

# Pair Assignment 3: Gathering, Cleaning, Merging and Exploring our Data

Course: Introduction to Collaborative Social Science Data Analysis

*Malte Berneaud-Kötz & Jonas Markgraf*

*Hertie School of Governance*

*14 April 2016*

## 1. Introduction<sup>1</sup>

## 2. Gathering the Data

### 2.1 Bank Board Data

We hand-collect a unique panel dataset on the composition of Boards of Directors in Bavaria's *Sparkassen*. This dataset includes detailed information on board member profiles which enables us to identify mayors on bank boards:

- name of board members;
- occupation of board members (identifier for mayors on board);
- position within board: normal board member, chairman, or vice chairman.

Annual information on Board of Directors is hand-collected from savings banks' annual reports available in PDF format on *Bundesanzeiger* for the years from 2006 to 2015; access to data prior 2006 is proprietary (Bureau van Dijk's *Bankscope* database), which restricts our observation period. The dataset on boardroom composition constitutes the first comprehensive and systematic investigation of Bavarian savings banks' corporate governance as information on German public banks' boards has not been systematically collected yet.

---

<sup>1</sup>This paper is based on and a part of a research project by Guillermo Rosas (Washington University in St. Louis; grosas@wustl.edu) and Jonas Markgraf (Hertie School of Governance; markgraf@hertie-school.org).

## 2.2 Municipal Election Data

A database on mayoral elections in Bavaria is available from the state statistical office upon request. It contains data on direct municipal elections between 1948 and 2014. With this database we are covering 79 of the 416 German *Sparkassen* (19%) and 2,099 municipalities (19% of all municipalities in Germany). The data for mayoral elections contains the following variables:

- election date;
- name of election winner and (at least) the first loser;
- party affiliation of candidates;
- vote shares of candidates;
- dummy for competitive elections (at least two candidates);
- dummy for ‘first-time mayor’;
- number of eligible voters in voting district (size of municipality).

After obtaining the raw data, we cleaned the data set and subsequently created additional variables needed in our analysis. The individual steps taken are outline in the

## 3. Cleaning the Data

### 3.1 *Sparkassen* Data

The data we obtained on the *Sparkassen* had really long and unwieldy names, which we changed to make them more manageable. Moreover, we standardized them to follow use underscores to seperate words and use lower case.

Some of the strings which we needed to use in our analysis contained unnecessary whitespace which we trimmed using the `str_trim` function from the `stringr` package.

### **3.2 Municipal Election Data**

### **3.3 Creating Additional Variables**

In addition to the variables available in the data set already, we created variables (1) distinguishing ‘first time mayors’, (2) identifying competitive elections where there were more than one candidate, (3) the number of times a single mayor was elected, which allowed us to identify mayor’s first re-election. This last point is important because our research is interested in the first re-election of mayors specifically.

## **4. Merging the Data**

## **5. Describing the Data**

## **6. First Inferences**