CHALLENGE:



Let us build "expert" on the fly



Proposed Work: Temporal Difference Learning & Apprenticeship Learning

Example: AlphaGo-Zero

[Silver, et.al, 17, Nature]

known & deterministic model



A slow Policy η (MCTS)



Example: AlphaGo-Zero

[Silver, et.al, 17, Nature]

known & deterministic model



CHALLENGE: What if model is unknown $\hat{P}(\cdot|s,a) \approx P(\cdot|s,a), \forall s \in S, a \in \mathcal{A}$ Not realistic



Model-based RL (e.g., iLQR) within a Trust-Region

[openAl Gym]

Dual Policy Iteration

[Sun et.al, 18, submitted to ICML]



[Sun et.al, 18, submitted to ICML]

Dual Policy Iteration



[Sun et.al, 18, submitted to ICML]

Dual Policy Iteration



Experiments

Conservative Policy Iteration [Kakade&Langford, 02] Synthetic discrete_MDP_1000 **Discrete MDPs** The lower the better (log-scale) 21.6 0.10 21.4 a1 0.20 21.2 0.95 0.05 21.0 0.30 0.40 AggreVaTeD_VI CPI 0.30 20.8 10 20 30 40 50 Garnet Problems **Batch** Iteration [Scherrer 14, ICML]

Our Approach: Frank-Wolf + VI

-1

Experiments



Summary

