

08/24/22 Wednesday.

Reinforcement Learning: Mathematical formalism for learning-based decision making.

Approach for learning decision making and control from experience.

→ Data Not i.i.d.

→ No ground truth.

Machine Learning: given $D = \{(x_i, y_i)\}$

learn to predict y from x $f(x) = y$.

Decisions: Action.

Consequence: Observation (State). Reward.

Why DEEP RL? able to Adapt. / handle unstructured environments \Rightarrow formalism for behavior.

↓
End-to-End: manual feature extraction \Rightarrow automatically learned features

(perception) (action)

Recognition + Control Separately \Rightarrow sensorimotor loop.

Learning reward functions from examples. (inverse reinforcement learning).

transfer knowledge between domains (Meta-learning, transfer learning).

learn to predict, use prediction to act. (predictive models).

imitation learning.

inferring intentions.

08/29/22 Monday.

s_t - state underlying true configuration.

o_t - observation.

a_t - action.

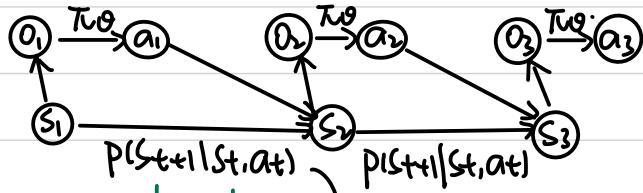
State results observation.

observation may or may not be enough to deduce state.

distribution of a_t given o_t .

$\pi_{\theta}(a_t | o_t)$ - policy.

$\pi(a_t | s_t)$ - policy (fully observed).



\Rightarrow Markov property.

Not for $P_{s+1}(s_{t+1} | o_t, a_t)$.

(past observations may give additional info. e.g. car).

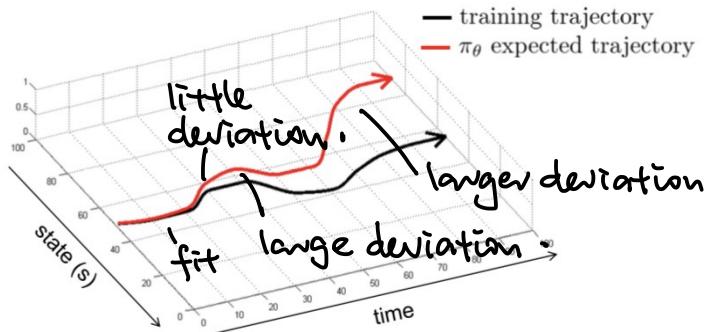
Imitation Learning.

behavior cloning. supervised learning.

o_t - left+center/right camera. \rightarrow map.

a_t - rotations by steering wheel.

left camera is supervised with steering angle that is a little bit to the right.



Distribution shift problem.

distribution of training data.

Make $P_{\text{data}}(o_t) = P_{\text{val}}(o_t)$? $T_{\pi \theta} \Rightarrow$ different o_t .

Make $T_{\pi \theta}$ perfect. Hard. \Rightarrow change data.

DAgger: Data Aggregation.

collect training data from $P_{\text{val}}(o_t)$ instead of $P_{\text{data}}(o_t)$.

- 1) train $T_{\text{ulat}}(t+1|t)$ from human data $D = \{O_1, a_1, \dots, O_N, a_N\}$.
- 2) run $T_{\text{ulat}}(t+1|t)$ to get dataset $D_{\pi} = \{O_1, \dots, O_M\}$.
- 3) Ask human to label D_{π} with actions a_t .
- 4) Aggregate: $D \leftarrow D \cup D_{\pi}$.

D_{π} will converge to the same distribution as D

Avoid distribution drift: mimic expert behavior very accurately, but don't overfit. 

reasons of failure:

① Non-Markovian behavior.

$T_{\text{ulat}}(t+1|t)$ behavior depends only on current observation.

→ unnatural for human beings.

if see the same thing twice, do the same thing twice, regardless of what happened before.

Reality: $T_{\text{ulat}}(O_1, O_2, \dots, O_t)$. Concatenate → input with various length, too large. X.

RNN: LSTM ✓.

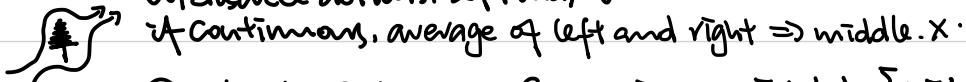
Causal Confusion: if brake is caused by the brake light or pedestrian.

possible solutions: include the history?

Augment?

② Multimodal behavior.

w discrete actions, softmax ✓.



① output mixture of Gaussians. $T_{\text{ulat}}(o) = \sum_i w_i N(\mu_i, \Sigma_i)$.

②: output parameters +++; high dimension data ⇒ challenge.

② Latent variable models.

prior distribution?

input: image + latent variable (noise). output: one gaussian.

③ Autoregressive discretization. discrete continuous ND data ⇒ exponential in n.

discretize one dimension at a time. ⇒ sample from Softmax ⇒ value 1. ⇒ feed into another NN. ⇒ (repeat...).

④ of imitation learning: Human need to provide data (only finite data; hard in some cases)

Human: unlimited data from own experience; gets self-improvement.

Cost / Reward Function

$\rightarrow \begin{cases} 0 & \text{Not eaten} \\ 1 & \text{eaten.} \end{cases}$

$$\min_{\theta} \mathbb{E}_{(s,a) \sim D_{\text{train}}, s' \sim p(s'|s,a)} [\delta'(s' = \text{eaten by tiger})].$$

$$= \min_{\theta} \mathbb{E}_{s \in T, a \in T} \left[\sum_t \delta(s_t = \text{eaten by tiger}) \right].$$

$C(s,a)$ - Cost function. $\max r(s_t, a_t)$.

$$C(s,a) = \begin{cases} 0 & \text{if } a = \pi^*(s) \\ 1 & \text{otherwise.} \end{cases}$$

Assume $t \in D_{\text{train}}$. $\Pr(a \neq \pi^*(s)|s) \leq \varepsilon$.

$$E \left[\sum_t C(s_t, a_t) \right] \leq \underbrace{\sum_t \cdot + (1-\varepsilon)(\varepsilon(T-1) + (1-\varepsilon)(\dots))}_{O(\varepsilon T^2)} \xrightarrow{\text{fall off} \Rightarrow \text{Training error}} \text{Training error, each } O(\varepsilon T).$$

$1-\varepsilon$ is large.

$\Pr(\text{mistake}) \leq \varepsilon \Rightarrow$ Supervised learning works.

Some probability distribution.

Assume $\Pr_{\text{train}}(s) \Pr(a \neq \pi^*(s)|s) \leq \varepsilon$. (enough if $E_{\text{train}}(s) \Pr(a \neq \pi^*(s)|s) \leq \varepsilon$.)

W/ DAgger, $\Pr_{\text{train}}(s) \rightarrow p_\theta(s)$ $E \left[\sum_t C(s_t, a_t) \right] \leq \varepsilon T$.

if $\Pr_{\text{train}}(s) \neq p_\theta(s)$, $p_\theta(s_t) = (1-\varepsilon)^t \Pr_{\text{train}}(s_t) + (-(1-\varepsilon)^t) \Pr(\text{mistake})(s_t)$.

$$|p_\theta(s_t) - \Pr_{\text{train}}(s_t)| = (1-(1-\varepsilon)^t) |\Pr(\text{mistake})(s_t) - \Pr_{\text{train}}(s_t)| \stackrel{\max(p)=1}{\leq} 2(1-(1-\varepsilon)^t)$$

$$\leq 2\varepsilon + ((1-\varepsilon)^t - 1 - \varepsilon t \text{ if } \varepsilon \in [0,1]).$$

$$\sum_t E_{p_\theta(s_t)}[C_t] = \sum_t \sum_s p_\theta(s_t) C_t(s_t) \leq \sum_t \sum_s \Pr_{\text{train}}(s_t) C_t(s_t) + |p_\theta(s_t) - \Pr_{\text{train}}(s_t)| \text{ max.}$$

$$\leq \sum_t (\varepsilon + 2\varepsilon t) \leq \varepsilon T + 2\varepsilon T^2 \quad O(\varepsilon T^2).$$

Goal-Conditioned behavioral cloning.

$\overrightarrow{p_i}$ $T_{\text{demos}}(s)$ policy for reaching p_i .

$\overrightarrow{p_i} \rightarrow p_i$ T_{demos, p_i} policy for reaching any p_i .

demo: $\{s_1, a_1, \dots, s_{T-1}, a_{T-1}, s_T\}$. $\xrightarrow{\text{successfully reach } s_T} \max \log \Pr(a_t | s_t, g=s_T)$.

Need demonstration? \Rightarrow start with random policy. \Rightarrow random data \Rightarrow label goals.

09/07/22 Wednesday.

reinforcement learning.

reward function. $r(s,a)$.

Markov chain. $M = \{S, T\}$.

S-state space.

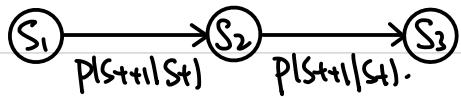
States $s \in S$. (discrete or continuous states).

T-transition operator.

\downarrow
matrix.

$P(s_{t+1}|s_t)$ probability.

let $\mu_{t,i} = P(s_t=i)$ $\bar{\mu}_t$ is a vector of probability.



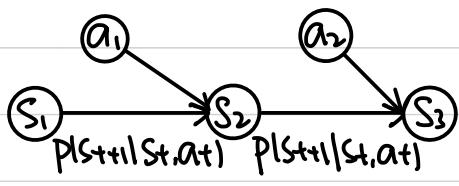
let $T_{i,j} = P(s_{t+1}=i|s_t=j)$, then $\bar{\mu}_{t+1} = T\bar{\mu}_t$.

Markov decision process. $M = \{S, A, T, r\}$.

S-state Space. States $s \in S$. (discrete or continuous states).

A-Action space. actions $a \in A$ (discrete or continuous actions).

T-transition operator. \Rightarrow tensor.³⁰



let $\mu_{t,j} = P(s_t=j)$. $\Sigma_{t,k} = P(a_t=k)$.

$T_{i,j,k} = P(s_{t+1}=i|s_t=j, a_t=k)$.

$\bar{\mu}_{t+1,i} = \sum_{j,k} T_{i,j,k} \mu_{t,j} \Sigma_{t,k}$.

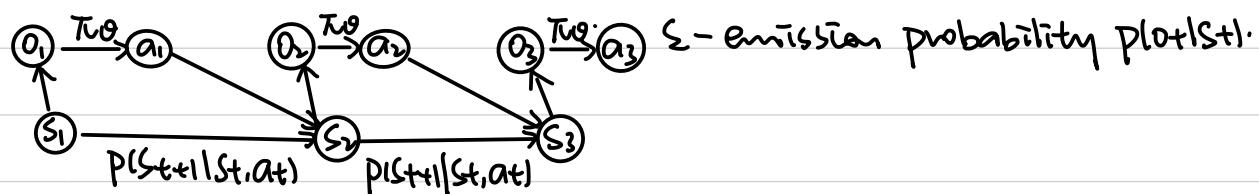
r-reward function. $r: S \times A \rightarrow \mathbb{R}$.

$r(s_t, a_t)$ - reward.

partially observed Markov decision process. $M = \{S, A, O, T, \Sigma, r\}$.

O-observation Space.

observation $\in O$ (discrete or continuous observation).



goal of reinforcement learning.

$$P(s_1, a_1, \dots, s_T, a_T) = P(s_1) \prod_{t=1}^T P(o_t|s_t) P(s_{t+1}|s_t, a_t).$$

$P(\tau)$

$$P((s_{t+1}, a_{t+1}) | (s_t, a_t)) = P(s_{t+1}|s_t, a_t) P(o_{t+1}|s_{t+1}).$$

$$\theta^* = \underset{\theta}{\operatorname{argmax}} E_{\gamma \sim P(\gamma)} \left[\sum_t r(s_t, a_t) \right].$$

$$= \underset{\theta}{\operatorname{argmax}} \sum_{t=1}^T E(s_t, a_t) \underset{\text{state-action marginal.}}{\uparrow} P(s_t, a_t) [r(s_t, a_t)].$$

finite horizon case. T.

$$(s_{t+1}, a_{t+1}) = \gamma(s_t, a_t). \Rightarrow (s_{t+k}, a_{t+k}) = \gamma^k(s_t, a_t).$$

γ state-action transition operator.

$t \rightarrow \infty$. if $P(s_t, a_t)$ converges to a stationary distribution μ , we can write $\mu = \gamma \mu$.

Assumptions: ① ergodicity.

every state can be reached from other states w/ nonzero probability.

② chain being aperiodic.

$$(\gamma - I)\mu = 0. \quad \mu \text{ is the eigenvector of } \lambda = 1. \text{ (always exist)} \quad \mu = P\theta(s, a).$$

$$\text{Infinite horizon case. } \theta^* = \underset{\theta}{\operatorname{argmax}} E(s, a) \underset{\text{e.g. GD.}}{\uparrow} P(s, a) [r(s, a)].$$

In RL, we always care about Expectations.

\rightarrow differentiable.

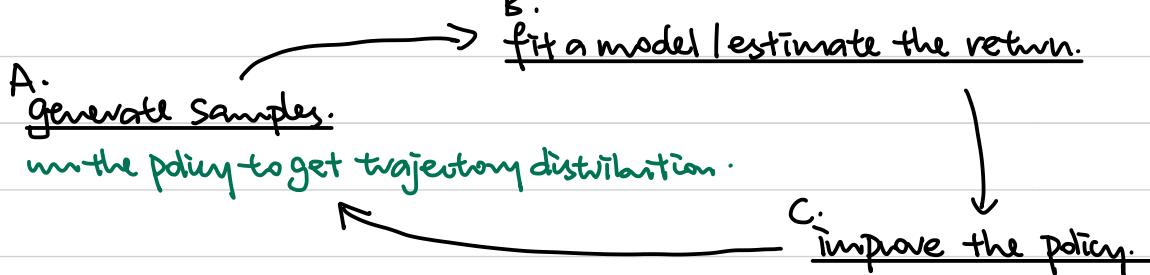
use smooth optimization algorithm for non-smooth reward function.

E.g. binary reward function. $r(x) = \begin{cases} 1 & \text{if not fall} \\ -1 & \text{if fall} \end{cases}$ Not Smooth.

$$Tu(a=\text{fall}) = 0.$$

$$E_{Tu}[r(x)] - \text{smooth in } \theta. \quad (-1)(0) + (1)(1-\theta).$$

Algorithms. Trial and Error.



Simple Example: B: $J(\theta) = E_{\pi} \left[\sum_t r_t \right] \approx \frac{1}{N} \sum_{t=1}^N \sum_i r_t^i$

$$C: \theta \leftarrow \theta + \alpha \nabla_{\theta} J(\theta).$$

expensive.

Another Example: backprop: B: learn f_{θ} s.t. $s_{t+1} \approx f_{\theta}(s_t, a_t)$.

C: backprop through f_{θ} and r to train $Tu(s_t) = a_t$.

Value Functions. Q-Functions. Evaluate Tu in B \Rightarrow C.

$$E_{\gamma \sim P(\gamma)} \left[\sum_t r(s_t, a_t) \right].$$

$$= \mathbb{E}_{S_1 \sim p(S_1)} [\mathbb{E}_{a_1 \sim \pi(a_1 | S_1)} [\underbrace{\mathbb{E}_{S_2 \sim p(S_2 | S_1, a_1)} [\mathbb{E}_{a_2 \sim \pi(a_2 | S_2)} [r(S_2, a_2) + \dots] | S_2]}_{Q(S_1, a_1)} | S_1].$$

$$= \mathbb{E}_{S_1 \sim p(S_1)} [\mathbb{E}_{a_1 \sim \pi(a_1 | S_1)} [Q(S_1, a_1)] | S_1].$$

if we know $Q(S_1, a_1)$, we can easily set $\pi(a_1 | S_1) = 1$ if $a_1 = \arg\max_a Q(S_1, a_1)$.

Q-function: $Q^\pi(s_t, a_t) = \sum_{t'=t}^T \mathbb{E}_\pi[r(s', a') | s_t, a_t]$ total expected future reward | st, at.

Value-function: $V^\pi(s_t) = \sum_{t'=t}^T \mathbb{E}_\pi[r(s', a') | s_t] = \mathbb{E}_{\pi^*}[r(s', a') | s_t] = Q^\pi(s_t, a_t)$.

RL-objective: $\mathbb{E}_{S_1 \sim p(S_1)} [V^\pi(S_1)]$.

why useful? ① if we know $Q^\pi(s, a)$, we can improve π . π' is at least as good as π .

② compute gradient to increase probability of good actions a .

if $Q^\pi(s, a) > V^\pi(s)$, a is better than average. \Rightarrow modify $\pi(a|s)$ to increase $\pi(a)$.

Types of RL Algorithms.

$$\pi^* = \arg\max_\pi \mathbb{E}_{\pi \sim \text{prior}} [\sum_t r(s_t, a_t)].$$

① policy gradient: directly differentiate obj. \Rightarrow GD.

② Value-Based: estimate Value fn / Q fn of optimal π^* to improve π . no explicit policy.

③ Actor-Critic: learn Value fn / Q fn of current π to improve π . ① + ②

④ Model-based RL: estimate the transition model for planning directly / improve policy | ...

④: $B \cdot P(s_{t+1} | s_t, a_t) \Rightarrow s_{t+1}$ or pro. dis. of s_{t+1} .

\rightarrow future actions without performing them.

c. a) use the model of the world to plan. No explicit policy.

e.g. (continuous) trajectory optimization / optimal control.

(discrete) Monte Carlo tree search.

b) back-propagate

c) learn a value function. (dynamic programming).

②: B : fit $V(s)$ or $Q(s, a)$ NN. C: set $\pi(s) = \arg\max_a Q(s, a)$.

①: B : $R_t = \sum_t r(s_t, a_t)$. C: $\theta \leftarrow \theta + \alpha \nabla_\theta E[\sum_t r(s_t, a_t)]$.

OR $V(s)$

③: B : fit $V(s)$ or $Q(s, a)$. C: $\theta \leftarrow \theta + \alpha \nabla_\theta E[Q(s, a)]$.

Different tradeoffs: A. Sample efficiency.

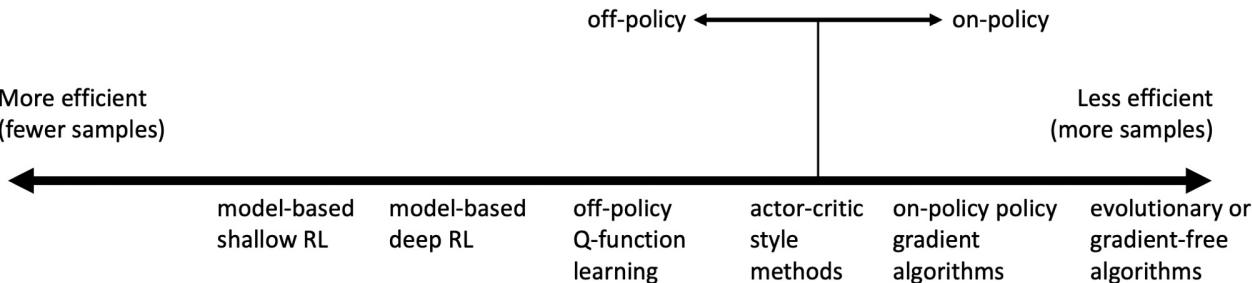
B. Stability and Ease of use.

BUT computation efficiency varies.

A. # samples we need to get a good policy.

off-policy algorithm: able to improve policy without generating new samples from that policy.

on-policy algorithm: generate samples each time the policy is changed.



B. Converge? Converge to what? Converge every time?

Though RL is often Not GD.

minimize error of fit.

> Model-based RL: optimized to be accurate model instead of obj.

> Policy gradient: GD. least efficient.

error of fit.

> Value function fitting: At best, minimize 'Bellman Error'

At worst, nothing.

fixed point iteration.

Different assumptions: #0 stochastic OR deterministic environment?

#1. full observability. (Value fitting).

#2. episodic learning. (policy gradient, model-based RL).
infinite horizon Some.

#3. Continuity | smoothness. (Value fitting, model-based RL).
discrete continuous Some.

easier to represent the policy or model.

π_{t+1} P_{t+1}

①: REINFORCE, Natural policy gradient. Trust region policy optimization.

②: Q-learning, Deep Q-Network (DQN), temporal difference learning. Fitted value iteration.

③: Asynchronous advantage actor-critic (A3C), Soft actor-critic (SAC)

④: Dyna. Guided policy search.

09/12/2022 Monday.

Policy Gradient

$$\theta^* = \underset{\theta}{\operatorname{argmax}} \left[\mathbb{E}_{\pi_\theta(z)} \left[\sum_t r(s_t, a_t) \right] \right]$$

$\downarrow J(\theta)$

$$J(\theta) = \mathbb{E}_{\pi_\theta(z)} [r(z)] = \int p_\theta(z) r(z) dz.$$

ith sample.

$$J(\theta) \approx \frac{1}{N} \sum_i r(s_{i,t}, a_{i,t}).$$

$$\begin{aligned} p_\theta(z) \nabla \log p_\theta(z) &= p_\theta(z) \frac{\nabla p_\theta(z)}{p_\theta(z)} = \nabla p_\theta(z). \\ \nabla J(\theta) &= \int \nabla p_\theta(z) r(z) dz = \int p_\theta(z) \nabla \log p_\theta(z) r(z) dz \\ &= \mathbb{E}_{\pi_\theta(z)} [\nabla \log p_\theta(z) r(z)]. \end{aligned}$$

By Obj. $\log p_\theta(z) = \log p(s_1) + \sum_{t=1}^T \log \pi_\theta(a_t | s_t) + \log p(s_{t+1} | s_t, a_t).$

$$\nabla \log p_\theta(z) = \sum_{t=1}^T \nabla \log \pi_\theta(a_t | s_t)$$

$$\nabla J(\theta) = \mathbb{E}_{\pi_\theta(z)} \left[\left(\sum_{t=1}^T \nabla \log \pi_\theta(a_t | s_t) \right) \left(\sum_{t=1}^T r(s_{i,t}, a_{i,t}) \right) \right].$$

known: $p(s_1), p(s_{t+1} | s_t, a_t)$.
unknown: $\pi_\theta(a_t | s_t)$.

REINFORCE Algorithm:

- ① Sample s_t^i from $\pi_\theta(a_t | s_t)$
- ② $\nabla J(\theta) \approx \sum_i \left(\sum_t \nabla \log \pi_\theta(a_t^i | s_t^i) \right) \left(\sum_t r(s_{i,t}, a_{i,t}) \right)$.
- ③ $\theta \leftarrow \theta + \lambda \nabla J(\theta)$.

→ formalize Trial and Error as GD.

Policy gradient: $\nabla J(\theta) \approx \frac{1}{N} \sum_i \left(\sum_{t=1}^T \nabla \log \pi_\theta(a_{i,t} | s_{i,t}) \right) \left(\sum_{t=1}^T r(s_{i,t}, a_{i,t}) \right)$.

a_t might be good/bad (since it's generated by our policy)
 \Rightarrow increase/decrease acc. to reward. \rightarrow weight.

Maximum likelihood: $\nabla J_{ML}(\theta) \approx \frac{1}{N} \sum_i \left(\sum_{t=1}^T \nabla \log \pi_\theta(a_{i,t} | s_{i,t}) \right)$

(Supervised learning). think a_t is good \Rightarrow Maximize/increase probability of a_t .

Gramssian.

Cont. case: $\pi_\theta(a_t | s_t) = N f_{NN}(s_t; \Sigma)$

↑ mean.

$$\log \pi_\theta(a_t | s_t) = -\frac{1}{2} \|f(s_t) - a_t\|_\Sigma^2 + \text{const.}$$

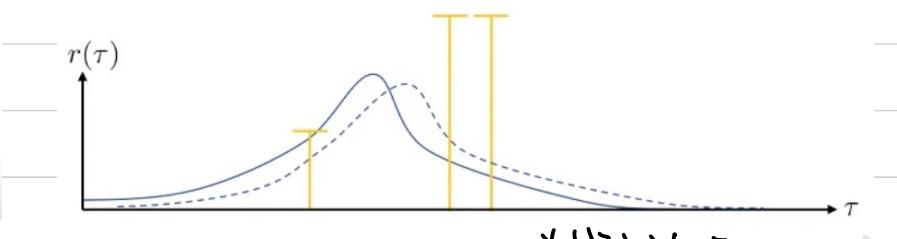
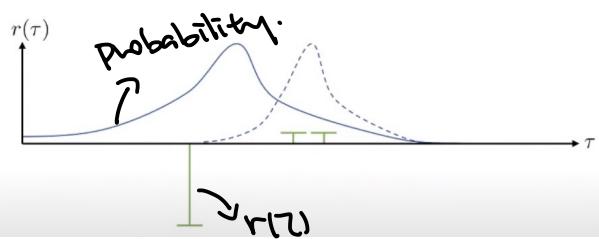
$$\nabla \log \pi_\theta(a_t | s_t) = -\frac{1}{2} \Sigma^{-1} (f(s_t) - a_t) \frac{df}{da}.$$

Partial observability.

s_t Markov
 a_t Non-Markov.

$\nabla J(\theta) \approx \frac{1}{N} \sum_i \left(\sum_{t=1}^T \nabla \log \pi_\theta(a_{i,t} | o_{i,t}) \right) \left(\sum_{t=1}^T r(s_{i,t}, a_{i,t}) \right)$. (I doesn't use Markov property).

Add a constant term to reward fn.



* High Variance.
if $r(z)=0$, $\text{DlogP}(r)$ doesn't matter at all.

w/ finite samples, largely depend on these random samples. \Rightarrow Need to reduce Variance.

w/ infinite samples, always converg. No

Reduce Variance

① Causality: policy at time t' cannot affect reward at time t if $t < t'$.

$$\mathbb{D}\bar{J}(\theta) \approx \frac{1}{N} \sum_{i=1}^N \sum_{t=1}^T \text{DlogP}(r_i(t)) \left(\sum_{t'=t}^T r_i(s_i(t'), a_i, t') \right). \quad \text{Smaller value expectation} \xrightarrow{\text{unbiased}} \text{smaller variance.}$$

"reward-to-go" $\hat{Q}_{i,t}$

② Baseline unbiased expectation, Smaller Variance.

$$\mathbb{D}\bar{J}(\theta) \approx \frac{1}{N} \sum_{i=1}^N \text{DlogP}(z_i) [r(z_i) - b]. \quad b = \frac{1}{N} \sum_{i=1}^N r(z_i) \quad \begin{array}{l} \text{(cannot increase } b \text{ for all } T \text{ with positive reward)} \\ \downarrow \\ \text{instead, higher than average.} \end{array}$$

$$\begin{aligned} \mathbb{E}[\text{DlogP}(z)b] &= \int p(z) \text{DlogP}(z) b dz \\ &= \int \text{DlogP}(z) b dz = b \int \text{DlogP}(z) dz = b \mathbb{D}\theta = 0. \end{aligned}$$

Optimal baseline to minimize variance:

$$\text{Var}(X) = \mathbb{E}[X^2] - \mathbb{E}[X]^2$$

$$\text{Var} = \mathbb{E}_z [p(z)] [\underbrace{(\text{DlogP}(z)(r(z)-b))^2}_{g(z)}] - \mathbb{E}_z [p(z)] [\text{DlogP}(z)(r(z)-b)]^2.$$

$$\frac{d\text{Var}}{db} = \frac{d}{db} \mathbb{E}_z [p(z)] [\underbrace{g(z)^2}_{r(z)}] = \frac{d}{db} (\mathbb{E}[g(z)^2] r(z)^2) - 2 \mathbb{E}[g(z)^2 r(z) b] + b^2 \mathbb{E}[g(z)^2] = 0.$$

$b = \frac{\mathbb{E}[g(z)^2 r(z)]}{\mathbb{E}[g(z)^2]}$ expected reward weighted by gradient magnitudes.

Policy gradient is on-policy. Since need generate data everytime we update θ .

NN only change a little bit each time \Rightarrow inefficient.

importance sampling. $\mathbb{E}_{x \sim p(x)} [f(x)] = \int p(x) f(x) dx = \int q(x) \frac{p(x)}{q(x)} f(x) dx = \mathbb{E}_{x \sim q(x)} \left[\frac{p(x)}{q(x)} f(x) \right].$

$$J(\theta) = \mathbb{E}_{z \sim p(z)} \left[\frac{p(z)}{q(z)} r(z) \right]$$

$$\frac{P(\pi)}{P(\pi')} = \frac{P(s_1) T_{\pi'}(a_1|s_1) P(a_1)}{P(s_1) T_{\pi}(a_1|s_1) P(a_1)} \quad \text{unknown}$$

$$\nabla_{\theta} J(\theta') = \mathbb{E}_{\pi' P(\pi)} \left[\frac{\nabla_{\theta} P(\pi)}{P(\pi)} r(t) \right] = \mathbb{E}_{\pi' P(\pi)} \left[\frac{P(\pi)}{P(\pi')} \nabla_{\theta} (\log P(\pi) / P(\pi')) r(t) \right].$$

$$\theta = \theta' := \mathbb{E}_{\pi' P(\pi)} [\nabla_{\theta} (\log P(\pi) / P(\pi')) r(t)].$$

$$\begin{aligned} \theta \neq \theta' &:= \mathbb{E}_{\pi' P(\pi)} \left[\left(\prod_{t=1}^T \frac{T_{\pi'}(a_t|s_t)}{T_{\pi}(a_t|s_t)} \right) \left(\sum_{t=1}^T \nabla_{\theta} \log T_{\pi'}(a_t|s_t) \right) \left(\sum_{t=1}^T r(s_t, a_t) \right) \right]. \\ &\stackrel{\text{by consistency}}{=} \mathbb{E}_{\pi' P(\pi)} \left[\sum_{t=1}^T \nabla_{\theta} \log T_{\pi'}(a_t|s_t) \left(\prod_{t'=t+1}^T \frac{T_{\pi'}(a_{t'}|s_{t'})}{T_{\pi}(a_{t'}|s_{t'})} \right) \left(\sum_{t'=t}^T r(s_{t'}, a_{t'}) \left(\prod_{t''=t'+1}^T \frac{T_{\pi'}(a_{t''}|s_{t''})}{T_{\pi}(a_{t''}|s_{t''})} \right) \right) \right]. \end{aligned}$$

\downarrow , T finite \Rightarrow exponent to 1
Variance expands to infinite

$$\text{Off-policy policy gradient} = \nabla_{\theta} J(\theta') \approx \frac{1}{T} \sum_{t=1}^T \frac{T_{\pi'}(s_t, a_t)}{T_{\pi}(s_t, a_t)} \nabla_{\theta} (\log T_{\pi'}(a_t|s_t)) r_t.$$

↑
ignore works if $\theta' \rightarrow \theta$

$$= \frac{T_{\pi'}(s_t)}{T_{\pi}(s_t)} \frac{\nabla_{\theta} (\log T_{\pi'}(a_t|s_t))}{T_{\pi}(s_t, a_t)}$$

implement policy gradient

Automatic differentiation:

"pseudo-loss": $\hat{J}(\theta) \approx \frac{1}{T} \sum_{t=1}^T \log(T_{\pi}(\cdot|s_t, a_t)) r_t \Rightarrow \nabla_{\theta} \hat{J}(\theta)$ is policy gradient.

Given:

```
# actions - (N*T) x Da tensor of actions
# states - (N*T) x Ds tensor of states
# q_values - (N*T) x 1 tensor of estimated state-action values
# Build the graph:
logits = policy.predictions(states) # This should return (N*T) x Da tensor of action logits
negative_likelihoods = tf.nn.softmax_cross_entropy_with_logits(labels=actions, logits=logits)
weighted_negative_likelihoods = tf.multiply(negative_likelihoods, q_values)
loss = tf.reduce_mean(weighted_negative_likelihoods)
gradients = loss.gradients(loss, variables)
```

* Consider large batches.

* tuning will be hard (Adam, Ok-ish)

Advanced Policy Gradient

$$\theta' \leftarrow \underset{\theta'}{\text{argmax}} \, (\theta' - \theta)^T \nabla_{\theta} J(\theta) \text{ s.t. } D(T_{\pi'}, T_{\pi}) \leq \epsilon.$$

parameter-independent divergence measure (e.g. KL-divergence)

$$D_{KL}(T_{\pi'} || T_{\pi}) = \mathbb{E}_{\pi'} [\log T_{\pi'} - \log T_{\pi}]$$

$$D_{KL}(T_{\pi'} || T_{\pi}) \approx (\theta' - \theta)^T F(\theta' - \theta) \quad \text{Fisher-information matrix. } F = \mathbb{E}_{\pi'} [\nabla_{\theta} \log T_{\pi}(\cdot|s) \nabla_{\theta} \log T_{\pi}(\cdot|s)^T]$$

$$\theta' \leftarrow \underset{\theta'}{\text{argmax}} \, (\theta' - \theta)^T \nabla_{\theta} J(\theta) \text{ s.t. } \| \theta' - \theta \|_F^2 \leq \epsilon. \quad \theta \leftarrow \theta + \lambda F^{-1} \nabla_{\theta} J(\theta)$$

natural policy gradient: pick 2.

Trust Region Policy optimization (TRPO): pick ϵ . solve for optimal λ while solving $F^{-1} \nabla_{\theta} J(\theta)$.

09/14/2022 Wednesday.

π^{θ} $V^{\pi}(s_t)$
Actor-Critic

policy-gradient: $\nabla J(\theta) \approx \frac{1}{N} \sum_{t=1}^T \sum_{a_t} \nabla \log \pi_\theta(a_t | s_t) \left(\sum_{t'=t+1}^T r(s_{t'}, a_{t'}) - b \right).$

unbiased, high variance, single-sample estimate.

$\hat{Q}_{i,t}$: estimate of expected future reward if we take $a_{i,t}$ in state $s_{i,t}$.

randomness $\Rightarrow \hat{Q}_{i,t} \approx \sum_{t'=t+1}^T E_{\pi^\theta}[r(s_{t'}, a_{t'}) | s_{i,t}, a_{i,t}]$. (better estimate).
 \uparrow decrease variance.

$\nabla J(\theta) \approx \frac{1}{N} \sum_{t=1}^T \sum_{a_t} \nabla \log \pi_\theta(a_t | s_t) \left(\hat{Q}_{i,t} - b \right)$. how much better $a_{i,t}$ is (avg) over all $a_{i,t}$ given $s_{i,t}$. (avg)
 \uparrow decrease variance.

$Q = \sum_{t'=t+1}^T E_{\pi^\theta}[r(s_{t'}, a_{t'}) | s_{i,t}, a_{i,t}]$ true expected reward-to-go.
 \uparrow decrease variance.

$b = \frac{1}{N} \sum_i Q(s_{i,t}, a_{i,t}) \Rightarrow b = V(s_t) = \underbrace{E_{\pi^\theta}[r(s_{t+1}) | s_t, a_t]}_{\text{all possible } a_t \text{ given } s_t}.$

$A^{\pi^*}(s_t, a_t) = Q^{\pi^*}(s_t, a_t) - V^{\pi^*}(s_t)$.

Actor-Critic: $\nabla J(\theta) \approx \frac{1}{N} \sum_{t=1}^T \sum_{a_t} \nabla \log \pi_\theta(a_t | s_t) A^{\pi^*}(s_t, a_t)$.

much lower variance, little bias.
 \uparrow fit Q^{π^*}, V^{π^*} or A^{π^*} in B.

fitting Q^{π^*}, V^{π^*} or A^{π^*} .

$Q^{\pi^*}(s_t, a_t) = r(s_t, a_t) + \sum_{t'=t+1}^T E_{\pi^\theta}[r(s_{t'}, a_{t'}) | s_t, a_t] = r(s_t, a_t) + \sum_{t'=t+1}^T E_{\pi^\theta}[p(s_{t+1} | s_t, a_t)] \bar{V}^{\pi^*}(s_{t+1}).$
 \uparrow Not Random \uparrow $\bar{V}^{\pi^*}(s_{t+1})$. Single-sample estimate at $t+1$ only.

$A^{\pi^*}(s_t, a_t) \approx r(s_t, a_t) + V^{\pi^*}(s_{t+1}) - V^{\pi^*}(s_t)$. V is easy to learn (depend only on state). \Rightarrow fit $V^{\pi^*}(s)$.

Monte Carlo policy evaluation (fit $V^{\pi}(s)$).

$V^{\pi}(s_t) = \sum_{t'=t}^T E_{\pi^\theta}[r(s_{t'}, a_{t'}) | s_t]$. $J(\theta) = E_{\pi^\theta}[V^{\pi}(s_1)]$.

model-free: $V^{\pi^*}(s_t) \approx \sum_{t'=t}^T r(s_{t'}, a_{t'})$ \uparrow GP

Want $V^{\pi^*}(s_t) = \frac{1}{N} \sum_{t'=t}^T \sum_{a_{t'}} r(s_{t'}, a_{t'})$
 \uparrow trajectory.

only at the beginning.

but we can't run multiple trials at a particular state.

model-based:

NN will learn to make it similar. \Rightarrow decrease variance.
 $y_{i,t}$: training data: $\{(s_{i,t}, \sum_{t'=t}^T r(s_{t'}, a_{t'}))\}$. \uparrow prevent overfit.

Supervised: $L(\phi) = \frac{1}{N} \sum_i \| \hat{V}_\phi(s_i) - y_i \|^2$

Ideal $y_{i,t} = \sum_{t'=t}^T E_{\pi^\theta}[r(s_{t'}, a_{t'}) | s_{i,t}] \approx r(s_{i,t}, a_{i,t}) + V^{\pi^*}(s_{i,t+1}) \approx r(s_{i,t}, a_{i,t}) + \hat{V}_\phi^*(s_{i,t+1})$
 \uparrow previous fitted.
 \uparrow bootstrapped estimate
 \uparrow high bias, low variance.

Actor-Critic Algorithm (batch)

- 1. Sample (S_i, a_i) from $T\theta(a|s)$
2. fit $\hat{V}_\phi^\pi(s)$ to sampled reward sums
3. evaluate $\hat{A}^\pi(S_i, a_i) = r(S_i, a_i) + \gamma \hat{V}_\phi^\pi(S'_i) - \hat{V}_\phi^\pi(S_i)$
4. $\nabla J(\theta) \approx \sum_i \nabla \log T\theta(a_i | S_i) \hat{A}^\pi(S_i, a_i)$
5. $\theta \leftarrow \theta + \alpha \nabla J(\theta)$.

Discount factor. (Reduce variance, introduce bias). ↗ Not counting all rewards.

if $T \rightarrow \infty$, $\hat{V}_\phi^\pi \rightarrow \infty$.

better to get rewards sooner than later.

$$y_{it} \approx r(S_{it}, a_{it}) + \gamma \hat{V}_\phi^\pi(S_{i,t+1}) + (1-\gamma) \dots \quad \text{1st step future.}$$

$$\begin{aligned} \text{to Monte Carlo GP: option 1: } \nabla J(\theta) &\approx \frac{1}{N} \sum_{i=1}^N \sum_{t=1}^T \nabla \log T\theta(a_{it} | S_{it}) \left(\sum_{t'=t}^T \gamma^{t-t'} r(S_{i,t'}, a_{i,t'}) \right) \\ \text{option 2: } \nabla J(\theta) &\approx \frac{1}{N} \sum_{i=1}^N \left(\sum_{t=1}^T \nabla \log T\theta(a_{it} | S_{it}) \left(\sum_{t'=1}^{T-1} \gamma^{t'-1} r(S_{i,t'}, a_{i,t'}) \right) \right) \quad \text{different} \\ &= \frac{1}{N} \sum_{i=1}^N \sum_{t=1}^T \nabla \log T\theta(a_{it} | S_{it}) \left(\sum_{t'=t}^T \gamma^{t'-1} r(S_{i,t'}, a_{i,t'}) \right) \\ &= \frac{1}{N} \sum_{i=1}^N \sum_{t=1}^T \gamma^{t-1} \nabla \log T\theta(a_{it} | S_{it}) \left(\sum_{t'=t}^T \gamma^{t'-t} r(S_{i,t'}, a_{i,t'}) \right) \end{aligned}$$

discount-future policy as well. \Rightarrow decide case by case.

option 1 > option 2: actually want infinite movements.

$$\text{to actor-critic: } \nabla J(\theta) \approx \frac{1}{N} \sum_{i=1}^N \sum_{t=1}^T \nabla \log T\theta(a_{it} | S_{it}) (r(S_{it}, a_{it}) + \gamma \hat{V}_\phi^\pi(S_{i,t+1}) - \hat{V}_\phi^\pi(S_{it}))$$

online actor-critic algorithm.

- 1. take action $a \sim T\theta(a|s)$, get (s, a, s', r)
2. update \hat{V}_ϕ^π using target $r + \gamma \hat{V}_\phi^\pi(s')$
3. evaluate $\hat{A}^\pi(s, a) = r(s, a) + \gamma \hat{V}_\phi^\pi(s') - \hat{V}_\phi^\pi(s)$
4. $\nabla J(\theta) \approx \nabla \log T\theta(a|s) \hat{A}^\pi(s, a)$
5. $\theta \leftarrow \theta + \alpha \nabla J(\theta)$.

Actor-Critic Design Decisions.

two network design.
shared network design.

Synchronized parallel update Ø using synchronized info. batch size = # parallel workers.
 Asynchronous parallel update Ø separately.

Off-policy actor-critic: sample batch from replay buffer. (all old transitions) $\rightarrow (s, a, s', r)$

Problem: T_θ is Not the latest \Rightarrow Step 3. Many y_i .

Step 5. Many $\nabla \log \pi_\theta(a_t | s_t)$

$$\text{Step 3: } \bar{V}^\pi(s_t) = \sum_{t'=t}^T \mathbb{E} V^\pi[r(s_t, a_t) | s_t].$$

$$\bar{Q}^\pi(s_t, a_t) = \sum_{t'=t}^T \mathbb{E} Q^\pi[r(s_t, a_t) | s_t, a_t].$$

$$3. \text{ update } \hat{Q}_\phi^\pi \text{ using targets } y_i = r_i + \gamma \hat{Q}_\phi^\pi(s'_i, a'_i) \quad a'_i \sim \pi_\theta(a_i | s_i) \quad L(\phi) = \frac{1}{N} \sum_i \| \hat{V}_\phi^\pi(s_i) - y_i \|^2.$$

$$\text{Step 5: } \nabla J(\theta) \approx \frac{1}{N} \sum_i \nabla \log \pi_\theta(a_t^\pi | s_t) \hat{Q}^\pi(s_t, a_t^\pi) \quad a_t^\pi \sim \pi_\theta(a_t | s_t).$$

\downarrow Convenience / high variance!

\uparrow from RB

\rightarrow Not Replay Buffer.

\rightarrow Not necessarily follow T_θ .

State-dependent baseline.

$$\nabla J(\theta) \approx \frac{1}{N} \sum_{t=1}^T \nabla \log \pi_\theta(a_t | s_t) \left(\left(\sum_{t'=t}^T \gamma^{t-t'} r(s_t, a_t, t') \right) - \hat{V}_\phi^\pi(s_t) \right). \quad \text{critics.}$$

Action-dependent baseline.

$$\text{Control variates: } \nabla J(\theta) \approx \frac{1}{N} \sum_{t=1}^T \nabla \log \pi_\theta(a_t | s_t) (\hat{Q}_{it} - \hat{Q}_\phi^\pi(s_t, a_t)) + \frac{1}{N} \sum_{t=1}^T \nabla \mathbb{E} \nabla \log \pi_\theta(a_t | s_t) (\hat{Q}_\phi^\pi(s_t, a_t)).$$

$$\text{eligibility traces, 1-h-step returns: } \hat{A}^\pi(s_t, a_t) = \sum_{t'=t}^{\text{term}} \gamma^{t-t'} r(s_t, a_t) - \hat{V}_\phi^\pi(s_t) + \gamma^{\text{term}} \hat{V}_\phi^\pi(s_{\text{term}}).$$

(cut before variance gets too big).

\uparrow Contribute Variance \uparrow usually $n > 1$; $n \uparrow$, variance \uparrow , bias \downarrow Contribute bias.

Generalized Advantage Estimation (GAE).

$$\hat{A}_{\text{GAE}}^\pi(s_t, a_t) = \sum_{n=1}^{\infty} \lambda^n \hat{A}_n^\pi(s_t, a_t).$$

$$\text{With } \alpha \propto \lambda^{n-1}. \quad \hat{A}_{\text{GAE}}^\pi(s_t, a_t) = r(s_t, a_t) + \gamma((1-\lambda)\hat{V}_\phi^\pi(s_{t+1}) + \lambda(r(s_{t+1}, a_{t+1}) + \gamma((1-\lambda)\hat{V}_\phi^\pi(s_{t+2}) + \lambda(r(s_{t+2}, a_{t+2}) + \dots))$$

PREFER CUT EARLIER.

$$= \sum_{t'=t}^{\infty} (\gamma \lambda)^{t'-t} \delta_t' \quad \delta_t' = r(s_t, a_t) + \gamma \hat{V}_\phi^\pi(s_{t+1}) - \hat{V}_\phi^\pi(s_t).$$

09/19/2022 Monday.

Value Function Methods.

$\underset{a_t}{\operatorname{argmax}} A^\pi(s_t, a_t) \rightarrow$ best action from s_t if we follow π regardless of $\pi(a_t | s_t)$.

C. policy iteration algorithm. \rightarrow 1. evaluate $\hat{A}^\pi(s, a) = \frac{r(s, a) + \gamma \mathbb{E}[V^\pi(s')]}{(Q^\pi(s, a))}$ ★

2. set $\pi' \leftarrow \pi'$

greedy policy. $\pi'(a_t | s_t) = \begin{cases} 1 & \text{if } a_t = \underset{a_t}{\operatorname{argmax}} A^\pi(s_t, a_t) = \underset{a_t}{\operatorname{argmax}} Q^\pi(s_t, a_t) \\ 0 & \text{otherwise.} \end{cases}$

S and a both discrete and small "tabular" \rightarrow guarantee to converge.

$T = |S| \times |S| \times |a|$ tensor. $\star V^T \leftarrow \underset{T^*(s)}{\text{E}}[r(s,a) + \gamma \underset{T^*(s')}{\mathbb{E}}[V^T(s')]]$

Value iteration algorithm \rightarrow 1. set $Q(s,a) \leftarrow r(s,a) + \gamma \underset{s'}{\mathbb{E}}[V(s')]$.

2. set $V(s) \leftarrow \max_a Q(s,a)$ from argmax.

> Continuous "Curse of dimensionality" $|S|$ too big.

fitted value iteration algorithm: 1. Set $y_i \leftarrow \max_{a_i} [r(s_i, a_i) + \gamma \underset{s'}{\mathbb{E}}[V_\phi(s')]]$.

2. Set $\phi \leftarrow \underset{\phi}{\text{argmin}} \frac{1}{N} \sum_i \|V_\phi(s_i) - y_i\|^2$.

if we don't know the transition dynamics: $\star Q^T \leftarrow r(s,a) + \gamma \underset{s'}{\mathbb{E}}[Q^T(s', T^*(s'))]$.

(transition model)

Model-free. Batch Mode.

using samples (s, a, s') instead of T^* .

parameters.

dataset size N , collection policy.

iteration k .

gradient step s .

1. Collect dataset $\{(s_i, a_i, s'_i, r_i)\}$ using some policy.

2. Set $y_i \leftarrow r(s_i, a_i) + \gamma \max_{a'_i} Q_\phi(s'_i, a'_i)$. approx. the value of T^* at s'

3. Set $\phi \leftarrow \underset{\phi}{\text{argmin}} \frac{1}{N} \sum_i \|Q_\phi(s_i, a_i) - y_i\|^2$

+: Works for off-policy; one NN $-$: no guarantee to converge.

iterate value instead of policy. \hookrightarrow input: s, a output: $Q_\phi(s, a)$

Batch Mode.

Online Q learning algorithm off-policy.

data available in sequential order.

1. take some action a_t and observe (s_t, a_t, s'_t, r_t)

2. $y_t = r(s_t, a_t) + \gamma \max_{a'_t} Q_\phi(s'_t, a'_t)$

3. $\phi \leftarrow \phi - \alpha \frac{d\phi}{da} (s_t, a_t) (Q_\phi(s_t, a_t) - y_t)$ *Not GD.

About the initial policy: not \heartsuit (fail to find high-reward states), need randomness.

$\rightarrow \pi_{\text{initial}}(s_t) = \begin{cases} 1-\epsilon & \text{if } a_t = \underset{a'_t}{\text{argmax}} Q_\phi(s_t, a'_t) \\ \epsilon & \text{otherwise.} \end{cases}$ \heartsuit epsilon-greedy. Varying ϵ .

$\rightarrow \pi_{\text{initial}}(s_t) \propto \exp(Q_\phi(s_t, a_t))$ Boltzmann Exploration.

+: if a_t have the highest Q , can keep both.

Wont waste time on a_t w/ extremely low Q .

Backup communicate.
 Consider operator $B: BV = \max_a (r_a + \gamma T V')$ \uparrow vector of rewards at all s for a .
 $|S| \times |S|; P(S' = i | S = j, a)$.

$V^*(s) = \max_a r(s, a) + \gamma E[V^*(s')]$ $\Rightarrow V^* = BV^*$. a fixed point of B . always exist / unique.

Value Iteration Algorithm converges to V^* since B is a contraction w.r.t. ∞ -norm.

$\hookrightarrow 1. V \leftarrow BV$

Habarau)

$$\|BV - B\bar{V}\|_\infty \leq \gamma \|V - \bar{V}\|_\infty \Rightarrow \|BV - V^*\|_\infty \leq \gamma \|V - V^*\|_\infty.$$

fitted Value Iteration Algorithm:

$\hookrightarrow 1. V \leftarrow \pi BV$ $\pi: \pi V = \underset{V' \in \Omega}{\operatorname{argmin}} \frac{1}{2} \sum \|V'(s) - V(s)\|^2$ \hookrightarrow all value-fn represented by (e. NN).
 \downarrow projection.
 contraction w.r.t. ℓ_2 -norm. $\|\pi V - \bar{V}\|^2 \leq \|V - \bar{V}\|^2$

πB is Not a contraction. \Rightarrow Not Converge (general).

fitted Q-Iteration Algorithm: Also Online Q-learning.

$\hookrightarrow 1. Q \leftarrow \pi BQ$. $B: BQ = r + \gamma \max_a Q$ $\pi: \pi Q = \underset{Q' \in \Omega}{\operatorname{argmin}} \frac{1}{2} \sum \|Q'(s, a) - Q(s, a)\|^2$.
 \downarrow contraction ∞ -norm. \downarrow contraction ℓ_2 -norm.

πB is Not a contraction. \Rightarrow Not Converge

Similarly, Actor-Critic using bootstrap target values. (i.e. $y = r + \gamma V(S')$) Not converge.

all algorithms use NN and bootstrap Not converge

01/21/2022 Wednesday.

Deep RL with Q-functions.

\ominus : Sequential states are strongly correlated.

Solution 1: parallel \leftarrow synchronized.

Solution 2: full Q-learning with replay buffer and target network.

\rightarrow 0. save target network parameters: $\phi' \leftarrow \phi$.

1. Collect dataset $\{(s_i, a_i, s'_i, r_i)\}$ using some policy, add it to B .

$\left[\begin{matrix} Nx \\ Kx \end{matrix} \right]$ 2. Sample a batch (s_i, a_i, s'_i, r_i) from B uniformly.

3. $\phi \leftarrow \phi - \alpha \sum_i \frac{\partial Q\phi}{\partial \phi} (s_i, a_i) (Q\phi(s_i, a_i) - \underbrace{[r_i s_i, a_i] + \gamma \max_{a'} Q\phi(s'_i, a'_i)]}_{\text{moving target}})$.

\ominus : moving target \Rightarrow doesn't change in the inner loop.

$K=1 \Rightarrow$ classic deep Q-learning Algorithm. (DQN.)

Alternative Target Network: update ϕ' each single step: $\phi' \leftarrow \gamma \phi' + (1-\gamma) \phi$.

(instead of every N steps).

e.g. $\gamma = 0.999$.

Supervised regression
SGD.

General View: Process 1: $(S, a, S', r) \xrightarrow{\text{finite}} \text{replay buffer} \xleftarrow{\text{all one step.}} T(\theta)(S)$

Process 2: $\phi \rightarrow \phi'$

process 3: Q-function regression.

Vanilla Q-learning: evict immediately. $\text{replay buffer} = 1$. \rightarrow all one step.

DQN: process 1/3 at the same speed. Process 2 slower.

fitted Q-iteration: process 3 in the inner loop of process 2. Process 2 in the inner loop of process 1.

Overestimate in Q-learning.

$$y_i = r_i + \gamma \max_{a'_j} Q\phi'(S'_j, a'_j) \quad E[\max(x_1, x_2)] \geq \max(E[x_1], E[x_2]).$$

$$= Q\phi'(S', \arg\max_{a'} Q\phi'(S', a'))$$

Solution:

$$\begin{cases} Q\phi_A(S, a) \leftarrow r + \gamma Q\phi_B(S', \arg\max_{a'} Q\phi_A(S', a')) \\ Q\phi_B(S, a) \leftarrow r + \gamma Q\phi_A(S', \arg\max_{a'} Q\phi_B(S', a')) \end{cases}$$

Double Q-learning. $y = r + \gamma Q\phi'(S', \arg\max_{a'} Q\phi(S', a')).$

N-step return estimator. $y_{i,t} = \sum_{t'=t}^{t+N-1} \gamma^{t-t'} r_{j,t'} + \gamma^N \max_{a'_j} Q\phi(S_j, t+N, a'_j, t+N).$
lower bias, higher variance.
* on policy only.

Q-learning with continuous actions

option A: stochastic optimization: cross-entropy Method (CEM)

CMA-ES

option B: easily maximized Q-functions: Normalized Advantage Functions (NAF).

$$Q\phi(S, a) = -\frac{1}{2} (\alpha - \mu\phi(S))^T P\phi(S) (\alpha - \mu\phi(S)) + V\phi(S).$$

$$\arg\max_{\alpha} Q\phi(S, a) = \mu\phi(S). \quad \max_{\alpha} Q\phi(S, a) = V\phi(S).$$

option C: ^{NN} learn an approximate maximizer: DDPG.

Compute $y_j = r_j + \gamma Q\phi'(S'_j, \mu\phi'(S'_j))$ using target nets $Q\phi'$ and $\mu\phi'$.

$$\phi \leftarrow \phi - \beta \sum_j \frac{d\phi}{d\theta}(S_j, a_j) (Q\phi(S_j, a_j) - y_j). \quad \theta \leftarrow \theta + \beta \sum_j \frac{d\mu}{d\theta}(S_j) \frac{dQ\phi}{da}(S_j, \mu(S_j)).$$