

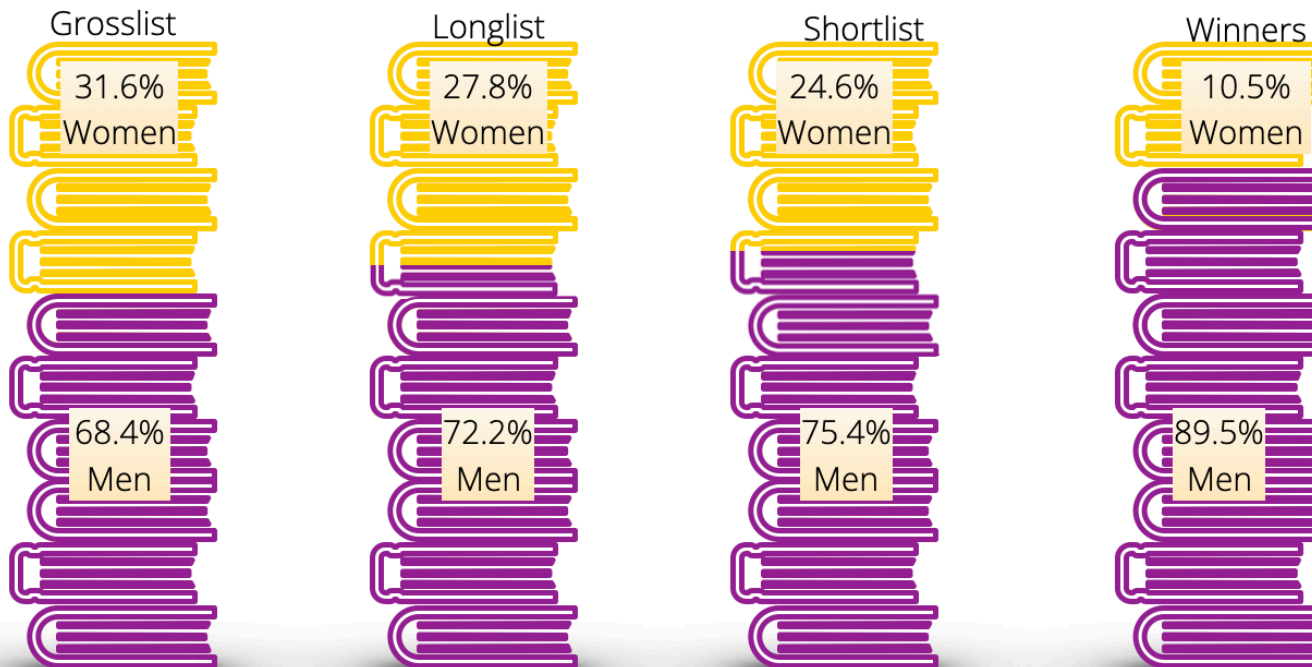
A Distant Reading of Gender Bias in Dutch Literary Prizes

Noa Visser MSc, dr. A. van Cranenburgh, dr. D. Nguyen

Introduction

Author gender inequality in Dutch literary prize nominations (Boekenbon Literatuurprijs and Libris Literatuur Prijs)

Inequality increases in the selection procedure



Percentages author gender in the selection procedure for the Libris Literatuur Prijs 1994 - 2013

Introduction

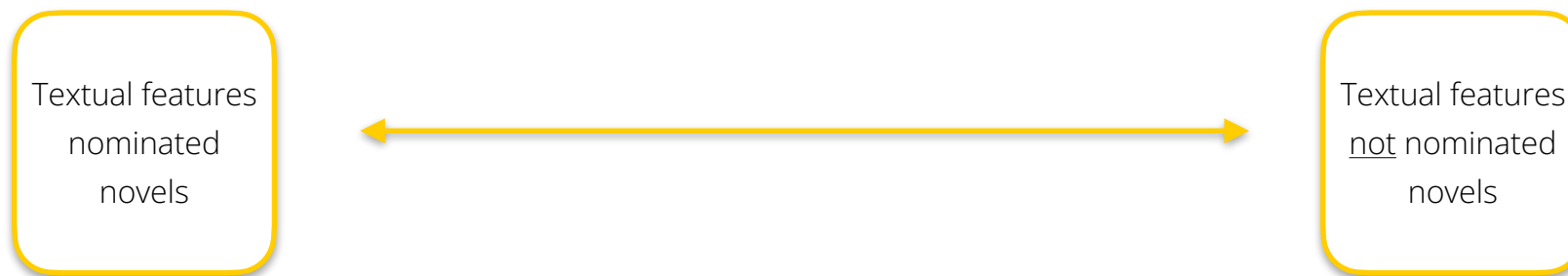
Dutch literary novels mainly written by (white) men

Perception of literary quality related to (white) men authors and publishers



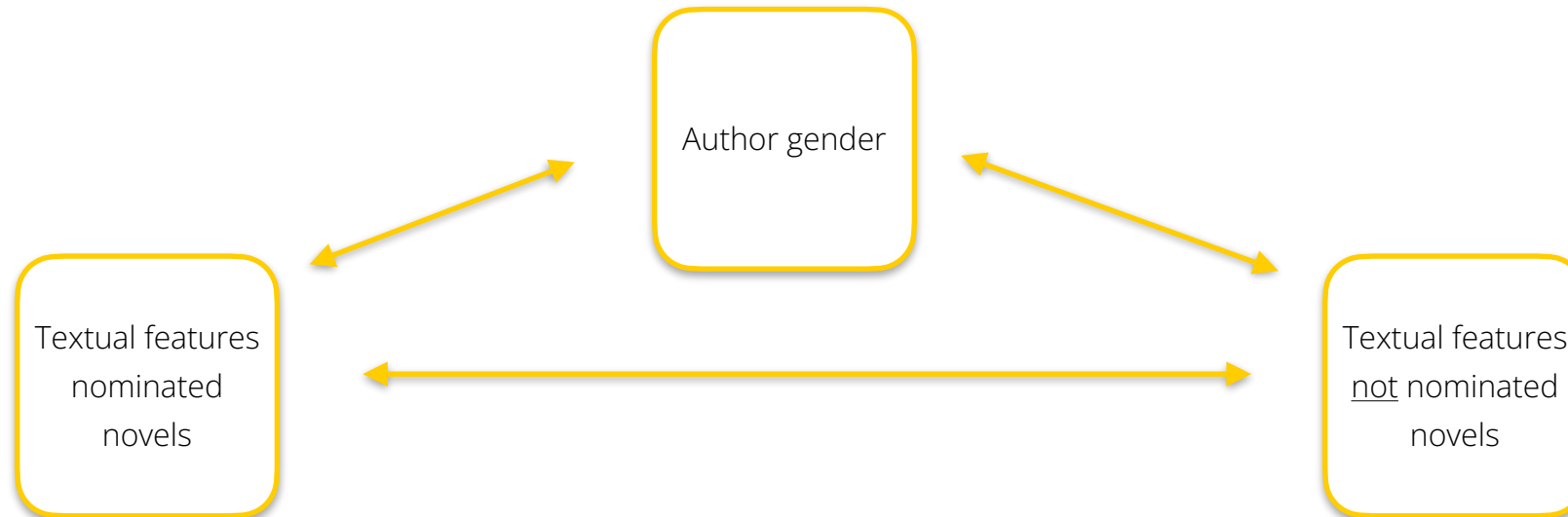
Research question

RQ1: Can nominated and not nominated novels be identified based on textual features alone?



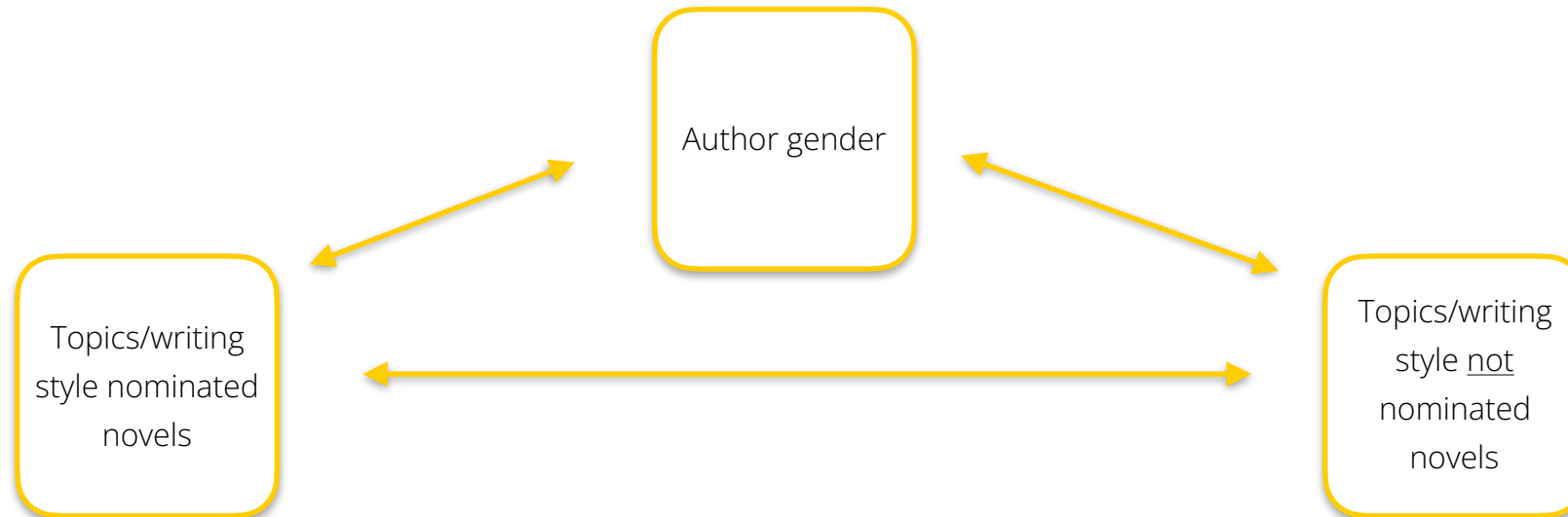
Research question

RQ2: Is there a relation between classifications on nominated and not nominated novels and author gender, where both classifications are based on textual features?



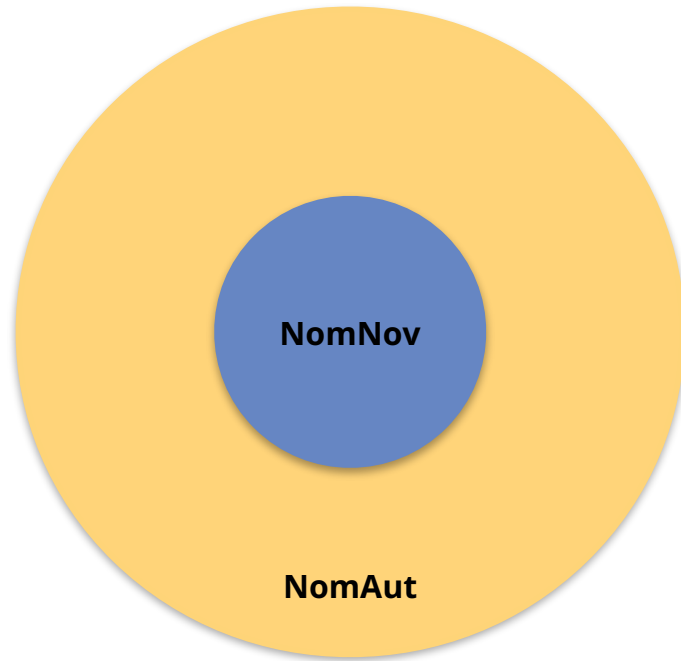
Research question

RQ3: Are the differences in topics/writing styles between books that are nominated for literary prizes and those that are not, related to author gender?

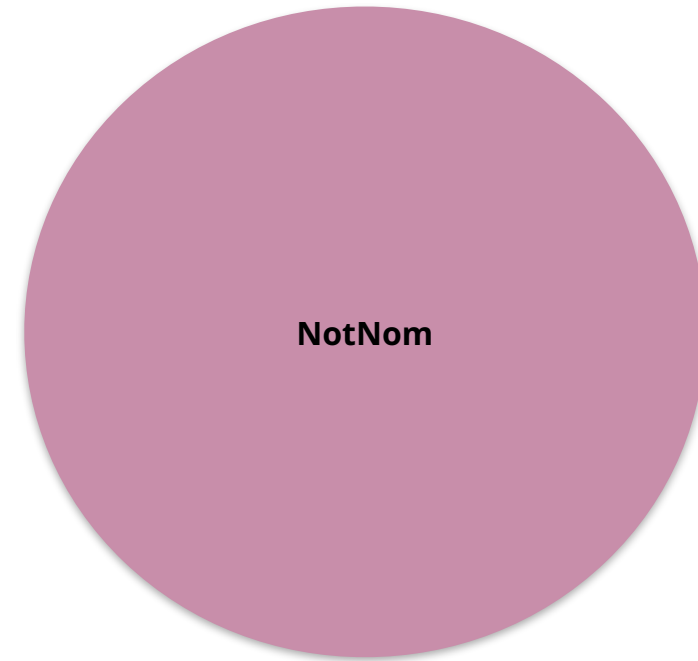


Dataset

Nominated authors



Not nominated authors

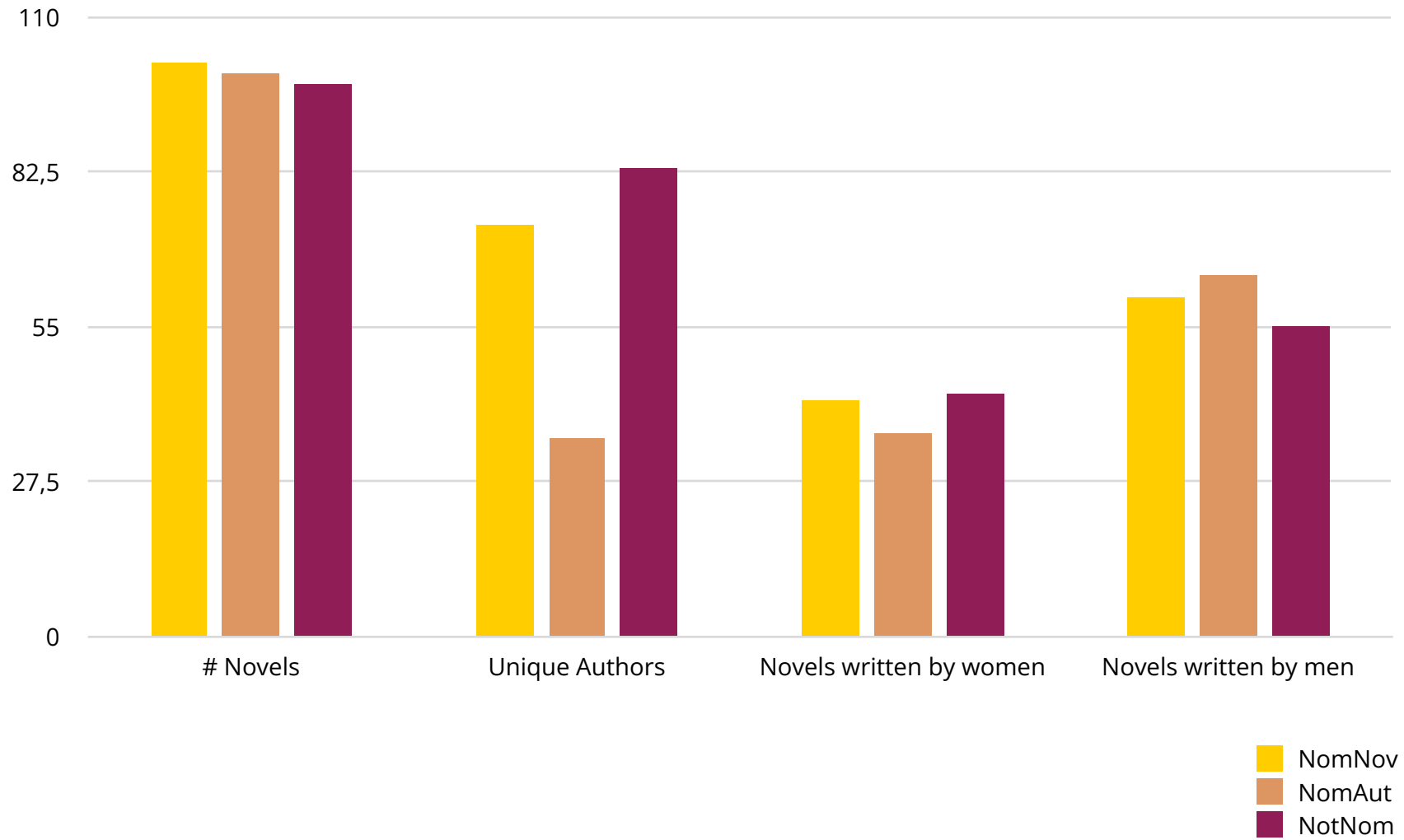


NomNov: Nominated novels

NomAut: Not nominated novels by nominated authors

NotNom: Not nominated novels by not nominated authors

Dataset





Method

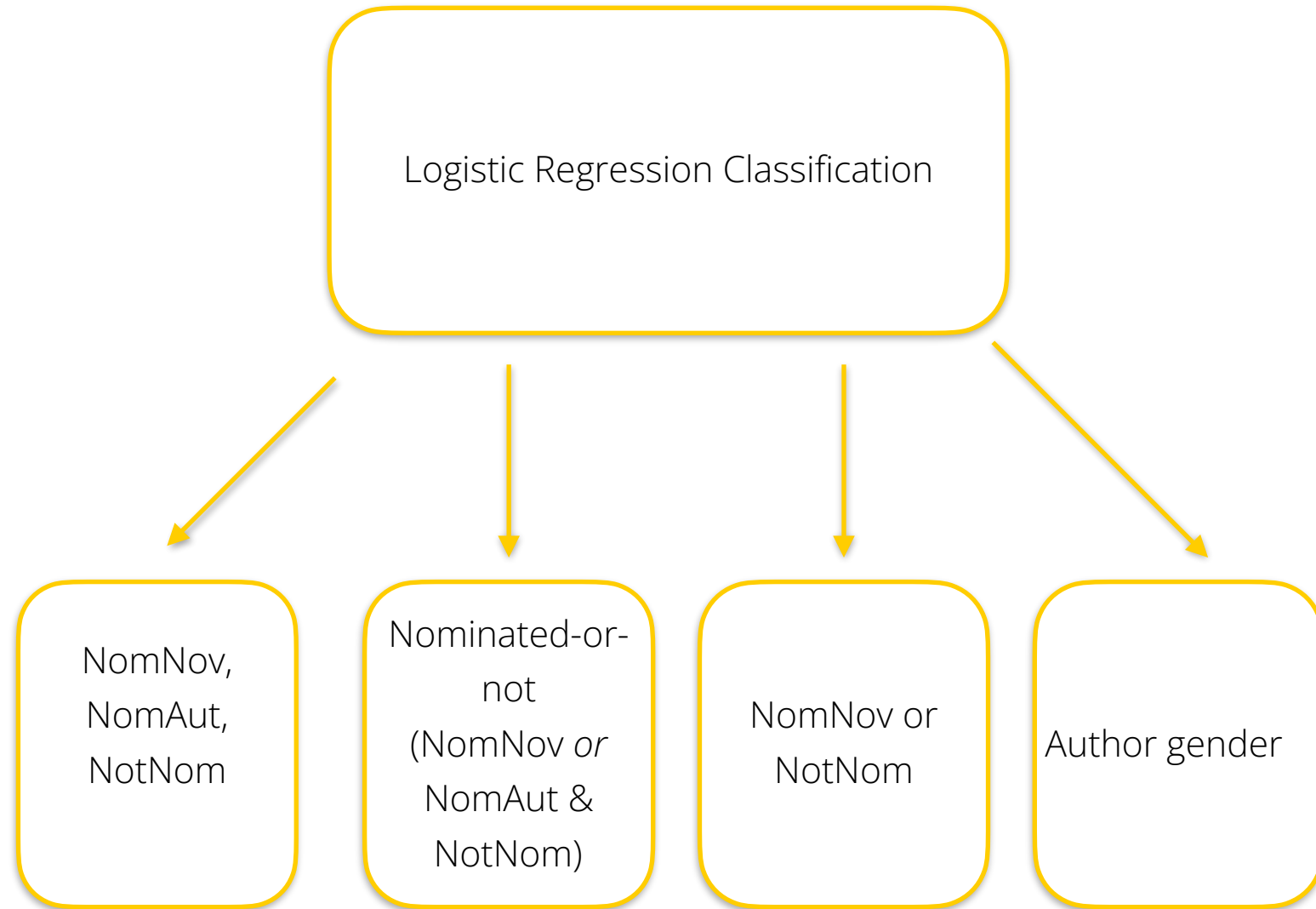
- Logistic Regression Classification
 - Tf-Idf vectoriser
 - Unigrams and bigrams
 - 5000 most frequent words
 - 5-fold cross-validation
- LDA Topic modelling
- Cosine delta

Method: Classification

Four different types of classification

All performed on complete dataset and balanced author gender subset

Precision, recall, F1-score and overall accuracy



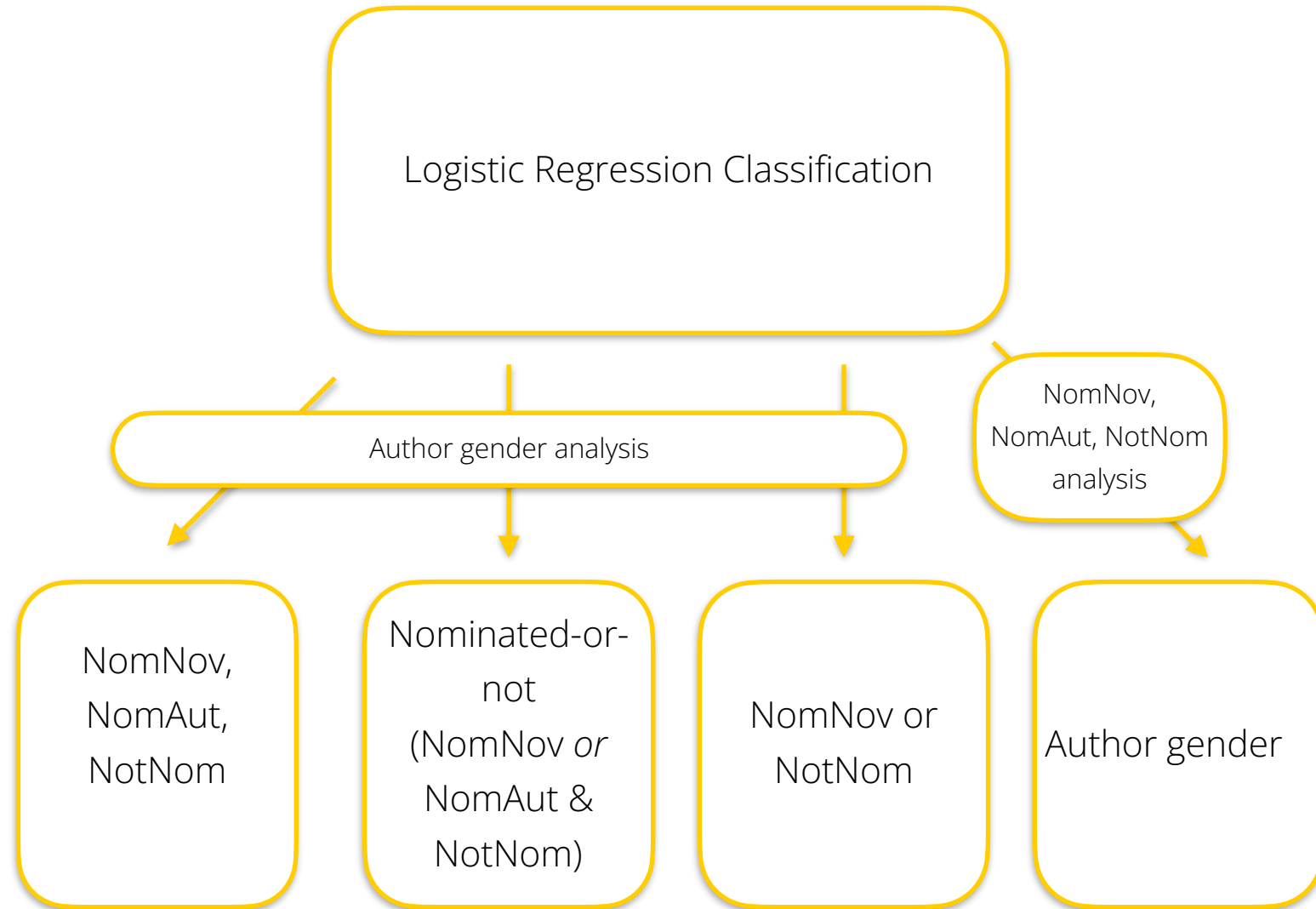
NomNov: nominated novels, NomAut: not nominated novels by nominated authors, NotNom: not nominated novels by not nominated authors

Method: Classification

Four different types of classification

All performed on complete dataset and balanced author gender subset

Precision, recall, F1-score and overall accuracy



NomNov: nominated novels, NomAut: not nominated novels by nominated authors, NotNom: not nominated novels by not nominated authors

Results: NomNov, NomAut, NotNom

COMPLETE CORPUS	Precision	Recall	F1-score	Standard deviation	Number of novels
NomNov	0.569	0.700	0.628	0.0134	100
NomAut	<u>0.567</u>	<u>0.333</u>	<u>0.420</u>	0.0285	102
NotNom	0.615	0.735	0.735	0.0284	98
Accuracy			0.587	0.0155	300

Overall accuracy better than
chance (0.306)

NomNov: nominated novels, NomAut: not nominated novels by nominated authors,
NotNom: not nominated novels by not nominated authors

Results: NomNov, NomAut, NotNom

NomAut worst performance,
regardless of author gender

Novels written by women lower
classification performance in
comparison to novels written by
men

Not nominated novels written by
women higher F1 scores than
nominated novels written by
women, for all three models

COMPLETE CORPUS				
Women	Precision	Recall	F1-score	Number of novels
NomNov	<u>0.500</u>	0.583	0.538	36
NomAut	0.517	<u>0.357</u>	<u>0.423</u>	42
NotNom	0.680	0.791	0.731	43
Accuracy			0.579	121
Men	Precision	Recall	F1-score	Number of novels
NomNov	0.605	0.766	0.676	64
NomAut	0.613	<u>0.317</u>	<u>0.418</u>	60
NotNom	<u>0.567</u>	0.691	0.623	55
Accuracy			0.592	179

NomNov: nominated novels, NomAut: not nominated novels by nominated authors,
NotNom: not nominated novels by not nominated authors

Results: NomNov, NomAut, NotNom

NomAut worst performance, regardless of author gender

Novels written by women lower classification performance in comparison to novels written by men

Not nominated novels written by women higher F1 scores than nominated novels written by women, for all three models

COMPLETE CORPUS				
Women	Precision	Recall	F1-score	Number of novels
NomNov	<u>0.500</u>	0.583	0.538	36
NomAut	0.517	<u>0.357</u>	<u>0.423</u>	42
NotNom	0.680	0.791	0.731	43
Accuracy			<u>0.579</u>	121
Men	Precision	Recall	F1-score	Number of novels
NomNov	0.605	0.766	0.676	64
NomAut	0.613	<u>0.317</u>	<u>0.418</u>	60
NotNom	<u>0.567</u>	0.691	0.623	55
Accuracy			<u>0.592</u>	179

NomNov: nominated novels, NomAut: not nominated novels by nominated authors, NotNom: not nominated novels by not nominated authors

Results: NomNov, NomAut, NotNom

NomAut worst performance,
regardless of author gender

Novels written by women lower
classification performance in
comparison to novels written by
men

Not nominated novels written by
women higher F1 scores than
nominated novels written by
women, for all three models

COMPLETE CORPUS				
Women	Precision	Recall	F1-score	Number of novels
NomNov	<u>0.500</u>	0.583	0.538	36
NomAut	0.517	<u>0.357</u>	<u>0.423</u>	42
NotNom	0.680	0.791	<u>0.731</u>	43
Accuracy			0.579	121
Men	Precision	Recall	F1-score	Number of novels
NomNov	0.605	0.766	0.676	64
NomAut	0.613	<u>0.317</u>	<u>0.418</u>	60
NotNom	<u>0.567</u>	0.691	0.623	55
Accuracy			0.592	179

NomNov: nominated novels, NomAut: not nominated novels by nominated authors,
NotNom: not nominated novels by not nominated authors

Results: Author Gender

COMPLETE CORPUS	Precision	Recall	F1-score	# novels
MAN	75.9	82.7	79.1	179
WOMAN	70.5	61.2	65.5	121
Accuracy			74.0	300

NOMNOV	Precision	Recall	F1-score	# novels
Man	77.1	84.4	80.6	64
Woman	<u>66.7</u>	<u>55.6</u>	<u>60.6</u>	36
Accuracy			74.0	100

NOMAUT	Precision	Recall	F1-score	# novels
Man	74.6	83.3	78.7	60
Woman	<u>71.4</u>	<u>59.5</u>	<u>64.9</u>	42
Accuracy			73.5	102

NOTNOM	Precision	Recall	F1-score	# novels
Man	75.9	80.0	77.9	55
Woman	<u>72.5</u>	<u>67.4</u>	<u>69.9</u>	43
Accuracy			74.5	98

Author gender prediction
surpasses chance (0.609)

NomNov: nominated novels, NomAut: not nominated novels by nominated authors,
NotNom: not nominated novels by not nominated authors

Results: Author Gender

COMPLETE CORPUS	Precision	Recall	F1-score	# novels
MAN	75.9	82.7	79.1	179
WOMAN	70.5	61.2	65.5	121
Accuracy			74.0	300

NOMNOV	Precision	Recall	F1-score	# novels
Man	77.1	84.4	80.6	64
Woman	66.7	55.6	60.6	36
Accuracy			74.0	100

NOMAUT	Precision	Recall	F1-score	# novels
Man	74.6	83.3	78.7	60
Woman	71.4	59.5	64.9	42
Accuracy			73.5	102

NOTNOM	Precision	Recall	F1-score	# novels
Man	75.9	80.0	77.9	55
Woman	72.5	67.4	69.9	43
Accuracy			74.5	98

Novels written by women lowest classification scores, on all classes

NomNov (nominated novels) written by women lowest classification score overall

NomNov: nominated novels, NomAut: not nominated novels by nominated authors, NotNom: not nominated novels by not nominated authors



Method: Exploration

- LDA Topic modelling
 - 50 topics
 - NomNov, NomAut, NotNom
 - Author gender
- Cosine delta
 - Exploration most frequent words
 - Correctly classified novels compared to misclassified novels

Results: LDA Topic Modelling

Men (NomNov and NotNom)

Topic 0: War

0.0174 Major (*majoor*)

0.0145 Soldier (*soldaat*)

0.0142 War (*oorlog*)

0.0141 Man (*man*)

0.0134 General officer (*generaal*)



NotNom (Men and Women)

Topic 23: Second World War

0.0223 German (*Duits*)

0.012 Prince (*prins*)

0.0103 Germany (*Duitsland*)

0.0091 War (*oorlog*)

0.0082 Jewish (*joods*)



NomNov Men, NomAut & NotNom Women

Topic 30: Health care

0.0312 Doctor (*dokter*)

0.0195 Patient (*patiënt*)

0.0172 Hospital (*ziekenhuis*)

0.0146 Doctor (*arts*)

0.0105 To say (*zeggen*)



NomNov: nominated novels, NomAut: not nominated novels by nominated authors,

NotNom: not nominated novels by not nominated authors

Results: Cosine Delta

Closely related writing style



Positive relation with novels
written by nominated authors
(NomNov and NomAut)





Conclusion

- Nominated and not nominated novels distinguishable
- Word use nominated novels further from women writers
- Author gender inequality rooted in homogenous writing style



Universiteit
Utrecht

Thank you for your attention
Feel free to ask any questions

Literature

Bamman, David, Jacob Eisenstein, and Tyler Schnoebelen (2014). "Gender identity and lexical variation in social media". In: *Journal of Sociolinguistics* 18.2, pp. 135–160. ISSN: 14679841. DOI: [10.1111/josl.12080](https://doi.org/10.1111/josl.12080). arXiv: [1210.4567](https://arxiv.org/abs/1210.4567).

Berkers, Pauwke (2009). Classification into the literary mainstream? Ethnic boundaries in the literary fields of the United States, the Netherlands and Germany, 1955-2005.

Blei, David M, Andrew Y Ng, and Michael I Jordan (2003). "Latent dirichlet allocation". In: the Journal of machine Learning research 3, pp. 993–1022.

Burrows, John (2002). "'Delta': a measure of stylistic difference and a guide to likely authorship". In: *Literary and linguistic computing* 17.3, pp. 267–287.

van Cranenburgh, Andreas and Rens Bod (2017). "A Data-Oriented Model of Literary Language". In: arXiv preprint arXiv:1701.03329.

Dera, Jeroen (2021). "De helaasheid der leeslijsten. Over diversiteit in het literatuuronderwijs". In: *De Lage Landen* 64 (1), pp. 115–121.

Dijkgraaf, Margot and René Appel (2013). *Vrouwen, mannen en de Libris Literatuur Prijs*. Stichting Libris Literatuur Prijs.

Herrmann, J Berenike, Arthur M Jacobs, and Andrew Piper (2021). "Computational Stylistics". In: *Handbook of Empirical Literary Studies*, p. 451.

Koolen, Corina et al. (2020). "Literary quality in the eye of the Dutch reader: The National Reader Survey". In: *Poetics* 79. February, p.101439. ISSN: 0304422X. DOI: [10.1016/j.poetic.2020.101439](https://doi.org/10.1016/j.poetic.2020.101439). URL: <https://doi.org/10.1016/j.poetic.2020.101439>.

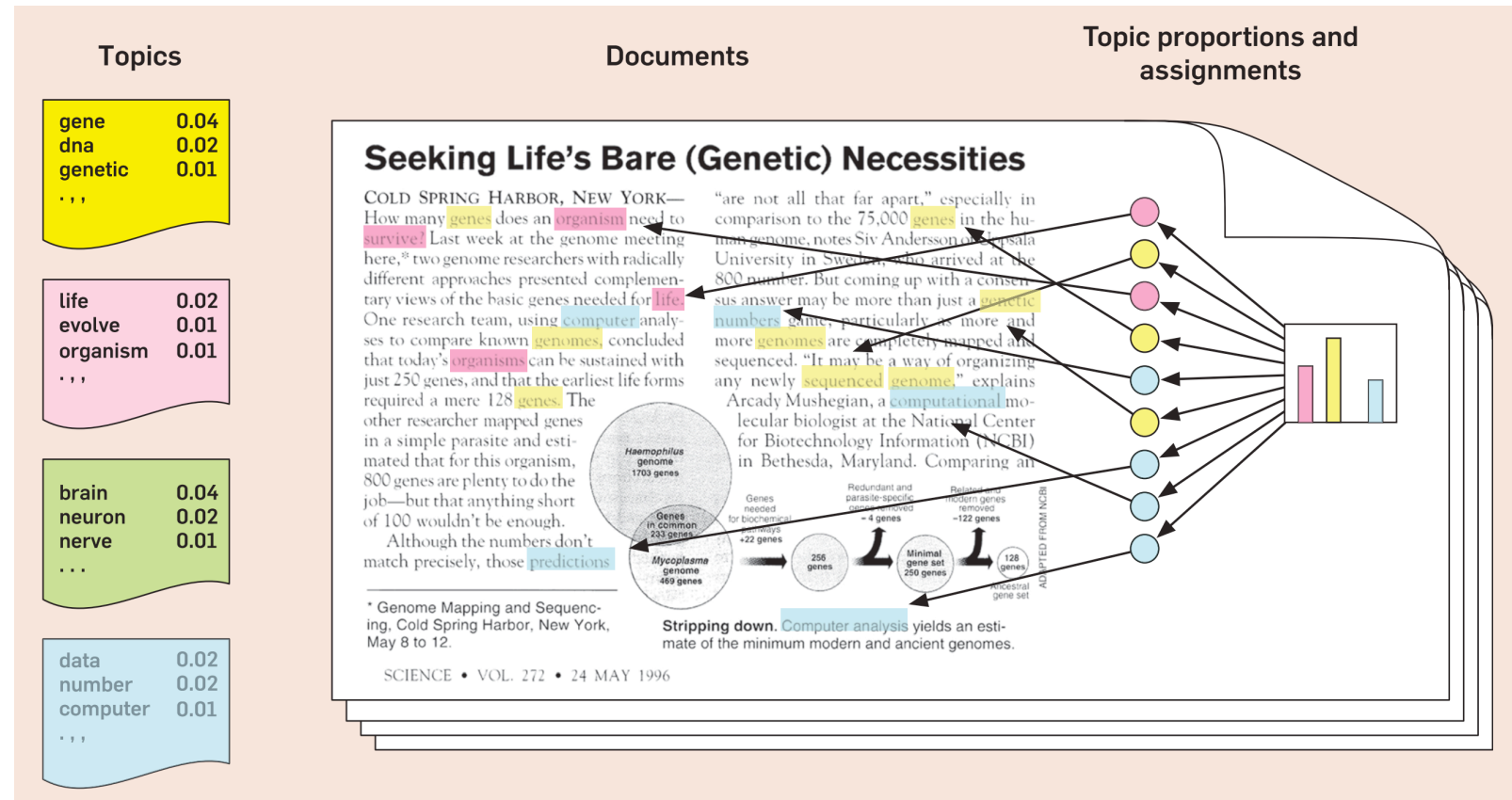
Koren, Timo and Christine Delhaye (2019). "Depoliticising literature, politicising diversity: ethno-racial boundaries in Dutch literary professionals' aesthetic repertoires". In: *Identities* 26.2, pp. 184–202. ISSN: 15473384. DOI: [10.1080/1070289X.2017.1391561](https://doi.org/10.1080/1070289X.2017.1391561). URL: <https://doi.org/10.1080/1070289X.2017.1391561>.

Background: LDA Topic modelling

LDA topic modelling

Unsupervised model to determine topics that occur in documents iteratively

Topics can be related to multiple documents



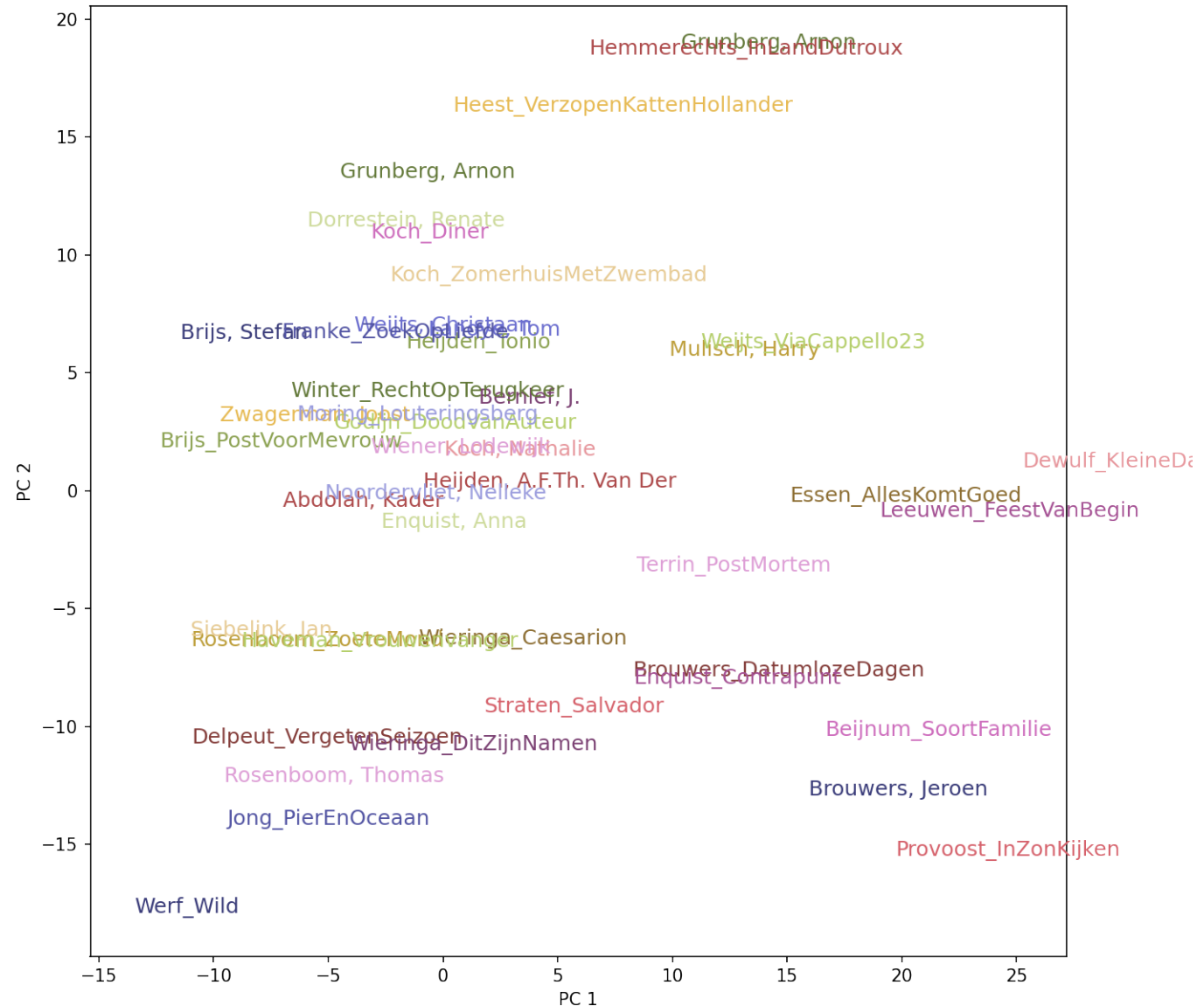
<https://pyro.ai/examples/proldlda.html>

Background: Cosine Delta

Identify authorship and writing style

Exploration 100-5000 most occurring words

Distance between words is calculated



Dataset

Collection of popular epub

DBNL

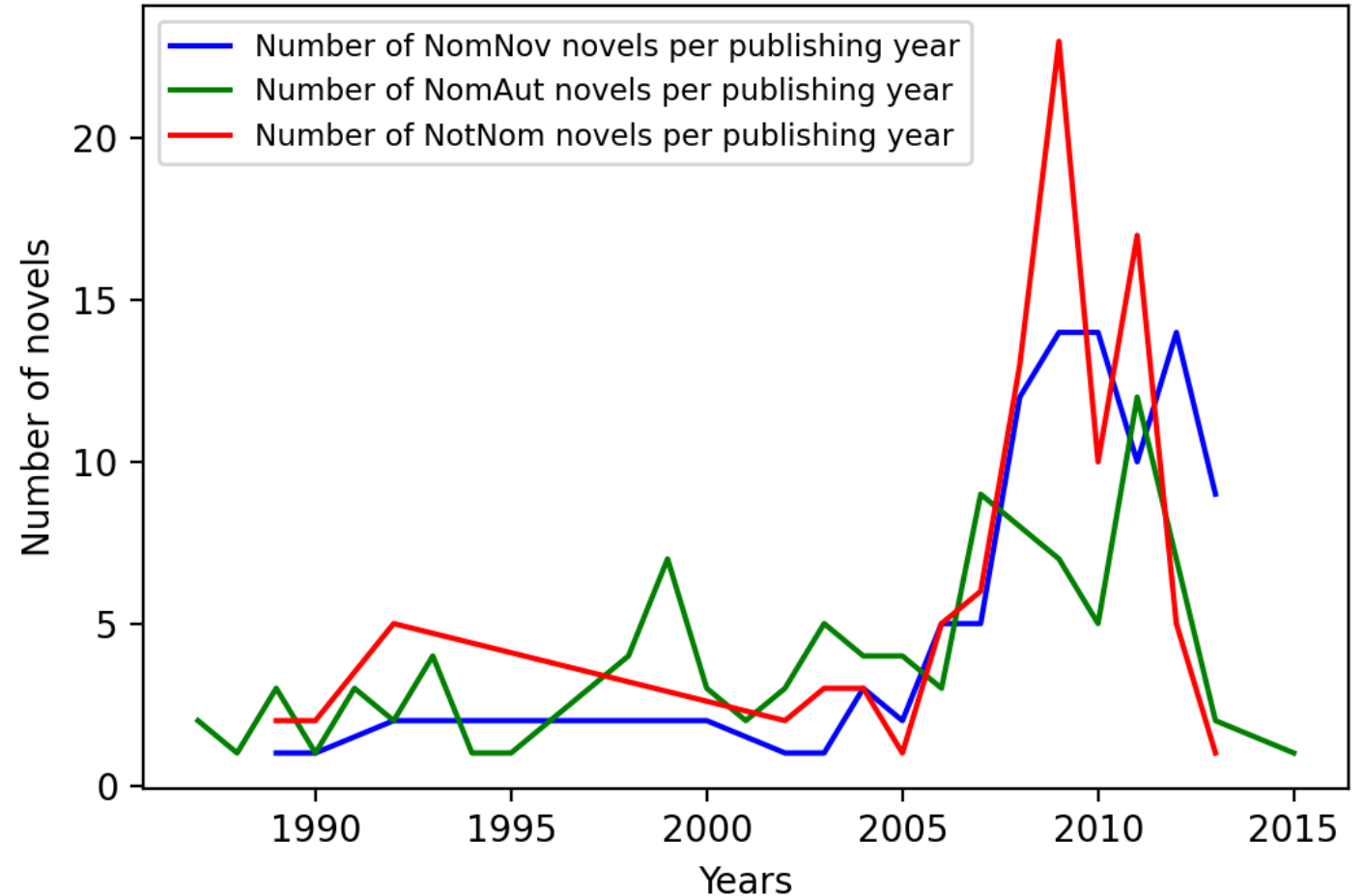
Riddle of Literary Quality

Dutch nominated novels

Estimation publishing year

Estimation author gender

Publishing years novels per category

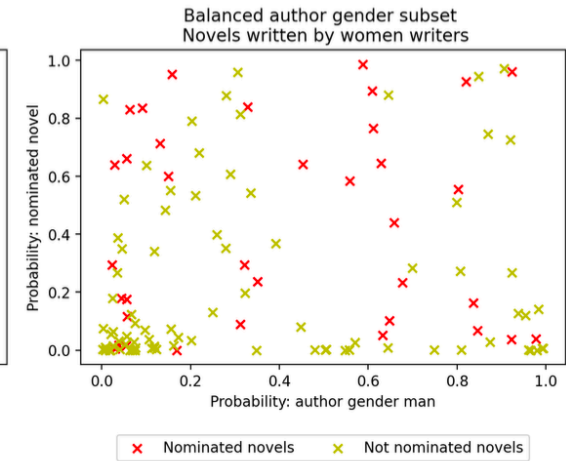
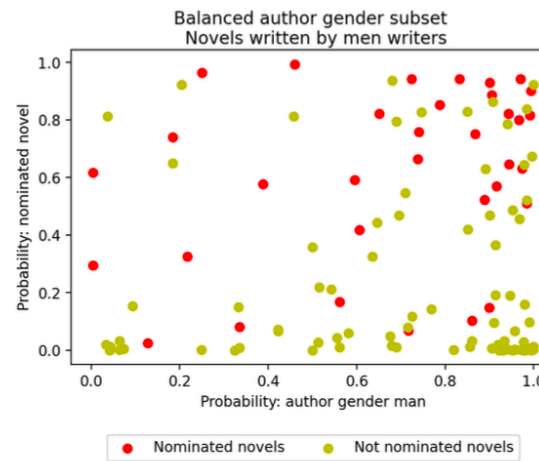
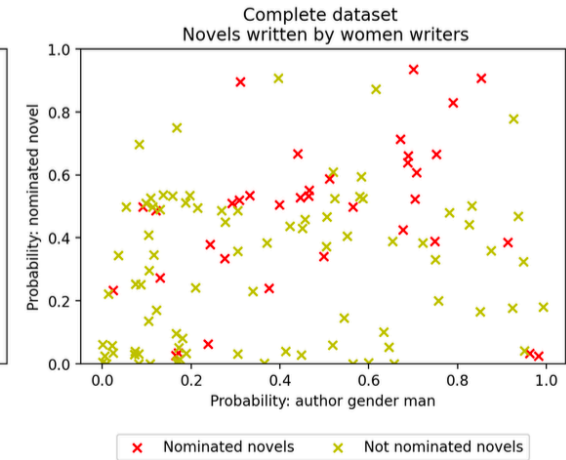
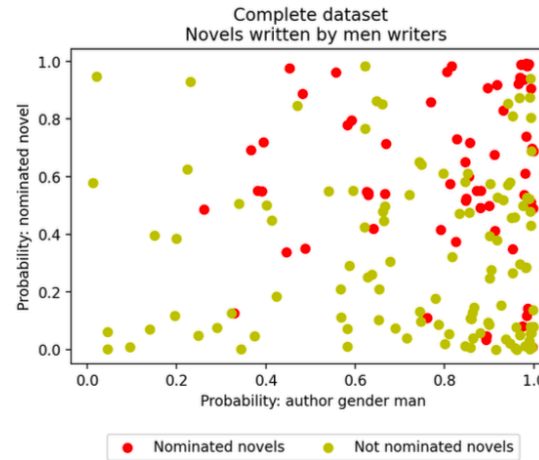


Results: Confidence Comparison

Nominated-or-not model and author gender prediction

Comparison confidence of classification

Relation between high probability to be nominated for a literary prize, and novels written by men



Results: highest weight features logistic regression

NomNov

Office (*bureau*)
I saw (*ik keek*)
Was not (*niet was*)
Swimming (*zwemmen*)

NomAut

Above (*boven*)
That still (*die nog*)
He went (*ging hij*)
To slide (*glijden*)
Her the (*haar de*)
He saw (*hij zag*)
Kilo (*kilo*)
Also be (*ook zijn*)
Party (*partij*)
When it (*toen het*)

NotNom

To happen (*gebeuren*)
No sense (*geen zin*)
Prison (*gevangenis*)
To slide (*glijden*)
Yes I (*ja ik*)
Can (*kan*)
Also from (*ook van*)
Stage (*podium*)
Affairs (*zaken*)
Sit (*zit*)