

---

# Supplementary Material: Variational Training for Large-Scale Noisy-OR Bayesian Networks

---

Geng Ji<sup>1,2</sup>   Dehua Cheng<sup>2</sup>   Huazhong Ning<sup>2,3</sup>   Changhe Yuan<sup>2,4</sup>  
Hanning Zhou<sup>2</sup>   Liang Xiong<sup>2</sup>   Erik B. Sudderth<sup>1</sup>

<sup>1</sup>UC Irvine   <sup>2</sup>Facebook AI Applied Research   <sup>3</sup>WeRide.ai   <sup>4</sup>CUNY Queens College

## A VARIATIONAL BOUND IS CONCAVE IN $r$

For each node  $i \in \{\mathcal{H} \cup \mathcal{O}^+\}$  of some document  $d$ , the subset of terms in the variational bound of Eq. (9) that depend on auxiliary variables  $r_i$  can be written as:

$$\mathcal{L}_{di}(r_i) = \sum_{k \in \mathcal{P}(i)} r_{k \rightarrow i} q_k [f(u_{k \rightarrow i}) - f(w_{0 \rightarrow i})].$$

The first partial derivative of this variational bound is

$$\frac{\partial \mathcal{L}_{di}}{\partial r_{k \rightarrow i}} = q_k \left( f(u_{k \rightarrow i}) - f(w_{0 \rightarrow i}) - \frac{w_{k \rightarrow i}}{r_{k \rightarrow i}} f'(u_{k \rightarrow i}) \right),$$

and its second partial derivatives equal

$$\frac{\partial^2 \mathcal{L}_{di}}{\partial r_{k \rightarrow i} \partial r_{\ell \rightarrow i}} = 0, \quad \frac{\partial^2 \mathcal{L}_{di}}{\partial r_{k \rightarrow i}^2} = q_k \frac{w_{k \rightarrow i}^2}{r_{k \rightarrow i}^3} f''(u_{k \rightarrow i}).$$

Here, the function

$$f''(a) = \frac{-\exp(a)}{(\exp(a) - 1)^2} < 0$$

is the second derivative of  $f(a)$ . Thus on the convex set of auxiliary parameters defined by Eq. (7), the (diagonal) Hessian matrix of  $\mathcal{L}_{di}$  is negative definite, and  $\mathcal{L}_{di}(r_i)$  is a strictly concave function of  $r_i$ .

## B INITIALIZATION OF $r$

We show that setting  $r_{k \rightarrow i} \propto w_{k \rightarrow i}$  globally optimizes our variational objective whenever the activation probabilities  $q_k$  for all parent nodes  $k \in \mathcal{P}(i)$  are equal. To prove this, note that optimizing Eq. (9) with respect to  $r_{k \rightarrow i}$  is equivalent to maximizing

$$\sum_{k \in \mathcal{P}(i)} r_{k \rightarrow i} \left[ f\left(w_{0 \rightarrow i} + \frac{w_{k \rightarrow i}}{r_{k \rightarrow i}}\right) - f(w_{0 \rightarrow i}) \right]. \quad (\text{B.1})$$

Given the non-negativity and normalization constraints in Eq. (7), we can apply Jensen's inequality in the opposite direction of typical variational derivations:

$$\begin{aligned} & \sum_{k \in \mathcal{P}(i)} r_{k \rightarrow i} \left[ f\left(w_{0 \rightarrow i} + \frac{w_{k \rightarrow i}}{r_{k \rightarrow i}}\right) - f(w_{0 \rightarrow i}) \right] \\ & \leq f\left(\sum_{k \in \mathcal{P}(i)} r_{k \rightarrow i} \left(w_{0 \rightarrow i} + \frac{w_{k \rightarrow i}}{r_{k \rightarrow i}}\right)\right) - f(w_{0 \rightarrow i}) \\ & = f\left(w_{0 \rightarrow i} + \sum_{k \in \mathcal{P}(i)} w_{k \rightarrow i}\right) - f(w_{0 \rightarrow i}). \end{aligned} \quad (\text{B.2})$$

The bound in the second line of Eq. (B.2) is achieved with equality if and only if  $w_{0 \rightarrow i} + \frac{w_{k \rightarrow i}}{r_{k \rightarrow i}}$  is constant for all parent nodes, which occurs when  $r_{k \rightarrow i} \propto w_{k \rightarrow i}$ .