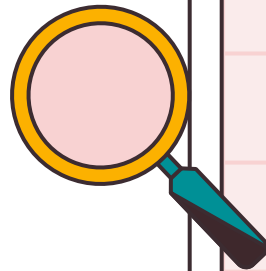
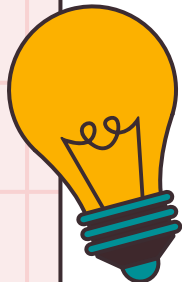




DS 5110/CS 5501

Scalable Ray: Sentiment Analysis of Amazon Reviews

By: Thomas Burrell, Ethan Assefa, Tatev Gomtsyan



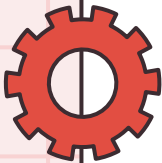


Table of Contents

01

**Project
Motivation**

02

Process

03

EDA & Analysis

04

**Results &
Conclusion**



Project Motivation



User-generated reviews have gained real momentum on e-commerce platforms like Amazon in the 21st century. They contain valuable insights regarding product quality, customer satisfaction, and market trends. Determining market sentiment and assessing unfilled needs is crucial for businesses interested in hidden niches, rich in potential customers. Big data analysis can provide a window into emerging trends and changing consumer preferences.



amazon



Our Data

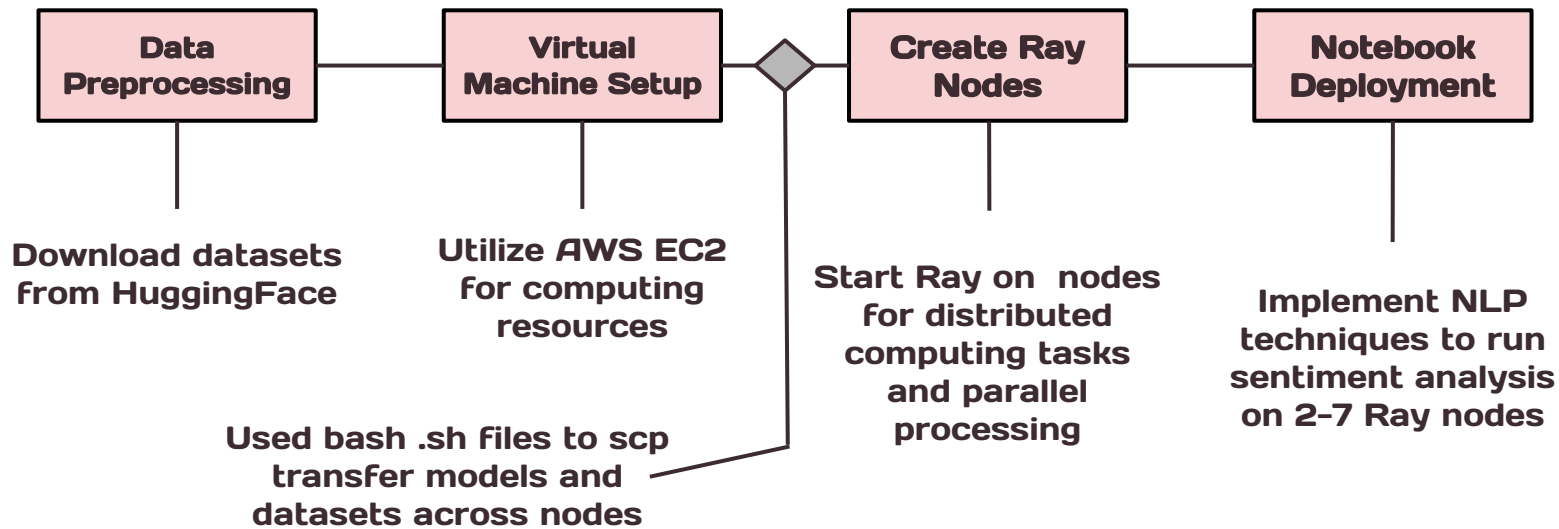
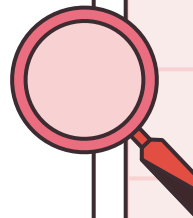
Year	#Review	#User	#Item	#R_Token	#M_Token	#Domain	Timespan
<u>2023</u>	571.54M	54.51M	48.19M	30.14B	30.78B	33	May'96 - Sep'23

Grouped by Category

Category	#User	#Item	#Rating	#R_Token	#M_Token	Download
All_Beauty	632.0K	112.6K	701.5K	31.6M	74.1M	review , meta
Amazon_Fashion	2.0M	825.9K	2.5M	94.9M	510.5M	review , meta
Appliances	1.8M	94.3K	2.1M	92.8M	95.3M	review , meta
Arts_Crafts_and_Sewing	4.6M	801.3K	9.0M	350.0M	695.4M	review , meta
Automotive	8.0M	2.0M	20.0M	824.9M	1.7B	review , meta
Baby_Products	3.4M	217.7K	6.0M	323.3M	218.6M	review , meta
Beauty_and_Personal_Care	11.3M	1.0M	23.9M	1.1B	913.7M	review , meta
Books	10.3M	4.4M	29.5M	2.9B	3.7B	review , meta
CDs_and_Vinyl	1.8M	701.7K	4.8M	514.8M	287.5M	review , meta
Cell_Phones_and_Accessories	11.6M	1.3M	20.8M	935.4M	1.3B	review , meta

- Text data
- # Rows: ~571 million
- We chose categories ~1-2.5 GB:
All_Beauty, Appliances,
Musical_Instruments, Video_Games
- Random sampling of each dataset

Process Workflow

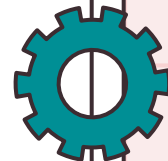


Using Ray

An open-source framework designed for scaling Python applications from a single computer to a large cluster with ease. Ray is particularly well-suited for machine learning, deep learning, and other data-intensive tasks.

Key features:

- Task Parallelism
- Actor Model
- Scalability and Flexibility



Resources

We used the AWS EC2 services to create our computing infrastructure

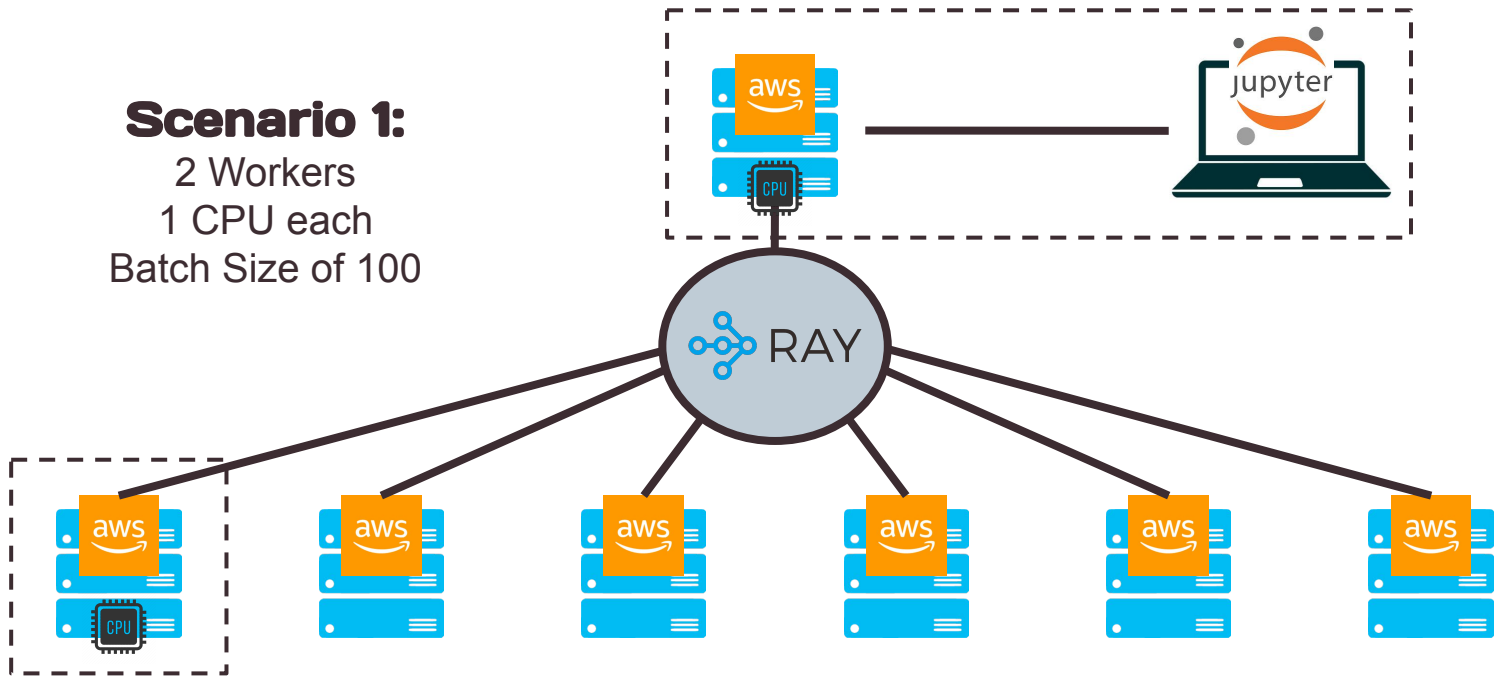
- The instance type we used was **t3.large**
 - On-Demand hourly rate: \$0.0832
 - vCPU: 2
 - Memory: 8 GiB
- We had **7 instances** total
- **Total cost was ~ \$60** at the end of all analysis



Instance Setup

Scenario 1:

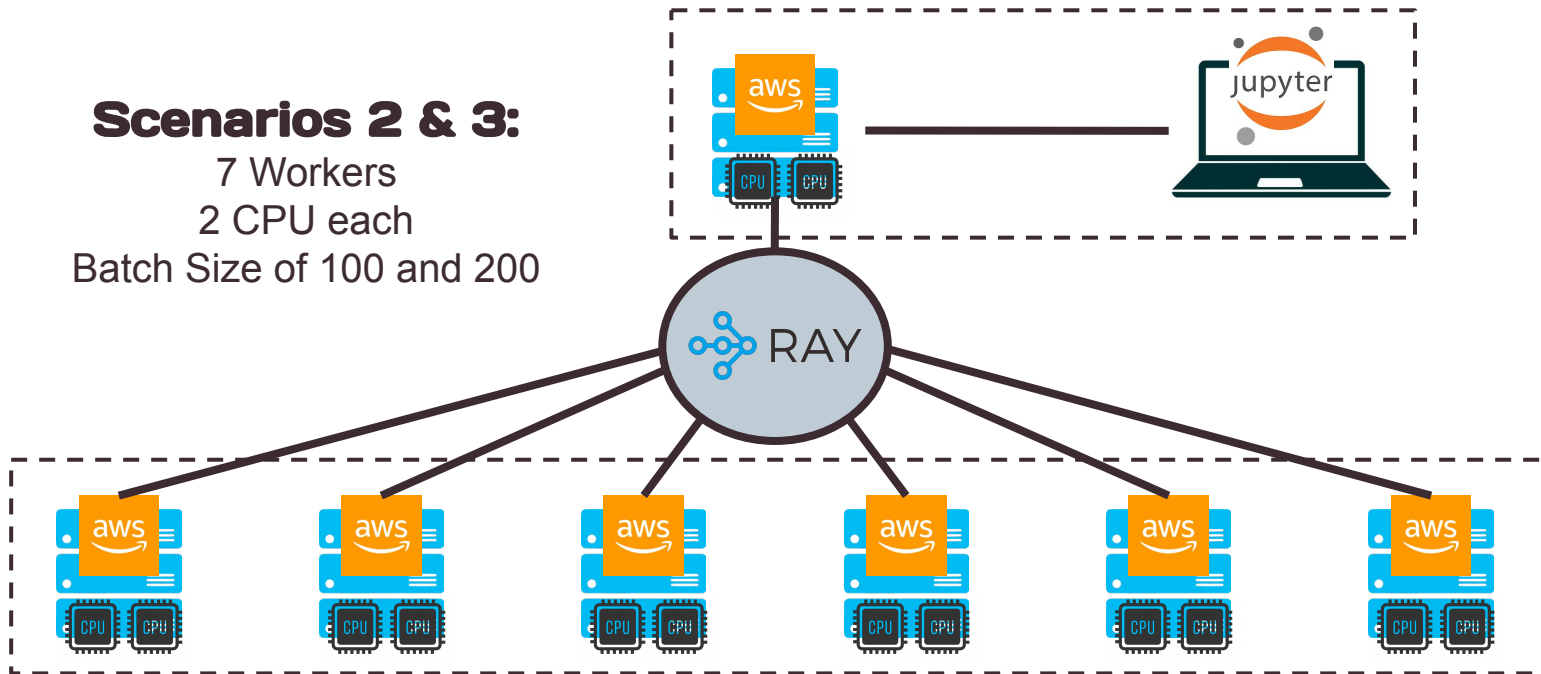
2 Workers
1 CPU each
Batch Size of 100



Instance Setup

Scenarios 2 & 3:

7 Workers
2 CPU each
Batch Size of 100 and 200



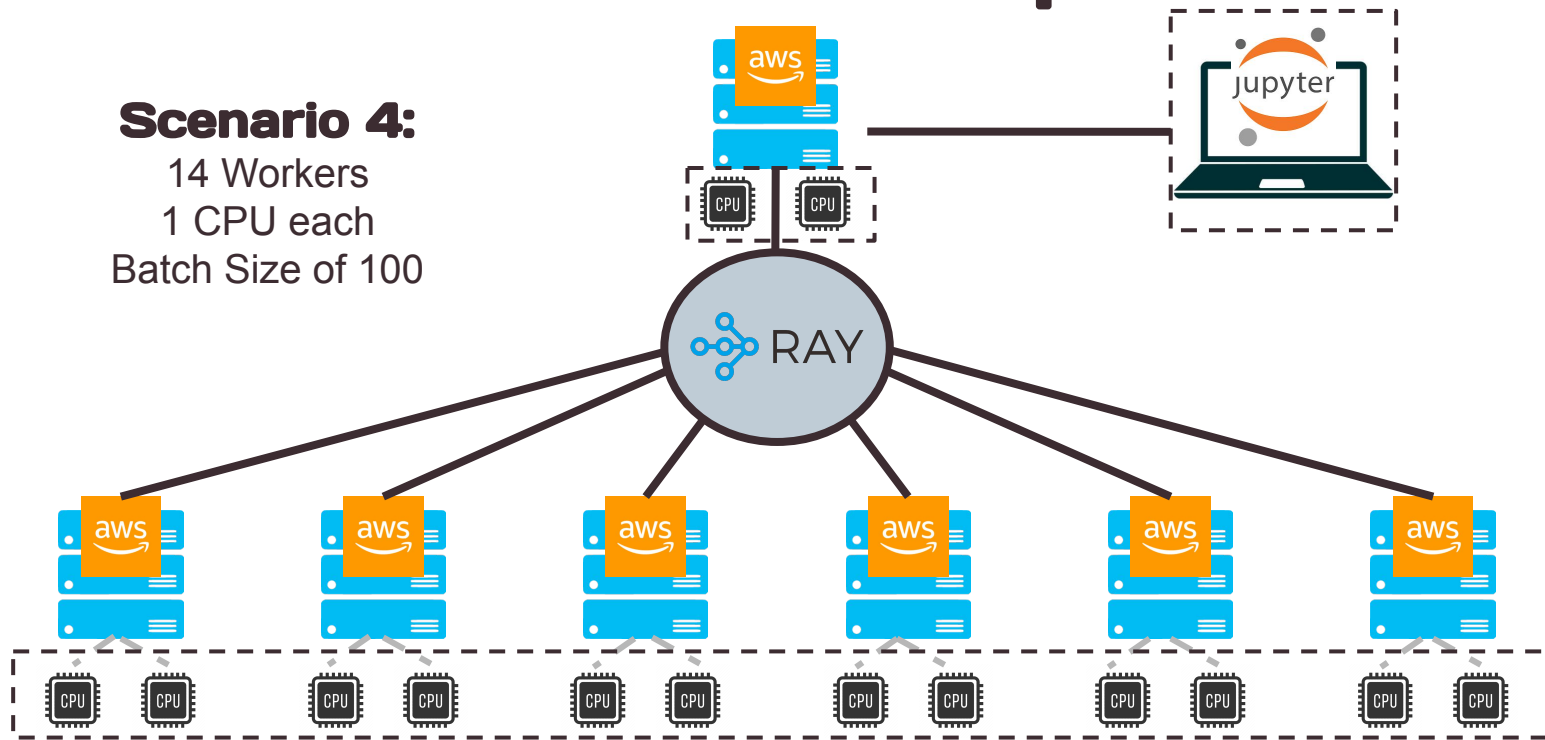
Instance Setup

Scenario 4:

14 Workers

1 CPU each

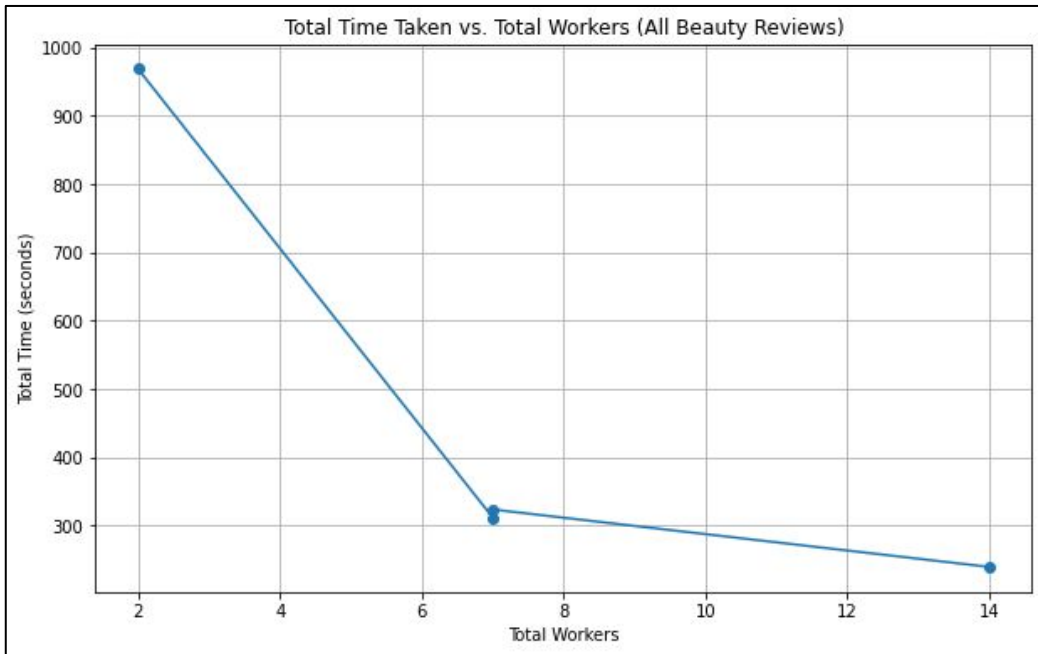
Batch Size of 100



'All Beauty' Cluster Performance

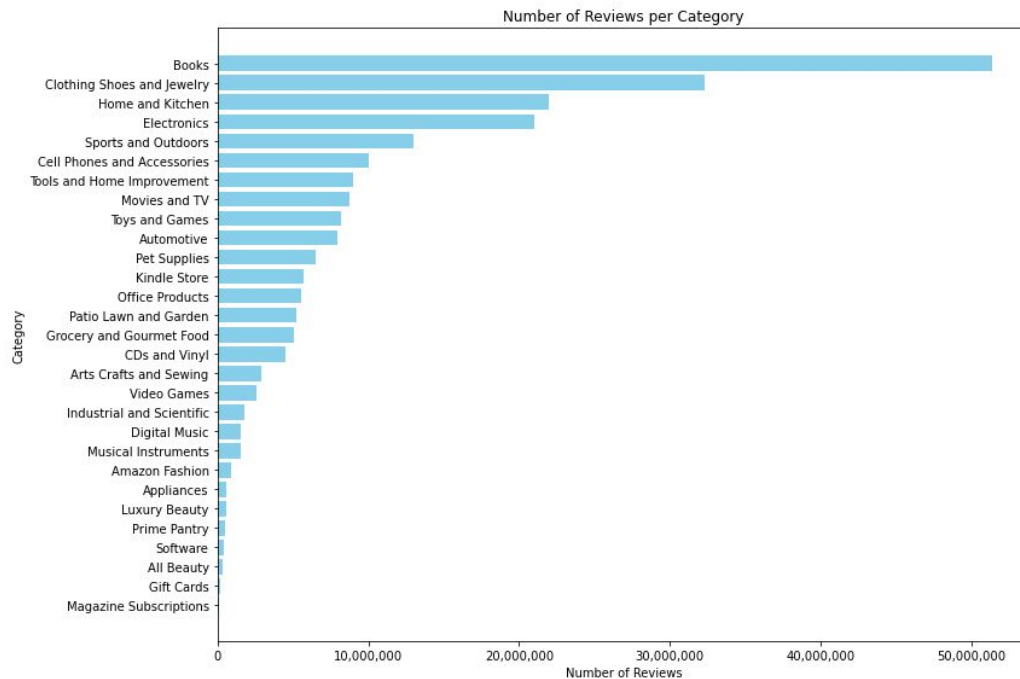
Total Workers	CPUs per Worker	Batch Size	Total Time (s)	Avg. Time per Batch (s)	Avg. Time per Text (s)
2	1	100	968.61	9.5	0.049
7	2	100	312.15	9.42	0.015
7	2	200	323.85	18.82	0.016
14	1	100	239.81	14.54	0.012

Ray Performance Results

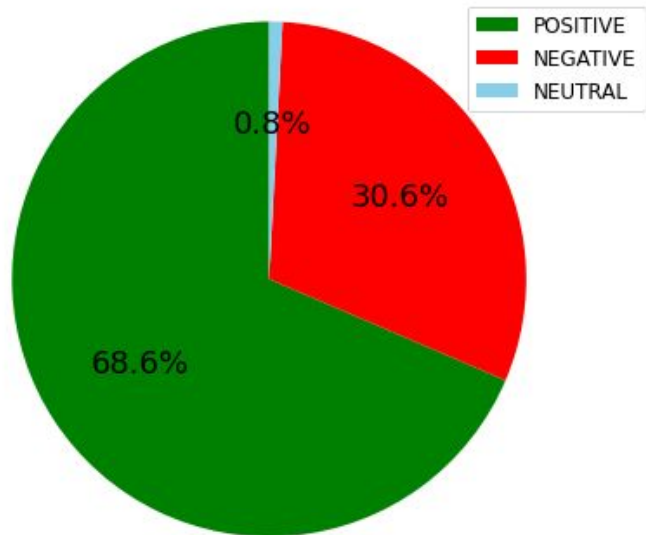


As we increased the number of workers, our **processing time decreased substantially**

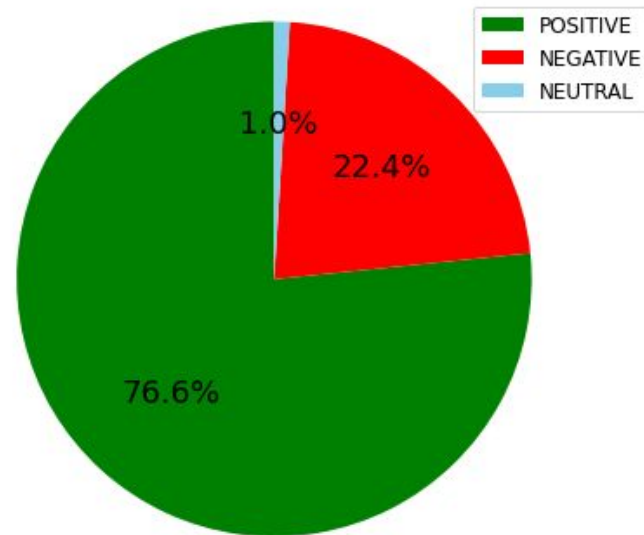
Exploratory Data Analysis



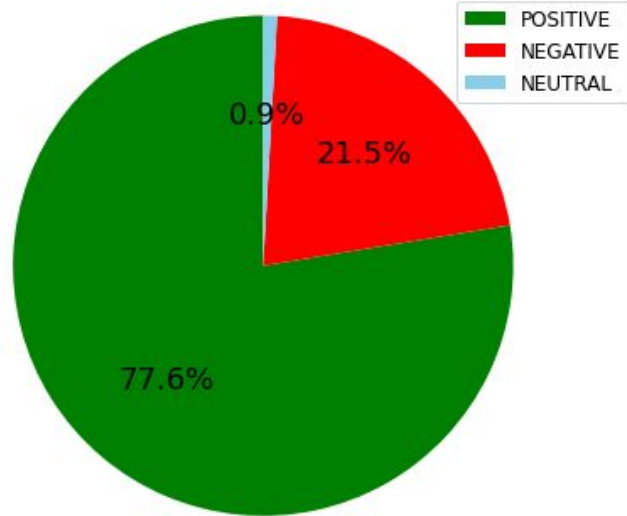
Distribution of Sentiment Categories for All Beauty Products



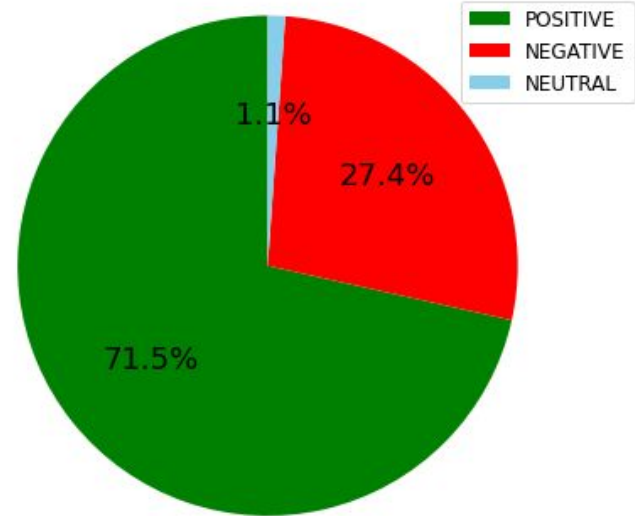
Distribution of Sentiment Categories for Appliances



Distribution of Sentiment Categories for Musical Instruments



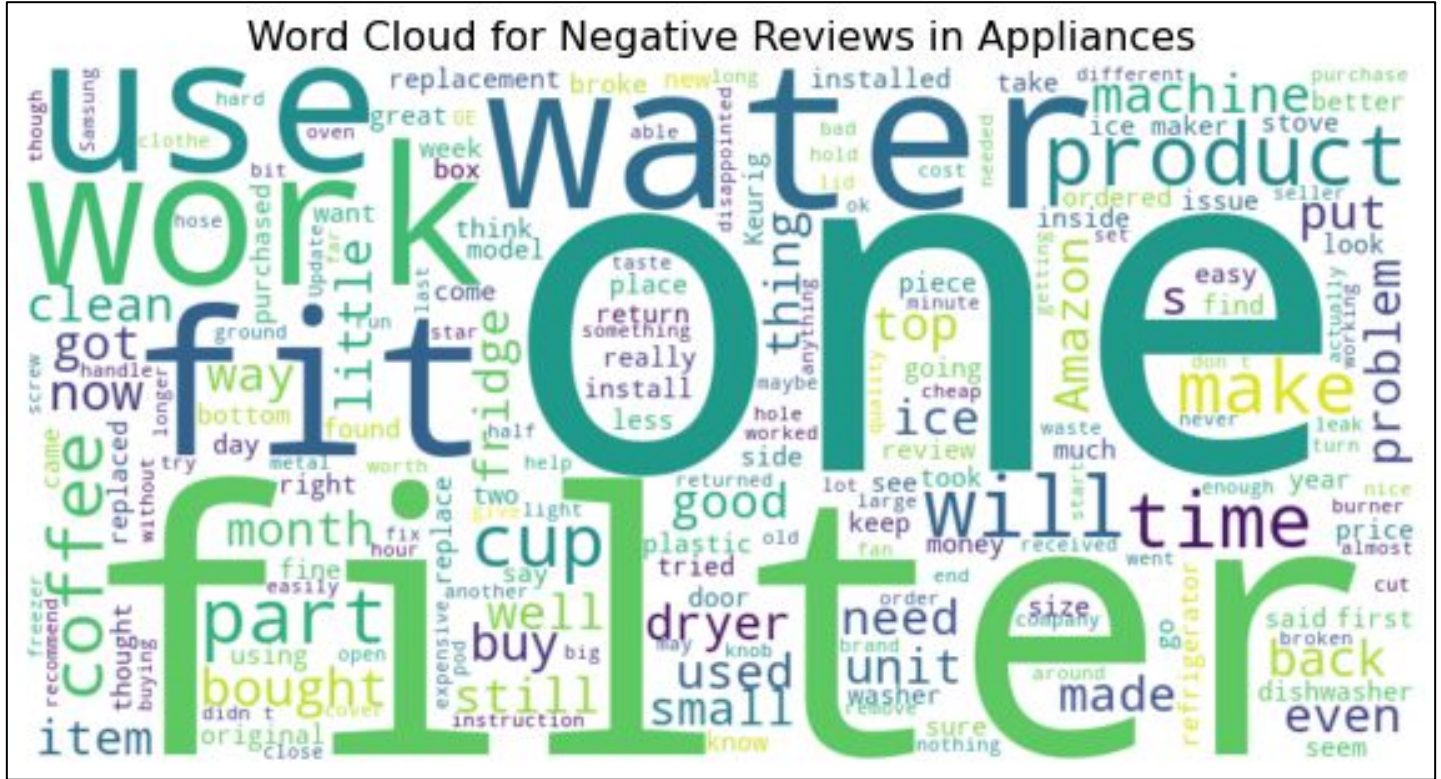
Distribution of Sentiment Categories for Video Games



[illegible]

[illegible]

[illegible][illegible]

[illegible]



Results

We ranked the **four** Amazon categories we investigated in terms of negative review percentage:

1. **Beauty Products** (30.6%)
2. **Video Games** (27.4%)
3. **Appliances** (22.4%)
4. **Musical Instruments** (21.5%)

Results Cont.

Common words/themes included in each negative review category:

1. **Beauty Products:** Hair, Color
2. **Video Games:** Controller, Keyboard
3. **Appliances:** Filter, Fit
4. **Musical Instruments:** Sound, Guitar



Conclusion

We identified a significant market opportunity within the **Beauty products** category, which had the highest percentage of negative reviews expressing dissatisfaction or criticism (30.6%). This indicates a clear demand for higher-quality products particularly in **hair-care**, and suggests a gap in the market that manufacturers can capitalize on.

Limitations/Future Work

- Computational Resources
- Non-text reviews (emojis)
- Year-by-year changes



Thanks!

Do you have any questions?

