# arXiv vs. snarXiv

**Anonymous Author(s)**
Affiliation
Address
email

## Abstract

We attempt to construct classifiers that can accurately distinguish between real
arXiv `hep-th` abstracts and fake abstracts generated from a context-free grammar
by the program snarXiv[1]. We study the bag-of-words and $n$-grams language models
as well as several different classifiers: naive Bayes, likelihood-ratio, and tf-idf. We
find the likelihood-ratio and tf-idf methods to perform extremely well, classifying
abstracts correctly with nearly 100% accuracy.

## 1   Introduction

The arXiv is a free online repository maintained by Cornell for preprints of scientific papers (many of
which go on to be published in journals) from fields including mathematics, physics, and computer
science. Submissions are moderated via an endorsement system to ensure that only legitimate
papers are posted, making it a valuable resource for scientists in many disciplines. In many areas of
theoretical physics, almost all new papers are posted to the arXiv before publication.

The snarXiv is a computer program created by physicist David Simmons-Duffin that generates
titles and abstracts in the style of arXiv's high-energy theory section (`hep-th`) using a context-free
grammar. Simmons-Duffin created a browser game, "arXiv vs. snarXiv,"[2] in which readers are
presented with two title/abstract pairs—one from `hep-th`, the other generated with snarXiv — and
are asked to choose which one is real. After 750,000 guesses, the success rate was only 59%; i.e.,
people mistakenly identified a randomly-generated snarXiv title/abstract as real 41% of the time.
This turns out to be a task that humans are remarkably (and amusingly) bad at, so we here investigate
whether machines are any better.

In the first part of this report, we discuss the principles of our text representation (§2) and classification
(§3) schemes. In the second, we present the results of parameter tuning and final classification
accuracies. Details of implementation are collected in Appendix A.

## 2   Language models

### 2.1   $n$-grams

Assume that each abstract $X$ has a probability $\mathbb{P}(X|Y)$ of being from source $Y \in \{-1, 1\}$, where
$Y = 1$ represents the label "snarXiv" and $Y = -1$ represents the label "arXiv." Suppose $X =
(w_1, w_2, \ldots, w_N)$ is an abstract containing $N$ words $\{w_i\}_{i=1}^{N}$, and let $w_i^{i+k}$ denote the word sequence

---

[1] `http://davidsd.org/2010/03/the-snarxiv/`

[2] Now sadly defunct, it appears, but some results are summarized at `http://davidsd.org/2010/09/the-arxiv-according-to-arxiv-vs-snarxiv/`

29   $(w_i, w_{i+1}, \ldots, w_{i+k})$. The $n$-grams language model approximates $\mathbb{P}(X|Y)$ as

$$\mathbb{P}(X|Y) = \mathbb{P}(w_1^N|Y) = \prod_{i=1}^{N} \mathbb{P}(w_i|w_1^{i-1}, Y) \tag{1}$$

$$\approx \prod_{i=1}^{N} \mathbb{P}(w_i|w_{i-n+1}^{i-1}, Y) \quad (n\text{-grams approx.}) \tag{2}$$

30 In other words, in the $n$-grams model we assume that the probability of a given word $w$ occurring in
31 some abstract only depends on what the previous $n-1$ words in the abstract are, where $n$ is some
32 (typically small) positive integer; the $n-1$ words on which $w$ depends are called the "context" of $w$.
33 In the simplest case $n = 1$, each word occurrence is assumed to be completely independent of any
34 context, which is why the 1-gram language model is often referred to as the "bag-of-words" (BoW)
35 model.

36 We define our $n$-grams estimate $\widehat{\mathbb{P}}(X|Y)_{n\text{-gram}}$ for a given training set from source $Y$ to be

$$\widehat{\mathbb{P}}(X|Y)_{n\text{-gram}} = \widehat{\mathbb{P}}(w_1^N|Y)_{n\text{-gram}} \equiv \prod_{i=1}^{N} \widehat{\mathbb{P}}(w_i|w_{i-n+1}^{i-1}, Y), \tag{3}$$

$$\text{where} \quad \widehat{\mathbb{P}}(w_i|w_{i-n+1}^{i-1}, Y) \equiv \frac{C_Y(w_{i-n+1}^i) + 1}{C_Y(w_{i-n+1}^{i-1}) + V}, \tag{4}$$

37 $C_Y(w_{i-n+1}^i)$ is the number of times that the $n$-gram $w_{i-n+1}^i$ occurs, $\mathcal{V}$ is a vocabulary of size
38 $V = |\mathcal{V}|$ containing all of the words that we wish to learn about, and $\sum_{w \in \mathcal{V}} C_Y(w)$ is the total
39 number of words in the training set, and $C_Y(w_{i-n+1}^{i-1})$ is the number of times that the $(n-1)$-gram
40 $w_{i-n+1}^{i-1}$ appears. The extra $+1$ on each word count $C_Y(w)$ comes from Laplace smoothing, which is
41 a precautionary measure taken to prevent zero probabilities from occurring in cases where we need to
42 classify an abstract that contains a word that did not appear in the training set. Other technical details
43 regarding our $n$-grams implementation (including parsing) can be found in Appendix A.

44 Note that since a typical arXiv/snarXiv abstract contains roughly 100 words, so the abstract probability
45 $\mathbb{P}(X|Y) = \prod_{i=1}^{N} \mathbb{P}(w_i|w_1^{i-1}, Y)$ is a product of $N \approx 100$ different word probabilities, each of which
46 is typically much less than one. As a result, the size of $\mathbb{P}(X|Y)$ is heavily dependent on the abstract
47 length and is usually extremely small, causing our estimates to hit numerical underflow in some cases.
48 To combat the issue of length dependence, we prefer to work with the perplexity

$$\text{PP}(X|Y) \equiv \sqrt[N]{\frac{1}{\mathbb{P}(X|Y)}} = \left[\prod_{i=1}^{N} \mathbb{P}(w_i|w_1^{i-1}, Y)\right]^{-1/N}, \tag{5}$$

49 which normalizes the abstract probability by taking the geometric mean of the individual word
50 probabilities (and inverting). To prevent numerical underflow from rendering our algorithms useless,
51 we always work with logs of probabilities in our code, which turns products into sums and keeps the
52 numbers reasonable. Note that the log of the perplexity is a simple arithmetic mean:

$$\log \text{PP}(X|Y) = -\frac{1}{N} \log \mathbb{P}(X|Y) = -\frac{1}{N} \sum_{i=1}^{N} \log \mathbb{P}(w_i|w_1^{i-1}, Y). \tag{6}$$

53 In practice, we use the logs of the classification conditions given below in Section 3 to classify our
54 test abstracts as arXiv or snarXiv; we only exponentiate back to perplexities for the purposes of
55 figures.

## 2.2   tf–idf

57 The family of $n$-gram models can be used to produce term probabilities for naive Bayes or likelihood-
58 ratio testing, but they also induce perfectly good $\frac{V!}{(V-n)!}$-dimensional vector representations of
59 documents, where each entry records the occurrence of a given $n$-gram in the document. These
60 vectors are large, but fairly sparse, and make rudimentary features for regression or clustering

algorithms. We enhance this model using the term-frequency–inverse-document-frequency (tf–idf) scheme as described in Jurafsky and Martin [2014, ch. 6].

Each entry $d_i$ in a document vector $d$ corresponds to a term $w_i$ in the vocabulary and receives a weight

$$d_i = \text{tf}_{i,d} \cdot \text{idf}_i. \tag{7}$$

We define

$$\text{tf}_{i,d} \equiv \begin{cases} 1 + \log_{10} C_{i,d} & : & C_{i,d} > 0 \\ 0 & : & \text{else} \end{cases}, \tag{8}$$

where $C_{i,d}$ is the number of occurrences of $w_i$ in the document $d$, and

$$\text{idf}_i \equiv \log_{10} \frac{N+1}{\text{df}_i + 1}, \tag{9}$$

where $N$ is the number of documents in the training corpus, and $\text{df}_i$ is the number of documents containing $w_i$. [3]

The principle of tf–idf is to weight vector components according to their discriminating power: our representation should emphasize the most unique parts of documents. Accordingly, the inverse document frequency metric completely ignores a word that occurs in every document, but gives great weight to a word occurring only once in the corpus. Further, our term-frequency metric grows sublinearly, since a word occurring ten times is not necessarily ten times as important as one occurring singly.

# 3 Classifiers

## 3.1 Naive Bayes classifier

Our first attempt at creating a program that distinguishes between arXiv and snarXiv abstracts utilized a naive Bayes (NB) classifier. The theory behind this classifier is simple: if $\mathbb{P}(Y = 1|X) > \mathbb{P}(Y = -1|X)$, then we classify $X$ as snarXiv; otherwise, we classify it as arXiv. We can use Bayes' theorem $\mathbb{P}(Y|X) = \frac{\mathbb{P}(Y)\mathbb{P}(X|Y)}{\mathbb{P}(X)}$ to rewrite this classification condition as

$$\widehat{Y}_{NB} \equiv \arg \max_{Y \in \{-1,1\}} \mathbb{P}(Y)\mathbb{P}(X|Y). \tag{10}$$

We estimate each probability in (10) by training on a large number of arXiv and snarXiv abstracts. We define our estimate[4] for $\mathbb{P}(Y)$ as the number of abstracts in training set $Y$ over the total number of abstracts in both the arXiv and snarXiv training sets. The conditional probabilities $\mathbb{P}(X|Y)$ are more complicated, and are approximated using a language model from Section 2.

## 3.2 Likelihood-ratio test

A slightly more sophisticated classification rule is the likelihood-ratio (LR) test. For a given abstract $X$, define the likelihood ratio $\Lambda(X) \equiv \frac{\mathbb{P}(X|Y=1)}{\mathbb{P}(X|Y=-1)}$. While we certainly want our classifier to correctly label all snarXiv abstracts as fakes, we also want to minimize the probability that an arXiv abstract gets incorrectly classified as fake. The Neyman-Pearson lemma states that the optimal classifier

$$\delta_{LR}^* \equiv \arg \max_\delta \mathbb{P}(\delta(X) = 1|Y = 1), \text{ subject to } \mathbb{P}(\delta(X) = 1|Y = -1) \leq \alpha \tag{11}$$

---

[3]Our definition of $\text{idf}_i$ varies from that given in the reference by smoothing terms which deal with the possibility of developing features from a corpus smaller than the one that produced the vocabulary and prevent infinite logs.

[4]The probabilities $\mathbb{P}(Y = \pm 1)$ are somewhat trivial since we generate the snarXiv abstracts ourselves and can control over how many arXiv vs. snarXiv abstracts will be in the train and test sets. Since the snarXiv was originally created to make abstracts that compete one-on-one with arXiv abstracts, we chose to set $\mathbb{P}(Y = \pm 1) = \frac{1}{2}$ for all tests reported in this paper.

for any fixed false-positive tolerance $\alpha$ takes the form

$$\mathbb{P}\left(\delta_{LR}^*(X) = 1\right) = \begin{cases} 1, & \Lambda(X) > \eta \\ \gamma, & \Lambda(X) = \eta \\ 0, & \Lambda(X) < \eta, \end{cases} \tag{12}$$

where $\eta$ and $\gamma$ can be treated as hyperparameters to be tuned in cross-validation. In practice, $\gamma$ is almost completely irrelevant since for any $\eta$, the odds that a language model like $n$-grams would yield an estimate for $\Lambda(X)$ that is *exactly* equal to $\eta$ are essentially zero.[5] We arbitrarily use $\gamma = 0.5$, meaning the only hyperparameter we tune is $\eta$. Because of the issues of small $\mathbb{P}(X|Y)$ probabilities discussed in Section 2.1, it is more practical to instead compare the perplexities $\text{PP}(X|Y) = \mathbb{P}(X|Y)^{-1/|X|}$. The relevant parameter in the LR-test is thus $\eta^{1/|X|}$. We would like our hyperparameter to not depend on the specifics of the input though, so we note that the average abstract length is about 120 words and define $\eta_{\text{PP}} \equiv \eta^{1/120}$ to be our new LR-hyperparameter.

## 3.3 Logistic regression

For weighted vector representations as produced by the tf–idf scheme, it is not obvious how to produce term or document likelihoods for use in a Bayes or LR classifier, so we must turn to other classification methods. Fortunately, there are a host of methods for binary classification of vector data, and it is convenient and familiar to choose regularized logistic regression. We adopt a model

$$\mathbb{P}(Y = y \mid x) = \frac{1}{1 + e^{-yw^T x}} \tag{13}$$

where $y \in \{-1, 1\}$ represents the class of a document (negative for arXiv, positive for snarXiv), $x$ is a tf–idf document vector, and $w$ and $b$ are parameters to be learned. We accomplish this by minimizing the convex regularized objective

$$J(w) = \sum_{i=1}^N \log\left[1 + \exp\left(-y_i(x_i^T w + b)\right)\right] + \lambda \|w\|_2^2, \tag{14}$$

where $\lambda$ is a hyperparameter. The optimum lacks a closed form, but it can be found easily in CVXPY even for the $V \sim 15000$-dimensional space created by our vocabulary. After training $w$ and $b$, we classify a document according to the sign of the exponent $(x_i^T w + b)$, since it can be related to the log-ratio of probabilities to be in each class.

# 4 Results

We report results for three different algorithms: bag-of-words model with naive Bayes classifer (BoW-NB), bag-of-words model with likelihood-ratio test (BoW-LR), and bigrams=(2-grams) model with likelihood-ratio test (bi-LR). In each case, we trained on an equal number $N_{\text{train}}^Y = 1000$ of arXiv and snarXiv abstracts, and we also tested on 1000 arXiv abstracts and 1000 snarXiv abstracts.

For BoW-NB, we found the classification accuracy to be 77%, with the only misclassifications being false positives (arXiv abstracts that the classifier mistakes as snarXiv); when a snarXiv abstract is tested, it is correctly identified as snarXiv 100% of the time. This is somewhat expected for a classifier which does nothing to constrain the false-positive rate.

For BoW-LR and bi-LR, we tuned the LR-hyperparameter $\eta_{\text{PP}}$ through cross-validation to maximize the classification accuracy on abstracts from both arXiv and snarXiv. We found $\eta_{\text{PP}} = 0.7$ to perform best in both models; Figure 1 shows histograms of our estimates for $\frac{\text{PP}(X|\text{snarXiv})}{\text{PP}(X|\text{arXiv})}$ for both the arXiv and snarXiv test corpora. The LR-classifier is clearly distinguishing between the two sets, giving a classification accuracy of 100%.

Logistic regression trained on 2000 tf–idf document vectors and tested on 8000 yielded error rates ranging from 0.05% to 0.375% depending on the regularization parameter (lower $\lambda$ tended to produce lower error, even down to $\lambda \sim 10^{-8}$), and exhibited the same pattern as the Bayes and LR classifiers of 100% accuracy at identifying snarXiv papers. The only errors at any regularization scale were misidentification of arXiv papers.

---

[5]The borderline case of the naive Bayes classifier where $\mathbb{P}(Y = \pm 1|X) = \frac{1}{2}$ is irrelevant for the same reason.
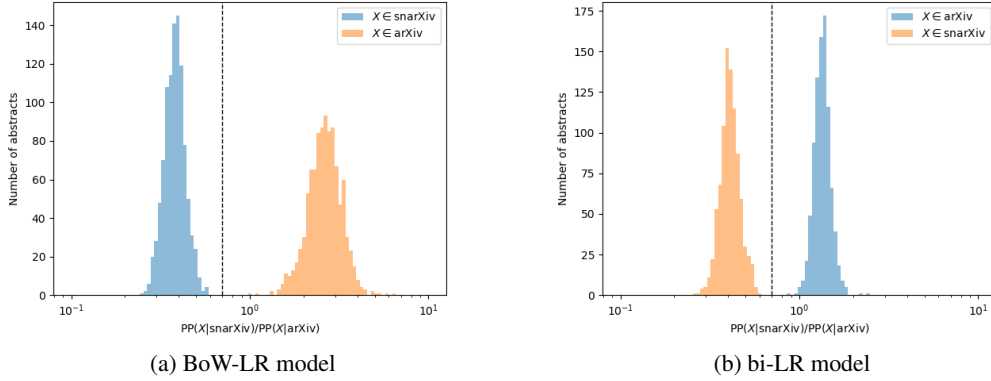
| (a) BoW-LR model | (b) bi-LR model |

Figure 1: Perplexity ratio histograms for both LR-models. The dotted black line corresponds to $\eta_{PP} = 0.7$.

## 5 Conclusion

We have demonstrated that the arXiv–snarXiv discrimination problem, while unusually difficult for humans, yields to even very simple analytical frameworks. This is most likely due to the profound different statistical structures of the corpora illustrated in our histograms, which are not evident to human inspection (which focuses on attempting to parse or impart semantic structure) but stand out readily to probabilistic classification methods.

## A Technical details

### A.1 Parsing abstract text

Our parser attempts to translate the raw abstract text (a string) to a list of formatted words that are easily identifiable between different abstracts. We first split the text on all spaces, newlines, and hyphens, and for uniformity we make all words lowercase. Next, we look for words that end in a period or question mark and insert the special word tokens `<e>` and `<s>` afterwards to mark the end of one sentence and the start of another. Lastly, we strip out all punctuation (except for the special tokens), prepend `<s>` to the list (marking the start of the first sentence), and delete the last word if the list ends `<e> <s>` (otherwise, we append `<e>`).

An example of parser input/output is

Input:  ``I'm taking a CSE class on machine-learning. It's a lot of work, but
pretty interesting.''

Output:   [<s>, im, taking, a, cse, class, on, machine, learning, <e>, <s>,
its, a, lot, of, work, but, pretty, interesting, <e>],

where each "word" in the output list is a string.

Our parsing scheme works well, but it leaves traces of some TeX commands from the raw abstract text; for example, it sends \mathbb{Z} → mathbbz. This would be fine if all arXiv and snarXiv abstracts have TeX commands written the same way, though it is possible in some cases for abstracts to write the same command in different ways. In practice, we find that these instances are rare enough for the concern to be negligible.

### A.2 Vocabulary

We defined our vocabulary $\mathcal{V}$ as the collection of $V = 15,315$ unique words that appeared at least twice in a training corpus of 12,000 randomly chosen parsed arXiv abstracts and 12,000 (randomly generated) parsed snarXiv abstracts. We decided not to include the 13,960 words that only appeared once since they seem to be relatively uncommon (often they are the residuals of TeX commands),

and including them would nearly double our vocabulary size. This pruning can be thought of as a feature-processing step.

## A.3 $n$-grams

It is fairly common for a word to appear in the test set that is not in the vocabulary. When this happens, we change the word to a special token <UNK>. We then treat <UNK> as we would any other word when counting $n$-gram occurrences $C_Y(\cdot)$.

We prepend copies of the start sentence token <s> to each parsed abstract to ensure that the first $n$-gram consists of $n-1$ copies of <s> and one "real" word. Similarly, we append copies of <e> to the end of each parsed abstract.

## References

Charu C. Aggarwal and ChengXiang Zhai. A survey of text classification algorithms. In *Mining Text Data*. Springer, Boston, Massachusetts, 2012. URL `https://doi.org/10.1007/978-1-4614-3223-4_6`.

Akshay Agrawal, Robin Verschueren, Steven Diamond, and Stephen Boyd. A rewriting system for convex optimization problems. *Journal of Control and Decision*, 5(1):42–60, 2018.

William Cavnar and John Trenkle. N-gram-based text categorization. 05 2001.

Steven Diamond and Stephen Boyd. CVXPY: A Python-embedded modeling language for convex optimization. *Journal of Machine Learning Research*, 17(83):1–5, 2016.

Dan Jurafsky and James H Martin. *Speech and language processing*, volume 3. Pearson London, 2014.

Mehran Sahami, Susan Dumais, David Heckerman, and Eric Horvitz. A bayesian approach to filtering junk e-mail. In *Learning for Text Categorization: Papers from the 1998 workshop*, volume 62, pages 98–105. Madison, Wisconsin, 1998.

Shrawan Kumar Trivedi. A study of machine learning classifiers for spam detection. In *Computational and Business Intelligence (ISCBI), 2016 4th International Symposium on*, pages 176–180. IEEE, 2016.