Definition of a problem statement and a short outline of the implementation

   Problem statement

Description of data acquisition / how it was collected (by you or the publisher of the data)

1. Data preparation (general)

2. Natural Language Processing

3. Supervised / Unsupervised ML

4. Network Analysis

# M2 - Exam - Network and NLP

`Code ▾`

## Malika Kuhlman & Michael Dybdahl

## 10/24/2019

**OBS** The notebook cannot be run in google Colab because of the *topicmodels* package.

Loading the *stop_words* that we will use in the tokenization.

`Hide`

```
set.seed(9)
data("stop_words")
```

# Definition of a problem statement and a short outline of the implementation

The link to our Github: https://github.com/michaeldybdahl/M2_exam (https://github.com/michaeldybdahl/M2_exam)

# Problem statement

How are the words in the text-reviews and based on a tidy format of the reviews, can we then compute any topics based on this? Can we compute a model that predict if a review is rated as a good review, based or linked to the NLP results?

# Description of data acquisition / how it was collected (by you or the publisher of the data)

The dataset contains reviews of some of amazons products, the dataset can be reached on this website https://data.world/datafiniti/consumer-reviews-of-amazon-products (https://data.world/datafiniti/consumer-reviews-of-amazon-products)

There are serveral datasets about the same topic, but with different amount of observations, we chose the dataset with 5000 observations. The dataset contains 24 variables and 5.000 observations (reviews).

<div style="text-align: right;">Hide</div>

```
reviews <- read_csv("https://query.data.world/s/foqg5o75hazenbwqdug534atoqiy
p3")
```

Have a look at the data structure

<div style="text-align: right;">Hide</div>

```
glimpse(reviews)
```

```
## Observations: 5,000
## Variables: 24
## $ id                   <chr> "AVqVGZNvQMlgsOJE6eUY", "AVqVGZNvQMlgsOJE6eU…
## $ dateAdded            <dttm> 2017-03-03 16:56:05, 2017-03-03 16:56:05, 2…
## $ dateUpdated          <dttm> 2018-10-25 16:36:31, 2018-10-25 16:36:31, 2…
## $ name                 <chr> "Amazon Kindle E-Reader 6\" Wifi (8th Genera…
## $ asins                <chr> "B00ZV9PXP2", "B00ZV9PXP2", "B00ZV9PXP2", "B…
## $ brand                <chr> "Amazon", "Amazon", "Amazon", "Amazon", "Ama…
## $ categories           <chr> "Computers,Electronics Features,Tablets,Elec…
## $ primaryCategories    <chr> "Electronics", "Electronics", "Electronics",…
## $ imageURLs            <chr> "https://pisces.bbystatic.com/image2/BestBuy…
## $ keys                 <chr> "allnewkindleereaderblack6glarefreetouchscre…
## $ manufacturer         <chr> "Amazon", "Amazon", "Amazon", "Amazon", "Ama…
## $ manufacturerNumber   <chr> "B00ZV9PXP2", "B00ZV9PXP2", "B00ZV9PXP2", "B…
## $ reviews.date         <dttm> 2017-09-03 00:00:00, 2017-06-06 00:00:00, 2…
## $ reviews.dateAdded    <lgl> NA, NA, NA, NA, NA, NA, NA, NA, NA, NA, NA, …
## $ reviews.dateSeen     <chr> "2018-05-27T00:00:00Z,2017-09-18T00:00:00Z,2…
## $ reviews.doRecommend  <lgl> FALSE, TRUE, TRUE, TRUE, TRUE, FALSE, TRUE, …
## $ reviews.id           <dbl> NA, NA, NA, 177283626, NA, NA, 187043823, NA…
## $ reviews.numHelpful   <dbl> 0, 0, 0, 3, 0, 0, 0, 0, 0, 0, 3, 0, 0, 0, 0,…
## $ reviews.rating       <dbl> 3, 5, 4, 5, 5, 5, 5, 4, 5, 5, 5, 5, 5, 5, 5,…
## $ reviews.sourceURLs   <chr> "http://reviews.bestbuy.com/3545/5442403/rev…
## $ reviews.text         <chr> "I thought it would be as big as small paper…
## $ reviews.title        <chr> "Too small", "Great light reader. Easy to us…
## $ reviews.username     <chr> "llyyue", "Charmi", "johnnyjojojo", "Kdperry…
## $ sourceURLs           <chr> "https://www.newegg.com/Product/Product.aspx…
```

# 1. Data preparation (general)

## 1.1. Data cleaning

We see that the data contains of 5000 observations and 24 variables, we are not going to use all of variables, therefore we now select the one we will use. Therefore we will only describe the variables we will use.

<div style="text-align: right;">Hide</div>

```
reviews <- reviews %>%
  select(reviews.username,
         dateAdded,
         name,
         primaryCategories,
         reviews.numHelpful,
         reviews.rating,
         reviews.text)
```

- **reviews.username:** contain the reviewsers username.
- **dateAdded:** coontain the date the reviews was written.
- **name:** contain the name of the product.
- **primaryCategories:** contain the PRimary category of the product, there are four different primary categories in our data; *Electronics*, *Electronics, Hardware*, *Electronics, Media*and *Office Supplies, Electronics*.
- **reviews.numHelpful:** contain a the value wether or not a reader found the review helpful.
- **reviews.text:** contain the reviews of the product.

Now we check for *NA* and drop them if any exists.

Hide

```
sum(is.na(reviews)) # First we check for NA values in the whole datasat.
```

```
## [1] 0
```

Hide

```
reviews = reviews %>%
  drop_na()
```

# 1.2. Recoding

We saw earlier that some of the variable are named *reviews.X*, this we want to change.

Hide

```
reviews <- reviews %>%
  rename(helpful = reviews.numHelpful,
         rating = reviews.rating,
         text = reviews.text,
         username = reviews.username,
         date = dateAdded)
```

Have a look at the dataset now.

Hide

```
glimpse(reviews)
```

```
## Observations: 5,000
## Variables: 7
## $ username         <chr> "llyyue", "Charmi", "johnnyjojojo", "Kdperry",…
## $ date             <dttm> 2017-03-03 16:56:05, 2017-03-03 16:56:05, 201…
## $ name             <chr> "Amazon Kindle E-Reader 6\" Wifi (8th Generati…
## $ primaryCategories <chr> "Electronics", "Electronics", "Electronics", "…
## $ helpful          <dbl> 0, 0, 0, 3, 0, 0, 0, 0, 0, 0, 3, 0, 0, 0, 0, 0…
## $ rating           <dbl> 3, 5, 4, 5, 5, 5, 4, 5, 5, 5, 5, 5, 5, 5, 5, 3…
## $ text             <chr> "I thought it would be as big as small paper b…
```

Now we have a dataset with 4987 observations and 9 variables

For later use, we add a new *ID* number to each username, since som users have reviewed more than one product.

Hide

```
username_df <- reviews %>%
  select(username) %>%
  distinct(username,
           .keep_all = TRUE)

username_df$ID <- seq.int(nrow(username_df))
```

Now we join the new ID with the dataset.

Hide

```
reviews <- reviews %>%
  left_join(username_df, by = "username") %>%
  select(ID, everything())

reviews$ID2 <- seq.int(nrow(reviews))
```

We change the rating from *doubble* to *factor*.

Hide

```
reviews$rating <- as.factor(reviews$rating)
```

We change the date from *dttm* to *date*.

Hide

```
reviews$date <- as.Date(as.POSIXct(reviews$date))
```

We create a new variable called *year*, where we put the year of every review.

Hide

```
reviews <- reviews %>%
  mutate(year = as.factor(year(date)))
```

Now we want to look into wich years the dataset has reviews from and how many reviews there are written in each year.

Hide

```
reviews %>%
  group_by(year) %>%
  count()
```

| year | n |
|------|---|
| <fctr> | <int> |
| 2015 | 44 |
| 2016 | 1026 |
| 2017 | 2495 |
| 2018 | 1435 |
| 4 rows | |

We see we have four years

- **2015** with 44 reviews
- **2016** with 1026 reviews
- **2017** with 2495 reviews
- **2018** with 1435 reviews

# 2. Natural Language Processing

We now do all the NLP here, and afterwards we have a section containing all the Network analysis.

# 2.1. NLP - preparation

## 2.1.1. Tokenization

In this section we are going to do a tokenization by removing *stop words*, *meaningless words*, *non-alphanumeric characters*, and so on.

First we seperate every word in each review text.

Hide

```
reviews_tidy <- reviews %>%
  select(ID2, ID, text) %>%
  unnest_tokens(output = word, input = text)
```

Hide

```
head(reviews_tidy)
```

| ID2 | ID | word |
|-----|-----|------|
| <int> | <int> | <chr> |
| 1 | 1 | i |
| 1 | 1 | thought |
| 1 | 1 | it |

| ID2 <int> | ID <int> | word <chr> |
|---|---|---|
| 1 | 1 | would |
| 1 | 1 | be |
| 1 | 1 | as |

6 rows

Then we remove *stop words*, using the stop_words we loaded in the begining.

Hide

```
reviews_tidy %<>% # I use "%<>%" which both assigns and pipe - instead of fi
rst assigning and then piping.
  anti_join(stop_words, by = "word") # I use anti_join instead of filter, be
cause it should be a bit faster.
```

Now we have a look at the most used words.

Hide

```
reviews_tidy %>%
  count(word, sort = TRUE) %>% # Count "word". "sort = TRUE" means that it w
ill sort the words in descending order of number words.
  head(20)
```

| word <chr> | n <int> |
|---|---|
| tablet | 1309 |
| love | 1090 |
| easy | 822 |
| bought | 785 |
| kindle | 764 |
| amazon | 694 |
| echo | 693 |
| alexa | 513 |
| loves | 506 |
| screen | 500 |

1-10 of 20 rows           Previous   **1**   2   Next

We see that words like *tablet*, *love* and *easy* are the most used words. Since the reviews are about *Electronics* we think thats why the word *Tablet* are the most used word togehter with *kindle*, *alexa* and *echo* wich all are names of some of the products.

After a closer look into the words, we found some words that could be removed.

Hide

```
own_stopwords <- tibble(word = c("bought", "yrs", "yokod", "yo", "xm", "woul
dnt", "wouldn", "woo", "withunderstandably", "wished", "wi", "wellthis", "we
llreasonably", "ample", "amazonso", "amazingly", "äîand"),
                         lexicon = "OWN")
```

Hide

```
reviews_tidy %<>% # I use "%<>%" which both assigns and pipe - instead of fi
rst assigning and then piping.
  anti_join(own_stopwords, by = "word") # I use anti_join instead of filter,
because it should be a bit faster.
```

Now we look at the top words again.

Hide

```
reviews_tidy %>%
  count(word, sort = TRUE) %>% # Count "word". "sort = TRUE" means that it w
ill sort the words in descending order of number words.
  head(20)
```

| word | n |
| --- | --- |
| <chr> | <int> |
| tablet | 1309 |
| love | 1090 |
| easy | 822 |
| kindle | 764 |
| amazon | 694 |
| echo | 693 |
| alexa | 513 |
| loves | 506 |
| screen | 500 |
| price | 477 |
| 1-10 of 20 rows | Previous **1** 2 Next |

We see that the top words are different, since we removed some of our own stop words.

Here we remowe numbers.

Hide

```
reviews_tidy %<>%
  mutate(word = trimws(gsub("[^\\s]*[0-9][^\\s]*", "", word, perl = T))) %>%
  filter(str_length(word) > 1) # this filter out the words that are blank.
```

Here we remowe non-alphanumeric characters.

Hide

```
reviews_tidy %<>%
  mutate(word = word %>% str_remove_all("[^[:alnum:]]") ) %>% # alnum = Alph
anumeric characters.
  filter(str_length(word) > 1) # filter out words with 1 character.
```

Now we look at the top $20$ of the least used words.

Hide

```
reviews_tidy %>%
  count(word, sort = TRUE) %>% # Count "word". "sort = TRUE" means that it w
ill sort the words in descending order of number words.
  top_n(-20)
```

| word | n |
| :--- | ---: |
| <chr> | <int> |
| abke | 1 |
| absent | 1 |
| absorbs | 1 |
| accelerometer | 1 |
| accept | 1 |
| acceptable | 1 |
| accesible | 1 |
| accesses | 1 |
| accessible | 1 |
| accident | 1 |

| 1-10 of 2,157 rows | Previous **1** 2 3 4 5 6 … 216 Next |

We see that the list contains words, that are only used once, this could indicate that some mispelling have happend or its not a very usefull word. Therefore we remove words used less than two times.

Hide

```
reviews_tidy = reviews_tidy %>%
  add_count(ID, word, name = "nword") %>%
  filter(nword > 1) %>%
  select(-nword)
```

Have a look again

Hide

```
reviews_tidy %>%
  count(word, sort = TRUE) %>% # Count "word". "sort = TRUE" means that it w
ill sort the words in descending order of number words.
  top_n(-20)
```

| word <chr> | n <int> |
| --- | --- |
| accessed | 2 |
| accessory | 2 |
| accidentally | 2 |
| accidently | 2 |
| accounts | 2 |
| actions | 2 |
| activate | 2 |
| activated | 2 |
| activities | 2 |
| acts | 2 |
| 1-10 of 872 rows | Previous **1** 2 3 4 5 6 … 88 Next |

We see some different words now, but we still find these words usefull and very informative.

We now make varible only containing the words in descending order of how much they are used.

Hide

```
topwords <- reviews_tidy %>%
  count(word, sort = TRUE)

head(topwords)
```

| word <chr> | n <int> |
| --- | --- |
| tablet | 748 |
| love | 477 |
| kindle | 462 |
| echo | 385 |
| easy | 337 |
| amazon | 320 |
| 6 rows | |

We use this new *topwords* variable to plot the topwords.

Hide

```
topwords %>%
  top_n(20, n) %>%
  ggplot(aes(x = word %>% fct_reorder(n), y = n)) + #fct_reorder factor the
 counts + it reorder the counts in a descending format.
  geom_col() +
  coord_flip() +
  labs(title = "Word Counts",
       x = "Frequency",
       y = "Top Words")
```



We see that the word *tablet* is way more than the rest. If we have a look at the other words, we see that overall the words are mainly the name of some electronics or the words are mainly about electronics.

*For later use, we save this reviews_tidy with a different name*

Hide

```
reviews_tidy_ML = reviews_tidy
```

# 2.2. Simple vectorization (Tf-idf)

In this section we add the *tf, term of frequency*, *idf, inverse document frequency* and *tf-idf, term frequency–inverse document frequency* to the tidy data.

Hide

```
reviews_tidy <- reviews_tidy %>%
  count(ID, word) %>%
  bind_tf_idf(ID, word, n) # I use this function to add the "tf", "idf" & "t
f_idf" values.
```

We have a look at the dwords with the highest *tf_idf* values

Hide

```
reviews_tidy %>%
  arrange(desc(tf_idf))
```

| ID | word | n | tf | idf | tf_idf |
|---:|:---|---:|---:|---:|---:|
| <int> | <chr> | <int> | <dbl> | <dbl> | <dbl> |
| 102 | äò | 2 | 1.000000000 | 7.510431 | 7.510430556 |
| 238 | äôre | 2 | 1.000000000 | 7.510431 | 7.510430556 |
| 1167 | response | 2 | 1.000000000 | 7.510431 | 7.510430556 |
| 1250 | bot | 2 | 1.000000000 | 7.510431 | 7.510430556 |
| 1334 | charges | 2 | 1.000000000 | 7.510431 | 7.510430556 |
| 1412 | mines | 2 | 1.000000000 | 7.510431 | 7.510430556 |
| 2150 | reliable | 2 | 1.000000000 | 7.510431 | 7.510430556 |
| 2534 | concept | 2 | 1.000000000 | 7.510431 | 7.510430556 |
| 2881 | launcher | 2 | 1.000000000 | 7.510431 | 7.510430556 |
| 3094 | lg | 2 | 1.000000000 | 7.510431 | 7.510430556 |

1-10 of 7,828 rows          Previous **1** 2   3   4   5   6 … 783 Next

# 2.3. Topic modelling / Clustering (LSA)

Now we will perform a LSA, which is stable when attempting to do dimensionality reduction as preprocessing for supervised ML workflows, or for visualization.

We have loaded the quanteda package, which is for corpus-token based text analysis.

First we have to create a document-feature-matrix

Hide

```
reviews_dfm <- reviews_tidy %>%
  count(ID, word) %>%
  cast_dfm(document = ID, term = word, value = n)

reviews_dfm
```

```
## Document-feature matrix of: 1,510 documents, 1,827 features (99.7% spars
e).
```

Now we get ready for the *LSA*, by choosing the number of dimensions.

Hide

```
reviews_lsa <- reviews_dfm %>%
  textmodel_lsa(nd = 5)
```

Let's have a look

Hide

```
reviews_lsa %>%
  glimpse()
```

```
## List of 5
##  $ sk            : num [1:5] 22.7 16.5 15 13.5 13
##  $ docs          : num [1:1510, 1:5] 0.0119 0.0191 0.0109 0.0199 0.015
...
##    ..- attr(*, "dimnames")=List of 2
##    .. ..$ : chr [1:1510] "3" "4" "6" "7" ...
##    .. ..$ : NULL
##  $ features      : num [1:1827, 1:5] 0.012388 0.004371 0.016513 0.000999
0.211243 ...
##    ..- attr(*, "dimnames")=List of 2
##    .. ..$ : chr [1:1827] "dark" "didnt" "happy" "im" ...
##    .. ..$ : NULL
##  $ matrix_low_rank: num [1:1510, 1:1827] 0.0209 0.0266 0.0103 0.0177 0.02
58 ...
##    ..- attr(*, "dimnames")=List of 2
##    .. ..$ : chr [1:1510] "3" "4" "6" "7" ...
##    .. ..$ : chr [1:1827] "dark" "didnt" "happy" "im" ...
##  $ data          :Formal class 'dgCMatrix' [package "Matrix"] with 6 slo
ts
##    .. ..@ i       : int [1:7828] 0 41 48 49 50 541 1395 0 959 1025 ...
##    .. ..@ p       : int [1:1828] 0 7 10 31 33 187 191 214 216 296 ...
##    .. ..@ Dim     : int [1:2] 1510 1827
##    .. ..@ Dimnames:List of 2
##    .. ..@ x       : num [1:7828] 1 1 1 1 1 1 1 1 1 1 ...
##    .. ..@ factors : list()
##  - attr(*, "class")= chr "textmodel_lsa"
```

Now we take the *LSA documents* and put them into a tibble, and adding the id's as row names.

Hide

```
reviews_lsa_loading <- reviews_lsa$docs %>%
  as.data.frame() %>%
  rownames_to_column(var = "ID") %>%
  as_tibble()

reviews_lsa_loading %>%
  head()
```

| ID<br><chr> | V1<br><dbl> | V2<br><dbl> | V3<br><dbl> | V4<br><dbl> | V5<br><dbl> |
|---|---|---|---|---|---|
| 3 | 0.01185456 | 0.0037554460 | 0.01558929 | 0.04261870 | 0.01507959 |
| 4 | 0.01913180 | 0.0023359416 | 0.01422975 | 0.04757917 | 0.02522116 |
| 6 | 0.01092350 | 0.0000638529 | 0.01024060 | 0.01088557 | 0.01348919 |
| 7 | 0.01985597 | 0.0063195407 | 0.01095280 | 0.02735569 | 0.01406400 |
| 8 | 0.01501553 | 0.0028959729 | 0.01667973 | 0.04820387 | 0.02481269 |
| 11 | 0.04093220 | 0.0084140442 | 0.02277632 | 0.03733989 | 0.05381831 |

6 rows

We can nicely visualize it using *UMAP* dimensionality reduction for optimizing the visualization of the feature space.

Hide

```
reviews_lsa_umap = umap(reviews_lsa_loading %>%
                        column_to_rownames("ID"),
                    n_neighbors = 15,
                    metric = "cosine",
                    min_dist = 0.01,
                    scale = TRUE,
                    verbose = TRUE,
                    n_threads = 8)
```
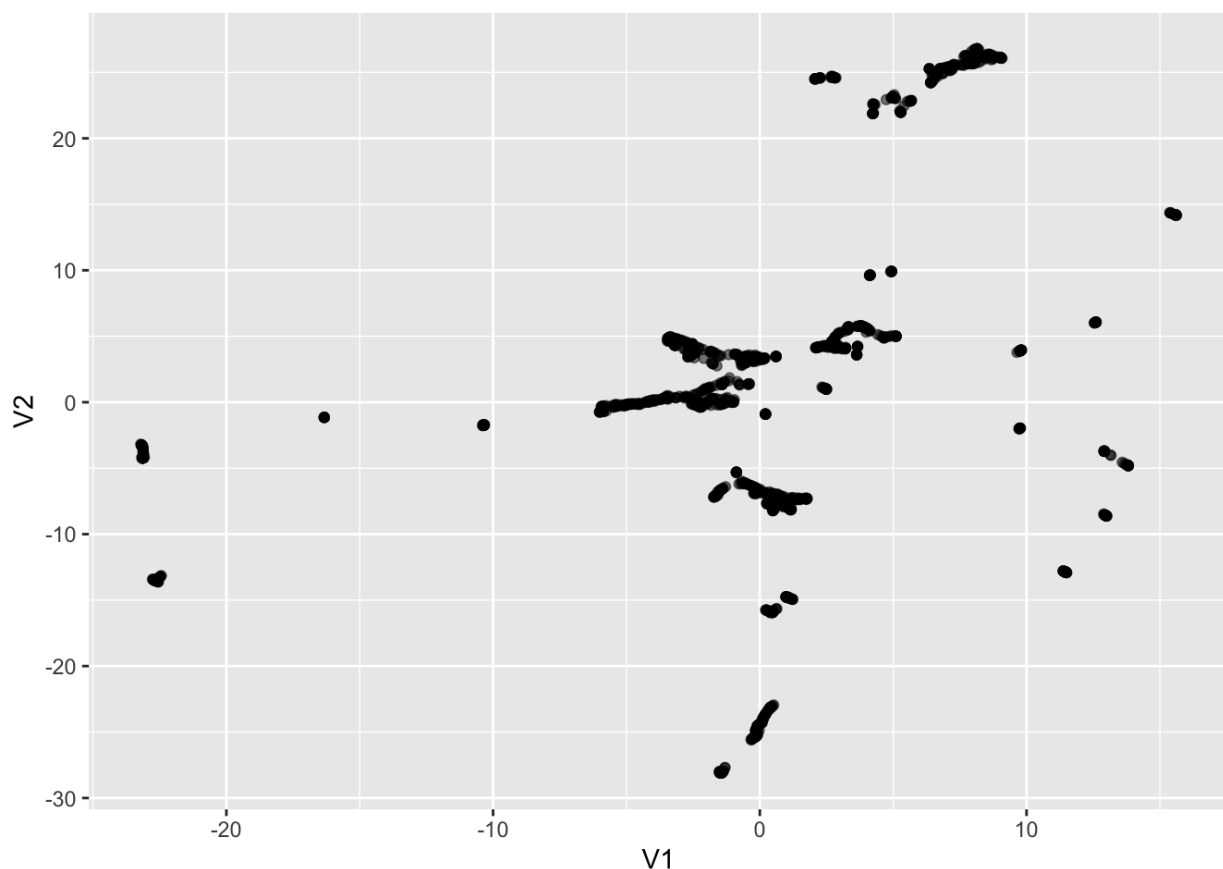
Then we put the result into a tibble for a nice layout.

Hide

```
reviews_lsa_umap  = reviews_lsa_umap %>%
    as.data.frame()
```

Now we can plot the result

Hide

```
reviews_lsa_umap %>%
    ggplot(aes(x = V1, y = V2)) +
    geom_point(alpha = 0.5)
```



We see that the some of the points tend to cluster togehter.

We now try to find clusters by using the *hdbscan*.
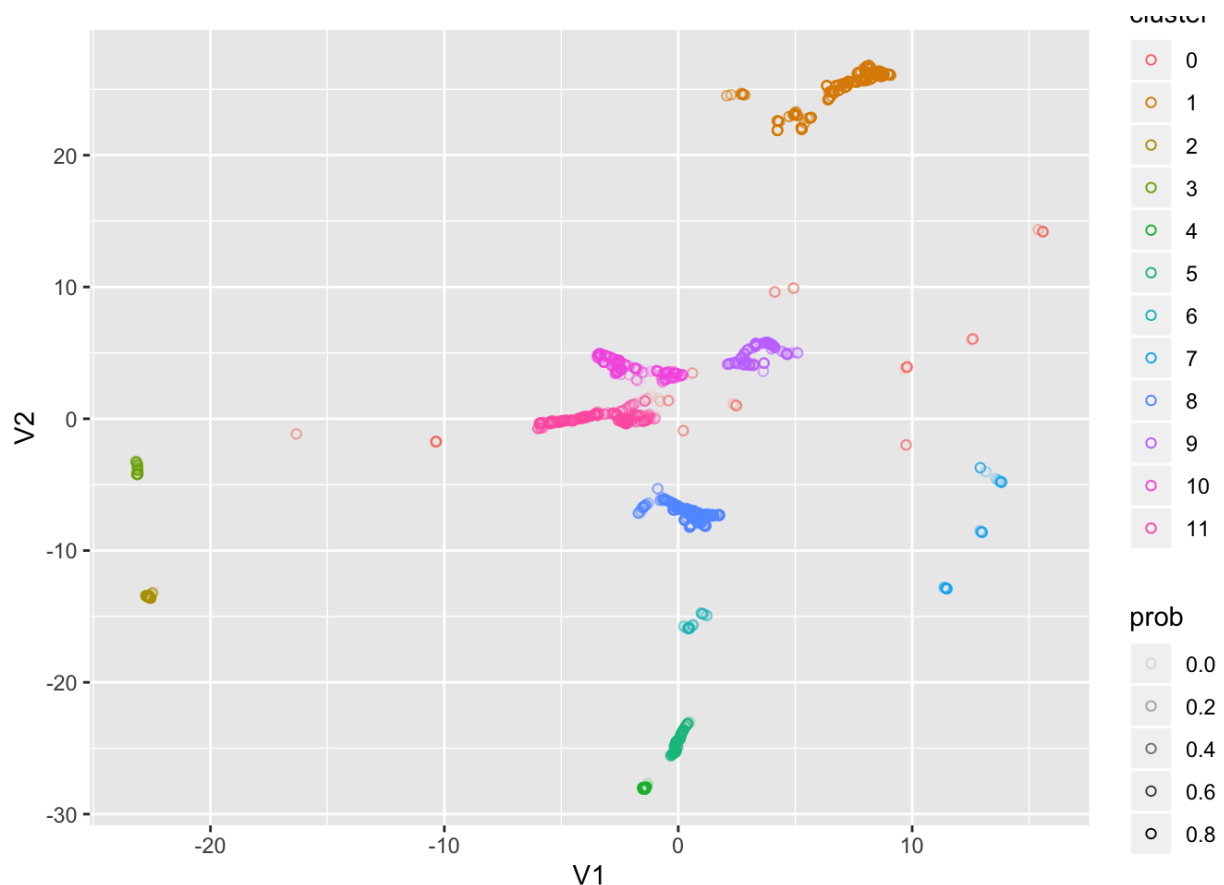
Hide

```
reviews_lsa_hdbscan <- reviews_lsa_umap %>%
  as.matrix() %>%
  hdbscan(minPts = 50)
```

Let plot again, but now we color by clusters.

Hide

```
reviews_lsa_umap %>%
  bind_cols(cluster = reviews_lsa_hdbscan$cluster %>% as.factor(),
            prob = reviews_lsa_hdbscan$membership_prob) %>%
  ggplot(aes(x = V1, y = V2, col = cluster)) +
  geom_point(aes(alpha = prob), shape = 21)
```



# 2.4. EDA / simple frequency-based analysis

Now we look closer into rating, we want to see of some words tend to show more in one rating then the others.

Hide

```
reviews_tidy3 = reviews_tidy %>%
  left_join(reviews, by = "ID")
```

Now we plot the 20 most used words in each rating.

Hide

```
reviews_tidy3 %>%
 count(rating, word, sort = TRUE) %>%
 group_by(rating) %>%
 top_n(10) %>%
 ungroup() %>%
 ggplot(aes(reorder_within(word, n, rating), n, fill = rating )) +
 geom_col(alpha = 0.8, show.legend = FALSE) +
 scale_x_reordered() +
 coord_flip() +
 facet_wrap(~rating, scales = "free") +
 scale_y_continuous(expand = c(0, 0)) +
 labs(x = NULL, y = "Word count", title = "20 Most frequent words in the 5 r
atings")
```

## 20 Most frequent words in the 5 ratings



we see that

- **Rating 1:** contains words like *constantly*, *returned*, *purchased* which could have a negative meaning and then there are a lot of words about the *electronics*.

- **Rating 2:** contains words like *responsive*, *time*, *fix* which also could have a negative meaning, and here we also see words describing the *electronics*

- **Rating 3:** contains words like the two other, but here the amount each words are used are almost the same for them all, except *tablet*. Therefore this looks loke the ones that don*t take this review seriously.

- **Rating 4:** containd words like *love*, *easy*, *price* and *product* this would usually have positive meaning in a text together with the describtion of the *electronic*.

- **Rating 5:** contain words like *love(s)*, *easy*, *price*, wich are way more positive words than the first two ratings. We also see that the amount each words are used are way higher here.

We now want to see why the amount each words are used differ that much through the ratings.

```
reviews %>%
  group_by(rating) %>%
  count()
```

| rating <fctr> | n <int> |
|---|---|
| 1 | 63 |
| 2 | 54 |
| 3 | 197 |
| 4 | 1208 |
| 5 | 3478 |
| 5 rows | |

We see that most of the reviewers gave the product a rating at $4$ or $5$, this could explain the differens in the amount of the word.

# 2.5. Topic modelling / Clustering (LDA)

**Preparing the Data**

For this application, we have to leave the tidy data, since the *topicmodels* package requires a document-term matrix as imput. We can easily produce it using the cast_dtm() function of tidytext. This matrix has to be term-frequency weighted, we do so using the *weightTf* function of the tm package for the weighting argument

```
reviews_dtm = reviews_tidy %>%
  count(ID, word) %>% #add a count of words in each id
  cast_dtm(document = ID, #Column containing document IDs as string or symbo
l
          term = word, #Column containing terms as string or symbol
          value = n, #Column containing values as string or symbol
          weighting = tm::weightTf) #The weighting function for the DTM/TDM

reviews_dtm
```

```
## <<DocumentTermMatrix (documents: 1510, terms: 1827)>>
## Non-/sparse entries: 7828/2750942
## Sparsity           : 100%
## Maximal term length: 18
## Weighting          : term frequency (tf)
```

We see that the matrix is $100\%$ sparse with $1510$ documents and $1827$ terms, which is an artefact of text data generally. Now we try to see if we could reduce that somewhat by deleting less often used terms.

```
reviews_dtm %>% removeSparseTerms(sparse = .99) #sparse is maximal allowed s
parsity.
```

```
## <<DocumentTermMatrix (documents: 1510, terms: 74)>>
## Non-/sparse entries: 3597/108143
## Sparsity            : 97%
## Maximal term length: 9
## Weighting           : term frequency (tf)
```

With max allowed sparsity at $99\%$ we only have $74$ terms out of $1827$, too little

Hide

```
reviews_dtm %>% removeSparseTerms(sparse = .999) #sparse is maximal allowed
 sparsity.
```

```
## <<DocumentTermMatrix (documents: 1510, terms: 806)>>
## Non-/sparse entries: 6807/1210253
## Sparsity            : 99%
## Maximal term length: 13
## Weighting           : term frequency (tf)
```

With max allowed sparsity at $99.9\%$ we only have $806$ terms out of $1827$, we might have to accept a high level of sparsity in order to still have a meaningful number of unique words.

Now we can perform an LDA, using the more accurate Gibbs sampling as method.

Hide

```
reviews_lda <- reviews_dtm %>%
  LDA(k = 5,
      method = "Gibbs",
      control = list(seed = 9))
```

## $\beta$ Word-Topic Association

$\beta$ is an output of the LDA model, indicating the propability that a word occurs in a certain topic. Therefore, loking at the top probability words of a topic often gives a good intuition regarding its properties.

Hide

```
# LDA output is defined for tidy(), so we can easily extract it
lda_beta <- reviews_lda %>%
  tidy(matrix = "beta") %>%
  group_by(topic) %>%
  arrange(topic, desc(beta)) %>% #the higher beta, the more a word is "used"
in a topic
  slice(1:15) %>%
  ungroup()

lda_beta %>% head()
```

| topic | term | beta |
|---|---|---|
| <int> | <chr> | <dbl> |

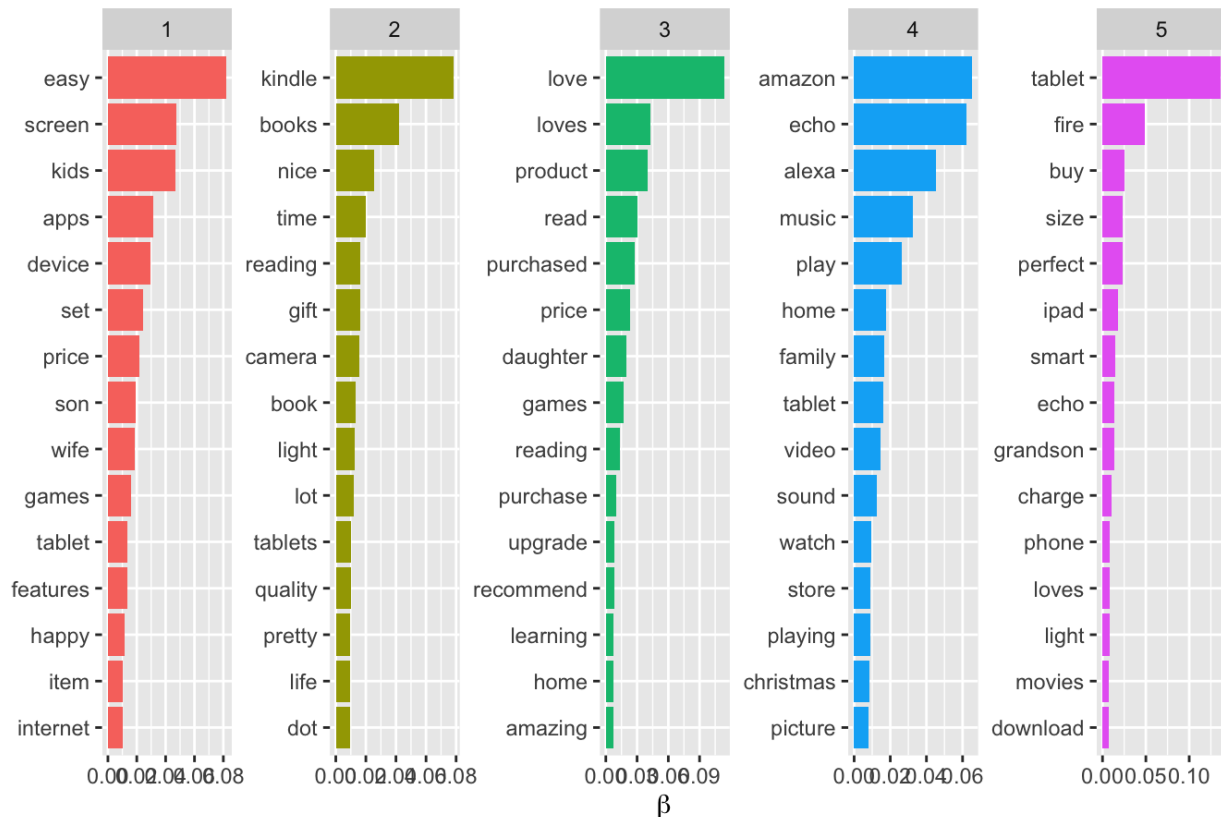| topic<br><int> | term<br><chr> | beta<br><dbl> |
|---:|:---|---:|
| 1 | easy | 0.08183222 |
| 1 | screen | 0.04752102 |
| 1 | kids | 0.04694916 |
| 1 | apps | 0.03150912 |
| 1 | device | 0.02922171 |
| 1 | set | 0.02407503 |

6 rows

We see that topic one have words like *easy*, *screen* and *kids* with high $\beta$ values.

Now we plot the top $10$ words in each topic.

Hide

```
lda_beta %>%
  mutate(term = reorder_within(term, beta, topic)) %>% #reordering the term
 by beta in each topic
  group_by(topic, term) %>% #first group by topic and then by term
  arrange(desc(beta)) %>% #rearrange data in descending order by beta
  ungroup() %>%
  ggplot(aes(term, beta, fill = as.factor(topic))) + #aes for the plot
  geom_col(show.legend = FALSE) + #Don't show legend
  coord_flip() + #flip the graph, so the bars will be horizontal
  scale_x_reordered() +
  labs(title = "Top 10 terms in each LDA topic",
       x = NULL,
       y = expression(beta)) +
  facet_wrap(~ topic, nrow = 1, scales = "free")
```

## Top 10 terms in each LDA topic



- **Topic 1:** We see that topic one is about quality for *HIFI* stuff.

- **Topic 2:** We see that topic two is mainly about reading. The word *Kindle* are top 1, but words like *books*, *read(ing)*, *light*, are also in here.

- **Topic 3:** We see that topic three is re mess of different types of words.

- **Topic 4:** We see that topic four is about *HIFI*, since words like *echo*, *alexa* and *music* are highly rated here.

- **Topic 5:** We see that topic five is mainly about tablets, since the words like *tablet*, *fire* and *ipad* are in this topic. together with some other tablet related words.

Overall we can say that the topics are mainly split into categories and not rating like we did earlier.

# 2.6. Embedding-model based vectorization (GloVe)

First we create a corpus

Hide

```
reviews_corpus <- reviews %>%
  corpus(docid_field = "ID",
         text_field = "text")
```

## Select features

Now we tokenize the corpus, we remove stopwords, punctuations, number and so on.

Hide

```
reviews_toks <- tokens(reviews_corpus, what = "word",
                       remove_numbers = TRUE, #Remove numbers
                       remove_punct = TRUE, #remove dots
                       remove_symbols = TRUE, #remove symbols
                       remove_separators = TRUE, #remove seperators
                       remove_url = TRUE, #remove url's
                       ngrams = 1) %>% #allowing unigrams
  tokens_tolower() %>% #only lower case letters
  tokens_remove(pattern = stopwords()) %>% #remove stopwords
  tokens_remove(pattern = own_stopwords$word) %>% #remove own words
  tokens_remove(pattern = c("[^[:alnum:]]")) #remove all non alpha numeric



reviews_toks %>% head(2)
```

```
## tokens from 2 documents.
## 1 :
##  [1] "thought"     "big"         "small"       "paper"       "turn"
##  [6] "just"        "like"        "palm"        "think"       "small"
## [11] "read"        "comfortable" "regular"     "kindle"      "definitely"
## [16] "recommend"   "paperwhite"  "instead"
##
## 2 :
## [1] "kindle"      "light"       "easy"        "use"         "especially"
## [6] "beach"
```

and then we keep the names of the features that occur five times or more.

Hide

```
feats <- dfm(reviews_toks, verbose = TRUE) %>%
    dfm_trim(min_termfreq = 5) %>%
    featnames()
```

## Constructing the feature co-occurence matrix

Hide

```
reviews_fcm <- fcm(reviews_toks,
             context = "window",
             count = "weighted",
             weights = 1 / (1:5),
             tri = TRUE)
```

## Fit word embedding model

Now we fit the *GloVe* model using *text2vec*

*GloVe* is an unsupervised learning algorithm for obtaining vector representations for words. Training is performed on aggregated global word-word co-occurrence statistics from a corpus, and the resulting representations showcase interesting linear substructures of the word vector space.

Hide

```
glove <- GlobalVectors$new(word_vectors_size = 50, #number of words in each
  vector
                            vocabulary = featnames(reviews_fcm),
                            x_max = 10)


reviews_main <- fit_transform(reviews_fcm,
                            glove,
                            n_iter = 20) #number of iterations
```

```
## INFO [2019-10-27 09:32:00] 2019-10-27 09:32:00 - epoch 1, expected cost
0.0806
## INFO [2019-10-27 09:32:01] 2019-10-27 09:32:01 - epoch 2, expected cost
0.0479
## INFO [2019-10-27 09:32:01] 2019-10-27 09:32:01 - epoch 3, expected cost
0.0389
## INFO [2019-10-27 09:32:01] 2019-10-27 09:32:01 - epoch 4, expected cost
0.0336
## INFO [2019-10-27 09:32:01] 2019-10-27 09:32:01 - epoch 5, expected cost
0.0300
## INFO [2019-10-27 09:32:01] 2019-10-27 09:32:01 - epoch 6, expected cost
0.0273
## INFO [2019-10-27 09:32:01] 2019-10-27 09:32:01 - epoch 7, expected cost
0.0251
## INFO [2019-10-27 09:32:01] 2019-10-27 09:32:01 - epoch 8, expected cost
0.0234
## INFO [2019-10-27 09:32:01] 2019-10-27 09:32:01 - epoch 9, expected cost
0.0220
## INFO [2019-10-27 09:32:01] 2019-10-27 09:32:01 - epoch 10, expected cost
0.0207
## INFO [2019-10-27 09:32:01] 2019-10-27 09:32:01 - epoch 11, expected cost
0.0197
## INFO [2019-10-27 09:32:01] 2019-10-27 09:32:01 - epoch 12, expected cost
0.0188
## INFO [2019-10-27 09:32:01] 2019-10-27 09:32:01 - epoch 13, expected cost
0.0180
## INFO [2019-10-27 09:32:01] 2019-10-27 09:32:01 - epoch 14, expected cost
0.0173
## INFO [2019-10-27 09:32:01] 2019-10-27 09:32:01 - epoch 15, expected cost
0.0167
## INFO [2019-10-27 09:32:02] 2019-10-27 09:32:01 - epoch 16, expected cost
0.0161
## INFO [2019-10-27 09:32:02] 2019-10-27 09:32:02 - epoch 17, expected cost
0.0156
## INFO [2019-10-27 09:32:02] 2019-10-27 09:32:02 - epoch 18, expected cost
0.0152
## INFO [2019-10-27 09:32:02] 2019-10-27 09:32:02 - epoch 19, expected cost
0.0148
## INFO [2019-10-27 09:32:02] 2019-10-27 09:32:02 - epoch 20, expected cost
0.0144
```

## Averagning learned word vectors

The two vectors are main and context. According to the *Glove* paper, averaging the two word vectors results in more accurate representation.

Hide

```
reviews_context <- glove$components
dim(reviews_context)
```

```
## [1]   50 5613
```

We see that the dimension of this model is $50 \times 5614$.

Hide

```
reviews_vectors <- as.dfm(reviews_main + t(reviews_context))
```

## Examining term representations

Now we can find the six closest word vectors for *kindle*.

Hide

```
new <- reviews_vectors["kindle", ]

# calculate the similarity
cos_sim <- textstat_simil(reviews_vectors, new,
                          margin = "documents", method = "cosine")

head(sort(cos_sim[, 1], decreasing = TRUE), 6)
```

```
##    kindle      fire    second       new      much     first
## 1.0000000 0.8208233 0.6757720 0.6563131 0.6323784 0.6260871
```

We see the words *fire*, *new*, *since*, *one* and *second* are the closest words to the word *kindle*

The six closest word vectors for *tablet*.

Hide

```
new <- reviews_vectors["tablet", ]

# calculate the similarity
cos_sim <- textstat_simil(reviews_vectors, new,
                          margin = "documents", method = "cosine")

head(sort(cos_sim[, 1], decreasing = TRUE), 6)
```

```
##    tablet     great      good     price      fire   perfect
## 1.0000000 0.7494200 0.7299844 0.7288418 0.6668515 0.6633848
```

We see that the words like *good*, *great*, *price* and *kids* are the closest to the word *tablet*.

As seen earlier the data caontains mostly of high rated reviews. Therefore it makes sence that the word *tablet* follows by positive words like *great* and *good*.

Now we make an average-vector-representations for the words

Hide

```
reviews_main_tibble = reviews_main %>%
  as.data.frame() %>% #making a dataframe
  rownames_to_column(var = "word") %>% #Row names
  as_tibble() #for nice layout

reviews_main_tibble %>% head()
```

| word | V1 | V2 | V3 | V4 | V5 | |
|------|-----|-----|-----|-----|-----|---|
| <chr> | <dbl> | <dbl> | <dbl> | <dbl> | <dbl> | |
| thought | -0.2731706 | -0.03957523 | 0.01527472 | 0.12807940 | -0.08243837 | -0.1331 |
| big | 0.1528246 | -0.25964165 | 0.21220364 | 0.51244879 | -0.30240378 | -0.2765 |
| small | -0.3029513 | 0.25653461 | 0.27191514 | 0.70423305 | 0.04605259 | 0.1813 |
| paper | 0.5156913 | 0.22531784 | -0.23414992 | 0.18228483 | 0.01700168 | -0.4546 |
| turn | -0.5437862 | -0.09267105 | 0.82648665 | -0.01936661 | 0.22087295 | 0.3563 |
| just | -0.4766783 | 0.01564204 | -0.27068293 | 0.25873142 | 0.14141901 | 0.0895 |

6 rows | 1-7 of 51 columns

Now we make a average-vector-representations for the reviews.

Hide

```
reviews_tidy2 = reviews_toks %>%
  dfm() %>%
  tidy()

reviews_vector2 = reviews_tidy2 %>%
  inner_join(reviews_main_tibble, by = c("term" = "word")) #joining the data
by word/term

head(reviews_vector2)
```

| docum... | term | co... | V1 | V2 | V3 | V4 | |
|----------|------|-------|-----|-----|-----|-----|---|
| <chr> | <chr> | <dbl> | <dbl> | <dbl> | <dbl> | <dbl> | < |
| 1 | thought | 1 | -0.2731706 | -0.03957523 | 0.01527472 | 0.1280794 | -0.08243 |
| 65 | thought | 1 | -0.2731706 | -0.03957523 | 0.01527472 | 0.1280794 | -0.08243 |
| 70 | thought | 1 | -0.2731706 | -0.03957523 | 0.01527472 | 0.1280794 | -0.08243 |
| 156 | thought | 1 | -0.2731706 | -0.03957523 | 0.01527472 | 0.1280794 | -0.08243 |
| 158 | thought | 1 | -0.2731706 | -0.03957523 | 0.01527472 | 0.1280794 | -0.08243 |
| 182 | thought | 1 | -0.2731706 | -0.03957523 | 0.01527472 | 0.1280794 | -0.08243 |

6 rows | 1-8 of 53 columns

Hide

```
reviews_vector2 = reviews_vector2 %>%
  select(-term, -count) %>% #remove term and count
  group_by(document) %>% #group the data by document
  summarise_all(mean) %>% #finding the mean of the coloumns for each documen
t
  ungroup()
```
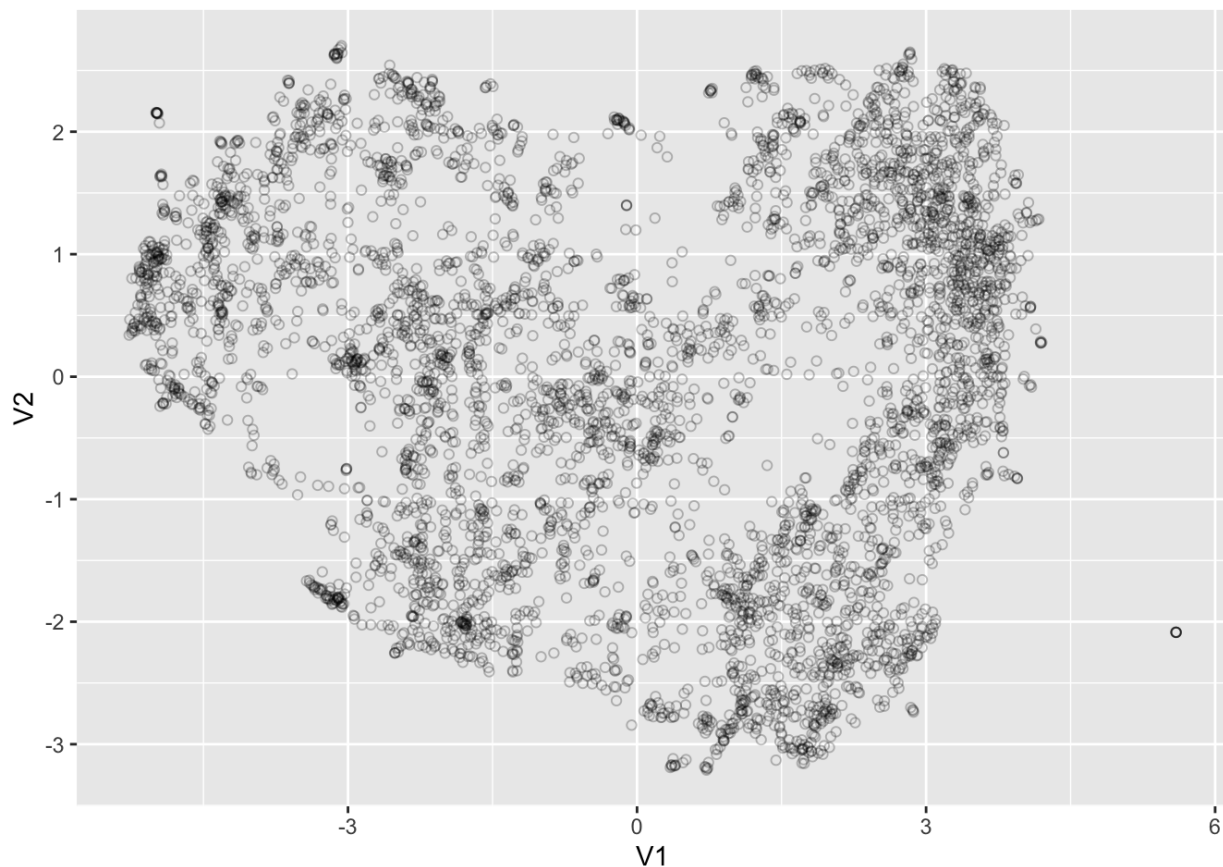
Let's have a look at the data

Hide

```
reviews_vector2_umap <- umap(reviews_vector2 %>% column_to_rownames("documen
t"),
                        n_neighbors = 15,
                        metric = "cosine",
                        min_dist = 0.01,
                        scale = TRUE,
                        verbose = TRUE) %>%
  as.data.frame()
```

Hide

```
reviews_vector2_umap %>%
  ggplot(aes(x = V1, y = V2)) +
  geom_point(shape = 21, alpha = 0.25)
```



It does not look like the data have any clusters right now. If we look closely we could argue about some few clusters at the plot, which might be *tablet* and/or *kindle* stuff.

# 3. Supervised / Unsupervised ML

# 3.1. Unsupervised ML

We do *kmeans* for the *tf_idf* values, and want $5$ clusters since we have $5$ different ratings in the data.

Hide

```
kmeans_tfidf = kmeans(reviews_tidy[,4:6], centers = 5, nstart = 20)
```
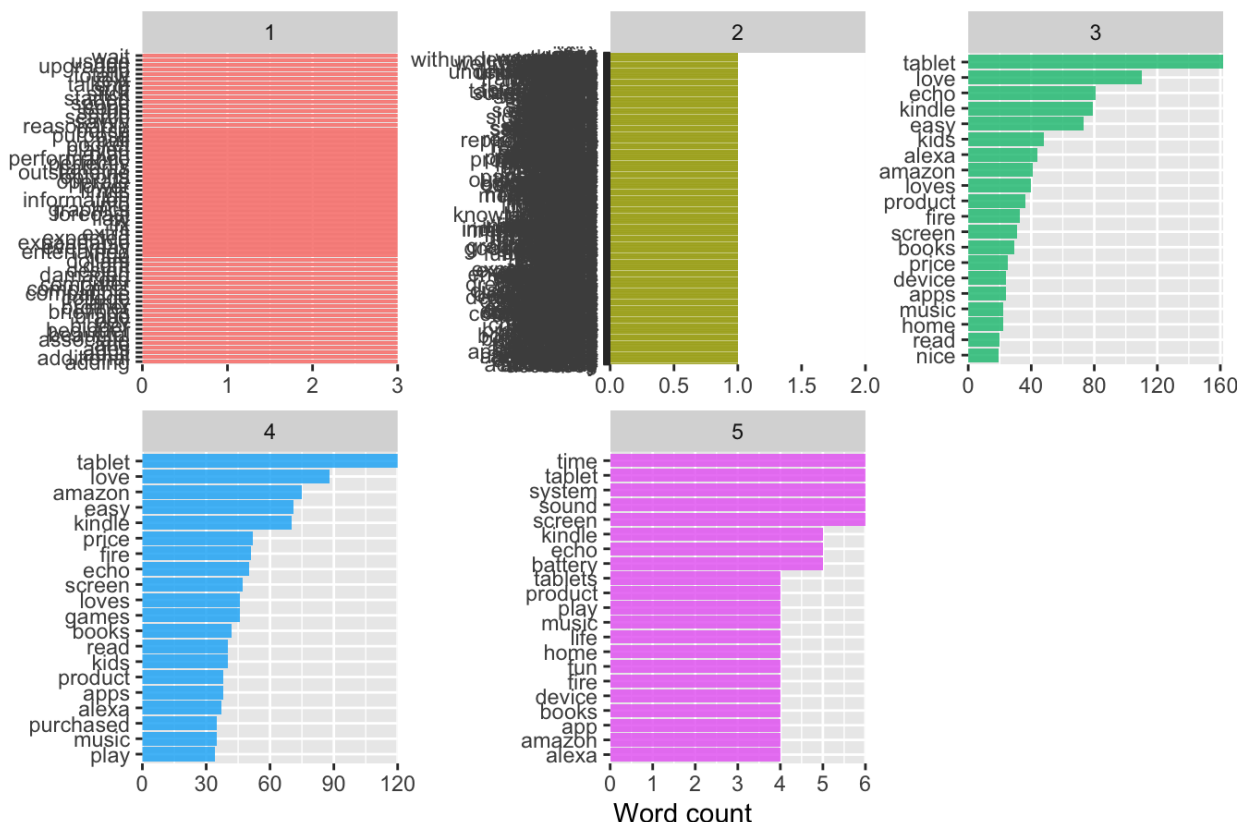
Hide

```
reviews_tidy = reviews_tidy %>%
  bind_cols(cluster = kmeans_tfidf$cluster)
```

Now we plot the reviews_tidy, and color by clusters.

Hide

```
reviews_tidy %>%
 mutate(cluster = as.factor(cluster)) %>%
 count(cluster, word, sort = TRUE) %>%
 group_by(cluster) %>%
 top_n(20,n) %>%
 ungroup() %>%
 ggplot(aes(reorder_within(word, n, cluster), n,
   fill = cluster
)) +
 geom_col(alpha = 0.8, show.legend = FALSE) +
 scale_x_reordered() +
 coord_flip() +
 facet_wrap(~cluster, scales = "free")  +
 scale_y_continuous(expand = c(0, 0)) +
 labs(x = NULL, y = "Word count",
   title = "Most frequent words in the 5 clusters")
```

## Most frequent words in the 5 clusters



Unfortunately the plots are messy. since some of the clusters contain words represented the same amount of times. Therefore we can only talk/discuss plot 1, 2 and 5.

The plots contain almost the same words, it is hard to find differences in these three plots.

# 3.2. Supervised ML

## 3.2.1. Method 1

Now we want to get ready for modeling. We split the data into training and testing sets, using *rsample*. We use the original data, reviews, because the reviews$text are our individual observations.

Hide

```r
reviews_split <- reviews %>%
  select(ID2) %>%
  initial_split() #by default it splits 3/4

train_data <- training(reviews_split) #train data

test_data <- testing(reviews_split) #test data
```

Now we transform our training data from a tidy structure to a sparse matrix, wich we will use for the machine learning algorithm.

Hide

```
sparse_words =  reviews_tidy_ML %>%
  count(ID2, word) %>%
  inner_join(train_data) %>% #keep rows from Twwets_tidy where there are mat
ching values in y, and all columns from x and y.
  cast_sparse(ID2, word, n) #Create a sparse matrix from row names, column n
ames, and values in a table
```

Hide

```
dim(sparse_words)
```

```
## [1] 1881 1754
```

We see we get a sparse matrix with $1885$ rows and $1764$ columns. This means we have $1885$ training observations and $1764$ features at this point.

Now we build a dataframe with response variable to associate each of the *rownames()* of the sparse matrix. We select the *rating*, since it will be our predictor later.

Hide

```
word_rownames = as.integer(rownames(sparse_words))
reviews_rating = reviews %>% select(ID2, rating)
```

Hide

```
reviews_joined = data_frame(ID2 = word_rownames) %>%
  left_join(reviews_rating, by = "ID2" )
```

Now we want to train our classification model. We use *glmnet* to fit a logistic regression model with *LASSO* regularization. It's a great fit for text classification because the variable selection that *LASSO* regularization performs can tell you which words are important for our prediction problem.

We make five models, since we have five types of ratings.

Hide

```r
#Model for rate 1
model_1 = cv.glmnet(sparse_words, # x variable = matrix .
                    reviews_joined$rating == 1, #response variable
                    family = "binomial",
                    keep = TRUE)

#Model for rate 2
model_2 <- cv.glmnet(sparse_words, #x variable = matrix .
                    reviews_joined$rating == 2,  #response variable
                    family = "binomial",
                    keep = TRUE)

#Model for rate 3
model_3 <- cv.glmnet(sparse_words, #x variable = matrix .
                    reviews_joined$rating == 3, #response variable
                    family = "binomial",
                    keep = TRUE)

#Model for rate 4
model_4 <- cv.glmnet(sparse_words, #x variable = matrix .
                    reviews_joined$rating == 4, #response variable
                    family = "binomial",
                    keep = TRUE)

#Model for rate 5
model_5 <- cv.glmnet(sparse_words, #x variable = matrix .
                    reviews_joined$rating == 5, #response variable
                    family = "binomial",
                    keep = TRUE)
```

Now we use *broom* to check the five models coeffecients. We also check for wich coefficients are largest in size, in each direction

Hide

```r
coef_1 <- model_1$glmnet.fit %>%
  tidy() %>%
  filter(lambda == model_1$lambda.1se) %>% #lambda with error within 1 stand
ard error
  group_by(estimate > 0) %>%
  top_n(10, abs(estimate)) %>%
  ungroup() %>%
  ggplot(aes(fct_reorder(term, estimate),  #plotting the top words in each d
irection
             estimate,
             fill = estimate > 0)) +
  geom_col(alpha = 0.8, show.legend = FALSE) +
  coord_flip() +
  scale_fill_brewer(palette = "Paired") +
  labs( title = "Rate 1",
      x = NULL)


coef_2 <- model_2$glmnet.fit %>%
  tidy() %>%
  filter(lambda == model_2$lambda.1se) %>% #lambda with error within 1 stand
ard error
  group_by(estimate > 0) %>%
  top_n(10, abs(estimate)) %>%
  ungroup() %>%
  ggplot( aes( fct_reorder( term, estimate),
               estimate,
               fill = estimate > 0)) +
  geom_col(alpha = 0.8, show.legend = FALSE) +
  coord_flip() +
  scale_fill_brewer(palette = "Paired") +
  labs( title = "Rate 2",
      x = NULL)


coef_3 <- model_3$glmnet.fit %>%
  tidy() %>%
  filter(lambda == model_3$lambda.1se) %>% #lambda with error within 1 stand
ard error
  group_by(estimate > 0) %>%
  top_n(10, abs(estimate)) %>%
  ungroup() %>%
  ggplot(aes(fct_reorder(term, estimate),
             estimate,
             fill = estimate > 0)) +
  geom_col(alpha = 0.8, show.legend = FALSE) +
  coord_flip() + #Horizontal barplot
  scale_fill_brewer(palette = "Paired") +
  labs(title = "Rate 3",
      x = NULL)


coef_4 <- model_4$glmnet.fit %>%
  tidy() %>%
  filter(lambda == model_4$lambda.1se) %>% #lambda with error within 1 stand
ard error
  group_by(estimate > 0) %>%
```
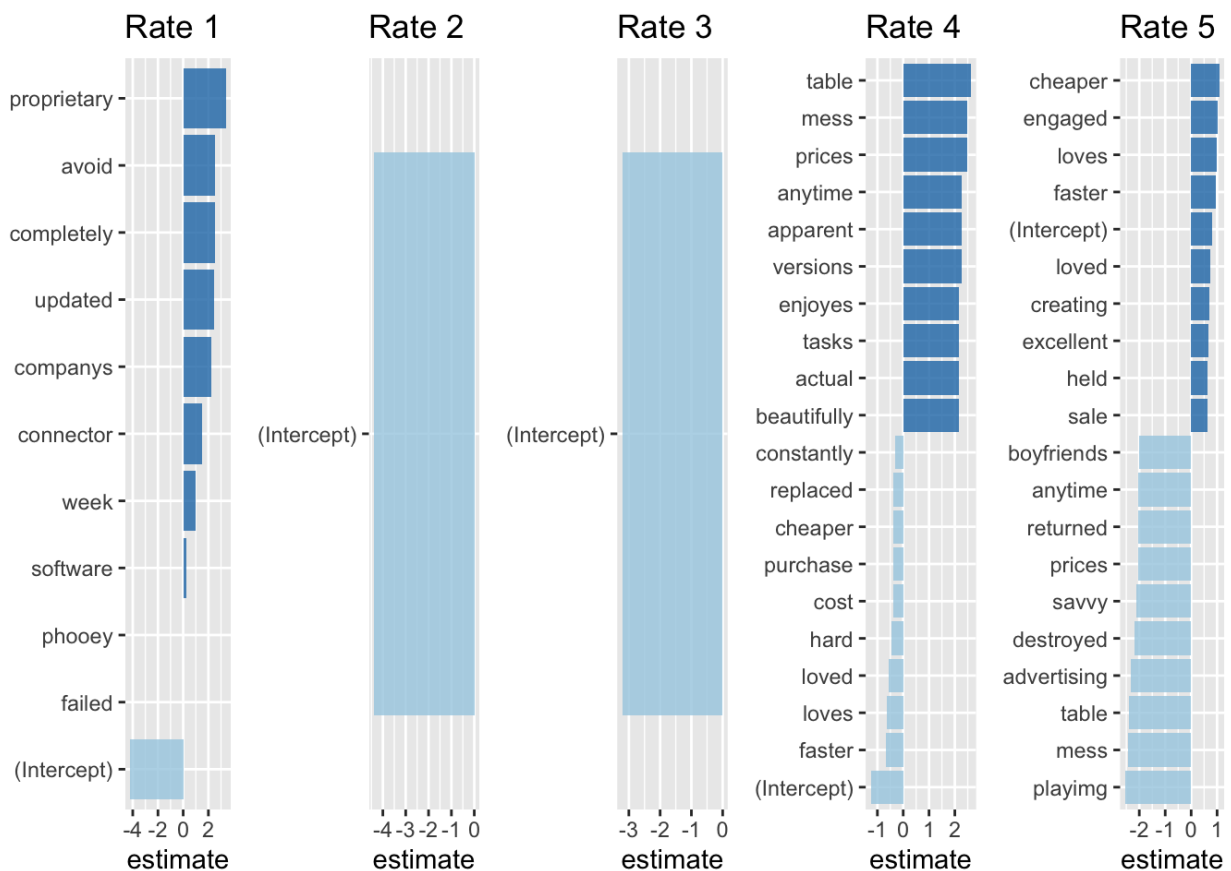
```
    top_n(10, abs(estimate)) %>%
    ungroup() %>%
    ggplot(aes(fct_reorder(term, estimate),
               estimate,
               fill = estimate > 0)) +
    geom_col(alpha = 0.8, show.legend = FALSE) +
    coord_flip() + #Horizontal barplot
    scale_fill_brewer(palette = "Paired") +
    labs(title = "Rate 4",
         x = NULL)


coef_5 <- model_5$glmnet.fit %>%
  tidy() %>%
  filter(lambda == model_5$lambda.1se) %>% #lambda with error within 1 stand
ard error
  group_by(estimate > 0) %>%
  top_n(10, abs(estimate)) %>%
  ungroup() %>%
  ggplot(aes(fct_reorder(term, estimate),
             estimate,
             fill = estimate > 0)) +
  geom_col(alpha = 0.8, show.legend = FALSE) +
  coord_flip() + #Horizontal barplot
  scale_fill_brewer(palette = "Paired") +
  labs(title = "Rate 5",
       x = NULL)

grid.arrange(coef_1, coef_2, coef_3, coef_4, coef_5, nrow = 1)
```



We see that not all the plots are very good, but we can see the meaning of some words for some topics.

Now we try to classify the rating for the test set.

Hide

```
model <- cv.glmnet(sparse_words,
                   reviews_joined$rating,
                   family = "multinomial",
                   keep = TRUE)

coefs <- model$glmnet.fit %>%
  tidy() %>%
  filter(lambda == model$lambda.1se)

intercept <- coefs %>%
  filter(term == "(Intercept)") %>%
  pull(estimate)

classifications <- reviews_tidy_ML %>%
  inner_join(test_data) %>% #Predicting the test data
  inner_join(coefs, by = c("word" = "term")) %>%
  group_by(class, ID2) %>%
  summarize(score = sum(estimate)) %>%
  mutate(probability = plogis(0 + score))
```

Let's have a look.

Hide

```
classifications
```

| class | ID2 | score | probability |
|:------|:---:|------:|------------:|
| <chr> | <int> | <dbl> | <dbl> |
| 1 | 227 | 1.320891929 | 0.78933006 |
| 1 | 305 | 2.540961206 | 0.92696393 |
| 1 | 516 | 1.333006916 | 0.79133758 |
| 1 | 545 | 0.666503458 | 0.66071978 |
| 1 | 664 | 4.404218312 | 0.98792200 |
| 1 | 665 | 1.270480603 | 0.78082501 |
| 1 | 878 | 1.333006916 | 0.79133758 |
| 1 | 879 | 0.666503458 | 0.66071978 |
| 1 | 1782 | 0.666503458 | 0.66071978 |
| 1 | 1839 | 3.472099305 | 0.96988340 |

1-10 of 604 rows                          Previous  **1**  2  3  4  5  6  …  61  Next

We need to add the original rating, before we can say anything.

Hide

```
res <- classifications %>%
  group_by(ID2) %>%
  filter(probability == max(probability)) %>%
  left_join(reviews, by = "ID2")
```

Now we put the result in a table, where we can see if the classification did good or not.

```
result <- table(res$rating, res$class)

result
```

```
##
##       1   2   3   4   5
##  1    5   0   0   0   1
##  2    0   1   1   0   2
##  3    0   0  11   2   7
##  4    2   0   1  53  53
##  5    4   5  11  58 236
```

We see that the model is not that good, since it predicts $\sim 50\%$ correct.

# 3.2.2. Method 2

Here we try do like we saw in M1.

We select our variables to the classification models.

```
reviews_tidy_ML <- reviews_tidy3 %>%
  select(tf, idf,
         tf_idf,
         helpful,
         rating,
         primaryCategories,
         year) %>%
  mutate(good_review = as.logical(as.integer(rating) >= 4, rating)) %>%
  select(-rating) # I set class == 0, so class 0 is TRUE and the rest FALSE:
```

## Splitting the data into train- and test set.

```
index <- createDataPartition(y = reviews_tidy_ML$good_review,
                             p = 0.75,
                             list = FALSE) # 75% to 25% split

training <- reviews_tidy_ML[index,]
test <- reviews_tidy_ML[-index,]
```

## Preprocessing using the recipes package.

```
reci <- recipe(good_review ~ ., data = training) %>%
  step_dummy(all_nominal(), -all_outcomes()) %>% # create dummy variables fo
r the nominal variables
  step_center(all_numeric(), -all_outcomes()) %>% # Centers all numeric vari
ables to mean = 0
  step_scale(all_numeric(), -all_outcomes()) %>% # this scales the numeric v
ariables
  step_zv(all_predictors())  # Removed predictors with zero variance

reci %<>%
  prep(data = train)
```

## Setting the x and y values.

Hide

```
# Predictors
x_train <- bake(reci, new_data = training) %>% select(-good_review) # I remo
ve class from the predictors.
y_train <- training %>% pull(good_review) %>% as.factor()
# test: split in y and x
x_test <- bake(reci, new_data = test) %>% select(-good_review)
y_test <- test %>% pull(good_review) %>% as.factor()
```

## Setting the workflow.

Hide

```
ctrl <- trainControl(method = "cv",
                     number = 4)
metric <- "Kappa" # I use Kappa because the class 0 (hate speech) is so low
 represented in the data.
```

## Fitting the model - logistic regression.

Hide

```
review_fit_log <- train(x = x_train,
              y = y_train,
              trControl = ctrl,
              metric = metric,
              method = "glm",
              family = "binomial")
```

Hide

```
review_fit_log
```

```
## Generalized Linear Model
##
## 17694 samples
##    10 predictor
##     2 classes: 'FALSE', 'TRUE'
##
## No pre-processing
## Resampling: Cross-Validated (4 fold)
## Summary of sample sizes: 13270, 13271, 13270, 13271
## Resampling results:
##
##   Accuracy   Kappa
##   0.9206511  0
```

## Prediction and evaluating.

Hide

```
pred_log <- predict(review_fit_log, newdata = x_test)
```

Hide

```
print("Confusion matrix for logistic model")
```

```
## [1] "Confusion matrix for logistic model"
```

Hide

```
confusionMatrix(pred_log, y_test, positive = 'TRUE')
```

```
## Confusion Matrix and Statistics
##
##           Reference
## Prediction FALSE TRUE
##      FALSE     0    0
##      TRUE    468 5429
##
##                Accuracy : 0.9206
##                  95% CI : (0.9134, 0.9274)
##     No Information Rate : 0.9206
##     P-Value [Acc > NIR] : 0.5123
##
##                   Kappa : 0
##
##  Mcnemar's Test P-Value : <2e-16
##
##             Sensitivity : 1.0000
##             Specificity : 0.0000
##          Pos Pred Value : 0.9206
##          Neg Pred Value :    NaN
##              Prevalence : 0.9206
##          Detection Rate : 0.9206
##    Detection Prevalence : 1.0000
##       Balanced Accuracy : 0.5000
##
##        'Positive' Class : TRUE
##
```

## Fitting the model - decision tree

Hide

```
reviews_fit_dt <- train(x = x_train,
              y = y_train,
              trControl = ctrl,
              metric = metric,
              method = "rpart")
```

Hide

```
reviews_fit_dt
```

```
## CART
##
## 17694 samples
##    10 predictor
##     2 classes: 'FALSE', 'TRUE'
##
## No pre-processing
## Resampling: Cross-Validated (4 fold)
## Summary of sample sizes: 13270, 13271, 13271, 13270
## Resampling results across tuning parameters:
##
##   cp          Accuracy   Kappa
##   0.01495726  0.9434272  0.4251322
##   0.02207977  0.9299197  0.1796601
##   0.02250712  0.9299197  0.1796601
##
## Kappa was used to select the optimal model using the largest value.
## The final value used for the model was cp = 0.01495726.
```

## Predicting and evaluating

Hide

```
pred_dt <- predict(reviews_fit_dt, newdata = x_test)
```

Hide

```
print("Confusion matrix for decision tree")
```

```
## [1] "Confusion matrix for decision tree"
```

Hide

```
confusionMatrix(pred_dt, y_test, positive = 'TRUE')
```

```
## Confusion Matrix and Statistics
##
##           Reference
## Prediction FALSE TRUE
##      FALSE   134    0
##      TRUE    334 5429
##
##                  Accuracy : 0.9434
##                    95% CI : (0.9372, 0.9491)
##       No Information Rate : 0.9206
##       P-Value [Acc > NIR] : 7.679e-12
##
##                     Kappa : 0.4249
##
##  Mcnemar's Test P-Value : < 2.2e-16
##
##               Sensitivity : 1.0000
##               Specificity : 0.2863
##            Pos Pred Value : 0.9420
##            Neg Pred Value : 1.0000
##                Prevalence : 0.9206
##            Detection Rate : 0.9206
##      Detection Prevalence : 0.9773
##         Balanced Accuracy : 0.6432
##
##          'Positive' Class : TRUE
##
```
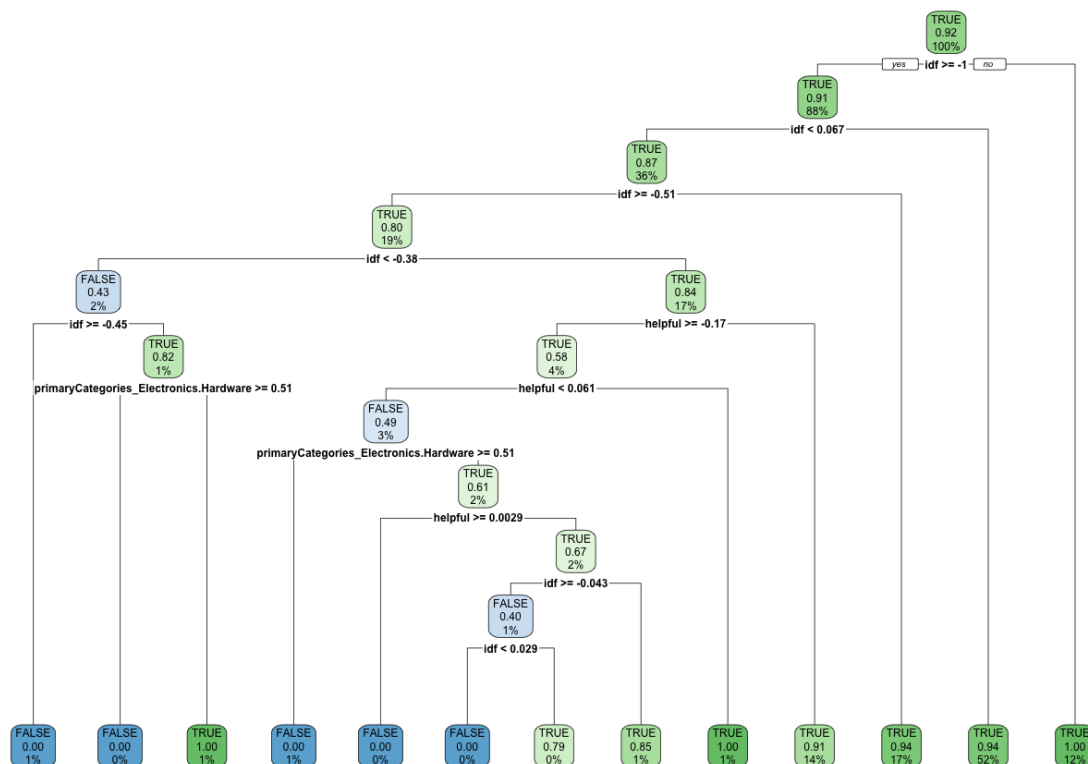
*We had some problems with the margins therefore we changed it.*

Hide

```
par(mar = c(1,1,1,1)) #The margins for the next plot are too large
```

Hide

```
reviews_fit_dt$finalModel %>%
  rpart.plot()
```

## Fitting the model - Random forest

Hide

```
reviews_fit_rf <- train(x = x_train,
                y = y_train,
                trControl = ctrl,
                metric = metric,
                method = "ranger",
                importance = "impurity", # To define how to measure variable
importantce (later more)
                num.trees = 25
                )
```

Hide

```
reviews_fit_rf
```

```
## Random Forest
##
## 17694 samples
##    10 predictor
##     2 classes: 'FALSE', 'TRUE'
##
## No pre-processing
## Resampling: Cross-Validated (4 fold)
## Summary of sample sizes: 13271, 13270, 13270, 13271
## Resampling results across tuning parameters:
##
##   mtry  splitrule   Accuracy   Kappa
##    2    gini        0.9255679  0.1072797
##    2    extratrees  0.9206511  0.0000000
##    6    gini        0.9540520  0.6316107
##    6    extratrees  0.9511698  0.5874154
##   10    gini        0.9499264  0.6216231
##   10    extratrees  0.9502091  0.6222944
##
## Tuning parameter 'min.node.size' was held constant at a value of 1
## Kappa was used to select the optimal model using the largest value.
## The final values used for the model were mtry = 6, splitrule = gini
##  and min.node.size = 1.
```

## Predicting and evaluating

Hide

```
pred_rf <- predict(reviews_fit_rf, newdata = x_test)
```

Hide

```
print("Confusion matrix for Random forest")
```

```
## [1] "Confusion matrix for Random forest"
```

Hide

```
confusionMatrix(pred_rf, y_test, positive = 'TRUE')
```

```
## Confusion Matrix and Statistics
##
##           Reference
## Prediction FALSE TRUE
##      FALSE   272   56
##       TRUE   196 5373
##
##                Accuracy : 0.9573
##                  95% CI : (0.9518, 0.9623)
##     No Information Rate : 0.9206
##     P-Value [Acc > NIR] : < 2.2e-16
##
##                   Kappa : 0.6613
##
##  Mcnemar's Test P-Value : < 2.2e-16
##
##             Sensitivity : 0.9897
##             Specificity : 0.5812
##          Pos Pred Value : 0.9648
##          Neg Pred Value : 0.8293
##              Prevalence : 0.9206
##          Detection Rate : 0.9111
##    Detection Prevalence : 0.9444
##       Balanced Accuracy : 0.7854
##
##        'Positive' Class : TRUE
##
```

# 4. Network Analysis

In this section we are going to construct a network based on the reviews data.

We did several constructions and tried interpret different attributes to the network, but for some reason we could not include any attributes. For example we wanted to try include *rating* and *primaryCategories* and use them for coloring or shaping the nodes

## Creating smaller sample of reviews.

We first tried to create a network inclusing the whole dataset, but because of the size, we where not be able to run the network/graph *(the system crashed)*.

Instead we resampled the review dataset to only 500 observations - so we instead got a datasat consiting of 500 reviews instead of 5.000.

Hide

```
reviews_sampled <- reviews[sample(nrow(reviews),500),]
```

Now we create a tidy data as we did earlier, just on the resampled set.

Hide

```
reviews_sampled_tidy <- reviews_sampled %>%
  select(ID, text) %>%
  unnest_tokens(output = word, input = text) %>%
  anti_join(stop_words %>% bind_rows(own_stopwords), by = "word") %>%
  mutate(word = trimws(gsub("[^\\s]*[0-9][^\\s]*", "", word, perl = T))) %>%
  filter(str_length(word) > 1) %>%# this filter out the words that are blan
k.
  mutate(word = word %>% str_remove_all("[^[:alnum:]]") ) %>% # alnum = Alph
anumeric characters.
  filter(str_length(word) > 1) # filter out words with 1 character.
```

To make a network not with too many words, we only use the 50 highest counted words in the tidy set.

Hide

```
sampled_50 <- reviews_sampled_tidy %>%
  count(word, sort = TRUE) %>% # we count the words
head(50) # we take the head of the 50 most frequent words.
```

Now we want the samlped tidy to only consist of the 50 most frequent words and here we use "left_join" to this.

First we get a big dataset, but this is because we have all the words with "NA" values, which is the words that where not a part of the 50 most frequent words. Therefore we use *drop_na()* to drop the rows with NA values.

Hide

```
reviews_sampled_tidy_joined <- reviews_sampled_tidy %>%
  left_join(sampled_50, by = "word") %>%
  drop_na()
```

Let's have a look.

Hide

```
glimpse(reviews_sampled_tidy_joined)
```

```
## Observations: 1,800
## Variables: 3
## $ ID   <int> 1825, 1825, 1825, 3594, 3594, 817, 817, 817, 817, 577, 577,…
## $ word <chr> "tablet", "tablet", "games", "recommend", "kindle", "alexa"…
## $ n    <int> 127, 127, 39, 22, 89, 40, 98, 19, 98, 21, 46, 25, 42, 16, 6…
```

We see that the data now consists of $1933$ observations and $3$ variables.

Now we can create the nodelist from the sampled tidy dataset.

Hide

```
nodes_reviews <- reviews_sampled_tidy_joined %>%
  group_by(ID) %>%
  select(word,ID) %>%
  ungroup() %>%
  distinct(word) # distinct does so we only retain the unique words from the
nodelist -> so we dont get duplicate words/nodes.
```

Then we create a dataframe that we are going to be use for *left_joining* the edges (word.x and word.y)

Hide

```
nodes4join <- reviews_sampled_tidy_joined %>%
  select(ID, word)
```

Hide

```
g_reviews <- reviews_sampled_tidy_joined %>%
  left_join(nodes4join, by = c("ID" = "ID"))
```

Hide

```
g_reviews <- g_reviews[c("word.x", "word.y")]
```

Here we transform the dataframe into a graph using *igraph*.

We set *directed = FALSE*, because it is an undirected network.

Hide

```
g_reviews <- graph_from_data_frame(d = g_reviews,
                                   directed = FALSE,
                                   vertices = nodes_reviews)
```
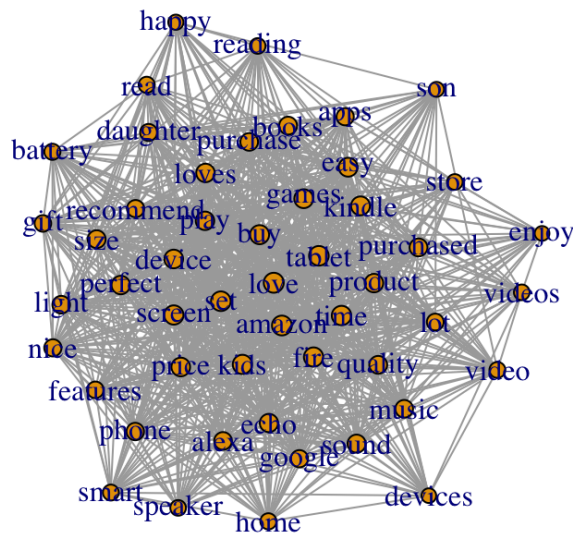
We use *simplify* to remove loops and multiple edges.

Hide

```
g_reviews <- simplify(g_reviews,
                      remove.multiple = TRUE,
                      remove.loops = TRUE,
                      edge.attr.comb = igraph_opt("edge.attr.comb"))
```

Now we can plot the network.

Hide

```
plot(g_reviews,
     vertex.size = 1 + sqrt(degree(g_reviews, mode = "all"))) # we assign th
e size to the degree and add "1+sqrt", so we dont get too big nodes.
```

At the graph we see that the size based on the degree is not that different between the words, but is seems as *device*, *tablet* & *love* is connected by *ID* to many words. This is not that surprising, because these where also some of the topwords from earlier.

Lets have a look at each nodes degree.

Hide

```
degree(g_reviews) %>%
  as.data.frame()
```

| | .<br><dbl> |
|---|---|
| tablet | 47 |
| games | 43 |
| recommend | 35 |
| kindle | 43 |
| alexa | 37 |
| love | 49 |
| features | 32 |
| video | 29 |
| screen | 45 |
| perfect | 41 |

| 1-10 of 50 rows | Previous 1 2 3 4 5 Next |
|---|---|

We can see that fx. *tap* only has degree $25$, but *love* has degree $49$. This makes sence, that not that many reviews has *tap* in the text, but rather more has *love* inside *(probably because the many high rated reviews)*.

# Alternative solution

We could have looked into bigrams and plot those as a network.

<div align="right">Hide</div>

```
bigrams <- reviews %>%
  unnest_tokens(bigram, text, token = "ngrams", n = 2) # n is the number of
 words to conside
```

<div align="right">Hide</div>

```
bigrams_separated <- bigrams %>%
  separate(bigram,c("word1","word2"),sep = " ")
```

<div align="right">Hide</div>

```
bigrams_filtered <- bigrams_separated %>%
  filter(!word1 %in% stop_words$word) %>%
  filter(!word2 %in% stop_words$word) %>%
  filter(!word1 %in% own_stopwords$word) %>%
  filter(!word2 %in% own_stopwords$word)
```

<div align="right">Hide</div>

```
bigram_counts <- bigrams_filtered %>%
  count(word1, word2, sort = TRUE)
```

<div align="right">Hide</div>

```
bigram_counts %>%
  head()
```

| word1 | word2 | n |
| :--- | :--- | ---: |
| <chr> | <chr> | <int> |
| kindle | fire | 145 |
| battery | life | 100 |
| amazon | fire | 78 |
| amazon | echo | 76 |
| amazon | prime | 66 |
| smart | home | 66 |

6 rows

<div align="right">Hide</div>

```
bigram_graph <- bigram_counts %>%
  filter(n > 15) %>%   #The occurence of the bigram is more than 15.
  graph_from_data_frame(directed = FALSE)
```

Hide

```
bigram_graph %>% ggraph(layout = "fr") +
  geom_edge_link() +
  geom_node_point(aes(size = degree(bigram_graph), colour = degree(bigram_
graph))) +
  scale_color_continuous(guide = "legend") +
  theme_graph() +
  geom_node_text(aes(label = name))
```