

Lecture 6 → RL. Value function approximations

→ scaling RL to real-world problems.

→ Incremental methods → take some function approximator (RVN) and when see some new data, online → update your value function.
→ Batch methods
→ look at set of data / history of data.

Large Scale RL → Backgammon → 10^{20} states
→ Go → 10^{17} states. → Can't build a table
→ want methods to work inspite of large # states. → not practical

eg don't want separate values for some position^x and another position at 1mm to x. → ideally value function should capture / understand that.

→ This is for value functions

→ Scale up for prediction & control.

Value function Approximators.

$Q(s, a)$ → used in the model free regime.

(X) Too many s/a to store in memory

(X) Slow to learn the state's value individually.

parametric function approximator

parameter.
 $\hat{V}(s, w) \approx V_{\pi}(s)$

gives estimate of $V_{\pi}(s)$ for any s .

$\hat{Q}(s, a, w) \approx Q(s, a) \rightarrow$ estimate for any state/action pair

① reduce memory

① generalize. → as states that we haven't seen.

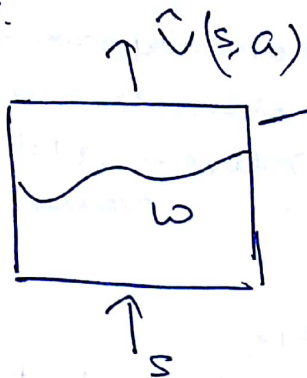
eg. $w \rightarrow$ weights of RVN

→ update using Bellman algo → Monte Carlo
TD learning

Meaning of function approx with a value function

Types

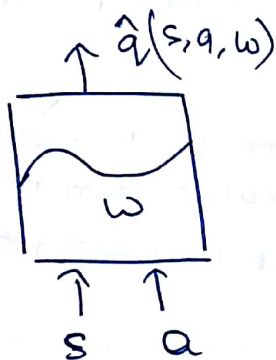
①



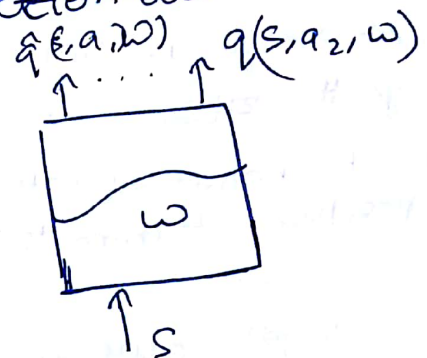
black box function approximator
with parameter $\rightarrow \omega$.
 \rightarrow value for a state s .

②

i) Action-in



ii) Action-out



Using the tunable parameter $\rightarrow \omega$, we get a universal function approximator for all states / (or/and) ~~all~~ action =

Many FAs:

\rightarrow NN, decision-tree, etc.

How to choose?

\rightarrow new focus on differentiable FAs \rightarrow

- Linear combination of features
- Neural Nets.

so that when we change ω , we know how $V(s, \omega)$ affected.

Challenge compared to supervised Learning

\hookrightarrow data is non-iid ^{like} as the $\hat{V}(s, \omega) \rightarrow$ approximating $V_{\pi}(s) \rightarrow$ but the π here is constantly being improved

Gradient Descent

$J(w) \rightarrow$ differentiable

use to find the minimum of $J(w)$

$$\nabla_w J(w) = \begin{bmatrix} \frac{\partial J(w)}{\partial w_1} \\ \vdots \\ \frac{\partial J(w)}{\partial w_n} \end{bmatrix}$$

adjust parameter in the downhill direction

$$\rightarrow \Delta w = -\eta \nabla_w J(w)$$

If \rightarrow by magic, we had actual

$V_\pi(s)$ that should

have been given \rightarrow simply do the mean square error as the cost function and find the minimum.

use SGD \rightarrow as a replacement for (Expectation).
(substitute)

Feature vector

$$\chi(s) = \begin{bmatrix} \chi_1(s) \\ \vdots \\ \chi_n(s) \end{bmatrix}$$

\rightarrow everything that says something about the state. (one thing).

Linear VFA $\rightarrow \hat{v}(s; w) = \chi(s)^T w = \sum_{j=1}^n \chi_j(s) w_j$

⊕ mean square error is quadratic

and hence, SGD converges to optimal global optimum

$$\Delta w = \alpha (V_\pi(s) - \hat{v}(s; w)) \chi(s) \rightarrow \text{adjust only where features are on/active}$$

update = step size \times error \times feature value

Table lookup \rightarrow essentially \rightarrow special case of linear VFA

$$\chi^{\text{table}}(s) = \begin{bmatrix} \mathbb{I}(s_1 = s) \\ \vdots \\ \mathbb{I}(s_n = s) \end{bmatrix} \rightarrow \text{only 1 in non-zero.}$$

$$\hat{v}(s; w) = \chi^{\text{table}}(s)^T w$$

\rightarrow gives the value for that specific state.

But we don't have V_π .

Use MC/ TD to give a target for our VFA.

→ For MC, target is G_t → unbiased estimator for V_π
$$\Delta W = \alpha (G_t - \hat{V}(s, w)) \nabla_w \hat{V}(s, w).$$

→ For TD(λ) → target is G_t^λ
eg. for TD(0) → $R_{t+1} + \gamma \hat{V}(s_{t+1}, w)$

~~MC~~ G_t → unbiased, noisy sample of $V_\pi(s)$.

→ we have training data like SL → $\langle s_1, G_1 \rangle \dots \langle s_T, G_T \rangle$

→ Linear MC policy evaluation

$$\Delta W = \alpha (G_t - \hat{V}(s_t, w)) \chi(s_t).$$

→ converges to local optimum (eg. with non linear VFA)

TD target → biased sample of true $V_\pi(s)$.

→ Still SL with data → $\langle s_1, R_2 + \gamma \hat{V}(s_2, w) \rangle, \dots$
 $\dots \langle s_{T-1}, R_T \rangle.$

$$\begin{aligned} \Delta W &= \alpha (\underbrace{\gamma \hat{V}(s', w) + R - \hat{V}(s, w)}_{\delta}) \nabla_w \hat{V}(s, w) \\ &= \alpha \delta \chi(s). \quad (\text{linear case}). \end{aligned}$$

[All these → incremental online updates].

Linear TD(0) $\xrightarrow{\text{close to}}$ global optimum.

TD(λ) → λ -return G_t^λ

$$\textcircled{1} \langle s_1, G_1^\lambda \rangle \dots \langle s_{T-1}, G_{T-1}^\lambda \rangle$$

$$\textcircled{2} \Delta W = \alpha (G_t^\lambda - \hat{V}(s_t, w)) \chi(s_t).$$

③ Here, Eligibility Trace → ET on the parameters, not the states.


$$\delta_t = R_{t+1} + \gamma \hat{V}(s_{t+1}, w) - \hat{V}(s_t, w)$$

$$E_t = \lambda \gamma E_{t-1} + \chi(s_t)$$

Doubt → gradient of $\hat{V}(s'; \omega)$ not included in optimizing?

$$\Delta \omega = \alpha (R + \hat{V}(s'; \omega) - \hat{V}(s, \omega)) \nabla_{\omega} \hat{V}(s, \omega)$$

↓ why not the gradient used here

But, it sort of means pulling both ends of the spring together → 
close

& it's shown to fail (But other modifications are there → residual gradient).
we need to ground stuff / something & we always do that in the future.

Control with VFA

- Approximate policy evaluator → $\hat{q}(s, a, \omega) \approx q_{\pi}$
 - Policy improvement - ϵ -greedy.
- update on every step → always use the most fresh experience.
- In practice, this goes to very close of the optimal π^* .

Action-VFA → same steps → just using Q instead of V .

Feature vector → $\chi(s, A) = \begin{bmatrix} \chi_1(s, A) \\ \vdots \\ \chi_n(s, A) \end{bmatrix}$

using the update eqn, we either pull up/down the function approximator → to what we actually see → feature coding

too many examples

Bootstrap → helps. $T_{\lambda}(\gamma) \geq TD(0)$. MLC → fails badly.
→ too variance.

TD → does not guarantee that TD would converge.

(Follow Table in slide).
to get where TD can be applied

Gradient TD → really follows the true gradient of the projected Bellman error.

fixes the problems with the TD stability

Why TD unstable?

→ TD does not follow the gradient g any objective function and hence, can diverge.

⇒ Control → problematic

	<u>Linear</u>	<u>Non-linear</u> (VFA)
Q-learning	X	X
Gradient Q-learning	✓	X

→ can keep chattering around the optimal policy.

Batch RL

- Problem with GD → throws away experience after one use. → not sample efficient.
- Find the best-fitting value to the entire batch.

Least Squares Prediction

→ minimize error over batch / entire dataset.

$$D = \{ \langle s_1, v_1^\pi \rangle, \dots, \langle s_T, v_T^\pi \rangle \}$$

Find best LS fit

↳ true value.

$$LS(w) = \sum_{t=1}^T [v_t^\pi - \hat{v}(s_t, w)]^2$$

$$= E_D [(v^\pi - \hat{v}(s, w))^2]$$

Easy way to find the LS → Experience Replay (ER)

- store the entire experience ^{Sol².} D
- sample from experience → $\langle s, v^\pi \rangle \sim D$
- Apply SGD to get that target. → global minimum obtained is the LS sol².
- $w^\pi = \arg \min_w LS(w)$

Idea → sampling repeatedly from experience → making more use of the data

Deep Q-Networks (DQN) → uses ER & fixed Q-targets

→ Big FA → NN

(Q-learning)
i.e. off policy

→ take action a_t → ε-greedy policy

→ store transition $(s_t, a_t, r_{t+1}, s_{t+1})$ in memory D .

→ sample random minibatches (s, a, r, s') from D .
(64).

→ compute Q-learning targets wrt old, fixed params w^- .

→ optimise MSE.

$$L(w_i) = E_{s,a,s',r \sim D} \left[(r + \gamma \max_{a'} Q(s', a', w_i^-) - Q(s, a, w_i))^2 \right]$$

→ variant - g SGD.

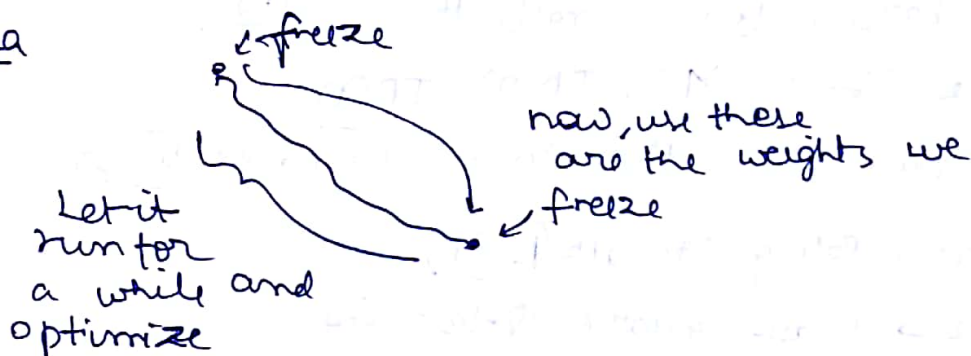
Generally, SARSA (and others) are unstable with SGD.

But this works. → Reasons
① ER. → stabilizes the NN as it decorrelates the trajectory (random sample).

② fixed Q-targets.

2 networks → freeze the old network & freeze the ~~old~~ bootstrapping

Idea



→ End-End Learning of $Q(s, a)$ from pixels s .

→ I/p → s → stack of raw pixels from last 4 frames.

→ O/p → $Q(s/a)$ for 18 possible actions

→ Reward → change in score.

SGD with ER \rightarrow works.
 \rightarrow For Linear FA, can jump directly to the minimum
 (closed form)
 policy evaluates in
1 step.

Σ Linear Least squares prediction:

\rightarrow Wouldn't want any updates at optimal point $\therefore E_D(\Delta w) = 0$

$$\alpha \sum_{t=1}^T x(s_t) (v_t^\pi - x(s_t)^T w) = 0$$

$$\sum_{t=1}^T x(s_t) v_t^\pi = \sum_{t=1}^T x(s_t) x(s_t)^T w$$

$$w = \left(\sum_{t=1}^T x(s_t) x(s_t)^T \right)^{-1} \sum_{t=1}^T x(s_t) v_t^\pi$$

Incremental matrix inverse $\rightarrow O(N^2)$.

Depends on the # parameters now, not the # states

\therefore if # parameter is small, this works!

Can replace v_t^π by MC, TD(0), TD(λ)

Performs much better than incremental \rightarrow for linear

Least squares Policy Iterates (LSPI).

- ① Evaluate \rightarrow Least square Q-learning
- ② Greedy policy ~~evaluation~~ improvement.