

# Machine Learning in Computational Literary Studies

## An Overview of methods and applications

Fotis Jannidis and Leonard Konle  
Würzburg University

### Introduction

Machine learning in its facets is an integral part of the methodological repertoire of Computational Literary Studies. Especially for the analysis of genre, period and authorship supervised and unsupervised learning has enabled researchers to base their theories on large text collections.

### Unsupervised Learning

Figure 1 shows the result of a combination of the unsupervised methods Topic Modeling (Blei et al. 2013) and Dimension Reduction to examine the genre of novels due to its publisher. Topic Modeling allows to discover distribution of hundreds of topics over large collections of text. By using the share of topics in documents as features, we can shrink this information into 2-dimensional space with Uniform Manifold Approximation (McInnes and Healy 2018).

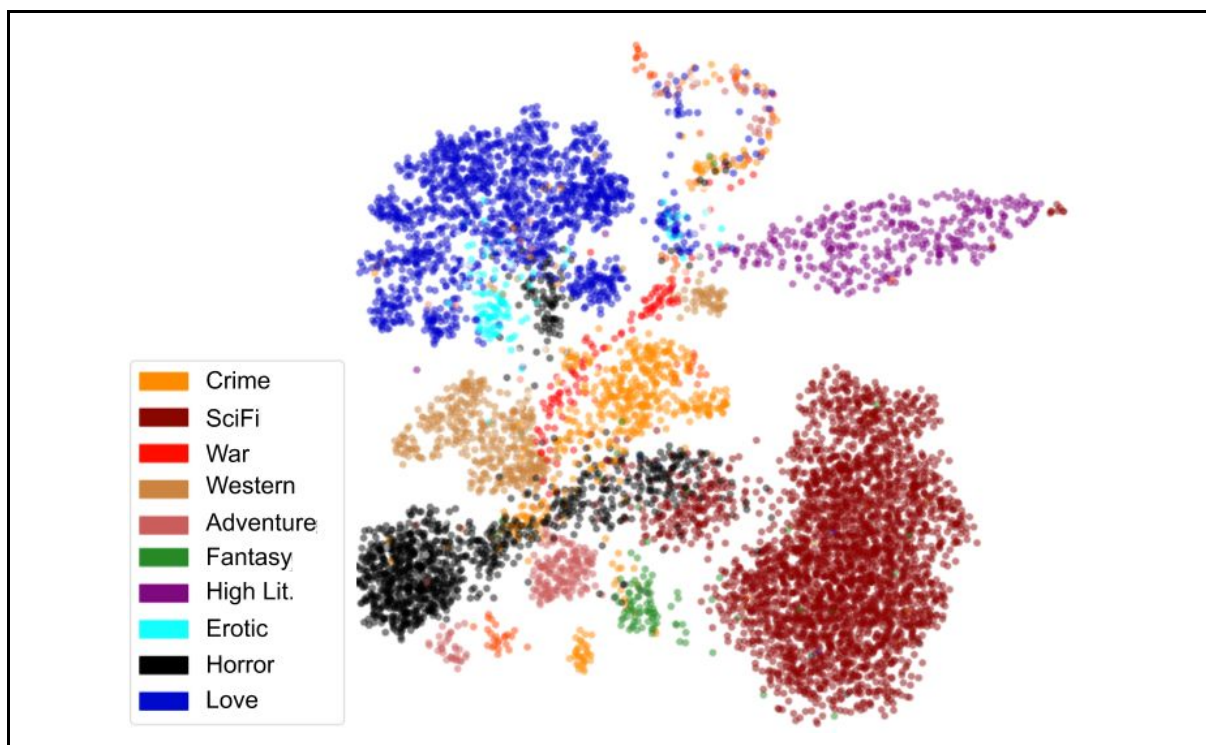


Figure 1: Genre novels plottet based on their topic proportions.

## Supervised Learning

Jannidis et al. 2018 aims at overcoming inconsistent use of quotation in novels to outline direct speech by creating annotated data from a collection of texts with uniform quotation to train a neural network. In the training phase quotation is removed from the data resulting in a model predicting direct speech without relying on quotation marks.

## Deep Learning and Word Embeddings

The combination of word embeddings (Mikolov 2013), Attention and Recurrent Neural Networks (Wang et. al. 2016) allows a more detailed view on the literary genres mentioned in Fig.2. While embeddings and recurrent blocks achieve a .96 f1-score on classifying the genre of segments with just 250 words, attention gives insights on which words are crucial for this decision (see Fig. 2).

[...] Before him lay not only the entire engine with the receiver mechanism, which picked up and processed the remote control signals, but also the generator for building the artificial gravity field in the cabin, the remote camera, which radiated its impulses onto the screen of the remote controller, and finally the field unit of two heavy disintegrators, which were rigidly built into the outer casing of the spaceship. NE saw that he had made a find. However, he also saw that it was now primarily a question of whether he could keep it. The distance from NE to NE was to be overcome for a vehicle of this kind in less than one hour. So if he didn't want to get in dangerous proximity to the enemy base, he had to act fast. With a few quick grips he disconnected the supply line to the remote control receiver and interrupted the contacts, so that the engine could no longer be influenced by any signal coming from outside. Then he examined the engine itself and found that it had stopped moving at the same moment. [...]

Figure 2: Segment from Kurt Mahr (1962): Perry Rhodan Band 47: *Gom antwortet nicht*. Darker red indicates higher attention of the neural network. Named Entities are replaced by NE. (Translated)

## Bibliography

- Blei, D. M., Ng, A. Y., & Jordan, M. I. (2003): Latent dirichlet allocation. *Journal of machine Learning research*, 3(Jan), 993-1022.
- Jannidis, F., Konle, L., Zehe, A., Hotho, A. and Krug, M. (2018): Analysing Direct Speech in German Novels. Dhd2018 Book of Abstracts.
- McInnes, L, Healy, J, (2018): *UMAP: Uniform Manifold Approximation and Projection for Dimension Reduction*, ArXiv e-prints 1802.03426
- Mikolov, T., Chen, K., Corrado, G., & Dean, J. (2013): Efficient estimation of word representations in vector space. *arXiv preprint arXiv:1301.3781.o*
- Wang, Y., Huang, M., & Zhao, L. (2016): Attention-based LSTM for aspect-level sentiment classification. In *Proceedings of the 2016 conference on empirical methods in natural language processing* (pp. 606-615).