

MATHEMATICAL FOUNDATIONS OF REINFORCEMENT LEARNING

LECTURE IV: APPROXIMATE DYNAMIC PROGRAMMING

- SUMMARY:
- COMPUTING OPTIMAL FEEDBACKS ONLINE
 - APPROXIMATE VALUE ITERATION
 - APPROXIMATE POLICY ITERATION
 - BELLMAN EQN. METHODS.

REFS: - CHAPTER 5, BERTSEKAS BOOK ON RL.

- A. ALLA, M. FALCONE, D.K. "AN EFFICIENT POLICY ITERATION ALGORITHM FOR DYNAMIC PROGRAMMING EQUATIONS", SISC, 2015.

PREVIOUS LECTURES:

INFINITE HORIZON COST:

$$J^{\mu}(i) = \lim_{N \rightarrow \infty} \mathbb{E} \left[\sum_{k=0}^N \alpha^k g(i_k, \mu, i_{k+1}) \mid i_0 = i \right]$$

discount (pointing to α^k)
running cost (pointing to $g(i_k, \mu, i_{k+1})$)

VALUE FUNCTION: $J^*(i) = \min_{\mu} J^{\mu}(i)$

V.F. SATISFIES BELLMAN EQUATION:

$$J^*(i) = \min_{\mu \in U(i)} \sum_{j \in X} p_{ij}(\mu) (g(i, \mu, j) + \alpha J^*(j))$$

valid $\forall i \in X$

WITH BELLMAN OPERATOR

$$TJ := \min_{\mu \in U(i)} \sum_{j \in X} \bar{p}_{ij}(\mu) (\bar{g}(i, \mu, j) + \alpha J(j)) \Rightarrow J^* = TJ^*$$

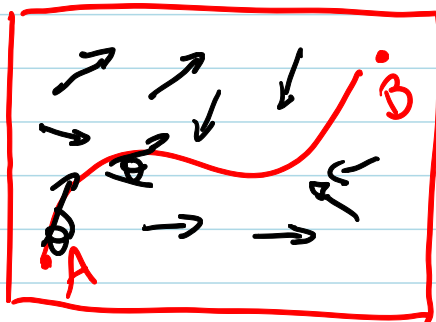
T IS A CONTRACTION MAP \Rightarrow FIXED POINT ITERATION

$$J^*(i) = \lim_{k \rightarrow \infty} J^k(i) \quad \text{WHERE } J^{k+1} = T J^k$$

ONCE J^* HAS BEEN COMPUTED, (μ^*) IS RECOVERED AS

$$\mu^*(i) \in \operatorname{argmin}_{u \in U(i)} \sum_j P_{ij}(u) (g(i; u, j) + \alpha \underline{J^*}(j))$$

EXAMPLE: ZERNELO NAVIGATION PROBLEM



GO FROM A TO B IN MIN TIME SUBJECT TO WIND

DISCRETE EQUS OF MOTION

$$x_{k+1} = x_k + \Delta t (WIND_x + \bar{V} \cos(\theta_k))$$

$$y_{k+1} = y_k + \Delta t (WIND_y + \bar{V} \sin(\theta_k))$$

CONTROL

STATE

(x_k, y_k)

SOME IDEAS FOR ANNS IN THE CONTROL FIELD

$$\mu^*(i) \in \arg \min_u \sum_j p_{ij}(u) (g(i, u, j) + \alpha \underbrace{J^*(j)}_{\downarrow})$$

REPLACE BY A LOW
COMPLEXITY ANN
 $\tilde{J}(j, \rho)$

EXAMPLE: LQ CONTROL

$$x_{k+1} = ax_k + bu_k, \quad x_k, u_k \in \mathbb{R}$$

$$g(x, u, y) = K u^2 + x^2, \quad K \text{ constant}$$

ASSUME A LINEAR ARCHITECTURE := $\tilde{J}(i, \rho) = \rho_0 + \rho_1 x + \rho_2 x^2$

$$\mu^*(x_i) \in \arg \min_u \underbrace{K u^2 + x_i^2 + \bar{\rho}_0 + \bar{\rho}_1 (ax_i + bu) + \bar{\rho}_2 (ax_i + bu)^2}_{g(u)} \Rightarrow g'(u) = 0$$

How do we incorporate the approx. of $J^*(i)$ by $\tilde{J}(i, \pi)$

↳ APPROXIMATE VALUE ITERATION:

$$\hat{J}_{k+1}(i) = \min_{\mu} \sum_j \varphi_{ij}(\mu) (g(i, \mu, j) + \alpha \tilde{J}(j, \pi_k))$$

$\forall i \in S_k$, THEN R_{k+1} IS UPDATED ACCORDING TO

$$\pi_{k+1} \in \arg \min_{\pi} \sum_{i \in S_k} \|\hat{J}_{k+1}(i) - \tilde{J}(i, \pi)\|^2$$

↳ DEFINE ARCHITECTURE $\tilde{J}(i, \pi)$ ✓

↳ INITIAL PARAMETRIC GUESS π_0 ✓

↳ UPDATE \hat{J}_{k+1} USING $\tilde{J}(i, \pi)$ OVER A SUBSET S_k

↳ USE $\hat{J}_{k+1}(i)$ TO TRAIN π_{k+1}

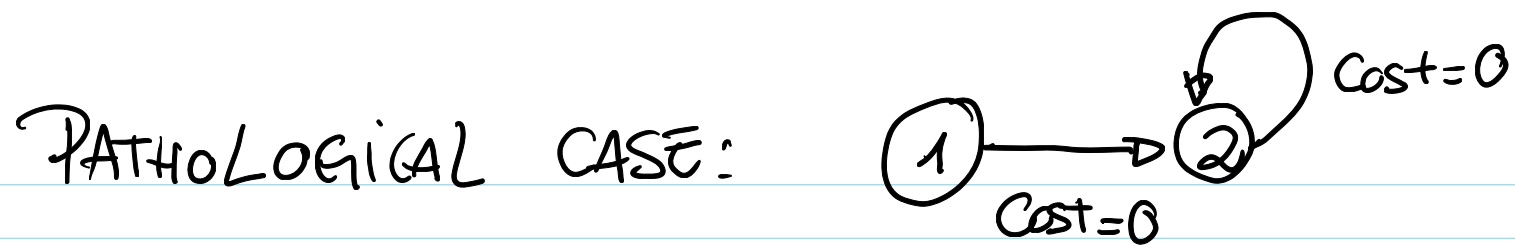
OBS: - THE ARCHITECTURE $\tilde{J}(i, r)$ REMAINS, ONLY r CHANGES AFTER EVERY ITERATION.

- THE EVALUATION OF P_{ij} 'S CAN BE COSTLY UNLESS $\# p_{ij} \neq 0$ IS SMALL.

OBS: AVI IS VERY SIMILAR VI, BUT IT IS CRITICAL TO ENSURE THAT

$$\| \tilde{J}(i, r_{k+1}) - \tilde{J}_{k+1}(i) \| \leq \delta$$

WITH δ SUFFICIENTLY SMALL, OTHERWISE THE CONTRACTIVITY OF THE BELLMAN OPERATOR IS LOST



$$\Rightarrow p_{11} = 0, p_{12} = 1, p_{22} = 1, p_{21} = 0$$

$g \equiv 0$, $\alpha > 0$ discount factor

It's clear that $J^*(1) = J^*(2) = 0$ (because $g \equiv 0$)

USE BELLMAN EQN. $J^*(i) = \sum_j p_{ij} (0 + \alpha J^*(j))$

BELLMAN

$$\begin{cases} J^*(1) = \alpha J^*(2) \\ J^*(2) = \alpha J^*(2) \end{cases}$$

VALUE ITERATION

$$J_{k+1}(1) = \alpha J_k(2)$$

$$J_{k+1}(2) = \alpha J_k(2)$$

ASSUME A LINEAR APPROXIMATION $\tilde{J}(i; r) = \underline{i \cdot r}$

$$\Rightarrow J_k(1) = r_k, J_k(2) = 2 \cdot r_k.$$

$$\Rightarrow T \tilde{J}_k = (2\alpha r_k, 2\alpha r_k) \Leftrightarrow \hat{J}_{k+1}$$

LINEAR LEAST SQUARES:

$$r_{k+1} \in \underset{r}{\operatorname{argmin}} \left(\underbrace{r - 2\alpha r_k}_{\tilde{J}(1, r)} \right)^2 + \left(\underbrace{2r - 2\alpha r_k}_{\tilde{J}(2, r)} \right)^2$$

$$\Rightarrow 2(r - 2\alpha r_k) + 2 \cdot (2r - 2\alpha r_k) \cdot 2 = 0$$

$$r_{k+1} = \frac{6}{5} \alpha r_k$$

\Rightarrow WE NEED $\alpha \stackrel{=}{=} < \frac{5}{6}$ TO CONVERGE ($r^* = 0$)

APPROXIMATE POLICY ITERATION

PI: GIVEN μ_k , SOLVE

$$J^{k+1} = T_{\mu^k} J^{k+1} \quad (\text{POLICY EVAL})$$

$$\mu^{k+1} = \underset{\mu}{\operatorname{argmin}} \sum p_{ij}(\mu) (g_{li, a_{ij}} + J^{k+1}(j))$$

(POLICY UPDATE)

↳ INCORPORATE ANN BY REPLACING J BY $\tilde{J}(l; z)$

BELLMAN EQN METHODS

Idea: Solve directly BELLMAN EQN BY USING A MU.

$$\min_{\mu} \sum_{i \in S} \left(\tilde{J}(i, \mu) - \min_{\mu} \sum_{u} p_{ij}(u) (g_{ij}(u) + \tilde{J}(i, \mu)) \right)^2$$

RESIDUAL OF BELLMAN EQN. ($J^* = T J^*$)

$$\text{RESIDUAL: } D(i, \mu) = \tilde{J}(i, \mu) - \min_{\mu} \sum_{u} p_{ij}(u) (g_{ij}(u) + \tilde{J}(i, \mu))$$

$$\Leftrightarrow \min_{\mu} \sum_{i \in S} D(i, \mu)^2$$

APPLYING SGD TO MINIMIZE THE RESIDUAL

? SAMPLE AT k

$$\pi_{k+1} = \pi_k - \gamma \underbrace{D(\pi, \pi_k)} \nabla_{\pi} \underbrace{D(\pi, \pi_k)}$$

$$= \pi_k - \gamma \underbrace{D(\pi, \pi_k)} \left(\sum_j p_{ij}(u) \nabla_{\pi} \tilde{J}(j, \pi_k) - \nabla_{\pi} \tilde{J}(i, \pi_k) \right)$$

WHERE

$$\bar{u} = \underset{u}{\operatorname{argmin}} \sum_j p_{ij}(u) (g(i, u, j) + \tilde{J}(j, \pi_k))$$







