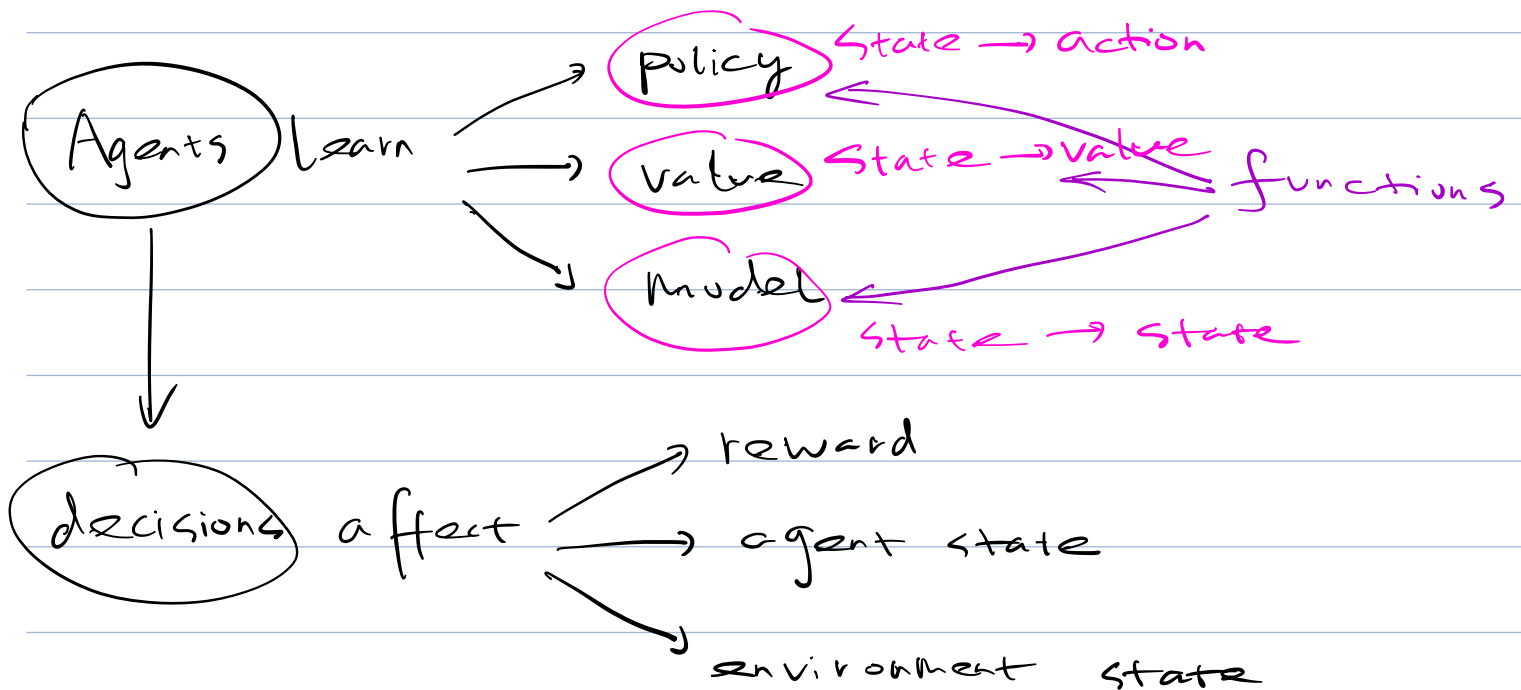


Deep Reinforcement Learning

(RL with function approximation)

Learning to make decisions

DNN



Value-based (Learn value function + implicit policy)

Learning value functions

G_0 : 10^{170} states

基于值的算法先评估每个 (s, a)

lookup tables

元组的Q值 $Q(s, a)$, 再根据Q值求最优策略.

State s : $V(s)$
State | s, a : $Q(s, a)$
action

$$V_{\theta}(s) \approx V_{\pi}(s)$$

$$Q_{\theta}(s, a) \approx Q_{\pi}(s, a)$$

Update θ using MC or TD learning

if the environment state is not fully observable:

→ Use agent state

→ Learn a state update function

$$S_{t+1} = u(S_t, D_{t+1})$$

Function Approximator:

1. ANN (DNN)
2. Decision Tree
3. Nearest Neighbor
4. Fourier / wavelet bases
5. Coarse coding

RL Notice:

- ① experience is not iid
- ② update non-stationary
- ③ feedback delayed

Classes of Function Approximation

Tabular

State Aggregation

Linear function approximation

Stochastic Gradient Descent

for value function Approximation

$$J(w) = \mathbb{E}_{\pi} [(V_{\pi}(s) - \hat{V}(s, w))^2]$$

$$\Delta w = -\frac{1}{\gamma} \alpha \nabla_w J(w)$$

$$= \alpha \mathbb{E}_{\pi} [(V_{\pi}(s) - \hat{V}(s, w)) \nabla_w \hat{V}(s, w)]$$

Sample:

$$\Delta w = \alpha (V_{\pi}(s) - \hat{V}(s, w)) \nabla_w \hat{V}(s, w)$$

Feature Vector

Represent state by a feature vector

$$\phi(s) = \begin{pmatrix} \phi_1(s) \\ \phi_2(s) \\ \vdots \\ \phi_n(s) \end{pmatrix}$$

$$\phi : \mathcal{S} \rightarrow \mathbb{R}^n$$

$$\phi_t = \phi(s_t)$$

Linear Approximation

$$V_{\theta}(s) = \theta^T \phi(s) = \sum_{j=1}^n \phi_j(s) \theta_j$$

$$J(\theta) = \mathbb{E}_{\pi} [(V_{\pi}(s) - \theta^T \phi(s))^2]$$

$$\nabla_{\theta} V_{\theta}(s_{t+1} = \phi(s_t) = \phi_t$$

$$\rightarrow \Delta \theta = \alpha (V_{\pi}(s_{t+1}) - V_{\theta}(s_t)) \phi_t$$

$$\text{update} = \text{step size} \times \underset{\text{error}}{\text{prediction}} \times \underset{\text{value}}{\text{feature}}$$

Table lookup features

Incremental Prediction Algorithms

MC:

$$\Delta \theta_t = \alpha (G_t - V_\theta(s)) \nabla_\theta V_\theta(s)$$

Linear TD (0):

$$\Delta \theta_t = \alpha (\underbrace{R_{t+1} + \gamma V_\theta(s_{t+1}) - V_\theta(s_t)}_{\delta_t \text{ "TD error"}}) \nabla_\theta V_\theta(s_t)$$

TD (λ):

$$\Delta \theta_t = \alpha (G_t^\lambda - V_\theta(s)) \nabla_\theta V_\theta(s_t)$$

Linear MC:

$$\begin{aligned} \Delta \theta_t &= \alpha (G_t - V_\theta(s_t)) \nabla_\theta V_\theta(s_t) \\ &= \alpha (G_t - V_\theta(s_t)) \phi_t \end{aligned}$$

Convergence of MC

$$\min_{\theta} \mathbb{E}[(G_t - V_{\theta}(s_{t+1}))^2]$$

$$= \mathbb{E}[\phi_t \phi_t^T]^{-1} \mathbb{E}[V_{\pi}(s_{t+1}) \phi_t]$$

Proof:

$$\nabla_{\theta} \mathbb{E}[(G_t - V_{\theta}(s_{t+1}))^2]$$

$$= \mathbb{E}[(G_t - V_{\theta}(s_{t+1})) \phi_t] = 0$$

$$\mathbb{E}[(G_t - \phi_t^T \theta) \phi_t] = 0$$

$$\mathbb{E}[G_t \phi_t - \phi_t \phi_t^T \theta] = 0$$

$$\mathbb{E}[\phi_t \phi_t^T] \theta = \mathbb{E}[G_t \phi_t]$$

$$\theta = \mathbb{E}[\phi_t \phi_t^T]^{-1} \mathbb{E}[V_{\pi}(s_{t+1}) \phi_t]$$

Convergence of TD

$$\min_{\theta} \mathbb{E} [(R_{t+1} + \gamma V_{\theta}(S_{t+1}) - V_{\theta}(S_t))^2]$$

$$= \mathbb{E} [\phi_t (\phi_t - \gamma \phi_{t+1})^T]^{-1} \mathbb{E} [R_{t+1} \phi_t]$$

TD method : faster , better

policy evaluation : prefer MC

Residual Bellman updates

TD:

$$\Delta \theta_t = \alpha \delta \nabla_{\theta} V_{\theta}(S_t)$$

$$\delta_t = R_{t+1} + \gamma V_{\theta}(S_{t+1}) - V_{\theta}(S_t)$$

Bellman residual gradient update:

$$\text{loss: } \mathbb{E}[\delta_t^2]$$

$$\text{update: } \Delta \theta_t = \alpha \delta_t \nabla_{\theta} (V_{\theta}(S_t) - \gamma V_{\theta}(S_{t+1}))$$

Work WORSE in practice

$$\frac{[R_{t+1} + \gamma V_{\theta}(S_{t+1})] - V_{\theta}(S_t)}{\text{TensorFlow}}$$

Action-Value Function Approximation

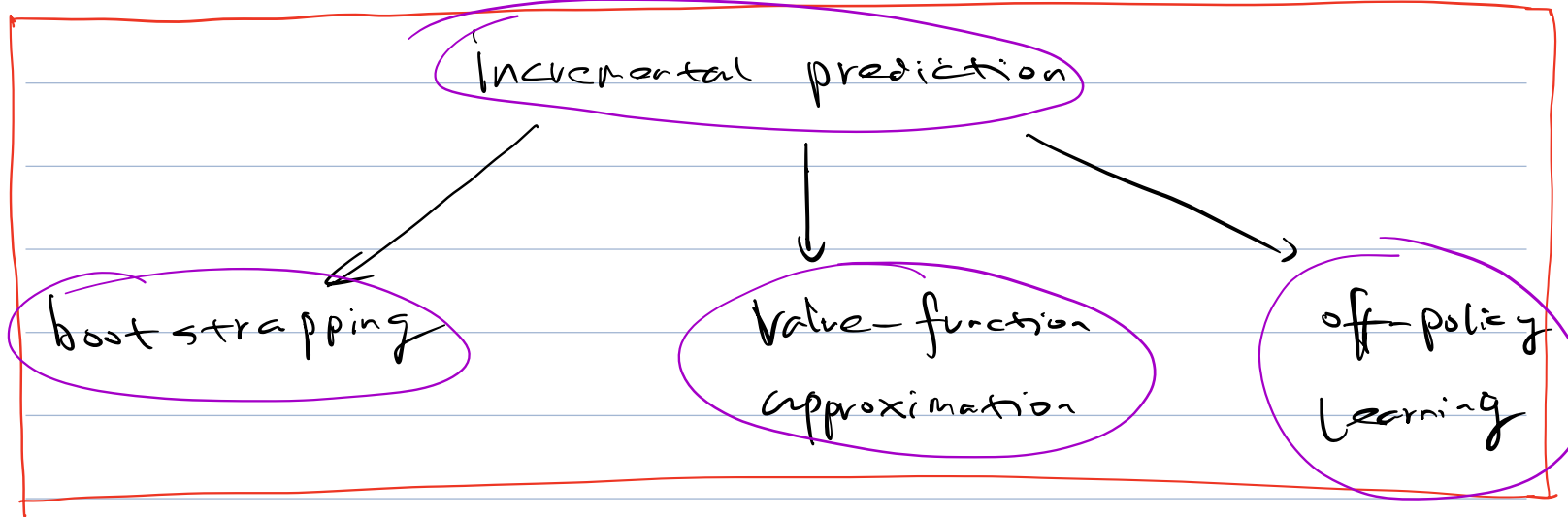
$$q_{\theta}(s, a) \approx q_{\pi}(s, a)$$

linear:

$$q_{\theta}(s, a) = \phi(s, a)^T \theta = \sum_{j=1}^n \phi_j(s, a) \theta_j$$

Stochastic Gradient Descent Update:

$$\begin{aligned} \Delta \theta &= \alpha (q_{\pi}(s, a) - q_{\theta}(s, a)) \nabla_{\theta} q_{\theta}(s, a) \\ &= \alpha (q_{\pi}(s, a) - q_{\theta}(s, a)) \phi(s, a) \end{aligned}$$



Tabular control learning algorithms

(Q-learning)

DQN (deep Q network)

Batch RL (Experience Replay)

Reinforcement Learning 5: Function Approximation and Deep Reinforcement Learning

Stochastic Gradient Descent with Experience Replay

Given experience consisting of $\langle \text{state}, \text{value} \rangle$ pairs

$$\mathcal{D} = \{ \langle S_1, \hat{v}_1^\pi \rangle, \langle S_2, \hat{v}_2^\pi \rangle, \dots, \langle S_T, \hat{v}_T^\pi \rangle \}$$

Repeat:

1. Sample state, value from experience

$$\langle s, \hat{v}^\pi \rangle \sim \mathcal{D}$$

2. Apply stochastic gradient descent update

$$\Delta \theta = \alpha (\hat{v}^\pi - v_\theta(s)) \nabla_\theta v_\theta(s)$$

Converges to least squares solution

$$\theta^\pi = \underset{\theta}{\operatorname{argmin}} \operatorname{LS}(\theta) = \underset{\theta}{\operatorname{argmin}} \mathbb{E}_{\mathcal{D}} [(\hat{v}_i^\pi - v_\theta(S_i))^2]$$

Linear Least Squares Prediction (2)

- At minimum of $\operatorname{LS}(\theta)$, the expected update must be zero

$$\mathbb{E}_{\mathcal{D}} [\Delta \theta] = 0$$

$$\alpha \sum_{t=1}^T \phi_t (\hat{v}_t^\pi - \phi_t^\top \theta) = 0$$

$$\sum_{t=1}^T \phi_t \hat{v}_t^\pi = \sum_{t=1}^T \phi_t \phi_t^\top \theta$$

$$\theta_t = \left(\sum_{t=1}^T \phi_t \phi_t^\top \right)^{-1} \sum_{t=1}^T \phi_t \hat{v}_t^\pi$$

- For N features, direct solution time is $O(N^3)$
- Incremental solution time is $O(N^2)$ using Sherman-Morrison

Convergence of Linear Least Squares Prediction Algorithms

On/Off-Policy	Algorithm	Table Lookup	Linear	Non-Linear
On-Policy	MC	✓	✓	✓
	LSMC	✓	✓	-
	TD	✓	✓	✗
	LSTD	✓	✓	-
Off-Policy	MC	✓	✓	✓
	LSMC	✓	✓	-
	TD	✓	✗	✗
	LSTD	✓	✓	-

Policy case the cross turns into a check mark for LS

{ TD / MC
Double Q-learning
Experience replay

Neural Q-learning

A network $q_\theta : Q_t \Rightarrow (q[1], \dots, q[M])$
(M actions)

An ϵ -greedy exploration policy

$$Q_t \Rightarrow \pi_t \Rightarrow A_t$$

Loss function:

Minimize

$$I(\theta) = \frac{1}{2} \left(R_{t+1} + \gamma \left[\max_a q_\theta(S_{t+1}, a) \right] - q_\theta(S_t, A_t) \right)^2$$

$$\nabla_\theta I(\theta) = \left(R_{t+1} + \gamma \max_a q_\theta(S_{t+1}, a) - q_\theta(S_t, A_t) \right) \nabla_\theta q_\theta(S_t, A_t)$$

Example: TF pseudo-code for Q-learning

```
# Compute Q values Q(S_t, .)
q = q_net(obs)

# Get action A_t
action = epsilon_greedy(q)

# Compute Q(S_t, A_t)
qa = q[action]

# Step in environment
reward, discount, next_obs = env.step(action)

# Get max of values at next state
max_q_next = tf.reduce_max(q_net(next_obs))

# Compute TD-error, do not to propagate into next state value
delta = reward + discount * tf.stop_gradient(max_q_next) - qa

# Define loss
q_loss = tf.square(delta)/2
```


Example: DQN

- ▶ DQN (Mnih et al. 2013, 2015) includes:
 - ▶ A **network** $q_\theta: O_t \mapsto (q[1], \dots, q[m])$ (m actions)
 - ▶ An ϵ -greedy **exploration policy**: $q_t \mapsto \pi_t \implies A_t$
 - ▶ A **replay buffer** to store and sample past transitions
 - ▶ A **target network** $q_{\theta^-}: O_t \mapsto (q^-[1], \dots, q^-[m])$
 - ▶ A Q-learning **loss function** on θ (uses replay and target network)

$$l(\theta) = \frac{1}{2} \left(R_{i+1} + \gamma \left[\max_a q_{\theta^-}(S_{i+1}, a) \right] - q_\theta(S_i, A_i) \right)^2$$

- ▶ An **optimizer** to minimize the loss (e.g., SGD, RMSprop, or Adam)
- ▶ Replay and target networks make RL look more like supervised learning
- ▶ It is unclear whether they are vital, but they help
- ▶ "DL-aware RL"

Multi-step updates

- ▶ When we bootstrap, updates use old estimates
- ▶ Information can propagate back quite slowly
- ▶ In MC information propagates faster, but the updates are noisier
- ▶ We can go in between TD and MC

case if you do TD with a single step you
actually only