

Amazon US Customer Reviews

TINGHUI XU, OUYANG XU, BOWEN TIAN, YIJIN GUAN, YIFAN DU

Introduction

Amazon is the online retailer with the largest variety of products in the world. It's meaningful for either customers or the business owners to know more about the reviews of the items. Our main goal is to explore which aspects are mostly mentioned. This can help sellers improve their stars. So we use CHTC to find the relationship between the high frequency words and rating stars.

Data Description

Source : Kaggle

<https://www.kaggle.com/cynthiarempel/amazon-us-customer-reviews-dataset>

Reviews of products of six categories (6 tsv files):

- Camera
- Book
- Digital Video Download
- Electronics
- Mobile Apps
- Digital Ebook

Variable Names	Description
marketplace	2 letter country code of the marketplace where the review was written.
customer_id	Random identifier that can be used to aggregate reviews written by a single author.
review_id	The unique ID of the review.
productid	The unique Product ID the review pertains to. In the multilingual dataset the reviews for the same product in different countries can be grouped by the same productid.
product_parent	Random identifier that can be used to aggregate reviews for the same product.
product_title	Title of the product.
product_category	Broad product category that can be used to group reviews.
star_rating	The 1-5 star rating of the review.
helpful_votes	Number of helpful votes.
total_votes	Number of total votes the review received.
vine	Review was written as part of the Vine program.
verified_purchase	The review is on a verified purchase.
review_headline	The title of the review.
review_body	The review text.
review_date	The date the review was written.

Data Processing

- **split.sh**: to split 6 tsv files into tens of small tsv files, each 100MB.
- **checklist.sh**: to find all the files that split from raw data and compile the file names to a list.
- **word_freq.R/.sh**: to calculate on these files in parallel and count the word frequency.
- If without parallel calculation, the word segmentation operation will take up a lot of memory and run for a long time

Data Processing

Example:

- Split the camera.tsv (1.1GB) into 11 small tsv files.
- Run `word_freq_array_camera.sh` to launch 11 small jobs. Each job does the tokenization and lemmatization to each csv file, and then calculates the frequency of each word, returning a csv file with Column word, star_ratings, and frequency.

Text Preprocessing

- Remove all of the punctuation and some html elements like `
` and `"`.
- Turn each and every letter to lower-case letter and do lemmatization, which can transform words like 'swims', 'swam', 'swimming' to 'swim'.
- Do tokenization to separate the strings into single words.
- Remove all the stop words to get our word boxes.

Computation for words frequency

- Group our word boxes to get words frequency results.
- Extract top 2000 words in each split file and merge them into one .csv file.
- Select all the nouns from the words we get.
- Do these on each of the 6 sub files.

Computation for words frequency with star rating

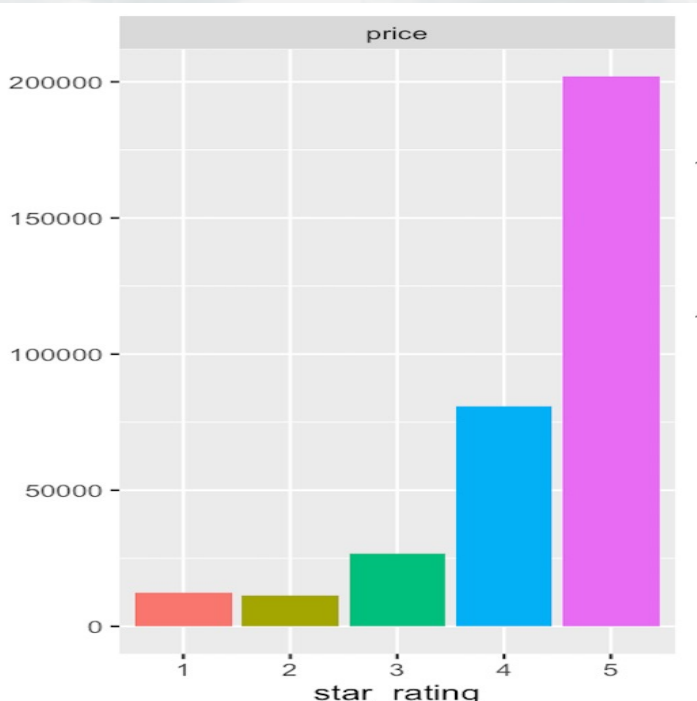
Part of our original result

	product_category	star_rating	word.stem	frequency
13	Camera	star_1	accident	453
14	Camera	star_1	accommod	201
15	Camera	star_1	accompani	8
16	Camera	star_1	accomplish	26
17	Camera	star_1	account	1152

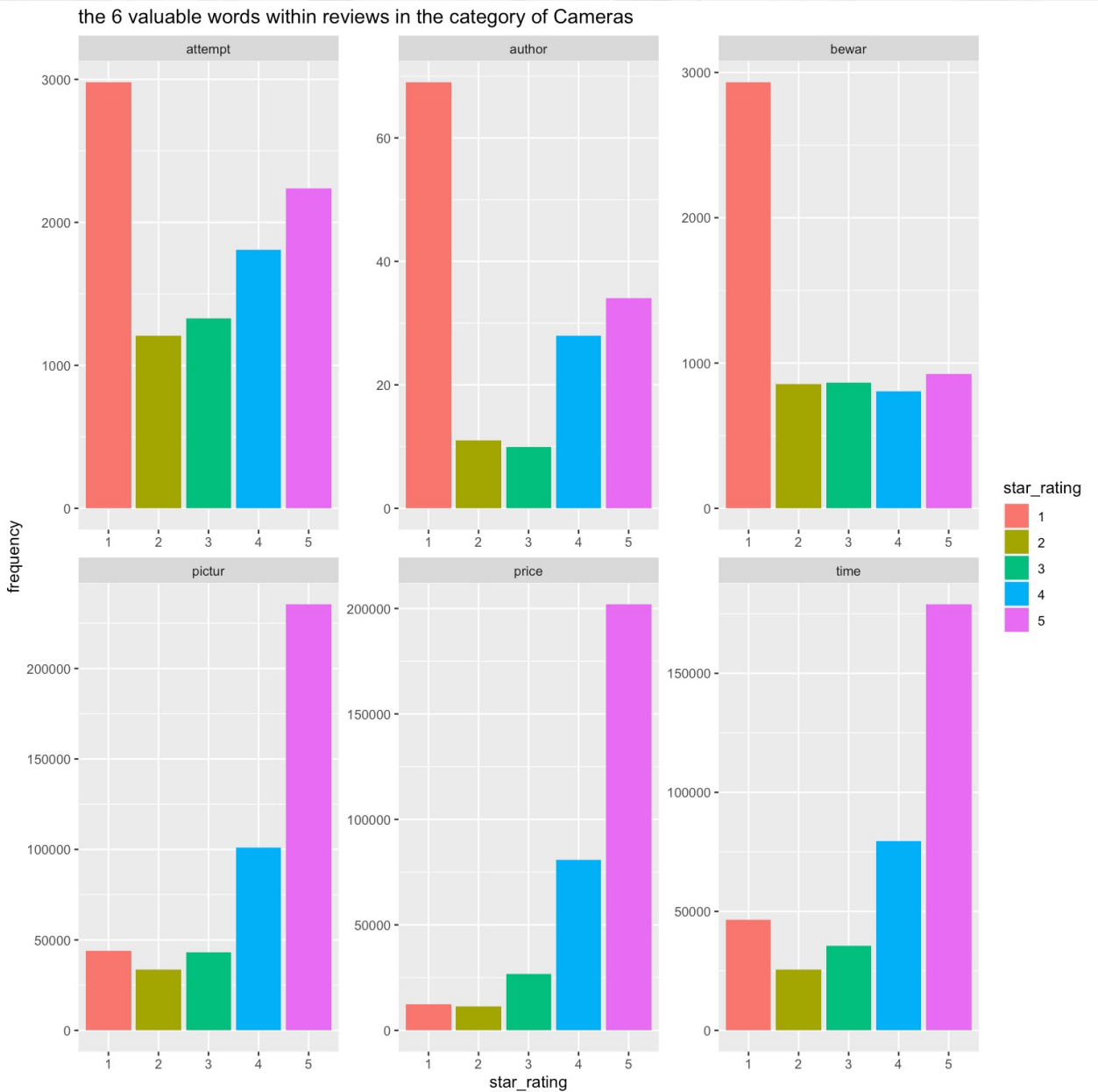
Word frequency v.s. Star rating plot

Part of our transformed result structure

	word.stem	star_1	star_2	star_3	star_4	star_5
13	accident	453	391	786	2085	4200
14	accommod	201	188	457	1190	2203
15	accompani	8	4	10	43	72
16	accomplish	26	40	60	181	339
17	account	1152	563	855	1717	2731

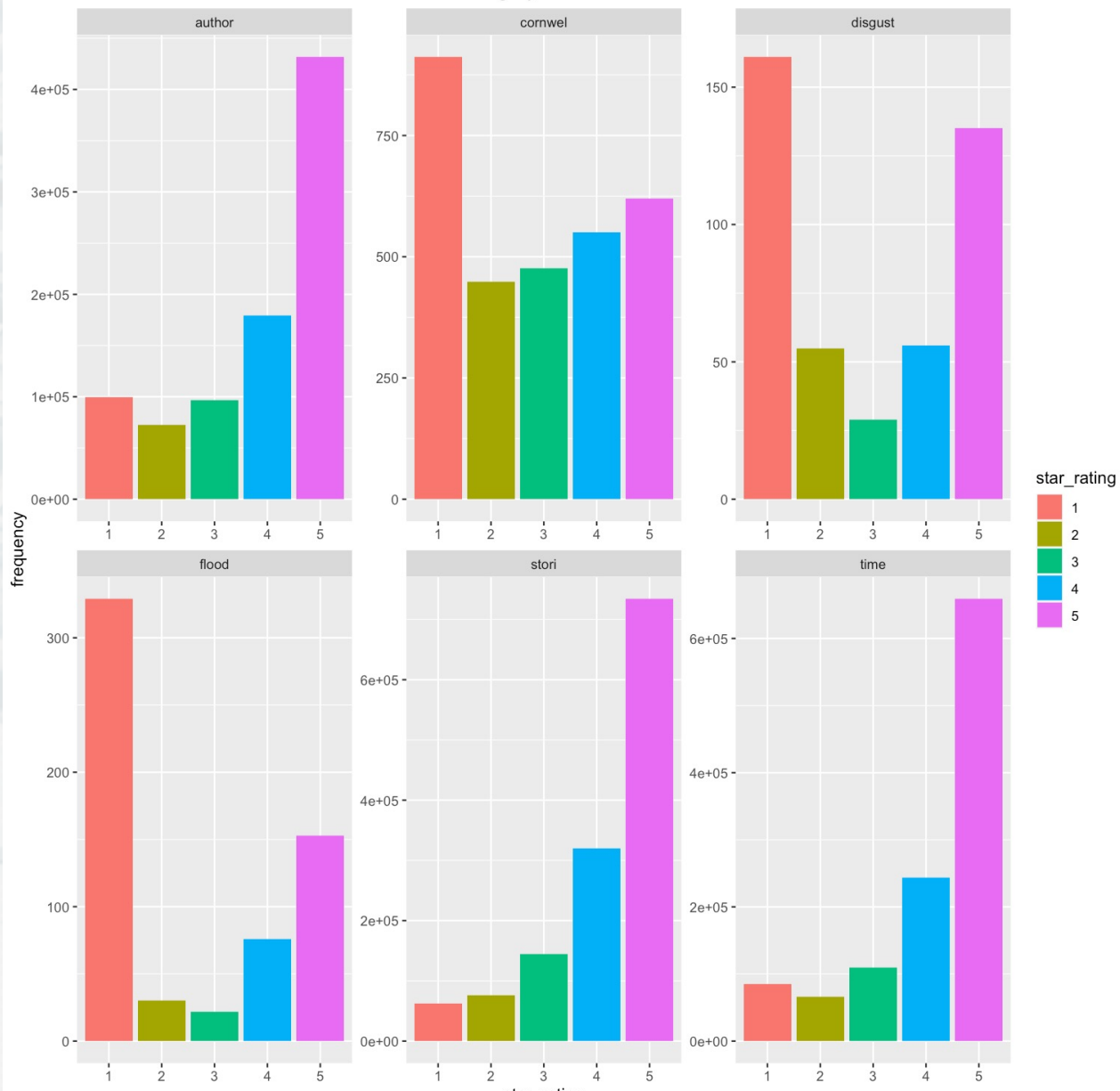


Conclusion

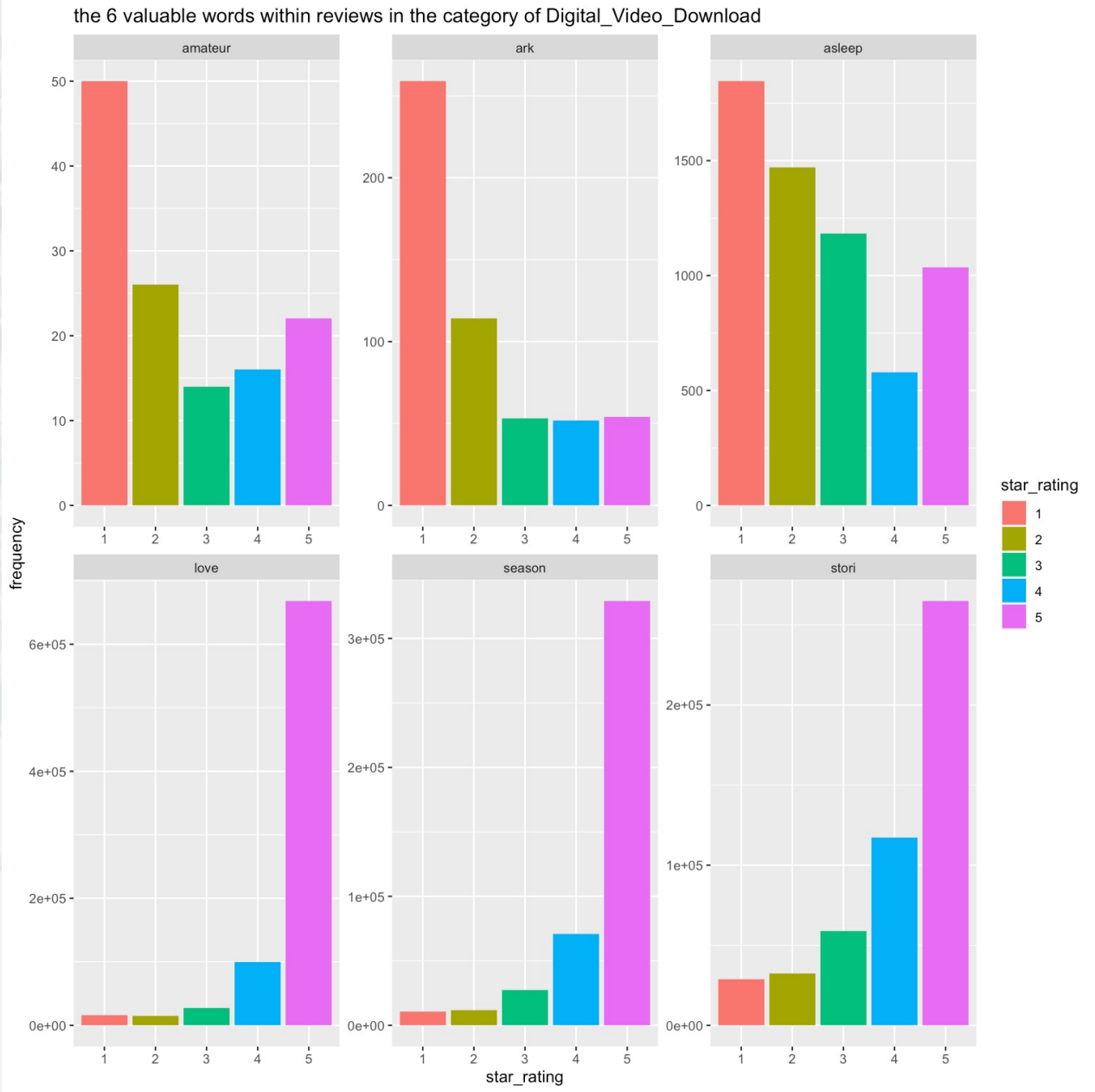


Conclusion

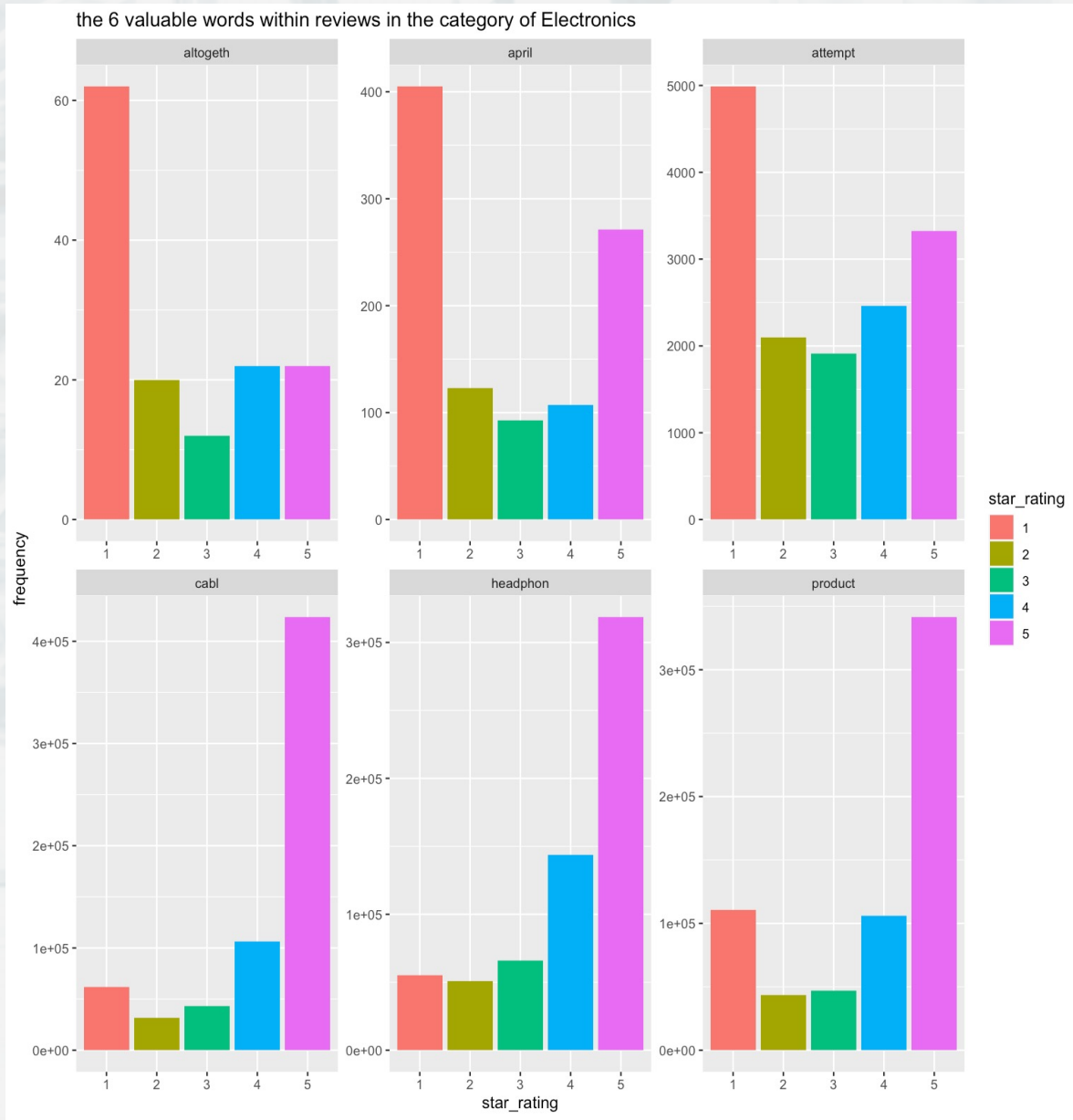
the 6 valuable words within reviews in the category of Books



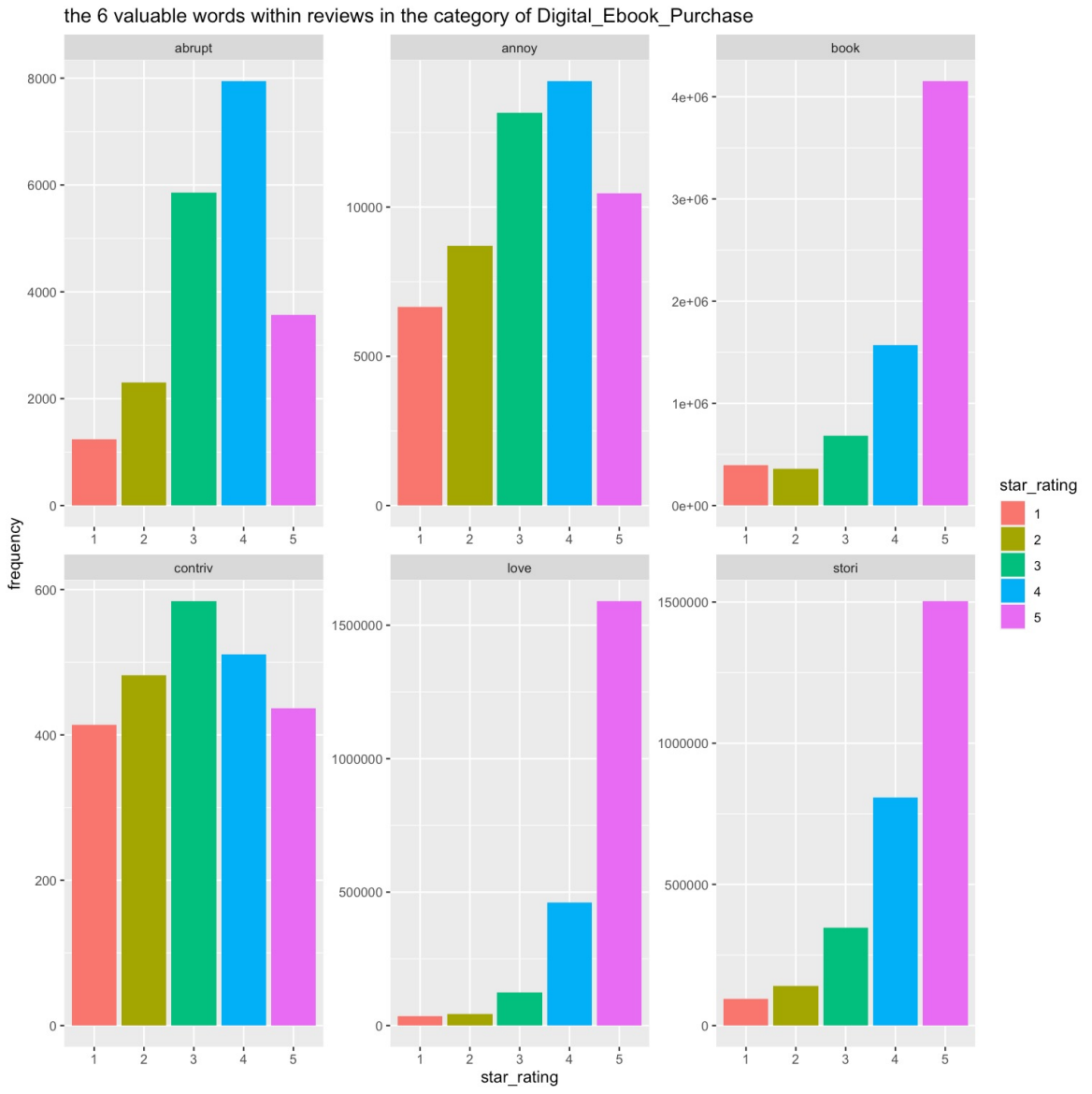
Conclusion



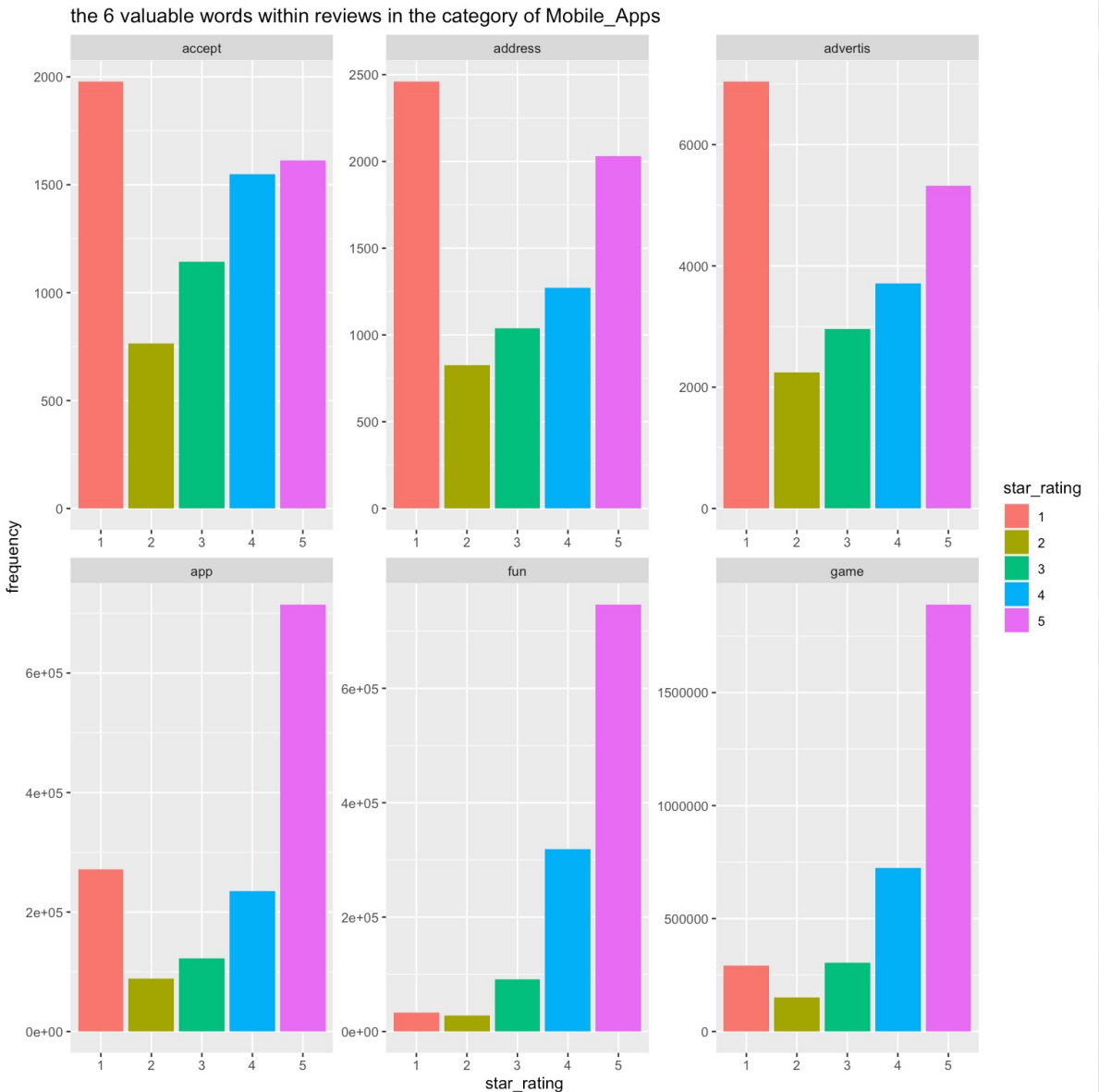
Conclusion



Conclusion



Conclusion





THANKS