# CS5344 Final Project - Group 21

JAMAL AHAMED, LALITHA RAVI, PARUL BANSAL, SREELAKSHMI
{A0268403E, A0268254X, A0268257R, A0268357N}
{e1101702@u.nus.edu , e1101553@u.nus.edu , e1101556@u.nus.edu , e1101656@u.nus.edu}
*National University of Singapore, Singapore*

*Abstract*—**This report demonstrates the application of the knowledge acquired in the CS5344 course, showcasing our implementation of big data technologies, data preprocessing, feature extraction, and sentimental analysis using the AL/DL approach to extract actionable insights from the Amazon Reviews dataset.**

## I. INTRODUCTION

In the fast-paced and competitive world of online shopping, it's crucial for businesses to know how customers feel. But it's not easy. There are so many online reviews, and customers express a wide range of emotions, from happy to unhappy. E-commerce sellers have a tough job figuring out what all these feelings mean and how they can use them to improve their products, marketing, and customer satisfaction efforts. To solve this challenge, we need to use advanced big data techniques because dealing with millions of reviews manually is impractical. The data includes sentiments, common issues, and comparisons, and sophisticated tools and technologies are necessary for extracting valuable insights efficiently.

Our project aims to help online sellers and brands succeed in the ever-changing world of e-commerce. Using advanced big data techniques and analytics, we want to make it easier for them to understand what buyers think. Our goal is to simplify the complex e-commerce landscape, giving sellers a complete understanding and helping them thrive in this dynamic marketplace. Singaporean housing market.

## II. MOTIVATION AND GOALS

This project aims to transform and strengthen the world of online shopping, especially for sellers and brands. We want to create an advanced system that analyzes Amazon reviews to provide valuable insights. Using sentiment analysis, we hope to understand buyer emotions and perspectives, helping sellers better meet consumer needs. By identifying common product issues, we can proactively solve problems, leading to higher product quality and customer satisfaction. The system will also excel in recognizing and highlighting key product features, giving sellers a competitive advantage. Additionally, it will offer comparative analysis, including sales metrics and features of similar products, providing a comprehensive overview for strategic decision-making. Ultimately, our goal is to empower sellers and brands with a cutting-edge solution that optimizes their offerings and enhances the overall online shopping experience. The code for the project can be found here: https://tinyurl.com/projectcodegrp21

## III. DATASET

Our project utilizes data from the 2018 Amazon reviews dataset featuring unique identifiers such as asin for each product, reviewerID for each reviewer, overall rating on a scale of 1 to 5 stars, reviewTime indicating the submission date and time, reviewText encompassing the full review content, title representing the product title, and brand denoting the product brand. Focusing specifically on this electronics-based dataset, which includes a diverse array of gadgets, appliances, and devices, our analysis incorporates a substantial size of 20,994,353 reviews and 766,868 products metadata. By leveraging the electronics-based information within this dataset, our project aims to unveil trends, monitor shifts in sentiment, pinpoint common issues, and spotlight product features unique to this industry. This intentional selection of the electronics-focused data source ensures that our analysis is not only comprehensive but also highly applicable to the nuanced dynamics of the electronics market. Figure 1 shows some features of reviews dataset.
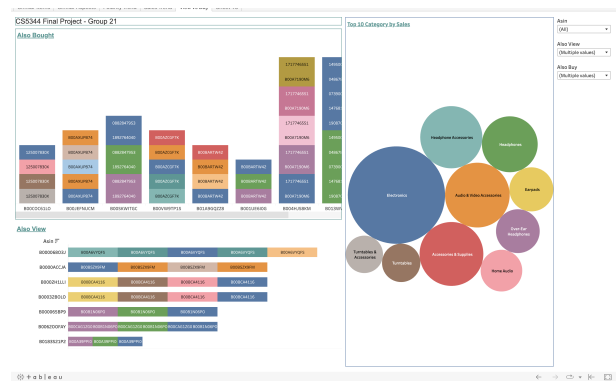


Fig. 1. Reviews Dataset

## IV. DATA PREPROCESSING

Data preprocessing is a fundamental and indispensable phase in big data analysis. It involves tasks such as handling missing values, standardizing data formats, removing duplicates, and transforming data into a consistent structure. The

primary goal of this preprocessing task is to intricately clean and refine the 'reviewText' field from the dataset so as to perform analyses and extract profound insights from it.

### A. Handling Null Values

- 'reviewText': Null values in the 'reviewText' column were removed from the dataset. The presence of text is crucial for any text-based analysis, and retaining records with null values in this context would not contribute meaningfully to downstream tasks.
- 'vote': The 'vote' column, representing the number of helpful votes, was handled by filling null values with zeros.

### B. Data Cleaning

The review text field from the data set had a lot of unclean data, so We created a cleaning script for the dataset.

- **Emojis**: Emojis found in the 'reviewText' column underwent conversion to text using the emoji Python library.
- **Special Characters, Links, and Mentions**: Special characters, links, mentions, and non-ASCII characters were systematically removed from the text. This process aimed to eliminate irrelevant or potentially disruptive elements from the review text.
- **Additional Blank Spaces**: Removed to enhance the overall consistency and readability of the text.
- **Text within Square Brackets**: Any text enclosed within square brackets was systematically removed.
- **Sequences with Three or More Slashes**: Sequences containing three or more slashes were systematically removed, promoting a cleaner and more structured representation of the text.
- **HTML Character Entities**: All HTML character entities were removed from the 'reviewText' column.
- **Text-Based Facial Expressions**: Text-based facial expressions or emoticons, such as ), were systematically removed.
- **Lowercasing**: In the final step of the cleaning process, the entire review text was converted to lowercase. This standardization facilitates uniformity in subsequent analyses and processing steps.

### C. Feature Generation

Feature generation was employed to create a sentiment label for each review, forming the basis for subsequent sentiment analysis In the label column generation process, we categorized reviews into three sentiment classes: Positive, Neutral, and Negative, based on the overall rating. The thresholds used for classification were set as follows:

- Reviews with an overall rating of 4 and above were labeled as Positive (Class 2.0).
- Reviews with an overall rating between 2 (exclusive) and 4 were labeled as Neutral (Class 1.0).
- Reviews with an overall rating of 2 and below were labeled as Negative (Class 0.0).

## V. ASPECT EXTRACTION

In our project, gaining insights into customer opinions about diverse product aspects is a pivotal step as Aspect extraction helps Brand owners/ dealers figure out the main themes and features that customers talk about. This is really helpful because it tells them what's good about products, what are the deciding factors, the general sentiment around them, and what needs improvement. By knowing these aspects, we can make our products better and make smart decisions based on what customers are saying. The dataset encompasses a wide array of products spanning different genres. As there are no predefined aspect labels for each product, proceeding with supervised aspect opinion extraction becomes challenging. To overcome this, we leveraged the dependency parser tree provided by Python's spaCy package.

Our methodology involved extracting pairs of words based on specific syntactic dependency paths. These paths served as indicators of how customers expressed their opinions regarding various product aspects. To streamline our analysis, we initiated by filtering out irrelevant reviews, retaining only those with a minimum of five votes – a criterion for relevance. Moreover to fully understand the intent of this aspect extraction we removed records of products about books as it mostly involved reviews about the content of the book like character description etc, which seemed irrelevant to our work.

Following this, we leveraged the dependency parser tree to establish seven rules that guided the extraction of aspects for each product. Our methodology allowed us to extract the top five aspects of every product, coupled with the accumulation of diverse opinions on various aspects from multiple reviewers. Moreover, we also accumulated the sentiment attached to each opinion/modifier using VADER Sentiment from the NLTK library.

Prior to aspect extraction, we aggregated all product reviews to generate a comprehensive reviews corpora, which served as the foundation for our aspect extraction process. The challenge arose when the size of the review text became substantial due to aggregation, surpassing the maximum length limit of 10,000,000 imposed by the spaCy NLP library. To address this, we strategically divided the reviews corpora into manageable chunks, each conforming to the specified maximum length restriction. This segmentation facilitated more efficient and effective processing of the extensive review data.

**Spacy for Aspect Extraction:** Spacy's dependency parsing serves as a foundation for identifying aspects (A) and their corresponding sentiment modifiers (M) which are customer opinions. A series of linguistic rules are systematically applied to extract these pairs

- Rule 1: Adjectival Modifier:
  - Description: This rule identifies aspects by capturing the relationship between adjectives and their corresponding nouns. It ensures that adjectives describing product features are correctly associated.
  - Example: In the phrase "Great product," the rule recognizes "product" as the aspect modified by the

adjective "great."

- Rule 2: Direct Object:
  - Description: This rule focuses on identifying aspects and sentiment modifiers by analyzing the direct object relationship. It captures nouns and adjectives that form the core components of a reviewer's opinion.
  - Example: In the sentence "I like the durable build," the rule extracts "build" as the aspect modified by the adjective "durable."

- Rule 3: : Adjectival Complement:
  - Description: By examining the adjectival complement relationship, this rule captures aspects and their associated sentiment modifiers. It is particularly useful in scenarios where adjectives modify nouns expressing opinions.
  - Example: In the statement "The sound of the speakers would be better," the rule identifies "sound" as the aspect enhanced by the adjective "better."

- Rule 4: Adverbial Modifier to a Passive Verb:
  - Description: This rule identifies aspects and sentiment modifiers when a verb is accompanied by an adverbial modifier. It ensures that the modifier is appropriately linked to the product feature.
  - Example: In the phrase "Well-designed product," the rule recognizes "product" as the aspect enhanced by the adverbial modifier "well-designed."

- Rule 5: Complement of a Copular Verb:
  - Description: Focusing on the complement of copular verbs, this rule captures aspects and sentiment modifiers. It is especially relevant when expressions like "is," "seems," or "appears" connect the subject and the complement.
  - Example: In the assertion "This is garbage," the rule identifies "this" as the subject and "garbage" as the complement representing the aspect.

- Rule 6: Interjections (INTJ):
  - Description: This rule targets interjections and captures aspects and sentiment modifiers in such expressions. It ensures that even in informal language, meaningful product aspects and sentiments are extracted.
  - Example: In the phrase "It's ok," the rule recognizes "it" as the aspect with "ok" as the sentiment modifier.

- Rule 7: Attribute (ATTR):
  - Description: Focused on attributes, this rule identifies aspects and sentiment modifiers by examining the syntactic relationship of attributes. It is valuable in scenarios where certain words attribute characteristics to products.
  - Example: In the statement "This seems interesting," the rule identifies "this" as the subject and "interesting" as the attribute expressing the aspect.

These rules guide us in understanding the most relevant aspects of the product. Figure **??** shows some of the rules formulated using the dependency tree
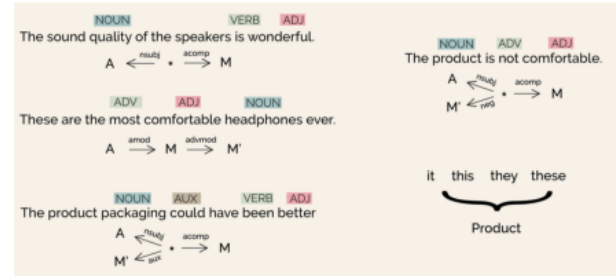


Fig. 2. Some of the rules formulated using Spacy's dependency tree

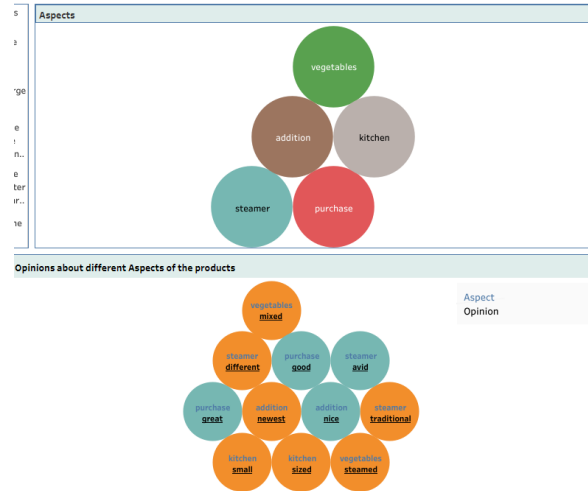Figure 3 shows the Top 5 aspects for each product and their opinions.



Fig. 3. Top 5 Aspects and their Opinion

## VI. SENTIMENT ANALYSIS USING ML/DL MODELS:

*1) Sentiment Classification for Product Reviews:* In the initial phase, we evaluated the results and did a sentiment analysis of product reviews using distinct machine learning models: Naive Bayes, Logistic Regression, and BERT. The models were evaluated using the following metrics: accuracy, precision, and recall. Prior to applying these models, the 'sentiment' column was transformed into a 'label' column where the sentiments were encoded as 0 for negative, 1 for neutral, and 2 for positive.

- **Naive Bayes:** Naive Bayes is a probabilistic classifier based on Bayes' theorem. It is a simple and efficient algorithm that can be effective for sentiment classification. Our model utilizes a hashing TF-IDF approach to extract features from product reviews. This technique combines the efficiency of hashing with the informative power of TF-IDF weighting, allowing us to capture the most salient

aspects of the reviews while maintaining computational efficiency.

- **Logistic regression:** Logistic regression is a statistical model that predicts the probability of a binary outcome. It works by fitting a logistic function to the data, which is a sigmoid function that squashes the input values between 0 and 1. The output of the logistic function is interpreted as the probability that the outcome is positive. Similar to Naive Bayes, we utilize a hashing TF-IDF approach to extract features. Logistic regression is a more powerful model than Naive Bayes, as it can capture some of the relationships between words in a review.

- **BERT-base**: BERT-base, which stands for Bidirectional Encoder Representations from Transformers, is a pre-trained natural language processing (NLP) model that has been shown to be effective in sentiment analysis. It is well-suited for processing sequential data, such as text. We fine-tuned the BERT model using our product review dataset. Due to the lack of official PySpark support for BERT, we employed Torch BERT to train a sentiment classification model on a subset of the product review data. The trained model was then effectively utilized to classify the entire dataset. BERT is well-suited for our task because it is efficient and can handle long-range dependencies in text.

After testing all the models, the following table in Figure 4 demonstrates the summary of experimental results along with the respective features of each model. In terms of accuracy, Logistic Regression with Hashing TF-IDF appears to perform the best among the three models, achieving the highest accuracy of 86%. BERT was expected to perform better than Naive Bayes and Logistic Regression for sentiment analysis of the product reviews as it can capture the nuances of human language and the relationships between words in a review. However, our experiments revealed that it performed slightly lower than Naive Bayes for sentiment classification of product reviews. This observation can be attributed to the limitations imposed by the chosen maximum sequence length of 512 tokens. While this length is sufficient for processing most product reviews, it may lead to the truncation of longer reviews, potentially discarding valuable contextual information that could improve sentiment prediction accuracy. Additionally, the computational overhead associated with further processing techniques was deemed excessive for the scope of this project. Therefore, a simpler approach was adopted to maintain computational efficiency, resulting in a slight trade-off in accuracy compared to Naive Bayes.

*2) Data Aggregation and Manipulation::* Aggregated metrics form the foundation for data-driven decision-making. Sellers and brands can use these metrics to optimize their strategies, allocate resources efficiently, and respond to chang-

| Model | Features | Accuracy | Precision | Recall |
|---|---|---|---|---|
| Naive-Bayes | Hashing TF-IDF(Tokenized word +Hashed Tokens +IDF-weighted features) | 78% | 83% | 78% |
| Logistic-Regression | Hashing TF-IDF(Tokenized words + Hashed Tokens +IDF-weighted features) | 86% | 83% | 83% |
| BERT-base | Text + BERT Embedding | 82% | 91% | 83% |

Fig. 4. Summary table of model results in different scenarios

ing market dynamics. For example, these average ratings and sales data to identify successful products or address quality issues.

In this phase of our project, we focused on organizing and transforming our raw data to extract meaningful insights. Here's a breakdown of the key steps:

- **Data Transformation:** We transformed our data to make it more useful. For instance, we converted the 'review-Time' to a date format, creating a new 'reviewDate' column. Additionally, we derived new columns representing the year, quarter, and month from the 'reviewDate'.

- **Time Calculations:** We calculated relevant time periods such as the previous year, quarter, and month based on the maximum date found in our data, stored as 'current_date'.

- **Rating Aggregation:** We aggregated our data to calculate average ratings for the previous year, month, and quarter. This process helped us understand the overall satisfaction levels expressed in the reviews.

- **Sales Aggregation:** To comprehend the sales trends, we computed total sales, average sales, and the number of sales for the previous year, month, and quarter.

- **Reviews Aggregation:** We analyzed the review patterns by determining the number of reviews for the previous year, month, and quarter.

- **Data Grouping and Metrics Calculation:** Grouping our data by product ('asin') and brand, we calculated various metrics, including average rating, total ratings, sales, and the number of reviews, offering a comprehensive view of each product's performance which is shown in the figure 5.

Fig. 5. Aggregated Data

In Our project, Tableau is used to make simple and interactive graphs that can help to understand complicated data better. In the Tableau figure 6 and figure 7, the top 10 products (Asin) and brands are showcased, with sales serving as the basis for ranking. The sales volume is represented by the height of each bar, allowing for a

visual comparison of the most successful products and brands. Notably, the brand ranking reveals Sony at the forefront, trailed by Samsung and Logitech in the second and third positions, respectively.
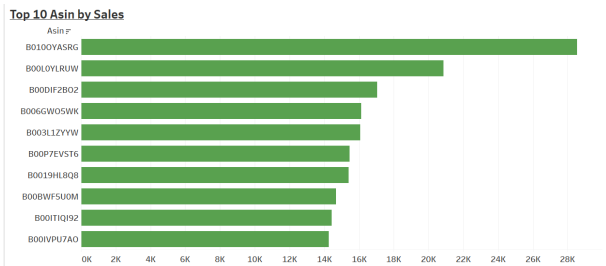


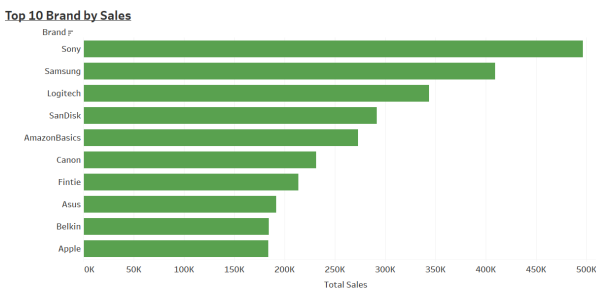Fig. 6. Top 10 Product ID's by Sales



Fig. 7. Top 10 brand by Sales

## VII. FIND SIMILAR ITEMS:

To find the Top 5 most similar products, we focused on using the Headphones category as the number of products in Electronics metadata is huge, and even using spark frameworks, we were getting out of memory error. So we first identified all the products in the metadata file which contains the words headphone, headphones, Headphone, or Headphones. Then as described above we aggregated the metadata with the extracted data such as to get num and total average reviews in the last year. After getting this information, we created a new column that had the product name, product category, and product description. We considered only active products, which had average ratings and average ratings in the last month.

We tried 2 methods to get the word embedings: TF-IDF and BERT. We chose the use Bert embeddings for clustering, as they gave better results. Then we normalized the Bert embedding and number of reviews to give equal weightage in clustering. We chose agglomerative clustering as the clustering mechanism, as a product can fall into various hierarchies, for example, headphone-¿ over the ear-¿ noise canceling, etc. Then for clustering, we clustered the products belonging to the same category since we don't want products like headphone extenders to be considered while clustering headphones. The normalized distance measure for clustering was chosen to be 10. We also removed those products, which

had less than 2 products in the same product category.

After clustering, the top 5 similar products for each category were chosen based on the minimum distance to the points in the same cluster. The figure below shows the clustering hierarchy for the category: over the ear headphones
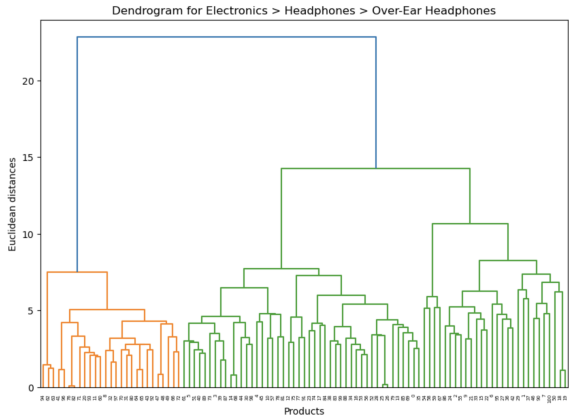


Fig. 8. Clustering Results of Over the year Headphones

## VIII. COMPETITOR ANALYSIS

Our project's competitor analysis offers Amazon vendors valuable insights into key performance indicators (KPIs) for both themselves and their competitors. Leveraging Tableau, we've created an interactive dashboard to provide vendors with a comprehensive view of this comparative data. Following is the link to our Tableau dashboard. https://tinyurl.com/tableureport

- **Performance of Similar Items:** The dashboard lets vendors see important info like last year's reviews, average ratings from the past month, and how their competitors did last year. This helps them understand better about the competitors which the shown in the figure 9.



Fig. 9. Competitor Product Analysis for top 5 Most Similar Products

- **Aspects of Similar Products:** This perspective highlights the top five aspects of each product along with corresponding opinions shown in the figure 10. This helps vendors by showing them the most important things about each product and what people think about them. It's a
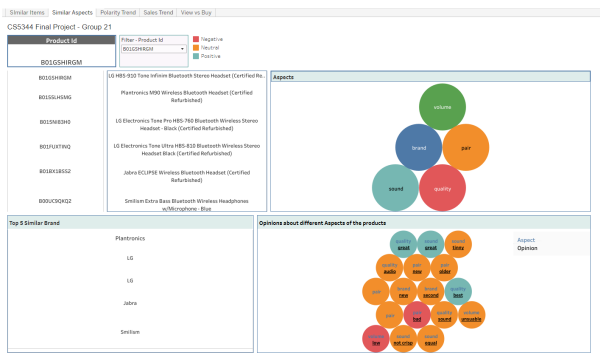
Fig. 10. Opinions about different Aspects of the products

quick way for vendors to understand what's working well and what might need improvement.

- **Polarity:** This dashboard shows how people feel about the product over time for vendors. Most of the reviews are neutral, with a good number being positive. Figure 11 shows a noticeable amount of negative reviews suggesting there might be some fake reviews. The line graph helps us see how opinions change over time.
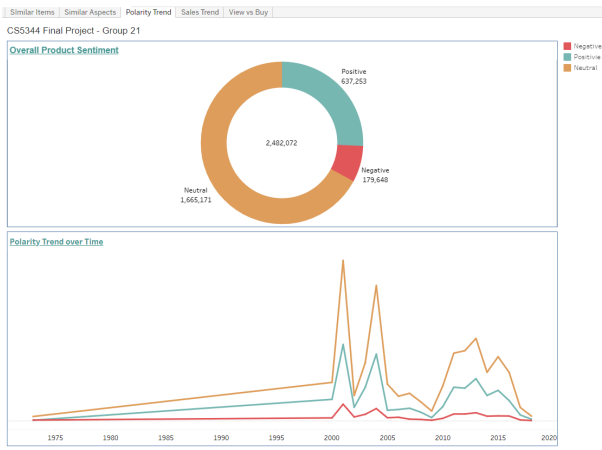


Fig. 11. Polarity analysis of the product over time

- **Sales Trend:** This dashboard assists vendors in studying the sales trends for each category and brand over time. Additionally, it provides insights into the correlation between product views and purchases, revealing valuable patterns that are visible in figure 12.

## REFERENCES

[1] Soujanya Poria, , Nir Ofek, Alexander Gelbukh, Amir Hussain, Lior Rokach, "Dependency Tree-based Rules for Concept-Level Aspect-based Sentiment Analysis"

[2] Wouter Bancken, Daniele Alfarone and Jesse Davis, "Automatically Detecting and Rating Product Aspects from Textual Customer Reviews"

[3] https://github.com/ishikaarora/Aspect-Sentiment-Analysis-on-Amazon-Reviews

[4] https://www.kaggle.com/code/phiitm/aspect-based-sentiment-analysis

[5] Devlin, J., Chang, M. W., Lee, K., & Toutanova, K. (2018). BERT: Pre-training of deep bidirectional transformers for language understanding. arXiv preprint arXiv:1810.04805
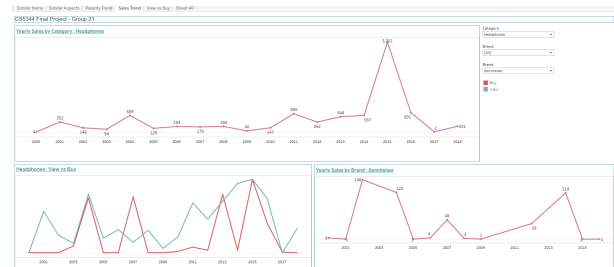
[6] Bollegala, D., D. Weir, and J. Carroll. Using multiple sources to construct a sentiment-sensitive thesaurus for cross-domain sentiment classification. In Proceedings of the 49th Annual Meeting of the Association for Computational Linguistics (ACL-2011), 2011.

[7] Subhabrata Mukherjee. Sentiment analysis of reviews.

[8] Bingwei Liu, Erik Blasch, Yu Chen, Dan Shen, and Genshe Chen. Scalable sentiment classification for big data analysis using naive bayes classifier. In Big Data, 2013 IEEE International Conference on, pages 99–104. IEEE, 2013.

[9] Sebastian Raschka. Naive bayes and text classification i-introduction and theory. arXiv preprint arXiv:1410.5329, 2014.

[10] Bing Liu. Sentiment analysis: Mining opinions, sentiments, and emotions. Cambridge University Press, 2015.

[11] Dataset Refence : https://jmcauley.ucsd.edu/data/amazon/

Fig. 12. Anlysis of Sales over Time