

## Removing Spam from Web Corpora Through Supervised Learning Using FastText

Vít Suchomel

Natural Language Processing Centre

Faculty of Informatics, Masaryk University, Brno, Czech Republic

xsuchom2@fi.muni.cz

### Abstract

Unlike traditional text corpora collected from trustworthy sources, the content of web based corpora has to be filtered. This study briefly discusses the impact of web spam on corpus usability and emphasizes the importance of removing computer generated text from web corpora.

The paper also presents a keyword comparison of an unfiltered corpus with the same collection of texts cleaned by a supervised classifier trained using FastText. The classifier was able to recognise 71 % of web spam documents similar to the training set but lacked both precision and recall when applied to short texts from another data set.

### 1 Web Spam in Text Corpora

It has been shown that boilerplate, duplicates, and spam skew corpus based analyses and therefore have to be removed, see nonsense examples of word use in an application for English learners based on a web corpus in figure 1. While the first two issues have been successfully addressed, e.g. by (Marek et al., 2007; Pomikálek, 2011; Versley and Panchenko, 2012; Schäfer and Bildhauer, 2013), spam might be still observed in web corpora as reported by (Kilgarriff and Suchomel, 2013). It was spam that represented the main difference between their 2008 and 2012 corpora crawled from the web. That is why a spam cleaning stage should be a part of the process of building web corpora.

The traditional definition of web spam is *actions intended to mislead search engines into ranking some pages higher than they deserve* (Gyöngyi and Garcia-Molina, 2005). The Google document ‘Fighting Spam’<sup>1</sup> describes the kinds of spam that

<sup>1</sup><https://www.google.com/insidesearch/>

The screenshot shows the SKELL search engine interface. At the top, there is a search bar with the word 'money' entered and a green 'Search' button. Below the search bar, there are three tabs: 'Examples', 'Word sketch', and 'Similar words'. The 'Examples' tab is selected. Below the tabs, the word 'money' is displayed in a large font, followed by '239.0 hits per million'. A list of 15 example sentences is shown, with the word 'money' highlighted in red in each sentence. The sentences are numbered 1 through 15. The first two sentences are: 1. The savings money account is really weak. 2. Car leasing dealer s money factor car leasing. The tenth sentence is: 10. Australian interest rates bank interest rates australia money interest.

Figure 1: Web spam in examples of use of word ‘money’ at [skell.sketchengine.co.uk](http://skell.sketchengine.co.uk) – see lines 2, 4 and 10.

Google finds, and what they do about it.

Text alteration techniques consist in changing the frequency properties of a web page content in favour of spam targeted words or phrases: *repetition of terms related to the spam campaign target, inserting a large number of unrelated terms, often even entire dictionaries, weaving of spam terms into contents copied from informative sites, e.g. news articles, glueing together sentences or phrases from different sources* as reported by (Gyöngyi and Garcia-Molina, 2005).

Automatically generated content does not provide examples of authentic use of a natural language. Nonsense, incoherent or any unnatural texts such as the following short instance have to be removed from a good quality web corpus: *Edmonton Oilers rallied towards get over the Montreal Canadiens 4-3 upon Thursday.Ryan Nugent-Hopkins completed with 2 aims, together with*

[howsearchworks/fighting-spam.html](http://howsearchworks.com/fighting-spam.html)

*the match-tying rating with 25 seconds remaining within just legislation.*<sup>2</sup>

The following types of automatically generated content are examples of documents penalised by Google:<sup>3</sup> *Text translated by an automated tool without human review or curation before publishing. Text generated through automated processes, such as Markov chains. Text generated using automated synonymizing or obfuscation techniques.* These kinds of spam should certainly be eliminated from web corpora while the other two examples given by Google may not present a harm to the corpus use: *Text generated from scraping Atom/RSS feeds or search results. Stitching or combining content from different web pages without adding sufficient value.*

In contrast to the traditional or search engine definitions of web spam, the corpus use point of view is not concerned with intentions of spam producers or the justification of the search engine optimisation of a web page. A text corpus built for NLP or linguistics purpose should contain coherent and consistent, meaningful, natural and authentic sentences in the target language. Only texts created by spamming techniques breaking those properties should be detected and avoided. The unwanted non-text is this: computer generated text, machine translated text, text altered by keyword stuffing or phrase stitching, text altered by replacing words with synonyms using a thesaurus, summaries automatically generated from databases (e.g. stock market reports, weather forecast, sport results – all of the same kind very similar), and finally any incoherent text. Varieties of spam removable by existing tools, e.g. duplicate content, link farms (quite a lot of links with scarce text), are only a minor problem.

Avoiding web spam by selecting trustworthy corpus sources such as Wikipedia, news sites, government and academic webs works well: (Baisa and Suchomel, 2014) show it is possible to construct medium sized corpora from URL whitelists and web catalogues. (Spoustová and Spousta, 2012) used a similar way of building a Czech web corpus. Also the BootCaT method (Baroni and Bernardini, 2004) indirectly avoids spam by relying on a search engine to find non-spam data. Despite the avoiding methods being successful, it is doubtful a huge web collection can be obtained

<sup>2</sup><http://masterclasspolska.pl/forum/>

<sup>3</sup>Google quality guidelines – <https://support.google.com/webmasters/answer/2721306>

just from trustworthy sources.

Furthermore, language independent methods of combating spam might be of use. (Ntoulas et al., 2006) reported web spamming was not only a matter of the English part of internet. Spam was found in their French, German, Japanese and Chinese documents as well.

## 2 Removing Spam Using a Supervised Classifier

This section describes training and evaluation of a supervised classifier to detect spam in web corpora.

We have manually annotated a collection of 1630 web pages from various web sources from years 2006 to 2015.<sup>4</sup> To cover the main topics of spam texts observed in our previously built corpora, we included 107 spam pages promoting medication, financial services, commercial essay writing and other subjects. Both phrase level and sentence level incoherent texts (mostly keyword insertions, n-grams of words stitched together or seemingly authentic sentences not conveying any connecting message) were represented. Another 39 spam documents coming from random web documents identified by annotators were included. There were 146 positive instances of spam documents altogether.

The classifier was trained using FastText (Joulin et al., 2016) and applied to a large English web corpus from 2015. The expected performance of the classifier was evaluated using a 30-fold cross-validation on the web page collection. Since our aim was to remove as much spam from the corpus as possible, regardless false positives, the classifier confidence threshold was set to prioritize recall over precision. The achieved precision and recall were 71.5 % and 70.5 % respectively. Applying this classifier to an English web corpus from 2015 resulted in removing 35 % of corpus documents still leaving enough data for the corpus use.

An inspection of the cleaned corpus revealed the relative count of usual spam related keywords dropped significantly as expected while general words not necessarily associated with spam were affected less as can be seen in table 1.

Another evaluation of the classifier was performed by manually checking 299 random web documents from the cleaned corpus and 25 ran-

<sup>4</sup>The collection is a part of another classification experiment by the same authors not covered by this paper.

dom spam documents removed by the classifier. The achieved precision was 40.0 % with the recall of 27.8 %. The error analysis showed the classifier was not able to recognise non-text rather than spam. 17 of 26 unrecognised documents were scientific paper references or lists of names, dates and places, i.e. *Submitted by Diana on 2013-09-25 and updated by Diana on Wed, 2013-09-25 08:32 or January 13, 2014 January 16, 2014 Gaithersburg, Maryland, USA*. Such web pages were not present in the training data since we believed it had been removed from the corpus sources by a boilerplate removal tool and paid attention to longer documents. Not counting these 17 non-text false negatives, the recall would reach 52.6 %.

To find out what was removed from the corpus, relative counts of lemmas<sup>5</sup> in the corpus were compared with the BNC<sup>6</sup> in figures 2 and 3. A list of lemmas in the web corpus with the most reduced relative lemma count caused by removing unwanted documents is presented in 4.

The inspection showed there were a lot of spam related words in the original web corpus and that spam words are no longer characteristic of the cleaned version of the corpus in comparison to the BNC.<sup>7</sup>

### 3 Conclusion

We view computer generated text as the main kind of spam decreasing the quality of web corpora. A classifier trained on spam documents was applied to remove unwanted content from a web corpus. Although the classifier significantly decreased the presence of spam related words in the corpus, it was not able to recognise short non-text documents. That remains to be addressed in the future.

### Acknowledgments

This work was partly supported by the Ministry of Education of CR within the LINDAT-Clarin project LM2015071. This publication was written with the support of the Specific University Research provided by the Ministry of Education, Youth and Sports of the Czech Republic.

<sup>5</sup>Corpora in the study were lemmatised by TreeTagger.

<sup>6</sup>The tokenisation of the BNC had to be changed to the same way the web corpus was tokenised in order to make the counts of tokens in both corpora comparable.

<sup>7</sup>The comparison with the BNC also revealed there are words related to the modern technology (e.g. *website, online, email*) and American English spelled words (*center, organization*) in the 2015 web corpus.

### References

- [Baisa and Suchomel2014] Vít Baisa and Vít Suchomel. 2014. Skell: Web interface for english language learning. In *Eighth Workshop on Recent Advances in Slavonic Natural Language Processing*, pages 63–70, Brno. Tribun EU.
- [Baroni and Bernardini2004] Marco Baroni and Silvia Bernardini. 2004. Bootcat: Bootstrapping corpora and terms from the web. In *LREC*.
- [Gyöngyi and Garcia-Molina2005] Zoltan Gyöngyi and Hector Garcia-Molina. 2005. Web spam taxonomy. In *First international workshop on adversarial information retrieval on the web (AIRWeb 2005)*.
- [Joulin et al.2016] Armand Joulin, Edouard Grave, Piotr Bojanowski, and Tomas Mikolov. 2016. Bag of tricks for efficient text classification. *arXiv preprint arXiv:1607.01759*.
- [Kilgarriff and Suchomel2013] Adam Kilgarriff and Vít Suchomel. 2013. Web spam. In Paul Rayson Stefan Evert, Egon Stemle, editor, *Proceedings of the 8th Web as Corpus Workshop (WAC-8) @ Corpus Linguistics 2013*, pages 46–52.
- [Marek et al.2007] Michal Marek, Pavel Pecina, and Miroslav Spousta. 2007. Web page cleaning with conditional random fields. In *Building and Exploring Web Corpora: Proceedings of the Fifth Web as Corpus Workshop, Incorporating CleanEval (WAC3), Belgium*, pages 155–162.
- [Ntoulas et al.2006] Alexandros Ntoulas, Marc Najork, Mark Manasse, and Dennis Fetterly. 2006. Detecting spam web pages through content analysis. In *Proceedings of the 15th international conference on World Wide Web*, pages 83–92. ACM.
- [Pomikálek2011] Jan Pomikálek. 2011. *Removing boilerplate and duplicate content from web corpora*. Ph.D. thesis, Masaryk University.
- [Schäfer and Bildhauer2013] Roland Schäfer and Felix Bildhauer. 2013. *Web Corpus Construction*, volume 6. Morgan & Claypool Publishers.
- [Spoustová and Spousta2012] Johanka Spoustová and Miroslav Spousta. 2012. A high-quality web corpus of Czech. In *LREC*, pages 311–315.
- [Versley and Panchenko2012] Yannick Versley and Yana Panchenko. 2012. Not just bigger: Towards better-quality web corpora. In *Proceedings of the seventh Web as Corpus Workshop (WAC7)*, pages 44–52.

Table 1: Comparison of the corpus before and after spam removal using the classifier. Corpus sizes and relative frequencies (number of occurrences per million words) of selected words are shown. Reducing the corpus to 55 % of the former token count, phrases strongly indicating spam documents such as “cialis 20 mg”, “payday loan” or “essay writing” were almost removed while innocent phrases from the same domains such as “oral administration”, “interest rate” or “pass the exam” were reduced proportionally to the whole corpus.

	Original corpus	Cleaned corpus	Kept in cleaned
<b>Document count</b>	58,438,034	37,810,139	64.7 %
<b>Token count</b>	33,144,241,513	18,371,812,861	55.4 %
“viagra”	229.71	3.42	0.8 %
“cialis 20 mg”	2.74	0.02	0.4 %
“aspirin”	5.63	1.52	14.8 %
“oral administration”	0.26	0.23	48.8 %
“loan”	166.32	48.34	16.1 %
“payday loan”	24.19	1.09	2.5 %
“cheap”	295.31	64.30	12.1 %
“interest rate”	14.73	9.80	36.7 %
“essay”	348.89	33.95	5.4 %
“essay writing”	7.72	0.32	2.3 %
“pass the exam”	0.34	0.36	59.4 %

Figure 2: Relative wordcount comparison of the original 2015 web corpus with British National Corpus, top 26 lemmas sorted by the keyword score.  $Score = \frac{f_{pm_1} + 100}{f_{pm_2} + 100}$  where  $f_{pm_1}$  is the count of lemmas per million in the focus corpus (3rd column) and  $f_{pm_2}$  is the count of lemmas per million in the reference corpus (5th column).

Lowercase lemma	Original English Web 2015		British National Corpus		Score
	frequency	frequency/mill	frequency	frequency/mill	
download	<a href="#">32,877,718</a>	992.0	<a href="#">35</a>	0.3	10.9
pdf	<a href="#">30,658,156</a>	925.0	<a href="#">37</a>	0.3	10.2
online	<a href="#">23,683,595</a>	714.6	<a href="#">596</a>	5.3	7.7
program	<a href="#">20,333,705</a>	613.5	<a href="#">5,814</a>	51.8	4.7
website	<a href="#">9,586,380</a>	289.2	0	0.0	3.9
center	<a href="#">9,903,586</a>	298.8	<a href="#">573</a>	5.1	3.8
essay	<a href="#">11,563,807</a>	348.9	<a href="#">2,317</a>	20.6	3.7
viagra	<a href="#">7,620,095</a>	229.9	0	0.0	3.3
url	<a href="#">7,168,836</a>	216.3	0	0.0	3.2
ebook	<a href="#">6,969,380</a>	210.3	0	0.0	3.1
web	<a href="#">7,206,520</a>	217.4	<a href="#">729</a>	6.5	3.0
internet	<a href="#">6,248,400</a>	188.5	<a href="#">97</a>	0.9	2.9
student	<a href="#">24,584,996</a>	741.8	<a href="#">22,133</a>	197.1	2.8
cialis	<a href="#">5,816,475</a>	175.5	0	0.0	2.8
blog	<a href="#">5,110,812</a>	154.2	0	0.0	2.5
email	<a href="#">5,074,946</a>	153.1	<a href="#">43</a>	0.4	2.5
cheap	<a href="#">9,787,744</a>	295.3	<a href="#">6,649</a>	59.2	2.5
epub	<a href="#">4,761,306</a>	143.7	0	0.0	2.4
video	<a href="#">10,278,042</a>	310.1	<a href="#">7,672</a>	68.3	2.4
free	<a href="#">20,406,767</a>	615.7	<a href="#">21,963</a>	195.6	2.4
u.s.	<a href="#">4,976,297</a>	150.1	<a href="#">458</a>	4.1	2.4
post	<a href="#">13,400,787</a>	404.3	<a href="#">12,576</a>	112.0	2.4
outlet	<a href="#">5,501,465</a>	166.0	<a href="#">1,375</a>	12.2	2.4
color	<a href="#">4,553,463</a>	137.4	<a href="#">143</a>	1.3	2.3
click	<a href="#">5,326,832</a>	160.7	<a href="#">1,273</a>	11.3	2.3
your	<a href="#">95,303,049</a>	2875.4	<a href="#">134,413</a>	1197.0	2.3

Figure 3: Relative wordcount comparison of the cleaned web corpus with British National Corpus

Lowercase lemma	<i>Cleaned English Web 2015</i>		<i>British National Corpus</i>		Score
	frequency	frequency/mill ☺	frequency	frequency/mill	
program	<a href="#">14,384,115</a>	782.9	<a href="#">5,814</a>	51.8	5.8
center	<a href="#">7,509,618</a>	408.8	<a href="#">573</a>	5.1	4.8
website	<a href="#">4,792,518</a>	260.9	0	0.0	3.6
student	<a href="#">16,973,541</a>	923.9	<a href="#">22,133</a>	197.1	3.4
online	<a href="#">4,753,580</a>	258.7	<a href="#">596</a>	5.3	3.4
u.s.	<a href="#">4,225,425</a>	230.0	<a href="#">458</a>	4.1	3.2
project	<a href="#">14,949,773</a>	813.7	<a href="#">21,742</a>	193.6	3.1
university	<a href="#">12,182,707</a>	663.1	<a href="#">18,899</a>	168.3	2.8
community	<a href="#">15,164,485</a>	825.4	<a href="#">26,564</a>	236.6	2.7
global	<a href="#">4,585,347</a>	249.6	<a href="#">3,529</a>	31.4	2.7
web	<a href="#">3,322,320</a>	180.8	<a href="#">729</a>	6.5	2.6
download	<a href="#">3,011,631</a>	163.9	<a href="#">35</a>	0.3	2.6
email	<a href="#">2,901,189</a>	157.9	<a href="#">43</a>	0.4	2.6
dr.	<a href="#">3,290,385</a>	179.1	<a href="#">1,215</a>	10.8	2.5
internet	<a href="#">2,753,028</a>	149.9	<a href="#">97</a>	0.9	2.5
our	<a href="#">39,914,081</a>	2172.6	<a href="#">93,457</a>	832.3	2.4
click	<a href="#">3,144,338</a>	171.2	<a href="#">1,273</a>	11.3	2.4
focus	<a href="#">6,345,601</a>	345.4	<a href="#">9,538</a>	84.9	2.4
technology	<a href="#">7,397,599</a>	402.7	<a href="#">12,865</a>	114.6	2.3
organization	<a href="#">5,944,514</a>	323.6	<a href="#">9,240</a>	82.3	2.3
research	<a href="#">12,854,262</a>	699.7	<a href="#">27,567</a>	245.5	2.3
update	<a href="#">3,452,461</a>	187.9	<a href="#">2,814</a>	25.1	2.3
datum	<a href="#">7,682,640</a>	418.2	<a href="#">14,212</a>	126.6	2.3
network	<a href="#">5,810,016</a>	316.2	<a href="#">9,291</a>	82.7	2.3
video	<a href="#">5,202,487</a>	283.2	<a href="#">7,672</a>	68.3	2.3
photo	<a href="#">3,054,229</a>	166.2	<a href="#">2,036</a>	18.1	2.3

Figure 4: Relative wordcount comparison of the original web corpus with the cleaned version

Lowercase lemma	<i>Original English Web 2015</i>		<i>Cleaned English Web 2015</i>		Score
	frequency	frequency/mill ☺	frequency	frequency/mill	
pdf	<a href="#">30,658,156</a>	925.0	<a href="#">1,851,347</a>	100.8	5.1
download	<a href="#">32,877,718</a>	992.0	<a href="#">3,011,631</a>	163.9	4.1
essay	<a href="#">11,563,807</a>	348.9	<a href="#">623,760</a>	34.0	3.4
viagra	<a href="#">7,620,095</a>	229.9	<a href="#">62,899</a>	3.4	3.2
ebook	<a href="#">6,969,380</a>	210.3	<a href="#">265,781</a>	14.5	2.7
cialis	<a href="#">5,816,475</a>	175.5	<a href="#">45,613</a>	2.5	2.7
url	<a href="#">7,168,836</a>	216.3	<a href="#">509,596</a>	27.7	2.5
buy	<a href="#">17,364,124</a>	523.9	<a href="#">2,867,958</a>	156.1	2.4
cheap	<a href="#">9,787,744</a>	295.3	<a href="#">1,180,506</a>	64.3	2.4
online	<a href="#">23,683,595</a>	714.6	<a href="#">4,753,580</a>	258.7	2.3
epub	<a href="#">4,761,306</a>	143.7	<a href="#">203,405</a>	11.1	2.2
prescription	<a href="#">4,646,919</a>	140.2	<a href="#">280,013</a>	15.2	2.1
outlet	<a href="#">5,501,465</a>	166.0	<a href="#">651,024</a>	35.4	2.0
book	<a href="#">29,921,305</a>	902.8	<a href="#">7,889,796</a>	429.5	1.9
generic	<a href="#">3,594,096</a>	108.4	<a href="#">257,090</a>	14.0	1.8
ugg	<a href="#">3,022,464</a>	91.2	<a href="#">93,591</a>	5.1	1.8
loan	<a href="#">5,512,504</a>	166.3	<a href="#">888,181</a>	48.3	1.8
jersey	<a href="#">4,873,552</a>	147.0	<a href="#">836,729</a>	45.5	1.7
insurance	<a href="#">7,150,681</a>	215.7	<a href="#">1,588,816</a>	86.5	1.7
pharmacy	<a href="#">2,941,876</a>	88.8	<a href="#">290,211</a>	15.8	1.6
sex	<a href="#">6,452,251</a>	194.7	<a href="#">1,502,817</a>	81.8	1.6
de	<a href="#">10,572,331</a>	319.0	<a href="#">2,986,557</a>	162.6	1.6
mg	<a href="#">2,776,320</a>	83.8	<a href="#">298,031</a>	16.2	1.6
you	<a href="#">195,234,032</a>	5890.4	<a href="#">68,409,350</a>	3723.6	1.6
binary	<a href="#">4,226,875</a>	127.5	<a href="#">839,475</a>	45.7	1.6
levitra	<a href="#">1,873,011</a>	56.5	<a href="#">34,646</a>	1.9	1.5