

---

# Supplementary material for "Graph Bernoulli Pooling"

---

Paper Id: 3310

## 1 Bernoulli Sampling Optimization Objective Derivation Details

2 **Notations.**  $\phi$  denotes the parameter set of the BernPool module,  $\psi$  is the parameter set of the other  
3 modules.

4 **Mutual Information Maximizing.** BernPool aims to maximize the mutual information between  
5 learned subgraph embeddings and corresponding labels, which can be formulated as:

$$\zeta_{MI} = MI(\mathbf{y}, f_{\psi, \phi}(\mathcal{G}, \mathcal{S})) = MI(\mathbf{y}, \mathbf{f}), \quad (1)$$

6 where  $f_{\psi, \phi}$  represents the graph embedding process. Moreover, we introduce  $\mathbf{f}$  to denote  $f_{\psi, \phi}(\mathcal{G}, \mathcal{S})$   
7 for simplification, which distributes in the embedding space  $\mathcal{F}$ .  $\mathcal{F}$  is spanned by the resulted  
8 embeddings of  $f_{\psi, \phi}$  inferred based on the input  $\mathcal{G}$  and reference set  $\mathcal{S}$ . Then, based on the connection  
9 between the mutual information and entropy, the objective can be further written as:

$$\begin{aligned} & \arg \max MI(\mathbf{y}, \mathbf{f}) \\ & = \arg \max H(\mathbf{y}) - H(\mathbf{y}|\mathbf{f}), \end{aligned} \quad (2)$$

10 where  $H(\mathbf{y})$  can be just omitted from the objective as it is independent from  $\tilde{\psi}, \tilde{\phi}$ . We have the  
11 following derivation:

$$\begin{aligned} & \arg \max -H(\mathbf{y}|\mathbf{f}) \\ & = \arg \max \sum_i -H(\mathbf{y}|\mathbf{f} = \mathcal{F}_i)p(\mathcal{F}_i) \\ & = \arg \max \sum_i p(\mathcal{F}_i)\mathbb{E}_{\mathbf{y}|\mathcal{F}_i}(\log p(\mathbf{y}|\mathbf{f} = \mathcal{F}_i)), \end{aligned} \quad (3)$$

12 where  $p(\mathcal{F}_i)$  means the probability of the  $i$ -th observation in the embedding space and can be  
13 rationally assumed to conform to the uniform distribution.  $L$  denotes the number of observations. As  
14 the observation  $\mathcal{F}_i$  means to be inferred based on an input sample  $\mathcal{G}_i$  with  $\mathcal{S}$ , we further denote  
15  $p(\mathbf{y}|\mathbf{f} = \mathcal{F}_i)$  equally  $p_{\psi, \phi}(\mathbf{y}|\mathcal{G}_i, \mathcal{S})$ . The objective can be further written as:

$$\begin{aligned} & \arg \max \sum_i \frac{1}{L} \mathbb{E}_{\mathbf{y}|\mathcal{G}_i, \mathcal{S}}(\log p_{\psi, \phi}(\mathbf{y}|\mathbf{f} = \mathcal{F}_i)) \\ & = \arg \max \sum_i \mathbb{E}_{\mathbf{y}|\mathcal{G}_i, \mathcal{S}}[\log \int p_{\psi}(\mathbf{y}|\mathcal{G}_i, \mathcal{S}, \mathbf{z})p_{\phi}(\mathbf{z}|\mathcal{G}_i, \mathcal{S})d\mathbf{z}] \\ & = \arg \max \sum_i \mathbb{E}_{\mathbf{y}|\mathcal{G}_i, \mathcal{S}}[\log \int q_{\phi}(\mathbf{z}|\mathcal{G}_i, \mathcal{S})p_{\psi}(\mathbf{y}|\mathcal{G}_i, \mathcal{S}, \mathbf{z})\frac{p_{\phi}(\mathbf{z}|\mathcal{G}_i, \mathcal{S})}{q_{\phi}(\mathbf{z}|\mathcal{G}_i, \mathcal{S})}d\mathbf{z}], \end{aligned} \quad (4)$$

16 where  $p_{\psi}(\mathbf{y}|\mathcal{G}_i, \mathcal{S}, \mathbf{z})$  is the conditional probability of label  $\mathbf{y}$ .  $p_{\phi}(\mathbf{z}|\mathcal{G}_i, \mathcal{S})$  denotes the conditional  
17 probability of the factor  $\mathbf{z}$ , which is usually intractable. Hence, we resort to the variational inference to  
18 approximate the intractable true posterior with  $q_{\phi}(\mathbf{z}|\mathcal{G}_i, \mathcal{S})$  that is the expected distribution. According

19 to the Jensen Inequality, the above formulation can be deduced as follows:

$$\begin{aligned}
&\geq \arg \max \sum_i \mathbb{E}_{\mathbf{y}|\mathcal{G}_i, \mathcal{S}} \left[ \int q_\phi(\mathbf{z}|\mathcal{G}_i, \mathcal{S}) \log(p_\psi(\mathbf{y}|\mathcal{G}_i, \mathcal{S}, \mathbf{z}) \frac{p_\phi(\mathbf{z}|\mathcal{G}_i, \mathcal{S})}{q_\phi(\mathbf{z}|\mathcal{G}_i, \mathcal{S})}) d\mathbf{z} \right] \\
&= \arg \max \sum_i \mathbb{E}_{\mathbf{y}|\mathcal{G}_i, \mathcal{S}} \left[ \int q_\phi(\mathbf{z}|\mathcal{G}_i, \mathcal{S}) \log p_\psi(\mathbf{y}|\mathcal{G}_i, \mathcal{S}, \mathbf{z}) + q_\phi(\mathbf{z}|\mathcal{G}_i, \mathcal{S}) \log \frac{p_\phi(\mathbf{z}|\mathcal{G}_i, \mathcal{S})}{q_\phi(\mathbf{z}|\mathcal{G}_i, \mathcal{S})} d\mathbf{z} \right] \quad (5) \\
&= \arg \max \sum_i \mathbb{E}_{\mathbf{y}|\mathcal{G}_i, \mathcal{S}} \left[ \mathbb{E}_{q_\phi(\mathbf{z}|\mathcal{G}_i, \mathcal{S})} [\log p_\psi(\mathbf{y}|\mathcal{G}_i, \mathcal{S}, \mathbf{z})] - D_{KL}(q_\phi(\mathbf{z}|\mathcal{G}_i, \mathcal{S}) || p_\phi(\mathbf{z}|\mathcal{G}_i, \mathcal{S})) \right].
\end{aligned}$$

20 As  $q_\phi(\mathbf{z}|\mathcal{G}_i, \mathcal{S})$  is predefined distribution,  $\mathbb{E}_{q_\phi(\mathbf{z}|\mathcal{G}_i, \mathcal{S})}$  can be regarded as a constant, the objective can  
21 be formulated as:

$$\begin{aligned}
&\arg \max \mathbb{E}_{\mathbf{y}|\mathcal{G}_i, \mathcal{S}} [\log p_\psi(\mathbf{y}|\mathcal{G}_i, \mathcal{S}, \mathbf{z})] - D_{KL}(q_\phi(\mathbf{z}|\mathcal{G}_i, \mathcal{S}) || p_\phi(\mathbf{z}|\mathcal{G}_i, \mathcal{S})) \\
&= \arg \max -\zeta_{CE} - D_{KL}(q_\phi(\mathbf{z}|\mathcal{G}_i, \mathcal{S}) || p_\phi(\mathbf{z}|\mathcal{G}_i, \mathcal{S})) \quad (6) \\
&= \arg \min \zeta_{CE} + D_{KL}(q_\phi(\mathbf{z}|\mathcal{G}_i, \mathcal{S}) || p_\phi(\mathbf{z}|\mathcal{G}_i, \mathcal{S})).
\end{aligned}$$

22 In addition, we can extend the above single-layer BernPool into multi-layer networks by deploying  
23 independent sampling factors in sequential graph pooling.

24 **Cross-entropy Loss Function  $\zeta_{CE}$  based on Subgraph Sampling.** Referring to the analysis of a  
25 random dropping method [1], we analyze the loss function of our proposed BernPool. We can derive  
26 two parts from  $\zeta_{CE}$ :

$$\zeta_{CE} = \mathcal{L}_{CE} + \sum_i \frac{1}{2} y_i (1 - y_i) \text{Var}(\tilde{h}_i), \quad (7)$$

27 where  $\mathcal{L}_{CE}$  is the original cross-entropy loss function, the second term tends the classification  
28 probability to 0 or 1 and reduces the variance of  $h_i$  in the training process.

29 Specifically, for analytical simplicity, we apply a single-layer graph convolution as the backbone  
30 model to perform the binary classification task. As mentioned in Section 3 of this paper,  $\mathbf{H} =$   
31  $\sigma(\hat{\mathbf{D}}^{-\frac{1}{2}} \hat{\mathbf{A}} \hat{\mathbf{D}}^{\frac{1}{2}} \mathbf{X} \mathbf{W})$  and  $\mathbf{y} = \text{sigmoid}(\mathbf{H})$  represents predicted probability. Thus the original cross-  
32 entropy loss function can be expressed as follows:

$$\mathcal{L}_{CE} = \sum_{j, y_j=1} \log(1 + e^{-h_j}) + \sum_{k, y_k=0} \log(1 + e^{h_k}). \quad (8)$$

33 When performing sampling in the original graph, the objective function can be regarded as adding a  
34 bias, which is expressed as follows:

$$E(\zeta_{CE}) = \sum_{j, y_j=1} [\log(1 + e^{-h_j}) + \mathbb{E}(u(\tilde{h}_j, h_j))] + \sum_{k, y_k=0} [\log(1 + e^{h_k}) + \mathbb{E}(v(\tilde{h}_k, h_k))]. \quad (9)$$

35

$$\begin{cases} u(\tilde{h}_j, h_j) = \log(1 + e^{-\tilde{h}_j}) - \log(1 + e^{-h_j}). \\ v(\tilde{h}_k, h_k) = \log(1 + e^{-\tilde{h}_k}) - \log(1 + e^{-h_k}). \end{cases} \quad (10)$$

36 We can approximate it with second-order Taylor expansion of  $u(\cdot)$  and  $v(\cdot)$  around  $h_j$  and  $h_k$ ,  
37 respectively. For instance:

$$\begin{aligned}
u(\tilde{h}_j, h_j) &= \frac{-e^{-h_j}}{1 + e^{-h_j}} (\tilde{h}_j - h_j) + \frac{1}{2} \frac{e^{-h_j}}{(1 + e^{-h_j})^2} (\tilde{h}_j - h_j)^2 \\
&= (-1 + y_j) (\tilde{h}_j - h_j) + \frac{1}{2} y_j (1 - y_j) (\tilde{h}_j - h_j)^2.
\end{aligned} \quad (11)$$

38 In the same way,  $v(\tilde{h}_k, h_k) = y_k (\tilde{h}_k - h_k) + \frac{1}{2} y_k (1 - y_k) (\tilde{h}_k - h_k)^2$ . So the above equation can be  
39 transformed as:

$$\begin{aligned}
E(\zeta_{CE}) &= \mathcal{L}_{CE} + E\left( \sum_{j, y_j=1} [(-1 + z_j) (\tilde{h}_j - h_j) + \frac{1}{2} y_j (1 - y_j) (\tilde{h}_j - h_j)^2] \right) \\
&\quad + E\left( \sum_{k, y_k=1} [z_k (\tilde{h}_k - h_k) + \frac{1}{2} y_k (1 - y_k) (\tilde{h}_k - h_k)^2] \right) \quad (12) \\
&= \mathcal{L}_{CE} + \sum_i \frac{1}{2} y_i (1 - y_i) \text{Var}(\tilde{h}_i).
\end{aligned}$$

40 **References**

- 41 [1] Taoran Fang, Zhiqing Xiao, Chunping Wang, Jiarong Xu, Xuan Yang, and Yang Yang. Dropmes-  
42 sage: Unifying random dropping for graph neural networks, 2023.