

# 1) Dataset Aggregation w/ DAgger

Setting: Discounted Infinite Horizon MDP

$$\mathcal{M} = \{ \mathcal{S}, \mathcal{A}, \mathcal{P}, \mathcal{R}, \gamma \}$$

unknown!  
possibly unobserved

Expert knows the optimal policy  $\pi^*$   
and we query the expert at any state  
during training

## Algorithm: DAgger

Initialize  $\pi^0$  and dataset  $\mathcal{D} = \emptyset$

For  $t=0, \dots, T-1$

1) Generate Dataset with  $\pi^t$  & expert

$$\mathcal{D}^t = \{ s_i, a_i^* \}_{i=1}^N \quad s_i \sim d_{\mathcal{M}}^{\pi^t} \quad a_i^* = \pi^*(s_i)$$

2) Data Aggregation:  $\mathcal{D} = \mathcal{D} \cup \mathcal{D}^t$

3) Update Policy via SL:

$$\pi^{t+1} = \underset{\pi \in \Pi}{\operatorname{argmin}} \left( \mathbb{E}_{s, a \sim \mathcal{D}} [\ell(\pi, s, a)] \right)$$

→ small if  $\pi(s) \approx a$

## 2) Online Learning

captures idea of learning from additional data over time

Iterative w/ 2 components

For  $t=0, 1, \dots, T-1$

1) learner chooses  $\theta_t$

2) Suffer the risk  $\mathcal{R}_t(\theta_t) = \mathbb{E}_{z \sim \mathcal{D}_t} [\ell(\theta_t, z)]$   
(expected loss)

We care about average regret

$$\frac{1}{T} R(T) = \frac{1}{T} \left[ \sum_{t=0}^{T-1} R_t(\theta_t) - \min_{\theta} \sum_{t=0}^{T-1} R_t(\theta) \right]$$

The baseline is the best parameter in hindsight

Difference from SL setting:

$D_t$  (&  $R_t$ ) can vary in many ways

Example: in DAgger, we choose  $\pi^t$   
and suffer  $\mathbb{E} \left[ \ell(\pi^t, S, \pi^*(S)) \right]$   
 $S \sim \mathcal{D}_M$

How should learner choose  $\theta_t$ ?

Algorithm: Follow the Regularized Leader

For  $t=0, 1, \dots, T-1$

$$\theta_t = \min_{\theta} \sum_{k=0}^{t-1} R_k(\theta) + \lambda f(\theta)$$

regularizer

$$\sum_{k=0}^{t-1} \mathbb{E}[\ell(\theta, z)] = \mathbb{E}[\ell(\theta, z)]$$

$z \sim \mathcal{D}_k$

$z \sim \frac{1}{t} \sum_{k=0}^{t-1} \mathcal{D}_k$

data aggregation

Theorem (FTL): if loss functions are convex  
and regularizer is strongly convex, then

$$\max_{R_0, \dots, R_{T-1}} \frac{1}{T} \left[ \sum_{t=0}^{T-1} R_t(\theta_t) - \min_{\theta} \sum_{t=0}^{T-1} R_t(\theta) \right] \leq O(1/\sqrt{T})$$

### 3) Analysis of Dagger

Corollary: if  $\ell(\pi^*, s, \pi^*(s)) = 0$ , then e.g.  $\ell(\pi(s) - a) \frac{1}{2}$   
and  $\ell > 0$

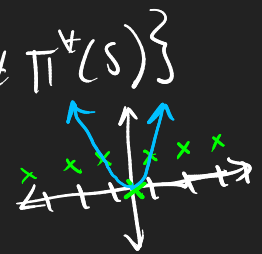
$$\min_{0 \leq t \leq T-1} \mathbb{E}_{\text{SMD}_{\pi^t}} \left[ \ell(\pi^t, s, \pi^*(s)) \right] \leq O(\sqrt{T}) = \epsilon_{FTL}$$

Proof:  $\pi^t$  plays role of  $\theta_t$ ,  $(s, \pi^*(s))$  is  $z$ ,  
 $d_M^{\pi^t}$  is  $D_t$

$$\begin{aligned} \min_{0 \leq t \leq T-1} R_t(\pi^t) &\leq \frac{1}{T} \sum_{t=0}^{T-1} R_t(\pi^t) \\ &= \frac{1}{T} \sum_{t=0}^{T-1} R_t(\pi^t) - \underbrace{R_t(\pi^*)}_{=0} \\ &\stackrel{\text{defined by } d_M^{\pi^t}}{=} \frac{1}{T} \left( \sum_{t=0}^{T-1} R_t(\pi^t) - \min_{\pi} \sum_{t=0}^{T-1} R_t(\pi) \right) \\ &\stackrel{\text{defined by arbitrary dist}}{\leq} \max_{R_0, \dots, R_{T-1}} \frac{1}{T} \left( \sum_{t=0}^{T-1} R_t(\pi^t) - \min_{\pi} \sum_{t=0}^{T-1} R_t(\pi) \right) \\ &\leq O(\sqrt{T}) = \epsilon_{FTL} \quad \square \end{aligned}$$

Key fact: accuracy guarantee on  $d_M^{\pi^t}$  instead

Theorem: if  $\ell(\pi, s, \pi^*(s)) \geq \underline{\epsilon} \mathbb{1}\{\pi(s) \neq \pi^*(s)\}$   
Then there is  $t \in \{0, \dots, T-1\}$  such that

$$\mathbb{E}_{\text{SMD}_{\pi^t}} [V^{\pi^*}(s) - V^{\pi^t}(s)] \leq \frac{(\max_{s,a} |A^{\pi^*}(s,a)|)}{1-\gamma} \cdot \epsilon_{FTL}$$


$$A^{\pi^*}(s,a) \leq 0 \quad \forall s,a \quad \left[ -Q^{\pi^*}(s,a) - V^{\pi^*}(s) \right]$$

$\max_{s,a} |A^{\pi^*}(s,a)|$  is the cost of messing up at one timestep.

Proof We apply PDL in other direction

$$\mathbb{E}_{S \sim \pi^*} [V^{\pi^*}(s) - V^{\pi^t}(s)] = \frac{1}{1-\gamma} \mathbb{E}_{S \sim d_{\pi^t}} [A^{\pi^*}(s, \pi^t(s))]$$

$$= \frac{1}{1-\gamma} \mathbb{E}_{S \sim d_{\pi^t}} [A^{\pi^*}(s, \pi^t(s)) - A^{\pi^*}(s, \pi^*(s))]$$

$$\geq \frac{1}{1-\gamma} \mathbb{E}_{S \sim d_{\pi^t}} \left[ \max_{s,a} |A^{\pi^*}(s,a)| \mathbb{1}_{\{\pi^t(s) \neq \pi^*(s)\}} \right]$$

$$\geq \frac{1}{1-\gamma} \max_{s,a} |A^{\pi^*}(s,a)| \cdot \overbrace{\mathbb{E}_{S \sim d_{\pi^t}} [\ell(\pi^t, s, \pi^*(s))]}^{R_t}$$

$$\mathbb{E}_{S \sim \pi^*} [V^{\pi^*}(s) - V^{\pi^t}(s)] \leq \frac{1}{1-\gamma} \max_{s,a} |A^{\pi^*}(s,a)| \cdot \epsilon_{FTL}$$