# Towards Evaluating Creativity in Language

*Matey Krastev*

A dissertation submitted in partial fulfilment
of the requirements for the degree of
**Bachelor of Science**
of the
**University of Aberdeen**.



Department of Computing Science

2022

# Declaration

No portion of the work contained in this document has been submitted in support of an application for a degree or qualification of this or any other university or other institution of learning. All verbatim extracts have been distinguished by quotation marks, and all sources of information have been specifically acknowledged.


Signed:


Date: 2022

# Abstract

An expansion of the title and contraction of the thesis.

# Acknowledgements

Much stuff borrowed from elsewhere.

# Contents

# Chapter 1

# Introduction

## 1.1 Inspiration

At the heart of the field of artificial intelligence is the concept of reproducing aspects defining human intelligence through rigourous examination and replication of the mechanisms that drive progress. This is a uniquely multi-faceted problem with a multitude of approaches, each tailored to a very specific aspect or manifestations of intelligence. Naturally, intelligence assumes following a logical pathway to arrive at sensible conclusions that interact with the real world in a beneficial way. This can be viewed through the lens of methodical, defined process that always follows a certain formula. An organism, assumed to be intelligent, might always follow such formulaic actions given a set of prerequisite conditions to accomplish a defined goal. Certain schools of thought theorise that there is such an order in every little action, and such thread of logicality interweaves every law of nature, known or not. Then, follows the question, can we recognise and define such a thread for the ambitious field of creativity? We seek not to properly define or constrain the subject of creativity, rather, we explore markers of what could only be a subset of the very broad field of human creativity.

The subject of the present document is exclusively the study of linguistic creativity. Henceforth, we seek to: confirm prior results of the research of psycholinguistics, affirm that hypotheses and conclusions drawn from them correlate highly with certain manifestations of creativity, and make firm the subject matter of creativity, that is, provide tools that may be used for exploration and analysis of specific creative features found in text.

**Chapter 2**

# Related Work

## 2.1 Challenge Landscape

### 2.1.1 Part of Speech Tagging

### 2.1.2 Word Sense Disambiguation

### 2.1.3 Sentiment Analysis

### 2.1.4 Phonetic Analysis

### 2.1.5 Natural Language Generation

- top KP sampling

- challenges

## 2.2 Creative Measures

- Burstiness of verbs and derived nouns: Patterns of language are sometimes 'bursty' Pierrehumbert (2012). This paper presents an analysis of text patterns for domain X. measures include XYZ...

- 

## 2.3 Tools

**Chapter 3**

# Methodology

## 3.1 Datasets

Some introductory text on the purpose and uses of the datasets...

### 3.1.1 Brown Corpus

The Brown Corpus (Francis and Kucera, 1979) is a widely used corpus in the field of computational linguistics, noted for the small variety of genres of literature it contains. The Corpus itself is founded on a compilation of American English literature from the year 1961. It is also small in terms of size, totalling around one million words, at least compared to modern corpora, which we also explore later on. The corpus also suffers from the issue of recency, as the works and language may be outdated for modern speakers of English.

Of interest is the fact that the corpus has been manually tagged for parts of speech, a process that tends to be error-prone. As we will see later on, this fact has implications in terms of the supervised learning algorithms we implement for creativity evaluation. Still, we opt to utilize it primarily for prototyping purposes and drawing preliminary conclusions about the effectiveness of the implemented algorithms, rather than in-depth analysis and publication of results.

### 3.1.2 Project Gutenberg

Project Gutenberg[1] is a large collection of more than 50,000 works available in the public domain. The collection contains literature from various years and various genres and thus is suitable for training and evaluation of the developed benchmarks in the context of creativity study.

As the Project does not offer an easy to process copy of its collection, we turn to the work of Gerlach and Font-Clos (2018). The team developed a catalogue for on-demand download of the entire set of books available on the Project Gutenberg website, intended for use in the study of computational linguistics. The tool avoids the overhead of writing a web-scraper or a manual parser for the downloadable collections of Project Gutenberg books made available by third parties, as well as enables easy synchronization of newly released literature. Instead, we are only required to develop a simple pipeline for the data to be fed into the utilized systems.

### 3.1.3 Hierarchical Neural Story Generation

In their work, Fan et al. (2018) trained a language model for text generation tasks on a dataset comprised of short stories submitted by multiple users given a particular premise (a prompt or a theme) by another user. [MK: Give an example for how one such short story would look like.] The dataset

---

[1] https://www.gutenberg.org/

in question is technically referred to a series of posts and comments (threads) to them on the popular social media platform REDDIT, and more tightly, the *subreddit* forum R/WRITINGPROMPTS. The authors of the work Fan et al. (2018) have made the dataset available for public use, and we have used it for the purpose of evaluating the performance of our creativity benchmarks. As described by the authors on their GitHub page[1], the paper models the first 1000 tokens (words) of each story.

[MK: How do we use this dataset? You should describe the process of how we use it. ]

### 3.1.4 WordNet

WordNet(Fellbaum, 1998) is a lexical database of semantic relations between words that links words into semantic relations including synonyms, hyponyms, and meronyms. The synonyms are grouped into sunsets with short definitions and usage examples. It can thus be seen as a combination and extension of a dictionary and thesaurus (Wikipedia contributors, 2023b).

For our specific use cases, we have identified it as a valuable resource in terms of relational representation of words in semantic space. In the given context, this enables us to traverse a semantic graph for synonyms and related words for the goal of enriching potential similarity between the set of creative parts of speech (i.e., nouns, adjectives, adverbs), which we narrow down our scope to in particular.

### 3.1.5 Numerical representations of semantic tokens

[MK: Potentially move this section to background work] The idea of representing words or lexical tokens as numerical vectors (or even scalars) is hardly new. For example the SimLex-999 dataset (Hill et al., 2015) gives values on a scale from 0 to 10, like the examples below, which range from near-synonyms (vanish, disappear) to pairs that scarcely seem to have anything in common (hole, agreement):

| word1 | word2 | score |
|-------|-------|-------|
| vanish | disappear | 9.8 |
| hole | agreement | 1.2 |

**Table 3.1:** Example Simlex-999 pairs

Early work on affective meaning by Osgood et al. (1957) found that words varied along three important dimensions of affective meaning:

- valence: the pleasantness of the stimulus

- arousal: the intensity of emotion provoked by the stimulus

- dominance: the degree of control exerted by the stimulus

Osgood et al. (1957) noticed that in using these 3 numbers to represent the meaning of a word, the model was representing each word as a point in a three-dimensional space, a vector whose three dimensions corresponded to the word's rating on the three scales. This revolutionary idea that word meaning could be represented as a point in space (e.g., that part of the meaning of

---

[1] `https://github.com/facebookresearch/fairseq/blob/main/examples/stories/README.md`

heartbreak can be represented as the point $[2.45, 5.65, 3.58]$) was the first expression of the vector semantics models that we introduce next. [MK: **You can paraphrase this**]

### Word2Vec

Mikolov et al. (2013) show in their work that words may be represented as dense vectors in $N$-dimensional space, and we can perform mathematical operations on them that may yield effective results in terms of word representation.

### Measuring distance in vector representations of semantic tokens

Intuition tells us that the dot product of vectors in $N$-dimensional space will grow when the set of vectors has similar values and decrease when the values are not similar. Thus, we can then construct the following metric for semantic similarity between vector representations of words:

$$D(v, w) = v \times w = \sum_{i=1}^{N} v_i w_i = v_1 w_1 + v_2 w_2 + \cdots + v_N w_N$$

The current metric, however, suffers from the problem that vectors of higher dimensions will inevitably be larger than vectors with lower dimensions. Furthermore, embedding vectors for words that occur frequently in text, tend to have high values in more dimensions, that is, they correlate with more words. The proposed solution is to normalize using the **vector length** as defined:

$$|v| = \sqrt{\sum_{i=1}^{N} v_i^2}$$

Therefore, we obtain the following:

$$\text{Similarity}(v, w) = \frac{v \times w}{|v||w|} = \frac{\sum_{i=1}^{N} v_i w_i}{\sqrt{\sum_{i=1}^{N} v_i^2} \sqrt{\sum_{i=1}^{N} w_i^2}}$$

This product turns out to be the same as the cosine of the angle between two vectors:

$$\frac{a \times b}{|a||b|} = \cos(\theta)$$

Therefore, we will call this metric the **cosine similarity** of two words. As mentioned, the similarity grows for vectors with similar features along the same dimensions. Note the boundaries of said cosine metric: we get $-1$ for vectors which are polar opposites, 0 for orthogonal vectors, and 1 for equivalent vectors. Of note is the fact that such learned vector embeddings only have values in the positive ranges, thus, it is impossible to have negative values for the cosine similarity ($\text{Similarity}(a, b) \in [0, 1]$).

Contrary to it, we also identify the metric of **cosine distance** between two vectors, as one minus the similarity of the vectors, or:

$$\text{Distance}(v, w) = 1 - \text{Similarity}(v, w)$$

The cosine distance may prove useful when dealing with minimisation problems as is often the case with machine learning.

## 3.2 Metrics

**Chapter 4**

# Design

In the following section, we introduce concepts behind the design of the developed library and how those will be implemented in the application. We also discuss the technology choices we have made and the reasoning behind them. Furthermore, we take a look at how the application will be structured and how it will be distributed, and how the users will be able to interact with it via a command-line interface, or apply declared methods and classes inside their own applications. Finally, we discuss the delivery of a documentation and a user guide, as well as the testing and validation of the application.

## 4.1 Functional Requirements

### 4.1.1 Accuracy, etc.

## 4.2 Non-functional Requirements

The package henceforth needs to satisfy a list of viable non-functional requirements, which we will list below.

### 4.2.1 Efficacy

The implemented algorithms must be viable to deploy both in small-scale and for large-scale applications. This means that the algorithms must be able to scale to large amounts of data, while also being able to run on a single machine. In our case, a user should be able to quickly evaluate their texts across several metrics, however, as we also strive to apply this benchmark to large-scale corpora, we must also ensure that the algorithms are scalable. We will therefore not only seek to reiterate on the existing literature and implement the most promising algorithms for the task of creativity evaluation, but also optimize them for deployment on HPC clusters.

### 4.2.2 Memory Requirements

The old adage that *memory is cheap* is not entirely true. While it is true that memory is cheap, it is also true that memory is not free (*and no, we cannot "just download more RAM"*). In fact, some LLMs simply tend to be too large to reliably fit within the memory constraints of a personal computer. [MK: We should probably cite some sources here.] Furthermore, model accuracy tends to grow with the size of the neural network and the size of the used vocabulary. Naturally, we then need to seek a compromise on the size of the models we use, as we cannot:

1. Use too large models during the research stage of the project, where we aim to process large corpora, evaluate the performance of the algorithms on them and make conclusions

about the data. If we do aim to speed up this process, it is very likely that we would benefit from parallel computing — but processing large sizes of text in parallel has a non-negligible likelihood of running out of allocated memory even on some HPC clusters.

2. Force users to run too large models on their personal computers, as this would be a very poor user experience. We do not plan to hardcode any models (large or small) in the application, however, the provided guides will reference certain smaller-scale pretrained LLMs. Naturally, we would provide a way for more experienced and more capable organizations or individuals to run larger models with minimal effort.

## 4.3 Technology Choices

### 4.3.1 Python

We will be using Python version 3.10.X as shipped by the Anaconda software package. We are aware that Python 3.11 brings non-negligible optimizations and faster execution speed for some Python scripts, however, in light of the fact that the Anaconda distribution is still shipping Python 3.10.X, we will be using that version for the time being. We will be using the Anaconda distribution as it is a very popular and mature distribution of Python, which is also very easy to install and use. It also comes with a large number of pre-installed packages, which will be very useful for the current developer experience.

### 4.3.2 PyTorch

PyTorch "is a machine learning framework based on the Torch library, used for applications such as computer vision and natural language processing, originally developed by Meta AI and now part of the Linux Foundation umbrella. It is free and open-source software released under the modified BSD licence", as described by Wikipedia contributors (2023a).

Any models used in the application will be implemented in PyTorch, as it is a very popular framework for deep learning and natural language processing. It is also a very flexible framework, which allows for easy implementation of new models and algorithms. Furthermore, it is a very popular framework, which means that there is a large community of developers and researchers who have already implemented many of the algorithms we plan to use. This means that we can easily reuse their code and adapt it to our needs.

**Comparison with TensorFlow**

### 4.3.3 NLTK

NLTK is a key Python library for natural language processing, primarily built for education purposes and managed as an open-source software, built to be relatively modular and lightweight. Commonly used by researchers and students for understanding and implementing algorithms for NLP tasks, it is a relatively popular and mature framework with a healthy extension ecosystem, where contributors are able to write their own modules and share them with the community.

NLTK

### 4.3.4 SpaCy

SpaCy is an open-source Python library for advanced natural language processing, designed to be easily used in production environments and implementing pipelines for enhanced NLP tasks.

Whereas NLTK is primarily used for research and education, SpaCy is commonly being applied in industry environments.

## 4.4 Code Style

### 4.4.1 PEP8

We will be using PEP8[1] (van Rossum et al., 2001) as our code style guide. PEP8 is a style guide for Python code, which is maintained by the Python Software Foundation. It is a very popular, and comprehensive style guide, widely used by many Python developers and organizations. It covers a wide range of topics, including naming conventions, indentation, line length, whitespace, comments, "docstrings" (short for documentation strings, or, more specifically, comments that explain the way a given procedure or a class works, inside the code itself), and so on.

### 4.4.2 Docstrings

### 4.4.3 Linting

### 4.4.4 Testing

### 4.4.5 Code Review

### 4.4.6 Version Control

The outlined project

Git and GitHub

## 4.5 Documentation

### 4.5.1 Documentation Framework

We use Sphinx in this household.

### 4.5.2 Sphinx

### 4.5.3 Hosting

---

[1]`https://peps.python.org/pep-0008`

**Chapter 5**

# Implementation

**Chapter 6**

# Evaluation

**Chapter 7**

# Discussion

# Bibliography

Fan, A., Lewis, M., and Dauphin, Y. (2018). Hierarchical Neural Story Generation. arXiv:1805.04833 [cs].

Fellbaum, C. (1998). *WordNet: An Electronic Lexical Database*. Bradford Books.

Francis, W. N. and Kucera, H. (1979). Brown corpus manual. *Letters to the Editor*, 5(2):7.

Gerlach, M. and Font-Clos, F. (2018). A standardized project gutenberg corpus for statistical analysis of natural language and quantitative linguistics. *CoRR*, abs/1812.08092.

Hill, F., Reichart, R., and Korhonen, A. (2015). SimLex-999: Evaluating semantic models with (genuine) similarity estimation. *Computational Linguistics*, 41(4):665–695.

Mikolov, T., Chen, K., Corrado, G., and Dean, J. (2013). Efficient Estimation of Word Representations in Vector Space. arXiv:1301.3781 [cs].

Osgood, C. E., Suci, G. J., and Tannenbaum, P. H. (1957). *The measurement of meaning*. University of Illinois press.

Pierrehumbert, J. B. (2012). Burstiness of Verbs and Derived Nouns. In Santos, D., Lindén, K., and Ng'ang'a, W., editors, *Shall We Play the Festschrift Game? Essays on the Occasion of Lauri Carlson's 60th Birthday*, pages 99–115. Springer Berlin Heidelberg, Berlin, Heidelberg.

van Rossum, G., Warsaw, B., and Coghlan, N. (2001). Style guide for Python code. PEP 8, Python Software Foundation.

Wikipedia contributors (2023a). Pytorch — Wikipedia, the free encyclopedia. `https://en.wikipedia.org/w/index.php?title=PyTorch&oldid=1146375871`. [Online; accessed 2-April-2023].

Wikipedia contributors (2023b). Wordnet — Wikipedia, the free encyclopedia. `https://en.wikipedia.org/w/index.php?title=WordNet&oldid=1143619785`. [Online; accessed 14-March-2023].