

# Unifying Knowledge Base Completion with PU Learning to Mitigate the Observation Bias

Jonas Schouterden, Jessa Bekker, Jesse Davis, Hendrik Blockeel

KU Leuven, Department of Computer Science, B-3000 Leuven, Belgium  
 Leuven.AI - KU Leuven Institute for AI, B-3000 Leuven, Belgium  
 {jonas.schouterden, jessa.bekker, jesse.davis, hendrik.blockeel}@kuleuven.be

## Abstract

Methods for Knowledge Base Completion (KBC) reason about a knowledge base (KB) in order to derive new facts that should be included in the KB. This is challenging for two reasons. First, KBs only contain positive examples. This complicates model evaluation which needs both positive and negative examples. Second, those facts that were selected to be included in the knowledge base, are most likely not an i.i.d. sample of the true facts, due to the way knowledge bases are constructed. In this paper, we focus on rule-based approaches, which traditionally address the first challenge by making assumptions that enable identifying negative examples, which in turn makes it possible to compute a rule’s confidence or precision. However, they largely ignore the second challenge, which means that their estimates of a rule’s confidence can be biased. This paper approaches rule-based KBC through the lens of PU learning, which can cope with both challenges. We make three contributions. (1) We provide a unifying view that formalizes the relationship between multiple existing confidences measures based on (i) what assumption they make about and (ii) how their accuracy depends on the selection mechanism. (2) We introduce two new confidence measures that can mitigate known biases by using propensity scores that quantify how likely a fact is to be included in the KB. (3) We show through theoretical and empirical analysis that taking the bias into account improves the confidence estimates, even when the propensity scores are not known exactly.

## 1 Introduction

Knowledge Bases (KBs) such as Wikidata (Vrandečić and Krötzsch 2014), YAGO (Rebele et al. 2016) and DBpedia (Auer et al. 2007) are large collections of knowledge about the world. They contain entities, such as *Audrey, Belgium* and *BestActress1960*, and facts about those entities such as  $\langle \text{Audrey, wasBornIn, Belgium} \rangle$  and  $\langle \text{Audrey, wonOscar, BestActress1960} \rangle$ . KBs are typically constructed through crowdsourcing or automatically extracting information from the web (Rebele et al. 2016). Consequently, these KBs are incomplete as they do not contain all facts.

Knowledge Base Completion (KBC) (Nickel et al. 2015) aims to address the issue of incompleteness by reasoning over the knowledge base in order to derive new facts that should be included in the KB. This is typically achieved by learning

a model from the initial incomplete KB. One common way to do this is to take a rule-based approach (Galárraga et al. 2013; Pellissier Tanon et al. 2017; Zupanc and Davis 2018). This would result in learning rules like  $\langle X, \text{wonOscar}, Y \rangle \wedge \langle X, \text{isMarriedTo}, Z \rangle \Rightarrow \langle Z, \text{livesIn}, \text{USA} \rangle$ , meaning that partners of Oscar winners usually live in the USA. A common measure for the quality of (intermediate) models is *confidence* (precision), which is the fraction of correctly predicted facts.

However, learning from the incomplete KB is challenging for two important reasons. First, KBs operate under *Open World* semantics which means that the truth value of any triple not in the KB is unknown. These triples could be true (i.e., they belong in the KB) or false (i.e., they should be excluded from the KB). Practically, this means the data only contains positive examples, whereas most learners require both positive and negative examples. It also implies that a rule’s confidence cannot be computed without making additional assumptions. Second, the way knowledge bases are constructed makes it highly unlikely that the facts included in the observed KB are an i.i.d. sample of the facts in the *ideal knowledge base*, aka the ground truth. Indeed, studies have shown that knowledge bases suffer from *observation biases*: They contain cultural biases, contain more facts about famous people and represents men and women differently (Callahan and Herring 2011; Wagner et al. 2015; Soulet et al. 2018). If knowledge base completion is applied while ignoring the observation bias, then the newly inferred facts are likely to strengthen the bias. Yet, to the best of our knowledge, this is what all current KBC methods do.

PU learning (learning from positive and unlabeled examples) (Bekker and Davis 2020), which is concerned with learning a binary classifier while only having access to positive and unlabeled examples, is well-equipped for addressing both these challenges. First, it perfectly matches the type of data available for KBC: the positive examples are the facts in the KB whereas the unlabeled data is any potential fact that is not included in the KB. Second, recent work in PU learning (Kato, Teshima, and Honda 2018; Bekker, Robberechts, and Davis 2019; Gong et al. 2021) has explicitly modeled the selection mechanism that determines the probability of observing a positive example’s label, i.e., the observation bias.

Motivated by this, we view the KBC task as a PU Learning

problem, which enables us to explicitly consider the selection mechanism. We consider rule-based approaches to KBC and make three contributions. (1) We provide a unifying view that formalizes the relationship between multiple existing confidences measures based on (i) what assumption they make about and (ii) how their accuracy depends on the selection mechanism. (2) We introduce two new confidence measures that can mitigate known biases by using propensity scores that quantify how likely a fact is to be included the KB. (3) We show through theoretical and empirical analysis that taking the bias into account improves the confidence estimates, even when the propensity scores are not known exactly.

## 2 Preliminaries

*Knowledge Bases (KBs)* store interlinked information about entities in the form of relations between the entities, often as RDF triple stores (WWW Consortium 2004). Using this format, the KB is a triple  $(\mathcal{E}, \mathcal{P}, F)$ , with  $\mathcal{E}$  the set of entities,  $\mathcal{P}$  the set of predicates and the  $F$  set of facts, denoted by  $\langle s, p, o \rangle$  triples with subject  $s \in \mathcal{E}$  and object  $o \in \mathcal{E}$  and predicate  $p \in \mathcal{P}$ , meaning that a relation of type  $p$  holds between entities  $s$  and  $o$ . The triples in a KB are a subset of the Cartesian product  $\mathcal{E} \times \mathcal{P} \times \mathcal{E}$  and each predicate and entity in  $\mathcal{E}$  and  $\mathcal{P}$  occurs at least once in a triple in the KB.

A knowledge base models a certain part of the world. We call a knowledge base *ideal* if it has a triple for each relevant fact, and *incomplete* if it contains only a subset of those triples. We will consistently use  $I$  to refer to the ideal knowledge in some context, and  $K$  to refer to a given "known" (incomplete) knowledge base. The task of *Knowledge Base Completion (KBC)* is then: given a knowledge base  $K$ , reconstruct the ideal knowledge base  $I$ .

The KBC task is often approached as follows: given an incomplete knowledge base, rules are derived of the form  $Body(s, o) \Rightarrow \langle s, p, o \rangle$ , with the semantics that if  $Body(s, o)$  (some condition on  $s$  and  $o$ ) is fulfilled in  $K$ , then  $\langle s, p, o \rangle$  is in  $I$ . These rules can then be used to derive new facts (facts that are not in  $K$ , but are in  $I$ ). We follow these semantics:  $Body(s, o)$  is always applied to  $K$ , predicting  $\langle s, p, o \rangle$  to be in  $I$ .  $Body(s, o)$  is typically a conjunctive condition (Galárraga et al. 2013; Fürnkranz, Gamberger, and Lavrač 2014; Pellissier Tanon et al. 2017; Zupanc and Davis 2018; Lajus, Galárraga, and Suchanek 2020), though this is not essential for this paper.

In the remainder of this paper, we will use the following notation. In the context of a specific rule,  $R$  refers to the rule itself, and  $p$  refers to the (fixed) predicate of the rule's prediction. We use the following indicator functions (which return 1 if the associated condition is true and 0 otherwise):

- $R(s, o) : \langle s, o \rangle$  fulfills the rule's conditions  $Body(s, o)$
- $y(s, o) = y(\langle s, p, o \rangle) : \langle s, p, o \rangle$  is in  $I$  ("is a fact")
- $l(s, o) = l(\langle s, p, o \rangle) : \langle s, p, o \rangle$  is in  $K$  ("is observed")
- $y(s) = y(\langle s, p \rangle) = \max_o y(s, o)$
- $l(s) = l(\langle s, p \rangle) = \max_o l(s, o)$

For readability, we use the short versions  $y(s, o)$ ,  $l(s, o)$ ,  $y(s)$ ,  $l(s)$  when  $p$  is implied (e.g., in the context of a single rule).

Note that  $y(s)$  and  $l(s)$  indicate that, for this specific  $s$ , at least one triple of the form  $\langle s, p, \cdot \rangle$  is respectively in  $I$  / in  $K$ .

Based on the above functions, we define the following sets:

- $\mathbf{R} = \{\langle s, o \rangle \mid R(s, o) = 1\}$
- $\mathbf{R}^+ = \{\langle s, o \rangle \in \mathbf{R} \mid y(s, o) = 1\}$
- $\mathbf{R}^l = \{\langle s, o \rangle \in \mathbf{R} \mid l(s, o) = 1\}$
- $\mathbf{R}_s^+ = \{\langle s, o \rangle \in \mathbf{R} \mid y(s) = 1\}$
- $\mathbf{R}_s^l = \{\langle s, o \rangle \in \mathbf{R} \mid l(s) = 1\}$

That is,  $\mathbf{R}$  is the rule's *coverage* containing all  $\langle s, o \rangle$  triples for which the rule fires (i.e. the rule's predictions); among those,  $\mathbf{R}^+$  and  $\mathbf{R}^l$  contain respectively the true and observed ones.  $\mathbf{R}_s^+$  and  $\mathbf{R}_s^l$  respectively restricts  $\mathbf{R}$  to triples with an  $s$  for which at least one  $\langle s, p, o \rangle$  is true or observed.

This paper focuses on *evaluating rules* of the format just described. *Confidence measures* are typically used to evaluate the quality of rules, during rule induction and model evaluation. The *confidence* of a rule is

$$\text{conf}(R) = \frac{\sum_{\langle s, o \rangle \in \mathbf{R}} y(s, o)}{|\mathbf{R}|} = \frac{|\mathbf{R}^+|}{|\mathbf{R}|}.$$

This definition reflects the fact that rules are executed on  $K$  but their predictions are considered correct if the predicted triple is in  $I$ .

When constructing a rule set from a knowledge base, a learner typically repeatedly tries to pick the highest-confidence rule from a number of options. As the learner has access to  $K$  but not  $I$ , it cannot compute  $\text{conf}(R)$ , so it must estimate it. Before discussing existing and novel estimators in section 4, we first discuss how current KBC approaches deal with this.

## 3 Assumptions on the Selection Mechanism

If  $K$  were an i.i.d. sample from the set of all triples, labeled positive or negative according to whether they are in  $I$ , and each triple had the same probability of being included in  $K$ , then simply counting how many of the predicted triples are labeled positive or negative would yield an unbiased estimator of  $\text{conf}(R)$  (just like test set accuracy is an unbiased estimator for population accuracy). But  $K$  contains no negative examples at all. This poses the following challenge:

How can one estimate the confidence of a rule without access to negative examples?

### 3.1 Typical Assumptions in KBC

In general, evaluating models without access to negative examples remains an open problem (Speranskaya, Schmitt, and Roth 2020; Pezeshkpour, Tian, and Singh 2020). A common approach in KBC is to make assumptions that allow deriving negative examples. Two prominent assumptions are:

**The closed-world assumption (CWA)** (naively) assumes that all facts not included in  $K$  are false. Hence, if a rule derives a fact not in  $K$ , that corresponds to a false positive.

**The partial-completeness assumption (PCA)** (a.k.a. *local closed-world assumption*) assumes that if a  $\langle s, p, o \rangle$  triple

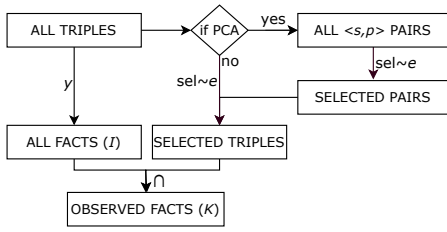


Figure 1: Which facts of  $I$  are observed in  $K$  depends the selection mechanism.

is observed, then all triples  $\langle s, p, \cdot \rangle$  in  $I$  with the same subject and predicate are observed (Galárraga et al. 2013; Dong et al. 2014). It follows from this that if  $\langle s, p, o \rangle$  is predicted and  $K$  contains  $\langle s, p, o' \rangle$  for some  $o' \neq o$ , but not  $\langle s, p, o \rangle$ , then this must be a false positive.

The key insight in our paper is that these assumptions fail to account for the underlying mechanism that determines how the KB is populated, which results in biased estimates of a rule’s confidence. In the following, we explicitly look at possible selection mechanisms and how the CWA and PCA assumptions connect with them.

### 3.2 Selection Mechanisms

Which facts are observed in  $K$  depends on how the KB was populated. Conceptually, this can be modeled by assuming that whether a fact is included or not in the KB depends on a *selection mechanism*  $\text{sel}(\langle s, p, o \rangle)$  which selects triples  $\langle s, p, o \rangle$  from  $\mathcal{E} \times \mathcal{P} \times \mathcal{E}$  (see Figure 1). If a selected triple is a fact in  $I$  then it becomes part of  $K$ :  $l(\langle s, p, o \rangle) = \text{sel}(\langle s, p, o \rangle)y(\langle s, p, o \rangle)$ .

The selection mechanism can operate in different ways. From a probabilistic point of view, the simplest version is that  $K$  is an independent and identically distributed (i.i.d.) sample from  $I$ . CWA is then the special case where each triple has probability 1 of being included in  $K$ . More realistically, groups of triples might be selected together (not independent) and some triples might be more likely to be selected than others (not identically distributed). PCA implies one particular type of dependence: It assumes a hierarchical selection mechanism that first selects pairs  $\langle s, p \rangle$ , then selects *all* triples  $\langle s, p, \cdot \rangle$  of the selected pairs, as depicted in Figure 1.

We next focus on the actual selection probabilities. While existing rule-based KBC work neither explicitly states nor considers the selection mechanism, the field of PU Learning makes such assumptions very explicit. Therefore, we follow their terminology. The probability with which a positive example is selected for inclusion is called its *propensity score*  $e$  (Bekker, Robberechts, and Davis 2019):

$$e(\langle s, p, o \rangle) = \Pr(\text{sel}(\langle s, p, o \rangle) = 1) \quad \# \text{ no PCA}$$

$$e(\langle s, p \rangle) = e(\langle s, p, o \rangle) = \Pr(\text{sel}(\langle s, p \rangle) = 1) \quad \# \text{ PCA}$$

From this, the probability that a triple appears in  $K$  follows:

$$\Pr(l(\langle s, p, o \rangle) = 1) = e(\langle s, p, o \rangle)y(\langle s, p, o \rangle)$$

We use shorthands  $e(s, o) = e(\langle s, p, o \rangle)$  and  $e(s) = e(\langle s, p \rangle)$ .

Assumptions in PU Learning about the selection probabilities range from *Selected Completely At Random (SCAR)*, where each positive example has the constant probability  $e(\cdot) = c$  to be selected, to *Selected At Random (SAR)*, where the propensity score is a function of the example’s features (Bekker, Robberechts, and Davis 2019).

Based on this, we propose the following taxonomy for assumptions about KBC selection probabilities:

*Closed World Assumption (CWA)*: All facts are observed:

$$e(\langle s, p, o \rangle) = 1.$$

*SCAR* All facts have the same probability to be selected:

$$e(\langle s, p, o \rangle) = c$$

*SCAR-per-predicate (SCAR<sub>p</sub>)*: All facts about the same predicate  $p$  have the same probability to be selected:

$$e(\langle s, p, o \rangle) = c_p$$

*SCAR-per-rule (SCAR<sub>R</sub>)*: All facts predicted by a rule  $R$  have the same probability to be selected:

$$R(s, o) = 1 \Rightarrow e(\langle s, p, o \rangle) = c_R$$

*SAR*: The probability that a fact gets selected, depends on its characteristics in the incomplete KB  $K$ .<sup>1</sup>

Note that this categorization is largely orthogonal to any dependence structures in the selection mechanism. In particular, all five levels are compatible with PCA, though additional restrictions may apply (e.g., SCAR-per-rule under PCA implies that rules with different  $c_R$  cannot cover the same subject  $s$ ).

Only SAR, the least strict assumption, can in general represent common observation biases such as higher propensity scores for famous entities, and certain predicates having different propensity scores for women and men (Callahan and Herring 2011; Wagner et al. 2015). SCAR-per-rule can include such biases, but only if each rule covers one group exclusively (famous or plebeian, man or woman).

Note that all except the SAR assumption consider the observed facts covered by a certain rule to be unbiased. The next section will show that this assumption is made implicitly by all existing confidence measures.

## 4 Confidence Estimators

We now return to the problem of estimating the confidence of a rule  $R$ ,  $\text{conf}(R)$ . First, we discuss and analyze estimators that have been used in the KBC field. Second, we will introduce new estimators that account for possible observation bias.

### 4.1 Existing Confidence Estimators

**CWA-based estimator** Under the closed-world assumption, a prediction is considered correct if it is known to be correct (i.e., the predicted triple is observed in  $K$ ), and incorrect otherwise. The confidence calculated as such is usually referred to as the *standard confidence* (Galárraga et al. 2013), but for clarity, we call it the CWA-based estimator.

$$\text{CWA}(R) = \frac{\sum_{\langle s, o \rangle \in \mathbf{R}} l(s, o)}{|\mathbf{R}|} = \frac{|\mathbf{R}^1|}{|\mathbf{R}|} = \text{conf}(R) \frac{|\mathbf{R}^1|}{|\mathbf{R}^+|}$$

<sup>1</sup>More realistically, the probability depends on its true (possibly unobserved) characteristics. However, in the KBC task, only the characteristic that are observable are relevant.

$CWA(R)$  generally *underestimates*  $\text{conf}(R)$  because  $K \subset I$  and therefore  $|\mathbf{R}^1| \leq |\mathbf{R}^+|$ , yielding  $CWA(R) \leq \text{conf}(R)$ .

Now consider different possible realizations of  $K$  given some  $I$ . Because the probability for a fact to be included in  $K$  is  $Pr(l(s, o)=1 \mid y(s, o)=1) = e(s, o)$ , the expected value for  $CWA(R)$  over all  $K$  is

$$\mathbb{E}_{\text{sel} \sim e} [CWA(R)] = \text{conf}(R) \frac{1}{|\mathbf{R}^+|} \sum_{\langle s, o \rangle \in \mathbf{R}^+} e(s, o).$$

Under SCAR-per-predicate, with a constant  $e(s, o) = c_p$  per predicate  $p$ , it has the same constant multiplicative bias  $c_p$  for all rules for a predicate  $p$ , meaning that the ranking of rules is still expected to be correct. To better analyse the problem for the setting where only the relative ranking of rules matter, we introduce an *inverse  $c_p$ -weighted* version of the CWA-based estimator  $ICW(R) = \frac{1}{c_p} CWA(R)$ , which is indeed unbiased under SCAR-per-predicate:  $\mathbb{E}_{\text{sel} \sim c_p} [ICW(R)] = \text{conf}(R)$ .

**PCA-based estimator** To solve the above-mentioned underestimation problem, the PCA-based estimator only considers the subset of predictions  $\mathbf{R}_s^1$  assumed to have a known truth value under the PCA assumption. That is, the PCA-based estimator only considers predictions  $\mathbf{R}_s^1$  for which the subject appears in an observed fact in  $K$  ( $l(s) = 1$ ). For all triples in  $\mathbf{R}_s^1$ , if the specific triple is observed ( $l(s, o) = 1$ ) then the prediction is considered correct, if it is not observed ( $l(s, o) = 0$ ) then the prediction is considered incorrect (Galárraga et al. 2013):

$$PCA(R) = \frac{\sum_{\langle s, o \rangle \in \mathbf{R}} l(s, o)}{\sum_{\langle s, o \rangle \in \mathbf{R}} l(s)} = \frac{|\mathbf{R}^1|}{|\mathbf{R}_s^1|}$$

While the PCA-based estimator is a commonly-used confidence estimator, we are, to the best of our knowledge, the first to study under which conditions it is expected to perform well. The PCA-based estimator can suffer from biases induced by three factors, which are mathematically derived and interpreted in Appendix A<sup>2</sup>:

$$\begin{aligned} \mathbb{E}_{\text{sel} \sim e} [PCA(R)] &\approx \frac{\mathbb{E}_{\text{sel} \sim e} [|\mathbf{R}^1|]}{\mathbb{E}_{\text{sel} \sim e} [|\mathbf{R}_s^1|]} \quad \text{first-order Taylor approximation} \\ &= \text{conf}(R) \cdot \text{bias}_{p \in \mathcal{A}}(R) \cdot \text{bias}_{y(s)=0}(R) \cdot \text{bias}_{e(s)}(R) \\ &= \text{conf}(R) \frac{\sum_{\langle s, o \rangle \in \mathbf{R}^+} e(s, o)}{\sum_{\langle s, o \rangle \in \mathbf{R}^+} e(s)} \frac{|\mathbf{R}|}{|\mathbf{R}^+|} \frac{\frac{1}{|\mathbf{R}^+|} \sum_{\langle s, o \rangle \in \mathbf{R}^+} e(s)}{\frac{1}{|\mathbf{R}_s^+|} \sum_{\langle s, o \rangle \in \mathbf{R}_s^+} e(s)}, \end{aligned}$$

The three biases can cancel each other out, making it hard to predict how  $PCA(R)$  will perform, when the problem is not well understood. We explain the three biases using an example rule that predicts which people won an Oscar, e.g.,  $\langle \text{Audrey}, \text{wonOscar}, \text{BestActress1954} \rangle$ . Under PCA,  $K$  is assumed to contain either all or none of the Oscars that each person has won.

<sup>2</sup>The appendices can be found at: <https://github.com/ML-KULeuven/KBC-as-PU-Learning>

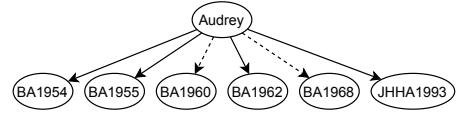


Figure 2: Correctly predicted facts  $\mathbf{R}^+$  for subject *Audrey* by the example rule. Dotted arrows indicate predicted facts not in  $K$ , which the PCA considers to be false positives.

To illustrate when  $\text{bias}_{p \in \mathcal{A}}$  arises, consider a KB that only contains facts denoting four of Audrey Hepburn’s six Oscars meaning the PCA assumption is violated for the  $\langle \text{Audrey}, \text{wonOscar} \rangle$  pair. Suppose a rule is learned that correctly derives all six of Audrey Hepburn’s Oscars (Figure 2). Making the PCA assumption results in the two Oscar wins not in the KB being incorrectly denoted as false positives, yielding an *underestimate* of the rule’s confidence just like in the CWA case. More formally, this bias arises whenever a rule fires for a  $\langle s, p \rangle$  pair for which the PCA assumption is violated and the rule derives a  $\langle s, p, o \rangle$  s.t.  $l(\langle s, p, o \rangle) = 0$  and  $y(\langle s, p, o \rangle) = 1$  (an unobserved fact). This bias never arises for functional predicates (where each subject appears in at most 1 fact) because the PCA trivially holds for such predicates.

To illustrate when  $\text{bias}_{y(s)=0}$  arises, consider a rule that predicts that Alan Rickman won an Oscar, which is a false positive since he has never won an Oscar. However, because there is no fact  $\langle \text{AlanRickman}, \text{wonOscar}, \cdot \rangle$  in  $K$ , since no such facts exist, the PCA estimator disregards the prediction in its confidence. In this case the PCA estimator *overestimates* the rule’s confidence. More generally, such an overestimate occurs whenever a rule predicts a triple  $\langle s, p, o' \rangle$  where  $\forall o : y(\langle s, p, o \rangle) = 0$  ( $s$  occurs in no facts). In these cases, these false positive predictions  $\mathbf{R} \setminus \mathbf{R}_s^+$  are ignored by the PCA.

The third bias factor  $\text{bias}_{e(s)}$  is the mean  $e(s)$  over all correct predictions  $\mathbf{R}^+$  made by the rule divided by the mean  $e(s)$  over all its predictions  $\mathbf{R}_s^+$  (restricting  $s$  to those  $s$  where  $y(s) = 1$ ) (Figure 3). This bias is  $> 1$  when correct predictions tend to have higher  $e(s)$ , or, put differently, when there are more correct predictions for high-propensity subjects. Vice versa, this bias is  $< 1$  when there are fewer correct predictions for high-propensity subjects. In our Oscars example, when a rule happens to give more accurate predictions for high-propensity Oscar winners than for low-propensity ones, the confidence of this rule is overestimated.

In our experiments in Section 5, **Q3** illustrates how  $PCA(R)$  can vary due to this bias when  $e(s)$  varies for different subjects.

## 4.2 Observation Bias Aware Confidence Estimators

In this section we propose two novel confidence estimators that can counteract observation biases by explicitly taking the selection mechanism into account. The difference between the estimators is whether or not they make the PCA assumption. The proposed estimators need propensity scores as input, therefore we additionally analyze their bias when using imperfect propensity scores and show that rough estimates are

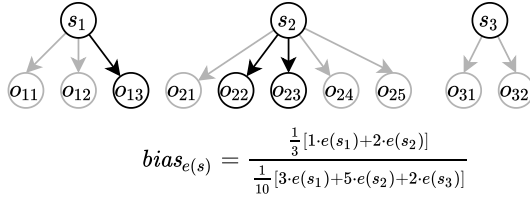


Figure 3: Example  $bias_{e(s)}$  calculation. The arrows indicate predictions  $\mathbf{R}_s^+$  for which the subject occurs in at least one fact. Black arrows are correct predictions  $\mathbf{R}^+$ , grey arrows are incorrect predictions  $\mathbf{R}_s^+ \setminus \mathbf{R}^+$ .

better than assuming that there is no observation bias.

The **Inverse Propensity Weighted estimator (IPW)** aims to debias  $CWA(R)$  by weighting the observed triples with inverse propensity score estimates  $\hat{e}(s, o)$ :

$$IPW(R) = \frac{1}{|\mathbf{R}|} \sum_{\langle s, o \rangle \in \mathbf{R}} \frac{l(s, o)}{\hat{e}(s, o)}$$

From its expected value over all possible  $K$ , it is clear that the estimator is *unbiased* when  $\hat{e}(s, o) = e(s, o)$ :

$$\mathbb{E}_{\text{sel} \sim e} [IPW(R)] = \text{conf}(R) \frac{1}{|\mathbf{R}^+|} \sum_{\langle s, o \rangle \in \mathbf{R}^+} \frac{e(s, o)}{\hat{e}(s, o)}.$$

Similarly, the **Inverse Propensity Weighted PCA-based estimator (IPW-PCA)** aims to debias  $PCA(R)$ :

$$IPW-PCA(R) = \frac{\sum_{\langle s, o \rangle \in \mathbf{R}} \frac{l(s, o)}{\hat{e}(s)}}{\sum_{\langle s, o \rangle \in \mathbf{R}} \frac{l(s)}{\hat{e}(s)}}$$

The first-order Taylor approximation of its expected value is:

$$\begin{aligned} & \mathbb{E}_{\text{sel} \sim e} [IPW-PCA(R)] \\ & \approx \text{conf}(R) \cdot bias_{PCA}^{IPW-PCA}(R) \cdot bias_{y(s)=0}(R) \cdot bias_{e(s)}^{IPW-PCA}(R) \\ & \approx \text{conf}(R) \frac{\sum_{\langle s, o \rangle \in \mathbf{R}^+} \frac{e(s, o)}{\hat{e}(s)}}{\sum_{\langle s, o \rangle \in \mathbf{R}^+} \frac{e(s)}{\hat{e}(s)}} \frac{|\mathbf{R}|}{|\mathbf{R}^+|} \frac{1}{|\mathbf{R}^+|} \sum_{\langle s, o \rangle \in \mathbf{R}^+} \frac{e(s)}{\hat{e}(s)} \end{aligned}$$

Here, 2 of the 3 bias factors are inverse propensity weighted versions of the corresponding  $PCA(R)$  biases. Note that when  $\hat{e}(s) = e(s)$ , the  $bias_{e(s)}^{IPW-PCA}(R)$  related to the selection mechanism completely disappears.

Most often, the exact propensity scores  $e(s, o)$  cannot be used for  $\hat{e}(s, o)$ , as they are unknown. When the propensity scores are not known exactly, using reasonable estimates for  $\hat{e}(s, o)$  can still result in a better confidence estimator than not using any  $\hat{e}(s, o)$  and making a SCAR assumption.

To investigate how accurate the propensity score estimates  $\hat{e}(\cdot)$  should be, we compare confidence estimators when the PCA does or does not hold, respectively:  $IPW-PCA(R)$  vs  $PCA(R)$  and  $IPW(R)$  vs  $CWA(R)$ <sup>3</sup>.

<sup>3</sup>To allow for rule comparison, we considered the calibrated version  $ICW(R) = \frac{1}{c_p} CWA(R)$ .

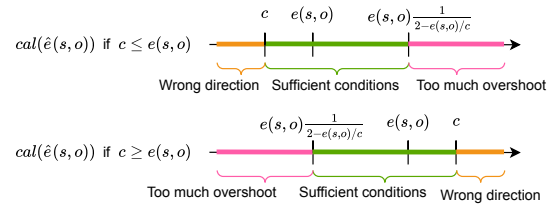


Figure 4: If  $cal(\hat{e}(s, o))$  is between  $c$  and  $e(s, o) \cdot \frac{1}{2 - e(s, o)/c}$ , then  $\langle s, o \rangle$  has a lower error contribution in the IPW(-PCA) estimator than in the CWA/PCA estimator.

In Appendix B, we show that an individual triple  $\langle s, o \rangle$  has a smaller error contribution in the IPW(-PCA) estimator than in the CWA/PCA estimators when its propensity score estimate  $\hat{e}(s, o)$  satisfies:

$$\begin{cases} c \leq cal(\hat{e}(s, o)) \leq e(s, o) \frac{1}{2 - e(s, o)/c} & , \text{ when } c \leq e(s, o), \\ c \geq cal(\hat{e}(s, o)) \geq e(s, o) \frac{1}{2 - e(s, o)/c} & , \text{ when } c \geq e(s, o), \end{cases}$$

with  $c = c_R$  under PCA and  $c = c_p$  otherwise. Here,  $cal(\hat{e}(s, o))$  multiplicatively calibrates  $\hat{e}(s, o)$ , so that  $\mathbb{E}[e(s, o) / cal(\hat{e}(s, o))] = \mathbb{E}[e(s, o) / c] = 1$  (note that  $\hat{e}(s, o)$  itself need not be calibrated). In other words, the IPW(-PCA) confidence estimator is preferable over the CWA/PCA confidence estimator as soon as  $cal(\hat{e}(s, o))$  deviates from  $c$  “in the right direction”, that is, towards  $e(s, o)$ , and this up to the point where it overshoots by a certain factor (see Figure 4). The allowed overshoot increases with  $e(s, o)/c$ .

As shown, reasonable estimates  $\hat{e}$  can be used in the IPW(-PCA) estimators that do not need be calibrated; only their relative values matter. In practice, these relative  $\hat{e}$  could be derived from domain knowledge, e.g., from research on KB bias (Callahan and Herring 2011; Wagner et al. 2015), or estimated through incompleteness estimation (Razniewski, Suchanek, and Nutt 2016; Galárraga et al. 2017). For example, in a movie-recommendation setting, Saito et al. (2020) use a movie’s popularity as its propensity score.

## 5 Experiments

We aim to empirically answer the following research questions: can we effectively account for observation biases (i.e., obtain more accurate confidence estimates) using the newly proposed propensity-based estimators, **(Q1)** when the propensities are known, **(Q2)** when propensities are guessed (“noisy” propensities), **(Q3)** even when the PCA assumption holds?

### 5.1 Experimental Setup

Evaluating a confidence estimator requires knowledge of  $I$ , which is generally unavailable for real-world KBs. We therefore equate  $I$  to a real-world KB from which  $K$ ’s are generated by applying different selection mechanisms. Our  $I$  is the popular KBC benchmark dataset Yago3-10 (Mahdisoltani, Biega, and Suchanek 2015). Rules predicting any  $p \in \mathcal{P}$  are mined from  $I$  with AMIE (Galárraga et al. 2015) with its default settings and a minimum  $CWA(R) \geq 0.1$ . This set

of rules serves as the testbed for our confidence estimators (thus, the same rules are used over all  $K$  and estimators). See Appendix E for the rule list.

The **applied selection mechanisms** differ in two ways. First, they either explicitly uphold the PCA (by selecting subjects  $s$  in **Q3**) or not (by selecting triples in **Q1**, **Q2**). For functional  $p$ , the PCA always holds by definition. Second, the mechanisms differ in which assumptions hold for the propensity scores: CWA,  $\text{SCAR}_p$  or SAR. Under  $\text{SCAR}_p$ ,  $c_p$  is varied. Two SAR mechanisms are considered: 1)  $\text{SAR}_{\text{group}}$  where the subjects of the triples are divided into two groups  $S_q, S_{-q}$  (e.g., actors and non-actors), each with a constant propensity score  $c_q, c_{-q}$ , and 2)  $\text{SAR}_{\text{pop}}$  where a triple’s propensity score is a logistic function of the number of facts in which the subject occurs, thus reflecting the subject’s *popularity*:

$$\begin{aligned} \#(s, p) &= |\{(s, q, \cdot) \in I\} \cup \{(\cdot, q, s) \in I\}|, q \neq p \\ e(\langle s, p, o \rangle) &= \max \left[ \frac{2}{1 + e^{-k \cdot \#(s, p)}} - 1, e_{\min} \right] \end{aligned}$$

More popular  $s$  have a higher  $e$ . The scaling factor  $k$  determines how often  $s$  must occur for a given  $e(\langle s, p, o \rangle)$ . Choosing  $e_{\min} > 0$  allows unpopular  $s$  to be selected.

When a rule’s coverage  $\mathbf{R}$  changes by applying the rule to different  $K$ , not only the estimators but also the rule’s actual confidence  $\text{conf}$  can change. In order to keep  $\text{conf}$  constant, the chosen selection mechanisms should not affect  $\mathbf{R}$ . Therefore, 1) only non-recursive rules are considered, and 2) each selection mechanism is applied to a single  $p \in \mathcal{P}$  at a time; the facts of  $\mathcal{P} \setminus \{p\}$  are completely included in  $K$  (cfr. CWA). This way, we can vary  $e(\cdot)$  for  $p$  while keeping the confidence  $\text{conf}(R)$  constant.

As **evaluation metric**, the Brier score  $\mathbb{E}_R[\widehat{\text{conf}}(R) - \text{conf}(R)]^2$  is chosen (with  $\widehat{\text{conf}}$  any estimator); this is a standard way of evaluating probability estimates.

For **Q1** and **Q2**, we compare  $\text{ICW}(R)$ <sup>4</sup> and  $\text{IPW}(R)$  to  $\text{CWA}(R)$  and  $\text{PCA}(R)$ . For **Q3**, we compare  $\text{IPW-PCA}(R)$  to  $\text{PCA}(R)$ , as the former modifies the latter to consider propensity scores.

Propensity scores are required to calculate  $\text{IPW-PCA}(R)$ . We use correct propensity scores  $e(\cdot)$  for the idealized scenarios in **Q1** and **Q3** and noisy versions  $\hat{e}$  for **Q2**.

More details about the exact setup can be found in Appendix C. Our source code is publicly available.<sup>5</sup>

## 5.2 Results

**(Q1)** Does using the ground truth propensity scores lead to a better confidence estimate? Table 1 shows the Brier scores for the estimators under  $\text{SCAR}_p$ ,  $\text{SAR}_{\text{group}}$  and  $\text{SAR}_{\text{pop}}$ . Only the leftmost  $\text{IPW}(R)$  column is relevant for **Q1**, i.e., the column with superscript **Q1**. The table shows that using correct propensity scores under  $\text{SCAR}_p$  and SAR results in a much better  $\text{conf}$  estimate: the Brier score for  $\text{IPW}(R)$  is often orders of magnitude lower than for the other estimators. (See also Tables 3, 4 and 5 in Appendix D.)

<sup>4</sup> For  $c_p$  the average  $e(\cdot)$  over all  $p$ -triples in  $K$  is used. Note that under  $\text{SCAR}_p$ ,  $\text{ICW}(R) = \text{IPW}(R)$ .

<sup>5</sup> <https://github.com/ML-KULEuven/KBC-as-PU-Learning>

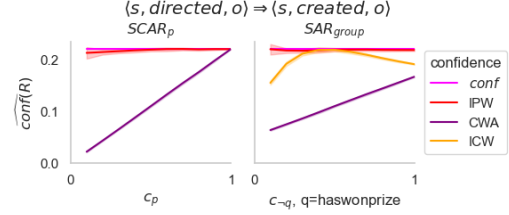


Figure 5: For a single rule under  $\text{SCAR}_p$  (left), dividing  $\text{CWA}(R)$  by  $c_p$  is (trivially) as good as  $\text{ICW}(R) = \text{IPW}(R) \approx \text{conf}(R)$  for all  $c_p$ . Under  $\text{SAR}_{\text{group}}$  (right), dividing  $\text{CWA}(R)$  by  $c_p$  fails to account for the bias for most  $c_{-q}$ <sup>4</sup>. Correct propensity scores are used, and  $c_q = 0.5$  for  $\text{SAR}_{\text{group}}$ .  $\text{conf}(R)$  is hidden by  $\text{IPW}(R)$ .

Figure 5 zooms in on a single rule. It shows how, under  $\text{SCAR}_p$ , both  $\text{IPW}(R)$  and  $\text{ICW}(R)$  almost perfectly compensate for  $\text{CWA}(R)$ ’s underestimation. However, under  $\text{SAR}_{\text{group}}$ ,  $\text{ICW}(R)$  does not recover  $\text{conf}(R)$  for most  $c_{-q}$ , while  $\text{IPW}(R)$  still does.

This illustrates how bad  $\text{CWA}(R)$  can be by ignoring the observation bias. If a learner merely *ranks* rules predicting the same predicate  $p$  under the simple  $\text{SCAR}_p$ , all  $\text{CWA}(R)$  are biased with the same constant  $c_p$ . Under this simple setting,  $\text{CWA}(R)$  works as well as  $\text{ICW}(R)$  (and hence  $\text{IPW}(R)$ ). However,  $\text{CWA}(R)$  fails under more complex settings (e.g.,  $\text{SAR}_{\text{group}}$ ), where using propensity scores allows  $\text{IPW}(R)$  to be clearly superior to  $\text{CWA}(R)$ .

Figure 6 (left) shows  $\text{PCA}(R)$  and  $\text{IPW-PCA}(R)$  for a single rule under  $\text{SCAR}_p$ , for both  $p$  (*person, diedin, place*) and its inverse  $p^{-1}$  (*place, wheredied, person*). Here, the PCA holds for  $p$  (a person dies in at most 1 place), but not for  $p^{-1}$  (only a fraction  $c_p$  of all the people who died somewhere are in  $K$ ). Therefore,  $\text{PCA}_p(R)$  remains constant for most  $c_p$ , differing from  $\text{conf}(R)$  with a constant factor  $\text{bias}_{y(s)=0}(R) = |\mathbf{R}|/|\mathbf{R}_s^+|$ . In contrast,  $\text{bias}_{p \in \mathcal{A}}$  causes  $\text{PCA}_{p^{-1}}(R)$  to vary with  $c_p$ .

The 3 factors that can cause  $\text{PCA}(R)$  to be biased interact; if they are unknown in advance, it is difficult to say how well  $\text{PCA}(R)$  will perform. For example in Figure 6 (left),  $\text{PCA}_{p^{-1}}(R)$  is equal to  $\text{conf}(R)$  at approximately  $c_p = 0.7$  due to an ‘accidental’ combination of these dimensions. However, if the PCA holds and  $|\mathbf{R}|/|\mathbf{R}_s^+|$  is close to 1,  $\text{PCA}(R)$  will be close to  $\text{conf}$  for all  $c_p$ , as shown with  $\text{PCA}_p(R)$ . Figure 6 (right) illustrates  $\text{PCA}(R)$  under SAR. The difference between  $\text{PCA}_p(R)$  and  $\text{conf}(R)$  is equal to  $|\mathbf{R}|/|\mathbf{R}_s^+|$  for the  $\text{SCAR}$  point ( $c_q = c_{-q}$ ), but changes when varying the relative number of triples in  $S_q$  and  $S_{-q}$  in  $K$  (see also the results for **Q3**).

This illustrates how complex the behavior of the  $\text{PCA}(R)$  is. Its disadvantage is its dependence on different interacting biases. In contrast, using propensity scores allows  $\text{IPW}(R)$  to be close to  $\text{conf}(R)$  in all settings.

**(Q2)** Can noisy propensity scores be used to improve confidence estimates? The rightmost IPW columns in Table 1 (with superscript **Q2**) show Brier scores for  $\text{IPW}(R)$  with respectively 0.1 and  $-0.1$  added as noise to the correct propen-

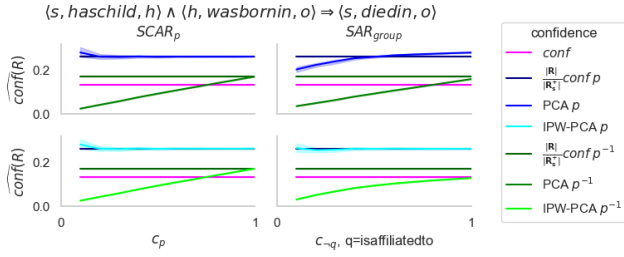


Figure 6: PCA holds for  $p = diedin$  but not for  $p^{-1} = wheredied$ . Under  $SCAR_p$  (left),  $PCA_p(R)$  for the example  $R$  differs from  $conf$  with  $|\mathbf{R}|/|\mathbf{R}_s^+|$  for most  $c_p$ , while  $PCA_{p^{-1}}(R)$  also varies with  $c_p$ . Under  $SAR_{group}$  with  $c_q = 0.5$  (right),  $PCA_p(R)$  also varies with  $c_{-q}$ , which  $IPW-PCA_p(R)$  corrects (see also figure 7).

sity scores for  $SCAR_p$  and  $SAR_{group}$ . For  $SAR_{pop}$ , the noisy propensity scores are obtained by increasing/decreasing  $k$  by 10%. The results show that  $IPW(R)$  is generally the best estimator for  $conf(R)$  if the noise is not too large. The exact differences between  $IPW(R)$ ,  $CWA(R)$  and  $PCA(R)$  depend on 1) the specific selection mechanism affecting  $CWA(R)$  and  $PCA(R)$  as seen in **Q1**, and 2) the noisy propensity scores affecting  $IPW(R)$ .

**(Q3)** When the PCA holds,  $PCA(R)$  becomes a better estimate as  $bias_{PCA}(R) = 1$ . Are there then still situations in which  $PCA(R)$  can be improved by using propensity scores? Here, we compare  $IPW-PCA(R)$  and  $PCA(R)$  under  $SAR_{group}$  and explicitly uphold the PCA for non-functional  $p$  by selecting subjects, e.g., if a person is a *citizen* of multiple countries, then either all or none of its triples are selected. We consider rules with predictions in both subject groups  $S_q$  and  $S_{-q}$ , with neither group dominating in size:  $0.3 \leq (|\mathbf{R} \cap S_q|)/|\mathbf{R}| \leq 0.7$ . We only include rules for which the group-local confidence (the confidence considering only the predictions in a group) differs by at least 0.1. The total confidence is the weighted mean of the group-local confidences where the weights are the fraction of predictions per group. By varying  $c_{-q}$  for a fixed  $c_q$ ,  $bias_{e(s)}(R)$  is varied (while  $bias_{y(s)=0}(R) = |\mathbf{R}|/|\mathbf{R}_s^+|$  remains constant): the relative number of subjects (and thus triples) in  $K$  belonging to each group varies, and is different from  $I$  for  $c_q \neq c_{-q}$ . Consequently, the total  $PCA(R)$  moves towards the group-local  $PCA(R)$  of the overrepresented group (Figure 7). However,  $IPW-PCA(R)$  remains relatively constant. Table 2 shows the Brier scores for this specific scenario. The results highlight how well  $PCA(R)$  works under PCA without needing propensity scores: although  $IPW-PCA(R)$  is mostly better, its improvement is rather small.

In conclusion,  $CWA(R)$  and  $PCA(R)$  fail to account for general observation biases in contrast to  $IPW(R)$  and  $IPW-PCA(R)$ , which make them explicit through propensity scores.

## 6 Related Work

Several **confidence measures for KBC** have been introduced to address the problem of dealing with the lack of

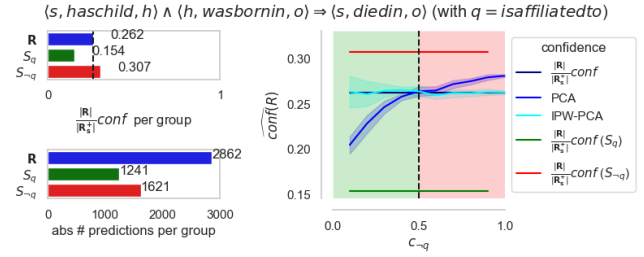


Figure 7: If the PCA holds under  $SAR_{group}$ ,  $PCA(R)$  for a rule with different group-local confidences (upper left) but a similar number of predictions per group (lower left) changes with  $c_{-q}$  towards  $|\mathbf{R}|/|\mathbf{R}_s^+| \cdot conf(R)$  of the overrepresented group, while  $IPW-PCA(R)$  remains more constant (right).  $c_q = 0.5$ .

negative examples (Galárraga et al. 2013; Pellissier Tanon et al. 2017; Zupanc and Davis 2018), but none of them handle general observation biases. The confidence measures from Zupanc and Davis and Pellissier Tanon et al. were omitted from the discussion, because, the former consists of an ad-hoc pipeline, and the latter consistently underestimates the confidence (see Appendix F). Other confidences were introduced with different goals:  $xconf$  (Zhou, Sadeghian, and Wang 2019) to limit computation effort and smooth confidence to cope with rules with low support (Meilicke et al. 2019).

Estimating propensity scores is closely related to estimating where the KB is more and less complete. The limited work on this topic combines several simple **completeness oracles** such as popularity, cardinality and change over time (Razniewski, Suchanek, and Nutt 2016; Galárraga et al. 2017). Soulet et al. (2018) estimate **representativeness** as deviation from an i.i.d. sample, but do not estimate which bias is causing the deviation nor propose a method for mitigating the bias at learning time.

Most of the work in **PU Learning** has been conducted under the  $SCAR$  assumption, where the labeled examples are an i.i.d. sample from the true positive examples (Elkan and Noto 2008). This is clearly violated by KBs. The more general  $SAR$  assumption allows for non-i.i.d. selection mechanisms, but needs additional assumptions to enable learning. Only a handful of such assumptions have been proposed (Kato, Teshima, and Honda 2018; Bekker, Robberechts, and Davis 2019; Gong et al. 2021). While none of these assumptions are sufficient for the KBC setting, the notion of explicitly modeling the selection mechanism inspired this paper.

**Recommender systems** solve a problem similar to KBC, but with only 1 predicate type ( $|\mathcal{P}| = 1$ ). Similar to our approach, Saito et al. (2020) and Gupta et al. (2021) adapt the PU learning loss function from Bekker, Robberechts, and Davis (2019). Amongst others, our paper differs from these works in 1) our unifying view on how KBs are constructed (including the taxonomy of assumptions) that allows analyzing the conditions under which evaluation metrics can be evaluated, 2) our analysis of the commonly used CWA and PCA estimators and 3) our newly proposed IPW(-PCA)

		# rules	CWA <sup>Q1,2</sup>	PCA <sup>Q1,2</sup>	ICW <sup>Q1,2</sup>	IPW <sup>Q1,2</sup>	IPW $-\Delta$ <sup>Q2</sup>	IPW $+\Delta$ <sup>Q2</sup>
SCAR <sub>p</sub>	$c_p = 0.3$	47	292.5	192.7	<b>4.6</b>	<b>4.6</b>	154.2	40.9
	$c_p = 0.7$	47	53.8	173.3	<b>1.1</b>	<b>1.1</b>	18.1	10.0
SAR <sub>group</sub>	$c_q = 0.5, c_{-q} = 0.3$	33	189.5	155.6	8.1	<b>3.4</b>	50.3	14.0
	$c_q = 0.5, c_{-q} = 0.7$	33	83.8	155.1	4.0	<b>1.8</b>	6.9	4.6
SAR <sub>pop</sub>	$k = 0.01$	47	458.6	264.2	168.5	62.8	81.8	<b>56.6</b>
	$k = 0.1$	47	172.3	182.7	51.0	<b>3.5</b>	7.4	6.2

Table 1: Results for **Q1** and **Q2** (see superscript).  $[\widehat{\text{conf}} - \text{conf}]^2 \cdot 10^4$  under SCAR<sub>p</sub>, SAR<sub>group</sub> and SAR<sub>pop</sub>. Results are averaged over  $p$ , the rules and (for SAR<sub>group</sub>)  $q$ . The 3 IPW confidence columns: 1 with correct  $\hat{e} = e$  (left) and 2 with noisy  $\hat{e} \neq e$  (middle and right). For SCAR<sub>p</sub>, noisy  $\hat{c}_p = c_p \pm 0.1$ . For SAR<sub>group</sub>, noisy  $\hat{c}_{-q} = c_{-q} \pm 0.1$ . For SAR<sub>pop</sub>, the noisy  $\hat{e}$  are obtained by using  $\hat{k} = k \pm 0.1k$ .

$p$	# R	$c_{-q} = 0.3$				$c_{-q} = 0.7$			
		PCA	IPW-PCA	IPW-PCA $-\Delta$	IPW-PCA $+\Delta$	PCA	IPW-PCA	IPW-PCA $-\Delta$	IPW-PCA $+\Delta$
dealwith	1	22.4	<b>16.9</b>	13.5	19.9	<b>9.1</b>	12.4	10.7	14.0
diedin	1	3.9	<b>1.6</b>	3.7	2.0	1.4	<b>0.5</b>	0.7	0.7
happenedin	1	6.6	<b>1.7</b>	4.0	3.3	2.3	<b>0.7</b>	1.0	1.1
iscitizenof	2	<b>13.1</b>	14.5	20.0	13.2	11.0	<b>9.9</b>	10.2	10.0
isleaderof	1	58.1	<b>55.5</b>	58.0	56.6	74.1	<b>72.0</b>	72.7	71.6
ispoliticianof	3	<b>8.0</b>	9.3	16.9	7.8	9.2	<b>8.3</b>	8.4	8.5
livesin	1	7.0	<b>6.8</b>	8.5	6.7	4.7	<b>4.0</b>	4.1	4.1
participatedin	1	16.0	<b>11.6</b>	8.1	14.1	<b>10.9</b>	13.2	12.1	14.2

Table 2: Results for **Q3**.  $[\widehat{\text{conf}} - |\mathbf{R}|/|\mathbf{R}_s^+| \cdot \text{conf}]^2 \cdot 10^4$  for SAR<sub>group</sub> with PCA upheld (avg. over  $q$  and rules predicting  $p$ ).  $c_{-q} \in \{0.3, 0.7\}$  with  $c_q = 0.5$  Bold is best per  $c_{-q}$  and  $p$ . Three *IPW-PCA*( $R$ ) columns for  $\hat{c}_{-q} = c_{-q} + \Delta$ ,  $\Delta \in \{0, \pm 0.1\}$ . Rules are included if  $0.3 \leq (|\mathbf{R} \cap S_q|)/|\mathbf{R}| \leq 0.7$  and the difference in group-local  $|\mathbf{R}|/|\mathbf{R}_s^+| \cdot \text{conf}(R)$  is at least 0.1.

confidence measures.

Biases are mostly studied in the **fairness** literature (Mehrabi et al. 2019; Barocas, Hardt, and Narayanan 2019), with the aim to learn bias-free models. This paper, in contrast, does not enforce certain ideals, but rather aims to increase model quality by being conscious of observation biases. This is a recent perspective in fairness literature (Blum and Stangl 2020).

## 7 Conclusion

We investigated rule evaluation, specifically confidence estimation, for knowledge base completion in the face of observation biases. Our theoretical and empirical analysis has shown that ignoring the observation bias results in biased confidence estimates. Yet, this is exactly what existing methods do. We have proposed two new confidence estimators that can mitigate known biases by using propensity scores that quantify how likely a fact is to be included in the KB. We have shown that these estimators are unbiased with respect to the observation bias. Our experiments showed that the Brier score of our *IPW*( $R$ ) measure is often orders of magnitude lower than those of the other estimators when observation biases are present. Our metric even outperforms the others when it has inexact values for the propensity scores.

## Acknowledgements

This research received funding from the Flemish Government under the ‘‘Onderzoeksprogramma Artificiële Intelligentie (AI) Vlaanderen’’ programme. Jonas Schouterden is supported by the KU Leuven Research Fund (C14/17/070). Jessa Bekker is also supported by the Research Foundation - Flanders under the Data- driven logistics project (FWO-S007318N), and Jesse Davis is also supported by the Research Foundation - Flanders (G0D8819N).

## References

- Auer, S.; Bizer, C.; Kobilarov, G.; Lehmann, J.; Cyganiak, R.; and Ives, Z. 2007. Dbpedia: A nucleus for a web of open data. In *The semantic web*, 722–735. Springer.
- Barocas, S.; Hardt, M.; and Narayanan, A. 2019. *Fairness and Machine Learning*. fairmlbook.org. <http://www.fairmlbook.org>.
- Bekker, J.; and Davis, J. 2020. Learning from positive and unlabeled data: A survey. *Machine Learning*, 109(4): 719–760.
- Bekker, J.; Robberechts, P.; and Davis, J. 2019. Beyond the Selected Completely At Random Assumption for Learning from Positive and Unlabeled Data. In *ECML PKDD: Joint European Conference on Machine Learning and Knowledge Discovery in Databases*. Springer.



- Blum, A.; and Stangl, K. 2020. Recovering from Biased Data: Can Fairness Constraints Improve Accuracy? In Roth, A., ed., *1st Symposium on Foundations of Responsible Computing (FORC 2020)*, volume 156 of *Leibniz International Proceedings in Informatics (LIPIcs)*, 3:1–3:20. Dagstuhl, Germany: Schloss Dagstuhl–Leibniz-Zentrum für Informatik. ISBN 978-3-95977-142-9.
- Callahan, E.; and Herring, S. 2011. Cultural bias in Wikipedia content on famous persons. *J. Assoc. Inf. Sci. Technol.*, 62: 1899–1915.
- Dong, X. L.; Gabrilovich, E.; Heitz, G.; Horn, W.; Lao, N.; Murphy, K.; Strohmman, T.; Sun, S.; and Zhang, W. 2014. Knowledge Vault: A Web-Scale Approach to Probabilistic Knowledge Fusion. In *Proceedings of the 20th ACM SIGKDD International Conference on Knowledge Discovery and Data Mining*, KDD '14, 601–610. Association for Computing Machinery. ISBN 9781450329569.
- Elkan, C.; and Noto, K. 2008. Learning classifiers from only positive and unlabeled data. In *Proceedings of the 14th ACM SIGKDD international conference on Knowledge discovery and data mining*, 213–220. ACM.
- Fürnkranz, J.; Gamberger, D.; and Lavrač, N. 2014. *Foundations of Rule Learning*. Springer Publishing Company, Incorporated. ISBN 3642430465, 9783642430466.
- Galárraga, L.; Razniewski, S.; Amarilli, A.; and Suchanek, F. M. 2017. Predicting Completeness in Knowledge Bases. *Proceedings of the Tenth ACM International Conference on Web Search and Data Mining*.
- Galárraga, L.; Teflioudi, C.; Hose, K.; and Suchanek, F. M. 2015. Fast rule mining in ontological knowledge bases with AMIE+. *The VLDB Journal*, 24: 707–730.
- Galárraga, L. A.; Teflioudi, C.; Hose, K.; and Suchanek, F. M. 2013. AMIE: Association Rule Mining Under Incomplete Evidence in Ontological Knowledge Bases. In *Proceedings of the 22nd International Conference on World Wide Web*, WWW '13, 413–422. New York, NY, USA: ACM. ISBN 978-1-4503-2035-1.
- Gong, C.; Wang, Q.; Liu, T.; Han, B.; You, J. J.; Yang, J.; and Tao, D. 2021. Instance-Dependent Positive and Unlabeled Learning with Labeling Bias Estimation. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 1–1.
- Gupta, S.; Wang, H.; Lipton, Z. C.; and Wang, Y. 2021. Correcting Exposure Bias for Link Recommendation. In *ICML*.
- Kato, M.; Teshima, T.; and Honda, J. 2018. Learning from positive and unlabeled data with a selection bias. In *International Conference on Learning Representations*.
- Lajus, J.; Galárraga, L.; and Suchanek, F. M. 2020. Fast and Exact Rule Mining with AMIE 3. *The Semantic Web*, 12123: 36 – 52.
- Mahdisoltani, F.; Biega, J. A.; and Suchanek, F. M. 2015. YAGO3: A Knowledge Base from Multilingual Wikipedias. In *CIDR*.
- Mehrabani, N.; Morstatter, F.; Saxena, N.; Lerman, K.; and Galstyan, A. 2019. A survey on bias and fairness in machine learning. *arXiv preprint arXiv:1908.09635*.
- Meilicke, C.; Chekol, M. W.; Ruffinelli, D.; and Stuckenschmidt, H. 2019. Anytime Bottom-Up Rule Learning for Knowledge Graph Completion. In *Proceedings of the Twenty-Eighth International Joint Conference on Artificial Intelligence, IJCAI-19*, 3137–3143. International Joint Conferences on Artificial Intelligence Organization.
- Nickel, M.; Murphy, K.; Tresp, V.; and Gabrilovich, E. 2015. A review of relational machine learning for knowledge graphs. *Proceedings of the IEEE*, 104(1): 11–33.
- Pellissier Tanon, T.; Stepanova, D.; Razniewski, S.; Mirza, P.; and Weikum, G. 2017. Completeness-Aware Rule Learning from Knowledge Graphs. In *The Semantic Web – ISWC 2017*, volume 1, 507–525. ISBN 978-3-319-68288-4.
- Pezeshkpour, P.; Tian, Y.; and Singh, S. 2020. Revisiting Evaluation of Knowledge Base Completion Models. In *Automated Knowledge Base Construction*.
- Razniewski, S.; Suchanek, F.; and Nutt, W. 2016. But what do we actually know? In *Proceedings of the 5th Workshop on Automated Knowledge Base Construction*, 40–44.
- Rebele, T.; Suchanek, F.; Hoffart, J.; Biega, J.; Kuzey, E.; and Weikum, G. 2016. YAGO: A multilingual knowledge base from wikipedia, wordnet, and geonames. In *International semantic web conference*, 177–185. Springer.
- Saito, Y.; Yaginuma, S.; Nishino, Y.; Sakata, H.; and Nakata, K. 2020. Unbiased Recommender Learning from Missing-Not-At-Random Implicit Feedback. In *Proceedings of the 13th International Conference on Web Search and Data Mining*, 501–509.
- Soulet, A.; Giacometti, A.; Bouchou-Markhoff, B.; and Suchanek, F. M. 2018. Representativeness of Knowledge Bases with the Generalized Benford’s Law. In *International Semantic Web Conference*.
- Speranskaya, M.; Schmitt, M.; and Roth, B. 2020. Ranking vs. Classifying: Measuring Knowledge Base Completion Quality. In *Automated Knowledge Base Construction*.
- Vrandečić, D.; and Krötzsch, M. 2014. Wikidata: a free collaborative knowledgebase. *Commun. ACM*, 57: 78–85.
- Wagner, C.; Garcia, D.; Jadidi, M.; and Strohmaier, M. 2015. It’s a man’s Wikipedia? Assessing gender inequality in an online encyclopedia. In *Proceedings of the International AAAI Conference on Web and Social Media*.
- WWW Consortium. 2004. RDF Primer (W3C Recommendation 2004-02-10).
- Zhou, X.; Sadeghian, A.; and Wang, D. 2019. Mining Rules Incrementally over Large Knowledge Bases. *ArXiv*, abs/1904.09399.
- Zupanc, K.; and Davis, J. 2018. Estimating rule quality for knowledge base completion with the relationship between coverage assumption. In *Proceedings of the Web Conference 2018*, 1–9.