

# T-REx: Table Repair Explanations

Daniel Deutch

Tel Aviv University  
danielde@post.tau.ac.il

Nave Frost

Tel Aviv University  
navefrost@mail.tau.ac.il

Amir Gilad

Tel Aviv University  
amirgilad@mail.tau.ac.il

Oren Sheffer

Tel Aviv University  
orensheffer@mail.tau.ac.il

## Abstract

Data repair is a common and crucial step in many frameworks today, as applications may use data from different sources and of different levels of credibility. Thus, this step has been the focus of many works, proposing diverse approaches. To assist users in understanding the output of such data repair algorithms, we propose T-REx, *a system for providing data repair explanations through Shapley values*. The system is generic and not specific to a given repair algorithm or approach: it treats the algorithm as a black box. Given a specific table cell selected by the user, T-REx employs Shapley values to explain the significance of each constraint and each table cell in the repair of the cell of interest. T-REx then ranks the constraints and table cells according to their importance in the repair of this cell. This explanation allows users to understand the repair process, as well as to act based on this knowledge, to modify the most influencing constraints or the original database.

## 1 Introduction

Multiple previous works have proposed algorithms for data repair using Denial Constraints (DCs) [2] or subsets thereof [5, 8, 3, 1]. These approaches employ algorithms that use the constraints to detect and change values in a database table. We propose a *system that provides explanations for data repairs by presenting the influence of each constraint and table cell*. An explanation for such a repair may be useful both as means of understanding the repair process and algorithm, and as a tool for debugging the qual-

ity of the constraints for the repair of this specific data.

T-REx<sup>1</sup> is a novel system for data repair explanations based on *Shapley values* [6]. The notion of Shapley values was originally suggested in the context of Game Theory as a measure of quantifying the contribution of each player in a cooperative game. It was later adopted by the Machine Learning (ML) community as a tool for evaluating the contribution of each feature in the model [4]. Given a repaired cell, T-REx computes and presents the Shapley values of the DCs and table cells that have influenced this repair. Our approach evaluates the contribution of the input directly rather than the contribution of hidden features which are used by a specific algorithm. *This allows our solution to treat the repair algorithm as a black box and only query it to compute the Shapley values of DCs and cells*. Explanations for the influence of DCs on the repair may assist users in correcting them and adapting them to the specific data and repair algorithm, while explanations about the influence of data cells can help in understanding the repair algorithm itself and changing specific cells to make the repair more accurate.

**Example 1.1** Consider the table in Figure 2a and the DCs in Figure 1 with the Shapley values of each DC on its left. C1 says that two tuples that share a team value must be in the same city, C2 says that if a pair of tuples share a city, they must have the same country, C3 says that two tuples that have the same league must have the same country, and C4 says that it is impossible for two different teams of

<sup>1</sup>Please refer to the video of the system at <https://youtu.be/xPVWzHP0uAk>

- $\frac{1}{6}$ : (C1)  $\forall t_1, t_2. \neg(t_1[Team] = t_2[Team] \wedge t_1[City] \neq t_2[City])$   
 $\frac{1}{6}$ : (C2)  $\forall t_1, t_2. \neg(t_1[City] = t_2[City] \wedge t_1[Country] \neq t_2[Country])$   
 $\frac{2}{3}$ : (C3)  $\forall t_1, t_2. \neg(t_1[League] = t_2[League] \wedge t_1[Country] \neq t_2[Country])$   
0: (C4)  $\forall t_1, t_2. \neg(t_1[Team] \neq t_2[Team] \wedge t_1[Year] = t_2[Year] \wedge t_1[League] = t_2[League] \wedge t_1[Place] = t_2[Place])$

Figure 1: Denial constraints with their Shapley value

---

**Algorithm 1:** Simple Repair Algorithm

---

**Input** : Set of constraints  $\mathcal{C}$ , a dirty database table  $T^d$

1. If tuple  $t$  has a contradiction according to  $C1$  then the  $City$  attribute will be modified to the most common one, i.e.,  $\arg \max_c \mathbb{P}[City = c]$ .
  2. If tuple  $t$  has a contradiction according to  $C2$  then the  $Country$  attribute will be modified to the most probable one given  $t[City]$ . i.e.,  $\arg \max_c \mathbb{P}[Country = c \mid City = t[City]]$ .
  3. If tuple  $t$  has a contradiction according to  $C3$  then the  $Country$  attribute will be modified to the most common one, i.e.,  $\arg \max_c \mathbb{P}[Country = c]$ .
  4. If tuple  $t$  has a contradiction according to  $C4$  then the  $Place$  attribute will be modified to the most probable one given  $t[Team]$ , i.e.,  $\arg \max_p \mathbb{P}[Place = p \mid Team = t[Team]]$ .
- 

the same league to finish in the same place in the same year. Consider the cell  $Country$  in the fifth row, denoted by  $t_5[Country]$ . For simplicity, assume that we have Algorithm 1 as a naïve repair algorithm<sup>2</sup>. T-REx computes the contribution of each DC and ranks them accordingly, where  $C3$  is the most influential DC. It contributed the most as the League value “La Liga” appears in 3 other tuples coupled with the value “Spain” in the attribute  $Country$ .  $C1$  and  $C2$  each contributed equally as  $C1$  caused the change of “Capital” to “Madrid” first and then  $C2$  caused the change of the value in the  $Country$  cell.  $C4$  is not involved in the repair so its contribution is 0.

Next, we measure the influence of different data

---

<sup>2</sup>In practice, the repair algorithm may be more sophisticated; our solution is agnostic to the complexity of the repair algorithm.

cells on this repair. Given Algorithm 1, observe that the value of  $t_1[Place]$  has no influence on the modification of  $t_5[Country]$  – as  $t_1$  has no contradictions with  $t_5$ , and the attribute  $Place$  does not affect  $Country$  in Algorithm 1. However, how can we determine if  $t_5[League]$  was more or less influential on the repair compared to  $t_6[City]$ ? Intuitively,  $t_5[League]$  is more influential than  $t_6[City]$ . This is because if  $t_5[League]$  had a different value, then tuple  $t_5$  would not have any contradictions according to  $C3$ . While if  $t_6[City]$  had a different value, then according to  $C1$  there would have been a contradiction between  $t_3$  and  $t_6$  (as both tuples would have  $Team$  value of “Real Madrid”, and an inconsistent  $City$ ) which would have been resolved by Algorithm 1. As a result T-REx will assign higher contribution to  $t_5[League]$  compared to  $t_6[City]$ .

T-REx takes as input the algorithm itself and its input which is a set of DCs and a dirty database table. Another input to the system is a specific table cell of interest whose repair requires explaining. The system then ranks the influencing DCs and table cells based on their Shapley value for this cell of interest. Generally, computing the Shapley value is exponential time in the number of DCs/table cells, and thus T-REx employs different algorithms to compute the Shapley value for DCs and for table cells. With DCs, the naïve approach is feasible as the number of DCs is usually small. Conversely, the number of cells in a table can be very large, so T-REx uses a sampling algorithm based on [7]. To compute the Shapley values, the system repeatedly changes the input of the repair algorithm and queries it, so it does not rely on the components or approach of a specific algorithm.

## 2 Technical Details

We give a short overview of the approach underlying T-REx.

### 2.1 Database Repair

$T$  will denote a database table with schema  $(A_1, \dots, A_m)$  where  $A_i$  is the  $i$ th attribute of  $T$ . For

	Team	City	Country	Year	League	Place
$t_1^d$	F.C. Barcelona	Barcelona	Spain	2019	La Liga	1
$t_2^d$	Atletico Madrid	Madrid	Spain	2019	La Liga	2
$t_3^d$	Real Madrid	Madrid	Spain	2019	La Liga	2
$t_4^d$	F.C. Barcelona	Barcelona	Catalonia	2018	La Liga	1
$t_5^d$	Atletico Madrid	Capital	España	2018	La Liga	2
$t_6^d$	Real Madrid	Madrid	Spain	2018	La Liga	3

(a) Dirty table (red cells are dirty)

	Team	City	Country	Year	League	Place
$t_1^c$	F.C. Barcelona	Barcelona	Spain	2019	La Liga	1
$t_2^c$	Atletico Madrid	Madrid	Spain	2019	La Liga	2
$t_3^c$	Real Madrid	Madrid	Spain	2019	La Liga	3
$t_4^c$	F.C. Barcelona	Barcelona	Spain	2018	La Liga	1
$t_5^c$	Atletico Madrid	Madrid	Spain	2018	La Liga	2
$t_6^c$	Real Madrid	Madrid	Spain	2018	La Liga	3

(b) Clean table (blue cells have been repaired)

Figure 2: Input dirty table and output clean table for La Liga standings

a tuple  $t \in T$ , the notation  $t[A_i] = v$  means that  $t$  has the value  $v$  in attribute  $A_i$ . We denote by  $T^d$  and  $T^c$  the database table prior to the repair and after it respectively. Extending this,  $t^d[A]$  and  $t^c[A]$  will also be used to denote a dirty and clean cell, respectively.

**Example 2.1** Consider the dirty and clean tables shown in Figures 2a, 2b, referred to as  $T^c$  and  $T^d$ . If we consider  $t_5$  in both tables, then the attribute  $t_5^d[\text{Country}]$  in  $T^d$  is changed in  $T^c$  from the value “España” to “Spain”.

We denote the repair algorithm by  $Alg$  and its input by (1)  $C$ , a set of DCs and (2)  $T^d$ , a dirty table. Also, denote  $Alg(C, T^d) = T^c$  as the output table of  $Alg$ . For our purposes, we will refer to  $Alg$  as a binary function as follows. Given a table cell  $t^d[A] \in T^d$ , the repair algorithm is a function  $Alg|_{t^d[A]} : (C, T^d) \rightarrow \{0, 1\}$ , where 1 signals that the value in  $t^d[A]$  is repaired to the value in  $t^c[A]$ , and 0 otherwise.

**Example 2.2** Consider the cell  $t_5[\text{City}]$  in Figures 2a and 2b. Without  $C1$  it would not have changed from “Capital” to “Madrid”, therefore:  $Alg|_{t_5[\text{City}]}(\{C1, C2, C3\}, T^d) = 1$  while  $Alg|_{t_5[\text{City}]}(\{C2, C3\}, T^d) = 0$ .

## 2.2 Shapley Value

In Cooperative Game Theory, Shapley value [6] is a way to distribute the worth of all players, assuming they cooperate. Let  $N$  be a finite set of players and  $v : 2^N \rightarrow \mathbb{R}$ ,  $v(\emptyset) = 0$  be a function (called a characteristic function).  $v$  maps sets of players to the joint worth they generate according to the game. The Shapley value of a player  $a$  is then defined as:

$$Shap(N, v, a) = \sum_{S \subseteq N \setminus \{a\}} \frac{|S|!(|N| - |S| - 1)!}{|N|!} \cdot (v(S \cup \{a\}) - v(S))$$

In our scenario, the model is a black box so the Shapley values are computed on the input itself, i.e., the constraints and the table. For constraints, we adapt the definition so that it reflects the contribution of a specific constraint to the repair of a cell, as follows.

$$Shap(C, Alg|_{t^d[A]}, C) = \sum_{S \subseteq C \setminus \{C\}} \frac{|S|!(|C| - |S| - 1)!}{|C|!} \cdot (Alg|_{t^d[A]}(S \cup \{C\}, T^d) - Alg|_{t^d[A]}(S, T^d))$$

Where  $t^d[A]$  is a specific cell of interest and  $C$  is a constraint whose contribution we want to determine. The “set of players” is the set of DCs while the table  $T^d$  remains constant.

**Example 2.3** Recall the tables in Figure 2 with the DCs in Figure 1 (Shapley values are on the left) and Algorithm 1. We now compute the contribution of each DC to the repair of the cell  $t_5[\text{Country}]$ , denoted  $t_5[C]$ . Algorithm 1 will repair  $t_5[C]$  only if we have the DCs  $\{C1, C2\}$ , or  $\{C3\}$ . According to the definition, we can compute the contribution of  $C1$  as follows: there are 8 subset of  $\{C2, C3, C4\}$ , and only for  $S = \{C2\}$  and  $S = \{C2, C4\}$  we have  $Alg|_{t_5[C]}(S \cup \{C1\}, T^d) = 1$  and  $Alg|_{t_5[C]}(S, T^d) = 0$ , so  $Shapley(C, T^d, C1) = \frac{2}{12}$ . The same computation applies to  $C2$ . For  $C3$  we have 6 out of 8 subsets  $S$  of  $\{C1, C2, C4\}$  that result in  $Alg|_{t_5[C]}(S \cup \{C3\}, T^d) = 1$  and  $Alg|_{t_5[C]}(S, T^d) = 0$ , including  $S = \emptyset$ . Thus,  $Shapley(C, T^d, C3) = \frac{2}{3}$ . As for  $C4$ , its presence or absence does not change the value of  $t_5[C]$ , so  $Shapley(C, T^d, C4) = 0$ .

Let us explain the intuition for the value of  $C3$  being double that of the pair  $\{C1, C2\}$ . Ignore for now

$C4$  since its contribution is 0. There are 5 subsets of the DCs  $\{C1, C2, C3\}$  for which we repair  $t_5[C]$ . These are  $\{C3\}$ ,  $\{C1, C2\}$ ,  $\{C1, C3\}$ ,  $\{C2, C3\}$ , and  $\{C1, C2, C3\}$ . Four of these sets contain  $C3$  while only two contain the pair  $\{C1, C2\}$  (for the subsets where one of these is present without its partner, the repair is due to  $C3$ ), thus, the contribution of  $C1$  and  $C2$ , as a pair, is half that of  $C3$ .

Similarly, we adjust the definition for the Shapley value of a cell. Given a repair of cell  $t^d[A]$  we define the formula for calculating the Shapley value of a cell  $t_i[B]$ , or intuitively, its contribution to the repair of  $t^d[A]$ .

$$Shap(D, Alg|_{t^d[A]}, t_i[B]) = \sum_{S \subseteq T^d \setminus \{t_i[B]\}} \frac{|S|!(|T^d| - |S| - 1)!}{|T^d|!} (Alg|_{t^d[A]}(C, S \cup \{t_i[B]\}) - Alg|_{t^d[A]}(C, S))$$

Where  $S \subseteq T^d$  means  $\forall t_j[C] \in T^d \setminus S. t_j[C] = \text{null}$ . Here, the “set of players” here is the set of cells in the table  $T^d$  while the set of constraints remains constant.

**Example 2.4** Reconsider our example with the DCs from Figure 1, Algorithm 1, and the tables in Figure 2. Consider the cell  $t_5[\text{Country}]$  whose value is changed from “España” to “Spain”. Among all the cells,  $t_5[\text{League}]$  has the highest Shapley value, next we will explain why. Notice that based on  $C3$  the inclusion of  $t_5[\text{League}]$  to any coalition that contains at least one of the pairs  $\{t_i[\text{Country}], t_i[\text{League}]\}$  for any  $i \in \{1, 2, 3, 6\}$  would result in the repair of  $t_5[\text{Country}]$  to “Spain”. Observe that there are  $175 \cdot 2^{27}$  such coalitions (since out of the relevant 8 cells there are  $2^8 - 3^4 = 175$  options to choose a coalition such that at least one pair exists, and excluding those cells and  $t_5[\text{League}]$  there are 27 remaining cells that can be either included or excluded from the coalition). Next, we will estimate the number of coalitions that are required for the fix based on  $C1$  and  $C2$ . According to these DCs, a coalition that contains  $\{t_3[\text{Team}], t_3[\text{City}], t_3[\text{Country}], t_5[\text{Team}]\}$  is required. There are  $2^{32}$  such coalitions. Since  $175 \cdot 2^{27}$  is more than five times larger than  $2^{32}$  we conclude that  $t_5[\text{League}]$  has the highest influence on the repair of  $t_5[\text{Country}]$  from “España” to

“Spain”. For simplicity we overlooked the coalition sizes, though they too play a role in the evaluation of Shapley values.

## 2.3 Computing Shapley Values

Shapley values can be computed from the definition, but the computation time may be exponential. For constraints, we can use the formula directly as their number is typically small. However, the number of table cells can be huge. Therefore, we use a novel algorithm based on probabilistic sampling [7] to approximate the contribution of a table cell.

**Example 2.5** Reconsider the table in Figure 2a. Suppose we are interested in the effect of the cell  $t_5[\text{City}]$  on the repair of the cell  $t_5[\text{Country}]$ . We initialize a variable  $\varphi = 0$ . We vectorize the table to get the vector  $x_T = (t_1[\text{Team}], t_1[\text{City}], \dots, t_2[\text{Team}], \dots, t_6[\text{Place}])$ . To sample a cell coalition, we take a random permutation of  $x_T$  – the coalition is the set of all of the cells that precede  $t_5[\text{City}]$ . Values of cells that are not part of the coalition will be replaced with a sample value from their column distribution. Once the cell coalition was formed we generate two instances of vectorized tables: one with the original value of  $t_5[\text{City}]$ , and the second where the  $t_5[\text{City}]$  value is replaced with random value. We then compute the difference in the result of  $Alg|_{t_5[\text{Country}]}$  for these two instances and add it to  $\varphi$ . We repeat this  $m$  times and output  $\frac{\varphi}{m}$ .

## 3 System Overview

T-REx is implemented in Python 3.6 and an underlying database engine in PostgreSQL 10.6. Its web-based GUI was built using JavaScript, CSS and HTML. The three screens of the system are shown in Figure 3 and the general architecture of T-REx is shown in Figure 4. Users first input a database table and a set of DCs to the HoLoClean system (Figure 3a and the arrow 1 in Figure 4). HoLoClean [5] is a holistic data repair system, that supports DCs, among other forms of constraints, and repairs the input table based on a probabilistic model involving



Figure 3: T-REx User Interface

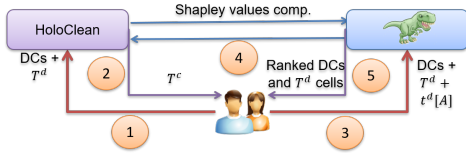


Figure 4: T-REx Architecture

machine learning techniques. After clicking the “Repair” button, users are presented with the repaired table, where repaired cells are highlighted (Figure 3b). Furthermore, when hovering over a repaired cell, the system shows its value before the repair. Now, T-REx allows users to choose any cell,  $t^d[A]$ , from the original table,  $T^d$ , whose value was changed, and mark it as a cell of interest and click the “Explain” button. The system then computes the Shapley values w.r.t. the chosen options by querying HoloClean as part of the computation. Once done, T-REx displays the DCs and table cells ranked from highest to lowest in terms of their Shapley value w.r.t.  $t^d[A]$ , where influencing DCs and cells are highlighted green and the darker the color, the more influencing the DC/cell is (Figure 3c). Again, when hovering over the DCs/cells users can also see their Shapley values. The user can continue the process by changing the DCs or values in  $T^d$ , and inputting it again to HoloClean to infer another repair, thus improving the repair iteratively.

## 4 Demo Scenario

Our demonstration will show that explaining repairs through Shapley values assists in understanding the repair process and debugging it. We will use a soccer database, scraped from Wikipedia, similarly to Fig-

ure 2a, and errors will be manually added into the table. We will start with an initial set of DCs. To get the repair, we will employ HoloClean that will output a clean table. Then, we will indicate a repaired cell of interest and show the most influential table cells and DCs involved in this repair, ranked according to their Shapley value. We will show how removing or changing the highest ranked DCs improves the repair of the specified table cell. We will use a similar scenario for table cells, where the DCs will be appropriate but some of the cells will cause a specific cell to be repaired in the wrong manner. After showing the obtained repair, we will invoke T-REx to rank the influencing table cells. We will then allow users to change values in the initial table and the DCs and choose different cells of interest to them. Users could then use T-REx to compute the Shapley value of the table cells and DCs that influenced the repair of their chosen cell and explore the system.

**Acknowledgements** This research has been funded by the European Research Council (ERC) under the European Union’s Horizon 2020 research and innovation programme (Grant agreement No. 804302), the Israeli Science Foundation (ISF) Grant No. 978/17, and the Google Ph.D. Fellowship. The contributions of Nave Frost and Amir Gilad are part of their respective Ph.D. thesis research conducted at Tel Aviv University.

## References

- [1] P. Bohannon, W. Fan, F. Geerts, X. Jia, and A. Kementsietsidis. Conditional functional dependencies for data cleaning. In *ICDE*, pages 746–755, 2007.
- [2] X. Chu, I. F. Ilyas, and P. Papotti. Discovering denial constraints. *PVLDB*, 6(13):1498–1509, 2013.
- [3] X. Chu, I. F. Ilyas, and P. Papotti. Holistic data cleaning: Putting violations into context. In *ICDE*, pages 458–469, 2013.
- [4] S. M. Lundberg and S. Lee. A unified approach to interpreting model predictions. In *NIPS*, pages 4765–4774, 2017.
- [5] T. Rekatsinas, X. Chu, I. F. Ilyas, and C. Ré. Holo-clean: Holistic data repairs with probabilistic inference. *PVLDB*, 10(11):1190–1201, 2017.
- [6] L. SHAPLEY. A value for n-person games. *Contributions to the Theory of Games*, (28):307–317, 1953.
- [7] E. Strumbelj and I. Kononenko. Explaining prediction models and individual predictions with feature contributions. *Knowl. Inf. Syst.*, 41(3):647–665, 2014.
- [8] M. Volkovs, F. Chiang, J. Szlichta, and R. J. Miller. Continuous data cleaning. In *ICDE*, pages 244–255, 2014.