# Value Veracity Estimation for Multi-Truth Objects via a Graph-Based Approach

Xiu Susie Fang[*], Quan Z. Sheng[*], Xianzhi Wang[†], Anne H.H. Ngu[‡]

[*]School of Computer Science, The University of Adelaide, SA 5005, Australia
{xiu.fang,michael.sheng}@adelaide.edu.au
[†]School of Computer Science and Engineering, UNSW Australia
xianzhi.wang@unsw.edu.au
[‡]Department of Computer Science, Texas State University, USA
angu@txstate.edu

## ABSTRACT

A fundamental issue with current *truth discovery* methods is that they generally assume only one true value for each object, while in reality objects may have multiple true values. We propose a graph-based approach, called *SmartMTD*, to relax this assumption in truth discovery. SmartMTD models and quantifies two types of source relations to estimate source reliability precisely and to detect malicious agreement among sources for multi-truth discovery. Two graphs are constructed based on the modeled source relations, which are further used to derive two aspects of source reliability via random walk computation.

## Keywords

Multi-Truth Discovery, Object Popularity, Copy Detection

## 1. INTRODUCTION

Considerable research efforts have been conducted to solve the *truth discovery* problem. Although these methods consider various factors to facilitate truth discovery, they commonly assume that each object has exactly one true value (i.e., *single-truth* assumption) [4, 2]. However, in the real world, multi-truth objects—such as the children of a person—widely exist. In this work, we study the problem of truth discovery for multi-truth objects, i.e., the multi-truth discovery (MTD) problem. We propose a graph-based model, called *SmartMTD*, as an overall solution, which incorporates two important implications, namely *source relations* and *object popularity*. We propose to model the two-sided relations among sources, based on which we construct graphs to capture source features. Specifically, source authority features and two-sided source precision are captured by $\pm$*supportive agreement graph*s, while source dependence degrees are quantified by $\pm$*malicious agreement graph*s. Random walk computation is applied on both graphs to estimate source re-

liability and independence degrees. We further propose to differentiate source reliability by the popularities of objects, to minimize the number of people misguided by false values.

## 2. THE SMARTMTD APPROACH

Given a set of multi-truth objects ($\mathcal{O}$), conflicting values $\mathcal{V}$ can be collected from a set of sources ($\mathcal{S}$). We denote the set of all values of an object $o$ provided by all sources in $\mathcal{S}$ as $\mathcal{U}_o$, the set of values provided by a source $s$ on $o$ as $\mathcal{V}_{s_o}$ (i.e., *positive claims*). By incorporating the *mutual exclusion assumption*, $s$ is believed to implicitly invalidate all the other values on $o$, the disclaimed values are denoted as $\tilde{\mathcal{V}}_{s_o}$ (i.e., *negative claims*), which is calculated by $\mathcal{U}_o - \mathcal{V}_{s_o}$. The goal of SmartMTD is to identify a set of true values ($\mathcal{V}_o^*$) from $\mathcal{V}$, for each object $o$, satisfying that $\mathcal{V}_o^*$ is as close to the ground truth $\mathcal{V}_o^g$ as possible, while estimating two aspects of source reliability, namely *positive precision* ($\tau(s)$, the probability of the positive claims of a source being true), and *negative precision* ($\tilde{\tau}(s)$, the probability of the negative claims of a source being false). The perfect truth discovery result satisfies $\mathcal{V}_o^* = \mathcal{V}_o^g$.

Intuitively, if the positive (resp., negative) claims of a source are agreed by most other sources, this source is likely to have high positive (resp., negative) precision. This means that the inter-source agreements (i.e., the common values claimed or disclaimed by two sources) indicate source reliability endorsement. This motivates us to measure sources' positive (resp., negative) precision by quantifying their +agreements (resp., –agreements), i.e., the agreements among sources regarding their positive (resp., negative) claims. In reality, sources may not only support one another by providing the same true claims, but also maliciously copy from others to provide the same false claims, which sometimes mislead the audience. Therefore, we identify two types of source relations. Specifically, sharing the same true values means one source supports/endorses the other source and indicates a *supportive relation* between two sources. We define the common values between sources as *supportive agreements* and measure source reliability by quantifying such supportive agreements. In contrast, if two sources share a significant number of false values, they are likely to copy from each other, indicating a *copying relation*. We define the sharing of false values as *malicious agreements*. In addition, the significance of knowing the truths of different objects may vary in reality. For example, the phone number of a restaurant

Table 1: Comparison of different methods: the best and second best performance values are in bold.

| Method | Book-Author Dataset | | | | | | | Parent-Children Dataset | | | | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | P | R | F1 | WP | WR | WF1 | T(s) | P | R | F1 | WP | WR | WF1 | T(s) |
| Voting | **0.84** | 0.63 | 0.72 | **0.83** | 0.64 | 0.72 | **0.07** | 0.88 | 0.85 | 0.87 | 0.69 | 0.68 | 0.69 | **0.56** |
| Sums | **0.84** | 0.64 | 0.73 | **0.83** | 0.64 | 0.72 | 0.85 | **0.90** | 0.89 | **0.90** | **0.88** | 0.86 | 0.87 | 1.13 |
| Avg-Log | **0.83** | 0.60 | 0.70 | **0.83** | 0.64 | 0.72 | 0.61 | **0.90** | 0.89 | 0.89 | **0.88** | 0.86 | 0.87 | **0.75** |
| TruthFinder | **0.84** | 0.60 | 0.70 | **0.83** | 0.60 | 0.70 | 0.74 | **0.90** | 0.89 | **0.90** | **0.88** | 0.85 | 0.86 | 1.24 |
| 2-Estimates | 0.81 | 0.70 | 0.75 | 0.80 | 0.68 | 0.74 | **0.38** | **0.91** | 0.89 | **0.90** | **0.88** | 0.86 | 0.87 | 1.34 |
| LTM | 0.82 | 0.65 | 0.73 | **0.82** | 0.62 | 0.71 | 0.98 | 0.87 | 0.90 | 0.88 | 0.86 | 0.89 | 0.87 | 0.99 |
| MBM | **0.83** | **0.74** | **0.78** | **0.82** | **0.71** | **0.76** | 0.67 | **0.90** | **0.92** | **0.91** | 0.87 | **0.90** | **0.88** | 2.17 |
| SmartMTD | **0.83** | **0.75** | **0.79** | **0.83** | **0.78** | **0.80** | 0.43 | **0.90** | **0.93** | **0.91** | **0.93** | **0.92** | **0.93** | 0.92 |

is more frequently used by customers and thus has bigger impact than the year when it is opened. Taking object popularity into consideration could help us better model MTD.

Algorithm 1 shows the procedure of SmartMTD. First, parameters including the iteration convergence threshold $\delta$, smoothing factor $\beta$, positive precision $pp_{max}$, negative precision $np_{max}$, two-sided dependence scores ($pc_{max}$ and $nc_{max}$) of sources with the highest visit probabilities in ±supportive agreement graphs and ±malicious agreement graphs are initialized, where both the confidence scores of each value $v$ being true or false ($\mathcal{C}_v$, $\mathcal{C}_{\bar{v}}$) are initialized by majority voting. The algorithm checks the *cosine similarities* of the two-sided source precision obtained by two successive iterations to determine whether it converges—the algorithm only terminated when the difference between such cosine similarities derived from two successive iterations becomes smaller than a certain threshold $\delta$.

## 3. EXPERIMENTS

We compared SmartMTD with two types of baselines: i) methods under single-truth assumption, including *Voting*, *Sums*, *Average-Log*, *TruthFinder*, and *2-Estimates*, ii) MTD methods, *LTM* [4] and *MBM* [2], on two real-world datasets, namely *book-author dataset* [3] and *Parent-Children dataset* (extracted from the *Biography dataset* [1]). All methods were evaluated in terms of precision (P), recall (R), $F_1$ score (F1), and execution time (T). Since we introduce a new concept of object popularity, we used object popularity weighted precision (WP), recall (WR) and $F_1$ score (WF1) as additional metrics. Table 1 shows the comparison results. SmartMTD consistently achieved the best results on all metrics except precision, on which SmartMTD still showed the second best performance on the experimental datasets. Among the three methods specially designed for MTD, our approach is the most efficient with the lowest execution time. This is because LTM includes complicated Bayesian inference over the complex probabilistic graphical model, and MBM conducts time-consuming copy detection, while our approach is based on a relatively simple graph model.

## 4. CONCLUSION

We propose a graph-based approach, *SmartMTD*, which incorporates two important concepts, *source relations* and *object popularity*, for MTD. In particular, we construct *supportive agreement graph*s to model the endorsement among sources, from which to derive two-sided source reliability, and *malicious agreement graph*s to capture copying relations among sources. We also consider and develop techniques to quantify object popularity based on object occurrences and source coverage. Empirical studies show the effectiveness of SmartMTD. Our future work involves exploring more implications such as the long-tail phenomenon on source coverage and source confidence on claims to improve SmartMTD.

---

**Algorithm 1:** The Algorithm of SmartMTD.

**Input:** $\mathcal{O}$, $\mathcal{S}$, and $\mathcal{V}$.
**Output:** $\mathcal{V}_o{}^*$.

1 Initialize $\delta$, $\beta$, $pp_{max}$, $np_{max}$, $pc_{max}$, $nc_{max}$
2 Initialize $\mathcal{C}_v$, $\mathcal{C}_{\bar{v}}$ for each $v \in \mathcal{V}$, $o \in \mathcal{O}$

// Object popularity quantification

3 **foreach** $o \in \mathcal{O}$ **do**
4     $\mathcal{P}_o = \sum_{s \in \mathcal{S}_O} \frac{1}{Cov(s)}$, $\mathcal{P}_o$ is the popularity degree of $o$, $\mathcal{S}_O$ is the set of sources provide values on $o$, $Cov(s)$ is the percentage of $s$'s provided objects over $\mathcal{O}$

5 **repeat**

    // Malicious agreement detection

6     **foreach** $o \in \mathcal{O}$ **do**
7         construct ±malicious agreement graphs for $\mathcal{S}_O$ (each node is a source belongs to $\mathcal{S}_O$) by quantifying the weights of each edge by
$$\omega_{c_O}(s_1 \rightarrow s_2) = \beta + (1-\beta) \cdot \frac{|A_o(s_1,s_2)|}{|\mathcal{V}_{s_{2_O}}|} \cdot (1 - \prod_{v \in A_o(s_1,s_2)} \mathcal{C}_v)$$
and $\tilde{\omega}_{c_O}(s_1 \rightarrow s_2) =$
$$\beta + (1-\beta) \cdot \frac{|\tilde{A}_o(s_1,s_2)|}{|\mathcal{V}_{s_{2_O}}|} \cdot (1 - \prod_{v \in \tilde{A}_o(s_1,s_2)} \mathcal{C}_{\bar{v}}), \text{ where}$$
$A_o(s_1, s_2) = \mathcal{V}_{s_{1_O}} \cap \mathcal{V}_{s_{2_O}}$, $\tilde{A}_o(s_1, s_2) = \tilde{\mathcal{V}}_{s_{1_O}} \cap \tilde{\mathcal{V}}_{s_{2_O}}$,
8         derive $\mathcal{D}(s,o)$, $\tilde{\mathcal{D}}(s,o)$ by applying random walk, where $\mathcal{D}(s,o)$ (resp., $\tilde{\mathcal{D}}(s,o)$) is the dependence score of $s$ providing positive (resp., negative) claims on $o$
9         compute $I(s,o)$ (resp., $\tilde{I}(s,o)$) by $1 - \mathcal{D}(s,o)$ (resp., $1 - \tilde{\mathcal{D}}(s,o)$), where $I(s,o)$ (resp., $\tilde{I}(s,o)$) is the independence score of $s$ providing positive (resp., negative) claims on $o$

    // ±Source Reliability computation

10     construct ±supportive agreement graphs (each node is a source belongs to $\mathcal{S}$) by quantifying the weights of each edge by $\mathcal{A}(s_1, s_2) =$
$$\sum_{o \in \mathcal{O}_{s_1} \cap \mathcal{O}_{s_2}} \frac{|A_o(s_1,s_2)|}{|\mathcal{V}_{s_{2_o}}|} \cdot (1 - \prod_{v \in A_o(s_1,s_2)} \mathcal{C}_v) \cdot \mathcal{P}_o \cdot I(s_1, o),$$
$$\omega(s_1 \rightarrow s_2) = \beta + (1-\beta) \cdot \frac{\mathcal{A}(s_1,s_2)}{|\mathcal{O}_{s_1} \cap \mathcal{O}_{s_2}|} \text{ and } \tilde{\mathcal{A}}(s_1,s_2) =$$
$$\sum_{o \in \mathcal{O}_{s_1} \cap \mathcal{O}_{s_2}} \frac{|\tilde{A}_o(s_1,s_2)|}{|\mathcal{V}_{s_{2_o}}|} \cdot (1 - \prod_{v \in \tilde{A}_o(s_1,s_2)} \mathcal{C}_{\bar{v}}) \cdot \mathcal{P}_o \cdot \tilde{I}(s_1, o),$$
$$\tilde{\omega}(s_1 \rightarrow s_2) = \beta + (1-\beta) \cdot \frac{\tilde{\mathcal{A}}(s_1,s_2)}{|\mathcal{O}_{s_1} \cap \mathcal{O}_{s_2}|}$$
11     derive $\tau(s)$, $\tilde{\tau}(s)$ by applying random walk

    // Value confidence score computation

12     **foreach** $v \in \mathcal{V}$, $o \in \mathcal{O}$ **do**
13         compute $\mathcal{C}_v$, $\mathcal{C}_{\bar{v}}$ by $\mathcal{C}_v = \frac{\sum_{s \in \mathcal{S}_v} \tau'(s) + \sum_{s \in \mathcal{S}_{\bar{v}}} (1 - \tilde{\tau}'(s))}{|\mathcal{S}_o|}$ and
$$\mathcal{C}_{\bar{v}} = \frac{\sum_{s \in \mathcal{S}_v} (1 - \tau'(s)) + \sum_{s \in \mathcal{S}_{\bar{v}}} \tilde{\tau}'(s)}{|\mathcal{S}_o|}, \text{ where } \mathcal{S}_v \text{ (resp., } \mathcal{S}_{\bar{v}})$$
is the set of sources claim (resp., disclaim) $v$ on $o$

14 **until** *convergence*;
15 **return** $\{(o,v)|v \in \mathcal{V} \wedge \mathcal{C}_v > \mathcal{C}_{\bar{v}} \wedge v \in \mathcal{U}_o, o \in \mathcal{O}\}$

---

## 5. REFERENCES

[1] J. Pasternack et al. Knowing what to believe (when you already know something). In *Proc. Intl. Conference on Computational Linguistics (COLING)*, pages 877–885, 2010.

[2] X. Wang et al. An integrated bayesian approach for effective multi-truth discovery. In *Proc. the 24th ACM Intl. Conference on Information and Knowledge Management (CIKM)*, pages 493–502, 2015.

[3] X. Yin, J. Han, and P. S. Yu. Truth discovery with multiple conflicting information providers on the web. *IEEE Transactions on Knowledge and Data Engineering (TKDE)*, 20(6):796–808, 2008.

[4] B. Zhao et al. A bayesian approach to discovering truth from conflicting sources for data integration. *Proc. the VLDB Endowment*, 5(6):550–561, 2012.