Probabilistic Logic Rule Learning from Small Imperfect Open-World Datasets using a Dempster-Shafer Theoretic Approach

Anonymous Author(s) Affiliation Address email

Abstract

We study the problem of learning *epistemic uncertainty measures* for probabilistic 1 logic rules from *small imperfect datasets*. While Bayesian approaches have had 2 tremendous success in learning probabilistic parameters for rules from complex 3 relational data, we lack good methods for handling small and incomplete datasets, 4 with imprecise and probabilistic data instances containing mutually-dependent 5 attributes, being obtained from multiple heterogeneous sources in the open world 6 where new attributes are introduced ad hoc. We propose a Dempster-Shafer ap-7 proach to address these challenges. 8

9 1 Introduction

There has been a growing interest in combining logic-based representations with probabilistic 10 reasoning mechanisms and machine learning techniques [1]. Under the rubrics of probabilistic 11 logic learning (PLL) [2] and statistical relational learning (SRL) [3] several approaches combine 12 probabilistic mechanisms (e.g., Bayesian Networks, Markov Networks, Stochastic Grammars), 13 with logical representation schemes (propositional logic, first-order logic), and machine learning 14 15 techniques that allow for automated learning of probabilistic parameters or relational structure from data. Some popular PLL and SRL paradigms include Bayesian Logic Programs (BLP), PRISM, ICL, 16 LPADs, ProbLog, P-Log, CP-Logic, PITA, Markov Logic Networks (MLN), Probabilistic Relational 17 Models (PRM), Bayesian Logic Networks (BLN) and Relational Dependency Networks (RDN). 18 These approaches adopt an underlying Bayesian probability framework expressed graphically as 19

Bayesian networks or Markov Networks and can learn probabilistic weights from data [4]. The Bayesian setting is an intuitive one that already has a number of off-the-shelf tools to make inferences, learn parameters and compute estimates relatively efficiently. While Bayesian approaches have enjoyed considerable success and shown potential in handling large datasets having a fairly complex relational structure such as NELL [5], there is little work on applying these approaches to learning from *small imperfect open-world datasets*.

Small imperfect open-world datasets can come from multiple distinct heterogeneous sources of 26 27 varying levels of capability and reliability, producing streams of incomplete, ignorant instances for known and unknown attributes that may be mutually dependent on each other. This type of data is 28 frequently encountered in autonomous agents and perceptual systems that learn normative behavior 29 or must learn in contextually-charged environments. Bayesian techniques are not well-suited in these 30 situations and can lead to non-intuitive results. Instead, we propose an alternative approach based on 31 Dempster-Shafer (DS) Theory of Belief Functions [6] together with an algorithm for automatically 32 learning belief-theoretic logic rules from small imperfect datasets in open-world contexts. 33

Submitted to 31st Conference on Neural Information Processing Systems (NIPS 2017). Do not distribute.

34 2 Dempster-Shafer Theory Background

DS-Theory is a measure-theoretic mathematical framework that allows for combining pieces of 35 uncertain evidential information to produce degrees of belief for the various events of interest. It 36 has been extensively used in sensor fusion networks, object tracking, and network security [7–9]. 37 In DS-Theory a set of elementary events of interest is called Frame of Discernment (FoD). The 38 FoD is a finite set of mutually exclusive events $\Theta = \{\theta_1, ..., \theta_N\}$. The power set of Θ is denoted 39 by $2^{\Theta} = \{A : A \subseteq \Theta\}$. Each set $A \subseteq \Theta$ has a certain weight, or *mass* associated with it. A *Basic Belief Assignment* (BBA) is a mapping $m_{\Theta}(\cdot) : 2^{\Theta} \to [0, 1]$ such that $\sum_{A \subseteq \Theta} m_{\Theta}(A) = 1$ 40 41 and $m_{\Theta}(\emptyset) = 0$. The BBA measures the support assigned to the propositions $A \subseteq \Theta$ only. The 42 subsets of A with non-zero mass are referred to as *focal elements* and comprise the set \mathcal{F}_{Θ} . The triple 43 $\mathcal{E} = \{\Theta, \mathcal{F}_{\Theta}, m_{\Theta}(\cdot)\}$ is called the *Body of Evidence* (BoE). For ease of reading, we sometimes omit 44 \mathcal{F}_{Θ} when referencing the BoE. Given a BoE $\{\Theta, \mathcal{F}_{\Theta}, m_{\Theta}(\cdot)\}$, the *belief* for a set of hypotheses A is 45 $Bel(A) = \sum_{B \subseteq A} m_{\Theta}(B)$. This belief function captures the total support that can be committed to A 46 without also committing it to the complement A^c of A. The *plausibility* of A is $Pl(A) = 1 - Bel(A^c)$. 47 Thus, Pl(A) corresponds to the total belief that does not contradict A. The uncertainty interval of 48 A is [Bel(A), Pl(A)], which contains the true probability P(A). In the limit case with no uncertainty, 49 we get Pl(A) = Bel(A) = P(A). 50

DS-Theory extends Bayesian theory in several ways. First, it allows for assigning probabilistic 51 measures to sets of these hypotheses (not just individual ones), including the set of all hypothesis. 52 This allows DS-Theory to consider ignorant and ambiguous information. Second, DS-theory does not 53 require assuming any prior distributions, which is useful when priors are difficult to justify, as is the 54 case with many open-world sensing and perception tasks. Third, DS-theoretic uncertainty generally 55 refers to epistemic uncertainty and corresponds to beliefs held by agents about the world. However, 56 probability theoretic uncertainty often refers to aleatory uncertainty as it relates to frequency of 57 occurrence, randomness and chance. Bayesian and DS-theories do share many commonalities and 58 DS-theory is often viewed as being a generalization of Bayesian theory. 59

The history of DS-Theory is not without controversy and has been criticized by some including Judea Pearl [10], who subsequently recalled this criticism [11]. Nevertheless, major strides have been made to address these concerns including approximation algorithms for reducing the time complexity of computation [12], decision theoretic aspects [13], graphical models [14], approaches to resolve conflicts that arose from Dempster's original rule of combination [15], and the development of DS-theoretic logical operators [16, 17].

One recent development in DS-theory, is an evidence filtering strategy that has upgraded Dempster's original rule of combination of evidence to accommodate the inertia of available evidence and address some challenges with respect to conflicting evidence [18]. In particular consider the BoEs $\mathcal{E}_1 = \{\Theta, \mathcal{F}_1, m_1(\cdot)\}$ and $\mathcal{E}_2 = \{\Theta, \mathcal{F}_2, m_2(\cdot)\}$, and a given $A \in \mathcal{F}_2$. The updated belief (from iteration t to t + 1) $Bel_{t+1} : 2^{\Theta} \to [0, 1]$ and the updated plausibility $Pl_{t+1} : 2^{\Theta} \to [0, 1]$ of an arbitrary proposition $B \subseteq \Theta$ are ¹:

$$Bel(B)_{t+1}^{\mathcal{E}_1} = \mu_t \cdot Bel(B)_t^{\mathcal{E}_1} + \nu_t \cdot Bel(B|A)_t^{\mathcal{E}_t}$$
$$Pl(B)_{t+1}^{\mathcal{E}_1} = \mu_t \cdot Pl(B)_t^{\mathcal{E}_1} + \nu_t \cdot Pl(B|A)_t^{\mathcal{E}_2}$$

where $\mu_t, \nu_t \ge 0, \mu_t + \nu_t = 1$. The conditional in the above equations are Fagin-Halpern conditionals which can be considered an extension of Bayesian conditional notions [19]. That is for a BoE $\mathcal{E} = \{\Theta, \mathcal{F}, m(\cdot)\}, A \subseteq \Theta$ and an arbitrary $B \subseteq \Theta$, the conditional beliefs and plausibility are given by:

$$Bel(B|A)^{\mathcal{E}} = Bel(A \cap B)^{\mathcal{E}} / [Bel(A \cap B)^{\mathcal{E}} + Pl(A \setminus B)^{\mathcal{E}}]$$
$$Pl(B|A)^{\mathcal{E}} = Pl(A \cap B)^{\mathcal{E}} / [Pl(A \cap B)^{\mathcal{E}} + Bel(A \setminus B)^{\mathcal{E}}]$$

66 We build on this and other developments to provide a unified probabilistic logic learning framework,

67 grounded on a Belief-Theoretic approach.

¹We specify the BoE superscript for Bel() and Pl() as needed to be precise, especially when we are combining two distinct BoEs.

3 **Proposed Belief-Theoretic Approach and its Unique Properties** 68

3.1 Rule System 69

Consider a propositional alphabet \mathcal{L} , in which we have all the standard symbols (variables, predicates, 70

functions) and logical connectives. In this alphabet, we define a belief theoretic rule, as follows: 71

Definition 1 (Belief-Theoretic Rule System). A belief-theoretic rule is an expression of the form: $\mathcal{R} := [\alpha, \beta] :: \psi \implies (\neg)\phi$

- where the ψ, ϕ are Uncertain Logic atoms, i.e., propositional formulas with an associated "uncer-72
- tainty interval" $[\alpha, \beta]$ defined under Dempster-ShaferTheory [20] as $\alpha = Bel(\mathcal{R}), \beta = Pl(\mathcal{R})$ with 73
- $0 \le \alpha \le \beta \le 1$. Thus, under the present formulation both the atoms and the rules have uncertainty 74
- intervals.² The (\neg) indicates that negation is optional in this rule. A Belief-Theoretic Rule System \mathcal{T} 75
- is a finite set of belief-theoretic rules \mathcal{R} .³ 76

Example 1 Consider an agent reasoning about actions it can perform in a car and in a house. We 77 can represent this scenario as a Belief-Theoretic Rule System, \mathcal{T} , as follows: 78

 $\begin{array}{ll} \mathcal{R}_1 := [0.8, 0.95] :: inCar \implies driving \\ \mathcal{R}_2 := [0.9, 1] :: inCar \implies texting \\ \mathcal{R}_3 := [0.9, 1] :: inHouse \implies texting \end{array} \begin{array}{ll} \mathcal{R}_4 := [0, 0.3] :: inHouse \implies running \\ \mathcal{R}_5 := [0.3, 0.6] :: inHouse \implies smoking \end{array}$ 79

The rules in this example have intuitive semantics, whereby the antecedents correspond to context and 80 the consequents correspond to actions taken. The location of the uncertainty interval between 0 and 1 81 suggests the degree of truth and falsity for the rule and the width of the interval generally suggests the 82 level of support or evidence for that rule. So, rules $\mathcal{R}_1, \mathcal{R}_2$, and \mathcal{R}_3 have tight uncertainty intervals 83 close to 1 indicating a confident support for their truth. 84

The rule system displays the agent's current level of belief or epistemic uncertainty about a certain 85 set of rules that are influenced by various pieces of evidence. However, it is typically the case that 86 87 the agent is not updating its beliefs about all the rules in a rule-system simultaneously, but instead it is considering only a certain subset that might pertain, for example, to a particular context that the 88 agent is in, and for which evidence arrives together. For instance, when an agent is collecting data by 89 observing cars, they may observe multiple actions being performed together: driving, texting, talking 90 etc. In this setting the agent might need to only consider those actions that are relevant in the context 91 at the time of data collection, i.e., those relevant to the "in car" context. 92

To formalize these intuitions, consider the rule system \mathcal{T} of Example 1. The subset of rules from the 93

rule system relevant to the context inCar is $\{\mathcal{R}_1, \mathcal{R}_2\} \subset \mathcal{T}$. We model the arriving evidence as well as the agent's growing body of beliefs as DS-theoretic frames of discernment $\Theta_{\mathcal{T}}^{inCar}$ comprising all possible combinations of the rule consequents (and negations) present in the selected subset of rules: 94

95

96 $\Theta_{\tau}^{inCar} = \{(texting, driving), (texting, \neg driving), (\neg texting, driving), (\neg texting, \neg driving)\}$

Modeling the frame of discernment in this exhaustive way ensures that elementary events in the 97 frame are mutually exclusive of each other. The subset of rules associated with this frame is called 98 a Rule frame $\mathcal{R}_{\mathcal{T}}^{inCar} = \{\mathcal{R}_1, \mathcal{R}_2\}$. Now, if we are interested in measuring the amount of support 99 in favor of, say rule \mathcal{R}_1 , based on our evidence that is captured in the frame $\Theta_{\mathcal{T}}^{inCar}$, we would 100 measure $Bel(\{(texting, driving), (\neg texting, driving)\})$ as this captures the level of support for 101 just *driving* irrespective of *texting*, from a body of evidence that captures both. 102

To generalize, a Rule Frame $\mathcal{R}^{\psi}_{\mathcal{T}}$ is a set of rules in rule system \mathcal{T} that share the same antecedent 103 ψ . Similarly, we can generalize the frame of discernment, Θ^{ψ}_{T} , by first defining a DS-theoretic 104 elementary event θ as a tuple of all the rule consequents ϕ (or their negations) present in a rule 105 frame $\mathcal{R}^{\psi}_{\mathcal{T}}$. A set of elementary events forms an indexed frame of discernment $\Theta^{\psi}_{\mathcal{T}}$. Defining an elementary event in this way allows us to represent, exhaustively, all possible combinations of the set 106 107 of consequents ϕ_1, \ldots, ϕ_k and their negations. 108

²The concepts presented in this paper are not limited to a propositional logic and are extendable to a first-order language. The proposed formulation does not preclude the possibility of multiple consequents and antecedents combined with logical operators.

³Rules \mathcal{R} in the rule system \mathcal{T} differ from each other either by antecedent and/or consequent.

109 3.2 Learning Uncertainty Intervals for Rules from Data

Data Format. For a rule system \mathcal{T} with n rules, consider an indexed FoD $\Theta_{\mathcal{T}}^{\psi}$ and a corresponding 110 rule frame $\mathcal{R}^{\psi}_{\mathcal{T}}$ comprising k rules. Consider a set $S = \{s_1, \ldots, s_m\}$ of m evidence sources. Let 111 the set of evidence sources provide a set of BoEs, defined as $\mathbb{E} = \{\mathcal{E}_1^{\Theta_{\mathcal{T}}^{\psi}}, \dots, \mathcal{E}_m^{\Theta_{\mathcal{T}}^{\psi}}\}$. Each BoE is a 112 DS-theoretic BoE \mathcal{E} and is associated with an indexed FoD $\Theta_{\mathcal{T}}^{\psi}$. For simplicity, we will assume that 113 all the BoEs in \mathbb{E} correspond to the same indexed frame $\Theta^{\psi}_{\mathcal{T}}$. That is, the sources S provide evidence 114 for the same rule frame $\mathcal{R}^{\psi}_{\mathcal{T}}$. We can then define an observation, the set $O_i = \{o_{i,1}, \ldots, o_{i,k}\}$, made 115 by a source s_i as a form of truth assignment where each $o_{i,j} \in \{0, 1, \epsilon\}$ indicates whether the source 116 observes, for a given antecedent ψ , whether a certain consequent ϕ_j , $1 \le j \le k$ is true (1), false (0) 117 or unknown (ϵ). For instance, an observation that a person in a car is texting, but not driving can 118 be represented as $O_i = \{1, 0\}$. We can combine the observation O_i with other information about 119 the source as well as a DS-theoretic mass assignment to form a data instance, defined as follows. 120 A data instance is a tuple $d = (s_i, O_i, m_{\Theta_{\tau}^{\psi}}(\cdot)_{s_i})$ comprising a specific source identifier $s_i \in S$, an 121 observation O_i , and a DS-theoretic mass assignment $m_{\Theta_{\mathcal{T}}^{\psi}}(\cdot)_{s_i}$ for source s_i per BoE, $\mathcal{E}_i^{\Theta_{\mathcal{T}}^{\psi}}$ provided 122 by that source. A dataset $\mathcal{D} = \{d_1, \dots, d_n\}$ is a finite set of n data instances. 123

Learning Problem. The learning problem can be defined as follows: Given an "unspecified" rule frame ${}^{\mathcal{R}}_{\mathcal{T}}^{\psi}$ ("-" suggests that parameter values are unspecified) with k rules, and a dataset \mathcal{D} , compute the parameters of the rule frame $\alpha_1, \ldots, \alpha_k, \beta_1, \ldots, \beta_k$.

Learning Algorithm. The rule learning algorithm (Algorithm 1) assigns uncertainty parameters to 127 each rule, updating those values as it considers each new data instance. Algorithm 1, displayed below, 128 achieves this form of rule learning. The algorithm iterates though each data instance d in the data set 129 \mathcal{D} (line 5) and, per instance, through each rule \mathcal{R} in the rule frame $\mathcal{R}^{\psi}_{\mathcal{T}}$ (line 6). For each iteration, we 130 first set the hyper-parameters μ and ν (line 7) that specify how much weight the algorithm will place 131 on previous learned knowledge (μ) and on each new data instance (ν). These hyper-parameters are 132 then used to compute a conditional belief and plausibility for a rule given that particular instance of 133 data (lines 8,9). The conditional beliefs and probabilities then yield an updated belief and plausibility 134 for each rule (lines 10, 11). Finally, the algorithm updates the uncertainty interval for each rule with 135 the new belief and plausibility values (lines 13,14). The result is a set of belief-theoretic rules (rules 136 accompanied with uncertainty intervals) (lines 16,17). 137

4 Comparing the Bayesian Approach with the Proposed Approach

To evaluate the proposed approach, we compare it to a Bayesian approach used in many PLL and SRL formulations (e.g, BLP, ProbLog). Consider a Bayesian clause of the form $\mathcal{R}^B := p :: \psi \implies (\neg)\phi$, where ψ, ϕ are Bayesian atoms and p is a point probability estimate. To learn p from data, we can establish a prior distribution and then update the distribution for each instance. We can assume an uninformative prior over each rule, as a uniform distribution: $\mathbf{p} \sim \text{Beta}(1,1) = \text{Uniform}(0,1)$. We now suppose that given a specific value of p_j for rule j, each individual source i will provide an observation $o_{i,j}$. We can compute the conditional, as follows: $P(o_{i,j}|p_j) = p_j^{o_{i,j}}(1-p_j)^{1-o_{i,j}}$. We can then compute the posterior, for n sources, where each source has a reliability measure or mass m'_i as follows:

$$P(p_j|o_{1,j},\ldots,o_{n,j}) = \prod_{i=1}^n \left[p_j^{o_{i,j}} (1-p_j)^{1-o_{i,j}} \right]^{m_i^{\prime}}$$

which simplifies to a Beta distribution:

$$\mathbf{p}_{\mathbf{j}}|o_{1,j},\ldots,o_{n,j} \sim \text{Beta}\left(\sum_{i=1}^{n} m'_{j}o_{i,j}+1,\sum_{i=1}^{n} m'_{j}(1-o_{i,j})+1\right)$$

In this case, the value of p in a Bayesian rule could be sampled from the distribution **p**. Since we do have a distribution, we can potentially estimate confidence intervals (or credible intervals) to generate

¹⁴¹ measures more akin to the proposed DS-based approach.

Algorithm 1 getParameters($\mathcal{D}, \mathcal{R}_{\mathcal{T}}^{\psi}$)

1: Input: $\mathcal{D} = \{d_1, \ldots, d_n\}$: Dataset containing *n* data instances 2: Input: ${}^{-}\mathcal{R}^{\psi}_{\mathcal{T}}$: An unspecified rule frame containing k rules \mathcal{R} 3: Initialize a DS Frame $\Theta^{\psi}_{T} = \{\theta_1, \dots, \theta_{2^k}\}$ 4: $m(\Theta_{\mathcal{T}}^{\psi}) \leftarrow 1$ 5: $t \leftarrow 0$ 6: for all $d \in \mathcal{D}$ do 7: Let \mathcal{E}_d be a BoE that corresponds to the data instance d8: Let \mathcal{E}_{Θ} be a BoE that corresponds to the indexed frame Θ_{τ}^{ψ} for all $\mathcal{R} \in \mathcal{R}^{\psi}_{\mathcal{T}}$ do 9: Set learning parameters μ_t and ν_t 10: bet rearring parameters μ_t allo ν_t $Bel(\mathcal{R}|d)^{\mathcal{E}_d} = Bel(\mathcal{R} \cap d)^{\mathcal{E}_d} / (Bel(\mathcal{R} \cap d)^{\mathcal{E}_d} + Pl(d \setminus \mathcal{R})^{\mathcal{E}_d})$ $Pl(\mathcal{R}|d)^{\mathcal{E}_d} = Pl(\mathcal{R} \cap d)^{\mathcal{E}_d} / (Pl(\mathcal{R} \cap d)^{\mathcal{E}_d} + Bel(d \setminus \mathcal{R})^{\mathcal{E}_d})$ $Bel(\mathcal{R})^{\mathcal{E}_{\Theta}}_{t+1} = \mu_t \cdot Bel(\mathcal{R})^{\mathcal{E}_{\Theta}}_t + \nu_t \cdot Bel(\mathcal{R}|d)^{\mathcal{E}_d}$ $Pl(\mathcal{R})^{\mathcal{E}_{\Theta}}_{t+1} = \mu_t \cdot Pl(\mathcal{R})^{\mathcal{E}_{\Theta}}_t + \nu_t \cdot Pl(\mathcal{R}|d)^{\mathcal{E}_d}$ 11: 12: 13: 14: end for 15: Set frame $\Theta^{\psi}_{\mathcal{T}}$ with $Bel(\mathcal{R})_{t+1}$ and $Pl(\mathcal{R})_{t+1}$ 16: 17: $t \leftarrow t + 1$ 18: end for 19: for all $\mathcal{R} \in \mathcal{R}^{\psi}_{\mathcal{T}}$ do $\alpha_{\mathcal{R}} \leftarrow Bel(\mathcal{R})^{\mathcal{E}_{\Theta}}$ 20: $\beta_{\mathcal{R}} \leftarrow Pl(\mathcal{R})^{\mathcal{E}_{\Theta}}$ 21: 22: end for 23: **Output:** $\alpha_1, \ldots, \alpha_k, \beta_1, \ldots, \beta_k$

142 4.1 Dataset and Experimental Conditions

We consider a small sample dataset \mathcal{D} shown in Table 1 containing ten data instances d from several 143 different sources, corresponding to an unspecified and modified version of the rule frame $-\mathcal{R}_{\psi}^{\psi}$ from 144 rules \mathcal{R}_1 and \mathcal{R}_2 in Example 1. For the evaluation, we consider four conditions that are variations 145 on the dataset in Table 1: (I) Reliable-Complete - the situation when there are no missing values 146 and the agent considers all the sources to be reliable and certain (ϵ are replaced with a 1 or 0, and 147 the masses all equal 1.0); (II) Reliable-Incomplete - while masses are still set to 1.0, several values 148 are missing in the data; (III) Unreliable-Complete - masses are set between 0 and 1, but there are 149 no missing values; and (IV) Unreliable-Incomplete – masses are set between 0 and 1, and there are 150 several missing values. We hypothesize that, even for such a simple dataset, the proposed approach 151 will be able to differentiate between these conditions and provide a richer sense for the uncertainty in 152 the data than the competing Bayesian Approach. 153

Source ID	driving	texting	mass
1	0	1	0.8
2	0	1	0.95
3	1	0	0.65
4	1	0	0.95
5	1	$\epsilon/0$	0.8
6	$\epsilon/0$	0	0.65
7	0	$\epsilon/1$	0.9
8	$\epsilon/1$	0	0.9
9	0	0	0.5
10	$\epsilon/0$	$\epsilon/1$	0.85

Table 1: Experimental Dataset \mathcal{D} . The ϵ refers to missing data entries. Variations of experimental conditions involve replacing ϵ with either a 1 or 0 as indicated, and mass values with 1.0.

We set the learning parameters $\nu_t = 1/|dataset|$, $\mu_t = 1 - \nu_t$ in the algorithm to emphasize that new information will be assigned a weight of 0.1, while the inertia of existing knowledge will be assigned a weight of 0.9. To remove order effects of the DS-based updating process, we ran the algorithm through 100 epochs, randomizing the order prior to each run. This allowed us to both reduce order effects and analyze convergence characteristics when compared to the Bayesian approach. On the Bayesian side, we looked at the maximum a posteriori (MAP) estimate and a 95% credible interval computed from the inverse CDF.

161 **5 Results and Discussion**

Ignorance, Reliability and Convergence. In this experiment, we set p of the Bayesian rules to the MAP of **p**. The general limitation with point estimates is that we cannot distinguish between the actual real-valued truth estimate and uncertainty in the truth value. Thus, for a value of 0.5 we do not know the amount of confidence in this estimate. We plot the MAP estimates (red circles) for both rules across the four conditions (Figure 1). Although there is some slight variation across various cases, the MAP estimates do not shed light on the reliability or completeness of the data.



Figure 1: Learned parameters provide richer sense of uncertainty. The proposed approach (blue) exposes more aspects of the epistemic uncertainty associated with small imperfect datasets than a traditional Bayesian approach (red). When the data is reliable and complete, the proposed approach converges to Bayesian. However, when the data displays unreliability or contains missing values, then the proposed approach allows for variable length intervals to express this uncertainty of evidence

We could also extend the definition of a Bayesian Clause to have an interval, corresponding to a 168 169 credible interval (CI) of p, allowing us to capture some of the richness present in the probability distribution. We computed a 95% CI for the probability distributions for each of the rules. The 170 CI suggests that there is a 95% chance that the calculated confidence interval from some future 171 experiment encompasses the true value of p. It does not, however, suggest that it contains the value 172 of true probability with 95% certainty, whereas the proposed Belief-theoretic uncertainty interval 173 states that the true value exists within the stated interval. Similar to the MAP estimate, the intervals 174 are also not informative, and although there is some widening in the presence of (un)reliability and 175 incompleteness, generally, there is not much variation between the conditions. Moreover, even if 176 the CI is suggestive of uncertainty more broadly, the upper and lower limits of the CI themselves, 177 do not provide any further information about the uncertainty of the rule. The belief-theoretic limits, 178 on the other hand, are well-defined and have specific meaning that pertain to the level of support 179 provided by the evidence. That is, the Bel() (lower limit) specifically represents the measure of 180 evidence supporting a proposition and the Pl() (upper limit) specifically measures evidence that 181 do not contradict the proposition. Thus, conversely, 1 - Bel() represents the level of doubt in the 182 evidence for the proposition and 1 - Pl() represents the level of disbelief in the evidence. This type 183 of information is not captured in a Bayesian CI. 184

The proposed approach captures variation in the data (conditions I-IV), while still converging to 185 Bayesian estimates when there is perfect data (condition I). Moreover, in one condition (condition 186 III, texting), the MAP estimate lies outside of the belief-theoretic interval suggesting that there 187 is a discrepancy between the different types of uncertainty being captured. We believe that the 188 MAP estimate captures aleatory uncertainty while the belief-theoretic approach captures epistemic 189 uncertainty, and this observed difference results from the selection of a potentially inappropriate prior 190 in the Bayesian approach. Not limited by a prior, the belief-theoretic approach allows for a more 191 dynamic update process and convergence to an estimate of epistemic uncertainty. 192

One advantageous feature of the proposed approach is that no matter how small the dataset, we can obtain a uncertainty interval based on the evidence received thus far. Although we do not show an instance-by-instance illustration of the algorithm, we can say that the algorithm begins with complete uncertainty [0, 1] and then with each input converges to either a point estimate (as is the case of ¹⁹⁷ condition I of reliable and complete data) or to an interval (as in the cases of conditions II, III, IV).

The rate and degree of convergence is also dependent on the selection of learning parameters m, n,

which roughly map to a learning rate.

Independence Between Relations. Another desirable feature of the proposed approach is that we 200 can ask a number of other questions of the indexed frame Θ^{ψ}_{T} that are not explicitly in the rule 201 system. For example, we might ask what is the uncertainty associated with (texting, driving). In 202 the Bayesian setting, if we assume these two actions are independent, we can multiply point estimates 203 (0.4 and 0.4 for condition I) and generate a non-zero probability of 0.16. In reality, these two actions 204 may not be independent of each other, and therefore it may be improper to make such an assumption. 205 Moreover, there is absolutely no evidence in Table 1 to support this non-zero probability as none of 206 the 10 data instances support both texting as well as driving. In contrast, in the proposed approach, 207 we do not make these assumptions, and instead directly query the same frame that was learned in the 208 experiment thus far. In doing so, we obtain an interval [0, 0]. This conforms to our intuitions about 209 the data as well as acceptable traffic norms. 210

211 6 Learning in an Open-World

Overall, we are interested in a numerical quantity that represents a degree to which the agent is 212 certain of something, or a degree to which the agent believes it, or a degree to which the evidence 213 supports it [21]. In the open world, measuring this sort of uncertainty requires the ability to process a 214 stream of information from multiple heterogeneous sources, say about a rule like those presented 215 above, and then incorporate and update uncertainty measures on this rule. The challenge is that one 216 information source may be quite different from another source not only in terms of reliability (as 217 discussed earlier) but also its repertoire of capabilities. For example, while one source can detect 218 both actions of texting and driving, another might only be able to detect the driving action, while still 219 another source might be able to detect a different action of "talking." 220

We elaborate this idea by extending the dataset in Table 1 and adding information from three new sources 11, 12 and 13, as shown in Table 2.

Source ID	texting	driving	talking	eating	mass
10	$\epsilon/0$	$\epsilon/1$	\geq	\succ	0.85
11	0	1	1	\geq	0.9
12	\geq	\ge	$>\!\!\!\!\!\!\!\!\!\!\!\!\!\!\!\!\!$	1	0.4
13	\geq	1	1	\times	0.6
12	\geq	\searrow	\geq	$\epsilon/0$	0.4

Table 2: Sample Dataset \mathcal{D}' : We extend Table 1 by first re-introducing source 10 and then sequentially incorporating sources 11, 12 and 13 and again source 10, each with different capabilities. Thus, when information from source 11 is received, a new attribute of "talking" must be incorporated into the frame and into the rule frame. When source 12 is processed, a new attribute of "eating" must be incorporated. Note, the agent does not see all these data instances simultaneously, but instead in sequence, which is more common in the open-world.

222

The first row of Table 2 includes the last entry of Table 1, which shows source 10 as being able to detect *texting* and *driving* and has provided observations $\{\epsilon/0, \epsilon/1\}$, accordingly. At this stage, the indexed frame $\Theta^{\psi}\tau$ is:

$$\Theta_{\mathcal{T}}^{\psi} = \{\theta_1 : (driving, texting), \theta_2 : (driving, \neg texting), \theta_3 : (\neg driving, texting), \theta_4 : (\neg driving, \neg texting)\}$$

Next, the agent receives evidence from source 11, which contains an additional previously unknown attribute *talking* and provides an observations for all three attributes *texting*, *driving*, and *talking*. The agent must now expand its indexed frame to incorporate this new previously unseen attribute and generate the following expanded frame $(\Theta_{\tau}^{\psi'})$:

 $\Theta_{\mathcal{T}}^{\psi\,'} = \{\theta_1': (driving, texting, talking), \dots \theta_5': (driving, texting, \neg talking), \dots\}$

²²³ DS-theory provides some useful notions to determine how $\Theta_{\mathcal{T}}^{\psi}$ relates to $\Theta_{\mathcal{T}}^{\psi'}$, which in turn will ²²⁴ allows us to dynamically grow and shrink the indexed frame as the agent sees new data. The notion

of "refinement" describes how one frame can be obtained from another by splitting some or all 225 of the elements of the initial frame. The canonical example of this operation is when a frame 226 $\Theta = \{animal, flowers\}$ is split into $\Theta' = \{dog, cat, rose, lily\}$. We can prove that expanding a 227 frame from $\Theta_{\mathcal{T}}^{\psi}$ to $\Theta_{\mathcal{T}}^{\psi'}$ is indeed such a refinement (although we cannot present the full proof here due to space limitations), by leveraging the exhaustive elaboration of consequents ϕ in the elementary 228 229 events $\hat{\theta}$ in Section 3.1, which in turn allows us to relate, for example $\{(\phi_1)\} \in \Theta$ with the expanded 230 $\{(\phi_1, \phi_2), (\phi_1, \neg \phi_2)\} \in \Theta'$. Showing that the two frames $\Theta_{\mathcal{T}}^{\psi}$ and $\Theta_{\mathcal{T}}^{\psi'}$ are a refinement also lets us establish that they are "compatible" in a DS sense. That means, we can show that the frames agree on the information defined in them, allowing us to prove that for $A \subseteq \Theta$, $Bel_{\Theta}(A) = Bel_{\Theta'}(\omega(A))$, 231 232 233 where ω is a refinement function $\omega : 2^{\Theta} \to 2^{\Theta'}$ mapping the two frames. 234 Next, when the agent receives information for a new attribute *eating* from source 12, the frame 235

is further refined. Because source 12 provides information for *only eating*, it doesn't make sense 236 to update the entire refined frame.⁴ Fortunately, DS-theory also defines the notion of "coarsening" 237 (inverse of refinement) which allows us to go in the opposite direction from $\Theta_{\mathcal{T}}^{\psi'}$ to $\Theta_{\mathcal{T}}^{\psi}$, an operation 238 that can be performed in $|\Theta| \log |\Theta'|$ time [12]. This ability to coarsen a frame allows us the possibility 239 of coarsening the frame to just one attribute (namely *eating*) and then incorporating the masses 240 assigned by source 12. The agent can then receive information from source 13, which does not 241 induce a refinement because there are no new attributes, ϕ , however, it does induce a coarsening as 242 it does not provide information for all attributes in the currently most-refined frame. Finally, the 243 agent might receive a new data instance from a previously known source, in this case source 12. The 244 task of coarsening the frame for source 12 is made easier this time because the agent already knows 245 the mapping between the frames from the prior computation. Thus, once an agent has encountered 246 a source, it remembers the capabilities of the source, and does not have to recompute these frame 247 mappings. 248

One of the most exciting aspects of the proposed approach is its ability to account for not just the 249 reliability of the sources, but as we have discussed, their capabilities as well. In doing so, we can 250 expand and grow the set of rules in real time, without always having to recompute joint distributions, 251 as we would need to do in a Bayesian approach. We can also update rules more efficiently as updating 252 a rule from a source with limited capacity does not impact existing knowledge about a capacity not 253 captured by the source. The DS-based approach allows us to speed up certain operations, especially 254 255 when new attributes are added ad-hoc and when sources provide information along a few dimensions; this is different from Bayesian approaches where even small open-world extension would require the 256 recalculation of the whole distribution. 257

7 General Discussion, Limitations and Conclusion

An advantage of learning belief-theoretic rules is to be able to apply existing DS-theoretic logic
formalisms (e.g., *Uncertain Logic* [20, 22]) to perform all manner of *inference* (e.g., modus ponens,
AND, OR). This sort of belief-theoretic inference has found applications in many AI and robotic
cognitive architectures [23, 24], so learning rule parameters from data would be beneficial.

DS-based operations are typically of exponential time complexity in the size of the frame since we consider all possible subsets of the frame. Although there are efficient implementations of DS-theoretic methods (graphical models and Monte-Carlo) [12], the proposed approach is generally intractable for large datasets. Thus, for large datasets, Bayesian approaches may be preferred.

Epistemic uncertainty is highly relevant to many open-world datasets. "Whereas the Bayesian language asks, in effect, that we think in terms of a chance model for the facts in which we are interested, the belief-function language asks that we think in terms of chance model for the reliability and meaning of our evidence." [25]. In this paper, we proposed a promising new probabilistic logic learning framework that uses a belief-theoretic logical representation combined with a learning methodology that allows for learning interval uncertainty for logical rules. The proposed approach offers several advantages over traditional Bayesian approaches when learning from *small imperfect*

274 datasets in the open world.

⁴Note the distinction between a source that is ignorant (when it is capable of providing an observation on an attribute but is unsure of its value) and the situation when the source is just incapable of providing any value as are the cases represented by red-crosses in Table 2.

275 **References**

- [1] De Raedt, L. Logical and relational learning (Springer Science & Business Media, 2008).
- [2] De Raedt, L. & Kersting, K. Probabilistic logic learning. ACM SIGKDD Explorations Newsletter
 5, 31–48 (2003).
- [3] Getoor, L. Introduction to statistical relational learning (MIT press, 2007).
- [4] Pearl, J. Probabilistic reasoning in intelligent systems: networks of plausible inference (Morgan Kaufmann, 2014).
- [5] Carlson, A. *et al.* Toward an architecture for never-ending language learning. In *AAAI*, vol. 5, 3 (2010).
- [6] Shafer, G. A Mathematical Theory of Evidence (Princeton University Press, 1976).
- [7] Delmotte, F. & Smets, P. Target identification based on the transferable belief model interpreta tion of dempster-shafer model. *IEEE Transactions on Systems, Man, and Cybernetics-Part A: Systems and Humans* 34, 457–471 (2004).
- [8] Seraji, H. & Serrano, N. A multisensor decision fusion system for terrain safety assessment.
 IEEE Transactions on Robotics 25, 99–108 (2009).
- [9] Mahoor, M. H. & Abdel-Mottaleb, M. A multimodal approach for face modeling and recognition.
 IEEE Transactions on Information Forensics and Security 3, 431–440 (2008).
- [10] Pearl, J. Reasoning with belief functions: a critical assessment. *Technical ReportR-136, UCLA* (1989).
- [11] Pearl, J. Reasoning with belief functions: An analysis of compatibility. *International Journal* of Approximate Reasoning 4, 363–389 (1990).
- [12] Wilson, N. Algorithms for dempster-shafer theory. In *Handbook of defeasible reasoning and uncertainty management systems*, 421–475 (Springer, 2000).
- [13] Smets, P. & Kennes, R. The transferable belief model. Artificial intelligence 66, 191–234 (1994).
- [14] Bergsten, U. & Schubert, J. Dempster's rule for evidence ordered in a complete directed acyclic
 graph. *International Journal of Approximate Reasoning* 9, 37–73 (1993).
- [15] Yager, R. R. On the dempster-shafer framework and new combination rules. *Information Sciences* 41, 93–137 (1987).
- [16] Nunez, R. C. *et al.* DS-based uncertain implication rules for inference and fusion applications.
 Information Fusion (FUSION), 2013 16th International Conference on 1934–1941 (2013).
- [17] Tang, Y., Hang, C. W., Parsons, S. & Singh, M. Towards argumentation with symbolic dempster shafer evidence. *Frontiers in Artificial Intelligence and Applications* 245, 462–469 (2012).
- [18] Dewasurendra, D. A., Bauer, P. H. & Premaratne, K. Evidence filtering. *IEEE Transactions on Signal Processing* 55, 5796–5805 (2007).
- [19] Fagin, R. & Halpern, J. Y. A new approach to updating beliefs. *arXiv preprint arXiv:1304.1119* (2013).
- [20] Núnez, R. C., Scheutz, M., Premaratne, K. & Murthi, M. N. Modeling uncertainty in first-order
 logic: A dempster-shafer theoretic approach. In *8th International Symposium on Imprecise Probability: Theories and Applications* (2013a).
- [21] Shafer, G. Non-additive probabilities in the work of bernoulli and lambert. In *Classic Works of the Dempster-Shafer Theory of Belief Functions*, 117–182 (Springer, 2008).

- [22] Núnez, R. C., Scheutz, M., Premaratne, K. & Murthi, M. N. Modeling uncertainty in first-order
 logic: A dempster-shafer theoretic approach. In *8th International Symposium on Imprecise Probability: Theories and Applications* (2013).
- ³²⁰ [23] Sarathy, V. & Scheutz, M. A logic-based computational framework for inferring cognitive ³²¹ affordances. *IEEE Transactions on Cognitive and Developmental Systems* **8** (2016).
- Williams, T., Briggs, G., Oosterveld, B. & Scheutz, M. Going beyond command- based
 instructions: Extending robotic natural language interaction capabilities. In *Proceedings of* AAAI (2015).
- [25] Shafer, G. & Tversky, A. Languages and designs for probability judgment. *Cognitive Science* 9, 309–339 (1985).