

# Detecting and Reporting Extensional Concept Drift in Statistical Linked Data

Albert Meroño-Peñuela<sup>1,2</sup>    Christophe Guéret<sup>2</sup>  
Rinke Hoekstra<sup>1,3</sup>    Stefan Schlobach<sup>1</sup>

<sup>1</sup>Department of Computer Science, VU University Amsterdam, NL

<sup>2</sup>Data Archiving and Networked Services, KNAW, NL

<sup>3</sup>Leibniz Center for Law, Faculty of Law, University of Amsterdam, NL

October 22th, 2013

First International Workshop on Semantic Statistics  
ISWC 2013

# Motivation

## Stability of Meaning of Concepts

As the world changes continuously, concepts also change their meaning over time



We call this *concept drift*

- Smooth transitions
- Radical transitions (*concept shift*)

*Concepts* are also present in SLD

- Variable meaning/semantics (what the variable is supposed to represent?)
- Variable values/factors (*RomschKatholik*, *RomsKatholic*, *KatholicChristelijk*)

To what extent stability of meaning of these concepts is guaranteed in data collected on *very long time ranges*?

If meaning is stable

- Old models are reusable
- Backwards comparisons always make sense

Else

- New models may be necessary
- Backwards comparisons may be incorrect

# Concept Drift

## Types of change of meaning

The meaning of a concept  $C$  can change in several ways<sup>1</sup>.

- **Intension drift** occurs when there is a difference in the *properties or attributes* of two variants of the same concept ( $sim_{int}(C', C'') \neq 1$ ).
- **Extension drift** occurs when there is a difference in the *individuals that belong* to two variants of the same concept ( $sim_{ext}(C', C'') \neq 1$ ).
- **Label drift** occurs when there is a difference in the *labels* of two variants of the same concept ( $sim_{label}(C', C'') \neq 1$ ).

---

<sup>1</sup>S. Wang, S. Schlobach, M. Klein, *What Is Concept Drift and How to Measure It?*, 2010

# Key ideas

If you remember nothing else, remember this

In this talk we present two key ideas

- A method to detect extensionally drifted concepts
- A recommendation to annotate such drifts in SLD

# Extensional drift detection

## Using statistics

- We define the *extension function*  $ext(C)$  of a concept  $C$  as the **number of individuals** that belong to  $C$

$$ext(C) = |\{a : C\}|$$

- We define the *extension similarity function*  $sim_{ext}(C', C'')$  between two variants  $C', C''$  of a concept  $C$  as the function that returns the **probability that  $C'$  and  $C''$  have identical populations**<sup>2</sup>

$$sim_{ext}(C', C'') = wilcox.test(ext(C'), ext(C''))$$

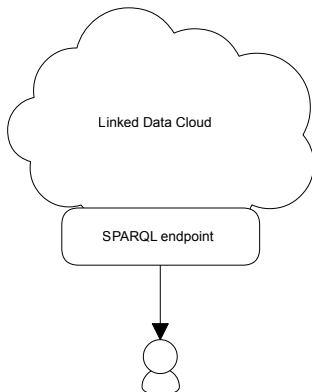
- Intuitive idea: concept variants with significantly different populations ( $p < 0.05$ ) suffer radical transitions

---

<sup>2</sup>F. Wilcoxon, *Individual comparisons by ranking methods*, 1945

# Extensional drift reporting

Using SPARQL UPDATES

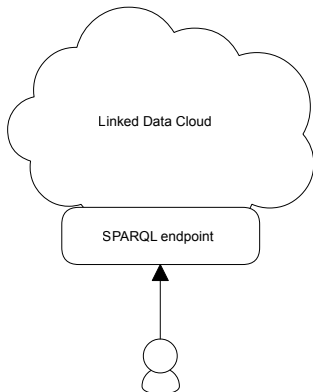


Typical analysis workflow. Users SELECT data from SLD, but analyses run offline and results are not pushed back



# Extensional drift reporting

Using SPARQL UPDATES



Proposal. After running analyses, UPDATE results (e.g. detected extensional drifts) back to the endpoint

## The Dutch historical censuses (1795 - 1791)

The image displays three components of historical Dutch census data:

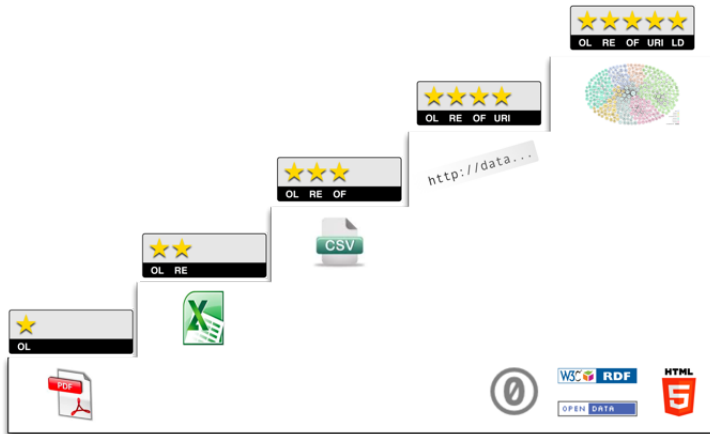
- Top Left:** A title page for the census of Opertseel in Deventer, featuring a population pyramid showing the distribution of the population by age and sex.
- Top Right:** A large, multi-column grid of census tables, likely representing individual households or individuals, with columns for names, ages, and other demographic information.
- Bottom Right:** A smaller table with handwritten entries, possibly a summary or a specific subset of the data, including names like 'Boudry Adriaan' and dates.

- Lots of statistical data (2,288 tables)
- Lots of concepts
- Big time span (176 years)

May contain lots of drifted concepts!

# Case-Study

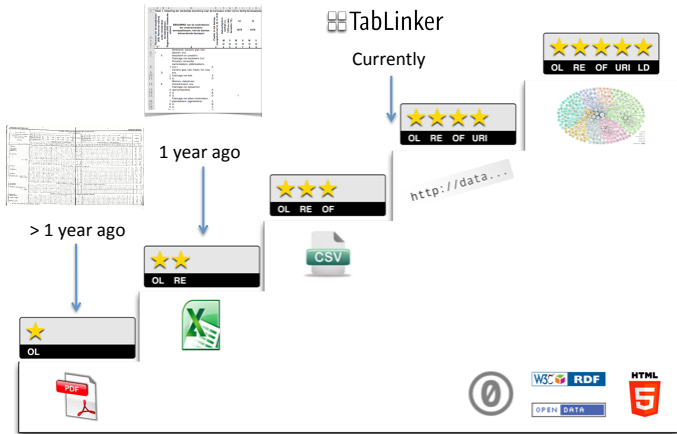
## Publishing Linked Census Data



Towards 5-star historical census data

# Case-Study

## Publishing Linked Census Data



TabLinker<sup>3</sup> : Supervised XLS2RDF converter

<sup>3</sup><https://github.com/Data2Semantics/TabLinker>

```
1 PREFIX qb: <http://purl.org/linked-data/cube#>
2 PREFIX d2s: <http://www.data2semantics.org/core/>
3 PREFIX skos: <http://www.w3.org/2004/02/skos/core#>
4 PREFIX ns: <row_property_URI>
5
6 SELECT ?d1label ... ?dnlabel ?p1label ... ?pmlabel ?population
7 FROM <named_graph_URI>
8 WHERE {
9   ?cell d2s:isObservation [ a qb:Observation ;
10                             qb:DimensionProperty ?d1 ... ?dn ;
11                             ns:property1 ?p1 ;
12                             ...
13                             ns:propertym ?pm ;
14                             qb:MeasureProperty ?population ] .
15   OPTIONAL {
16     ?cell d2s:isObservation [ns:propertyk ?pk ] .
17     ?pk skos:prefLabel ?pklabel .
18     ...
19   }
20   OPTIONAL {
21     ?cell d2s:isObservation [qb:DimensionProperty ?di ] .
22     ?di skos:prefLabel ?dilabel .
23     ?pr skos:broader ?pu .
24     ?pu skos:broader ?pv .
25     d1 ... dn skos:prefLabel ?d1label ... ?dnlabel .
26     p1 ... pm skos:prefLabel ?p1label ... ?pmlabel .
27   }
28   FILTER (?d1 IN (v1, ..., vr)) ...
29   FILTER (?dn IN (w1, ..., ws))
30 }
```



## SPARQL template for unfolding RDF Data Cubes

Census contains very heterogeneous data

Subset selection: occupation census of 1889 and 1899

Age range	Gender	Marital status	Municipality
36-50	M	G	Velsen
51-60	V	O	Zaandam
23-35	VROUWEN	O.	Haarlem
36-50	MANNEN	G.	Weesp

Class	Subclass	Occupation	Position	Population
II	b	Aanemers	A	3
IV	a	Agenten	B	1
XXI	d	Ambtenaren en beambten	C	2
I	h	Afwerken van huizen	A	5

### 1889 HISCO mappings

Arbeiders	D	99900
Kooplieden	A	41025
Winkeliers	A	41030
Handel in voorwerpen van kleding.		41025
Kooplieden	A	41025
Winkeliers	A	41030
Winkelbedienden	C	45130
Handel in voorwerpen van voeding en genot.		41025
Arbeiders	C	99900
Arbeiders	D	99900
Dépothouders	B	45130
Kooplieden	A	41025
Winkeliers	A	41030
Winkelbedienden	C	45130
Handel in voorwerpen van woning.		44130
Kooplieden	A	41025
Winkeliers	A	41030
Handel in boek- en kunstwerken (incl. dagbladen)		41030
Uitegevers	A	21110
Winkeliers	A	41030
Handel in luxe artikelen.		41025
Kooplieden	B	41025
Winkeliers	A	41030
Handel in levend vee en gevogelte		41025
Kooplieden	A	41025
Handel in andere waren.		41025
Arbeiders	D	99900
Kooplieden	A	41025
Kramers en rondventers	A	45220
Kramers en rondventers	D	45220
Magazijn- en pakhuisnechts	D	97145
Winkeliers	A	41030

### 1899 HISCO mappings

Schaapherders	D	62430
Hoenderfokkers	A	61260
Hoenderfokkers	D	61260
Vogelkweekers	A	61290
Bijenhouders	A	61290
Bijenhouders	D	61290
Boterboeren	A	77530
Boterboeren	B	77530
Boterboeren	D	77530
Melkboeren (niet melkslijters)	B	41030
Beestensnijders	A	77330
Hoenderparkhouders	A	61260
Tuinlieden	B	62740
Bloembollenkweekers	B	61270
Boomkweekers	B	61230
Boomsnoeiers	C	62730
Nettenboeters	D	75465
Scheepslossers	A	97120
Scheepslossers	C	97120
Scheepslossers	D	97120
Schelpenvisschers	A	64990
Schelpenvisschers	B	64990
Schelpenvisschers	C	64990
Schelpenvisschers	D	64990
Wiervisschers	A	64990
Vischsnijders	D	77940
Visschers	B	64100
Visschers	C	64100
Jagers	D	64960
Eendenkooihouders	A	64960
Eendenkooihouders	D	64960
Pers. in algem. dienst	C	30000

## HISCO normalization

# Case-Study

## Extensional drift detection in R

Is there extensional concept drift between two variants of the same occupation (i.e. same HISCO code)?

E.g. is there concept drift between *ship loaders* of 1889 and 1899?



Are the occupations in the two censuses comparable?

- 217 common HISCO codes
- 72.2% of all 1889 HISCO codes
- 71.1% of all 1899 HISCO codes

# Case-Study

## Extensional drift detection in R

For all  $h$  common HISCO codes

- Query for population distributions of  $h$  in 1889
- Query for population distributions of  $h$  in 1899
- Do `wilcox.test` between the two and get *p-values*

# Case-Study

## Results

HISCO Occupation	p-value
97125 Loader of ship, truck, wagon or airplane	1.83e-10
21110 General manager	4.23e-09
41025 Working proprietor (wholesale, retail trade)	1.52e-08
79100 Tailor	7.75e-07
57030 Barber, hairdresser	1.17e-04
88010 Jeweller	1.84e-04

(a) Occupations with stronger ext. drift.

Group Type	p-value
7, 8, 9 Production, transport, operators	2.03e-19
5 Service workers	1.88e-12
4 Sales workers	2.94e-08
2 Administrative and managerial	4.20e-08

(c) Major groups with stronger ext. drift.

HISCO Occupation	p-value
53190 Other cooks	1.00
75452 Lace weaver	1.00
75490 Other weavers	1.00
75990 Other spinners, weavers, knitters, dyers	1.00
77690 Other bakers, pastry cooks and confectionery makers	1.00

(b) Occupations with greater ext. stability.

Group Type	p-value
6 Agriculture, animal husbandry, fishermen, hunters	0.38
0, 1 Professional and technical	0.16
3 Clerical	1.40e-04

(d) Major groups with greater ext. stability.

Table 2: Wilcoxon test p-values per HISCO code ((a),(b)) and major group ((c),(d)).

- Late industrialization of the Netherlands (late 19th century)
- 19.35% of occupations show radical extensional drifts

# Case-Study

## Extensional drift reporting

```
1 PREFIX d2s: <http://www.data2semantics.org/core/>
2 PREFIX d2s1889: <urn:nbn:nl:ui:13-m4k-4lp>
3 PREFIX d2s1899: <urn:nbn:nl:ui:13-988-0dq>
4
5 INSERT DATA {
6   GRAPH <named_graph_URI> {
7     d2s1889:Sjouwerlieden d2s:isDrift [
8       d2s:extDrift d2s1899:Expeditie_bevrachters_bestellers_sjouwerlieden ,
9         d2s1899:Personeel_voor_laden_en_lossen ,
10        d2s1899:Personeel_voor_lading_en_lossing ,
11        d2s1899:Sjouwerlieden ;
12    d2s:weight 1.83e-10 ] . } }
```

Listing 1.3: Excerpt of the SPARQL query reporting back extensionally drifted occupational concepts. Only the drift for one occupational concept is shown. Inverse drifts from the second graph to the first are also issued.

## Closing the pull-push cycle<sup>4</sup>

<sup>4</sup><https://github.com/albertmeronyo/ConceptDrift/>

# Conclusions & Future work

- Concept drift: concepts change over time
  - It affects model reusability and backwards querying
  - Extensional drifts can be detected with `wilcox.test`
  - We SPARQL UPDATE drifts to let others know
- 
- Scale up variables, tables
  - Parametrization of the `wilcox.test` depending on the time gap
  - Integration of intension and label drift measures
  - RDF HISCO will be released soon
  - HISCO mappings will be published as separate named graphs (so that you can link yours)

Thank you  
Questions, suggestions?

---

@albertmeronyo  
<http://www.cedar-project.nl>  
<http://www.data2semantics.org>