

Robust Machine-Learning Approaches for Efficient Functional Dependency Approximation

Philipp Jung

Beuth University of Applied Sciences

philippjung@posteo.de

November 27, 2019

Overview

- 1 Motivation
- 2 Objectives
- 3 FD Imputer
- 4 DepDetector
- 5 Prospects

Motivation

Example of an FD

left hand side				right hand side
Id	First Name	Surname	Zip	Town
1	Alice	Smith	19139	Munich
2	Peter	Meyer	19139	Munich
3	Ana	Parker	19139	Munich
4	John	Pick	12055	Berlin
5	John	Pick	19139	Munich

Table: Example of the non-minimal FD
 $\{\text{Id, First Name, Surname, Zip}\} \rightarrow \text{Town}.$

FD Fields of Application

FDs are constraints on a relational scheme commonly used for

- schema normalization of relational databases
- data cleaning, e.g. HoloClean¹
- data exploration

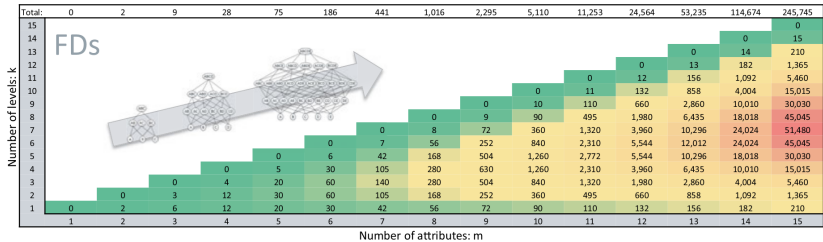
¹Heidari et al. 2019.

Challenges of FD Detection

Remark

FD detection is a particularly complex problem to solve.

FD Detection Search Space Size



The number of FD candidates for m attributes is $\mathcal{O}\left(\frac{m}{2} \cdot 2^m\right)$.

¹Image from Abedjan et al. 2019

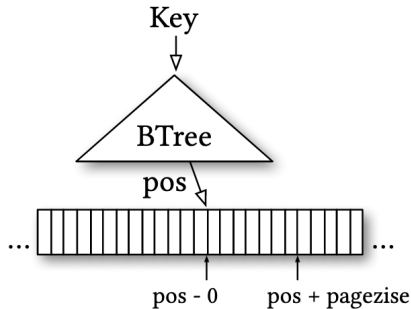
Learned Algorithms

Kraska et al. showed in 2018 that a Binary Tree can be interpreted as a learned index structure.²

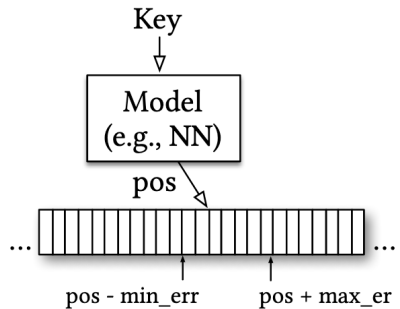
²Kraska et al. 2018.

BTree as a Learned Model

(a) B-Tree Index



(b) Learned Index



²Image from Kraska et al. 2018

Learning FDs

Remark

Learned algorithms offer a new research-directory when solving old algorithmic problems.

Objectives

Research Objectives

- Interpret FDs as *learned* constraints
- Detect FDs with machine-learning techniques
- Compare results to existing algorithms

DataWig

- DataWig is a framework for learning models to impute missing values in tables
- Data imputation: Replace missing or faulty data
- Models trained by DataWig use either regression or multi-label classification
- DataWig models can be benchmarked against FD Imputer

FD Imputer

FD Imputer: FDs as Models

- Interpret FDs as rules for data imputation
- Write FD Imputer: An imputation model entirely based on FDs
- Measure *Robustness*: Either the F1-Score or the MSE that FD Imputer obtains for a FD

How FD Imputer works

- 1 Split dataset in train-set and test-set
- 2 Detect FDs on train-set using HyFD³
- 3 For each FD, impute the right hand side for each tuple in the test-set by looking for a tuple with an identical left hand side in the train-set
- 4 Compute *Robustness* by evaluating FD Imputer's performance for each FD (F1-Score or Mean Squared Error)

³Papenbrock and Naumann 2016.

FD Imputer Functionality Example

A	B	C	D
Green	Bus	Portugal	Lisbon
Yellow	Car	Portugal	Lisbon

(a) Train set

A	B	C	D
Yellow	Bus	Spain	?
Blue	Bus	Portugal	?

(b) Test set

In this example, FD Imputer uses the FD $C \rightarrow D$.

FD Imputer Functionality Example

A	B	C	D
Green	Bus	Portugal	Lisbon
Yellow	Car	Portugal	Lisbon

(a) Train set

A	B	C	D
Yellow	Bus	Spain	-
Blue	Bus	Portugal	Lisbon

(b) Test set

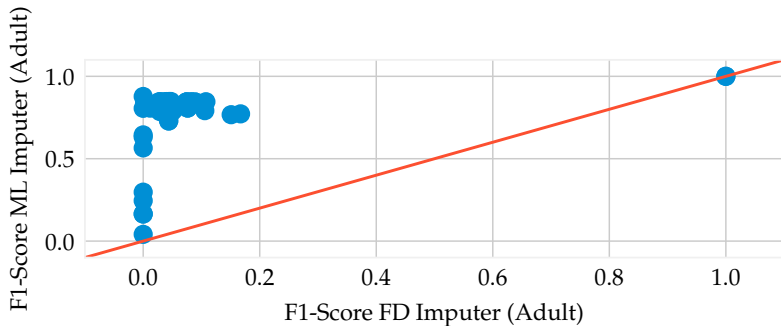
In this example, FD Imputer uses the FD $C \rightarrow D$.

Benchmarking FD Imputer

Dataset	#FDs _{train}	#FD (F1 = 0)	F1 _{mean}	F1 _{max}
Abalone	193	45	0.0008	0.0048
Adult	88	10	0.0669	1.0000
Balance S.	7	6	0.0000	0.0000
Breast C. W.	77	10	0.2198	0.7539
Chess	9	8	0.0000	0.0000
Iris	8	1	0.1274	0.2252
Letter	80	17	0.2347	0.3737
Nursery	11	10	0.0000	0.0000

Table: Performance of the FD Imputer on a selection of UCI datasets.

Benchmarking FD Imputer



Benchmarking FD Imputer

- Some FDs are more robust than others
- FD Imputer performs generally worse than the model trained with DataWig
- This concerns only classifiable data – it is generally impossible to impute continuous numerical data with FD Imputer

FD Imputer is Overfitting

Haykin 2008

“[Overfitting] is essentially a ‘look-up table’, which implies that the input-output mapping [...] is not smooth.”⁴

⁴Haykin 2008, p. 165.

FD Imputer is Overfitting

Haykin 2008

“[Overfitting] is essentially a ‘look-up table’, which implies that the input-output mapping [...] is not smooth.”⁴

- Due to the implementation of FD Imputer, the train-set is merely a table to look up imputation values
- No generalization takes place whatsoever
- No empirical risk minimization (ERM) is applied!

⁴Haykin 2008, p. 165.

Deriving FDs from Trained Models

When overfitting DataWig models, one cannot ensure that overfitting takes place on FD-attributes. Thus, it does not appear to be possible to derive FDs with trained models based on empirical risk minimization (ERM).

- It does not appear to be possible to derive FDs with ERM-based imputation models
- But what about Relaxed Functional Dependencies (RFDs)?

DepDetector

RFDs as Models

- An RFD is based on the definition of an FD
- It alters that definition to serve a specific purpose
- There are many RFDs defined, such as Metric Functional Dependencies, Conditional Functional Dependencies or Approximate Functional Dependencies.

Example of an RFD

Example

Koudas et al. introduce Metric Functional Dependencies (MFDs) to find constraints on tables that contain rows with slightly different formatting or slightly deviating values.⁵

⁵Koudas et al. 2009.

Example for Noisy Data

Id	First name	Last name	Zip	Town
1	Alice	Smith	19139	Munich
2	Peter	Meyer	19139	Muinch
3	Ana	Parker	19139	Munich
4	John	Pick	12055	Berlin

The FD Zip \rightarrow Town is violated by the second row. A MFD that considers the Levenshtein distance between two entries rather than exact equality can still detect the constraint.

RFDs as Models

- Due to their training being ERM-based, DataWig models are capable of learning constraints on noisy data as well
- During training, DataWig models “learn” relaxations – there is no need to manually set up a threshold of a Levenshtein distance as in the previous example

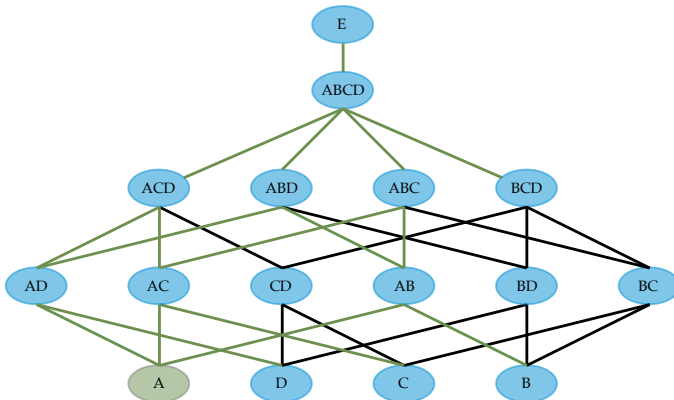
RFDs as Models

- Due to their training being ERM-based, DataWig models are capable of learning constraints on noisy data as well
- During training, DataWig models “learn” relaxations – there is no need to manually set up a threshold of a Levenshtein distance as in the previous example

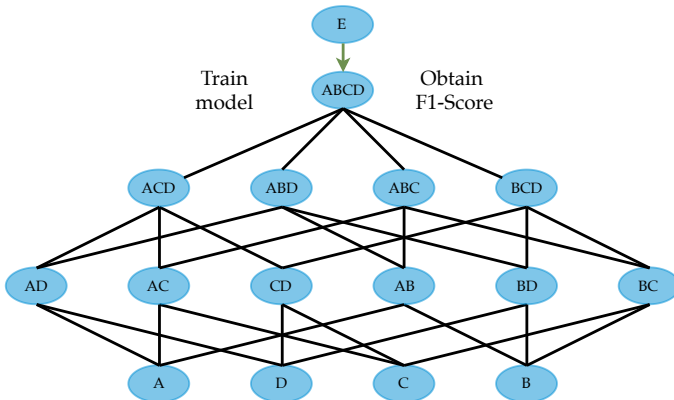
Question

How does one find a *minimal* dependency with this approach?

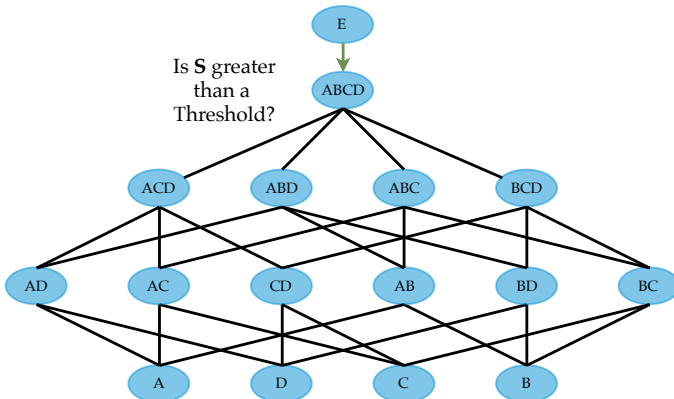
Functionality of DepDetector



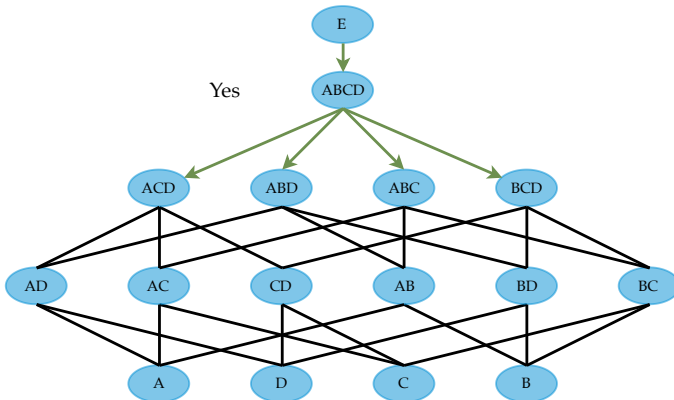
Functionality of DepDetector



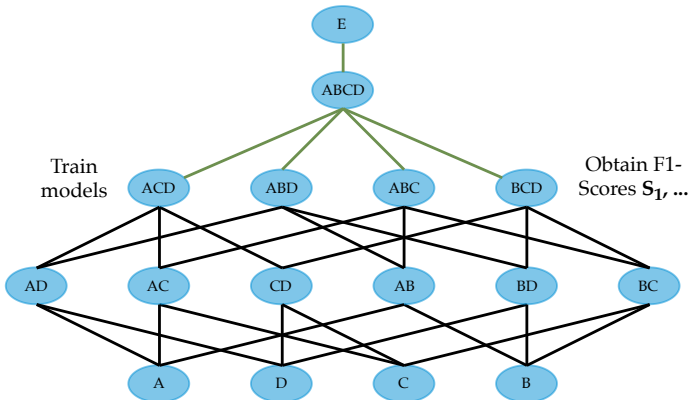
Functionality of DepDetector



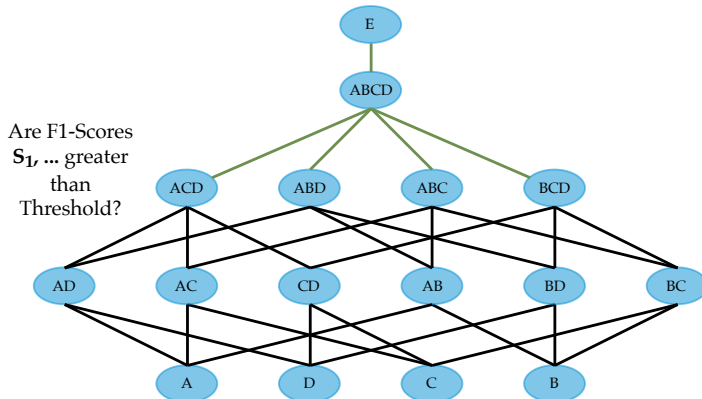
Functionality of DepDetector



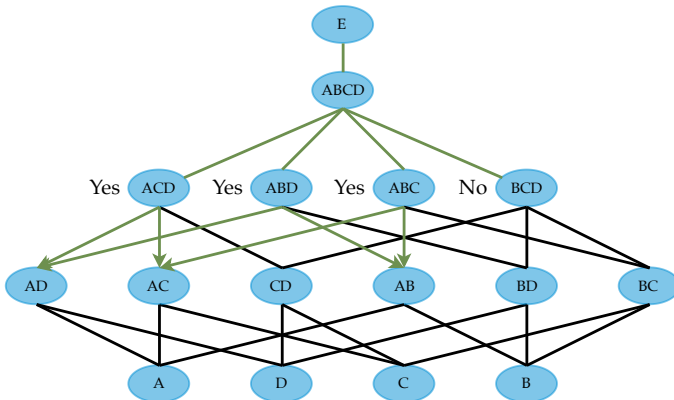
Functionality of DepDetector



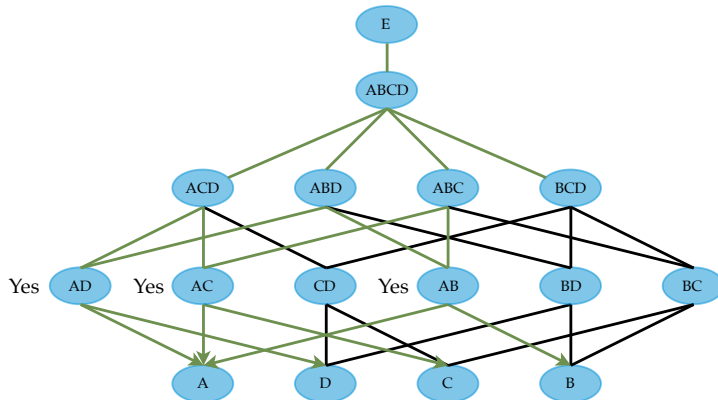
Functionality of DepDetector



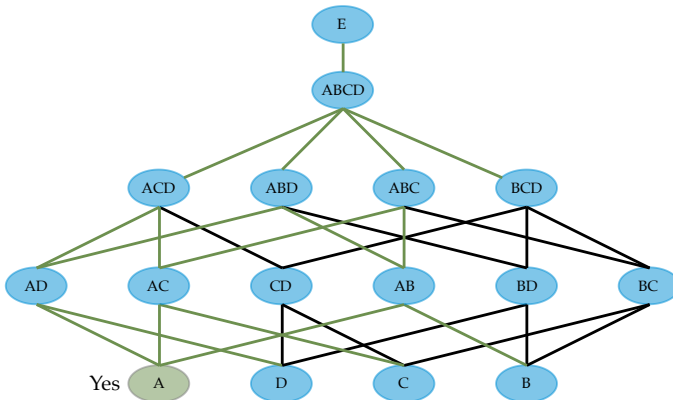
Functionality of DepDetector



Functionality of DepDetector



Functionality of DepDetector



DepDetector Properties

- DepDetector solves an optimization problem on a directed graph, whereas in FD detection, tuple-comparison is performed
- DepDetector depends on one threshold for classifiable data and continuous numerical data respectively
- Machine learning classifier/regressor models have the potential to unify RFDs

DepDetector Dependency Results

Dataset	Cols	Rows	# FDs	Greedy dependencies	Complete dependencies
Abalone	10	4177	175	7 (42 min)	TL
Adult	16	32561	93	TL	TL
Balance-S.	6	625	7	3 (67 s)	3 (80 s)
Chess	8	28056	9	1 (117 min)	1 (340 min)
Iris	6	150	9	5 (38 s)	8 (43s)
Letter	18	20000	78	TL	TL
Nursery	11	12960	11	3 (110 min)	TL

‘TL’ indicates a time limit of 350 min.

Prospects

Prospects for further Research

- Minimal dependency detection algorithms for learned relaxations can be further optimized (Concurrency, applying Graph-Theory)
- A complexity analysis for the for such algorithms can be performed
- Approaches to calculate thresholds from the data can be introduced

Thank you for your attention!

References I

Ziawasch Abedjan et al. *Data Profiling*. 2019. isbn: 9781681734477. doi: <https://doi.org/10.2200/S00878ED1V01Y201810DTM052>.

Simon Haykin. *Neural Networks and Learning Machines Third Edition*. Pearson Prentice Hall, 2008. isbn: 9780131471399.

References II

Alireza Heidari et al. “HoloDetect: Few-Shot Learning for Error Detection”. In: *Proceedings of the 2019 International Conference on Management of Data*. SIGMOD '19. ACM, 2019, pp. 829–846. isbn: 978-1-4503-5643-5. doi: 10.1145/3299869.3319888. url: <http://doi.acm.org/10.1145/3299869.3319888>.

N. Koudas et al. “Metric Functional Dependencies”. In: (Mar. 2009), pp. 1275–1278. issn: 1063-6382. doi: 10.1109/ICDE.2009.219.

References III

Tim Kraska et al. “The Case for Learned Index Structures”. In: *Proceedings of the 2018 International Conference on Management of Data*. SIGMOD '18. Houston, TX, USA: ACM, 2018, pp. 489–504. isbn: 978-1-4503-4703-7. doi: 10.1145/3183713.3196909. url: <http://doi.acm.org/10.1145/3183713.3196909>.

Thorsten Papenbrock and Felix Naumann. “A Hybrid Approach to Functional Dependency Discovery”. In: SIGMOD '16 (2016), pp. 821–833. doi: 10.1145/2882903.2915203. url: <http://doi.acm.org/10.1145/2882903.2915203>.

Benchmarking FD Imputer Continuous Data

Dataset	# cFDs _{train}	# 0-Coverage cFDs	Coverage (%)
Abalone	139	84	0.1277
Adult	11	5	0.1217
Balance S.	1	1	0.0000
Breast C. W.	1	1	0.0000
Chess	1	1	0.0000
Iris	4	4	0.0000
Letter	0	0	-
Nursery	1	1	0.0000

Table: Imputation coverage of FD Imputer on all UCI datasets for which FDs with continuous data in the RHS were detected.

FD Imputer Continuous Data Definitions

$$\text{mean missing values per cFD} = \sum_i \frac{\text{missing imputations}}{\text{cFD}_i} \cdot (\# \text{ cFDs})^{-1}$$
$$\text{mean coverage} = \left(1 - \frac{\text{mean missing values per cFD}}{\# \text{ rows in } r_{\text{test}}} \right) \cdot 100$$

Machine Learning Models

data representation



model class



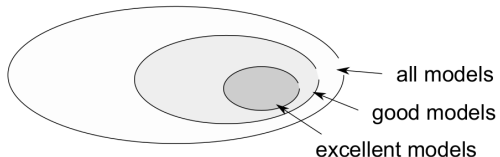
performance measure



optimization



validation



⁵Image from Prof. Obermeyer, Neural Information Processing Group TU Berlin

Definition Robustness Continuous Data

If A contains continuous numerical data, the MSE is determined to measure robustness:

$$\text{robustness}_{\text{MSE}} = \frac{1}{p - m} \sum_{i=m+1}^p (t_i^*[A] - t'_i[A])^2. \quad (1)$$

Here, $p \in \mathbb{N}$ denotes the number of tuples in a relational instance and m is the number of imputed tuples.

Definition Robustness Classifiable Data

If A contains classifiable data, the F1-Score is calculated to measure robustness:

$$\text{robustness}_{\text{F1-Score}} = \left(\frac{\text{Recall}(r_{\text{imp}}, r_{\text{test}})^{-1} + \text{Precision}(r_{\text{imp}}, r_{\text{test}})^{-1}}{2} \right)^{-1} \quad (2)$$