



# Provenance Language Dynamics in the Dutch Golden Age

Marjo van Koppen

DANS, Den Haag, 2018-09-03

Meertens  
Instituut





# Overview

- Input data
- Editing data
- Availability of the data
- Provenance for researchers



# Input Data

- Dbnl-web version of ‘The letters of P.C. Hooft’.

[https://dbnl.org/tekst/hoof001hwva02\\_01/colofon.php](https://dbnl.org/tekst/hoof001hwva02_01/colofon.php)

- FoLiA-data files from Nederlab: no version number, just a date (19-09-2017).

- Separated actual letters from notes, foreword etcetera.



# Editing data

- Tokenization and Tagging with *Adelheid* (versie 1.0).
- *Adelheid* was bound to give a lot of mistakes.
- We used it as a tool to preprocess the data for manual correction (maar veel fouten en onvolledig).



# Editing data

(Semi) Manual  
enrichment: pos-  
tagging and socio-  
linguistic  
information

By 10  
(R)MAstudents

<i>Document characteristics</i>	
Category	business, personal
Type	regular, appendix
Goal	express thanks, compliment, excuse, ask a favour, ask information, ask advice, admonish, inform, remember, persuade, order, allow, invite
Topic	business, literature, domestic affairs, love, death, news, religion/ethics

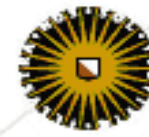
  

<i>Correspondent characteristics</i>	
Group	name
Individual	name, birth/death date, gender, occupation, literary author, relation to P.C. Hooft

<i>Letter segmentation</i>	
Introductory greeting, opening (optional), narrative, closing (optional), final greeting	

Table 1: Sociolinguistic annotation set.



# Editing data

Adelheid controle tagging

http://localhost:3035/adelheid\_check.html

557 / 1219

wonder vruchtbare akker des verstands . zo ghij wakker verhoet datter de ikker gheen onkruid meer in en zaaije (die

Gebruik pijljestoetsen

*datterde*    *dat | ter*

	<i>huidig</i>	<i>controle</i>	
<i>lemma</i>	dat+er	<input type="text" value="dat+er"/>	<input type="checkbox"/> modern alternatief
<i>pos</i>	VNW	<input type="radio"/> N <input type="radio"/> ADJ <input type="radio"/> WW <input type="radio"/> BW <input checked="" type="radio"/> VNW <input type="radio"/> LID <input type="radio"/> TW <input type="radio"/> VZ <input type="radio"/> VG <input type="radio"/> SPEC <input type="radio"/> LET	<input type="checkbox"/> pos/features onduidelijk
<i>features</i>	betr	<input type="checkbox"/> pers <input type="checkbox"/> aanw <input type="checkbox"/> onbep <input type="checkbox"/> +negonb <input type="checkbox"/> vrag <input type="checkbox"/> bez <input type="checkbox"/> refl <input checked="" type="checkbox"/> betr <input type="checkbox"/> 1 <input type="checkbox"/> 2 <input type="checkbox"/> 3 <input type="checkbox"/> ev <input type="checkbox"/> mv <input type="checkbox"/> +nom <input type="checkbox"/> +nonnom <input type="checkbox"/> +forme <input type="checkbox"/> +formn <input type="checkbox"/> +formr <input type="checkbox"/> +forms	
	<input type="button" value="-"/>		
<i>pos</i>	BW	<input type="radio"/> N <input type="radio"/> ADJ <input type="radio"/> WW <input checked="" type="radio"/> BW <input type="radio"/> VNW <input type="radio"/> LID <input type="radio"/> TW <input type="radio"/> VZ <input type="radio"/> VG <input type="radio"/> SPEC <input type="radio"/> LET	
<i>features</i>	+pers	<input type="checkbox"/> +aanw <input type="checkbox"/> +gener <input type="checkbox"/> +onbep <input type="checkbox"/> +vrag <input checked="" type="checkbox"/> +pers <input type="checkbox"/> +vz <input type="checkbox"/> +neg <input type="checkbox"/> +negcl <input type="checkbox"/> +betr <input type="checkbox"/> +comp <input type="checkbox"/> +super <input type="checkbox"/> +prtcl	





# POS-tagging data: Gustave tool

Adelheid controle tagging

http://localhost:3035/adelheid\_check.html

557 / 1219

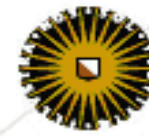
wonder vruchtbare akker des verstands . zo ghij wakker verhoet datter de ikker gheen onkruid meer in en zaaije (die

Gebruik pijljestoetsen

*datterde*    *dat | ter*

	<i>huidig</i>	<i>controle</i>	
<i>lemma</i>	dat+er	<input type="text" value="dat+er"/>	<input type="checkbox"/> modern alternatief
<i>pos</i>	VNW	<input type="radio"/> N <input type="radio"/> ADJ <input type="radio"/> WW <input type="radio"/> BW <input checked="" type="radio"/> VNW <input type="radio"/> LID <input type="radio"/> TW <input type="radio"/> VZ <input type="radio"/> VG <input type="radio"/> SPEC <input type="radio"/> LET	<input type="checkbox"/> pos/features onduidelijk
<i>features</i>	betr	<input type="checkbox"/> pers <input type="checkbox"/> aanw <input type="checkbox"/> onbep <input type="checkbox"/> +negonb <input type="checkbox"/> vrag <input type="checkbox"/> bez <input type="checkbox"/> refl <input checked="" type="checkbox"/> betr <input type="checkbox"/> 1 <input type="checkbox"/> 2 <input type="checkbox"/> 3 <input type="checkbox"/> ev <input type="checkbox"/> mv <input type="checkbox"/> +nom <input type="checkbox"/> +nonnom <input type="checkbox"/> +forme <input type="checkbox"/> +formn <input type="checkbox"/> +formr <input type="checkbox"/> +forms	
	<input type="button" value="-"/>		
<i>pos</i>	BW	<input type="radio"/> N <input type="radio"/> ADJ <input type="radio"/> WW <input checked="" type="radio"/> BW <input type="radio"/> VNW <input type="radio"/> LID <input type="radio"/> TW <input type="radio"/> VZ <input type="radio"/> VG <input type="radio"/> SPEC <input type="radio"/> LET	
<i>features</i>	+pers	<input type="checkbox"/> +aanw <input type="checkbox"/> +gener <input type="checkbox"/> +onbep <input type="checkbox"/> +vrag <input checked="" type="checkbox"/> +pers <input type="checkbox"/> +vz <input type="checkbox"/> +neg <input type="checkbox"/> +negcl <input type="checkbox"/> +betr <input type="checkbox"/> +comp <input type="checkbox"/> +super <input type="checkbox"/> +prtcl	





# Availability of the corpus

- The corpus, including the manual of the tagset and the way in which the sociolinguistics enrichment is defined will become part of Nederlab
- Problem: the letters of P.C. Hooft are from an edited volume (part of DBNL) which is copyrighted.





# Provenance for researchers:

- Metadata about the source file (preprocessing):
  - which version of the text (which print? Date of original tekst? Author? Editor? Etc.)
- Information of the input text into the tools:
  - Editorial matter separated or not?
  - Notes separated or not?



# Provenance for researchers:

- Information about the enrichment
  - What information has been added?
  - How has that information been added?  
Manually or automatically?
    - If automatically: with which tool?
    - If manually: according to what protocol?
  - How is that information defined?



# Provenance for researchers:

- Information about the query:
  - What was the query used to search for the data?
  - Are (all or part of) the data manually checked or not?



Thanks for Your Attention!