

# Vorbesprechung/Introductory Meeting: Text Mining for Historical Documents

Caroline Sporleder and Martin Schreiber

Computational Linguistics &  
Kultur- und Mediengeschichte  
Universität des Saarlandes

Wintersemester 2010/11

18.01.2011

## Organisational Stuff

# What is it?

## Project Seminar

- a theoretical part (class presentations)
- a practical part (group work: design and implementation of a small or not so small project)
- interdisciplinary (history students and CS/Coli students)

## Course Objectives

- overview of text mining techniques
- awareness of challenges of the cultural heritage domain
- hands-on experience with text mining techniques
- working in an interdisciplinary environment and communicating with non-experts

No previous knowledge of text mining or NLP is necessary, however:

- interest in language (processing) and
- interest in history
- programming experience

. . . are essential!

There will also be **serious group work!**

# Credits

## Coli Students

- 5 CPs (150 hours)
- presentation (20-30 minutes)
- participation in group work (practical work / implementation)
- written report (one per group)

## Computer Science Students

- 7 CPs (210 hours)
- presentation (20-30 minutes)
- written summary of the presentation or another topic from the theoretical part (around 5 pages)
- participation in group work (larger amount of practical work than Coli students)
- written report (one per group)

⇒ presentations should be prepared for the first day of the course  
⇒ the written reports/summaries are due on 1.4.2011

# Composition of Grades

- 70% project work and report (one grade per group, unless this would lead to serious unfairness)
- 30% for presentation (and summary)

## Presentations

### Content

- main aspects of the paper(s) covered
- good explanations of content (examples, graphics)
- **main points of presentation understandable to non-experts**
- discussion of advantages/disadvantages of presented method
- use of additional (printed) references (optionally)

### Form

- slides contain references
- slides are well-structured
- the topic was well-presented (interaction with audience etc.)
- the presentation is not too short
- the presentation is not (much) too long

## Group Work and Report

- project is well thought through (relevance for historians, use of existing methodology, design and implementation)
- report is well-written (motivation/relevance, references to published work, description of design)
- how well did the group work together?

# When and Where

When: 21.2.-04.03.2011/11.3.2011  
9:30(?) - 18:00

Where: to be decided

# Background

## A growing field:

- cultural heritage institutes possess large collections of data (primary data and metadata)
- more and more digitisation projects (national and international):
  - digitisation as a safeguard against data loss
  - governments have come to see CH as a valuable asset
  - digitised data can be accessed more easily

## Unlocking the value of CH collections

- efficient retrieval of information is crucial (e.g. automatic searching)
- but more is needed for sophisticated information access (e.g., automatic structuring of unstructured documents, linking of related documents, disambiguation of person names etc.)
- much primary and nearly all metadata are textual  
⇒ text mining and natural language processing (NLP) play a big role

## CH as a testbed for NLP technology

- noise introduced during digitisation
- archaic language, 'free' orthography
- few NLP resources for these domains
- multi-lingual and multi-modal content
- variety of data formats and knowledge representation standards (hamper interoperability between collections)

⇒ tools need to be very robust and make use of all available resources

# Example: From paper to information (1)

## Original Manuscript

1031.	<i>Mesophorus phaeostictus</i>	x	Metzner
1032.	<i>Dactyloceras</i>	x	
1033.	<i>Ostolites magnificatus</i>	x	Metzner
1034.	<i>Agyneta</i>	x	
1035.	<i>Calopoda elongata</i>	x	
1036.	<i>Agyneta</i>	x	
Fleck: B. 46			
1037.	<i>Choristoneura fumiferana</i>	x	
1038.	<i>Thaumetopoea pityocampa</i>	x	
1039.	<i>Choristoneura fumiferana</i>	x	
1040.	<i>Zelotinus nocturnus</i>	x	
1041.	<i>Hypena leucophaea</i>	x	
1042.	<i>Cochylis rosaceana</i>	x	
1043.	<i>Cochylis rosaceana</i>	x	
1044.	<i>Lycophantis pyralis</i>	x	
1045.	<i>Pyralis malvae</i>	x	
1046.	<i>Pyralis malvae</i>	x	
1047.	<i>Pyralis malvae</i>	x	
1048.	<i>Pyralis malvae</i>	x	
1049.	<i>Pyralis malvae</i>	x	
1050.	<i>Pyralis malvae</i>	x	
1051.	<i>Pyralis malvae</i>	x	
Fleck: B. 47			
1052.	<i>Sphingopus syriacus</i>	x	
1053.	<i>Bombyx Mori</i>	x	
1054.	<i>Stathmox monacha</i>	x	
1055.	<i>Solar melan</i>	x	
1056.	<i>Polyphemus argus</i>	x	
1057.	<i>Saurida tibialis</i>	x	
1058.	<i>Solar melan</i>	x	
1059.	<i>Amata phegea</i>	x	
1060.	<i>Pollina hebraica</i>	x	
1061.	<i>Polyphemus marginatus</i>	x	
1062.	<i>Polyphemus marginatus</i>	x	
1063.	<i>Egista angulifera</i>	x	
1064.	<i>Rhyacia minuta</i>	x	
1065.	<i>Hydriomena nevada</i>	x	
1066.	<i>Pyrgula pomona</i>	x	
1067.	<i>Polyommatus amandus</i>	x	
1068.	<i>Charaxes offrenda</i>	x	
1069.	<i>Paraceraspilus leda</i>	x	
1070.	<i>Lycaenidae</i>	x	
1071.	<i>Lycaenidae</i>	x	
1072.	<i>Lycaenidae</i>	x	
1073.	<i>Lycaenidae</i>	x	
1074.	<i>Lycaenidae</i>	x	
1075.	<i>Lycaenidae</i>	x	
1076.	<i>Lycaenidae</i>	x	
1077.	<i>Lycaenidae</i>	x	
1078.	<i>Lycaenidae</i>	x	
1079.	<i>Lycaenidae</i>	x	
1080.	<i>Lycaenidae</i>	x	
1081.	<i>Lycaenidae</i>	x	
1082.	<i>Lycaenidae</i>	x	
1083.	<i>Lycaenidae</i>	x	
1084.	<i>Lycaenidae</i>	x	
1085.	<i>Lycaenidae</i>	x	
1086.	<i>Lycaenidae</i>	x	
1087.	<i>Lycaenidae</i>	x	
1088.	<i>Lycaenidae</i>	x	
1089.	<i>Lycaenidae</i>	x	
1090.	<i>Lycaenidae</i>	x	
1091.	<i>Lycaenidae</i>	x	
1092.	<i>Lycaenidae</i>	x	
1093.	<i>Lycaenidae</i>	x	
1094.	<i>Lycaenidae</i>	x	
1095.	<i>Lycaenidae</i>	x	
1096.	<i>Lycaenidae</i>	x	

# Example: From paper to information (2)

## Digital Copy

Film II, pagina 18 blm 37 van de no's 15 t/m 22		5
	RMNH 15034	
	RMNH 15480	
Reg.nr. 15139		15. <u>Gonatodes humeralis</u> , Cultuurtuin, Paramaribo, op de stammen van bomen, juist boven grond o.a. kantkrab, palmtakken, boom met paarse vruchten aan stam. ± 9.30 - 16.30. Etige aan basis stam. <del>met</del> cacao onder bladeren.
	RMNH 15690	16. <u>Ameiva ameiva</u> , Cultuurtuin, Paramaribo, 3-V-1968, wegberm, tussen gras, onder stam. ± 10 u.
Uitgenomen rex.: gezield med Paris Voorzien van A. B. L. F. en f.	→ RMNH 15245	17. <u>Chemidophorus lemniscatus</u> , Cultuurtuin, Paramaribo, 3-V-1968, wegberm tussen gras, onder stam. ± 10 u.
Reg.nr. 15230	top lang blauw, rest rase	18. <u>Anolis auratus</u> , Cultuurtuin, Paramaribo, 3-V-1968, wegberm, lage keurheid vegetatie, tot 50 cm hoog. ± 10 u.
	Reg.nr. 15700.	19. <u>Gymnophthalmus speciosus</u> , Cultuurtuin, Paramaribo, 3-V-1968, (Cacao boom) onder afgevallen blad van boom met paarse vruchten aan stam, onder en vermoedelijk ook in vermoeide boomstammen. 9.30 u - 12.30 u zeer talrijk, 15.00 - 17.00 u minder talrijk, zonnend op o.a. plankwortsels kantkrab.
	Reg.nr. 16738	20. <u>Cercosaura ocellata</u> , Cultuurtuin, Paramaribo, 3-V-1968. onderafgevallen blad en gemerde <u>Tradescantia</u> , licht bos, boom.
	RMNH 16207	21. <u>Hylos rubra</u> , Cultuurtuin, Paramaribo, 3-V-1968, op stam van boom met paarse vruchten aan stam (=cacao)
289, onderkant stam oranje	RMNH 16293	22. <u>Leptodeiralus apodusoides</u> (det. Dr. A.R. Meyer, 1972), Cultuurtuin, Paramaribo, 3-V-1968, tussen afgevallen blad van boom met paarse vruchten op stam 11.30 u. Naar grote (± 2 m) sprongen.
		23. <u>Hemidactylus mabouia</u> , Paramaribo, op buitenmuur huis o.a. Waterkant, 3-V-1968, l.d. F. Westerling.
		24. <u>Gonatodes humeralis</u> , Paramaribo, Cultuurtuin, 4-V-1968. op onderkant boomstam.

## Example: From paper to information (3)

### Transcript

1 ex. *Phyllobates femoralis* At base of tree on small island, primary forest, 20.45-22.00 u. RMNH 23865

*Lithodytes lineatus*, Brownsberg, aan voet, onder stuk rot hout, 13.07.1968, 8.45 u., RMNH 26076

Dorsolateraal strepen heldergeel, tekening op dijen vuurrood, veel feller als bij *Phyllobates femoralis*.

*Gonocephalus auritus* Meyer, 3 ex. (1 juv.), Misool. Hoedt 1867.

RMNH 17656 *Eleutherodactylus zeuctotylus* 1 [vrouw] Lelygebergte, 4 km N.O. van airstrip, distr. Marowijne, Suriname, 19-VIII-1975, onder stuk hout, 610m, 1 [plus] d M. S. Hoogmoed.

# Finding Information: Problems with keyword search

**Aim:** find all specimens of “*Phyllobates femoralis*”

1 ex. *Phyllobates femoralis* At base of tree on small island, primary forest, 20.45-22.00 u. RMNH 23865

*Lithodytes lineatus*, Brownsberg, aan voet, onder stuk rot hout, 13.07.1968, 8.45 u., RMNH 26076

Dorsolateraal strepen heldergeel, tekening op dijen vuurrood, veel feller als bij *Phyllobates femoralis*.

*Gonyocephalus auritus* Meyer, 3 ex. (1 juv.), Misool. Hoedt 1867.

RMNH 17656 *Eleutherodactylus zeuctotylus* 1 [vrouw] Lelygebergte, 4 km N.O. van airstrip, distr. Marowijne, Suriname, 19-VIII-1975, onder stuk hout, 610m, 1 [plus] d M. S. Hoogmoed.

# Finding Information: Problems with keyword search

**Aim:** find all specimens of “*Phyllobates femoralis*”

1 ex. *Phyllobates femoralis* At base of tree on small island, primary forest, 20.45-22.00 u. RMNH 23865

*Lithodytes lineatus*, Brownsberg, aan voet, onder stuk rot hout, 13.07.1968, 8.45 u., RMNH 26076

Dorsolateraal strepen heldergeel, tekening op dijen vuurrood, veel feller als bij *Phyllobates femoralis*.

*Gonyocephalus auritus* Meyer, 3 ex. (1 juv.), Misool. Hoedt 1867.

RMNH 17656 *Eleutherodactylus zeuctotylus* 1 [vrouw] Lelygebergte, 4 km N.O. van airstrip, distr. Marowijne, Suriname, 19-VIII-1975, onder stuk hout, 610m, 1 [plus] d M. S. Hoogmoed.

# Segmenting Fieldbook Entries

⇒ Number, Genus, Species, Subspecies, Biotope, Location,  
Date/Time, Remarks,

1 ex. *Phyllobates femoralis* At base of tree on small island, primary forest, 20.45-22.00 u. RMNH 23865

*Lithodytes lineatus*, Brownsberg, aan voet, onder stuk rot hout, 13.07.1968, 8.45 u., RMNH 26076  
Dorsolateraal strepen heldergeel, tekening op dijen vuurrood, veel feller als bij *Phyllobates femoralis*.

*Gonyocephalus auritus* Meyer, 3 ex. (1 juv.), Misool.  
Hoedt 1867.

RMNH 17656 *Eleutherodactylus zeuctotylus* 1  
Lelygebergte, 4 km N.O. van airstrip, distr.  
Marowijne, Suriname, 19-VIII-1975, onder stuk hout,  
610m, 1 [plus] d M. S. Hoogmoed.

# Beyond keyword search

**Aim:** find all specimens of “*Phyllobates femoralis*”

**Query:** Genus=*Phyllobates* and Species=*femoralis*

1 ex. *Phyllobates femoralis* At base of tree on small island, primary forest, 20.45-22.00 u. RMNH 23865

*Lithodytes lineatus*, Brownsberg, aan voet, onder stuk rot hout, 13.07.1968, 8.45 u., RMNH 26076

Dorsolateraal strepen heldergeel, tekening op dijen vuurrood, veel feller als bij *Phyllobates femoralis*.

*Gonyocephalus auritus* Meyer, 3 ex. (1 juv.), Misool. Hoedt 1867.

RMNH 17656 *Eleutherodactylus zeuctotylus* 1 [vrouw] Lelygebergte, 4 km N.O. van airstrip, distr. Marowijne, Suriname, 19-VIII-1975, onder stuk hout, 610m, 1 [plus] d M. S. Hoogmoed.

# Beyond keyword search

**Aim:** find all specimens of “*Phyllobates femoralis*”

**Query:** Genus=*Phyllobates* and Species=*femoralis*

**1 ex. *Phyllobates femoralis*** At base of tree on small island,  
primary forest, **20.45-22.00 u.** RMNH 23865

*Lithodytes lineatus*, Brownsberg, aan voet, onder  
stuk rot hout, **13.07.1968**, **8.45 u.**, RMNH 26076  
Dorsolateraal strepen heldergeel, tekening op dijen  
vuurrood, veel feller als bij *Phyllobates femoralis*.

*Gonyocephalus auritus* Meyer, **3 ex.** (1 juv.), Misool.  
Hoedt 1867.

RMNH 17656 *Eleutherodactylus zeuctotylus* 1 [vrouw]  
Lelygebergte, 4 km N.O. van airstrip, distr.  
Marowijne, Suriname, **19-VIII-1975**, onder stuk hout,  
610m, 1 [plus] d M. S. Hoogmoed.

# Topics

- Pre-processing
- Preservation of Digital Data
- Non-Standard Language
- Semantic Web
- Metadata
- Information Extraction
- Text Mining
- Multi-Modal Data
- Personalisation

A detailed list of topics can be found on the web <http://www.coli.uni-saarland.de/courses/tm-hist11/readings.html>.