

# RL with function approximation

RLFA-1

$\Rightarrow n_t, g_t \quad t \geq 0$  = state/cst realizations for  $n_{t+1} = f(n_t, \mu(n_t), w_t)$   
 when MDP is operated with policy  $\mu(\cdot)$  (stochastic)  
 $g_t = g(n_t, \mu(n_t), w_t)$

we want to estimate

$$V^{\mu}(n), \quad n \in \mathcal{X} \quad Q^{\mu}(n, u) \quad n, u \in \mathcal{X} \times \mathcal{U} \quad \xrightarrow{\text{then}} \text{use GPI}$$

however:  $|\mathcal{X}| = m$  and  $|\mathcal{U}| = m$  are too big for an exact representation  
 of  $V^{\mu}(n)$  and  $Q^{\mu}(n, u)$   $\hookrightarrow$  chess no. of board configurations  
 $> 10^{100}$

$\Rightarrow$  use function approximation

$$V^{\mu}(n) \approx \hat{V}(n, \vartheta) \quad Q^{\mu}(n, u) \approx \hat{Q}(n, u, \vartheta) \quad \vartheta \in \mathbb{R}^9$$

requirements: store  $\vartheta$  and evaluate  $\hat{V}(n, \vartheta)$  and  $\hat{Q}(n, u, \vartheta)$   
 for given  $n$  and  $u$

linear:  $\hat{V}(n, \vartheta) = \varphi(n)^T \vartheta$        $\varphi(n) \quad n \in \mathcal{X}$  = basis function

$$\hat{Q}(n, u, \vartheta) = \varphi(n, u)^T \vartheta \quad \varphi(n, u) \quad n, u \in \mathcal{X} \times \mathcal{U}$$

$\hookrightarrow$  less modeling capabilities, strong convergence properties

tabular case as a particular case:

$$\vartheta \in \mathbb{R}^m \quad \varphi(i) = \begin{bmatrix} 1 \\ 0 \\ \vdots \\ 0 \\ 0 \end{bmatrix} - i\text{-th} \Rightarrow \varphi(i)^T \vartheta = \vartheta_i = V^{\mu}(i)$$

mutatis mutandis for  $Q^{\mu}(n, u)$

non-linear

e.g.  $\hat{V}(n, \vartheta)$  and  $\hat{Q}(n, u, \vartheta)$  = Neural Network

$\hookrightarrow$  convergence critical, strong modeling power

$V^{\mu}(n)$  to begin with ( $Q^{\mu}(n, u)$  conceptually identical)

Approximation achieved via projection onto  $\mathcal{M} = \{\hat{V}(n, \vartheta) : \vartheta \in \mathbb{R}^q\}$

$\Rightarrow$  Given  $V(n)$ , let

$$\Pi[V(n)] = \underset{\hat{V} \in \mathcal{M}}{\operatorname{argmin}} \mathbb{E}\left[\left(V(n_t) - \hat{V}(n_t, \vartheta)\right)^2\right]$$

expectation taken w.r.t  
to the stationary  
distribution of  $n_t$  when  
MDP operated with  $\mu(\cdot)$

$$\hat{V}(n, \bar{\vartheta}) \text{ where } \bar{\vartheta} = \underset{\vartheta}{\operatorname{argmin}} \sum_{i=1}^M \xi_i (V(i) - \hat{V}(i, \vartheta))^2$$

$\hookrightarrow \xi_i = \Pr\{n_t = i\} \quad t \rightarrow \infty$

### DIRECT (BATCH) APPROACH

goal: compute  $\Pi[V^{\mu}(n)]$

$$\mathbb{E}\left[\left(V^{\mu}(n_t) - \hat{V}(n_t, \vartheta)\right)^2\right] \approx \frac{1}{T} \sum_{t=0}^{T-1} \left(V^{\mu}(n_t) - \hat{V}(n_t, \vartheta)\right)^2$$

↑ available visitation  $n_1, n_2, \dots$

$$V^{\mu}(n_e) \approx \sum_{e=t}^{T-1} \gamma^{e-t} g_e \approx \sum_{i=0}^{+\infty} \gamma^i g_{t+i}$$

$$\Rightarrow \bar{\vartheta} = \underset{\vartheta}{\operatorname{argmin}} \frac{1}{T} \sum_{t=0}^{T-1} \left(\hat{V}(n_t, \vartheta) - \sum_{i=0}^{T-1-t} \gamma^i g_{t+i}\right)^2$$

► can be solved via a gradient method

$$\hat{\vartheta}_{n+1} = \hat{\vartheta}_n - \alpha \sum_{t=0}^{T-1} \nabla_{\vartheta} \hat{V}(n_t, \hat{\vartheta}_n) \left( \hat{V}(n_t, \hat{\vartheta}_n) - \sum_{i=0}^{T-1-t} \gamma^i g_{t+i} \right) \rightsquigarrow \begin{array}{l} \text{stochastic gradient} \\ \text{can be used} \end{array}$$

↑ to be modified

stochastic gradient: case

$$\hat{\vartheta}_{K+1} = \hat{\vartheta}_n - \alpha \nabla_{\vartheta} \hat{V}(n_n, \hat{\vartheta}_n) \left( \hat{V}(n_n, \hat{\vartheta}_n) - \sum_{i=0}^{T-1-K} \gamma^i g_{K+i} \right)$$

↑ the whole batch is still needed → we short T and change batch memory

OBS: we can write

$$\begin{aligned} \hat{\vartheta}_{n+1} &= \hat{\vartheta}_n - \alpha \left( \sum_{t=0}^{T-1} \nabla_{\vartheta} \hat{V}(n_t, \hat{\vartheta}_n) \hat{V}(n_t, \hat{\vartheta}_n) - \sum_{t=0}^{T-1} \sum_{e=t}^{T-1} \nabla_{\vartheta} \hat{V}(n_t, \hat{\vartheta}_n) \gamma^{e-t} g_e \right) \\ &= \hat{\vartheta}_n - \alpha \left( \sum_{t=0}^{T-1} \nabla_{\vartheta} \hat{V}(n_t, \hat{\vartheta}_n) \hat{V}(n_t, \hat{\vartheta}_n) - \sum_{t=0}^{T-1} \sum_{e=0}^t \nabla_{\vartheta} \hat{V}(n_e, \hat{\vartheta}_n) \gamma^{t-e} g_t \right) \end{aligned}$$

↑ incremental approach  
incremental update:  
can be used with very a long batch

stochastic gradient: case

$$\hat{\vartheta}_{n+1} = \hat{\vartheta}_n - \alpha \left( \nabla_{\vartheta} \hat{V}(n_n, \hat{\vartheta}_n) \hat{V}(n_n, \hat{\vartheta}_n) - \sum_{e=0}^K \underbrace{\nabla_{\vartheta} \hat{V}(n_e, \hat{\vartheta}_n)}_{\hat{\vartheta} \rightarrow \hat{\vartheta}_0} \gamma^{n-e} g_n \right)$$

↑ incremental update

# INDIRECT APPROACHES $\Rightarrow$ projected Bellman equation

RLFA - 3

$$V^{\mu}(u) : V^{\mu}(u) = T^{\mu}[V^{\mu}(u)] = E\left[g(u, \mu(u), w) + \gamma V^{\mu}(f(u, \mu(u), w))\right]$$

We consider a projection of this equation over the space  $\mathcal{M}$

$$\hat{V}(u, \bar{\theta}) = \Pi_{\downarrow} [T^{\mu}[\hat{V}(u, \bar{\theta})]] \quad \hat{V}(u, \bar{\theta}) \in \mathcal{M} \quad T^{\mu}[\hat{V}(u, \bar{\theta})] \notin \mathcal{M}$$

projection over  $\mathcal{M}$

more explicitly  $\hat{V}(u, \bar{\theta}) \in \mathcal{M}$  w.r.t.  $u_t$  w.r.t.  $w_t$

$$\hat{V}(u, \bar{\theta}) = \underset{\hat{V}(u, \bar{\theta}) \in \mathcal{M}}{\operatorname{argmin}} \mathbb{E}_{u_t} \left[ \left( \hat{V}(u_t, \bar{\theta}) - \mathbb{E}_{w_t} \left[ g(u_t, \mu(u_t), w_t) + \gamma \hat{V}(f(u_t, \mu(u_t), w_t), \bar{\theta}) \right] \right)^2 \right]$$

$\Updownarrow$  if  $\hat{V}(\cdot, \bar{\theta})$  one to one ... if  $\hat{V}(u, \bar{\theta}) = \varphi(u)^T \bar{\theta}$  then  $\varphi(u)$  linear ind.

$$\bar{\theta} = \underset{\varphi}{\operatorname{argmin}} \mathbb{E} \left[ \left( \hat{V}(u_t, \bar{\theta}) - \left( g(u_t, \mu(u_t), w_t) + \gamma \hat{V}(f(u_t, \mu(u_t), w_t), \bar{\theta}) \right) \right)^2 \right]$$

[ $\varphi(u)^T \cdot \varphi(w)$   
full name

assume the projected equation has one and only one solution  
(under certain condition, it does  $\rightarrow$  see earlier)

then  $\bar{\theta}$  := solution to  $\nabla_{\theta} (\mathbb{E}[(\dots)^2]) = 0$  i.e. it must be

$$\mathbb{E} \left[ \nabla_{\theta} \hat{V}(u_t, \bar{\theta}) \cdot \left( \hat{V}(u_t, \bar{\theta}) - \left( g(u_t, \mu(u_t), w_t) + \gamma \hat{V}(f(u_t, \mu(u_t), w_t), \bar{\theta}) \right) \right) \right] = 0$$

$\Downarrow$

can be written as a fixed point problem: find  $\bar{\theta}$  such that

$$\bar{\theta} = \bar{\theta} - \mathbb{E} \left[ \nabla_{\theta} \hat{V}(u_t, \bar{\theta}) \cdot \left( \hat{V}(u_t, \bar{\theta}) - \left( g(u_t, \mu(u_t), w_t) + \gamma \hat{V}(f(u_t, \mu(u_t), w_t), \bar{\theta}) \right) \right) \right]$$

$\Downarrow$   
can be obtained via stochastic approximation  
+ bootstrapping

$u_t, g_t \quad t=0, 1, \dots$  realization of state and cost

$\hat{\theta}_t$  = estimate of  $\bar{\theta}$  at time  $t$

given  $\hat{\theta}_t$ , new target becomes

bootstrapped

stochastic realization of  
the right-hand side

$$\hat{\theta}_t + \nabla_{\theta} \hat{V}(u_t, \hat{\theta}_t) \left( g_t + \gamma \hat{V}(u_{t+1}, \hat{\theta}_t) - \hat{V}(u_t, \hat{\theta}_t) \right)$$

or

and direction from  $\hat{\vartheta}_t$  towards target:

$$\nabla_{\vartheta} \hat{V}(x_t, \hat{\vartheta}_t) (g_t + \gamma \hat{V}(x_{t+1}, \hat{\vartheta}_t) - \hat{V}(x_t, \hat{\vartheta}_t)) \rightarrow \text{Temporal Difference}$$

TD(0) algorithm:

$$\begin{cases} \delta_t = g_t + \gamma \hat{V}(x_{t+1}, \hat{\vartheta}_t) - \hat{V}(x_t, \hat{\vartheta}_t) \\ \hat{\vartheta}_{t+1} = \hat{\vartheta}_t + \alpha_t \delta_t \cdot \nabla_{\vartheta} \hat{V}(x_t, \hat{\vartheta}_t) \end{cases}$$

In the linear case  $\hat{V}(x, \vartheta) = \varphi(x)^T \vartheta \Rightarrow \nabla_{\vartheta} \hat{V}(x, \vartheta) = \varphi(x)$

$$\begin{cases} \delta_t = g_t + \gamma \varphi(x_{t+1})^T \hat{\vartheta}_t - \varphi(x_t)^T \hat{\vartheta}_t \\ \hat{\vartheta}_{t+1} = \hat{\vartheta}_t + \alpha_t \delta_t \varphi(x_t) \end{cases} \quad \begin{array}{l} \text{when } \varphi(i) = \begin{bmatrix} 0 \\ 0 \\ \vdots \\ 1_{x_i=1} \\ 0 \end{bmatrix} - i\text{-th} \\ \varphi(x_t) = \begin{bmatrix} 1_{x_1=1} \\ \vdots \\ 1_{x_t=1} \end{bmatrix} \Rightarrow \text{back to the} \\ \text{tabular} \\ \text{case!} \end{array}$$

In the linear case, projected Bellman equation has indeed one and only one solution!  $\Rightarrow$  convergence in TD(0) well-established

►  $\varphi(x)^T \bar{\vartheta} = \Pi[\bar{T}^u(\varphi(x)^T \bar{\vartheta})]$  idea:  $\Pi[\bar{T}^u(\cdot)]$  is a contraction  
 $\downarrow$   
 projection onto a linear subspace = orthogonal projection

$$\begin{aligned} \Pi[V(x)] &= \varphi(x)^T \cdot \underset{\vartheta}{\operatorname{argmin}} \mathbb{E}[(\varphi(x)^T \vartheta - V(x))^2] = && \text{linear operator} \\ &= \varphi(x)^T \cdot \underbrace{\mathbb{E}[\varphi(x) \varphi(x)^T]}_{\text{exists and is unique if } \varphi(x) \text{ are linear independent}}^{-1} \cdot \mathbb{E}[\varphi(x) V(x)] && \text{linear in } V(x) \end{aligned}$$

exist and is unique if  $\varphi(x)$  are linear independent + some assumptions on the stationary distribution of  $x_t$

Let  $\|V(x)\|_S^2 = \mathbb{E}[V(x)^2]$  = same norm (and scalar product) used to define the projection  $\Pi$   
 $\uparrow$   
 stationary distribution of  $x_t$

Then, given  $V_1(x)$  and  $V_2(x)$  we have:

$$\|\Pi[V_1(x)] - \Pi[V_2(x)]\|_S^2 = \|\Pi[V_1(x) - V_2(x)]\|_S^2 \leq$$

$$\leq \|\pi[V_1(n) - V_2(n)]\|_S^2 + \|V_1(n) - V_2(n) - \pi[V_1(n) - V_2(n)]\|_S^2 \quad \text{RLFA.5}$$

$$= \|V_1(n) - V_2(n)\|_S^2$$

orthogonal  $\Rightarrow$  use Pythagora's Th.

Hence,  $\|\pi[V_1(n)] - \pi[V_2(n)]\|_S \leq \|V_1(n) - V_2(n)\|_S$  and

$$\|\pi[T^u(V_1(n))] - \pi[T^u(V_2(n))]\|_S \leq \|T^u(V_1(n)) - T^u(V_2(n))\|_S$$

$\leq \beta \|V_1(n) - V_2(n)\|_S$  for issue  $T^u(\cdot)$  contractive w.r.t. to  
 $\beta < 1$   
 the max norm, also for  $\|\cdot\|_S$ ,  
 under some condition on the stationary  
 distribution + linear independence of  $\varphi(n)$

$\Downarrow$  projected Bellman eq always has a unique solution  
 and convergence of TD( $\alpha$ ) is guaranteed

In the linear case, alternatives to solve approximately  
 the projected Bellman equation exist!

$$\bar{\vartheta} = \underset{\vartheta}{\operatorname{argmin}} \mathbb{E} \left[ (\varphi(u_t)^\top \vartheta - (g(u_t, \mu(u_t), w_t) + \gamma \varphi(f(u_t, M(u_t), w_t))^\top \bar{\vartheta})^2 \right]$$

$$\approx \underset{\vartheta}{\operatorname{argmin}} \frac{1}{T} \sum_{t=0}^{T-1} \left( \varphi(u_t)^\top \vartheta - g_t - \gamma \varphi(u_{t+1})^\top \bar{\vartheta} \right)^2$$

Least Squares  
 Temporal Difference       $\underbrace{\text{data driven approximation}}$ , over a realization  $u_t, g_t \in \mathbb{R}^{n_u, n_g}$

LSTD: solve the Least Square problem (quadratic in  $\vartheta$ )

$$\hat{\vartheta} := \frac{1}{T} \sum_{t=0}^{T-1} \varphi(u_t) \cdot \left[ (\varphi(u_t)^\top - \gamma \varphi(u_{t+1})^\top) \bar{\vartheta} - g_t \right] = 0$$

$$\Rightarrow \hat{\vartheta}_T = \left[ \sum_{t=0}^{T-1} \varphi(u_t) (\varphi(u_t)^\top - \gamma \varphi(u_{t+1})^\top) \right]^{-1} \times \sum_{t=0}^{T-1} \varphi(u_t) \cdot g_t$$

recursive implementation are possible so as to avoid  
 memory saturation and computational issues

easiest

$$\begin{cases} C_t = C_{t-1} + \varphi(\pi_t) (\varphi(\pi_t)^T - \gamma \varphi(\pi_{t+1})) \\ d_t = d_{t-1} + \varphi(\pi_t) g_t \\ \hat{\vartheta}_{t+1} = C_t^{-1} \cdot d_t \end{cases}$$

$\rightarrow$  refined version that avoid inversion exist

LSPE: use value iteration for the projection operator

$$\begin{aligned} \hat{\vartheta}_{t+1} &= \underset{\vartheta}{\operatorname{argmin}} \mathbb{E} \left[ \left( \varphi(\pi_t)^T \vartheta - (g(\pi_t, \mu(\pi_t), w_t) + \gamma \varphi(f(\pi_t, \mu(\pi_t), w_t))^T \hat{\vartheta}_t) \right)^2 \right] \\ &\approx \underset{\vartheta}{\operatorname{argmin}} \frac{1}{t+1} \sum_{j=0}^t \left( \varphi(\pi_j)^T \vartheta - g_j - \gamma \varphi(\pi_{j+1})^T \hat{\vartheta}_t \right)^2 \\ &= \left[ \sum_{j=0}^t \varphi(\pi_j) \varphi(\pi_j)^T \right]^{-1} \cdot \sum_{j=0}^t \varphi(\pi_j) \left[ g_j + \gamma \varphi(\pi_{j+1})^T \hat{\vartheta}_t \right] \\ &\quad \pm \varphi(\pi_j)^T \hat{\vartheta}_t \\ &= \hat{\vartheta}_t + \left[ \sum_{j=0}^t \varphi(\pi_j) \varphi(\pi_j)^T \right]^{-1} \cdot \sum_{j=0}^t \varphi(\pi_j) \left( g_j + \gamma \varphi(\pi_{j+1})^T \hat{\vartheta}_t - \varphi(\pi_j)^T \hat{\vartheta}_t \right) \end{aligned}$$

TD

$\downarrow$   
possibly  
changed to  $\hat{\vartheta}_j$   
 $\rightarrow$  to favor an  
incremental  
implementation

Least Square  
Policy Evolution

TD( $\lambda$ ), LSTD( $\lambda$ ), LSPE( $\lambda$ ) are also possible

$\rightarrow$  repeat everything with  $(1-\lambda) \sum_{n=1}^{+\infty} \lambda^{n-1} (\mathbf{T}^n)^n$  in place  
of  $\mathbf{T}^n$   $\rightarrow$  cumbersome calculations, yet easy to  
implement via eligibility traces

TD( $\lambda$ ):

$$\begin{cases} \delta_t = g_t + \gamma \hat{V}(\pi_{t+1}, \hat{\vartheta}_t) - \hat{V}(\pi_t, \hat{\vartheta}_t) \\ z_t = \nabla_{\vartheta} \hat{V}(\pi_t, \hat{\vartheta}_t) + \gamma \lambda z_{t-1} \\ \hat{\vartheta}_{t+1} = \hat{\vartheta}_t + \alpha_t \delta_t z_t \end{cases}$$

LSTD( $\lambda$ ):

$$z_t = \varphi(\pi_t) + \gamma \lambda z_{t+1}$$

$$\hat{\vartheta}_{t+1} = \left[ \sum_{j=0}^t z_j \left( \varphi(\pi_j)^\top - \gamma \varphi(\pi_{j+1})^\top \right) \right]^{-1} \cdot \sum_{j=0}^t z_j \cdot g_j$$

LSPE( $\lambda$ ):

$$z_t = \varphi(\pi_t) + \gamma \lambda z_{t+1}$$

$$\hat{\vartheta}_{t+1} = \hat{\vartheta}_t + \left[ \sum_{j=0}^t \varphi(\pi_j) \varphi(\pi_j)^\top \right]^{-1} \cdot \sum_{j=0}^t z_j \left( g_j + \gamma \varphi(\pi_{j+1}) \hat{\vartheta}_t - \varphi(\pi_j)^\top \hat{\vartheta}_t \right)$$

Everything can be repeated for  $\hat{Q}(\pi, u, \vartheta)$   
without conceptual changes

$$\hat{V}(\pi_t, \vartheta) \rightarrow \hat{Q}(\pi_t, u_t, \vartheta)$$

$$\hat{V}(\pi_{t+1}, \vartheta) \rightarrow \hat{Q}(\pi_{t+1}, \mu(\pi_{t+1}), \vartheta) = \hat{Q}(\pi_{t+1}, u_{t+1}, \vartheta) \quad (\text{on-policy})$$

$$\varphi(\pi_t) \rightarrow \varphi(\pi_t, u_t)$$

$$\varphi(\pi_{t+1}) \rightarrow \varphi(\pi_{t+1}, \mu(\pi_{t+1})) = \varphi(\pi_{t+1}, u_{t+1}) \quad (\text{on policy})$$

$\Rightarrow$  estimating  $\hat{Q}(\pi, u, \vartheta)$  for model-free policy improvement