Generalization in RL

Sham M. Kakade (and Wen Sun)

- 2 Today: SL vs. RL
- Supervised Learning (SL) : Let's review

RL and generalization

- Is Agnostic Learning Possible?
- Lower bounds



The need for strategic exploration



- agent starts at s₀
- Iength of chain is H
- chance of hitting goal state in H steps is $(1/3)^H$ with a random policy

UCBVI: Optimistic Model-based Learning

Inside iteration *n* :

Use all previous data to estimate transitions $\widehat{P}_{1}^{n}, ..., \widehat{P}_{H-1}^{n}$

Design reward bonus $b_h^n(s, a), \forall s, a, h$

Optimistic planning with learned model: $\pi^n = \text{Value-Iter}\left(\left\{ \widehat{P}_h^n, r_h + b_h^n \right\}_{h=1}^{H-1} \right)$

Collect a new trajectory by executing π^n in the real world $\{P_h\}_{h=0}^{H-1}$ starting from s_0

UCBVI: Put All Together

For
$$n = 1 \rightarrow N$$
:
1. Set $N_h^n(s, a) = \sum_{i=1}^{n-1} \mathbf{1}\{(s_h^i, a_h^i) = (s, a)\}, \forall s, a, h$
2. Set $N_h^n(s, a, s') = \sum_{i=1}^{n-1} \mathbf{1}\{(s_h^i, a_h^i, s_{h+1}^i) = (s, a, s')\}, \forall s, a, a', h$
3. Estimate $\widehat{P}^n : \widehat{P}_h^n(s'|s, a) = \frac{N_h^n(s, a, s')}{N_h^n(s, a)}, \forall s, a, s', h$
4. Plan: $\pi^n = VI\left(\{\widehat{P}_h^n, r_h + b_h^n\}_h\right)$, with $b_h^n(s, a) = cH\sqrt{\frac{\ln(SAHN/\delta)}{N_h^n(s, a)}}$
5. Execute $\pi^n : \{s_0^n, a_0^n, r_0^n, \dots, s_{H-1}^n, a_{H-1}^n, r_{H-1}^n, s_H^n\}$

Theorem: UCBVI Regret Bound

$$\mathbb{E}\left[\mathsf{Regret}_{N}\right] := \mathbb{E}\left[\sum_{n=1}^{N} \left(V^{\star} - V^{\pi^{n}}\right)\right] \leq \widetilde{O}\left(H^{2}\sqrt{S^{2}AN}\right)$$

Remarks:

Note that we consider expected regret here (policy π^n is a random quantity). High probability version is not hard to get (need to do a martingale argument)

Dependency on H and S are suboptimal; but the **same** algorithm can achieve $H^2 \sqrt{SAN}$ in the leading term [Azar et.al 17 ICML]



Supervised Learning (SL) : Let's review

4 RL and generalization

- Is Agnostic Learning Possible?
- Lower bounds

Maze example: r = -1 per time-step and policy





[David Silver. Advanced Topics: RL]

What we want to solve: (the large state space case)



TD GAMMON [Tesauro 95]



[AlphaZero, Silver et.al, 17]



[OpenAl Five, 18]

To what extent is generalization in RL similar to (or different from) that in supervised learning?

- Up to now, we have focussed on "tabular" MDPs (theoretically important)
 - We ultimately seek learnability results where number of states is large (or $|\mathcal{S}| = \infty$).
 - This is a question of generalization.
- Supervised Learning: two lines of thinking
 - Optimal learning: try to learn the Bayes optimal classifier. need very strong assumptions.
 - Agnostic learning: try to do as well best classifier in some (restricted) class $\mathcal{H}.$
- If rather than trying to be 'optimal' in RL, does trying to do agnostic learning make our task easier?

2 Today: SL vs. RL

Supervised Learning (SL) : Let's review

RL and generalization

- Is Agnostic Learning Possible?
- Lower bounds

Binary Classification

• *n* labeled examples: $(x_i, y_i)_{i=1}^n$, with $x_i \in \mathcal{X}$ and $y_i \in \{0, 1\}$. A set \mathcal{H} of binary classifiers, where for $h \in \mathcal{H}$, $h : \mathcal{X} \to \{0, 1\}$. Define the empirical error and the true error as:

$$\widehat{\operatorname{err}}(h) = \frac{1}{N} \sum_{i=1}^{N} \mathbf{1}(h(x_i) \neq y_i), \quad \operatorname{err}(h) = \mathbb{E}_{(X,Y) \sim D} \mathbf{1}(h(X) \neq Y).$$

where $\mathbf{1}(h(x) \neq y)$ is 0 if h(x) = y and 1 otherwise.

If the samples are drawn i.i.d. according to a joint distribution *D* over (*x*, *y*), then, by Hoeffding's inequality, for a fixed *h* ∈ *H*, with probability at least 1 − δ:

$$|\operatorname{err}(h) - \widehat{\operatorname{err}}(h)| \leq \sqrt{\frac{1}{2N}\log\frac{2}{\delta}}.$$

- Binary classification is special case of RL. Consider learning in an MDP, with two actions where the effective horizon is 1.
- $|\mathcal{A}| = 2$, $\gamma = 0$, and the reward function is $r(s, a) = \mathbf{1}(\text{label}(s) = a)$.
- Note in SL, we rarely make restrictions that \mathcal{X} (i.e. \mathcal{S}) is finite.
- Note that $\mu(s_0) \leftrightarrow D(x)$ (*D* is the distribution of our data)

Your HW0: This and the union bound give rise to what is often referred to as the "Occam's razor" bound:

Proposition

(The "Occam's razor" bound) Suppose \mathcal{H} is finite. Let $\widehat{h} = \arg \min_{h \in \mathcal{H}} \widehat{err}(h)$ and $h^* = \arg \min_{h \in \mathcal{H}} err(h)$. With probability at least $1 - \delta$:

$$\textit{err}(\widehat{h}) - \textit{err}(h^{\star}) \leq \sqrt{rac{2}{N} \log rac{2|\mathcal{H}|}{\delta}}$$

(The logarithmic dependence is the most naive complexity measure of $\ensuremath{\mathcal{H}}$, yet the bound is strong.)

- |*H*|: a set of Boolean functions on *X*.
 Even though |*H*| may be infinite, the number of possible behaviors of on a finite set of states is not necessarily exhaustive.
- We say that the set {*x*₁, *x*₂, . . . *x_d*} is shattered if there exists an *h* ∈ *H* that can realize any of the possible 2^{*d*} labellings.
- The Vapnik–Chervonenkis (VC) dimension is the size of the largest shattered set.
- Let $d = VC(\mathcal{H})$. The Sauer–Shelah lemma: the number of possible labellings on a set of *n* points by functions in \mathcal{H} is at most $\left(\frac{en}{d}\right)^d$. For $d \ll m$, this is much less than 2^n .

Review: Half Spaces



- Let \mathcal{H}_{half} be the set of halspaces on $\mathcal{X} = \mathbb{R}^d$.
- The VC dimension is $VC(\mathcal{H}_{half}) = d + 1$

The following classical bound highlights how generalization is possible on infinite hypothesis classes with bounded complexity.

Proposition

(VC dimension and generalization) Let $\hat{h} = \arg \min_{h \in \mathcal{H}} \sum_{i=1}^{n} \mathbf{1}(h(x_i) \neq y_i)$ and $h^* = \arg \min_{h \in \mathcal{H}} err(h)$. Suppose \mathcal{H} has a bounded VC dimension. For $m \geq VC(\mathcal{H})$, we have that with probability at least $1 - \delta$:

$$\textit{err}(\widehat{h}) - \textit{err}(h^{\star}) \leq \sqrt{rac{c}{n} \left(\textit{VC}(\mathcal{H}) \log rac{2n}{\textit{VC}(\mathcal{H})} + \log rac{2}{\delta}
ight)},$$

where c is an absolute constant

- 2 Today: SL vs. RL
- 3 Supervised Learning (SL) : Let's review
 - RL and generalization
 - Is Agnostic Learning Possible?
 - Lower bounds

RL and Agnostic Learning

- We have a set of policies Π (either finite or infinite).
 - Π could be a parametric set.
 - Π could be greedy policys on a a set of parametric value functions
 V = {f_θ : S × A → ℝ| θ ∈ ℝ^d}.
 - Π may not contain π^* .
- in agnostic learning, we have the optimization problem:

(

$$\max_{\pi\in\Pi}\mathbb{E}_{\boldsymbol{s}_0\sim\mu}\boldsymbol{V}^{\pi}(\boldsymbol{s}_0)$$

We want to (approx) solve this with a small number of sample trajectories.

- analogous to agnostic learning in SL
 - binary classification: |A| = 2, $\gamma = 0$, $r(\cdot)$ being the labeling reward.
 - relevant dependencies for RL:

Complexity(
$$\Pi$$
), $|S|$, $|A|$, N

RL Sampling Model

- Assume sampling access to the MDP in a μ -reset model:
 - start at a state $s_0 \sim \mu$
 - we can rollout a policy π of our choosing
 - we can terminate the trajectory at will.

(weaker model than generative model)

Lemma

(Effective Horizon and Truncation) We have that:

$$|V^{\pi}(s_0) - \mathbb{E}_{\pi}\left[\sum_{t=0}^{H} \gamma^t r(s_t, a_t) \mid s_0\right]| \leq \gamma^H / (1 - \gamma),$$

For $H = \frac{\log (1/(\epsilon(1-\gamma)))}{1-\gamma}$ we will have an ϵ approximation to $V^{\pi}(s_0)$.

Lemma

(Near unbiased estimation of $V^{\pi}(s_0)$) Let π_{uar} denote the policy which chooses actions uniformly at random at every state. We have that:

$$\begin{aligned} |\mathcal{A}|^{H} \mathbb{E}_{\pi_{uar}} \left[\mathbf{1} \left(\pi(s_{0}) = a_{0}, \ldots, \pi(s_{H}) = a_{H} \right) \sum_{t=0}^{H} \gamma^{t} r(s_{t}, a_{t}) \right] \\ = \mathbb{E}_{\pi} \left[\sum_{t=0}^{H} \gamma^{t} r(s_{t}, a_{t}) \right]. \end{aligned}$$

- the estimated the reward of π on a trajectory is nonzero only when π takes precisely the same actions as the π_{uar} on the trajectory then estimated reward of $|\mathcal{A}|^{H}$ is equal to that of π_{uar} .
- the factor of |A|^H which is due to this being a high variance estimate.
 We will return to this point in the next section.

Denote the *n*-th sample by $(s_0^n, a_0^n, r_1^n, s_1^n, \dots, s_H^n)$, where *H* is a cutoff time where the trajectory ends. A nearly, unbiased estimate of the γ -discounted reward of a given policy π is given by:

$$\widehat{V}^{\pi}(s_0) = \frac{|\mathcal{A}|^H}{N} \sum_{n=1}^N \mathbf{1}\Big(\pi(s_0^n) = a_0^n, \dots, \pi(s_H^n) = a_H^n\Big) \sum_{t=0}^H \gamma^t r(s_t^n, a_t^n).$$

Proposition

(Generalization in RL) Suppose Π is finite. Let $\hat{\pi} = \arg \max_{\pi \in \Pi} \hat{V}^{\pi}(s_0)$ and $\pi^* = \arg \max_{\pi \in \Pi} V^{\pi}(s_0)$. With probability at least $1 - \delta$:

$$V^{\widehat{\pi}}(s_0) \geq rg\max_{\pi\in\Pi} V^{\pi}(s_0) - rac{\epsilon}{2} - |\mathcal{A}|^H \sqrt{rac{2}{N}\lograc{2|\Pi|}{\delta}}.$$

A VC Theorem for RL

- Suppose $|\mathcal{A}| = 2$. Each $\pi \in \Pi$ can be viewed as Boolean function.
- VC(Π) is well defined.
- For *N* trajectories, the Sauer–Shelah lemma bounds the number of possible labellings on a set of *N* trajectories (of length *H*) by $\left(\frac{eNH}{d}\right)^d$, where $d = VC(\Pi)$.
- this leads to the following proposition:

Proposition

(A VC Theorem for RL) Suppose $|\mathcal{A}| = 2$. Let $\widehat{\pi} = \arg \max_{\pi \in \Pi} \widehat{V}^{\pi}(s_0)$ and $\pi^* = \arg \max_{\pi \in \Pi} V^{\pi}(s_0)$. With probability at least $1 - \delta$:

$$V^{\widehat{\pi}}(s_0) \geq rg\max_{\pi\in\Pi} V^{\pi}(s_0) - 2^{\mathcal{H}} \; \sqrt{rac{c}{n} \left(\mathit{VC}(\Pi) \log rac{2n}{\mathit{VC}(\Pi)} + \log rac{2}{\delta}
ight)},$$

where c is an absolute constant.

- 2 Today: SL vs. RL
- Supervised Learning (SL) : Let's review

RL and generalization

- Is Agnostic Learning Possible?
- Lower bounds

- What we want, for our agnostic sample complexity:
 - no dependence on $|\mathcal{S}|$ (or logarithmic)
 - poly H dependence
 - to depend reasonably on a complexity measure of H e.g. poly log |H| dependence
- Is this possible?
 No :(
- This is why RL is hard!
 - it is hard in practice...
 - how should we study it?

Proposition

(Lower Bound for The Complete Policy Class) Suppose $|\mathcal{A}| = 2$ and $|\mathcal{S}| = 2^{H}$, where $H = \lfloor \frac{\log(2)}{1-\gamma} \rfloor$. Let Π be the set of all 2^{H} policies. There exists a family of MDPs such that if a deterministic algorithm \mathcal{A} is guaranteed to find a policy π such that:

$$V^{\widehat{\pi}}(\mathit{s}_{0}) \geq rg\max_{\pi \in \Pi} V^{\pi}(\mathit{s}_{0}) - 1/4.$$

then A must use $N \ge 2^H$ trajectories.

Observe that $\log |\Pi| = H \log(2)$, so this already rules out the possibility of logarithmic dependence on the size of the policy class, without having an exponential dependence on *H*.

A General Lower Bound

before: we had a complete policy class. practice: Π is a restricted (and smaller) class.

Proposition

(Lower Bound for an Arbitrary Policy Class) Suppose |A| = 2 and Π is an arbitrary policy class. There exists a family of MDPs s.t. any algorithm A that is guaranteed to find a policy $\hat{\pi}$ s.t.:

$$\mathbb{E}\left[V^{\widehat{\pi}}(\textbf{\textit{s}}_{0})
ight] \geq rg\max_{\pi\in\Pi}V^{\pi}(\textbf{\textit{s}}_{0})-\epsilon.$$

then \mathcal{A} must use an expected number of trajectories N where

$$N \ge c \frac{\min\{2^H, 2^{VC(\Pi)}\}}{\epsilon^2}$$

in the worst case, we (nearly) have to do exhaustive search (trying $2^{VC(\Pi)}$ polices, which is the effective the number of policies in Π)

- 2 Today: SL vs. RL
- 3 Supervised Learning (SL) : Let's review
- 4 RL and generalization
 - Is Agnostic Learning Possible?
 - Lower bounds

- The tabular case:
 - we can understand fundamental issues of exploration vs exploitation
 - we can't get at generalization.
- need stronger assumptions/side info.
- this course (and the field) take the following approaches:
 - Structural (and Modelling) Assumptions: By making stronger assumptions about the world, we can move away from agnostic learning and escape the curse of dimensionality. We will see examples of this in Part 2.
 - Distribution Dependent Results (and Distribution Shift): When we move to policy gradient methods (in Part 3), we will consider results which depend on given distribution of how we obtain samples. Here, we will make connections to transfer learning.
 - Imitation learning and behavior cloning: here will consider models where the agent has input from, effectively, a teacher, and we will see how this alleviates the problem of curse of dimensionality.