

# Active Imitation Learning (+ other IL ideas)

August 7, 2024

## Contents

<b>1</b>	<b>Deterministic expert</b>	<b>1</b>
1.1	Setup	1
1.2	Mirror descent (trajectories)	1
1.3	DPOing the policy update	2
1.4	Mirror descent (history-dependent policies)	3
<b>2</b>	<b>DPO in sharpening</b>	<b>4</b>
<b>3</b>	<b>Misspecification</b>	<b>5</b>
3.1	Insufficiency of log loss	5
3.2	Misspecification in Hellinger distance for deterministic experts	5
3.3	Stochastic experts, starting point: Scheffe with TV distance	5

## 1 Deterministic expert

### 1.1 Setup

Let  $\Pi$  be a given deterministic policy class. Denote by  $\Pi_{\text{RNS}}$  the set of randomized nonstationary Markovian policies, and by  $\Pi_{\text{RNM}}$  the set of randomized non-Markovian policies.

For a fixed  $\hat{\pi}$ , define  $\Pi_\varepsilon(\hat{\pi}) := \{\pi \in \Pi : \rho(\pi, \hat{\pi}) \geq \varepsilon\}$  to be the set of policies from  $\Pi$  that disagree with  $\hat{\pi}$  with probability at least  $\varepsilon$ . Our objective is to find for another policy class  $\tilde{\Pi}$

$$\sup_{p \in \tilde{\Pi}} \mathcal{L}(p) := \inf_{\pi \in \Pi_\varepsilon(\hat{\pi})} \mathbb{P}^p[\exists h : \pi(x_h) \neq \hat{\pi}(x_h)]. \quad (1)$$

The class  $\tilde{\Pi}$  may be  $\Pi, \Pi_{\text{RNS}}, \Pi_{\text{RNM}}$ , in order of increasing generality. Let  $\tau = \{x_1, a_1, \dots, x_H, a_H, x_{H+1}\}$  denote a trajectory.

### 1.2 Mirror descent (trajectories)

Let  $\mathcal{P} = \{\mathbb{P}^\pi : \pi \in \tilde{\Pi}\}$  be the set of admissible laws of trajectories induced by rolling  $\pi \in \tilde{\Pi}$  out in the MDP. We can equivalently consider the problem in [Eq. \(1\)](#) as

$$\sup_{p \in \mathcal{P}} \mathcal{L}(p) = \inf_{\pi \in \Pi_\varepsilon(\hat{\pi})} \sum_{\tau} p(\tau) \cdot \mathbb{I}[\exists h : \pi(x_h) \neq \hat{\pi}(x_h)].$$

If  $\tilde{\Pi}$  is a convex set in policy space, then  $\mathcal{P}$  is also a convex set in trajectory space.

For a fixed  $p$ , define  $\pi_p := \min_{\pi \in \Pi} \mathbb{P}^p[\exists h : \pi(x_h) \neq \hat{\pi}(x_h)]$ . The derivative is

$$\frac{\partial \mathcal{L}(p)}{\partial p(\tau)} = \mathbb{I}[\exists h : \pi_p(x_h) \neq \hat{\pi}(x_h)].$$

From this it can be observed that  $\sup_{p \in \mathcal{P}} \|\nabla \mathcal{L}(p)\|_\infty \leq 1$ .

**Procedure.** The mirror descent procedure is as follows. We abbreviate  $\mathbb{P}^t \equiv \mathbb{P}^{p^t}$ ; similarly, the best-response policy with respect to  $p^t$  is  $\pi^t \equiv \pi_{p^t}$ . Initialize  $p^1$ , a law over trajectories. Then for  $t = 1, \dots, T$ :

1. Obtain the best-response policy w.r.t the current  $p^t$ ,

$$\pi^t = \operatorname{argmin}_{\pi \in \Pi_\varepsilon(\hat{\pi})} \mathbb{P}^t[\exists h : \pi(x_h) \neq \hat{\pi}(x_h)].$$

2. Solve the mirror descent update

$$p^{t+1} = \operatorname{argmin}_{p \in \mathcal{P}} \mathbb{P}^p[\exists h : \tilde{\pi}^t(x_h) \neq \hat{\pi}(x_h)] - \beta D_{\text{KL}}(p \| p^t). \quad (2)$$

**Analysis.** Let  $p^*$  be such that  $\mathcal{L}(p^*) = \sup_{p \in \mathcal{P}} \mathcal{L}(p)$ . The mirror descent guarantee states that

$$\begin{aligned} \sum_t \mathcal{L}(p^*) - \sum_t \mathcal{L}(p^t) &\leq \frac{D_{\text{KL}}(p^* \| p^1)}{\beta} + \frac{\beta}{2} \sum_t \|\nabla \mathcal{L}(p^t)\|_\infty^2 \\ &\leq \frac{\log(1 + D_{\chi^2}(p^* \| p^1))}{\beta} + \frac{T\beta}{2} \end{aligned}$$

Choosing  $\beta = \sqrt{\frac{2 \log(1 + D_{\chi^2}(p^* \| p^1))}{T}}$ ,

$$\sum_t \mathcal{L}(p^*) - \sum_t \mathcal{L}(p^t) \leq \sqrt{2 \log(1 + D_{\chi^2}(p^* \| p^1)) T}.$$

Put another way, there exists  $t \in [T]$  such that

$$\mathcal{L}(p^*) - \mathcal{L}(p^t) \leq \sqrt{\frac{2 \log(1 + D_{\chi^2}(p^* \| p^1))}{T}}.$$

This pays for log-coverability of  $p^*$  over  $p^1$ . However, we cannot actually solve for the update in [Eq. \(2\)](#) without constructing  $\mathcal{P}$ .

### 1.3 DPOing the policy update

Consider again the trajectory-level formulation. To make the DPO substitution, we first need to recover the maximizer to each MD update. Suppose  $p^t$  is fixed and admissible, in the sense that  $p^t = \mathbb{P}^{\tilde{\pi}^t}$  for some (possibly non-Markovian) policy  $\tilde{\pi}^t$ . Let

$$p_\star^{t+1} = \sup_{p \in \Delta(\tau)} \{ \mathbb{P}^p[\exists h : \pi^t(x_h) \neq \hat{\pi}(x_h)] - \beta D_{\text{KL}}(p \| p^t) \}. \quad (3)$$

For each  $\tau$ , this takes the closed form

$$\begin{aligned} \mathbb{I}[\exists h : \pi^t(x_h) \neq \hat{\pi}(x_h)] &= \beta \log \left( \frac{p_\star^{t+1}(\tau)}{p^t(\tau)} \right) + Z \\ &= \beta \log \left( \frac{\tilde{\pi}_\star^{t+1}(a_{1:H} \mid x_{1:H})}{\tilde{\pi}^t(a_{1:H} \mid x_{1:H})} \right) + Z, \end{aligned}$$

where  $\tilde{\pi}_*^{t+1}$  is the policy that induces  $p_*^{t+1}$ , i.e.,  $p_*^{t+1} = \mathbb{P}^{\tilde{\pi}_*^{t+1}}$ . **[TODO: Needs to be more exact. Does  $\tilde{\pi}_*^{t+1}$  always exist, e.g., by factoring out the transition probabilities?]**

With this substitution, we can consider the following alternative procedure. Given a policy class  $\Pi_{\text{DPO}}$ , data in the form of pairs of trajectories drawn as  $(\tau, \tau') \sim \pi_{\text{ref}}$ , and initial policy  $\tilde{\pi}^1$ ,

1. Obtain the best-response policy w.r.t the current  $\tilde{\pi}^t$ ,

$$\pi^t = \underset{\pi \in \Pi_\epsilon(\tilde{\pi})}{\operatorname{argmin}} \mathbb{P}^{\tilde{\pi}^t}[\exists h : \pi(x_h) \neq \hat{\pi}(x_h)].$$

2. Let  $g^t(\tau) := \mathbb{I}[\exists h : \pi^t(x_h) \neq \hat{\pi}(x_h)]$ . Solve the DPO update

$$\tilde{\pi}^{t+1} = \underset{\tilde{\pi} \in \Pi_{\text{DPO}}}{\operatorname{argmin}} \mathbb{E}_{\tau, \tau' \sim \pi_{\text{ref}}} \left[ \left( g^t(\tau) - g^t(\tau') - \beta \log \left( \frac{\tilde{\pi}(a_{1:H} \mid x_{1:H})}{\tilde{\pi}^t(a_{1:H} \mid x_{1:H})} \right) + \beta \log \left( \frac{\tilde{\pi}(a'_{1:H} \mid x'_{1:H})}{\tilde{\pi}^t(a'_{1:H} \mid x'_{1:H})} \right) \right)^2 \right] \quad (4)$$

**Assumption 1.1** (Policy completeness). *Given  $\beta$ , for any  $\pi \in \Pi_{\text{DPO}}$  and trajectory-level reward function  $g \in \{g^\pi(\tau) = \mathbb{I}[\exists h : \pi(x_h) \neq \hat{\pi}(x_h)] : \pi \in \Pi\}$ , there exists  $\pi' \in \Pi_{\text{DPO}}$  such that*

$$g(\tau) = \beta \log \left( \frac{\pi'(a_{1:H} \mid x_{1:H})}{\pi(a_{1:H} \mid x_{1:H})} \right) + Z, \quad \forall \tau.$$

**Sketch.**

- The update in Eq. (4) approximately solves Eq. (3), and pays for all-policy coverage over  $\pi_{\text{ref}}$
- **[TODO: Mirror descent guarantee involves  $p^*$  computed from what policy class? What kind of policy class is  $\Pi_{\text{DPO}}$ ?]**

## 1.4 Mirror descent (history-dependent policies)

Here we consider mirror descent in (possibly history-dependent) policy space. Let  $\tilde{x}_h = \{x_1, a_1, \dots, x_{h-1}, a_{h-1}, x_h\}$  be the history up until time  $h$ . Now  $p \in \tilde{\Pi}$  maps  $\tilde{x}_h \rightarrow \Delta(\mathcal{A})$ . First we calculate the gradient for a fixed  $p$ .

$$\frac{\partial \mathcal{L}(p)}{\partial p(a_h \mid \tilde{x}_h)} = \mathbb{P}^p(\tilde{x}_h) \cdot \mathbb{P}^p[\exists h : \pi_p(x_h) \neq \hat{\pi}(x_h) \mid \tilde{x}_h, a_h].$$

**Procedure.** The mirror descent procedure is as follows. Initialize  $p^1$ . Then for  $t = 1, \dots, T$ :

1. Obtain the best-response policy w.r.t the current  $p^t$ ,

$$\pi^t = \underset{\pi \in \Pi_\epsilon(\tilde{\pi})}{\operatorname{argmin}} \mathbb{P}^t[\exists h : \pi(x_h) \neq \hat{\pi}(x_h)].$$

2. Compute the surrogate value function

$$Q^t(\tilde{x}_h, a_h) = \mathbb{P}^t[\exists h : \pi^t(x_h) \neq \hat{\pi}(x_h) \mid \tilde{x}_h, a_h].$$

3. Solve the mirror descent update

$$\begin{aligned} p^{t+1} &= \sup_{p \in \tilde{\Pi}} \sum_h \mathbb{E}^{p^t} [Q^t(\tilde{x}_h, p(\tilde{x}_h))] - \beta \underbrace{D_{\text{KL}}(\mathbb{P}^p \parallel \mathbb{P}^t)}_{= \mathbb{E}^p[\sum_h D_{\text{KL}}(p(\tilde{x}_h) \parallel p^t(\tilde{x}_h))]} \\ &= \mathbb{E}^p[\sum_h D_{\text{KL}}(p(\tilde{x}_h) \parallel p^t(\tilde{x}_h))] \end{aligned}$$

**Analysis.** **[TODO: Conjugate norm?]** The regularization term is strongly convex with respect to  $p(a_h \mid \tilde{x}_h)$  in the  $\|\cdot\|_{1, d^p(\tilde{x}_h)}$  norm...? But there is a mismatch between distributions over which expectations are taken in the two terms.

## 2 DPO in sharpening

Dropping context dependence for simplicity, for each  $t$  the objective is to solve

$$\Phi^t(\pi) = \mathbb{E}_\pi[\log \pi^1(a)] - \beta D_{\text{KL}}(\pi \parallel \pi^t). \quad (5)$$

Given data collected from a data collection policy  $\mu^t$ , the DPO-style update is

$$\mathcal{L}^t(\pi) = \mathbb{E}_\mu \left[ \left( \log \pi^1(a) - \log \pi^1(b) - \beta \log \frac{\pi(a)}{\pi^t(a)} + \beta \log \frac{\pi(b)}{\pi^t(b)} \right)^2 \right]. \quad (6)$$

Our goal is to derive the number of samples for  $\max_{\pi \in \Pi} \Phi^t(\pi) - \Phi^t(\hat{\pi}) \leq \varepsilon$ , where  $\hat{\pi} = \max_{\pi \in \Pi} \mathcal{L}^t(\pi)$ .

**Lemma 2.1.** *For any comparator policy  $\pi$ ,*

$$\Phi^t(\pi) - \Phi^t(\hat{\pi}) \leq \mathbb{E}_{a \sim \pi, b \sim \hat{\pi}} \left[ \log \pi^1(a) - \log \pi^1(b) - \beta \log \left( \frac{\hat{\pi}(a)}{\pi^t(a)} \right) + \beta \log \left( \frac{\hat{\pi}(b)}{\pi^t(b)} \right) \right].$$

**Sketch.** Suppose we collect data with  $\pi^t$  (another more advanced option is  $\pi^1$ ). Then we expect that, roughly,

$$\Phi^t(\pi^{t+1,*}) - \Phi^t(\hat{\pi}) \lesssim \left\| \frac{\pi^{t+1,*}}{\pi^t} \right\|_\infty \cdot \left\| \frac{\hat{\pi}}{\pi^t} \right\|_\infty \cdot \varepsilon_{\text{stat}},$$

where  $\varepsilon_{\text{stat}}^2 = \mathbb{E}_{(a,b) \sim \pi^t} \left[ \left( \log \pi^1(a) - \log \pi^1(b) - \beta \log \left( \frac{\hat{\pi}(a)}{\pi^t(a)} \right) + \beta \log \left( \frac{\hat{\pi}(b)}{\pi^t(b)} \right) \right)^2 \right]$ .

[TODO:

- **Concentrability of  $\pi^{t+1,*}$ :** Need to check bound on  $\frac{\pi^{t+1,*}}{\pi^t}$  given  $R_{\max} \in (\infty, 0]$ , or at least something very small as the lower limit.
- **Concentrability of  $\hat{\pi}$ :** This is all-policy. How to control?
- **Dependence on ‘reward range’:** May need  $\left| \log \frac{\pi^1(a)}{\pi^1(b)} \right| \leq C$  for all pairs of actions  $(a, b)$ , is there something more refined? Clipping if we’re willing to throw away actions with very low probability?

]

**Proof of Lemma 2.1.** For a reward function  $r$  define  $\Phi_r^t(\pi) := \mathbb{E}_\pi[r(a)] - \beta D_{\text{KL}}(\pi \parallel \pi^t)$ , and let  $\hat{r} = \beta \log \left( \frac{\hat{\pi}}{\pi^t} \right)$ . Since  $\Phi_{\hat{r}}(\hat{\pi}) = 0$ , we have  $\beta D_{\text{KL}}(\hat{\pi} \parallel \pi^t) = \mathbb{E}_{\hat{\pi}}[\hat{r}]$ .

Further, from Lemma E.2 of our paper, we know that  $\hat{\pi} = \arg\max_{\pi \in \Delta(\mathcal{A})} \mathbb{E}_\pi[\hat{r}] - \beta D_{\text{KL}}(\pi \parallel \pi^t)$ . Since  $\Phi_{\hat{r}}(\pi) \leq \Phi_{\hat{r}}(\hat{\pi}) = 0$ , we observe that  $-\beta D_{\text{KL}}(\pi \parallel \pi^t) \leq -\mathbb{E}_\pi[\hat{r}]$ .

Decomposing the error and using the above inequalities, we obtain

$$\begin{aligned} \Phi^t(\pi) - \Phi^t(\hat{\pi}) &= \mathbb{E}_\pi[\log \pi^1(a)] - \beta D_{\text{KL}}(\pi \parallel \pi^t) - \mathbb{E}_{\hat{\pi}}[\log \pi^1(a)] + \beta D_{\text{KL}}(\hat{\pi} \parallel \pi^t) \\ &\leq \mathbb{E}_\pi[\log \pi^1(a) - \hat{r}(a)] - \mathbb{E}_{\hat{\pi}}[\log \pi^1(a) - \hat{r}(a)] \\ &= \mathbb{E}_{a \sim \pi, b \sim \hat{\pi}} \left[ \log \pi^1(a) - \log \pi^1(b) - \beta \log \left( \frac{\hat{\pi}(a)}{\pi^t(a)} \right) + \beta \log \left( \frac{\hat{\pi}(b)}{\pi^t(b)} \right) \right]. \end{aligned}$$

□

### 3 Misspecification

#### 3.1 Insufficiency of log loss

#### 3.2 Misspecification in Hellinger distance for deterministic experts

Is it possible to learn a policy where the error scales with misspecification in Hellinger distance? This is a natural goal to pursue since we care about outputting a policy that's close to  $\pi^*$  in Hellinger distance. Further,  $D_H^2(P, Q) \leq D_{KL}(P \parallel Q) \leq \log(1 + D_{\chi^2}(P \parallel Q))$ , so this would be a strict improvement on the misspecification error under log loss.

For deterministic expert policies, this is possible with the  $L_{\max}$  loss, and it doesn't require known transitions. Here, the  $L_{\max}$  loss is equivalent to Hellinger distance up to a constant, and doesn't suffer from the same blow-up issues when a candidate policy is off-support relative to  $\pi^*$ .

Recall that for a deterministic policy  $\pi^*$ ,

$$L_{\max}(\pi) = \mathbb{E}^{\pi^*} \mathbb{E}_{a'_{1:H} \sim \pi(x_{1:H})} [\mathbb{I}[\exists h : a'_h \neq a_h]],$$

and  $\hat{L}_{\max}$  is the empirical version. It can be observed that

$$\frac{1}{2} L_{\max}(\pi) \leq D_H^2(\mathbb{P}^\pi, \mathbb{P}^{\pi^*}) \leq 2 L_{\max}(\pi).$$

#### 3.3 Stochastic experts, starting point: Scheffe with TV distance

Let  $\mathcal{P} = \{\mathbb{P}^\pi : \pi \in \Pi\}$ . For any  $P, Q \in \mathcal{P}$ , define the witness function

$$g_{P,Q} = \operatorname{argmax}_{|g| \leq \frac{1}{2}} \mathbb{E}_P[g] - \mathbb{E}_Q[g]$$

and the set of discriminator functions as

$$\mathcal{G} = \{g_{P,Q} : P, Q \in \mathcal{P}, P \neq Q\}.$$

Output the policy

$$\hat{\pi} = \operatorname{argmin}_{\pi \in \Pi} \max_{g \in \mathcal{G}} \hat{\mathbb{E}}[g] - \mathbb{E}_{\mathbb{P}^\pi}[g]. \quad (7)$$

**Proposition 3.1.** *The output of Eq. (7) satisfies*

$$D_{TV}(\mathbb{P}^{\pi^*}, \mathbb{P}^{\hat{\pi}}) \leq 3 \min_{\pi \in \Pi} D_{TV}(\mathbb{P}^{\pi^*}, \mathbb{P}^\pi) + 2\varepsilon_{\text{stat}},$$

where  $\varepsilon_{\text{stat}} := \max_{g \in \mathcal{G}} |\hat{\mathbb{E}}[g] - \mathbb{E}_{\pi^*}[g]|$ .

**[TODO: comparison to hellinger bound; incomparable]**

**Proof of Proposition 3.1.** Fix any  $\bar{\pi} \in \Pi$ . Using the triangle inequality,

$$D_{TV}(\mathbb{P}^{\pi^*}, \mathbb{P}^{\hat{\pi}}) \leq D_{TV}(\mathbb{P}^{\pi^*}, \mathbb{P}^{\bar{\pi}}) + D_{TV}(\mathbb{P}^{\bar{\pi}}, \mathbb{P}^{\hat{\pi}}).$$

Let  $\tilde{g} = g_{\mathbb{P}^{\bar{\pi}}, \mathbb{P}^{\hat{\pi}}}$ . By construction,  $\tilde{g} \in \mathcal{G}$  so

$$\begin{aligned} D_{TV}(\mathbb{P}^{\bar{\pi}}, \mathbb{P}^{\hat{\pi}}) &= \mathbb{E}_{\bar{\pi}}[\tilde{g}] - \mathbb{E}_{\hat{\pi}}[\tilde{g}] \\ &= \mathbb{E}_{\bar{\pi}}[\tilde{g}] - \hat{\mathbb{E}}[\tilde{g}] + \hat{\mathbb{E}}[\tilde{g}] - \mathbb{E}_{\hat{\pi}}[\tilde{g}] \\ &\leq \mathbb{E}_{\bar{\pi}}[\tilde{g}] - \hat{\mathbb{E}}[\tilde{g}] + \max_{g \in \mathcal{G}} \{\hat{\mathbb{E}}[g] - \mathbb{E}_{\hat{\pi}}[g]\} \end{aligned}$$

$$\leq \mathbb{E}_{\bar{\pi}}[\tilde{g}] - \widehat{\mathbb{E}}[\tilde{g}] + \max_{g \in \mathcal{G}} \left\{ \widehat{\mathbb{E}}[g] - \mathbb{E}_{\bar{\pi}}[g] \right\},$$

since  $\hat{\pi} = \operatorname{argmin}_{\pi \in \Pi} \max_{g \in \mathcal{G}} \left\{ \widehat{\mathbb{E}}[g] - \mathbb{E}_{\pi}[g] \right\}$ . Next, define  $\varepsilon_{\text{stat}} = \max_{g \in \mathcal{G}} \left| \widehat{\mathbb{E}}[g] - \mathbb{E}_{\pi^*}[g] \right|$ . Letting  $\bar{g} = \max_{g \in \mathcal{G}} \widehat{\mathbb{E}}[g] - \mathbb{E}_{\bar{\pi}}[g]$ , we have

$$\begin{aligned} D_{\text{TV}}(\mathbb{P}^{\bar{\pi}}, \mathbb{P}^{\hat{\pi}}) &\leq \mathbb{E}_{\bar{\pi}}[\tilde{g}] - \widehat{\mathbb{E}}[\tilde{g}] + \widehat{\mathbb{E}}[\tilde{g}] - \mathbb{E}_{\bar{\pi}}[\tilde{g}] \\ &\leq \mathbb{E}_{\bar{\pi}}[\tilde{g}] - \mathbb{E}_{\pi^*}[\tilde{g}] + \mathbb{E}_{\pi^*}[\tilde{g}] - \mathbb{E}_{\bar{\pi}}[\tilde{g}] + 2\varepsilon_{\text{stat}} \\ &\leq 2 \sup_{|g| \leq \frac{1}{2}} \{ \mathbb{E}_{\bar{\pi}}[g] - \mathbb{E}_{\pi^*}[g] \} + 2\varepsilon_{\text{stat}} \\ &= 2D_{\text{TV}}(\mathbb{P}^{\bar{\pi}}, \mathbb{P}^{\pi^*}) + 2\varepsilon_{\text{stat}}. \end{aligned}$$

□

For the [Proposition 3.1](#) to hold for general  $f$ -divergences, we need (1) a general version of the triangle inequality to isolate the misspecification term in the first step; (2) a concentration inequality from  $\widehat{\mathbb{E}}$  to  $\mathbb{E}_{\pi^*}$  that holds for all  $g \in \mathcal{G}$ , which means that  $\mathcal{G}$  must be bounded; and (3) possibly symmetry of the discriminator set. These properties should be satisfied by the Hellinger distance and triangular discrimination, in addition to TV.

#### Questions.

- **Imitation learning (known dynamics):** What objective should we use for the triangular discrimination metric to get a fast rate? What happens when you try to use Hellinger directly?
- **Imitation learning (unknown dynamics):** The objective in [Eq. \(7\)](#) requires known dynamics. How many queries to the model are required if the dynamics are not known? Can we lower bound the number of queries?
- **Distribution learning:** For the analysis in [Proposition 3.1](#) to go through, we need the divergence to be symmetric and bounded. Can either of these be relaxed to extend this result to general  $f$ -divergences?
- **Online imitation learning:** What questions can we generate in this setting?