A brief introduction to imitation learning

Reinforcement Learning



Reinforcement learning (RL):

Aim to finding optimal strategies through interaction with the environment



The interaction process can be modeled as MDP $< S, A, R, P, \gamma, (D) >$

Reinforcement Learning





$$egin{aligned} G_t &\doteq R_{t+1} + \gamma R_{t+2} + \gamma^2 R_{t+3} + ... = \sum_{k=0}^\infty \gamma^k R_{t+k+1} \ v_\pi(s) &= \mathbb{E}_\pi[G_t | S_t = s] \ q_\pi(s,a) &\doteq \mathbb{E}_\pi[G_t | S_t = s, A_t = a] \ V_\pi &= \mathbb{E}_\pi[G_t] = \mathbb{E}_\pi[\sum_{t=0}^\infty \gamma^t R(s_t,a_t)] \end{aligned}$$

$$egin{aligned} &
ho_{\pi}(s,a) = \pi(a|s)\sum_{t=0}^{\infty}\gamma^{t}P(S_{t}=s|\pi) \ &V_{\pi} = \sum_{s,a}
ho(s,a)R(s,a) \end{aligned}$$

Cons of RL





- RL requires a large amount of samples. (data inefficient)
- It's hard to design proper reward function for each particular task.
- However, it is usually esay to obtain expert-level demonstrations

Imitation Learning





- problem setting
 - Learner is given samples of trajectories from the expert
 - Learner can interact with envrioment, but can't get reinforcement signal of any kind, which can be modeled as MDP/R
 - objective: $\max_{\pi} V^{\pi} \Leftrightarrow \min_{\pi} [V^{\pi_E} V^{\pi}]$
- Application
 - Warmup for RL
 - When it's hard to define reward function

(e.g. AlphaGo)
(e.g. Self-driving cars)

Methods of IL

- Behavior Cloning (BC)
- Inverse Reinforcement Learning (IRL)
- Adversarial Structured Imitation Leaning





Behavior Cloning



mimic by matching actions distribution with supervised learning method

$$\mathcal{D} = \{ au_1, au_2, ...\}, au = \{(s_1, a_1), (s_2, a_2), ...\}$$



Algorithm 1 Abstract of behavioral cloning

Collect a set of trajectories demonstrated by the expert \mathcal{D} Select a policy representation π_{θ} Select an objective function \mathcal{L} Optimize \mathcal{L} w.r.t. the policy parameter θ using \mathcal{D} return optimized policy parameters θ

Behavior Cloning



• In the simplest case, BC learn a policy to minimize the KL divergence

$$\min_{\pi} \mathbb{E}_{s \sim d_{\pi_E}} \left[D_{KL}(\pi_E(\cdot|s), \pi(\cdot|s))
ight] := \mathbb{E}_{(s,a) \sim
ho_{\pi_E}} \left[\log \! \left(rac{\pi_E(a|s)}{\pi(a|s)}
ight)
ight]$$

• In practice, we optimize the objective with finite samples

 $\max_{\pi \in \Pi} \sum_{(s,a) \in \tau_E} \log(\pi(a|s)).$

- For discrete action spaces, it reduces to learn a classifier
- For continuous action spaces, it reduces to learn a regressor

Behavior Cloning



- some papers
 - Alvinn: An autonomous land vehicle in a neural network. 1989
 - **DAVE-2**: Bojarski M , Testa D D , Dworakowski D , et al. *End to End Learning for Self-Driving Cars*. 2016.
 - **DAgger**: Stéphane Ross, Geoffrey J. Gordon, and Drew Bagnell. A reduction of *imitation learning and structured prediction to no-regret online learning*. In Proceedings of the 14th InternationalConference on Artificial Intelligence and Statistics (AISTATS'11), pages 627–635, 2011.
 - **DAgger by coaching:** He He, Hal Daumé, III, and Jason Eisner. *Imitation learning by coaching*. In Proceedings of the 25th International Conference on Neural Information Processing Systems Volume 2, NIPS?12, pages 3149?3157, USA, 2012. Curran Associates Inc.
 - AggreVaT: Stéphane Ross and J. Andrew Bagnell. *Reinforcement and imitation learning via interactive no-regret learning*. CoRR, abs/1406.5979, 2014.

ALVINN: An Autonomous Land Vehicle In a Neural Network

Dean A. Pomerleau January 1989 CMU-CS-89-107-



- use a 3-layer back-propagation network
- Train the network by having it observe live sensor data as a human drives the





• Loss

- fit a Gaussian to the network's output vector.
- Measure distance between Gaussian's peak and human steering direction.





- Learning to Correct Steering Errors
 - If the human drives perfectly, the network never learns to make corrections when it drifts off the desired track
 - Crude solution
 - Turn learning off temporarily, and drive off course.
 - Turn learning back on, and let the network observe the human making the necessary corrections.
 - Repeat





- Simulating the Steering Errors
 - Let humans drive as best they can.
 - Increase training set variety by artificially shifting and rotating the video images, so that the vehicle appears at different orientations relative to the road.
 - Generate 14 random shift/rotations for each image.
 - A simple steering model is used to predict how a human driver would react to each transformation





- After a long right turn, the network will be biased toward turning right, since recent training data focused on right turns
- Balanced Training Images: Keep a buffer of 200 training images. Replace 15 old exemplars with new ones derived from the current camera image. Replacement strategies:
 - Replace the image with the lowest error
 - Replace the image with the closest steering direction
- Online training details
 - 1. Take current camera image plus 14 shifted/rotated variants, each with computed steering direction.
 - 2. Replace 15 old exemplars in the 200 element training exemplar buffer with these 15 new ones.
 - 3. Perform one epoch of backpropagation learning on the training exemplar buffer.
 - 4. Repeat steps 1-3 until the network's predicted steering direction reliably matches the person's steering direction.

End to End Learning for Self-Driving Cars

Mariusz Bojarski NVIDIA Corporation Holmdel, NJ 07735 Davide Del Testa NVIDIA Corporation Holmdel, NJ 07735 Daniel Dworakowski NVIDIA Corporation Holmdel, NJ 07735 Bernhard Firner NVIDIA Corporation Holmdel, NJ 07735

Mathew Monfort

NVIDIA Corporation

Holmdel, NJ 07735

Beat Flepp NVIDIA Corporation Holmdel, NJ 07735 Prasoon Goyal NVIDIA Corporation Holmdel, NJ 07735 Lawrence D. Jackel NVIDIA Corporation Holmdel, NJ 07735

Urs Muller NVIDIA Corporation Holmdel, NJ 07735 Jiakai Zhang NVIDIA Corporation Holmdel, NJ 07735 Xin Zhang NVIDIA Corporation Holmdel, NJ 07735 Jake Zhao NVIDIA Corporation Holmdel, NJ 07735

Karol Zieba NVIDIA Corporation Holmdel, NJ 07735



• Overview of the DAVE-2 System



Figure 1: High-level view of the data collection system.

- We represent the steering command as 1/r to make our system independent of the car geometry, where r is the turning radius in meters.
- We use 1/r instead of r to prevent a singularity when driving straight (the turning radius for driving straight is infinity).



• Training data was collected by driving on a wide variety of roads and in a diverse set of lighting and weather conditions



- Training with data from only the human driver is not sufficient. The network must learn how to recover from mistakes.
 - The left and the right camera —> Images for two specific off-center shifts
 - viewpoint transformation of the image from the nearest camera —> Additional shifts between the cameras
 - The steering label for transformed images is adjusted to one that would steer the vehicle back to the desired location and orientation in two seconds.



Network Architecture



Output: vehicle control

Fully-connected layer Fully-connected layer Fully-connected layer

Convolutional feature map 64@1x18

Convolutional feature map 64@3x20

Convolutional feature map 48@5x22

Convolutional feature map 36@14x47

Convolutional feature map 24@31x98

Normalized input planes 3@66x200

Input planes 3@66x200

as a controller for steering

feature extraction

hard-coded and not adjusted in the learning process

•



• Training with three cameras







• Simulation & on-road test



 When the off-center distance exceeds one meter, a virtual human intervention is triggered to do interventions. We estimate what percentage of the time the network could drive the car (autonomy)

autonomy =
$$(1 - \frac{(\text{number of interventions}) \cdot 6 \text{ seconds}}{\text{elapsed time [seconds]}}) \cdot 100$$

 After a trained network has demonstrated good performance in the simulator, the network is loaded on the test car and taken out for a road test. For a typical drive, we are autonomous approximately 98% of the time.

Problem of BC



- There are three common problems with supervised imitation learning
 - 1. copies unnecessay action (SL take all errors equally)
 - 2. compounding error because of fit single-timestep decisions
 - 3. mismatch/covriate shift (not satisfy the i.i.d assumption)





$$P(s,a) = P(a|s)P(s)$$

A Reduction of Imitation Learning and Structured Prediction to No-Regret Online Learning

Stéphane Ross

Robotics Institute Carnegie Mellon University Pittsburgh, PA 15213, USA stephaneross@cmu.edu Geoffrey J. Gordon Machine Learning Department Carnegie Mellon University Pittsburgh, PA 15213, USA ggordon@cs.cmu.edu

J. Andrew Bagnell

Robotics Institute Carnegie Mellon University Pittsburgh, PA 15213, USA dbagnell@ri.cmu.edu



- DAgger attempts to collect expert demonstrations under the state distribution induced by the learned policy.
- It can be seen most naturally as an on-policy to imitation learning: the expert provides the correct actions to take, but the input distribution of examples comes from the learner's own behavior.





• The general DAGGER algorithm.



- We typically use $\beta_1 = 1, \beta_i = p^{i-1} (i \ge 2)$ or $\beta_i = I(i = 1)$, optionally allowing the algorithm queries the expert to choose controls for a fraction of the time to better leverage the presence of the expert
- in general, $\{eta_i\}$ be a sequence such that $areta_N=rac{1}{N}\sum_{i=1}^Neta_i o 0$ as N o 0



- This algorithm can be interpreted as a *Follow-The-Leader* algorithm that at each iteration we pick the best policy under all trajectories seen so far over the iterations.
- As a special case of FTL, DAgger enjoys the property of being no-regret

$$\frac{1}{N} \sum_{i=1}^{N} \ell_i(\pi_i) - \min_{\pi \in \Pi} \frac{1}{N} \sum_{i=1}^{N} \ell_i(\pi) \le \gamma_N$$
$$\lim_{N \to \infty} \gamma_N = 0$$



Sup



human expert training laps are all very similar expert is a planning algorithm which hacked the game

8

9

10 $x 10^{4}$

Sm0.1



- $J(\pi)$ the expected total reward of trajectories starting with the initial state I
- $d_{\pi} = \frac{1}{T} \sum_{t=1}^{T} d_{\pi}^{t}$ the empirical mean of state distribution induced over each time step
- $C_{\pi}(s) = \mathbb{E}_{a \sim \pi(s)}[R(s, a)]$ the total reward in a T-step trajectory
- $l(s,\pi)$ the observed surrogate loss
- Assuming l(s, π) is the 0-1 loss (or upper bound on the 0-1 loss) implies the following performance guarantee with respect to any task reward function C bounded in [0, 1]:

Theorem 1 Let denote $\epsilon = \mathbb{E}_{s \sim d_{\pi^*}}[l(s, \pi, \pi^*)]$, then there exists $\pi \in \pi_{1:N}$ such that $J(\pi) \leq J(\pi^*) + T^2 \epsilon$

Theorem 5 Let $\epsilon_N = \min_{\pi \in \Pi} \frac{1}{N} \sum_{i=1}^N \mathbb{E}_{s \sim d_{\pi_i}}[l(s,\pi)]$ be the true loss of the best policy, then if $N = O(\frac{T}{\epsilon})$ there exists $\pi \in \pi_{1:N}$ and u such as $J(\pi) \leq J(\pi^*) + uT\epsilon_N + O(1)$

Attia A, Dayan S. Global overview of Imitation Learning. 2018. arXiv preprint arXiv:arXiv:1801.06503

Summary of BC



- Advantages
 - 1. easy to implement
 - 2. efficient and stable to train
- Disadvantage
 - 1. distribution mistmatch
 - 2. compounding error

- When to use
 - 1. 1-step deviation not too bad
 - 2. expert episodes cover state space
- When not to use
 - 1. 1-step deviation lead to catastrophic error
 - 2. optimze long-term objective

Summary of BC



- other problem setting ?
 - Using ambiguous demonstration data to imitate learning, thus reducing the cost of collecting demonstration
 - ...



Inverse RL



• The idea is to learn the optimal reward function that can explain the expert's behavior most appropriately



Inverse RL



- However, since a policy can be optimal for multiple reward functions, the problem of determining the reward function is "ill-posed"
- To obtain the unique solution, many studies have proposed additional objective functions to be optimized
 - maximize the margin
 - maximize the entropy
 -

Inverse RL



- some papers
 - ALVIL: P. Abbeel and A. Y. Ng. *Apprenticeship learning via inverse reinforcement learning*. In Proceedings of the international conference on Machine learning (ICML), 2004.
 - MaxEnt: B. D. Ziebart, A. Maas, J. A. Bagnell, and A. K. Dey. *Maximum entropy inverse reinforcement learning*. In Proceedings of the Twenty-Second Con_x0002_ference on Artificial Intelligence (AAAI), pages 1433–1438, 2008.
 - **MaxCausalEnt**: B. D. Ziebart, J. A. Bagnell, and A. K. Dey. *Modeling interaction via the principle of maximum causal entropy*. in Proc. of International Conference on Machine Learning, Haifa, Israel, 2010.
 - N. D. Ratliff, J. A. Bagnell, and M. A. Zinkevich. *Maximum margin planning*. In Proceedings of the international conference on Machine learning (ICML), pages 729–736, 2006b.
 - **Guided cost learning:** Chelsea Finn, Sergey Levine, and Pieter Abbeel. 2016. *Guided cost learning: deep inverse optimal control via policy optimization*. In Proceedings of the 33rd International Conference on International Conference on Machine Learning - Volume 48 (ICML'16). JMLR.org, 49–58.
 - Boularias, A., Kober, J., and Peters, J. *Relative entropy inverse reinforcement learning*. In International Conference on Artificial Intelligence and Statistics (AISTATS), 2011.

Apprenticeship Learning via Inverse Reinforcement Learning

Pieter Abbeel Andrew Y. Ng Computer Science Department, Stanford University, Stanford, CA 94305, USA

PABBEEL@CS.STANFORD.EDU ANG@CS.STANFORD.EDU



- finite state MDP (S, A, P, γ, D, R) , the reward function assume to be bounded in absolute value by 1
- features mapping function: $\phi: \mathcal{S}
 ightarrow [0,1]^k$
 - feature vector: $\mathbf{f}_s = \phi(s)$
 - feature expectations: $\mu(\pi) = \mathbb{E}_{\pi,s_0 \sim D}[\sum_{t=0}^\infty \gamma^t \phi(s_t)] \in \mathbb{R}^k$
- assume there is some "true" reward function: $R^*(S) = w^{*\top}\phi(S), \ w^* \in \mathbb{R}^k$, In order to ensure that the rewards are bounded by 1, we also assume $\|w^*\|_1 \leq 1$
 - $R(s) = \boldsymbol{w}^{ op} \phi(s)$ $V^{\pi} = \mathbb{E}_{\pi, s_0 \sim D} [V^{\pi}(s_0 \otimes D)]$

•
$$G_{\tau} = \sum_{t=0}^{\infty} \gamma^t R(s_t) = \sum_{t=0}^{\infty} \gamma^t \boldsymbol{w}^{\top} \phi(s_t)$$

• $V^{\pi}(s) = \mathbb{E}_{\pi,s_0=s}[G_{\tau}]$

$$egin{aligned} &\pi = \mathbb{E}_{\pi,s_0 \sim D} \left[V^{\pi}(s_0)
ight] \ &= \mathbb{E}_{\pi,s_0 \sim D} \left[\sum_{t=0}^\infty \gamma^t R(s_t)
ight] \ &= oldsymbol{w}^ op \mathbb{E}_{\pi,s_0 \sim D} \left[\sum_{t=0}^\infty \gamma^t \phi(s_t)
ight] \ &= oldsymbol{w}^ op \mu(\pi) \end{aligned}$$



- if we have found some set of policies $\pi_1, \pi_2, ..., \pi_n$, we can mix them to get a new policy
 - At the start of a trajectory, choose one policy with probability γ_i , then we always acts according to the selected policy until trajectory end
 - the feature expectations of the mixed policy is

$$\sum_{i=1}^{n} \lambda_{i} \mu(\pi_{i}) \ (\lambda_{i} \ge 0, \sum_{i} \lambda_{i} = 1)$$

• Use demonstration data $\{s_0^{(i)}, s_1^{(i)}, \ldots\}_{i=1}^m$ to estimate the expert's feature expectations $\mu_E = \mu(\pi_E)$

$$\hat{\mu}_E = \frac{1}{m} \sum_{i=1}^m \sum_{t=0}^\infty \gamma^t \phi(s_t^{(i)})$$



37/5

to find a policy whose performance is close to expert's on the unknown reward function, just need to find a policy π̃ such that ||μ(π̃) − μ_E||₂ ≤ ε

$$\begin{split} |E[\sum_{t=0}^{\infty} \gamma^{t} R(s_{t}) | \pi_{E}] - E[\sum_{t=0}^{\infty} \gamma^{t} R(s_{t}) | \tilde{\pi}]| &\longleftarrow \text{gap of policy value} \\ &= |w^{T} \mu(\tilde{\pi}) - w^{T} \mu_{E}| & | \wedge \\ &\leq ||w||_{2} ||\mu(\tilde{\pi}) - \mu_{E}||_{2} & | \wedge \\ &\leq 1 \cdot \epsilon = \epsilon &\longleftarrow \text{gap of feature expection} \end{split}$$

 $egin{aligned} & oldsymbol{x}^ opoldsymbol{y}| = |oldsymbol{x}||_2|oldsymbol{y}||_2\cos(heta) \leq |oldsymbol{x}||_2|oldsymbol{y}||_2 \ & |oldsymbol{w}||_2 \leq |oldsymbol{w}||_1 \leq 1 \end{aligned}$



- 1. Randomly pick some policy $\pi^{(0)}$, compute (or approximate via Monte Carlo) $\mu^{(0)} = \mu(\pi^{(0)})$, and set i = 1.
- 2. Compute $t^{(i)} = \max_{w:||w||_2 \leq 1} \min_{j \in \{0., (i-1)\}} w^T (\mu_E \mu^{(j)})$, and let $w^{(i)}$ be the value of w that attains this maximum.
- 3. If $t^{(i)} \leq \epsilon$, then terminate.
- 4. Using the RL algorithm, compute the optimal policy $\pi^{(i)}$ for the MDP using rewards $R = (w^{(i)})^T \phi$.
- 5. Compute (or estimate) $\mu^{(i)} = \mu(\pi^{(i)})$.
- 6. Set i = i + 1, and go back to step 2.

t



$$egin{aligned} &(i) = \max_{w: ||w||_2 \leq 1} \min_{j \in \{0, 1, 2, ..., (i-1)\}} m{w}^ op (\mu_E - \mu^{(i)}) \end{aligned}$$

- min: fix reward function as $w^{\top}\phi(s)$, find a policy (by RL) to minimize the gap of feature expection
- max: fix policy, find parameter w^{\top} (i.e. reward function) to maximize te gap of policy value

 $\begin{array}{ll} \max_{t,w} & t \leftarrow \text{margin of policy value} \\ \text{s.t.} & w^T \mu_E \ge w^T \mu^{(j)} + t, \ j = 0, \dots, i-1 \\ ||w||_2 \le 1 \end{array}$

 The key idea is to match the feature expectations (FEM), and then use additional objective functions (max margin here) to deal with the "ill-posed" problem

Maximum Entropy Inverse Reinforcement Learning

Brian D. Ziebart, Andrew Maas, J.Andrew Bagnell, and Anind K. Dey

School of Computer Science Carnegie Mellon University Pittsburgh, PA 15213

bziebart@cs.cmu.edu, amaas@andrew.cmu.edu, dbagnell@ri.cmu.edu, anind@cs.cmu.edu





• Objective: maximize the the entropy of trajectory distribution under the condition of feature matching

Problem Formulation			
$arg max_p$	H	I(p)	(entropy)
subject to	$\mathbb{E}_{\pi^{\mathbb{E}}}\left[oldsymbol{\phi}(au) ight] = \mathbb{E}_{\pi^{\mathbb{L}}}\left[oldsymbol{\phi}(au) ight],$		(feature matching)
	$\sum_{\tau\in\mathcal{T}}p(\tau)=1,$	$\forall \tau \in \mathcal{T} : p(\tau) > 0$	(prob. distr.)

 choose the distribution that does not exhibit any additional preferences beyond matching feature expectations

MaxEnt



• for Deterministic MDP

$$P(au| heta) = rac{1}{Z(heta)} exp(heta^ op \mathbf{f}_ au) = rac{1}{Z(heta)} exp(heta^ op \sum_{s_j \in au} \mathbf{f}_{s_j})$$

• for Non-Deterministic MDP, paths are determined by the action choices of the agent and the random outcomes of the MDP

$$P(\tau|\theta,T) = \sum_{o \in O} P_T(o) \frac{exp(\theta^{\top} \mathbf{f}_{\tau})}{Z(\theta,o)} I_{\tau \in o}$$

$$\approx \frac{1}{Z(\theta,T)} exp(\theta^{\top} \mathbf{f}_{\tau}) \prod_{s_{t+1},a_t,s_t \in \tau} P_T(s_{t+1}|a_t,s_t)$$

transition distribution

deterministic transition

• Stochastic Policy $P(a|\theta,T) \propto \sum_{\tau:a \in \tau_{t=1}} P(\tau|\theta,T)$

MaxEnt



• Maximizing the entropy of the distribution over paths subject to the feature constraints from observed data implies that we maximize the likelihood of the observed data under the maximum entropy (exponential family) distribution derived above

$$heta^* = arg \max_{\theta} L(\theta) = arg \max_{\theta} \sum_{examples \; \tilde{\tau}} log P(\tilde{\tau}|\theta, T)$$

$$\nabla L(\theta) = \tilde{\mathbf{f}} - \sum_{\tau} P(\tau|\theta, T) \mathbf{f}_{\tau} = \tilde{\mathbf{f}} - \sum_{s_i} D_{s_i} \mathbf{f}_{s_i}$$

state visitation frequencies

45/57

MaxCausalEnt

 MaxCausalEnt: B. D. Ziebart, J. A. Bagnell, and A. K. Dey. *Modeling interaction* via the principle of maximum causal entropy. in Proc. of International Conference on Machine Learning, Haifa, Israel, 2010.

$$P(\mathbf{A}^{T}||\mathbf{S}^{T}) \triangleq \prod_{t=1}^{T} P(A_{t}|\mathbf{S}_{1:t}, \mathbf{A}_{1:t-1}).$$

$$P(\mathbf{A}|\mathbf{S}) = \prod_{t=1}^{T} P(A_{t}|\mathbf{S}_{1:T}, \mathbf{A}_{1:t-1}).$$

$$H(\mathbf{A}^{T}||\mathbf{S}^{T}) \triangleq E_{\mathbf{A},\mathbf{S}}[-\log P(\mathbf{A}^{T}||\mathbf{S}^{T})]$$

$$= \sum_{t=1}^{T} H(A_{t}|\mathbf{S}_{1:t}, \mathbf{A}_{1:t-1}),$$

 $\begin{aligned} \operatorname*{argmax}_{\{P(A_t|\mathbf{S}_{1:t},\mathbf{A}_{1:t-1})\}} H(\mathbf{A}^T || \mathbf{S}^T) \\ \mathrm{such that:} \ E_{\mathbf{S},\mathbf{A}}[\mathcal{F}(\mathbf{S},\mathbf{A})] = \tilde{E}_{\mathbf{S},\mathbf{A}}[\mathcal{F}(\mathbf{S},\mathbf{A})] \\ \mathrm{and} \ \forall_{\mathbf{S}_{1:t},\mathbf{A}_{1:t-1}} \sum_{A_t} P(A_t|\mathbf{S}_{1:t},\mathbf{A}_{1:t-1}) = 1, \\ \mathrm{and given:} \ P(\mathbf{S}^T || \mathbf{A}^{T-1}). \end{aligned}$



Summary of IRL



- advantages:
 - 1. does not need interactive expert
 - 2. very efficient when trained (in some cases can outperform the demonstrator)
 - 3. has long-term planning
- disadvantages:
 - 1. can be difficult to train
- use when:
 - 1. an interactive expert is not available
 - 2. it might be easier to learn the reward functions than the expert's policy

Generative Adversarial Imitation Learning

Jonathan Ho Stanford University hoj@cs.stanford.edu Stefano Ermon Stanford University ermon@cs.stanford.edu



- IRL methods aim to recover the expert's cost function and then extract a policy from that cost function with reinforcement learning. This approach is indirect and can be slow
- We propose a new general framework for directly extracting a policy from data, as if it were obtained by reinforcement learning following inverse reinforcement learning
- IRL:

$$\underset{c \in \mathcal{C}}{\operatorname{maximize}} \left(\underset{\pi \in \Pi}{\min} - H(\pi) + \mathbb{E}_{\pi}[c(s,a)] \right) - \mathbb{E}_{\pi_{E}}[c(s,a)]$$
$$\operatorname{RL}(c) = \underset{\pi \in \Pi}{\operatorname{arg\,min}} - H(\pi) + \mathbb{E}_{\pi}[c(s,a)]$$





- IRL methods aim to recover the expert's cost function and then extract a policy from that cost function with reinforcement learning. This approach is indirect and can be slow
- We propose a new general framework for directly extracting a policy from data, as if it were obtained by reinforcement learning following inverse reinforcement learning
- general form of IRL:

$$\underset{c \in \mathcal{C}}{\operatorname{maximize}} \left(\underset{\pi \in \Pi}{\min} - H(\pi) + \mathbb{E}_{\pi}[c(s,a)] \right) - \mathbb{E}_{\pi_{E}}[c(s,a)]$$

$$\underset{\pi \in \Pi}{\operatorname{RL}(c)} = \underset{\pi \in \Pi}{\operatorname{arg\,min}} - H(\pi) + \mathbb{E}_{\pi}[c(s,a)]$$

$$\underset{\pi \in \Pi}{\operatorname{E}_{\pi}[c(s,a)]} \triangleq \mathbb{E}\left[\sum_{t=0}^{\infty} \gamma^{t} c(s_{t},a_{t})\right]$$

 $H(\pi) \triangleq \mathbb{E}_{\pi}[-\log \pi(a|s)]$ is the γ -discounted causal entropy



• in case of overfit, we incorporate a (closed, proper) convex cost function regularizer $\psi : \mathbb{R}^{S \times A} \to \overline{\mathbb{R}}$ into our study, it must be convex as a function defined on all of $\mathbb{R}^{S \times A}$

$$\operatorname{IRL}_{\psi}(\pi_{E}) = \underset{c \in \mathbb{R}^{S \times \mathcal{A}}}{\operatorname{arg\,max}} - \psi(c) + \left(\underset{\pi \in \Pi}{\min} - H(\pi) + \mathbb{E}_{\pi}[c(s,a)] \right) - \mathbb{E}_{\pi_{E}}[c(s,a)]$$
$$\tilde{c} \in \operatorname{IRL}_{\psi}(\pi_{E}).$$

one-to-one correspondence between policy set and occupancy measures set

$$\rho_{\pi}(s,a) = \pi(a|s) \sum_{t=0}^{\infty} \gamma^{t} P(s_{t} = s|\pi)$$

$$\mathcal{D} \triangleq \{\rho_{\pi} : \pi \in \Pi\}$$

$$= \{\rho: \rho \ge 0 \text{ and } \sum_{a} \rho(s,a) = p_{0}(s) + \gamma \sum_{s',a} P(s|s',a)\rho(s',a) \ \forall s \in \mathcal{S} \}.$$

distribution of starting states

Proposition 3.1 (Theorem 2 of Syed et al. [29]). If $\rho \in D$, then ρ is the occupancy measure for $\pi_{\rho}(a|s) \triangleq \rho(s,a) / \sum_{a'} \rho(s,a')$, and π_{ρ} is the only policy whose occupancy measure is ρ .



• for a function $f : \mathbb{R}^{S \times A} \to \overline{\mathbb{R}}$, its convex conjugate is $f^*(x) = \sup_{y \in \mathbb{R}^{S \times A}} x^T y - f(y)$

Proposition 3.2. RL \circ IRL $_{\psi}(\pi_E) = \arg \min_{\pi \in \Pi} -H(\pi) + \psi^*(\rho_{\pi} - \rho_{\pi_E})$

Lemma 3.1. Let $\overline{H}(\rho) = -\sum_{s,a} \rho(s,a) \log(\rho(s,a) / \sum_{a'} \rho(s,a'))$. Then, \overline{H} is strictly concave, and for all $\pi \in \Pi$ and $\rho \in D$, we have $H(\pi) = \overline{H}(\rho_{\pi})$ and $\overline{H}(\rho) = H(\pi_{\rho})$.

Proof of Lemma 3.7 First, we show strict concavity of \overline{H} . Let ρ and ρ' be occupancy measures, and suppose $\lambda \in [0, 1]$. For all s and a, the log-sum inequality [6] implies:

$$-(\lambda\rho(s,a) + (1-\lambda)\rho'(s,a))\log\frac{\lambda\rho(s,a) + (1-\lambda)\rho'(s,a)}{\sum_{a'}(\lambda\rho(s,a') + (1-\lambda)\rho'(s,a'))}$$
(19)

$$= -(\underbrace{\lambda\rho(s,a)}_{a} + \underbrace{(1-\lambda)\rho'(s,a)}_{a})\log \frac{\lambda\rho(s,a) + (1-\lambda)\rho'(s,a)}{\underbrace{\lambda\sum_{a'}\rho(s,a')}_{a'} + \underbrace{(1-\lambda)\sum_{a'}\rho'(s,a')}_{a'}}$$
(20)

$$\sum a_i \log \frac{a_i}{b_i} \ge \sum a_i \log \frac{\sum a_i}{\sum b_i} \ge -\lambda \rho(s, a) \log \frac{\lambda \rho(s, a)}{\lambda \sum_{a'} \rho(s, a')} - (1 - \lambda) \rho'(s, a) \log \frac{(1 - \lambda) \rho'(s, a)}{(1 - \lambda) \sum_{a'} \rho'(s, a')} \mathsf{D}^2 \tag{21}$$

$$= \lambda \left(-\rho(s,a) \log \frac{\rho(s,a)}{\sum_{a'} \rho(s,a')} \right) + (1-\lambda) \left(-\rho'(s,a) \log \frac{\rho'(s,a)}{\sum_{a'} \rho'(s,a')} \right),$$
(22)

Summing both sides over all s and a shows that

$$\bar{H}(\lambda\rho + (1-\lambda)\rho') \ge \lambda\bar{H}(\rho) + (1-\lambda)\bar{H}(\rho')$$



Lemma 3.1. Let $\bar{H}(\rho) = -\sum_{s,a} \rho(s,a) \log(\rho(s,a) / \sum_{a'} \rho(s,a'))$. Then, \bar{H} is strictly concave, and for all $\pi \in \Pi$ and $\rho \in D$, we have $H(\pi) = \bar{H}(\rho_{\pi})$ and $\bar{H}(\rho) = H(\pi_{\rho})$.

$$H(\pi) = \mathbb{E}_{\pi}[-\log \pi(a|s)]$$

= $-\sum_{s,a} \rho_{\pi}(s,a) \log \pi(a|s)$
= $-\sum_{s,a} \rho_{\pi}(s,a) \log \frac{\rho_{\pi}(s,a)}{\sum_{a'} \rho_{\pi}(s,a')}$
= $\bar{H}(\rho_{\pi}),$

$$\bar{H}(\rho) = -\sum_{s,a} \rho(s,a) \log \frac{\rho(s,a)}{\sum_{a'} \rho(s,a')}$$
$$= -\sum_{s,a} \rho_{\pi_{\rho}}(s,a) \log \pi_{\rho}(a|s)$$
$$= \mathbb{E}_{\pi_{\rho}}[-\log \pi_{\rho}(a|s)]$$
$$= H(\pi_{\rho}).$$



Proposition 3.2.
$$\operatorname{RL} \circ \operatorname{IRL}_{\psi}(\pi_{E}) = \operatorname{arg\,min}_{\pi \in \Pi} -H(\pi) + \psi^{*}(\rho_{\pi} - \rho_{\pi_{E}})$$

Let $\tilde{c} \in \operatorname{IRL}_{\psi}(\pi_{E}), \tilde{\pi} \in \operatorname{RL}(\tilde{c}) = \operatorname{RL} \circ \operatorname{IRL}_{\psi}(\pi_{E}), \text{ and } f^{*}(x) = \sup_{y \in \mathbb{R}^{S \times A}} x^{T}y - f(y)$
 $\pi_{A} \in \operatorname{arg\,min}_{\pi} -H(\pi) + \psi^{*}(\rho_{\pi} - \rho_{\pi_{E}})$

$$= \operatorname{arg\,min\,max}_{\pi} -H(\pi) - \psi(c) + \sum_{s,a} (\rho_{\pi}(s,a) - \rho_{\pi_{E}}(s,a))c(s,a) \quad (32)$$

We wish to show that $\pi_A = \tilde{\pi}$. To do this, let ρ_A be the occupancy measure of π_A , let $\tilde{\rho}$ be the occupancy measure of $\tilde{\pi}$, and define $\bar{L} : \mathcal{D} \times \mathbb{R}^{S \times A} \to \mathbb{R}$ by

$$\bar{L}(\rho, c) = -\bar{H}(\rho) - \psi(c) + \sum_{s,a} \rho(s, a)c(s, a) - \sum_{s,a} \rho_{\pi_E}(s, a)c(s, a).$$
(33)

The following relationships then hold, due to Proposition 3.1

$$\rho_A \in \arg\min_{\rho \in \mathcal{D}} \max_{c} \bar{L}(\rho, c), \tag{34}$$

$$\tilde{c} \in \arg \max_{c} \min_{\rho \in \mathcal{D}} \bar{L}(\rho, c),$$
(35)

$$\tilde{\rho} \in \underset{\rho \in \mathcal{D}}{\operatorname{arg\,min}} \bar{L}(\rho, \tilde{c}). \tag{36}$$

$$\operatorname{IRL}_{\psi}(\pi_{E}) = \underset{c \in \mathbb{R}^{S \times \mathcal{A}}}{\operatorname{argmax}} - \psi(c) + \left(\underset{\pi \in \Pi}{\min} - H(\pi) + \mathbb{E}_{\pi}[c(s,a)] \right) - \mathbb{E}_{\pi_{E}}[c(s,a)]$$
$$\mathbb{E}_{\pi}[c(s,a)] = \sum_{s,a} \rho_{\pi}(s,a)c(s,a) \quad \min_{\rho} \max_{c} \bar{L}(\rho,c) = \max_{c} \min_{\rho} \bar{L}(\rho,c) \quad 53/5$$



- Proposition 3.2 tells us that ψ -regularized inverse reinforcement learning, implicitly, seeks a policy whose occupancy measure is close to the expert's, as measured by the convex function ψ^*
- conclusion
 - IRL is a dual of an occupancy measure matching problem
 - The induced optimal policy is the primal optimum
- Enticingly, this suggests that various settings of ψ lead to various imitation learning algorithms that directly solve the optimization problem given by Proposition 3.2
 - Corollary 3.2.1. If ψ is a constant function, $\tilde{c} \in \operatorname{IRL}_{\psi}(\pi_E)$, and $\tilde{\pi} \in \operatorname{RL}(\tilde{c})$, then $\rho_{\tilde{\pi}} = \rho_{\pi_E}$.
 - apprenticeship learning:

 $\begin{aligned} & \underset{\pi}{\mininimize} \max_{c \in \mathcal{C}} \mathbb{E}_{\pi}[c(s, a)] - \mathbb{E}_{\pi_{E}}[c(s, a)] \\ \delta_{\mathcal{C}} : \mathbb{R}^{\mathcal{S} \times \mathcal{A}} \to \overline{\mathbb{R}}, \text{ defined by } \delta_{\mathcal{C}}(c) = 0 \text{ if } c \in \mathcal{C} \text{ and } +\infty \text{ otherwise,} \\ & \underset{c \in \mathcal{C}}{\max} \mathbb{E}_{\pi}[c(s, a)] - \mathbb{E}_{\pi_{E}}[c(s, a)] = \underset{c \in \mathbb{R}^{\mathcal{S} \times \mathcal{A}}}{\max} - \delta_{\mathcal{C}}(c) + \sum_{s, a} (\rho_{\pi}(s, a) - \rho_{\pi_{E}}(s, a))c(s, a) = \delta_{\mathcal{C}}^{*}(\rho_{\pi} - \rho_{\pi_{E}}) \\ & \underset{\pi}{\mininimize} - H(\pi) + \underset{c \in \mathcal{C}}{\max} \mathbb{E}_{\pi}[c(s, a)] - \mathbb{E}_{\pi_{E}}[c(s, a)] \end{aligned}$



• Generative adversarial imitation learning

$$\psi_{\text{GA}}(c) \triangleq \begin{cases} \mathbb{E}_{\pi_E}[g(c(s,a))] & \text{if } c < 0\\ +\infty & \text{otherwise} \end{cases} \text{ where } g(x) = \begin{cases} -x - \log(1 - e^x) & \text{if } x < 0\\ +\infty & \text{otherwise} \end{cases}$$

$$\psi_{\text{GA}}^{*}(\rho_{\pi} - \rho_{\pi_{E}}) = \max_{D \in (0,1)^{S \times \mathcal{A}}} \mathbb{E}_{\pi}[\log(D(s,a))] + \mathbb{E}_{\pi_{E}}[\log(1 - D(s,a))]$$

$$\underset{\pi}{\text{minimize }} \psi_{\text{GA}}^*(\rho_{\pi} - \rho_{\pi_E}) - \lambda H(\pi) = D_{\text{JS}}(\rho_{\pi}, \rho_{\pi_E}) - \lambda H(\pi)$$



Algorithm 1 Generative adversarial imitation learning

- 1: Input: Expert trajectories $\tau_E \sim \pi_E$, initial policy and discriminator parameters θ_0, w_0
- 2: for $i = 0, 1, 2, \dots$ do
- 3: Sample trajectories $\tau_i \sim \pi_{\theta_i}$
- 4: Update the discriminator parameters from w_i to w_{i+1} with the gradient

$$\hat{\mathbb{E}}_{\tau_i}[\nabla_w \log(D_w(s,a))] + \hat{\mathbb{E}}_{\tau_E}[\nabla_w \log(1 - D_w(s,a))]$$
(17)

5: Take a policy step from θ_i to θ_{i+1} , using the TRPO rule with cost function $\log(D_{w_{i+1}}(s, a))$. Specifically, take a KL-constrained natural gradient step with

$$\hat{\mathbb{E}}_{\tau_i} \left[\nabla_\theta \log \pi_\theta(a|s) Q(s,a) \right] - \lambda \nabla_\theta H(\pi_\theta),$$
where $Q(\bar{s}, \bar{a}) = \hat{\mathbb{E}}_{\tau_i} [\log(D_{w_{i+1}}(s,a)) \mid s_0 = \bar{s}, a_0 = \bar{a}]$
(18)

6: end for





RL platform



https://github.com/wxc971231/RL_Platform-CrossTheWall

