

Identifying the relation between author gender and Dutch suspense novels

Noa Visser, Utrecht University

Corina Koolen, Huygens ING

In this research, the relation between author gender and the text of suspense novels is explored. Novels have been classified on author gender before, but with a main focus on the identification of author gender (Koppel et al., 2002). It is important to control variables in the corpus and not to overgeneralise results when using NLP techniques to research the relation between gender and text (Koolen & Van Cranenburgh, 2017). Therefore, this research will focus on novels within one genre, namely Dutch suspense novels.

Dutch suspense novels are an interesting genre, as they contain the subgenre of ‘literary thriller’, which has a very specific content and style. Another signature characteristic of this genre is that most authors are women. Therefore the author Paul Goeken, a Dutch man, published his literary thrillers under the female pen name Suzanne Vermeer as a marketing strategy (Van Lieshout, 2013). His novels under the name of Suzanne Vermeer are bestselling novels, whereas the suspense novels he has published under his own name are less successful. This suggests that in Dutch suspense novels subgenres are influenced by gender. Therefore, this research focuses on one main question: What is the relation between author gender and suspense novels? This short paper will be a first step in this research. In order to answer this question, gender will not be treated as a binary variable, as gender is established to be a more complex construct than a zero or one-question (Butler, 1998).

The corpus exists of 134 readily available recent Dutch-language suspense novels published between 2007 and 2012. As can be seen in Table 1, the corpus contains works written by men, women, Paul Goeken under the pen name Suzanne Vermeer and authors of unknown gender and duos of mixed gender, with a maximum of three works per author. Both translated and originally Dutch novels were used, in order to have a sufficient amount of novels and different authors of different genders. Each novel was vectorized with a Term-frequency - Inverse document-frequency Vectorizer, which tokenized the documents per word pair. An initial experiment was a bottom-up clustering, using K-means to see if meaningful groupings related to gender could be identified, which proved unsuccessful.

	Women	Men	Unknown/multiple
Translated	31	59	7
Original	23	11	3

Table 1: Division corpus

We decided to explore to what extent author gender could be identified by an SVM, a known well-performing algorithm on the classification of author gender (Koppel et al., 2002, Peersman et al., 2011, Vicente et al., 2015). Author gender was classified in three ways: male/female, female/non-female and male/non-male¹, using cross validation with a minimum of three different authors per fold. The three models were necessary for an accurate analysis of the influence of gender. The addition of female/non-female and male/non-male classification enabled us to analyse the classification of author gender without treating female and male authors as opposites and also to include authors of unknown and multiple gender.

	Male/female	Female/non-female	Male/non-male
Accuracy	79.0%	82.1%	78.4%
F1 score	75.5	75.5	78.8
Precision	76.9	84.1	80.6
Recall	74.1	68.5	77.1

Table 2: Results SVM author gender classification

The results in Table 2 indicate that a highly distinguishable selection of novels by female authors is represented in the corpus, as the female/non-female SVM has the highest accuracy and precision. This contains works in the subgenre literary thriller. Surprisingly, the novels of Lieneke Dijkzeul, a Dutch woman, were misclassified by all the models. A keyword analysis comparing works of Lieneke Dijkzeul to works of other Dutch women showed that the word *geweer* (rifle) differentiates her works, suggesting that storyline related word choice influences the classification on author gender. However, altering *geweer* in *pistool* (pistol), which is used in the novels of the other Dutch women, does not change the results. A thorough keyword and vector analysis is needed to interpret the difference between the novels of Lieneke Dijkzeul and other Dutch women. Remarkably, the works by Suzanne Vermeer were all classified as non-female, even though they were specifically marketed as to be written by a woman.

In conclusion, the high accuracy and precision score of the female/non-female SVM shows that the word pairs within these novels are closely related. This indicates that suspense novels written by women have a distinctive style and word use. However, the misclassification of the works of

¹ The words female, male, non-female and male are used to clearly state the distinction between the different classifications. As this research is focussed on gender and not on sex, the classification female/non-female should be interpret as classified as written by a woman/not classified as written by a woman etc.

Suzanne Vermeer and Lieneke Dijkzeul suggests that marketing strategies are more influential on labelling certain suspense novels as literary thrillers, than a writing style which is characteristic to women. Further research is needed to investigate whether writing style and word use within suspense novels are influenced by author gender or whether gendered subgenres such as literary thriller are created as a marketing strategy. Therefore, performing a cross genre evaluation would be insightful. As it is suspected that the inclusion of both original and translated works in the corpus influences the results, we aim to only include original novels in the dataset.

References

- Butler, J. (1998). *Imitation and Gender Insubordination*, pages 722–730. Blackwell Publishers.
- Koolen, C., & van Cranenburgh, A. (2017, April). These are not the Stereotypes You are Looking For: Bias and Fairness in Authorial Gender Attribution. In *Proceedings of the First ACL Workshop on Ethics in Natural Language Processing* (pp. 12-22).
- Koppel, M., Argamon, S., and Shimoni, A. R. (2002). Automatically Categorizing Written Texts by Author Gender. *Literary and Linguistic Computing*, 17(4):401–412.
- Peersman, C., Daelemans, W., and Van Vaerenbergh, L. (2011). Predicting age and gender in online social networks. In *Proceedings of the 3rd International Workshop on Search and Mining User-generated Contents*, SMUC '11, pages 37–44, New York, NY, USA. ACM.
- Van Lieshout, M. (2013, June 4). Mysterieuze schrijfster met een uitgekiende strategie. *De Volkskrant*. Retrieved from <https://www.volkskrant.nl/nieuws-achtergrond/mysterieuze-schrijfster-met-een-uitgekiende-strategie~b70deaef/>
- Vicente, M., Carvalho, J. P., and Batista, F. (2015). Using unstructured profile information for gender classification of portuguese and english twitter users. In Sierra-Rodríguez, J.-L., Leal, J.-P., and Simões, A., editors, *Languages, Applications and Technologies*, pages 57–64, Cham. Springer International Publishing.