What makes modern machine learning work?









What about reinforcement learning?



Mnih et al. '13



Schulman et al. '14 & '15



Levine*, Finn*, et al. '16





enormous gulf







3

Can we develop data-driven RL methods?



Levine, Kumar, Tucker, Fu. Offline Reinforcement Learning: Tutorial, Review, and Perspectives on Open Problems. '20

What does offline RL mean?

off-policy RL

 π_k

 \mathbf{s}, r

а

rollout(s)

rollout data $\{(s_i, a_i, s'_i, r_i)\}$

 π_{k+1}

 $\mathcal{D}^{\mathsf{buffer}}$

update

 π_{k+1}



offline reinforcement learning



Formally:

$$\mathcal{D} = \{(\mathbf{s}_i, \mathbf{a}_i, \mathbf{s}'_i, r_i)\}$$

$$\mathbf{s} \sim d^{\pi_{\beta}}(\mathbf{s})$$

$$\mathbf{a} \sim \pi_{\beta}(\mathbf{a}|\mathbf{s})$$

$$\mathbf{s}' \sim p(\mathbf{s}'|\mathbf{s}, \mathbf{a})$$

$$r \leftarrow r(\mathbf{s}, \mathbf{a})$$

RL objective:
$$\max_{\pi} \sum_{t=0}^{T} E_{\mathbf{s}_{t} \sim d^{\pi}(\mathbf{s}), \mathbf{a}_{t} \sim \pi(\mathbf{a}|\mathbf{s})} [\gamma^{t} r(\mathbf{s}_{t}, \mathbf{a}_{t})]$$

How is this even possible?

1. Find the "good stuff" in a dataset full of good and bad behaviors

2. Generalization: good behavior in one place may suggest good behavior in another place

3. "Stitching": parts of good behaviors can be recombined





Why should we care?











In some settings there exist very good decision policies and we would like to automate them

- One idea: humans provide reward signal when RL algorithm makes decisions
- Good: simple, cheap form of supervision
- Bad: High sample complexity

Alternative: imitation learning

Rewards that are **dense in time** closely guide the agent. How can we supply these rewards?

- Manually design them: often brittle
- Implicitly specify them through demonstrations



Emma Brunskill (CS234 Reinforcement Learn

- Expert provides a set of demonstration trajectories: sequences of states and actions
- Imitation learning is useful when it is easier for the expert to demonstrate the desired behavior rather than:
 - Specifying a reward that would generate such behavior,
 - Specifying the desired policy directly

• Input:

- State space, action space
- Transition model $P(s' \mid s, a)$
- No reward function R
- Set of one or more teacher's demonstrations $(s_0, a_0, s_1, s_0, ...)$ (actions drawn from teacher's policy π^*)
- Behavioral Cloning:
 - Can we directly learn the teacher's policy using supervised learning?
- Inverse RL:
 - Can we recover R?
- Apprenticeship learning via Inverse RL:
 - Can we use R to generate a good policy?

- Formulate problem as a standard machine learning problem:
 - Fix a policy class (e.g. neural network, decision tree, etc.)
 - Estimate a policy from training examples $(s_0, a_0), (s_1, a_1), (s_2, a_2), \ldots$
- Two notable success stories:
 - Pomerleau, NIPS 1989: ALVINN
 - Summut et al., ICML 1992: Learning to fly in flight simulator

ALVINN



30x32 Video Input Retina

< □ > < □ > < □ > < □ > < □ > < □ >

- Often behavior cloning in practice can work very well, especially if use BCRNN
- See What Matters in Learning from Offline Human Demonstrations for Robot Manipulation. Mandlekar et al. CORL 2021

Potential Problem with Behavior Cloning: Compounding Errors

Supervised learning assumes iid. (s, a) pairs and ignores temporal structure Independent in time errors:

Error at time t with probability $\leq \epsilon \in \mathbb{E}[\text{Total errors}] \leq \epsilon T$

Problem: Compounding Errors



Data distribution mismatch!

In supervised learning, $(x, y) \sim D$ during train and test. In MDPs:

- Train: $s_t \sim D_{\pi^*}$
- Test: $s_t \sim D_{\pi_{\theta}}$

A Reduction of Imitation Learning and Structured Prediction to No-Regret Online Learning, Ross et al. 2011

Emma Brunskill (CS234 Reinforcement Learn

Problem: Compounding Errors



- Error at time t with probability ϵ
- Approximate intuition: $\mathbb{E}[\text{Total errors}] \le \epsilon(T + (T 1) + (T 2) \dots + 1) \propto \epsilon T^2$
- Real result requires more formality. See Theorem 2.1 in http://www.cs.cmu.edu/~sross1/publications/ Ross-AIStats10-paper.pdf with proof in supplement: http://www.cs.cmu.edu/~sross1/publications/ Ross-AIStats10-sup.pdf

Emma Brunskill (CS234 Reinforcement Learn

A Reduction of Imitation Learning and Structured Prediction to No-Regret Online Learning, Ross et al. 2011

DAGGER: Dataset Aggregation

```
Initialize \mathcal{D} \leftarrow \emptyset.

Initialize \hat{\pi}_1 to any policy in \Pi.

for i = 1 to N do

Let \pi_i = \beta_i \pi^* + (1 - \beta_i) \hat{\pi}_i.

Sample T-step trajectories using \pi_i.

Get dataset \mathcal{D}_i = \{(s, \pi^*(s))\} of visited states by \pi_i

and actions given by expert.

Aggregate datasets: \mathcal{D} \leftarrow \mathcal{D} \bigcup \mathcal{D}_i.

Train classifier \hat{\pi}_{i+1} on \mathcal{D}.

end for

Return best \hat{\pi}_i on validation.
```

- Idea: Get more labels of the expert action along the path taken by the policy computed by behavior cloning
- Obtains a stationary deterministic policy with good performance under its induced state distribution
- Key limitation?

- Given state space, action space, transition model P(s' | s, a)
- No reward function R
- Set of one or more expert's demonstrations (s₀, a₀, s₁, s₀,...) (actions drawn from teacher's policy π^{*})
- Goal: infer the reward function R
- Assume that the teacher's policy is optimal. What can be inferred about *R*?

- Recall linear value function approximation
- Similarly, here consider when reward is linear over features

•
$$R(s) = oldsymbol{w}^{ op} x(s)$$
 where $oldsymbol{w} \in \mathbb{R}^n, x: S o \mathbb{R}^n$

- Goal: identify the weight vector \boldsymbol{w} given a set of demonstrations
- The resulting value function for a policy π can be expressed as

$$V^{\pi}(s_0) = \mathbb{E}_{s \sim \pi} \left[\sum_{t=0}^{\infty} \gamma^t R(s_t) | s_0 \right]$$

- Recall linear value function approximation
- Similarly, here consider when reward is linear over features

•
$$R(s) = oldsymbol{w}^{ op} x(s)$$
 where $w \in \mathbb{R}^n, x: S o \mathbb{R}^n$

- Goal: identify the weight vector \boldsymbol{w} given a set of demonstrations
- The resulting value function for a policy π can be expressed as

$$V^{\pi}(s_0) = \mathbb{E}_{s \sim \pi} \left[\sum_{t=0}^{\infty} \gamma^t R(s_t) \mid s_0 \right] = \mathbb{E}_{s \sim \pi} \left[\sum_{t=0}^{\infty} \gamma^t \boldsymbol{w}^T x(s_t) \mid s_0 \right]$$
$$= \boldsymbol{w}^T \mathbb{E}_{s \sim \pi} \left[\sum_{t=0}^{\infty} \gamma^t x(s_t) \mid s_0 \right]$$
$$= \boldsymbol{w}^T \mu(\pi)$$

 where μ(π)(s) is defined as the discounted weighted frequency of state features under policy π, starting in state s₀.

Relating Frequencies to Optimality

- Assume $R(s) = \boldsymbol{w}^T x(s)$ where $w \in \mathbb{R}^n, x : S \to \mathbb{R}^n$
- Goal: identify the weight vector \boldsymbol{w} given a set of demonstrations
- $V^{\pi} = \mathbb{E}_{s \sim \pi} [\sum_{t=0}^{\infty} \gamma^t R^*(s_t) \mid \pi] = \mathbf{w}^T \mu(\pi)$ where $\mu(\pi)(s) =$ discounted weighted frequency of state *s* under policy π .

$$V^* \ge V^{\pi}$$

Relating Frequencies to Optimality

- Recall linear value function approximation
- Similarly, here consider when reward is linear over features

•
$$R(s) = oldsymbol{w}^{ op} x(s)$$
 where $w \in \mathbb{R}^n, x: S o \mathbb{R}^n$

- Goal: identify the weight vector **w** given a set of demonstrations
- The resulting value function for a policy π can be expressed as

$$V^{\pi} = \boldsymbol{w}^{T} \mu(\pi)$$

• $\mu(\pi)(s) = \text{discounted weighted frequency of state } s$ under policy π .

$$\mathbb{E}_{s \sim \pi^*} [\sum_{t=0}^{\infty} \gamma^t R^*(s_t) \mid \pi^*] = V^* \ge V^{\pi} = \mathbb{E}_{s \sim \pi} [\sum_{t=0}^{\infty} \gamma^t R^*(s_t) \mid \pi] \quad \forall \pi$$

 Therefore if the expert's demonstrations are from the optimal policy, to identify *w* it is sufficient to find *w*^{*} such that

$$w^{*T}\mu(\pi^*) \geq w^{*T}\mu(\pi), \forall \pi \neq \pi^*$$

- Want to find a reward function such that the expert policy outperforms other policies.
- For a policy π to be guaranteed to perform as well as the expert policy π*, sufficient if its discounted summed feature expectations match the expert's policy [Abbeel & Ng, 2004].
- More precisely, if

$$\|\mu(\pi) - \mu(\pi^*)\|_1 \le \epsilon$$

then for all w with $||w||_{\infty} \leq 1$:

$$|w^T \mu(\pi) - w^T \mu(\pi^*)| \leq \epsilon$$

- There is an infinite number of reward functions with the same optimal policy.
- There are infinitely many stochastic policies that can match feature counts
- Which one should be chosen?

- Many different approaches
- Two of the key papers are:
 - Maximumum Entropy Inverse Reinforcement Learning (Ziebart et al. AAAI 2008)
 - Generative adversarial imitation learning (Ho and Ermon, NeurIPS 2016)

- Again assume a linear reward function $R(s) = \boldsymbol{w}^T x(s)$
- Define the total feature counts for a single trajectory τ_j as: $\mu_{\tau_j} = \sum_{s_i \in \tau_j} x(s_i)$
 - Note that this is a slightly different definition that we saw earlier
- The average feature counts over *m* trajectories is: $\tilde{\mu} = \frac{1}{m} \sum_{i=1}^{m} \mu_{\tau_i}$

Deterministic MDP Path Distributions



- Consider all possible H-step trajectories in a deterministic MDP
- For a linear reward model, a policy is completely specified by its distribution over trajectories
- Which policy/distribution should we choose given a set of *m* demonstrations?

• Principle of max entropy: choose distribution with no additional preferences beyond matching the feature expectations in the demonstration dataset

$$\max_{P} - \sum_{\tau} P(\tau) \log P(\tau) s.t. \quad \sum_{\tau} P(\tau) \mu_{\tau} = \tilde{\mu} \qquad \sum_{\tau} P(\tau) = 1$$
(1)

In the linear reward case, this is equivalent to specifying the weights
 w that yield a policy with the max entropy constrained to matching
 the feature expectations

Ziebart et al., 2008

Emma Brunskill (CS234 Reinforcement Learn

Max Entropy Principle

 Maximizing the entropy of the distribution over the paths subject to the feature constraints from observed data implies we maximize the likelihood of the observed data under the maximum entropy (exponential family) distribution¹.

$$P(\tau_j \mid w) = \frac{1}{Z(w)} \exp\left(w^T \mu_{\tau_j}\right) = \frac{1}{Z(w)} \exp\left(\sum_{s_i \in \tau_j} w^T x(s_i)\right)$$
$$Z(w, s) = \sum_{\tau_s} \exp\left(w^T \mu_{\tau_s}\right)$$

• Strong preference for low cost paths, equal cost paths are equally probable.

¹Jaynes 1957

Emma Brunskill	(CS234	Reinforcement	Learr
----------------	--------	---------------	-------

- Many MDPs of interest are stochastic
- For these the distribution over paths depends both on the reward weights and on the stochastic dynamics

$$P(\tau_j \mid w, P(s' \mid s, a)) \approx \frac{\exp\left(w^T \mu_{\tau_j}\right)}{Z(w, P(s' \mid s, a))} \prod_{s_i, a_i \in \tau_j} P(s_{i+1} \mid s_i, a_i)$$

• Select w to maximize likelihood of data:

$$w^* = rg\max_w L(w) = rg\max_w \sum_{ ext{examples}} \log P(\tau \mid w)$$

• The gradient is the difference between expected empirical feature counts and the learner's expected feature counts, which can be expressed in terms of expected state visitation frequencies

$$abla L(w) = ilde{\mu} - \sum_{\tau} P(\tau \mid w) \mu_{ au} = ilde{\mu} - \sum_{s_i} D(s_i) x(s_i)$$

- where $D(s_i)$: state visitation frequency
- Do we need to know the transition model to compute the above?

MaxEnt IRL Algorithm

Backward pass

- 1. Set $Z_{s_i,0} = 1$
- 2. Recursively compute for N iterations

$$\begin{split} Z_{a_{i,j}} &= \sum_{k} P(s_k | s_i, a_{i,j}) e^{\text{reward}(s_i | \theta)} Z_{s_k} \\ Z_{s_i} &= \sum_{a_{i,j}} Z_{a_{i,j}} \end{split}$$

Local action probability computation

3.
$$P(a_{i,j}|s_i) = \frac{Z_{a_{i,j}}}{Z_{s_i}}$$

Forward pass

- 4. Set $D_{s_i,t} = P(s_i = s_{\text{initial}})$
- 5. Recursively compute for t = 1 to N

$$D_{s_{i},t+1} = \sum_{a_{i,j}} \sum_{k} D_{s_{k},t} P(a_{i,j}|s_{i}) P(s_{k}|a_{i,j},s_{i})$$

Summing frequencies

$$6. \ D_{s_i} = \sum_t D_{s_{i,t}}$$

→ ∃ →

1

- Max entropy approach has been hugely influential
- Provides a principled way for selecting among the (many) possible reward functions
- The original formulation requires knowledge of the transition model or the ability to simulate/act in the world to gather samples of the transition model
 - Check your understanding: was this needed in behavioral cloning?

- Inverse RL approaches provide a way to learn a reward function
- Generally interested in using this reward function to compute a policy whose performance equals or exceeds the expert policy
- One approach: given learned reward function, use with regular RL
- Can we more directly learn the desired policy?

- Imitation learning can greatly reduce the amount of data need to learn a good policy
- Challenges remain and one exciting area is combining inverse RL / learning from demonstration and online reinforcement learning
- For a look into some of the theory between imitation learning and RL, see Sun, Venkatraman, Gordon, Boots, Bagnell (ICML 2017)