

# Supplementary Material to “Augmenting and Tuning Knowledge Graph Embeddings”

**Robert Bamler\***

Department of Computer Science  
University of California, Irvine  
Irvine, CA 92617

**Farnood Salehi\***

École Polytechnique Fédérale  
de Lausanne (EPFL),  
Switzerland

**Stephan Mandt**

Department of Computer Science  
University of California, Irvine  
Irvine, CA 92617

## The Role of Parameter Uncertainty in the Proposed Hyperparameter Optimization

Bayesian inference and the idea of measuring uncertainty are somewhat uncommon in the literature for knowledge graph embeddings. To clarify why uncertainty is important in the proposed hyperparameter optimization (Section 3.1 of the main text), we compare the method here to a more naive approach that will turn out to fail because it ignores uncertainty. We discuss the failure of the naive approach and the benefit of estimating parameter uncertainty first intuitively and then more formally.

**Intuitive picture.** The variational EM algorithm maximizes (a lower bound on) the marginal likelihood  $p(\mathcal{S}'|\lambda)$ , Eq. 14 of the main text. In a more naive attempt to hyperparameter tuning without cross validation, one might be tempted to skip the marginalization over model parameters  $\mathbf{E}$  and  $\mathbf{R}$ . Instead, one might try to directly maximize the log joint probability  $\log p(\mathbf{E}, \mathbf{R}, \mathcal{S}'|\lambda)$ , i.e., minimize the loss  $L = -\log p(\mathbf{E}, \mathbf{R}, \mathcal{S}'|\lambda)$ , over  $\mathbf{E}$ ,  $\mathbf{R}$ , and hyperparameters  $\lambda$ . This approach, however, would lead to divergent solutions because the log joint probability is unbounded.

The log joint probability contains the log priors (first two terms on the right-hand side of Eq. 10 of the main text). Figure S1 shows the prior  $p(E_{ek}|\lambda_e^E)$ , Eq. 8 of the main text, of a single component  $E_{ek}$  of an entity embedding assuming, for simplicity, a real embedding space. With growing regularizer strength (increasing  $\lambda_e^E$ ), the prior becomes narrower and narrower. As the peak narrows, it also grows higher due to the normalization constraint. In the limit  $\lambda_e^E \rightarrow \infty$ , the prior collapses to an infinitely narrow and high  $\delta$ -peak at zero.

A hyperparameter tuning method that ignores posterior uncertainty can exploit the unbounded growth of the maximum of the prior to send the log joint distribution to infinity (i.e., the loss  $L \rightarrow -\infty$ ). Without any posterior uncertainty, one can set the model parameter  $E_{ek}$  precisely to

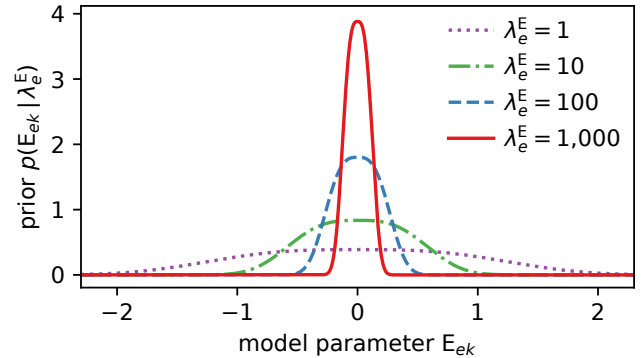


Figure S1: Prior (Eq. 8 of the main text with  $p = 3$ ) for a single model parameter  $E_{ek}$ . As the prior gets more peaked for growing regularizer strength  $\lambda_e^E$ , the height of the peak grows unboundedly. An (incorrect) hyperparameter optimization method that ignores parameter uncertainty could end up exploiting this unboundedness and diverge to  $\lambda_e^E \rightarrow \infty$ .

zero and then make the value of  $p(E_{ek}|\lambda_e^E)$  arbitrarily large by sending  $\lambda_e^E \rightarrow \infty$ .

We prevent this collapse of the prior to a  $\delta$ -peak by keeping track of parameter uncertainty. Admitting a nonzero uncertainty for  $E_{ek}$  no longer allows us to set  $E_{ek}$  precisely and deterministically to zero. Any slightly nonzero value of  $E_{ek}$  would have no support under a  $\delta$ -peaked prior.

**Formal derivation.** We now formalize the above intuitive picture and show that the specific variational approximation chosen in Eqs. 12-13 indeed suffices to prevent any divergent solutions.

Assuming again a real embedding space, the log prior of a single model parameter  $E_{ek}$  is given by (cf., Eq. 16 of the main text)

$$\log p(E_e|\lambda_e^E) = \frac{1}{p} \left[ \log \lambda_e^E - \lambda_e^E |E_{ek}|^p \right] + \text{const.} \quad (\text{S1})$$

As discussed above, setting  $E_{ek} = 0$  and sending  $\lambda_e^E \rightarrow \infty$  sends the right-hand side of Eq. S1 to infinity. This can

\*joint first authorship.

even be relaxed: the log prior diverges for  $\lambda_e^E \rightarrow \infty$  as long as we keep  $E_{ek}$  small enough, i.e., as long as  $E_{ek} = O((\lambda_e^E)^{-1/p})$ . This is why maximizing the log joint distribution over  $\lambda$  leads to divergent solutions.

Instead of maximizing the log joint distribution, the variational EM algorithm maximizes the ELBO, Eq. 14 of the main text. Note first that the ELBO itself is bounded from above by zero: the ELBO is a lower bound on the marginal log likelihood  $\log p(\mathbb{S}'|\lambda)$ , where  $p(\mathbb{S}'|\lambda) \leq 1$  since it is a discrete probability distribution.

Further, the ELBO has a maximum at finite values for the variational parameters and hyperparameters. Maximizing the ELBO ensures that each model parameter is associated with a nonzero uncertainty  $\sigma_{e/r}^{E/R} > 0$  since the entropy term

$$H[q_{\mu,\sigma}] = \sum_{e \in [N_e]} \log \sigma_e^E + \sum_{r \in [N_r]} \log \sigma_r^R + \text{const.} \quad (\text{S2})$$

imposes an infinite penalty if any  $\sigma_{e/r}^{E/R} \rightarrow 0$ . The entropy term in the ELBO thus has an additional regularizing effect. Combined with the other regularizing term in the ELBO, the expected log prior, we obtain for a given model parameter  $E_{ek}$  using Eq. S1, up to an additive constant,

$$\begin{aligned} & \mathbb{E}_{q_{\mu,\sigma}} \left[ \log p(E_{ek}|\lambda_e^E) - \log (q_{\sigma_{ek}^E, \mu_{ek}^E}(E_{ek})) \right] \\ &= \frac{1}{p} \log \lambda_e^E - \frac{\lambda_e^E}{p} \mathbb{E}_{q_{\sigma,\mu}} [ |E_{ek}|^p ] + \log \sigma_{ek}^E \quad (\text{S3}) \\ &= \frac{1}{p} \left[ \log \left( \lambda_e^E (\sigma_{ek}^E)^p \right) - \lambda_e^E \mathbb{E}_{\epsilon \sim \mathcal{N}(0,1)} \left[ |\mu_{ek}^E + \sigma_{ek}^E \epsilon|^p \right] \right] \end{aligned}$$

where, in the last equality, we made the dependency on  $\sigma_{ek}^E$  explicit by reparameterizing the normal distributed random variable  $E_{ek} = \mu_{ek}^E + \sigma_{ek}^E \epsilon$  in terms of a standard-normal distributed variable  $\epsilon$ .

Maximizing the right-hand side of Eq. S3 over  $\mu_{ek}^E$  yields  $\mu_{ek}^E = 0$  by symmetry, thus simplifying the objective to

$$\begin{aligned} & \frac{1}{p} \left[ \log \left( \lambda_e^E (\sigma_{ek}^E)^p \right) - \lambda_e^E (\sigma_{ek}^E)^p \mathbb{E}_{\epsilon \sim \mathcal{N}(0,1)} [ |\epsilon|^p ] \right] \\ &= \frac{1}{p} \left[ \log \left( \lambda_e^E (\sigma_{ek}^E)^p \right) - \lambda_e^E (\sigma_{ek}^E)^p c_p \right] \quad (\text{S4}) \end{aligned}$$

with some numerical constant  $c_p > 0$ . The right-hand side of Eq. S4 is structurally similar to the right-hand side of Eq. S1, which had divergent solutions: both are a difference between a logarithmic and a linear term in  $\lambda_e^E$ . However, while in Eq. S1, we were able to send the logarithmic term to infinity and still keep the linear term bounded, this is not possible in Eq. S4, which has the same argument (up to a constant  $c_p$ ) for the logarithmic and the linear term. The right-hand side of Eq. S4 has a maximum at a finite value for  $\lambda_e^E (\sigma_{ek}^E)^p$ . Thus, the variational EM algorithm avoids divergent solutions.