

**TỔNG LIÊN ĐOÀN LAO ĐỘNG VIỆT NAM  
TRƯỜNG ĐẠI HỌC TÔN ĐỨC THẮNG  
KHOA CÔNG NGHỆ THÔNG TIN**



**ĐỒ ÁN GIỮA KÌ MÔN  
HỆ THỐNG THƯƠNG MẠI THÔNG MINH**

## **MIDTERM PROJECT**

*Người hướng dẫn:* **ThS DƯƠNG HỮU PHÚC**

*Người thực hiện:* **PHẠM NGUYỄN – 52000092**

**NGÔ HOÀNG KHÔI – 52000676**

**TRẦN TRUNG HIẾU – 52000888**

**Lớp : 20050261**

**Khoá : 24**

**THÀNH PHỐ HỒ CHÍ MINH, NĂM 2023**

**TỔNG LIÊN ĐOÀN LAO ĐỘNG VIỆT NAM  
TRƯỜNG ĐẠI HỌC TÔN ĐỨC THẮNG  
KHOA CÔNG NGHỆ THÔNG TIN**



**ĐỒ ÁN GIỮA KÌ MÔN  
HỆ THỐNG THƯƠNG MẠI THÔNG MINH**

## **MIDTERM PROJECT**

*Người hướng dẫn:* **ThS DƯƠNG HỮU PHÚC**

*Người thực hiện:* **PHẠM NGUYỄN – 52000092**

**NGÔ HOÀNG KHÔI – 52000676**

**TRẦN TRUNG HIẾU – 52000888**

**Lớp : 20050261**

**Khoá : 24**

**THÀNH PHỐ HỒ CHÍ MINH, NĂM 2023**

## LỜI CẢM ƠN

Chúng em xin chân thành gửi lời cảm ơn này đến thầy Dương Hữu Phúc giảng viên phụ trách giảng dạy bộ môn Hệ thống thương mại thông minh. Nhờ có sự tận tình giảng dạy, truyền đạt kiến thức của quý thầy mà chúng em mới đủ kiến thức để hoàn thành đồ án giữa kỳ này.

Song song với đó, chúng em cũng xin gửi lời cảm ơn đến Khoa Công Nghệ Thông Tin, trường Đại học Tôn Đức Thắng vì đã tạo điều kiện cho chúng em học tập, nghiên cứu trong suốt quá trình học tập môn học này nói riêng và cả quá trình học tại môi trường Đại học nói chung. Một lần nữa chúng em xin gửi lời cảm ơn chân thành đến mọi người và chúc tất cả thật nhiều sức khỏe.

## **ĐỒ ÁN ĐƯỢC HOÀN THÀNH TẠI TRƯỜNG ĐẠI HỌC TÔN ĐỨC THẮNG**

Tôi xin cam đoan đây là sản phẩm đồ án của riêng tôi / chúng tôi và được sự hướng dẫn của ThS Dương Hữu Phúc;. Các nội dung nghiên cứu, kết quả trong đề tài này là trung thực và chưa công bố dưới bất kỳ hình thức nào trước đây. Những số liệu trong các bảng biểu phục vụ cho việc phân tích, nhận xét, đánh giá được chính tác giả thu thập từ các nguồn khác nhau có ghi rõ trong phần tài liệu tham khảo.

Ngoài ra, trong đồ án còn sử dụng một số nhận xét, đánh giá cũng như số liệu của các tác giả khác, cơ quan tổ chức khác đều có trích dẫn và chú thích nguồn gốc.

**Nếu phát hiện có bất kỳ sự gian lận nào tôi xin hoàn toàn chịu trách nhiệm về nội dung đồ án của mình.** Trường đại học Tôn Đức Thắng không liên quan đến những vi phạm tác quyền, bản quyền do tôi gây ra trong quá trình thực hiện (nếu có).

*TP. Hồ Chí Minh, ngày tháng năm*

*Tác giả*

*(ký tên và ghi rõ họ tên)*

*Nguyễn*

*Phạm Nguyễn*

*Khôi*

*Ngô Hoàng Khôi*

*Hiếu*

*Trần Văn C*

## PHẦN XÁC NHẬN VÀ ĐÁNH GIÁ CỦA GIẢNG VIÊN

### Phần xác nhận của GV hướng dẫn

---

---

---

---

---

---

---

Tp. Hồ Chí Minh, ngày      tháng      năm  
(kí và ghi họ tên)

### Phần đánh giá của GV chấm bài

---

---

---

---

---

---

---

Tp. Hồ Chí Minh, ngày      tháng      năm  
(kí và ghi họ tên)

## **SUMMARY**

The midterm project questions will be organized into two types

- QUESTION 1: Exploratory data analysis (EDA).

The midterm project questions will be organized into two types

- QUESTION 2: Augumented the dataset.

The midterm project questions will be organized into two types

## MỤC LỤC

LỜI CẢM ƠN .....	i
PHẦN XÁC NHẬN VÀ ĐÁNH GIÁ CỦA GIẢNG VIÊN .....	iii
SUMMARY .....	iv
MỤC LỤC.....	1
DANH MỤC CÁC BẢNG BIỂU, HÌNH VẼ, ĐỒ THỊ .....	3
QUESTION 1 – EXPLORATORY DATA ANALYSIS (EDA) .....	5
1.1 Topic .....	5
1.2 Answer .....	5
1.2.1 Hypothesis: .....	5
1.2.2 Problem a. ....	6
1.2.3 Problem b. ....	9
1.2.4 Problem c. ....	12
1.2.5 Problem d. ....	13
1.2.6 Problem e. ....	15
1.2.7 Problem f. ....	19
2.1 Topic .....	20
2.2 Answer .....	21
2.2.1 Problem a .....	21
2.2.2 Problem b .....	22
2.2.3 Problem c .....	24

## DANH MỤC KÍ HIỆU VÀ CHỮ VIẾT TẮT

### CÁC KÝ HIỆU

$f$  Tần số của dòng điện và điện áp (Hz)

$p$  Mật độ điện tích khối (C/m<sup>3</sup>)

### CÁC CHỮ VIẾT TẮT

EDA Exploratory data analysis



## DANH MỤC CÁC BẢNG BIỂU, HÌNH VẼ, ĐỒ THỊ

### DANH MỤC HÌNH

<i>Hình 1. 1 Filter [Avg.Price]</i> .....	6
<i>Hình 1. 2 Set [best-selling scales N]</i> .....	7
<i>Hình 1. 3 Top 10 best selling.</i> .....	7
<i>Hình 1. 4 Top 1000 best selling.</i> .....	8
<i>Hình 1. 5 Top 10000 best selling.</i> .....	9
<i>Hình 1. 6 Set [worst-selling scales N]</i> .....	9
<i>Hình 1. 7 Top 10 worst-selling</i> .....	10
<i>Hình 1. 8 Top 100 worst-selling</i> .....	10
<i>Hình 1. 9 Top 1000 worst-selling</i> .....	11
<i>Hình 1. 10 Top 1000 worst-selling</i> .....	11
<i>Hình 1. 11 Relationships between variable in dataset.</i> .....	13
<i>Hình 1. 12 Condition of set 1D.</i> .....	14
<i>Hình 1. 13 Top of set 1D.</i> .....	14
<i>Hình 1. 14 Top 10 product which are most expensive and ratings lower than 3.0</i> .....	14
<i>Hình 1. 15 Top 100 product which are most expensive and ratings lower than 3.0</i> .....	15
<i>Hình 1. 16 Top 1000 product which are most expensive and ratings lower than 3.0</i> .....	15
<i>Hình 1. 17 Condition of Set 1E</i> .....	16
<i>Hình 1. 18 Top of 1E</i> .....	17
<i>Hình 1. 19 Top 10 the cheapest product with ratings more than 4.0</i> .....	17
<i>Hình 1. 20 Top 100 the cheapest product with ratings more than 4.0</i> .....	18
<i>Hình 1. 21 Top 1000 the cheapest product with ratings more than 4.0</i> .....	18
<i>Hình 1. 22 Top of set 1F</i> .....	19
<i>Hình 1. 23 Products with the greatest preferential prices.</i> .....	19
<i>Hình 1. 24 The condition of promotion.</i> .....	23
<i>Hình 1. 25 Top products should be appear in a promotion.</i> .....	24

<i>Hình 1. 26 Top n worst profit city .....</i>	<i>24</i>
<i>Hình 1. 27 Top n best profit city .....</i>	<i>25</i>
<i>Hình 1. 28 Top best-profit and worst-profit cities in scales N.....</i>	<i>25</i>

## **DANH MỤC BẢNG**

<i>Table 1. Dataset variables 1 .....</i>	<i>12</i>
<i>Table 2. Dataset variables 2. ....</i>	<i>21</i>

## **QUESTION 1 – EXPLORATORY DATA ANALYSIS (EDA)**

### **1.1 Topic**

Exploratory data analysis (EDA) [3] is used to analyze and investigate datasets and summarize their main characteristics, often employing data visualization methods. It helps determine how best to manipulate data sources to get the answers you need, making it easier for data scientists to discover patterns, spot anomalies, test a hypothesis, or check assumptions. EDA is primarily used to see what data can reveal beyond the formal modeling or hypothesis testing task and provides a better understanding of data set variables and the relationships between them. Based on the idea of EDA, using Tableau software to answer the following questions:

- a. Find best-selling products in the scale of 10, 100, 1K, 10K items, and visualize them.
- b. Find worst-selling products in the scale of 10, 100, 1K, 10K items, and visualize them.
- c. Find and explain the relationships between dataset variables.
- d. Find products which are most expensive but have ratings lower than 3.0, in the scale of 10, 100, 1K items.
- e. Find products which are cheapest and have ratings more than 4.0, in the scale of 10, 100, 1K items.
- f. Find products which have the largest discounted price.

### **1.2 Answer**

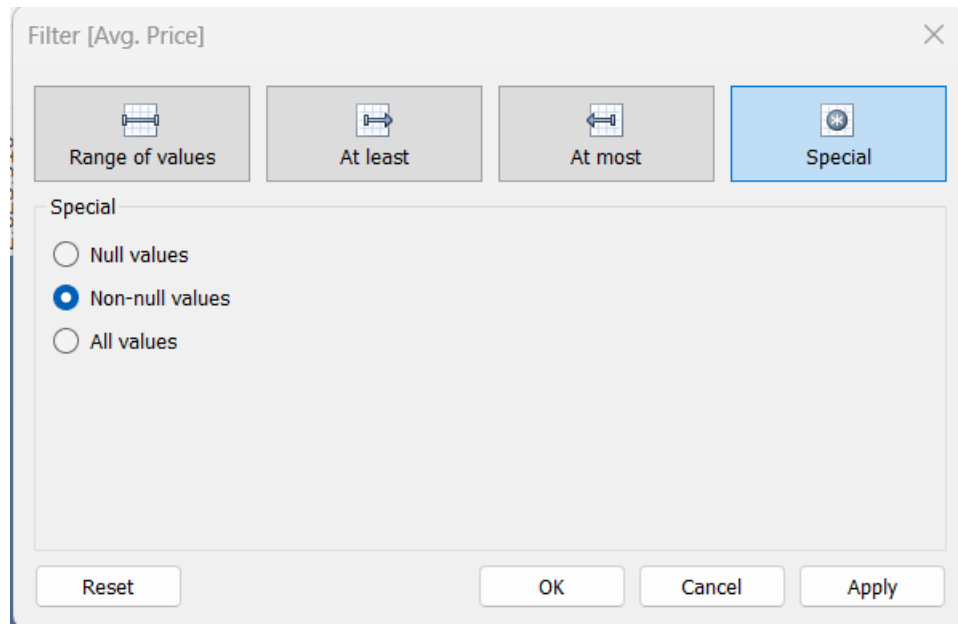
This chapter presents the reason for choosing the topic, purpose, object and scope of research, scientific and practical significance of the topic; Scientific basis for choosing topics...;

#### ***1.2.1 Hypothesis:***

According to my group's hypothesis:

- Because there are products that appear many times in the data set, we should be objective. We'll use the average price of those occurrences as that product's price and remove the currency character for ease of calculation, then remove the null value for that reason are products that have not yet been brought to market.

$$\text{Price} = \text{INT}(\text{RIGHT}([\text{Actual Price}], \text{LEN}([\text{Actual Price}]) - 1))$$



*Hình 1. 1 Filter [Avg.Price]*

- Because we do not have a representative value for purchases, we use the numerical value of ratings as a substitute. After making a purchase, people can provide feedback, so the number of purchases and the number of ratings are proportional to each other.

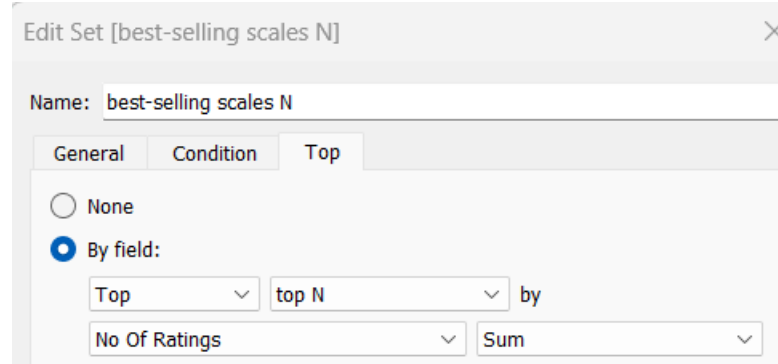
Our team hypothesizes that the 'Discount Price' column represents the price of the product after a discount. Therefore, we have:

$$\text{Discount amount} = \text{AVG}([\text{Price}]) - [\text{Discounted Price (AVG)}]$$

$$\text{Discount}(\%) = ((\text{AVG}([\text{Price}]) - [\text{Discounted Price (AVG)}]) * 100) / \text{AVG}([\text{Price}])$$

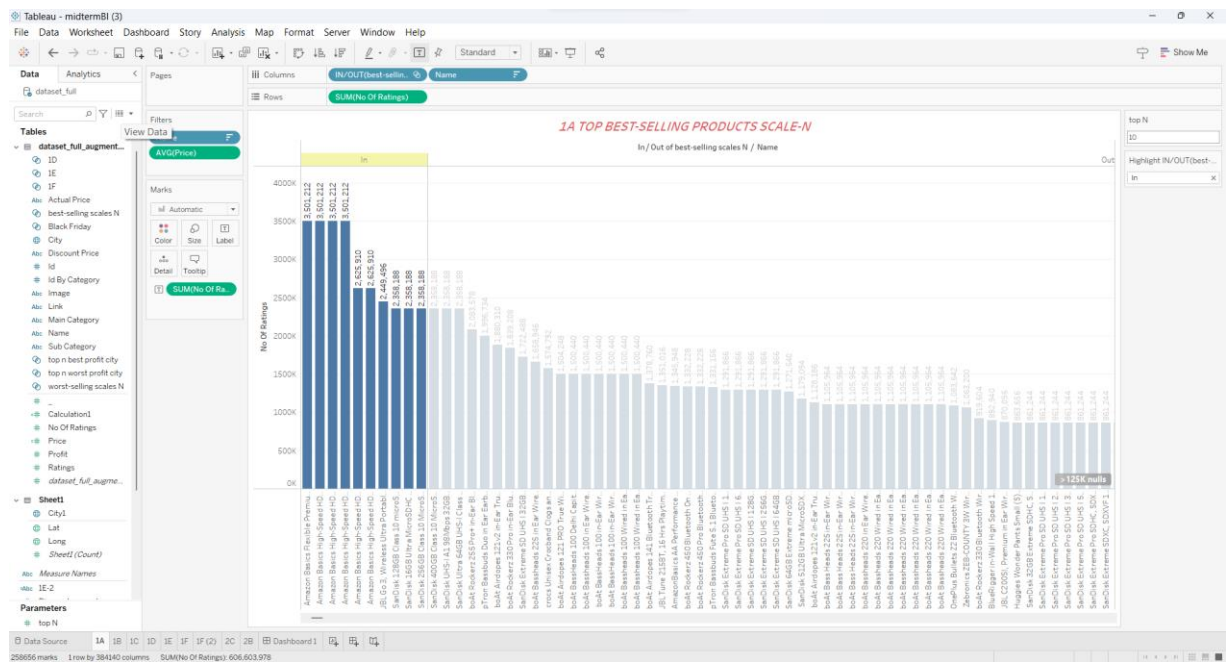
### **1.2.2 Problem a.**

We hypothesize that the products with the highest number of purchases are also the ones with the highest number of reviews:

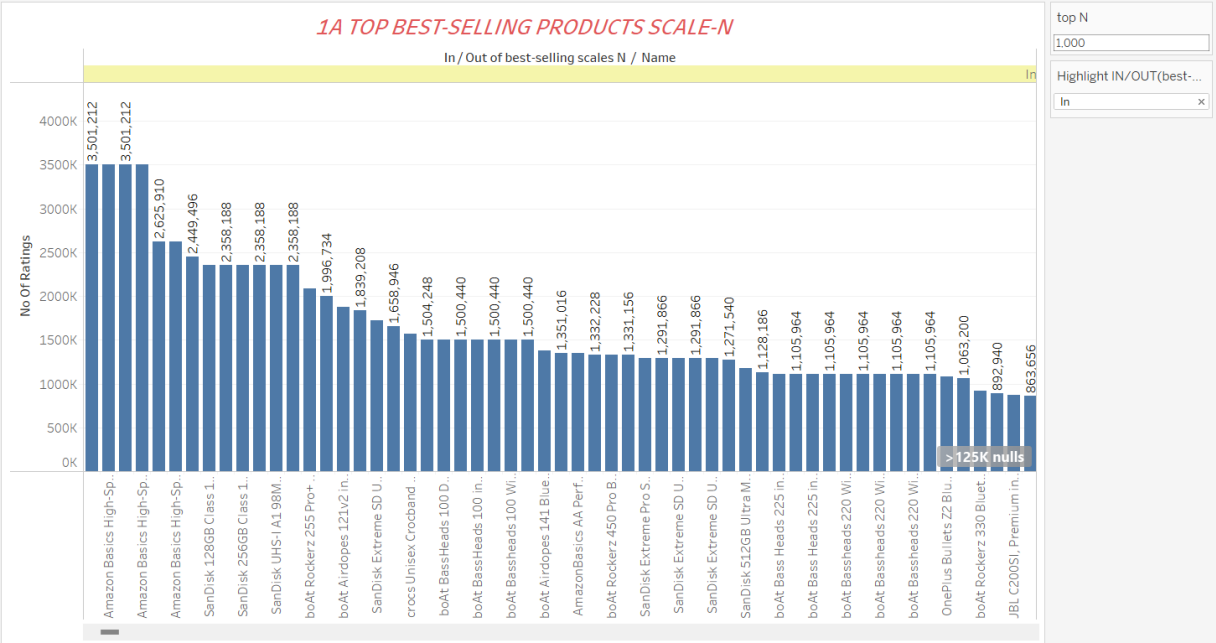


Hình 1. 2 Set [best-selling scales N]

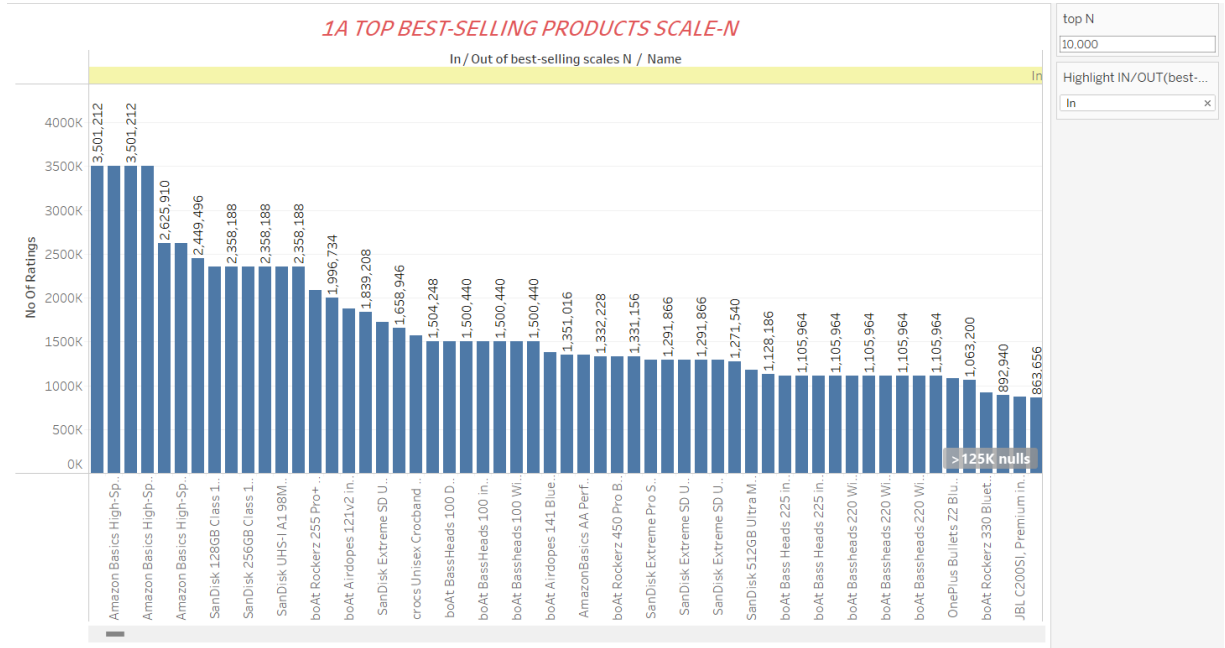
The following will be a presentation of data visualization using Tableau. Here, we will create a set to be able to count the number of best-selling products in the top 10, 100, 1000, 10000



Hình 1. 3 Top 10 best selling.



Hình 1. 4 Top 1000 best selling.



Hình 1. 5 Top 1000 best selling.

### 1.2.3 Problem b.

We hypothesize that the products with the lowest number of purchases are also the ones with the lowest number of reviews.

Edit Set [worst-selling scales N]

Name: worst-selling scales N

General Condition Top

☐ None

☒ By field:

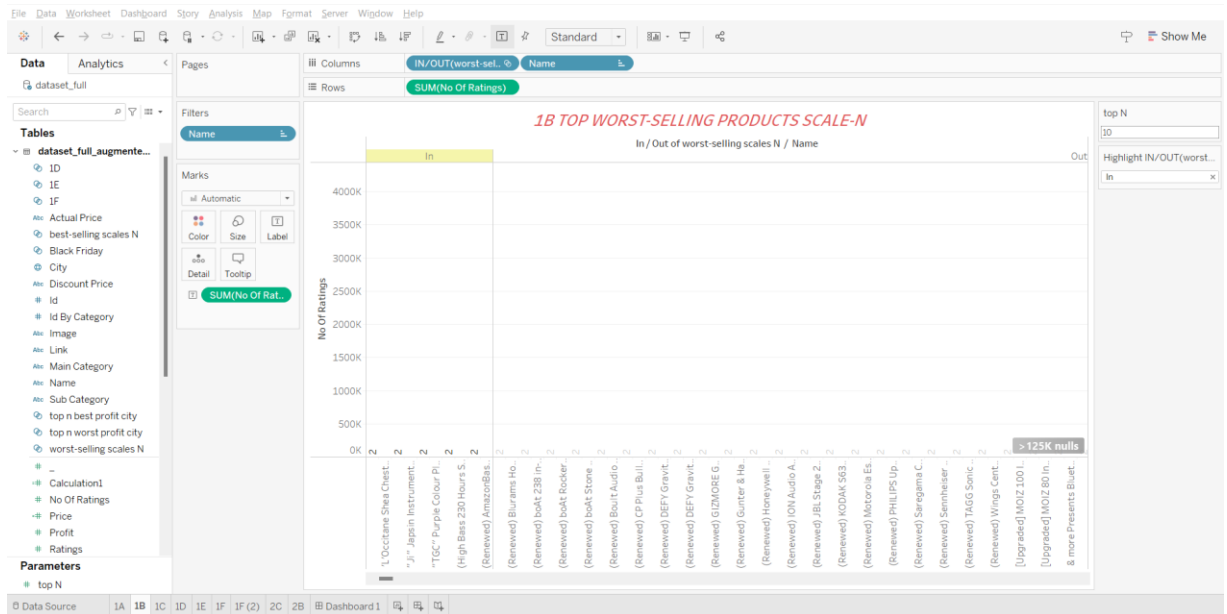
Bottom top N by

No Of Ratings Sum

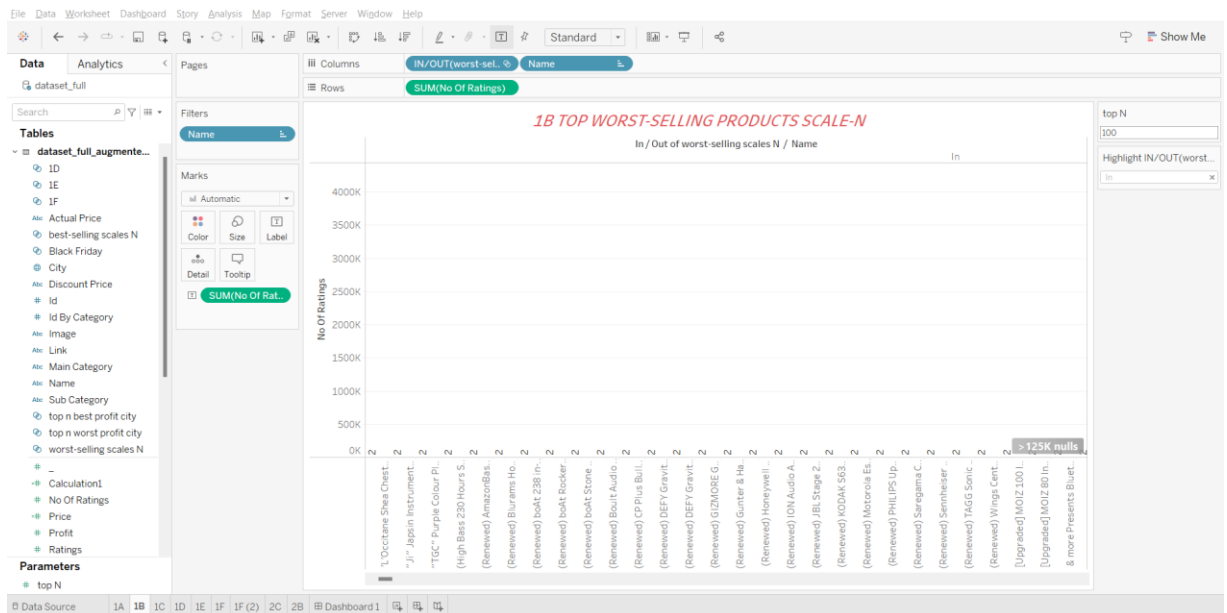
Hình 1. 6 Set [worst-selling scales N]

The following will be a presentation of data visualization using Tableau. Here, we will create a set to be able to count the number of worst-selling products in the top 10, 100,

1000,10000.

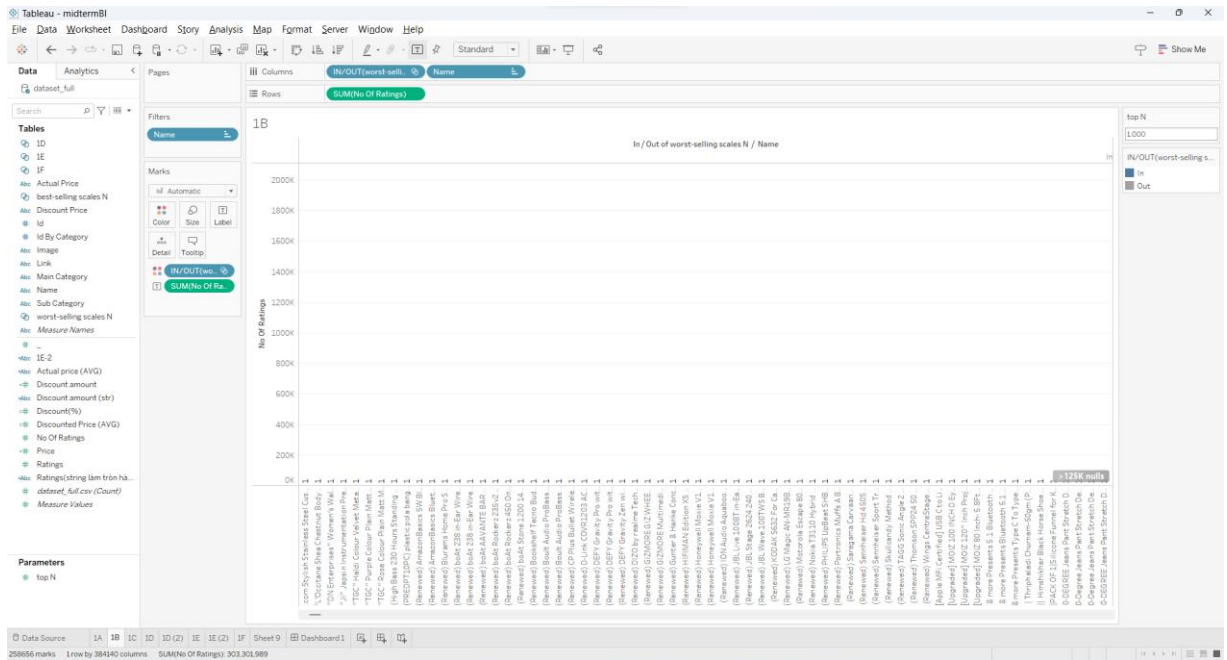


Hình 1. 7 Top 10 worst-selling

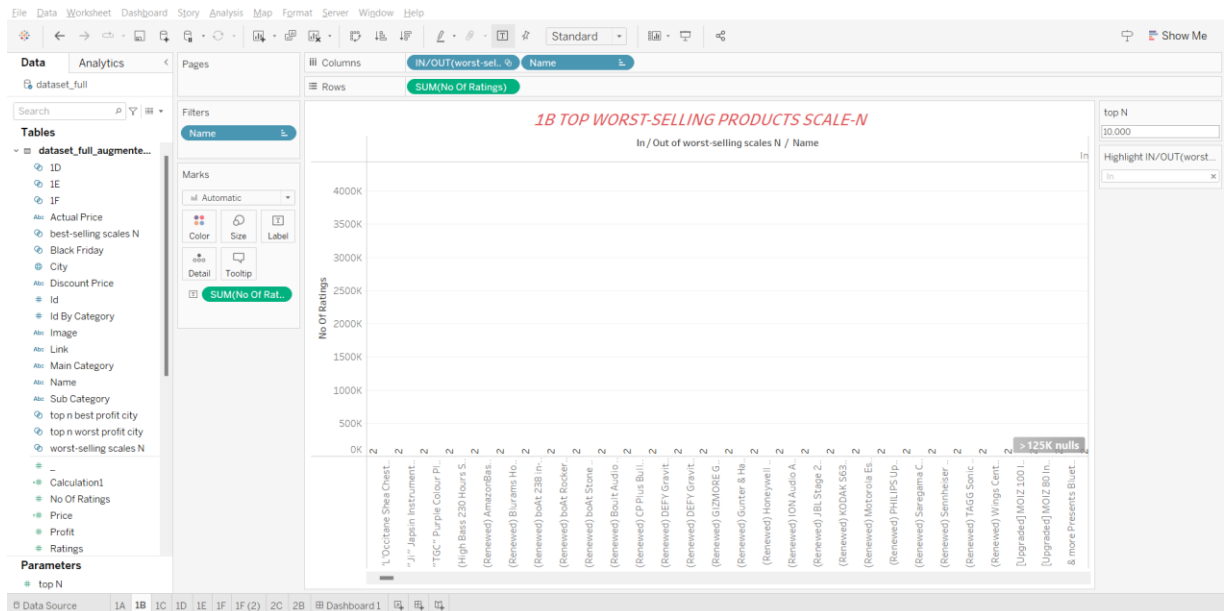


Hình 1. 8 Top 100 worst-selling.





Hình 1. 9 Top 1000 worst-selling.



Hình 1. 10 Top 1000 worst-selling.

### 1.2.4 Problem c.

Variables

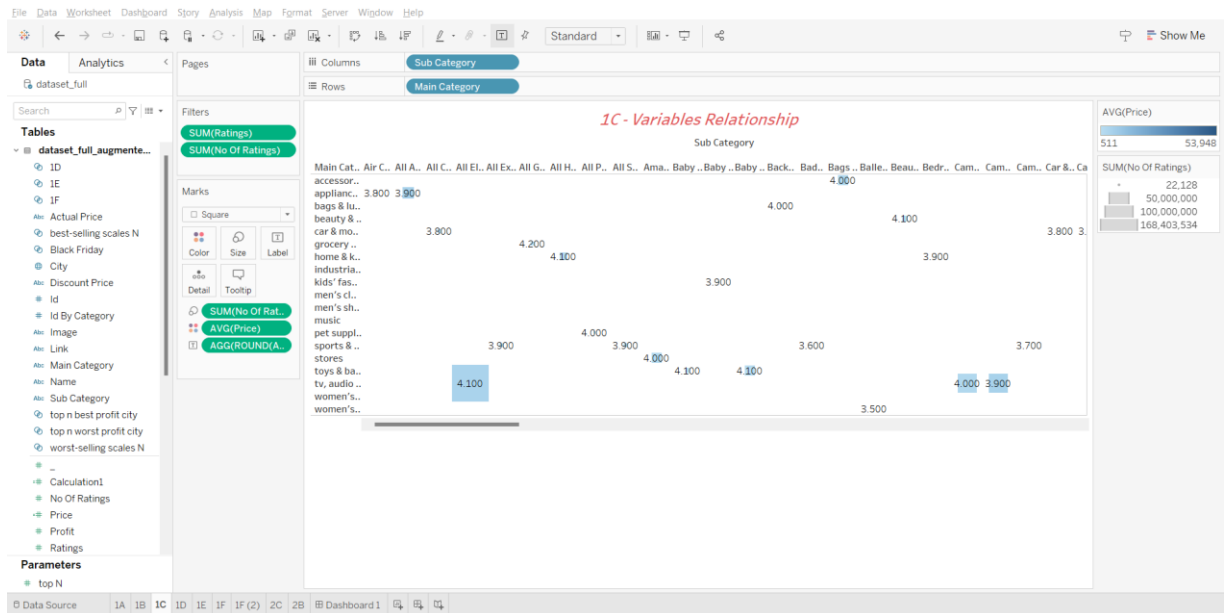
Attribute	Variable	Description
Name	numerical variable	The name of the product
Main_category	categorical variable	The main category of the product belong
Sub_category	categorical variable	The sub category of the product belong
Image	No numerical variable	
Link	No numerical variable	
ratings	numerical variable	Ratings score
No of ratings	numerical variable	the number of ratings
Discount_price	numerical variable	the price after discounted
Actual_price	numerical variable	Origin price

*Table 1. Dataset variables 1*

Find and explain relationship

1. Name (numerical variable) and Main\_category (categorical variable):  
The variable "Name" represents the name of the product, while "Main\_category" indicates the main category to which the product belongs. This relationship allows us to understand the association between the product's name and its main category. Each product has a unique name, and it is categorized into a specific main category.
2. Main\_category (categorical variable) and Sub\_category (categorical variable):  
The variable "Main\_category" represents the main category of the product, while "Sub\_category" denotes the sub-category to which the product belongs. This relationship helps us understand the hierarchical structure of the categories. Each product falls under a main category and can further be classified into a specific sub-category within that main category.
3. Ratings (numerical variable) and No of ratings (numerical variable):  
The variable "Ratings" represents the score or rating given to the product, and "No of ratings" indicates the number of ratings received by the product. This relationship enables us to analyze the correlation between the rating score and the number of ratings. Generally, higher-rated products tend to attract more ratings, indicating their popularity or customer satisfaction.
4. Discount\_price (numerical variable) and Actual\_price (numerical variable):  
The variable "Discount\_price" represents the price of the product after applying a discount, while "Actual\_price" indicates the original price of the product. This relationship helps us understand the impact of discounts on the product's pricing. The

discount price is derived from the original price by reducing it based on the discount offered.



Hình 1. 11 Relationships between variable in dataset.

### 1.2.5 Problem d.

Edit Set [1D]

Name: 1D

General Condition Top

☐ None

☒ By field:

Ratings Average

< 3

Range of Values

Min:  Load

Max:

Hình 1. 12 Condition of set 1D.

Edit Set [1D]

Name: 1D

General Condition Top

☐ None

☒ By field:

Top top N by

Price Average

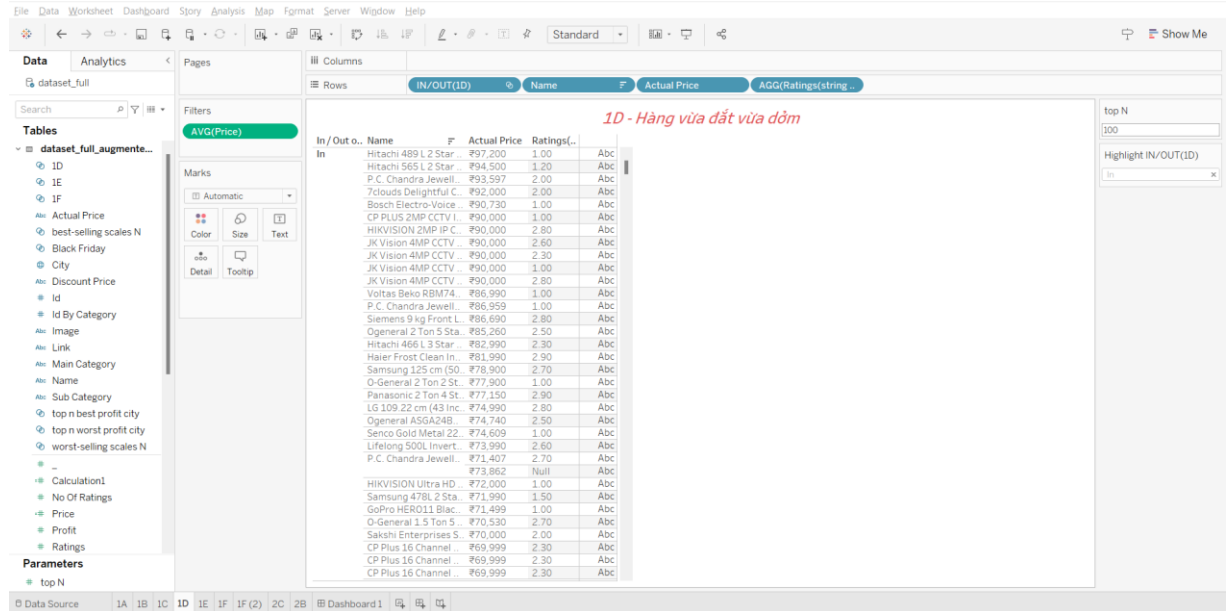
Hình 1. 13 Top of set 1D.

To find products that are most expensive and have a rating below 3.0. Use the Name and Ratings columns from the given data then create a set to set the conditions and list the scale of items.

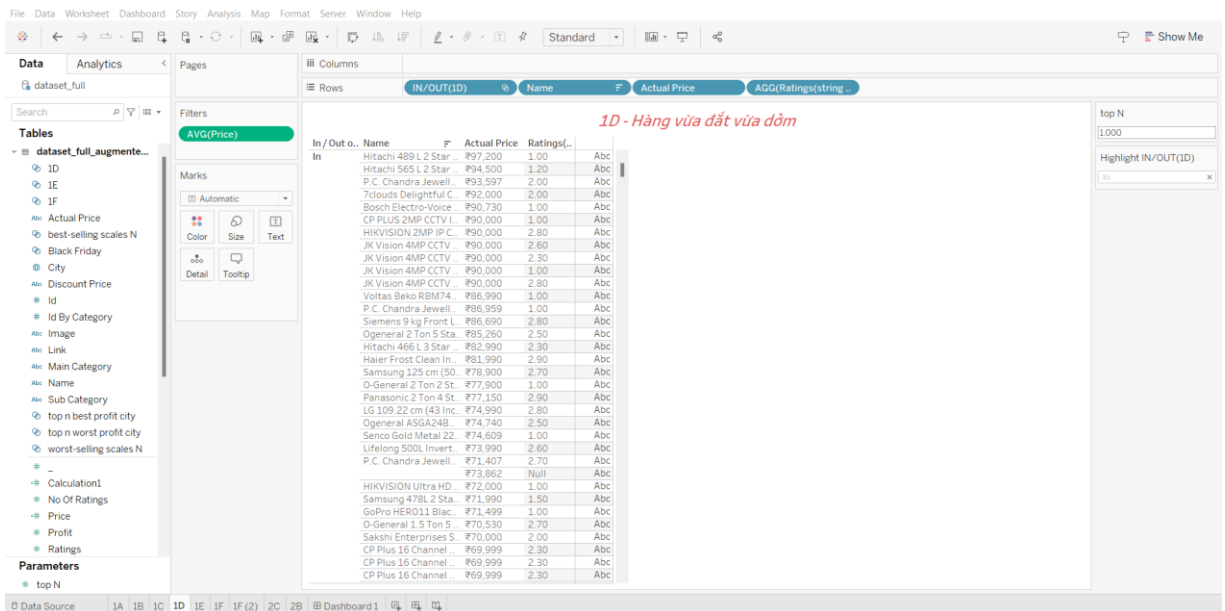
1D - Hàng vừa đắt vừa dởm

In / Out o.	Name	P	Actual Pr.	Ratings	AGG(Ratings(string...))
In	Hitachi 489 L 2 Star ..	¥97,200	1.00	Abc	
	Hitachi 565 L 2 Star ..	¥94,500	1.20	Abc	
	P. C. Chandra Jewell ..	¥93,597	2.00	Abc	
	7clouds Delightful C. ..	¥92,000	2.00	Abc	
	Bosch Electro-Voice ..	¥90,730	1.00	Abc	
	CP PLUS 2MP CCTV I. ..	¥90,000	1.00	Abc	
	WIKIVISION 2MP IP C. ..	¥90,000	2.80	Abc	
	JK Vision 4MP CCTV ..	¥90,000	2.60	Abc	
	JK Vision 4MP CCTV ..	¥90,000	2.30	Abc	
	JK Vision 4MP CCTV ..	¥90,000	1.00	Abc	
Out	Basics 670 L French ..	¥99,999	Null	Abc	
	coccaa 164 cm (65 i. ..	¥99,999	4.20	Abc	
	Denon AVR-S760H 7 ..	¥99,999	4.50	Abc	
	ESC Medicams Endo. ..	¥99,999	Null	Abc	
	Fujifilm X-T30 II Bod. ..	¥99,999	4.50	Abc	
	Google Pixel 6 Pro 5. ..	¥99,999	3.80	Abc	
	Hometronics Mastic. ..	¥99,999	5.00	Abc	
	JBL Bar 800 Pro, 7.1 ..	¥99,999	4.10	Abc	
	LETSPRAY® LP-100A. ..	¥99,999	4.10	Abc	
	MEKE Meike 85mm ..	¥99,999	4.20	Abc	
	Sonos Arc - The Pre. ..	¥99,999	4.70	Abc	
	Sonos Arc: Arcg1Ux1. ..	¥99,999	4.70	Abc	
	Acer 178 cm (70 inch. ..	¥99,990	4.40	Abc	
	BenQ TH575 4K Com. ..	¥99,990	4.50	Abc	
	ELECTROLUX 453L F. ..	¥99,990	Null	Abc	
	Havells-Lloyd 146cm. ..	¥99,990	3.00	Abc	
	Hisense 164 cm (65 i. ..	¥99,990	4.30	Abc	
	Power Guard 165 cm. ..	¥99,990	4.10	Abc	
	Sony SRS-XV900 X-S. ..	¥99,990	4.40	Abc	
	usha misty 25ltr. st. ..	¥99,990	3.70	Abc	
	HP 245 G8 Laptop P. ..	¥99,977	Null	Abc	
	PC Jeweller The Frig. ..	¥99,905	Null	Abc	
	Apple 2020 MacBoo. ..	¥99,900	4.70	Abc	
	Apple iPhone 14 Plu. ..	¥99,900	4.50	Abc	

Hình 1. 14 Top 10 product which are most expensive and ratings lower than 3.0



Hình 1. 15 Top 100 product which are most expensive and ratings lower than 3.0



Hình 1. 16 Top 1000 product which are most expensive and ratings lower than 3.0

### 1.2.6 Problem e.

Edit Set [1E] X

Name: 1E

General Condition Top

☐ None

☐ By field:

Ratings(string làm tròn hàng đơn vị) Custom

=

Range of Values

Min: Load

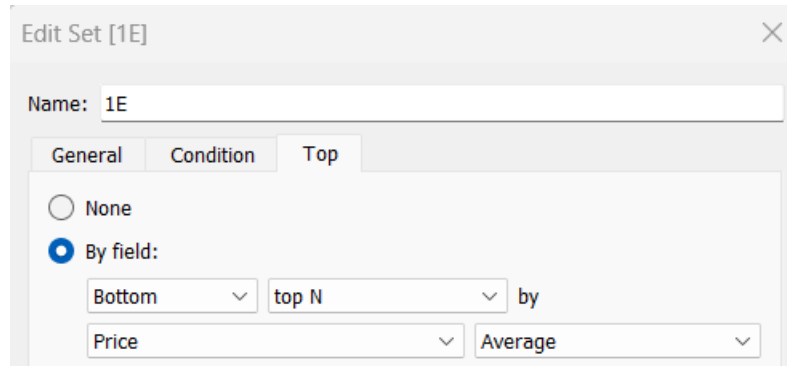
Max:

☒ By formula:

```
AVG([Price]) >= 0
AND
AVG([Ratings]) > 4
```

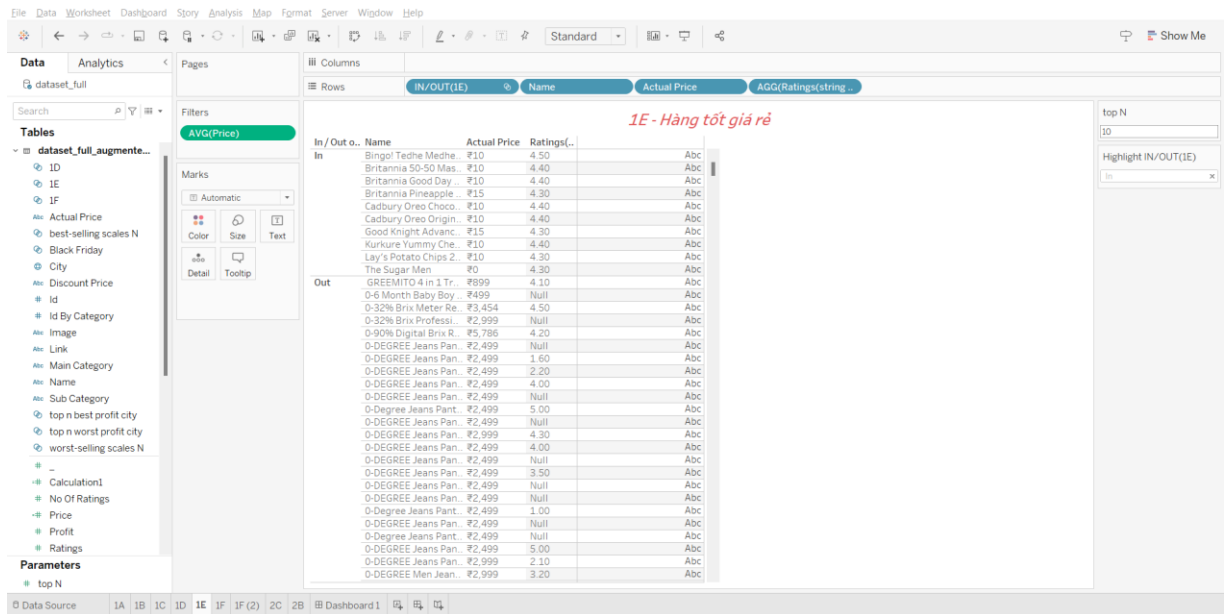
Reset OK Cancel Apply

Hình 1. 17 Condition of Set 1E



Hình 1. 18 Top of 1E

The cheapest products have a rating above 4.0. Use 2 data columns: Ratings and Name for processing.



Hình 1. 19 Top 10 the cheapest product with ratings more than 4.0

Tableau Desktop interface showing a table of products. The table is filtered by AVG(Price) and sorted by top N (100). The table is titled "1E - Hàng tốt giá rẻ".

IN / Out o..	Name	Actual Pr..	Ratings(.
In	Alpenliebe Pop. 64g	¥50	4.20
In	Amul Kool Koko Can...	¥40	4.20
In	APLUS Super Saver	¥36	4.20
In	Appara White Dustie	¥20	4.20
In	Aquafina Drinking W...	¥20	4.30
In	Bingöl Tedhe Medhe	¥10	4.50
In	Britannia 50-50 Mas...	¥10	4.40
In	Britannia Good Day	¥10	4.40
In	Britannia Milk Bikis	¥25	4.30
In	Britannia Pineapple	¥15	4.30
In	Britannia Toastea Pr...	¥50	4.30
In	Britannia Toastea Pr...	¥40	4.10
In	Britannia Vita Marie	¥30	4.30
In	Brooke Bond, Red La...	¥30	4.30
In	Cadbury Dairy Milk	¥30	4.40
In	Cadbury Oreo Choco	¥10	4.40
In	Cadbury Oreo Origin	¥10	4.40
In	Cadbury Oreo Vanill	¥35	4.40
In	Camin Kokuyo Fine	¥30	4.30
In	Cavin's Vanilla Milk	¥40	4.30
In	Centre Fresh Mint C...	¥50	4.30
In	Ching's Dark Soy Sa...	¥25	4.20
In	Classmate Octane N...	¥50	4.30
In	Classmate Octane-B	¥50	4.30
In	Coca-Cola Zero Suga...	¥40	4.30
In	Cornitos Nachos Cri...	¥35	4.40
In	Cup Noodles Mazed	¥50	4.20
In	Cup Noodles Veggi	¥50	4.10
In	Del Monte Tomato K...	¥40	4.20
In	Dettol Disinfectant	¥49	4.30
In	Dettol Liquid Handw...	¥40	4.30
In	Dettol Original Ger...	¥25	4.50
In	Dukes Bourbon Pre...	¥36	4.30
In	Dukes Waffly - Choco	¥50	4.20

Hình 1. 20 Top 100 the cheapest product with ratings more than 4.0

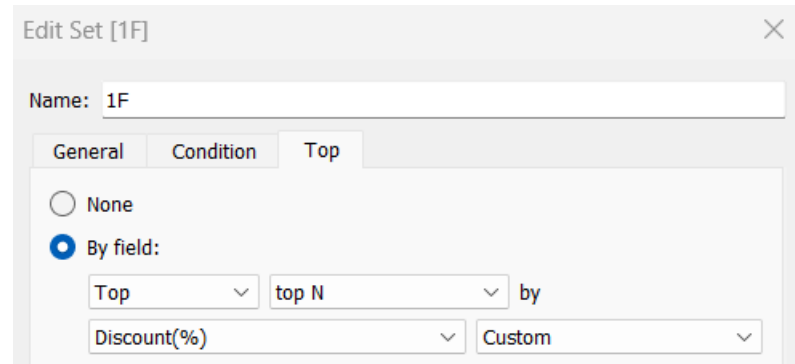
Tableau Desktop interface showing a table of products. The table is filtered by AVG(Price) and sorted by top N (1,000). The table is titled "1E - Hàng tốt giá rẻ".

IN / Out o..	Name	Actual Pr..	Ratings(.
In	1K ohm 1/4 Watt Re...	¥99	5.00
In	1Pc Cosmetic Puff P...	¥149	4.50
In	1st Step Baby Wet...	¥99	4.30
In	2 Spool x 300m Orga...	¥100	4.20
In	2K ohm Variable Res...	¥99	4.10
In	3M 1270 Ear Plugs S...	¥128	4.30
In	3M Command, Large	¥149	4.20
In	3M Post-it Sticky No...	¥135	4.40
In	3M Scotch Kids Sciss...	¥149	4.40
In	3M Scotch Magic Ta...	¥140	4.30
In	3V PRODUCTS, Data	¥110	4.10
In	3V PRODUCTS, Jathi	¥125	5.00
In	4.7K ohm Variable R...	¥99	5.00
In	7 Up Soft Drink - 2.2	¥99	4.40
In	01 Super Shop Extra	¥149	4.40
In	22nF/630V (0.022uF	¥99	4.50
In	22uF/25V Electrolyt...	¥99	5.00
In	24 Mantra Organic 7	¥95	4.30
In	24 Mantra Organic B	¥55	4.30
In	24 Mantra Organic S	¥145	4.30
In	24 Mantra Organic	¥95	4.40
In	24 Mantra Organic	¥80	4.40
In	24 Mantra Organic	¥145	4.30
In	47nF/630V (0.047uF	¥99	5.00
In	1000uF/16V Electrol...	¥99	5.00
In	3386P- 2K Ohm 0.5	¥99	5.00
In	Amritamehari Chur...	¥140	4.50
In	Aswagandha Table	¥125	4.80
In	Eaneer Kuzhampou...	¥130	4.40
In	Gandharvahasthad...	¥125	5.00
In	Talisapatradi Chur...	¥100	4.60
In	A D Food & Herbs Or...	¥129	5.00
In	A.W. Faber-Castell I...	¥125	4.40
In	Aachi Coriander Pow...	¥129	4.20

Hình 1. 21 Top 1000 the cheapest product with ratings more than 4.0

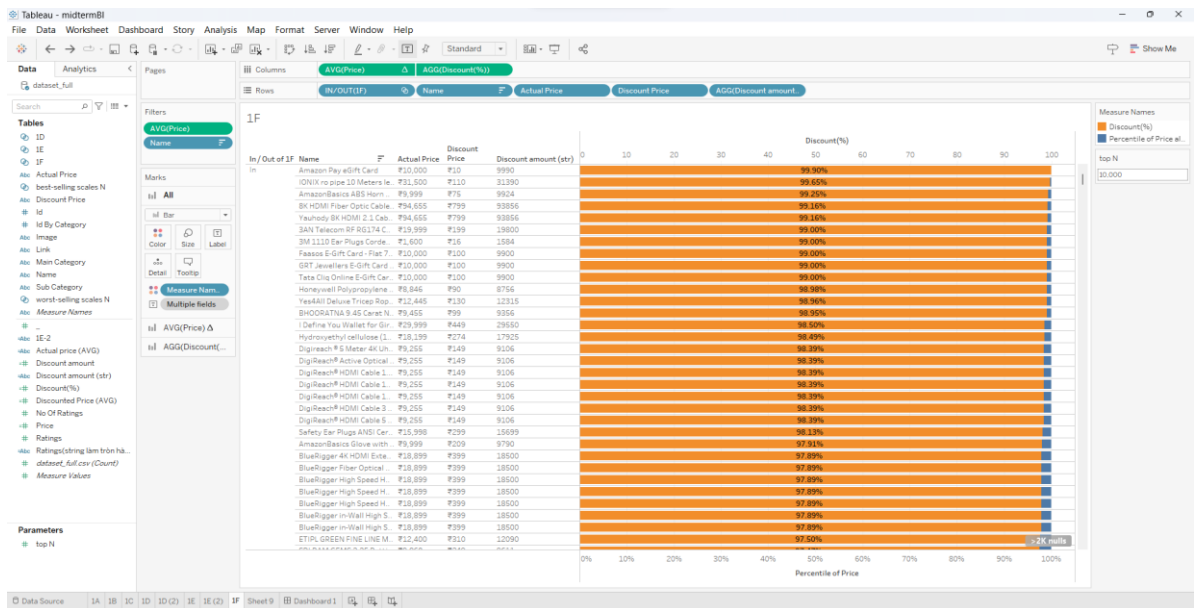


### 1.2.7 Problem f.



Hình 1. 22 Top of set 1F

Products with the greatest preferential prices.



Hình 1. 23 Products with the greatest preferential prices.

## QUESTION 2 – AUGMENTED THE DATASET

### 2.1 Topic

To employ more power of Tableau software, in this section, you need to augment the given dataset in both  $x$  and  $y$  axes. For  $x$  axis, you will insert two new columns, ie.. "profit" and "city", where "profit" indicates the number of products which were sold in the corresponding "city" of US. For  $y$  axis, it's quite easy to imagine, you need to duplicate each instance of the dataset to at least 100 new ones, and randomly generate the "profit" and "city" values for each of them. Table 2 illustrates the new dataset after augmenting in both  $x$  and  $y$  axes. You should note that the "id" and "id by category" columns need to be altered to make sure these identification numbers are unique.

Table 2. Illustrate a new dataset after augmenting $x$ and $y$ axes				
<i>id</i>	<i>name</i>	<i>...</i>	<i>profit</i>	<i>city</i>
0	Lloyd 1.5 Ton 3 Star Inverter Split	...	126	South Bend
1	Lloyd 1.5 Ton 3 Star Inverter Split	...	32	Springfield
2	Lloyd 1.5 Ton 3 Star Inverter Split	...	147	Savannah
3	Lloyd 1.5 Ton 3 Star Inverter Split	...	296	Knoxville
4	Lloyd 1.5 Ton 3 Star Inverter Split	...	13	Salt Lake City
5	Lloyd 1.5 Ton 3 Star Inverter Split	...	58	Cincinnati

The dataset after augmenting is now having enough information to support the decision-making activities. You need to answer the following questions:

- Find and explain the relationships between dataset variables
- Get a list of products that should be included in a promotion campaign based on the selling profits by city. For example, you can define a measure such as "if a product has a good selling profit at a city, it should be discounted as  $x$  percent to improve its profit", or "we should discount that jacket's price because it is going to be outdated soon". Whatever the measure you define, you need to define a threshold, for example, to determine whether a product has a good selling profit or not.
- Visualize the best-and-worst-selling-profit on world map.

## 2.2 Answer

### 2.2.1 Problem a

Variables

Attribute	Variable	Description
Name	numerical variable	The name of the product
Main_category	categorical variable	The main category of the product belong
Sub_category	categorical variable	The sub category of the product belong
Image	No numerical variable	
Link	No numerical variable	
ratings	numerical variable	Ratings score
No of ratings	numerical variable	the number of ratings
Discount_price	numerical variable	the price after discounted
Actual_price	numerical variable	Origin price
Profit	numerical variable	The profit of the product
City	numerical variable	The city where the product be sold

*Table 2. Dataset variables 2.*

Find and explain the relationship:

1. Profit and Main\_category: Different main categories may have varying profitability levels.
2. Profit and Sub\_category: Profitability can differ across sub-categories within a main category.
3. Profit and Ratings: Higher ratings may indicate increased profitability, but other factors also influence it.
4. Profit and No of ratings: More ratings may suggest higher profitability, considering other factors.
5. Profit and Discount\_price: Discounts can attract more customers, affecting profitability.
6. Profit and City: Profitability can vary across cities due to market factors and preferences.

7. City and Main\_category: Different cities may have varying preferences for specific product categories.
8. City and Sub\_category: Sub-categories may have different popularity levels across cities.
9. City and Ratings: Customer ratings may vary across different cities, indicating regional sentiment.
10. City and No of ratings: The number of ratings may differ based on city population and market activity.

### ***2.2.2 Problem b***

Promotion details:

- Name: Black Friday.
- Condition: The products with ratings higher than 3.5, a discount rate equal to or less than 30%, and profits exceeding the average profit of the city where they are sold.
- Describe: Discounting popular and high-profit products in cities while ensuring that these products have not been discounted by more than 30% of their original value.

Edit Set [Black Friday] ✕

Name: Black Friday

General Condition Top

☐ None

☐ By field:

1E-2 Custom

=

Range of Values

Min: Load

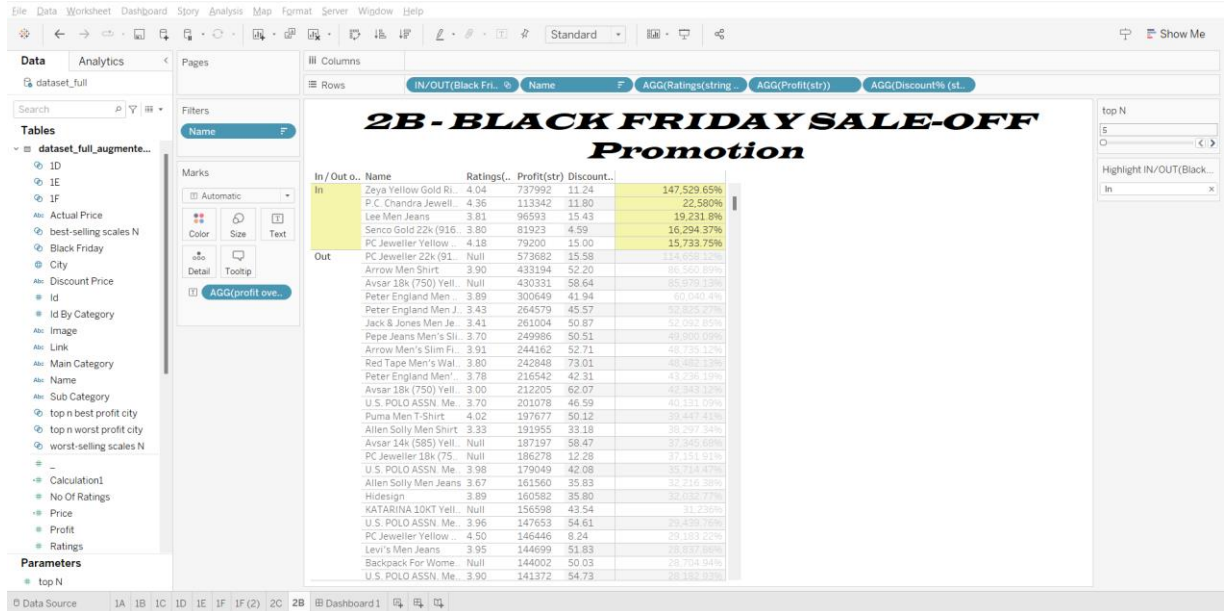
Max:

☒ By formula:

```
(AVG([Ratings]) >= 3.5)
AND
[Discount(%)] <= 30 AND [Discount(%)] > 0
AND
AVG([Profit]) >= AVG({FIXED [City] : AVG(
AND
AVG([Price]) >= 0
```

Reset OK Cancel Apply

Hình 1. 24 The condition of promotion.



Hình 1. 25 Top products should be appear in a promotion.

### 2.2.3 Problem c

The screenshot shows the "Edit Set [top n worst profit city]" dialog box. The "Name:" field is set to "top n worst profit city". The "General" tab is selected, and the "By field:" radio button is chosen. The "Bottom" dropdown is set to "top N" and the "Sum" dropdown is set to "Sum".

Hình 1. 26 Top n worst profit city

Edit Set [top n best profit city] ✕

Name:

General Condition Top

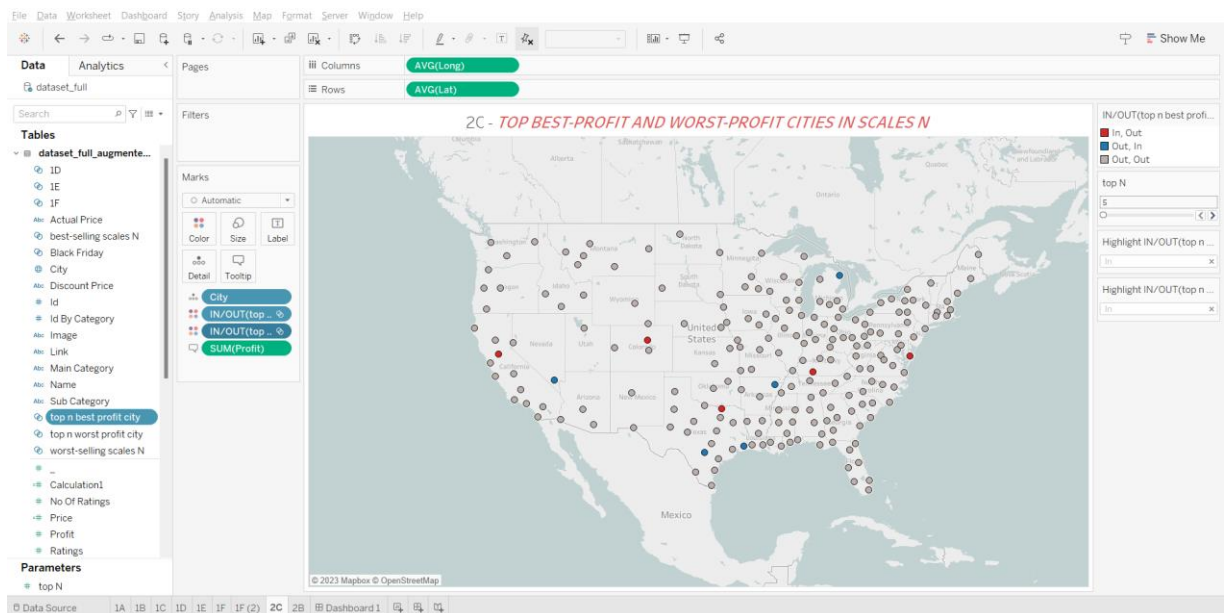
☐ None

☒ By field:

Top  by

Profit

Hình 1. 27 Top n best profit city



Hình 1. 28 Top best-profit and worst-profit cities in scales N

## TÀI LIỆU THAM KHẢO

### Tiếng Anh

- [1] Midterm Project Description - <https://cs504049.fastai.dev/midterm.html>
- [2] pandas.read\_csv-[https://pandas.pydata.org/docs/reference/api/pandas.read\\_csv.html](https://pandas.pydata.org/docs/reference/api/pandas.read_csv.html)
- [3] Exploratory data analysis - [https://en.wikipedia.org/wiki/Exploratory\\_data\\_analysis](https://en.wikipedia.org/wiki/Exploratory_data_analysis)