

Exploiting Human-AI Dependency for Learning to Defer

Anonymous Authors¹

In our paper, we proposed a surrogate loss called DCE for the vanilla Learning to Defer (L2D) setting. Here, we further provide an extension for the proposed DCE loss to the setting of L2D with deferral costs and L2D with multiple experts.

1. L2D with Deferral Costs

In this part, we would like to introduce the extension of DCE with deferral cost (DCEc). In L2D with deferral cost $c(\mathbf{x})$. The performance of an L2D system could be formulated as follow loss function:

$$L_{01c}^\perp(f(\mathbf{x}), y, m) = \mathbb{I}_{f(\mathbf{x}) \neq y} \mathbb{I}_{f(\mathbf{x}) \neq \perp} + (\mathbb{I}_{m \neq y} + c(\mathbf{x})) \mathbb{I}_{f(\mathbf{x}) = \perp},$$

Previous studies showed that the Bayes optimal classifier for L2D with deferral cost $c(\mathbf{x})$ can be formulated as follows:

$$f^*(\mathbf{x}) = \begin{cases} \perp, & \mathbb{P}(Y = M|\mathbf{x}) > c(\mathbf{x}) + \max_{y \in \mathcal{Y}} \eta_y(\mathbf{x}), \\ \arg \max_{y \in \mathcal{Y}} \eta_y(\mathbf{x}), & \text{otherwise.} \end{cases} \quad (1)$$

Then, the generalized version of DCE for L2D with deferral costs could be defined as:

$$\begin{aligned} L_{DCEc}^\perp(g(\mathbf{x}), y, m) = & -\mathbb{I}_{m \neq y} \log(\psi_{\mathcal{Y}^\perp}^y(g(\mathbf{x}))) \\ & - \mathbb{I}_{m=y} (\log(\psi_{\mathcal{Y}}^y(g(\mathbf{x}))) + \log(\psi_{\mathcal{Y}^\perp/q}^\perp(g(\mathbf{x})))) \\ & - c(\mathbf{x}) \log(\psi_{\mathcal{Y}^\perp}^q(g(\mathbf{x}))), \end{aligned}$$

Where $q = \arg \max_{y \in \mathcal{Y}} g_y(\mathbf{x})$ denotes the predicted class of the model. Then we could train a model by minimize the risk $R_{DCEc}^\perp(g) = \mathbb{E}_{p(\mathbf{x}, y, m)} [L_{01c}^\perp(g(\mathbf{x}), y, m)]$.

Theorem 1.1. *Let us define $g^*(\mathbf{x}) \in \arg \min_g R_{DCEm}^\perp(g)$ the optimized scoring function, then $\phi(g^*(\mathbf{x}))$ meets the condition in Eq. 1, which further means L_{DCEc}^\perp is a consistent loss for L2D with deferral costs.*

The proof of Theorem 1.1 is provided in Appendix A.

¹Anonymous Institution, Anonymous City, Anonymous Region, Anonymous Country. Correspondence to: Anonymous Author <anon.email@domain.com>.

2. L2D with Multiple Experts

In this part, we show how to extend the DCE loss to L2D with multiple experts. We call this extended version DCEm. L2D with multiple experts could be formulated as minimizing the following loss function

$$\begin{aligned} L_{01m}^\perp(f(\mathbf{x}), y, m_1, m_2, \dots, m_J) = & \mathbb{I}_{f(\mathbf{x}) \neq y} \mathbb{I}_{f(\mathbf{x}) \in \mathcal{Y}} \\ & + \sum_{j=1}^J (\mathbb{I}_{m_j \neq y}) \mathbb{I}_{f(\mathbf{x}) = \perp_j}, \end{aligned}$$

Where J is used to denote the number of experts, and m_j denotes the prediction made by the j -th expert. When $f(\mathbf{x}) = \perp_j$, the L2D system would defer the prediction to the j -th expert. Previous studies have shown that the Bayes optimal classifier for L2D with multiple experts should meet the following condition:

$$f^*(\mathbf{x}) = \begin{cases} \perp_t, & \mathbb{P}(Y = M_t|\mathbf{x}) > \max_{y \in \mathcal{Y}} \eta_y(\mathbf{x}), \\ \arg \max_{y \in \mathcal{Y}} \eta_y(\mathbf{x}), & \text{otherwise,} \end{cases}$$

where $t = \arg \max_{0 \leq j \leq J} \mathbb{P}(Y = M_j|\mathbf{x})$.

Then the extended DCE for L2D with multiple experts could be formulated as:

$$\begin{aligned} L_{DCEm}^\perp(g(\mathbf{x}), y, m_1, m_2, \dots, m_J) = & -\log(\psi_{\mathcal{Y}}^y(g(\mathbf{x}))) - \log(\psi_{\mathcal{Y} \cup \{j | \mathbb{I}[m_j \neq y]\}}^\perp(g(\mathbf{x}))) \\ & - \sum_{j=1}^J \mathbb{I}[m_j = y] \log(\psi_{\mathcal{Y}^\perp/q}^{\perp_j}(g(\mathbf{x}))), \end{aligned}$$

Where $\mathcal{Y}^{\perp_j} = \mathcal{Y} \cup \{\perp_j\}$. However, L_{DCEm} is not a consistent loss for L2D with multiple experts. Since DCE considers the dependency patterns between different M_j and Y during training, various M_j are connected through Y . These diverse M_j mutually influence each other during the training process, resulting in increased complexity. How to utilize dependency patterns in scenarios with multiple experts poses an interesting challenge.

A. Proof of Theorem 1.1

To begin with proof, let us define

$$C_{\text{DCEc}}^{\perp}(g(\mathbf{x})) = \sum_{y \in \mathcal{Y}} \left[-\mathbb{P}(Y = y, M = y | \mathbf{x}) (\log(\psi_y^y(g(\mathbf{x}))) + \log(\psi_{\mathcal{Y}^{\perp}/q}^{\perp}(g(\mathbf{x})))) \right. \\ \left. - (\mathbb{P}(Y = y, M \neq y | \mathbf{x}) + c(\mathbf{x})) \log(\psi_y^y(g(\mathbf{x}))) \right]$$

The conditional surrogate risk w.r.t. \mathbf{x} .

Let g^* be the optimal scoring function for DCEc.

We first proof that $\arg \max_{\tilde{y} \in \mathcal{Y}} g_{\tilde{y}}^*(\mathbf{x}) = \arg \max_{\tilde{y} \in \mathcal{Y}} \eta_{\tilde{y}}(\mathbf{x})$ by contradiction. Let us use $r = \arg \max_{\tilde{y} \in \mathcal{Y}} \eta_{\tilde{y}}(\mathbf{x})$, $q = \arg \max_{\tilde{y} \in \mathcal{Y}} g_{\tilde{y}}^*(\mathbf{x})$ to denote the index of maximum dimension in posterior distribution and scoring function respectively. For simplicity in notation, we represent the score vector outputted by the scoring function as $\mathbf{s} = g(\mathbf{x})$, and $\mathbf{s}^* = g^*(\mathbf{x})$

Suppose $r \neq q$, i.e. $\mathbf{s}_q^* > \mathbf{s}_r^*$. Then we show that we could obtain a lower value of conditional surrogate risk by switching the value between \mathbf{s}_q^* and \mathbf{s}_r^* . Let us define $\tilde{\mathbf{s}}$ the score vector obtained by switching the value between \mathbf{s}_q^* and \mathbf{s}_r^* , i.e. $\tilde{\mathbf{s}}_r = \mathbf{s}_q^*$, $\tilde{\mathbf{s}}_q = \mathbf{s}_r^*$ and $\tilde{\mathbf{s}}_i = \mathbf{s}_i^*$ for $i \in \mathcal{Y}^{\perp}$, $i \neq r, q$. Thus $\arg \max_{\tilde{y} \in \mathcal{Y}} \tilde{\mathbf{s}} = r$ Then $C_{\text{DCEc}}^{\perp}(\mathbf{s}^*) - C_{\text{DCEc}}^{\perp}(\tilde{\mathbf{s}})$ could be expressed as:

$$\begin{aligned} & C_{\text{DCEc}}^{\perp}(\mathbf{s}^*) - C_{\text{DCEc}}^{\perp}(\tilde{\mathbf{s}}) \\ &= \mathbb{P}(Y = q, M = q | \mathbf{x}) (\log(\psi_y^q(\tilde{\mathbf{s}})) - \log(\psi_y^q(\mathbf{s}^*)) + \log(\psi_{\mathcal{Y}^{\perp}/r}^{\perp}(\tilde{\mathbf{s}})) - \log(\psi_{\mathcal{Y}^{\perp}/q}^{\perp}(\mathbf{s}^*))) \\ &+ \mathbb{P}(Y = q, M \neq q | \mathbf{x}) (\log(\psi_{\mathcal{Y}^{\perp}}^q(\tilde{\mathbf{s}})) - \log(\psi_{\mathcal{Y}^{\perp}}^q(\mathbf{s}^*))) \\ &+ \mathbb{P}(Y = r, M = r | \mathbf{x}) (\log(\psi_y^r(\tilde{\mathbf{s}})) - \log(\psi_y^r(\mathbf{s}^*)) + \log(\psi_{\mathcal{Y}^{\perp}/r}^{\perp}(\tilde{\mathbf{s}})) - \log(\psi_{\mathcal{Y}^{\perp}/q}^{\perp}(\mathbf{s}^*))) \\ &+ \mathbb{P}(Y = r, M \neq r | \mathbf{x}) (\log(\psi_{\mathcal{Y}^{\perp}}^r(\tilde{\mathbf{s}})) - \log(\psi_{\mathcal{Y}^{\perp}}^r(\mathbf{s}^*))) \\ &\stackrel{(a)}{=} \mathbb{P}(Y = q, M = q | \mathbf{x}) (\log(\frac{\exp(\mathbf{s}_r^*)}{\exp(\mathbf{s}_q^*)}) + \log(\frac{\exp(\mathbf{s}_r^*) + \sum_{i \neq r, q} \exp(\mathbf{s}_i^*)}{\exp(\mathbf{s}_r^*) + \sum_{i \neq r, q} \exp(\mathbf{s}_i^*)}) \\ &+ \mathbb{P}(Y = q, M \neq q | \mathbf{x}) (\log(\frac{\exp(\mathbf{s}_r^*)}{\exp(\mathbf{s}_q^*)})) \\ &+ \mathbb{P}(Y = r, M = r | \mathbf{x}) (\log(\frac{\exp(\mathbf{s}_q^*)}{\exp(\mathbf{s}_r^*)}) + \log(\frac{\exp(\mathbf{s}_r^*) + \sum_{i \neq r, q} \exp(\mathbf{s}_i^*)}{\exp(\mathbf{s}_r^*) + \sum_{i \neq r, q} \exp(\mathbf{s}_i^*)}) \\ &+ \mathbb{P}(Y = r, M \neq r | \mathbf{x}) (\log(\frac{\exp(\mathbf{s}_q^*)}{\exp(\mathbf{s}_r^*)})) \\ &= (\mathbb{P}(Y = r | \mathbf{x}) - \mathbb{P}(Y = q | \mathbf{x})) \log(\frac{\exp(\mathbf{s}_q^*)}{\exp(\mathbf{s}_r^*)}) \\ &> 0 \end{aligned}$$

The equation (a) holds since $\sum_{i \in \mathcal{Y}} \exp(\mathbf{s}_i^*) = \sum_{i \in \mathcal{Y}} \exp(\tilde{\mathbf{s}}_i)$ and $\sum_{i \in \mathcal{Y}^{\perp}} \exp(\mathbf{s}_i^*) = \sum_{i \in \mathcal{Y}^{\perp}} \exp(\tilde{\mathbf{s}}_i)$. Because $\mathbb{P}(Y = r | \mathbf{x}) > \mathbb{P}(Y = q | \mathbf{x})$ and $\mathbf{s}_q^* > \mathbf{s}_r^*$, then $(\mathbb{P}(Y = r | \mathbf{x}) - \mathbb{P}(Y = q | \mathbf{x})) \log(\frac{\exp(\mathbf{s}_q^*)}{\exp(\mathbf{s}_r^*)}) > 0, C_{\text{DCE}}^{\perp}(\mathbf{s}^*) > C_{\text{DCE}}^{\perp}(\tilde{\mathbf{s}})$. This contradicts to \mathbf{s}^* is the minimizer of C_{DCE}^{\perp} , thus $r = q$. Then we can conclude the proof of $r = q$.

Then we prove that $\mathbf{s}_r^* > \mathbf{s}_{\perp}^*$ when $\mathbb{P}(Y = r, Y \neq M | \mathbf{x}) + c(\mathbf{x}) > \mathbb{P}(Y \neq r, Y = M | \mathbf{x})$. Suppose $\mathbf{s}_{\perp}^* > \mathbf{s}_r^*$ when $\mathbb{P}(Y = r, M \neq r | \mathbf{x}) + c(\mathbf{x}) > \mathbb{P}(Y \neq r, Y = M | \mathbf{x})$ and we use \mathbf{s}' to represent the score vector obtained by switching the value between \mathbf{s}_r^* and \mathbf{s}_{\perp}^* , i.e. $\mathbf{s}'_r = \mathbf{s}_{\perp}^*$, $\mathbf{s}'_{\perp} = \mathbf{s}_r^*$ and $\mathbf{s}'_i = \mathbf{s}_i^*$ for all $i \in \mathcal{Y}$, $i \neq r$. We can directly derive that

$r = \arg \max_{\tilde{y} \in \mathcal{Y}} s'_r$. Then $C_{\text{DCE}}^\perp(s^*) - C_{\text{DCE}}^\perp(s')$ could be expressed as:

$$\begin{aligned}
 & C_{\text{DCE}}^\perp(s^*) - C_{\text{DCE}}^\perp(s') \\
 &= \mathbb{P}(Y = r, M = r | \mathbf{x}) (\log(\psi_Y^r(s')) - \log(\psi_Y^r(s^*)) + \log(\psi_{Y^\perp/r}^\perp(s')) - \log(\psi_{Y^\perp/r}^\perp(s^*))) \\
 &+ \mathbb{P}(Y = r, M \neq r | \mathbf{x}) (\log(\psi_{Y^\perp}^r(s')) - \log(\psi_{Y^\perp}^r(s^*))) \\
 &+ \sum_{i \in \mathcal{Y}, i \neq r} \mathbb{P}(Y = i, M = i | \mathbf{x}) (\log(\psi_Y^i(s')) - \log(\psi_Y^i(s^*)) + \log(\psi_{Y^\perp/r}^\perp(s')) - \log(\psi_{Y^\perp/r}^\perp(s^*))) \\
 &+ \sum_{i \in \mathcal{Y}, i \neq r} \mathbb{P}(Y = i, M \neq i | \mathbf{x}) (\log(\psi_{Y^\perp}^i(s')) - \log(\psi_{Y^\perp}^i(s^*))) \\
 &+ c(\mathbf{x}) (\log(\psi_{Y^\perp}^q(s')) - \log(\psi_{Y^\perp}^q(s^*))) \\
 &= \mathbb{P}(Y = r, M \neq r | \mathbf{x}) (\log(\psi_{Y^\perp}^r(s')) - \log(\psi_{Y^\perp}^r(s^*))) \\
 &+ \sum_{i \in \mathcal{Y}, i \neq r} \mathbb{P}(Y = i, M = i | \mathbf{x}) (\log(\psi_Y^i(s')) - \log(\psi_Y^i(s^*)) + \log(\psi_{Y^\perp/r}^\perp(s')) - \log(\psi_{Y^\perp/r}^\perp(s^*))) \\
 &+ c(\mathbf{x}) \log\left(\frac{\exp(s_\perp^*)}{\exp(s_r^*)}\right) \\
 &= \mathbb{P}(Y = r, M \neq r | \mathbf{x}) \log\left(\frac{\exp(s_\perp^*)}{\exp(s_r^*)}\right) \\
 &+ \sum_{i \in \mathcal{Y}, i \neq r} \mathbb{P}(Y = i, M = i | \mathbf{x}) \left(\log\left(\frac{\sum_{j \in \mathcal{Y}, j \neq r} \exp(s_j^*) + \exp(s_r^*)}{\sum_{j \in \mathcal{Y}, j \neq r} \exp(s_j^*) + \exp(s_\perp^*)}\right) + \log\left(\frac{\exp(s_r^*)}{\exp(s_\perp^*)} \frac{\sum_{j \in \mathcal{Y}, j \neq r} \exp(s_j^*) + \exp(s_\perp^*)}{\sum_{j \in \mathcal{Y}, j \neq r} \exp(s_j^*) + \exp(s_r^*)}\right) \right) \\
 &+ c(\mathbf{x}) \log\left(\frac{\exp(s_\perp^*)}{\exp(s_r^*)}\right) \\
 &= \mathbb{P}(Y = r, M \neq r | \mathbf{x}) \log\left(\frac{\exp(s_\perp^*)}{\exp(s_r^*)}\right) \\
 &- \mathbb{P}(Y \neq r, M = Y | \mathbf{x}) \log\left(\frac{\exp(s_\perp^*)}{\exp(s_r^*)}\right) \\
 &+ c(\mathbf{x}) \log\left(\frac{\exp(s_\perp^*)}{\exp(s_r^*)}\right) \\
 &= (\mathbb{P}(Y = r, M \neq r | \mathbf{x}) + c(\mathbf{x}) - \mathbb{P}(Y \neq r, M = Y | \mathbf{x})) \log\left(\frac{\exp(s_\perp^*)}{\exp(s_r^*)}\right) > 0
 \end{aligned}$$

Thus $C_{\text{DCE}}^\perp(s^*) > C_{\text{DCE}}^\perp(s')$, which is contradictory to s^* is a minimizer of C_{DCE}^\perp , then we conclude the proof that $s_r^* > s_\perp^*$ if $\mathbb{P}(Y = r, Y \neq M | \mathbf{x}) + c(\mathbf{x}) > \mathbb{P}(Y \neq r, Y = M | \mathbf{x})$.

Lastly, we prove that $s_\perp^* > s_r^*$ when $\mathbb{P}(Y \neq r, Y = M | \mathbf{x}) > \mathbb{P}(Y = r, Y \neq M | \mathbf{x}) + c(\mathbf{x})$. Suppose $s_r^* > s_\perp^*$ when $\mathbb{P}(Y \neq r, Y = M | \mathbf{x}) > \mathbb{P}(Y = r, Y \neq M | \mathbf{x}) + c(\mathbf{x})$ and we still use s' to represent the score vector obtained by switching the value between s_r^* and s_\perp^* , i.e. $s'_r = s_\perp^*$, $s'_\perp = s_r^*$ and $s'_i = s_i^*$ for all $i \in \mathcal{Y}, i \neq r$. Let $t = \arg \max_{\tilde{y} \in \mathcal{Y}} s'_t$. We define

$\epsilon = \sum_{i \in \mathcal{Y}, i \neq r} \exp(s_i^*)$ and $\epsilon' = \sum_{i \in \mathcal{Y}, i \neq r} \exp(s'_i)$. We could derive that $\epsilon \geq \epsilon'$. Then $C_{\text{DCE}}^\perp(s^*) - C_{\text{DCE}}^\perp(s')$ could be expressed as:

$$\begin{aligned}
 & C_{\text{DCEc}}^{\perp}(\mathbf{s}^*) - C_{\text{DCEc}}^{\perp}(\mathbf{s}') \\
 &= \mathbb{P}(Y = r, M = r|\mathbf{x}) (\log(\psi_{\mathbf{y}}^r(\mathbf{s}')) - \log(\psi_{\mathbf{y}}^r(\mathbf{s}^*)) + \log(\psi_{\mathbf{y}^{\perp}/t}^{\perp}(\mathbf{s}')) - \log(\psi_{\mathbf{y}^{\perp}/r}^{\perp}(\mathbf{s}^*))) \\
 &+ \mathbb{P}(Y = r, M \neq r|\mathbf{x}) (\log(\psi_{\mathbf{y}^{\perp}}^r(\mathbf{s}')) - \log(\psi_{\mathbf{y}^{\perp}}^r(\mathbf{s}^*))) \\
 &+ \sum_{i \in \mathcal{Y}, i \neq r} \mathbb{P}(Y = i, M = i|\mathbf{x}) (\log(\psi_{\mathbf{y}}^i(\mathbf{s}')) - \log(\psi_{\mathbf{y}}^i(\mathbf{s}^*)) + \log(\psi_{\mathbf{y}^{\perp}/t}^{\perp}(\mathbf{s}')) - \log(\psi_{\mathbf{y}^{\perp}/r}^{\perp}(\mathbf{s}^*))) \\
 &+ \sum_{i \in \mathcal{Y}, i \neq r} \mathbb{P}(Y = i, M \neq i|\mathbf{x}) (\log(\psi_{\mathbf{y}^{\perp}}^i(\mathbf{s}')) - \log(\psi_{\mathbf{y}^{\perp}}^i(\mathbf{s}^*))) \\
 &+ c(\mathbf{x}) (\log(\psi_{\mathbf{y}^{\perp}}^q(\mathbf{s}')) - \log(\psi_{\mathbf{y}^{\perp}}^q(\mathbf{s}^*))) \\
 &= \mathbb{P}(Y = r, M = r|\mathbf{x}) (\log(\frac{\exp(\mathbf{s}_{\perp}^*)}{\epsilon + \exp(\mathbf{s}_{\perp}^*)}) - \log(\frac{\exp(\mathbf{s}_r^*)}{\epsilon + \exp(\mathbf{s}_r^*)}) + \log(\frac{\exp(\mathbf{s}_r^*)}{\epsilon' + \exp(\mathbf{s}_r^*)}) - \log(\frac{\exp(\mathbf{s}_{\perp}^*)}{\epsilon + \exp(\mathbf{s}_{\perp}^*)})) \\
 &+ \mathbb{P}(Y = r, M \neq r|\mathbf{x}) \log(\frac{\exp(\mathbf{s}_{\perp}^*)}{\exp(\mathbf{s}_r^*)}) \\
 &+ \sum_{i \in \mathcal{Y}, i \neq r} \mathbb{P}(Y = i, M = i|\mathbf{x}) (\log(\frac{\exp(\mathbf{s}_i^*)}{\epsilon + \exp(\mathbf{s}_{\perp}^*)}) - \log(\frac{\exp(\mathbf{s}_i^*)}{\epsilon + \exp(\mathbf{s}_r^*)}) + \log(\frac{\exp(\mathbf{s}_r^*)}{\epsilon' + \exp(\mathbf{s}_r^*)}) - \log(\frac{\exp(\mathbf{s}_{\perp}^*)}{\epsilon + \exp(\mathbf{s}_{\perp}^*)})) \\
 &+ c(\mathbf{x}) (\log(\frac{\exp(\mathbf{s}'_t)}{\exp(\mathbf{s}_r^*)}) \\
 &= \mathbb{P}(Y = r, M = r|\mathbf{x}) (\log(\frac{\epsilon + \exp(\mathbf{s}_r^*)}{\epsilon' + \exp(\mathbf{s}_r^*)})) \\
 &+ \mathbb{P}(Y = r, M \neq r|\mathbf{x}) \log(\frac{\exp(\mathbf{s}_{\perp}^*)}{\exp(\mathbf{s}_r^*)}) \\
 &+ \mathbb{P}(Y \neq r, M = Y|\mathbf{x}) (\log(\frac{\exp(\mathbf{s}_r^*)}{\exp(\mathbf{s}_{\perp}^*)}) + \log(\frac{\epsilon + \exp(\mathbf{s}_r^*)}{\epsilon' + \exp(\mathbf{s}_r^*)})) \\
 &+ c(\mathbf{x}) (\log(\frac{\exp(\mathbf{s}'_t)}{\exp(\mathbf{s}_r^*)}) \\
 &\geq \mathbb{P}(Y = r, M = r|\mathbf{x}) (\log(\frac{\epsilon + \exp(\mathbf{s}_r^*)}{\epsilon' + \exp(\mathbf{s}_r^*)})) \\
 &+ (\mathbb{P}(Y \neq r, M = Y|\mathbf{x}) - \mathbb{P}(Y = r, M \neq r|\mathbf{x}) - c(\mathbf{x})) \log(\frac{\exp(\mathbf{s}_r^*)}{\exp(\mathbf{s}_{\perp}^*)}) \\
 &+ \mathbb{P}(Y = r, M \neq r|\mathbf{x}) \log(\frac{\epsilon + \exp(\mathbf{s}_r^*)}{\epsilon' + \exp(\mathbf{s}_r^*)}) > 0
 \end{aligned}$$

Which conclude the proof $\mathbf{s}_{\perp}^* > \mathbf{s}_r^*$ when $\mathbb{P}(Y \neq r, Y = M|\mathbf{x}) > \mathbb{P}(Y = r, Y \neq M|\mathbf{x}) + c(\mathbf{x})$ by contradiction.