



UNIVERSITY OF AMSTERDAM

BACHELOR THESIS

YouTube's bias towards conspiracy content: an analysis of the YouTube algorithm

May 12, 2021

Student:
Roan Schellingerhout

Supervisor:
Dr. Maarten Marx

Abstract

YouTube is the second largest website on the internet. It brings in over 34.6 billion page views each month, most of which consist of videos being watched. 70% of those videos are found through YouTube's recommendation algorithm. Unfortunately, this recommender system, like most, is sensitive to filter bubbles. In recent years, conspiracy content has been gaining popularity on the website, partially because of the way in which the algorithm handles it. To gain a better understanding of how the algorithm deals with conspiracies, an experiment was setup wherein brand new YouTube accounts were made to watch conspiracy content adhering to different watch strategies. After watching 15 conspiracy videos, all accounts ended up with significantly more conspiracy recommendations than the baseline. Accounts that watched conspiracy videos that were recommended to them by YouTube ended up in filter bubbles especially quickly. While accounts ended up in a filter bubble after watching between five and seven videos, getting out of a filter bubble took considerably longer. This suggests that YouTube's algorithm is indeed susceptible to the creation of filter bubbles of conspiracy content.

Keywords: Recommender systems, Conspiracy theories, Filter bubbles, YouTube, Machine learning

1 Introduction

- ¹ YouTube attracts an average of 34.6 billion page views per month, making it
- ² the world's largest video-sharing website and the second largest website on the
- ³ entire internet (Neufeld, 2021). The overwhelming majority of those page views

4 come from users watching videos, 70% of which are recommended to users by
5 YouTube's algorithm (Cooper, 2020). All types of content get produced and
6 consumed on the website. However, conspiracy content has been booming on
7 YouTube (Donzelli et al., 2018). Alt-right (or far-right) and conspiracy chan-
8 nels are starting to grow their audiences, which could have many negative con-
9 sequences for society at large. For example, the number of people who are
10 distrustful of science is increasing, a development in which conspiracy content
11 on YouTube plays a role. Whenever this increased distrust relates to impor-
12 tant topics, such as believing in the efficacy of vaccines, it can create genuine
13 dangers to the public. As it turns out, more than half of the American popula-
14 tion has doubts about - or is definitely against - taking the COVID-19 vaccine
15 (Rosenbaum, 2021).

16 To better understand how YouTube's algorithm (and recommender
17 algorithms in general) allow(s) conspiracy content to thrive, this research will
18 investigate how quickly the algorithm develops a preference for conspiracy video;
19 in other words: how many videos a user needs to watch before they get sent
20 *down the rabbit hole*.

21 This research is based on the assumption that YouTube's recom-
22 mender system is susceptible to the creation of filter bubbles. This concept,
23 coined by Pariser (2011), has been studied in-depth on many social media web-
24 sites, YouTube included. While results vary slightly, the common finding is that
25 YouTube recommendations do indeed lead to filter bubbles and that extremist
26 and conspiracy content is more likely to do so (Bryant, 2020; O'Callaghan et al.,
27 2013; Ledwich and Zaitsev, 2019). This effect can lead to the radicalization of
28 impressionable users, with deleterious consequences. However, an important
29 factor herein is how quickly a user's recommendations turn into a bubble. If
30 the user has enough time to be exposed to other types of content, they might
31 stray away from the more extreme, preventing them from adopting a potentially
32 harmful view (Bozdag and Van Den Hoven, 2015). That is why this research
33 looks at how *quickly* a user could end up in a filter bubble on YouTube. To
34 do so, brand-new accounts will be made to watch YouTube videos according
35 to different watch strategies; after each video watched, the recommendations of
36 the user will be labeled as being either conspiracy content or regular content
37 by a machine learning classifier, to determine whether or not the user is in a
38 bubble.

1.1 Research question

39 For this research the following research question has been formulated: *What*
40 *is the impact of different watch strategies on the number of conspiracy videos*
41 *that have to be watched until a user's YouTube recommendations start preferring*
42 *conspiracy content?* In this scenario, 'preferring' will be defined as the situation
43 in which the amount of conspiracy videos present in the recommendations is
44 significantly higher than that of the baseline.

45 To assist in answering the research question, the following sub-questions
46 will be answered:

- 47 • How do different watch strategies on YouTube influence the type of con-
48 spiracy content that is recommended to a user?
- 49 • How long does it take for YouTube recommendations to stop preferring
50 conspiracy videos, once they have started doing so?
- 51 • What type of classifier is suitable for labeling conspiracy videos on YouTube?

2 Theoretical Framework

2.1 Filter bubbles on social media

52 Whenever the user of a website finds themselves in their own information uni-
53 verse, in which the content and recommendations play into the user's preexisting
54 opinions and beliefs, they are in a filter bubble (Pariser, 2011). Users are by
55 themselves in such bubbles and each bubble is unique. Different bubbles can
56 have overlap, but each bubble is precisely tuned to an individual. In traditional
57 media, a user makes a conscious *choice* what types of opinions they want to
58 hear, for example by choosing to watch a broadcaster with a specific political
59 opinion. Online this decision is implicit: based on the user's behavior, their
60 content is filtered automatically by an algorithm, without explicit consent.

2.2 Filter bubbles on YouTube

61 Previous research has found that YouTube's recommendation algorithm runs
62 the risk of creating filter bubbles. Roth et al. (2020) came to this conclusion
63 after they analysed YouTube recommendations based on content. YouTube has
64 two distinct types of recommendations: recommendations based on the user's
65 viewing behavior and recommendations based on the content of the current
66 video a user is watching. In their research, Roth et al. focused predominantly
67 on recommendations based on content. They found that such recommendations
68 could quickly lead to a decrease in information diversity (read: filter bubbles)
69 and that this decrease happened sooner for videos with a lot of views; the
70 more views a video had, the less diverse its related recommendations. They
71 speculate that this can be explained by the fact that YouTube tends to store
72 more information about videos with a high view count, allowing the algorithm to
73 give better recommendations for such videos. They also predict that, whenever
74 the algorithm has more information about a user to its disposal, it can combine
75 said information with the information it has about a certain video, which could
76 lead to an even stronger limitation of recommendations. According to Ledwich
77 and Zaitsev (2019), a user's viewing behavior is responsible for approximately
78 70% of their recommendations; this behavior could therefore play a big role in
79 the creation of filter bubbles on YouTube.

80 Once a YouTube user has entered a filter bubble, it can be difficult
81 to escape it. The most common way to help users get out of a filter bubble is
82 by exposing them to content covering viewpoints other than their own (Bozdag

83 and Van Den Hoven, 2015). However, for YouTube, this could form an issue.
84 YouTube makes its money by displaying advertisements to a user. The longer
85 a user stays on the website, the more profit YouTube can make. As a result,
86 YouTube's algorithm prefers recommending videos that are likely to generate a
87 lot of watch time (Maack, 2019). As it turns out, controversial content (such as
88 conspiracies) tends to have a higher audience retention: people keep watching
89 controversial content for longer (Birch, 2019). Whenever content is surprising
90 (which conspiracy theories often are), it is more likely to capture and keep a
91 user's attention. Thus, by showing the user more diverse content, the algorithm
92 would actively hinder its own goal. Because of this design, filter bubbles are
93 commonplace on YouTube.

2.3 Conspiracy content on YouTube

94 YouTube has limited rules with regards to the spread of conspiracy videos
95 (YouTube, 2021). As long as the content does not directly incite violence or
96 endangers the public health (e.g. misinformation about the COVID-19 virus),
97 objectively incorrect ideas are allowed to be shared on YouTube. As a result,
98 YouTube is a home to multiple conspiracy communities. Conspiracy theories
99 such as 'the earth is flat and the government is hiding it from us', 'the world
100 will end soon and only followers of this specific religion will be spared', and 'the
101 world is ruled by cannibalistic, satanic pedophiles' (better known as QAnon)
102 gather millions of views on the platform (Paolillo, 2018; Miller, 2021). Though
103 such videos could be considered harmful to society, they are not suppressed by
104 YouTube. Whenever a user shows interest in this type of content, they will
105 be recommended similar videos, even when YouTube is aware of their harmful
106 nature (Ledwich and Zaitsev, 2019; Maack, 2019).

2.4 The YouTube algorithm

107 The YouTube algorithm tries to recommend videos based on the expected watch
108 time they will generate, rather than the probability of a user clicking on them
109 (Covington et al., 2016). This decision was made in order to decrease the
110 likelihood of misleading videos (also known as *clickbait*) being recommended.
111 However, gathering feedback about videos through their watch time can cause
112 a lot of noise, making it difficult to measure user satisfaction. As it turns
113 out, even when a users enjoy a certain video, they are unlikely to watch it
114 completely. On average, users watch around 50-60% of a video before they
115 switch it off (Park et al., 2016). Though, videos that are well-structured, or
116 especially interesting, can improve this percentage up till 70-80%, where nearly
117 half of the viewers actually finish the video in its entirety (Lang, 2018). After a
118 video has been watched, there is a 41.6% chance that the user decides to watch
119 a recommended video. Which recommendation the user will choose, follows a
120 Zipf-distribution ($\alpha = 0.78$) with regards to the position of the video in the list
121 of recommendations (Zhou et al., 2010).

122 All in all, previous research has found that YouTube's algorithm is
123 sensitive to filter bubbles and that it has a tendency to recommend conspiracy
124 content. It is also speculated that the algorithm makes decisions based on the
125 user's viewing behavior, which it combines with the content of videos. In order to
126 keep the user on the website as long as possible, which is profitable for YouTube,
127 the algorithm prefers recommending videos that it suspects the user will watch
128 for a longer period of time, even when they may contain harmful content. Based
129 on this information, further research can be done on the origination of filter
130 bubbles and the spread of conspiracy content on YouTube. For example, little is
131 known about how quickly a user's recommendations adapt to a user's behavior,
132 even though this is a critical aspect when it comes to the creation of so-called
133 *rabbit holes*. Furthermore, no research has been done into the way different
134 types of videos (recommendations in different locations, random videos, etc.)
135 influence YouTube's algorithm. For example, whenever a user primarily watches
136 recommended videos, the algorithm could see this as implicit positive feedback,
137 which could cause a snowball-effect.

3 Methodology

3.1 Watching conspiracy videos

138 In order to determine how different watch strategies affect the YouTube algo-
139 rithm, a python script was created to automatically log into a Google account
140 and proceed to watch YouTube videos. The script was made using Selenium
141 WebDriver: a suite of tools used for browser automation.

142 In the following sections, each aspect of the initial experiment will be
143 explained. Firstly, an explanation will be given about how to log into Google
144 accounts using a python bot. Then, the watch strategies as mentioned in the
145 research question will be defined, followed by a description of their video watch-
146 ing behavior. Afterwards, some restrictions about the videos being watched will
147 be mentioned. Additionally the actual script allowing the bots to be run will be
148 described. Lastly, the precise form of the output of the script will be explained.

3.1.1 Google login

149 Due to Google's strict policy regarding automation within their ecosystem, many
150 obstacles are put into place to prevent users from logging into a Google account
151 using automated software such as a selenium script. To circumvent this restric-
152 tion, two steps had to be taken. Firstly, the selenium WebDriver had to be
153 accompanied by the selenium-stealth package, which removes metadata about
154 the current browser, so that it is less obvious that a WebDriver is being used.
155 Additionally, because this metadata was removed, the Google login service was
156 unable to check what browser the client was using. This results in a warning
157 to the user that their current browser may be insecure, which prohibits them
158 from logging in. To avoid this warning, the Google account needs to have been

159 created within a WebDriver, such as Google's ChromeDriver or Mozilla's Gecko-
160 Driver. Therefore, all twenty accounts were manually created in ChromeDriver.
161 Since Google accounts require a phone number verification upon creation, six
162 free (prepaid) SIM cards were ordered from various providers in order to create
163 the accounts. Each SIM card could create two to three accounts before it was
164 blocked due to being used too many times.

3.1.2 The watch strategies

165 After all accounts had been created, they were subdivided into four distinct
166 watch strategies, making for a total of five accounts per strategy.

- 167 **1. Random videos (baseline)** The first watch strategy is the simplest one.
168 The bots following it will watch random, non-conspiracy videos from a
169 dataset. This watch strategy is used as the baseline to compare the other
170 three strategies to.
- 171 **2. Random conspiracies** The second strategy is similar to the first: the ad-
172 hering bots watch random conspiracy videos from a dataset.
- 173 **3. Video recommendations** The second-to-last strategy starts off in the same
174 way as strategy 2: it chooses a random conspiracy video from a dataset
175 to watch. However, it then watches the four most similar videos in the
176 dataset (based on cosine similarity) in order to allow the algorithm to
177 get a feel for the user's interests. After watching those five initial videos,
178 it starts looking at the recommended videos displayed next to the cur-
179 rent video and chooses the recommendation that is most likely to be a
180 conspiracy video. These recommendations consist of a combination of
181 recommendations based on the content of the current video and the per-
182 sonalized recommendations of the user. By using this strategy, bots are
183 likely to go *down the rabbit hole* and eventually end up in a filter bubble.
- 184 **4. Homepage recommendations** Finally, the last strategy is similar to the
185 previous one, though with one alteration: rather than choosing a rec-
186 ommended video from the list of recommendations next to the current
187 video, it will choose a recommended video from the YouTube homepage
188 of the account. Compared to the third strategy, this will lead to the user
189 watching more personalized recommendations rather than content-based
190 recommendations, possibly speeding up the creation and/or increasing the
191 strength of the filter bubble.

192 Strategy 2 was not added in order to study the formation of filter bubbles, since
193 it was unlikely that a filter bubble would come from this strategy. Considering
194 the video choices would be all over the place, it would be difficult for the algo-
195 rithm to determine the specific interest of the user. However, comparing how
196 the algorithm responds to conspiracy content in general as opposed to regular
197 content might still yield interesting results.

For strategy 3 and 4, the likelihood of a recommendation being a conspiracy video was estimated by a neural network using the title, description, transcript, channel description, and channel keywords of the specific video. Though this is more information than a regular user would have to their disposal, it has to be kept in mind that humans are able to interpret the thumbnail of the recommendations, can have foreknowledge about the channel uploading the video or the subject being mentioned in the title, etc. Thus, the choice was made to allow the neural network to consider the transcript and look at general information about the uploader to balance the scales.

Each strategy was executed by five different accounts in order to decrease the probability of a random streak of videos altering the result. The individual accounts watched a total of fifteen videos as described by their watch strategy for a total of three hundred videos watched by the script.

Additionally, to simulate real-world user behavior, the average watch time for the videos was normally distributed with a mean of 55% and a standard deviation of 25% (Park et al., 2016; Lang, 2018). The watch time was not be able to exceed a value of 100 or subceed a value of 0, as a video cannot be watched for more than 100% or less than 0%. In the same vein, the clicking behavior of users was simulated as accurately as possible. Whenever none of the recommendations were predicted to be conspiracy videos, the probability of a user clicking on a video at position k within a given list of recommendations (its click-through rate: CTR), was determined using the following formula:

$$CTR(k; N, \alpha) = \frac{1/k^\alpha}{\sum_{n=1}^N (1/n^\alpha)} \quad (1)$$

Wherein N is the total number of recommendations and α is the distribution's exponent value ($\alpha = 0.78$) (Zhou et al., 2010). Using this formula, when considering the first twenty recommendations, the first recommendation will have a click-through rate of approximately 20.6%, after which the CTR quickly decreases, until a probability of 1.9% at the twentieth recommendation.

3.1.3 Running the bots

After the accounts were logged in, they started watching YouTube videos according to their watch strategy. However, some restrictions were put into place to make sure the bots did not take too long (considering three hundred videos had to be watched in total, some limitations had to apply). For example, the bots were not allowed to watch videos over an hour long, nor were they allowed to watch live streams, as those could theoretically go on infinitely. Additionally, the random videos at the start of the third and fourth strategy were first manually inspected to make sure the bots would not start the experiment by watching a falsely flagged conspiracy video. Considering the way in which the dataset was created, it is possible that some videos that are flagged as conspiracy videos are, in reality, normal videos. This could happen whenever a conspiracy channel uploads a regular video for once (e.g. a holiday video, promotion of some product, etc.). These *false positives* are far and few between, however,

238 having one of them be selected as the first video for strategy three or four had
239 to be prevented, as that could have greatly altered the final results. With these
240 restrictions in mind, the following script was created and run for all twenty bots,
241 keeping track of the videos they watched and the homepage recommendations
242 they had after each video:

Algorithm 1: Watch YouTube videos according to a watch strategy

Data: User information and a video dataset

Result: The watched videos and homepage recommendations of the user

```
1 for twenty bots do
2   initialize WebDriver;
3   log into Google account;
4   for fifteen videos do
5     if there is a recommendation to be watched then
6       | go to the link;
7     else
8       | pick a random video to watch based on usertype;
9       | determine how long it will get watched;
10      | go to the link;
11    get video metadata and store for overview of watched videos;
12    watch video for given amount of time;
13    if usertype == 3 then
14      | pick recommendation next to current video to watch next;
15      | determine watch time for found recommendation;
16    go to YouTube homepage;
17    store current recommendations for overview;
18    if usertype == 4 then
19      | pick homepage recommendation to watch next;
20      | determine watch time for found recommendation;
21 return watched videos and homepage recommendations;
```

243 Running the script for all twenty bots resulted in two different datasets: the first
244 containing the videos watched by the bots and the second containing the home-
245 page recommendations for all bots, after each number of videos watched. To
246 determine the influence of the watch strategies on the algorithm, the recommen-
247 dations were labeled as being either conspiracy or non-conspiracy videos by the
248 classifier. By then grouping all recommendations by their watch strategy and
249 the number of videos watched before them (e.g. the recommendations after the
250 third video watched by all bots with strategy one), it was possible to calculate

Example DataFrame

bot	vid_no	video_id	views	likes	dislikes	vid_len	title	description	transcript	channel_desc	keywords	conspiracy
1	5	u6sN6K-4_qo	49778	2913	16	536	All I Want for Christmas (Is a Pardon from the...	It's almost Christmas time. Saxophones by Jon ...	It's almost christmas time and i'm all by myse...	Nick Lutsko is a songwriter, producer, and per...	Nick Lutsko Puppets Etc. ALL SHOOK UP Greezy R...	False
1	5	p212D-WvKGO	71645	3073	38	639	How A Narcissist Reacts When You Behave Badly	How a narcissist reacts when you do something ...	[Music] today I'm going to talk about the diff...	Pop Culture -- Narcissism & Sociopathy -- Phil...	Narcissism Sociopathy Philosophy "Trauma Infor...	True
1	5	7vLbNZs2sFk	25137	132	18	439	FORECAST: Phoenix will stay in the 70s though ...	High pressure over the Pacific is producing a ...	let's take a look at what we're tracking weath...	Get all the news of the day, plus investigativ...	cbs5az 3tv "cbs 5 phoenix" "cbs 5 az" ktvk kph...	False

Figure 1: An example of the experiment's output

aggregates about general statistics of the recommendations, such as view count and video duration, and the percentage of conspiracy videos present amongst them. This lead to four groups with fifteen entries of different statistics (one for each video watched). In order to find out whether any of the differences between the groups were significant at any point, a number of independent sample t-tests were performed.

3.1.4 Outputs of the experiment

To create the initial recommendation dataset, the bots stored the following information about their top twenty recommendations after each video watched:

- the id of the bot that had gotten the recommendations;
- the amount of videos watched before the recommendation was given;
- the URL of the video being recommended;
- the URL of the channel that uploaded the recommendation.

To then allow the classifier to label the videos, and to calculate the aggregates per strategy, additional data was gathered using YouTube's API. For each recommendation, the title, description, transcript, (dis)likes, views, video duration, channel description, and channel keywords were collected. The title, description, transcript, channel description, and channel keywords were collected in order for the classifier to label the videos, as will be further explained in section 3.3. The other features were collected in order to calculate the aggregates used for the aforementioned ANOVAs. An example of the output can be seen in figure 1.

3.2 Leaving the filter bubble

To find out how quickly different users can get out of a filter bubble after they have gotten into one, an additional experiment similar to the one described before was set up. This experiment however, was significantly simpler. Rather than having different bots behave differently, all bots behaved the exact same

276 way. After the bots adhering to strategy 2, 3, or 4 had gotten into a filter bubble,
277 this experiment was performed on them in order to see how long it takes for
278 users adhering to different watch strategies to leave a filter bubble once they
279 find themselves in one.

3.2.1 The setup

280 Considering the experiment studies the way in which users leave a filter bubble,
281 only the bots that had actually gotten into a filter bubble were used. This
282 means that the first five bots, which were used as a baseline, were ignored.
283 The remaining bots went through the same starting procedure as they did in
284 the earlier experiment: the WebDriver was initialized and the bots logged into
285 their corresponding accounts. The same restrictions for videos (i.e. maximum
286 watch time) from the other experiment applied. However, rather than the bots
287 watching videos adhering to their original strategy, all bots watched random
288 non-conspiracy videos from the dataset. Once again, after each video, the bots
289 stored their homepage recommendations. Through doing so, it became possible
290 to see how many videos would have to be watched, for each bot, before their
291 recommendations started looking similar to that of the baseline again. Thus,
292 the following script was created:

Algorithm 2: Getting out of a filter bubble

Data: User information and a video dataset

Result: The watched videos and homepage recommendations of the user

```
1 for fifteen bots do
2   initialize WebDriver;
3   log into Google account;
4   for fifteen videos do
5     pick a random non-conspiracy video to watch;
6     determine how long it will get watched;
7     go to the link;
8     get video metadata and store for overview of watched videos;
9     watch video for given amount of time;
10    go to YouTube homepage;
11    store current recommendations for overview;
12 return watched videos and homepage recommendations;
```

293 After the script had been run, the same steps were taken as in the previous
294 experiment: the title, description, transcript, channel description, and channel
295 keywords of each recommendation on were downloaded, after which the classifier
296 predicted whether or not each recommendations was a conspiracy video. Then,

the recommendations were grouped by watch strategy and number of videos watched. In doing so, the results could be used to determine how quickly the recommendations for each watch strategy returned back to normal. By then analysing the results again using independent sample t-tests, the *strength* (or *inescapability*) of the filter bubbles created by different watch strategies could be measured.

3.3 Machine learning

3.3.1 Data gathering

To answer the research question, it is necessary to determine which YouTube videos can be considered conspiracy videos. Considering the large amount of videos getting recommended, determining each video manually is simply not possible. There are two possible ways to solve this problem. Firstly, there is a dataset which contains nearly 7000 YouTube channels that have been manually labeled based on their political view - almost 3000 of which were labeled as conspiracy channels (Ledwich and Zaitsev, 2019); whenever a video is made by one such channel, it can be considered a conspiracy video. However, due to the enormous amount of existing YouTube channels, the odds of a video being uploaded by a channel that is not present in this dataset are very large. For those videos, a supervised machine learning classifier was used. To optimize performance, five different classifiers have been trained and compared: k-nearest neighbors, support-vector machine, neural network, logistic regression, and ridge regression.

In order to train these machine learning algorithms, a training dataset was created. To get a labeled dataset of conspiracy and non-conspiracy videos, use was made of the aforementioned channel dataset made by Ledwich and Zaitsev (2019). For each channel in that dataset, the title, description, and transcript of the ten most recently uploaded videos were downloaded using YouTube's API. Videos uploaded by a conspiracy channel were then labeled as conspiracy videos, and videos uploaded by a channel from a different category were labeled as normal videos. Additionally, the channel description and channel keywords (which are used for targeted advertising on YouTube) were added to each video. The final dataset contains 65.683 unique YouTube videos, 22.156 of which are considered as conspiracy videos.

3.3.2 Data cleaning

However, this dataset was not yet suitable for machine learning, as the data was still messy. Therefore, multiple steps were taken in order to clean the data. Firstly, the two classes (conspiracy and non-conspiracy) were balanced, so that the classifier would not develop a bias for non-conspiracy videos. Rather than opting for balancing the two classes through the use of class-weights (a technique where weights are attributed to classes, thereby telling the classifier that getting a prediction correct for a certain, underrepresented class is more important), the

choice was made to under-sample the data in order to equalize both classes (both containing 22.156 videos, for a total of 44.312 videos) (Lemaître et al., 2017; Sun et al., 2006). As there was plenty of data in the dataset, under-sampling was more convenient than implementing class-weights. After both classes had been balanced, the text for each video had to be translated into English. Since the original dataset by Ledwich and Zaitsev (2019) also contained channels by non-English speakers, these videos had to be automatically translated. Then, a few common cleaning methods were applied: all text was converted to lowercase, after which special characters, such as emojis were removed, whereafter stop words were removed and all words were stemmed using the porter stemmer (Karaa, 2013). Finally, each video was TF-IDF vectorized to allow the classifiers to function.

3.3.3 Performance optimization

After splitting the dataset into a training, test, and validation set, the hyperparameters of each algorithm were tuned to get the optimal performance (Feurer and Hutter, 2019). Performance was measured using four distinct metrics: the accuracy, which shows the share of correct predictions; the recall, which shows what fraction of truly positive samples were correctly labeled as such; the precision, which shows what part of the positive predictions were correct; and the F1-score, which is the harmonic mean of the recall and precision (Sokolova and Lapalme, 2009). For each classifier, different configurations of hyperparameters (such as the kernel and the penalty-parameter) were systemically tested - each possible combination was tried. The classifiers were trained on the training set and the optimal hyperparameters were determined based on the performance of the classifiers on the validation set. By saving these performance measures for every configuration, for every classifier, the optimal configuration for every classifier could be determined. Lastly, the classifiers were equipped with their optimal hyperparameters and then tested for the final time on the test set. By comparing the performance of every optimally configured classifier on the test set, the best-performing classifier could be chosen (Reitermanova, 2010).

Additionally, the added value of using a machine learning ensemble was measured. By having each classifier make a prediction for all videos in the dataset, a new dataset was created, wherein the features were the predictions made by the different classifiers. By using all possible combinations of classifiers, and then having different neural networks use those features as input, a machine learning ensemble was created. This ensemble was then optimized in a similar way to the classifiers individually.

4 Results

4.1 The recommended content

- All strategies except for the baseline lead to filter bubbles;

Conspiracy recommendations

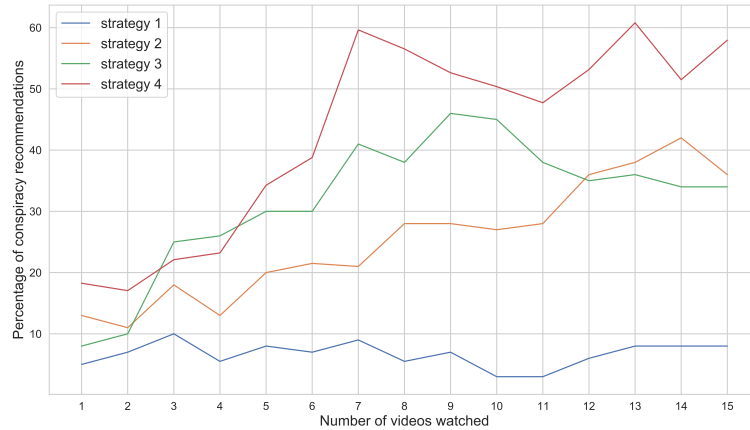


Figure 2: The percentage of conspiracy recommendations after each number of videos per strategy

- strategy 4 had the biggest bubble;
- strategy 2 and 3 do not differ all too much, which is interesting;
- seems like YouTube sort of 'caps' the amount of conspiracy videos on the homepage, as each time a new height is reached, the number starts decreasing;
- All strategies except the baseline lead to a lower diversity of content;
- Average views drop rapidly when a user starts watching content, regardless of the strategy;
- typical example of a cold-start problem;
- The average duration of videos is higher for the conspiracy strategies, but it seems somewhat random.

4.2 Leaving the filter bubble

- The number of conspiracy recommendations drops quickly, but keeps hovering a fair amount above baseline;
- This means that YouTube quickly adapts to the new content being watched, but still keeps the previous content in mind;

Conspiracy recommendations

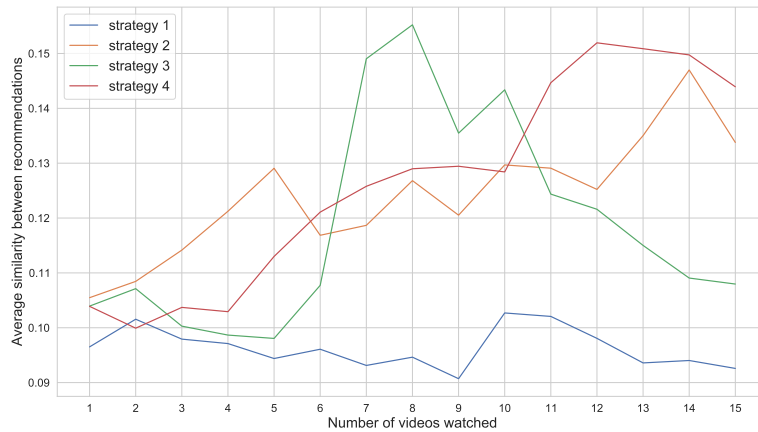


Figure 3: The similarity of recommended videos after each number of videos watched

387

- The amount at which it hovers, seems to be proportional to the initial

Conspiracy recommendations

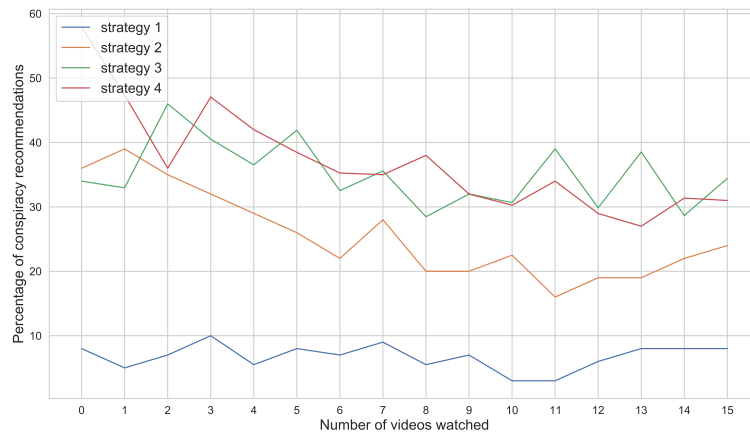


Figure 4: The percentage of conspiracy recommendations after each number of videos per strategy while attempting to leave the filter bubble

value of the filter bubble: if a user gets recommended more conspiracy content at first, more of that will stay once they stop watching that type of content;

- The number of videos until the algorithm stops recommending conspiracy content is well over 15;
- This is interesting, as a user ends up in a bubble after ± 5 videos watched.

4.3 Machine learning

Classifier performance

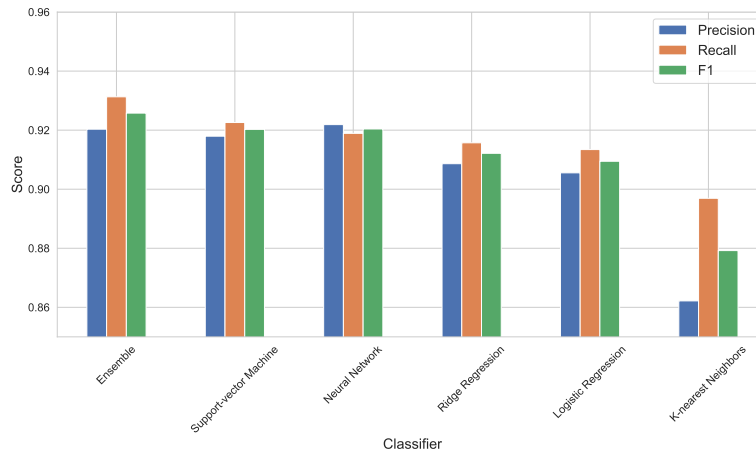


Figure 5: Metrics for each classifier with optimized hyperparameters

The hyperparameter tuning lead to impressive scores for all classifiers. When making predictions for the test set, the best-performing classifier is the support-vector machine making use of the Radial Basis Function (RBF) kernel and a penalty parameter (C-value) of 10. The SVM is tied for F1-score with the neural network using the identity activation function, with 10 hidden layers of 10 neurons. Ridge regression with a sparse-cg solver and penalty (alpha) value of 0.1 takes the third place, very closely followed by logistic regression with an L2 penalty, a penalty (C) value of 20 and a newton-cg solver. The worst-performing classifier is also the simplest of the bunch: the k-nearest neighbors classifier (K=1). Although its performance is still formidable, it does substantially worse than the others. An overview of all metrics for each classifier can be seen in figure 5. The ten best-performing configurations for each classifier can be found in appendix A.

Noteworthy is the fact that the optimal ensemble actually outperforms the support-vector machine by a slight margin. This ensemble, consisting of the

409 SVM, the neural network, and surprisingly, the k-nearest neighbor classifiers,
410 gets slightly higher scores than the runner-up across the board. The ensemble
411 had a 16-way tie for best-performing parameters, all of which contained at least
412 the SVM, neural network, and k-NN classifiers.

413 Though the ensemble outperforms the other classifiers, it has a signif-
414 icant drawback: its training time is significantly larger than that of the individ-
415 ual classifiers. Support-vector machines are infamous for their slowness when
416 there is a lot of training data, and neural networks can require a lot of training
417 time whenever the number of neurons gets large (Burges and Schölkopf, 1997;
418 Kamarathi and Pittner, 1999). Requiring both algorithms to run will therefore
419 require a lot of additional training time. Considering the marginal performance
420 increase, the cost outweighs the benefit. As a result, when taking everything into
421 account, the support-vector machine is the best classifier for labeling conspiracy
422 videos on YouTube.

423 Additionally, for predicting the likelihood of recommendations being
424 conspiracy videos in real-time, as was done for strategy 3 and 4, the best-
425 performing classifier is the neural network. Here, the neural network is preferred
426 over the support-vector machine, as neural networks are better optimized for
427 providing probabilities of samples belonging to a certain class (Specht, 1990).

5 Discussion

428 To do.

6 Conclusion

429 To do.

7 Future work

430 To do.



8 Planning

Table 1: Planning

Week	Handelingen	Afgehandeld
28/03-03/04	Hyperparameters optimaliseren	Ja
03/04-10/04	Classifier-ensemble optimaliseren	Ja
11/04-17/04	Google accounts maken, deelvraag 3 maken	Ja
18/04-24/04	Inleiding uitbreiden	Ja
25/04-01/05	Uitvoeren experiment, labelen met classifier	Ja
02/05-08/05	Beginnen deelvraag 1	Ja
09/05-15/05	Deelvraag 1 afschrijven	
16/05-22/05	Beginnen deelvraag 2	
23/05-29/05	Deelvraag 2 afschrijven	
30/05-05/06	Discussie schrijven	
06/06-12/06	Tekst proof-readen, laatste aanpassingen	
13/06-19/06	Inleveren scriptie en verdediging	



References

- 431 Birch, M. K. S. (2019). White rabbit. the logic and proportion of conspiracy
432 theory videos on youtube: a foucauldian discourse analysis. *Malmö univer-*
433 *sitet*.
- 434 Bozdag, E. and Van Den Hoven, J. (2015). Breaking the filter bubble: democ-
435 racy and design. *Ethics and information technology*, 17(4):249–265.
- 436 Bryant, L. V. (2020). The youtube algorithm and the alt-right filter bubble.
437 *Open Information Science*, 4(1):85–90.
- 438 Burges, C. J. and Schölkopf, B. (1997). Improving the accuracy and speed
439 of support vector machines. *Advances in neural information processing*
440 *systems*, pages 375–381.
- 441 Cooper, P. (2020). How does the youtube algorithm work? a guide to getting
442 more views.
- 443 Covington, P., Adams, J., and Sargin, E. (2016). Deep neural networks for
444 youtube recommendations. In *Proceedings of the 10th ACM conference on*
445 *recommender systems*, pages 191–198.
- 446 Donzelli, G., Palomba, G., Federigi, I., Aquino, F., Cioni, L., Verani, M., Car-
447 ducci, A., and Lopalco, P. (2018). Misinformation on vaccination: A quan-
448 titative analysis of youtube videos. *Human vaccines & immunotherapeutics*,
449 14(7):1654–1659.
- 450 Feurer, M. and Hutter, F. (2019). Hyperparameter optimization. In *Automated*
451 *Machine Learning*, pages 3–33. Springer, Cham.
- 452 Kamarathi, S. V. and Pittner, S. (1999). Accelerating neural network training
453 using weight extrapolations. *Neural networks*, 12(9):1285–1299.
- 454 Karaa, W. B. A. (2013). A new stemmer to improve information retrieval.
455 *International Journal of Network Security & Its Applications*, 5(4):143.
- 456 Lang, P. (2018). Youtube average view duration - the 50% rule.
- 457 Ledwich, M. and Zaitsev, A. (2019). Algorithmic extremism: Examining
458 youtube’s rabbit hole of radicalization. *arXiv preprint arXiv:1912.11211*.
- 459 Lemaître, G., Nogueira, F., and Aridas, C. K. (2017). Imbalanced-learn: A
460 python toolbox to tackle the curse of imbalanced datasets in machine learn-
461 ing. *The Journal of Machine Learning Research*, 18(1):559–563.
- 462 Maack, M. M. (2019). ‘youtube recommendations are toxic,’ says dev who
463 worked on the algorithm.
- 464 Miller, D. T. (2021). Characterizing qanon: Analysis of youtube comments
465 presents new conclusions about a popular conservative conspiracy. *First*
466 *Monday*.



- 467 Neufeld, D. (2021). The 50 most visited websites in the world.
- 468 O’Callaghan, D., Greene, D., Conway, M., Carthy, J., and Cunningham, P.
469 (2013). The extreme right filter bubble. *arXiv preprint arXiv:1308.6149*.
- 470 Paolillo, J. C. (2018). The flat earth phenomenon on youtube. *First Monday*.
- 471 Pariser, E. (2011). *The filter bubble: What the Internet is hiding from you*.
472 Penguin UK.
- 473 Park, M., Naaman, M., and Berger, J. (2016). A data-driven study of view
474 duration on youtube. In *Proceedings of the International AAAI Conference*
475 *on Web and Social Media*, volume 10.
- 476 Reitermanova, Z. (2010). Data splitting. In *WDS*, volume 10, pages 31–36.
- 477 Rosenbaum, L. (2021). Escaping catch-22—overcoming covid vaccine hesitancy.
- 478 Roth, C., Mazières, A., and Menezes, T. (2020). Tubes and bubbles topological
479 confinement of youtube recommendations. *PloS one*, 15(4):e0231703.
- 480 Sokolova, M. and Lapalme, G. (2009). A systematic analysis of performance
481 measures for classification tasks. *Information processing & management*,
482 45(4):427–437.
- 483 Specht, D. F. (1990). Probabilistic neural networks. *Neural networks*, 3(1):109–
484 118.
- 485 Sun, Y., Kamel, M. S., and Wang, Y. (2006). Boosting for learning multiple
486 classes with imbalanced class distribution. In *Sixth international conference*
487 *on data mining (ICDM’06)*, pages 592–602. IEEE.
- 488 YouTube (2021). Youtube community guidelines & policies - how youtube works.
- 489 Zhou, R., Khemmarat, S., and Gao, L. (2010). The impact of youtube rec-
490 ommendation system on video views. In *Proceedings of the 10th ACM*
491 *SIGCOMM conference on Internet measurement*, pages 404–410.

Appendices

A Hyperparameter tuning



Ensemble	Activation	Layers	Neurons	Accuracy	Precision	Recall	F1
svm, nn, knn	logistic	1	20	0.939318	0.948923	0.931204	0.93998
svm, nn, knn	relu	1	20	0.939318	0.948923	0.931204	0.93998
svm, nn, knn	identity	1	1	0.939318	0.948923	0.931204	0.93998
svm, nn, knn	identity	10	1	0.939318	0.948923	0.931204	0.93998
svm, nn, knn	logistic	1	10	0.939318	0.948923	0.931204	0.93998
ridge, svm, nn, knn	tanh	10	1	0.939318	0.948923	0.931204	0.93998
svm, nn, knn	tanh	1	10	0.939318	0.948923	0.931204	0.93998
svm, nn, knn	tanh	1	20	0.939318	0.948923	0.931204	0.93998
svm, nn, logr, knn	identity	1	1	0.939318	0.948923	0.931204	0.93998
svm, nn, logr, knn	identity	1	10	0.939318	0.948923	0.931204	0.93998

(a) Ensemble.

Kernel	C	Accuracy	Precision	Recall	F1
rbf	10.0	0.936309	0.945473	0.928747	0.937035
rbf	100.0	0.935557	0.942289	0.930713	0.936465
rbf	1.0	0.925276	0.930163	0.922850	0.926492
poly	10.0	0.916499	0.946017	0.886978	0.915547
poly	100.0	0.915246	0.944940	0.885504	0.914257
linear	1.0	0.913741	0.917944	0.912531	0.915229
poly	1.0	0.909729	0.935065	0.884521	0.909091
linear	10.0	0.904965	0.907882	0.905651	0.906765
linear	100.0	0.898195	0.905830	0.893366	0.899555
rbf	0.1	0.878887	0.878906	0.884521	0.881705

(b) Support-vector machine.

Activation	Layers	Neurons	Accuracy	Precision	Recall	F1
identity	10	10	0.923019	0.935484	0.912039	0.923613
identity	25	10	0.921013	0.933031	0.910565	0.921661
relu	10	10	0.919007	0.917561	0.924324	0.920930
identity	10	20	0.916750	0.906056	0.933661	0.919652
relu	10	20	0.915998	0.912221	0.924324	0.918233
tanh	10	10	0.915747	0.914592	0.920885	0.917728
relu	1	1	0.915747	0.931876	0.900737	0.916042
tanh	25	20	0.915496	0.919052	0.914988	0.917016
tanh	10	20	0.914744	0.920178	0.912039	0.916091
logistic	1	1	0.913741	0.925516	0.903686	0.914470

(c) Neural network.



Solver	Alpha	Accuracy	Precision	Recall	F1
auto	0.1	0.918506	0.919118	0.921376	0.920245
sparse_cg	0.1	0.918506	0.919118	0.921376	0.920245
sag	0.1	0.918255	0.919902	0.919902	0.919902
auto	1.0	0.917252	0.923497	0.913514	0.918478
sparse_cg	1.0	0.917252	0.923497	0.913514	0.918478
sag	1.0	0.917252	0.923497	0.913514	0.918478
sag	10.0	0.878385	0.893002	0.865356	0.878962
auto	10.0	0.878134	0.892549	0.865356	0.878743
sparse_cg	10.0	0.878134	0.892549	0.865356	0.878743
auto	100.0	0.812437	0.854545	0.762162	0.805714

(a) Ridge regression.

Penalty	C	Solver	Accuracy	Precision	Recall	F1
l2	20	newton-cg	0.918506	0.924107	0.915479	0.919773
l2	20	saga	0.918506	0.924107	0.915479	0.919773
l2	20	sag	0.918506	0.924107	0.915479	0.919773
l2	10	sag	0.916249	0.920831	0.914496	0.917653
l2	10	newton-cg	0.916249	0.920831	0.914496	0.917653
l2	10	saga	0.916249	0.920831	0.914496	0.917653
l2	10	lbfgs	0.915747	0.919506	0.914988	0.917241
l2	20	lbfgs	0.914744	0.921432	0.910565	0.915966
none	1	sag	0.913992	0.923848	0.906143	0.914909
none	10	saga	0.913240	0.922461	0.906143	0.914229

(b) Logistic regression.

K	Accuracy	Precision	Recall	F1
1	0.889669	0.888456	0.896314	0.892368
3	0.888415	0.882212	0.901720	0.891859
4	0.879137	0.908899	0.848157	0.877478
5	0.873370	0.858482	0.900246	0.878868
6	0.873119	0.882441	0.866830	0.874566
2	0.872618	0.935043	0.806388	0.865963
7	0.868355	0.848891	0.902703	0.874970
8	0.867603	0.867382	0.874201	0.870778
9	0.861585	0.835672	0.907125	0.869934
10	0.859579	0.854397	0.873710	0.863946

(c) K-nearest neighbors.

Table 3: Results of the classifiers on the validation set with different hyperparameters. The best score per metric is written in bold.