

**OUTLIERS AND ANOMALY DETECTION**  
**USING**  
**POWERBI DAX**  
**AND**  
**PYTHON MACHINE LEARNING**

Table of Contents

1. MACHINE LEARNING ANOMALY DETECTION IN POWERBI .....	2
2. STATISTICAL OUTLIER DETECTION IN TIME SERIES DATA .....	4
3. ANOMALY DETECTION IN SUPERSTORE SALE DATA .....	6

## 1. MACHINE LEARNING ANOMALY DETECTION IN POWERBI

This project detects anomaly using Isolation Forest method using Machine Learning with PowerBI DAX formulas and Python scrips to analyze the temperature parameter of a machine.

The main idea, which is different from other popular outlier detection methods, is that Isolation Forest explicitly identifies anomalies instead of profiling normal data points. Isolation Forest, like any tree ensemble method, is built on the basis of decision trees. In these trees, partitions are created by first randomly selecting a feature and then selecting a random split value between the minimum and maximum value of the selected feature.

The Isolation Forest 'isolates' observations by randomly selecting a feature and then randomly selecting a split value between the maximum and minimum values of the selected feature.

### **Backgroud of the Dataset used:**

The dataset contains day wise temperature recording for a machine. The anomaly detection is performed for the temperature column in the dataset. A sample of the data used is provided below.

Sample source database table: Source: Github, Type: CSV, Columns: 2, Rows: 1769

Date	Temperature
03/01/2014	45.868
03/02/2014	47.606
03/03/2014	42.58
03/04/2014	46.03
03/05/2014	44.992
03/06/2014	45.238
03/07/2014	45.752
03/08/2014	46.476
03/09/2014	42.752
03/10/2014	46.156
03/11/2014	43.35
03/12/2014	46.292
03/13/2014	43.428
03/14/2014	46.314
03/15/2014	43.688
03/16/2014	45.4

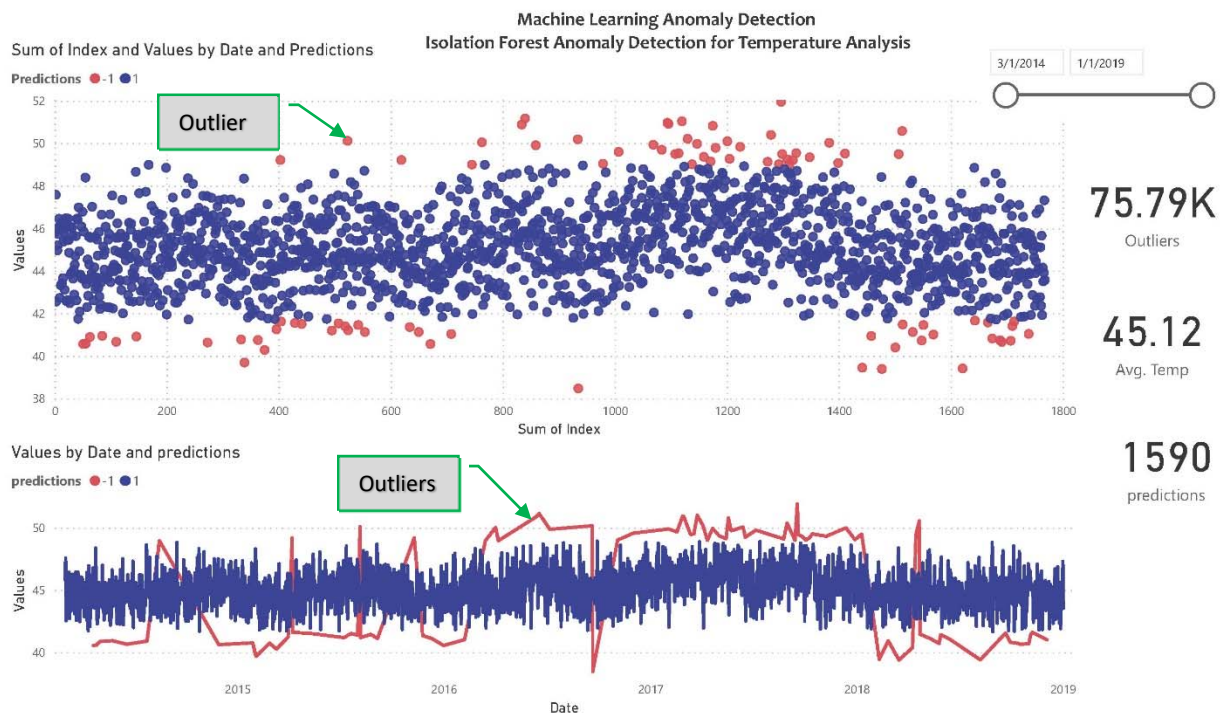
The above raw data was extracted and transformed in PowerBI for further analysis using Python code. The SKlearn Machine Learning module was used to determine the outliers using the Isolation Forest method. The outliers identified are represented as (-1) under the Predictions column in the below table.

The screenshot shows the PowerBI Desktop interface. The main view is a data table with columns: Date, Values, Index, and Predictions. The table contains data from December 2016 to April 2015. A Python script window is open, showing a script for Isolation Forest anomaly detection. The script uses the sklearn library to create a model and fit it to the data. Callouts highlight specific features: 'Predictions from Machine Learning model with Outliers' points to the Predictions column, 'PowerBI DAX operations using Power Query Editor' points to the Fields pane, and 'Python script for ML Model' points to the Python script window.

Date	Values	Index	Predictions
Saturday, December 22, 2016	46.924	1757	1
Sunday, December 23, 2018	42.792	1758	1
Monday, December 24, 2018	45.134	1759	1
Tuesday, December 25, 2018	42.382	1760	1
Wednesday, December 26, 2018	45.68	1761	1
Thursday, December 27, 2018	41.948	1762	1
Friday, December 28, 2018	44.178	1763	1
Saturday, December 29, 2018	43.526	1764	1
Sunday, December 30, 2018	43.944	1765	1
Monday, December 31, 2018	47.344	1766	1
Tuesday, January 01, 2019	43.594	1767	1
Sunday, April 20, 2014	40.586	50	-1
Friday, April 25, 2014	40.61	55	-1
Friday, May 02, 2014	40.916	62	-1
Saturday, May 24, 2014	40.97	84	-1
Wednesday, June 18, 2014	40.69	109	-1
Thursday, July 24, 2014	40.938	145	-1
Friday, November 28, 2014	40.658	272	-1
Tuesday, January 27, 2015	40.798	332	-1
Monday, February 02, 2015	39.718	338	-1
Thursday, February 26, 2015	40.774	362	-1
Tuesday, March 10, 2015	40.306	374	-1
Tuesday, March 31, 2015	41.278	395	-1
Tuesday, April 07, 2015	49.238	402	-1

The predicted output values are plotted in PowerBI to visually represent the outliers in the dataset. The below dashboard provides a dynamic visualization of the complete dataset.

### MACHINE LEARNING ANOMALY DETECTION

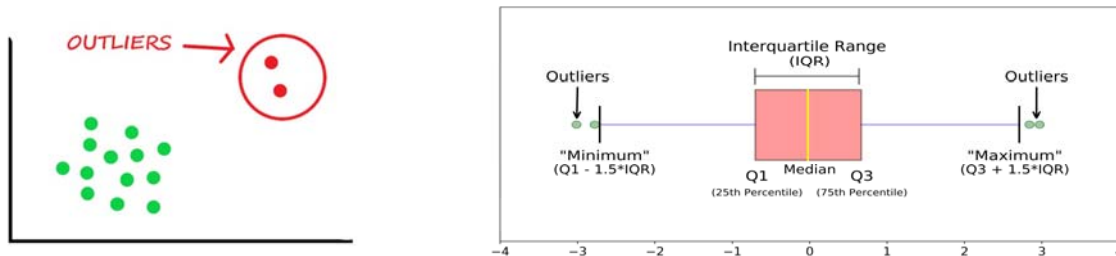


## 2. STATISTICAL OUTLIER DETECTION IN TIME SERIES DATA

Outlier detection or Anomaly detection is the identification of rare items, events or observations which rise suspicions by differing significantly from the majority of the data. Outlier detection is the process of detecting and excluding them from a given set of data. Typical anomalous problems in a real world could be:

- A bank fraud
- Medical problems
- A structural defect
- Errors in a text, measurement, etc

Outliers can be explained as the data point that is more than 1.5 times the Inter Quartile Range (IQR) where Q1 is the median of the lower half of the data, and Q3 is the median of the upper half of the data.



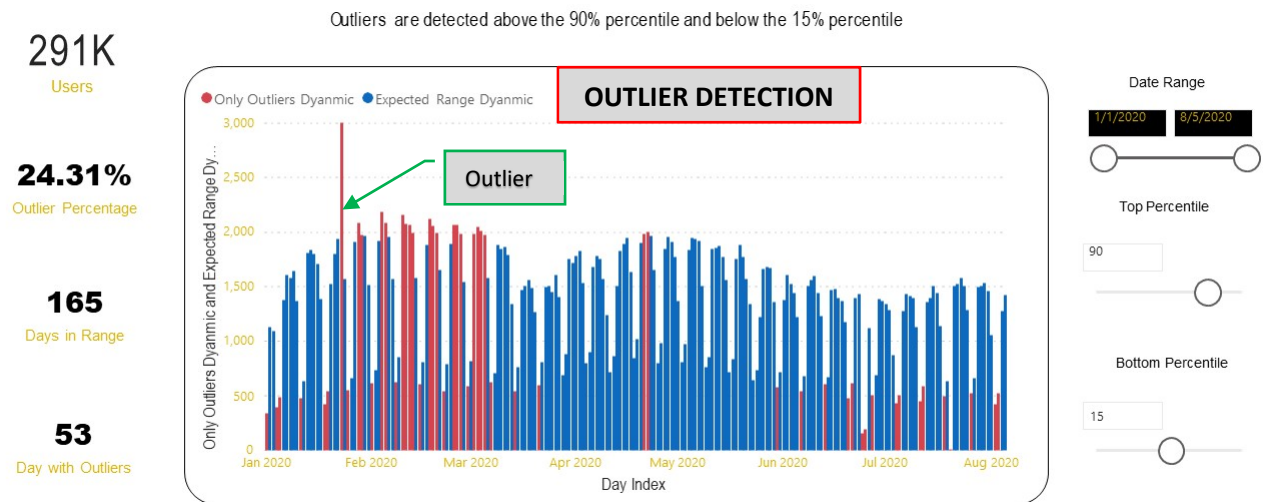
Microsoft PowerBI provides anomaly detection and statistical analysis methods that can be applied to a dataset. The result of the analysis can be translated as a report and a dashboard visual for easy detection. The analysis also provides a deeper insight into the business. A typical data was used to identify the outliers in a dataset using PowerBI and the results are presented below.

### **Background of the Dataset used:**

The dataset contains day wise users index for a newly launched website. The number of users visiting the website are plotted against the day index. A sample of the data used is provided below.

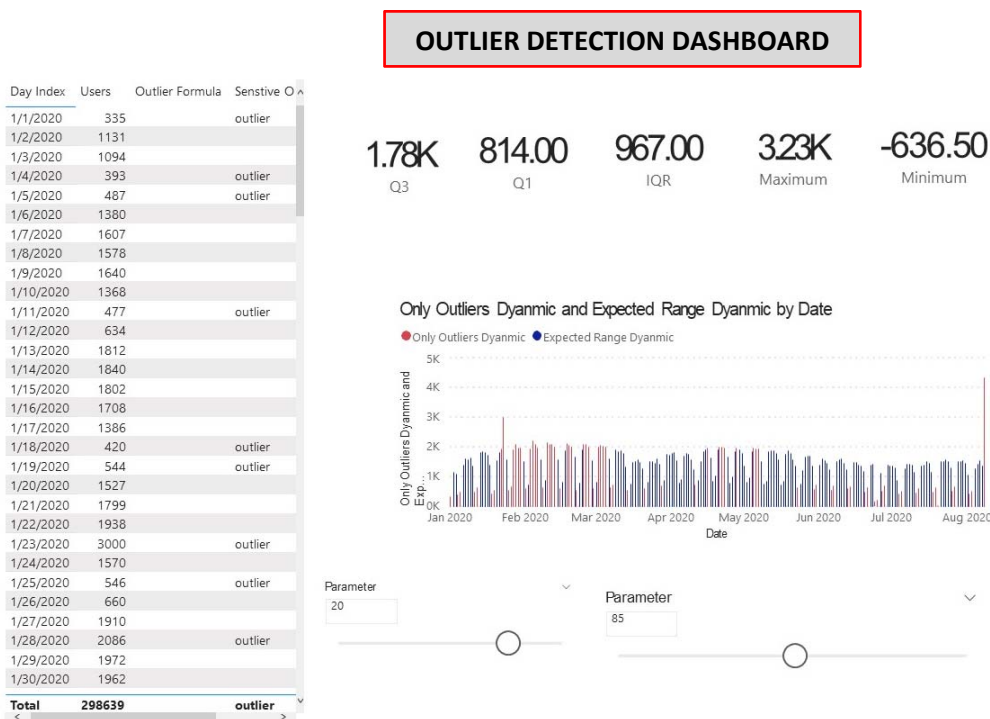
Sample database table: Source: Github, Type: CSV, Columns: 1, Rows: 222

DayIndex Users
01/01/2020 0:00 335
01/02/2020 0:00 1131
01/03/2020 0:00 1094
01/04/2020 0:00 393
01/05/2020 0:00 487
01/06/2020 0:00 1380
01/07/2020 0:00 1607
01/08/2020 0:00 1578
01/09/2020 0:00 1640



In the above chart, the outliers are calculated in PowerBI from the given dataset and they are represented as red coloured bars.

Dynamic dashboard with table view developed in PowerBI.



## Conclusion:

Outlier information is very useful when data is compared with the original data. The above identification of outliers from the original dataset provides an opportunity to further investigate the causation using PowerBI.

### 3. ANOMALY DETECTION IN SUPERSTORE SALE DATA

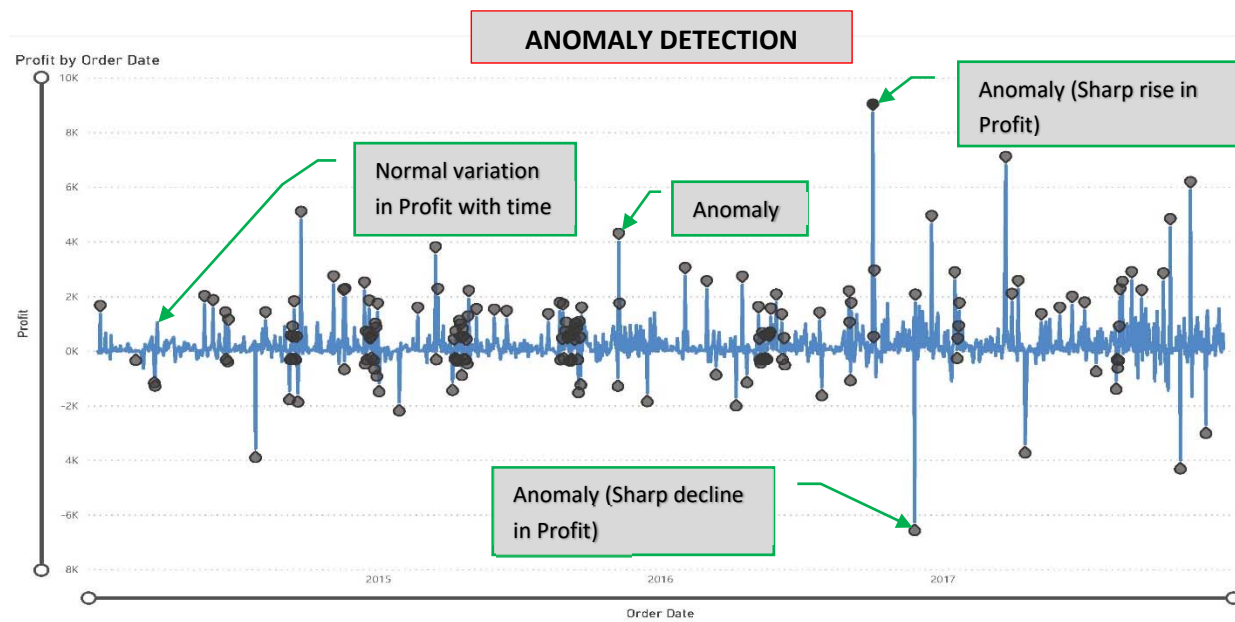
**Anomaly detection** (or outlier analysis) is referred as something that deviates from what is standard, normal, or expected. Anomalous data can indicate critical incidents, such as a fraudulent transaction, technical glitch, or potential opportunities, or a change in consumer behavior. Identifying outliers and correlating with various influencing factors can deliver insights to business decision makers.

#### **Background of the Dataset used:**

The dataset contains the sale of consumer products by a superstore in USA. The Order date and Profits of the superstore was analysed to understand the anomalies in the rise and fall of profits with respect to time of goods purchased. Two sharp in the data points, one for sharp increase in profit one for sharp decline in profit were selected for further analysis.

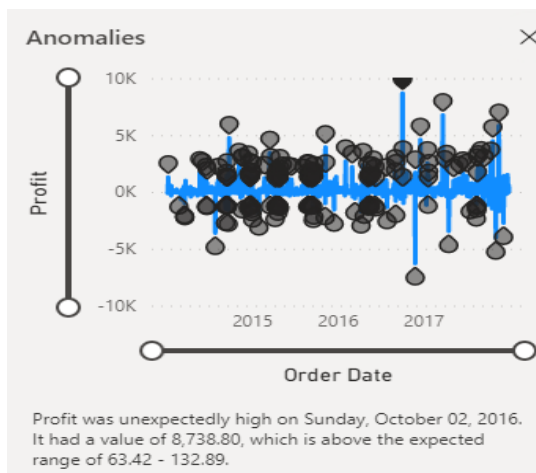
Sample database table: Source: Github, Type: Excel, Columns: 21, Rows: 9995, Columns analysed: Order date and Profit

Row ID	Order Date	Ship Date	Postal Code	Region	Category	Sales	Quantity	Discount	Profit
43	Sunday, July 17, 2016	Friday, July 22, 2016	90049	West	Office Supplies	77.88	2	0	3.894
514	Thursday, December 21, 2017	Monday, December 25, 2017	90049	West	Office Supplies	6.63	3	0	1.7901
....	....	....	....	....	....	....	....	....	....



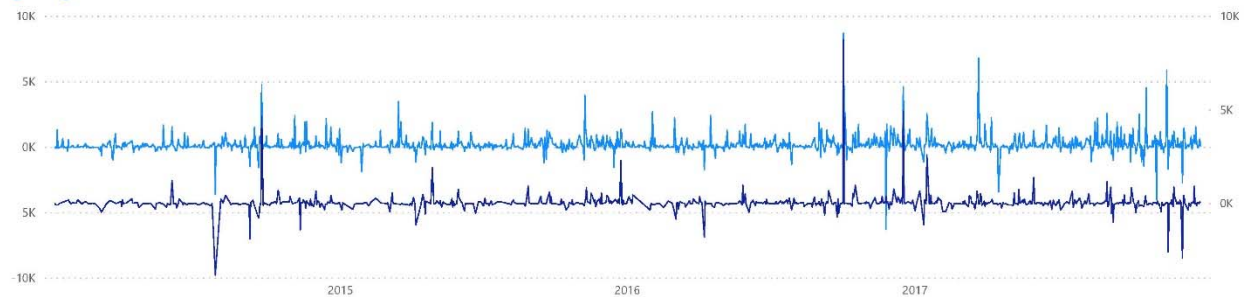
## Rise in Profit – A Deeper Analysis using PowerBI

When the data point was selected, the reasons contributing to the anomaly were available for further investigation. The sharp rise in Profit for the identified above was found to be contributed by two major reasons (Central region and Copiers product sub-category). Their contribution is more than 50% for the rise in profit. The detailed graphs and the findings are presented below.



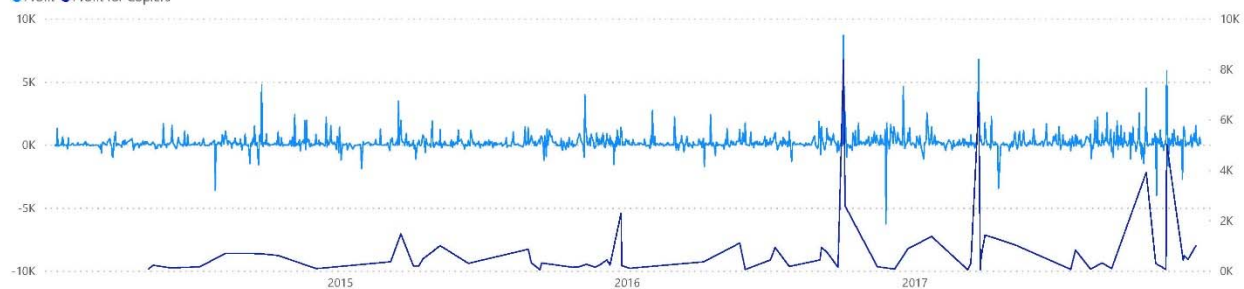
Profit and Profit for Central by Order Date

● Profit ● Profit for Central



Profit and Profit for Copiers by Order Date

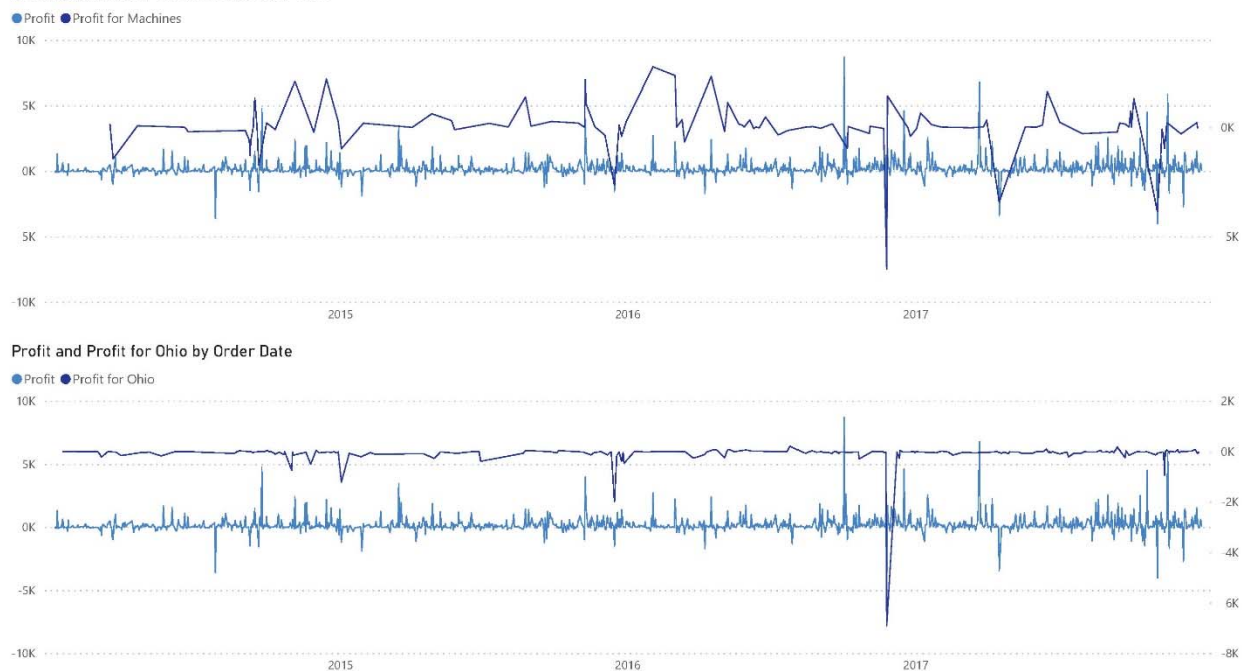
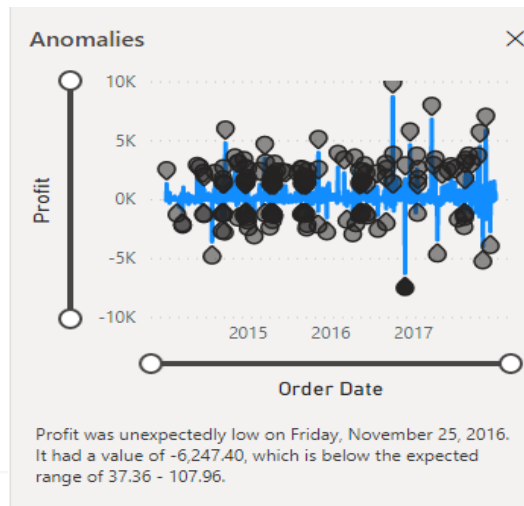
● Profit ● Profit for Copiers





### Profit decline by Sub-category and State-wise Comparison:

Similar to the above the reasons two factors were found as the reasons for sharp decline in Profit in 2016. These two factors, Ohio state and product Machines under the Sub-Category were responsible for more than 50% of the total Profit decline. The detailed graphs and the findings are presented below.



### Conclusion:

From the above charts and analysis, the underlying reasons and their significant impact to the business can be well correlated. The knowledge gained from this analysis could be further leveraged to promote the business interest.