

054 filtering techniques. These involved looking for certain keywords in the website and searching for
055 associated websites [3]. Machine learning methods began to be applied including analysis of linked
056 documents and embedded URLs.

057 Text classification has many different applications, from suggesting keywords for documents, to
058 classically detecting spam emails. It has been applied in various ways for content filtering where the
059 problem is modelled as a binary classification problem. This approach then makes use of the vast
060 research in classification algorithms in machine learning.

061 The Reddit website is a website designed for user-generated content with over eight million regular
062 users. A user can post a topic with a title, and either a link to a webpage or a short piece of text to
063 elicit discussion [14]. A user can also tag their post as NSFW or later another user can do so. The
064 website receives a huge amount of new volume constantly. This means that not every NSFW post
065 will be flagged appropriately either because the posting user decided not to, or it was only viewed by
066 a small number of users who also failed to flag it. A user can click on a link unaware of where it will
067 take them, and be presented with content that could cause disciplinary action in many workplaces.
068 Due to the huge volume of traffic, it is important that the classification can be successful from only
069 the details in the Reddit post and not necessarily content at the associated URL. Therefore a fast and
070 accurate algorithm for detecting and flagging NSFW posts would be a greatly valuable addition to
071 the Reddit infrastructure.

072 This paper examines text classification techniques and appropriate machine learning classifiers ap-
073 plied to this binary classification problem. Firstly various feature extraction methods are tested based
074 on the 'bag of words' concept to generate feature vectors for each post. This involves examining
075 which fields of a Reddit post are the most important for successful classification. Then several dif-
076 ferent binary classifiers are examined the various benefits. In the end we present a reliable method
077 for NSFW classification of new Reddit posts.

078

079 **1.2 Related Work**

080

081 Many existing commercial software tools including NetNanny [8] have been developed for filtering
082 internet content. These use a mix of heuristics and internally administered black-lists for blocking
083 content such as [1]. Early approaches [7] for examining the content of website included examining
084 linked pages, the content of images and specific keywords. Newer machine learning approaches use
085 a full text and image classification strategy on the content of the website such as [11] and [2].

086 Most text classification problems involve large corpora such as full news articles. The recent rise
087 of Twitter has changed researchers focuses to using text classification techniques on short text se-
088 quences of less than 50 words. Much research has involved sentiment prediction from messages on
089 Twitter such as [6]. Other researchers have specifically examined Reddit as an interesting source of
090 predictive machine learning data such as [13] which attempted to predict post popularity. However
091 the area of filtering for inappropriate content using short texts remain an interesting field of research.

092

093 **2 Testing Methodology**

094

095 **2.1 Data**

096

097 Selecting the appropriate data was very important for the study. A Reddit scraper was created to
098 pull all posts with associated meta-information from the Reddit website. The scraper collected posts
099 between 21st December 2012 and 8th January 2013. The scraper was executed almost three months
100 later in early April. This allowed the Reddit community sufficient time to have viewed the posts and
101 tagged for NSFW if appropriate. Over 1 million reddit posts were scraped during those dates.

102 Reddit data can also be scored by users with either an up-vote or a down-vote. This allowed for a
103 metric of a minimum number of users who had viewed the page. By thresholding posts that had a
104 minimum number of total votes, a more robust set of data was generated with a higher likelihood that
105 posts had been appropriately tagged as NSFW. Figure 2 shows the quality of the classifier increase
106 as the minimum total votes is increases which shows the theory. If the threshold was set too high,
107 a large proportion of posts would be omitted. After analysis of the data, the threshold was decided

108
109
110
111
112
113
114
115
116
117
118
119
120
121
122
123
124
125
126
127
128
129
130
131
132
133
134
135
136
137
138
139
140
141
142
143
144
145
146
147
148
149
150
151
152
153
154
155
156
157
158
159
160
161

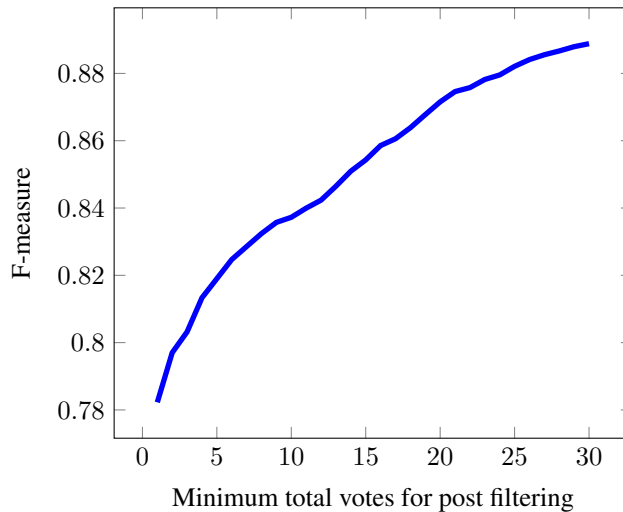


Figure 2: More popular posts are more likely to be correctly annotated, thereby creating a more robust test set

to be the median of the data which was 10. This meant that only posts of above average popularity would be tested, and thereby give a more reliable dataset.

500,000 posts with this threshold were selected. Only 8.6% of these posts were tagged as NSFW which was representative of the normal data. This huge class skew creates an additional challenge in both classification and also proper testing.

2.2 Testing

Because of the high class skew, a basic accuracy metric would not be appropriate. This is due to the simple idea that if a classifier tags all posts as Safe for Work, it would be 95% accurate as only 5% of posts are NSFW. It is most important that posts that are NSFW are tagged appropriately, as the potential cost to a user clicking on incorrectly labelled Safe for Work material could be large. On the other hand, Safe for Work posts cannot be tagged NSFW to often to affect the quality of posts on Reddit.

$$F_{measure} = 2 \frac{precision * recall}{precision + recall}$$

The F-measure (shown above), trades off precision and recall (also known as sensitivity). Due to the need to balance the success of positive classifications as well as negative classifications, we selected the F-measure as our target metric.

Furthermore the classifiers were testing using 5-fold cross-validation. This meant that in each test-case, the training data had 400,000 posts and the test data had 100,000 posts.

3 Features

3.1 Extraction

Several properties of each Reddit post were captured. The most visible to the user on Reddit is the title and would be a key indicator for the metric. The subreddit, which is the category of the post decided by the user, and the author's username of the post were also captured.

In order to extract numerical vectors from these text features, the 'bag of words' concept was used. This method finds each unique word in the data set, and then for each post counts the number of each word present to generate a very large and normally very sparse feature vector. The idea for 'bag of

162
163
164
165
166
167
168
169
170
171
172
173
174
175
176
177
178
179
180
181
182
183
184
185
186
187
188
189
190
191
192
193
194
195
196
197
198
199
200
201
202
203
204
205
206
207
208
209
210
211
212
213
214
215

words' comes from the simple idea that an email with the word 'viagra' anywhere in it would more than likely be spam.

The 'bag of words' method has been adapted in several ways. Firstly, the technique can be extended to bigrams. This method searches for all word pairings instead of single words and is able to add more contextual information to the feature vector. It has been suggested that often bi-grams are enough to capture the core concept of phrases, e.g. 'United States of America' can be captured by the bigram 'United States'. Therefore the use of tri-grams rarely gives additional gain. This was tested and is also shown in Figure 3.

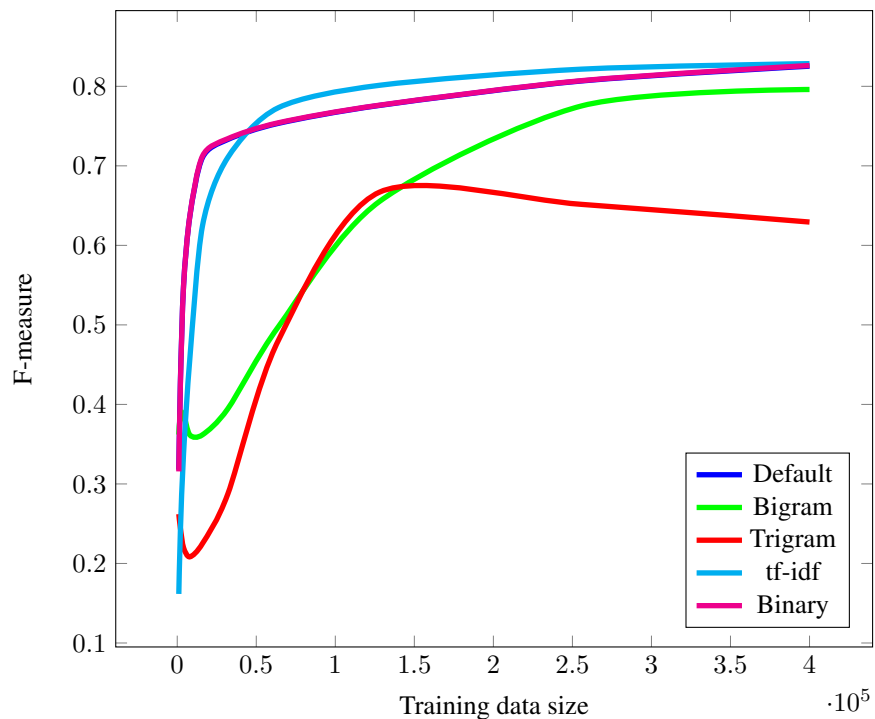


Figure 3: A comparison of different text extraction techniques

The frequency of words and varied text lengths can also skew the feature vector. The tf-idf method normalises the data for length of text and frequency of common words [9]. This is beneficial to this problem so that certain words are not over-weighted and the result of this is shown in Figure 3.

A binary feature extractor was also tested. This would only test for a word appearance and not count word occurrence. In the figure, the binary classifier overlaps closely with the default single word tokenizer and gives no additional benefit. It should be noted that in all these cases, the lower-case text was used. Also stop-words (common words in the English language) and punctuation were also removed.

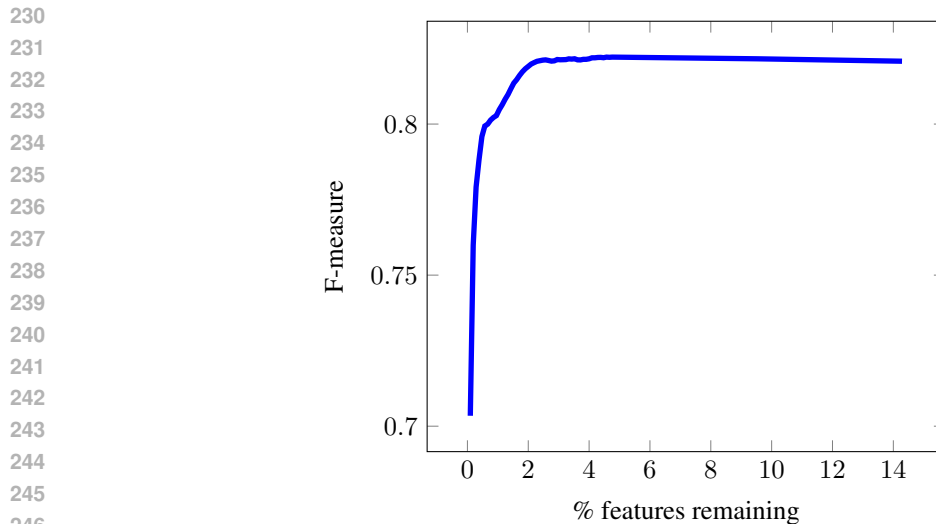
After cross-validation the tf-idf tokenizer was selected as the best feature extractor. Because of the very large data-sets used, the feature hashing approach was applied [15]. This uses the hash of words to calculate the column ID and does not require on a dictionary reducing the memory and computation requirements. The extraction methods were implemented using the Scikit Learn libraries [12].

3.2 Selection

The 'bag of words' methodology creates huge numbers of features. Specifically using our data set, a basic single word 'bag of words' algorithm creates over one million features.

216 Traditionally feature selection is an excellent way to prune the number of features to a manageable
217 level. This is beneficial for several reasons. Often classifiers will not be as successful with a very
218 high number of classifiers. This is due to the very high dimensionality of the problem, and the chal-
219 lenge of creating a relevant fit around the given data. It can also be useful for greater understanding
220 in the problem to identify which features are important in the classification and which are not. How-
221 ever it is more challenging in text classification. This is both due to the extremely large number of
222 vectors and also the incredible sparseness of the data.

223 It has been noted that the results of feature selection in a text classification study vary [4] and de-
224 pends heavily on the data-set used. We tested the chi-squared technique for feature selection in
225 order to reduce the number of features. The results are shown in Figure 4. These interesting results
226 show the reliance on a small subset of the features for the majority of successful classifications.
227 Further analysis using a non-hashing vectorizer revealed that a significant proportion of the remain-
228 ing features were related to subreddits. This shows that subreddits are very important in successful
229 classification.



248 Figure 4: Results of chi-squared removal of different proportions of the feature space

251 4 Classification

252 Various binary classifiers were tested on the dataset given using the reduced feature set. Cross-
253 validation was used to adjust the various parameters for the best results for this data set. The results
254 of the different optimised strategies outlined in this section are shown in Figure 5.

255 The Multinomial Naive bayes method and Bernoulli Naive Bayes method use the basic probabilities
256 of word occurrences as well as the class frequencies to calculate the probability that a post is NSFW.
257 These techniques used the Scikit Learn python library [12].

258 Neural Networks mimic the behaviour of neurons in the human brain. Each neuron takes in multiple
259 inputs and only fires (giving a particular output) when certain constraints are met on the inputs. A
260 multi-layer network of neurons is built, where each input feature is linked to a neuron and neurons
261 are interlinked on several layers until a single output is given. This output decides whether the
262 feature data given should be classified as SFW or NSFW. Support vector machines, a method for
263 splitting a data-set by transforming data and splitting with a hyperplane, and logistic regression are
264 also tested.

265 The Vowpal Wabbit tool [10] from Yahoo Research and Microsoft Research was used to test logistic
266 regression, support vector machine and neural networks approaches on the very large data set
267 used. Using cross-validation the size of the hidden layer in the neural network (i.e. the number of
268 additional neurons between the inputs and the output used) is adjusted for optimal results.
269

270
271
272
273
274
275
276
277
278
279
280
281
282
283
284
285
286
287
288
289
290
291
292
293
294
295
296
297
298
299
300
301
302
303
304
305
306
307
308
309
310
311
312
313
314
315
316
317
318
319
320
321
322
323

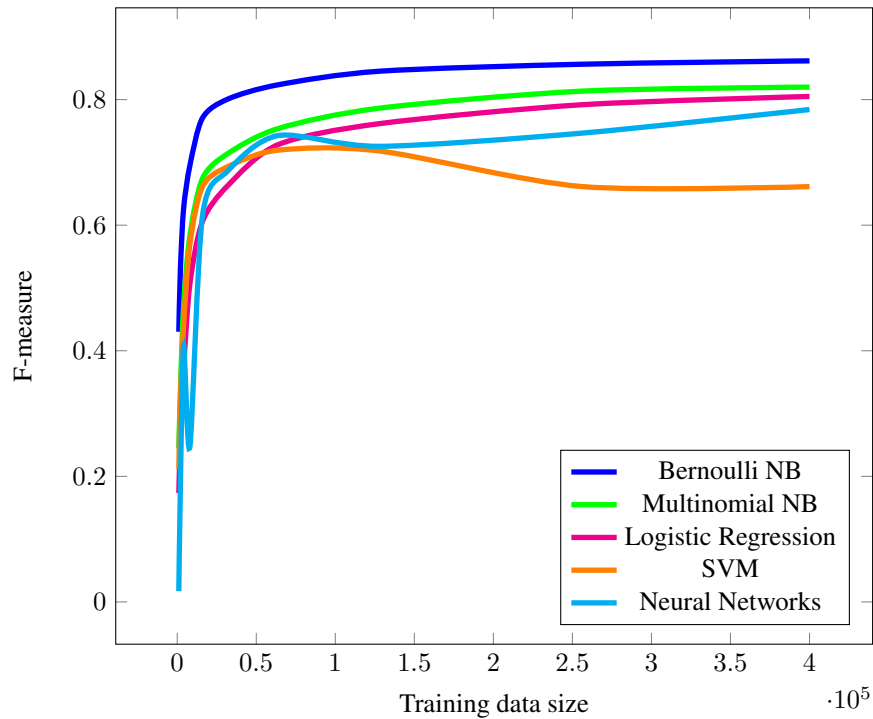


Figure 5: Results of each classifier using tf-idf feature extractor

Classifier	
Feature extractor	tf-idf
Feature selection limit	22000
Model	Bernoulli NB
Results	
True Positive	35088
False Positive	6636
True Negative	453650
False Negative	4626
Sensitivity	0.883
Specificity	0.986
F-Measure	0.862

Table 1: Detailed results for the best classifier for full data size

5 Results and Discussion

As had been shown the use of tf-idf tokenizer had given the best set of features for this data set. This feature set was further pruned using chi-squared-based feature selection. The results in Figure 5 show that the Bernoulli Naives Bayes implementation gives the highest success. The specific results from this classifier are shown in Table 1.

This set of results together show the significance of the subreddit as a key predictor of a post’s NSFW/SFW class. This is re-inforced by the evidence that bigrams do not offer improved performance, which suggests that the single words of the subreddit name are more important than word

324
325
326
327
328
329
330
331
332
333
334
335
336
337
338
339
340
341
342
343
344
345
346
347
348
349
350
351
352
353
354
355
356
357
358
359
360
361
362
363
364
365
366
367
368
369
370
371
372
373
374
375
376
377

Post Titles

My friends friend had an accident
Totally legit
Unimpressed Kitchen
So I googled reddit and down votes and this came up.
I don't think they are the reason why it sank ..
Somebody left this book at work

Table 2: Examples of false negative posts with ambiguous titles

combinations inside the title. This was also proposed by the large number of features that could be removed through feature selection. Interestingly a simpler Naive Bayes model gives better results than SVMs and neural networks which suggest that in this data-set these more complex models suffer from over-fitting.

5.1 The Challenging Subset

The results from the various classifier highlight that there remains a subset of posts that are difficult to classify. A cross-analysis of these posts shows several reasons for their challenging classification.

An initial look at the post titles highlight that many are linguistically ambiguous and would be very difficult for a human to identify the possible content. Some examples of the titles are shown in Table 5.1. Furthermore an analysis of the subreddit and author of the posts highlight the difficulty in using these metrics. A large proportion of these posts are from deleted users which causes the author name to be "[Deleted]". This author name has over 9000 posts attributed to it, 22% of which are tagged NSFW. Because of the large variability and high proportion of this author name, the author features become significantly less predicitive of NSFW posts. Many posts are also from subreddits with no clearly defined bias towards SFW or NSFW which makes the subreddit a more limiting predictor in these cases.

It may be possible to identify these challenging posts through the same metrics and flag them for further analysis. Then a deeper analysis of the linked URL or linked images could be done for a better classification. This would be interesting area of further research to improve the quality of classification for this challenging subset of data.

5.2 Title Only Classification

In order that this classification system could be used on a more diverse data set than just Reddit posts, we tested whether only using the title of the Reddit post would be sufficient to successfully classify posts. The same method of feature detection was used on the title only. This caused the F-measure to drop to 0.42.

After analysis of the failing posts, it can be shown that the problem with the challenging subset has been enlarged. The text in the title can contain very ambiguous language which causes the much lower success rate of the classifier. With the given success rate, the classifier could not be used a reliable metric for tagging posts. Furthermore the author and subreddit fields are certainly very important to the success of the classifier.

6 Conclusion

This paper introduces an effective classifier for NSFW posts on Reddit. It tests the various features that may be used and showed that the most effective results were gained from using a Bernouilli Naives Bayes classifier with a tf-idf feature extractor. While the sensitivity and specificity is not high enough to surpass human intervention, the classifer could be used as an excellent complement to suggest tags for posts.

378
379
380
381
382
383
384
385
386
387
388
389
390
391
392
393
394
395
396
397
398
399
400
401
402
403
404
405
406
407
408
409
410
411
412
413
414
415
416
417
418
419
420
421
422
423
424
425
426
427
428
429
430
431

References

- [1] Brenda S Baker and Eric Grosse. Local control over filtered www access. In *Proceedings of the 4th World Wide Web Conference*, 1995.
- [2] Thomas Deselaers, Lexi Pimenidis, and Hermann Ney. Bag-of-visual-words models for adult image classification and filtering. In *Pattern Recognition, 2008. ICPR 2008. 19th International Conference on*, pages 1–4. IEEE, 2008.
- [3] Rongbo Du, Reihaneh Safavi-Naini, and Willy Susilo. Web filtering using text classification. In *Networks, 2003. ICON2003. The 11th IEEE International Conference on*, pages 325–330. IEEE, 2003.
- [4] George Forman. Feature selection for text classification. *Computational methods of feature selection*, pages 257–276, 2008.
- [5] Wordle generating tool. Net Nanny: content-control software, <http://www.wordle.com>, 2013.
- [6] Alec Go, Richa Bhayani, and Lei Huang. Twitter sentiment classification using distant supervision. *CS224N Project Report, Stanford*, pages 1–12, 2009.
- [7] Paul Greenfield, Peter Rickwood, Huu Cuong Tran, and Australian Broadcasting Authority. *Effectiveness of Internet filtering software products*. Australian Broadcasting Authority, 2001.
- [8] ContentWatch Inc. Net Nanny: content-control software, <http://www.netnanny.com>, 2013.
- [9] Thorsten Joachims. A probabilistic analysis of the rocchio algorithm with tfidf for text categorization. Technical report, DTIC Document, 1996.
- [10] J Langford, L Li, and A Strehl. Vowpal wabbit online learning project, 2007.
- [11] Pui Y Lee, Siu C Hui, and Alvis Cheuk M Fong. Neural networks for web content filtering. *Intelligent Systems, IEEE*, 17(5):48–57, 2002.
- [12] F. Pedregosa, G. Varoquaux, A. Gramfort, V. Michel, B. Thirion, O. Grisel, M. Blondel, P. Prettenhofer, R. Weiss, V. Dubourg, J. Vanderplas, A. Passos, D. Cournapeau, M. Brucher, M. Perrot, and E. Duchesnay. Scikit-learn: Machine Learning in Python . *Journal of Machine Learning Research*, 12:2825–2830, 2011.
- [13] Jordan Segall and Alex Zamoshchin. Twitter sentiment classification using distant supervision. *CS229 Project Report, Stanford*.
- [14] Troy Steinbauer. Information and social analysis of reddit. *CS224N Project Report, Stanford*.
- [15] Kilian Weinberger, Anirban Dasgupta, John Langford, Alex Smola, and Josh Attenberg. Feature hashing for large scale multitask learning. In *Proceedings of the 26th Annual International Conference on Machine Learning*, pages 1113–1120. ACM, 2009.