# Mining Paper Catalogues

A Multilingual Solution to Reduce Verbose Fields to Consistent Terminology

Dr Tim Evans

**UNIVERSITY** *of York*

Felix Kußmaul

**University of Cologne**

31 August 2017

EAA Maastricht 2017

# MOTIVATION

**Figure 1:** Sample from *Conspectus* catalogue.

# Oh dear!

### Problem

Running texts contain a lot of *irrelevant information* (for machine processing).

This makes database lookups without keywords <span style="color:orange">extremely inefficient</span>.

## What we **have**:

### Production

Subform 23.1 is probably made in many parts of Italy; examples in Padana ware do not show applied decoration. Subform 23.2 is made in Italy but apparently not in the Padana region.

### Distribution

Subform 23.2 is very common throughout the Mediterranean region, with sporadic examples found in the North and in North-Italy; Subform 23.1 is relatively uncommon.

UNSTRUCTURED DATA

## What we **want**:

```
{
  "form": "23.1",
  "origin": "Italy",
  "decoration": "none",
  "occurs": "uncommon"
},
{
  "form": "23.2",
  "origin": "Italy, not Padana",
  "occurs": "Mediterranean region;
            North-Italy"
}
```

STRUCTURED DATA

## What we have:

### Production
Subform 23.1 is probably made in many parts of Italy; examples in Padana ware do not show applied decoration. Subform 23.2 is made in Italy but apparently not in the Padana region.

### Distribution
Subform 23.2 is very common throughout the Mediterranean region, with sporadic examples found in the North and in North-Italy; Subform 23.1 is relatively uncommon.

UNSTRUCTURED DATA

## What we want:

```
{
    "form": "23.1",
    "origin": "Italy",
    "decoration": "none",
    "occurs": "uncommon"
},
{
    "form": "23.2",
    "origin": "Italy, not Padana",
    "occurs": "Mediterranean region;
               North-Italy"
}
```

STRUCTURED DATA

## What we **have**:

### Production

Subform 23.1 is probably made in many parts of Italy; examples in Padana ware do not show applied decoration. Subform 23.2 is made in Italy but apparently not in the Padana region.

### Distribution

Subform 23.2 is very common throughout the Mediterranean region, with sporadic examples found in the North and in North-Italy; Subform 23.1 is relatively uncommon.

**UNSTRUCTURED DATA**

## What we **want**:

```
{
    "form": "23.1",
    "origin": "Italy",
    "decoration": "none",
    "occurs": "uncommon"
},
{
    "form": "23.2",
    "origin": "Italy, not Padana",
    "occurs": "Mediterranean region;
              North-Italy"
}
```

**STRUCTURED DATA**

# TEXT MINING

# Classification

## Definition: Text Mining

**Text Mining** is a general term covering several different ideas, e. g.:

- Information retrieval
- Statistical analysis
- Information extraction
- ...

# Classification

## Definition: Text Mining

**Text Mining** is a general term covering several different ideas, e. g.:

- Information retrieval
- Statistical analysis
- Information extraction
- …

# Classification

## Definition: Text Mining

**Text Mining** is a general term covering several different ideas, e. g.:

- Information retrieval
- Statistical analysis
- Information extraction
- ...

# Classification

## Definition: Text Mining

**Text Mining** is a general term covering several different ideas, e. g.:

- Information retrieval
- Statistical analysis
- **Information extraction**
- ...

# Classification

## Definition: Text Mining

**Text Mining** is a general term covering several different ideas, e. g.:

- Information retrieval
- Statistical analysis
- **Information extraction**
- …

# Information Extraction

## Definition: Information Extraction (IE)

*"[IE] is the task of automatically extracting structured information from unstructured […] documents."*

Some other facts about Information Extraction:

- Computer scientists have a hard time with IE (for over 30 years now!)
- IE is really super difficult and often inaccurate.

**Definition: Information Extraction (IE)**

*"[IE] is the task of automatically extracting structured information from unstructured […] documents."*

Some other facts about Information Extraction:

- Computer scientists have a hard time with IE (for over 30 years now!)
- IE is really super difficult and often inaccurate.

## Definition: Information Extraction (IE)

*"[IE] is the task of automatically extracting structured information from unstructured [...] documents."*

Some other facts about Information Extraction:

- Computer scientists have a hard time with IE (for over 30 years now!)
- IE is really super difficult and often inaccurate.

## DISCLAIMER

In this presentation, we show **preliminary** results, as this project is still work in progress.

**Figure 2:** IE Process Pipeline.

# POS-Tagging

jumps

| The | quick | brown | fox | jump | over | the | lazy | dog | . |
|-----|-------|-------|-----|------|------|-----|------|-----|---|
| DT | JJ | JJ | NN | VBD | IN | DT | JJ | NN | . |

Figure 3: POS-tagging examples after lemmatisation.

**Figure 3:** POS-tagging examples after lemmatisation.

## Rule-based approach

- High precision, but lower recall
  ⇒ Many many rules?!

## Machine-learning approach

- Needs to be trained!

- Lower precision, but high recall

## Rule-based approach

- High precision, but lower recall
  ⇒ **Many many rules?!**

**K 612**                                    Abb. 117,2

*Amphoriskos* mit breitem, annähernd zylindrischem Hals und davon abgesetzter Schulter. Schmaler, aussen vorkragender Wulstrand mit an der Innenseite umlaufender breiter Riefe. Zwei Bandhenkel vom Mündungsrand zur Schulter. Form des Gefässkörpers unbekannt.

Ton   I,B mit rotem Überfang

M.    Dm 9.5–10 cm

Vk.   Selten

Die Form fand sich mit Scherbenmaterial des späten 3. Jhs. vergesellschaftet im Nilmesserbecken (Schnitt 2710, Schicht E, s. Abb. 1; vgl. dazu JARITZ, *Elephantine* III, S. 51 mit Anm. 334). Ein Mündungsfragment mit Henkel stammt aus der das 3. und 4. Jh. umfassenden Schicht C/D (Schnitt 2710, s. Abb. 1).

Dat.   Spätes 3.–4. Jh. n. Chr.

## Machine-learning approach

- Needs to be trained!

- Lower precision, but high recall

## Rule-based approach

- High precision, but lower recall
  ⇒ Many many rules?!

**K 612**         Abb. 117,2
*Amphoriskos* mit breitem, annähernd zylindrischem Hals und davon abgesetzter Schulter. Schmaler, aussen vorkragender Wulstrand mit an der Innenseite umlaufender breiter Riefe. Zwei Bandhenkel vom Mündungsrand zur Schulter. Form des Gefässkörpers unbekannt.
Ton    I,B mit rotem Überfang
M.      Dm 9.5–10 cm
Vk.    Selten
Die Form fand sich mit Scherbenmaterial des späten 3. Jhs. vergesellschaftet im Nilmesserbecken (Schnitt 2710, Schicht E, s. Abb. 1; vgl. dazu Jaritz, *Elephantine* III, S. 51 mit Anm. 334). Ein Mündungsfragment mit Henkel stammt aus der das 3. und 4. Jh. umfassenden Schicht C/D (Schnitt 2710, s. Abb. 1).
Dat.    Spätes 3.–4. Jh. n. Chr.

## Machine-learning approach

- Needs to be trained!

- Lower precision, but high recall

## Rule-based approach

- High precision, but lower recall
  $\Rightarrow$ Many many rules?!

**K 612**                                    Abb. 117,2
*Amphoriskos* mit breitem, annähernd zylindrischem Hals und davon abgesetzter Schulter. Schmaler, aussen vorkragender Wulstrand mit an der Innenseite umlaufender breiter Riefe. Zwei Bandhenkel vom Mündungsrand zur Schulter. Form des Gefässkörpers unbekannt.
Ton   I,B mit rotem Überfang
M.      Dm 9.5–10 cm
Vk.    Selten
Die Form fand sich mit Scherbenmaterial des späten 3. Jhs. vergesellschaftet im Nilmesserbecken (Schnitt 2710, Schicht E, s. Abb. 1; vgl. dazu JARITZ, *Elephantine* III, S. 51 mit Anm. 334). Ein Mündungsfragment mit Henkel stammt aus der das 3. und 4. Jh. umfassenden Schicht C/D (Schnitt 2710, s. Abb. 1).
Dat.   Spätes 3.–4. Jh. n. Chr.

## Machine-learning approach

- Needs to be trained!
- Lower precision, but high recall

MISSING FIGURE

# Relation Extraction

| Subject | Relation | Object |
|---|---|---|
| quick brown fox | jump over | lazy dog |

## Adapting the NER

Stanford CoreNLP only recognises 8 entities types:

| | |
|---|---|
| PERSON | DATE |
| ORGANIZATION | TIME |
| LOCATION | MONEY |
| PERCENT | MISC |

So we have to add the custom type FORM. Adjusting DATE also necessary.

## Temporal Expressions

With *HEIDELTIME* temporal expressions are mapped to TIMEX3 standard

$$\text{around 140 B.C.} \longmapsto \text{APPROX BC0140}$$

$$\text{second quarter first century B.C.} \longmapsto \text{XXXX-Q2 BC00}$$

$$\text{first half third century A.D.} \longmapsto \text{XXXX-H1 02}$$

HEIDELTIME supports many other languages, e.g. German, Italian, French, …

HEIDELPLACE?!

## Temporal Expressions

With *HEIDELTIME* temporal expressions are mapped to TIMEX3 standard

$$\text{around 140 B.C.} \longmapsto \text{APPROX BC0140}$$

$$\text{second quarter first century B.C.} \longmapsto \text{XXXX-Q2 BC00}$$

$$\text{first half third century A.D.} \longmapsto \text{XXXX-H1 02}$$

HEIDELTIME supports many other languages, e.g. German, Italian, French, …

HEIDELPLACE?!

# MULTILINGUALISM

# Background

Two problems:
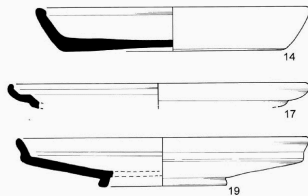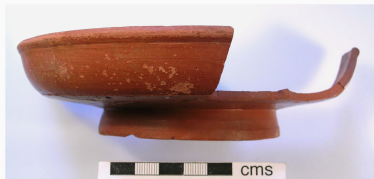
- Linguistic
- Conceptual

Figure 4: Plate, platter or dish?

# Creating controlled vocabularies

- Sherd type (e.g. rim)
- Form (e.g. plate)
- Decoration form (e.g. burnished)
- Decoration color (e.g. yellow)
- Fabric (e.g. bla)

Using tools developed for the ARIADNE project by the Hypermedia Research
Group at the University of South Wales



Creation of a neutral spine based on the Getty Institute's Art and Architecture
Thesaurus (AAT)

more bla

more bla

more bla

# Preliminary results: felix

challenges:

choice of tools, coreferences in text, eloquence of archaeologists, maybe calculating F-value?

# Mining Paper Catalogues

A Multilingual Solution to Reduce Verbose Fields to Consistent Terminology

Dr Tim Evans

Felix Kußmaul

**UNIVERSITY** *of York*

**University of Cologne**

31 August 2017

EAA Maastricht 2017