ETH
Eidgenössische Technische Hochschule Zürich
Swiss Federal Institute of Technology Zurich

geoinformation engineering
ikg    eth zurich

IKG
Institut für Kartografie und Geoinformation

# Semantic description of spatial analysis data sets on the Web

## Problem statement and Main Goal

A problem, which is often observed on the Web is that there is lack of provenance information about spatial and statistical data sets, such as their origin, the operations performed for their generation and the people involved in the procedure. Due to this lack of information, data analysts do not know which analysis method can be meaningfully applied to a specific data set and often implement methods, which are inappropriate. As a result, they draw wrong conclusions about the data set.

The main objective of this Master thesis was to investigate how geographic data sets can be described with linked data for improved analysis. In this line of research the following research questions have been developed in order to reach the main goal of the Thesis:

• Which are the available linked data vocabularies for the main purpose of the Thesis?
• Which are the differences between these vocabularies and meta-data standards?
• In how far do they serve this purpose, and what is missing?
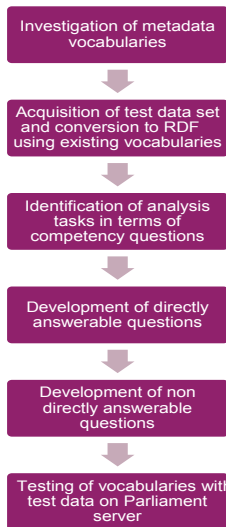• How could a linked data vocabulary for geographic and statistical analysis look like?



Figure 1: Approach and workflow of the Thesis

• Translation of each competency question into a query
• Request of a particular data set and its retrieval were possible through the Parliament server
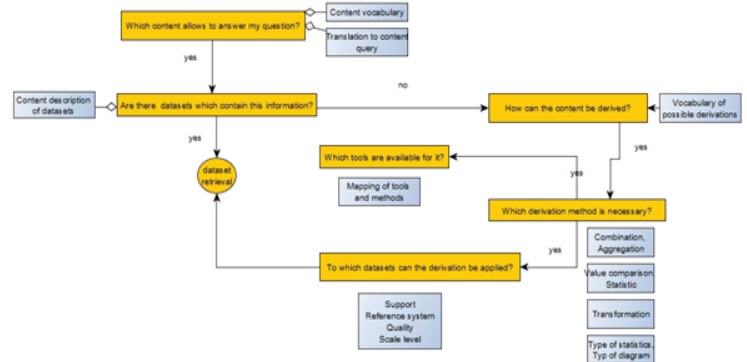


Figure 2: General methodology for retrieving data for analysis tasks



Figure 3: Example of the resulting RDF model

## Data and Methodology

• Two json files providing meta-information about the data sources and data sets included in the AURIN portal were used as data
• The principal objective of the methodology was to test vocabularies and whether they allow retrieval of data sets corresponding to answers to the competency questions that can be asked for analysis purposes
• Direclty answeable questions can be answered from the content of information in a database using content vocabularies
• Indirectly answerable questions cannot be answered directly from the content of information in the database, since there are some restrictions on the data to be retrieved based on some operations that should be carried out
• These operations are combination, comparison, transformation and calculation of statistics relevant to a particular data set.
• There is a kind of meta-information determining to which data sets the derivation can be meaningfully applied

## Conclusions

• Development of general methodology to answer analysis questions based on the content of information included in the data sets and other tools enabling users to conduct analysis operations
• Research questions were answered to a great extent
• Only two provenance properties were used, since no specific information provided in the json metadata files about the procedures, through which the data sets were generated or whether any transformations were conducted

## Outlook

• Integration of the resulting RDF model in a data portal
• Creation of metadata descriptions by the users or by automatic translation tools
• Possibility for the users to annoatate their own data with the proposed vocabularies
• Development of a running recommender system for retrieval of the most appropriate data set for analysis puproses

Author: Ariadni Gaki
Head: Prof.Dr.Martin Raubal,    Advisor: Dr.Simon Scheider

Master Thesis
Autumn Semester 2015