
Text Ranking and Classification using Data Compression

Nitya Kasturi
Meta
nityakasturi@fb.com

Igor L. Markov
Meta
imarkov@fb.com

Abstract

A well-known but rarely used approach to text categorization uses conditional entropy estimates computed using data compression tools. Text affinity scores derived from compressed sizes can be used for classification and ranking tasks, but their success depends on the compression tools used. We use the `Zstandard` compressor and strengthen these ideas in several ways, calling the resulting language-agnostic technique `Zest`. In applications, this approach simplifies configuration, avoiding careful feature extraction and large ML models. Our ablation studies confirm the value of individual enhancements we introduce. We show that `Zest` complements and can compete with language-specific multidimensional content embeddings in production, but cannot outperform other counting methods on public datasets.

1 Motivation

The idea of comparing texts using off-the-shelf lossless data compression tools goes back to [1], which in 2002 linked *entropy estimation* and using `gzip` on text with text similarity metrics. Given two strings X and Y , one compresses each of them individually and also the string $X + Y$. Similar strings compress better after being concatenated. An affinity score for X and Y is computed from the three resulting bytesizes. This computation is simple, requires little infrastructure, works for any language, and naturally handles similar words, word forms, typos, etc. It can be easily applied to multi-class classification (e.g., binning news articles by category) and ranking relative to known examples.

Following up on [1], the 2004 work in [2] showcased the versatility of compression-based estimators by demonstrating their practical advantages for (i) clustering time-series and text, (ii) anomaly detection, (iii) classification for various non-text datasets. These techniques are attractive in practice because they can handle many data formats without the need to understand their structure, are easy to configure, and are reasonably efficient in terms of computation. The authors of [2] also emphasize the “parameter-free” nature of compression-based estimators, in comparison with various ML techniques that may produce different results depending on parameter settings. In practice, some compression tools are better than others on certain kinds of inputs (making a choice of a compression tool is a type of fitting to data), and hyperparameter settings affect compression performance: compression block size and word size (8/12/16 bits), level of effort, parallel execution, efficient use of CPU instruction sets, etc. Moreover, compression headers and dictionaries embedded in each compressed file spoil entropy estimates, especially for small files. In fact, a controversy ensued after the publication of [1] because zipping could not outperform Naive Bayes [3].

Despite the past controversies and apparent lack of broad adoption, the idea to use data compression for text classification has become so mainstream today that it appears in a key AI textbook [4, Chapter 23]. There, Russel and Norvig point out that “In effect, compression algorithms are creating a language model. The LZW algorithm in particular directly models a maximum-entropy probability distribution.” Presumably, compression ratio then estimates the probability that given input was

generated by the language model. Motivated by this discussion, Halford offers [5] a recent (2021) re-evaluation of the idea from [1], using a recent version of `gzip` compressor on the public News data set available with `sklearn` [12]. Indeed, Multinomial Naive Bayes from `sklearn` clearly outperforms `gzip`-based classification by F1 score and "It's also dramatically faster". Halford evaluates several other modern compressions — `zlib`, `bz2`, and `lzma`, which show wildly different results in terms of F1 and runtime. All of these tools are much slower than Multinomial Naive Bayes, and only `lzma` shows performance on par with it but at the cost of being several times slower than other data compressors. In contrast to these recent pessimistic results, we show how to significantly improve the ML performance and runtime of compression-based text classification.

2 Insights and Overview

We revisit compression-based text affinity scores [1] because modern compressors are faster than tools from 2002 and produce better entropy estimates. The lossless open-source `zstandard` compressor [6, 7, 8], developed by Yann Collet at Facebook around 2016, produces results not far from those of (slow) arithmetic coding that are considered close to the Shannon bound. The `zstandard` compressor offers a *dictionary interface*, which allows us to improve the approach of [1] and make it more practical, especially for small and medium-sized strings. `Zstandard` can build a compression dictionary and use it to compress many small files, to avoid the overhead of separate dictionaries. We use the dictionary interface to compress texts to be classified or ranked. This sharpens entropy estimates and improves speed versus compressing concatenated files (where the same file would be concatenated with multiple other files). The resulting text affinity scores can be used as inputs to multimodal classifiers and rankers. The simplicity of this language-agnostic method is attractive when building ML platforms, especially for product engineers without ML background.

We leverage text affinity estimation in text ranking and classification. For example, given positive and negative examples for 2-class classification, we first build compression dictionaries for each class. As an option, the texts can be normalized by removing punctuation within sentences (but not spaces) and lowercasing the remaining letters (for languages without upper/lower cases, this is a no-op). With `Zstandard`, there is no need to concatenate files as in [1], thanks to the dictionary interface.

The original approach has several major weaknesses that we are able to address. To illustrate them, consider classifying news articles in topics — Politics, Celebrities, and Sports. Some important words and phrases appear in multiple topics, but with different frequencies, for example, "Arnold Schwarzenegger". For a sufficiently large set of examples, such words compress equally well for each class, and do not contribute useful information. This is particularly detrimental when classifying or ranking short texts. Downsampling the examples would help with common words and phrases, but undermine the handling of rare words and phrases (which would not be compressed for any class). To address this challenge, we use *a set of dictionaries of telescoping sizes* — this way, common words are differentiated by smaller dictionaries and rare words are differentiated by larger dictionaries. Another challenge is the heavier impact of longer words on compression ratios. We address it by *word padding* to fixed length, e.g., "hello" padded to 10 characters becomes "hellohello". We configure `zstandard` to minimize headers in compressed files and, furthermore, subtract the compressed size of an empty string from compressed sizes of evaluated pieces of text. For each classification class, we average the byte compression ratio over multiple dictionaries. Subtracting this number from 1.0 produces an affinity score, for which "greater is better". In particular, a sentence that was seen in some class examples may return a value close to 1.0, whereas a sentence in a different script (e.g., Greek vs Cyrillic) would not compress well, resulting in a value close to 0.0. For multiclass classification, we subtract the min class score from all scores. This handles words present in many classes. For ranking applications, affinity scores can be sorted to produce an ordering.

3 Implementation and Empirical Evaluation

Our PyTorch implementation is based on an untrained `Torchscript` module for `Zest`. The Python `zstandard` library was not supported with `TorchScript`, so we implemented the `Zstandard` interface in C++ as a `TorchScript` module using the original C library. For comparison with other linear models and public text datasets, we used the Python `zstd` library directly. The `Zest` transformer takes in lists of text containing the examples per class to train separate dictionaries, as

well as a list of text features for evaluation. It produces affinity scores per class, which can be used separately or combined (e.g., by subtraction) into one score.

3.1 Deployment in a production ML Platform

We onboarded Zest to a company-internal ML platform that hosts hundreds of ML models for prediction, classification, etc. Using production data, we evaluate the use of Zest text affinity scores as additional features to these models, where we can judge features by their importance values in the context of various other features. The specific application discussed in this paper uses text features to rank search results in an internal tool (posts) produced for a given search query. Some model features provide context: `user_id`, `group_id`, and `post owner_id`. The most useful features in the deployed ML model include the number of characters, query proximity to the post text, and components of a multi-dimensional content embedding. In this study, we check if compression-based affinity scores are comparable in their utility to language-specific content embeddings, which require additional training and cannot handle text in many languages, or mixed-language text. In contrast, Zest can be useful where content embeddings are not easily available.

To prepare input for text affinity computation, we feature-engineer the positive and negative post examples to evaluate the post text. For each user on the internal tool, we fetch the 12 most recent posts from the 5 groups with which the user interacted most recently. We split these examples between positive and negative by whether the user viewed them more than once. As most posts end up being negative examples for a user, we added posts that are trending on the internal tool as positive examples. On average, a user has 39 negative examples with 851 characters each and 20 positive examples with 1659 characters each. This way, the number of characters passed into the positive and negative compression dictionaries balances out.

3.2 Empirical evaluation within a larger model

We ran a baseline gradient-boosted decision tree (GBDT) workflow with no additional features, a GBDT workflow with Zest scores as features, and a GBDT workflow with 3-gram affinity scores (percentage of matched 3-grams) as features. The Zest transformer took roughly 50% longer to compute affinity scores in our workflows compared to the n-gram transformer (both implemented by us). However, the Zest classifier code ran much faster with the Python implementation, taking 200-300 seconds to build dictionaries on 25k examples, whereas the transformer would take 40-50 minutes for a few thousands of examples. This could be due to the memory limits placed on TorchScript modules when running in production. Feature importance of the top Zest feature was #16 (compared to #201 for the n-gram models), ahead of hundreds of content-embedding dimensions and behind of only 7 of them.

We then tested removing embedding features to see how well Zest can fill in the missing information. Based on random removal of embedding dimensions, the new features allow the model to drop 150 dimensions of both the post and query embeddings, while improving normalized entropy (NE) by 7% compared to the baseline model. With half of the embedding features removed, the top Zest feature was at #8 with only 4 embedding features ahead. An important distinction between Zest and the embedding features is that the Texas SIF embeddings [9, 10] used in this application can only support English and Spanish, whereas the Zest transformer is language-agnostic and can handle mixed-language text. Incidentally, the Texas SIF embeddings are shown to be a "Tough-to-Beat Baseline for Sentence Embeddings" [9]. We also checked feature quality of the Zest transformer by removing features representing an ID, given that the ID features had high feature importance and could be affecting model quality. After removing the ID features, feature importance jumped from #16 to #9, again being outperformed by only 4 embedding dimensions.

Table 1 compares different workflow runs with different configurations along with the offline results for comparison. There is a distinction in the score distribution between the True and False examples, with most True examples around 0.03 - 0.11, and most False examples around 0-0.07, indicating that the feature values can be strong enough to help with classification on its own.

3.3 Standalone comparisons versus linear models

As demonstrated on production data, Zest features can be useful when passed into a model along with other features. To compare Zest to other text transformers while avoiding non-text features,

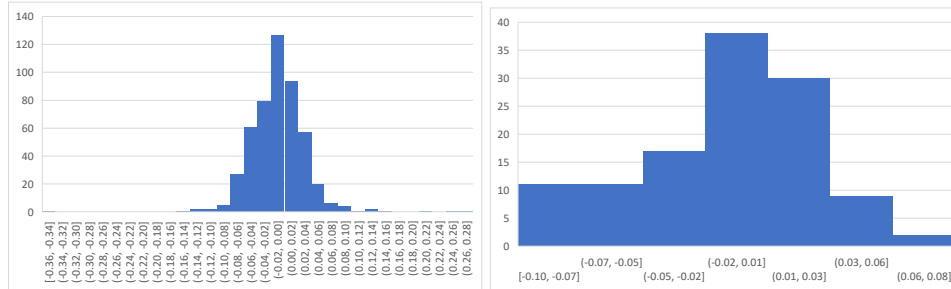


Figure 1: Zest score distribution for True and False examples respectively.

to perform ablation studies, and experiment with ensembling, we worked with various public text classification and sentiment analysis datasets. We ran Logistic Regression (LR) trained on Bag of Words (BoW) as a feature, LR trained on Facebook’s InferSent sentence embeddings [11], and the multi-class version of Zest on a dataset with various categories of news headlines [12] and various sentiment datasets (Stanford Movie sentiment [13] and IMDB movie review sentiment [14]). We ran Zest with 1, 2, and 4 telescoping dictionaries each to check performance on the news headlines dataset. The Zest model with 4 telescoping dictionaries performed significantly better than the rest. Word padding generally improved accuracy by 0.5-1% based on the dataset. Compared to sophisticated language models — BERT and such, — Zest has a much smaller resource footprint and is easier to work with, yet customizable.

The ensemble (averaged) prediction of Zest and LogisticRegression with a BoW performs the best on the news headline categories dataset. However, BoW as a feature worked well with Logistic Regression — it ran faster and outperformed standalone Zest by 1-5% in accuracy, depending on the dataset. The sentiment analysis datasets both had an accuracy of $\approx 87\%$ with LR and BoW versus $\approx 78\text{-}80\%$ for Zest, no matter how it was ensembled.

Table 1: Comparisons in a production setting

model	post emb.	query emb.	feat. imp.	NE
Baseline	All	All	—	0.196
Zest	All	All	0.29%	0.212
3-gram	All	All	0.06%	0.216
Zest	None	None	3.57%	0.448
Zest	150 dims	150 dims	1.54%	0.189
Baseline	All	All	—	0.173
Baseline	None	All	—	0.161
Zest	None	All	1.14%	0.151
3-gram	None	All	0.17%	0.151

Table 2: Comparisons for news headlines [12]

model	size (MB)	train sec.	acc.
NB/BoW	3.34	0.170	0.926
LR/BoW	1.67	26.53	0.947
LR/InferSent	777	822.3	0.874
Zest 4D	3.04	74.43	0.924
Zest 4D/LR-BoW	4.71	100.9	0.951
Zest 2D	1.02	60.49	0.886
Zest 2D/LR-BoW	2.69	87.02	0.946
Zest 1D	0.02	49.83	0.741
Zest 1D/LR-BoW	1.69	76.36	0.938

The baseline model is a GBDT. Alternatively, we pass either Zest or the 3-gram counter in as additional features to the baseline model. BoWs are also compared to Zest and also passed in as features, to a Logistic Regression (LR) or Naive Bayes (NB). Zest is ablated with 1, 2 and 4 telescoping dictionaries. Additional sentiment analysis datasets (not shown) exhibit similar trends.

3.4 Comparisons with other compression algorithms

As discussed in Section 1, Halford [5] compares various compression algorithms by speed and quality of results on a categorical newsgroups dataset. While `zstd` is not included, the LZMA algorithm performs the best but takes over 30 minutes for 4,000 documents. We extend Halford’s work by including `zstd` in two ways — as a regular compressor (without using its dictionary mode) and via Zest. Such comparisons can clarify if (1) other compression algorithms are competitive and (2) if our improvements help with performance and speed.

Halford uses the Newsgroups dataset [16] that contains 20,000 documents from 20 newsgroups. GZip, Zlib, and BZ2 all underperform, with an average of 0.65-0.75 for both precision and recall. The best-performing algorithm, LZMA, attains both a precision and recall of 0.897 and takes 30 minutes to complete. zstd used as a regular compressor achieves a precision 0.857 and completes in under a minute. Zest takes only 40 seconds to train and test the documents, with FOM 0.899. Zest is a significant improvement over using compression algorithms on its own, likely due to the telescoping dictionaries.

4 Conclusions and Future Work

We have demonstrated model-free language-agnostic text features based on data compression that can be useful to text rankers and classifiers. In addition to using a modern data compression tool, our implementation goes beyond the ideas in [1] by leveraging the *dictionary mode* in zstandard, using *telescoping dictionaries* and performing *word padding*. Empirical performance of TorchScript module on a production ML platform is competitive with that of content-embedding features. However, for some simpler datasets with clear distinction between text classes, BoW shows better ML performance, while being simple and fast.

Evidence from Section 3 allows us to conclude the following:

1. Compression-based methods can be significantly improved by telescoping dictionaries and word padding.
2. Counting methods can achieve strong performance in some cases.
3. Ability to handle synonyms via embeddings offers no advantage on some practical datasets.
4. Different datasets and classification types favor different models.
5. Zest is competitive with trained sentence embeddings in production settings.
6. Zest outperforms other counting methods such as n-grams in production datasets.

Although Zest performance is strong and can be an adequate replacement for content embeddings in a production environment, BoWs as features consistently performed the best on all public datasets that we used to evaluate Zest. cursory analysis suggests that these datasets allow identifying each class by a small set of words, making explicit counting more accurate than compression-based methods and trained word/content embeddings. However, in a production environment as the one described in Section 3, frequent class-specific words are less common, allowing methods like Zest to be on par with competition but with no training. We also evaluated (Markov-chain) language models that estimate the probability that a given text was generated by a given language [4, 15]. In our experiments (not shown), they produce slightly more accurate classifiers than Zest, but tend to be more complex and consume greater resources. Compression-based methods using compression algorithms out of the box do not perform well either, with the exception of lzma. zstd slightly underperforms in comparison with lzma, but can cut the compression time significantly. Zest with its telescoping dictionaries is efficient and performs well in our compression-based experiments.

Our findings are useful when designing ML platforms that need to deal with text features without asking product engineers to write ML code. They suggest maintaining several lightweight, language-agnostic text features including compression-based ones, and letting the model-training process choose which features are helpful. In many applications, such low-hanging features provide performance that is close to or better than more sophisticated word/content embeddings, while using a much smaller resource and latency footprint. Unlike more sophisticated methods, lightweight methods tend to be language agnostic and can be implemented without language detection.

It is straightforward to generalize our methods to the applications explored in [2], and we believe that many of our conclusions are going to generalize as well.

Our code and experiments are available at <https://github.com/facebookresearch/zest>.

References

- [1] D. Benedetto, E. Caglioti, V. Loreto, "Language Trees and Zipping," *Phys. Rev. Lett.*, 2002.
- [2] E. J. Keogh, S. Lonardi, Ch. (Ann) Ratanamahatana, "Towards parameter-free data mining," *KDD 2004*: 206-215.
- [3] D. Benedetto, E. Caglioti, V. Loreto "On J. Goodman's comment to "Language Trees and Zipping"", <https://arxiv.org/abs/cond-mat/0203275>
- [4] S. Russell and P. Norvig, *Artificial Intelligence: A Modern Approach*, 4th ed., 2020. <http://aima.cs.berkeley.edu/>
- [5] M. Halford, "Text classification by data compression," 2021, <https://maxhalford.github.io/blog/text-classification-by-compression/>
- [6] Zstandard, *Wikipedia*, <https://en.wikipedia.org/wiki/Zstandard>
- [7] Y. Collet and C. Turner, "Smaller and faster data compression with Zstandard," *Facebook Engineering Blog*, 2016, <https://engineering.fb.com/2016/08/31/core-data/smaller-and-faster-data-compression-with-zstandard/>
- [8] Yann Collet, M. Kucherawy, "Zstandard Compression and the application/zstd Media Type," RFC 8478, 2018, DOI: 10.17487/RFC8478.
- [9] S. Arora, Y. Liang, T. Ma, "A Simple but Tough-to-beat Baseline for Sentence Embeddings," *ICLR 2017*.
- [10] S. Arora, M. Khodak, N. Saunshi, K. Vodrahalli, "A Compressed Sensing View of Unsupervised Text Embeddings, Bags-of-n-grams, and LSTMs", *ICLR 2018*.
- [11] A. Conneau, D. Kiela, H. Schwenk, L. Barrault, A. Bordes, "Supervised Learning of Universal Sentence Representations from Natural Language Inference Data," *EMNLP 2017*.
- [12] M. Lichman, "UCI Machine Learning Repository," *Irvine, CA: University of California, School of Information and Computer Science 2013*.
- [13] R. Socher, A. Perelygin, J. Wu, J. Chuang, C. Manning, A. Ng, C. Potts, "Recursive Deep Models for Semantic Compositionality Over a Sentiment Treebank", *EMNLP 2013*.
- [14] A. Maas, R. Daly, P. Pham, D. Huang, A. Ng, C. Potts, "Learning Word Vectors for Sentiment Analysis", *ACL 2011*.
- [15] D. Jurafsky and J. Martin, "Speech and Language Processing: An Introduction to Natural Language Processing, Computational Linguistics, and Speech Recognition," *Prentice Hall 2000*.
- [16] K. Lang, "Newsweeder: Learning to filter netnews," *ICML 1995*.