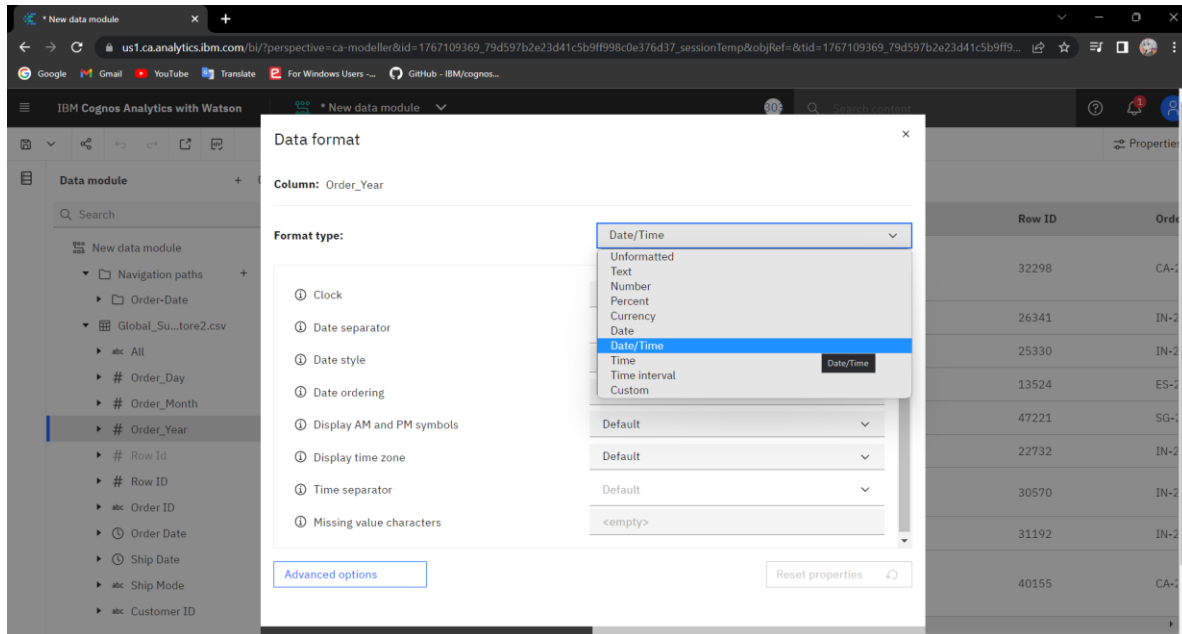# Creating "All" calculation for navigation path

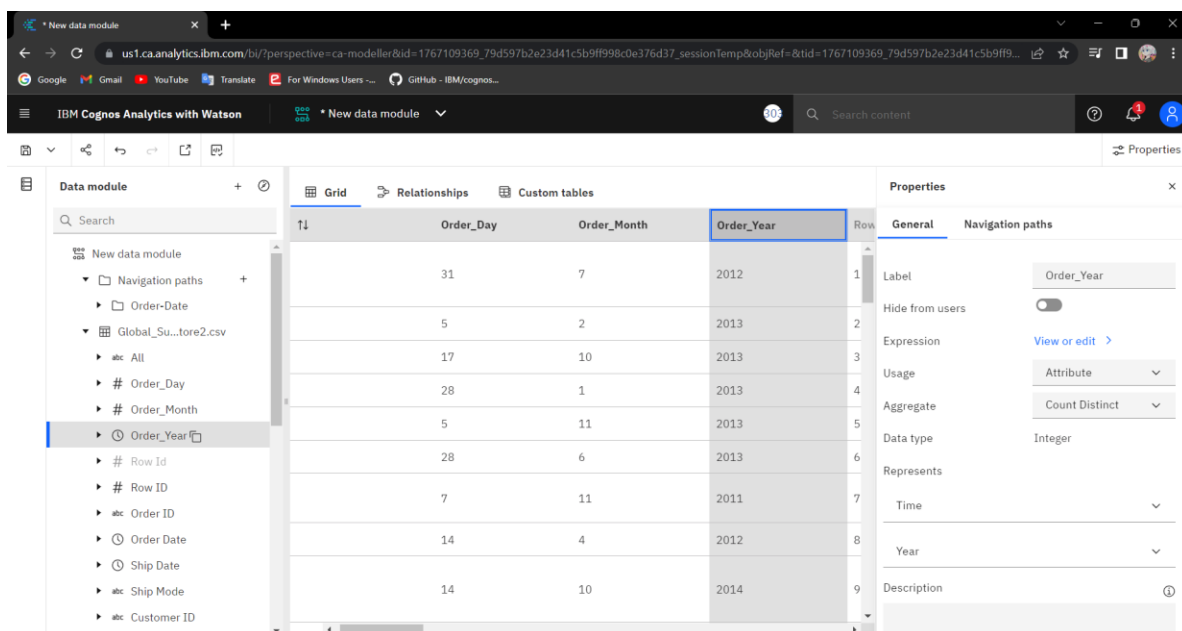# Creating "Order_Date" navigation path





The order year, month and day are in numerical value, so we have to change it to date values.

The following steps are used to change from numerical value to date values

**Step 1:** Click on **Order_Year** and --> Format data and change the format type to "Date/Time" then click ok.



**Step 2:** Go to **Order_Year** properties and change the Usage to 'Attribute', Aggregate to 'Count distinct', Represents to 'Time' and change the display option to 'Show members'.

Repeat steps 1 and 2 for **Order_Month** and **Order_Day** To change it to date values.



Now all the numerical values are changed into date value

Then save the dataset