

On consistency scores in text data with an implementation in R *

Ke-Li Chiu *University of Toronto*
Rohan Alexander *University of Toronto*

In this paper, we introduce a reproducible cleaning process for the text extracted from PDFs using n-gram models. Our approach compares the originally extracted text with the text generated from, or expected by, these models using earlier text as stimulus. To guide this process, we introduce the notion of a consistency score, which refers to the proportion of text that is expected by the model. This is used to monitor changes during the cleaning process, and across different corpuses. We illustrate our process on text from the book *Jane Eyre* and introduce both a Shiny application and an R package to make our process easier for others to adopt.

Keywords: text-as-data; natural language processing; quantitative analysis; optical character recognition

1 Introduction

When we think of quantitative analysis, we may like to think that our job is to ‘let the data speak’. But this is rarely the case in practise. Datasets can have errors, be biased, incomplete, or messy. In any case, it is the underlying statistical process, of which the dataset is an artifact, that is typically of interest. In order to use statistical models to understand that process, we typically need to clean and prepare the dataset in some way. This is especially the case when we work with text data because they are typically not designed for analysis. This cleaning and preparation requires us to make many decisions. To what extent should we correct obvious errors? What about slightly-less-obvious errors? Again, this is especially the case for text datasets because: they encode a large amount of information; an aspect may only be an error in context; and the high-dimensionality of the data at even its most basic level (26 English characters, compared with 10 numbers). Although cleaning and preparation is a necessary step, we may be concerned about the extent to which we have introduced new errors, and the possibility that we have made decisions that have affected, or even driven, our results.

In this paper we introduce the concept of consistency and implement it for a text corpus. Consistency refers to the proportion of data that are able to be forecast by a statistical model, based on preceding and surrounding data. Further, we define internal consistency as when the model is trained on the dataset itself, and external consistency as when the model is trained on a more general dataset. Together, these concepts provide a guide to the cleanliness and reasonableness of a dataset. This can be important when deciding whether a dataset is fit for purpose; when carrying out data cleaning and preparation tasks; and as a comparison between datasets. For instance, if the Canadian population is known to be 38 million from a census, but the sum of age-groups is initially 35 million, rising to 38 million only after an inputation error is corrected, then that post-

*Thank you to Monica Alexander and Alex Luscombe for helpful comments. A Shiny application for interactive and small-scale examples is available: <https://kelichiu.shinyapps.io/Arianna/>. An R package for larger scale applications is available: <https://github.com/RohanAlexander/arianna>. Our code and datasets are available: https://github.com/RohanAlexander/consistency_scores_in_text_datasets. Comments on the 13 January 2021 version of this paper are welcome at: rohan.alexander@utoronto.ca.

cleaning dataset would have higher consistency than the pre-cleaning dataset. When applied to a text dataset, this means that we construct a model that forecasts the next word, given the previous words. If the forecast word is always the same as the actual word, then the consistency is 1. That said, it is not the absolute level of consistency that is our focus, but instead measuring how it changes during the cleaning and preparation process.

To provide an example, consider the sentence, ‘the cat in the...’. A child who has read this book could tell you that the next word should be ‘hat’. Hence if the sentence was actually ‘the cat in the bat’, then that child would know something was likely wrong. The consistency score would likely be lower than if the sentence were ‘the cat in the hat’ because we expected ‘hat’, but actually found ‘bat’. After we correct this error, the consistency score would likely increase. By examining how consistency scores evolve in response to changes made to the text during the data preparation and cleaning stages we can better understand the effect of the changes. Including consistency scores and their evolution when datasets are shared allows researchers to be more transparent about the effect of their decisions in these steps of the analysis workflow. And finally, the use of consistency scores allows for an increased level of automation in the cleaning process.

We employ n-gram models to calculate consistency scores and generate word candidates for text corrections. The n-gram approach involves identifying words that are commonly found together. In our application, we use these sequences to calculate consistency scores and generate word candidates for text corrections. For instance, if we are using a sequence of three words, or tri-grams, then we identify the first two words in the sequence and evaluate if the last word in the sequence is ‘as expected’ by comparing it to the corresponding tri-gram in a comparison dataset. The sequences we employ in this paper are a mixture of three, four, and five, consecutive words also known as tri-grams, 4-grams, and 5-grams.

The comparison dataset can be internal, meaning that the dataset is from the same source as the text being evaluated. On the other hand, a comparison dataset can be external, which means the dataset is from a different source, typically a larger and more universal dataset. The external dataset we use in the application discussed in this paper is constructed using a subset of text data sourced from Open WebText Corpus (Gokaslan and Cohen, 2019). Since the content of Open WebText Corpus is from the internet, the variety of the text data is promising for meeting the generalizability and universality that an external dataset should possess, although it will also have extensive bias.

The resources that we have developed to support this paper include a Shiny app available at: <https://kelichiu.shinyapps.io/aRianna/>. That app computes internal and external consistency scores for corpus excerpts. We have also developed an R Package, aRianna, that allows our approach to be used on larger datasets, which is available at: <https://github.com/RohanAlexander/aRianna>.

In this paper, we introduce consistency scores and an R package that implements them. In Section 2, we provide an overview of language models and n-grams in particular. In Section 3, we deconstruct the package functions and provide demonstration example of the package usage. Finally, in Section 4, we discuss our findings, limitations, and future directions.

2 Background

2.1 Language model

Given the nature of human languages, some combinations of words tend to occur more frequently than others. Think of ‘good’, which is more often followed by ‘morning’, than ‘duck’. As such,

we could consider English text production as a conditional probability, $\Pr(w_k | w_1^{k-1})$, where k is the number of words in a sentence, w_k is the predicted word, and w_1^{k-1} is the history of the word occurring in a sequence (Brown et al., 1992). In this way, the generation of some prediction, w_k , is based on the history, w_1^{k-1} . This is the underlying principle of all language models. Essentially, the goal of statistical language modeling is to estimate probability distributions over different linguistic units — terms, words, sentences, and even documents (Bengio et al., 2003). However, this is difficult as language is categorical. If we consider each word in a vocabulary as a category, then the dimensionality of a language dataset often becomes large quickly (Rosenfeld, 2000). One reason there is such a variety of statistical language models is that there are various ways of dealing with this fundamental problem. N-gram models are statistical language models that work by taking the co-occurrence of words in a sequence into account.

2.2 N-gram models

The foundation of an n-gram language model is the conditional probability set-up introduced above. An n-gram model is a probabilistic language model that predicts the next term in a sequence of terms (Bengio et al., 2003). The n in n-gram refers to the number of terms in that sequence. Consider the following excerpt from *Jane Eyre*: ‘We had been wandering’. ‘We’ is a uni-gram, ‘We had’ is a bi-gram, ‘We had been’ is a tri-gram, and ‘We had been wandering’ is a 4-gram. Notice that the two tri-grams in this excerpt, ‘We had been’, and ‘had been wandering’, overlap. The use of n-gram models enables us to assign probabilities to both the next sequence of words and just the next word. For instance, consider the two sentences: ‘We had been wandering’, and ‘We had been wangling’. The former is likely to be more frequently encountered in a training corpus. Thus, a 4-gram would assign a higher probability to the next word being ‘wandering’ than ‘wangling’, given the sequence ‘We had been’.

To predict the next word, we have to take the sequence of preceding words into account, which requires knowing the probability of the sequence of words. The probability of a sequence appearing in a corpus follows the chain rule:

$$\begin{aligned} \Pr(\text{We, had, been, wandering}) &= \Pr(\text{We}) \times \Pr(\text{had} | \text{We}) \times \Pr(\text{been} | \text{We, had}) \\ &\quad \times \Pr(\text{wandering} | \text{We, had, been}). \end{aligned}$$

However, the likelihood that more and more words will occur next to each other in an identical sequence becomes smaller and smaller, making prediction difficult. Alternatively, we can approximate the probability of a word depending on only the previous word. This is known as the ‘Markov assumption’ and it allows us to approximate the probability using only the last n words (Brown et al., 1992):

$$\Pr(\text{We, had, been, wandering}) \approx \Pr(\text{We}) \times \Pr(\text{had} | \text{We}) \times \Pr(\text{been} | \text{had}) \times \Pr(\text{wandering} | \text{been}).$$

As a bi-gram model only considers the immediately preceding word, under the Markov assumption, an n -gram model can be reduced to a bi-gram model with n being any number:

$$\Pr(w_n | w_1^{n-1}) \approx \Pr(w_n | w_{n-1})$$

In our application, we use a model that has a mixture of tri-grams, 4-grams, and 5-grams to provide term replacement options for unexpected terms. Longer n-grams such as 5-grams capture more context and potentially generate more accurate predictions. However, there is less chance that a longer sequence of words will fit a specific combination which can lead to sparsity. Therefore, shorter n-grams should still be included to preserve the generalizability of the predictions.

Consider the idiom ‘elephant in the room’. If by mistake, this is written as ‘elephant in the roon’ and correction is needed for the last word, then we can apply both a tri-gram model or a 4-gram model to predict the candidate word:

Using tri-gram model: ‘in the <prediction>’

Using 4-gram model: ‘elephant in the <prediction>’

In the 4-gram model, one more word, ‘elephant’, is included, which provides more context and increased chance of ‘room’ being accurately predicted. However, there are times when the literal location of an elephant may be being described; for instance, ‘elephant in the pond’. The combination of this sequence as a 4-gram might not be in a training dataset and no prediction would be available if only 4-gram models are used. Hence the use of tri-grams to provide other word candidates including ‘in the pond’.

2.3 Text preparation and cleaning

In general, text requires a great deal of preparation and cleaning before it can be analysed. For instance, `compositr` (Luscombe, 2021) is an R package that runs through 13 common steps: remove non-UTF8 characters; insert missing comma spaces; replace misspelled words with most likely equivalent; remove HTML tags; convert case; remove URLs; remove numbers; remove punctuation; replace elongated words; remove extra white space; tokenization; stopwords removal; and finally stem/lemmatize. And these are just some of the steps that may be taken.

Text preparation and cleaning steps and the uncertainty they bring, assume a fairly clean text dataset in the first instance. For instance, if a dataset has been created from Twitter then hashtags and emoji would need to be considered. And if the dataset had been created from OCR then there would likely be systematic errors to be addressed. One benefit of using packages such as `compositr` is that these decisions become repeatable and so their effect on subsequent analysis can be easily understood through repeated running of the full workflow. The notion of consistency provides a measure of this without the need to run the full workflow.

2.4 Consistency

Here we define, mathematically, what we mean by consistency in the case of text. Consider some text extract for which we want a consistency score. If there are n words, of which k are forecast correctly by a model, then the consistency score is the proportion: k/n . We implement this in a dataframe using a binary column. Each word is a row, and a column, ‘correctly_forecast’, tracks the forecast, classifying each word as 1 if it was correctly forecast, and 0 otherwise. Then the internal consistency estimator is given by Equation (1):

$$c_i = \frac{\text{sum 'correctly_forecast' column}}{\text{length of dataset}}. \quad (1)$$

We distinguish between a model trained on the dataset for which consistency is being estimated, which results in an estimate of internal consistency, c_i , and a model trained on another dataset, which results in an estimate of external consistency, c_e , although the procedure is the same once the forecast has been made. We consider the forecast of each word independently, allowing for correlation in future work.

3 The aRianna R package

3.1 Package dependencies

There are a variety of ways to implement an n-gram model within R (R Core Team, 2019) including packages such as `quanteda` (Benoit et al., 2018), `tidytext` (Silge and Robinson, 2016) and `tm` (Feinerer and Hornik, 2019). We use `quanteda` because it has a comprehensive set of functions for conducting text analysis. We also use `tidyverse` (Wickham et al., 2019) for data manipulation:

```
install.packages("quanteda")
install.packages("tidyverse")
library(quanteda)
library(tidyverse)
```

3.2 Function to make the internal dataset

To obtain the internal consistency score, we first construct a dataset from the text that is itself the text being evaluated. The function `aRianna::make_internal_consistency_dataset()` is created for this purpose. The function turns a `body_of_text` into a dataset to compare the text with. To do so, the function first transforms the `body_of_text` into word tokens with the punctuation removed and the letters turned to lowercase:

```
make_internal_consistency_dataset <-
  function(body_of_text) {
    tokens_from_example <-
      quanteda::tokens(body_of_text, remove_punct = TRUE)
    tokens_from_example <-
      quanteda::tokens_tolower(tokens_from_example)
    # code continues...
```

Next, the `aRianna::make_internal_consistency_dataset()` function turns the individual tokens into tri-grams, 4-grams, and 5-grams, keeping only the n-grams that appear in the data for more than once. These common n-grams are stored in a tibble:

```
# ...continue from previous
# Create ngrams from the tokens
toks_ngram <-
  quanteda::tokens_ngrams(tokens_from_example, n = 5:3)

# Convert to tibble so we can use our familiar verbs
all_tokens <-
  tibble::tibble(tokens = toks_ngram[[1]])

# We only want the common ones, not every one.
all_tokens <-
  all_tokens %>%
  dplyr::group_by(tokens) %>%
  dplyr::count() %>%
```

```
dplyr::filter(n > 1) %>%
dplyr::ungroup()
# code continues...
```

Finally, the `aRianna::make_internal_consistency_dataset()` function splits the n-grams into two parts: the first $n - 1$ words and the last word. As a result, a tibble is created with each n-gram containing its original sequence, the first words, and the last word:

```
# ...continue from previous
# Create a tibble that has the first two words in one column then the third
all_tokens <-
  all_tokens %>%
  dplyr::mutate(tokens = stringr::str_replace_all(tokens, "_", " "),
               first_words = stringr::word(tokens, start = 1, end = -2),
               last_word = stringr::word(tokens, -1),
               tokens = stringr::str_replace_all(tokens, " ", "_"),
               first_words = stringr::str_replace_all(first_words, " ", "_")
  ) %>%
  dplyr::rename(last_word_expected = last_word) %>%
  dplyr::select(-n)
}
```

3.3 Function to generate internal consistency

After we have the internal consistency dataset, we need a function to compare the text being evaluated and the internal dataset in order to retrieve the internal consistency score. We define the `aRianna::generate_internal_consistency_score` function to do this. Similar to the `aRianna::make_internal_consistency_dataset` function, this function first transforms the text being evaluated into tri-grams, 4-grams, and 5-grams tokens, then splits the n-grams to their first words and last words:

```
generate_internal_consistency_score <-
function(text_to_check, consistency_dataset){
  # Create tokens with errors
  tokens_from_example_with_errors <-
    quanteda::tokens(text_to_check, remove_punct = TRUE)
  tokens_from_example_with_errors <-
    quanteda::tokens_tolower(tokens_from_example_with_errors)

  # Create ngrams from the tokens with errors
  toks_ngram_with_errors <-
    quanteda::tokens_ngrams(tokens_from_example_with_errors, n = 5:3)

  all_tokens_with_errors <-
    tibble::tibble(tokens = toks_ngram_with_errors[[1]])

  all_tokens_with_errors <-
    all_tokens_with_errors %>%
```

```
dplyr::mutate(ngram = sapply(strsplit(tokens, "_"), length),
             tokens = stringr::str_replace_all(tokens, "_", " "),
             first_words = stringr::word(tokens, start = 1, end = -2),
             last_word = stringr::word(tokens, -1),
             tokens = stringr::str_replace_all(tokens, " ", "_"),
             first_words = stringr::str_replace_all(first_words, " ", "_"))
# code continues...
```

Next, the function uses `dplyr::left_join()` to combine the tibble of text tokens and the internal consistency tibble. The `dplyr::left_join()` function returns all the rows from the text tokens tibble and all the columns from both tibbles (Wickham et al., 2018). By combining the two tibbles, the function generates two additional columns — the `last_word` column are the last words of all n-grams, and the `last_word_expected` contains the last words that are in the internal consistency dataset. By comparing the words in `last_word` and `last_word_expected`, we can identify the ‘unexpected’ words. If a last word is in `last_word` but not `last_word_expected`, it is regarded as an unexpected word and the corresponding word in `last_word_expected` serves as a replacement candidate:

```
# ...continue from previous
all_tokens_with_errors <-
  all_tokens_with_errors %>%
  dplyr::left_join(dplyr::select(consistency_dataset, -tokens),
                  by = c("first_words"))
# code continues...
```

The consistency score is calculated by the number of words that have been predicted by the model divided by the total number of words in the data. The function identifies the counts of words that are predicted by the model and divided by the total number of words in the input text:

```
# ...continue from previous
# Calculate the internal consistency score:
false_count <-
  length(internal_consistency[which(internal_consistency$as_expected == FALSE &
                                   internal_consistency$ngram == 3)])
word_count <- sapply(strsplit(text_to_check, " "), length)
true_count <- word_count - false_count

internal_consistency <- tibble::tibble(
  "as_expected" = true_count,
  "unexpected" = false_count,
  "consistency" = true_count/word_count
)
# code continues...
```

The `aRianna::generate_internal_consistency_score` function also lists the identified text errors and generates replacement candidates to the text errors:


```
# ...continue from previous
# Identify which words were unexpected
unexpected <-
  all_tokens_with_errors_only %>%
  dplyr::mutate(as_expected = last_word == last_word_expected) %>%
  dplyr::filter(as_expected == FALSE) %>%
  dplyr::select(-tokens, -ngram, -as_expected)
}
```

3.4 Function to generate external consistency

The `aRianna::generate_external_consistency_score` function works identically. The only difference is that it compares the text being evaluated with an external consistency dataset that is larger and more general. The external consistency dataset we employed in the application is constructed using a subset of text data from Open WebText Corpus ([Gokaslan and Cohen, 2019](#)). In general, any dataset could be used.

3.5 Demonstration

To install the package and load the library:

```
devtools::install_github("RohanAlexander/arianna")
library(aRianna)
```

For this demonstration, we use the first paragraph of Jane Eyre as the internal text data. A sentence within the paragraph is modified with the intention to contain an error and serves as the text to be evaluated:

```
body_of_text <- "There was no possibility of taking a walk that day.
We had been wandering, indeed, in the leafless shrubbery an hour in
the morning; but since dinner (Mrs. Reed, when there was no company,
dined early) the cold winter wind had brought with it clouds so sombre,
and a rain so penetrating, that further out-door exercise was now out
of the question."

text_to_evaluate <- "when there was na company"
```

The first step is to turn the body of text into the internal consistency dataset. The generated internal consistency dataset is a tibble that contains the tri-grams, 4-grams and 5-grams that appear in the internal text data more than once. Since only 'there_was_no' has more than one occurrence, it is the only n-gram in the internal consistency dataset:

```
internal_consistency_dataset <-
  aRianna::make_internal_consistency_dataset(body_of_text)
internal_consistency_dataset
## # A tibble: 1 x 3
##   tokens      first_words last_word_expected
##   <chr>      <chr>      <chr>
## 1 there_was_no there_was  no
```


The next step is to compare the text to be evaluated with the internal consistency dataset that we created in the previous step. The function `aRianna::generate_internal_consistency_score()` takes two arguments: the text to evaluate, and the internal consistency dataset. The function identifies the word 'na' as an unexpected word and generates the internal consistency score as 0.8. The function also lists 'no' as the replacement of 'na':

```
aRianna::generate_internal_consistency_score(
  text_to_evaluate, internal_consistency_dataset)

## $'internal consistency'
##   as_expected unexpected consistency
## 1           4           1         0.8

## $'unexpected words'
## # A tibble: 1 x 3
##   first_words last_word last_word_expected
##   <chr>      <chr>      <chr>
## 1 there_was  na          no
```

To get the external consistency score, we use the `aRianna::generate_external_consistency_score()` function. The function takes only one argument, the text to evaluate, and compares it with the external consistency dataset. Here, the text to evaluate is 'there was no possibiliti', with the word "possibility" being wrongly spelled. The function identifies the word 'possibiliti' as an unexpected word, and generates the external consistency score 0.75. Because the external dataset is larger, it is capable of providing more replacement candidates for 'possibiliti'. There are no replacements based on 5-grams, so the replacements based on 4-grams are displayed first, followed by the replacements based on tri-grams:

```
text_to_evaluate <- "there was no possibiliti"
aRianna::generate_external_consistency_score(text_to_evaluate)

## aRianna::generate_external_consistency_score(text_to_evaluate)
## $'external consistency'
## # A tibble: 1 x 3
##   as_expected unexpected consistency
##   <int>      <int>      <dbl>
## 1           3           1         0.75

## $'unexpected words'
## # A tibble: 11 x 3
##   first_words last_word last_word_expected
##   <chr>      <chr>      <chr>
## 1 there_was_no possibiliti evidence
## 2 there_was_no possibiliti immediate
## 3 there_was_no possibiliti infrastructure
## 4 there_was_no possibiliti one
## 5 there_was_no possibiliti possibility
## 6 there_was_no possibiliti sound
```

```
## 7 there_was_no possibiliti way
## 8 was_no        possibiliti good
## 9 was_no        possibiliti longer
## 10 was_no       possibiliti more
## 11 was_no       possibiliti wonder
```

The choice of what do now depends on a variety of factors including the purpose of the analysis, and the size of the corpus. One option is to consider only a subset of potential replacements, or to consider only replacements over a certain probability threshold. These are parameters that can be tuned in more sophisticated language models. The choice of which word to change ‘possibili’ to could affect the consistency score if it is not what the model expects. Future work, with more sophisticated language models, could allow tuning in this regard.

As the dataset increases in size, the number of possible replacements will also likely increase. Various options for decreasing the number of possible replacements exist, starting with string distance algorithms and increasing in complexity. It is important to recognise that we are focused on the fact that they were not forecast, rather than what they should be replaced with. However, improving this aspect helps to make the consistency score more relevant and useful, especially in text settings where the cleaning and preparation aspect is likely to be iterative.

4 Discussion

Outside of textbook examples, almost all real-world datasets need to be cleaned and prepared before they can be analyzed. While there are many tools to help with reproducibility of analysis, options for ensuring reproducibility in earlier stages of a typical data science workflow are more limited. Our approach discussed in this paper can help. Our notion of consistency can be tracked at many points while the dataset is being cleaned and prepared. Stages in which it changes substantially can be investigated. While there may be good reasons for such as change, by identifying them it is easier for other researchers to understand the effects of cleaning and preparation on eventual analysis results.

Our approaches are applicable to any dataset, but in this paper we applied them to one based on text. The high-dimensionality of text, and the messiness of it when it comes from real-world sources, means that often considerable cleaning and preparation is done to text datasets. We demonstrated how sequences of words can be used to forecast the very next word. That forecast can be compared with the actual word observed and this difference expressed as a proportion, that we call a ‘consistency score’.

The language model that we used is n-grams. Language models underpinned by n-grams are widely applied in text prediction, spelling-correction, and machine translation (Brown et al., 1992). As demonstrated in our application, we use a n-gram based model to calculate internal and external consistency scores and generate a list of words as text correction candidates. However, n-gram models do not take the linguistic structure of language, nor broader context, into account. For instance, Rosenfeld (2000, p. 1) discusses language in this context, saying that ‘...it may as well be a sequence of arbitrary symbols, with no deep structure, intention or thought behind’. Hence next-word prediction using n-gram-based language models can be limited. Here think of a two-gram involving the word ‘good’ such as ‘good morning’. At scale, these can identify missing or unusual words, and work quickly, but they lack nuance. For instance, an equally reasonable two-gram involving the word ‘good’ is ‘good work’. However, because n-gram models do not

take the surrounding text into account, they do not have the capacity to judge if ‘morning’ or ‘work’ should be suggested as the next word after ‘good’.

In future work we intend to investigate more advanced language models such as word embeddings and Transformers. We will consider pre-trained word embedding models including Word2Vec (Mikolov et al., 2013a) (Mikolov et al., 2013b) and GloVe (Pennington et al., 2014), which place each word in a multi-dimensional space such that distance between words can illustrate their relationship (Bengio et al., 2003). This approach of representing words in vectors has a long history, but the implementation of Bengio et al. (2003) has become the foundation for much subsequent work. For instance, Stoltz and Taylor (2019) use word embeddings to identify how documents relate to concepts.

We will also consider incorporating generative pre-trained Transformer models such as GPT-2 (Radford et al., 2019), GPT-3 (Brown et al., 2020), and BERT (Devlin et al., 2018). These models are based on the Transformer architecture (Vaswani et al., 2017). Proposed by Vaswani et al. (2017) in 2017, the Transformer network architecture for neural networks has been found to outperform both recurrent neural network-based and convolutional neural network-based models in computational efficiency (Vaswani et al., 2017). Since the emergence of the Transformer model, most of the representative pre-trained models are built on this architecture (Hanretty et al., 2018).

Although the level of the consistency score itself is not of major importance compared with its changes, improving the score by using more sophisticated models would be beneficial in terms of identifying incorrect words. In our future work we are interested in implementing different language models to compare their effectiveness in the task of generating consistency scores and providing suggested text corrections. We intend to use our future findings to continuously improve the aRianna package. We are also interested in expanding the notion of consistency beyond text and applying it to other datasets.

References

- Bengio, Y., Ducharme, R., Vincent, P., and Jauvin, C. (2003). A neural probabilistic language model. *Journal of Machine Learning Research*, 3(Feb):1137–1155.
- Benoit, K., Watanabe, K., Wang, H., Nulty, P., Obeng, A., Müller, S., and Matsuo, A. (2018). quanteda: An r package for the quantitative analysis of textual data. *Journal of Open Source Software*, 3(30):774.
- Brown, P. F., Della Pietra, V. J., Desouza, P. V., Lai, J. C., and Mercer, R. L. (1992). Class-based n-gram models of natural language. *Computational linguistics*, 18(4):467–480.
- Brown, T. B., Mann, B., Ryder, N., Subbiah, M., Kaplan, J., Dhariwal, P., Neelakantan, A., Shyam, P., Sastry, G., Askell, A., et al. (2020). Language models are few-shot learners. *arXiv preprint arXiv:2005.14165*.
- Devlin, J., Chang, M.-W., Lee, K., and Toutanova, K. (2018). Bert: Pre-training of deep bidirectional transformers for language understanding. *arXiv preprint arXiv:1810.04805*.
- Feinerer, I. and Hornik, K. (2019). *tm: Text Mining Package*. R package version 0.7-7.
- Gokaslan, A. and Cohen, V. (2019). Openwebtext corpus. <http://Skylion007.github.io/OpenWebTextCorpus>.
- Hanretty, C., Lauderdale, B. E., and Vivyan, N. (2018). Comparing strategies for estimating constituency opinion from national survey samples. *Political Science Research and Methods*, 6(3):571–591.
- Luscombe, A. (2021). *compositr: Efficient tools for preprocessing text data in R*. R package version 0.0.0.9000.
- Mikolov, T., Chen, K., Corrado, G., and Dean, J. (2013a). Efficient estimation of word representations in vector space. *arXiv preprint arXiv:1301.3781*.
- Mikolov, T., Sutskever, I., Chen, K., Corrado, G. S., and Dean, J. (2013b). Distributed representations of words and phrases and their compositionality. In *Advances in neural information processing systems*, pages 3111–3119.
- Pennington, J., Socher, R., and Manning, C. D. (2014). Glove: Global vectors for word representation. In *Proceedings of the 2014 conference on empirical methods in natural language processing (EMNLP)*, pages 1532–1543.
- R Core Team (2019). *R: A Language and Environment for Statistical Computing*. R Foundation for Statistical Computing, Vienna, Austria.
- Radford, A., Wu, J., Child, R., Luan, D., Amodei, D., and Sutskever, I. (2019). Language models are unsupervised multitask learners. *OpenAI Blog*, 1(8):9.
- Rosenfeld, R. (2000). Two decades of statistical language modeling: Where do we go from here? *Proceedings of the IEEE*, 88(8):1270–1278.
- Silge, J. and Robinson, D. (2016). tidytext: Text mining and analysis using tidy data principles in r. *JOSS*, 1(3).

- Stoltz, D. S. and Taylor, M. A. (2019). Concept mover’s distance: measuring concept engagement via word embeddings in texts. *Journal of Computational Social Science*, 2(2):293–313.
- Vaswani, A., Shazeer, N., Parmar, N., Uszkoreit, J., Jones, L., Gomez, A. N., Kaiser, Ł., and Polosukhin, I. (2017). Attention is all you need. In *Advances in neural information processing systems*, pages 5998–6008.
- Wickham, H., Averick, M., Bryan, J., Chang, W., McGowan, L. D., François, R., Grolemund, G., Hayes, A., Henry, L., Hester, J., Kuhn, M., Pedersen, T. L., Miller, E., Bache, S. M., Müller, K., Ooms, J., Robinson, D., Seidel, D. P., Spinu, V., Takahashi, K., Vaughan, D., Wilke, C., Woo, K., and Yutani, H. (2019). Welcome to the tidyverse. *Journal of Open Source Software*, 4(43):1686.
- Wickham, H., François, R., Henry, L., and Müller, K. (2018). *dplyr: A Grammar of Data Manipulation*. R package version 0.7.6.