

4. Proposed Work: Temporal Difference Learning & Apprenticeship Learning

Imitation Learning

Expert Feedback

Machine Learning Algorithm

Policy π



• SVM

- Gaussian Process
- Kernel Estimator
- Deep Networks
- Random Forests
- LWR

Maps states to <u>actions</u>

Efficient RL via Imitation



Extra Assumptions: Planner & Controller (robotics) (reward) signal Optimal Planner, MPC) [Choudhury, et.al, ICRA, RSS, 17, Pan et.al, 17]

Ground Truth Labels + Utility (NLP)

Search Algorithm (e.g. A*) as expert

Why bother imitating when you have cost/reward signals?



Formalizing Advantages

ADVANTAGE #1

Less sensitive to local optimality

DAgger (Data Aggregation) [Ross et.al, 11, AISTATS]

AggreVaTe (Aggregate with Values) [Ross&Bagnell14, arxiv]

ADVANTAGE #2 More Sample Efficient (i.e., Learns faster)

There exist problems, s.t. with access to optimal expert, i.e., $\pi^e=\pi^*$

IL learns exponentially faster than RL

[Sun et.al, 17,ICML]

Consider a simple, tree like MDP

let us assume $\pi^e = \pi^*$, and we can query $Q^*(s, a)$



 $Q^{\pi}(s,a)\,$: the total cost of taking action $a\,$ at $s\,$ then following policy $\pi\,$



Halving: Eliminate half of the nodes at every iteration

[Sun et.al, 17,ICML]

In time logarithmic in #states, we know an optimal policy



In time logarithmic in #states, we know an optimal policy

Now if we can only query **unbiased but noisy** Q^*





But For Pure RL



Proof uses a reduction from Multi-Armed Bandit

Ex: AggreVaTe



Cost-Sensitive classification dataset



Towards Differentiable AggreVaTe (AggreVaTeD)

[Sun et.al, 17, ICML]

Stochastic parameterized policy:

 $\pi_{\theta} = \pi(\cdot|s;\theta)$



[Sun et.al, 17, ICML]

Differentiable AggreVaTe (AggreVaTeD)

$$\pi_{\theta_{n}} \longrightarrow \{s, \begin{bmatrix} Q^{*}(s, a_{1}) \\ \dots \\ Q^{*}(s, a_{A}) \end{bmatrix}\}_{N}$$

$$\nabla_{\theta} \ell_{n}(\theta)|_{\theta_{n}} \longrightarrow \ell_{n}(\theta) = \sum_{s} \sum_{a} \pi(a|s; \theta)Q^{*}(s, a)$$
Cost-Sensitive loss on
$$Cost-Sensitive loss on$$

AggreValeD-GD (Gradient Descent) $\theta_{n+1} = \theta_n - \mu \nabla_{\theta} \ell_n(\theta)|_{\theta = \theta_n}$ Cost-Sensitive loss on the *new batch (no aggregation)*

AggreVaTeD-NG (Natural Gradient):

$$\theta_{n+1} = \theta_n - \eta_n I(\theta_n)^{-1} \nabla_{\theta_n} \ell_n(\theta_n)$$

[Sun et.al, 17, ICML]

Differentiable AggreVaTe (AggreVaTeD)

Discrete MDP, Tabular Policy:



GD ensures *Convergence*

34

Dependency Parsing

Dependency Parsing on Handwritten Algebra Data



Dependency Parsing

Dependency Parsing <=> Sequential Decision Making



Performance of AggreVaTeD, RL, and DAgger



37

RL: Natural Policy Gradient [Kakade02, Bagnell 04] DAgger result from Duyck & Gordon, 15

