

When Can Proxies Improve the Sample Complexity of Preference Learning?

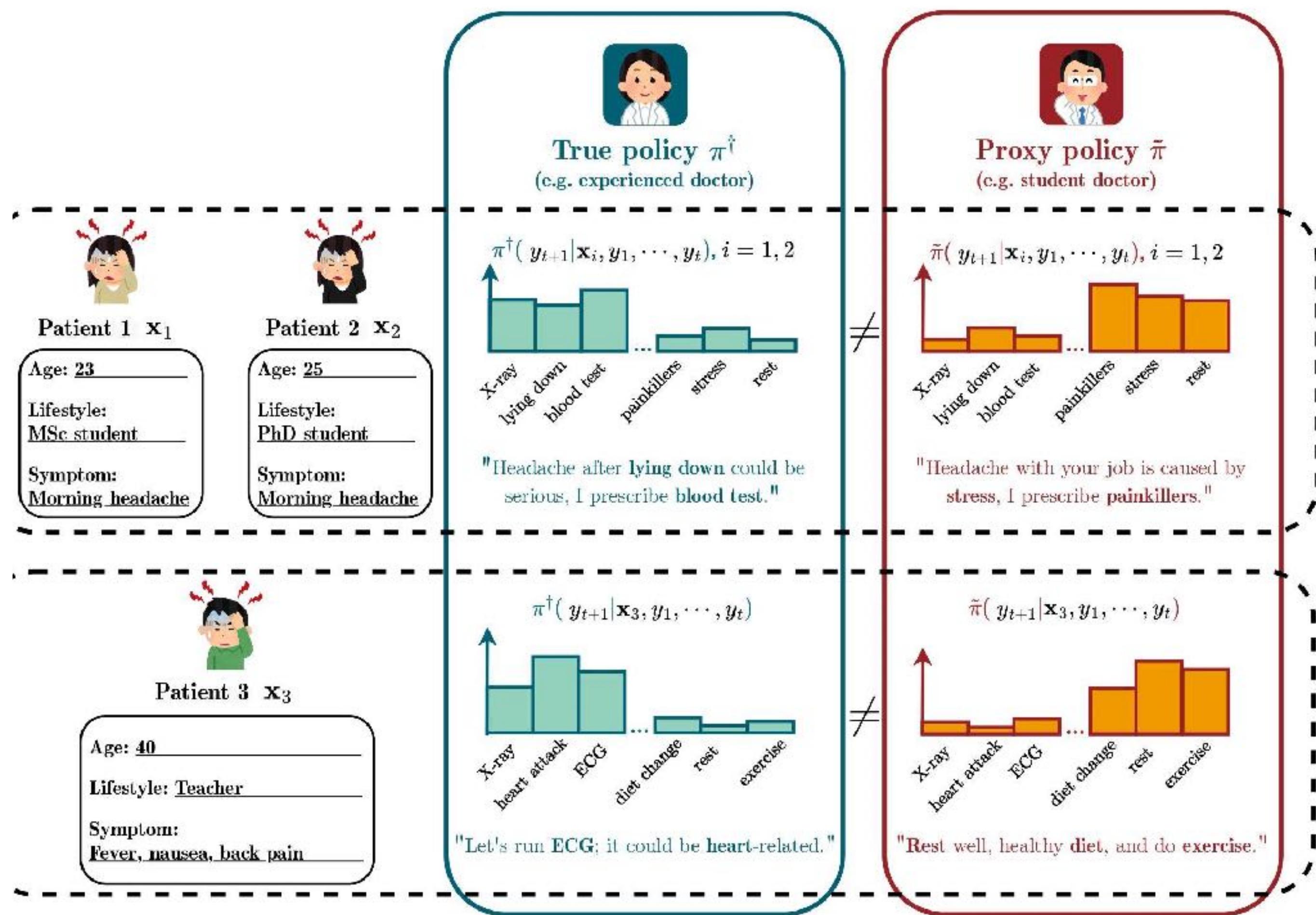


Yuchen Zhu, Daniel Augusto de Souza, Zhengyan Shi, Mengyue Yang,
Pasquale Minervini, Matt J. Kusner, Alexander D'Amour



Sufficient Conditions for Low-Rank/Dimensional Adaptation

I. Motivating Example



II. Mathematical set-up

$$\mathbf{x} \sim p_{\mathcal{X}}; \quad \mathbf{y}_1, \mathbf{y}_2 \stackrel{\text{i.i.d.}}{\sim} \pi_{\text{ref}}(\cdot | \mathbf{x}); \quad (\mathbf{y}_w, \mathbf{y}_l) = (\mathbf{y}_1, \mathbf{y}_2) \quad \text{if } b = 1$$

$$b \sim \text{Bern}(\sigma(r(\mathbf{x}, \mathbf{y}_1) - r(\mathbf{x}, \mathbf{y}_2))); \quad (\mathbf{y}_w, \mathbf{y}_l) = (\mathbf{y}_2, \mathbf{y}_1) \quad \text{if } b = 0$$

$$r(\mathbf{x}, \mathbf{y}) = \beta \log \frac{\pi(\mathbf{y} | \mathbf{x})}{\pi_{\text{ref}}(\mathbf{y} | \mathbf{x})}.$$

$$\arg \max_{\pi} \mathbb{E}_{(\mathbf{x}, \mathbf{y}_w, \mathbf{y}_l) \sim G} \left[\log \sigma \left(\beta \log \frac{\pi(\mathbf{y}_w | \mathbf{x})}{\pi_{\text{ref}}(\mathbf{y}_w | \mathbf{x})} - \beta \log \frac{\pi(\mathbf{y}_l | \mathbf{x})}{\pi_{\text{ref}}(\mathbf{y}_l | \mathbf{x})} \right) \right]$$

III. Sufficient Conditions

Condition 1. For a given metric $d_{\mathcal{P}_{\mathcal{Y}}}$ on the space of distributions on \mathcal{Y} , there is some positive scalar L such that

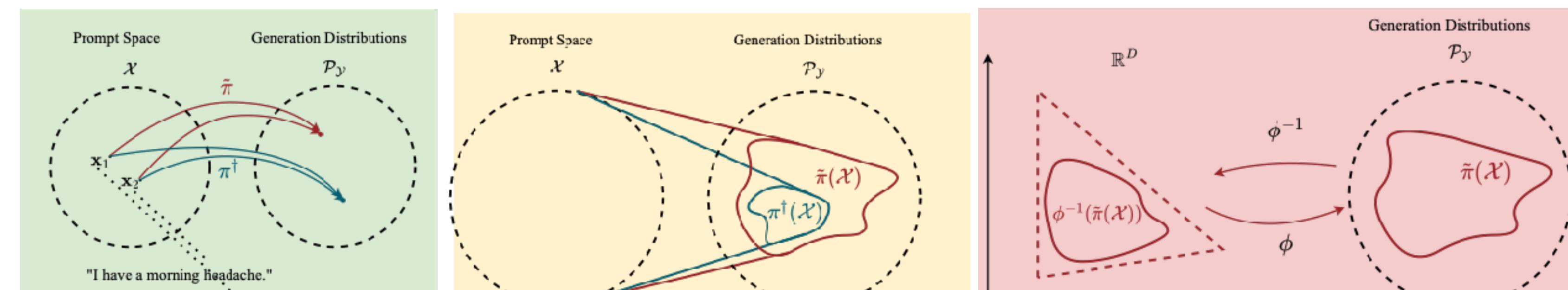
$$d_{\mathcal{P}_{\mathcal{Y}}}(\pi^{\dagger}(\cdot | x_1), \pi^{\dagger}(\cdot | x_2)) \leq L d_{\mathcal{P}_{\mathcal{Y}}}(\tilde{\pi}(\cdot | x_1), \tilde{\pi}(\cdot | x_2)) \quad (1)$$

Condition 2 (Image Inclusion).

$$\pi^{\dagger}(\mathcal{X}) \subseteq \tilde{\pi}(\mathcal{X}) \quad (2)$$

Condition 3 (Low-dimensional encoding). There exists an injective function $\phi: \mathcal{V} \rightarrow \mathcal{P}_{\mathcal{Y}}$, where $\mathcal{V} \subset \mathbb{R}^D$ is a bounded convex polytope with $D + 1$ vertices, such that:

1. The image $\phi(\mathcal{V})$ contains the image of the policies: $\pi^{\dagger}(\mathcal{X}) \subseteq \tilde{\pi}(\mathcal{X}) \subseteq \phi(\mathcal{V})$;
2. It is $(L_{\phi}, L_{\phi^{-1}})$ -bi-Lipschitz with its left inverse $\phi^{-1}: \mathcal{P}_{\mathcal{Y}} \rightarrow \mathcal{V}$: $\frac{1}{L_{\phi^{-1}}} \|v_1 - v_2\|_p \leq d(\phi(v_1), \phi(v_2)) \leq L_{\phi} \|v_1 - v_2\|_p$, where d is a metric on $\mathcal{P}_{\mathcal{Y}}$.



IV. Low-dimensional Adaptation

Theorem. We work under Conditions 1–3. For some D , there exists a Lipschitz invertible function $\tilde{\phi}: \mathcal{V} \rightarrow \tilde{\pi}(\mathcal{X})$ satisfying Condition 4, $\tilde{\Theta} \in \mathbb{R}^{N \times (D+1)}$ and $\tilde{\tau}^{\circ}: \mathcal{X} \rightarrow \Delta^D$ such that $\tilde{\pi} = \tilde{\phi} \circ \tilde{\Theta} \tilde{\tau}^{\circ}$. Moreover, for any $(\tilde{\phi}, \tilde{\Theta}, \tilde{\tau}^{\circ})$ such that $\tilde{\pi} = \tilde{\phi} \circ \tilde{\Theta} \tilde{\tau}^{\circ}$, there exists a Lipschitz continuous function $\bar{\pi}^{\dagger}: \Delta^D \rightarrow \Delta^D$ such that $\pi^{\dagger} = \tilde{\phi} \circ \tilde{\Theta} \circ \bar{\pi}^{\dagger} \circ \tilde{\tau}^{\circ}$.

$$\tilde{\pi} = \tilde{\phi} \circ \tilde{\Theta} \circ \tilde{\tau}^{\circ}$$

$$\pi^{\dagger} = \tilde{\phi} \circ \tilde{\Theta} \circ \bar{\pi}^{\dagger} \circ \tilde{\tau}^{\circ}$$

V. Improved Sample Complexity

Theorem. Assuming we have trained to convergence on proxy data, ...

we need

$$n(\epsilon, \omega) = \Omega \left(\frac{D}{\epsilon^2} \left(\frac{96L_{\phi}\|\tilde{\Theta}\|_p L_{\bar{\pi}} \sqrt{D}}{\epsilon} \right)^D \log \left(\frac{96L_{\phi}\|\tilde{\Theta}\|_p \sqrt{D}}{\epsilon} \right) - \log \omega \right)$$

samples to generalise.

Theorem. Otherwise ...

We need

$$n(\epsilon, \omega) = \Omega \left(\frac{D'}{\epsilon^2} \left(\frac{48L_{\phi}\|\tilde{\Theta}\|_p L_{\bar{\pi}} E'(p, D') \sqrt{D'}}{\epsilon} \right)^{D'} \log \left(\frac{48L_{\phi}\|\tilde{\Theta}\|_p L_{\bar{\pi}} E'(p, D') \sqrt{D'}}{\epsilon} \right) - \log \omega \right)$$

VI. Real world application

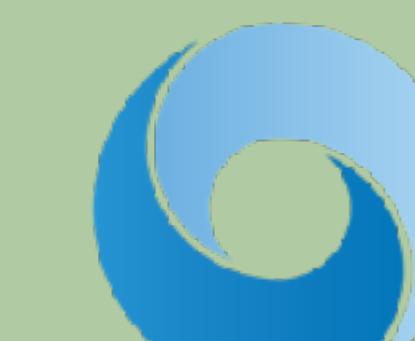
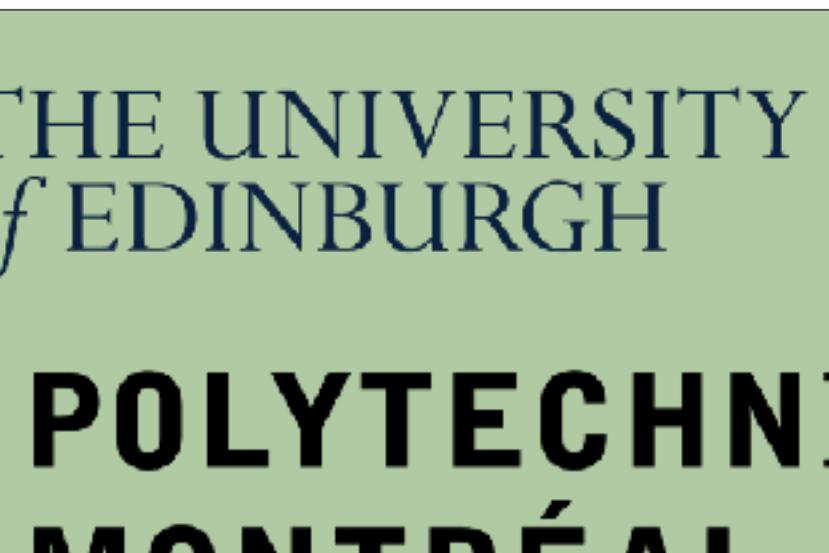
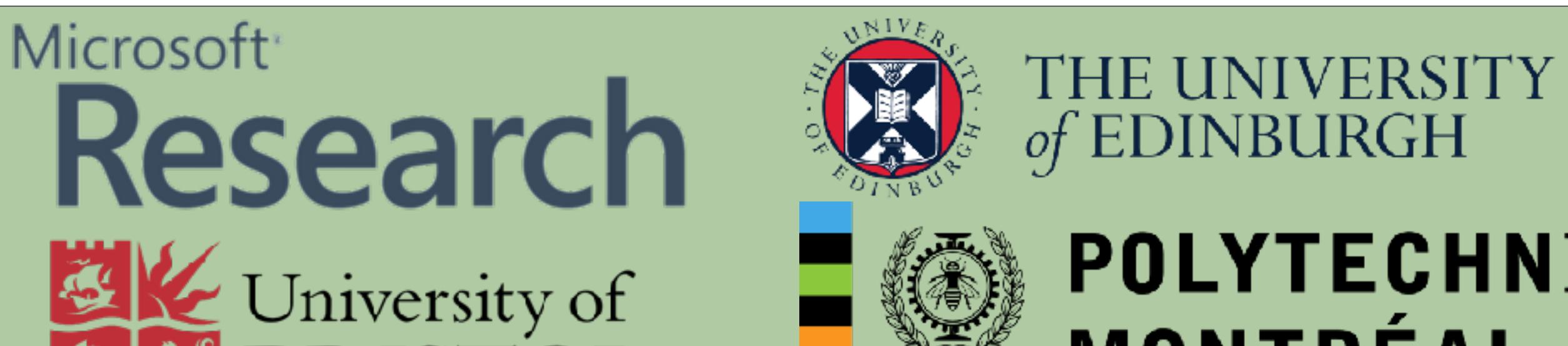
Tempered softmax gives a real-world example which satisfies our conditions.

- $\mathcal{X} = \mathbb{R}^5, \mathcal{P}_{\mathcal{X}} = \mathcal{N}(\mathbf{0}, I_5)$
- $\mathcal{Y} = \{1, 2, 3\}$, so $\mathcal{P}_{\mathcal{Y}} = \Delta^2$, a two-simplex.
- $D = 1$. Condition 3 is satisfied.
- $\pi^{\dagger}: \mathbb{R}^5 \rightarrow \mathbb{R} \rightarrow \Delta^2$.
- $\log(\tilde{\pi}(y_k | x)) = \frac{\log(\pi^{\dagger}(y_k | x))}{T}$. Temperature $T = 5$.
- $\pi_{\text{ref}} = \text{Uniform}\{1, 2, 3\}$

Division by a constant T satisfies conditions 1 and 2 where $d_{\mathcal{P}_{\mathcal{Y}}}$ is the log difference.

| | π_{ref} | π^{\dagger} | $\tilde{\pi}$ | $\tilde{\pi}_{\theta}$ | π_{θ}^{\dagger} |
|------|--------------------|-----------------|---------------|------------------------|--------------------------|
| mean | 0.63 | 0.0 | 0.33 | 0.34 | 0.32 |
| std | 0.00 | 0.0 | 0.00 | 0.014 | 0.096 |

Table 1: Results for the Tempered Reward experiment



Google DeepMind