
Algorithm 1: Training algorithm for explaining graph classification

- 1: **Input:** A set of input graphs (normalized $V \times V$ symmetric matrix observations) with i -th graph represented by $G_o^{(i)}$, node features $\mathbf{X}^{(i)}$ (if available), graph features $\mathbf{H}^{(i)}$, and a binary outcome $Y^{(i)}$, a trained GNN model: $\text{GNNE}_{\Phi_0}(\cdot)$ and $\text{GNNC}_{\Phi_1}(\cdot)$.
 - 2: **for** each graph $G_o^{(i)}$ **do**
 - 3: $\mathbf{Z}^{(i)} \leftarrow \text{GNNE}_{\Phi_0}(G_o^{(i)}, \mathbf{X}^{(i)})$.
 - 4: $Y_o^{(i)} \leftarrow \text{GNNC}_{\Phi_1}(\mathbf{Z}^{(i)})$.
 - 5: **end for**
 - 6: **for** each epoch **do**
 - 7: **for** each graph $G_o^{(i)}$ **do**
 - 8: $\Omega \leftarrow$ latent variables parameterized by neural networks.
 - 9: **for** $k \leftarrow 1$ **to** K **do**
 - 10: $\hat{G}_s^{(i,k)} \leftarrow$ sampled from binary edge variables (preserving weights if needed).
 - 11: $\hat{Y}_s^{(i,k)} \leftarrow \text{GNNC}_{\Phi_1}(\text{GNNE}_{\Phi_0}(\hat{G}_s^{(i,k)}, \mathbf{X}^{(i)}))$
 - 12: **end for**
 - 13: **end for**
 - 14: Compute loss with the objective function.
 - 15: Update parameters Ψ with backpropagation.
 - 16: **end for**
-