

# Exploiting Human-AI Dependency for Learning to Defer

Anonymous Authors<sup>1</sup>

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 follows:

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

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

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

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

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

Where  $q = \arg \max_{y \in \mathcal{Y}} g_y(\mathbf{x})$  denotes the predicted class of the model.

The consistency proof for DCEc is provided in Appendix A.

## 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

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

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

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

Where  $t = \arg \max_{0 \leq i \leq J} \mathbb{P}(Y = M_i|\mathbf{x})$

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

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

It is noteworthy that when  $J = 1$  and  $m \neq y$ .

$$\begin{aligned} L_{DCE}^\perp(g(\mathbf{x}), y, m) = & -\log(\psi_{\mathcal{Y}}^y(g(\mathbf{x}))) - \log(\psi_{\mathcal{Y} \cup \perp}^{\mathcal{Y}}(g(\mathbf{x}))) \\ = & -\log(\psi_{\mathcal{Y} \cup \perp}^y(g(\mathbf{x}))), \end{aligned}$$

which is identical to the single expert version.

This generalized version of surrogate loss holds consistency for L2D with multiple experts. The proof is provided in Appendix B.

## References

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

Preliminary work. Under review by the International Conference on Machine Learning (ICML). Do not distribute.

## A. Consistency Proof of DCEc

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  $C_{\text{DCE}}^{\perp}$ , 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  $C_{\text{DCE}}^\perp$ , 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  $t = \arg \max_{\tilde{y} \in \mathcal{Y}} s'_y$ . 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_Y^r(\mathbf{s}')) - \log(\psi_Y^r(\mathbf{s}^*)) + \log(\psi_{Y^{\perp}/t}^{\perp}(\mathbf{s}')) - \log(\psi_{Y^{\perp}/r}^{\perp}(\mathbf{s}^*))) \\
 &+ \mathbb{P}(Y = r, M \neq r|\mathbf{x}) (\log(\psi_{Y^{\perp}}^r(\mathbf{s}')) - \log(\psi_{Y^{\perp}}^r(\mathbf{s}^*))) \\
 &+ \sum_{i \in \mathcal{Y}, i \neq r} \mathbb{P}(Y = i, M = i|\mathbf{x}) (\log(\psi_Y^i(\mathbf{s}')) - \log(\psi_Y^i(\mathbf{s}^*)) + \log(\psi_{Y^{\perp}/t}^{\perp}(\mathbf{s}')) - \log(\psi_{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_{Y^{\perp}}^i(\mathbf{s}')) - \log(\psi_{Y^{\perp}}^i(\mathbf{s}^*))) \\
 &+ c(\mathbf{x}) (\log(\psi_{Y^{\perp}}^q(\mathbf{s}')) - \log(\psi_{Y^{\perp}}^q(\mathbf{s}^*))) \\
 &= \mathbb{P}(Y = r, M = r|\mathbf{x}) \left( \log\left(\frac{\exp(\mathbf{s}_{\perp}^*)}{\epsilon + \exp(\mathbf{s}_{\perp}^*)}\right) - \log\left(\frac{\exp(\mathbf{s}_r^*)}{\epsilon + \exp(\mathbf{s}_r^*)}\right) + \log\left(\frac{\exp(\mathbf{s}_r^*)}{\epsilon' + \exp(\mathbf{s}_r^*)}\right) - \log\left(\frac{\exp(\mathbf{s}_{\perp}^*)}{\epsilon + \exp(\mathbf{s}_{\perp}^*)}\right) \right) \\
 &+ \mathbb{P}(Y = r, M \neq r|\mathbf{x}) \log\left(\frac{\exp(\mathbf{s}_{\perp}^*)}{\exp(\mathbf{s}_r^*)}\right) \\
 &+ \sum_{i \in \mathcal{Y}, i \neq r} \mathbb{P}(Y = i, M = i|\mathbf{x}) \left( \log\left(\frac{\exp(\mathbf{s}_i^*)}{\epsilon + \exp(\mathbf{s}_{\perp}^*)}\right) - \log\left(\frac{\exp(\mathbf{s}_i^*)}{\epsilon + \exp(\mathbf{s}_r^*)}\right) + \log\left(\frac{\exp(\mathbf{s}_r^*)}{\epsilon' + \exp(\mathbf{s}_r^*)}\right) - \log\left(\frac{\exp(\mathbf{s}_{\perp}^*)}{\epsilon + \exp(\mathbf{s}_{\perp}^*)}\right) \right) \\
 &+ c(\mathbf{x}) \left( \log\left(\frac{\exp(\mathbf{s}_t')}{\exp(\mathbf{s}_r^*)}\right) \right) \\
 &= \mathbb{P}(Y = r, M = r|\mathbf{x}) \left( \log\left(\frac{\epsilon + \exp(\mathbf{s}_r^*)}{\epsilon' + \exp(\mathbf{s}_r^*)}\right) \right) \\
 &+ \mathbb{P}(Y = r, M \neq r|\mathbf{x}) \log\left(\frac{\exp(\mathbf{s}_{\perp}^*)}{\exp(\mathbf{s}_r^*)}\right) \\
 &+ \mathbb{P}(Y \neq r, M = Y|\mathbf{x}) \left( \log\left(\frac{\exp(\mathbf{s}_r^*)}{\exp(\mathbf{s}_{\perp}^*)}\right) + \log\left(\frac{\epsilon + \exp(\mathbf{s}_r^*)}{\epsilon' + \exp(\mathbf{s}_r^*)}\right) \right) \\
 &+ c(\mathbf{x}) \left( \log\left(\frac{\exp(\mathbf{s}_t')}{\exp(\mathbf{s}_r^*)}\right) \right) \\
 &\geq \mathbb{P}(Y = r, M = r|\mathbf{x}) \left( \log\left(\frac{\epsilon + \exp(\mathbf{s}_r^*)}{\epsilon' + \exp(\mathbf{s}_r^*)}\right) \right) \\
 &+ (\mathbb{P}(Y \neq r, M = Y|\mathbf{x}) - \mathbb{P}(Y = r, M \neq r|\mathbf{x}) - c(\mathbf{x})) \log\left(\frac{\exp(\mathbf{s}_r^*)}{\exp(\mathbf{s}_{\perp}^*)}\right) \\
 &+ \mathbb{P}(Y = r, M \neq r|\mathbf{x}) \log\left(\frac{\epsilon + \exp(\mathbf{s}_r^*)}{\epsilon' + \exp(\mathbf{s}_r^*)}\right) > 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.

## B. Consistency Proof of DCEm

To begin with proof, let us define  $L = 2^{\{1,2,\dots,J\}}$  the powerset of  $\{1, 2, 3, \dots, J\}$ ,  $L_i = \{E \in L | i \in E\}$ . Let us use  $\mathbb{P}(Y = r, E|\mathbf{x})$  to denote the probability that  $Y = r$  and  $\forall i \in E, m_i = r, \forall j \notin E, m_j \neq r$ , for  $E \in L$ .  $\mathcal{Y}^{\perp, E} = \mathcal{Y} \cup \{\perp, i | i \notin E\}$ .

$$\begin{aligned}
 C_{\text{DCEm}}^{\perp}(g(\mathbf{x})) &= \sum_{y \in \mathcal{Y}} \sum_{E \in L} \left[ -\mathbb{P}(Y = y, E|\mathbf{x}) (\log(\psi_Y^y(g(\mathbf{x}))) + \sum_{i \in E} \log(\psi_{Y^{\perp}, i/q}^{\perp}(g(\mathbf{x})))) \right. \\
 &\quad \left. + \log(\psi_{Y^{\perp}, E}^{\mathcal{Y}}) \right]
 \end{aligned}$$

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

$$\begin{aligned} & C_{\text{DCEm}}^\perp(s^*) - C_{\text{DCEm}}^\perp(\tilde{s}) \\ &= \sum_{E \in L_i} \mathbb{P}(Y = r, E|x) (\log(\psi_{\mathcal{Y}}^y(g(\mathbf{x}))) + \sum_{i \in E} \log(\psi_{\mathcal{Y}^\perp, i/q}^\perp(g(\mathbf{x}))) + \log(\psi_{\mathcal{Y}^\perp, E}^\mathcal{Y})) \\ &> 0 \end{aligned}$$

Then we can conclude the proof of  $r = q$  by contradiction.

Then we prove that  $s_r^* > s_{\perp, i}^*$  when  $\mathbb{P}(Y = r, Y \neq M_i|x) > \mathbb{P}(Y \neq r, Y = M_i|x)$ . Suppose  $s_{\perp, i}^* > s_r^*$  when  $\mathbb{P}(Y = r, M_i \neq r|x) > \mathbb{P}(Y \neq r, Y = M_i|x)$  and we use  $s'$  to represent the score vector obtained by switching the value between  $s_r^*$  and  $s_{\perp, i}^*$ , i.e.  $s'_r = s_{\perp, i}^*$ ,  $s'_{\perp, i} = s_r^*$  and  $s'_j = s_j^*$  for all  $j \in \mathcal{Y}, j \neq r$ . We can directly derive that  $r = \arg \max_{\tilde{y} \in \mathcal{Y}} s'$ . Then  $C_{\text{DCEm}}^\perp(s^*) - C_{\text{DCEm}}^\perp(s')$  could be expressed as:

$$\begin{aligned} & C_{\text{DCEm}}^\perp(s^*) - C_{\text{DCEm}}^\perp(s') \\ &= \sum_{E \in L} -\mathbb{P}\left[(Y = y, E|x) (\log(\psi_{\mathcal{Y}}^y(g(\mathbf{x}))) + \sum_{i \in E} \log(\psi_{\mathcal{Y}^\perp, i/q}^\perp(g(\mathbf{x}))) \right. \\ &\quad \left. + \log(\psi_{\mathcal{Y}^\perp, E}^\mathcal{Y})\right] \end{aligned}$$

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

Lastly, we prove that  $s_{\perp}^* > s_r^*$  when  $\mathbb{P}(Y \neq r, Y = M|x) > \mathbb{P}(Y = r, Y \neq M|x)$ . Suppose  $s_r^* > s_{\perp}^*$  when  $\mathbb{P}(Y \neq r, Y = M|x) > \mathbb{P}(Y = r, Y \neq M|x)$  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  $t = \arg \max_{\tilde{y} \in \mathcal{Y}} s'_y$ . We define

$\epsilon = \sum_{i \in \mathcal{Y}, i \neq r} \exp(s^*)$  and  $\epsilon' = \sum_{i \in \mathcal{Y}, i \neq r} \exp(s')$ . 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{DCE}}^{\perp}(\mathbf{s}^*) - C_{\text{DCE}}^{\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}^*))) \\
 &= \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}^*)})) \\
 &= \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^*)})) \\
 &= \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})) \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  $\mathbb{P}(Y \neq r, Y = M|\mathbf{x}) > \mathbb{P}(Y = r, Y \neq M|\mathbf{x})$  by contradiction.