



Historical legislation:
segmentation, datafication, and
automatic metadating

Dr. Annemieke Romein

Datum: 17 Nov. 2020

Annemieke.Romein@Huygens.knaw.nl



Contents



intro

LA:
segmentation

Datafication

Automatic
metadating

Conclusion



Introduction

- Early modernist (political/ legal/ institutional) focus on political terminology (training in Rotterdam)
- Turned DH about 3 years ago... (although I did statistics)
 - Old-fashioned way of doing paleography... but now *the* go-to person on Transkribus (according to colleagues).
- Postdoc at UGent (2017-2020: Law and Order?);
- Postdoc at KNAW Huygens ING (2020-2023: A Game of Thrones?);
- Researcher-in-Residence at the KB National Library of the Netherlands (2019) to automatically metadata texts.

intro

LA:
segmentation

Datafication

Automatic
metadataing

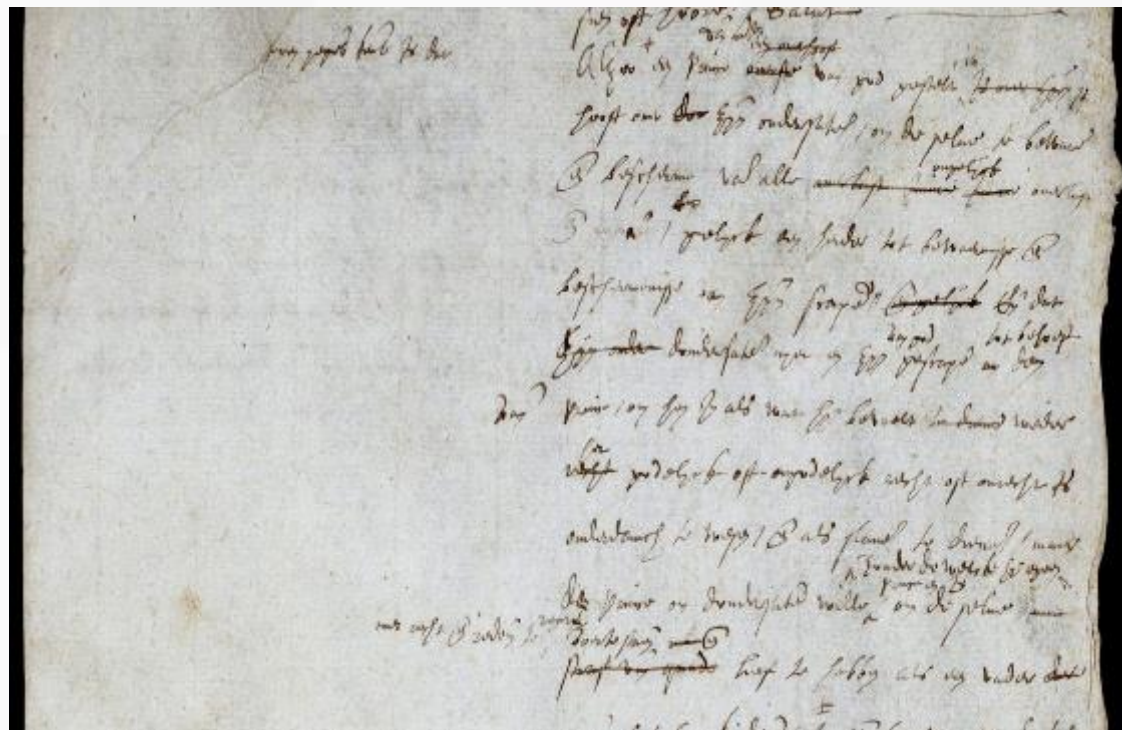
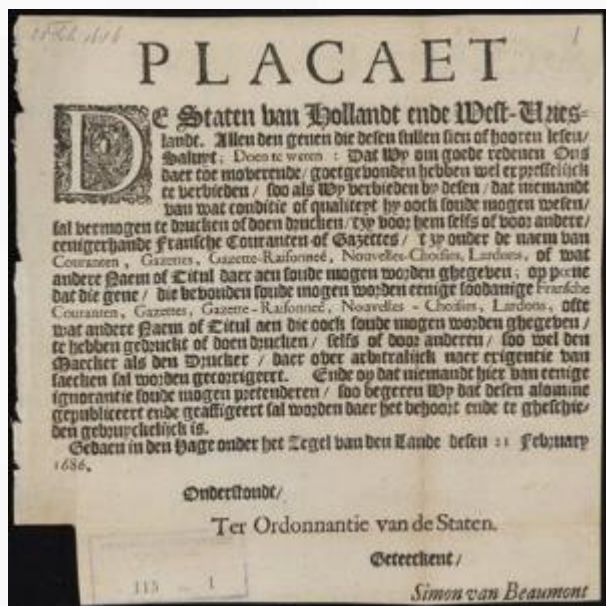
Conclusion

Project/ program: early modern ordinances



Paalhuis,
Amsterdam
(ca. 1660)

Appearances



Left: <https://library.maastrichtuniversity.nl/collections/databases/digitale-plakkaatzoeke/>

Right: Plakkaat van Verlatinghe, pagina 1 (NA)

intro

LA:
segmentation

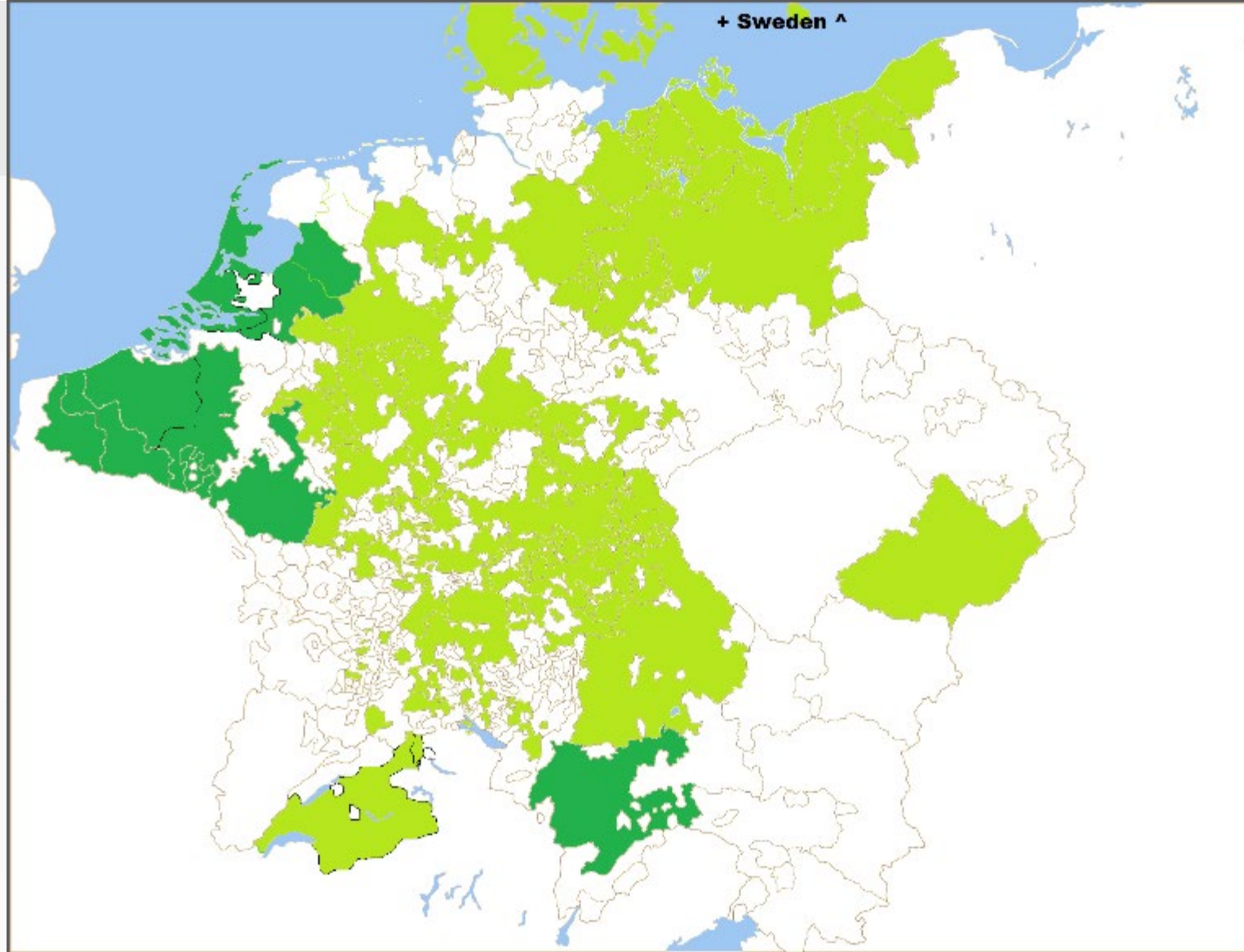
Datafication

Automatic
metadating

Conclusion

Police-ordinances include much information





MPIeR Repertorium and other initiatives (dark green), period 1500-1800.



The 'smoking gun' and the RQ

'Early Modern European states struggled for survival, making it impossible to 'reinvent the wheel' each time a problem arose. Hence, it was of tremendous importance to copy, adapt and implement normative rules (aka legislation) that were already proven succesful elsewhere.'

But...

- Where to start?
- How to prove this?
- Data?

intro

LA:
segmentation

Datafication

Automatic
metadating

Conclusion



intro

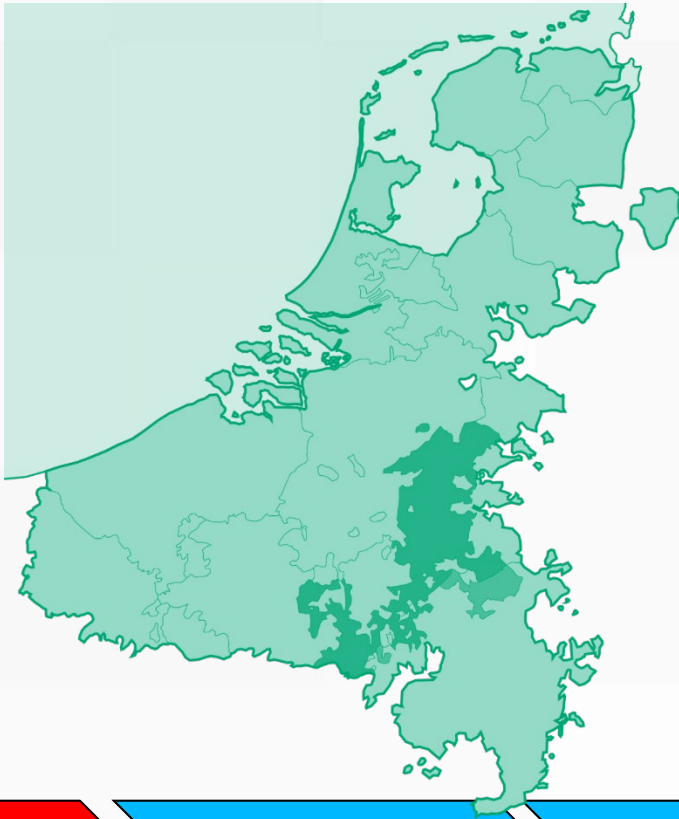
LA:
segmentation

Datafication

Automatic
metadating

Conclusion

Narrowing down the focus



intro

LA:
segmentation

Datafication

Automatic
metadating

Conclusion

From going narrow to broadening up... DH

‘Entangled Histories’ as a DH-pilot study... to:

- Work on automatic segmentation
- Working on the readability of early modern print (Dutch Gothic/ Dutch Romantype/ French Romantype)
- Automatic metadating early modern laws

A horizontal navigation bar consisting of five chevron-shaped segments pointing to the right. The first segment is red and contains the text 'intro'. The remaining four segments are blue and contain the text 'LA: segmentation', 'Datafication', 'Automatic metadating', and 'Conclusion' respectively.

intro

LA:
segmentation

Datafication

Automatic
metadating

Conclusion



Team 'Entangled Histories'

Core team (6 months)

... Annemieke Romein (PI)

... Sara Veldhoen (Scientific Programmer)

... Michel de Gruijter (Project Manager)

And help of many others!



Part 1. Segmentation of texts

intro

LA:
segmentation

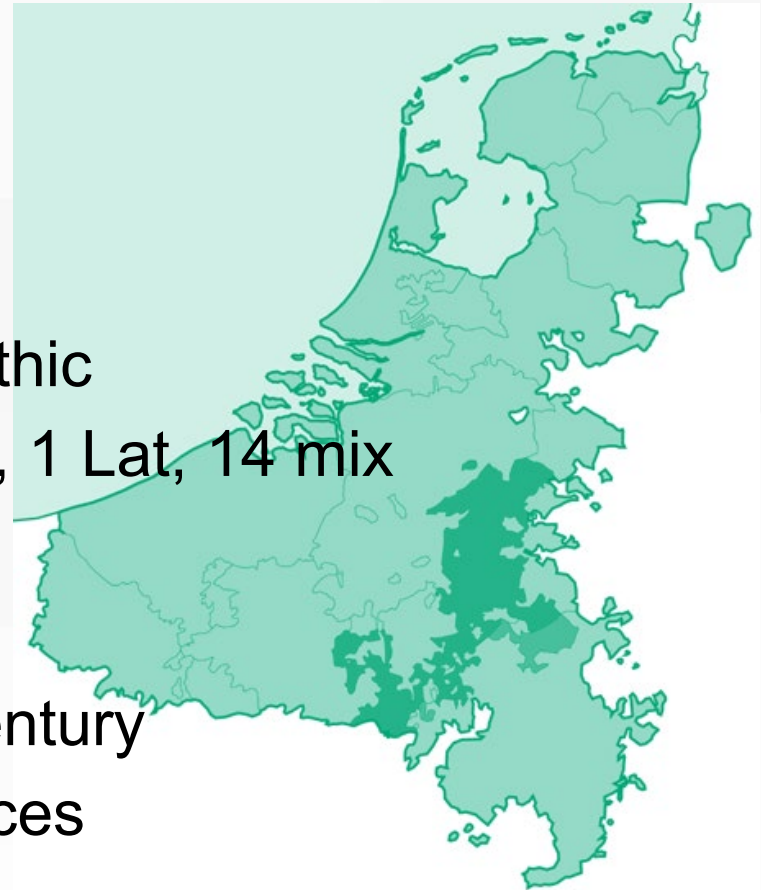
Datafication

Automatic
metadating

Conclusion

Digitisation: Sources in numbers

Sources #	108
Font	88 roman, 20 gothic
Language	67 Dutch, 26 FR, 1 Lat, 14 mix
Pages#	75.000
Characters#	550 mln
Publication date	17th and 18th century
Region	Nearly all provinces



intro

LA:
segmentation

Datafication

Automatic
metadating

Conclusion



intro

LA:
segmentation

Datafication

Automatic
metadating

Conclusion



Teaching the computer to read and to segment texts

Transkribus

page-number 35 header Groot Gelders paragraph

Placaet aangaende de Facht,
en dat op Veluwe nyemant sal mogen huys-
honden houden, ten waere die twee voorste
clawwen tottet leste lit toe syn affgeknypt.

By den Keyser.

WY laeten weten, dat nyemant geen huys-
 honden off bastaert-wynden op de Ve-
 luwe halde, ten sy dat den selven twee voire-
 ste principae clawwen wes tot den lesten lit
 toe van ene hoer voireste voeten affgeknypt sijn,
 brynnen veertien daezhen nae publicatie van de

te bysulder van wegen Key. Maj^r, als Hertoge
 van Gelre ende Grave van Zutphen gelast ende
 bevolen, Und eerst den tollenaeren voorten een
 bouck alleene to halden ende daer inne to schrij-
 ven ende aentcykenen alle goeden die sy vertol-
 len fullen, in wyns schip die goeden sijn, wiese
 toekommen, in wat Steden off plaetsen die schip-
 pers en cooplyden off eygenaers van den goeden
 woenen ende burgerschap houden.

Die affcommende goeden by sich selven, en-
 de die opvaerende goeden oick by sich selven,
 in wat maenden en op wat daghen van elcker
 maent die vertolt sullen worden.

Van gelijken sullen die tolschryvers haer con-
 trebouck alleen halden, en daer inne op 't aenge-

page-number 36

E Staten
 van Stadt Gronin-
 gen ende Ommelan-
 den: DOEN TE
 WETEN. Alhoewel
 wy verhoopt hadden
 dat dooz het redres en
 correctie van de Lijste
 van verpachtinge by
 ons in het laest verlo-
 pen Jaer 1622. gedaen
 ende alomme gepubli-
 ceert / eenighsins gheremedieert ende belet sulden zijn
 geweest / de veelvuldige frauden ende dieverpen / Item
 listige ende schadelijcke practiquen / composities en an-
 dere abuyten / streckende tot groten aff-breuck van de
 Generale Middelen / en merckelijcke prejuditie van alle
 goede getrouwe Ingesetenen ende oprechte Liefhebbers
 van 't Vaderlant / dewelcke deur alfulcke publicque Die-
 ven ende Wyanden van 't Gemeene beste / niet alleene in
 haere neeringe becoztet worden / maer oock de schade
 ende 't achterhepde van 't Gemeene by den selven veroor-
 sacckt



Segmentation (the thing you are supposed to do first...)

- Methods: Abbyy FineReader v.11 (✘/✓)
- P2PaLA (α-version) (✘/ ✓)
- NLE Document Understanding (AI-based) (✓)
- Rule-based approach (✓!)

intro

LA:
segmentation

Datafication

Automatic
metadating

Conclusion

Part 2. Datafication/ Text Recognition



intro

LA:
segmentation

Datafication

Automatic
metadating

Conclusion

Improving the OCR- results throu

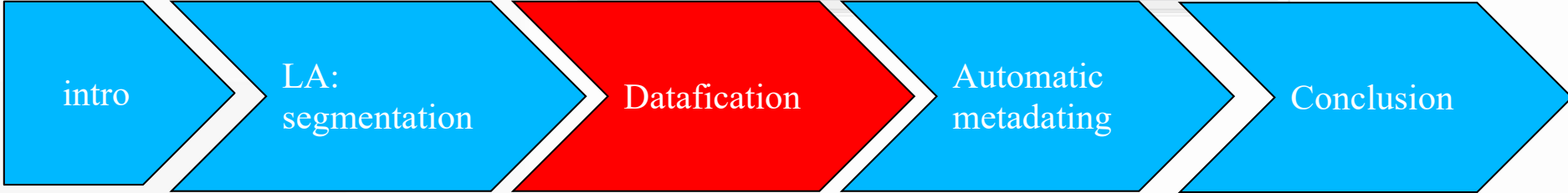


Nr. of Words:	Nr. of Lines:	
2134	309	
Show Train Set	Show Test Set	Show Character Set

Learning Curve

Epoch	CER Train (%)	CER Test (%)
0	~70	~95
5	0.02	1.69
75	0.02	1.69

CER on Train Set:	0.02%	CER on Test Set:	1.69%
-------------------	-------	------------------	-------



Transkribus... what's that?



<https://www.youtube.com/watch?v=gS3GvveIN4>



Transcribus v1.3.0.7-SNAPSHOT (31_11_2019_13:12): Loaded doc: KBNLB170078181.png, ID: 169342, Page 1, file: 0000001.png [Image Meta Info: (Resolution: 1.0, w/h: 1610 * 2164)] [current line: w/h: 1316 * 78]

Server Overview Layout Metadata Tools

Logout info@camelin.nl

Document... Find

Document Manager User Manager


Versions Jobs

Recent Documents... User activity

Collections: Plakkaten (40652, Owner) Col-ID

Documents: HTR Model Data

ID	Title	Pages	Uploader	Uploaded
174...	GENT900000157721.png	1038	sara.veltho...	Mon Jun 17...
174...	GENT900000122888.png	532	sara.veltho...	Mon Jun 17...
169...	KBNLB030060241.png	560	sara.veltho...	Tue Jun 04...
169...	KBNLB0300194313.png	50	sara.veltho...	Tue Jun 04...
169...	KBNLB0300168073.png	504	sara.veltho...	Tue Jun 04...
169...	KBNLB0300145425.png	230	sara.veltho...	Tue Jun 04...
169...	KBNLB0300113811.png	758	sara.veltho...	Tue Jun 04...
169...	KBNLB170078181.png	34	sara.veltho...	Tue Jun 0...
169...	KBNLB013024231.png	324	sara.veltho...	Tue Jun 04...
169...	USA000017480.png	246	sara.veltho...	Tue Jun 04...
169...	USA000005624.png	336	sara.veltho...	Tue Jun 04...
169...	USA000118944.png	328	sara.veltho...	Tue Jun 04...
169...	USA000166296.png	1000	sara.veltho...	Tue Jun 04...
169...	UBL000022736.png	300	sara.veltho...	Tue Jun 04...
169...	UBL000034544.png	328	sara.veltho...	Tue Jun 04...
169...	UBL000045373.png	626	sara.veltho...	Tue Jun 04...
169...	UBL000048244.png	805	sara.veltho...	Tue Jun 04...
168...	UBL000035604.png	228	sara.veltho...	Mon Jun 03...
168...	UBL000035548.png	498	sara.veltho...	Mon Jun 03...
168...	UBL000031728.png	214	sara.veltho...	Mon Jun 03...
168...	USA000168294.png	942	sara.veltho...	Mon Jun 03...
168...	USA000136173.png	178	sara.veltho...	Mon Jun 03...



1-1 O.nantien en Statuten

2-1 die de Keyserlycke Maiesteyt in zijnder ieghe-

2-2 woordicheyt op den .vij. dach Octobris Int iaer

3-1 M.CCCCC.xxi.heeft doen lesen en verclaren den Sta-

3-2 ten van sinen landen van herwerts over- en de welke al-

4-1 daer wtgeroepen en gepubliceert zijn geweest opden xv

intro

LA:
segmentation

Datafication

Automatic
metadating

Conclusion



Automatic Text Recognition

- Basically “Handwriting Text Recognition” but then applied to printed texts.
- It performed better than we had expected.

TABLE 1. RESULTS PER CREATED MODEL (CER).

Model name [ID]	Training (CER)	Test (CER)	# Words (training)	# Lines (training)
Dutch_Gothic_Print [Model ID18944]	0.22%	1.71%	51143	7143
French_18thC_printed [Model ID19166]	0.33%	0.65%	38487	3883
Romantype_Dutch_Print [Model ID19423]	1.26%	1.17%	88105	13013



Bonus-output

- We created three HTR+-models within Transkribus that have been made public:
 - Dutch Gothic Print
 - French 18thC Print
 - Romantype Dutch Print(since October 2020, these are also available in PyLaia)

intro

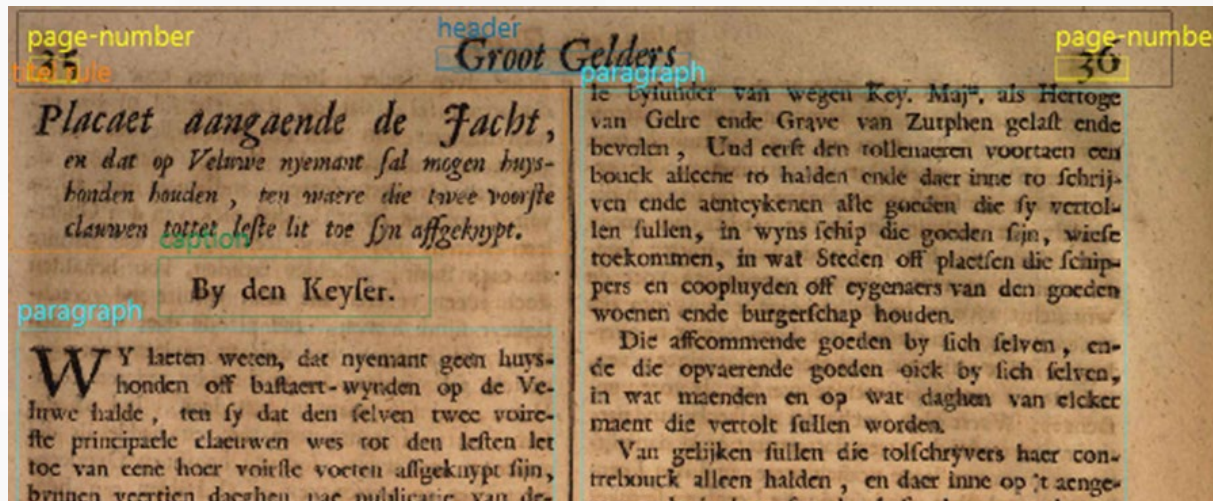
LA:
segmentation

Datafication

Automatic
metadating

Conclusion

Part 3. Automatic metadating



intro

LA:
segmentation

Datafication

Automatic
metadating

Conclusion

Method for automatic metadating

- Use of a controlled subject vocabulary
- This was developed within the projects 'Repertorium der Polizeyordnungen' and 'Gute Polizey und Polizeywissenschaft'.
 - Not perfect for all situations, but additions are possible

intro

LA:
segmentation

Datafication

Automatic
metadating

Conclusion

Method for automatic metadating

- 4 level deep hierarchical categorisation, but we needed to add a 5th level (police ordinance or not)
- So, at level 1: yes police ordinance or not
- Level 2: 5 options
- Level 3: 25 options
- Level 4: 163
- Level 5: 1584 options (currently)

intro

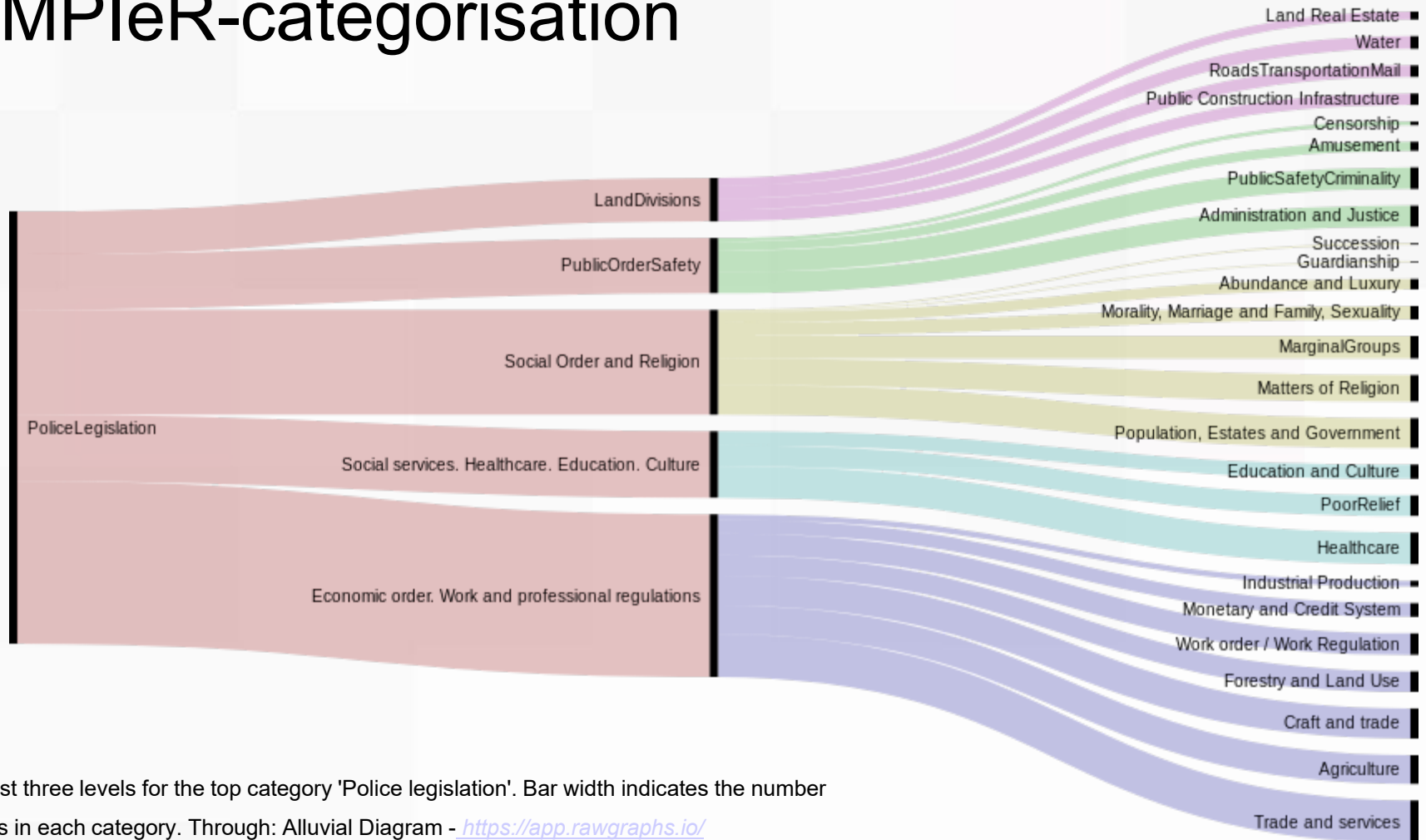
LA:
segmentation

Datafication

Automatic
metadating

Conclusion

MPIeR-categorisation



The first three levels for the top category 'Police legislation'. Bar width indicates the number of texts in each category. Through: Alluvial Diagram - <https://app.rawgraphs.io/>

SKOS: <https://doi.org/10.5281/zenodo.3564586>



Categorisation

Manually labelled ± 3000 laws.

Used in the casestudy: 470.

Hierarchical categorisation – 4(+1) layers (MPLeR).

Converted to SKOS format (Short Knowledge Organisation System)

<https://doi.org/10.5281/zenodo.3564586>

intro

LA:
segmentation

Datafication

Automatic
metadating

Conclusion



annif

intro

LA:
segmentation

Datafication

Automatic
metadating

Conclusion



Training Annif

SET-UP

- Proof of concept: 470 labeled 'documents' (laws)
- <10 subjects each, 3.3 on average
- 69% and 28% of annotations at level 5 (deepest) and 4
- Training per category level
- 10-fold Cross Validation
- Backend: TF-IDF
- Limit = 4, threshold = 0.4

intro

LA:
segmentation

Datafication

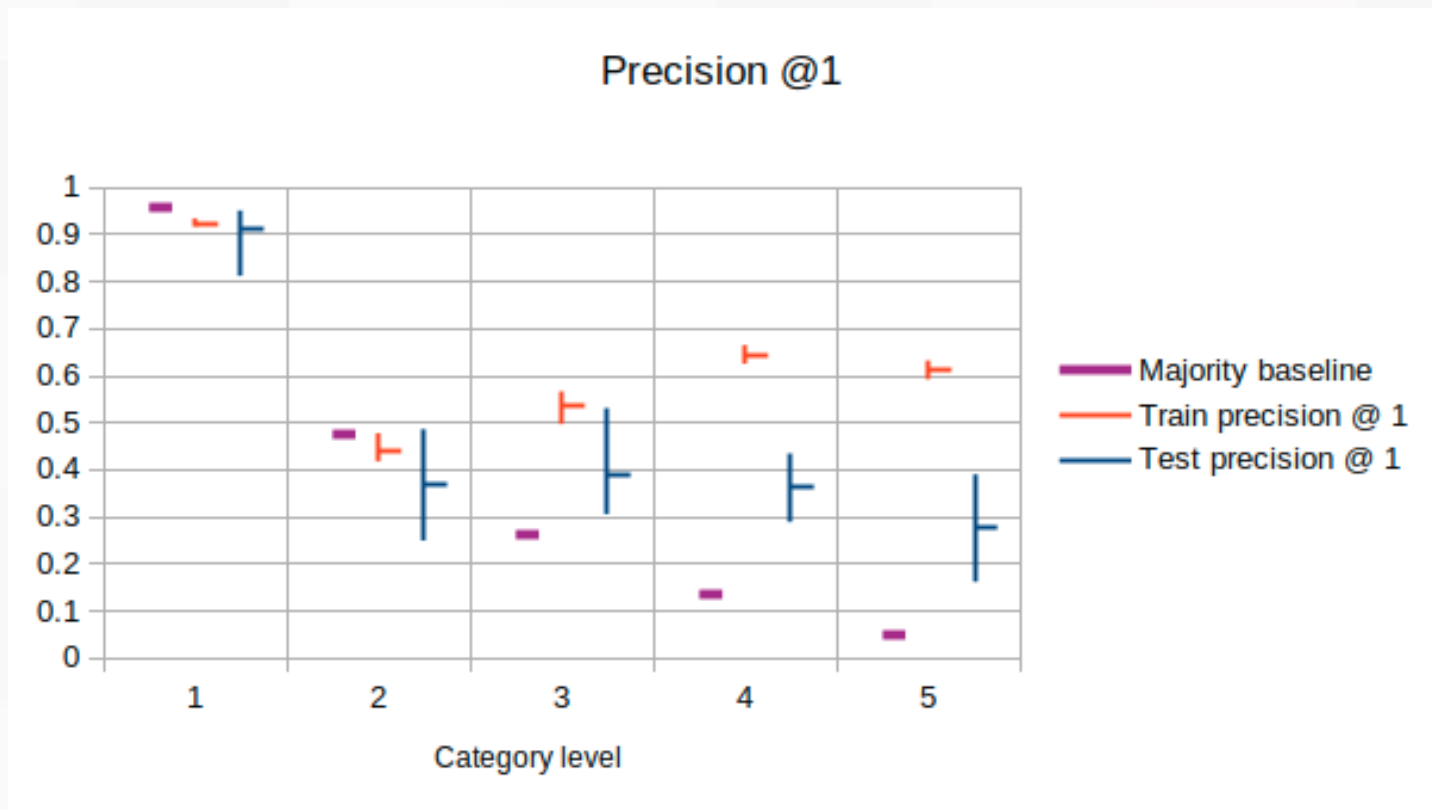
Automatic
metadating

Conclusion



Training Annif

RESULTS



intro

LA:
segmentation

Datafication

Automatic
metadating

Conclusion

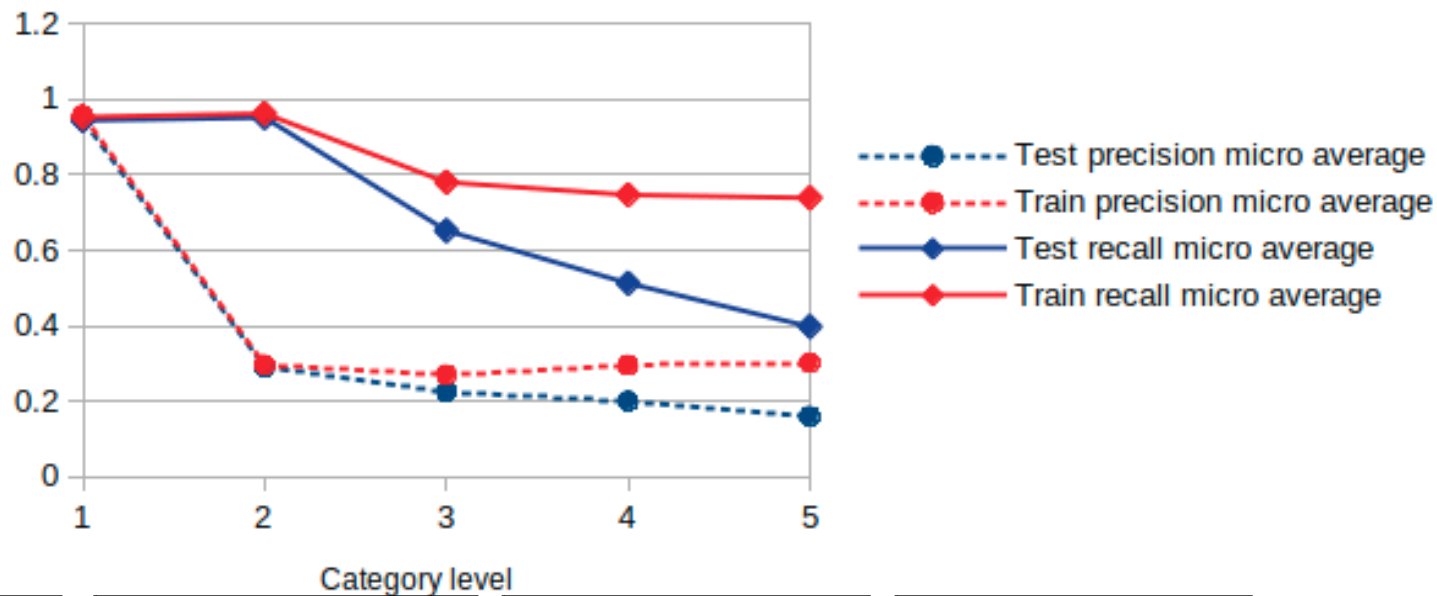


Training Annif

RESULTS

Micro averaged recall & precision

Limit=4, Threshold=0.4



intro

LA:
segmentation

Datafication

Automatic
metadating

Conclusion



Conclusions

intro

LA:
segmentation

Datafication

Automatic
metadating

Conclusion



Conclusions

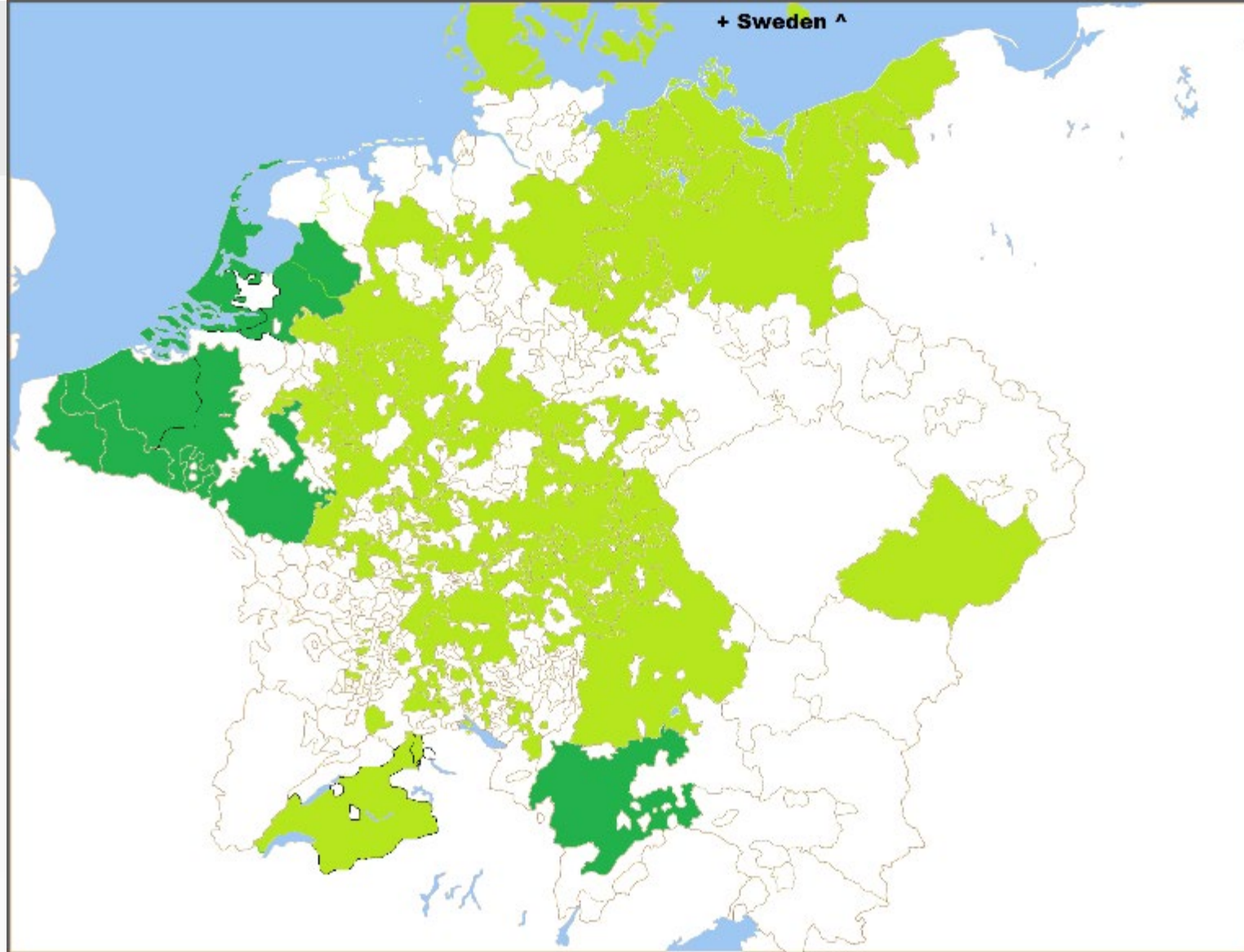
- Yes, the computer can automatically metadata the early modern legislation and it can be trained to become really good (but manual training will be necessary).
- No, we couldn't automatically segment the texts just yet.
- No, we couldn't prove the hypothesis to the fullest yet – application on all of the books was just too time consuming.



Training Annif

DISCUSSION

- Limited amount of data (470texts)
- Only 'dumb' TF-IDF – curious about other backends
- Impact of OCR quality?
- Stemmer: Snowball Dutch
Impact of historical language & spelling variation?
- Making use of the hierarchy in categories



MPIeR Repertorium and other initiatives (dark green), period 1500-1800.



... future/ further research...

- In my current project I will try to test the use of Annif on German texts (Berne).
- I hope to apply for funding for future projects that will start connecting 'the dots' through Linked Data.
- If someone is interested in testing out Annif on early modern Swedish ordinances... let me know!



Further Reading

POINTERS

- DH Benelux journal publication:
<http://journal.dhbenelux.org/journal/issues/002/article-23-romein/article-23-romein.html>
- Info on KB lab:
Dataset: <https://lab.kb.nl/dataset/entangled-histories-ordinances-low-countries>
Blogs: <https://lab.kb.nl/about-us/blog/entangled-histories-bumps-road-and-bursts-success>
- Other work with Annif at KB:
<https://zenodo.org/record/3899723#.X2B2N1bRZJd>



Further Reading

- Annemieke continues to work on early-modern legislation
(<https://en.huygens.knaw.nl/projecten/game-of-thrones/>)
- At the KB, they're experimenting with Annif in a tool to aid the cataloguing department:
'Demosaurus'.

See <https://zenodo.org/record/3899723#.X2B6Q4bRZJc>



Thank you for your attention.