

# Analyzing Amazon Dataset Proposing Better Metrics Unstructured Data Management

Harish Bharatham  
(worked with 3 other teammates)

# TOPIC INTRODUCTION

## What is Amazon

Amazon is one of the largest online shopping websites in the world. The site is widely known for its wide selection of books, although the site has expanded to sell electronics, music, furniture, and apparel. Similar to eBay, users can also purchase and sell items using Amazon's online marketplace system. Amazon was founded in 1995 by Jeff Bezos and is based out of Seattle, Washington.

## Amazon Reviews

Reviews by ordinary people have become an essential mechanism for selling almost anything online; they are used for resorts, dermatologists, neighborhood restaurants, high-fashion boutiques, churches, parks, astrologers and healers — not to mention products like garbage pails, tweezers, spa slippers and cases for tablet computers. In many situations, these reviews are supplanting the marketing department, the press agent, advertisements, word of mouth and the professional critique. But not just any kind of review will do. They have to be somewhere between enthusiastic and ecstatic. — NY Times

Amazon allows users to submit reviews to the web page of each product. Reviewers must rate the product on a rating scale from one to five stars. Amazon provides a badging option for reviewers which indicate the real name of the reviewer (based on confirmation of a credit card account) or which indicate that the reviewer is one of the top reviewers by popularity. Customers may comment or vote on the reviews, indicating whether they found a review helpful to them. If a review is given enough "helpful" hits, it appears on the front page of the product. In 2010, Amazon was reported as being the largest single source of Internet consumer reviews.

However, off late Amazon has been receiving a lot of backlash for having fake reviews on the website. The sellers seek shoppers on dozens of networks, including Amazon Review Club and Amazon Reviewers Group, to give glowing feedback in exchange for money or other compensation. The practice artificially inflates the ranking of thousands of products, experts say, misleading consumers. This is one aspect we address in this project by using SAS Text Miner.

In this project, we have focussed on the recommendation/ review system of Amazon.com. We have utilized the knowledge gained in this course to create a pseudo database schema in MongoDB as well as Neo4j (in the form of graph data model). Next, we have presented an overview of the dataset available by executing MongoDB queries on the dataset to get a general overview of the data. Based on these queries (but not limited to them) we have come up with some findings/ patterns in the data and have executed MongoDB queries for them as well. We have then attempted to formulate some interesting analytics metrics for the Amazon website specifically related to the product reviews and reviewers. We have also identified the queries required to run for identifying helpfulness of reviews and also identifying common set of reviewers who review a common set of products. Lastly, we have used SAS text miner to identify fake reviews among the dataset and then we have designed a template of the webpage using the analytics metrics.

# DESIGN

## 1. MongoDB Pseudo Database Schema

Before defining the pseudo database, we will define the constraints and indexes that we will be using for the Amazon dataset. Following the constraints defined:

- a. Constraint: The overall of a review must be either 1, 2, 3, 4 or 5.

```
db.createCollection("review",
{
  validator:
  { "overall": { $in: [1,2,3,4,5] }
  }
})
```

- b. Constraint: The price of a product must be greater than \$0.

```
db.product.aggregate("product",
{
  validator:
  { "price" : {$gt : 0}}
})
})
```

- c. Index: We can create a textual index on the reviewText of reviews.

```
db.review.createIndex( { reviewText: "text" } )
```

## Pseudo Database Schema

For the Amazon dataset provided, we propose the following pseudo database schema.

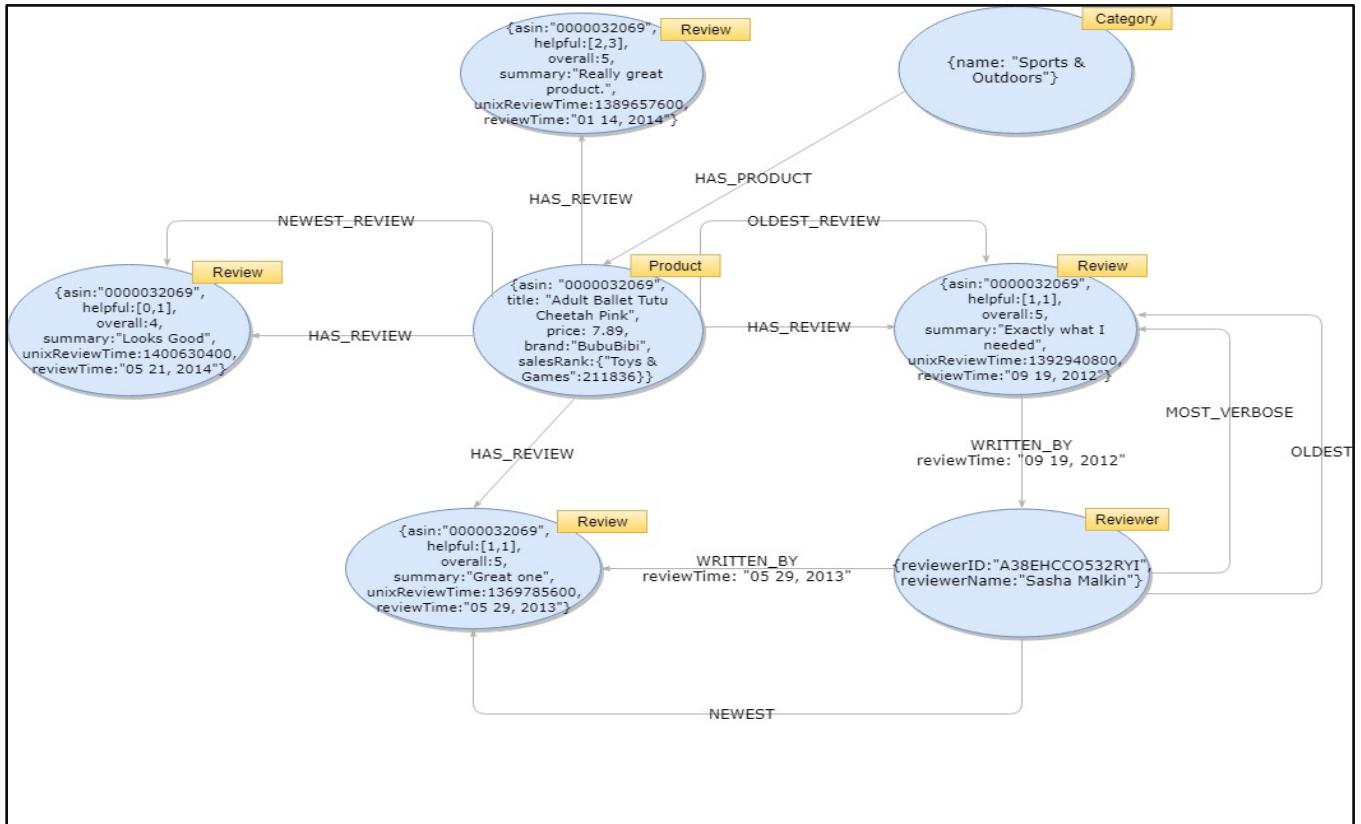
### Embedding Reviews in Products

Placing orders and recommendation system for products are the two pillars of the Amazon website. Since, we are dealing with reviews in this project, we think that reviews of a particular product should be embedded within the specific product document inside the product collection. Whenever a reviewer clicks on a product for purchasing, the product page opens and all reviews for that product are displayed along with a review histogram. Therefore, queries involving retrieving the reviews for a product will run faster and more efficiently since we will get all the data in the same collection. We will not need to reference multiple collections resulting in increase in the running time of the query. Basically, we will be having an array of the review documents inside the corresponding product document. Below is a sample of the product collection after embedding reviews.

```
{
  "asin" : "0001048791",
  "salesRank" : { "Books" : 6334800 },
```

```
"imUrl" : "http://ecx.images-amazon.com/images/I/51MKPOT4DBL.jpg",
"categories" : [ [ "Books" ] ], "title" : "The Crucible: Performed by Stuart Pankin, Jerome Dempsey & Cast"
reviews:[
{
    "reviewerID" : "AA8M6331NI1EN",
    "asin" : "B00000JBLH",
    "reviewerName" : "ZombieMom",
    "helpful" : [0,0],
    "reviewText" : "Bought this for my boss because he lost his. He loves this calculator & would not be caught without it. It really helps him with his day to day work & he is the CEO of my company. If the CEO swears by them then they must be a great little calculator.",
    "overall" : 5,
    "summary" : "Fast shipping & great price for this awesome calculator",
    "unixReviewTime" : 1384387200,
    "reviewTime" : "11 14, 2013",
    "category" : "Office"
},
{
    "reviewerID" : "AA8M6331NI1EN",
    "asin" : "B00000JBLH",
    "reviewerName" : "ZombieMom",
    "helpful" : [ 0, 0 ],
    "reviewText" : "Bought this for my boss because he lost his. He loves this calculator & would not be caught without it. It really helps him with his day to day work & he is the CEO of my company. If the CEO swears by them then they must be a great little calculator.",
    "overall" : 5,
    "summary" : "Fast shipping & great price for this awesome calculator",
    "unixReviewTime" : 1384387200,
    "reviewTime" : "11 14, 2013",
    "category" : "Office"
},
...  
]}
```

## 2. Neo4j Graph Data Model



### Description of Graph Data Model

We have identified the following entities based on our analysis of the Amazon dataset.

1. Product
2. Review
3. Reviewer
4. Category

We have taken these entities to be the Node Types which are represented in the labels for each node in our model. Based on the Cypher query language, each node consists the respective attributes/ properties separated by commas and enclosed in curly braces.

The relationships between the nodes are represented using directed lines. The direction of the relationship basically gives more semantics to the relationship and tells us as that we can only read the relationship in that direction. Few of the relationships also have some properties. For example, WRITTEN\_BY relationship has the date property which gives the date on which that particular user wrote that particular review.

Note: In our illustration of graph data model, we have listed only a few properties for each node due to space constraints and to keep the clarity of the model intact. We have mentioned all the properties of each node in the following paragraphs.

While preparing our graph data model, we have considered that the queries that we have written so far can be answered using our model. Few examples of the queries and their corresponding relationships in our data model are:

1. Query: Identify date of oldest review per reviewer.

This query is answered by the OLDEST relationship between the User and Review nodes. The query groups the User collection on the basis of reviewerID. Now, taking each reviewerID we can easily find its oldest review using OLDEST relationship.

Note: We have further extended our model to identify the newest review given by a reviewer, which can be answered using the NEWEST relationship between Reviewer and Review nodes.

2. Query: Top 10 most verbose reviewers.

To answer this query, we have created a relationship MOST\_VERBOSE between User and Review nodes. This will give us the user who has written the longest text for a review. In future, if there comes another user with an even longer text for that particular review, this relationship will change to something like PREVIOUS or SECOND\_MOST\_VERBOSE (not shown in the model) and the relationship between the new User node and the Review node will become MOST\_VERBOSE.

In order to determine the oldest or newest reviewer for any product, we have the WRITTEN\_BY relationship with 'date' property. We can easily sort the reviews based on their date and then retrieve the corresponding reviewer.

The following tables explain the details about the nodes and relationships used in our data model and the thought process behind this approach.

#### Nodes:

Node No	Label	Properties	Justification
Node 1	Review	reviewerID asin reviewerName helpful reviewText overall summary unixReviewTime reviewTime category	Review is an entity that we want to model and store information about. We want to use our graph data model for queries involving reviews, so we have made it as a Node. <u>Example of query</u> : Display top most useful reviews per business.
Node 2	Product	asin related title price salesRank imUrl brand categories	Product is an entity that we want to model and store information about. We want to use our graph data model for queries involving Product, so we have made it as a Node. <u>Example of query</u> : Display average number of reviews per product.
Node 3	Reviewer	reviewerID reviewerName	Reviewer is an entity that we want to model and store information

			about. We want to use our graph data model for queries involving Reviewer, so we have made it as a Node. <u>Example of query</u> : Identify the reviewer who has been reviewing for the longest period of time.
Node 4	Category	category_name	Category is an entity that we want to model and store information about. We want to use our graph data model for queries involving Category, so we have made it as a Node. <u>Example of query</u> : Display average review ratings per category. Note: This is not a previously asked query but might be a useful statistic to determine the credibility of a business.

#### Relationships:

Relationship Name	Nodes Involved	Properties	Justification
HAS_PRODUCT (Category HAS_PRODUCT Product)	Category1 Product1		Our graph data model includes product and the category to which that product belongs. We can create this relationship by matching the category name (of Category node) to the categories (of Product node). During graph traversal, we can simply go to the particular category whose products we are looking for and then traverse through only those relationships which emerge from that Category node and end at any Product node. This way we will be able to access all those products that belong to that particular category.
WRITTEN_BY (Review WRITTEN_BY Reviewer)	Review3 Reviewer1	reviewTime	Each review is written by a reviewer. This relationship will have an attribute known as reviewTime which denotes the date on which the reviewer posted the review.
HAS REVIEW (Product HAS REVIEW Review)	Product1 Review2		Every product has multiple reviews. A product may or may not have any reviews.

NEWEST_REVIEW (Product NEWEST_REVIEW Review)	Product1 Review2		In the graph data model, every product will have newest review. This traversal will help us answer queries such as: Identify the reviewer with the newest review for each product.
OLDEST_REVIEW (Product OLDEST_REVIEW Review)	Product1 Review3		Every product has an oldest review. This traversal will help us answer queries such as: Identify the reviewer with the oldest review for each product.
OLDEST (Reviewer OLDEST Review)	Reviewer1 Review1		Review is written by a Reviewer. Each reviewer will have an oldest review which will be the first review that he/she writes. In order to find such reviews faster, we have a direct relationship between Reviewer and Review named, OLDEST. This way the user will not need to traverse through all Review nodes for a Reviewer and check the reviewTime property and then determine the oldest Review. This approach of traversal will help us answer queries such as: Identify the oldest review for a reviewer.
NEWEST (Reviewer NEWEST Review)	Reviewer1 Review2		Review is written by a Reviewer. Each reviewer will have a newest review which will be the most recent or last review that he/she writes. In order to find such reviews faster, we have a direct relationship between Reviewer and Review named, NEWEST. This way the user will not need to traverse through all Review nodes for a Reviewer and check the reviewTime property and then determine the newest Review. This approach of traversal will help us answer queries such as: Identify the newest review for a reviewer.
MOST_VERBOSE (Reviewer MOST_VERBOSE Review)	Reviewer1 Review3		In our graph data model, every reviewer will have one review for which he/she has written the maximum text. Every time the reviewer enters another review with more text, this new review becomes the most verbose review and the previous one becomes the second most verbose review. By modelling in this way, we can

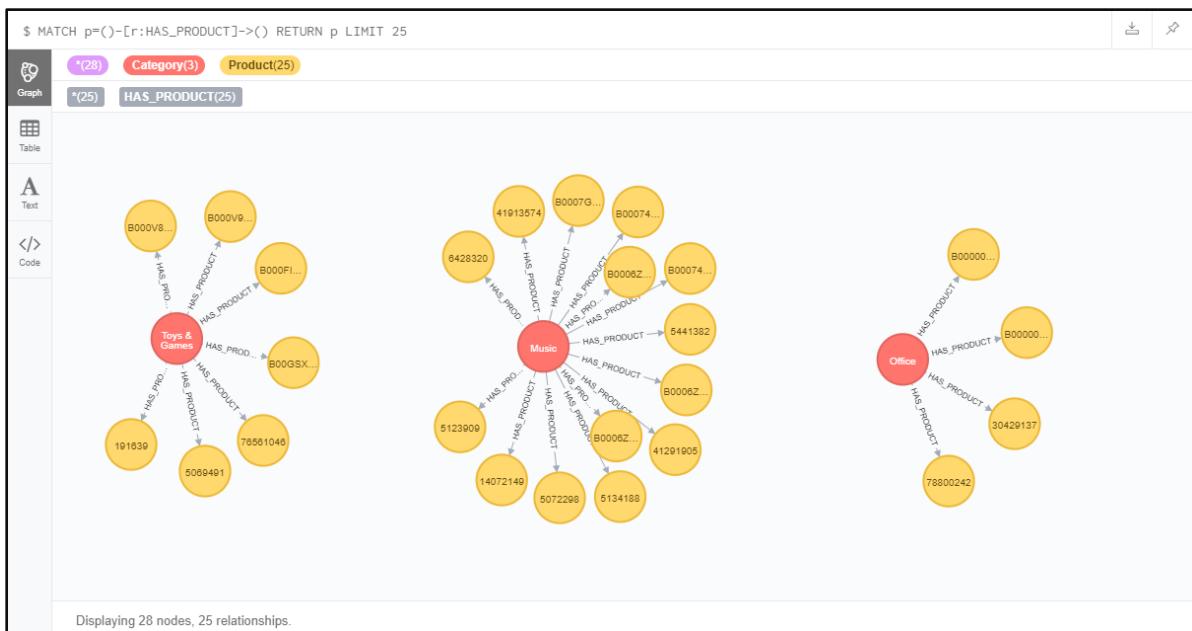
			answer queries such as: Most verbose reviewers.
--	--	--	---

Following are the CREATE queries for the three major relationships of our graph data model.

### 1. Category HAS\_PRODUCT Product

```
MATCH (category:Category), (product:Product)
WHERE product.categories CONTAINS category.category_name
CREATE (category)-[:HAS_PRODUCT]->(product)
```

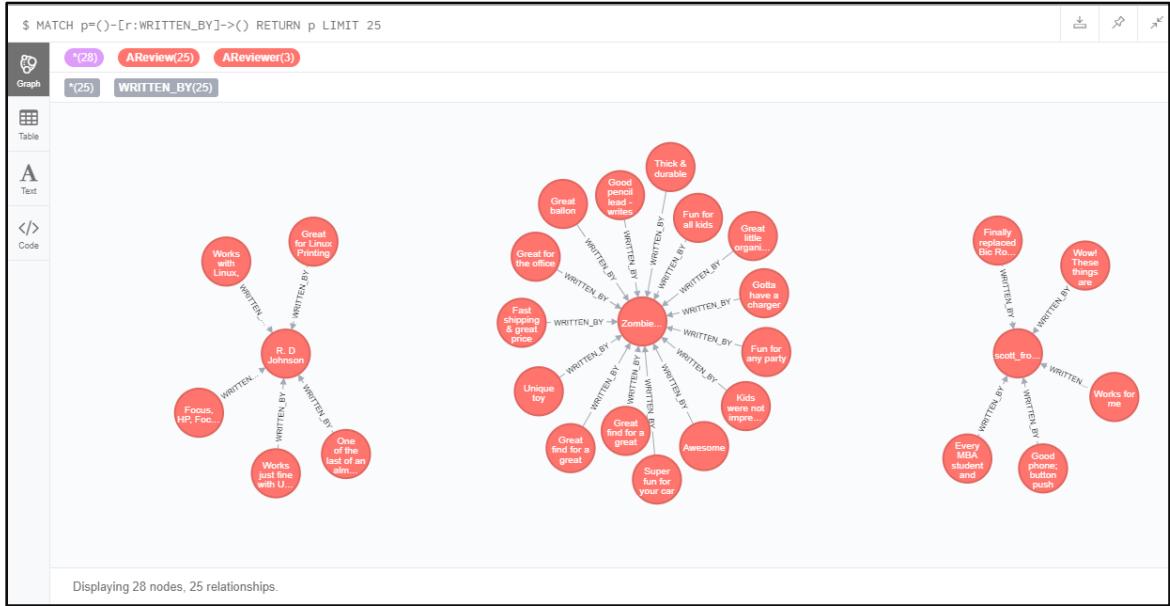
#### Representative Sample:



### 2. Review WRITTEN\_BY Reviewer

```
MATCH (review:Review), (reviewer:Reviewer)
WHERE review.reviewerID = reviewer.reviewerID
CREATE (review)-[:WRITTEN_BY]->(reviewer)
```

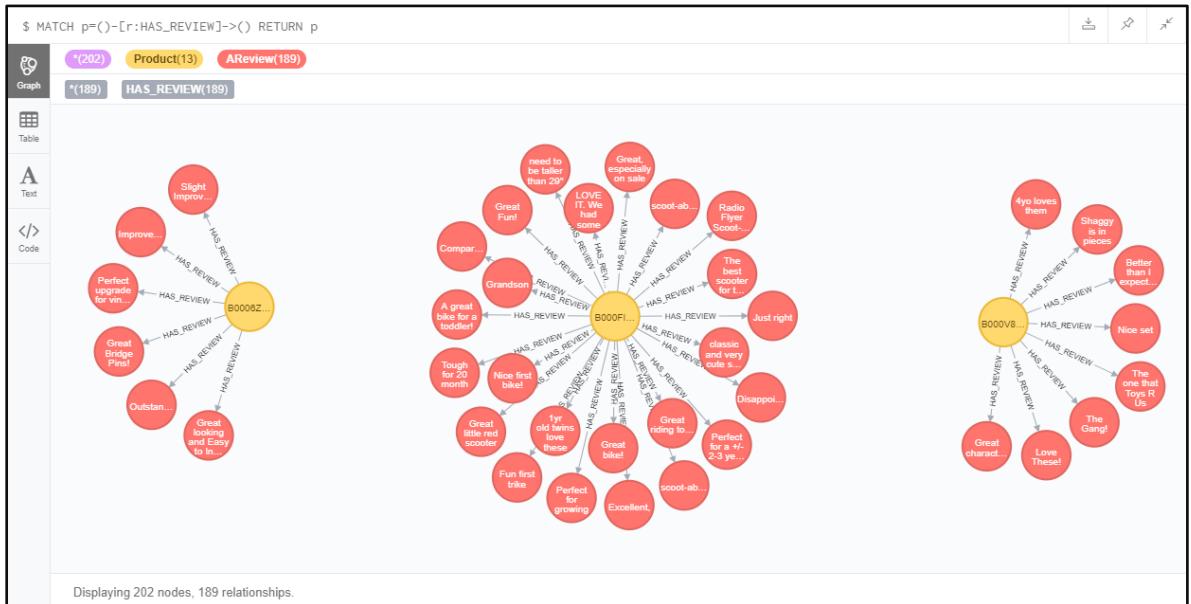
## Representative Sample:



### 3. Product HAS\_REVIEW Review

```
MATCH (review:Review), (product:Product)
WHERE review.asin = product.asin
CREATE (product)-[:HAS REVIEW]->(review)
```

## Representative Sample:



# BASIC UNDERSTANDING OF THE DATA

## 1. MongoDB Queries

Instead of going with all the reviews for all categories, we have selected the “Office Products”, “Musical Instruments” and “Toys & Games” categories in order to limit our dataset to a reasonable size. We have also downloaded and imported the product metadata which gives all information about a particular product.

After importing the JSON files for the reviews of each of the categories, we have added a “category” attribute to each of these collections. We have the value of this attribute as “Office”, “Music” and “Toys & Games” respectively. Next, we have merged these individual review collections into a single “review” collection. We have used the following queries to do so.

For adding category attribute:

```
db.office_products.update({},{$set:{'category': "Office"}},{multi:true})  
db.musical_instruments.update({},{$set:{'category': "Music"}},{multi:true})  
db.toys_and_games.update({},{$set:{'category': "Toys & Games"}},{multi:true})
```

For merging the review collections for each category into a single review collection:

```
db.office_products.copyTo("review")  
db.musical_instruments.copyTo("review")  
db.toys_and_games.copyTo("review")
```

Following are the queries we have executed to get an overview of the dataset under consideration:

- What is the overall number of products?

### Wording of query

```
db.product.distinct("asin").length
```

### Snapshot showing successful execution

```
MongoDB Enterprise > db.product.distinct("asin").length  
67291
```

### Summary of results

To find the total number of products in the Amazon dataset, we have simply applied the distinct function on the asin field and then calculated the length of the collection. The length function returns the total number of documents in that collection.

- b. What is the overall number of reviews?

**Wording of query**

```
db.review.count()
```

**Snapshot showing successful execution**

```
MongoDB Enterprise > db.review.count()
231116
```

**Summary of results**

In this query, we have used the count function on the review collection to calculate the total number of reviews in our dataset.

- c. What is the overall number of reviewers?

**Wording of query**

```
db.review.distinct("reviewerID").length
```

**Snapshot showing successful execution**

```
MongoDB Enterprise > db.review.distinct("reviewerID").length
23998
```

**Summary of results**

In order to get the total number of reviewers, we have identified the distinct reviewer IDs in the review collection. After applying the length function on this, we get the total number of these distinct reviewers. It is important to apply the distinct function here because a reviewer can review more than one products also. So doing a simple count on the review collection will not give the correct number of reviewers for the products.

- d. What is the overall number of reviews with ratings less than 3?

**Wording of query**

```
db.review.find({overall:{$lt:3}}, {_id:0, asin:1, summary:1, overall:1}).count()
```

**Snapshot showing successful execution**

```
MongoDB Enterprise > db.review.find({overall:{$lt:3}}, {_id:0, asin:1, summary:1, overall:1}).count()
14328
```

**Summary of results**

In order to get the total number of reviews having a rating less than 3, we have used overall:{\$lt:3} as the where criteria for the find function. We have then applied the count function on the result obtained to get the total number of matched reviews. We have displayed the asin, summary and overall fields in the final output and we have also suppressed the \_id field.

- e. What is the overall number of reviews with ratings more than 3?

#### Wording of query

```
db.review.find({overall:{$gt:3}}, {_id:0, asin:1, summary:1, overall:1}).count()
```

#### Snapshot showing successful execution

```
MongoDB Enterprise > db.review.find({overall:{$gt:3}}, {_id:0, asin:1, summary:1, overall:1}).count()
194599
```

#### Summary of results

In order to get the total number of reviews having a rating greater than 3, we have used overall:{\$gt:3} as the where criteria for the find function. We have then applied the count function on the result obtained to get the total number of matched reviews. We have displayed the asin, summary and overall fields in the final output and we have also suppressed the \_id field.

- f. What is the total number of reviews per product?

#### Wording of query

```
db.review.aggregate([{$group: {_id: "$asin", review_count: {$sum:1}}}] )
```

#### Snapshot showing successful execution

```
MongoDB Enterprise > db.review.aggregate([{$group: {_id: "$asin", review_count: {$sum:1}}}] )
{ "_id" : "B00JBJIVXGC", "review_count" : 13 }
{ "_id" : "B00FXKIG5I", "review_count" : 7 }
{ "_id" : "B00EPWAQS0", "review_count" : 5 }
{ "_id" : "B00C5B20QE", "review_count" : 5 }
{ "_id" : "B00BTGMI5O", "review_count" : 6 }
{ "_id" : "B00BLQTZDA", "review_count" : 6 }
{ "_id" : "B00BL6IYW8", "review_count" : 7 }
{ "_id" : "B00BHK2UHI", "review_count" : 6 }
{ "_id" : "B00B5OY8EE", "review_count" : 5 }
{ "_id" : "B00B1N06PO", "review_count" : 6 }
{ "_id" : "B00AZUAORE", "review_count" : 24 }
{ "_id" : "B00A716FB0", "review_count" : 11 }
{ "_id" : "B0095814U0", "review_count" : 5 }
{ "_id" : "B009MDMQJ4", "review_count" : 5 }
{ "_id" : "B0097WDJZA", "review_count" : 11 }
{ "_id" : "B0087UPSLQ", "review_count" : 8 }
{ "_id" : "B007T80GLK", "review_count" : 16 }
{ "_id" : "B007T8KI6M", "review_count" : 8 }
{ "_id" : "B007K3EGXW", "review_count" : 8 }
{ "_id" : "B007K3EGWI", "review_count" : 5 }
Type "it" for more
MongoDB Enterprise >
```

## Summary of results

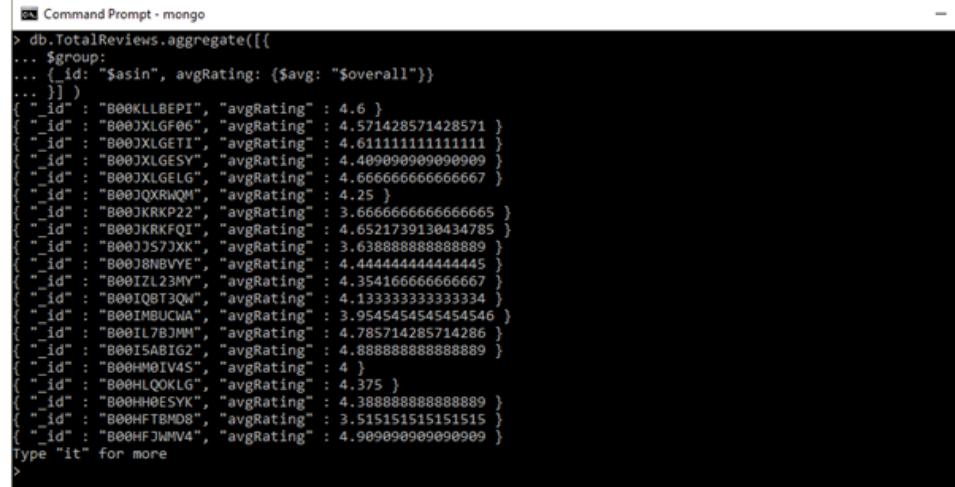
Here we are the total number of reviews for each product and for that first we are grouping all the reviews by asin (product ID) field and then calculating the total number of documents for each product by using \$sum function.

- g. What is the average rating per product?

## Wording of query

```
db.review.aggregate([{
  $group:
  {_id: "$asin", avgRating: {$avg: "$overall"}}
}])
```

## Snapshot showing successful execution



```
Command Prompt - mongo
> db.TotalReviews.aggregate([
... $group:
... {_id: "$asin", avgRating: {$avg: "$overall"}}
...])
{
  "_id": "B00KLLBEP1", "avgRating": 4.6
}
{
  "_id": "B00JXLGF06", "avgRating": 4.571428571428571
}
{
  "_id": "B00JXLGET1", "avgRating": 4.611111111111111
}
{
  "_id": "B00JXLGESY", "avgRating": 4.409090909090909
}
{
  "_id": "B00JXLGELG", "avgRating": 4.6666666666666667
}
{
  "_id": "B00JQXRWQM", "avgRating": 4.25
}
{
  "_id": "B00JKRP22", "avgRating": 3.6666666666666665
}
{
  "_id": "B00JKRKFQ1", "avgRating": 4.6521739130434785
}
{
  "_id": "B00JJ573XK", "avgRating": 3.6388888888888889
}
{
  "_id": "B00JNBVYE", "avgRating": 4.4444444444444445
}
{
  "_id": "B00IZL23MV", "avgRating": 4.3541666666666667
}
{
  "_id": "B00IQBT3QW", "avgRating": 4.1333333333333334
}
{
  "_id": "B00IMBUUWA", "avgRating": 3.9545454545454546
}
{
  "_id": "B00IL7BJMM", "avgRating": 4.785714285714286
}
{
  "_id": "B00ISAB1G2", "avgRating": 4.8888888888888889
}
{
  "_id": "B00HW0IV45", "avgRating": 4
}
{
  "_id": "B00HLQOKLG", "avgRating": 4.375
}
{
  "_id": "B00HH0ESYK", "avgRating": 4.3888888888888889
}
{
  "_id": "B00HFTBMD8", "avgRating": 3.5151515151515
}
{
  "_id": "B00HFJWNVA", "avgRating": 4.999090909090909
}
> Type "it" for more
```

## Summary of results

Here we are querying out the top 10 products which received most positive reviews. We are grouping the total reviews by asin (product id) which is unique for each product and then we are calculating the average of overall for each product by using \$avg function. After that we are sorting the results in descending order and limiting the results to top 10.

- h. What is the date of the first review per category?

## Wording of query

```
db.review.aggregate([
  $group:
  {_id: "$category", firstReviewDate: { $min: "$reviewTime" }}
])
```

### Snapshot showing successful execution

```
Command Prompt - mongo
> db.TotalReviews.aggregate([
... {$group: {_id: "$category", firstReviewDate: { $min: "$reviewTime" }}}
... ]
{
  "_id" : "Music", "firstReviewDate" : "01 1, 2009"
}
{
  "_id" : "Office", "firstReviewDate" : "01 1, 2005"
}
{
  "_id" : "Toys & Games", "firstReviewDate" : "01 1, 2001"
}>
```

### Summary of results

Here we are querying the first review per category and the date of first review per category. We are grouping the total reviews by category field and then taking the minimum of reviewTime field by using \$min function.

- i. What is the top 10 most prolific reviewers?

### Wording of query

```
db.review.aggregate([
{$group:{_id: "$reviewerID", review_count: {$sum:1}}},
{$sort: {review_count:-1}},
{$limit:10}
])
```

### Snapshot showing successful execution

```
Command Prompt - mongo
> db.TotalReviews.aggregate([
... {$group:{_id: "$reviewerID", review_count: {$sum:1}}},
... {$sort: {review_count:-1}},
... {$limit:10}
...
{
  "_id" : "AJGU56YG8G1DQ", "review_count" : 550
}
{
  "_id" : "A22CW0ZHY3NJH8", "review_count" : 181
}
{
  "_id" : "A1POFVXXUZR3IQ", "review_count" : 180
}
{
  "_id" : "A1M8AYAL3L8ACP", "review_count" : 168
}
{
  "_id" : "A1DDOGXEYECQ8", "review_count" : 167
}
{
  "_id" : "A1II2ZRPKZAQD", "review_count" : 165
}
{
  "_id" : "A3OXHLG6DI8RW8", "review_count" : 165
}
{
  "_id" : "A1EVV74UQYYVKRY", "review_count" : 163
}
{
  "_id" : "A2WW57XX2UQLM6", "review_count" : 163
}
{
  "_id" : "A2QDOJFFLFGF18", "review_count" : 141
}>
```

### Summary of results

Here we are querying out the 10 most prolific reviewers i.e. with the most number of reviews. We are grouping the total reviews by reviewer id which is unique for each reviewer and then

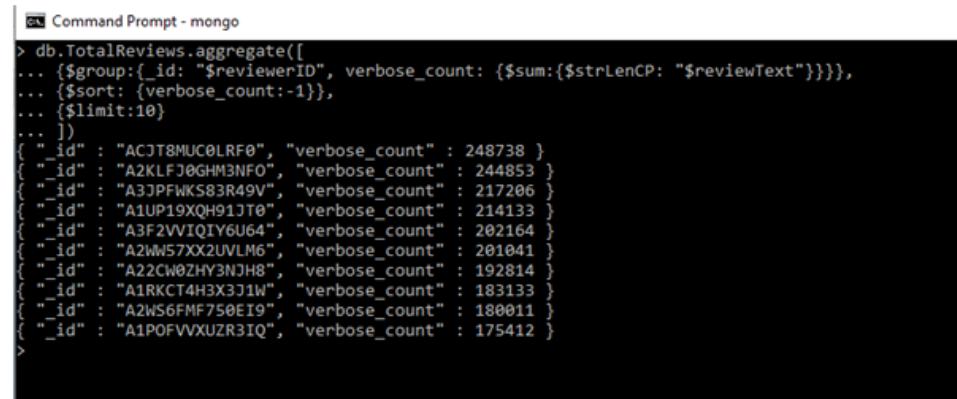
counting the number of documents for each reviewer. After that we are sorting the results in descending order and limiting the results to top 10.

- j. What is the top 10 most verbose reviewers?

#### Wording of query

```
db.review.aggregate([
  {$group:{_id: "$reviewerID", verbose_count: {$sum:{$strLenCP: "$reviewText"}}}},
  {$sort: {verbose_count:-1}},
  {$limit:10}
])
```

#### Snapshot showing successful execution



```
> db.TotalReviews.aggregate([
... {$group:{_id: "$reviewerID", verbose_count: {$sum:{$strLenCP: "$reviewText"}}}},
... {$sort: {verbose_count:-1}},
... {$limit:10}
... ])
{
  "_id" : "ACJT8MUC0LRF0",
  "verbose_count" : 248738
},
{
  "_id" : "A2KLFJ0GHM3NFO",
  "verbose_count" : 244853
},
{
  "_id" : "A3JPFWK583R49V",
  "verbose_count" : 217206
},
{
  "_id" : "A1UP19XQH91JT0",
  "verbose_count" : 214133
},
{
  "_id" : "A3F2VVIQIY6U64",
  "verbose_count" : 202164
},
{
  "_id" : "A2WM57XX2UVLM6",
  "verbose_count" : 201841
},
{
  "_id" : "A22CW0ZHY3NJHQ",
  "verbose_count" : 192814
},
{
  "_id" : "A1RKCT4H3X3J1W",
  "verbose_count" : 183133
},
{
  "_id" : "A2WS6FMF750EI9",
  "verbose_count" : 180011
},
{
  "_id" : "A1POFVVXUZR3IQ",
  "verbose_count" : 175412
}>
```

#### Summary of results

Here we are querying out the 10 most verbose reviewers i.e. with the max amount of review text. We are grouping the total reviews by reviewer id which is unique for each reviewer and then we are calculating the length of review text for each reviewer. After that we are sorting the results in descending order and limiting the results to top 10.

- k. What is the top 10 most positive reviewers?

#### Wording of query

```
db.review.aggregate([
  {$group:{_id: "$reviewerID", avgRating: {$avg: "$overall"}}},
  {$sort: {avgRating:-1}},
  {$limit:10}
])
```

### Snapshot showing successful execution

```
cmd Command Prompt - mongo
... {$group:{_id: "$reviewerID", avgRating: {$avg: "$overall"}}},
... {$sort: {avgRating:-1}},
... {$limit:10}
...
{
  "_id" : "AA75INOMAHULX", "avgRating" : 5 },
  "_id" : "A1B3GN09C8YX0N", "avgRating" : 5 },
  "_id" : "A17USEQCX31NS2", "avgRating" : 5 },
  "_id" : "A361M14PU2GUEG", "avgRating" : 5 },
  "_id" : "A13SAVYJLEH4D1", "avgRating" : 5 },
  "_id" : "A27DDHWHE540I5", "avgRating" : 5 },
  "_id" : "A3LWTCX08GMD9N", "avgRating" : 5 },
  "_id" : "AY5IL05TZIPZO", "avgRating" : 5 },
  "_id" : "A3VD1ELQVO02DW", "avgRating" : 5 },
  "_id" : "A2X2CJK2C5WZVW", "avgRating" : 5 }
```

### Summary of results

Here we are querying out the 10 most positive reviewers. We are grouping the total reviews by reviewer id which is unique for each reviewer and then we are calculating the average of overall field for each reviewer. After that we are sorting the results in descending order and limiting the results to top 10.

- I. What are the top 10 products with most number of reviews?

### Wording of query

```
db.review.aggregate([
  {$group:{_id: "$asin", review_count: {$sum:1}}},
  {$sort: {review_count:-1}},
  {$limit:10}
])
```

### Snapshot showing successful execution

```
cmd Command Prompt - mongo
... {$group:{_id: "$asin", review_count: {$sum:1}}},
... {$sort: {review_count:-1}},
... {$limit:10}
...
{
  "_id" : "B0010T3QT2", "review_count" : 311 },
  "_id" : "B004S8F7QM", "review_count" : 309 },
  "_id" : "B0089RPUHO", "review_count" : 253 },
  "_id" : "B0039N7ELS", "review_count" : 227 },
  "_id" : "B0039S7NO6", "review_count" : 215 },
  "_id" : "B001W30D20", "review_count" : 207 },
  "_id" : "B0027CTFB0", "review_count" : 205 },
  "_id" : "B002NPBT50", "review_count" : 196 },
  "_id" : "B000LSZVKA", "review_count" : 194 },
  "_id" : "B000N5QNSK", "review_count" : 192 }
```

## Summary of results

Here we are querying out the top 10 products with most number of reviews. We are grouping the total reviews by asin (product id) which is unique for each product and then we are calculating the total number of documents for each product. After that we are sorting the results in descending order and limiting the results to top 10.

- m. What is the top 10 most negative reviewers?

## Wording of query

```
db.review.aggregate([
  {$group:{_id:{reviewerId: "$reviewerID", name: "$reviewerName"}, avgRating: {$avg: "$overall"}}, 
  {$sort: {avgRating:1}}, 
  {$limit:10}
])
```

## Snapshot showing successful execution



```
Command Prompt - mongo
> db.TotalReviews.aggregate([
... {$group:{_id:{reviewerId: "$reviewerID", name: "$reviewerName"}, avgRating: {$avg: "$overall"}}, 
... {$sort: {avgRating:1}}, 
... {$limit:10}
... ])
{
  "_id": {
    "reviewerId": "A3Q798TJMRUIRU",
    "name": "doll lady \\\"doll lady\\\"",
    "avgRating": 1
  },
  "_id": {
    "reviewerId": "A24TVPGG6JGH4",
    "avgRating": 1
  },
  "_id": {
    "reviewerId": "A2BB4DGBRUVGKMM",
    "avgRating": 1
  },
  "_id": {
    "reviewerId": "A3U521S7IPSFZG",
    "avgRating": 1
  },
  "_id": {
    "reviewerId": "AHFLWC1G4OBW1",
    "name": "Steven Davidson",
    "avgRating": 1
  },
  "_id": {
    "reviewerId": "A1STG5612WHLX",
    "name": "Andriy K.",
    "avgRating": 1
  },
  "_id": {
    "reviewerId": "A6FIAB28IS79",
    "avgRating": 1
  },
  "_id": {
    "reviewerId": "A5NIQWTH7QVM",
    "avgRating": 1
  },
  "_id": {
    "reviewerId": "A2LB4YHGBWOLSU",
    "name": "Kelly",
    "avgRating": 1
  },
  "_id": {
    "reviewerId": "A3IFRQETHPZ5GJ",
    "name": "Al L.",
    "avgRating": 1
  }
}
```

## Summary of results

Here we are querying out the 10 most negative reviewers. We are grouping the total reviews by reviewer id which is unique for each reviewer and then we are calculating the average of overall field for each reviewer. After that we are sorting the results in ascending order and limiting the results to top 10.

- n. What are the top 10 products with most number of positive reviews?

## Wording of query

```
db.review.aggregate([
  {$group:{_id: "$asin", avgRating: {$avg: "$overall"}}, 
  {$sort: {avgRating:-1}}, 
  {$limit:10}
])
```

```
{$limit:10}  
])
```

### Snapshot showing successful execution



```
ca: Command Prompt - mongo  
> db.TotalReviews.aggregate([  
... {$group:{_id: "$asin", avgRating: {$avg: "$overall"} }},  
... {$sort: {avgRating:-1}},  
... {$limit:10}  
... ])  
{ "_id" : "B006LU77M2", "avgRating" : 5 }  
{ "_id" : "B0075LDFNQ", "avgRating" : 5 }  
{ "_id" : "B005YVMIQO", "avgRating" : 5 }  
{ "_id" : "B0089INHUU", "avgRating" : 5 }  
{ "_id" : "B009FFXVQQ", "avgRating" : 5 }  
{ "_id" : "B006LNU00K", "avgRating" : 5 }  
{ "_id" : "B005S0YXEG", "avgRating" : 5 }  
{ "_id" : "B005TLVNMY", "avgRating" : 5 }  
{ "_id" : "B005VC8CK2", "avgRating" : 5 }  
{ "_id" : "B0062C2CGU", "avgRating" : 5 }  
>
```

### Summary of results

Here we are querying out the top 10 products which received most positive reviews. We are grouping the total reviews by asin (product id) which is unique for each product and then we are calculating the average of overall for each product by using \$avg function. After that we are sorting the results in descending order and limiting the results to top 10.

- o. What are the top 10 products with most number of negative reviews?

### Wording of query

```
db.review.aggregate([  
{$group:{_id: "$asin", avgRating: {$avg: "$overall"} }},  
{$sort: {avgRating:1}},  
{$limit:10}  
])
```

### Snapshot showing successful execution

```
> db.TotalReviews.aggregate([
... {$group:{_id: "$asin", avgRating: {$avg: "$overall"}},},
... {$sort: {avgRating:1}},
... {$limit:10}
...])
{ "_id" : "B000AR84B0", "avgRating" : 1 },
{ "_id" : "B0017867QK", "avgRating" : 1 },
{ "_id" : "B002VUCA00", "avgRating" : 1.111111111111112 },
{ "_id" : "B003B1WF7Q", "avgRating" : 1.333333333333333 },
{ "_id" : "B009F70UE6", "avgRating" : 1.375 },
{ "_id" : "B001TE5C1S", "avgRating" : 1.4 },
{ "_id" : "B001FVSZAK", "avgRating" : 1.4 },
{ "_id" : "B000SDU7ZW", "avgRating" : 1.4 },
{ "_id" : "B003AN6H42", "avgRating" : 1.4 },
{ "_id" : "B001P06IL0", "avgRating" : 1.5 }
```

### Summary of results

Here we are querying out the top 10 products which received most negative reviews. We are grouping the total reviews by asin (product id) which is unique for each product and then we are calculating the average of overall for each product by using \$avg function. After that we are sorting the results in ascending order and limiting the results to top 10.

## FINDINGS AND INTERPRETATION

We have the following findings from the Amazon dataset provided.

- Comparison between the overall number of positive and negative reviews as well as per category.

First, we calculated the total number of reviews having a rating greater than or equal to 3, which gives us the total number of positive reviews. Next, we calculated the overall number of reviews having a rating less than 3 i.e. negative reviews.

#### Query:

```
db.review.count({overall:{$gte:3}}) => 216788
db.review.count({overall:{$lt:3}}) => 14328
```

```
MongoDB Enterprise > db.review.count({overall:{$gte:3}})
216788
MongoDB Enterprise > db.review.count({overall:{$lt:3}})
14328
MongoDB Enterprise >
```

Next, we have calculated the number of positive and negative reviews per category. We have done this for all 3 categories which we have considered.

For category “Office”,

```
db.review.find({category:"Office", overall:{$gte:3}}).count() => 50402
```

```
db.review.find({category:"Office", overall:{$lt:3}}).count() => 2856
```

```
MongoDB Enterprise > db.review.find({category:"Office", overall:{$gte:3}}).count()
50402
MongoDB Enterprise > db.review.find({category:"Office", overall:{$lt:3}}).count()
2856
MongoDB Enterprise >
```

For category “Music”,

```
db.review.find({category:"Music", overall:{$gte:3}}).count() => 9794
db.review.find({category:"Music", overall:{$lt:3}}).count() => 467
```

```
MongoDB Enterprise > db.review.find({category:"Music", overall:{$gte:3}}).count()
9794
MongoDB Enterprise > db.review.find({category:"Music", overall:{$lt:3}}).count()
467
MongoDB Enterprise >
```

For category “Toys& Games”,

```
db.review.find({category:"Toys& Games", overall:{$gte:3}}).count() => 156592
db.review.find({category:"Toys& Games", overall:{$lt:3}}).count() => 11005
```

```
MongoDB Enterprise > db.review.find({category:"Toys& Games", overall:{$gte:3}}).count()
156592
MongoDB Enterprise > db.review.find({category:"Toys& Games", overall:{$lt:3}}).count()
11005
MongoDB Enterprise >
```

- b. Identification of the product category which is being reviewed for the longest period of time.

We identified the category of products which is being reviewed for the longest period of time. From this we can interpret that this particular category has been around for the longest period of time and has been constantly receiving reviews. This may indicate that the popularity and sales of that product category because only if the product (belonging to that category) is being bought by customers will it receive more number of reviews. The dataset available consists of reviews only till the year 2014; the newest review date we have in our review collection is 23<sup>rd</sup> July 2014.

Just on a side note, the reviews go back to 28<sup>th</sup> July 2014. This implies our dataset consists of reviews from the year 2000 to 2014 for the categories, Office Products, Musical Instruments and Toys & Games.

#### Query:

1. We have first converted the review time for each review from String to Date format using the below function. This was required to be done because we preferred to display the oldest and newest review times in date format instead of the unix review times (epochs).

```

db.review.find().forEach(function(element){
    element.reviewTime = new Date(element.reviewTime);
    db.review.save(element);
})

```

This query simply converts the strings to dates and has no visible output.

2. Next, we simply grouped the review collection by category and then calculated the oldest review date and newest review date using the \$min and \$max functions respectively.

```

db.review.aggregate([
{$group: {"_id": "$category", oldestReviewTime: {$min: "$reviewTime"}, newestReviewTime: {$max: "$reviewTime"}}}
])

```

Output:

```

MongoDB Enterprise > db.review.aggregate([
... {$group: {"_id": "$category", oldestReviewTime: {$min: "$reviewTime"}, newestReviewTime: {$max: "$reviewTime"}}}
... ])
{ "_id" : "Music", "oldestReviewTime" : ISODate("2004-09-18T04:00:00Z"), "newestReviewTime" : ISODate("2014-07-22T04:00:00Z") }
{ "_id" : "Toys& Games", "oldestReviewTime" : ISODate("2000-07-28T04:00:00Z"), "newestReviewTime" : ISODate("2014-07-23T04:00:00Z") }
{ "_id" : "Office", "oldestReviewTime" : ISODate("2000-09-29T04:00:00Z"), "newestReviewTime" : ISODate("2014-07-23T04:00:00Z") }
MongoDB Enterprise >

```

- c. Identification of the product category which has the most number of recommended reviews.

We have formulated a definition for “recommended reviews” for a particular product. This definition is that if a reviewer precisely writes something like “I recommend this product to anyone.”, we consider it as a review which actually recommends buying the product. Similarly, if the review has words “not recommend”, then we will consider this review as one which discourages buying the product. Next for each category, we calculate the total number of reviews which recommends buying the overall products in that category. It is evident from the below query that the Toys & Games category has the most number of recommended products based on the corresponding reviews.

Query:

```

db.review.aggregate([
{$match: {$text: {$search:"recommend Recommend -\"not recommend\"}}},
{$group:{_id:$category, recommendedReviews:{$sum:1}}}
])

```

```

MongoDB Enterprise > db.review.aggregate([
... {$match: {$text: {$search:"recommend Recommend -\"not recommend\""}}, 
... {$group:{_id:"$category", recommendedReviews:{$sum:1}}}
...
])
{
  "_id" : "Office", "recommendedReviews" : 7901 }
{
  "_id" : "Music", "recommendedReviews" : 1162 }
{
  "_id" : "Toys& Games", "recommendedReviews" : 19555 }
MongoDB Enterprise >

```

- d. Identification of the product which has the most number of “also bought” and “also viewed” products.

For each product, Amazon shows some related products that were either viewed or bought by other customers. Amazon basically uses this strategy to increase its sales. If the user buys the current product as well as one or more of the “also viewed” and/or “also bought” products, this results in an increase in number of sales from Amazon’s perspective. Therefore, it will be helpful if we are able to determine the products which have the highest number of related (also bought/ also viewed) products.

In the below query, we assume that if the product has at least 10 also viewed and at least 10 also bought products, we consider it to have a high number of related products. The output shows the asin for the products which are considered to have a high number of related products.

Amazon can probably highlight the products which have high number of related products in order to lead the users on those product’s pages and increase their probability of buying their related products as well.

#### Query:

```

db.product.find( {"related.also_viewed" : {$exists:true},
$where:'this.related.also_viewed.length>10', "related.also_bought" : {$exists:true},
$where:'this.related.also_bought.length>10'}, {asin:1, _id:0} )

```

```

MongoDB Enterprise > db.product.find( {"related.also_viewed" : {$exists:true}, $where:'this.related.also_viewed.length>10', "related.also_bought" : {$exists:true}, $where:'this.related.also_bought.length>10'}, {asin:1, _id:0} )
{
  "asin" : "0000032069" }
{
  "asin" : "0000032034" }
{
  "asin" : "0000031852" }
{
  "asin" : "0000032050" }
{
  "asin" : "0000031909" }
{
  "asin" : "0000031895" }
{
  "asin" : "0000031887" }
{
  "asin" : "0005123909" }
{
  "asin" : "0005119367" }
{
  "asin" : "0005134188" }
{
  "asin" : "0061997102" }
{
  "asin" : "0071436588" }
{
  "asin" : "0000031909" }
{
  "asin" : "0000031852" }
{
  "asin" : "0000032034" }
{
  "asin" : "0000032050" }
{
  "asin" : "0000031895" }
{
  "asin" : "0000031887" }
{
  "asin" : "0005123909" }
{
  "asin" : "0005119367" }
Type "it" for more
MongoDB Enterprise >

```

# ANALYTICS

## MongoDB Metrics

Based on our analysis, we have come up with a few metrics which could be useful to include on Amazon's webpages for individual products and reviewers. These metrics are not available in Amazon's website currently.

### A) Performance over time metric

#### Reason for choosing the metric:

This metric is shows the performance of each product with time. One of the problem with amazon reviews is the date of the reviews is not properly displayed to the user to show how old a particular review is. So we decided to show a graph that displays the rating of the product from the time it was launched.

#### Wording of the query:

```
db.Reviews.aggregate([{$sort:[{"unixReviewTime":1}]},{$group:{_id:"$asin",PerformanceOverTime:{$push:{stars:"$overall",date:"$reviewTime"}}}},{$project:{ProductID:"$_id",PerformanceOverTime:{$$slice:[ "$PerformanceOverTime",10]}},_id:0}],{$limit:5}],{allowDiskUse:true}).pretty()
```

#### Summary:

We used aggregate pipeline to group the data based on productID, and \$push function to create an array of ratings with their time. We have sorted the data based on the review date in ascending order, in this way we will show the ratings of the product from the first review. This data can be used to display a graph on the product page.

#### Screenshot of Mongodb output:

```
{  
  "ProductID" : "B00H3QRFAC",  
  "PerformanceOverTime" : [  
    {  
      "stars" : 5,  
      "date" : "07 8, 2014"  
    },  
    {  
      "stars" : 5,  
      "date" : "07 8, 2014"  
    },  
    {  
      "stars" : 4,  
      "date" : "07 10, 2014"  
    },  
    {  
      "stars" : 4,  
      "date" : "07 11, 2014"  
    },  
    {  
      "stars" : 3,  
      "date" : "07 12, 2014"  
    },  
    {  
      "stars" : 5,  
      "date" : "07 13, 2014"  
    }  
  ]  
}
```

## Screenshot of proposed metric on amazon page:



## B) Sort by Helpful votes

### Reason for choosing the metric:

With this metric we have decided to sort the reviews based on helpful votes of reviews. One of the feature that amazon lacks in sorting is that the reviews cannot be sorted based on the usefulness of the reviews.

### Wording of the query:

```
db.Reviews.aggregate([{$group:{"_id":{"Product":"$asin"}, "Review":{$push:{Percentage:$multiply:[{$cond:[{$eq:[{"$arrayElemAt":["$helpful",1]},0]},0,{${divide:[{"$arrayElemAt":["$helpful",0],"$arrayElemAt":["$helpful",1]}]}]}},100]}, "ReviewText": "$reviewText"}}}], {"$project": {"Review": "$Review"}, {"$sort": {"Review.Percentage": -1}}, {"$limit": 5}], {allowDiskUse: true}).pretty()
```

### Summary:

We have decided to use aggregate pipeline to group reviews based on the productID and using \$push to create an array of the review text and a percentage of helpfulness field that we calculated based on the data helpful votes data. We sorted the data based on this percentage, thereby displaying the most helpful review on the top.

### Screenshot of Mongodb output:

```
{  
    "_id" : {  
        "Product" : "0786955708"  
    },  
    "Review" : [  
        {  
            "Percentage" : 50,  
            "ReviewText" : "I got this game from a local shop and I must say I wish I would've known this was on amazon because its way cheaper. The game is super fun and to anyone that is even just a little bit interested in buying it DO IT! Wrath of Ashardalon is that fun and even more so when you have a full 6 people to play."  
        },  
        {  
            "Percentage" : 23.076923076923077,  
            "ReviewText" : "Hated this product. Predictable. Not fun. It attempts to be a fast setup table top game but feels as a meshed together mess of crap. Often times I thought this is like d&d; battle.... Just kidding! We never played it again after one play through. Would not recommend this for an older gaming group that enjoys board/card games like talisman and 7 wonders and ascension.F"  
        },  
        {  
            "Percentage" : 0,  
            "ReviewText" : "I have loved all the D&D board games from Wizards of the Coast so far (I have the other two already) and this one was no exception. Great pieces and sturdy construction for a lot of fun. I especially liked the more challenging Encounter cards in this game in particular. Very fun and relatively easy to learn. There is a bit of reading but once you get though that it is very simple and once you get it you don't need to really re-read the instructions. Maybe we went back for clarification once or twice in the first few games, but that was it. If you loved old school D&D but wished it was a bit more tactile (you know, game pieces and stuff) as opposed to all pencil and paper; like I did, then you will LOVE this!"  
        },  
        {  
            "Percentage" : 0,  
            "ReviewText" : "I own this and Castle Ravenloft, which are the two I would always put together. WoA is what most recommend as the one to get if you just get one -- it's difficulty is a average (unlike CR which is a real challenge and heroes are good). These D&D games are so much fun because it truly is a system of sorts. You will end up having a ton of little figures (sadly not painted the way some other games have them, but then again you get so many of it) that makes sense in a jillion little towns, and an equally well illustrated rulebook along with a really nice collection of scenarios to play a lot of 'em!'. This game, while you have the 1-2hrs needed for some of the scenarios with friends, just plays so well with the whole group. It's truly cooperative, every character is different and everyone can contribute into the fun. I do suggest that if you buy this game, play it for a bit w/ the base rules, and if you ever start feeling like you are finding it a bit mechanically repetitive; then pause, go to boardgamegeek and find the variants (one I published in fact) that allow for an initiative system - it will change everything and you'll be back to enjoy this game in no time! Absolutely worth every penny when it is ~ $40."  
        },  
        {  
            "Percentage" : 0,  
            "ReviewText" : "This is in my opinion the best dungeon crawler out there. It can be setup, played and torn down all in under 2 hours. It makes a fun ride, and is fun with your kids too (mine are 5 and 7)."  
        }  
    ]  
}
```

### Screenshot of proposed metric on amazon page:

The screenshot shows a list of five reviews from the Amazon product page, each highlighted with a red box. The reviews are as follows:

- pascal thibon** **★★★★★ NICE!** The shape of the toy is good, the child likes it very much  
92% of the voters found this review helpful.
- David** **★★★★★ Very interesting toys** Very interesting toys  
90% of the voters found this review helpful.
- samantha0424** **★★★★★ Pass!** This broke after about a half hour of use and seemed to leave scratches on the floor somehow too.  
90% of the voters found this review helpful.
- Christina Schulte** **★★★★★ Five Stars** Like it very much  
85% of the voters found this review helpful.
- Linda Thorpe** **★★★★★ Five Stars** My grandson loved this soccer ball. Quick delivery.  
84% of the voters found this review helpful.

### C) Most Popular products under the same Category:

#### Reason for choosing this metric:

Amazon's website shows similar items for a product. These items being shown may not be very popular or may not be highly rated. We want to recommend the best products to the user. Therefore, we sorted the products in the same category based on the highest rating and then sorted again based on the most number of reviews. This will allow the user to view the most popular products under the same category. In this way, the user can make better informed decision while purchasing any product. According to us, rating and number of reviews will define popularity of product.

#### Wording of Query:

```
db.allreviews.aggregate([{$match:{'category':'toys'}},{$group:{'_id':{'Category':'$category','prod':'$asin'},'Avg':{$avg:'$overall'},'Total':{$sum:1}}},{$sort:{'Avg':-1}},{$sort:{'Total':-1}},{$limit:5}])
```

## Summary:

For a given Category "Toys", we have found out the average rating for each product and then calculated the total reviews available for each product. We sorted the results based on the average rating and the average number of reviews available for a particular product in the given category.

## Screenshot from MongoDB:

```
{ "_id" : { "Category" : "toys", "prod" : "B00458F7QM" }, "Avg" : 4.86084142394822, "Total" : 309 }
{ "_id" : { "Category" : "toys", "prod" : "B008RPUH0" }, "Avg" : 4.600790513833992, "Total" : 253 }
{ "_id" : { "Category" : "toys", "prod" : "B003957N0G" }, "Avg" : 4.716279069767442, "Total" : 215 }
{ "_id" : { "Category" : "toys", "prod" : "B001W30D20" }, "Avg" : 4.5893719806763285, "Total" : 207 }
{ "_id" : { "Category" : "toys", "prod" : "B002NPBT50" }, "Avg" : 4.566326530612245, "Total" : 196 }
```

## Screenshot from our proposed Amazon webpage:



## D) Most sold products under the same Category:

### Reason for choosing this metric:

We would like the users to view the most sold products under the category that they are browsing. This will enable the users to make better informed decisions while choosing a product for buying.

### Wording of Query:

```
db.products.find({"salesRank.Books":{$exists:true}}, {"salesRank.Books":1}).sort({"salesRank.Books":1}).limit(5)
```

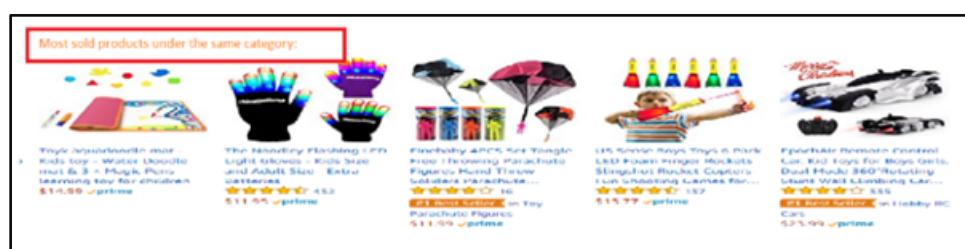
### Summary:

For a given category, we have displayed top 5 products with the least sales rank. Some of the documents in our collection did not have a sales rank attribute, so we excluded those documents from our query.

## Screenshot from MongoDB:

```
{ "_id" : ObjectId("5adcaf924745c1259a62a82d"), "asin" : "0062345214", "salesRank" : { "Books" : 6 } }
{ "_id" : ObjectId("5adcaf8b4745c1259a628fb"), "asin" : "0061122416", "salesRank" : { "Books" : 14 } }
{ "_id" : ObjectId("5adcaf8e4745c1259a629b0a"), "asin" : "0061950726", "salesRank" : { "Books" : 25 } }
{ "_id" : ObjectId("5adcaf904745c1259a6299e6"), "asin" : "0062024843", "salesRank" : { "Books" : 31 } }
{ "_id" : ObjectId("5adcaf924745c1259a62a7f4"), "asin" : "006232005X", "salesRank" : { "Books" : 39 } }
```

## Screenshot from our proposed Amazon webpage:



## E) Percentage of helpfulness of each review:

### Reason for choosing this metric:

Every review has helpful votes. These are displayed in Amazon. But in our webpage, we would like to display percentage of helpfulness of a review because there might be cases when a review has got 100 votes out of which there are only 5 helpful votes. So, the percentage of helpfulness is only 5 %, while there can be a review which received totally 5 votes out of which all 5 are Helpful. So, this review has 100% helpfulness. Hence, we would like the user to view this metric to judge whether to take a review into consideration or not. This will prevent the user from getting biased by a highly positive or highly negative review, which was otherwise not considered Helpful by majority of the users.

### Wording of Query:

```
db.allreviews.aggregate([{$group:{"_id":{"Review":"$reviewText","Total":{"$arrayElemAt": [ "$helpful", 1 ]}, "helpfulVotes":{"$arrayElemAt": [ "$helpful", 0 ]}}}, {$project:{ "Decimal":{ $cond: [ { $eq: [ "$_id.Total", 0 ] }, 0,{ $divide:[ "$_id.helpfulVotes", "$_id.Total"]]} ]}}, {$project:{ "NotHelpful":{ $subtract:[ "$_id.Total", "$_id.helptfulVotes"]}, "Percentage":{ $multiply:[ "$Decimal",100]} }}, {allowDiskUse:true}])
```

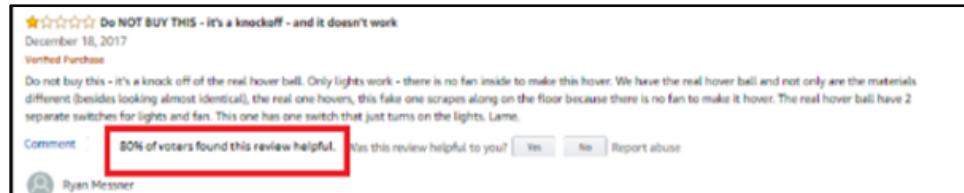
### Summary:

In the json collection we have an array called helpful. The first element of the array is number of helpful votes and the second element is the total number of votes for a particular review. So, for each review, we have calculated helpful votes divided by total votes, multiplied by 100 to calculate the percentage. If a review did not have any votes, then we excluded those reviews from division to avoid divide by zero error and displayed percentage of helpfulness as zero for those reviews.

### Screenshot from MongoDB:

```
[{"_id": {"Review": "", "Total": 0, "helpfulVotes": 0}, "NotHelpful": 0, "Percentage": 0}, {"_id": {"Review": "", "Total": 1, "helpfulVotes": 0}, "NotHelpful": 1, "Percentage": 0}, {"_id": {"Review": "", "Total": 1, "helpfulVotes": 1}, "NotHelpful": 0, "Percentage": 100}, {"_id": {"Review": "", "Total": 2, "helpfulVotes": 1}, "NotHelpful": 1, "Percentage": 50}, {"_id": {"Review": "", "Total": 2, "helpfulVotes": 2}, "NotHelpful": 0, "Percentage": 100}, {"_id": {"Review": "", "Total": 3, "helpfulVotes": 1}, "NotHelpful": 2, "Percentage": 33.3333333333333}, {"_id": {"Review": "", "Total": 3, "helpfulVotes": 2}, "NotHelpful": 1, "Percentage": 66.66666666666666}, {"_id": {"Review": "", "Total": 3, "helpfulVotes": 3}, "NotHelpful": 0, "Percentage": 100}]
```

### Screenshot from our proposed Amazon webpage:



## F) Top 5 recommended reviews for a product:

### Reason for choosing this metric:

We have come up with our own definition of recommended reviews. We have defined a list of words and searched for those words in all the reviews for a particular product. Wherever we found a strong match based on high metascore, we displayed those reviews as recommended reviews. The reviews with the lowest metascore are considered as non-recommended reviews. We feel this will be a helpful criterion for the users to judge the quality of a product by taking into account the experience of other reviewers.

### **Wording of Query:**

We first created an Index in reviewText field using the query: db.allreviews.createIndex({reviewText: "text"})

```
db.allreviews.find({ $text: { $search: "good bad worse excellent price quality worst great delivery fast slow service time cost poor" }, "asin": "1384719342" }, { score: { $meta: "textScore" } }).sort( { score: { $meta: "textScore" } } ).limit(5).pretty()
```

### **Summary:**

For a given product, we searched reviews with the following words: "good bad worse excellent price quality worst great delivery fast slow service time cost poor". We sorted the obtained results based on metascore. So top 5 reviews with the highest relevance to these words are displayed.

### **Screenshot from MongoDB:**

```
{  
    "_id" : ObjectId("5ada2e5123f6ae830d9e0295"),  
    "reviewerID" : "A2IBPI20UZIR0U",  
    "asin" : "1384719342",  
    "reviewerName" : "cassandra tu \\"Yeah, well, that's just like, u...",  
    "helpful" : [  
        0,  
        0  
    ],  
    "reviewText" : "Not much to write about here, but it does exactly what it's supposed to. filters out  
    one of the lowest prices pop filters on amazon so might as well buy it, they honestly work the same despite  
    overall" : 5,  
    "summary" : "good",  
    "unixReviewTime" : 1393545600,  
    "reviewTime" : "02 28, 2014",  
    "category" : "music",  
    "score" : 0.8125  
}
```

### **Screenshot from our proposed Amazon webpage:**



### **G) Average rating given by a reviewer irrespective of products:**

#### **Reason for choosing this metric:**

When a user reads a review and is strongly influenced by it, we feel it will help the user to also check whether the reviewer is a positive or negative reviewer in general. A reviewer may give only positive reviews or only negative reviews. So, this metric will help a prospective buyer to not be influenced by just one such review and to take other reviews also into consideration.

### **Wording of Query:**

```
db.allreviews.aggregate([{$group:{"_id": "$reviewerID", "Avg": {$avg: "$overall"}}, {$project: {"Average": "$Avg", "PositiveOrNegativeReviewer": {$cond: [{"$gte": ["$Avg", 3]}, "Positive", "Negative"]}}}]})
```

### **Summary:**

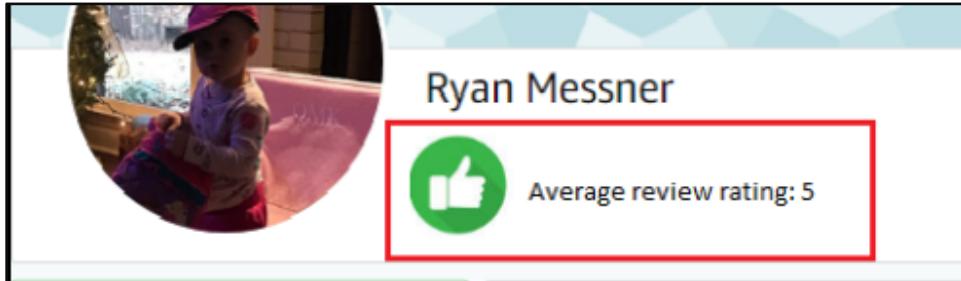
In this query, we have found out the average review given by each reviewer and based on that average, we have determined whether the reviewer is a positive or negative reviewer in general. If average rating

is greater than 3, the reviewer is considered to be a Positive reviewer. If the average rating is less than 3, the reviewer is considered to be a Negative reviewer in general.

#### Screenshot from MongoDB:

```
{
  "_id": "AOE0XMITAJWDT", "Average": 4.2, "PositiveOrNegativeReviewer": "Positive"
  "_id": "A3RNSX8VCK38MA", "Average": 4.333333333333333, "PositiveOrNegativeReviewer": "Positive"
  ...
  "_id": "A3E07UI5GJ2BT7", "Average": 4.2, "PositiveOrNegativeReviewer": "Positive"
  ...
  "_id": "A1CUORA1W0E33B", "Average": 4.6, "PositiveOrNegativeReviewer": "Positive"
  ...
  "_id": "A2YDAGENHKY5FX", "Average": 1.833333333333333, "PositiveOrNegativeReviewer": "Negative"
  ...
  "_id": "AP42XSKS2MMYZ", "Average": 2.4, "PositiveOrNegativeReviewer": "Negative"
  ...
  "_id": "A3MKON3VM0W1TH", "Average": 4.4, "PositiveOrNegativeReviewer": "Positive"
}
```

#### Screenshot from our proposed Amazon webpage:



#### H) Percentage of helpfulness of each reviewer:

##### Reason for choosing this metric:

In Amazon, only the total number of Helpful votes received by a reviewer is displayed. But in our proposed webpage, we would like to show the total number of Not Helpful votes and the percentage of helpfulness of a particular reviewer. This is because a reviewer may have received 5 helpful votes but may have received 20 not helpful votes, while another reviewer might have received 5 helpful votes and 0 not helpful votes. Hence, we would like to display the percentage of helpfulness by calculating total helpful votes received by a user divided by total votes received by a user.

##### Wording of Query:

```
db.allreviews.aggregate([
  {
    $group: {
      "_id": {
        "Reviewer": "$reviewerID",
        "Total": "$arrayElemAt": [
          "$helpful", 1
        ],
        "helpfulVotes": "$arrayElemAt": [
          "$helpful", 0
        ]
      }
    },
    $project: {
      "Decimal": {
        $cond: [
          { $eq: [
            "$_id.Total", 0
          ] },
          0,
          {
            $divide: [
              "$_id.helpfulVotes",
              "$_id.Total"
            ]
          }
        ]
      },
      "NotHelpful": {
        $subtract: [
          "$_id.Total",
          "$_id.helpfulVotes"
        ]
      },
      "Percentage": {
        $multiply: [
          "$Decimal",
          100
        ]
      }
    }
  }
], {allowDiskUse: true})
```

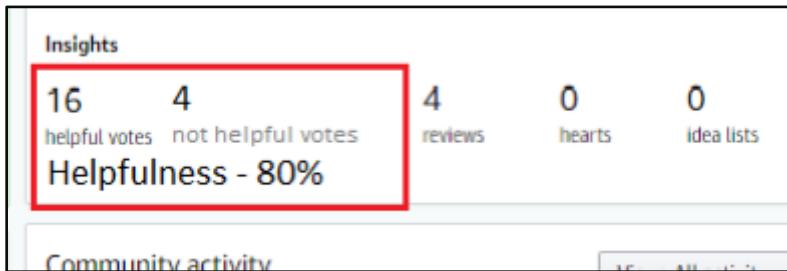
##### Summary:

We have calculated the total helpful, total not helpful and total votes received by every reviewer. Then we have divided the total helpful votes by total votes received and multiplied it by 100 to obtain the percentage. If any reviewer received 0 total votes, we excluded those reviewers to avoid divide by zero error and assigned zero as the percentage of helpfulness for those reviewers.

#### Screenshot from MongoDB:

```
{
  "_id": {
    "Reviewer": "A2F9E9MKV3MT02",
    "Total": 1,
    "helpfulVotes": 0,
    "NotHelpful": 1,
    "Percentage": 0
  },
  ...
  "_id": {
    "Reviewer": "A255V6KHM31VFD",
    "Total": 3,
    "helpfulVotes": 0,
    "NotHelpful": 3,
    "Percentage": 0
  },
  ...
  "_id": {
    "Reviewer": "A2Q8GFSEUIG151",
    "Total": 9,
    "helpfulVotes": 9,
    "NotHelpful": 0,
    "Percentage": 100
  },
  ...
  "_id": {
    "Reviewer": "A31N8XY2UTB25C",
    "Total": 7,
    "helpfulVotes": 7,
    "NotHelpful": 0,
    "Percentage": 100
  },
  ...
  "_id": {
    "Reviewer": "A3NM1MT3Q2FHXX",
    "Total": 9,
    "helpfulVotes": 7,
    "NotHelpful": 2,
    "Percentage": 77.77777777777779
  },
  ...
  "_id": {
    "Reviewer": "A3205FZH994CNY",
    "Total": 2,
    "helpfulVotes": 2,
    "NotHelpful": 0,
    "Percentage": 100
  },
  ...
  "_id": {
    "Reviewer": "A2I980OPNIJNYY",
    "Total": 3,
    "helpfulVotes": 2,
    "NotHelpful": 1,
    "Percentage": 66.66666666666666
  },
  ...
  "_id": {
    "Reviewer": "A1SCAMM0TFeG9I",
    "Total": 6,
    "helpfulVotes": 6,
    "NotHelpful": 0,
    "Percentage": 100
  }
}
```

Screenshot from our proposed Amazon webpage:



## HELPFULNESS OF REVIEWS

Text based definition:

We have come up with our own definition of recommended reviews. We have defined a list of words and searched for those words in all the reviews for a particular product. Wherever we found a strong match based on high metascore, we considered those reviews as helpful. The reviews with lower metascores are considered as not helpful. We feel this will be a helpful criterion for the users to judge the quality of a product by taking into account the experience of other reviewers. We searched reviews with the following words: "good bad worse excellent price quality worst great delivery fast slow service time cost poor". While length can also be considered a criterion to define helpfulness of a review, we considered metascores as the parameter to define helpfulness of reviews. We have considered reviews having metascore greater than 5 as helpful.

```
db.allreviews.find({ $text: { $search: "good bad worse excellent price quality worst great delivery fast slow service time cost poor" } },{ "_id":0,"asin":1,score: { $meta: "textScore" } }).sort( { score: { $meta: "textScore" } } ).limit(5).pretty()
```

NOTE: For ease of comparison with the non-text-based approach, we have displayed only the top 5 products having the most helpful reviews.

```
{ "asin" : "B001YTK3XK", "score" : 10.278137943352384 }
{ "asin" : "B004H058VQ", "score" : 8.778493804372712 }
{ "asin" : "B00CMQTVUA", "score" : 8.5383478303362 }
{ "asin" : "B00AVWKLWU", "score" : 8.533585379464286 }
{ "asin" : "B003YL412U", "score" : 8.388176037901129 }
```

Non-Text based definition:

In Amazon, every review gets helpful votes. We have calculated the percentage of helpfulness of each review. We consider percentage of helpfulness as a criterion for judging helpfulness of a review because there might be cases when a review has got 100 votes out of which there are only 5 helpful votes. So, the percentage of helpfulness is only 5 %, while there can be a review which received totally 5 votes out of which all 5 are Helpful. So, this review has 100% helpfulness.

```

[{"id": {"Review": "", "ProductID": "B00005UFC1", "Total": 1, "helpfulVotes": 1}, "NotHelpful": 0, "Percentage": 100}, {"id": {"Review": "", "ProductID": "B00005AWBS", "Total": 2, "helpfulVotes": 2}, "NotHelpful": 0, "Percentage": 100}, {"id": {"Review": "", "ProductID": "B0000C9WI1", "Total": 3, "helpfulVotes": 3}, "NotHelpful": 0, "Percentage": 100}, {"id": {"Review": "", "ProductID": "B00000DMER", "Total": 2, "helpfulVotes": 2}, "NotHelpful": 0, "Percentage": 100}, {"id": {"Review": "", "ProductID": "B000087L1L", "Total": 8, "helpfulVotes": 8}, "NotHelpful": 0, "Percentage": 100}]

```

From the above 2 screenshots, we can see that the results obtained from text based and non-text-based approaches do not match. This is because in the text-based approach, we have considered our own defined set of words as a measure of helpfulness. In the non-text-based approach, we have taken helpful votes as the measure of helpfulness. As helpfulness varies according to users' perception, therefore, the two approaches will result in different outcomes.

## NEO4J – REVIEWERS REVIEWING COMMON SET OF PRODUCTS

### Wording of query

```

MATCH (reviewer:AResponse)-[:WRITTEN_BY]-(:AReview)-[:HAS_REVIEW]-(product:Product)
WITH reviewer, COLLECT(product.asin) AS commonProducts
WITH commonProducts, COLLECT(reviewers) as commonReviewers
WHERE SIZE(commonReviewers) > 1
RETURN commonProducts, commonReviewers

```

### Snapshot showing successful execution

\$ MATCH (reviewer:AResponse)-[:WRITTEN_BY]-(:AReview)-[:HAS_REVIEW]-(product:Product) WITH reviewer, COLLECT(product.asin) AS commonProducts, COLLECT(reviewers) as commonReviewers WHERE SIZE(commonReviewers) > 1 RETURN commonProducts, commonReviewers	
commonProducts	commonReviewers
["B00000JBLH"]	["A1F1A0QQ2PVKH5", "A32T2H81500JLU", "A2JFOHC9W629IE", "A2XRMQA6PJ5ZJ8", "AA8M6331NI1EN", "A49R5DBXXQD5E", "A38NELQT98S4H8", "A3MAFS04ZABRG0"]
["B00GSXJIMO"]	["ATM0311QH86HN", "A2H667ZUNLNTF1", "A2I4A4VGWII18Z", "A397CWU6DM305G", "A1JOANG98KBFKU", "A4EBYIKOMYFF9", "ALNFHVS3SC4FV", "A17C7XZHFOS12W", "AOM63D4VL32WM", "A3MKPAIK06U60K", "A14X244VGHWPXS", "ALQ4USPEQ9L5N", "A1GNYV0RA0EQSS", "ALDAF4VFLRH", "A2AYYKD40ODSLY", "A2NBCFDISGUQT", "A30QAOLNZ4415R", "ACR4HKUT808U1", "AWH2AY17ZU7W2", "A1Q08H4B9QGT2N", "A2OJCTT5WL2HR", "A2IFKH3T10387", "A27GITTN6AVW5I", "A2RVCF21L2NPUS", "A1ZVFCPHCWVF71", "A24RCBRDXRXR0Y", "A2503LT8PZIHAD", "A1E5WXNT7E02A", "A2OG83SEU9HL6D", "A3H5QGH50D2XGX", "AGFGY4EJ37VS2", "A2974R9BTZP0J", "ASEYZYO85D0TA", "A2B3A8NQIQUIT08", "A12E0Y0J6584RT", "A2LAS014NSRW72", "A2V9DTXTQ5YIMK", "A3W45GP0UKU8L", "A33V944SO2YCN", "A1GUN2J1FDYNF3", "A27BKXQAR305QL", "A34GM1776WTJDK", "AHNHLGKIZ2N3B", "A8INT3NMHRUS4", "A1II2ZRPKZAQQD"]
["B00000JBO8"]	["A2PWUD52M4VHGT", "A2ZX137EP1NT", "A0M1MEG343AP", "A3RRWUJU333R49", "A2Z9S2RQD542CP", "A346MAIT1GXHOH", "A3QMJMTLJC34QC", "A81JUYW2SL24", "A39XUL8UWT4015", "A304GUEPCPYM3Z"]
["B000FIOB2A", "B00GSXJIMO"]	["A3Q9HY0LU7MN6E", "A1QT3N5FUQ4I7"]
["B0006ZXFWO"]	["A32BCTL0HLOXRV", "A6A03NM42KUXL", "A300S5QIIMBNQH", "A1VOONTYMM0SDA", "A149F91FB9WTW6", "A3PGQWCJSJPCYDH", "A3A4UKJUP34E18", "A1FHOW9NV0H8XR", "A36IOSZXG6HOAU", "A2M0T1A400GS72", "A18P5I03P4UJA", "AUH2DKYUUT7K9", "A3ARMRKUJUQGLS", "A2RIQKV400XME1", "A3PQULEC090AV6", "A3QTFU2ZWH7WS"]
["B00074B67A"]	["A2MIP3AQVSF2SS", "A23HJIRAT5PIMG0", "A10GLT9EBDTK3N", "A3K6M4WVNNTYQEM"]
["B00074WF08"]	["A21DER2RMJDX7", "A25REOQJ9EDELE", "AKALIDNS3NNMZ", "AEHDGUE3L90ZS", "ALV01A5UB8DG0"]
["B000V825HM"]	["A12MB10SS6KE2", "AKJHHD5VEH7VG", "A3M606TB6E6C5N4", "A1JF919PT186ZI", "A1GJHYOS9T56GR", "AXARTAQU05PFL", "A1HVR12V7N8QUG", "A83Y6TKNOM5DQ"]
["B0007GGUGA"]	["A298INOY52JEH8", "A1LCG14GZJ5LO", "A27DR1V0079F1V", "A2GA2GM4LWESIC", "A3BTH66WVQCAM9", "A2G12DY50U700V", "A2QMNA0ZT6XMRR", "A39RL01WOC0Q0Q", "A2D9QW4TY8CFB0", "A2B5LSF5V8SSZA", "A16UCWYGY4924HH", "A3IS1RBMOOULWI", "A1BN3DBBWLMT6V", "A30J7WQV0ZNRXG", "A2E08GHKXEF141", "A46D3MTB5LIU", "AAGD3GA9ZVPLQ", "A3G6KKDDMMFQQ5K"]
["B0006ZXFZQ"]	["A2CJVLER896Q7L", "A3LBEIXA1JL28I", "A13UHRU7ZBXFL", "A3HIXX5UIDJ403N", "A9RAHOMYQAB51"]

Started streaming 14 records after 434 ms and completed after 434 ms.

### Summary of results

Out of all the paths in our graph model, we have matched the ones where we have a review being written by a reviewer and the product to which that review belongs. After that, given a reviewer, we have collected the list of products they have reviewed. Next, using this collection of products,

we have further collected the list of reviewers who have written these reviews. Then we have applied a condition that the size of the reviewers collection should be greater than one since we want to find the set of products being reviewed by a common set of reviewers and not just one reviewer. Finally, we return/ display the collections/ lists of the reviewers and the lists of products that they reviewed.

## SAS – IDENTIFYING FAKE REVIEWS

### **SAS Text mining:**

#### **Issue of fake reviews:**

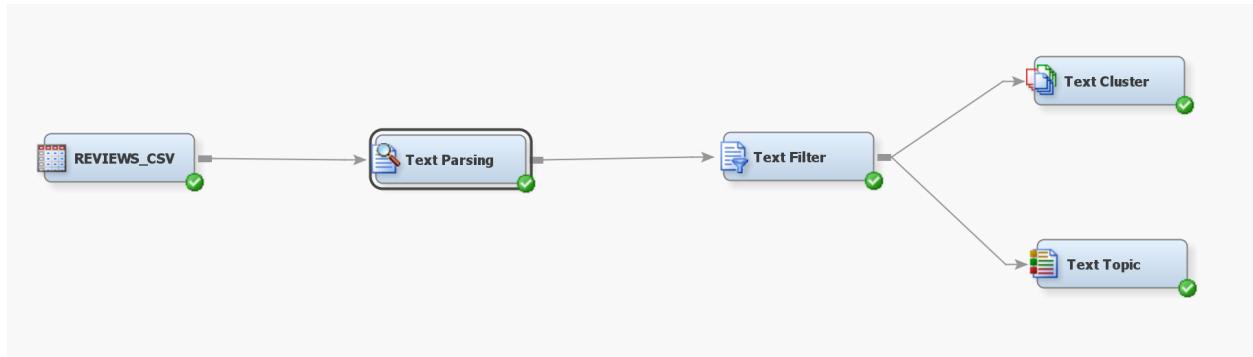
As the review systems are being extensively used by all the online businesses to provide the users best experience possible, some businesses, sellers and vendors of services and products are taking advantage of this system. In our analysis of the amazon review system, reportedly several sellers have been gaining fake positive reviews by offering free products and discounts to the customers. There are also cases where businesses target competitions and pay customers to provide fake negative reviews to their products. This is deteriorating the entire usability of the review system, providing bad experiences to users which, when reviewed are getting lost in the fraudulent reviews of that product. Reportedly, almost 100 facebook groups were found split up in several geographic locations and by categories, where merchants are continuously asking customers to provide fraud reviews.

To analyze the fake reviews in some of the amazon's reviews, we chose to use some extreme emotion carrying adjectives like excellent, worst, etc. Our explanation of doing this is that fake reviews generally tend to lack details of the product but just use words like this to represent the positiveness and high ratings. To implement our idea, we have converted the CSV file into a SAS table using SAS Enterprise Guide 7.1. In SAS Enterprise, we created a library and data source for this file and imported the table.

#### **Data Dictionary:**

<b>id</b>	<b>ID</b>	<b>Nominal</b>	<b>No</b>	<b>No</b>	.	.
<b>asin</b>	<b>Inout</b>	<b>Nominal</b>	<b>No</b>	<b>No</b>	.	.
<b>helpful</b>	<b>Inout</b>	<b>Nominal</b>	<b>No</b>	<b>No</b>	.	.
<b>overall</b>	<b>Inout</b>	<b>Interval</b>	<b>No</b>	<b>No</b>	.	.
<b>reviewText</b>	<b>Text</b>	<b>Nominal</b>	<b>No</b>	<b>No</b>	.	.
<b>reviewTime</b>	<b>Inout</b>	<b>Interval</b>	<b>No</b>	<b>No</b>	.	.
<b>reviewerID</b>	<b>Inout</b>	<b>Nominal</b>	<b>No</b>	<b>No</b>	.	.
<b>summary</b>	<b>Text</b>	<b>Nominal</b>	<b>No</b>	<b>No</b>	.	.
<b>unixReview</b>	<b>Inout</b>	<b>Interval</b>	<b>No</b>	<b>No</b>	.	.

## Methodology:



To implement our idea, we have converted the CSV file into a SAS Enterprise Guide 7.1. In SAS Enterprise, we created a library and data source for this file and imported the table.

### Text Parsing:

After we imported the data, we use the text parsing node to clean the data and make some necessary modifications for further processing. We used the properties panel of the text parsing node to:

1. Set the 'find entities' option to standard.
2. Set 'Detect parts of speech' to NO.
3. Ignore all the entities except Miscellaneous proper nouns (prop\_misc).
4. Ignore all the parts of speech except adjectives.
5. Ignore numbers and punctuations.

After running the text parsing node, we generated the list of terms with their frequencies, to look and understand how it can be better processed by identifying the terms that are rarely used.

### Output of text parsing node:

Term	Role	Attribute	Freq	# Docs	Keep	Parent/Child Status	Parent ID	Rank for Variable numdocs
+ be	Alpha		16202	38615	+		270861	
+ have	Alpha		79468	262000	+		270838	
+ like	Alpha		54760	55717	+		270819	
+ do	Alpha		32892	53871	+		270848	
+ good	Alpha		32892	53871	+		270887	
s work	Alpha		19303	39925	+		270893	
very	Alpha		20253	42084	+		270900	
not	Alpha		19303	39925	+		270907	
last	Alpha		19303	39925	+		270920	
most	Alpha		19303	39925	+		270924	
less	Alpha		19303	39925	+		270926	
+ easier	Alpha		16202	11417	+		270930	
+ think	Alpha		16202	11417	+		270934	
+ little	Alpha		16202	10317	+		270936	
also	Alpha		16202	10680	+		270937	
well	Alpha		16202	10680	+		270938	
+ product	Alpha		16202	11317	+		270939	
+ time	Alpha		32892	56021	+		270942	
+ color	Alpha		32892	56021	+		270950	
real	Alpha		17950	35924	+		270953	
self	Alpha		14710	33234	+		270954	
point	Alpha		14710	33234	+		270955	
+ come	Alpha		13102	31819	+		270957	
+ that	Alpha		12956	31610	+		270958	
one	Alpha		10361	26101	+		270960	
+ under	Alpha		9310	20000	+		270961	
+ buy	Alpha		9310	20000	+		270962	
+ want	Alpha		16202	77987	+		190128	
now	Alpha		16202	77987	+		270963	
+ price	Alpha		16202	71187	+		270978	
+ think	Alpha		16202	71187	+		270979	
+ quality	Alpha		16202	70242	+		270980	
+ print	Alpha		16202	70242	+		270981	
+ take	Alpha		8474	17373	+		270987	
try	Alpha		8474	16884	+		270988	
too	Alpha		8474	16884	+		270989	
+ recommend	Alpha		11186	21187	+		270990	
+ write	Alpha		11186	21187	+		270991	
+ see	Alpha		8956	20101	+		270992	
+ even	Alpha		8956	20101	+		270993	
+ go	Alpha		8956	20101	+		270994	
there	Alpha		8956	20101	+		270995	
+ feel	Alpha		26917	61117	+		270996	
+ keep	Alpha		26917	61117	+		270997	
seem	Alpha		26917	61117	+		270998	
easy	Alpha		7956	14444	+		270999	
hard	Alpha		7956	14444	+		271000	
over	Alpha		6317	8289	+		271002	
overcom	Alpha		6317	8289	+		271003	
back	Alpha		6317	8289	+		271004	
out	Alpha		6317	8289	+		271005	
all	Alpha		1516	3693	+		271006	
+ case	Alpha		1516	3693	+		271007	
+ office	Alpha		1516	3693	+		271008	
+ very	Alpha		1516	3693	+		271009	
+ nice	Alpha		1516	3693	+		271010	
+ sit	Alpha		1516	3693	+		271011	
+ first	Alpha		1516	3693	+		271012	
through	Alpha		1516	3693	+		271013	
still	Alpha		1516	3693	+		271014	
role	Alpha		1516	3693	+		271015	
box	Alpha		10498	20101	+		271016	
new	Alpha		10498	20101	+		271017	
long	Alpha		6066	20101	+		271018	
late	Alpha		6066	20101	+		271019	
what	Alpha		5806	20101	+		271020	
name	Alpha		15440	20101	+		185920	
row	Alpha		15440	20101	+		112389	
set	Alpha		9546	44793	+		271021	
get	Alpha		9546	44793	+		271022	
when	Alpha		7334	46217	+		271023	
when	Alpha		5573	45097	+		69724	
fine	Alpha		5573	45097	+		15010	
none	Alpha		5573	45097	+		271024	
decade	Alpha		5573	45097	+		271025	

## Text filter:

We added a text filter node, to drop the terms that have less weight, by entering the number of documents that word should be present in order to be considered. We used the ‘filter viewer’ in the properties of text filter node to keep words that would best fit for this analysis manually and combined some words and grouped them as synonyms. We performed spell check by selecting the ‘check spelling’ option as yes from the properties.

We also performed concept linking on some of the words, just to see if we can analyze something from the data. The outputs of the text filtering node and the concept linking are as below.

## Text Filtering Output:

Term	Role	Attribute	Status	Weight	Imported Frequency	Freq	Number of Imported Documents	# Docs	Rank	Parent/Child Status	Parent ID
+ had				0.117	1583	138	1479	108	1+		1327723
+ little				Alpha	100	17850	1171	13070	103	1++	240496
+ good				Alpha	100	8500	54	7167	43	1++	2207721
+ other				Alpha	100	8500	54	7167	43	1++	2207721
+ much				Alpha	100	8500	54	7167	43	1++	2207721
+ us				Alpha	100	8500	54	7167	43	1++	2207721
+ only				Alpha	100	8500	54	7167	43	1++	2207721
+ new				Alpha	100	8500	54	7167	43	1++	2207721
+ more				Alpha	100	8500	54	7167	43	1++	2207721
+ now				Alpha	100	8500	54	7167	43	1++	2207721
+ never				Alpha	100	8500	54	7167	43	1++	2207721
+ travel				Alpha	100	8500	54	7167	43	1++	2207721
+ travel				Alpha	100	8500	54	7167	43	1++	2207721
+ man				Alpha	100	8500	54	7167	43	1++	2207721
+ man				Alpha	100	8500	54	7167	43	1++	2207721
+ man				Alpha	100	8500	54	7167	43	1++	2207721
+ man				Alpha	100	8500	54	7167	43	1++	2207721
+ man				Alpha	100	8500	54	7167	43	1++	2207721
+ different	Miscellaneous Prop.	Entity		Alpha	100	293	24	361	298	1+	117452
+ different				Alpha	100	293	24	361	298	1+	117452
+ hard				Alpha	100	1845	30	1560	227	2++	159132
+ fast				Alpha	100	1845	30	1560	227	2++	159132
+ several				Alpha	100	289	23	2173	203	2++	159132
+ several				Alpha	100	289	23	2173	203	2++	159132
+ expensive				Alpha	100	289	23	2173	203	2++	159132
+ expensive	Miscellaneous Prop.	Entity		Alpha	100	289	23	2173	203	2++	159132
+ expensive				Alpha	100	289	23	2173	203	2++	159132
+ bid				Alpha	100	483	30	4073	222	2++	117452
+ bid				Alpha	100	483	30	4073	222	2++	117452
+ all				Alpha	100	3186	24	2848	224	2++	210552
+ unless				Alpha	100	3186	24	2848	224	2++	210552
+ sure				Alpha	100	3186	24	2848	224	2++	210552
+ public				Alpha	100	3186	24	2848	224	2++	210552
+ happy				Alpha	100	3186	24	2848	224	2++	210552
+ clean				Alpha	100	3186	24	2848	224	2++	210552
+ how				Alpha	100	3186	24	2848	224	2++	210552
+ ink cartridges	Noun Group			Alpha	0.440	2668	21	2438	20	3++	88695
+ simple				Alpha	100	3561	24	1967	1930	4++	250501
+ simple				Alpha	100	3561	24	1967	1930	4++	250501
+ fast				Alpha	100	3548	29	1817	1800	4++	120673
+ fast				Alpha	100	3548	29	1817	1800	4++	120673
+ home				Alpha	100	2634	24	2240	177	4++	81788
+ home				Alpha	100	2634	24	2240	177	4++	81788
+ easier				Alpha	100	1500	14	1463	165	4++	23111
+ difficult				Alpha	100	1500	14	1463	165	4++	23111
+ difficult				Alpha	100	1500	14	1463	165	4++	23111
+ available				Alpha	100	1941	21	1970	1510	4++	270890
+ available				Alpha	100	1941	21	1970	1510	4++	270890
+ quality				Alpha	0.491	2059	20	1871	1727	4++	10113
+ perfect				Alpha	100	1941	20	1871	1727	4++	10113
+ perfect				Alpha	100	1941	20	1871	1727	4++	10113
+ standard				Alpha	0.490	2157	18	1843	1424	5++	126144
+ standard				Alpha	100	2157	18	1843	1424	5++	126144
+ real				Alpha	0.489	1174	14	1163	143	5++	123389
+ rain				Alpha	100	1174	14	1163	143	5++	123389
+ minute				Alpha	0.487	2710	14	2472	1265	5++	861193
+ excellent				Alpha	100	2710	14	2472	1265	5++	1260443
+ excellent				Alpha	100	2710	14	2472	1265	5++	1260443
+ right				Alpha	0.487	1508	14	1324	1150	5++	1260443
+ right				Alpha	100	1508	14	1324	1150	5++	1260443
+ clear				Alpha	0.487	2770	12	2069	172	5++	54728
+ clear				Alpha	100	2770	12	2069	172	5++	54728
+ full				Alpha	0.486	1424	14	1323	1151	5++	150573
+ full				Alpha	100	1424	14	1323	1151	5++	150573
+ accurate				Alpha	0.486	1424	14	1323	1151	5++	150573
+ accurate				Alpha	100	1424	14	1323	1151	5++	150573
+ thus				Alpha	0.486	1526	14	1366	1152	5++	144064
+ thus				Alpha	100	1526	14	1366	1152	5++	144064
+ light				Alpha	0.486	1860	14	1611	1123	5++	150767
+ light				Alpha	100	1860	14	1611	1123	5++	150767
+ enough				Alpha	0.486	1109	14	993	1124	5++	150767
+ enough				Alpha	100	1109	14	993	1124	5++	150767
+ frame				Alpha	0.486	1109	14	993	1124	5++	150767
+ frame				Alpha	100	1109	14	993	1124	5++	150767
+ perfect				Alpha	0.486	1109	14	993	1124	5++	150767
+ extra				Alpha	0.486	1420	14	1301	111	5++	68554
+ extra				Alpha	100	1420	14	1301	111	5++	68554
+ acceptable				Alpha	0.485	1253	14	1062	111	5++	116373
+ acceptable				Alpha	100	1253	14	1062	111	5++	116373
+ bad thing	Noun Group			Alpha	0.485	1111	14	883	111	5++	120529
+ bad thing				Alpha	100	1111	14	883	111	5++	120529
+ second				Alpha	0.484	1110	14	1016	10	5++	141773
+ second				Alpha	100	1110	14	1016	10	5++	141773
+ comfortable				Alpha	0.484	573	14	1173	100	5++	150367
+ comfortable				Alpha	100	573	14	1173	100	5++	150367
+ still				Alpha	0.484	1379	14	1289	100	5++	220729
+ still				Alpha	100	1379	14	1289	100	5++	220729
+ possible				Alpha	0.484	1624	14	860	99	5++	220729
+ possible				Alpha	100	1624	14	860	99	5++	220729
+ wonderful				Alpha	0.484	1624	14	860	99	5++	220729
+ wonderful				Alpha	100	1624	14	860	99	5++	220729
+ count quality	Noun Group			Alpha	0.484	631	14	620	99	5++	11424
+ count quality				Alpha	100	631	14	620	99	5++	11424
+ black and white				Alpha	0.484	631	14	620	99	5++	11424
+ black and white				Alpha	100	631	14	620	99	5++	11424
+ back				Alpha	0.484	1216	14	1093	99	5++	104923
+ back				Alpha	100	1216	14	1093	99	5++	104923
+ reason	Miscellaneous Prop.	Entity		Alpha	0.484	204	14	162	99	5++	117452

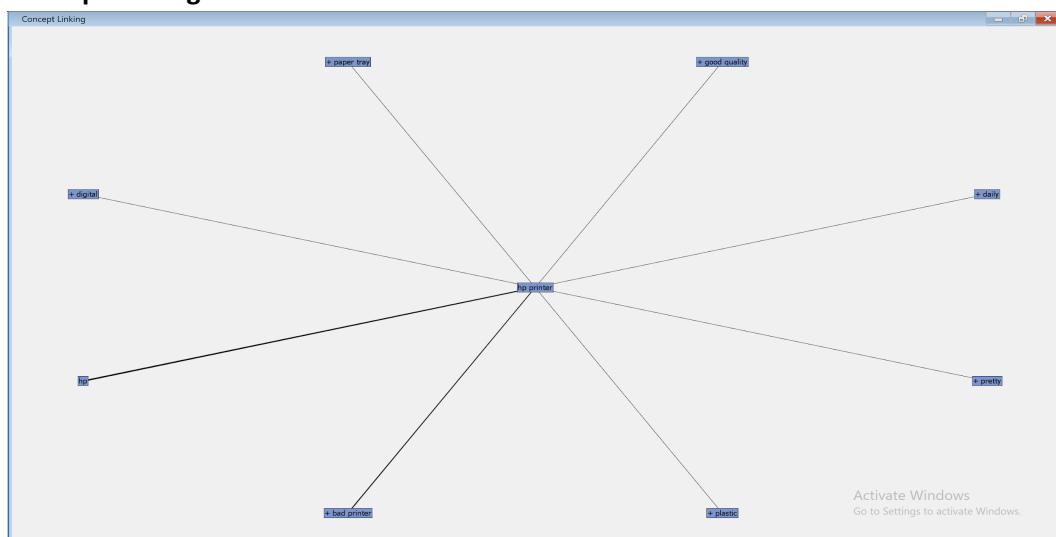
## Manually keeping and dropping terms:

Terms							
TERM	FREQ	# DOCS	KEEP	WEIGHT ▾	ROLE	ATTRIBUTE	Alpha
expensive printer	6	5	✓	0.854	Noun Group	Alpha	
expensive scanner	6	5		0.854	Noun Group	Alpha	
file storage	6	5		0.854	Noun Group	Alpha	
inexpensive laser printer	6	5		0.854	Noun Group	Alpha	
mock	6	5	✓	0.854	Noun Group	Alpha	
excellent photo quality	6	5	✓	0.854	Noun Group	Alpha	
mj	18	6		0.854	Miscellaneous ...	Entity	
wf-7010	22	7		0.854	Miscellaneous ...	Entity	
epson wp-4540	15	7		0.852	Miscellaneous ...	Entity	
epson expression home xp-4...	9	6		0.852	Miscellaneous ...	Entity	
lubricant	9	6		0.852		Alpha	
introductory	9	6		0.852		Alpha	
staple puller	9	6		0.852	Noun Group	Alpha	
mp990	9	6		0.852	Miscellaneous ...	Entity	
acrobail	20	6		0.852	Miscellaneous ...	Entity	
old at&t	15	6		0.852	Noun Group	Mixed	
nx625	17	8		0.852	Miscellaneous ...	Entity	
transitional notebook paper	12	7		0.851	Noun Group	Alpha	
firmware version	11	6		0.851	Noun Group	Alpha	
black leatherette	11	6		0.851	Noun Group	Alpha	
claria ink	11	6		0.851	Miscellaneous ...	Entity	
hl-5370dw	11	6		0.851	Miscellaneous ...	Entity	
dear	15	8		0.851		Alpha	
cl-241xl	15	6		0.851	Miscellaneous ...	Entity	
metal arm	13	6		0.85	Noun Group	Alpha	
battery pack	10	6		0.85	Noun Group	Alpha	
workforce series	10	6		0.85	Miscellaneous ...	Entity	
w x d	10	6		0.85	Miscellaneous ...	Entity	
orange light	10	6		0.85	Noun Group	Alpha	
unbroken	5	5		0.85		Alpha	
avery ticket	5	5		0.85	Noun Group	Alpha	
stress test	5	5		0.85	Noun Group	Alpha	
everyday usage	5	5		0.85	Noun Group	Alpha	
cheap kind	5	5	✓	0.85	Noun Group	Alpha	
printer unit	5	5	✓	0.85	Noun Group	Alpha	
fairly long time	5	5	✓	0.85	Noun Group	Alpha	

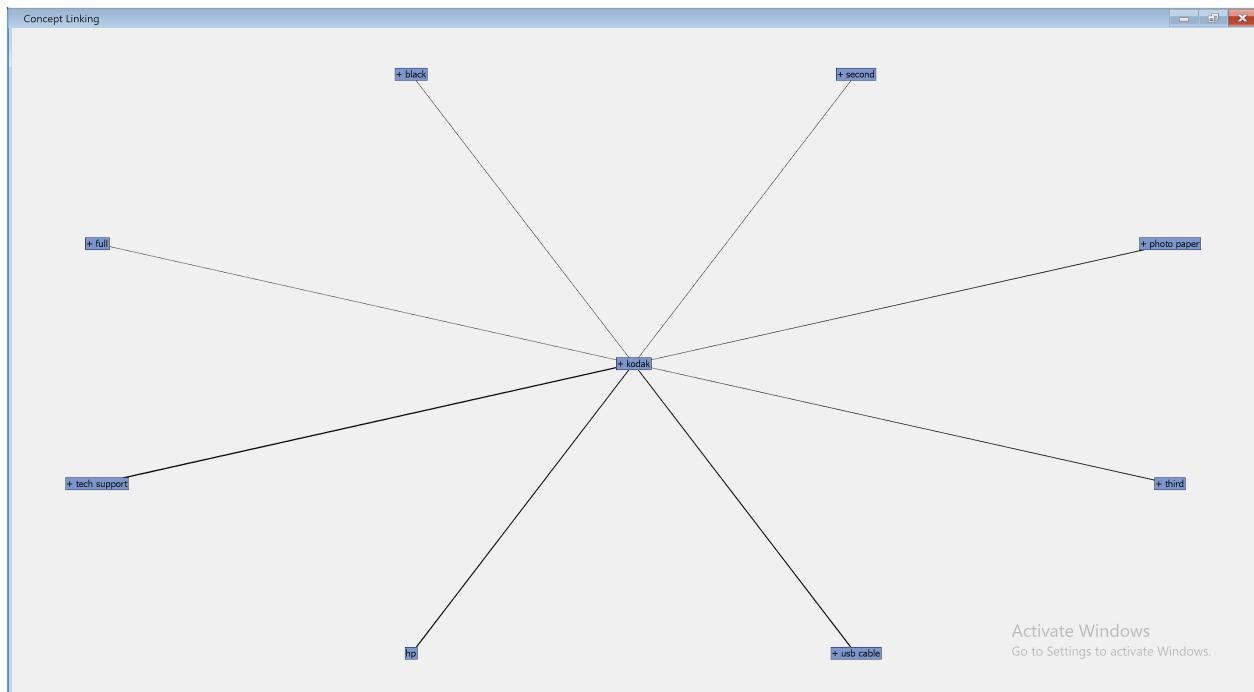
## Concept linking:

We performed concept linking on the word 'HP Printer' and it shows that there's a darker line towards the words 'bad printer' and 'plastic'. We found this analysis interesting and performed this on other printer brands Epson and Kodak and inferred that epson had more stronger connection to words 'wireless' and 'automatic', indicating its features. Kodak had stronger connection to the word 'tech support', from which we inferred that the product might have more issues.

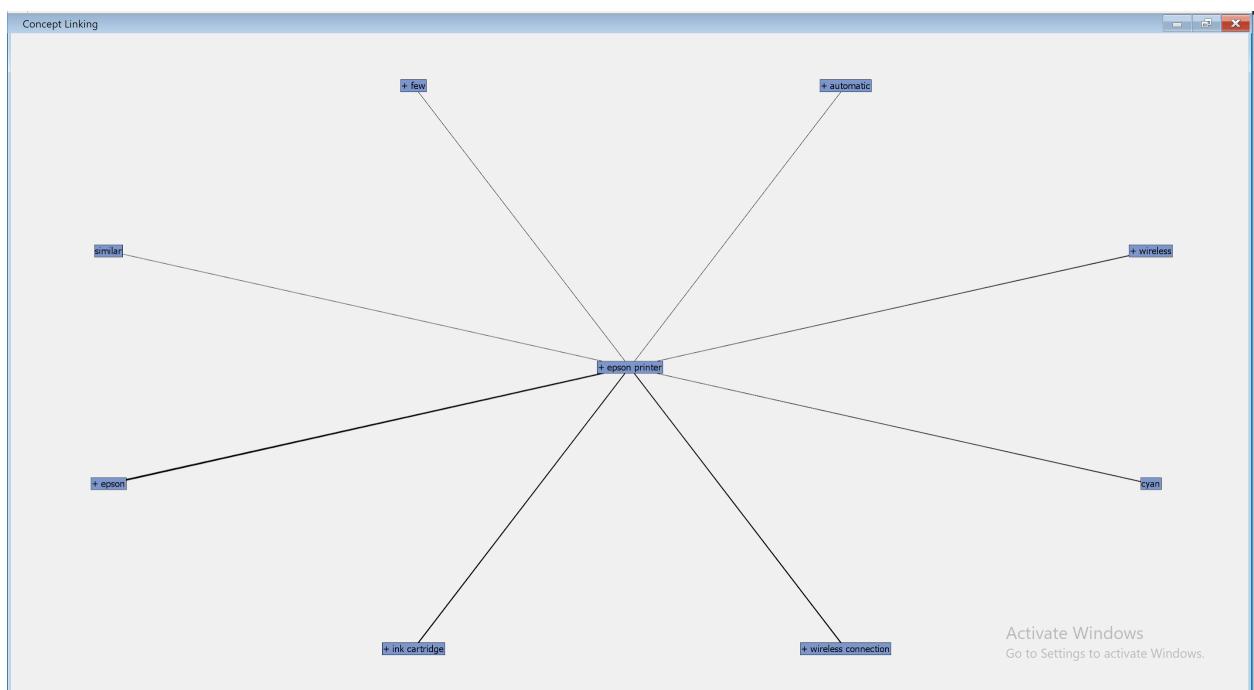
## Concept Linking for HP:



## Concept Linking for KODAK:

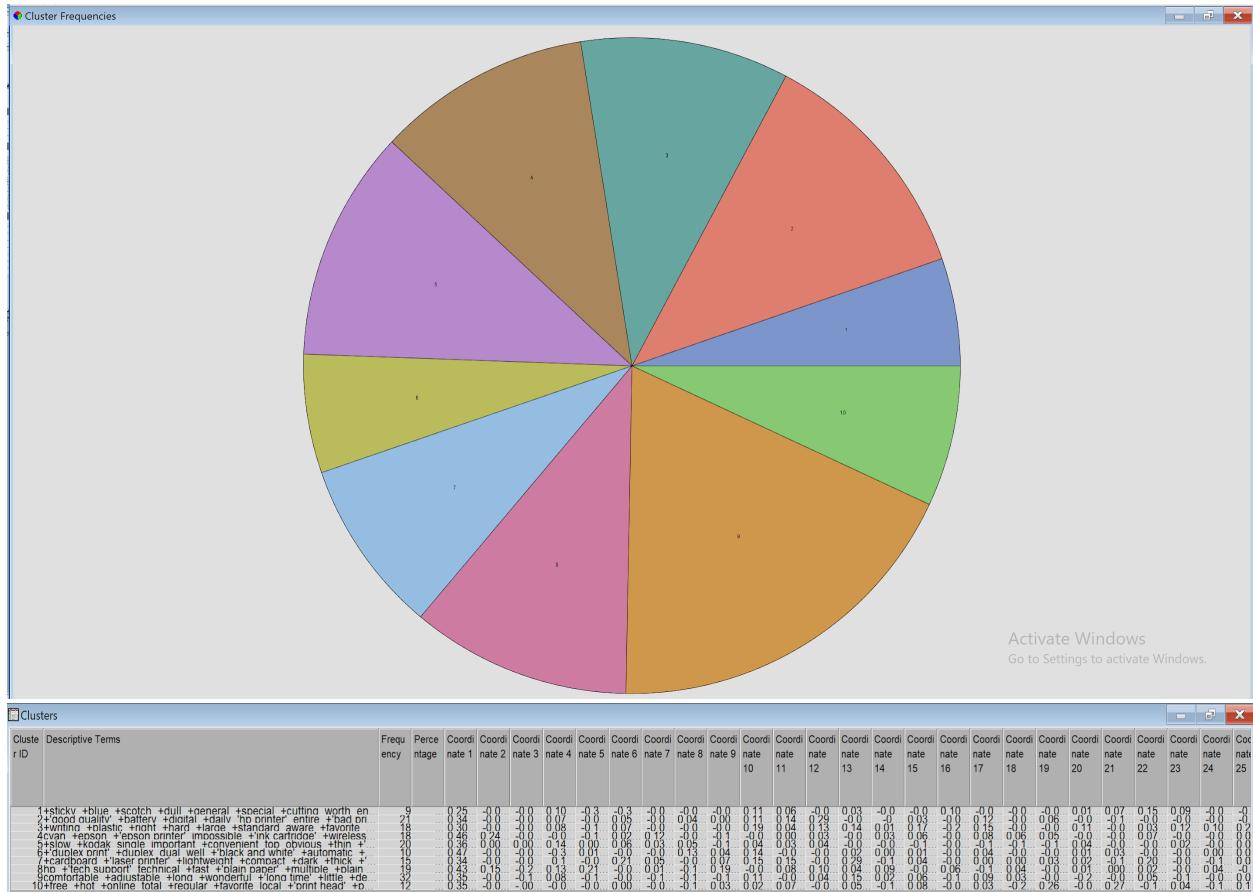


## Concept linking for EPSON:



## Text Clustering:

We attached a text clustering node to display the text in clusters to better understand the trend. We used the properties panel of the text clustering node to generate exactly 10 clusters. The outputs of this node are as below.



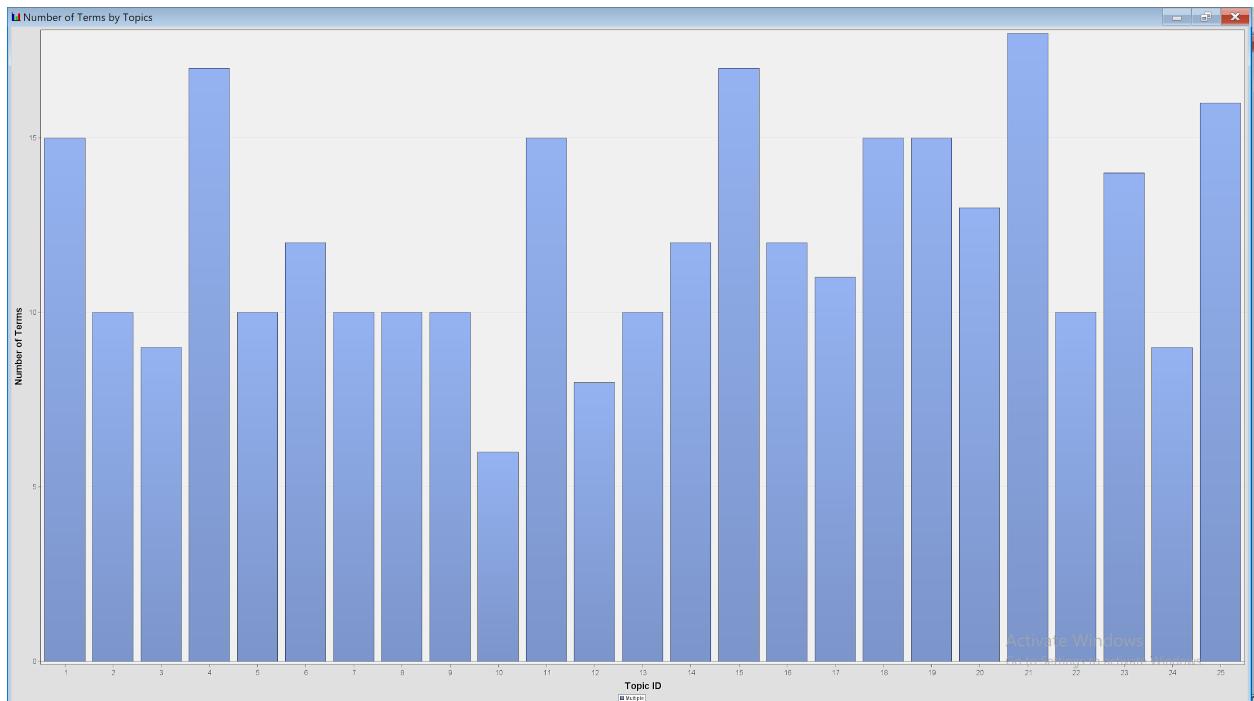
## Text Topic:

We used text topic node to combine words into topics that can help analyze the data further. We used the outputs to plot graphs between the topics and the documents, topics and terms.

## Topics:

Category	Topic ID	Document Cutoff	Term Cutoff	Topic	Number of Terms	# Docs
Multiple	1	0.192	0.111	+low +heavy +dark actual +soft	15	17
Multiple	2	0.172	0.112	+black +white +white h	19	15
Multiple	3	0.172	0.112	+sticky +scotch +blue worth +second	19	15
Multiple	4	0.143	0.113	+duplex +duplex print +automatic +	10	10
Multiple	5	0.163	0.114	+laser printer +cardboard +fax mac	10	10
Multiple	6	0.131	0.115	+top +bottom +cardboard running +	10	10
Multiple	7	0.130	0.116	+ink cartridge +toner printer single	10	10
Multiple	8	0.130	0.117	+battery +cutting entire +dull +thin	19	15
Multiple	9	0.130	0.118	+expensive impossible total +hand	19	15
Multiple	10	0.130	0.119	+bottom +real +cheap separate +front	19	15
Multiple	11	0.130	0.120	+bad part +little +bad +lightweight	21	20
Multiple	12	0.130	0.121	+print head +complicated +head +tie	20	15
Multiple	13	0.130	0.122	+scotch +difficult +desk +thick +tie	20	15
Multiple	14	0.130	0.123	+no no center +bad center +ink car	20	15
Multiple	15	0.130	0.124	+lightweight +break important +elec	20	15

### Terms by Topics:



### Terms by Documents:

