



ERASMUS UNIVERSITY ROTTERDAM

Erasmus School of Economics

Master thesis Data Science and Marketing Analytics

THE EFFECT OF SOCIAL MEDIA CONTENT ON CONSUMER ENGAGEMENT

Does social media content on trending topics improve online
consumer engagement?

Supervisor: Karpienko, R.
Second assessor: Dekker, R.

Sabine Perquin
425700

Table of content

1	Introduction	3
2	Theoretical framework	5
2.1	<i>Defining owned social media</i>	5
2.2	<i>Defining earned social media</i>	5
2.3	<i>The effect of OSM on business performance</i>	6
2.4	<i>The effect of ESM on business performance</i>	6
2.5	<i>Brand versus trend consistent content</i>	7
2.6	<i>The effect of emotions in OSM on ESM engagement</i>	8
2.7	<i>The effect of trend consistent content on ESM for different brand types</i>	9
3	Data	12
3.1	<i>Dataset with the selection of brands</i>	12
3.2	<i>Dataset on the timelines of the brands</i>	15
3.3	<i>Dataset on brand related content posted by consumers</i>	17
3.4	<i>Data cleaning and pre-processing</i>	18
4	Methodology	18
4.1	<i>Measuring the four different brand types by performing sentiment analysis</i>	18
4.1.1	<i>Measuring brand love with sentiment analysis</i>	18
4.1.2	<i>First dimension for measuring brand love: emotional engagement</i>	20
4.1.3	<i>Second dimension for measuring brand love: respect</i>	20
4.1.4	<i>The four brand types</i>	21
4.2	<i>Classifying brands with hierarchical clustering</i>	22
4.3	<i>Measuring owned social media</i>	25
4.4	<i>Measuring trend consistent content with Latent Dirichlet Allocation</i>	25
4.5	<i>Measuring earned social media</i>	28
4.6	<i>Linear regression to determine the effect of OSM on ESM</i>	29
4.7	<i>Random Forest to determine the importance of emotions on ESM</i>	30
4.8	<i>LIME to determine the effect of sentiment on predictions</i>	32
5	Results	32
5.1	<i>Sentiment analysis on the brand related content</i>	32
5.2	<i>Hierarchical clustering results for classification of the brands</i>	33
5.2.1	<i>Classification of the brands on the basis of the hierarchical clustering results</i>	34
5.3	<i>LDA results for classification of owned social media content</i>	36
5.4	<i>Regression analysis results to determine the effect of OSM on ESM</i>	40
5.4.1	<i>Variables included in the linear regression</i>	41
5.4.2	<i>Descriptive statistics on the variables used in the linear regression</i>	41
5.4.3	<i>Linear regression results</i>	43
5.4.4	<i>Hypothesis testing</i>	45
5.5	<i>Random forest to analyse the importance of emotions on the prediction of ESM</i>	46

5.6	<i>LIME to determine the effect of joy, trust and anger on ESM</i>	46
6	Conclusion	47
	References	50
	Appendix A	53
	Appendix B	57

Abstract

In this paper, the effect of jumping on social media trends on online consumer engagement is studied by looking at two types of social media content that is posted by brands on Twitter. It is expected that consumer engagement is influenced by the loving relationship that brands have with their consumers, which is referred to as brand love. Therefore, the brands used in this analysis are classified into one of four brand types of the basis of the level of brand love. Using a unique dataset that was scraped from Twitter, I find that the extent to which consumers engage in online social media depends on the type of content that was posted and on the degree of brand love. Only for loved brands, jumping on social media trends benefits consumer engagement. For brands that do not have a loving relationship with their consumers, consumer engagement is improved by posting online marketing content that is consistent with the traditional marketing message, instead of jumping on media trends.

1 Introduction

Many brands use their social media channels for posting marketing content on social movements that are trending among the public in order to engage more with their consumers. Social media marketing has become a very common and important part of the media mix of firms. Posting marketing content on social media channels results in benefits for the firm, among which increased consumer engagement and improved brand performance (Kumar et al., 2016; Stephen & Galak, 2012). Whereas marketers traditionally strive for consistency in their marketing message to create strong and loved brands, heaps of firms deviate from their conventional marketing message nowadays to follow media trends in their online content (Kay, 2006). However, where some marketing campaigns on media trends increase engagement among consumers, other social media content fails and does not catch on. Thus, understanding how consumers conceive social media content is key for firms when determining whether to jump on media trends or not if they wish to improve their consumer engagement.

Existing literature has examined how social media content posted by firms can increase consumer engagement. Consumer engagement is important for brands as it leads to increased satisfaction and more loyalty. Furthermore, stronger consumer-brand relationships are established and purchases are increased (Brodie et al., 2013; Laroche et al., 2013; van Doorn et al., 2010). However, literature lacks on determining whether different types of online marketing content have different effects on consumer engagement. This research determines this differential effect by distinguishing two different types of social media content: brand consistent content and trend consistent content. *Brand consistent content* is defined as content on brands' social media pages that is consistent with the brand message. Social media content that deviates from the brand message and that is on topics that receive major media attention, such as societal concerns or crises, is referred to as *trend consistent content*. It is expected that consumers react differently to social media content concerning societal matters and crises than to brand consistent content, since social and emotional factors influence social media engagement of consumers heavily (Lovett et al., 2013). Thus, it is important for brands to understand how consumers react to their social media content so that they can improve consumer engagement via their social media channels. Once firms understand their customers, they can create strong and loved brands and engage with their customers (Vernuccio et al., 2015). Therefore, the following question is proposed:

To what extent does brand versus trend consistent social media content encourage consumers to engage in social media content creation?

Trend consistent content is beneficial in improving consumer engagement because of two reasons. First, consumers pay more attention to posts on trends in the media because consumers are already aware of these topics, which in return leads to increased consumer engagement (Brodie et al., 2013). Secondly, consistency in posting brand consistent content is increasingly difficult to pursue, because brands should connect with consumers on an emotional level, which are subject to change. However, in contrast to these two allegations, consumers understand a brand better if brands employ consistency in their branding decisions (Kay, 2006). Thus, existing literature is inconclusive on whether jumping on trends is more beneficial for firms than pursuing consistency in the marketing message. This study fills this gap by examining the differential consumer response to brand versus trend consistent content.

The effect of the two different content types is examined by using social media content from the social media platform Twitter. A unique characteristic of Twitter is, that compared to other social media channels, a large volume of content is created everyday (Zhang et al., 2011), which was the main motivator for using Twitter in this analysis. A unique Twitter dataset was extracted including 82 brands, for which the effect of different types of social media content on consumer engagement is studied. Based on existing theory, it is believed that consumers perceive social media content from brands they love differently than brands that they do not love (Berger & Milkman, 2012; Lovett et al., 2013). Therefore, this research distinguishes four different brand types on the basis of brand love, in which the 82 brands are classified. The tweets posted by these brands were classified into brand or trend consistent content based on the topics present in the tweets. If the tweet is on either COVID-19 or the Black Lives Matters movement, the tweet is classified as trend consistent content. Once the brands and tweets have been classified, the effect of brand versus trend consistent content on consumer engagement was examined. Findings suggest that only brands that have already established a loving relationship with their consumers gain from posting on social media trends. Other brands should first create a loved and strong brand by posting brand consistent content.

This research paper is structured in the following way. The next section discusses existing literature on social media content, after which four different hypotheses are formed. Section 3 outlines the data collection and preparation process, and Section 4 includes the methodology of

this research. The results of the analyses are discussed and interpreted in Section 5, after which a conclusion and discussion is given in Section 6.

2 Theoretical framework

Previous studies categorize social media content into two categories: owned and earned social media. This study examines the effect of owned social media on earned social media. Therefore, this section first elaborates on owned and earned social media content and their relevance. Thereafter, two types of owned social media are distinguished: brand and trend consistent content. It is expected that trend consistent content triggers emotions among consumers, and therefore the influence of emotions on consumer engagement is discussed. Lastly, based on existing theory on earned social media creation, it is believed that brand love affects online consumer engagement greatly. Therefore, four different types of brands are distinguished on the basis of the literature on brand love.

2.1 Defining owned social media

Owned social media (OSM) is content created and controlled by the brand itself on their social media network sites (Colicev et al., 2018). An example of OSM is brand-initiated marketing communication on its social media network site such as their Instagram page or Twitter account. In this research, OSM is defined as all content posted by brands on their official social media page. More specifically, OSM are all organic tweets, which are all tweets excluding replies and retweets, created and posted by the brands on their Twitter timeline.

2.2 Defining earned social media

The second type of media content, *earned social media* (ESM), is published content on brands that is created by consumers on social media networks (Smith et al., 2012). Creating brand content is no longer only controlled by the brands themselves. With the help of social media channels everyone is able to generate and publish online content nowadays. Existing literature uses several different terms when discussing ESM content, such as user-generated content and electronic word-of-mouth. These terms all refer to social media content on brands that is created by and shared among users of social media platforms (Choi & Lee, 2017; Dhar & Chang, 2009; Goh et al., 2013). In this study, ESM content is defined as all content created by consumers in response to OSM posted by brands. More specifically, ESM content is the number retweets that OSM posted by brands receive.

2.3 The effect of OSM on business performance

Posting OSM content is mutually beneficial for brands and consumers. Consumers gain from OSM as they learn about the products and services offered by the brand, and brands post OSM to improve their customer relationships (Kivetz & Simonson, 2000). This results in more peer-to-peer recommendations and increased consumer engagement, which leads to increased business performance (Kumar et al., 2016). However, where some OSM succeeds in this and is catching on, other content fails tremendously. Thus, it is crucial for companies to get a better understanding of how consumers react to OSM.

Colicev et al. (2018) study the effect of OSM on three consumer mindset metrics: purchase intent, brand awareness and consumer satisfaction. Their results show the benefits of OSM, as it leads to more customer satisfaction and brand awareness. This implies that marketers should use their social media marketing strategies to increase brand awareness and improve customer satisfaction rather than for other purposes. They also find that OSM activity has a positive effect on ESM, which they measure by the number of Twitter followers and ESM volume engagement. This implies that marketers can improve ESM by engaging in OSM on their social media channels. Improved ESM has positive effects for firms, which are discussed in the following section.

2.4 The effect of ESM on business performance

The need for ESM content has been growing and many brands have, for example, set up brand communities on online social networks to increase customer engagement (Goh et al., 2013). The following studies have examined the effects of ESM on business performance and show how ESM can be beneficial for firms.

First, Chevalier & Mayzlin (2006) study the effect of ESM on sales performance by studying online reviews. They focus on the sentiment in reviews and how this affects book sales. Results show that an additional, positive review on a book has a positive effect on sales whereas an additional negative review decreases sales. Moreover, an additional negative review is more powerful in decreasing sales of a book than a positive review is in increasing the sales. Thus, sentiment in earned media affects sales, where the effect of negative earned media on sales is larger than the effect of positive earned media.

The influence of ESM on future sales has also been examined by Dhar & Chang (2009). However, in contrast to Chevalier & Mayzlin (2006), they use blog posts instead of reviews to determine the effect of ESM on sales. In their research, they use online music sales and hypothesize that social networks matter for these sales. Thus, they do not examine the sentiment in ESM but rather the volume. Their results show that the volume of blog posts is positively correlated with future sales. This implies that brands can increase their future sales by increasing ESM volume.

Lastly, Kim & Johnson (2016) examine the effect that ESM has on consumer decision-making on the social media platform Facebook. They investigate whether ESM triggers an emotional and/or cognitive response within consumers that translates into certain consumer behaviour. Within their framework, consumers' reaction to ESM was examined on the following aspects: pleasure, arousal, information quality, information pass-along, impulse buying and future-purchase intention. Brand awareness was used as a control variable to control for any pre-existing knowledge that can influence a consumer's attitude towards a brand. Findings revealed that emotional and cognitive responses caused by ESM posted by other consumers significantly influenced consumer behaviour and their attitude towards a company. This implies that ESM on social media platforms influences consumer behaviour and their opinions on a brand.

All in all, existing literature shows the importance of ESM for brands. ESM does not only positively affect sales performance, but it also influences consumer behaviour and their attitude towards a company. Therefore, it is important for brands to optimize the ESM creation among consumers. This study will examine how brands can achieve this through their OSM content by analysing consumer responses to brand versus trend consistent content.

2.5 Brand versus trend consistent content

Existing literature determines the importance of OSM and ESM for brands. However, existing literature does not take into consideration that brands post OSM on different topics and how this influences ESM. Therefore, this research identifies two different types of OSM and determines their effect on ESM.

The two types of OSM used in this research are: brand consistent content and trend consistent content. Brand consistent content is OSM content that is in line with the brand message. Trend consistent content deviates from this traditional brand message and is on topics that do not

regard the brand itself but rather regard societal matters and concerns. In this research, brand consistent content is distinguished from trend consistent content by considering two global societal concerns. First, the situation around the Corona virus (COVID-19) is considered. COVID-19 is an infectious disease that caused a pandemic end of 2019, which induced a global health crisis (World Health Organisation 2020). This crisis led to a movement in the marketing industry, as many brands deferred planned marketing campaigns and dedicated new campaigns to the virus (Warc.com, 2020). The second social movement that is considered in this research is the Black Lives Matter (BLM) movement. BLM is a global organization with the mission to exterminate white supremacy and racism (Black Lives Matter 2020). After a violent incident in the United States between a black man and a police officer on the 25th of May 2020, the BLM movement caught people's attention globally. Where some brands suddenly devoted their entire marketing campaigns on the movement, other brands remained silent. The difference in reaction of brands towards these two concerns makes the two matters suitable for this research. Thus, tweets that are on COVID-19 or BLM are classified as trend consistent content. The remainder of the tweets are classified as brand consistent content.

2.6 The effect of emotions in OSM on ESM engagement

The trend consistent content on the Corona crisis and BLM movement triggers emotions among consumers as the content is about tragic real-life concerns. Existing literature shows that emotions present in social media content have an effect on the drive for ESM creation. It is expected that trend consistent content leads to more ESM. However, because people conceive media from brands that they love differently than from brands that they do not love, this expectation only holds for brands that have already established a loving relationship with their consumers.

Berger & Milkman (2012) examine New York Times articles to discuss how emotion in online content has an influence on why consumers share content and how often they share this. They state that practical information is mostly shared because of altruistic reasons and to help others. Sharing this practical, useful information has social exchange value and may generate reciprocity. Emotional content is shared by consumers to reduce dissonance, make sense of their experience or to deepen the connections with their social networks. This implies that consumers share trend consistent content in order to make sense of the two topics COVID-19 and BLM. Their findings also reveal that positive content is more viral than negative content.

Reasoning for this is that consumers often share content to present themselves, and consequently, positive content may reflect consumers more positively.

Lovett et al. (2013) their research is in line with this. They also state that emotional drivers cause consumers to create ESM. Additionally, they identify social and functional drivers. Social drivers relate to self-enhancement and the desire to socialize with other online consumers and the functional driver represents the need to obtain and share practical information. This suggests that brand consistent content, which is more practical information about brands, is shared because of functional drivers and that trend consistent content is shared because of emotional drivers. Furthermore, their results show that positive content is shared more than negative content. The reasoning behind this proposition is in line with that of Berger & Milkman (2012), which is that consumers share positive content to improve their image and identity. However, their research goes beyond the valence of content and they also examine how specific emotions present in OSM evoke different reactions among consumers. Results show that marketing content that evokes emotions characterized by arousal (e.g., amusement and anger) is more likely to be shared than content with emotion that is characterized by deactivation (sadness), regardless of the valence of the emotion.

The study of Daugherty et al. (2008) finds similar results as the previous two studies. Besides emotions playing a role, consumers create content because it helps them understand their environment and the topic at hand and because of self-enhancement. By creating content, they feel a sense of substantial wisdom and they feel part of a community that share the same beliefs and values. This means that consumers create ESM on trend consistent topics to understand these topics and feel part of, for example, the BLM community.

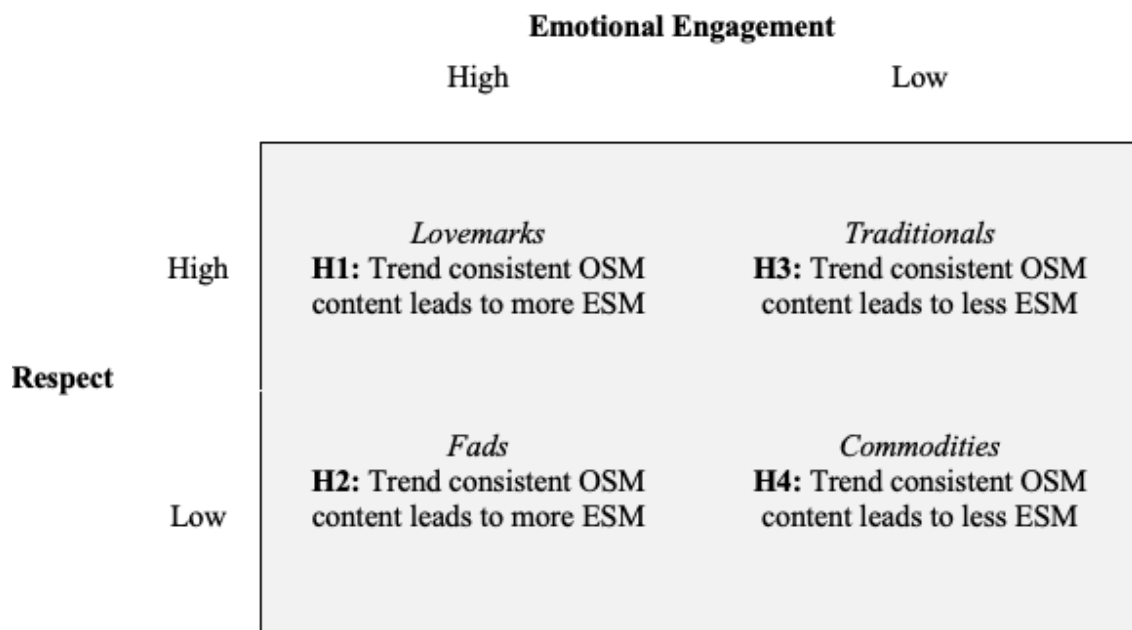
2.7 The effect of trend consistent content on ESM for different brand types

It is expected that trend consistent content stimulates ESM engagement more than brand consistent content. Reasoning for this is that content that evokes arousal is shared more by consumers. Furthermore, consumers create ESM to get a deeper understanding of the topic, to deepen their connections with their network and for self-enhancement. However, consumers perceive content posted by loved brands differently than content posted by brands with which they do not have a loving relationship. Therefore, it is expected that posting trend consistent content has different effects on consumer engagement for brands with different levels of brand love.

Brand love goes beyond simply liking a brand. In this research, brand love is seen as a broad and long-term relationship between the consumer and the brand. It is the degree of emotional attachment that a consumer has towards a brand (Batra et al., 2012; Wallace et al., 2014). Based on the level of brand love, four types of brands are distinguished in this research along two dimensions, emotional engagement and respect: 1) Lovemarks, 2) Fads, 3) Traditionals and 4) Commodities. An overview of the four brand types with corresponding hypotheses can be obtained in Figure 2.1.

Figure 2.1

Framework of different brand types



The four different brand types with their score on the two dimensions emotional engagement and respect and their corresponding hypotheses.

The first brand type is the type of brand that is loved most by consumers and is called *Lovemarks*. Batra et al. (2012) state that loved brands use social betterment in their marketing campaigns, have existential brand meaning and bond with consumers emotionally. Additionally, loved brands should connect to deeper meanings and important values in life. Furthermore, brand love has a positive effect on word-of-mouth creation and engagement in ESM (Carroll & Ahuvia, 2006). Combining these findings on brand love with the literature that states that consumers create content to feel part of a community and to understand their environment and the topic at hand (Daugherty et al., 2008; Hennig-Thurau et al., 2004), it is expected that consumers engage more in ESM if Lovemarks post on trend consistent topics. Therefore, the first hypothesis is developed:

Hypothesis 1: For Lovemarks (brands with the strongest loving relationship with consumers), trend consistent owned social media content leads to more earned social media.

The second type of brands is called *Fads*. Consumers are enthusiastic and positive about Fads and feel love for these brands. However, consumers have a lack of respect for these brands. Therefore, the brand will eventually disappear, and the consumer-relationship is short lived. Thus, consumers feel brand love for Fads but to a lesser extent than for Lovemarks. Expected is that, because these brands are loved by consumers, consumers create an adequate amount of ESM for Fads (Wallace et al., 2014). Furthermore, because Fads are a hype among consumers, it is relevant for these brands to respond to topics that consumers are already aware of, making it more likely for them to pay attention and engage in ESM. Therefore, posting trend consistent content rather than brand consistent content will create more ESM for Fads, which leads to the second hypothesis:

Hypothesis 2: For Fads (brands that have a loving but short lived relationship with consumers), trend consistent owned social media content leads to more earned social media.

Traditionals is the third type of brands distinguished in this research. These brands are respected by consumers but lack in creating a loving relationship with their consumers. Therefore, these brands are seen as more traditional brands that do not create a form of passion or attachment with their consumers. Because consumers do not feel any emotional engagement towards these brands, they are less triggered to create ESM. Thus, for Traditionals it is important to create brand meaning and value in order to establish brand love and improve ESM (Batra et al., 2012). This can be achieved by getting more understanding among consumers for the brand, which is done by pursuing consistency in the marketing content (Kay, 2006). Therefore, it is expected that Traditionals should stick to brand consistent content and should not jump on media trends in their OSM. By pursuing brand consistent OSM, a loved brand can be created which leads to more ESM engagement. Hence, it is expected that trend consistent content has a negative effect on ESM and the third hypothesis was developed:

Hypothesis 3: For Traditionals (brands that have no loving relationship with consumers but are respected), trend consistent owned social media content leads to less earned social media.

The last brand type that is considered is *Commodities*. These are brands that are not respected by consumers nor loved. The products or services sold by these brands are merely seen as commodities and the brands do not have any type of relationship with their consumers. Similarly, as for Traditionals, it is important for Commodities to create a strong brand meaning to create more brand love and, in return, more ESM engagement. Moreover, consistency in all aspects of the brand strategy is essential as they have to develop trust and respect among consumers, which is best done by repeated experiences with the brand (Elliott & Yannopoulou, 2007). Therefore, these brands should focus on content that is consistent with their brand message and brand meaning. If Commodities post on trending topics rather than brand consistent content, this may lead to misunderstanding among consumers and less engagement. This leads to the last hypothesis:

Hypothesis 4: For Commodities (brands that have no loving relationship with consumers nor are respected), trend consistent owned social media content leads to less earned social media.

3 Data

The data needed for testing the hypotheses requires an extensive and challenging data collection process. In total, three different datasets are used. The first dataset contains the selection of brands, their brand value and the industry in which they operate. The other two datasets are collected by scraping data from Twitter with the use of a Twitter Application Programming Interface (API), which enables the reading and scraping of Twitter data. The second dataset contains data on the timelines of the brands. This dataset contains the OSM and ESM on the tweets posted by the brands. Furthermore, this dataset contains the date on which the tweet is posted, the number of favourites, number of followers of the brand and a unique user- and status ID. The third dataset is composed of tweets about the brands that are posted by consumers. This dataset is used for distinguishing the four brand types among the brands.

3.1 Dataset with the selection of brands

The first data frame used in this study contains the selection of brands. From these brands, the OSM and ESM are extracted later on. The selection of brands is done on the basis of the Forbes' annual list of the World's Most Valuable Brands of 2019, which ranks the top 100 brands based on their brand value. This valuation of brand value is performed on the basis of the brand's revenue and earnings before interest and taxes, corrected for the brand's capital and the role the

brand plays in each industry. These 100 most valuable brands are collected in a data frame, where each row corresponds to a brand. For every brand, the brand value and industry are given as stated by Forbes. Thereafter, the main Twitter account for each brand is looked up in order to extract their Twitter data. If brands have multiple Twitter accounts, the Twitter account that has the most followers and that posts tweets in English is used. From the selection of brands, three brands are omitted. Marlboro is omitted because this brand does not have a Twitter account. Additionally, Costco and Apple are omitted because they do not post any tweets on their Twitter accounts. This resulted in collecting Twitter accounts for a selection of 97 brands (Table 3.1). The top five brands are the Big Five tech giants: Apple (omitted from final selection), Google, Microsoft, Amazon and Facebook. Furthermore, there is a variety of industries in the selection, such as the beverages industry, automotive industry and apparel industry.

Table 3.1

Brand selection

Rank	Brand	Industry	Brand Value	Number of tweets	Time frame
2	Google	Technology	\$ 207.5 B	3200	44
3	Microsoft	Technology	\$ 162.9 B	3198	514
4	Amazon	Technology	\$ 135.4 B	3199	557
5	Facebook	Technology	\$ 70.3 B	3199	698
6	Coca-Cola	Beverage	\$ 64.4 B	3200	149
7	Disney	Leisure	\$ 61.3 B	3200	778
8	Samsung	Technology	\$ 50.4 B	3200	1854
9	Louis Vuitton	Luxury	\$ 47.2 B	3198	1623
10	McDonald's	Restaurants	\$ 46.1 B	3200	10
11	Toyota	Automotive	\$ 41.5 B	3200	54
12	Intel	Technology	\$ 39.5 B	3200	1061
13	Nike	Apparel	\$ 39.1 B	3200	919
14	AT&T	Telecom	\$ 37.3 B	3200	55
15	Cisco	Technology	\$ 36.0 B	3196	600
16	Oracle	Technology	\$ 35.7 B	3199	738
17	Verizon	Telecom	\$ 32.3 B	3198	139
18	Visa	Financial Services	\$ 31.8 B	3196	2409
19	Walmart	Retail	\$ 29.5 B	3200	29
20	GeneralElectric	Diversified	\$ 29.5 B	3185	1570
21	Budweiser	Alcohol	\$ 28.9 B	1555	2389
22	SAP	Technology	\$ 28.6 B	3200	490
23	Mercedes-Benz	Automotive	\$ 28.5 B	3200	117
24	IBM	Technology	\$ 28.2 B	3200	296
26	Netflix	Technology	\$ 26.7 B	3149	237
27	BMW	Automotive	\$ 25.9 B	3200	419

28	AmericanExpress	Financial Services	\$ 25.1 B	3199	328
29	Honda	Automotive	\$ 24.5 B	3200	863
30	L'Oréal	Consumer Packaged Goods	\$ 22.8 B	3190	1162
31	Gucci	Luxury	\$ 22.6 B	3200	947
32	Hermès	Luxury	\$ 21.6 B	351	1849
33	Nescafé	Beverage	\$ 20.4 B	3200	1257
34	Home Depot	Retail	\$ 19.2 B	3200	1484
35	Accenture	Business Services	\$ 19.1 B	3169	1373
36	Pepsi	Beverage	\$ 18.2 B	3200	240
37	Starbucks	Restaurants	\$ 17.8 B	3200	38
38	Mastercard	Financial Services	\$ 17.3 B	3199	598
39	Frito-lay	Consumer Packaged Goods	\$ 16.3 B	3200	491
40	IKEA	Retail	\$ 15.8 B	3200	1621
41	Zara	Retail	\$ 14.7 B	3200	237
42	Gillette	Consumer Packaged Goods	\$ 14.5 B	3200	573
43	HSBC	Financial Services	\$14.4 B	3082	2751
44	Audi	Automotive	\$ 13.8 B	3200	451
45	J.P. Morgan	Financial Services	\$ 13.7 B	3199	1481
46	Deloitte	Business Services	\$ 13.5 B	3197	674
47	Sony	Technology	\$ 13.3 B	3151	475
48	UPS	Transportation	\$ 13.3 B	3200	124
49	Bank of America	Financial Services	\$ 13.2 B	3194	1799
50	Chase	Financial Services	\$ 13.1 B	3199	1462
51	Adidas	Apparel	\$ 12.9 B	3199	1070
52	Chanel	Luxury	\$ 12.8 B	3106	3297
53	Siemens	Diversified	\$ 12.7 B	3200	1031
54	Nestle	Consumer Packaged Goods	\$ 12.3 B	3196	802
55	CVS	Retail	\$ 12.3 B	3200	98
56	Cartier	Luxury	\$ 12.2 B	2382	2511
57	Porsche	Automotive	\$ 12.1 B	3196	2769
58	ESPN	Media	\$ 11.9 B	3200	119
59	Citi	Financial Services	\$ 11.8 B	3200	988
60	Wells Fargo	Financial Services	\$ 11.8 B	3194	412
61	Adobe	Technology	\$ 11.5 B	3200	661
62	Pampers	Consumer Packaged Goods	\$ 11.5 B	3200	542
63	Corona	Alcohol	\$ 11.4 B	3185	2689
64	T-Mobile	Telecom	\$ 11.4 B	3200	274
65	eBay	Technology	\$ 11.3 B	3198	1305
66	Chevrolet	Automotive	\$ 11.3 B	3200	216
67	PayPal	Technology	\$ 11.3 B	3198	1811
68	Ford	Automotive	\$ 11.2 B	3200	216
69	Red Bull	Beverage	\$ 11.1 B	3175	937
70	PwC	Business Services	\$ 11.0 B	3200	838
71	HP	Technology	\$ 11.0 B	3197	1804
72	Colgate	Consumer Packaged Goods	\$ 10.7 B	3200	961
73	Fox	Media	\$ 10.6 B	3179	374

74	Lowe's	Retail	\$ 10.5 B	3200	51
75	Lancôme	Consumer Packaged Goods	\$ 10.4 B	3190	1532
76	H&M	Retail	\$ 10.4 B	3199	1264
77	Lexus	Automotive	\$ 10.3 B	3200	917
78	Santander	Financial Services	\$ 9.7 B	3200	1080
80	Rolex	Luxury	\$ 9.5 B	225	659
81	Hyundai	Automotive	\$ 9.5 B	3200	348
82	Danone	Consumer Packaged Goods	\$ 9.3 B	3199	2630
83	Heineken	Alcohol	\$ 9.3 B	3200	2016
84	Uniqlo	Apparel	\$ 9.2 B	3199	933
85	Goldman Sachs	Financial Services	\$ 8.9 B	3200	607
86	Hennessy	Alcohol	\$ 8.9 B	3154	2592
87	Nintendo	Technology	\$ 8.8 B	3198	594
88	AXA	Financial Services	\$ 8.8 B	3198	1637
89	Allianz	Financial Services	\$ 8.8 B	3196	1452
90	Dell	Technology	\$ 8.7 B	3199	389
91	Caterpillar	Heavy Equipment	\$ 8.6 B	3200	791
92	LEGO	Leisure	\$ 8.6 B	3200	196
93	Huawei	Technology	\$ 8.5 B	3200	107
94	John Deere	Heavy Equipment	\$ 8.4 B	3200	853
95	UBS	Financial Services	\$ 8.3 B	3199	1228
96	KFC	Restaurants	\$ 8.3 B	3200	41
97	Burger King	Restaurants	\$ 8.2 B	3200	111
98	EY	Business Services	\$ 8.0 B	3200	552
99	FedEx	Transportation	\$ 7.9 B	3200	555
100	Volkswagen	Automotive	\$ 7.9 B	3200	933

Selection of brands with their rank, industry and brand value in billion US dollars.

3.2 Dataset on the timelines of the brands

The second dataset that is collected in this analysis contains data on the timelines of the brands. This dataset is scraped from Twitter with the use of the Twitter API and contains the OSM and ESM of the brands.

The Twitter API allows to retrieve a maximum of 3200 most recent tweets on the brands' Twitter timelines. Thus, for every brand up to 3200 of their latest tweets are collected and a dataset was composed where every row corresponds to a tweet of one of the brands. For every tweet, the original tweet text is obtained, and a distinction is made between organic tweets, retweets and replies to users. In this research, the organic tweets are the OSM content of the brands. The number of replies by each brand is the measure for frequency that is used to classify brands into different brand types. The retweets that brands post on their timelines are not used in this study. Besides the text of the tweet, the number of retweets is collected. These retweets are the measure of ESM in this analysis. Other variables retrieved by the Twitter API are the

date on which the tweet is posted, the number of favourites, number of followers of the brand and a unique user- and status ID.

After retrieving the data on the timelines of the brands, inactive brands are omitted from the selection of brands. Brands are regarded as inactive when they have not posted the 3200 tweets within the last five years. This results in omitting ten brands from the selection. Furthermore, the two brands Rolex and Hermès are omitted since they have only posted around 300 tweets within the past five years and are therefore also seen as inactive brands. The final selection of brands includes 85 brands. This selection, together with the variables retrieved by the Twitter API, can be obtained in Table 1, Appendix A.

The 85 brands have posted an average of 3197 tweets per brand over an average period of 721 days, which is equivalent to almost two years (Table 3.2). The total timeframe of collected tweets is approximately five years and ranges from August 2015 until October 2020. Brands devote their Twitter account mostly to replying to users as, on average, half of the total number of tweets posted on the timeline of a brand are replies. The fraction of OSM content posted by brands out of the total number of tweets posted is slightly below this, namely 35.60% on average. This means that, on average, 35.60% of the content posted by brands is OSM content. The fraction of retweets is fairly low, only 9.35% on average. This implies that brands use their Twitter accounts mostly to interact with consumers followed by posting OSM content.

Table 3.2

Descriptive statistics

	Number of tweets	Percentage of replies	Percentage of organic tweets	Percentage of retweets	Timeframe in days
Minimum	3149	0.03%	0.25%	0.00%	10
Mean	3197	55.06%	35.60%	9.35%	721
Maximum	3200	99.75%	99.12%	45%	1854

Descriptive statistics on the number of tweets, the percentage of replies, organic tweets (OSM) and retweets out of the total number of tweets posted by a brand and the timeframe.

The brands that interact most with their consumers through Twitter are brands from the food and beverages industry such as McDonalds, KFC and Coca-Cola (Table 1, Appendix A). The majority of tweets posted by these brands are replies and only a small percentage of their total number of tweets is OSM. Furthermore, McDonalds is the most active brand on Twitter, as they

have posted 3200 tweets within a timeframe of only 10 days. In contrast to these brands, luxury fashion brands such as Gucci and Louis Vuitton, and banks like Goldman Sachs and UBS post almost only OSM and rarely reply to consumers.

3.3 Dataset on brand related content posted by consumers

In order to classify the brands into one of the four different brand types, a third dataset is collected with the use of the Twitter API. This dataset contains tweets about the brands that are posted by consumers. From these tweets, the opinions and thoughts of consumers on the brands can be extracted to distinguish different types of brands. Thus, this third dataset is used for performing sentiment analysis on, after which the brands are clustered and classified into one of the four brand types.

The Twitter API can only retrieve search results from six to nine days ago. A maximum of 3000 English tweets is retrieved per brand, resulting in a dataset with 214,694 observations. This number of tweets is chosen based on previous research, where sentiment analysis on tweets was performed with roughly the same number of tweets (Pak & Paroubek, 2010; Zhang et al., 2011). The tweets that are collected are tweets that contain mentions to one of the brands. These tweets can be replies that consumers specifically address towards the brand but also tweets in which the brand is mentioned casually. Instead of extracting tweets that include mentions to the brand, the Twitter API can also retrieve tweets with a hashtag of the brand name. However, the choice of using brand mentions instead of hashtags is made to be certain that consumers talk about the brand instead of other topics. An example where using a hashtag will not result in useful tweets is the brand Oracle as many tweets with #oracle are on other topics than this brand.

After retrieving the brand related content, brands that consumers do not talk actively about are omitted from the brand selection to make sure that there is enough data to perform analyses on. This results in omitting three brands, Pampers, Lancôme and Santander. Therefore, the final brand selection contains 82 brands (Table 2, Appendix A). The total volume of brand related content retrieved for these 82 brands is 196,677 tweets, which results in 2,396 brand related tweets per brand on average. The timeframe of this third dataset is nine days, from the 23rd of October 2020 until the 2nd of November 2020.

3.4 Data cleaning and pre-processing

The second and the third dataset contain tweets that are used in the analyses. Hence, the text in these tweets needs cleaning and pre-processing to reduce noise. First, non-alphabetic signs, excess spaces, excess punctuation, URLS and pictures are removed from the text using regular expression techniques. These regular expression techniques are also used to remove Twitter specific text such as retweet headers, usernames and the hash for hashtags. Furthermore, repeated characters are replaced with single characters, for example: ‘hellooooo’ is replaced with ‘hello’. In order to remove stop words, a list with stop words is created which includes standard words. Twitter related stop words and the 82 brand names are added to this list, since the brand names will not add any explanatory value. Lastly, because Twitter is a social media platform, many tweets contain emoticons. These emoticons can be punctuation-based emoticons or emojis, which are images. Both types of emotions are replaced with the corresponding meaning.

4 Methodology

4.1 Measuring the four different brand types by performing sentiment analysis

Based on the existing literature, four different brand types are distinguished from one another on the basis of the level of brand love (Roberts, 2005; Thomson et al., 2005). This research defines two dimensions of brand love: emotional engagement and respect. The dimensions are measured by the frequency of interaction that brands have with consumers and the extent to which consumers post emotional content on the brands.

4.1.1 Measuring brand love with sentiment analysis

In order to determine what emotions are present in the content posted by consumers and, based on these emotions, measure brand love, sentiment analysis is performed. This method is performed on the third dataset, which includes the brand related tweets posted by consumers. From the results of the sentiment analysis, every brand is given an average score for the emotions joy, trust and anger which are used to measure brand love and to classify the brands into one of the four brand types.

Sentiment analysis is a method used to extract an author’s emotional intent from a piece of text with the use of natural language processing algorithms. It can be applied to any type of textual form and can be performed on a document level, sentence level or even word level, depending

on how detailed the analysis needs to be. One approach to sentiment analysis is to state the polarity of a document by using a polarity score, which indicates whether the text is positive, neutral or negative (Jurek et al., 2015). Beyond sentiment analysis on the polarity, analysis on the emotional states can be performed. One of these frameworks to classify emotions was created by Robert Plutchik, who believed that there are eight primary emotions: anger, fear, sadness, disgust, surprise, anticipation, trust and joy (Plutchik, 1980). This study uses Plutchik's framework to detect the emotions joy, trust and anger in the brand related content.

Two main approaches for performing sentiment analysis are machine learning methods, which are supervised methods and methods that rely on dictionaries, unsupervised methods.

The machine-learning sentiment analysis approach uses pre-classified texts for training a classifier. In this pre-classified training set, each text is labelled to a pre-set class. Then, with this trained classifier, unseen texts are classified into the classes positive, negative or neutral or in different classes of emotions (Kolchyna et al., 2015). An advantage of this method is that it does not rely on a dictionary and it often has higher accuracy in prediction than opposed to the lexicon-based method (Jurek et al., 2015). However, finding a pre-labelled training set for the machine learning analysis is a challenging task.

The lexicon-based method relies on a dictionary of words, with a score assigned to each word or sentence to indicate the valence or emotion. As mentioned, a main advantage of this method is that it does not require training of the model. However, one of the criticisms on the lexicon-based approach is that the dictionaries may be unreliable and may not include all words necessary (Taboada et al., 2011). Nevertheless, since an unlabelled dataset is used in this research, the lexicon-based method is preferred.

This research uses the NRC sentiment dictionary set up by Mohammad & Turney (2018) which includes the eight primary emotions of Plutchik's wheel of emotion. Sentiment analysis in this study is performed on word-level. Thus, tweet text from the brand related content posted by consumers is represented as a bag-of-words. In a bag-of-words, the order of the words and grammatical rules are ignored. Then, every word in the data frame is scored based on the emotions in the NRC dictionary. If a word corresponds to the emotion of interest, this word is given a 1. If the word does not correspond to the emotion of interest, a value of 0 is given (Jurek et al. 2015).

4.1.2 First dimension for measuring brand love: emotional engagement

The first dimension on which brand love is measured is emotional engagement. Emotional engagement is defined as the emotional connection that brands have with their consumers that results in brand love. Emotional engagement is measured by two criteria: the frequency of interaction that brands have with consumers and the extent to which consumers post content on the brands that contain the emotion joy, as measured by sentiment analysis.

The first criterium for emotional engagement is the frequency of interactions. This is measured by the number of replies a brand posts relative to the number of OSM tweets and retweets. A relationship between a brand and consumers develops based on the interaction between them. Relationships require frequent, interactive behaviour. Thus, in order for a consumer to love a brand, brands should have frequent interactions with them (Batra et al., 2012; Thomson et al., 2005).

The second criterium for measuring emotional engagement is the extent to which consumers post content on brands that elicit the emotion joy. Consumers describe loved brands in positive terms (Batra et al., 2012; Vernuccio et al., 2015). Furthermore, loved brands provoke a strong feeling of pleasure and arousal. Thus, brands that have more emotional engagement with consumers evoke a positive emotion, pleasure and arousal. Three positive emotions that evoke pleasure and arousal are amusement, happiness and joy (Berger & Milkman, 2012; Bottenberg, 1975; Russell & Mehrabian, 1977). However, this study uses the NRC sentiment and emotion dictionary, which only contains eight primary emotions: anger, fear, anticipation, trust, surprise, sadness, joy and disgust (Mohammad & Turney, 2018). Other secondary and tertiary emotions, such as amusement and happiness, are not present in the NRC lexicon but have been categorized among these primary emotions in previous research. Both amusement and happiness have been categorized as part of the primary emotion joy, which is hence the second measure for emotional engagement (Shaver et al., 1987).

4.1.3 Second dimension for measuring brand love: respect

The second dimension is respect. Respect represents perceptions of the performance and the more functional attributes of a brand. It is defined as the positive attitude that consumers have towards a certain brand based on the evaluation of brand performance, rather than based on emotions or feelings. Furthermore, respect is the only emotion related to cognition. Respect is

measured by the degree to which consumers post tweets that contain the emotions trust and anger.

The first measure for respect is the extent to which consumers elicit the emotion trust in their posts on the brands. Existing literature states that respect is, among other things, reflected by brand performance and the trust that consumers have in a brand (Pawle & Cooper, 2006; Roberts, 2005). Furthermore, Elliott & Yannopoulou (2007) see trust as a critical component for a brand to build a loving relationship with consumers.

The second measure for the dimension respect is the degree to which consumers post angry content on the brand. Consumers lose respect in a brand when they are angry or upset (Elliott & Yannopoulou, 2007). Furthermore, consumers are least trusting, and thus respecting, when they are angry (Dunn & Schweitzer, 2005)

4.1.4 The four brand types

These two dimensions with the corresponding measures are used to classify the brand into one of the four brand types (Figure 4.1).

The first brand type, Lovemarks, are brands that are loved most by consumers. Lovemarks score high on emotional engagement and on respect. Therefore, they are characterized by frequent interaction with their consumers and a high level of the emotions joy and trust. These brands do not evoke anger among consumers.

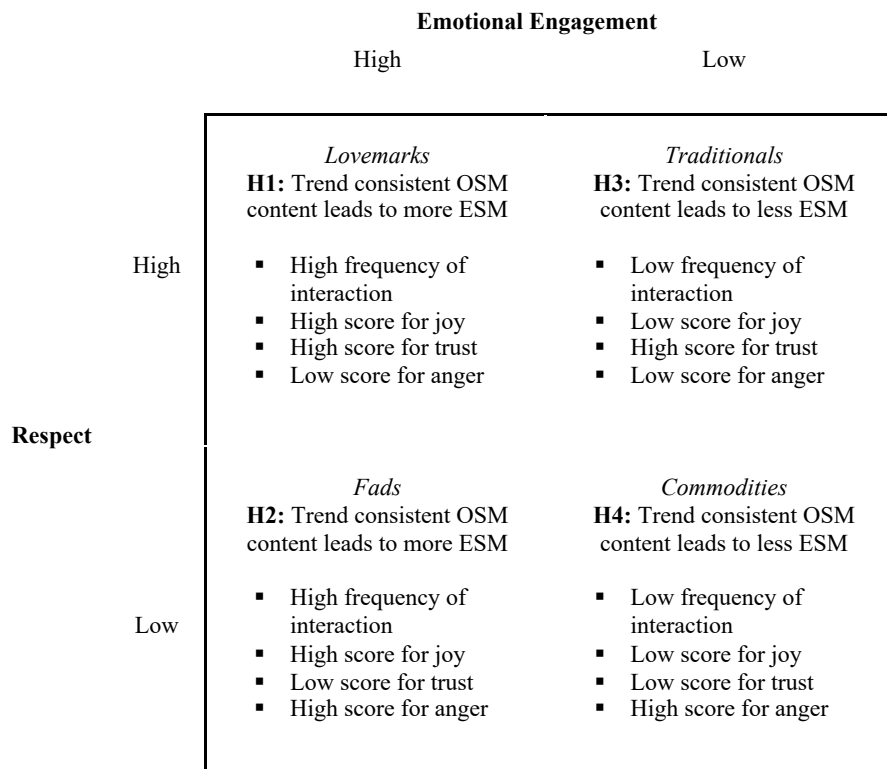
Fads, which are loved but short lived brands, score high on emotional engagement but low on respect. These brands interact frequently with their consumers and consumers post tweets on these brands that contain the emotion joy. However, these brands are not respected by consumers and thus they have no trust in the brands. Lastly, consumers talk angrily about these brands more often than they do about Lovemarks.

Traditionals, which are brands that are not loved but are respected, score high on the dimension respect but not on the dimension emotional engagement. This implies that these brands do not interact frequently with consumers nor do consumers tweet posts that contain the emotion joy. Consumers do post tweets that contain the emotion trust, as the brands are respected, and they do not talk angrily about these brands.

Lastly, brands that score low on respect and low on emotional engagement are the Commodities. Characteristics of this brand type are that consumers post tweets on these brands that do not contain the emotions joy and trust. Moreover, consumers are likely to talk angrily about these brands. Lastly, brands that are classified as Commodities do not interact with their consumers frequently.

Figure 4.1

Framework of different brand types



Four different types of brands with corresponding hypothesis and the four different measures for emotional engagement and respect.

4.2 Classifying brands with hierarchical clustering

Once sentiment analysis is performed, the brands can be classified into the four brand types. This is done by performing hierarchical clustering on the four measures for brand love: frequency of interaction, the emotion joy, the emotion trust and the emotion anger. The emotions joy, trust and anger are extracted from the results of the sentiment analysis.

Hierarchical clustering is an unsupervised learning method that is used for finding subgroups, or clusters, in an unlabelled dataset. The clustering methods seeks to split the dataset into distinct groups in such a way that observations within a group are similar to one another, whilst the observations between groups are rather different from each other (James et al., 2000).

Hierarchical clustering is preferred as opposed to other clustering methods, such as k-means clustering, since it works well with numerical data and on small datasets (Rusch, 2015). Furthermore, hierarchical clustering produces clusters of higher quality and offers a more natural way to organize real-world observations (Bouguettaya et al., 2015).

The method groups observations on the basis of their similarity. The clustering results are represented in a tree-based structure, which is called a dendrogram. This enables one to clearly interpret the clustering structure of hierarchical clustering. In hierarchical clustering, every observation starts in its own cluster. Then, each observation is paired with neighbouring observations or clusters one at the time, which continues until all clusters are linked to one another. Thus, the method does not consider the global structure of the data but makes decisions considering the local pattern. This method is referred to as the ‘bottom-up’ approach, or the agglomerative method. Another method is divisive clustering, which is a ‘top-down’ method, where all data points start in one cluster and are split into smaller clusters until all observations are its own cluster (Rusch, 2015). However, because clustering all the way up from individual data points identifies smaller clusters, agglomerative clustering is performed in this research. After the choice between agglomerative or divisive clustering, three other diagnostics have to be determined: the similarity measure, linkage method and the optimal number of clusters.

The similarity measure is used for the pairing of the neighbouring observations. Different similarity measures exist, which can be based on the correlation between observations or on the distance. When using correlation-based similarity measures, a linear relationship between the variables is required which is not applicable for the measures of brand love (Jain et al., 1999). Thus, a distance-based similarity measure is preferred. The measures anger and joy have some outliers, so the Manhattan distance measure is used in this research, which is more robust and less sensitive to outliers (Fowlkes & Mallows, 1983). The Manhattan distance (d) is calculated by taking the absolute difference between the coordinates (x, y) of the individual data objects (i), as noted in Formula 1.

$$d(x, y) = \sum_i |x_i - y_i| \quad (1)$$

Once individual observations have been paired into clusters on the basis of the Manhattan distance, these clusters need to be linked to other clusters. For linking these clusters, the

Manhattan distance does not suffice, and another measure should be used, leading to the next diagnostic that is decided on, the linkage method. This is the measure that is used for pairing groups of observations. Linkage measures merge clusters on the basis of the distance between them. Different linkage methods exist, among which single, complete, average and Ward's, which differ in characterizing the similarity between clusters. The single linkage method merges clusters on the basis of minimal inter-cluster dissimilarity. This means that the distance between two clusters is the minimum of the distance of all pairs of observations in the two clusters. However, this method is not suitable for this research as it performs poorly when there is noise between clusters and there are outliers in the data (Patel et al., 2015). The complete linkage method computes pairwise dissimilarities on the basis of the maximum distance of all pairs of observations between two clusters. It produces more compact and more balanced clusters. However, it is biased towards globular clusters (Jain et al., 1999). The average clustering method computes all pairwise distances between observations in the two clusters and takes the average to determine the distance between the clusters. However, similar to the single linkage method, the average linkage method does not perform well if there is noise between clusters (James et al., 2000). The last linkage method discussed is the Ward's method, which is used in this research. The method minimizes the total cluster variance. In other words, the distance between two clusters is determined by how much the sum of squares will increase once the clusters are merged. The distance between clusters is calculated with the use of the squared Euclidean distance. This method does well if there is noise between clusters and produces compact, even-sized clusters (Ward, 1963).

The last decision is the optimal number of clusters. This can be determined with the average silhouette method. This measure determines the quality of a cluster by measuring the distance between the resulting clusters. This is done by comparing the mean intra cluster distance, $a(i)$, which is the average Euclidean distance between observations within a cluster, to the mean distance to the nearest cluster, $b(i)$ (Formula 2). The mean distance to the nearest cluster is the average Euclidean distance from the data points of the cluster to data points from other clusters. The result is a silhouette coefficient, $s(i)$, which can range from -1 to 1.

$$s(i) = \frac{b(i) - a(i)}{\max(a(i), b(i))} \quad (2)$$

A silhouette coefficient close to 1 means that it is likely that the method clustered the observations in the correct cluster. On the contrary, if the coefficient is close to -1, this means that the observation is probably in the wrong cluster. Thus, a higher silhouette coefficient indicates a higher cluster quality. The optimal number of clusters is therefore obtained at the maximum silhouette coefficient.

4.3 Measuring owned social media

After the brands have been classified into the four brand types, the OSM content has to be classified into trend or brand consistent content. Thus, a measure for OSM has to be defined. Existing literature has distinguished several dimensions in OSM. Berger & Milkman (2012) examine newspaper articles on two dimensions: the valence of the article and the specific emotions in the content. Dobeles et al. (2007) categorize owned media content in six primary emotions: surprise, joy, anger, fear, sadness and disgust.

In this study, OSM is classified in either brand consistent content or trend consistent content. This is measured by the topics present in the tweets as determined by LDA topic modelling. This means that if tweets are on topics that concern COVID-19 or BLM, these tweets are classified as trend consistent content. All other tweets are classified as brand consistent content.

4.4 Measuring trend consistent content with Latent Dirichlet Allocation

Latent Dirichlet Allocation (LDA) is a latent topic modelling method that is designed to analyse documents (Blei et al, 2002). LDA is an unsupervised learning method, looking for previously undetected topics present in text. In this research, it is used to classify OSM into brand consistent content or trend consistent content. Thus, LDA is performed on the OSM tweets that are posted by the brands, which are extracted from the second dataset.

LDA is a generative probabilistic model used to find clusters within documents, where a document is represented by a bag-of-words. The method first determines whether specific terms are part of a particular topic or not. Then, the method observes the words present in the text and assigns each text a probability of being part of a topic. Thus, LDA assigns topics to terms and, dependent on the frequency of terms present in the text, assigns these texts to the topics. The method is considered a soft-clustering method because a term can belong to multiple topics and a document can be about multiple topics (Blei et al., 2002).

The main idea is that documents are denoted as mixtures over latent topics, where every topic is represented by a distribution over words. First of all, the topic probability distribution is determined, which gives the share of each document belonging to a certain topic. The hidden topic distribution needs to be rendered in order to uncover the topic structure. Thus, LDA observes the probabilities over all words (β_k) in all documents and defines k number of topics. Every document is then given a probability for each topic (θ_n). Thus, first β_k needs to be generated for each topic and then for the given documents the topic probabilities θ_n can be determined. LDA does this by computing the posterior distribution of the hidden variables, which is the conditional distribution. LDA uses a Dirichlet distribution prior on the per topic word distribution as well as on the per document distribution. A Dirichlet distribution is a distribution that is defined over a vector of probabilities (Reisenbichler & Reutterer, 2019). In contrast to classical clustering methods, where membership of a cluster is a binary variable, every term belongs to all clusters but with different probabilities. The probabilities that a certain term belongs to topic k are expressed in a vector with the probabilities over all words that adds up to 1 (Formula 3). The shape of the prior per topic word distribution of β_k is governed by parameter δ . Similarly, every topic belongs partially to all documents (n) of which probabilities are also expressed in a vector (Formula 4). This vector contains the probabilities over all topics, which indicate the activeness of certain topics within a document. The shape of the prior per document distribution of θ_n is also governed by a parameter, α (Kwartler, 2017).

$$\beta_k \sim \text{Dirichlet}(\delta_1, \dots, \delta_k) \quad (3)$$

$$\theta_n \sim \text{Dirichlet}(\alpha_1, \dots, \alpha_k) \quad (4)$$

Once the topic distribution has generated the topic probabilities, these topic probabilities determine topics for each word (i) using a multinomial distribution. The multinomial distribution is a distribution that assigns a probability to discrete outcomes. Thus, the exact allocation of topics to the word can be performed, denoted by z_{in} in Formula 5. Then, for each and every word, a corresponding topic is assigned, denoted by w_{in} in Formula 6. Again, this is done using the multinomial distribution (Kwartler, 2017).

$$z_{in} \sim \text{Multinomial}(\theta_n) \quad (5)$$

$$w_{in} \sim \text{Multinomial}(\beta_{z_{in}}) \quad (6)$$

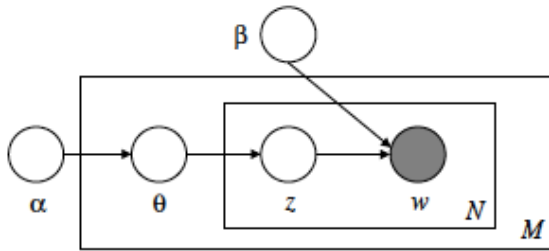
Once α and β are known, the joint distribution of the topic distribution θ_n , the set of N topics z and the set of N words w is given by:

$$p(\theta, z, w | \alpha, \beta) = p(\theta | \alpha) \prod_{n=1}^N p(z_n | \theta) p(w_n | z_n, \beta). \quad (7)$$

This process of LDA is visualized in Figure 4.2, where the outer box, M , represents the documents and the inner box, N , represents the repeated choice of topics and words within a document.

Figure 4.2

Visualization of the process of LDA



Representation of the process of the LDA model where M represents the documents and N represents the repeated choice of topics and words within a document.

LDA can be performed with different techniques. In this research, the Gibbs sampling method is used. This is a sampling method that generates draws out of the distribution. The final result is the average of those draws. This method is preferred because it works well on short documents, and thus it is expected that it works well on the tweets of the brands (Kwartler, 2017).

Before performing LDA on the Twitter data, two parameters have to be predetermined. First, the number of topics, k , should be determined. A higher value of k leads to a better fit to the data. However, the method can become computationally expensive, and a high value of k may lead to overfitting of the data. Thus, in making the choice of the number of topics, a trade-off should be made between model complexity and in-sample fit (Kwartler, 2017).

The optimal number of topics can be determined with the perplexity value, which evaluates the fit of the model to the word distribution of a corpus. The perplexity value is calculated by first

determining the log likelihood of the model. The log likelihood calculates the probability that the process described by the model fits the data that is actually observed. In other words, the log-likelihood determines the probability of the fit of the model of a set of unseen documents, given the topics and the topic-distribution as estimated by the LDA model. This means that the higher the log likelihood, the higher the probability of observing the actual data with the model (Heinrich, 2005). Then, the perplexity is calculated by dividing the log likelihood by the total number of words and taking the exponent of this (Formula 8). The maximum log likelihood determines the best model fit, which means that we want to minimize the perplexity value. Thus, a lower perplexity value indicates a better model fit (Wang et al., 2012).

$$Perplexity = \exp\left(-\frac{\text{loglikelihood}}{\text{Total number of words}}\right) \quad (8)$$

Secondly, the parameter α should be chosen, which controls the sparseness of the topic distribution via the variance. If α is set to a very small value, the variance will be rather large. Thus, a small α creates a very sparse topic distribution. If α is very large, this means that the variance of the topic probability is small and all realizations of the Dirichlet distribution will be quite similar to each other. This generates a very centred distribution around the probability equal to 0.5. An α equal to exactly 1 creates a uniform distribution, where all the shares are equally likely. In LDA modelling with Gibbs sampling the same α is assumed for every topic. Nevertheless, the size of this α has to be determined (Kwartler, 2017). Again, the optimal value is determined by the minimum perplexity value.

After tuning the parameters, the final LDA model is performed and the OSM tweets are classified into trend consistent content and brand consistent content.

4.5 *Measuring earned social media*

Now that the brands have been classified into the brand types and the OSM tweets into brand or trend consistent content, the effect of the different types of OSM on online consumer engagement for the different brand types can be examined. The observed output on social media network sites of online consumer engagement is ESM. Hence, consumer engagement is measured by looking at brand-related ESM. Earned media has been characterized by a couple of measures. Colicev et al. (2018) distinguish three dimensions in ESM: volume, valence and brand fan following. Volume captures the amount of earned media content that is created and shared for brands. The valence of ESM captures the positive and negative sentiment of the

social media content and total brand following which shows how popular a brand is among consumers on social media networks. This is in line with the research of Tirunillai & Tellis (2012), who use product reviews and product ratings as forms of user-generated content and characterize these by three metrics: volume, valence and ratings.

Based on the existing literature, ESM is measured by the number of retweets, which indicates how much the OSM is shared by consumers.

4.6 Linear regression to determine the effect of OSM on ESM

To examine the effect of brand and trend consistent content on ESM, a multiple linear regression analysis is performed. This is a supervised learning technique that predicts a linear relationship between the continuous dependent variable, ESM, and eight independent variables. Content type and brand type are the independent variables of interest. Furthermore, six control variables are added which are expected to influence ESM creation among consumers. First of all, brand specific control variables are added, which are the number of followers of the brand account, the percentage of OSM relative and the brand value. It is expected that the number of followers influences ESM creation positively and therefore the number of followers per brand was added (Colicev et al., 2018). The second control variable is the percentage of OSM. This variable shows how many of the tweets posted by a brand are OSM tweets, relative to the number of retweets and replies. Kumar et al. (2016) state that active social media communication increases favorable brand attitudes, which is expected to increase engagement in ESM. Thus, the higher the percentage of OSM, the higher the number of marketing content a brand has posted which is expected to positively influence ESM. Third, brand value is added to control for brand awareness that influences a consumer's intention to create ESM (Kim & Johnson, 2016). Moreover, three tweet-specific control variables are included: the length of the tweets, the number of favourites the tweet received and the day on which the tweet was posted. Tweet length is expected to affect ESM positively as well as the number of favourites that a tweet receives (Batra et al., 2012). To determine whether there are any time-varying effects, the day on which the tweet was posted is added as predictor variable. Last, because we are interested in the effect of trend consistent ESM on the different brand types, the interaction effect between content type and brand type is added (Formula 9).

$$ESM_i = \beta_0 + \beta_1 Brand\ type_i + \beta_2 Content\ type_{i3} + \beta_3 Content\ type * Brand\ type_i + \beta_4 Followers_i + \beta_5 Brandvalue_i + \beta_6 OSM_i + \beta_7 Favourite_i + \beta_8 Tweet\ length_i + \beta_9 Day_i + \epsilon \quad (9)$$

Where β_0 is the intercept with the y-axis and a constant term, and β_p is the slope coefficient for each independent variable that indicates the size and direction of the relationship. ϵ is the error of the model.

The regression fits the best linear relationship between the dependent and independent variable by minimizing the least squares, which are the distances between the actual observations and the predicted observations. The quality of the linear regression is assessed by the R^2 of the regression (James et al., 2000). This is a goodness-of-fit measure that indicates how well the model fits the data. It indicates the percentage of the variance in the dependent variable that can be explained by the independent variables in the model. This means that the larger the R^2 , the better the model fits the observations and the better the quality is of the model.

A linear regression is performed in this analysis for two reasons. First, because of its easy implementation and simplicity, it can occasionally outperform more advanced predictive modelling techniques (van der Heide et al., 2019). Furthermore, the output of a linear regression is very informative, as the coefficients give the direction of the relationship between the predictor and response variable along with the size of this relationship. Therefore, the linear regression is a suitable method for testing the four hypotheses formulated in this research. One of the disadvantages of a linear regression is that the linear regression relies on some strong assumptions, such as the absence of multicollinearity among the independent variables and a normal distribution of the data. Furthermore, the residuals should be homoscedastic, meaning that the variance of the residuals throughout the data should be similar, and the residuals should be independent. Lastly, linearity between the dependent and independent variables is assumed (Uyanık & Güler, 2013). If one of these assumptions is not met, performance of the linear regression will be suboptimal.

4.7 Random Forest to determine the importance of emotions on ESM

According to previous studies, sentiment and emotions, which influence brand love, are very important in stimulating ESM creation among consumers (Berger & Milkman, 2012; Lovett et

al., 2013). By performing a random forest and examining the variable importance, it is determined whether sentiment and emotions are indeed important predictors in predicting ESM as compared to seven other variables that have an effect on ESM. Thus, a random forest is performed on the emotions trust, anger and joy that were obtained in sentiment analysis and on seven predictor variables used in the linear regression: brand type, content type, followers, percentage of OSM, brand value, tweet length and the number of favourites.

Random forest is an ensemble method, which is a method that combines the predictions from multiple other methods to make more accurate predictions. Random forest provides an improvement over other ensemble methods, such as bagging, because the method introduces an additional force of randomness by decorrelating the trees. When building B decision trees on bootstrapped samples, the random forest introduces randomness by only considering m predictors at each split in the decision trees. This force of randomness prevents the random forest from being dominated by one strong predictor variable and reduces overfitting. The final prediction of the dependent variable is done by averaging the predictions of all B trees (James et al., 2000).

Before performing a random forest, three important decisions have to be made to maximize prediction performance. First, the splitting criteria used in each tree should be considered. The decision trees are created by using binary recursive partitioning, which means that each split depends on the previous splits. In regression problems, these splits can be made by minimizing the variance in each node or, in other words, minimizing the residual sum of squares (Formula 10). This is done by summing the squared difference between the predicted value (\hat{y}_i) and the actual value (y_i).

$$RSS = \sum_{i \in X_m} (y_i - \hat{y}_i)^2 \quad (10)$$

Secondly, the hyper parameter m should be tuned correctly, which controls the balance between decorrelation of the trees and the predictive strength of the random forest (Breiman, 2001). The hyper parameter is tuned by performing five-fold cross validation. In this method, the data is randomly divided into five folds, where each fold is used as a test set once. The four remaining folds are used as a training set. Then, the model is fit on the training folds, and this fitted model is used on the test fold. The best value for m is the value with the lowest cross validation error. The last decision that is considered before performing the random forest is the number of trees.

The number of trees is technically not a hyper parameter, but performance of the model depends highly on this choice. The outcome of the random forest becomes more stable the more trees there are, however, computations can get very expensive. Thus, choosing the number of trees is a trade-off between stable estimates and computational time (James et al., 2000). The number of trees is determined by five-fold cross validation. Thereafter, the random forest is performed. The performance of the random forest is assessed by examining the root mean square error (RMSE). The RMSE calculates the difference between the prediction and the actual target, squares these, takes the average and lastly takes the square root of this value. The lower the RMSE, the better the performance of the random forest (James et al., 2000).

After performing the random forest, the variable importance in predictions is analysed to determine whether the three emotions trust, joy and anger are important in predicting ESM. This is done by permuting one variable at a time and predict on the test set, with the permuted dataset. The mean squared error (MSE) of the original predictions and the permuted predictions are then compared. The predictor variables are ordered from most important variable to least important variable by the change in MSE. Thus, if permutation of a predictor variable results in a large change in MSE, the variable is considered an important variable (Breiman, 2001).

4.8 LIME to determine the effect of sentiment on predictions

The last method performed in this research is Local Interpretable Model-Agnostic Explanations (LIME) on the random forest. This method is performed to examine the effect of the three emotions trust, anger and joy on the predictions of ESM. LIME creates a new dataset that consists of permuted samples, after which it tests what happens to predictions with these changes in the dataset. The output of LIME is a new prediction for these permuted samples and the contribution of each feature to the prediction. Thus, by zooming in on a few observations, LIME provides local model interpretability from which the influence of features on the prediction can be obtained (James et al., 2000).

5 Results

5.1 Sentiment analysis on the brand related content

Sentiment analysis is performed on the brand related tweets posted by consumers. For every brand, a separate analysis is done to determine the score for the emotions joy, trust and anger that are present in the content posted by consumers per brand. Thus, from the second data frame,

a bag-of-words per brand is created. From this bag-of-words, sentiment analysis is performed on word-level, where a score for the three emotions is given per word on the basis of a 1 when the emotion conforms with the word and a 0 otherwise. Per brand, all the scores per emotion are added up together and divided by the number of words that are collected in order to get the average sentiment scores per brand. Thus, for every brand, one average sentiment score for the three emotions joy, trust and anger is obtained. Table 5.1 contains the minimum, mean and maximum score per emotion. These statistics show that consumers mostly feel the emotion trust towards the brand. The emotion joy has the widest range. Furthermore, the mean average score for anger is the lowest of all three emotions. This implies that, on average, consumers express the emotion anger the least in their Twitter content on the brands compared to the other two emotions.

Table 5.1

Descriptive statistics on the emotions that are present in brand related content posted by consumers

	Joy	Anger	Trust
Minimum	0.54	0.66	1.23
Mean	4.42	2.83	6.30
Maximum	11.50	10.07	10.59

The minimum, mean and maximum score for the emotions joy, anger and trust present in brand related content

5.2 Hierarchical clustering results for classification of the brands

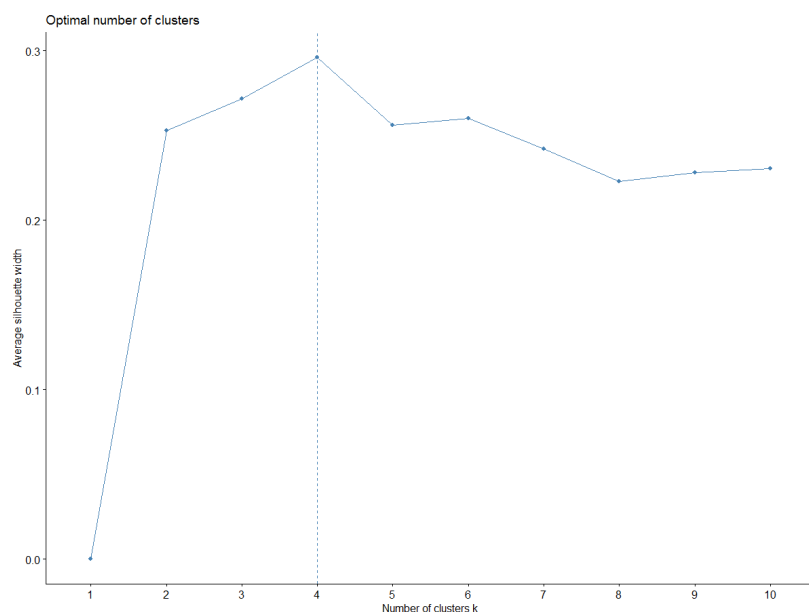
After performing sentiment analysis, hierarchical clustering is performed on the sentiment scores that are obtained in the sentiment analysis and an additional measure, frequency of interaction. The sentiment score for joy and the frequency of interaction are used to measure the dimension emotional engagement. The sentiment scores for the emotions anger and trust are used to measure the dimension respect. Thus, the three sentiment scores joy, anger and trust are merged into a data frame together with the variable frequency of interaction. In this data frame, every row corresponds to one of the 82 brands. Then, hierarchical clustering is performed on these four measures, after which the brands are classified into one of the four brand types: Lovemarks, Fads, Traditionals and Commodities. All rows of the data frame are scaled to be able to compare the four measures in the cluster analysis.

As mentioned before, the data contains outliers for the emotions anger and joy, and there is no linear relationship between most variables (Figures 1-4, Appendix B). Therefore, hierarchical

clustering with the Manhattan distance measure and Ward's linkage method is used. The optimal number of clusters is determined on the basis of the Silhouette plot. Based on the Silhouette plot, six clusters is determined to be the optimal number of clusters and the dendrogram is cut off at k equal to six (Figure 5.1; Figure 5, Appendix B).

Figure 5.1

Silhouette plot



Silhouette plot to determine the optimal number of clusters. The optimal number of clusters is at the maximum of the silhouette plot.

5.2.1 Classification of the brands on the basis of the hierarchical clustering results

For all six clusters obtained in the hierarchical clustering, the average values for the four measures are given (Table 5.2). Based on these values, the brands within a cluster are classified into one of the four brand types.

Table 5.2

Averages per cluster

Cluster	Percentage of replies	Joy	Anger	Trust	Size	Brand Type
1	14.87%	5.73	2.88	6.88	23	Lovemarks
2	85.57%	3.62	2.91	7.91	12	Lovemarks
3	83.40%	3.30	2.43	5.19	21	None
4	87.47%	1.53	7.56	2.84	3	Commodities
5	87.60%	7.07	2.63	6.88	10	Lovemarks
6	23.31%	3.29	2.41	5.91	13	Traditionals

Table with the averages per variable for each of the six clusters. This table contains the non-scaled values for the purpose of easier interpretation.

Cluster one has a fairly low share of frequency of interaction, only 14.87% of all tweets posted by these brands are replies to consumers. This implies that they do not interact with consumers on Twitter frequently but devote most of their tweets to OSM and retweets. However, the average value for joy in this cluster is very high compared to other clusters. Therefore, it may be concluded that this brand scores well on the dimension emotional engagement. Moreover, the average value for trust is also high and, compared to the values of joy and trust, the cluster scores low on the emotion anger. This means that this cluster also scores well on the dimension respect and it is best classified into the brand type Lovemarks.

The second cluster interacts with their consumers frequently as they have one of the highest values for percentage of replies. This cluster also has a high average value for the emotion joy and thus this cluster performs well on the dimension emotional engagement. Furthermore, cluster two has the highest average value for trust and compared to the values for joy and trust, a low average value for anger. Therefore, this cluster also performs well on the dimension respect and best fits the brand type Lovemarks.

When analysing the values for cluster three, none of the brand types can be identified. This is because cluster three has a high frequency of interaction and low values for the three emotions: joy, anger and trust. Thus, the values do not follow a pattern that fits any of the brand types and therefore this cluster is not used in further analyses.

Cluster four has a high value for frequency of interaction, 87.47% of the tweets posted by these brands are replies to consumers. However, this cluster has the lowest average value for joy. Furthermore, this cluster has a very high value for anger and a low value for trust. Thus, the high frequency of interaction may be explained by the fact that these brands have to reply to angry consumer very often in order to satisfy them. Therefore, these brands are classified into the brand type Commodities despite the high frequency of interaction.

The values for cluster five show a clear pattern for Lovemarks. This cluster has the highest frequency of interaction, 87.60% of the tweets are replies to consumers. Furthermore, this cluster has the highest average value for joy which means that the cluster scores high on the dimension emotional engagement. The cluster also scores high on the dimension respect because their average value for trust is high and their average value for anger is low.

Last, cluster six has a low percentage of replies, only 23.31%, which implies that their frequency of interaction is low. Moreover, they have the second lowest average value for joy which means that they do not score well on the dimension emotional engagement. This cluster also has a low average score for the emotion anger and compared to the scores for joy and anger a high average value for trust. Therefore, cluster six scores well on the dimension respect and is classified as Traditionals.

The final classification of the brands into the six clusters can be obtained in Table 1, Appendix B.

5.3 *LDA results for classification of owned social media content*

The OSM content posted by the brands is classified into trend consistent content and brand consistent content by performing LDA topic modelling. The OSM content is extracted from the second dataset, which is the dataset on the timeline of the brands. Then, the OSM tweets are converted to a Document-Term Matrix (DTM), where the rows correspond to the tweets and the columns to the individual words present in the tweets. Stemming is applied to the DTM, after which the DTM was split into a training set, containing 80% of the data, and a validation set with 20% of the data.

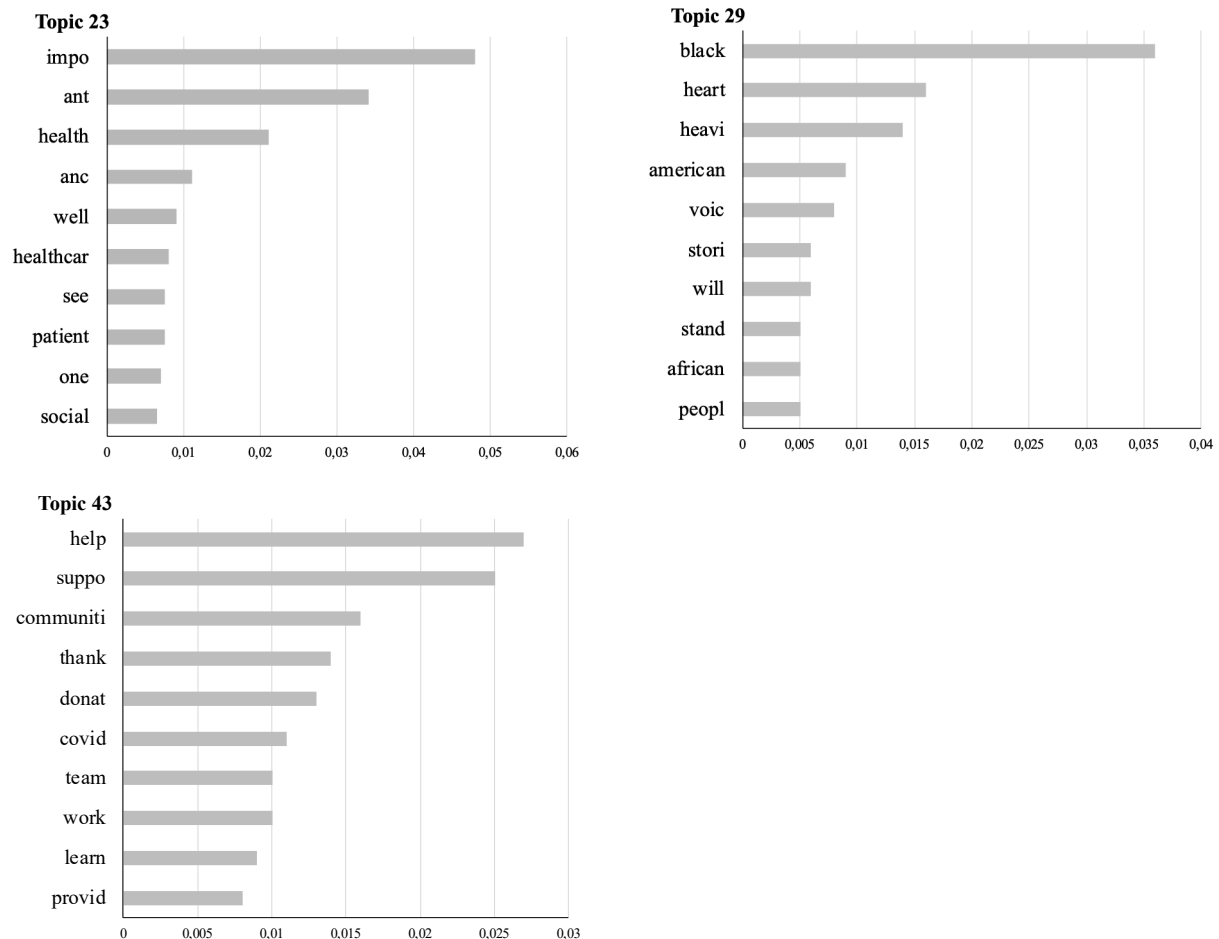
Before detecting the topics present in the tweets and running the LDA method on the data, the parameters k and α have to be optimized. This is done by running the LDA model on the validation set and determining the optimal number of topics on the basis of the perplexity value first. With this optimal number of topics, the model is performed with different values for α and the optimal α is determined on the basis of the perplexity. Then, with this optimal α , different values for k are tried again and vice versa until both parameters do not change anymore. After performing this method of tuning the parameters, the optimal value for k is 50 and the optimal value for α 0.1 (Figure 6 and 7, Appendix B).

From the results of the LDA model, the 50 topics present in the Twitter data are examined on the basis of the ten most important words. Out of the 50 topics, three of the topics can be categorized as trend consistent topics (Figure 5.2). The other 47 topics are categorized as brand

consistent topics. Visualization of the ten most important words for all 50 topics can be obtained in Figure 8, Appendix B.

Figure 5.2

Topics 23, 29 and 43 with the 10 corresponding terms



Visualization of the ten most important terms for topics 23, 29 and 43 concerning the COVID-19 pandemic and the BLM movement.

Topic 23 contains terms such as health, healthcare and patient. Furthermore, the world social is in the top ten words for this topic, which can refer to the new concept of social distancing that arose during the pandemic. Therefore, topic 23 is on COVID-19. Topic 29 includes the terms black, American, African and people and is thus on the BLM movement. Moreover, the terms heart, voice and stand are in line with the many protests and statements made by people and brands supporting the BLM movement. Topic 43 is also on the COVID-19 pandemic since the terms help, support, community, donation and covid are included in the ten most important words for this topic.

After determining which topics correspond to trend consistent matters, the most likely topic is determined for every tweet. This is the topic that has the highest probability of belonging to each tweet. All tweets that have topic 23, 29 or 43 as most likely topic are classified as trend consistent content and all other tweets as brand consistent content. In total, 4618 tweets are classified as trend consistent content. The majority of trend consistent content is on COVID-19 and approximately one quarter of the trend consistent content is on the BLM movement (Table 5.3).

Table 5.3

Number of tweets for the three trend consistent topics.

	Topic 23	Topic 29	Topic 43
Topic	COVID-19	BLM	COVID-19
Number of tweets	959	1138	2521

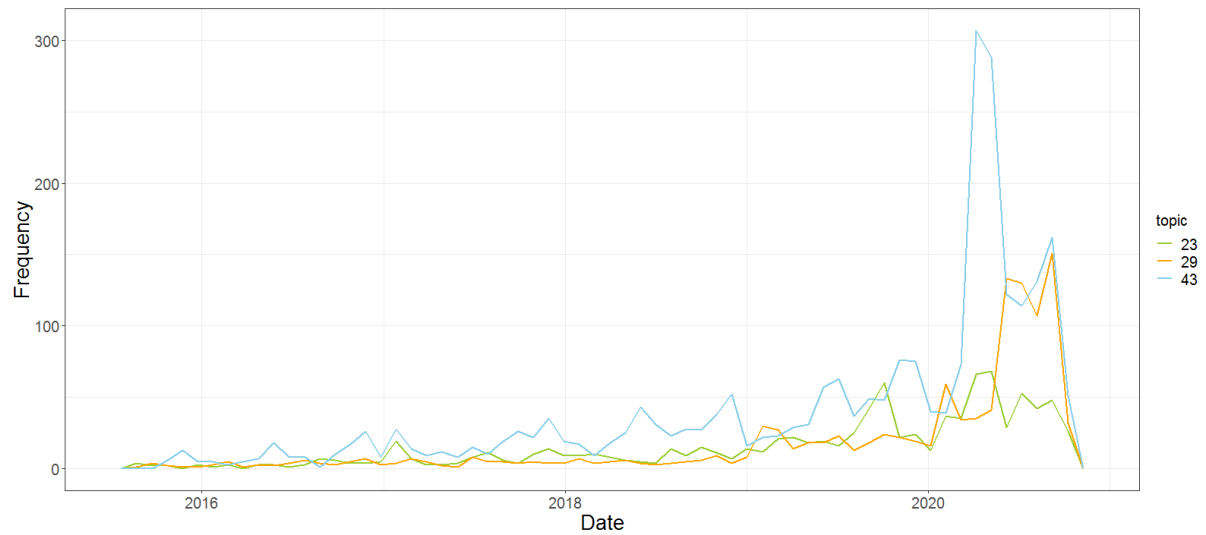
Topic 23, 29 and 43 with the corresponding trend consistent topic and number of tweets on this topic.

The frequency of trend consistent content posted over time can be observed in Figure 5.3. From this graph, it can be observed that brands have posted most of their trend consistent content in 2020. This is in line with expectations, as the two societal concerns went viral end 2019 and during 2020. The blue and green line show the frequency of posts on the COVID-19 virus. Both lines exhibit a peak at the end of 2019. This peak corresponds to the discovery of the virus. Moreover, a large peak can be observed in the posting frequency on the COVID-19 topic around March-April 2020. This corresponds to the start of the pandemic (World Health Organisation, 2020). The orange line in the graph visualizes the posting frequency for tweets on the BLM movement. In line with expectations, there is an increase in trend consistent posts on the BLM movement from May 2020 onwards, which is when the BLM movement went viral on social media.

Six examples of tweets on the trend consistent topics can be observed in Table 5.4. Not all tweets that are classified as trend consistent content are on the two topics. However, when analysing the classification, the majority of tweets are on COVID-19 and the BLM movement.

Figure 5.3

Frequency plot over time of the trend consistent content posted by brands



Visualization of when brands posted trend consistent content

Table 5.4*OSM content that is classified under trend consistent content*

Topic	Tweet	Brand	Date
23	With the rise in COVID-19 cases in the US, we're putting an alert at the top of @facebookapp and @instagram to remind everyone to wear face coverings and find more prevention tips from the CDC in our COVID-19 Information Center. #WearAMask https://t.co/ID29oACjoB	Facebook	02-07-2020
29	From black denim to vintage washes to matching sets! #HMMagainze has charted the front rows and compiled the three best denim styles sported by fashion month's it girls. Which look is your favourite? https://t.co/7DfWPwU01e https://t.co/tDcore3iL2	H&M	14-09-2018
29	To be silent is to be complicit. Black lives matter. We have a platform, and we have a duty to our Black members, employees, creators and talent to speak up.	Netflix	30-05-2020
29	"I am not going to stand up to show pride in a flag for a country that oppresses Black people and people of color." —Colin Kaepernick addressing the media, four years ago today, about his decision to sit for the national anthem https://t.co/QuD2eSkunF	ESPN	26-08-2020
43	Our team of volunteers travel around the world to help with non-profits. See how they're making an impact. https://t.co/5Ry9qTBNOW	JP Morgan	10-07-2017
43	Gucci supports One World #TogetherAtHome to celebrate #COVID19 frontline workers. #HealthForAll Watch the live show now: https://t.co/drB5370Els	Gucci	04-18-2020

Six examples of owned social media content that are classified under trend consistent content on the basis of the most likely topic of the tweet. For every tweet the topic, original tweet text, brand and date the tweet was posted are given.

5.4 Regression analysis results to determine the effect of OSM on ESM

After the brands have been classified into brand types by the clustering method and the OSM postings of the brands have been classified into brand or trend consistent content with the use of LDA, a regression analysis is performed to test the hypotheses.

5.4.1 Variables included in the linear regression

For the linear regression analysis, an additional dataset is composed. This dataset is on tweet level, which means that every row corresponds to an OSM tweet of one of the brands. The tweets of the brands in cluster 3 are omitted as this cluster could not be classified into one of the brand types. This results in a dataset with a total of 67,136 tweets. The overview of the variables in this dataset is given in Table 5.5.

Table 5.5

Overview of the variables used in the regression analysis

Variable	Data type	Variable type	Explanation	Source
Retweets	Numerical	Response	Measure for ESM	Twitter API
Brand type	Categorical	Predictor	The four brand types: 1) Lovemarks 2) Fads 3) Traditionals 4) Commodities	Sentiment analysis and clustering analysis
Content type	Categorical	Predictor	Content type of the tweet is 1) trend consistent or 2) brand consistent	LDA analysis
Followers	Numerical	Control	Number of followers of a brand	Twitter API
Percentage of OSM	Numerical	Control	Percentage of OSM posted relative to the retweets and replies	Number of OSM tweets over total tweets posted by a brand
Brand value	Numerical	Control	The brand value in USD	Forbes
Tweet length	Numerical	Control	Number of characters used in a tweet	Twitter API
Favourites	Numerical	Control	Number of favourites a tweet received	Twitter API
Day	Date	Control	Day on which the tweet was posted	Twitter API

The variables used in the linear regression. For every variable, the data type, variable type, a short explanation and the source for the variable is given.

For every tweet, the number of retweets is given, which is the measure of ESM and thus the response variable in the regression analysis. Two predictor variables of interest are added to the dataset, which are the content type of the tweet and the brand type. Furthermore, as mentioned previously, six control variables are added because it is believed that these variables have an effect on ESM. These six control variables are the number of followers, the percentage of OSM out of the total tweets that a brand has posted, the brand value, tweet length, number of favourites and the day on which the tweet was posted.

5.4.2 Descriptive statistics on the variables used in the linear regression

Out of the 67,136 tweets, only 3965 tweets are on trend consistent topics. This is in line with expectations since not all brands post trend consistent content. Furthermore, the brands that do

post trend consistent content do not post on these topics all the time. The timeframe of the dataset is five years, from August 2015, to October 2020. Brands have an average of 2,766,015 followers and, on average, post more replies and retweets than OSM content (Table 5.6). The brand value ranges from 8 million USD to around 160 million USD. Moreover, a tweet receives more favourites on average than retweets. This may be because favourites are only saved on a user's Twitter profile whereas a retweet is shared with all the followers of this user. Thus, favouriting may have a lower threshold for consumers than retweeting a tweet leading to more favourites than retweets.

Table 5.6

Descriptive statistics on the data frame containing all OSM tweets

Variable	Minimum	Mean	Maximum
ESM	0	204	216,381
Followers	46,868	2,766,015	35,912,234
Percentage of OSM	0.25	42.62	99.12
Brand value	7.90	22.16	162.90
Tweet length	5	131	280
Favourites	0	959	1,087,873

The statistics on ESM, the number of followers, the percentage of OSM posted, the brand value in billion USD, the number of favourites and the tweet length.

Besides analysing descriptive statistics on the entire data frame, the descriptive statistics per brand type are analysed to get a first impression of the characteristics of the different brand types. Out of the tweets posted by Lovemarks, 6.24% is trend consistent content. For Traditionals this percentage is 5.06% and for Commodities 7.46%. Thus, Commodities post on trend consistent most as compared to the other brand types. The brand type that gets the most ESM on average are the Traditionals, with 358 retweets per tweet on average (Table 5.7). Traditionals also receive most favourites on their tweets. These high levels may be explained by Traditionals posting the longest tweets, the most OSM and have the highest average number of followers as compared to Lovemarks and Commodities.

Table 5.7*Descriptive statistics per brand type*

Variable	Minimum	Mean	Maximum
Lovemarks			
ESM	0	141	119,111
Followers	46,868	2,543,926	35,912,234
Percentage of OSM	1.41	39.83	99.12
Brand value	7.90	22.86	162.90
Tweet length	15	134	280
Favourites	0	769	236,837
Traditionals			
ESM	0	358	216,381
Followers	161,137	2,818,018	12,233,722
Percentage of OSM	24.45	59.50	96.44
Brand value	8.30	17.38	50.50
Tweet length	5	122	280
Favourites	0	1,522	1,087,873
Commodities			
ESM	1	57	7,831
Followers	447,234	1,673,677	3,669,310
Percentage of OSM	0.25	11.46	23.84
Brand value	13.80	32.40	46.10
Tweet length	12	154	280
Favourites	3	479	33,165

The statistics on ESM, the number of followers, the percentage of OSM posted, the brand value in billion USD, the number of favourites and the tweet length.

5.4.3 Linear regression results

Before performing a linear regression, the log of all numerical variables is taken in order to reduce skewness in the data. A linear regression model including all variables is performed. For the two categorical variables, content type and brand type, reference levels are determined. Since we are interested in the effect of trend consistent content as opposed to the base level, brand consistent content, the reference level for content type is set to brand consistent content. For the variable brand type, the brand type Traditionals is set as a reference level. This is because it is expected that trend consistent content leads to more ESM for Lovemarks as opposed to Traditionals whereas for Commodities it is expected to lead to less ESM as opposed to Traditionals. Thus, the brand type Traditionals serves as a reference category as opposed to the other two brand types. Lastly, the interaction effect of content type and brand type was

included in the regression. The regression has an R-squared of 0.9017, which means that 90.17% of the data fits the model. Results of the regression can be obtained in Table 5.8.

Table 5.8

Results of the linear regression model

Variable	Coefficient
Trend content	0.003 *
Brand type Lovemarks	-0.005 ***
Brand type Commodities	-0.046 ***
Followers	-0.008 ***
Percentage of OSM	0.086 ***
Brand value	0.117 ***
Tweet length	0.003 ***
Favourites	0.856 ***
Day	-3.427 e-04 ***
Trend content * Lovemarks	0.005 **
Trend content * Commodities	0.007

Coefficients of the predictor variables included in the linear regression.

**** indicates significance at the 1% level*

*** indicates significance at the 5% level*

** indicates significance at the 10% level*

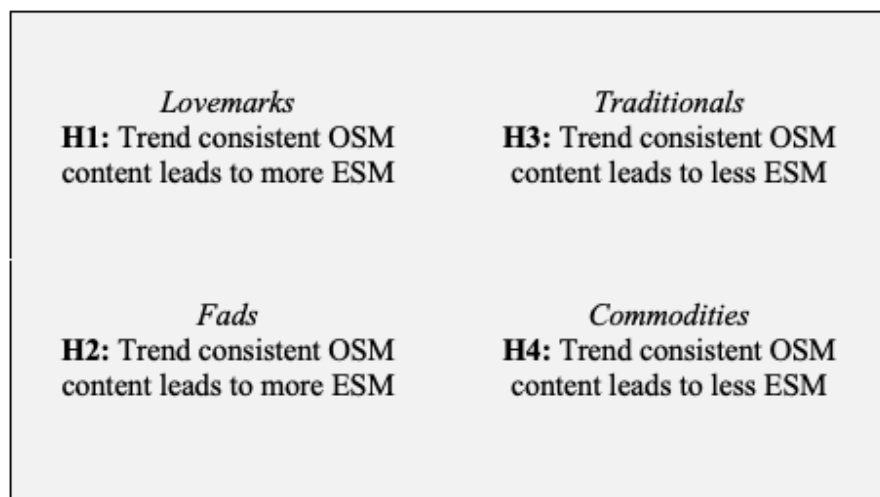
First of all, as can be obtained from the results, we generally find that the results indicate that there is a different relationship for trend consistent content for the different brand types. The effect of both the brand types, Lovemarks and Commodities, indicate a significant and negative relationship to ESM. However, the regression coefficient for the interaction effect between content type and brand type indicate different results. The positive coefficient for the interaction effect of trend content and Lovemarks implies that if Lovemarks increase their trend consistent content with 1 post as opposed to brand consistent content, their ESM increases with 0.5% more than that of Traditionals. However, the coefficient for trend type times Commodities is insignificant, and thus no conclusion can be drawn on the effect of trend consistent content for Commodities. Furthermore, against expectations, an increase in followers leads to a decrease in ESM. This means that if brands wish to improve their ESM, their focus should not be on increasing their followers base as, keeping all else constant, this leads to a decrease in ESM. They should rather focus on posting more OSM relative to replies and retweets and increasing their brand value, as both variables have a significant and positive coefficient. The coefficients for the number of favourites and tweet length are also significant and positive. The coefficient for the time trend variable, days, is significant and negative but very close to zero. This means that time has almost no effect on the number of retweets that a tweet receives.

5.4.4 Hypothesis testing

Based on the results from the regression analysis, the four hypotheses are tested (Figure 5.4).

Figure 5.4

Hypotheses for different brand types



The four different brand types with corresponding hypotheses.

The first hypothesis states that trend consistent content leads to more ESM for Lovemarks. From the regression coefficients, it can be concluded that Lovemarks receive more ESM than Traditionals when posting trend consistent content. Therefore, there is not enough evidence to reject the first hypothesis and it is concluded that trend consistent content leads to more ESM for Lovemarks.

The second hypothesis is on the effect of trend consistent content on Fads. No conclusion can be drawn for this hypothesis since Fads are not present in the data.

The third hypothesis states that trend consistent content leads to less ESM for Traditionals. The results from the regression show that Traditionals receive less ESM when posting trend consistent content than Lovemarks. Additionally, the effect of trend consistent content on Commodities as opposed to Traditionals is insignificant. Thus, there is not enough evidence to reject the third hypothesis. This leads to the conclusion that trend consistent content leads to less ESM for Traditionals.

Lastly, the fourth hypothesis on Commodities states that trend consistent content leads to less ESM. However, the interaction coefficient for trend content and Commodities is insignificant.

This implies that there is not enough evidence to reject the hypothesis which means that trend consistent content leads to less ESM for Commodities.

5.5 Random forest to analyse the importance of emotions on the prediction of ESM

Based on existing literature, it is expected that emotions play an important role in the creation of ESM. Furthermore, the four brand types build on the expectation that brand love, and thus emotions, play a key role in consumer engagement. To test these expectations, a random forest is performed on a dataset with ESM as response variable and ten predictor variables: trend content, brand type, followers, percentage of OSM, brand value, tweet length, favourites, joy, trust and anger. The dataset is split into a training set containing 70% of the data, a validation set with 15% of the data and a test set with 15% of the data.

Before training the random forest, two hyperparameters are tuned on the validation set: the value of the number of predictors, m , and the number of trees used in the random forest. Both values are tuned with the use of five-fold-cross validation. The optimal value of m is equal to 8 and the optimal value of the number of trees is equal to 250 trees. Then this random forest model is trained on the training set after which predictions were made on the test data. The RMSE of the random forest is 2096.64 retweets.

The variable importance is computed through the permutation method to test the importance of the emotions joy, trust and anger. All variables are ranked according to their mean increase in MSE (Figure 9, Appendix B). In line with expectations, the emotions trust, joy and anger are amongst the five most important predictor variables in the random forest. Trust is the second most important variable, joy the fourth and anger the fifth most important variable in predictions. This implies that the emotions that consumers feel towards a brand are of great importance in predicting ESM.

5.6 LIME to determine the effect of joy, trust and anger on ESM

The last method performed in this research is the LIME method. This method is used to determine the sign of the effect of the three emotions joy, trust and anger on ESM predictions. LIME is performed on five random predictions of the random forest. From the results, it can be observed that for all five observations, regardless of the type of brand or type of content, the emotion trust has the largest, positive effect on the prediction of ESM (Figure 10, Appendix B).

This means that if consumers feel more trust towards a brand, the ESM of this brand increases. The emotion joy also has a positive effect on all five ESM predictions and the emotion anger has a negative effect. This implies that joyful feelings among consumers towards a brand increase a brand's ESM. However, consumers talking negatively and angrily about the brand has a negative effect on ESM. These findings are in line with expectations, as more joy and trust lead to more brand love, which in return is expected to increase ESM (Batra et al., 2012; Thomson et al., 2005).

6 Conclusion

In this study, it is determined to what extent brand versus trend consistent content encourages consumers to engage in ESM. This is done by classifying social media content into brand and trend consistent content and distinguishing four brand types on the basis of brand love. I apply a linear regression analysis to determine the effect that trend consistent content has on the ESM of the different brand types. Furthermore, I establish whether emotions are indeed important in predicting ESM, as suggested by the existing literature.

The findings suggest that the extent to which consumers engage in ESM does depend on the social media type that brands post and, additionally, on the degree of brand love. From the regression results, it is concluded that for brands that have a strong, loving relationship with consumers, which are the Lovemarks, trend consistent content is beneficial in terms of ESM as opposed to posting brand consistent content. Thus, for Lovemarks, online consumer engagement is improved by posting on trend consistent topics. However, for Traditionals and Commodities posting trend consistent content has a negative effect on ESM. This means that if these brands post trend consistent content, consumers are less encouraged to engage in ESM than when these brands post brand consistent content.

To test the expectation that emotions in social media content have a strong influence on ESM, the variable importance in predictions of a random forest is examined. I find that the three emotions trust, joy and anger are among the five most important predictors. Furthermore, as expected by the literature on brand love, the emotions trust and joy have a positive effect on prediction of ESM and anger has a negative effect.

These findings imply that managers should not just follow social media trends and post on trend consistent topics as it does not lead to benefits for all brands. Dependent on the level of brand love that brands have among their consumers, brands should decide whether to post trend

consistent content or not. Managers can assess the level of brand love on the basis of the frequency of interaction they have with consumers, whether consumers feel joy and trust and whether they have to deal with many angry consumers or not. If a manager observes that they interact frequently with their consumers, consumers feel joy and trust and the brand gets little complaints from angry customers, the brand can classify itself as a loved brand. Managers of loved brands can improve their online consumer engagement by jumping on trend consistent topics in their social media content. However, if a brand does not interact frequently with their consumers and consumers do not have positive opinions about a brand, but rather post social media posts containing anger and complaints, there is no loving relationship between the brand and consumers. If a brand is not loved by consumers, the manager should stick to brand consistent content to first increase brand love. Only once a strong and loved brand is created, these brands can improve consumer engagement by posting on trend consistent topics.

The current study gives a good basis for determining when posting trend consistent content is beneficial for brands. However, the study has three limitations which give room for improvement by future research. First, there was a large variability in the timeframe of the OSM content posted by brands. The Twitter API only allows to retrieve the 3200 most recent posts. For some brands this covered a timeframe of five years, whereas other brands had posted 3200 tweets in only ten days. This led to performing analyses on different timeframes which can influence the results as tweets posted in 2020 were compared to tweets posted five years ago. Hence, future research can improve results by collecting tweets within one timeframe for all brands.

A second limitation in this study also concerns constraints in collecting Twitter data with the Twitter API. Because the Twitter API only allows to collect brand mentions within a timeframe of six to nine days, sentiment analysis was performed on tweets that were posted only within this short timeframe. Therefore, the classification of brands was performed on consumers' opinions within a short time frame, which may not represent the general consumer opinion on brands as opinions are due to change. Thus, future research should perform classification of the brands over a larger timeframe to get a better representation of the relationship between brands and consumers.

Lastly, this research only examines the volume of ESM by using the number of retweets. Results can be extended by also investigating the valence of ESM, as an increase in ESM does not necessarily indicate a positive attitude towards a brand (Lovett et al., 2013). Unfortunately,

the valence and emotions present in ESM could not be retrieved by the Twitter API, as it does not allow for extracting replies to OSM.

References

- Batra, R., Ahuvia, A., & Bagozzi, R. P. (2012). Brand love. *Journal of Marketing*, 76(2), 1–16. <https://doi.org/10.1509/jm.09.0339>
- Berger, J., & Milkman, K. L. (2012). What makes online content viral? *Journal of Marketing Research*, 49(2), 192–205. <https://doi.org/10.1509/jmr.10.0353>
- Blei, D. M., Ng, A. Y., & Jordan, M. T. (2002). Latent dirichlet allocation. *Advances in Neural Information Processing Systems*, 3, 993–1022.
- Bottenberg, E. H. (1975). *Phenomenological and operational characteriation OF FACTOR-ANALYTICALLY DERIVED DIMENSIONS OF EMOTION*. 1253–1254.
- Bouguettaya, A., Yu, Q., Liu, X., Zhou, X., & Song, A. (2015). Efficient agglomerative hierarchical clustering. *Expert Systems with Applications*. <https://doi.org/10.1016/j.eswa.2014.09.054>
- Breiman, L. (2001). Random forests. *Machine Learning*. <https://doi.org/10.1023/A:1010933404324>
- Brodie, R. J., Ilic, A., Juric, B., & Hollebeek, L. (2013). Consumer engagement in a virtual brand community: An exploratory analysis. *Journal of Business Research*, 66(1), 105–114. <https://doi.org/10.1016/j.jbusres.2011.07.029>
- Carroll, B. A., & Ahuvia, A. C. (2006). Some antecedents and outcomes of brand love. *Marketing Letters*, 17(2), 79–89. <https://doi.org/10.1007/s11002-006-4219-2>
- Chevalier, J. A., & Mayzlin, D. (2006). The effect of word of mouth on sales: Online book reviews. *Journal of Marketing Research*, 43(3), 345–354. <https://doi.org/10.1509/jmkr.43.3.345>
- Choi, B., & Lee, I. (2017). Trust in open versus closed social media: The relative influence of user- and marketer-generated content in social network services on customer trust. *Telematics and Informatics*, 34(5), 550–559. <https://doi.org/10.1016/j.tele.2016.11.005>
- Colicev, A., Malshe, A., Pauwels, K., & O'Connor, P. (2018). Improving consumer mindset metrics and shareholder value through social media: The different roles of owned and earned media. *Journal of Marketing*, 82(1), 37–56. <https://doi.org/10.1509/jm.16.0055>
- Daugherty, T., Eastin, M. S., & Bright, L. (2008). Exploring Consumer Motivations for Creating User-Generated Content. *Journal of Interactive Advertising*, 8(2), 16–25. <https://doi.org/10.1080/15252019.2008.10722139>
- Dhar, V., & Chang, E. A. (2009). Does Chatter Matter? The Impact of User-Generated Content on Music Sales. *Journal of Interactive Marketing*, 23(4), 300–307. <https://doi.org/10.1016/j.intmar.2009.07.004>
- Dobele, A., Lindgreen, A., Beverland, M., Vanhamme, J., & van Wijk, R. (2007). Why pass on viral messages? Because they connect emotionally. *Business Horizons*, 50(4), 291–304. <https://doi.org/10.1016/j.bushor.2007.01.004>
- Dunn, J. R., & Schweitzer, M. E. (2005). Feeling and believing: The influence of emotion on trust. *Journal of Personality and Social Psychology*, 88(5). <https://doi.org/10.1037/0022-3514.88.5.736>
- Elliott, R., & Yannopoulou, N. (2007). The nature of trust in brands: A psychosocial model. *European Journal of Marketing*, 41(9–10), 988–998. <https://doi.org/10.1108/03090560710773309>
- Fowlkes, E. B., & Mallows, C. L. (1983). A method for comparing two hierarchical clusterings. *Journal of the American Statistical Association*. <https://doi.org/10.1080/01621459.1983.10478008>
- Goh, K. Y., Heng, C. S., & Lin, Z. (2013). Social media brand community and consumer behavior: Quantifying the relative impact of user- and marketer-generated content. *Information Systems Research*, 24(1), 88–107. <https://doi.org/10.1287/isre.1120.0469>
- Heinrich, G. (2005). Parameter estimation for text analysis. *Bernoulli*.

- Hennig-Thurau, T., Gwinner, K. P., Walsh, G., & Gremler, D. D. (2004). Electronic word-of-mouth via consumer-opinion platforms: What motivates consumers to articulate themselves on the Internet? *Journal of Interactive Marketing*, 18(1), 38–52. <https://doi.org/10.1002/dir.10073>
- Jain, A. K., Murty, M. N., & Flynn, P. J. (1999). Data clustering: A review. *ACM Computing Surveys*. <https://doi.org/10.1145/331499.331504>
- James, G., Witten, D., Hastie, T., & Tibshirani, R. (2000). An introduction to Statistical Learning. In *Current medicinal chemistry*. <https://doi.org/10.1007/978-1-4614-7138-7>
- Jurek, A., Mulvenna, M. D., & Bi, Y. (2015). Improved lexicon-based sentiment analysis for social media analytics. *Security Informatics*, 4(1). <https://doi.org/10.1186/s13388-015-0024-x>
- Kay, M. J. (2006). Strong brands and corporate brands. *European Journal of Marketing*, 40(7–8), 742–760. <https://doi.org/10.1108/03090560610669973>
- Kim, A. J., & Johnson, K. K. P. (2016). Power of consumers using social media: Examining the influences of brand-related user-generated content on Facebook. *Computers in Human Behavior*, 58, 98–108. <https://doi.org/10.1016/j.chb.2015.12.047>
- Kivetz, R., & Simonson, I. (2000). The effects of incomplete information on consumer choice. *Journal of Marketing Research*, 41(2), 57–63. <https://doi.org/10.2501/JAR-41-2-57-63>
- Kolchyna, O., Souza, T. T. P., Treleaven, P., & Aste, T. (2015). *Twitter Sentiment Analysis: Lexicon Method, Machine Learning Method and Their Combination*. <http://arxiv.org/abs/1507.00955>
- Kumar, A., Bezawada, R., Rishika, R., Janakiraman, R., & Kannan, P. K. (2016). From social to sale: The effects of firm-generated content in social media on customer behavior. *Journal of Marketing*, 80(1), 7–25. <https://doi.org/10.1509/jm.14.0249>
- Kwartler, T. (2017). Text Mining in Practice with R. In *Text Mining in Practice with R*. <https://doi.org/10.1002/9781119282105>
- Laroche, M., Habibi, M. R., & Richard, M. O. (2013). To be or not to be in social media: How brand loyalty is affected by social media? *International Journal of Information Management*, 33(1), 76–82. <https://doi.org/10.1016/j.ijinfomgt.2012.07.003>
- Lovett, M. J., Renana, P., & Shachar, R. O. N. (2013). On brands and word of mouth. *Journal of Marketing Research*, 50(4), 427–444. <https://doi.org/10.1509/jmr.11.0458>
- Mohammad, S., & Turney, P. (2018). Emotions Evoked by Common Words and Phrases: Using Mechanical Turk to Create an Emotion Lexicon, Proceedings of the {NAACL}-HLT 2010 Workshop on Computational Approaches to Analysis and Generation of Emotion in Text. *Proceedings of the NAACL HLT 2010 Workshop on Computational Approaches to Analysis and Generation of Emotion in Text*, 14(June), 26–34. <http://www.wjh.harvard.edu/%0Ahttp://saifmohammad.com/WebPages/lexicons.html>
- Pak, A., & Paroubek, P. (2010). Twitter as a corpus for sentiment analysis and opinion mining. *Proceedings of the 7th International Conference on Language Resources and Evaluation, LREC 2010*, 1320–1326. <https://doi.org/10.17148/ijarcce.2016.51274>
- Patel, S., Sihmar, S., & Jatain, A. (2015). A study of hierarchical clustering algorithms. *2015 International Conference on Computing for Sustainable Global Development, INDIACom 2015*, 3(10), 537–541.
- Pawle, J., & Cooper, P. (2006). Measuring emotion - Lovemarks, the future beyond brands. *Journal of Advertising Research*, 46(1), 38–48. <https://doi.org/10.2501/S0021849906060053>
- PLUTCHIK, R. (1980). A GENERAL PSYCHOEVOLUTIONARY THEORY OF EMOTION. In *Theories of Emotion*. <https://doi.org/10.1016/b978-0-12-558701-3.50007-7>
- Reisenbichler, M., & Reutterer, T. (2019). Topic modeling in marketing: recent advances and

- research opportunities. *Journal of Business Economics*. <https://doi.org/10.1007/s11573-018-0915-7>
- Roberts, K. (2005). *Lovemarks: The future beyond brands*. Powerhouse Books.
- Rusch, T. (2015). R for Marketing Research and Analytics . *Journal of Statistical Software*. <https://doi.org/10.18637/jss.v067.b02>
- Russell, J. A., & Mehrabian, A. (1977). Evidence for a three-factor theory of emotions. *Journal of Research in Personality*, 11(3), 273–294. [https://doi.org/10.1016/0092-6566\(77\)90037-X](https://doi.org/10.1016/0092-6566(77)90037-X)
- Shaver, P., Schwartz, J., Kirson, D., & O'Connor, C. (1987). Emotion Knowledge: Further Exploration of a Prototype Approach. *Journal of Personality and Social Psychology*, 52(6). <https://doi.org/10.1037/0022-3514.52.6.1061>
- Smith, A. N., Fischer, E., & Yongjian, C. (2012). How Does Brand-related User-generated Content Differ across YouTube, Facebook, and Twitter? *Journal of Interactive Marketing*, 26(2), 102–113. <https://doi.org/10.1016/j.intmar.2012.01.002>
- Stephen, A. T., & Galak, J. (2012). The effects of traditional and social earned media on sales: A study of a microlending marketplace. *Journal of Marketing Research*, 49(5), 624–639. <https://doi.org/10.1509/jmr.09.0401>
- Taboada, M., Brooke, J., & Voll, K. (2011). *Lexicon-Based Methods for Sentiment Analysis*. December 2009.
- Thomson, M., MacInnis, D. J., & Park, C. W. (2005). The ties that bind: Measuring the strength of consumers' emotional attachments to brands. *Journal of Consumer Psychology*, 15(1), 77–91. https://doi.org/10.1207/s15327663jcp1501_10
- Tirunillai, S., & Tellis, G. J. (2012). Does chatter really matter? Dynamics of user-generated content and stock performance. *Marketing Science*, 31(2), 198–215. <https://doi.org/10.1287/mksc.1110.0682>
- Uyanık, G. K., & Güler, N. (2013). A Study on Multiple Linear Regression Analysis. *Procedia - Social and Behavioral Sciences*. <https://doi.org/10.1016/j.sbspro.2013.12.027>
- van der Heide, E. M. M., Veerkamp, R. F., van Pelt, M. L., Kamphuis, C., Athanasiadis, I., & Ducro, B. J. (2019). Comparing regression, naive Bayes, and random forest methods in the prediction of individual survival to second lactation in Holstein cattle. *Journal of Dairy Science*. <https://doi.org/10.3168/jds.2019-16295>
- van Doorn, J., Lemon, K. N., Mittal, V., Nass, S., Pick, D., Pirner, P., & Verhoef, P. C. (2010). Customer engagement behavior: Theoretical foundations and research directions. *Journal of Service Research*, 13(3), 253–266. <https://doi.org/10.1177/1094670510375599>
- Vernuccio, M., Pagani, M., Barbarossa, C., & Pastore, A. (2015). Antecedents of brand love in online network-based communities. A social identity perspective. *Journal of Product and Brand Management*, 24(7), 706–719. <https://doi.org/10.1108/JPBM-12-2014-0772>
- Wallace, E., Buil, I., & de Chernatony, L. (2014). Consumer engagement with self-expressive brands: Brand love and WOM outcomes. *Journal of Product and Brand Management*, 23(1), 33–42. <https://doi.org/10.1108/JPBM-06-2013-0326>
- Wang, Y., Agichtein, E., & Benzi, M. (2012). TM-LDA: Efficient online modeling of latent topic transitions in social media. *Proceedings of the ACM SIGKDD International Conference on Knowledge Discovery and Data Mining*, 123–131. <https://doi.org/10.1145/2339530.2339552>
- Ward, J. H. (1963). Hierarchical Grouping to Optimize an Objective Function. *Journal of the American Statistical Association*. <https://doi.org/10.1080/01621459.1963.10500845>
- Zhang, L., Ghosh, R., Dekhil, M., Hsu, M., & Liu, B. (2011). Combining lexicon-based and learning-based methods for twitter sentiment analysis. *HP Laboratories Technical Report*, 89.

Appendix A

Table 1

Selection of the 85 brands after omitting inactive brands

Brand	Industry	Number of tweets	% of organic tweets	% of retweets	% of replies	Time frame (days)
Accenture	Business Services	3169	70,4	20,07	9,53	1373
Adidas	Apparel	3199	5,1	1,22	93,69	1070
Adobe	Technology	3200	64,34	4	31,66	661
Allianz	Financial Services	3196	64,67	33,79	1,53	1452
Amazon	Technology	3199	8,47	2,66	88,87	557
AmericanExpress	Financial Services	3199	11,19	0,69	88,12	328
AT&T	Telecom	3200	10,28	1,28	88,44	55
AudiOfficial	Automotive	3200	23,84	1,94	74,22	451
AXA	Financial Services	3198	35,21	36,55	28,24	1637
Bank of America	Financial Services	3194	42,42	39,04	18,53	1799
BMW	Automotive	3200	23,53	2,91	73,56	419
BurgerKing	Restaurants	3200	1,78	0,03	98,19	111
Caterpillar	Heavy Equipment	3200	57,91	20,5	21,59	791
Chase	Financial Services	3199	48,95	8,85	42,2	1462
Chevrolet	Automotive	3200	5,69	1,41	92,91	216
Cisco	Technology	3196	74,59	14,89	10,51	600
Citi	Financial Services	3200	77,38	21,69	0,94	988
CocaCola	Beverage	3200	0,75	0,03	99,22	149
Colgate	Consumer Packaged Goods	3200	2,22	1,03	96,75	961
CVS	Retail	3200	1,41	0,59	98	98
Dell	Technology	3199	16,41	4,72	78,87	389
Deloitte	Business Services	3197	85,83	13,32	0,84	674
Disney	Leisure	3200	85,38	14,12	0,5	778
eBay	Technology	3198	65,2	10,29	24,52	1305
ESPN	Media	3200	69,06	30,63	0,31	119
EY	Business Services	3200	63,25	26,53	10,22	552
Facebook	Technology	3199	9,22	5	85,78	698
FedEx	Transportation	3200	11,69	1,41	86,91	555
Ford	Automotive	3200	5,59	0,66	93,75	216
Fox	Media	3179	48,13	26,86	25,01	374
Frito-Lay	Consumer Packaged Goods	3200	4,53	0,75	94,72	491
GeneralElectric	Diversified	3185	43,99	3,86	52,15	1570
Gillette	Consumer Packaged Goods	3200	7,41	0,62	91,97	573
GoldmanSachs	Financial Services	3200	92,62	7,34	0,03	607
Google	Technology	3200	2,97	1,81	95,22	44
Gucci	Luxury	3200	99,12	0,34	0,53	947
H&M	Retail	3199	74,84	1,44	23,73	1264
Home Depot	Retail	3200	73	2,88	24,12	1484
Honda	Automotive	3200	26,75	2,53	70,72	863
HP	Technology	3197	33,75	3,69	62,56	1804
Huawei	Technology	3200	18,62	7,78	73,59	107

Hyundai	Automotive	3200	14,62	0,91	84,47	348
IBM	Technology	3200	6,5	4,06	89,44	296
IKEA	Retail	3200	74,56	2,25	23,19	1621
Intel	Technology	3200	35,72	9,56	54,72	1061
J.P. Morgan	Financial Services	3199	82,96	13,44	3,59	1481
JohnDeere	Heavy Equipment	3200	30,88	3,06	66,06	853
KFC	Restaurants	3200	0,44	0,03	99,53	41
L'Oréal	Consumer Packaged Goods	3190	61,16	26,24	12,6	1162
Lancôme	Consumer Packaged Goods	3190	72,04	14,01	13,95	1532
LEGO	Leisure	3200	12,31	2	85,69	196
Lexus	Automotive	3200	26,06	1,75	72,19	917
LouisVuitton	Luxury	3198	98,87	0,59	0,53	1623
Lowes	Retail	3200	1,44	0,44	98,12	51
Mastercard	Financial Services	3199	18,32	18,88	62,8	598
McDonalds	Restaurants	3200	0,25	0	99,75	10
MercedesBenz	Automotive	3200	17,06	2,12	80,81	117
Microsoft	Technology	3198	45,37	26,58	28,05	514
Nescafe	Beverage	3200	7,16	0	92,84	1257
Nestle	Consumer Packaged Goods	3196	45,12	13,02	41,86	802
Netflix	Technology	3149	24,45	45	30,55	237
Nike	Apparel	3200	2,94	0,31	96,75	919
Nintendo	Technology	3198	73,89	25,58	0,53	594
Oracle	Technology	3199	59,58	27,26	13,16	738
Pampers	Consumer Packaged Goods	3200	1,22	0,06	98,72	542
PayPal	Technology	3198	43,65	8,35	48	1811
Pepsi	Beverage	3200	4,03	0,12	95,84	240
PwC	Business Services	3200	96,44	3,5	0,06	839
Red Bull	Beverage	3175	21,61	10,39	68	937
Samsung	Technology	3200	49,16	12,53	38,31	1854
Santander	Financial Services	3200	77,94	16,28	5,78	1080
SAP	Technology	3200	69,44	2,81	27,75	490
Siemens	Diversified	3200	38,66	31,37	29,97	1031
Sony	Technology	3151	25,01	40,62	34,37	475
Starbucks	Restaurants	3200	0,94	0,62	98,44	38
T-Mobile	Telecom	3200	11,16	2,12	86,72	274
Toyota	Automotive	3200	2,06	0,16	97,78	54
UBS	Financial Services	3199	86,25	9,47	4,28	1228
UniqloUSA	Apparel	3199	62,36	4,03	33,6	933
UPS	Transportation	3200	4,19	1,22	94,59	124
Verizon	Telecom	3198	13,95	19,32	66,73	139
Volkswagen	Automotive	3200	15,41	2,66	81,94	933
Walmart	Retail	3200	1,97	0,28	97,75	29
Wells Fargo	Financial Services	3194	14,43	11,58	73,98	412
Zara	Retail	3200	6,53	0	93,47	237

Selection of brands with the industry, number of tweets, the percentage of replies, organic tweet and retweets out of the total number of tweets posted by a brand and the timeframe.

Table 2

Final selection of the 82 brands for retrieving the mentions posted by users.

Brands	Industry	Number of mentions
Accenture	Business Services	1957
Adidas	Apparel	3000
Adobe	Technology	3000
Allianz	Financial Services	757
Amazon	Technology	3000
AmericanExpress	Financial Services	1022
AT&T	Telecom	3000
AudiOfficial	Automotive	621
AXA	Financial Services	3000
Bank of America	Financial Services	3000
BMW	Automotive	3000
BurgerKing	Restaurants	3000
Caterpillar	Heavy Equipment	371
Chase	Financial Services	2787
Chevrolet	Automotive	3000
Cisco	Technology	3000
Citi	Financial Services	1668
CocaCola	Beverage	3000
Colgate	Consumer Packaged Goods	941
CVS	Retail	2135
Dell	Technology	1913
Deloitte	Business Services	2574
Disney	Leisure	3000
eBay	Technology	3000
ESPN	Media	3000
EY	Business Services	1056
Facebook	Technology	3000
FedEx	Transportation	3000
Ford	Automotive	3000
Fox	Media	3000
Frito-Lay	Consumer Packaged Goods	497
GeneralElectric	Diversified	483
Gillette	Consumer Packaged Goods	635
GoldmanSachs	Financial Services	2038
Google	Technology	3000
Gucci	Luxury	3000
H&M	Retail	1983
Home Depot	Retail	3000
Honda	Automotive	2642
HP	Technology	2270
Huawei	Technology	3000
Hyundai	Automotive	3000

IBM	Technology	3000
IKEA	Retail	1193
Intel	Technology	3000
JohnDeere	Heavy Equipment	1418
J.P. Morgan	Financial Services	3000
KFC	Restaurants	3000
LEGO	Leisure	3000
Lexus	Automotive	3000
L'Oréal	Consumer Packaged Goods	701
LouisVuitton	Luxury	3000
Lowes	Retail	3000
Mastercard	Financial Services	3000
McDonalds	Restaurants	3000
MercedesBenz	Automotive	3000
Microsoft	Technology	3000
Nescafe	Beverage	341
Nestle	Consumer Packaged Goods	3000
Netflix	Technology	3000
Nike	Apparel	3000
Nintendo	Technology	3000
Oracle	Technology	2545
PayPal	Technology	3000
Pepsi	Beverage	3000
PwC	Business Services	1703
Red Bull	Beverage	3000
Samsung	Technology	3000
SAP	Technology	3000
Siemens	Diversified	1848
Sony	Technology	3000
Starbucks	Restaurants	3000
T-Mobile	Telecom	3000
Toyota	Automotive	3000
UBS	Financial Services	754
UniqloUSA	Apparel	309
UPS	Transportation	3000
Verizon	Telecom	3000
Volkswagen	Automotive	1088
Walmart	Retail	3000
Wells Fargo	Financial Services	2492
Zara	Retail	935

Final selection of brands with the industry and the number of mentions from 23rd of October until the 2nd of November 2020.

Appendix B

Figure 1
Histogram for all values of the variable anger in which outliers can be detected

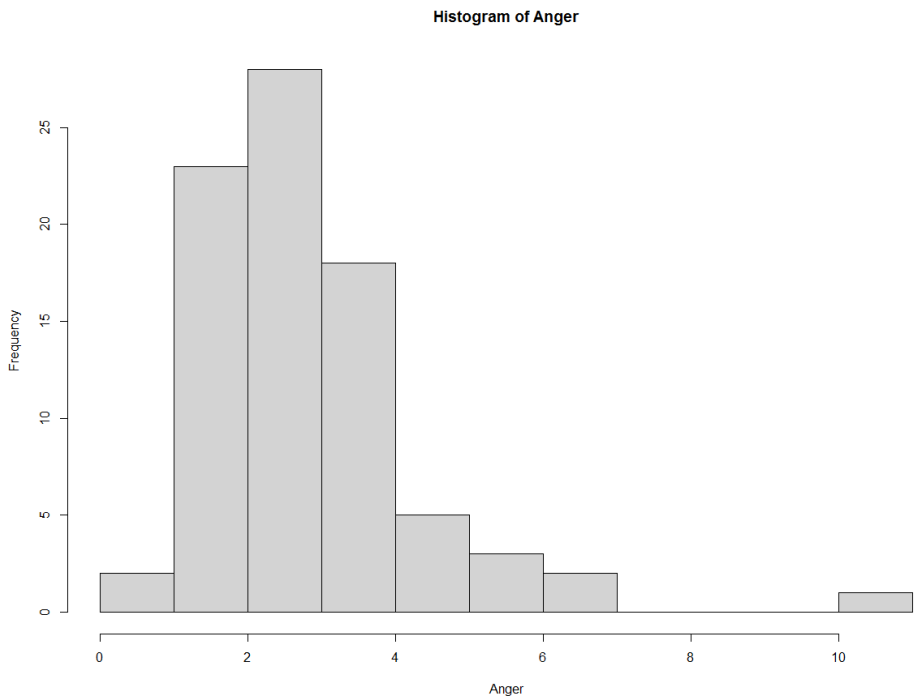


Figure 2
Histogram for all values of the variable joy in which outliers can be detected

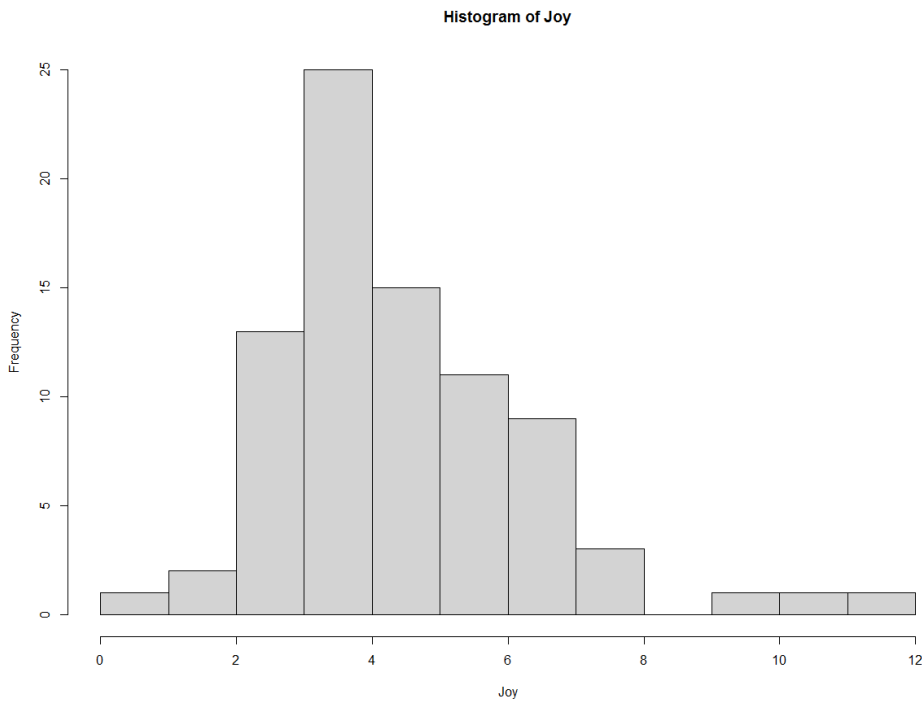


Figure 3

Scatterplot that shows that there is no linear relationship between the variables % of replies and Joy

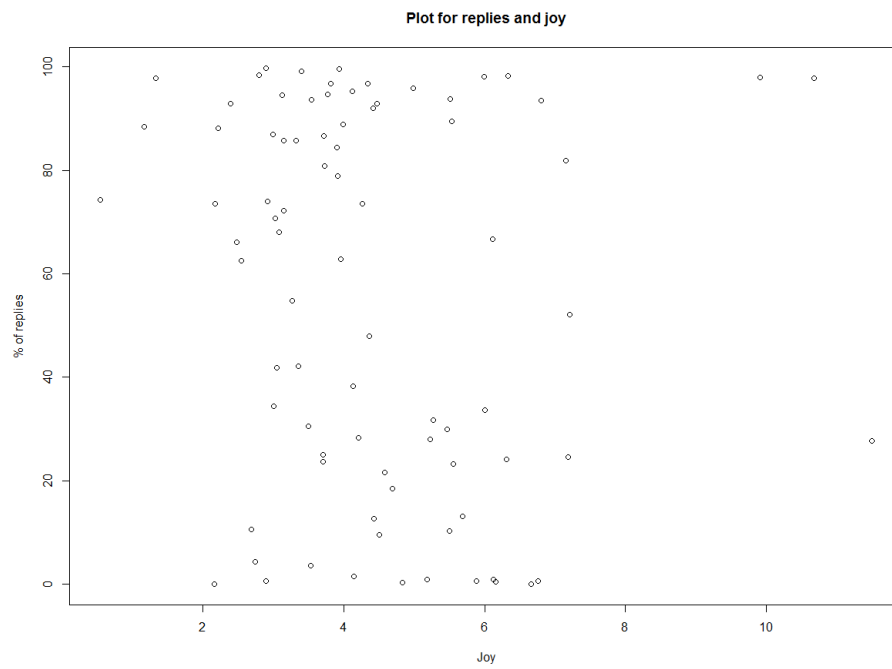


Figure 4

Scatterplot that shows that there is no linear relationship between the variables Anger and Trust

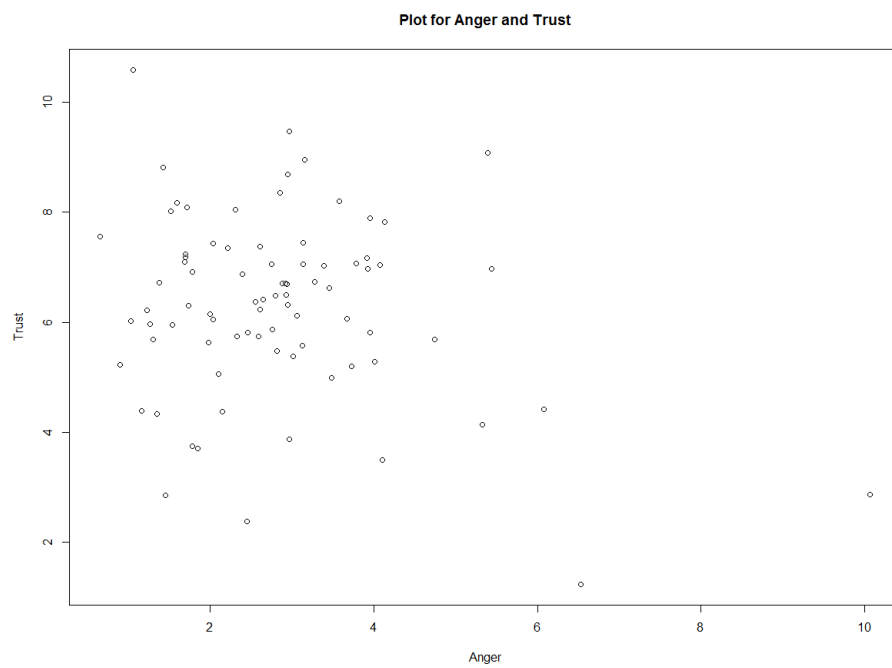


Table 1
Clustering results of the final selection of brands

Brand	% of replies	Joy	Anger	Trust	Cluster	Brand type
Accenture	9,53	4,50099	2,946648	8,689912	1	Loved brands
Adobe	31,66	5,2655	2,936814	6,694354	1	Loved brands
Allianz	1,53	4,139467	3,455956	6,620075	1	Loved brands
AXA	28,24	4,2121	4,748505	5,689413	1	Loved brands
Bank of America	18,53	4,692031	3,387078	7,029347	1	Loved brands
Caterpillar Inc	21,59	4,585439	3,275313	6,731549	1	Loved brands
Citi	0,94	6,126377	5,442719	6,975972	1	Loved brands
Deloitte	0,84	5,186737	3,012353	5,374558	1	Loved brands
Disney	0,5	6,156996	3,952445	7,897016	1	Loved brands
eBay	24,52	7,186858	2,310062	8,042437	1	Loved brands
ESPN	0,31	4,831063	2,799308	6,479043	1	Loved brands
EY	10,22	5,498362	1,263453	5,966308	1	Loved brands
GoldmanSachs	0,03	6,655974	2,646351	6,415397	1	Loved brands
Gucci	0,53	6,755946	2,038777	7,435539	1	Loved brands
HomeDepot	24,12	6,307922	1,375187	6,726457	1	Loved brands
IKEAUSA	23,19	5,558452	2,964942	9,474782	1	Loved brands
Loreal	12,6	4,427583	4,077378	7,037439	1	Loved brands
Louis Vuitton	0,53	5,880197	1,518484	8,027019	1	Loved brands
Microsoft	28,05	5,230386	1,691027	7,098381	1	Loved brands
Oracle	13,16	5,688718	3,139981	7,454458	1	Loved brands
SAP	27,75	11,50458	2,147683	4,367761	1	Loved brands
Siemens	29,97	5,464262	3,668291	6,064131	1	Loved brands
Uniqlo	33,6	6,006297	1,537903	5,957859	1	Loved brands
Adidas	93,69	3,541693	3,578974	8,208028	2	Loved brands
Colgate	96,75	4,339754	1,429566	8,81411	2	Loved brands
Dell	78,87	3,908689	1,700443	7,235034	2	Loved brands
Facebook	85,78	3,322011	1,71875	8,084239	2	Loved brands
FedEx	86,91	2,997118	1,786744	6,916427	2	Loved brands
Gillette	91,97	4,41408	3,7843	7,068151	2	Loved brands
Hyundai	84,47	3,904358	3,135593	7,058111	2	Loved brands
IBM	89,44	5,530606	3,155952	8,951092	2	Loved brands
Intel	54,72	3,267974	5,388526	9,084967	2	Loved brands
LEGO	85,69	3,147658	2,85763	8,359635	2	Loved brands
Mercedes Benz	80,81	3,727736	4,131629	7,82626	2	Loved brands
Walmart	97,75	1,328558	2,218115	7,347505	2	Loved brands
Amazon	88,87	3,983983	3,128818	5,578933	3	Fads
AmericanExpress	88,12	2,213542	2,445652	2,377717	3	Fads
BMW	73,56	4,264168	1,302801	5,682217	3	Fads
Chevrolet	92,91	2,386727	1,984346	5,633337	3	Fads
Coca-Cola	99,22	3,394256	2,930084	6,498404	3	Fads
Fritolay	94,72	3,764304	2,966913	3,873039	3	Fads
Google	95,22	4,123387	2,612528	6,232295	3	Fads
Honda	70,72	3,024449	2,461504	5,80606	3	Fads

HP	62,56	2,541364	2,886561	6,701583	3	Fads
Huawei	73,59	2,169486	1,456767	2,843045	3	Fads
JohnDeere	66,06	2,481048	0,903591	5,222452	3	Fads
KFC	99,53	3,931139	3,95776	5,812406	3	Fads
Lexus	72,19	3,141178	2,755111	7,063262	3	Fads
Mastercard	62,8	3,952209	1,35244	4,330612	3	Fads
Nescafe	92,84	4,474273	4,013686	5,27701	3	Fads
Nike	96,75	3,808909	1,786099	3,744351	3	Fads
Redbull	68	3,075103	3,725606	5,204021	3	Fads
Starbucks	98,44	2,795193	2,594915	5,738419	3	Fads
T-Mobile	86,72	3,713428	1,162791	4,388597	3	Fads
UPS	94,59	3,121853	1,02719	6,022155	3	Fads
WellsFargo	73,98	2,909776	3,481962	4,982206	3	Fads
ATT	88,44	1,164386	6,533321	1,231378	4	Low love brands
Audi	74,22	0,54107	10,06577	2,86394	4	Low love brands
McDonalds	99,75	2,889069	6,079541	4,418143	4	Low love brands
BurgerKing	98,19	6,337136	3,929024	6,978305	5	Loved brands
Cvs pharmacy	98	9,913773	1,055459	10,59228	5	Loved brands
Ford	93,75	5,50653	2,915637	6,713731	5	Loved brands
General Electric	52,15	7,212395	3,064582	6,122309	5	Loved brands
Lowe's	98,12	5,989566	2,549033	6,372581	5	Loved brands
Pepsi	95,84	4,984038	2,039221	6,045996	5	Loved brands
Toyota	97,78	10,68096	1,695729	7,188183	5	Loved brands
Verizon	66,73	6,114637	2,760149	5,872872	5	Loved brands
VW	81,94	7,149966	3,919467	7,165276	5	Loved brands
ZARA	93,47	6,808369	2,324809	5,745599	5	Loved brands
Chase	42,2	3,348921	1,226932	6,215115	6	Traditional brands
Cisco	10,51	2,689671	1,847334	3,700476	6	Traditional brands
FOX TV	25,01	3,708841	0,657078	7,556399	6	Traditional brands
hm	23,73	3,70249	2,817824	5,471822	6	Traditional brands
JP Morgan	3,59	3,526841	5,330836	4,138577	6	Traditional brands
Nestle	41,86	3,052091	1,738811	6,294938	6	Traditional brands
Netflix	30,55	3,493089	2,944712	6,317608	6	Traditional brands
Nintendo America	0,53	2,896572	4,104075	3,493136	6	Traditional brands
PayPal	48	4,357342	2,099085	5,053721	6	Traditional brands
PwC	0,06	2,157551	2,389614	6,874059	6	Traditional brands
Samsung Mobile	38,31	4,12758	2,611007	7,379612	6	Traditional brands
Sony	34,37	3,003364	1,994234	6,150889	6	Traditional brands
UBS	4,28	2,740353	1,596421	8,180385	6	Traditional brands

Figure 5

Dendrogram for the cluster results

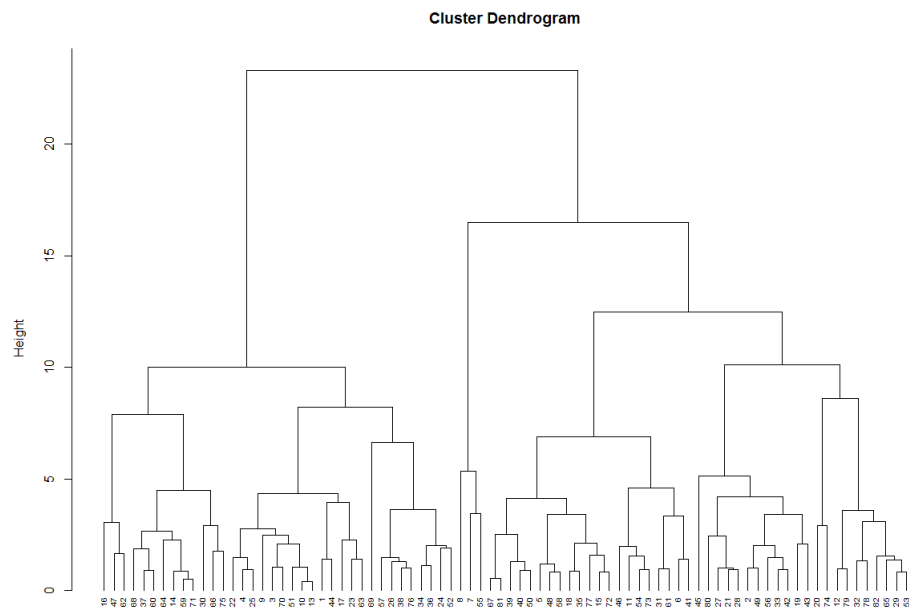


Figure 6

Perplexity plot for the optimal value of k

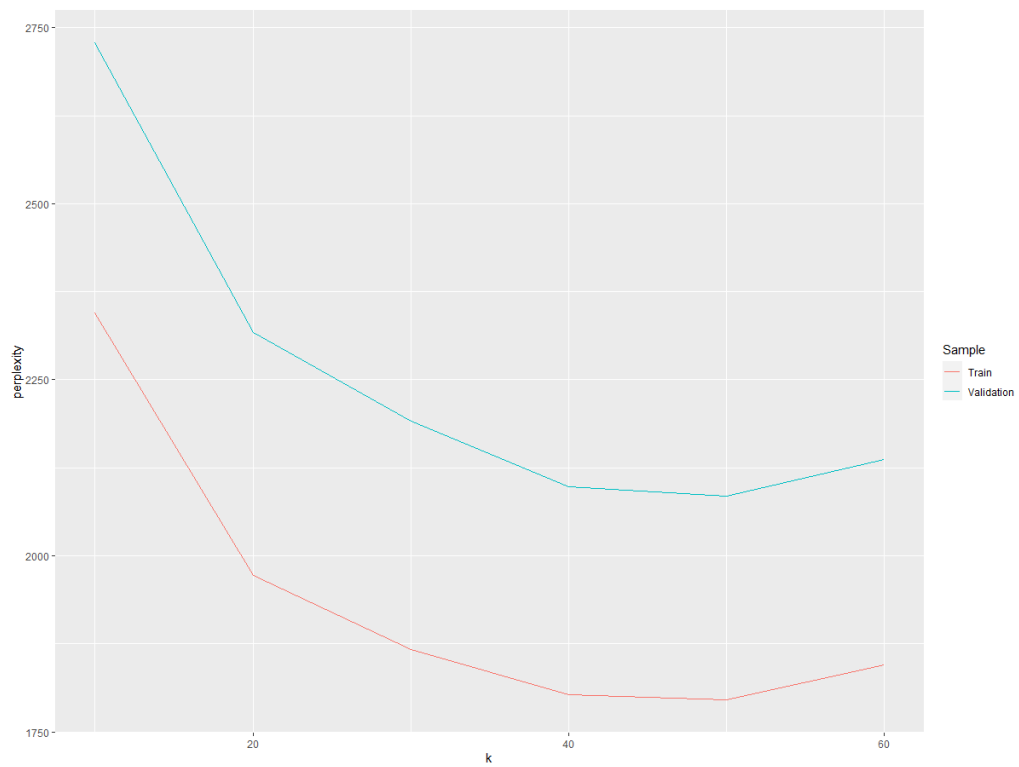


Figure 7

Perplexity plot for the optimal value of α

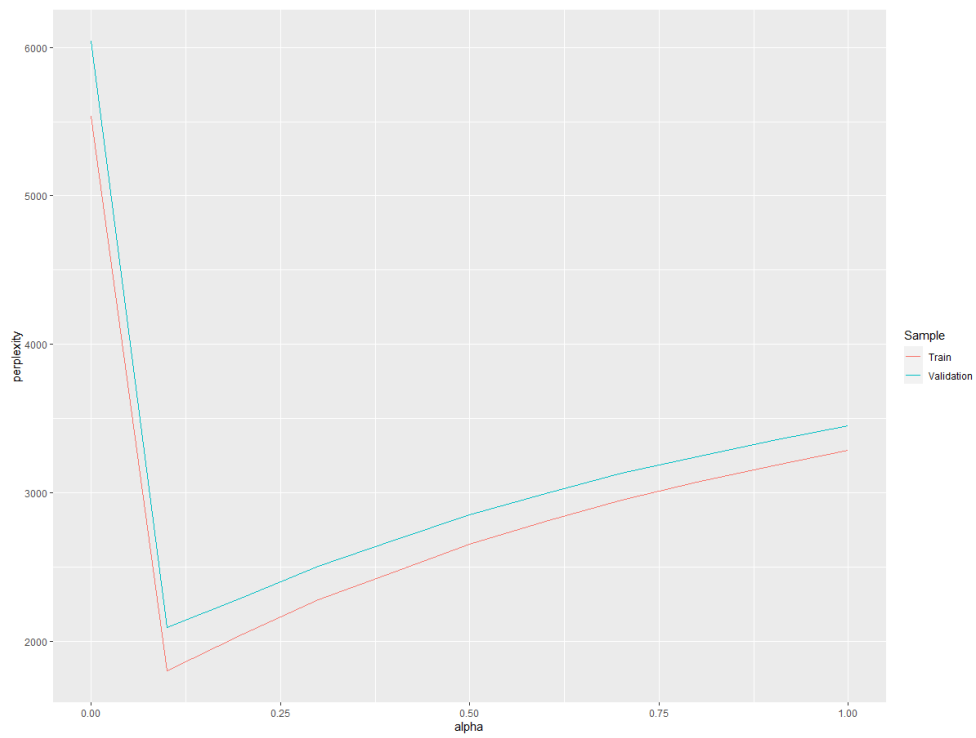
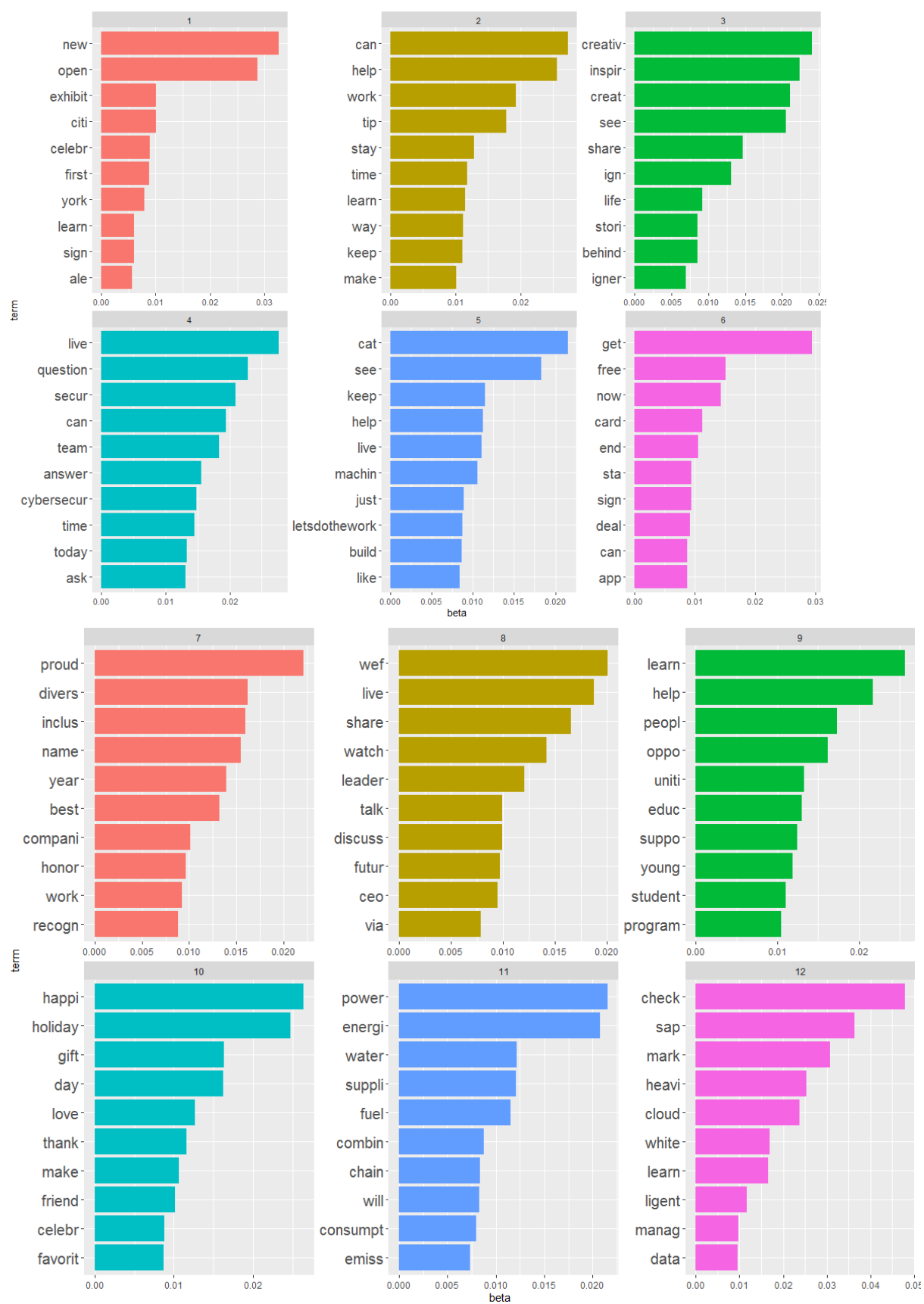
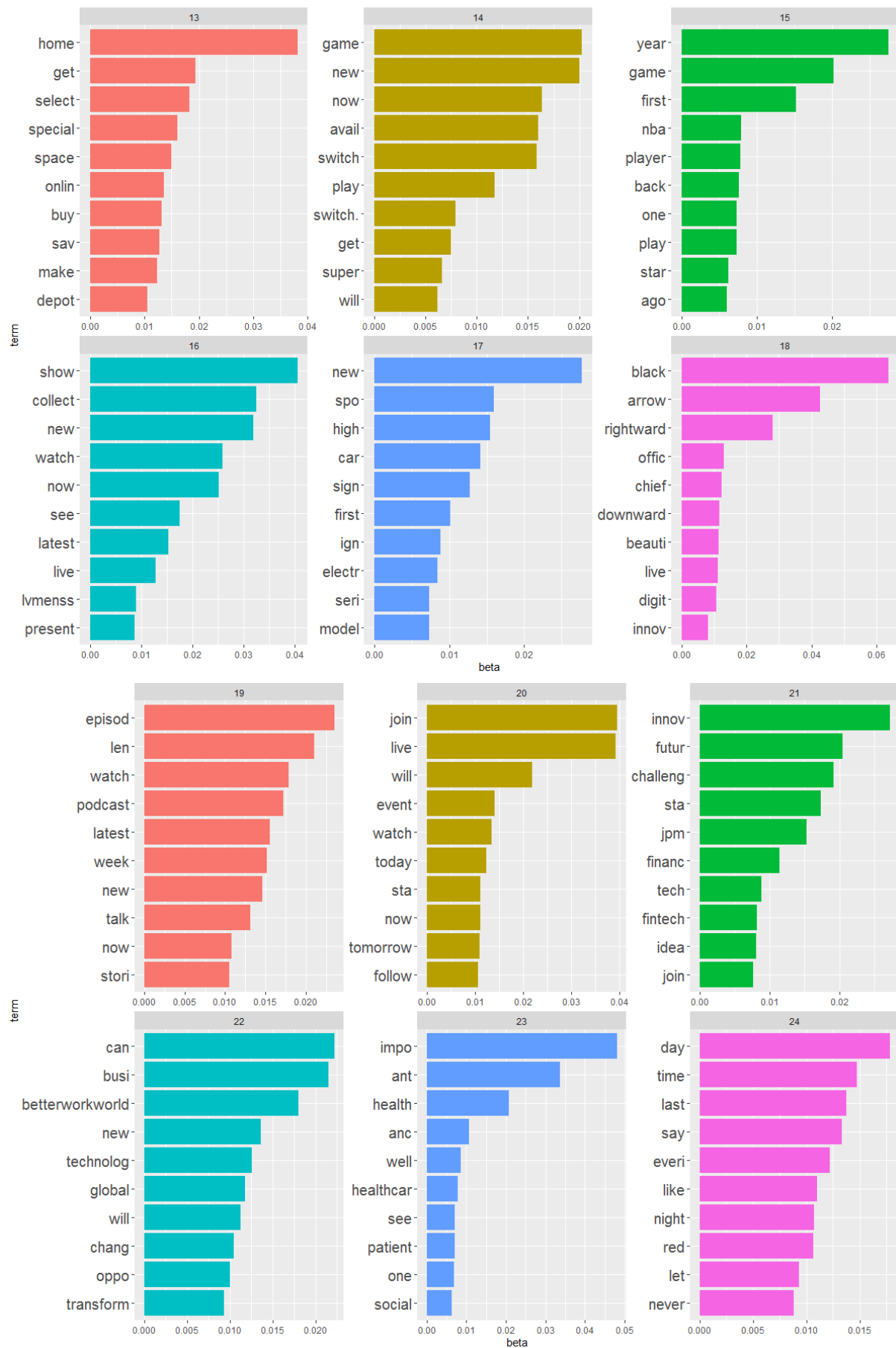
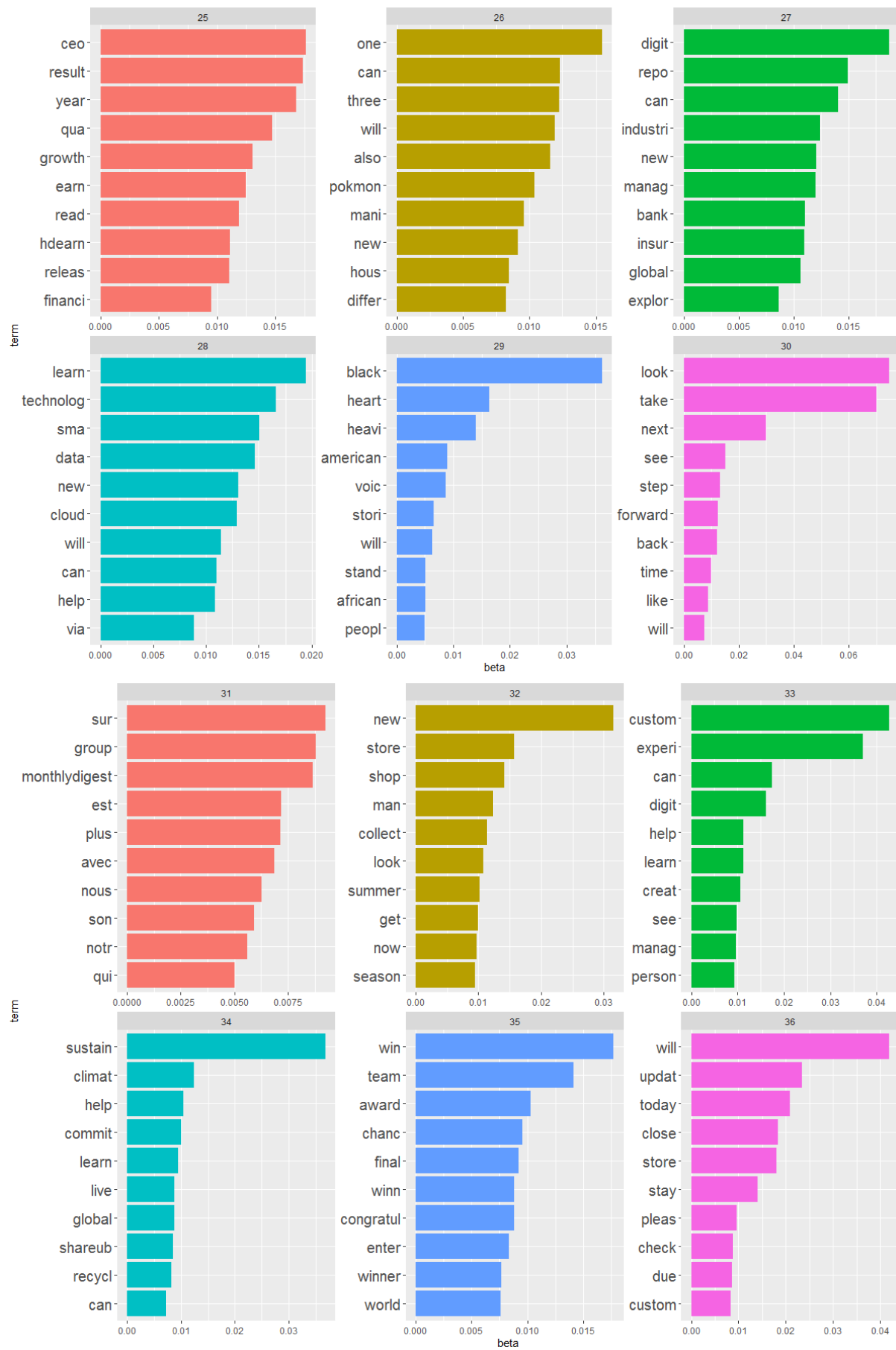


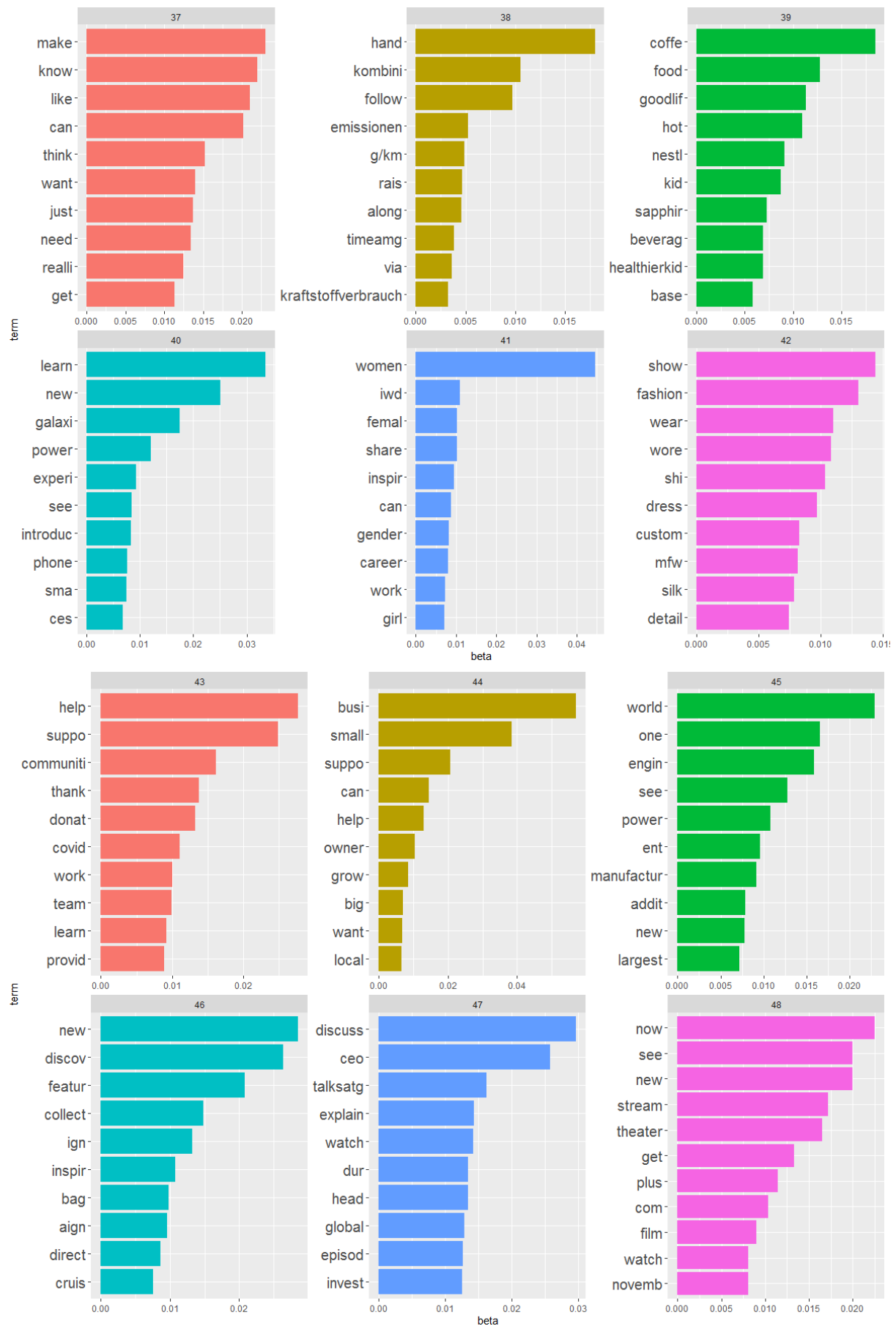
Figure 8

The 50 topics from the LDA analysis with corresponding ten words









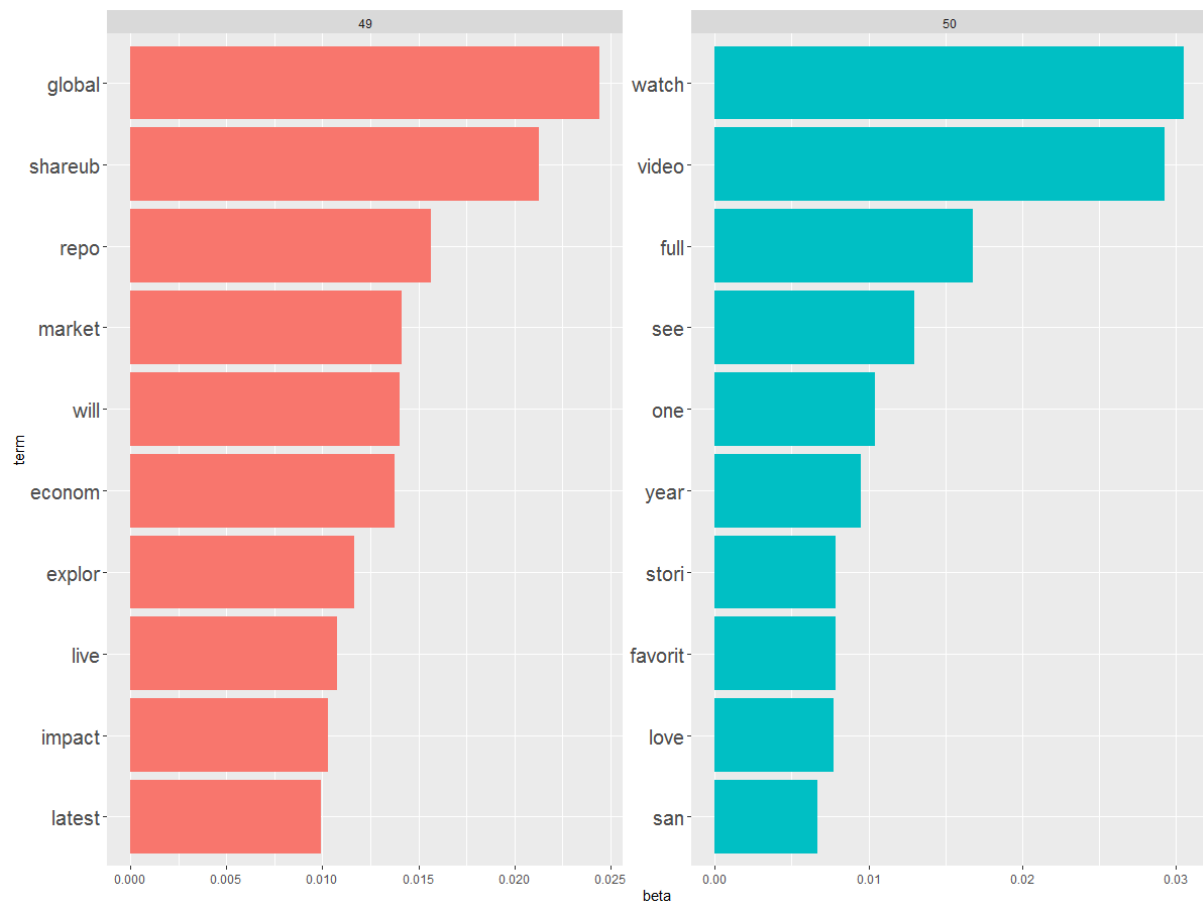


Figure 9
Variable importance in predicting ESM

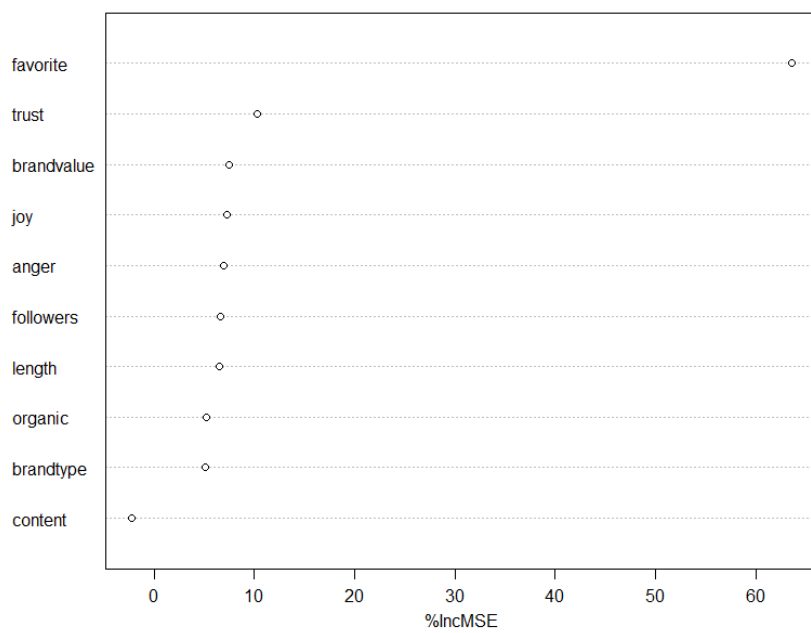
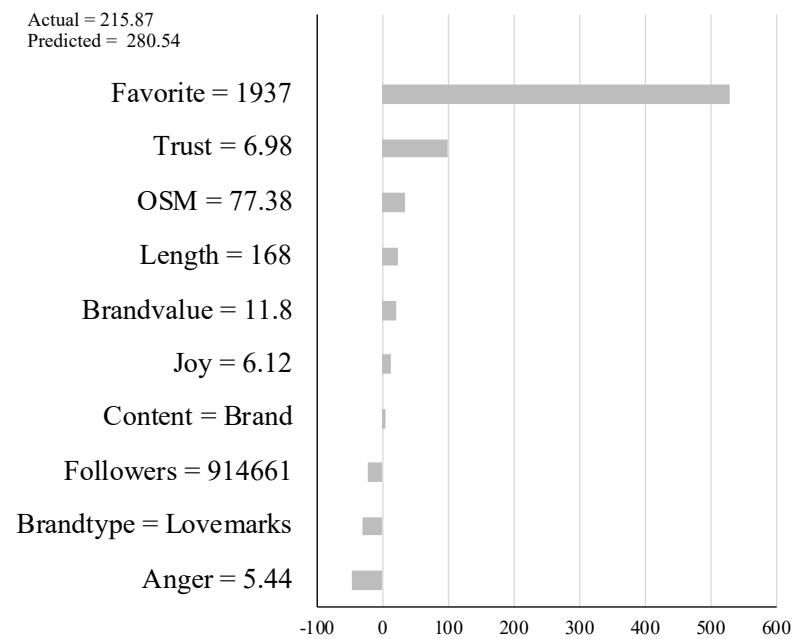
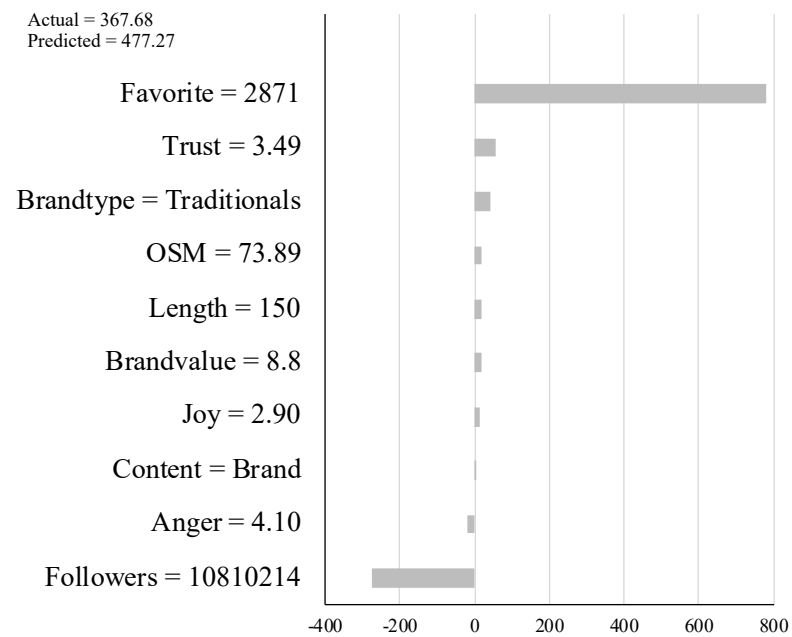
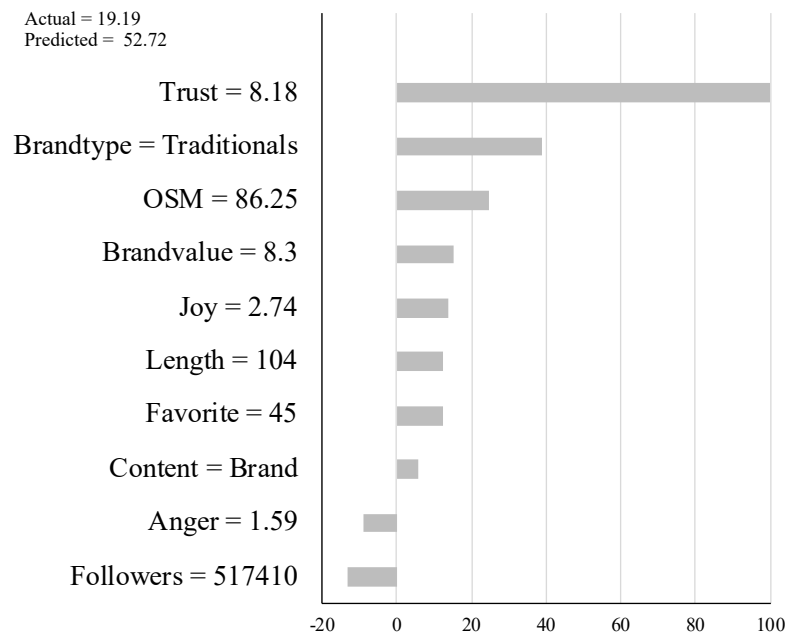
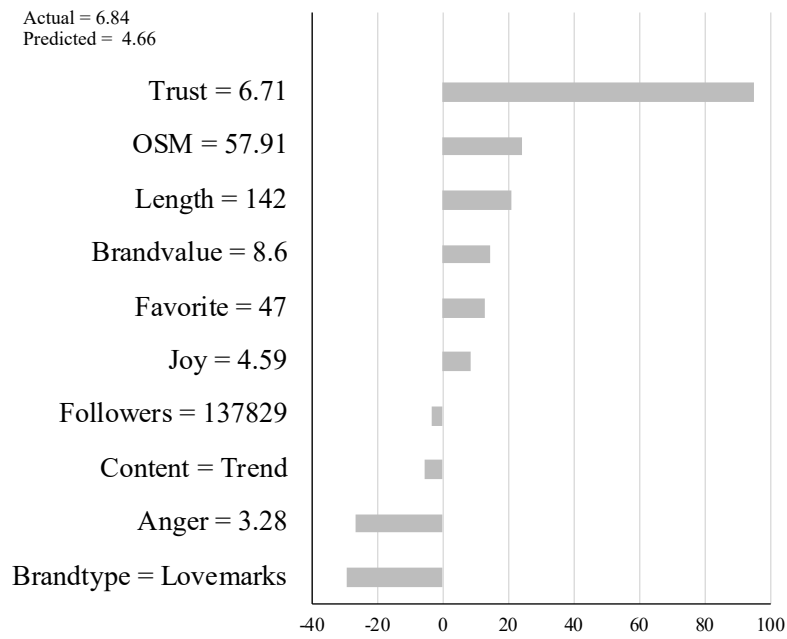


Figure 10

LIME results for five random observations in the prediction of the random forest





Actual = 32.50
Predicted = 13.03

