



BEUTH HOCHSCHULE FÜR TECHNIK BERLIN
University of Applied Sciences

Master's Thesis

Robust Machine-Learning Approaches for Efficient Functional Dependency Approximation

by

Philipp Jung

Business Administration and Engineering - Project Management M.A.

Fachbereich I

Philipp Jung
Matriculation Number: 872855
01.10.2019

Advisors:
Prof. Dr. Biessmann
Prof. Dr. Schoeneberg

Abstract

In this thesis the applicability of empirical risk minimization techniques for approximate functional dependency detection is explored. A measure called ‘robustness’ is introduced to assess the ability of a functional dependency when used for value imputation. Robustness is measured for functional dependencies from a selection of datasets and functional dependencies are ranked by robustness. Imputations derived from functional dependencies are compared to imputations derived from a model trained using the DataWig framework. It is shown that trained imputation models share some properties with two relaxed functional dependencies. Finally, an algorithm for relaxed functional dependency detection based on learned imputation models is presented.

Contents

1	Introduction	1
2	Theory	2
2.1	Introduction of Functional Dependencies	2
2.2	Additional Relational Database Theory	4
2.3	Database Normalization	5
2.4	Relaxed Functional Dependencies	7
2.4.1	FDs Relaxing on the Extent	8
2.4.2	FDs Relaxing on the Attribute Comparison	9
2.5	Model Performance Measures	10
2.5.1	Classification Performance Measure	11
2.5.2	Regression Performance Measure	13
2.6	Robustness of FDs	13
2.7	Datawig Imputer	15
2.8	FD Imputer: Measuring Robustness	17
2.9	ML Imputer: Dependency Based Imputation	20
2.10	DepDetector: RFD Discovery	20
2.10.1	DepDetector Algorithm	21
3	Experiments	23
3.1	FD Imputer	23
3.1.1	Imputation of Classifiable Data	24
3.1.2	Imputation of Continuous Data	25
3.2	ML Imputer	27
3.2.1	Imputation of Classifiable Data	28
3.2.2	Imputation of Continuous Data	29
3.3	Overfitting ML Imputer	30
3.4	Dependency detection	33
3.4.1	DepDetector Results on Various Datasets	33
3.4.2	Dependency stability	34
4	Discussion	36
4.1	Robustness	36
4.2	Comparing ML Imputer with FD Imputer	36
4.2.1	Imputing Sequential RHSs	36
4.2.2	Imputing Classifiable RHSs	37
4.2.3	Overfitting ML Imputer	38
4.3	Dependency Detection with DepDetector	38
5	Conclusion	39

References	40
6 Appendix	42

List of Figures

1	Venn diagram representing the extent relaxation criterion.	8
2	Illustration of a binary confusion matrix. ‘Prediction’ refers to predicted labels $y_{pred}(x)$ while ‘Ground Truth’ represents the actual labels $y(x)$. . .	11
3	Each predicted label y_x is represented by a circle. Hollow circles stand for negative labels and full circles for positive labels.	12
4	State diagram showing the histogram established by Datawig Imputer to encode example 18	16
5	Hasse diagram of the search lattice when discovering minimal, non-trivial FDs on a relational schema with five columns.	20
6	F1-Scores of the seven most-robust FDs on the Adult dataset when imputing classifiable values.	23
7	Robustness of the seven top-performing FDs on the Letter dataset when imputing classifiable RHSs.	24
8	The figure compares the F1-Score of the FD Imputer compared to the F1-Score of the ML Imputer. Each point represents one FD.	28
9	F1-Scores of the seven top-performing FDs on the Adult dataset.	29
10	Comparison of an overfitted model with a conventionally trained model on the Adult dataset.	31
11	Overtrained models are compared with conventionally trained models on the Abalone dataset. Continuous RHSs were imputed to get MSEs.	31
12	Overfitted imputation-models are trained with FD $[2,3,6,7] \rightarrow 5$ on the Abalone dataset with varying dataset-size.	32
13	Analysis on the Iris dataset reveals that results become continuously more stable for larger τ	34
14	On the dataset Balance-scale, the analysis of the result shows a peak at $\tau = 4$	35
15	The figure compares the F1-Score of the FD Imputer compared to the F1-Score of ML Imputer. Each point represents one FD.	37
16	Mean robustness of all FDs with a classifiable RHS in function of the total number of FDs with a classifiable RHS divided by the total number of attributes in a dataset. No functional relation between these two values can be stated.	43

List of Tables

1	Example of an FD.	3
2	The compound attributes ADDRESS and NAME can be split into their atomic components TOWN and STREET as well as FIRST NAME and SURNAME, respectively.	6
3	Even though column ZIP determines the content of column TOWN, an FD is not capable of displaying this fact due a typing error in column TOWN.	8
4	The relation r is split into train- and test sets with a split-ratio $s = 0.5$	18
5	Performance of the FD Imputer on a selection of UCI datasets.	25
6	Imputation coverage of FD Imputer on all UCI datasets for which FDs with continuous data in the RHS were detected.	26
7	Statistical analysis of the computed MSEs when running FD Imputer with 84 different cFDs on the Abalone dataset.	27
8	Performance of ML Imputer on eight different datasets when training on FDs with classifiable RHS entries.	29
9	Statistical analysis of the Abalone dataset obtained by running ML Imputer with the same 84 cFDs as used in the generation of table 7.	30
10	Result of running dependency detection on selected datasets. Values in brackets are the runtime until convergence. 'TL' indicates that a time limit of 350 min is reached.	33

1 Introduction

Data-driven methods change the way computer scientists approach algorithmic problems. Rather than designing and implementing complex algorithms themselves, recent advances in machine learning have allowed for learned algorithms.

Kraska et al. showed in their 2018 publication “The case for Learned Index Structures” that different index structures can be replaced by learned ones, greatly improving performance. [Kra+18] In the field of data cleaning and data enrichment, HoloClean lead the way for machine-learning approaches in the domain of data cleaning. [Hei+19] HoloClean is agnostic of the way data is structured, making it versatile for many different domains of application.

In this work, machine-learning techniques are applied to the field of relational database theory — more precisely, functional dependency detection. Stemming from the beginnings of relational database theory, functional dependencies were introduced to formalize normalization of relational schemata.

In the theoretical part of this thesis, the basic relational database terminology is introduced. Furthermore, limitations of canonical functional dependencies are mentioned and relaxed functional dependencies are introduced. With reference to Koudas et al., functional dependencies’ robustness is discussed. A method for measuring robustness is proposed. Mean square error and F1-Score are introduced to measure the performance of imputation models.

Machine-learning classification theory necessary for understanding the basic functionality of the Datawig framework [Bie+18] is discussed. It is described how robustness of FDs is measured using an algorithm called FD Imputer. FD Imputer’s machine-learning counterpart, called ML Imputer, is presented. The theoretical part of this thesis is concluded with a discussion of how to detect dependencies using learned imputation models using an algorithm called DepDetector.

Several experiments are conducted to explore machine-learning techniques applied to functional dependency analysis. FD Imputer is run on a range of datasets and FDs are ranked according to their robustness. The experiment is repeated with ML Imputer. ML Imputer is overfitted and implications of overfitting on the FD robustness-ranking are examined. The final building block of the experimental section is formed by dependency detection with DepDetector.

Results are discussed and recommendations for future research in learned dependency detection are issued in the final section of this work.

2 Theory

2.1 Introduction of Functional Dependencies

FDs are a way of expressing “a priori knowledge of restrictions or constraints on permissible sets of data”. [Mai83, p. 42] To be able to define FDs, they must be placed in the context of the domain from which they originate: relational database theory.

Definition (Relation Scheme, Attribute Names, Domain, Relation, Tuples) A *relation scheme*¹ R is a finite set of *attribute names* $R = \{A_1, A_2, \dots, A_n\}$, where to each attribute name A_i corresponds a set D_i , called *domain* of A_i , $1 \leq i \leq n$. Let $D = D_1 \cup D_2 \cup \dots \cup D_n$, then a *relation* r on relation scheme R is a finite set of mappings $\{t_1, t_2, \dots, t_p\}$ from R to D :

$$t_i : R \rightarrow D$$

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

$$t(A_i) \subseteq D_i.$$

In application, attribute names are commonly called *column name* or *column attribute*. One can think of them as labels of columns in which data is stored.

Definition (Functional Dependency) Consider a relation r on scheme R with subset $X \subseteq R$ and a single attribute $A_i \in R$. A FD $X \rightarrow A$ is said to be *valid* in r , if and only if

$$t_i[X] = t_j[X] \Rightarrow t_i[A] = t_j[A] \quad (1)$$

holds for 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 write $X \rightarrow A$. X is called the *left hand side* (LHS) of an FD, whilst A is called the *right hand side* (RHS) of the same FD. A FD $X \rightarrow Y$ is called *trivial*, if $Y \subseteq X$. [Stu16, p. 163]

Example Considering table 1, one can see that every tuple in the left hand side subset of the relation uniquely determines the right hand side. We say that ID, FIRST NAME, SURNAME, TOWN *functionally determines* ZIP or write $\{\text{ID, FIRST NAME, SURNAME, TOWN}\} \rightarrow \text{ZIP}$. [Mai83, p. 43]

If inspected closely, one can discover even more FDs in table 1. For example, $\text{TOWN} \rightarrow \text{ZIP}$ and $\text{ID} \rightarrow \text{ZIP}$. Since TOWN and ID are subsets of $\{\text{ID, FIRST NAME, SURNAME, TOWN}\}$, we call the FD $\{\text{ID, FIRST NAME, SURNAME, TOWN}\} \rightarrow \text{ZIP}$ *non-minimal*.

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

left hand side				right hand side
ID	FIRST NAME	SURNAME	TOWN	ZIP
1	Alice	Smith	Munich	19139
2	Peter	Meyer	Munich	19139
3	Ana	Parker	Munich	19139
4	John	Pick	Berlin	12055
5	John	Pick	Munich	19139

Table 1: Example of an FD.

Definition (Minimal FD) A FD $X \rightarrow A$ is *minimal*, if no subset of X functionally determines A . [Pap+15, p. 2] Thus, $ID \rightarrow ZIP$ and $TOWN \rightarrow ZIP$ are *minimal FDs*.

Definition (Logical Implication) Let F be the set of FDs on R . $X \rightarrow Y$ is *logically implied* by F , or algebraically expressed

$$F \models X \rightarrow Y, \quad (2)$$

if every relational instance r of R , which satisfies all dependencies in F , also satisfies $X \rightarrow Y$. [Stu16, p. 166]

Thus, if $F = \{A \rightarrow B, A \rightarrow C, BC \rightarrow D\}$, the following logical implications are true:

$$\begin{aligned} F &\models A \rightarrow B \\ F &\models A \rightarrow BC \\ F &\models A \rightarrow D \end{aligned}$$

Definition (Transitive Closure) Let F be a set of FDs. The *transitive closure* F^+ of F is defined by

$$F^+ := \{X \rightarrow Y \mid F \models X \rightarrow Y\} \quad (3)$$

Definition (Full Dependency, Determinant) Let R be a relation scheme and let $Y \in R$ be an attribute on R . Furthermore, let $X \subseteq R$ be a set of attributes on R . The attribute Y is called *fully dependent* on X , if Y is functionally determined by X and X is minimal. [Sch17, p. 61] We call a *determinant* a set of attributes that fully functionally determines another attribute.

2.2 Additional Relational Database Theory

When real-world data used by an application is stored on a machine according to the relational model, it is usually stored in a relational database. *Database normalization* is the original domain of application for FDs. [Cod70, p. 381] In the following section, concepts necessary to introduce databases are defined. In addition, definitions of terms used in database normalization are given.

Definition (Superkey) Let r be a relation and let $R = \{A_1, A_2, \dots, A_n\}$, $n \in \mathbb{N}$, be a relation scheme on which r is defined. Let K be a subset of R , such that $K = \{A_1, A_2, \dots, A_m\} \subseteq R$, where $m \leq n$, $m \in \mathbb{N}$. The subset K is called *superkey*, if for any tuple $t_i \in r$ the relation

$$t_i(A_k) = t_j(A_k) \Rightarrow t_i \equiv t_j$$

holds for any single $A_k \in K$. [Mai83, p. 4] In other words, if K is a superkey, any K -value of a tuple identifies that tuple uniquely. [Sch17, p. 32]

Definition (Candidate Key) A superkey K is called *candidate key*, if it is *minimal*. [Sch17, p. 32] The notion of K being minimal means that K is no longer a superkey, if any attribute $A_k \in K$ is removed from K .

Definition (Prime Attributes) An attribute A on a relational scheme R is called *prime*, if A is part of a key of R . Otherwise, A is called *non-prime*.

Definition (Relational Database) Following the definition of a relation scheme R , one can formally introduce databases and database schemes:

We assume that R is composed of two parts, S and K . We call S a *set of attributes* and K a *set of designated keys* and describe this composition by writing $R = (S, K)$. A *relational database scheme* \mathcal{R} over \mathbf{U} can now be defined as a collection of relation schemes $\mathcal{R} = \{R_1, R_2, \dots, R_p\}$, where $R_i = (S_i, K_i)$, $1 \leq i, j \leq p$,

$$\mathbf{U} := \bigcup_{i=1}^p S_i$$

We demand that $S_i \neq S_j$ if $i \neq j$.

A *relational database* d on a database scheme \mathcal{R} is a collection of relations $d = \{r_1, r_2, \dots, r_p\}$ such that for each relation scheme $R = (S, K)$ in \mathcal{R} there is a relation r in d such that r is a relation on S that satisfies every *key* in K . [Mai83, p. 94]

2.3 Database Normalization

In 1970, Edgar F. Codd defined the *relational model* [Cod70] and pioneered many important concepts in database theory. When introducing the relational database model in his 1970 article “A relational model of data for large shared data banks”, Codd formalized database normalization alongside. [Cod70]

Describing what will be known to academia as *First Normal Form* (1NF), Codd states that “problems treated [when normalizing databases] are those of *data independence*”, aiming to protect future users of large databases “from having to know how the data is organized in the machine”. [Cod70, p. 1]

Databases at the time were structured hierarchically or navigationally. This design was centered on efficiency, optimized for handling queries as fast as possible. While this design yielded good performance in an era when computing time was cost-intensive, it came with a heavy cost of complexity: “Teams of programmers were needed to express queries to extract meaningful information. [...] Such databases [...] were absolutely inflexible[y]”. [IBM03]

Codd’s relational model shifted the focus of database architecture away from efficiency towards a new design, centered around the user of a database. However, this new approach of database architecture came with new challenges. How does one create a database-scheme with favorable properties?

Generally, redundantly stored data is an indicator [Wat14, p. 61] that so-called *anomalies* might occur, leading to inconsistency when operating the database. [Stu16, p. 162] Kleuker writes that “with a certain ‘database-sense’ ” one could sense when splitting tables up is necessary to prevent anomalies. [Kle11, p. 76]

This ‘database-sense’ was formalized when Codd introduced database normalization as part of the relational model. [Cod70, p. 381] Different *normal forms* were defined. Codd used three normal forms, which he called First-, Second- and Third Normal Form.²

Subsequently, Fourth-, Fifth- and even higher-order Normal Forms were defined by other authors. However, these normal forms find little reception in real-world applications. [Sch17, p. 58] Today, FDs are primarily used in database normalization. [CDP16, p. 1] Thus, the derivation of the most commonly used normal-forms is of interest within the scope of this work.

²The Third Normal Form was later newly formulated by Boyce and Codd and named Boyce-Codd Normal Form (BCNF). Due to its more elegant definition, BCNF is introduced rather than Codd’s original 3NF.

Definition (First Normal Form) A relation scheme R is in *First Normal Form* (1NF), if values in $\text{dom}(A)$ are atomic for every attribute $A \in R$. [Mai83, p. 96] Likewise, a database scheme \mathbf{R} is in 1NF if every relation scheme $R \in \mathbf{R}$ is in 1NF.

Example Consider table 2 which represents two relational instances on two different relational schemes. It serves as an example of what is called atomic- and compound data in the Relational Database model. [Cod90, p. 6]

Compound Scheme			Atomic Scheme			
NAME	ADDRESS		FIRST NAME	SURNAME	TOWN	STREET
1	Alice Smith	Munich, Flurstr.	Alice	Smith	Munich	Flurstr.
2	Peter Smith	Munich, Flurstr.	Peter	Smith	Munich	Flurstr.
3	Ana Parker	Munich, Anastr.	Ana	Parker	Munich	Anastr.
4	John Pick	Berlin, Flurstr.	John	Pick	Berlin	Flurstr.

Table 2: The compound attributes ADDRESS and NAME can be split into their atomic components TOWN and STREET as well as FIRST NAME and SURNAME, respectively.

Note that the compound scheme's attributes can be decomposed into several other attributes. The atomic scheme's attributes however cannot be further split into any meaningful smaller components.

1NF is the very foundation of the Relational Model. It demands, that the only type of compound data is the relation itself. [Cod90, p. 6]

Definition (Second Normal Form) A relation scheme R is said to be in *Second Normal Form* (2NF) in respect to a set of FDs F , if it is in 1NF and every nonprime attribute is fully dependent on every key of R . [Mai83, p. 99] This definition can be extended for databases: A database scheme \mathbf{R} is in 2NF with respect to F if every relation scheme $R \in \mathbf{R}$ is in 2NF with respect to F .

It can be shown that a database scheme in 2NF is also in 1NF. [Sch17, p. 58]

Definition (Boyce-Codd Normal Form) A relation scheme R is in *Boyce-Codd Normal Form* (BCNF), if every determinant in R is a candidate key. [Sch17, p. 65] Equally, a

relational database scheme \mathbf{R} is in BCNF if every relation scheme $R \in \mathbf{R}$ is in BCNF.

If a database scheme is in BCNF, it is also in 2NF. [Sch17, p. 58] A database in BCNF effectively eliminates redundancies apart from candidate keys, freeing a database in BCNF from anomalies. [Sch17, p. 67]

2.4 Relaxed Functional Dependencies

When Edgar F. Codd proposed the relational model in 1970, the concept of keys and thus the concept of FDs were already present. [Mai83, p. 70] Pursuing the goal of database normalization, Codd separated the definition of an FD from keys in 1972.

In the ensuing time, researchers abstracted the concept of FDs and transferred the approach. FDs were used in research to solve problems other than schema normalization. [CDP16, p. 161] Many of these use-cases, such as ‘approximate classification’ or ‘source merging’, arose due to the fact that real-world datasets are almost always *noisy*.

Noise in databases has many faces: Entries might be corrupted by missing data, wrongly entered data or incomplete data. A database might have been merged with identical entries formatted in different ways. [Kou+09, p. 1] In these cases, functionally dependent column-combinations are not detected as such by an FD detection algorithm searching for FDs as defined in 1. This may result in misleading insights when searching for FDs or normalizing a database.

Taking these limitations of *canonical FDs*³ into account, a multitude of new dependencies was defined, “aiming to solve specific problems”. [CDP16, p. 147] Since all of these new kinds of dependencies *relax* the conditions in equation 1, they are called *relaxed functional dependencies* (RFDs).

Caruccio et al. classified and compared 35 different RFDs. [CDP16, p. 151] They state that each RFD was “based on its underlying relaxation criteria”. [CDP16, p. 149] They proceed to define two relaxation criteria by which they distinguish RFDs.

The first relaxation criterion is called *attribute comparison*. It refers to “the type of tuple comparison used on the LHS and RHS” in equation 1.

The other relaxation criterion *extent* implies a relaxation of the set of tuples for which an FD is valid.

³To accentuate the difference between FD and RFD, FDs are also called canonical FDs. This notion refers to FDs as defined in equation 1.

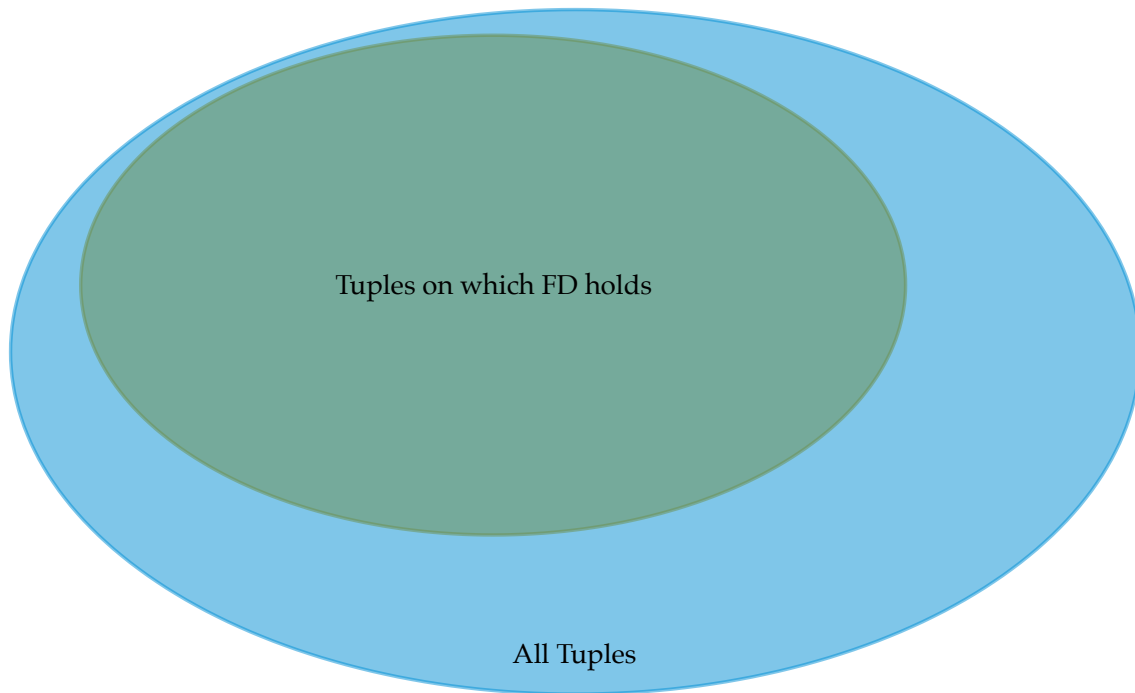


Figure 1: Venn diagram representing the extent relaxation criterion.

2.4.1 FDs Relaxing on the Extent

The term *extent* is used to describe how an RFD relaxes on the subset of tuples for which the RFD is satisfied. While the definition of an FD demands that the condition of an FD must be valid for all tuples in a relational table, RFDs that relax on the extent only hold on a subset of all tuples.

This idea is implemented by different RFDs in different ways. In the following section, a selection of RFDs relaxing the extent of the FD is briefly presented.

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

Table 3: Even though column ZIP determines the content of column TOWN, an FD is not capable of displaying this fact due a typing error in column TOWN.

Approximate Functional Dependency (AFD) Applicability of FDs is improved by AFDs by “holding on *almost*⁴ every tuple” [CDP16, p. 151]. To illustrate this, table 3 shows an example of noisy data. The potential FD $\text{TOWN} \rightarrow \text{ZIP}$ is not captured by the definition of a canonical FD. Due to a typing error, the potential FD is invalidated in the row where $\text{ID} = 2$. To still capture the relation, a different dependency-measure than given in the definition of the canonical FD is required.

Tuples that do not correspond to the canonical FD are measured as fraction of the total of tuples on a relation r as follows:

$$\Psi(X, Y) = \frac{\min(|r_1| \mid r_1 \subseteq r \text{ and } X \rightarrow Y \text{ hold in } r \setminus r_1)}{|r|} \quad (4)$$

Here, function $\Psi(X, Y)$ is called *coverage measure* of an RFD $X \rightarrow Y$. Ψ “quantifies the satisfiability degree of an RFD [...] on r ” [CDP16, p. 150] and is used in the definition of an AFD to be compared to a threshold $\epsilon \in [0, 1]$.

If $\Psi(X, Y)$ is smaller or equal to ϵ , the AFD is said to hold on a relation r . Applied to table 3, the AFD $\text{TOWN} \rightarrow \text{ZIP}$ holds, if $\epsilon \geq 0.25$.

Conditional Functional Dependencies (CFDs) *Conditional Functional Dependencies* employ conditions to define the subset on which a dependency holds. Originally, those conditions exclusively allow the definition of constraints using the equality operator. [CDP16, p. 152] Applied to table 3, a CFD $\text{TOWN} \rightarrow \text{ZIP}$ holds under the condition that a) entries in column $\text{TOWN} = \text{Munich}$ or b) entries in column $\text{TOWN} = \text{Berlin}$.

Other RFDs based on the definition of CFDs include *extended conditional functional dependencies* (ECFDs) by Bravo et al. that allow the disjunction-operator as well as the inequality-operator for defining conditions. [Bra+08]

Also, Chen et al. defined CFD^ps, which introduce $<, \leq, >, \geq$ and \neq for defining conditions. [CFM09]

2.4.2 FDs Relaxing on the Attribute Comparison

Instead of specifying subsets of tuples for which an FD is valid, FDs relaxing on the attribute comparison alter the condition under which a dependency is said to be valid.

⁴Highlighting by the author.

One common RFD relaxing on the attribute comparison is the metric functional dependency.

Metric Functional Dependencies (MFDs) First introduced by Koudas et al. in 2009, *Metric Functional Dependencies* were defined to address violations of canonical FDs due to “small variations [...] in data format and interpretation” [Kou+09, p. 1].

According to equation 1, an FD $ZIP \rightarrow TOWN$ is said to be valid if

$$t[ZIP] = t'[ZIP] \Rightarrow t[TOWN] = t'[TOWN] \quad (5)$$

holds for all pairs of distinct tuples $t_i, t_j \in r$, where r is a relation on scheme R . The definition of an MFD replaces this equation by demanding that if and only if for all pairs of distinct tuples $t_i, t_j \in r$ the metric condition

$$t[ZIP] = t'[ZIP] \Rightarrow d(t[TOWN], t'[TOWN]) \leq \delta \quad (6)$$

holds, the MFD $ZIP \rightarrow TOWN$ is said to be valid on r . [Kou+09, p. 2] Here, δ is called a *tolerance parameter*, indicating how far the result of a *metric* $d(t[Y], t'[Y])$ can deviate from exact equality.

One such metric d would be the absolute difference between two values, such that $d(t[Y], t'[Y]) = |t[Y] - t'[Y]|$. A popular metric used to measure the difference between two sequences is the Levenshtein distance, where $d(t[Y], t'[Y]) = \text{lev}_{t,t'}(|t|, |t'|)$ and

$$\text{lev}_{t,t'} = \begin{cases} \max(i, j), & \text{if } \min(i, j) = 0, \\ \min \begin{cases} \text{lev}_{a,b}(i-1, j) + 1 \\ \text{lev}_{a,b}(i, j-1) + 1 \\ \text{lev}_{a,b}(i-1, j-1) + 1_{(a_i \neq b_j)} \end{cases} & \text{otherwise.} \end{cases}$$

Regarding the example given in table 3, the Levenshtein distance between Munich and Muinch is 2. Thus the MFD $ZIP \rightarrow TOWN$ is valid, if $d = \text{lev}$ and $\delta = 2$.

2.5 Model Performance Measures

In machine learning, models are tested in order to determine whether they accurately predict new data points. [Bis06, p. 32] This is usually done by measuring the model's performance on a separate dataset not involved in training, the so called test set X_{test} . [Bis06, p. 6] Depending on the kind of task a model is trained to perform, different metrics need to be employed.

2.5.1 Classification Performance Measure

To visualize the performance of a classifier, a *confusion matrix* can be created. [Tha18, p. 2] The simplest case of a confusion matrix can be constructed when measuring the performance of a binary classifier, that is a classifier that separates data points into two classes. Figure 2 shows such a 2×2 binary confusion matrix. Here, ‘Ground Truth’

Prediction	Ground Truth	
	Positive	Negative
Positive	True Positive	False Positive
Negative	False Negative	True Negative

Figure 2: Illustration of a binary confusion matrix. ‘Prediction’ refers to predicted labels $y_{pred}(x)$ while ‘Ground Truth’ represents the actual labels $y(x)$.

describes the label $y(x)$ of some data point $x \in X_{test}$, where $y \in \{0,1\}$. ‘Prediction’ identifies the predicted labels $y_{pred}(x)$ that the model generates after it has been executed on the test-set X_{test} .

Whenever $y_{pred}(x) = y(x)$, $x \in X_{test}$ holds, the predicted label can be assigned to be either a *True Positive* (TP) or a *True Negative* (TN). The opposite holds as well, such that a falsely predicted label will be either a *False Negative* (FN) or a *False Positive* (FP). [Tha18, p. 2]

Using the classification introduced by the binary confusion matrix, all predicted labels y_{pred} are assigned to the four sets TP, TN, FN and FP. Using these four sets, we can introduce measures for classification performance.

Precision is a measure that depicts the proportion of correctly classified positive samples to the total amount of samples classified as positive. [Tha18, p. 4] This can be formulated

as

$$Precision = \frac{|TP|}{|TP| + |FP|} \quad (7)$$

where $|A|$ represents the cardinality of a set A . Precision measures how many elements classified as positive are True Positives.

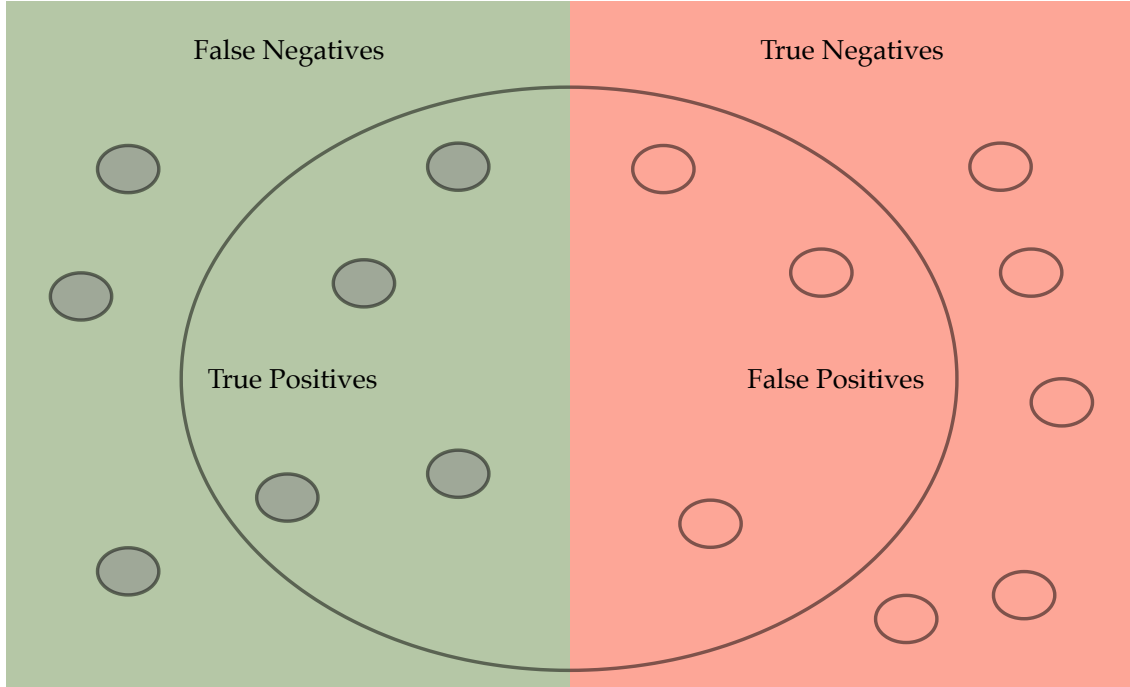


Figure 3: Each predicted label y_x is represented by a circle. Hollow circles stand for negative labels and full circles for positive labels.

Recall, also called *sensitivity*, represents the share of correctly classified positive samples to the total amount of positive samples. [Tha18, p. 3] This can be formalized as

$$Recall = \frac{|TP|}{|TP| + |FN|} \quad (8)$$

Recall measures how many of the elements labelled as positives were actually selected.

The harmonic mean of precision and recall is called *F1-Score* or *F1-Measure*:

$$F1\text{-Score} = \left(\frac{Recall^{-1} + Precision^{-1}}{2} \right)^{-1} \quad (9)$$

Selecting models based on the highest F1-Score thus maximizes both precision and recall, which is a desired property of classifiers.

The F1-Score was derived under the name of *MUC-4* by Nancy Chinchor in 1992. [Chi92] Chinchor based her work on the book 'Information Retrieval' by Van Rijsbergen from 1979. [Rij79]

2.5.2 Regression Performance Measure

For regression problems, the distance between a predicted label $y_{pred}(x)$ and the actual label $y(x)$ of a data-point specifies the performance. In supervised learning, the *Mean Square Error* (MSE) or the *sum of squared errors* over the training sample are commonly used as performance measure. [Hay08, p. 35]

In this work, the MSE is used as performance measure for models that perform regression tasks. This choice is motivated by favorable properties of the MSE. For the one hand, the MSE can be easily computed.

Let $Y \in \mathbb{R}^n$ be the n -dimensional vector of imputed values and $\hat{Y} \in \mathbb{R}^n$ be the vector of actual RHS values. Then the MSE is defined as [Moo+11, p. 597]

$$\text{MSE} = \frac{1}{n} \sum_{i=1}^n (Y_i - \hat{Y}_i)^2. \quad (10)$$

On the other hand, the MSE of an estimator Θ “is a function of both its variance and its bias”. [WMS08, p. 393] It can be expressed as

$$\text{MSE}(\Theta) = \text{Var}(\Theta) + (\text{Bias}(\Theta))^2.$$

Thus, if DepDetector optimizes by searching for the dependency with the smallest MSE, it thereby selects for small variance and small bias.

2.6 Robustness of FDs

Koudas et al., when introducing the MFD, write that the definition of the canonical FD was “not *robust* enough to capture functional relationships on data obtained from merging heterogeneous sources [...]”⁵. [Kou+09, p. 1] However, they do neither provide a definition of robustness, nor a way to measure it.

Bleifuß et al. [Ble+16, p. 3] define a measure similar to robustness called *correctness*. They do so by comparing a set of RFDs *Out* to a set of FDs *Gold* on the same relational instance r as follows:

$$\text{correctness} = \frac{|Out \cap Gold^+|}{|Out|}$$

$Gold^+$ represents the transitive closure of *Gold*, meaning that it includes all non-minimal FDs implied by *Gold* as defined in equation 3.

⁵Highlighting added by the author.

By this definition, the notion ‘correctness’ considers *Out* to be more or less ‘correct’ depending on how many members it has in common with $Gold^+$. If an RFD is used to find relations on noisy data, this definition of ‘correctness’ does not meaningfully measure dependencies that attempt to take the noise into account. This means that an RFD that takes the typing error in table 3 into account would lead to a less ‘correct’ set of RFDs than an RFD that does not.

In conclusion, Koudas et al. do not provide a definition of robustness and ‘correctness’ cannot be used for measuring robustness. Inspired by the theory of machine learning, the application of cross-validation techniques to FDs to measure robustness is proposed.

Definition (Train Set, Test Set, Split Ratio) Let $r = \{t_1, t_2, \dots, t_p\}$ be a relation on a relation scheme R where $p \in \mathbb{N}$ is the number of tuples in r . Let $s \in [0, 1]$ be the *split ratio* and let $m = \lfloor s \cdot p \rfloor$, then we define:

$$r_{train} = \{t_1, t_2, \dots, t_m\} \quad (11)$$

$$r_{test} = \{t_{m+1}, t_{m+2}, \dots, t_p\}, \quad (12)$$

such that $r = r_{train} \dot{\cup} r_{test}$ is split into two disjunct subsets. We call r_{train} the *train set* and r_{test} the *test set*.⁶ [SV08, p. 56]

Definition (Imputation Derived by an FD) Let r be a relational instance on a relational scheme R . In addition, let $X \subseteq R$ and $A \in R, A \notin X$. Furthermore, let $X \rightarrow A$ be a non-trivial, minimal FD of r_{train} , $r = r_{train} \dot{\cup} r_{test}$. For all tuples $t \in r_{train}$ where

$$t[X] = t'[X], \quad (13)$$

we call the RHS $t[A]$ the *imputation* of $t' \in r_{test}$ derived by the FD $X \rightarrow A$.

Definition (Set of Imputations) Let r be a relational instance on a relational scheme R . Moreover, let $X \subseteq R$ and $A \in R, A \notin X$. Also, let $X \rightarrow A$ be a non-trivial, minimal FD of r_{train} , where $r = r_{train} \dot{\cup} r_{test}$. We then call

$$r_{imp}^i = \{t[A] \mid t[X] = t'_i[X], t'_i \in r_{test}, t \in r_{train}\} \quad (14)$$

set of imputations of the i th tuple in the test set.

⁶The naming ‘test set’ is ambiguous in literature. Bishop [Bis06] and Smola et al. [SV08] call ‘validation set’ what is called ‘test set’ by Burkov [Bur19, ch. 5, p. 8-9]. In this work, the notation from [Bur19] is adopted.

Definition (Robustness) Let $r = \{t_1, t_2, \dots, t_p\}$ be a relational instance on a relational scheme \mathbf{R} , $p \in \mathbb{N}$. Moreover let $X \subseteq \mathbf{R}$ and $A \in \mathbf{R}$, $A \notin X$. Let $X \rightarrow A$ be a non-trivial, minimal FD of r_{train} . Furthermore, let $r = r_{train} \cup r_{test}$ be split with split ratio $s \in [0, 1]$ and let $m = \lfloor s \cdot p \rfloor$.

For each tuple $t'_i \in r_{test}$, let r_{test}^i be the set of imputations with regard to $X \rightarrow A$. For all $r_{test}^i \neq \emptyset$, we arbitrarily choose one tuple $t_i^* \in r_{test}^i$. We call

$$r_{imp} = \{t_i^* \mid i \in \{m+1, m+2, \dots, p\}\} \quad (15)$$

set of imputations of the test split.

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

$$\text{robustness}_{\text{MSE}} = \frac{1}{p-m} \sum_{i=m+1}^p (t_i^*[A] - t'_i[A])^2. \quad (16)$$

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

$$\text{robustness}_{\text{F1-Score}} = \left(\frac{\text{Recall}(r_{imp}, r_{test})^{-1} + \text{Precision}(r_{imp}, r_{test})^{-1}}{2} \right)^{-1} \quad (17)$$

2.7 Datawig Imputer

Datawig Imputer is a framework to train models that perform classification- and regression tasks used to impute missing values in tables. [Bie+18, p. 1] Models are generated using *supervised learning*.

In this work, Datawig Imputer is used to train classification-models and regression-models on subsets of relational instances based on sets of attributes of a relation scheme. This function is referred to throughout this work by the name ‘ML Imputer’. The DepDetector algorithm later introduced in this work uses ML Imputer to train models for dependency detection.

The way Datawig Imputer interprets inputs, observations are a set of tuples $\{(x_1, y_1), (x_2, y_2), \dots\}$. [SV08, p. 10] We call $x_i \in \mathbb{F}$ *feature-vector* in a *feature-space* $\mathbb{F} \subseteq \mathbb{R}^n$ and $y_i \in S$ *output-attributes* or *labels*. [DHS00, p. 7]

If Datawig Imputer solves a classification problem, a function

$$f : \mathbb{F} \rightarrow S$$

is being approximated, where \mathbb{F} is the feature-space and $S = \{a_1, a_2, \dots, a_M\}$ is a set of labels.

Datawig Imputer assumes that the column to be imputed, the so-called *label column*, contains data that can be transformed to obtain attributes a_1, \dots, a_M . [Bie+18, p. 2018] All columns on a database table that are *not* the label column are called *feature columns*. These feature columns are transformed by Datawig Imputer to obtain feature-vectors.

The transformation performed to obtain learnable inputs assign to the string data of each column c a numerical representation x^c . [Bie+18, p. 2020] These functions are called *encoders*. When solving a classification problem, Datawig Imputer separates data into two categories: *categorical data* and *sequential data*. [Bie+18, p. 2017] Different encoders are chosen to transform categorical data and sequential data.

Categorical data are data that consist of *categorical variables*. “A categorical variable places a case into one of several groups of categories,” define Moore et al. [Moo+11, p. 4] An example-column c_{cat} containing categorical variables can be represented by the following set of colors:

$$c_{cat} = \{\text{blue, red, yellow, blue, blue, red}\} \quad (18)$$

Datawig Imputer creates a histogram of column c and uses the index of the histogram $x^c \in \{1, 2, \dots, M_c\}$ to generate a numerical representation. Thus the encoded column

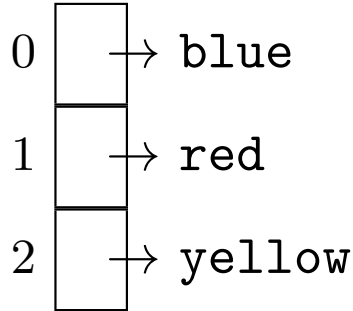


Figure 4: State diagram showing the histogram established by Datawig Imputer to encode example 18

c_{cat} would have the following form:

$$c_{cat}^c = \{0, 1, 2, 0, 0, 1\}$$

Sequential data is data where the sequence in which individual data points are stored in contains information. An example for sequential data would be a column containing non-categorical strings, like usernames:

$$c_{seq} = \{\text{ItalyPaleAle, sisou, primos63}\}$$

The numerical representation $c_{seq}^c \in \{0, 1, 2, \dots, A_c\}^{S_c}$ of the sequential column “is a vector of length S_c , where S_c denotes the length of the sequence or string in column c_{seq} .” [Bie+18, p. 2020] This numerical representation is called *n-gram representation*. For further details on the n-gram implementation, consider [Bie+18, p. 2020].

Finally, the Datawig-Imputer learns to predict the label distribution of $y \in 1, 2, \dots, D_y$ from the feature-vector x . It therefore models $p(y|x, \Theta)$, the D_y - dimensional probability vector over all possible values in the to-be imputed column in function of feature-vector x and *learned model parameters* Θ . [Bie+18, p. 2021]

Datawig uses a logistic regression type output layer to achieve this:

$$p(y|x, \Theta) = \text{softmax}[Wx + b]$$

Here, the learned model parameters $\Theta = (W, z, b)$ depends on the learned *weights* W and *biases* b as well as z , representing all parameters of the column-specific feature extraction. Parameters Θ are learned by minimizing cross-entropy loss between predicted and observed labels y which is computed as follows:

$$\Theta = \min_{\Theta} \sum_1^N -\log(p(y|x, \Theta))^{\top} \text{onehot}(y) \quad (19)$$

Here, $\log()$ represents the element-wise logarithm. $\text{onehot}(y) \in 0, 1^{D_y}$ stands for a one-hot encoding of one label y .

Equation 19 applies the principle of *Empirical Risk Minimization*. [Vap92, p. 832] The term right of the equality sign is called *empirical risk-function* for the classification problem. By minimizing risk, in this case cross-entropy loss, weights Θ that lead to an optimized classification performance are ‘learned’.

Training is done using standard backpropagation and stochastic gradient-descent on mini-batches. The exact network layout is described in the paper “Deep Learning for Missing Value imputation in Tables with Non-Numerical Data” [Bie+18, p. 2022] in full detail.

If a regression model is learned, DataWig uses a “standard regression loss function” [Bie+18, p. 2025] in an approach similar to the one described for the classification case.

2.8 FD Imputer: Measuring Robustness

In order to measure robustness of FDs as defined in equations 17 and 16, an algorithm called *FD Imputer* is introduced. To describe how FD Imputer functions, an example

follows. In table 4 a splitting is exemplified for $s = 0.5$ and $p = 4$. The approach and terminology involved in the definition of robustness are based on methods from statistical cross-validation used in model selection. [Hay08, p. 172] When measuring robustness, FDs are first detected on the train set using an FD detection algorithm. In this work, the HyFD algorithm is chosen. [PN16] FD Imputer is then used for measuring robustness of each FD detected on the train set.

A	B	C	D
Blue	Car	Portugal	Lisbon
Yellow	Car	Portugal	Lisbon
Green	Bus	Spain	Madrid
Grey	Car	Portugal	Lisbon

(a) Original relation r .

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

(b) Train set r_{train} .

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

(c) Test set r_{test} .

Table 4: The relation r is split into train- and test sets with a split-ratio $s = 0.5$.

FD Imputer iterates over all tuples in the test set and searches for a tuple in the train set with the exact same LHS. If one or more such tuples in the train set are found, FD Imputer randomly selects one of these tuples and saves the RHS of that tuple. As defined in the previous section, this RHS value is called *imputation* of the RHS value on the train set.

This procedure is repeated for all tuples in the train set. In a last step, FD Imputer compares the imputations with the actual RHS values in the train set. This enables a performance measurement. If the RHS of the FD contains categorical data, the F1-measure can be calculated for each FD. Otherwise, in case the FD's RHS contains continuous data, the MSE is chosen as measure of robustness.

Consider table 4. The FD $C \rightarrow D$, when used to find a RHS value on r_{test} , correctly returns 'Lisbon'. However, FD $A \rightarrow D$ for example will lead to FD Imputer yielding 'Madrid', which is an incorrect imputation.

The implementation of FD Imputer leverages SQL join clauses to retrieve imputations. FD Imputer performs an inner join on all LHS columns, joining train-set and test-set. A

successive second left join clause concatenates the original test-set with the column of imputed tuples that originate from the first join. The result is a set of imputations as described in equation 15.

Robustness as a F1-Score The imputed values retrieved by FD Imputer can be leveraged to create a measure for robustness. First, each imputed value is interpreted as a label. These labels are then each classified as in the binary classification case (see figure 2). Next, precision and recall are calculated and the F1-Score is derived according to equation 9.

Lastly, the F1-Score weighted according to the number of true instances for each label is calculated.⁷

This means that to a set $L = \{l_1, l_2, \dots, l_n\}$ of $n \in \mathbb{N}$ labels, a set $X = \{x_1, x_2, \dots, x_n\}$ of n F1-Scores corresponds. Let $W = \{w_1, w_2, \dots, w_n\}$ be a set of weights, where each weight represents the number of true instances for each label. Then the weighted F1-Score of a dependency can be calculated as follows:

$$\text{F1-Score} = \frac{\sum_{i=1}^n x_i w_i}{\sum_{i=1}^n w_i} \quad (20)$$

This weighted F1-Score, calculated for each FD on a relation r , is called robustness of an FD.

Robustness as a Mean Square Error Depending on the kind of data imputed by FD Imputer, the measure used to express robustness is adapted. As shown in the previous paragraph, performance when imputing categorical data can be measured by the F1-Score. If the RHS's content is of a continuous data type, using a classification performance measure is pointless.

Therefore, FD Imputer distinguishes between continuous and categorical data. When imputing continuous data, the *mean square error* (MSE) is chosen as measure of robustness. The MSE can be computed according to equation 10.

⁷See [sklearn documentation](#) for `sklearn.metrics.f1_score()` with weighted average for more details on the implementation.

2.9 ML Imputer: Dependency Based Imputation

The ML Imputer algorithm adapts the DataWig framework so that it can be used to run on FDs. For a FD $X \rightarrow A$, ML Imputer interprets attribute-members of X as feature columns and the attribute A as label column.

This modification is made to DataWig in order to compare the results yielded by FD Imputer to the results of a learned model for imputation. Both imputers use cross-validation to find impute-values on a test set. They derive either a F1-Score as in equation 20 for classifiable data or a MSE as in equation 10.

2.10 DepDetector: RFD Discovery

In the field of data profiling, an extensive stock of theory and algorithms for FD discovery has been created in recent decades. [Abe+19, p. 39] New techniques such as massive parallelization or hybrid algorithms combining different discovery strategies continuously shortened the time needed to detect FDs on a relational table. [Abe+19, p. 40]

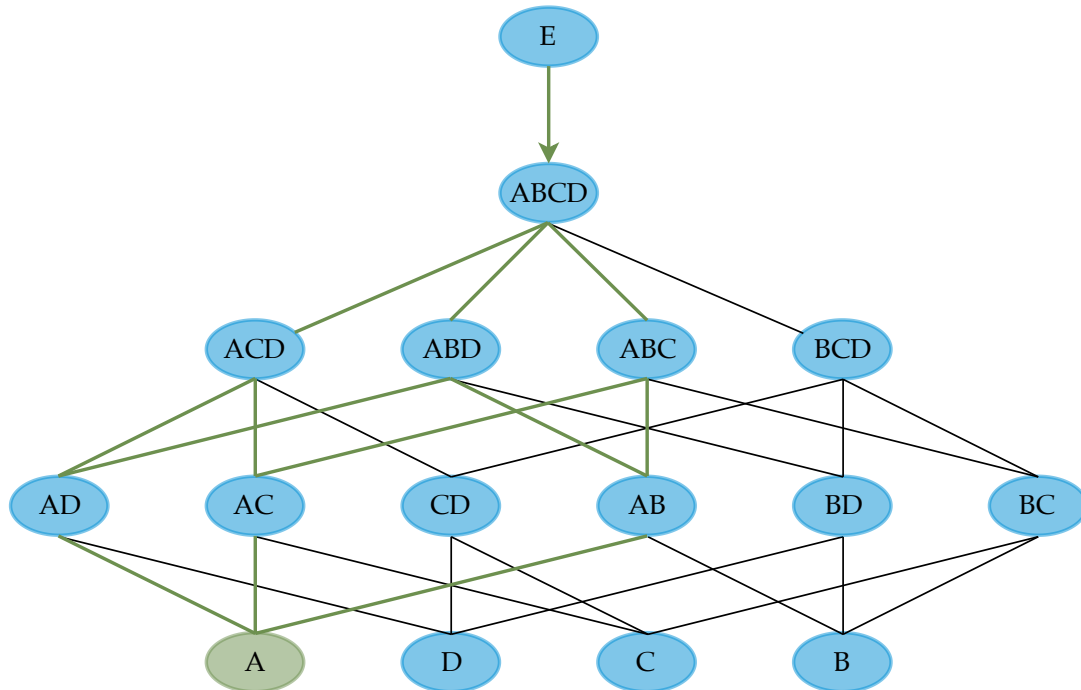


Figure 5: Hasse diagram of the search lattice when discovering minimal, non-trivial FDs on a relational schema with five columns.

Notably the TANE algorithm developed by Huhtala et al. discovered FDs faster than any other algorithm at the time. [Huh+99, p. 108] It paved the way for other algo-

rithms that use the same search strategy called ‘lattice traversal’, such as FUN [NC08], FD_MINE [YHB02] and DFD [ASN14]. [Abe+19, p. 39]

In an effort to leverage the imputation-capabilities of ML Imputer for dependency discovery, the algorithm called *DepDetector* is introduced. *DepDetector* finds RFDs that relax both on the attribute comparison as well as on the extent. When discovering dependencies, *DepDetector* searches all non-trivial, minimal dependencies. Similar to FD discovery algorithms, the search space can be modeled as a graph labeling problem. [Abe+19, p. 24] Here, every possible combination of attributes, called *power set* of attribute combinations, from the nodes of the graph. Due to the properties of the graph, the problem can be modelled as a *lattice*. [Abe+19, p. 24] This has been visualized using a Hasse diagram in figure 5: Column E is the RHS for which a minimal LHS is searched. Green edges visualize paths to the minimal LHS A. On this path, a number of candidate-LHSs are passed through.

2.10.1 DepDetector Algorithm

DepDetector searches for dependencies by traversing the search lattice as depicted in figure 5. In a first step, ML Imputer is run, imputing the RHS with all other columns on the relational scheme as LHS. This yields the model’s robustness. For the algorithm to continue traversing the search lattice in search of dependencies, this model’s robustness must exceed an initial threshold value. The threshold value depends on the data type of the potential RHS. The model is required to either achieve an F1-Score equal or greater than 0.9 or an MSE of less than $0.2 \cdot \bar{y}$, where \bar{y} is the arithmetic mean of all RHS values in the dataset. If the initial model does not exceed the threshold, we can assume that there is no dependency present on the relational instance.

However, if the initial threshold is surpassed, the search-algorithm traverses the search lattice in a ‘top-down’ manner as depicted in figure 5. For each node in the search lattice a model is trained and evaluated. Again, the search only continues when the model exceeds a threshold value: either a smaller MSE than the parent-node’s or a F1-Score greater than 98% of the F1-Score of the parent node. This is repeated until the threshold is not exceeded by any new model or until a minimal LHS is reached.

DepDetector Relaxing on the Extent The threshold-driven search-strategy is effectively a relaxation on the extent of the canonical FD when a potential RHS is of a classifiable data type. It is similar to an AFDs approach of “holding on almost every tuple” — the imputation needs to be correct on almost every instance, but not on every instance. In addition, ML Imputer is capable of learning conditions, similar to CFDs. As shown in the experimental section of this work, ML Imputer learns the order of entries in a

relational instance as well as thresholds, such as a CFD^P would do. Compared to CFDs, ECFDs and CFD^P s, it is important to note that these conditions are never set manually, but learned.

The behavior of DepDetector as described in the paragraph above is also displayed in listing 1. Note that DepDetector is implemented with two search-modes called ‘greedy’ and ‘complete’. The ‘complete’ search-mode behaves as described above. It is optimized to find all dependencies on a relational scheme. Meanwhile the ‘greedy’ search-mode is optimized to detect just *one* dependency. This approach converges faster than ‘complete’ and can be used if a single minimal dependency per RHS is sufficient. However, it neither guarantees to find minimal dependencies, nor does it guarantee to find the most robust dependency.

```

1 def get_complete_candidates(root):
2     convergence = True
3     most_recent_nodes = root.get_newest_children()
4     for node in most_recent_nodes:
5         if root.is_continuous: # continuous data
6             if (node.score < node.parent.score) and (len(node.name) > 1):
7                 convergence = False
8                 pot_lhs = node.name
9                 for col in pot_lhs:
10                     candidate_lhs = [c for c in pot_lhs if c != col]
11                     add_node(candidate_lhs,
12                             parent=node,
13                             score=None)
14
15         elif not root.is_continuous: # classifiable data
16             if (node.score > node.parent.score*0.98) \
17                 and (len(node.name) > 1):
18                 convergence = False
19                 pot_lhs = node.name
20                 for col in pot_lhs:
21                     candidate_lhs = [c for c in pot_lhs if c != col]
22                     add_node(candidate_lhs,
23                             parent=node,
24                             score=None)

```

Listing 1: ‘Complete’ candidate generation in the DepDetector algorithm

DepDetector Relaxing on Attribute Comparison A fundamental difference that manifests itself when comparing dependencies discovered by DepDetector with MFDs or FDs is due the fact that DepDetector does not derive the existence of a dependency based on attribute comparison as in equations 1 or 6. Instead, DepDetector uses ERM-strategies when executing ML Imputer to measure a dependency's robustness. Thus, DepDetector optimizes for maximum robustness and a minimum number of LHS attributes.

3 Experiments

To examine robustness as it is defined in the previous chapter, a number of experiments are conducted. Then, in order to evaluate the capabilities of ERM-techniques for FD discovery, the DepDetector algorithm is run on a number of datasets. For the experiments, datasets from the UCI Machine Learning Repository are used. [DG17]

Since the algorithms analyzed in this section handle missing values differently, rows containing missing values are excluded from the datasets due to possible inconsistencies when comparing results.

3.1 FD Imputer

FD Imputer is run for every FD found on a train set of a dataset. The measured performance-score is called robustness of the FD. First, experiments on the imputation of classifiable data with the FD Imputer are presented. Then, results obtained when imputing continuous data are discussed.

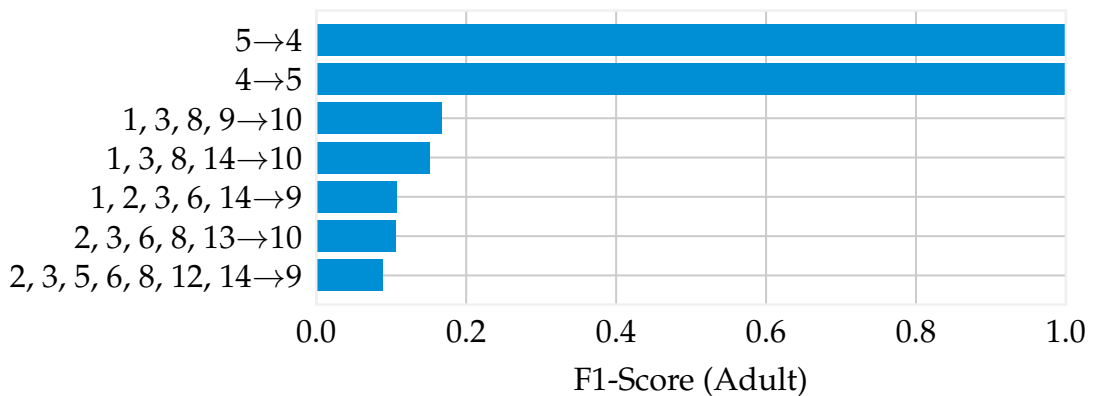


Figure 6: F1-Scores of the seven most-robust FDs on the Adult dataset when imputing classifiable values.

3.1.1 Imputation of Classifiable Data

Figure 6 shows the performance of FD Imputer on the Adult dataset. The two top performing FDs have a perfect F1-Score of 1 each. An explanation for this circumstance can be found when analyzing the content of columns 4 and 5.

Column 4 contains information about the highest educational level achieved. There are 16 different categories of educational level defined. Each category is assigned an integer in a range spanning from 0 to 15. This integer is the content of column 5. Thus, the relation between column 4 and column 5 can be modeled by a bijective function, projecting the domains of each attribute onto the other. In consequence, FDs between column 4 and 5 are perfectly robust.

All other FDs found on the Adult dataset lead to F1-Scores smaller than 0.2, being substantially less robust than the two top-performing FDs. It can be derived that only the two top-performing FDs can be used to meaningfully impute data. If new data was added to the dataset, it can thus safely be assumed that these two FDs still held. In reverse, no value in a column other than 4 or 5 can be meaningfully imputed using FDs detected on the Adult dataset due to the small robustness of these FDs. Figure 7

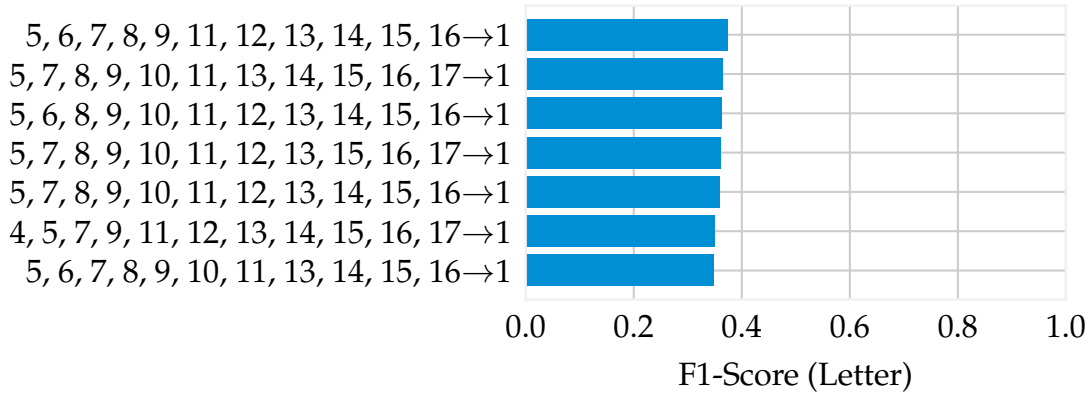


Figure 7: Robustness of the seven top-performing FDs on the Letter dataset when imputing classifiable RHSs.

displays the result of running FD Imputer on the Letter dataset. Of the 80 FDs found on the train split, the seven most robust ones share the same RHS, column 1. No F1-Score greater than 0.37 is achieved on this dataset.

All analyzed datasets other than Adult yield no FD with a perfect score. Only an FD on the Breast Cancer Wisconsin dataset achieves a score higher than 0.75. Table 5 provides a summary of how FD Imputer performs on eight different datasets.

In Table 5, Column #FDs indicates how many FDs were found on the complete dataset. #FDs_{train} contains the number of FDs that were detected on the train set. $F1_{\text{mean}}$ and $F1_{\text{max}}$ indicate the arithmetic mean and maximal F1-Score achieved on each dataset respectively. The last column named #(F1 = 0) provides the number of FDs that scored a F1-Score of 0.

Dataset	Attributes	#FDs	#FDs _{train}	Classification Performance		
				#FD (F1 = 0)	$F1_{\text{mean}}$	$F1_{\text{max}}$
Abalone	10	175	193	45	0.0008	0.0048
Adult	16	93	88	10	0.0669	1.0000
Balance S.	6	7	7	6	0.0000	0.0000
Breast C. W.	12	57	77	10	0.2198	0.7539
Chess	8	9	9	8	0.0000	0.0000
Iris	6	9	8	1	0.1274	0.2252
Letter	18	78	80	17	0.2347	0.3737
Nursery	10	11	11	10	0.0000	0.0000

Table 5: Performance of the FD Imputer on a selection of UCI datasets.

The results displayed in table 5 show that robustness depends strongly on the dataset considered. However, no relation between #FDs and robustness can be identified (see appendix figure 16). Analysis on a bigger body of datasets seems necessary to draw further conclusions.

3.1.2 Imputation of Continuous Data

The majority of columns for which FDs are detected in this work contain classifiable data. However, some FDs are found whose RHSs contains continuous data. For brevity, we call such an FD continuous FD or *cFD*. Table 6 displays an overview of how FD Imputer performs when imputing continuous RHS values. Continuous FDs were detected on seven train sets of the eight datasets considered — column ‘# continuous FDs_{train}’ in table 6 shows how many cFDs were detected.

Column ‘0-Coverage cFDs’ displays for how many cFDs not a single imputation was found. If all cFDs found on the train set are ‘0-coverage cFDs’, the resulting ‘Coverage’ is thus 0% with missing minimum and maximum values.

The column named ‘Coverage’ indicates the mean imputation coverage. Mean imputation coverage is the mean percentage of rows in the test set for which imputations were found per cFD. It was computed according to the following formula:

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

Note that the computed mean coverage is a percentage.

On the Letter dataset, no cFD is detected. Five out of seven datasets for which cFDs were detected have a mean imputation coverage of 0%. Concerning the two datasets with nonzero mean imputation coverage, the mean imputation coverage is approximately one per mill.

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

Since the MSE is an absolute error measure, comparisons between cFDs concerning different RHSs do not provide any insights. Table 7 gives an overview of all cFDs examined on the Abalone dataset. Mind that every cFD that is used by FD Imputer to return imputations yields exactly one MSE.

Column ‘RHS’ indicates the RHS, for which cFDs were examined. The next column ‘Var(MSE)’ contains the experimental variance of the mean MSE. The latter is called $\overline{\text{MSE}}$ and is displayed in the following column. The last two columns ‘MSE_{min}’ and ‘MSE_{max}’ contain the minimum and maximum MSEs.

From the eight RHSs displayed in table 7, six RHSs show similar statistical results. For these six RHSs, the variance of MSEs is very small and lies between 10^{-7} and 10^{-9} . In addition, the mean MSEs are distributed within two orders of magnitude, the maximum MSE within one order of magnitude.

All values observed for RHS 0 are between 10^5 and 10^{21} bigger than values observed for other RHSs. Column 0 of the Abalone dataset contains the row-ID. Since there is no dependency between any column-combination in the dataset and the row-ID, a high MSE seems plausible.

RHS	Var(MSE)	$\overline{\text{MSE}}$	MSE_{\min}	MSE_{\max}
0	$2.0712 \cdot 10^{13}$	$7.9385 \cdot 10^6$	$8.2705 \cdot 10^5$	$1.3418 \cdot 10^7$
2	$1.3039 \cdot 10^{-7}$	$2.3393 \cdot 10^{-4}$	0.0000	$9.2500 \cdot 10^{-4}$
3	$3.5634 \cdot 10^{-7}$	$3.5412 \cdot 10^{-5}$	0.0000	$1.6250 \cdot 10^{-4}$
4	$3.4298 \cdot 10^{-9}$	$2.0972 \cdot 10^{-4}$	$1.1250 \cdot 10^{-4}$	$3.1250 \cdot 10^{-4}$
5	$3.0624 \cdot 10^{-9}$	$6.9775 \cdot 10^{-5}$	$9.0000 \cdot 10^{-6}$	$1.6700 \cdot 10^{-4}$
6	$2.2867 \cdot 10^{-9}$	$1.5017 \cdot 10^{-4}$	$8.4500 \cdot 10^{-5}$	$1.9700 \cdot 10^{-4}$
7	0.0000	0.0000	0.0000	0.0000
8	$7.4015 \cdot 10^{-9}$	$1.0366 \cdot 10^{-4}$	0.0000	$2.2250 \cdot 10^{-4}$

Table 7: Statistical analysis of the computed MSEs when running FD Imputer with 84 different cFDs on the Abalone dataset.

If column 7 of the Abalone dataset is imputed with cFDs, the task results in a perfect MSE of 0. On closer inspection, column 7 is the RHS of eleven cFDs. Only one of those eleven cFDs can be used by FD Imputer to impute exactly one value — which happens to be an exact match.

Half of the RHSs are determined by at least one cFD with a perfect MSE of 0. This shows that perfect imputations are common when running FD Imputer with cFDs.

3.2 ML Imputer

ML Imputer is executed with the same set of FDs as FD Imputer. Analogously to the discussion of experiments executed with FD Imputer, the evaluation will be separated into two sections, one discussing FDs with classifiable RHSs, the other examining FDs with continuous RHSs.

3.2.1 Imputation of Classifiable Data

The first dataset analyzed is Adult. Figure 8 displays the F1-Scores of the seven best performing FDs. The two top-scoring FDs are the same for FD Imputer and ML Imputer. The third most top-performing FD is $0 \rightarrow 14$, a relation between row-ID and nationality. There is no apparent dependency between row-ID and nationality. Inspect-

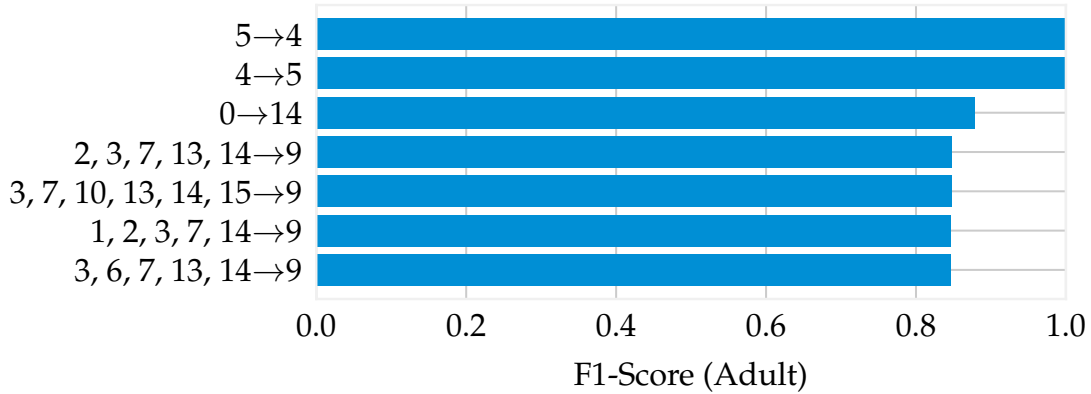


Figure 8: The figure compares the F1-Score of the FD Imputer compared to the F1-Score of the ML Imputer. Each point represents one FD.

ing the distribution of values in column 14, one cause for the performance of FD $0 \rightarrow 14$ can be derived. Of 24154 rows in the train set, 22020 entries contain the value 'United-States'. ML Imputer learns the same imputation for all 2987 entries in the test set, namely 'United-States'. This perfectly explains the performance of FD $0 \rightarrow 14$.

The bar plot in figure 9 shows the performance of ML Imputer on the Chess dataset. The FD with the biggest F1-Score is $0 \rightarrow 7$. Column 0 contains the row-ID, whereas column 7 contains the outcome of a chess-game, assuming both players play perfectly. Again, there does not seem to exist a relation between the two columns. On closer inspection however it can be observed that the Chess-dataset is sorted. For tuples with row-ID of less than 2796, column 7 contains always the value 'draw'. Tuples that have a row-ID bigger than 2795 but smaller than 2824 always contain the value 'zero' and so on.

This showcases the capabilities of ML Imputer to learn relaxation on the extent, similar to what is used in the definition of CFDs, ECFDs and CFD^ps: ML Imputer approximately learns a pattern where a dependency $0 \rightarrow 7$ holds, if $0 \leq \text{row-ID} \leq 2795$ or $2796 \leq \text{row-ID} \leq 2823$, and so on.

Table 8 shows a generally higher performance of ML Imputer compared to FD Imputer. There are no FDs for which ML Imputer scores 0. This can be explained by the property of ML Imputer of always returning some imputation value, even if the probability with

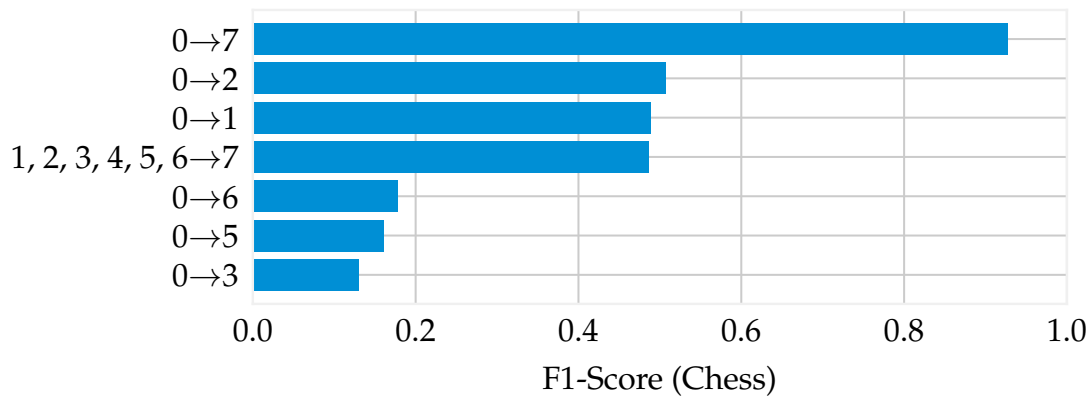


Figure 9: F1-Scores of the seven top-performing FDs on the Adult dataset.

which ML Imputer predicts a value is low. On the other hand there is FD Imputer, which frequently does not return any imputation value at all.

Dataset	#FDs	#FDs _{train}	Classification Performance		
			F1 _{mean}	F1 _{max}	F1 _{min}
Abalone	175	193	0.3697	0.5664	0.0619
Adult	93	88	0.7720	1.0000	0.0409
Balance Scale	7	7	0.4670	0.9443	0.0255
Chess	9	9	0.3744	0.9276	0.1162
Iris	9	8	0.9243	1.0000	0.1275
Letter	78	80	0.7207	0.9134	0.0058
Nursery	11	11	0.9915	0.4531	0.1138

Table 8: Performance of ML Imputer on eight different datasets when training on FDs with classifiable RHS entries.

3.2.2 Imputation of Continuous Data

When ML Imputer is run with cFDs, internally a regression is performed. To examine how this mechanism behaves, imputation generation based on cFDs was performed. ML Imputer always returns imputations, thus there are no missing values as listed in table 6 — ML Imputer coverage is always 100%.

As done for the FD Imputer in table 7, an in-depth analysis is performed on the Abalone dataset in table 9. The biggest mean MSE is returned for RHS 0, which contains the row-ID. This can be explained with the observation that the data is neither sorted, nor are there any other dependencies between the contents of the Abalone dataset and the row-ID.

RHS	Var(MSE)	$\overline{\text{MSE}}$	MSE_{\min}	MSE_{\max}
0	$3.7736 \cdot 10^9$	$1.4656 \cdot 10^6$	$1.3884 \cdot 10^6$	$1.6045 \cdot 10^6$
2	$1.3150 \cdot 10^{-5}$	$1.5828 \cdot 10^{-3}$	$3.0536 \cdot 10^{-4}$	$1.4630 \cdot 10^{-2}$
3	$5.3958 \cdot 10^{-6}$	$9.6910 \cdot 10^{-4}$	$2.0273 \cdot 10^{-4}$	$1.0238 \cdot 10^{-2}$
4	$9.2724 \cdot 10^{-8}$	$3.2361 \cdot 10^{-4}$	$2.2714 \cdot 10^{-4}$	$1.5399 \cdot 10^{-3}$
5	$3.5212 \cdot 10^{-3}$	$2.1342 \cdot 10^{-2}$	$1.6185 \cdot 10^{-3}$	$2.5844 \cdot 10^{-1}$
6	$1.5414 \cdot 10^{-4}$	$6.6988 \cdot 10^{-3}$	$1.2499 \cdot 10^{-3}$	$5.2755 \cdot 10^{-2}$
7	$1.1345 \cdot 10^{-5}$	$1.9497 \cdot 10^{-3}$	$4.8120 \cdot 10^{-4}$	$1.2581 \cdot 10^{-2}$
8	$1.7069 \cdot 10^{-5}$	$2.8521 \cdot 10^{-3}$	$9.7318 \cdot 10^{-4}$	$1.9730 \cdot 10^{-2}$

Table 9: Statistical analysis of the Abalone dataset obtained by running ML Imputer with the same 84 cFDs as used in the generation of table 7.

All other RHSs lead to performance measures that seem plausible. Variances are between one and four orders of magnitude smaller than the corresponding arithmetic means, minimum and maximum MSEs are separated by up to two orders of magnitude.

3.3 Overfitting ML Imputer

If the neural network trained by ML Imputer learns too many examples, the resulting model is just a memorization of the whole training-set — a phenomenon referred to as overfitting. [Hay08, p. 164] To further investigate ML Imputer’s capabilities of overfitting models, experiments are conducted to compare imputation performance between a overfitted model and a normally trained model. Overfitting was achieved by first uniting r_{train} and r_{test} and then training and testing the model on the same unified relational instance. Internally, this unification of train set and test set implies that the resulting train-error and test-error are the same entity, effectively preventing cross-validation.

Figure 10 compares an overfitted model with a cross-validated model on the Adult dataset. Every blue data point represents one FD, providing the data with which the models were trained. The red graph in figure 10 is the graph of the identity function.

All data points that lie below the graph of the identity function indicate that the overfitted model scores higher than the cross-validated model. Data points that lie above the graph of the identity function represent the opposite. Both models score equally well, if a data point lies on the graph of the identity function. The distribution of data points

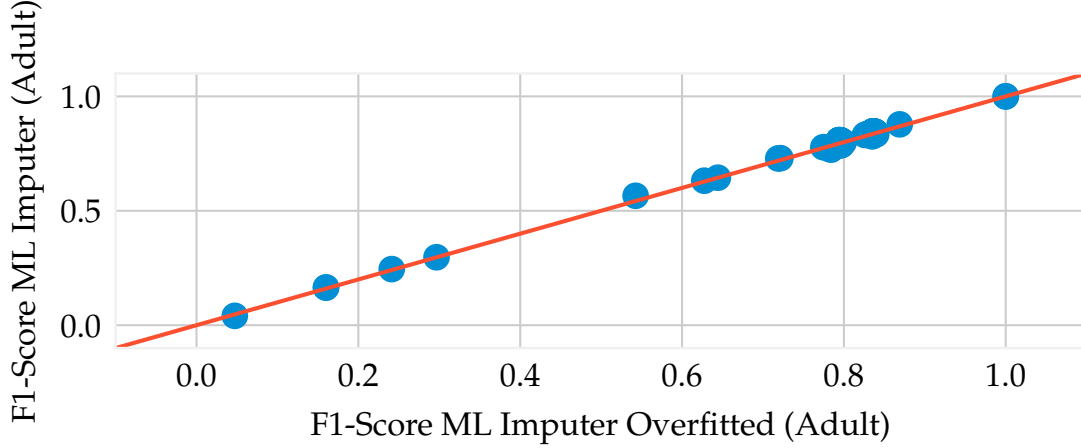


Figure 10: Comparison of an overfitted model with a conventionally trained model on the Adult dataset.

in figure 10 is well approximated by the identity function. It can thus be stated that the classification models trained by ML Imputer are not capable of overfitting: If the model complexity cannot be arbitrarily increased, overfitting is prevented by the network's simple design. [DHS00, ch. 6, p. 51-52] When continuous data is imputed, ML Imputer uses

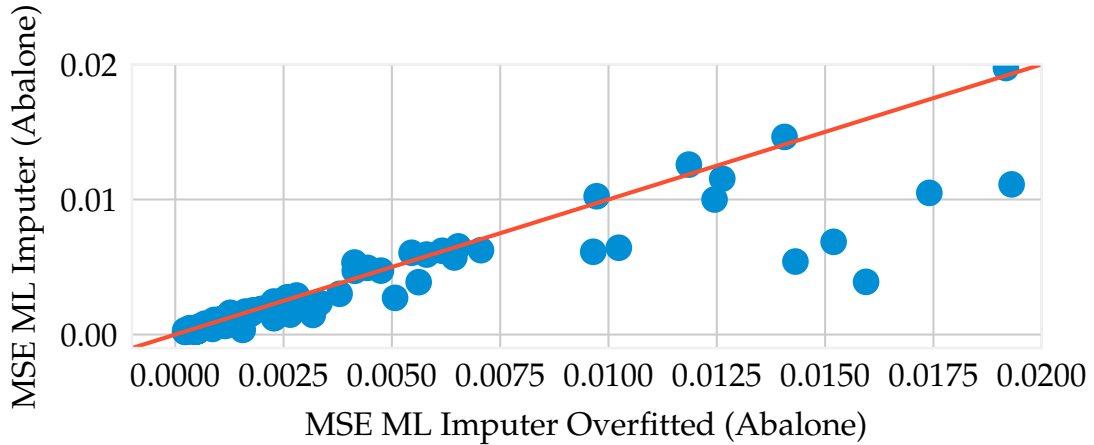


Figure 11: Overtrained models are compared with conventionally trained models on the Abalone dataset. Continuous RHSs were imputed to get MSEs.

another network design than when performing a classification task. [Bie+18, p. 2025] The

behavior regarding overfitting of the numerical regression loss function used in training to generate models to impute continuous data is represented in figure 11. Note that, while in figure 10 a greater score indicates a better performance, in figure 11 the inverse is true.

It can be observed that the network design used for regression on continuous data is more vulnerable to overfitting than the network used to impute classifiable data. The majority of cFDs results in nearly equally big MSEs of overfitted models and conventionally trained models. This is represented in figure 11 by data points that are positioned in close proximity of the graph of the identity function.

However, a considerable number of data points indicate that the overfitted models perform *worse* than the conventionally trained ones. At first glance, this result appears to be counterintuitive. One would expect a smaller MSE if more data is used in training.

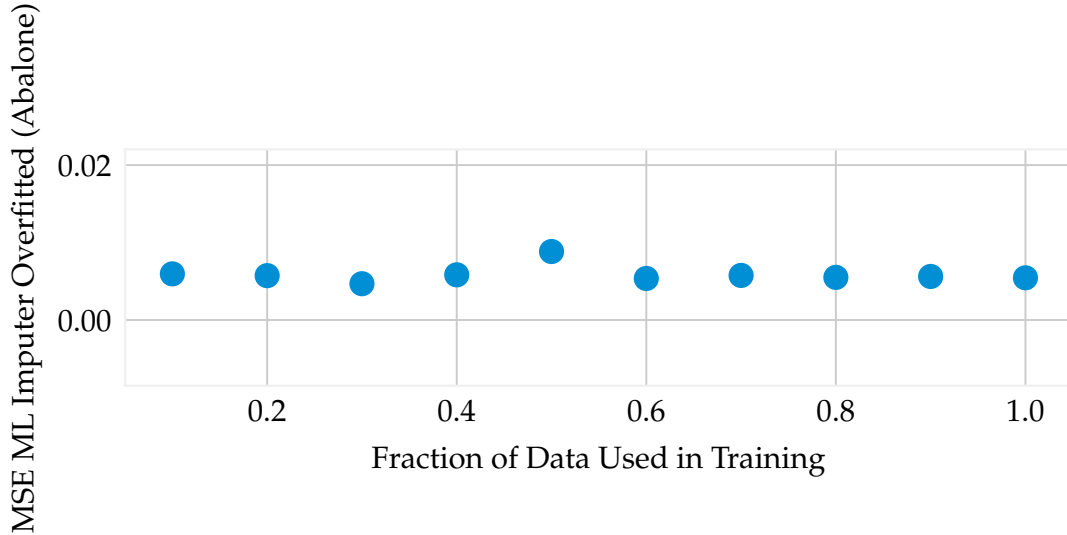


Figure 12: Overfitted imputation-models are trained with FD $[2, 3, 6, 7] \rightarrow 5$ on the Abalone dataset with varying dataset-size.

However, this phenomenon is well-known in practice. For example, Meng and Xie discuss the effect of training models with differently sized datasets. They show that more data during training does not necessarily imply a better model. [MX14, p. 1] On one side, additional data only leads to better models if it provides additional information to solve a problem. On the other side, as discussed previously, a neural network needs to be sufficiently complex to be able to leverage new data.

Figure 12 shows the MSE achieved by an overfitted ML Imputer model trained with different fractions of the total dataset. One can observe that the model trained on 10% of the total dataset scores a lower MSE as the model trained on 50% of the total dataset. Thus, when an overfitted model in figure 11 performs worse than the conventionally

trained counterpart, it can be anticipated that either the additional data contains no or little additional information regarding the regression problem, or that the design of the regression network does not provide enough complexity to incorporate additional information.

3.4 Dependency detection

DepDetector is run on number of datasets to detect generalized imputation dependencies. Since the models used by DepDetector to find minimal dependencies are results of stochastic processes, they are subjects to fluctuations. To further investigate how these fluctuations influence the result of the DepDetector algorithm, the number of training cycles used in training is varied and performance measures are compared in a section called *stability analysis*.

3.4.1 DepDetector Results on Various Datasets

Dependencies were detected on the eight datasets analyzed throughout this work. Models used by DepDetector were trained on a CPU. The machine on which the experiment is conducted is equipped with 7.2 GiB of RAM and an Intel Core i5-6300U CPU, a processing unit that runs four cores with clock speeds of 2.40GHz each. Table 10 lists the

Dataset	Cols	Rows	# FDs	Greedy dependencies	Complete dependencies
Abalone	10	4177	175	7 (42 min)	TL
Adult	16	32561	93	TL	TL
Balance-Scale	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

Table 10: Result of running dependency detection on selected datasets. Values in brackets are the runtime until convergence. ‘TL’ indicates that a time limit of 350 min is reached.

results. The column named ‘Greedy dependencies’ numbers how many minimal depen-

dependencies were discovered using the greedy search-strategy. Column ‘Complete dependencies’ indicates how many minimal dependencies were detected using the complete search-strategy.

For six of the eight datasets, the greedy search-strategy converged within 350 minutes. The complete search-strategy converged for three datasets within 350 minutes. The complexity of FD discovery is known to be $\mathcal{O}(n^2(\frac{m}{2})^2 2^m)$, where m is the number of attributes and n is the number of records. [Pap+15, p. 2] The findings in table 10 appear to be according to this rule, since the time limit was reached only for datasets with a great number of columns and/or rows.

3.4.2 Dependency stability

To investigate the stability of detected minimal dependencies, DepDetector is run for a fixed RHS with varying numbers of training cycles τ . The experiment was conducted on the Iris dataset and the Balance-Scale dataset. Minimal dependencies were searched following the ‘complete’ search-strategy.

The number of training cycles τ is increased stepwise from 3 to 15. For each τ the number of undetected minimal LHSs is calculated. This is done by first selecting the result of the run with $\tau_{max} = 15$ training cycles as a reference. Then, the LHSs obtained when $\tau < \tau_{max}$ are compared to the ones obtained when $\tau = \tau_{max}$. For each LHS found when DepDetector is run with τ_{max} training-cycles that is not contained in the result of a run where $\tau < \tau_{max}$, the number of undetected minimal LHSs is increased by one.

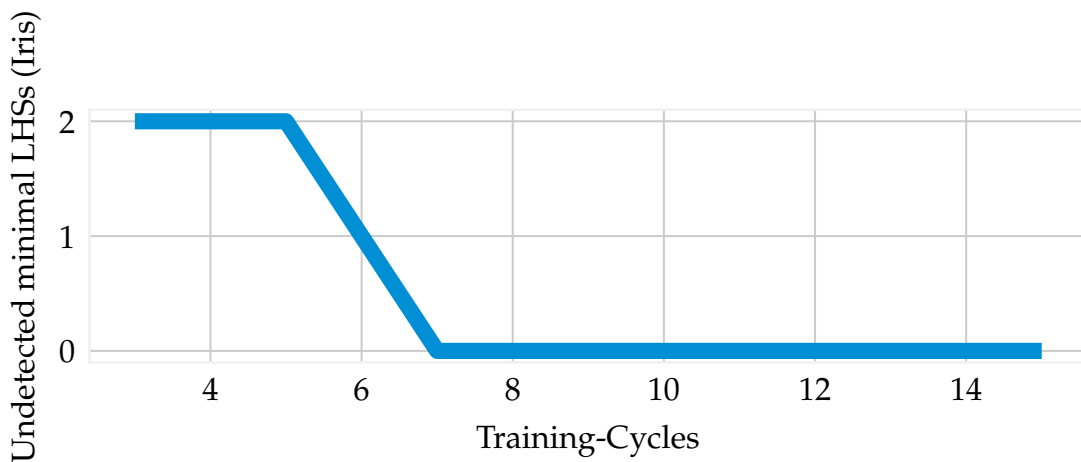


Figure 13: Analysis on the Iris dataset reveals that results become continuously more stable for larger τ .

Figure 13 shows the result of this procedure on the Iris dataset. When dependencies are searched with τ_{max} training cycles, DepDetector finds two minimal LHS combinations and thus two Dependencies. If $\tau \in [3, 5]$, both minimal dependencies remained undetected. In the run where $\tau = 6$, one of the two minimal LHSs was found. For $\tau \in [7, 15]$ training cycles, DepDetector discovered both minimal dependencies.

There are no fluctuations in the amount of undetected minimal LHSs once seven or more training cycles were used in dependency detection. This indicates that the minimal dependencies detected by DepDetector do not fluctuate every time new models are trained during detection.

The experiment is repeated on the Balance-Scale dataset. Figure 14 displays a non-linear behavior when detecting minimal LHSs in function of τ . Even though the result obtained for $\tau = 3$ is minimal, the result found when $\tau = 4$ misses one minimal dependency. For all DepDetector-runs performed with five or more training-cycles, all

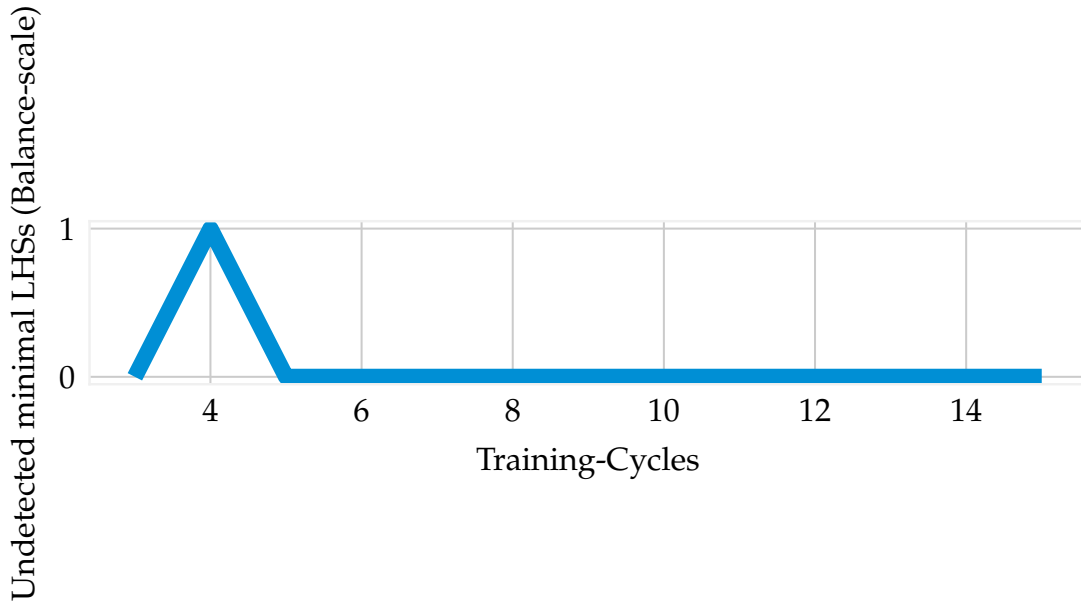


Figure 14: On the dataset Balance-scale, the analysis of the result shows a peak at $\tau = 4$.

minimal LHSs are detected. This is another hint that minimal dependencies identified with DepDetector converge to a stable solution.

While these two examples are no proof that DepDetector always yields a stable set of minimal LHSs, they motivate a default value of $\tau = 10$ training-cycles during dependency detection.

4 Discussion

4.1 Robustness

Robustness is introduced in order to measure the imputation-capabilities of an FD using FD Imputer. It is successfully demonstrated that FDs with a high robustness also contain humanly explainable meanings, for example when discussing figure 6.

However, it was possible to demonstrate that FDs do not offer any insights about the profile of data. They are fit for the specific task of schema normalization. However, the notion that “any dependency [could] be turned into a rule to check for errors in the data” [Abe+19, p. 9] does not seem to be true in general, but only for highly robust FDs with classifiable values in the RHS.

If an FD contains continuous values in the RHS, FD Imputer generally performs poorly. As shown in table 6, FD Imputer does not retrieve imputations for most RHSs in the test-set. In future works, implementing FD Imputer with a selection of RFDs might lead to further insights regarding robustness of dependencies containing continuous RHS values.

4.2 Comparing ML Imputer with FD Imputer

FD Imputer is implemented to probe the feasibility of value imputation using FDs. The usage of FDs for value implementation also makes a comparison with learned classifiers and regression-models possible.

The behavior of FD Imputer is similar to what can be observed when analyzing a case of overfitting. [SV08, p. 56] Haykins writes that “[Overfitting] is essentially a ‘look-up table’, which implies that the input-output mapping [...] is not smooth.” [Hay08, p. 165] The way FD Imputer functions is *literally* by using the test-set as a look-up table.

4.2.1 Imputing Sequential RHSs

Since FDs are established by exact equality of values, imputations of continuous values are very rare and, if they exist, either extremely accurate or very wrong (see table 7. FD Imputer cannot approximate numerical values, due to the definition of an FD. Data is always assumed to be classifiable.

In contrast to this, ML-Imputer is able to perform regression, predicting a continuous label for a given input with a specific uncertainty. This circumstance leads to a far superior performance of ML Imputer when imputing continuous values. Although MSEs for the FD Imputer model are generally smaller than MSEs measured for models trained with DepDetector, this effect is just a manifestation of the overfitting that takes place when FD Imputer looks up values. This result is emphasized by the findings in table 6: On the datasets considered in this work, FD Imputer cannot retrieve imputations from the test-set for an average of 99.8% of RHS values.

4.2.2 Imputing Classifiable RHSs

Figure 15 compares the F1-Scores of both ML Imputer and FD Imputer on the Adult dataset. One can observe that for almost all FDs, ML Imputer performs better than the FD Imputer. FD Imputer performance and ML Imputer performance appear to be partially proportional. If the score achieved by ML Imputer is lower than 0.7, the FD Imputer's F1-Score for the same FD is 0.

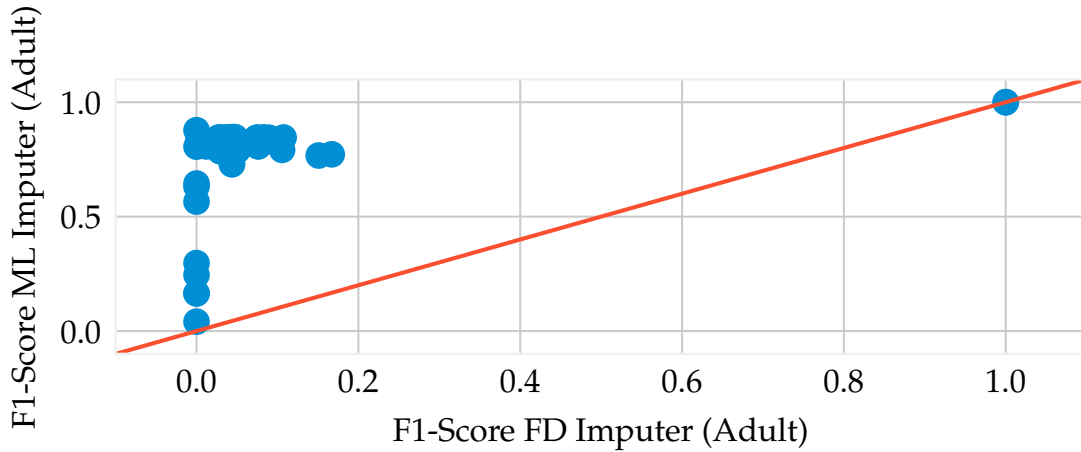


Figure 15: The figure compares the F1-Score of the FD Imputer compared to the F1-Score of ML Imputer. Each point represents one FD.

However, for FDs where ML Imputer's scores are greater than 0.7, the FD Imputer's scores are greater than 0. The two FDs for which the FD Imputer performs as well as the ML Imputer were identified and discussed in the previous sections.

4.2.3 Overfitting ML Imputer

In order to mimic the overfitting that takes place when FD Imputer looks up imputations on the train-set, it might be interesting to overfit ML Imputer models on purpose. As shown in figure 10 and figure 12, neither the classifier-models, nor the regression-models trained by ML Imputer are capable of overfitting. Future research-effort could be put into overfitting the models trained by ML Imputer with the goal of approximating the properties of FD Imputer.

4.3 Dependency Detection with DepDetector

The DepDetector algorithm was implemented to demonstrate dependency detection using models trained with ERM-strategies. The ‘greedy’ search-strategy was successfully used to detect minimal dependencies that are robust towards noise when used for training new models. The dependencies detected this way are a blend of RFDs as discussed in the theory-section of this work.

Dependencies found with DepDetector can be used to save computation time when training a network for imputation. Since minimal dependencies are often a subset of the whole relational instance, future imputation tasks can thus be performed more quickly by training models with less data.

Furthermore, minimal dependency detection helps explain the way a trained imputation-model works internally. By reducing the amount of columns with which a model is trained, the way imputations are derived can often be reconstructed by inspecting the remaining columns.

It was not possible to derive FDs using ERM-strategies. As discussed in the previous paragraphs, the behavior of FD Imputer might be approximated by overfitting models generated by ML Imputer. In a future effort to derive robust FDs with DepDetector, it might be suitable to run the DepDetector-algorithm training overfitted models on purpose.

Computing results, especially using the ‘complete’-strategy, is computationally expensive. Compared with FD detection algorithms, DepDetector is several orders of magnitude slower when detecting minimal dependencies. For the Chess-dataset, the ‘complete’-strategy takes about 10^5 times longer until it converges, the ‘greedy’-strategy takes approximately 10^4 times longer.

DepDetector was implemented as a proof of concept. There is a big margin for performance improvement. Better performance can be achieved by training models on GPUs

instead of CPUs. Leveraging asynchronous programming techniques, models on the same level of the search lattice can be trained in parallel. Furthermore, DepDetector internally models the search lattice as a tree rather than a lattice as in figure 5. This leads to redundant model training.

In general, the DepDetector approach is different from classical FD discovery algorithms in that minimal dependencies are detected solving an optimization problem on a graph. The score of a parent-node has direct implications for the performance of its child-nodes. Future research might investigate if this property is advantageous when searching for dependencies on big datasets.

5 Conclusion

In pursuit of identifying robust FDs, it is shown that FD Imputer can be used to measure robustness as defined in this work. Experiments with FD Imputer lead to the conclusion that only FDs with RHSs containing classifiable data can be robust. Generally, FDs are highly susceptible to noise in data.

A comparison between ML Imputer and FD Imputer suggests that imputation results yielded using FD Imputer are similar to using an overfitted model. It is not possible to use FDs to impute continuous numerical data.

Dependencies with learned relaxation on the attribute comparison were successfully detected on datasets using machine learning techniques with an algorithm called DepDetector. Two search-strategies were implemented and benchmarked on eight datasets.

References

- [Abe+19] Ziawasch Abedjan et al. *Data Profiling*. 2019. ISBN: 9781681734477. DOI: <https://doi.org/10.2200/S00878ED1V01Y201810DTM052>.
- [ASN14] Ziawasch Abedjan, Patrick Schulze, and Felix Naumann. "DFD: Efficient Functional Dependency Discovery". In: *CIKM*. 2014.
- [Bie+18] Felix Biessmann et al. "'Deep' Learning for Missing Value Imputation in Tables with Non-Numerical Data". In: *ACM International Conference on Information and Knowledge Management* 27 (2018), p. 9. DOI: <https://doi.org/10.1145/3269206.3272005>.
- [Bis06] Christopher M. Bishop. *Pattern Recognition and Machine Learning*. Springer, 2006.
- [Ble+16] Tobias Bleifuß et al. "Approximate Discovery of Functional Dependencies for Large Datasets". In: *Proceedings of the 25th ACM International on Conference on Information and Knowledge Management*. *CIKM '16*. Indianapolis, Indiana, USA: ACM, 2016, pp. 1803–1812. ISBN: 978-1-4503-4073-1. DOI: [10.1145/2983323.2983781](https://doi.org/10.1145/2983323.2983781). URL: <http://doi.acm.org/10.1145/2983323.2983781>.
- [Bra+08] Loreto Bravo et al. "Increasing the Expressivity of Conditional Functional Dependencies without Extra Complexity". In: *Proceedings - International Conference on Data Engineering* (May 2008), pp. 516–525. DOI: [10.1109/ICDE.2008.4497460](https://doi.org/10.1109/ICDE.2008.4497460).
- [Bur19] Andriy Burkov. *The Hundred-Page Machine Learning Book*. Andriy Burkov, 2019. ISBN: 978-1999579500.
- [CDP16] Loredana Caruccio, Vincenzo Deufemia, and Giuseppe Polese. "Relaxed Functional Dependencies - A Survey of Approaches". In: *IEEE Transactions on Knowledge and Data Engineering* (2016), pp. 147–165.
- [CFM09] Wenguang Chen, Wenfei Fan, and Shuai Ma. "Analyses and Validation of Conditional Dependencies with Built-in Predicates". In: (2009), pp. 576–591.
- [Chi92] Nancy Chinchor. "MUC-4 Evaluation Metrics". In: *Proceedings of the 4th Conference on Message Understanding*. *MUC4 '92*. McLean, Virginia: Association for Computational Linguistics, 1992, pp. 22–29. ISBN: 1-55860-273-9. DOI: [10.3115/1072064.1072067](https://doi.org/10.3115/1072064.1072067). URL: <https://doi.org/10.3115/1072064.1072067>.
- [Cod70] Edgar F Codd. "A relational model of data for large shared data banks". In: *Communications of the ACM* 13.6 (1970), pp. 377–387.
- [Cod90] Edgar F. Codd. *The Relational Model for Database Management: Version 2*. Addison-Wesley Publishing Company, 1990. ISBN: 0-201-14192-2. URL: <https://dl.acm.org/citation.cfm?id=77708>.
- [DG17] Dheeru Dua and Casey Graff. *UCI Machine Learning Repository*. 2017. URL: <http://archive.ics.uci.edu/ml>.
- [DHS00] Richard O. Duda, Peter E. Hart, and David G. Stork. *Pattern Classification (2nd Edition)*. Wiley-Interscience New York, 2000. ISBN: 0471056693.

- [Hay08] Simon Haykin. *Neural Networks and Learning Machines Third Edition*. Pearson Prentice Hall, 2008. ISBN: 9780131471399.
- [Hei+19] 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](https://doi.org/10.1145/3299869.3319888). URL: <http://doi.acm.org/10.1145/3299869.3319888>.
- [Huh+99] Ykä Huhtala et al. “TANE: An Efficient Algorithm for Discovering Functional and Approximate Dependencies”. In: *The Computer Journal* 42 (1999), pp. 100–111.
- [IBM03] Research News IBM. “Former IBM Fellow Edgar (Ted) Codd passed away on April 18”. In: (Apr. 23, 2003). URL: https://web.archive.org/web/20190425094107/https://www.ibm.com/ibm/history/exhibits/builders/builders_codd.html (visited on 06/26/2019).
- [Kle11] Stephan Kleuker. *Grudkurs Datenbankentwicklung*. Vieweg+Teubner Verlag, 2011. ISBN: 978-3-8348-9925-5. URL: <https://link.springer.com/book/10.1007/978-3-8348-9925-5>.
- [Kou+09] N. Koudas et al. “Metric Functional Dependencies”. In: (Mar. 2009), pp. 1275–1278. ISSN: 1063-6382. DOI: [10.1109/ICDE.2009.219](https://doi.org/10.1109/ICDE.2009.219).
- [Kra+18] 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](https://doi.org/10.1145/3183713.3196909). URL: <http://doi.acm.org/10.1145/3183713.3196909>.
- [Mai83] David Maier. *The Theory of Relational Databases*. Computer Science Pr, 1983. ISBN: 0914894420. URL: <http://web.cecs.pdx.edu/~maier/TheoryBook/TRD.html>.
- [Moo+11] David S. Moore et al. *The Practice of Statistics for Business and Economics*. W. H. Freeman and Company, 2011.
- [MX14] Xiao-Li Meng and Xianchao Xie. “I Got More Data, My Model is More Refined, but My Estimator is Getting Worse! Am I Just Dumb?” In: *Econometric Reviews* 33.1-4 (2014), pp. 218–250. DOI: [10.1080/07474938.2013.808567](https://doi.org/10.1080/07474938.2013.808567). eprint: <https://doi.org/10.1080/07474938.2013.808567>. URL: <https://doi.org/10.1080/07474938.2013.808567>.
- [NC08] Noel Novelli and Rosine Cicchetti. “FUN: An Efficient Algorithm for Mining Functional and Embedded Dependencies”. In: vol. 1973. Jan. 2008, pp. 189–203. DOI: [10.1007/3-540-44503-X_13](https://doi.org/10.1007/3-540-44503-X_13).
- [Pap+15] Thorsten Papenbrock et al. “Functional Dependency Discovery: An Experimental Evaluation of Seven Algorithms”. In: *Proc. VLDB Endow.* 8.10 (2015), pp. 1082–1093. ISSN: 2150-8097. DOI: [10.14778/2794367.2794377](https://doi.org/10.14778/2794367.2794377). URL: <https://doi.org/10.14778/2794367.2794377>.

- [PN16] Thorsten Papenbrock and Felix Naumann. “A Hybrid Approach to Functional Dependency Discovery”. In: SIGMOD ’16 (2016), pp. 821–833. doi: [10.1145/2882903.2915203](https://doi.org/10.1145/2882903.2915203). URL: <http://doi.acm.org/10.1145/2882903.2915203>.
- [Rij79] C. J. Van Rijsbergen. *Information Retrieval*. 2nd. Newton, MA, USA: Butterworth-Heinemann, 1979. ISBN: 0408709294.
- [Sch17] Edwin Schicker. *Datenbanken und SQL*. Springer Vieweg, 2017.
- [Stu16] Thomas Studer. *Relationale Datenbanken*. Springer Vieweg, 2016.
- [SV08] Alex Smola and S.V.N. Vishwanathan. *Introduction to Machine Learning*. Cambridge University Press, 2008. ISBN: 521825830. URL: [svn://smola@repos.stat.purdue.edu/thebook/trunk/Book/thebook.tex](http://svn.stat.purdue.edu/thebook/trunk/Book/thebook.tex).
- [Tha18] Alaa Tharwat. “Classification assessment methods”. In: *Applied Computing and Informatics* (2018). ISSN: 2210-8327. doi: <https://doi.org/10.1016/j.aci.2018.08.003>. URL: <http://www.sciencedirect.com/science/article/pii/S2210832718301546>.
- [Vap92] V. Vapnik. “Principles of Risk Minimization for Learning Theory”. In: *Advances in Neural Information Processing Systems 4*. Ed. by J. E. Moody, S. J. Hanson, and R. P. Lippmann. Morgan-Kaufmann, 1992, pp. 831–838. URL: <http://papers.nips.cc/paper/506-principles-of-risk-minimization-for-learning-theory.pdf>.
- [Wat14] Adrienne Watt. *Database Design*. BCcampus, 2014. URL: <https://opentextbc.ca/dbdesign01/>.
- [WMS08] Dennis D. Wackerly, William Mendenhall III, and Richard L. Scheaffer. *Mathematical Statistics with Applications, Seventh Edition*. Thomson Learning Inc., 2008. ISBN: 978-0-495-38508-0.
- [YHB02] Hong Yao, Howard Hamilton, and Cory Butz. “FD_Mine: Discovering Functional Dependencies in a Database Using Equivalences.” In: Jan. 2002, pp. 729–732. doi: [10.1109/ICDM.2002.1184040](https://doi.org/10.1109/ICDM.2002.1184040).

6 Appendix

```

1 def get_greedy_candidates(root):
2     convergence = True
3     threshold = root.most_recent_nodes()[0].parent.score
4     highscore_node = None
5
6     for node in most_recent_nodes():
7         if root.is_continuous: # sequential RHS
8             if (node.score <= 1.02*threshold) and (len(node.name) > 1):
9                 highscore_node = node
10        elif not root.is_continuous: # classifiable RHS
11            if (node.score >= 0.98*highscore) and (len(node.name) > 1):
12                highscore_node = node

```

```

13
14     if highscore_node is not None:
15         convergence = False
16         pot_lhs = highscore_node.name
17         for col in pot_lhs:
18             candidate_lhs = [c for c in pot_lhs if c != col]
19             add_node(candidate_lhs,
20                     parent=highscore_node,
21                     score=None)

```

Listing 2: ‘Greedy’ Candidate Generation in DepDetector

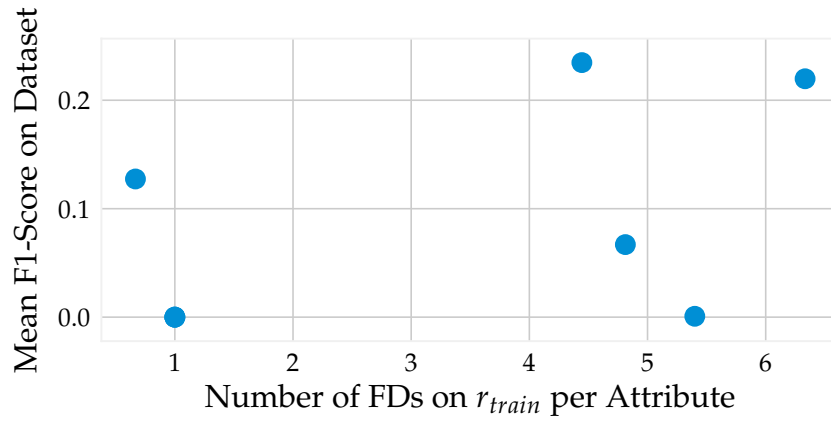


Figure 16: Mean robustness of all FDs with a classifiable RHS in function of the total number of FDs with a classifiable RHS divided by the total number of attributes in a dataset. No functional relation between these two values can be stated.