# The Mediascape of Dutch Chroniclers (1500-1850)
## Labeling Media Mentions in Early Modern Chronicles Using CRF

Alie Lassche & Roser Morante

Leiden University, VU Amsterdam

*DH Benelux. June 3, 2021*

# Outline I

1. Chronicles and the early modern mediascape

2. Corpus and annotation task

3. Machine learning media mentions

4. Discussion

# Outline

# Chronicles and the early modern mediascape

*Chronicling Novelty. New knowledge in the Netherlands, 1500-1850*

- Project managers: Judith Pollmann (Leiden University) & Erika Kuijpers (VU Amsterdam)
- About: Circulation and evaluation of new knowledge, ideas and technologies among a non-specialist public
- `www.chroniclingnovelty.com`



## Changing mediascapes and the collection of knowledge

- Explore how the use of computational methods allow me to get more insight in the media early modern chroniclers used and the information they received

# Chronicles and the early modern mediascape

# Chronicles and the early modern mediascape

- The whole of sources of information available to the early modern chronicler
- Two elements:
  - **The media that are mentioned by the authors of the chronicle**
  - The information that the authors are reporting on, which has been obtained from the media

# Outline

1. Chronicles and the early modern mediascape

2. **Corpus and annotation task**

3. Machine learning media mentions

4. Discussion

# Corpus

## Corpus characteristics

- 350 manuscripts, 70,000 pages
- Written in the Dutch language in the Low Countries between 1500 and 1850
- Transcribed and annotated with Transkribus (HTR) and the help of volunteers (VeleHanden)

Three labels

- `receiver`, `source`, `perception` (oral/heard, written/read, seen, else)

This morning `<source>`mayor Vorsterman`</source>` came `<perception: oral/heard>` telling `</perception>` `<receiver>`us`</receiver>` that because of the disease, no one was allowed to be buried in the church

# Annotation task

Inter Annotator Agreement shows difficulty of the task

|            | F-score 1 | F-score 2 |
|------------|-----------|-----------|
| *all labels*  | 0.589     | 0.742     |
| *source*      | 0.208     | 0.764     |
| *receiver*    | 0.777     | 0.619     |
| *perception*  | 0.707     | 0.727     |

Table: Inter Annotator Agreement in the media annotation task.

Used data in first experiments: 12 volume chronicle by Jozef Van Walleghem about Bruges (1779 - 1800)

| chronicle      | p (% labeled) | n sources | n receivers | n perceptions |
|----------------|---------------|-----------|-------------|---------------|
| Van Walleghem  | 1165 (17%)    | 519       | 272         | 510           |

Table: Characterics of the data.

# Outline

# Machine learning media mentions

- Classifier should assign a label to every token
  - `O`, `source-B`, `source-I`, `receiver-B`, `receiver-I`, `perception-B`, or `perception-I`
- Word vector model was trained on the data, using `fastText`
- 100 dimensional vector were used as 100 features for the classifier
- Model was trained using Conditional Random Fields (CRF) in `sklearn-crfsuite`
  - Able to deal with sequential data implicitly
  - Useful when working with sequence labeling tasks
  - Calculates features during training
- Information about the previous and next word as features (onsets, offsets, digits, vectors)
- Train set and test set were manually split (70% training and 30% testing of every volume)

# Machine learning media mentions

|  | precision | recall | f1-score | support |
|---|---|---|---|---|
| perception-B | 0.643 | 0.552 | 0.594 | 134 |
| receiver-B | 0.828 | 0.578 | 0.681 | 83 |
| source-B | 0.513 | 0.488 | 0.500 | 121 |
| source-I | 0.670 | 0.448 | 0.537 | 145 |
| perception-I | 0.500 | 0.143 | 0.222 | 7 |
| receiver-I | 0.125 | 0.111 | 0.118 | 9 |

Table: F-scores per label.

# Preliminary results

| label n | label n+1 | weight |
|---|---|---|
| O | O | 4.163518 |
| receiver-I | receiver-I | 4.029864 |
| receiver-B | receiver-I | 3.232955 |
| source-I | source-I | 2.888268 |
| perception-I | perception-I | 2.374591 |
| source-B | source-I | 1.727776 |

Table: Top likely transitions.

| weight | label | feature |
|---|---|---|
| 1.988054 | source-I | +1:word.isdigit() |
| 1.744118 | perception-B | word[-5:]: lesen |
| 1.472060 | perception-B | word[-5:]: ndigt |
| 1.466354 | O | bias |
| 1.145304 | source-B | -1:word[-5:]: nde |

Table: Top positive state features.

# Preliminary results

## Three main types of errors

- The token with which the label starts: which token gets the suffix –B?
  - ▶ **de** Gendsche Gazette
  - ▶ **alderverschrikkelijkste** berichten
- Lonely tokens: model is better in predicting when there is a sequence of labeled tokens
  - ▶ decreet
  - ▶ gerugt
- Consistency of the model
  - ▶ The word *men* often appears as `receiver`, but not always
  - ▶ Model labels media mentions that are overlooked by the annotators
  - ▶ Annotator confused `source` with `receiver`, but model assigned the correct labels

# Outline

# Discussion

- **Removal of punctuation**: our data consisted of one long sentence
  - Problematic when trying to optimize the model
  - Early modern chroniclers were very inconsistent in using punctuation
- **Lowering all characters**: mentions of oral sources or newspaper could not be recognized by its capital
  - Early modern chroniclers were very inconsistent in using capitals
- **Cluster word vectors**
- **Add a lexicon**: using a lexicon with frequently used words might improve the labeling of lonely tokens
- **Goal**: train a model that can be used to label media mentions in the whole corpus, and that can be used as a method to facilitate the close reading

Thanks for your attention!