

A Supplementary Materials

A.1 Proof of Theorem 1

Proof of Theorem 1.

$$\min_{\hat{P}(\hat{\mathbf{y}}|\mathbf{x})} \max_{\check{P}(\check{\mathbf{y}}|\mathbf{x})} \mathbb{E}_{\mathbf{X} \sim \tilde{P}; \hat{\mathbf{Y}}|\mathbf{X} \sim \hat{P}; \check{\mathbf{Y}}|\mathbf{X} \sim \check{P}} [\text{loss}(\hat{\mathbf{Y}}, \check{\mathbf{Y}})] \quad (11)$$

$$\begin{aligned} & \text{subject to: } \mathbb{E}_{\mathbf{X} \sim \tilde{P}; \check{\mathbf{Y}}|\mathbf{X} \sim \check{P}} [\Phi(\mathbf{X}, \check{\mathbf{Y}})] = \mathbb{E}_{\mathbf{X}, \mathbf{Y} \sim \tilde{P}} [\Phi(\mathbf{X}, \mathbf{Y})] \\ & \stackrel{(a)}{=} \max_{\check{P}(\check{\mathbf{y}}|\mathbf{x})} \min_{\hat{P}(\hat{\mathbf{y}}|\mathbf{x})} \mathbb{E}_{\mathbf{X} \sim \tilde{P}; \hat{\mathbf{Y}}|\mathbf{X} \sim \hat{P}; \check{\mathbf{Y}}|\mathbf{X} \sim \check{P}} [\text{loss}(\hat{\mathbf{Y}}, \check{\mathbf{Y}})] \end{aligned} \quad (12)$$

$$\begin{aligned} & \text{subject to: } \mathbb{E}_{\mathbf{X} \sim \tilde{P}; \check{\mathbf{Y}}|\mathbf{X} \sim \check{P}} [\Phi(\mathbf{X}, \check{\mathbf{Y}})] = \mathbb{E}_{\mathbf{X}, \mathbf{Y} \sim \tilde{P}} [\Phi(\mathbf{X}, \mathbf{Y})] \\ & \stackrel{(b)}{=} \max_{\check{P}(\check{\mathbf{y}}|\mathbf{x})} \min_{\theta} \min_{\hat{P}(\hat{\mathbf{y}}|\mathbf{x})} \mathbb{E}_{\mathbf{X}, \mathbf{Y} \sim \tilde{P}; \hat{\mathbf{Y}}|\mathbf{X} \sim \hat{P}; \check{\mathbf{Y}}|\mathbf{X} \sim \check{P}} [\text{loss}(\hat{\mathbf{Y}}, \check{\mathbf{Y}}) + \theta^T (\Phi(\mathbf{X}, \check{\mathbf{Y}}) - \Phi(\mathbf{X}, \mathbf{Y}))] \end{aligned} \quad (13)$$

$$\stackrel{(c)}{=} \min_{\theta} \max_{\check{P}(\check{\mathbf{y}}|\mathbf{x})} \min_{\hat{P}(\hat{\mathbf{y}}|\mathbf{x})} \mathbb{E}_{\mathbf{X}, \mathbf{Y} \sim \tilde{P}; \hat{\mathbf{Y}}|\mathbf{X} \sim \hat{P}; \check{\mathbf{Y}}|\mathbf{X} \sim \check{P}} [\text{loss}(\hat{\mathbf{Y}}, \check{\mathbf{Y}}) + \theta^T (\Phi(\mathbf{X}, \check{\mathbf{Y}}) - \Phi(\mathbf{X}, \mathbf{Y}))] \quad (14)$$

$$\stackrel{(d)}{=} \min_{\theta} \mathbb{E}_{\mathbf{X}, \mathbf{Y} \sim \tilde{P}} \max_{\check{P}(\check{\mathbf{y}}|\mathbf{x})} \min_{\hat{P}(\hat{\mathbf{y}}|\mathbf{x})} \mathbb{E}_{\hat{\mathbf{Y}}|\mathbf{X} \sim \hat{P}; \check{\mathbf{Y}}|\mathbf{X} \sim \check{P}} [\text{loss}(\hat{\mathbf{Y}}, \check{\mathbf{Y}}) + \theta^T (\Phi(\mathbf{X}, \check{\mathbf{Y}}) - \Phi(\mathbf{X}, \mathbf{Y}))] \quad (15)$$

$$\begin{aligned} & \stackrel{(e)}{=} \min_{\theta_e, \theta_v} \mathbb{E}_{\mathbf{X}, \mathbf{Y} \sim \tilde{P}} \max_{\check{P}(\check{\mathbf{y}}|\mathbf{x})} \min_{\hat{P}(\hat{\mathbf{y}}|\mathbf{x})} \left[\sum_i^n \text{loss}(\hat{Y}_i, \check{Y}_i) \right. \\ & \quad \left. + \theta_e \cdot \sum_{(i,j) \in E} [\phi(\mathbf{X}, \check{Y}_i, \check{Y}_j) - \phi(\mathbf{X}, Y_i, Y_j)] + \theta_v \cdot \sum_i^n [\phi(\mathbf{X}, \check{Y}_i) - \phi(\mathbf{X}, Y_i)] \right] \end{aligned} \quad (16)$$

$$\begin{aligned} & \stackrel{(f)}{=} \min_{\theta_e, \theta_v} \mathbb{E}_{\mathbf{X}, \mathbf{Y} \sim \tilde{P}} \max_{\check{P}(\check{\mathbf{y}}|\mathbf{x})} \min_{\hat{P}(\hat{\mathbf{y}}|\mathbf{x})} \sum_{\check{\mathbf{y}}, \check{\mathbf{y}}} \hat{P}(\hat{\mathbf{y}}|\mathbf{x}) \check{P}(\check{\mathbf{y}}|\mathbf{x}) \left[\sum_i^n \text{loss}(\hat{y}_i, \check{y}_i) \right. \\ & \quad \left. + \theta_e \cdot \sum_{(i,j) \in E} [\phi(\mathbf{x}, \check{y}_i, \check{y}_j) - \phi(\mathbf{x}, y_i, y_j)] + \theta_v \cdot \sum_i^n [\phi(\mathbf{x}, \check{y}_i) - \phi(\mathbf{x}, y_i)] \right] \end{aligned} \quad (17)$$

$$\begin{aligned} & \stackrel{(g)}{=} \min_{\theta_e, \theta_v} \mathbb{E}_{\mathbf{X}, \mathbf{Y} \sim \tilde{P}} \max_{\check{P}(\check{\mathbf{y}}|\mathbf{x})} \min_{\hat{P}(\hat{\mathbf{y}}|\mathbf{x})} \left[\sum_i^n \sum_{\hat{y}_i, \check{y}_i} \hat{P}(\hat{y}_i|\mathbf{x}) \check{P}(\check{y}_i|\mathbf{x}) \text{loss}(\hat{y}_i, \check{y}_i) \right. \\ & \quad \left. + \sum_{(i,j) \in E} \sum_{\hat{y}_i, \check{y}_j} \check{P}(\check{y}_i, \check{y}_j|\mathbf{x}) [\theta_e \cdot \phi(\mathbf{x}, \check{y}_i, \check{y}_j)] - \sum_{(i,j) \in E} \theta_e \cdot \phi(\mathbf{x}, y_i, y_j) \right. \\ & \quad \left. + \sum_i^n \sum_{\check{y}_i} \check{P}(\check{y}_i|\mathbf{x}) [\theta_v \cdot \phi(\mathbf{x}, \check{y}_i)] - \sum_i^n \theta_v \cdot \phi(\mathbf{x}, y_i) \right]. \end{aligned} \quad (18)$$

The transformation steps above are described as follows:

- (a) We flip the min and max order using minimax duality [36]. The domains of $\hat{P}(\hat{\mathbf{y}}|\mathbf{x})$ and $\check{P}(\check{\mathbf{y}}|\mathbf{x})$ are both compact convex sets and the objective function is bilinear, therefore, strong duality holds.
- (b) We introduce the Lagrange dual variable θ to directly incorporate the equality constraints into the objective function.
- (c) The domain of $\check{P}(\check{\mathbf{y}}|\mathbf{x})$ is a compact convex subset of \mathbb{R}^n , while the domain of θ is \mathbb{R}^m . The objective is concave on $\check{P}(\check{\mathbf{y}}|\mathbf{x})$ for all θ (a non-negative linear combination of minimums of affine functions is concave), while it is convex on θ for all $\check{P}(\check{\mathbf{y}}|\mathbf{x})$. Based on Sion's minimax theorem [37], strong duality holds, and thus we can flip the optimization order of $\check{P}(\check{\mathbf{y}}|\mathbf{x})$ and θ .
- (d) Since the expression is additive in terms of $\check{P}(\check{\mathbf{y}}|\mathbf{x})$ and $\hat{P}(\hat{\mathbf{y}}|\mathbf{x})$, we can push the expectation over the empirical distribution $\mathbf{X}, \mathbf{Y} \sim \tilde{P}$ outside and independently optimize each $\check{P}(\check{\mathbf{y}}|\mathbf{x})$ and $\hat{P}(\hat{\mathbf{y}}|\mathbf{x})$.
- (e) We apply our description of loss metrics which is additively decomposable into the loss for each node, and the features that can be decomposed into node and edge features. We also

separate the notation for the Lagrange dual variable into the variable for the constraints on node features (θ_v) and and the variable for the edge features (θ_e).

- (f) We rewrite the expectation over $\hat{P}(\hat{\mathbf{y}}|\mathbf{x})$ and $\check{P}(\check{\mathbf{y}}|\mathbf{x})$ in terms of the probability-weighted average.
- (g) Based on the property of the loss metrics and feature functions, the sum over the exponentially many possibilities of $\hat{\mathbf{y}}$ and $\check{\mathbf{y}}$ can be simplified into the sum over individual nodes and edges values, resulting in the optimization over the node and edge marginal distributions.

□

A.2 Proof of Theorem 2

Proof of Theorem 2.

$$\max_{\mathbf{Q} \in \Delta} \min_{\mathbf{p} \in \Delta} \sum_i^n \left[\mathbf{p}_i \mathbf{L}_i(\mathbf{Q}_{pt(i);i}^T \mathbf{1}) + \langle \mathbf{Q}_{pt(i);i}, \mathbf{B}_{pt(i);i} \rangle + (\mathbf{Q}_{pt(i);i}^T \mathbf{1})^T \mathbf{b}_i \right] \quad (19)$$

subject to: $\mathbf{Q}_{pt(pt(i));pt(i)}^T \mathbf{1} = \mathbf{Q}_{pt(i);i} \mathbf{1}, \forall i \in \{1, \dots, n\}$

$$\stackrel{(a)}{=} \max_{\mathbf{Q} \in \Delta} \min_{\mathbf{u}} \min_{\mathbf{p} \in \Delta} \sum_i^n \left[\mathbf{p}_i \mathbf{L}_i(\mathbf{Q}_{pt(i);i}^T \mathbf{1}) + \langle \mathbf{Q}_{pt(i);i}, \mathbf{B}_{pt(i);i} \rangle + (\mathbf{Q}_{pt(i);i}^T \mathbf{1})^T \mathbf{b}_i \right] \\ + \sum_i^n \mathbf{u}_i^T (\mathbf{Q}_{pt(pt(i));pt(i)}^T \mathbf{1} - \mathbf{Q}_{pt(i);i} \mathbf{1}) \quad (20)$$

$$\stackrel{(b)}{=} \min_{\mathbf{u}} \max_{\mathbf{Q} \in \Delta} \min_{\mathbf{p} \in \Delta} \sum_i^n \left[\mathbf{p}_i \mathbf{L}_i(\mathbf{Q}_{pt(i);i}^T \mathbf{1}) + \langle \mathbf{Q}_{pt(i);i}, \mathbf{B}_{pt(i);i} \rangle + (\mathbf{Q}_{pt(i);i}^T \mathbf{1})^T \mathbf{b}_i \right] \\ + \sum_i^n \mathbf{u}_i^T (\mathbf{Q}_{pt(pt(i));pt(i)}^T \mathbf{1} - \mathbf{Q}_{pt(i);i} \mathbf{1}) \quad (21)$$

$$\stackrel{(c)}{=} \min_{\mathbf{u}} \max_{\mathbf{Q} \in \Delta} \min_{\mathbf{p} \in \Delta} \sum_i^n \left[\mathbf{p}_i \mathbf{L}_i(\mathbf{Q}_{pt(i);i}^T \mathbf{1}) + \langle \mathbf{Q}_{pt(i);i}, \mathbf{B}_{pt(i);i} \rangle + \langle \mathbf{Q}_{pt(i);i}, \mathbf{1} \mathbf{b}_i^T \rangle \right] \\ + \sum_i^n \left[\langle \mathbf{Q}_{pt(pt(i));pt(i)}, \mathbf{1} \mathbf{u}_i^T \rangle - \langle \mathbf{Q}_{pt(i);i}, \mathbf{u}_i \mathbf{1}^T \rangle \right] \quad (22)$$

$$\stackrel{(d)}{=} \min_{\mathbf{u}} \max_{\mathbf{Q} \in \Delta} \min_{\mathbf{p} \in \Delta} \sum_i^n \left[\mathbf{p}_i \mathbf{L}_i(\mathbf{Q}_{pt(i);i}^T \mathbf{1}) + \langle \mathbf{Q}_{pt(i);i}, \mathbf{B}_{pt(i);i} + \mathbf{1} \mathbf{b}_i^T - \mathbf{u}_i \mathbf{1}^T + \sum_{k \in ch(i)} \mathbf{1} \mathbf{u}_k^T \rangle \right]. \quad (23)$$

The transformation steps above are described as follows:

- (a) We introduce the Lagrange dual variable \mathbf{u} , where \mathbf{u}_i is the dual variable associated with the marginal constraint of $\mathbf{Q}_{pt(pt(i));pt(i)}^T \mathbf{1} = \mathbf{Q}_{pt(i);i} \mathbf{1}$.
- (b) Similar to the analysis in Theorem 1, strong duality holds due to Sion's minimax theorem. Therefore we can flip the optimization order of \mathbf{Q} and \mathbf{u} .
- (c) We rewrite the vector multiplication over $\mathbf{Q}_{pt(i);i} \mathbf{1}$ or $\mathbf{Q}_{pt(i);i}^T \mathbf{1}$ with the corresponding Frobenius inner product notations.
- (d) We regroup the terms in the optimization above by considering the parent-child relations in the tree for each node. Note that $ch(i)$ represents the children of node i .

□