# Value Function Approximation [2]

DOROZHKO Anton

Novosibirsk State University

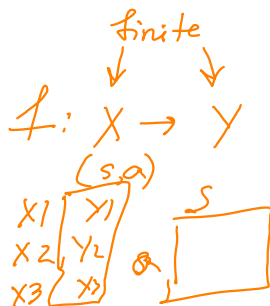May 27, 2020

N\* Novosibirsk
State
University
\*THE REAL SCIENCE

---

[2]David Silver's Lecture 6, Stanford CS231 lecture 5

# Outline

- Previous lecture : Control (makind decision) in Model-Free case
- **This time: Value function approximation**

# Last time: Model-Free Control

- How to learn policy from experience ?
- Tabular representation of $Q(s, a)$ and $V(s)$
- In real world:
  - Backgammon $\sim 10^{20}$
  - Go $\sim 10^{100}$
  - Robotic control - continuous
- Tabular is insufficient



N* Novosibirsk
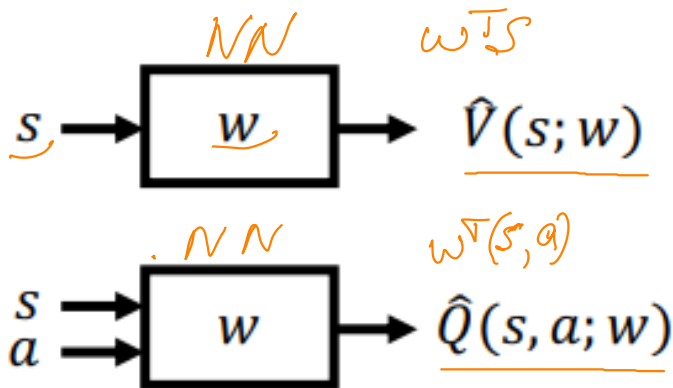State
University
*THE REAL SCIENCE

# Recall: RL Involves

- Optimization
- Delayed consequences
- Exploration
- Generalization

$target, \ opt \ alg$

$G_t = R_t + \gamma$

$\lim_{t \to \infty} N(s, a) \to \infty \quad \forall_{s,a}$

$R_{t+2} \cdots$

# Recall: RL Involves

- Optimization
- Delayed consequences
- Exploration
- **Generalization**

# Value Function Approximation

Represent a (state-action / state) value function with parametrized function
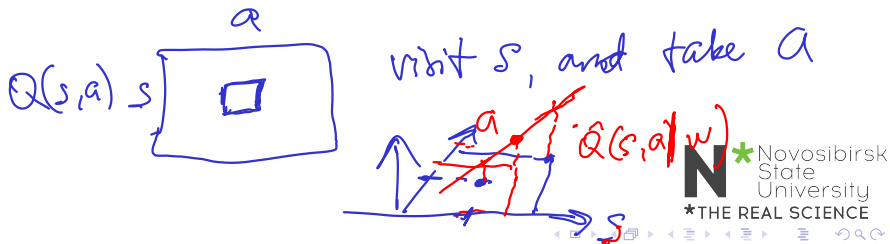
# Motivation for VFA

✓ model based

✗ model free

$$MDP \ (S, \ A, \ P, \ R, \ \gamma)$$

dynamics

- Don't want to explicitly store/learn for every single state
  - Dynamics or reward model  $\hat{P}, \ \hat{R}$
  - Value  $\vee$
  - State-action  $Q$
  - Policy  $\pi(a|s) \sim \pi_\theta(a|s)$
- Want a compact representation that generalizes across states, states-actions
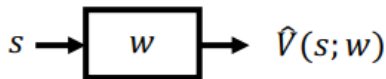
# Benefits of generalization

- Reduce memory needed to store $(P, R)/V/Q/\pi$
- Reduce computation of $(P, R)/V/Q/\pi$
- Reduce experience needed to learn a good $(P, R)/V/Q/\pi$

# Value Function Approximation

- Represent a (state-action / state) value function with parametrized function



$$s \rightarrow \boxed{w} \rightarrow \hat{V}(s; w)$$

$$\begin{matrix} s \\ a \end{matrix} \rightarrow \boxed{w} \rightarrow \hat{Q}(s, a; w)$$

- **Which function approximation ?**

*linear*          *deep RLs*      | ML alg
                   *NN*          RF SVM |

# Function Approximatiors

- Linear combinations of features
- Neural networks
- Decision trees
- Nearest neighbors $\quad kNN$
- Fourier / wavelets

N* Novosibirsk
State
University
*THE REAL SCIENCE

# Gradient Descent

- $J(w)$ - a differentiable function
- Find $w$ that minimizes $J$
- The gradient

$$\nabla_w J(w) = \left[ \frac{\partial J(w)}{\partial w_1}, ..., \frac{\partial J(w)}{\partial w_N} \right]$$

$$w \leftarrow w - \alpha \nabla_w J(w)$$

learning rate

# Stochastic Gradient Descent

*sample*

- Find $w$ that minimizes loss between $V^\pi(s)$ and $\hat{V}(s; w)$

FA    $Q^\pi(s, a)$

- MSE

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

- use gradient descent to find local minumum

$$\Delta w = -\frac{1}{2}\alpha\nabla_w J(w)$$

- SGD samples the gradient

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

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

$\delta$ error

# Model Free VFA Prediction / Evaluation

- Model-free policy evaluation
  - Follow a fixed policy $\pi$
  - Estimate $V^{\pi}(s)$ and/or $Q^{\pi}$
- Maintain lookup table for $V^{\pi}(s)$ and/or $Q^{\pi}$
- Update with MC or TD
- **With VFA: update is fitting of the VFA**

# Feature Vectors


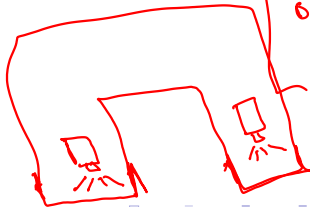
Use feature vector to represent state $s$

$$x(s) = \begin{pmatrix} x_1(s) \\ x_2(s) \\ ... \\ x_n(s) \end{pmatrix}$$

Example: laser sensor, state aliasing

# Linear VFA

- 

$$\hat{V}(s; w) = \sum_{j=1}^{n} x_j(s)w_j = x(s)^T w$$

- Objective function

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

$$\frac{\partial \hat{V}(s, w)}{\partial w} = = X(s)$$

- Weight update:

$$\nabla w \sim -\alpha \nabla_w J(w)$$

- Update = step-size $\times$ prediction error $\times$ feature value

$$\alpha \qquad [V^\pi - \hat{V}] \; X(s)$$

# MC VFA

*cum reward at the end of the episode*

- $G_t$ is unbiased noisy sample of $V^\pi(s_t)$
- Can be used in updates with pairs

$$(s_1, G_1), \quad (s_2, G_2), ..., (s_T, G_T)$$

- Update with linear VFA

$$\Delta w = \alpha(G_t - \hat{V}(s_t; w)\nabla_w \hat{V}(s_t; w)$$
$$= \alpha(G_t - \hat{V}(s_t; w)x(s_t)$$
$$= \alpha(G_t - x(s_t)^T w)x(s_t)$$

$$J(w)$$

- Note: $G_t$ can be very noisy

# MC linear VFA

1: Initialize $\mathbf{w} = \mathbf{0}$, $k = 1$
2: **loop**
3:     Sample $k$-th episode $(s_{k,1}, a_{k,1}, r_{k,1}, s_{k,2}, \ldots, s_{k,L_k})$ given $\pi$
4:     **for** $t = 1, \ldots, L_k$ **do**
5:         **if** First visit to $(s)$ in episode $k$ **then**
6:             $G_t(s) = \sum_{j=t}^{L_k} r_{k,j}$
7:             Update weights:

$$w = w - \alpha \left[ \overset{\hat{v}}{G_t(s)} - X(s)^T w \right] X(s)$$

$$\underbrace{\phantom{G_t(s) - X(s)^T w}}_{error}$$

8:         **end if**
9:     **end for**
10:    $k = k + 1$
11: **end loop**

# Baird(1995)-like example with MC policy Evaluation

$W_0 = [1, 1, 1, ..., 1]$

$\Delta w = [-3, 0, .., 0, -1.5]$

all $r = 0$

$S_2$ $a_1$ $r_1 = 0$

$S_7$ $a_1$ $0$

$S_7$ $a_1$ $0$

terminate

$S1 [1, 0, 0, 0, 0, 0, 0, 0, 1]$ ←



$a_1$ ——

$a_2$ - - -

$w_0 = [1 1 1 1 | 1 1 1 1]$

$w_7 + 2w_8$  $S7$  .99

.01 Terminate

$\alpha = \frac{1}{2}$

$G_t = 0$

- MC update $\Delta w = \alpha(G_t - x(s_t)^T w)x(s_t)$
- With samll prob $s_7$ goes to terminal.
  $x(s_7)^T = [0, 0, 0, 0, 0, 0, 1, 2]$

$V(S_1) = X(S_1)^T W_0 =$

$= 3$

$0.5[0 - 3](2 0 .. 0]$

w1 w2w3      w7 w8

# Convergence Guarantees for linear VFA



- <u>Markov Chain</u> defined by <u>MDP</u> + <u>policy</u> (fixed) will converge to some distribution over states $d(s)$
- $d(s)$ - stationary distribution of $\pi$
- $\sum_s d(s) = 1$
- $d(s)$ satisfies

$$d(s') = \sum_s \left[ \sum_a \pi(a|s) p(s'|s, a) \right] d(s)$$

$$d(\bar{s}) = P\, d(\bar{s})$$

$$p(s \to s')$$
$$p(s'|s)$$

$$d(s) = p\left(\begin{array}{c} \text{being in} \\ \text{state } s \end{array}\right)$$

# Convergence Guarantees for linear VFA [1]

- MSE for linear VFA for $\pi$

$$\mathbb{E}_{\pi}\left[\left(V^{\pi}(s) - \hat{V}^{\pi}(s;w)\right)^2\right]$$

$$MSVE(w) = \sum_{s \in S} d(s)(V^{\pi}(s) - \hat{V}^{\pi}(s;w))^2$$

- where
  - $d(s)$ - stationary distribution
  - $\hat{V}^p i(s;w) = x(s)^T w$ linear VFA
- MC will converge to minimum MSE

$$MSVE(w_{MC}) = min_w \sum_{s \in S} d(s)(V^{\pi}(s) - \hat{V}^{\pi}(s;w))^2$$

---

[1]Tsitsiklis and Van Roy. An Analysis of Temporal-Difference Learning with Funciton Approximation 1997

# Batch Monte-Carlo VFA

- logged episodes from a policy
- Solve analytically

MC error

$$arg \min_w \sum_{i=1}^{N} (\overbrace{G(s_i) - x(s_i)^T w)^2}$$

- result

$$w = (X^T X)^{-1} X^T \mathbf{G}$$

$$G \begin{bmatrix} G(s_1) \\ \vdots \\ G(s_N) \end{bmatrix}$$

$$\begin{bmatrix} x(s_1) \\ \vdots \\ x(s_N) \end{bmatrix}$$

- Note: computationally costly
- Note: no Markov assumptions

# Recall: TD with tables

min err to prediction $G_t :\rightarrow r + \gamma V^{TD}(s')$

↑

→ play some games

- Bootsrap + sampling to approximate $V^\pi$
- Updates $V^\pi(s)$ after each transition $(s, a, r, s')$:

$$V^\pi(s) = V^\pi(s) + \alpha(r + \gamma V^\pi(s') - V^\pi(s))$$

$\Delta V^\pi(s)$

- Target is $r + \gamma V^\pi(s')$ - baised estimate of $V^\pi(s)$
- Represent $V^\pi(s)$ as table

# Recall: TD(0) with VFA

- Bootsrap + sampling to approximate $V^\pi$
- Updates $V^\pi(s)$ after each transition $(s, a, r, s')$:

$$V^\pi(s) = V^\pi(s) + \alpha(r + \gamma V^\pi(s') - V^\pi(s))$$

- Target is $r + \gamma V^\pi(s')$ - baised estimate of $V^\pi(s)$
- In VFA target is $r + \gamma \hat{V}^\pi(s'; w)$
- 3 forms of approximation:

# Recall: TD(0) with VFA

- Bootsrap + sampling to approximate $V^\pi$
- Updates $V^\pi(s)$ after each transition $(s, a, r, s')$:

$$V^\pi(s) = V^\pi(s) + \alpha(r + \gamma V^\pi(s') - V^\pi(s))$$

- Target is $r + \gamma V^\pi(s')$ - baised estimate of $V^\pi(s)$
- In VFA target is $r + \gamma \hat{V}^\pi(s'; w)$
- 3 forms of approximation:
    - function approximation

---

[26] David Silver's Lecture 6, Stanford CS231 lecture 5

# Recall: TD(0) with VFA

- Bootsrap + sampling to approximate $V^\pi$
- Updates $V^\pi(s)$ after each transition $(s, a, r, s')$:

$$V^\pi(s) = V^\pi(s) + \alpha(r + \gamma V^\pi(s') - V^\pi(s))$$

- Target is $r + \gamma V^\pi(s')$ - baised estimate of $V^\pi(s)$
- In VFA target is $r + \gamma \hat{V}^\pi(s'; w)$
- 3 forms of approximation:
    - function approximation
    - bootstrapping

---

[26] David Silver's Lecture 6, Stanford CS231 lecture 5

# Recall: TD(0) with VFA

- Bootsrap + sampling to approximate $V^\pi$
- Updates $V^\pi(s)$ after each transition $(s, a, r, s')$:

$$V^\pi(s) = V^\pi(s) + \alpha(r + \gamma V^\pi(s') - V^\pi(s))$$

*bootstrap*

- Target is $r + \gamma V^\pi(s')$ - baised estimate of $V^\pi(s)$
- In VFA target is $r + \gamma \hat{V}^\pi(s'; w)$
- 3 forms of approximation:

*deadly triad* [
- function approximation ←
- bootstrapping ]
- sampling ] ← some games

---

[26] David Silver's Lecture 6, Stanford CS231 lecture 5

# Recall: TD(0) with VFA

- Bootsrap + sampling to approximate $V^\pi$
- Updates $V^\pi(s)$ after each transition $(s, a, r, s')$:

$$V^\pi(s) = V^\pi(s) + \alpha(r + \gamma V^\pi(s') - V^\pi(s))$$

- Target is $r + \gamma V^\pi(s')$ - baised estimate of $V^\pi(s)$
- In VFA target is $r + \gamma \hat{V}^\pi(s'; w)$
- 3 forms of approximation:
  - function approximation
  - bootstrapping
  - sampling
- Note: we are still on-policy

---

[26] David Silver's Lecture 6, Stanford CS231 lecture 5

# TD(0) with VFA

- Reduce TD(0) to supervised learning with

$$(s_1, \underbrace{r_1 + \gamma \hat{V}^\pi(s_2; w)}_{G_t}), (s_2, r_2 + \gamma \hat{V}^\pi(s_3; w))...$$

- Minimize

$$J(w) = \mathbb{E}_\pi[(r_j + \gamma \hat{V}^\pi(s_{j+1}; w) - \hat{V}^\pi(s_j; w))^2]$$

- Update:

$$\Delta w = \alpha(r + \gamma \hat{V}^\pi(s'; w) - \hat{V}^\pi(s; w)\nabla_w \hat{V}^\pi(s; w)$$
$$= \alpha(r + \gamma \hat{V}(s'; w) - \hat{V}^\pi(s_t; w)x(s)$$
$$= \alpha(\underbrace{r + \gamma x(s')^T w - x(s_t)^T w}_{G_t})x(s)$$

# TD(0) with VFA

1: Initialize $\mathbf{w} = \mathbf{0}$, $k = 1$
2: **loop**
3:     Sample tuple $(s_k, a_k, r_k, s_{k+1})$ given $\pi$     1 step
4:     Update weights:

$$\mathbf{w} = \mathbf{w} + \alpha(r + \gamma \underline{\mathbf{x}(s')^T \mathbf{w}} - \underline{\mathbf{x}(s)^T \mathbf{w}})\mathbf{x}(s)$$

                                           NN           NN

5:     $k = k + 1$
6: **end loop**

# Convergence Guarantees for linear VFA [1]

- MSE for linear VFA for $\pi$

$MDP + \pi \Rightarrow$ Markov chain

$$MSVE(w) = \sum_{s \in S} d(s)(V^\pi(s) - \hat{V}^\pi(s;w))^2$$

$\downarrow$

$d(s)$

- where
  - $d(s)$ - stationary distribution
  - $\hat{V}^p i(s;w) = \underline{x(s)}^T w$ linear VFA
- MC will converge to minimum MSE

TD

$\gamma = 0$

$G_t = R_t$

$$MSVE(w_{TD}) \le \frac{1}{1-\gamma} \min_w \sum_{s \in S} d(s)(V^\pi(s) - \hat{V}^\pi(s;w))^2$$

---

[1] Tsitsiklis and Van Roy. An Analysis of Temporal-Difference Learning with Funciton Approximation 1997

# Control Using VFA

- use VFA to represent state-action values $\hat{Q}^\pi(s, a; w) \approx Q^\pi$
- Interleave:
  - Approximate policy eval using VFA
  - $\epsilon$-greedy improvement ?
- Can be unstable (Deadly Triad)
  - Function approximation — DRL
  - Bootstrapping
  - Off-policy learning

$$(s, a, r, s') \sim \beta$$

$$a \sim \beta \qquad a' \text{ at } s' \sim \pi$$

greedy policy impr.

$$\pi = \underset{a}{\arg\max} \, Q(s,a)$$

$\epsilon$-greedy

$$\pi = \begin{cases} \underset{a}{\arg\max} Q & p = 1-\epsilon \\ \text{random } a & p = \epsilon \end{cases}$$

# Action-Value approximation

- Find $w$ that minimizes loss between $Q^\pi(s, a)$ and $\hat{Q}(s, a; w)$
- MSE

$$J(w) = \mathbb{E}_\pi[(Q^\pi(s, a) - \hat{Q}(s, a; w))^2]$$
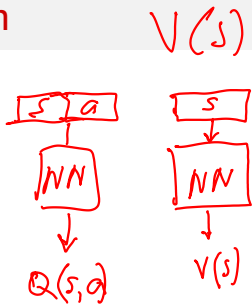
- SGD samples the gradient

$$-\frac{1}{2}\nabla_w J(w) = \mathbb{E}_\pi[(Q^\pi(s, a) - \hat{Q}(s, a; w))\nabla_w \hat{Q}^\pi(s, a; w)]$$

$$\Delta w = -\frac{1}{2}\alpha \nabla_w J(w)$$

# Linear $Q(s, a)$ function approximation

$V(s)$

- Freatures represet both state and action

$S$   $A$

$$x(s, a) = \begin{pmatrix} x_1(s, a) \\ x_2(s, a) \\ ... \\ x_n(s, a) \end{pmatrix}$$



- State-action function as weighted linear combination of features

$$\hat{Q}(s, a; w) = x(s, a)^T w = \sum_{j=1}^{n} x_j(s, a) w_j$$

- SGD

$$\nabla_w J(w) = \nabla_w \mathbb{E}_\pi[(Q^\pi(s, a) - \hat{Q}^\pi(s, a; w))^2]$$

# Incremental Model-Free Control

- MC

  MC target

  $$\Delta w = \alpha(G_t - \hat{Q}^\pi(s_t, a_t; w))\nabla_w \hat{Q}(s_t, a_t; w)$$

- SARSA

  on-policy

  $$\Delta w = \alpha(r + \gamma \hat{Q}(s', a'; w) - \hat{Q}^\pi(s_t, a_t; w))\nabla_w \hat{Q}(s_t, a_t; w)$$

- Q-learning

  off-policy

  $$\Delta w = \alpha(r + \gamma \max_{a'} \hat{Q}(s', a'; w) - \hat{Q}^\pi(s_t, a_t; w))\nabla_w \hat{Q}(s_t, a_t; w)$$

  TD target / bootstrapped

# Convergence of TD with VFA

$$\hat{Q}(s, a) = r + \gamma \overset{max}{\hat{Q}(s', a)}$$

reward signal from env

target

- Sutton and Barto Ch. 11
- TD is not following gradient of $J(w) = \mathbb{E}\left[(Q - \hat{Q}(s, a|w))^2\right]$
- Bellman Operator are contractions, but VFA fitting can be an expansion

# Convergence of Prediction Algorithms

$\widehat{V}(s)$

form of function $\widehat{V}$

| On/Off-Policy | Algorithm | Table Lookup | Linear | Non-Linear |
|---|---|---|---|---|
| On-Policy | MC | ✓ | ✓ | ✓ |
| | TD | ✓ | ✓ | ✗ |
| | Gradient TD | ✓ | ✓ | ✓ |
| Off-Policy | MC | ✓ | ✓ | ✓ |
| | TD | ✓ | ✗ | ✗ |
| | Gradient TD | ✓ | ✓ | ✓ |

$X(s)^{\top} w$

# Convergence of Control Algorithms

$$Q(s,a) = X(s,a)^T W$$

| Algorithm | Table Lookup | Linear | Non-Linear |
|---|---|---|---|
| Monte-Carlo Control | ✓ | (✓) | ✗ |
| Sarsa | ✓ | (✓) | ✗ |
| Q-learning | ✓ | ✗ | ✗ |
| Gradient Q-learning | ✓ | ✓ | ✗ |

(✓) = chatters around near-optimal value function

off policy

# What You Should Understand

*Tabular V, Q $\rightarrow$ func approx*

- Implement TD(0) and MC on policy evaluation with <u>linear function approximation</u>

- Be able to implement <u>Q-learning</u>, SARSA and MC control with function approximation

- List the 3 issues that can cause instability
  - Function approximation  *NN*
  - Bootstrapping  *fit on prediction*
  - Off-policy learning

*deadly triad*

*easier to collect more experience*

Explain project details