

BDA Project

Amazon Reviews (2022-23)

Presentation | 27.05.2025

UC | Big Data Analytics

Professor | Márcia L. Baptista (NOVA IMS)

Group W

André Silvestre	20240502
Filipa Pereira	20240509
João Henriques	20240499
Umeima Mahomed	20240543

TP1 & P2

DATA COLLECTION & PREPROCESSING

Collection

Filtering the data collected from *Amazon Reviews's 23 dataset* collected in 2023 by *McAuley Lab*



ETL

Data Pre-Processing

- Deduplicate
- Convert to the correct data types
- Clear timestamp



EDA

Exploratory Analysis & Verify Collection Quality

Visualize and Analyse the data obtained



ANALYSIS | TEXT MINING

Text Mining

- Tokenization
- Stop Words Removal
- N-Grams



Sentiment Analysis

Transformers

Text Sentiment Classification

- *citizenlab/twitter-xlm-roberta-base-sentiment-finetuned*
- *MoritzLaurer/mDeBERTa-v3-base-xnli-multilingual-nli-2mil7*



Topic Analysis

Transformers

Classification of Text Topics

- *MoritzLaurer/mDeBERTa-v3-base-xnli-multilingual-nli-2mil7*



+ EXTRA

Clustering

K-Means



Graph Analysis

Graph Analysis User-Product



Streaming Simulation

Batch and Advanced Streaming Simulation in

databricks





Figure 1 | Amazon box.
Source: Google Images



Figure 2 | Amazon headquarters
Source : SustainableJapan

Objective

Modelling Amazon Tech Reviews for Consumer Insight

Analysis of Amazon Electronics and Computers reviews (2022-2023) using Sentiment and Topic Modelling, Clustering, and Graph Analysis to uncover Consumer Purchasing Behaviour.

Project Management

Control Instruments

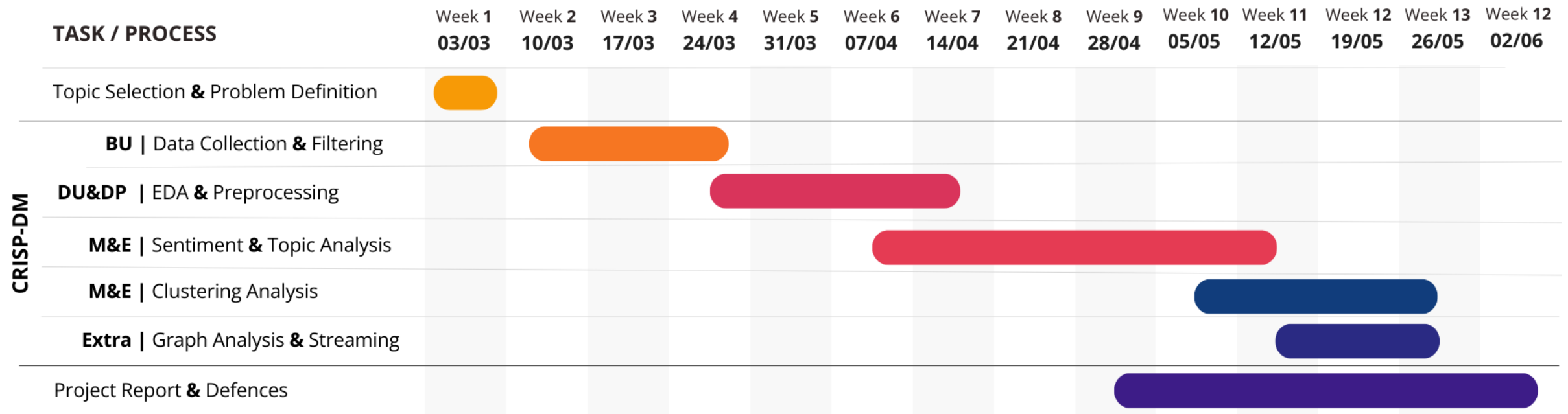
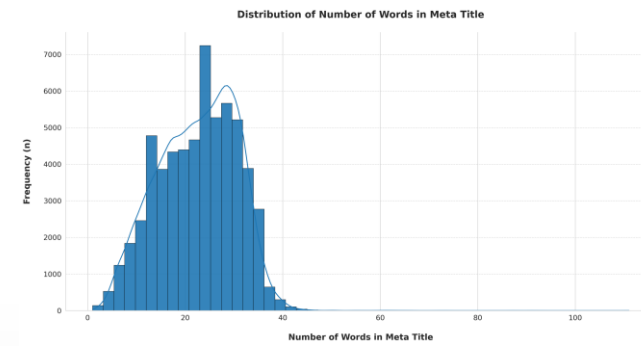
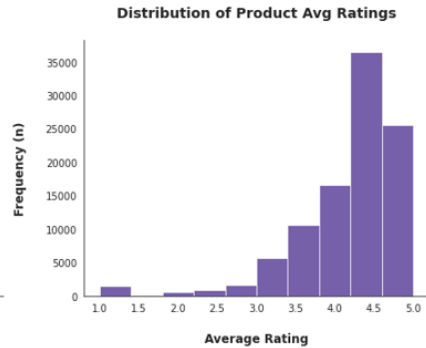


Figure 1 | Gantt Chart with Task Distribution.

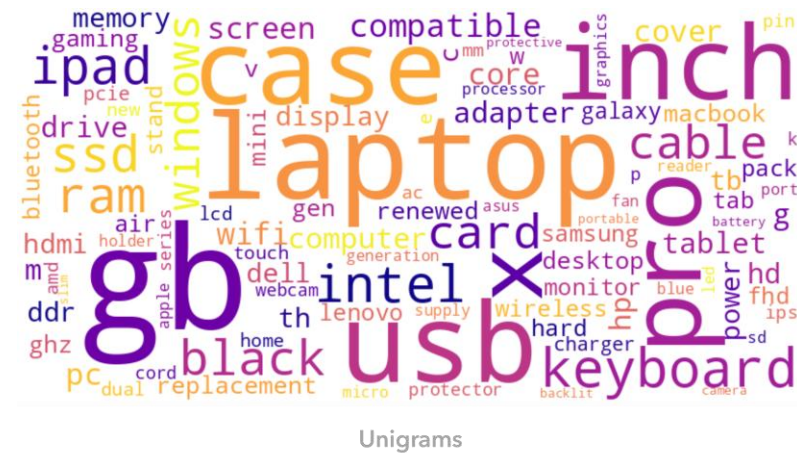


Products



	n	%
Store		
Amazon Renewed	6402	6.4
HP	3699	3.7
Lenovo	2040	2.0
ASUS	1886	1.9
Dell	1444	1.4
Generic	948	1.0
SAMSUNG	830	0.8
MOSISO	744	0.7
SanDisk	725	0.7
Logitech	473	0.5

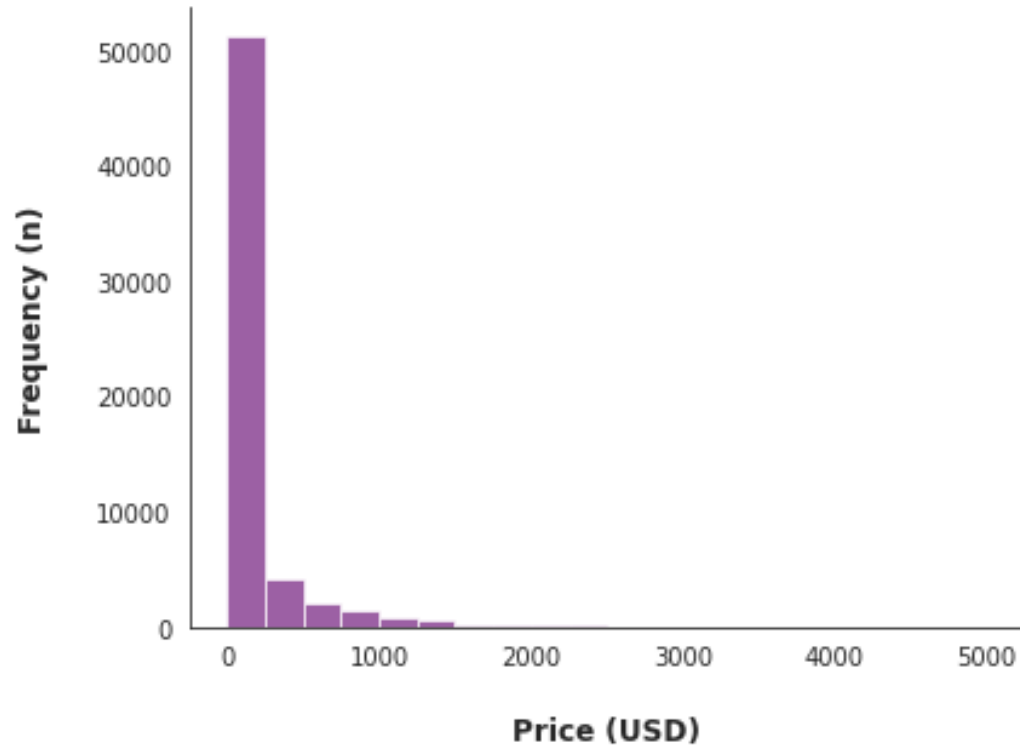
Price Category	n	%	Cumulative %
0 - 10	8774	8.8	8.8
10 - 50	28197	28.3	37.1
50 - 100	6958	7.0	44.1
100 - 250	7332	7.4	51.5
250 - 500	4205	4.2	55.7
500+	5612	5.6	61.3
Unknown	38557	38.7	100.0



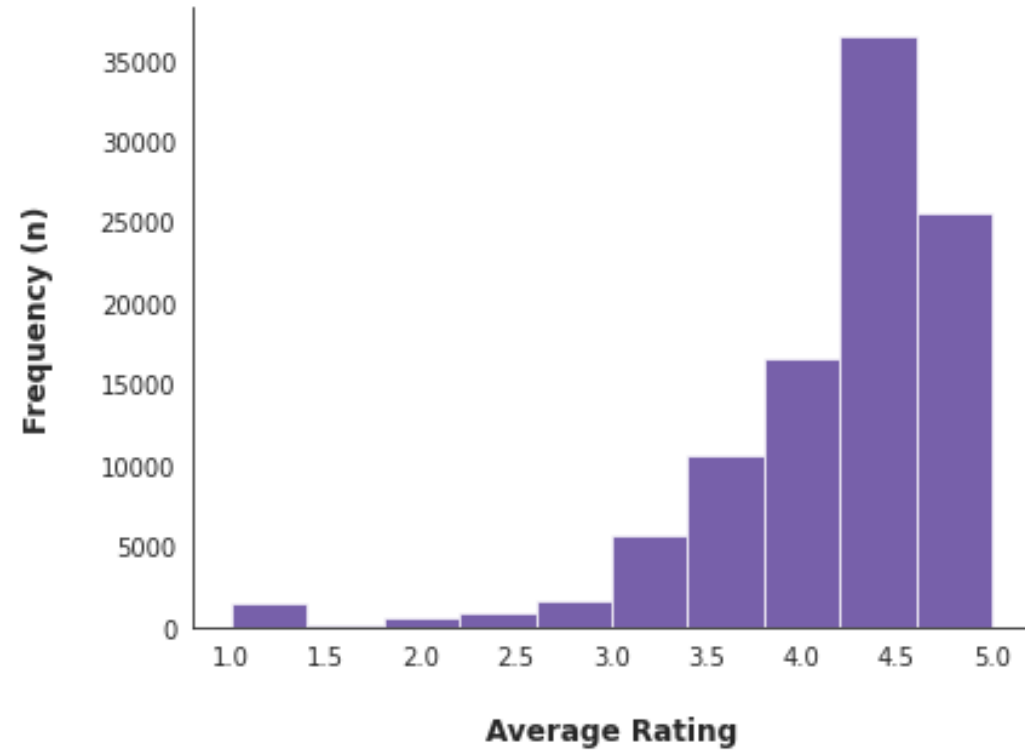


Products

Distribution of Product Prices

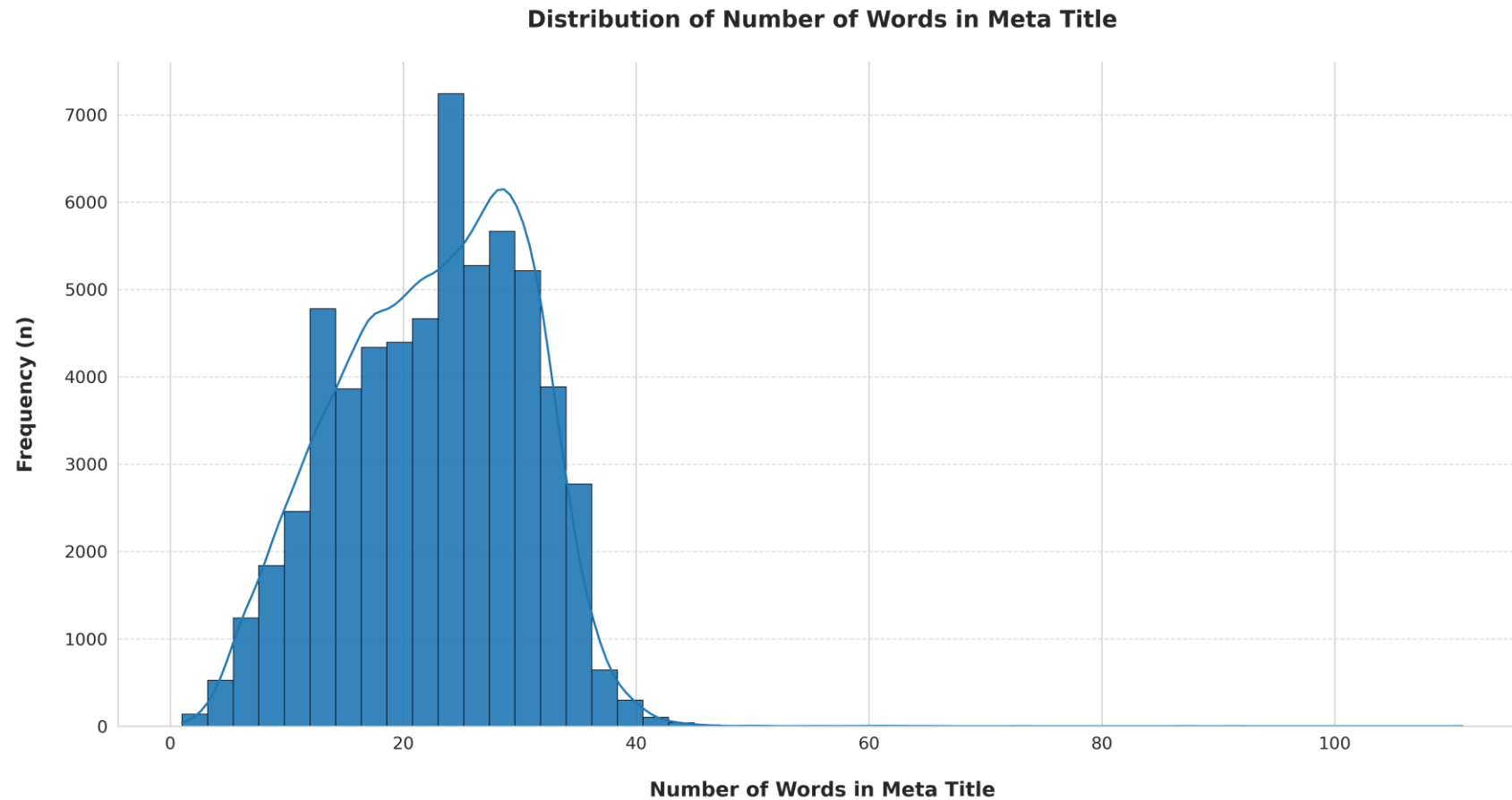


Distribution of Product Avg Ratings



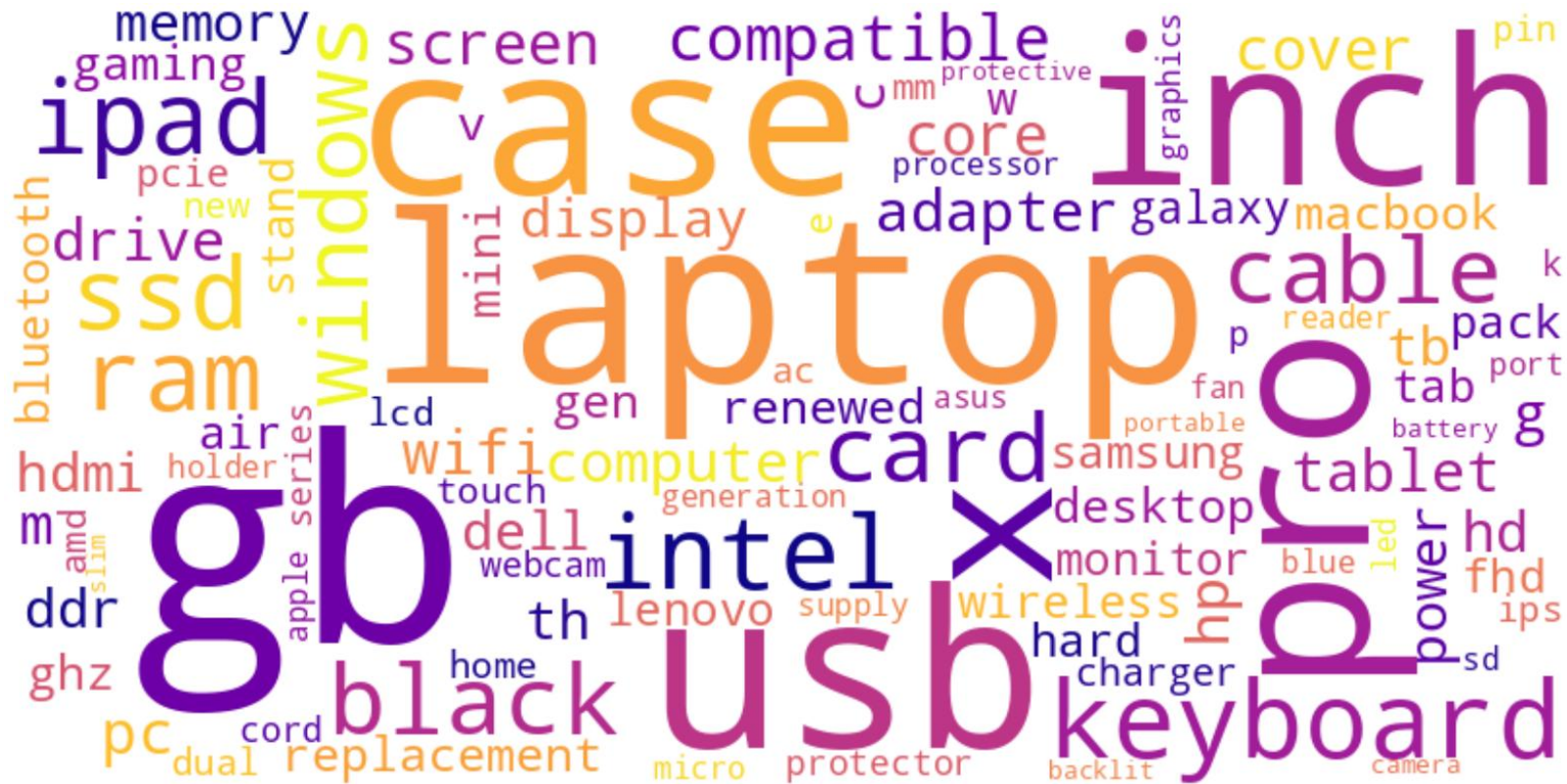


Products





Products



Unigrams

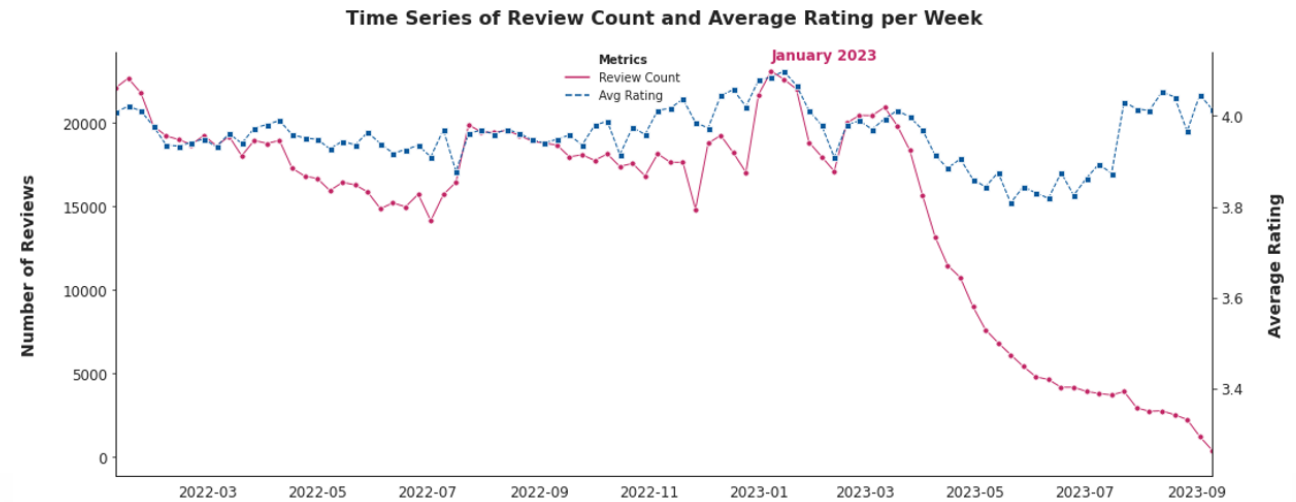
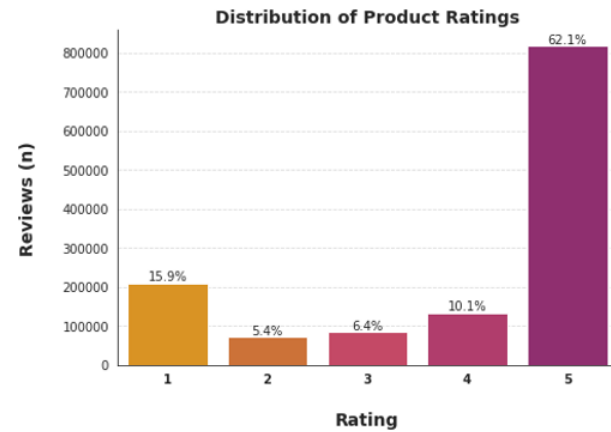


Products

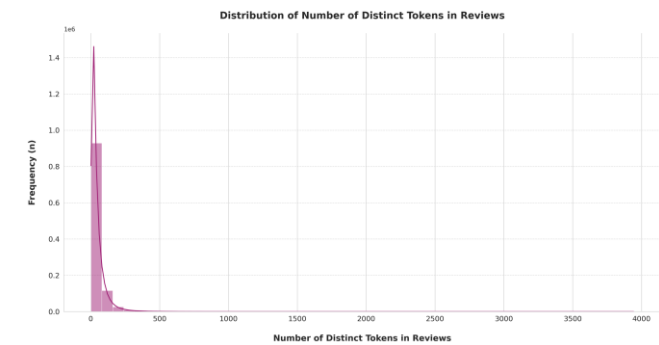
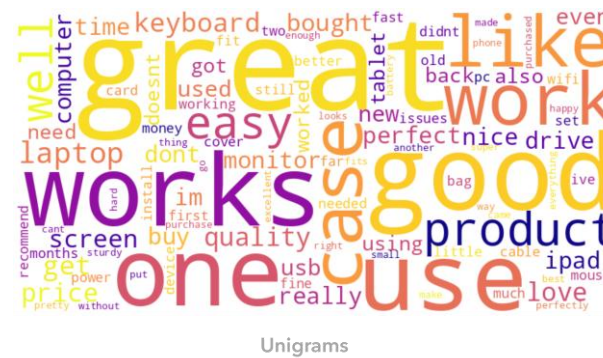
	n	%
Store		
Amazon Renewed	6402	6.4
HP	3699	3.7
Lenovo	2040	2.0
ASUS	1886	1.9
Dell	1444	1.4
Generic	948	1.0
SAMSUNG	830	0.8
MOSISO	744	0.7
SanDisk	725	0.7
Logitech	473	0.5

	n	%	Cumulative %
Price Category			
0 - 10	8774	8.8	8.8
10 - 50	28197	28.3	37.1
50 - 100	6958	7.0	44.1
100 - 250	7332	7.4	51.5
250 - 500	4205	4.2	55.7
500+	5612	5.6	61.3
Unknown	38557	38.7	100.0

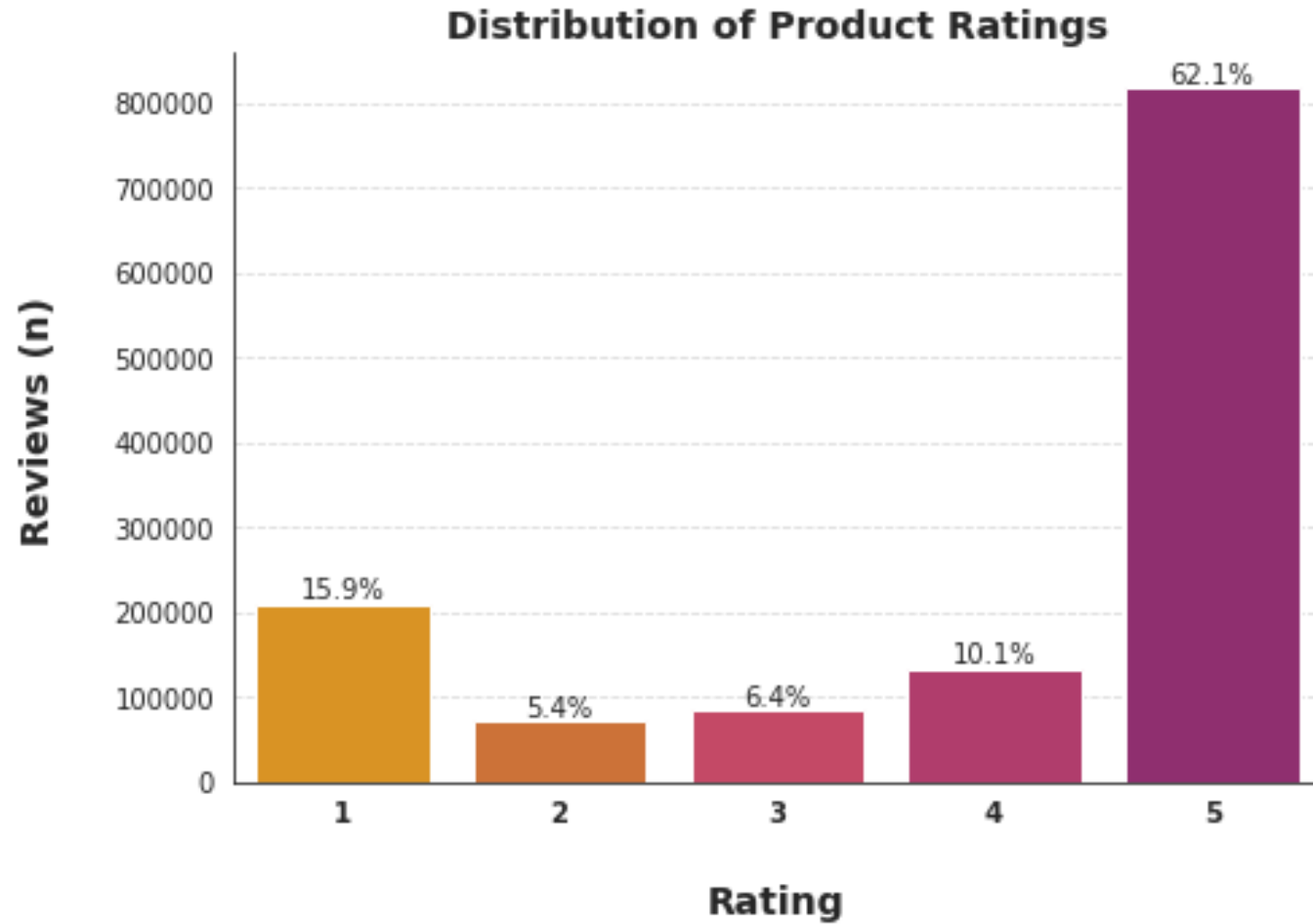
Reviews



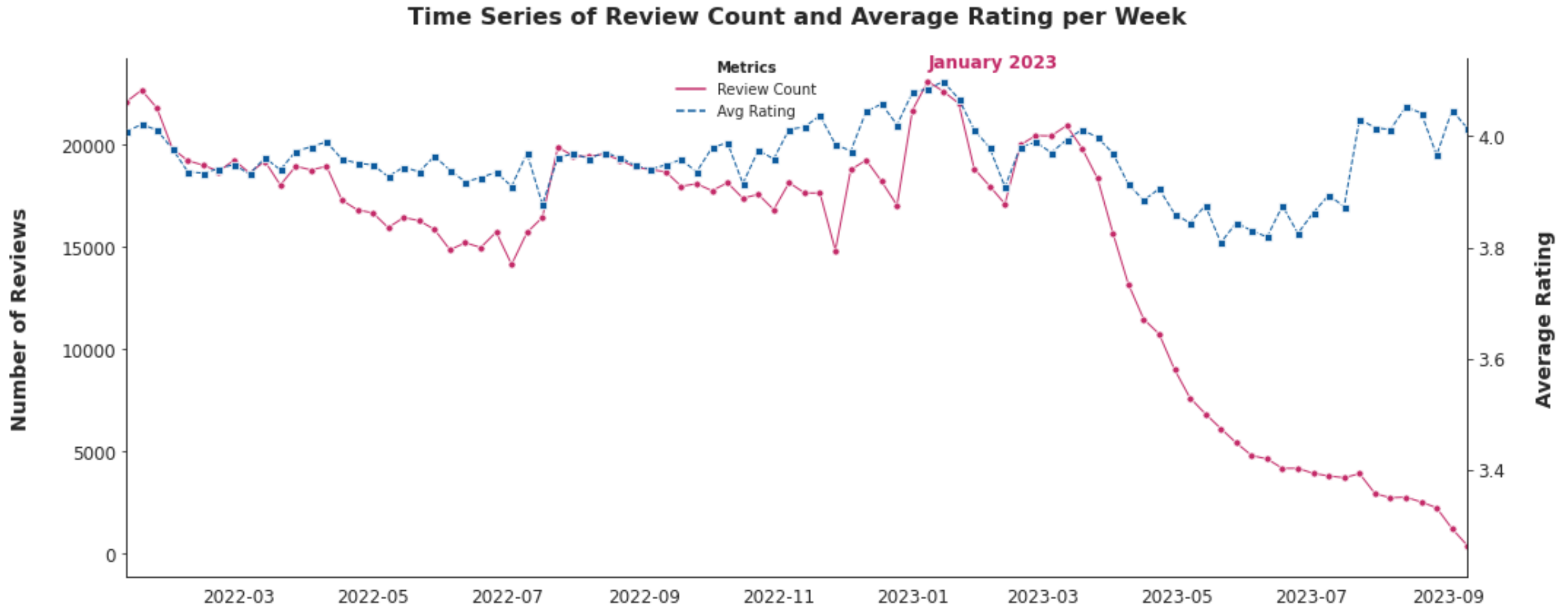
	n	%	Cumulative %
Verified Purchase			
True	1226888	93.21	93.21
False	89405	6.79	100.00



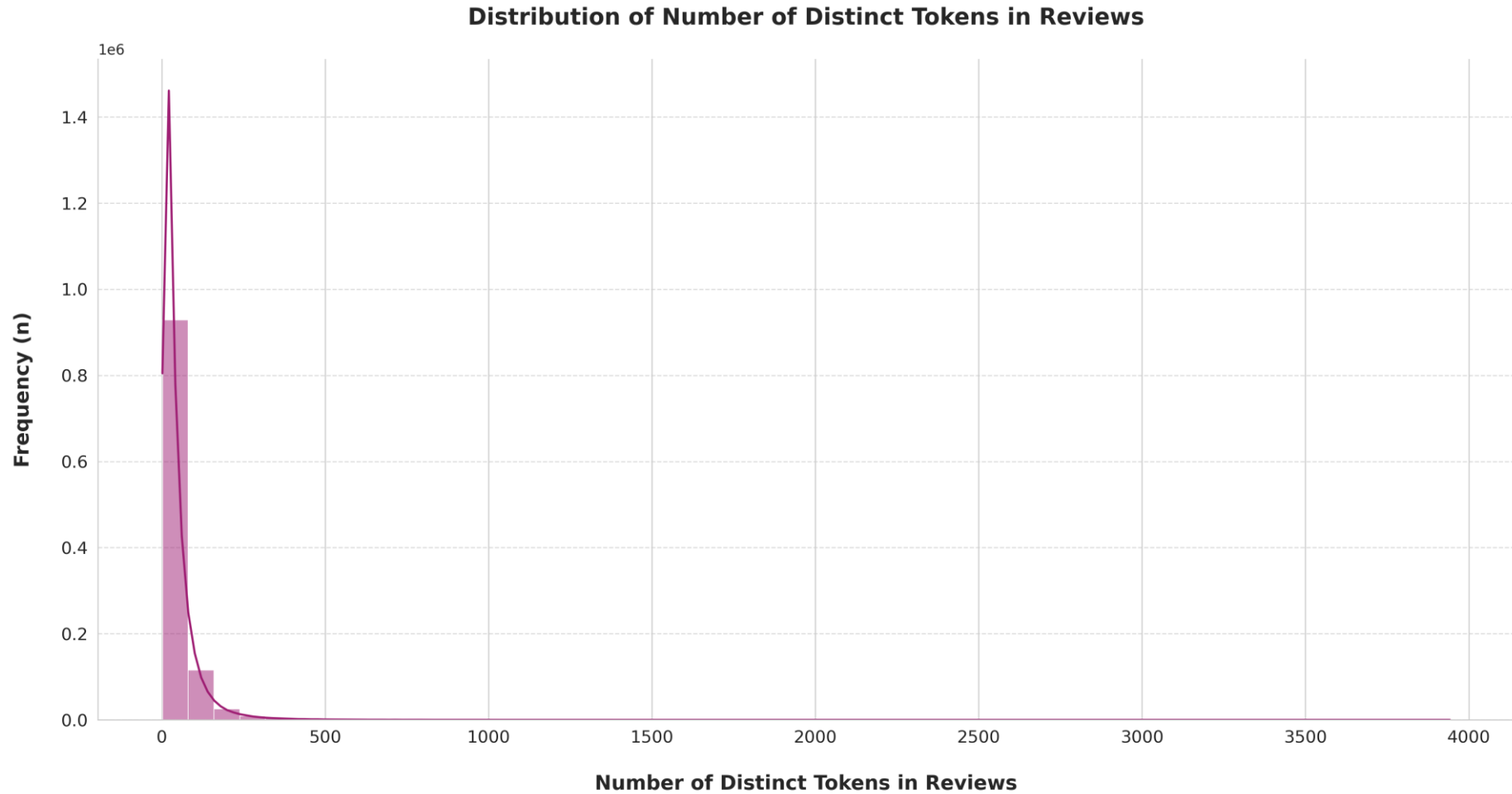
Reviews



Reviews



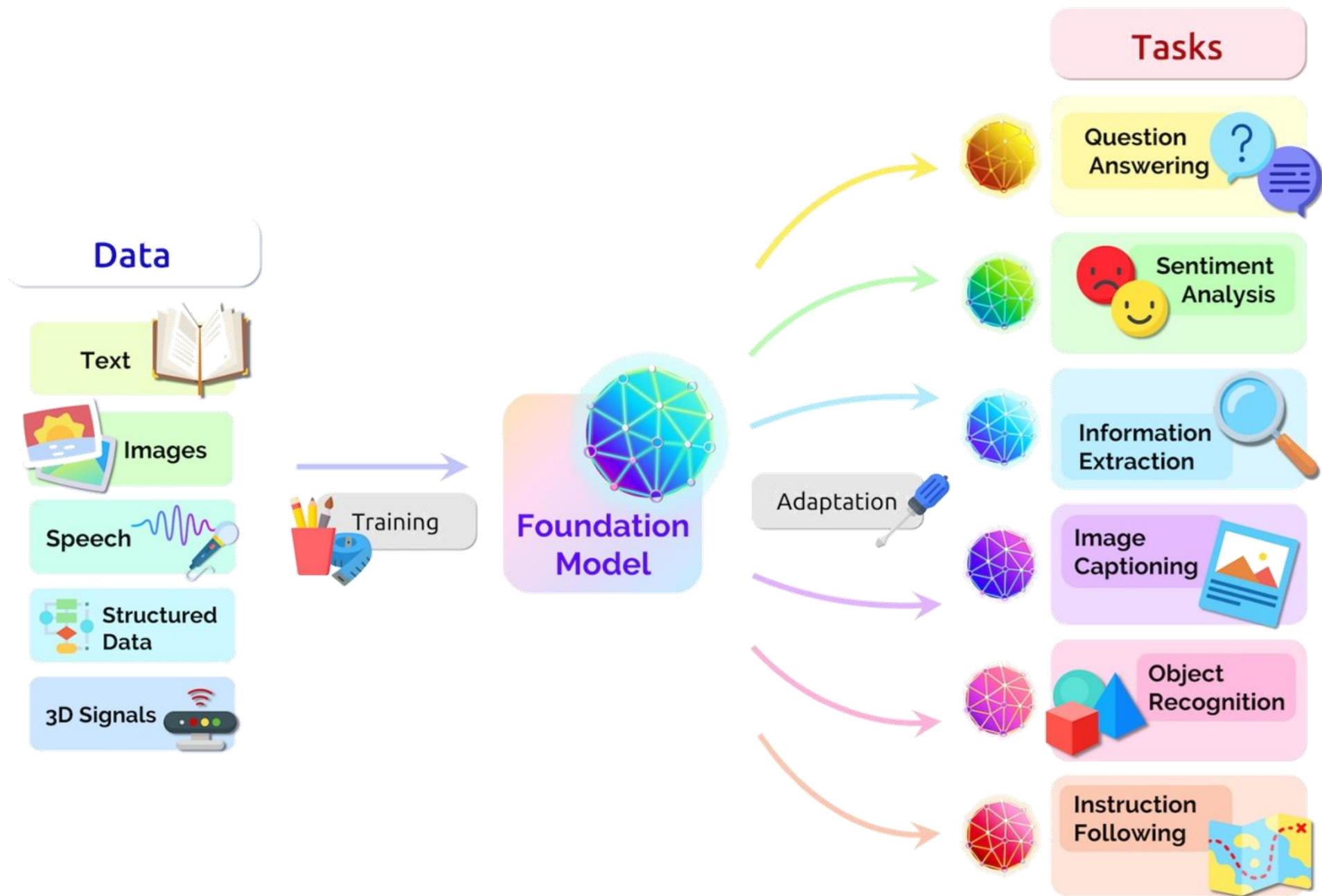
Reviews



Reviews



Unigrams



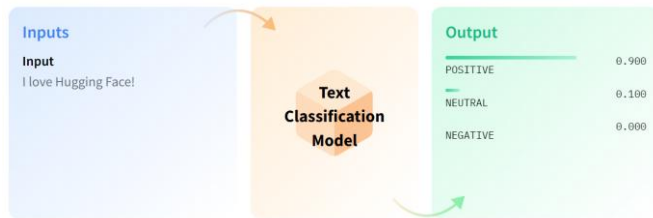


Hugging Face



Text Classification

Task of assigning a label or class to a given text, based on **predefined classes** trained by the model.



citizenlab/[twitter-xlm-roberta-base-sentiment-finetuned](#)
facebook + Google

- It is an **MLLM** (Multi-Lingual Language Model)
- Trained with **2.5TB** of *CommonCrawl* data filtered
- Adjusted in **~58M tweets** for *Sentiment Analysis*



Source: [Cardiff NLP Group](#) (2020)



Zero-Shot Classification

A task in which a model is trained on a set of labelled examples but is **then able to classify new examples** of previously unseen classes.



MoritzLaurer/[mDeBERTa-v3-base-xnli-multilingual-nli-2mil7](#)
Microsoft + facebook + Google

- It is an **MLLM** (Multi-Lingual Language Model) trained by *Microsoft* in +100 languages.
- Trained with **2.7M hypothesis-premise** pairs

Source: [Microsoft](#) (2023)



Hugging Face

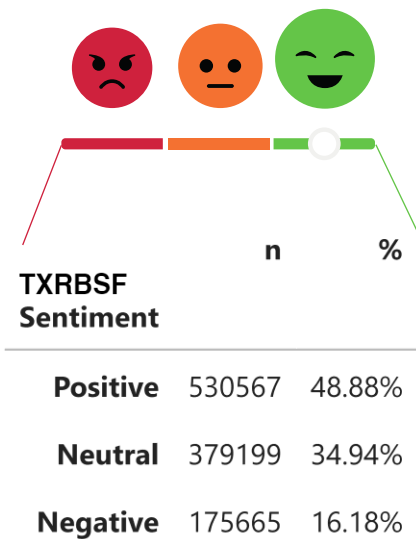
“

Beautiful, I really recommend it I
liked everything no complains and I
bought it for my son

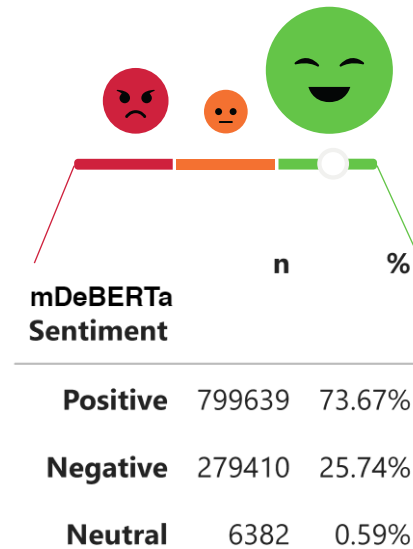
Amazon review

”

citizenlab/ twitter-xlm-roberta-base-sentiment-finetuned



MoritzLaurer/ mDeBERTa-v3-base-xnli-multilingual-nli-2mil7



Classification Matrix of the Results of Both Models

mDeBERTa Model Positive Neutral Negative

TXRBSF Model

Positive	525061	341	5165
Neutral	257855	5453	115891
Negative	16723	588	158354

n %

Combined Sentiment

Positive	525061	48.37
Positive Tendency	258196	23.79
Neutral	27341	2.52
Negative Tendency	116479	10.73
Negative	158354	14.59



Topic Analysis

Amazon Products



Hugging Face

🔗 Predefined Topics [12]

- 📺 Laptops
- 📺 Desktops
- 📺 PC Gaming
- 📺 Monitors
- 📺 Tablets
- 📺 Computer Components
- 📺 Computer Accessories
- 📺 Networking
- 📺 Drives & Storage
- 📺 Printers & Ink
- 📺 Software
- 📺 Others

“

**AOPEN by Acer 16PM1Q Bbmuiux
15.6" Full HD 1920 x 1080 IPS ...**

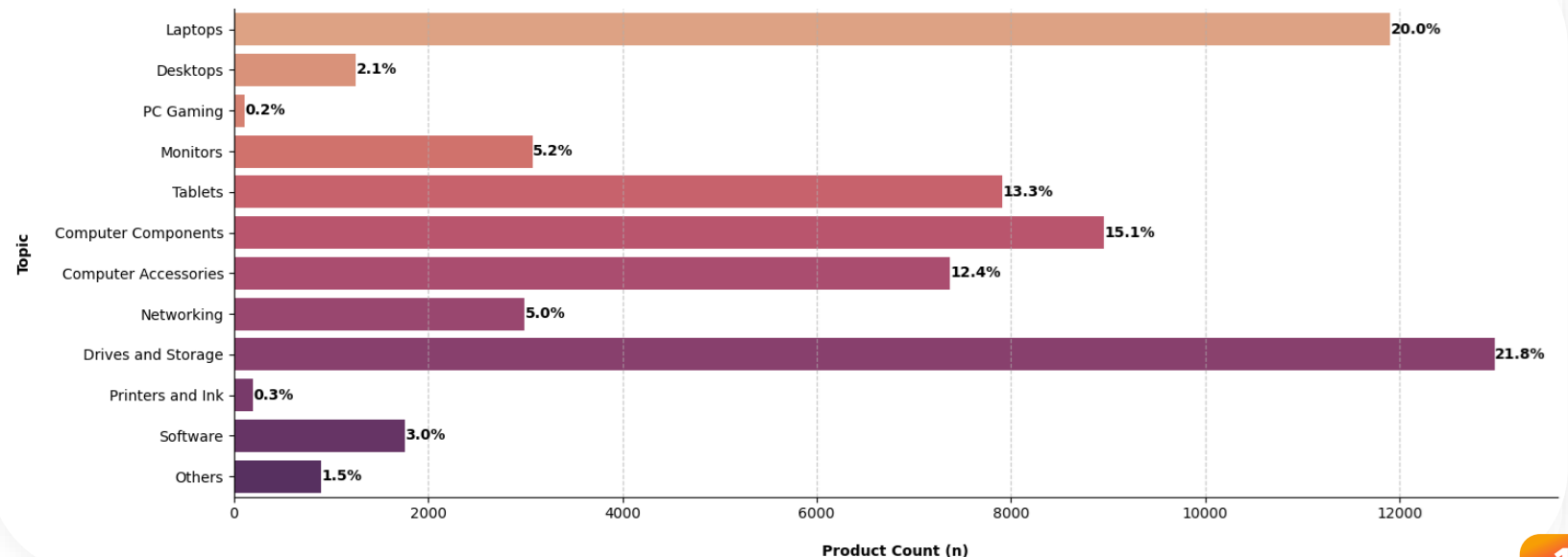
Example of an excerpt of Amazon electronic product

”

1st

Monitors
Score: 0.945

Topic Distribution of Products



Clustering Analysis

PCA & K-Means

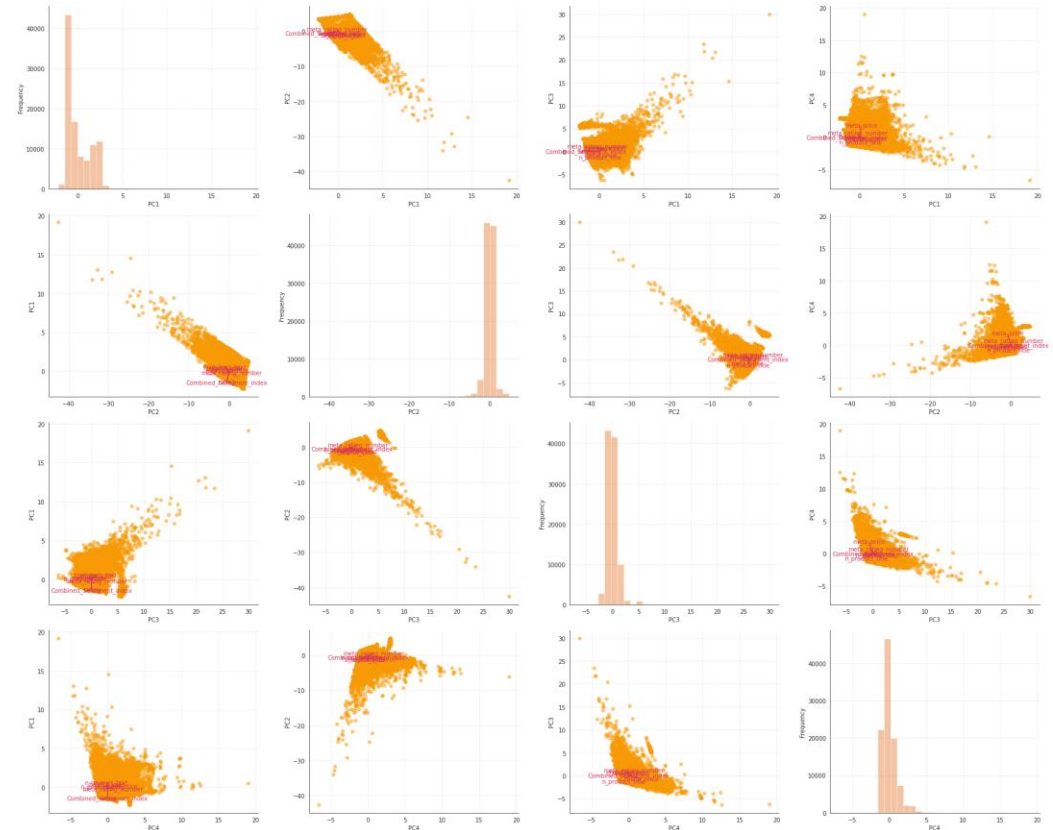
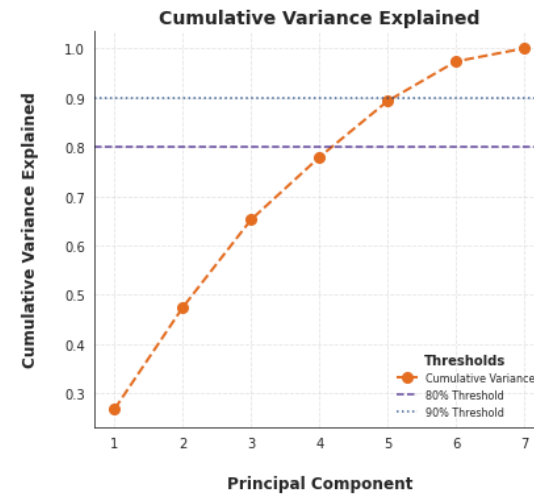
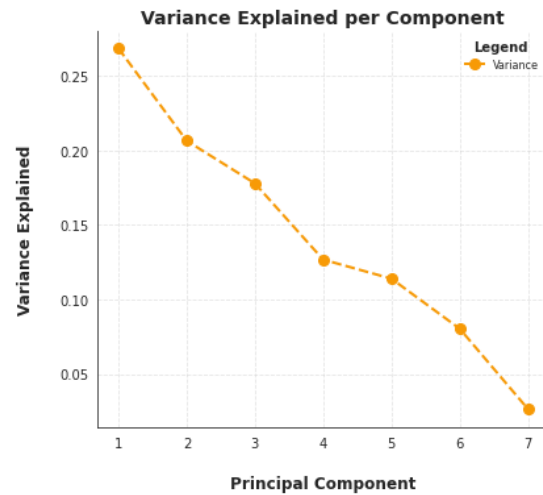
Combined_Sentiment	Negative	Negative Tendency	Neutral	Positive Tendency	Positive
Price Bucket					
0 - 10	15967 (14.7%)	11152 (10.3%)	2523 (2.3%)	28919 (26.7%)	49877 (46.0%)
10 - 50	83618 (14.1%)	57187 (9.7%)	14362 (2.4%)	137872 (23.3%)	299363 (50.5%)
50 - 100	20333 (15.5%)	15430 (11.7%)	3467 (2.6%)	32399 (24.7%)	59795 (45.5%)
100 - 250	22864 (15.0%)	18923 (12.4%)	4145 (2.7%)	36604 (24.0%)	70155 (45.9%)
250 - 500	9232 (15.2%)	7755 (12.8%)	1571 (2.6%)	13864 (22.9%)	28172 (46.5%)
500+	6340 (15.9%)	6032 (15.1%)	1273 (3.2%)	8538 (21.4%)	17699 (44.4%)

Products: ~10k
Reviews: ~100k

PCA

(Principal Component Analysis) - 7 Features

PCA Variance Explained



Clustering Analysis

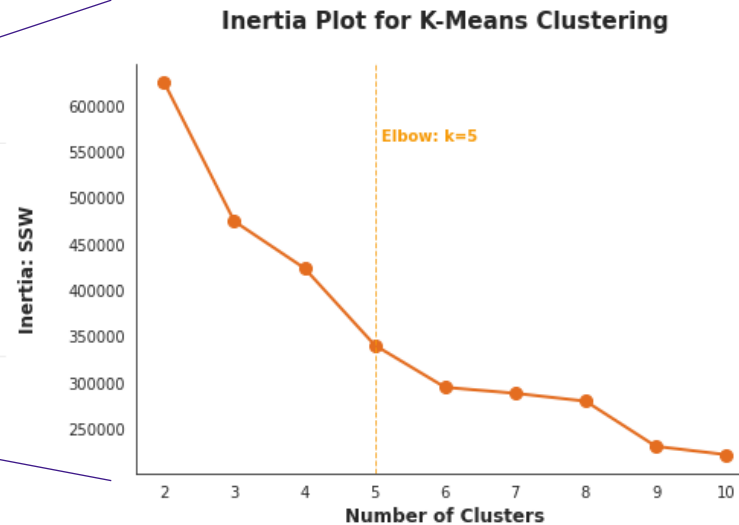
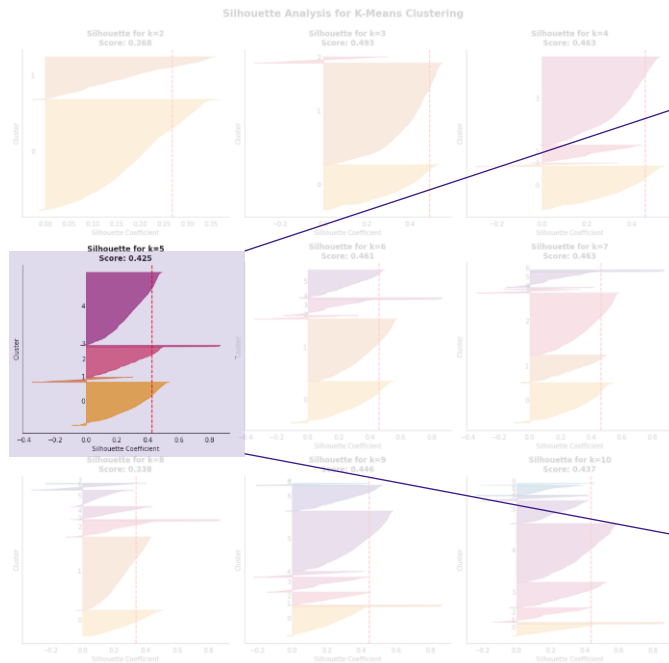
PCA & K-Means

Combined_Sentiment	Negative	Negative Tendency	Neutral	Positive Tendency	Positive
Price Bucket					
0 - 10	15967 (14.7%)	11152 (10.3%)	2523 (2.3%)	28919 (26.7%)	49877 (46.0%)
10 - 50	83618 (14.1%)	57187 (9.7%)	14362 (2.4%)	137872 (23.3%)	299363 (50.5%)
50 - 100	20333 (15.5%)	15430 (11.7%)	3467 (2.6%)	32399 (24.7%)	59795 (45.5%)
100 - 250	22864 (15.0%)	18923 (12.4%)	4145 (2.7%)	36604 (24.0%)	70155 (45.9%)
250 - 500	9232 (15.2%)	7755 (12.8%)	1571 (2.6%)	13864 (22.9%)	28172 (46.5%)
500+	6340 (15.9%)	6032 (15.1%)	1273 (3.2%)	8538 (21.4%)	17699 (44.4%)

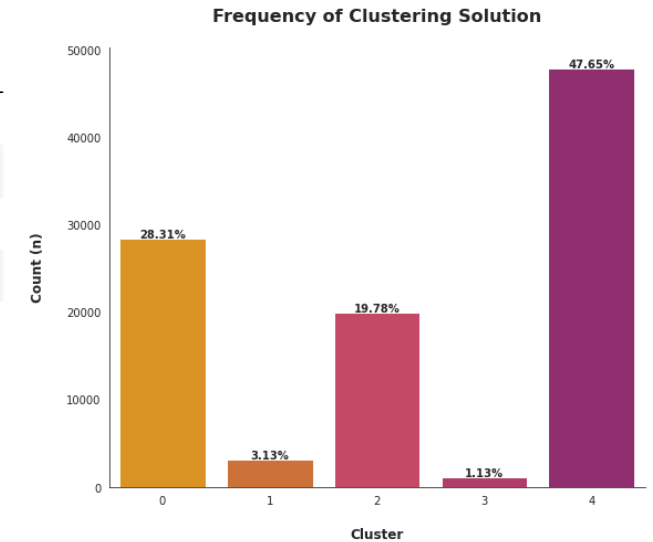
Products: ~10k
Reviews: ~100k

Clustering

K-Means & Original 7 Features



Cluster	n	%
0	28444	28.31
1	3143	3.13
2	19874	19.78
3	1134	1.13
4	47881	47.65



Confusion Matrix: True Sentiment vs Clusters

	0	1	2	3	4
Positive	89	179	12329	591	32683
Positive Tendency	1165	1923	5806	327	13181
Neutral	951	43	607	18	1225
Negative Tendency	11051	922	929	105	780
Negative	15188	76	203	93	12

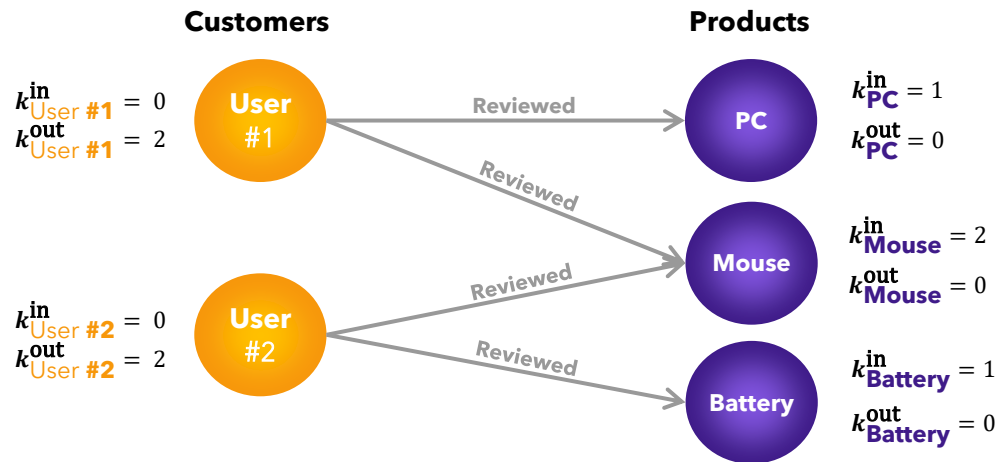
True Sentiment

Cluster

Combined_Sentiment	Negative	Negative Tendency	Neutral	Positive Tendency	Positive
Price Bucket					
0 - 10	15967 (14.7%)	11152 (10.3%)	2523 (2.3%)	28919 (26.7%)	49877 (46.0%)
10 - 50	83618 (14.1%)	57187 (9.7%)	14362 (2.4%)	137872 (23.3%)	299363 (50.5%)
50 - 100	20333 (15.5%)	15430 (11.7%)	3467 (2.6%)	32399 (24.7%)	59795 (45.5%)
100 - 250	22864 (15.0%)	18923 (12.4%)	4145 (2.7%)	36604 (24.0%)	70155 (45.9%)
250 - 500	9232 (15.2%)	7755 (12.8%)	1571 (2.6%)	13864 (22.9%)	28172 (46.5%)
500+	6340 (15.9%)	6032 (15.1%)	1273 (3.2%)	8538 (21.4%)	17699 (44.4%)

Products: ~10k
Reviews: ~100k

Graph ( GraphFrames)
Users-Products



Algorithms Applied

Page Rank

Higher PageRank Scores: For products, indicate popularity or frequent reviews.

Community Detection Label Propagation

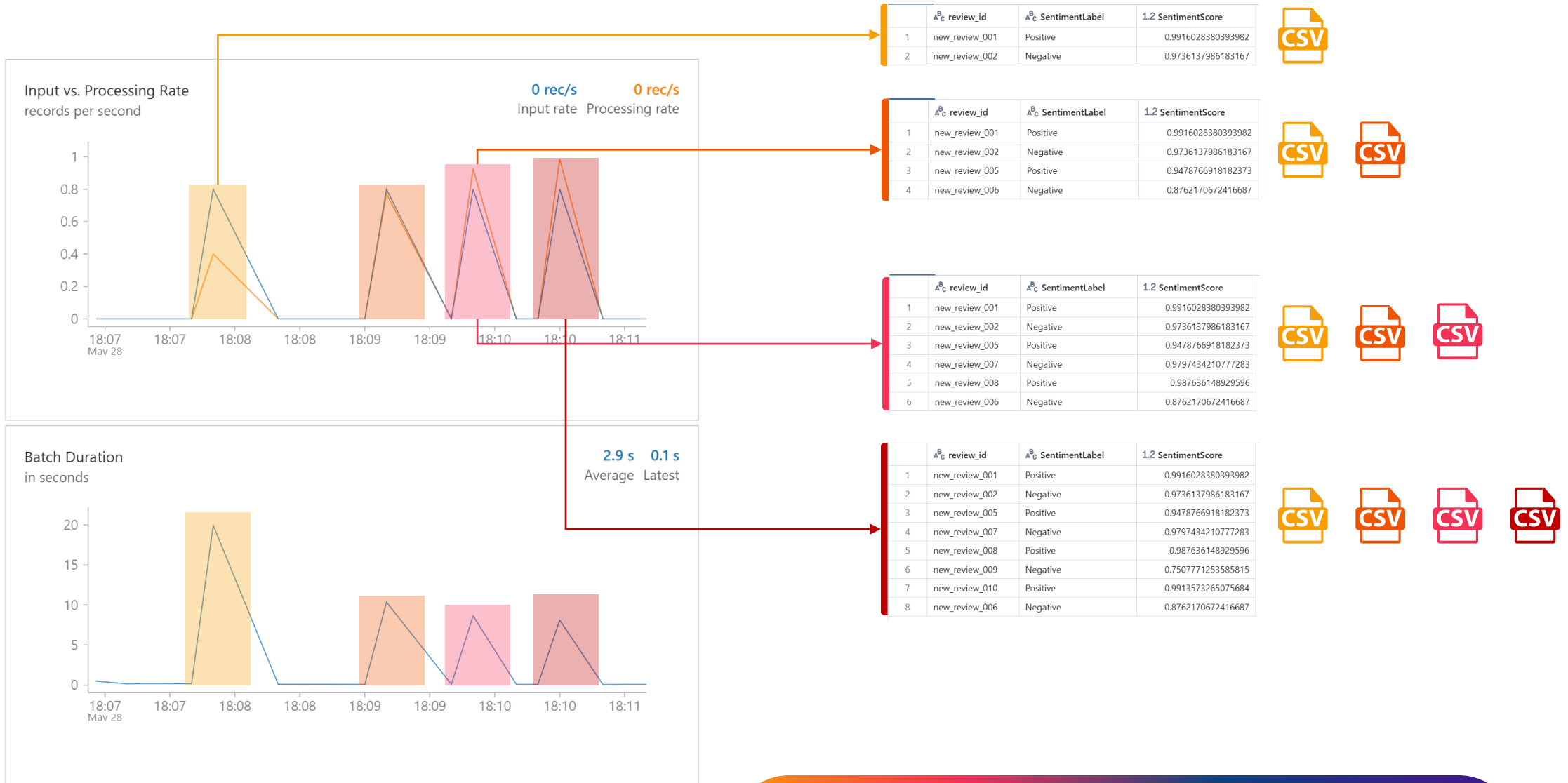
Larger Communities: indicate groups of users frequently reviewing similar high-value products, suggesting shared interests or preferences.

Smaller Communities: may represent niche markets or specialized product categories, which could be targeted for personalized marketing.

Degree Distribution Analysis

In-Degree: Products with high in-degrees are reviewed by many users, indicating popularity or visibility.

Out-Degree: Users with high out-degrees are prolific reviewers, contributing significantly to the review ecosystem.



Conclusion

Time Series of Average True Rating vs. Predicted Sentiment per Week

