# From Flat Lists to Taxonomies: Bottom-up Concept Scheme Generation in Linked Statistical Data

Albert Meroño-Peñuela
@albertmeronyo

SemStats @ ISWC 19-10-2014

Data Archiving and Networked Services

DANS  VU  UNIVERSITY AMSTERDAM  eHumanities
Royal Netherlands Academy of Arts and Sciences

# Motivation

## Dutch Historical Censuses (1795-1971)
## [Public Historical Statistical Data]



http://lod.cedar-project.nl/cedar/

# Dimension Reusability

```
cedar:BRT_1889_02_T1-S0-K17-h a qb:Observation ;
    cedar:population "12"^^xml:integer ;
    maritalstatus:maritalStatus
        maritalstatus:single ;
    cedarterms:occupationPosition cedarterms:job-D ;
    sdmx-dimension:sex sdmx-code:sex-F ;
    cedarterms:occupation hisco:88030 ;
    sdmx-dimension:refArea gg:11150 ;
    cedarterms:belief cedar:118 ;
    cedarterms:houseType cedar:Klooster ;
    prov:wasDerivedFrom
        cedar:BRT_1889_08_T1-S0-K17 ;
    prov:wasGeneratedBy
        cedar:BRT_1889_08_T1-S0-K17-activity .
```

# LSD Dimensions



## LSD Dimensions

Counting **2461** dimensions in **588** SPARQL endpoints in Linked Statistical Data

| | Dimension URI | Label | References ▾ |
|---|---|---|---|
| 👁 | http://purl.org/linked-data/cube#measureType | measure type | 20 |
| 👁 | http://purl.org/linked-data/sdmx/2009/dimension#civilStatus | Civil Status | 19 |
| 👁 | http://purl.org/linked-data/sdmx/2009/dimension#refArea | Reference Area | 19 |
| 👁 | http://purl.org/linked-data/sdmx/2009/dimension#freq | Frequency | 19 |
| 👁 | http://purl.org/linked-data/sdmx/2009/dimension#sex | Sex | 19 |
| 👁 | http://purl.org/linked-data/sdmx/2009/dimension#refPeriod | Reference Period | 19 |
| 👁 | http://purl.org/linked-data/sdmx/2009/dimension#age | Age | 19 |
| 👁 | http://purl.org/linked-data/sdmx/2009/dimension#timePeriod | Time Period | 19 |
| 👁 | http://purl.org/linked-data/sdmx/2009/dimension#currency | Currency | 19 |
| 👁 | http://purl.org/linked-data/sdmx/2009/dimension#education… | Education Level | 19 |

Showing 1 to 10 of 2461 rows    10 ▴    records per page

<<  <  **1**  2  3  4  5  >  >>

http://lsd-dimensions.org/
https://github.com/albertmeronyo/LSD-Dimensions
Hourly JSON-LD dumps

# What if dimensions aren't out there?

- Need to build them
- Input: flat lists of non-standard values
- Output: standard concept scheme
- Knowledge intensive problem

# Dutch Historical Building Types

```
 1    Huis van Arrest
 2    Weeshuis
 3    Overige huizen
 4    St. Anthony Gasthuis
 5    Overige huizen
 6    Hotel "de Gouden Leeuw"
 7    Overige huizen
 8    Klooster der Franciscanen
 9    Overige huizen
10    Gesticht "de Goede Herder"
11    Overige huizen
12    Gesticht Ommerschans
13    Overige huizen
14    Ned. Herv. Weeshuis
15    Overige huizen
16    Karmelieten klooster
17    Overige huizen
18    St. Nicolaas-Gesticht
```
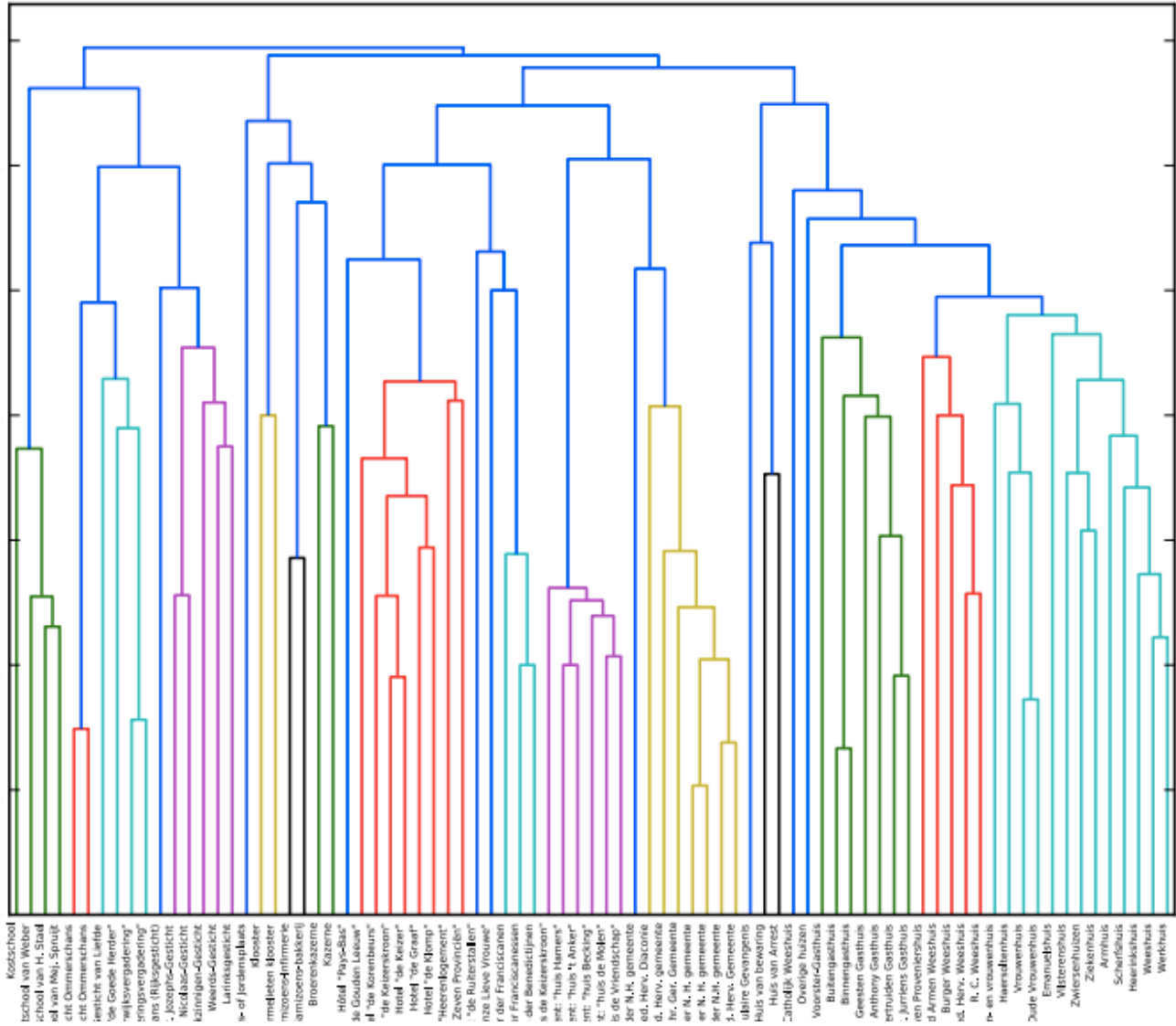
# TabCluster: data-driven concept scheme generation

https://github.com/CEDAR-project/TabCluster

Leverages lexical and semantic term properties

- Lexical Hierarchical Clustering
  - ScyPy python implementation
  - Similarity transitivity
  - Clustering threshold & other parameters
- Semantic Tagging of Clusters (WordNet & DBPedia) w/ skos:Concept
  - Term frequency-based
  - Less common broader concept category

# Output & Evaluation

# Thank you

## Questions, suggestions, comments most welcome

@albertmeronyo

https://github.com/CEDAR-project/TabCluster
http://lsd-dimensions.org/