
Detective Novel Text Generation using RNN LSTM

Annysia Dupaya, Pranav Rajan, Safi Benyahia
KTH Royal Institute of Technology
dupaya, pranavr, safib@kth.se

Abstract

1 This project focuses on Natural Language Processing from simple recurrent-neural
2 networks(RNN) to LSTMs and if time permitting, some of the modern approaches
3 in the NLP toolkit including word embeddings (Word2Vec, Bert, Glove) and the
4 HuggingFace Transformer Library.

5 0.1 Project Goal

6 We plan to implement default Project 3 + the extensions.

7 0.2 Overview

8 The goal of this project is to investigate and implement different neural network architectures
9 including recurrent-neural networks(RNNs), LSTMS, GRU, Word Embeddings and Transformers for
10 the text generation problem. Traditional Natural Language Processing has been based on statistical
11 and computational linguistic approaches including NER, Logistic Regression that involve significant
12 feature engineering to produce useful results. With the success of Deep Learning and the Transformer
13 architecture, text generation and natural language processing has been transformed with the way
14 the models are able to find and learn patterns in the data resulting in tools such as ChatGPT, highly
15 accurate language translation tools, search tools and more.

16
17 In this project we will be extending Assignment 4 to experiment and understand the tradeoffs
18 and implementations of different neural network approaches to text generation. Using the RNN
19 architecture as a baseline, we will train neural network models on text chosen from Project Gutenberg
20 to learn patterns and generate unique text with certain characteristics such as writing style, theme and
21 language. We will compare the baseline RNN model with one and two-layer LSTM implementations
22 both qualitatively and quantitatively using evaluation techniques such as prediction loss, n-gram
23 frequency and quality of generated text.

24 0.3 Dataset Used

25 The group would like to propose the works of Conan Doyle (Sherlock Holmes) and Maurice Leblanc
26 (Arsène Lupin) from the Project Gutenberg site. We will do an 80-10-10 split for training, testing,
27 and validation. For example in our proposed base project using the Sherlock Holmes dataset with
28 12 books, to account for writing style changes that can happen over time, we will split each book
29 80-10-10.

30 For example, in one book, we use 80% of the first portion, then the remaining 20% with the latter
31 portion of the book.

32 0.4 Technology Used

33 The group is considering using the following technologies:

- 34 • PyTorch
- 35 • PyTorch Lightning
- 36 • HuggingFace Transformers Library
- 37 • Spacy
- 38 • NLTK

39 0.5 Implementation

40 The group plans to implement the code from scratch using the technologies previously mentioned.
 41 The group will be extending Assignment 4, and following closely with the architecture used by
 42 Vaswani et al. [1] Some papers that we may reference include the original ChatGPT paper by Radford
 43 et al. [2] and the ChatGPT2 paper. All the members of this group have taken/are currently taking a
 44 class taught by Professor Johan Boye from the Speech, Hearing and Language Division and thus are
 45 familiar with the basics of how to set this project up.

46 0.6 Initial Set of Experiments

47 The group's first set of initial experiments will be to investigate how the different training parameters
 48 affect the model (ex. batch size and learning rate). Another set of experiments the group plans to do is
 49 checking the quality of the generated text of greater length. We plan on using strategies taught during
 50 the language engineering course such as using the BERTScore, entropy and the ideas described in [3].

51 0.7 Milestones Timeline

- 52 • E grade:
 - 53 – Complete the basic assignment as specified for Default Project 3
- 54 • D-C grade:
 - 55 – Train on two datasets: Conan Doyle + Maurice Leblanc
 - 56 – Check the performance when training the two datasets separately
 - 57 – Create a mixed dataset and compare the performance on the mixed and initial datasets
- 58 • B-A grade:
 - 59 – Use words as the basic entry in the network and use Glove for word embeddings and
 - 60 investigate Byte-Pair Encoding (BPE) tokenization.
 - 61 – Compare the quality of generated text when using word embeddings or BPE tokeniza-
 - 62 tions compared to the base project.

63 0.8 Member Skills and Knowledge

- 64 • Annysia - has knowledge of RNNs in the course Scalable Machine Learning, mostly in
 65 their usage in an encoder-decoder configuration. She also has experience working with
 66 information systems and textual information. She hopes to gain deeper understanding of
 67 RNNs by implementing them with LSTM architecture.
- 68 • Pranav - first time taking a deep learning class but has completed DD2477 and currently
 69 taking DD2417 with Professor Johan Boye and is familiar with the language engineering
 70 techniques needed for the assignment. Also had exposure to machine learning/basic deep
 71 learning theory from DD1420 and DD2434. By the end of this project, Pranav hopes to
 72 understand how to implement the different architectures in the assignment, the tradeoffs and
 73 have a deeper understanding of neural network approaches to language engineering.
- 74 • Safi - Has completed the Machine Learning class and is currently taking a Language
 75 Engineering course with exposure to word embedding techniques like word2vec and GloVe.
 76 Safi is familiar with the basic concepts of neural networks and aims to gain practical
 77 understanding of implementing RNN architectures with LSTM layers, learning how they
 78 process sequential data for text generation, and developing skills in hyperparameter tuning
 79 to improve model performance.
 80 Also familiar with the language engineering techniques,

81 **0.9 Target Project Grade**

82 The group aims to get an A for this project.

83 **References**

- 84 [1] Ashish Vaswani et al. *Attention Is All You Need*. 2023. arXiv: 1706.03762 [cs.CL]. URL:
85 <https://arxiv.org/abs/1706.03762>.
- 86 [2] Alec Radford and Karthik Narasimhan. “Improving Language Understanding by Generative Pre-
87 Training”. In: 2018. URL: <https://api.semanticscholar.org/CorpusID:49313245>.
- 88 [3] Ari Holtzman et al. “The curious case of neural text degeneration”. In: *arXiv preprint*
89 *arXiv:1904.09751* (2019).