# A Robust Approach for Discovering Functional Dependencies using Machine Learning Approaches

von

Philipp Jung

Philipp Jung
Matrikelnummer: 872855
16.03.2019

Gutachter:
Prof. Felix Biessmann
Dr. Zweit Gutachterin

ABSTRACT. Lorem ipsum dolor sit amet, consetetur sadipscing elitr, sed diam nonumy eirmod tempor invidunt ut labore et dolore magna aliquyam erat, sed diam voluptua. At vero eos et accusam et justo duo dolores et ea rebum. Stet clita kasd gubergren, no sea takimata sanctus est Lorem ipsum dolor sit amet.

# Contents

# 1 Introduction



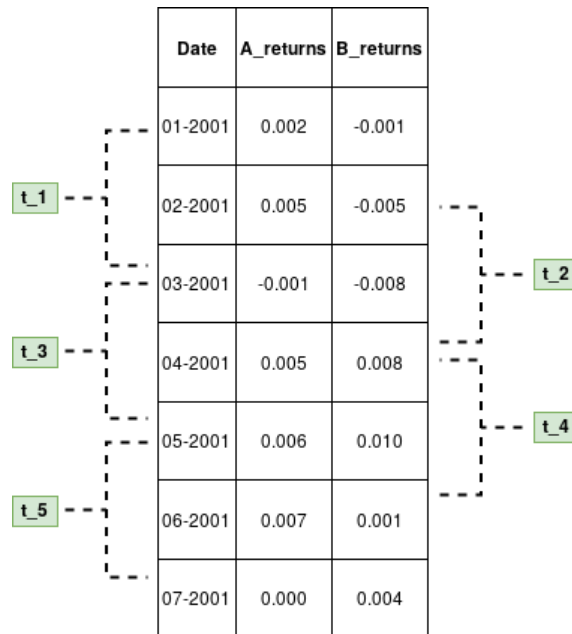| Date | A_returns | B_returns |
|------|-----------|-----------|
| 01-2001 | 0.002 | -0.001 |
| 02-2001 | 0.005 | -0.005 |
| 03-2001 | -0.001 | -0.008 |
| 04-2001 | 0.005 | 0.008 |
| 05-2001 | 0.006 | 0.010 |
| 06-2001 | 0.007 | 0.001 |
| 07-2001 | 0.000 | 0.004 |

Figure 1: Illustration of the rolling-window approach for a time-series containing seven time-steps filled with mock-data. Five subsets of length 3 divide the time-series.

This approach is schematically described in figure 1.

## 2 Theory

*Functional dependencies*(FD) are a way of expressing "*a priori* knowledge of restrictions or constraints on permissible sets of data"[Mai83, p.42] in relational database theory. In order to give a definition of FDs, some concepts stemming from relational database theory need to be introduced beforehand.

### 2.1 Relational Database Theory

A *relation scheme*[1] $R$ is a finite set of *attribute names* $\{A_1, A_2, \ldots, A_n\}$, where to each attribute name $A_i$ corresponds a set $D_i$, called *domain* of $A_i$, $1 \leq i \leq n$.

Let $\boldsymbol{D} = D_1 \cup D_2 \cup \cdots \cup D_n$, then a *relation* $r$ on relation scheme $R$ is a finite set of mappings $\{t_1, t_2, \ldots, t_p\}$ from $R$ to $\boldsymbol{D}$:

$$t_i : R \to \boldsymbol{D}, \tag{1}$$

where we call those mappings *tuples* under the constraint that [Mai83, p.2]

$$t(A_i) \subseteq D_i. \tag{2}$$

### 2.2 Definition of a Functional Dependency

For giving a definition of a FD, relation $r$ on scheme $R$ with subset $X \subseteq R$ and a single attribute $A_i \in R$ are considered. A FD $X \to A$ is said to be *valid* in $r$, if and only if

$$t_i[X] = t_j[X] \Rightarrow t_i[A] = t_j[A] \tag{3}$$

holds for all all pairs of distinct tuples $t_i, t_j \in r$.[Abe+19, p.21] We say that $X$ *functionally determines* $A$[Mai83, p.43] and name $X$ the *left side*, whilst calling $A$ the *right side*.

### 2.3 Approximate Functional Dependencies

In the field of data profiling an extensive body of theory and algorithms for FD detection has been created in the past decades. These mainly consider FDs as defined in equation 3. However, the strict detection of FDs yields results that are solely applicable in a strictly controlled environment. Real-world datasets faced by data-scientists or database engineers are often *noisy*. Entries might be spelled incorrectly and inconsistencies are to be expected.

Here goes a nice example explaining table 1 *Approximate FDs (AFDs)* relax the strict definition of FDs and introduce an error-measure.
The error measure for

---

[1]also called *relational schema* in literature[Abe+19, p.21]

| First name | Data | |
| --- | --- | --- |
| | Last name | ZIP |
| Alice | Smith | 19139 |
| Pencil | 1 | big |
| Marker | 4 | |
| Fountain Pen | 43 | green |

Table 1: AFD example I have to work out.

# 3 Execution

Lorem ipsum dolor sit amet, consetetur sadipscing elitr, sed diam nonumy eirmod tempor invidunt ut labore et dolore magna aliquyam erat, sed diam voluptua. At vero eos et accusam et justo duo dolores et ea rebum. Stet clita kasd gubergren, no sea takimata sanctus est Lorem ipsum dolor sit amet.

## 3.1 Begriffsdiskussion

Lorem ipsum dolor sit amet, consetetur sadipscing elitr, sed diam nonumy eirmod tempor invidunt ut labore et dolore magna aliquyam erat, sed diam voluptua. At vero eos et accusam et justo duo dolores et ea rebum. Stet clita kasd gubergren, no sea takimata sanctus est Lorem ipsum dolor sit amet.

# 4 Discussion

Lorem ipsum dolor sit amet, consetetur sadipscing elitr, sed diam nonumy eirmod tempor invidunt ut labore et dolore magna aliquyam erat, sed diam voluptua. At vero eos et accusam et justo duo dolores et ea rebum. Stet clita kasd gubergren, no sea takimata sanctus est Lorem ipsum dolor sit amet.

## 4.1 Begriffsdiskussion

Lorem ipsum dolor sit amet, consetetur sadipscing elitr, sed diam nonumy eirmod tempor invidunt ut labore et dolore magna aliquyam erat, sed diam voluptua. At vero eos et accusam et justo duo dolores et ea rebum. Stet clita kasd gubergren, no sea takimata sanctus est Lorem ipsum dolor sit amet.

# References

[Abe+19]   Ziawasch Abedjan et al. *Data Profiling*. 2019. ISBN: 9781681734477. DOI: `https://doi.org/10.2200/S00878ED1V01Y201810DTM052`.

[Mai83]   David Maier. *The Theory of Relational Databases*. Computer Science Pr, 1983. ISBN: 0914894420. URL: `http://web.cecs.pdx.edu/~maier/TheoryBook/TRD.html`.