

[ICML2016](#)**International Conference on Machine Learning**

June 19 – June 24, 2016, New York, United States

**Reviews For Paper****Paper ID** 491**Title** Thompson Sampling and Compositions in Knowledge Bases with Uncertainty**Masked Reviewer ID:** Assigned\_Reviewer\_3**Review:**

Question	
Clarity (Assess the clarity of the presentation and reproducibility of the results.)	Poor (Hard to follow)
Clarity - Justification	The combination of frequent grammatical errors, typos, and dense technical content make this paper rather difficult to read and follow.
Significance (Does the paper contribute a major breakthrough or an incremental advance?)	Below Average
Significance - Justification	This paper appears to be a reasonable, although incremental, advance upon previous work on knowledge base completion.
Detailed comments. (Explain the basis for your ratings while providing constructive feedback.)	I am not familiar with work in this area, and as this paper is presented in a rather dense and technical fashion, it is difficult for me to comment directly on the approach taken. However I do have a couple of questions after reading this paper. Firstly, what is the use case for the active KB completion model? Given that any useful KB will contain millions of triples, it seems unlikely that an approach requiring a human in the loop will scale. Second, it is unclear to me what the compositional triples add to the model, as they do not appear to contain any information which is not already available from the transitive semantics of the base triples. Does this merely make inference easier?
Overall Rating	Weak reject
Reviewer confidence	Reviewer's evaluation is an educated guess

**Masked Reviewer ID:** Assigned\_Reviewer\_4**Review:**

Question	
Clarity (Assess the clarity of the presentation and reproducibility of the results.)	Above Average
Clarity -	Overall, I believe this is a clearly written paper, which is easy to follow and that I enjoyed reading.

Justification	I think that a better discussion of previous work could be added, in particular the relation between the proposed method and other compositional models.
Significance (Does the paper contribute a major breakthrough or an incremental advance?)	Below Average
Significance - Justification	Since I am not an expert in Bayesian modeling, I kind of fail to see what are the main contributions of the paper, and their significance. Indeed, the models are not particularly new, neither is Thomson sampling. I believe the main contribution is the compositional aspect of the model, which should be emphasized by the authors. (Compositional models of knowledge bases have already been proposed, using non Bayesian approach).
Detailed comments. (Explain the basis for your ratings while providing constructive feedback.)	I have a couple of questions for the authors: - Most experiments were performed on rather small datasets. How does this approach would scale to larger knowledge bases, such as Freebase? - While the logit model seems more natural, it is outperformed by the Gaussian model on most experiments. Do you have an explanation for that?
Overall Rating	Weak accept
Reviewer confidence	Reviewer's evaluation is an educated guess
Please acknowledge that you have read the author rebuttal. If your opinion has changed, please summarize the main reasons below.	I have read the author rebuttal.

**Masked Reviewer ID:** Assigned\_Reviewer\_5

**Review:**

Question	
Clarity (Assess the clarity of the presentation and reproducibility of the results.)	Excellent (Easy to follow)
Clarity - Justification	Most of the paper is easy to follow. Some of the high level ideas described in the introduction, such as the connection between explore-exploit trade-off and knowledge completion/extraction, are a bit ambiguous at first, but become clear later as more details come.
Significance (Does the paper contribute a major	Above Average

breakthrough or an incremental advance?)	
Significance - Justification	The proposed approach is a nice combination of Bayesian inference, tensor modeling and Thompson sampling. Based on the empirical results, it may have a chance to have wider practical impact. However, it is not as significant as a major breakthrough.
Detailed comments. (Explain the basis for your ratings while providing constructive feedback.)	<p>This paper nicely combines Bayesian, tensor modeling, and Thompson sampling to derive an active learning/inference method for knowledge graph completion. Although the individual components are not new, the proposed combination seems novel and empirically successful. One question is on the value of the two extensions for handling compositionality. They substantially increase modeling and computational complexity, but the empirical gain is only marginal. Also, only short compositions were considered in the experiments, which may be due to computational reasons. Overall, this paper seems to be a nice contribution to Bayesian inference/active learning in knowledge graphs. Some detailed comments are below.</p> <p>* Related work: In addition to those already discussed, two lines of work seem relevant: (1) active search [1,2], a form of active learning that aims to maximize the number of true positive instances returned to users. (2) Theoretical analysis of Thompson or Posterior sampling [3].</p> <p>* In Section 6.3: How was ROC-AUC computed exactly, or how was the ROC curve generated?</p> <p>* In Figure 4, the x-axis seems to extend beyond the total number of triples in the data (Table 3).</p> <p>References:  [1] Garnett et al., Bayesian Optimal Active Search and Surveying, ICML 2012  [2] Wang et al., Active Search on Graphs, KDD 2013  [3] Russo and Van Roy, Learning to Optimize Via Posterior Sampling, Mathematics of Operations Research. Vol. 39. No. 4, pp. 1221-1243, 2014.</p>
Overall Rating	Weak accept
Reviewer confidence	Reviewer is knowledgeable
Please acknowledge that you have read the author rebuttal. If your opinion has changed, please summarize the main reasons below.	I've read the rebuttal.

**Masked Reviewer ID:** Assigned\_Reviewer\_6

**Review:**

Question	
Clarity (Assess the clarity of the presentation and reproducibility of the results.)	Below Average

Clarity -  
Justification

The manuscript is difficult to read due to numerous grammatical errors. There is clearly clarity of thought behind what the authors are proposing, but the ideas are significantly muddled by the presentation. The text is riddled with simple mechanical problems such as subject-verb agreement errors. I do not believe the text is up to the standards of an elite international conference such as ICML. Below I list incomplete corrections for the first two sections; I would recommend the authors seek copyediting help before resubmitting.

#### Abstract:

line 12 (12): Human -> A Human (or Humans)

20: tensors, that explicitly models -> tensors that explicitly model [delete comma, correct verb]

27-28: sampling provide effective exploit-explore -> sampling provides effective exploitation--exploration [correct verb, rephrase "exploit-explore"] (I would suggest replacing all instances of "exploit-explore" with "exploitation--exploration")

#### Introduction:

35: contact facts: "contact" is incorrect; "contain facts" or "store facts" perhaps?

42: facts, the first learning -> facts: the first is learning

63: propose -> proposed

64-65 triples, however the algorithm find -> triples; however, the algorithm finds

65: high reconstruction: "high" is incorrect; perhaps "faithful?"

65: I propose recasting "find[s] it difficult to achieve..." as "the algorithm had problems simultaneously achieving high recall and faithful reconstruction"

66-67: Having an active...: This sentence is unclear in that I don't really follow the reasoning. \_Why\_ would an exploration/exploitation strategy help marry these two concepts? This is especially confusing because it is not yet clear at this point in the manuscript what exactly "exploitation" or "exploration" entails in the context of knowledge base completion/extraction, and whether it is sensible to interpret them the same in both contexts. As this is effectively the thesis of this work, I would suggest expanding and clarifying your discussion on this point.

70-73: The transition to the sentence "Also known as paths in knowledge graphs" is rather jarring.

77: techniques are recently developed -> techniques were recently developed

87-92: "The probabilistic model...": This sentence is pretty sprawling and I would suggest recasting as two sentences. Perhaps: "...knowledge completion. Here we use Thompson sampling, an approach...."

97: with and without and -> with and without

102: sampling provide -> sampling provides

103: triples, we also -> triples; we also

106: In passive learning setting -> In the passive learning setting

108: help when training set -> help when the training set

113-114: outperforms active learning strategies... -> outperforming active learning strategies that focus solely on exploitation or exploration

116: bridging the gaps -> bridging the gap

117: knowledge compositions -> knowledge composition

#### Related Work:

130-134: The description of "Bayesian/Non-Bayesian" is somewhat arbitrary and confusing. Point estimates are often used in Bayesian settings. A frequentist would probably disagree that they have no method

	<p>of quantifying uncertainty. In fact, I'd suggest leaving the entire discussion out; it's unnecessary for the ICML audience.</p> <p>140: Relational learning problems for: is there a word or phrase missing after "for?"</p> <p>141: is common when dataset -&gt; is common when a dataset</p> <p>142: such as in -&gt; such as</p> <p>143: in recommenders systems setting -&gt; in the recommender system setting</p> <p>144: edges in the graph has labels -&gt; edges in the graph have labels</p> <p>149-198: These two paragraphs are quite unclear, especially the description of the missing N, A, C entry and the sprawling "Note that we re-formulate..." sentence, which contains two semicolons connecting three somewhat disconnected ideas.</p> <p>162: compositions objectives -&gt; composition objective</p> <p>163: in the probabilistically -&gt; probabilistically</p>
Significance (Does the paper contribute a major breakthrough or an incremental advance?)	Below Average
Significance - Justification	<p>The main contributions of this paper are:</p> <ol style="list-style-type: none"> <li>1. A Bayesian reformulation of the RESCAL tensor factorization method for modeling relational knowledge bases, allowing for a quantification of uncertainty in missing entries and the latent features, BRESCAL. Two variations are proposed, one using a Gaussian observation model and one using a Bernoulli observation model.</li> <li>2. A Thompson sampling approach for actively querying missing entries in the relational tensor, given the above model.</li> <li>3. A particle filtering approach to avoid constantly reperforming the inference over the latent parameters of the BRESCAL model.</li> <li>4. An extension to the above reason about compositional relations.</li> </ol> <p>I will handle these in turn. The first contribution (1) is relatively minor, in my opinion, as it is a straightforward and predictable extension of (Bayesian) probabilistic matrix factorization (BPMF) to tensor factorization. The Gibbs sampling updates are very similar to the original BPMF paper (Salakhutdinov/Mnih, ICML 2008). The second contribution (2) is perhaps somewhat interesting but not a huge leap methodologically over what is currently done and not established empirically by the authors to be useful or needed. The third contribution (3) is similar in that it applies straightforward, simple SMC methods to the proposed model, which is only marginally novel. The authors also never establish that this is necessary. Finally, the fourth contribution (4) is somewhat interesting and perhaps the most-novel contribution of the paper. The weight of this contribution, however, is hard to judge as it is hard to build intuition from either the discussion or the experiments about whether the additive or multiplicative model should be preferred (the experiments prefer the latter).</p>
	<p>I believe the focus of this paper, active relational knowledge base completion, is an interesting problem, although I wish the authors had done a better job at motivating the importance of the problem (e.g., how is the oracle implemented or realized in typical applications?). In addition to concerns about the clarity and novelty of the methods presented in this paper, I have a number of higher-level comments and questions regarding the work.</p> <p>One of the main contributions of the paper is a Thompson sampling</p>

Detailed  
comments.  
(Explain the basis  
for your ratings  
while providing  
constructive  
feedback.)

approach to actively querying a relational knowledge base, wherein we first sample from the posterior distribution of the latent features in the BRESICAL model (E and R), then select the currently unlabeled tuple with the highest expected value given the sampled E and R (line 328) to query. Thompson sampling is not novel; however, I have not seen it applied to the tensor completion task. That said, the authors never actually establish that this is necessary at all, nor that the gains in performance they see are due to the Thompson sampling approach or simply the model used. How much does the performance degrade if we forgo sampling E and R and instead simply use MAP estimates of these? At the moment, the Thompson sampling feels like extra complexity that is possibly not even needed. A second, more minor issue I have with the discussion surrounding the Thompson sampling proposal is that it is an instance of somehow achieving an exploration/exploitation trade off. How? The argmax in 328 looks an awful lot like exploitation. Is the exploration supposed to be coming from the sampling over E and R? If so, I would consider that a highly nonstandard use of the term.

Another nominal contribution is the proposed sequential Monte Carlo technique wherein we maintain a set of (E, R) particles that we propagate through time. The authors suggest this is necessary to avoid costly retraining after every query. However, the datasets considered in this paper are all quite small, and I can't imagine running the inference takes a terribly long time, especially if we begin the MCMC chain where we left off (which, presumably, should still be a reasonable location in the parameter space; the posterior can't change too much after a single observation). Furthermore, if the oracle is costly (e.g., I have to ask a human whether the relation exists), and presumably it should be for active learning to be worthwhile at all, then shouldn't I have plenty of time to do as much exhaustive MCMC as I'd like? Even a scenario existed where the particle filtering were actually necessary, due to a high-bandwidth oracle perhaps, there is no investigation in this paper into the effect of this choice in the proposed framework. How much do the results degrade using the particle filtering framework from a method using extensive sampling?

I feel the manuscript is incomplete without an investigation into the effects and importance of these two choices.

The idea of including compositional relations in the model struck me as intriguing, but I was not able to come away from section 5 with a good intuition about whether the additive or multiplicative model was more appropriate. Experiments suggest the latter is better, so why propose the additive model at all? There is also some sloppy math/notation in this section; for example, equation (13) is the same as equation (5) but should probably be closer to equation (11). A final point of confusion here was why the authors claim the triples in the expanded tensor cannot be queried. Why not? Why can't a system propose a compositional relation and have it checked by an oracle? Furthermore, if that is not a possible query, why allow it in section 6.2?

Another issue with the compositional model is that the experiments seem to suggest that it hurts performance on the completion task across a variety of datasets (although it, paradoxically, helps on pure prediction tasks). If that's true, why include this material at all? It really detracts from the main message of the paper.

Another point of confusion I had with this paper was the inclusion of both the Gaussian and the logistic observation models. The latter is clearly more suited to the task, although I can believe that the former works better in practice. That said, the way the Gaussian model is treated in the experiments is confusing. Regret in Figure 1(c) and 1(d) is clearly being

	<p>measured against the highest-valued entry in the sampled X matrix (around 33). The magnitude of this entry is clearly meaningless, so why measure against it? It makes absolutely no sense in the context of this setting. Regarding Figure 1, is the passive curve also indicating performance across 10 repetitions? A brief experiment conducted on my part suggests there should be larger error bars.</p> <p>What were the priors used in section 6.3?</p> <p>Why are the regrets of the additive and multiplicative models in Figure 2 so vastly different? Was this experiment repeated? That it wasn't is the only possible explanation I can think of.</p> <p>Are the results in Figure 4 a single run of the experiment? If so, why? We could sample different test sets to get some variability in the results. If the results are for a single run, it's unclear whether the proposed method is truly as amazing as it appears or whether it just got lucky.</p> <p>Minor corrections/comments:</p> <ul style="list-style-type: none"> <li>- The first column of table 2 should probably have conditioning bars after the variables</li> <li>- The formula around line 318 is rendered more confusing than it needs to be by taking an expectation of an indicator; why not cast this as a more-readable probability?</li> </ul>
Overall Rating	Strong reject
Reviewer confidence	Reviewer is an expert
Please acknowledge that you have read the author rebuttal. If your opinion has changed, please summarize the main reasons below.	No.