netherlands Science center







Erik Tjong Kim Sang, Kim Smeenk, Aysenur Bilgin, Tom Klaver, Laura Hollink, Jacco van Ossenbruggen, Frank Harbers and Marcel Broersma

CLIN29, Groningen, 31/01/2019

Task: automatically predict genres of Dutch newspaper articles

Data: 2,930 Dutch newspaper articles with 16 different genre labels

Examples of genre labels: news, column, editorial, interview



Previous work: Harbers and Lonij (2017) obtained 65% accuracy on this task

Our method: machine learning (MLP, NB, RF, SVM)

Result: 70% accuracy with SVM (interannotator agreement: 77%)



We want to use the distribution of genres over time (1955-1995) to study the effects of depillarization of Dutch newspapers

The quality of the proposed genre labels should be very good, in particular: their predicted distributions should be excellent



Question

Can you convince us that the genre prediction system works well enough to base our future studies on?

- 1. Open the genre classification system
- 2. Look for components that could introduce bias
- 3. Improve the **transparency** of the system with data visualizations

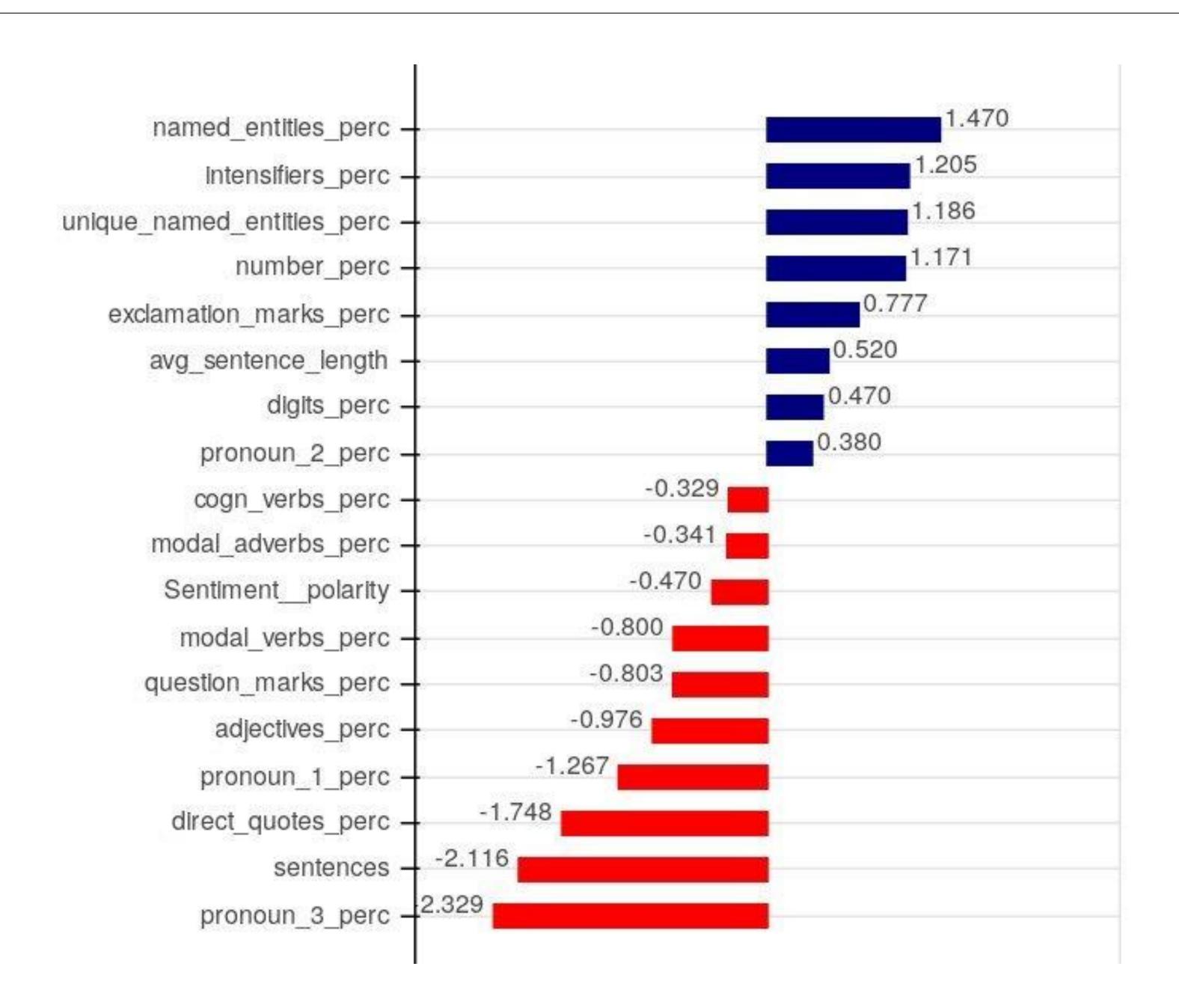
We have built a platform supporting step 3



VOOR AAN DE RADIO
t TWEEDE DIVISIE A i Portu"a__Psv Hilversum-EDO__
Enschede rviv RCH—Graafschap 5 Go ZFC—Zwolse
Boys f ADO—Telstii. Heerenveen—Wageningen ... i
DWS—Sitter__, Zwartemeer—AGOVV i VIVV—HerVclés
Vitesse—Spel. Cambuur 5 Sparta—Nac PEC-FC
Zaanstreek EERSTE- DIVISIE ' ' Haarlem~Tubantia i SSar.™ TWEEDE DIVISIE B 'f Willen, H-Ive'lov Fortuna VI.Xerxes •' VW—Blanw » Baronie—'t Gooi tSSB-S&»
gfcfZe.DvS ■:.:::. i 'SS:3S""'U"■'■>'* """■'■"■ &£-__e*i_■:::::::: !• Helmondia—Limburgia «t zijn opgenomen in de
sport-toto. De curfl.--.j_. ' '" drukte z'.I') reserve-wedstrijden.
j"""v «__. ' *A- - - -v"-'^"-"JV-_--__r_-^-».---I^v-"--__nj_-

Paper version

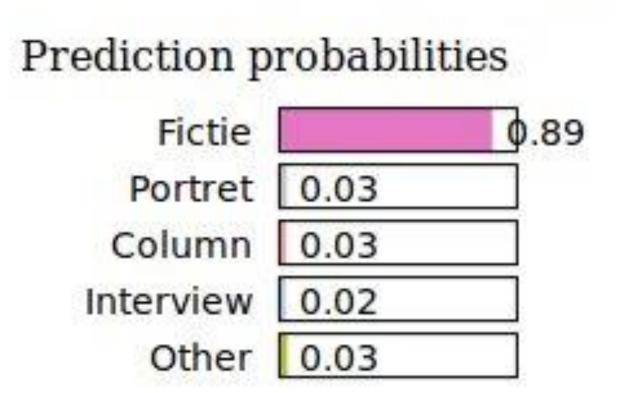
Digital version



Explanation for Article 79 using SVC BOW N3BGS

NOT Fictie

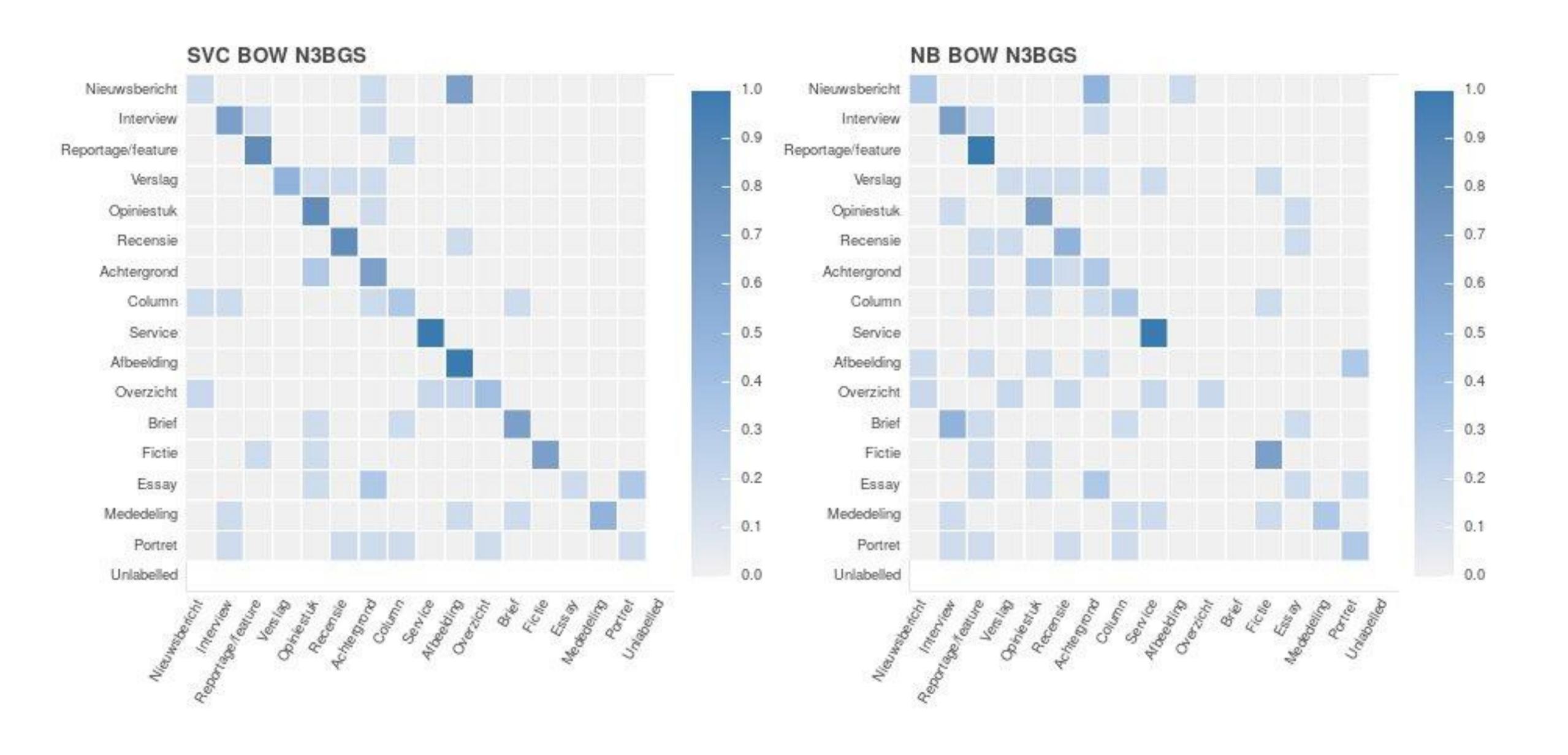
SVC BOW N3BGS predicts Fictie. Actual genre is Fictie.





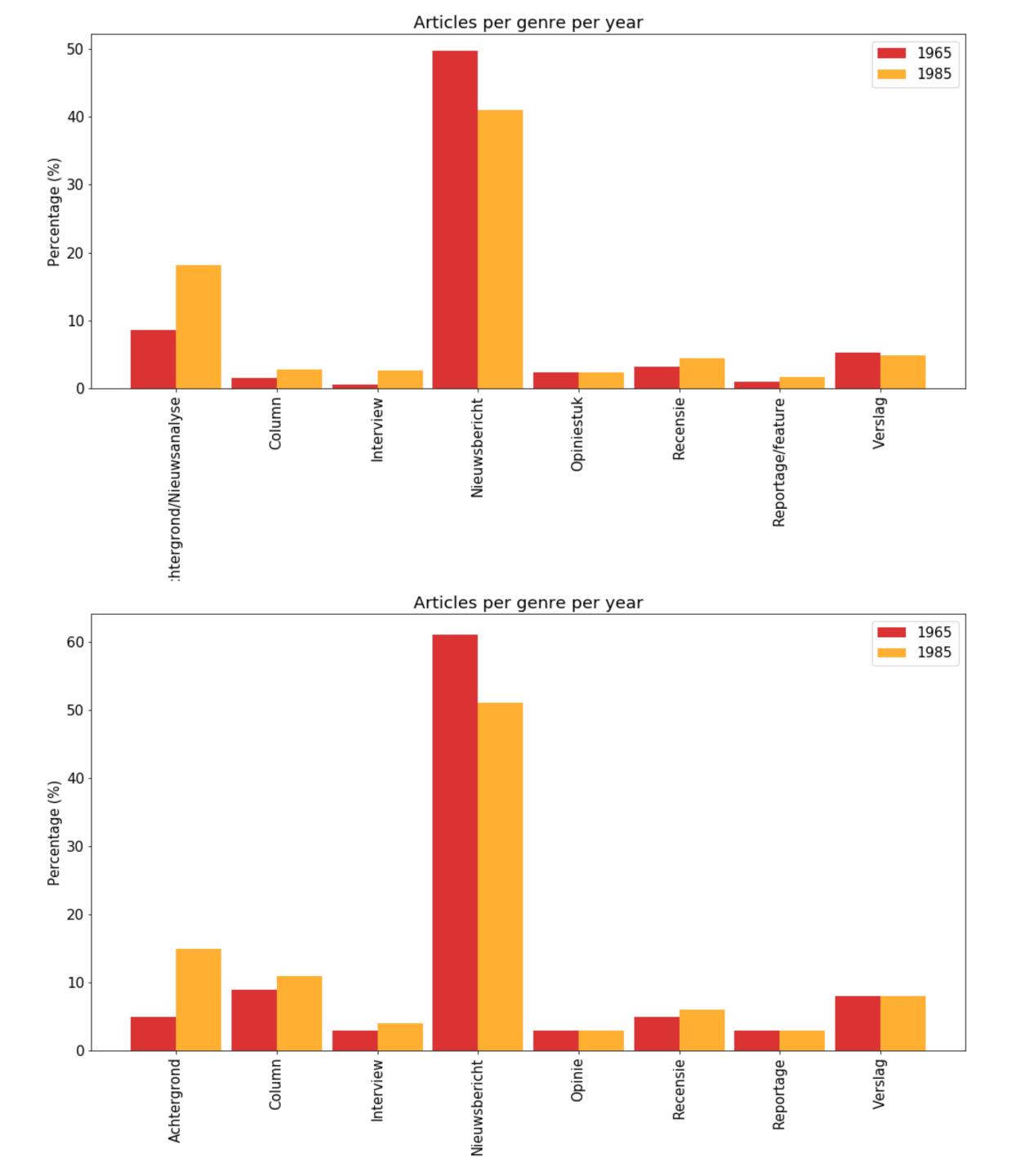
, * dan* wii ontdekt hebben . Gistere opmerking. Slr Tristram zei:. "He proberen je over te halen . Denk je v ernstl Het ls er ; jde t aangewezen p nuis door te dringen/tzou jij dan het twijfel, "antwoordde. Shield, / « n eens door te maken . " -* . ; | • « ' Ludovic . " Ik zal het paneel ; zëker sprak erover, d » * htf deze week e Thane kuchte . " En hoe — neem Jn 🤄 kunnen een raam openbreken , " 'ar , bang . dat je Si # Tristram nooit to beantwoordde haar blik door hsarï ... 'r ' (WORDT VERVOLGD)

% " heb precies bet soort miv/dat prima



Gold standard data

Machine labeled data



The domain scientists regard the current quality of the predicted genre labels as too low to be used as a basis for further study

This involves both the label accuracy and the provided explanations for the labels

1. Collect more training data to improve model accuracy

2. Employ word vectors to overcome lack of training data

3. Look for better features, to generate better explanations

4. Evaluate alternative more advanced machine learners

Concluding remark

Improving the transparency of our classifier has improved the insights in the classification task, both for domain scientists and computer scientists



