

Could language be the key to detecting fake news?

David Shariatmadari, Mon 2 Sep 2019 09.00, theguardian.com

Purveyors of disinformation can be caught out by the particular words they use, according to new research

5 The internet represents the biggest explosion of data in human history. There's more out there, and more access to it than ever before. The information ecosystem is a bit like a tropical rainforest: luxuriant, dense and fiercely competitive. As such, it contains its fair share of predators and poisonous plants.

10 Deliberately misleading articles, websites and social media posts can come about for lots of different reasons: they might be trying to influence elections or policies; they might represent a form of cyberwarfare between states; they might be aimed at raising someone's profile and influence, or discrediting their opponents. Or they might simply be about making money, relying on the attention-grabbing nature of outrageous lies to generate ad revenue, as in the case of the "digital gold rush" that saw a small Macedonian town register more than 150 pro-Trump websites during the 2016 presidential
15 race.

One thing they may have in common, however, is the language they use.

20 Having a reliable way of identifying fake news is important. The whole reason it's a problem is that it mimics reliable reporting – and people can't always tell the difference. That's why, for the past few years, researchers have been trying to work out what the linguistic characteristics of fake news are. Computers that are fed material already classified as misleading are able to identify patterns in the language used. They're then able to apply that knowledge to new material, and flag it as potentially dubious.

25 One such project, led by Fatemeh Torabi Asr at Simon Fraser University in Canada, recently found that "on average, fake news articles use more ... words related to sex, death and anxiety". "Overly emotional" language is often deployed. In contrast, "Genuine news ... contains a larger proportion of words related to work (business) and money (economy)."

30 Another group of researchers analysed the relationship of various grammatical categories to fake news. They concluded that words which can be used to exaggerate are all found more often in deliberately misleading sources. These included superlatives, like "most" and "worst", and so-called subjectives, like "brilliant" and "terrible". They noted that propaganda tends to use abstract generalities like "truth" and "freedom", and intriguingly showed that use of the second-person pronoun "you" was closely linked to fake news.

35 Some of these approaches have their problems. Jack Grieve, at the University of Birmingham, cautions that scholars don't always control for genre – so the differences in language seen above might just come down to the difference between a more formal news article, and a more casual Facebook post.

40 To get around this problem, Grieve's team has compared 40 retracted and 41 non-retracted articles by Jayson Blair, who resigned from the New York Times in disgrace in 2003. These were produced in a single genre – national newspaper writing – but they still displayed subtle, probably unconscious differences in register, related, according to Grieve, to the different communicative purposes they served (on the one hand to inform, on the other to deceive). Even though he was trying to pass his work off as factual, there were subtle tells that only become evident when the data is crunched. For example, there were more emphatics like "really" and "most" in Blair's retracted articles. He used

45 shorter words and his language was less “informationally dense”. The present tense cropped up more often and he relied on the third person pronouns “he” and “she” rather than full names – something that’s typical of fiction.

50 So what does all this tell us? Clearly, we don’t have a foolproof means of telling fact from fake yet, but there are certain features that should put us on our guard. Is the writing more informal than you’d expect? Does it contain lots of superlatives and emphatic language? Does it make subjective judgments or read more like narrative than reportage? Ultimately, we may have to rely on artificial intelligence to do the heavy lifting for us – and it should be able to tell us whether those telltale linguistic patterns seen in large datasets of fake news, invisible to the “naked eye”, are present.

55 For me there’s an interesting correspondence with certain kinds of political rhetoric here. The language of fakery, with its powerful subjective statements and focus on anxiety, has something in common with that used by populist leaders. Their style, which often involves “adversarial, emotional, patriotic and abrasive speech” should put us on our guard too. Cooler heads make for a more boring read, but they might get you a little closer to the truth.

• *David Shariatmadari is a Guardian editor and writer, and author of Don’t Believe A Word: The Surprising Truth About Language*

60