

MATHEMATICAL FOUNDATIONS OF REINFORCEMENT LEARNING

LECTURE II: DYNAMIC PROGRAMMING ALGORITHMS

- * DYNAMIC PROGRAMMING
 - * FINITE HORIZON CONTROL
 - * INFINITE HORIZON CONTROL
 - * STATIONARY POLICIES
 - * BELLMAN OPERATOR
 - * VALUE ITERATION
 - * POLICY ITERATION

DYNAMIC PROGRAMMING

INGREDIENTS:

a) A DISCRETE-TIME EVOLUTION,
 x_1, \dots, x_n STATES.

b) x_0 , A **TERMINATION** STATE.

c) A CONTROL SET, $u(i) \in U(i)$

d) TRANSITION PROBABILITIES $P_{ij}(u)$

$$x_i \xrightarrow{u} x_j$$

e) A COST THAT ACCUMULATES OVER TIME:

$$\alpha^k g(i, u, j)$$

g : **running cost**; k = ITERATION

$0 < \alpha < 1$: DISCOUNT FACTOR (DEPRECIATION)

THE GOAL: TO FIND AN OPTIMAL POLICY MAP

$$\Pi := \{ \mu_0, \mu_1, \dots \}$$

$\mu_k(i) :=$ OPTIMAL FEEDBACK AT TIME k AT STATE i
 $\mu_k(i) \in \bar{U}(i)$

FOR A FIXED POLICY, THE SEQUENCE OF STATES BECOMES A MARKOV CHAIN WITH TRANSITION PROBABILITIES

$$P(i_{k+1} = j \mid i_k = i) = P_{ij}(\mu_k(i))$$

↳ Obs: EQUIVALENT TO $\frac{dx}{dt} = f(x(t), u(t))$ WITH $u(t)$ A GIVEN SIGNAL

DIFFERENT DP EQUATIONS FOR DIFFERENT HORIZONS K AND COST g

FINITE HORIZON CONTROL

Fix A CONTROL HORIZON : $0, \dots, N$

GIVEN A POLICY $\pi = \{\mu_0, \dots, \mu_{N-1}\}$ AND AN INITIAL STATE x_i ,

THE EXPECTED TOTAL COST IS GIVEN BY

$$J_N^\pi(i) := \mathbb{E} \left[\underbrace{\alpha^N G(i_N)}_{\text{TERMINAL COST}} + \sum_{k=0}^{N-1} \alpha^k g(i_k, \mu_k(i_k), i_{k+1}) \mid i_0 = i \right]$$

→ W/RESPECT TO THE PROB. OF THE MC $\{x_0, x_1, \dots\}$

THE OPTIMAL COST-TO-GO (VALUE FUNCTION) IS GIVEN BY

$$J_N^*(i) = \min_{\pi} J_N^\pi(i)$$

INFINITE HORIZON CONTROL

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

OPTIMAL COST-TO-GO: $J^*(i) = \min_{\pi} J^\pi(i)$

FOR ∞ -HORIZON PROBLEMS, WE LOOK FOR A **STATIONARY** POLICY,

$$\pi = \{ \mu, \mu, \mu, \dots, \mu \} \rightarrow J^\mu(i)$$

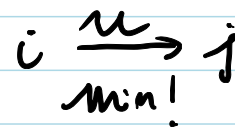
DYNAMIC PROGRAMMING EQUATIONS

FOR FINITE HORIZON CONTROL:

I. THE CASE WITH $N=1$ (**INSTANTANEOUS CONTROL**)

$$J_1^*(i) = \min_{\mu_0} \sum_{j=1}^m p_{ij}(\mu_0(i)) (g(i, \mu_0(i), j) + \alpha G(j))$$

$$= \min_{\mu \in \mathcal{U}(i)} \sum_{j=1}^m \underbrace{p_{ij}(\mu)}_{E[\dots]} \left(\underbrace{g(i, \mu, j)}_{\text{EXPECTED PRESENT COST}} + \alpha \underbrace{G(j)}_{\text{FUTURE COST}} \right)$$



→ SOLVING 1 OPTIMIZATION FOR EACH i

II. DP FOR AN ARBITRARY HORIZON N

BELLMAN $\Rightarrow J_N^*(i) = \min_{u \in U(i)} \sum_j P_{ij} (g(i, u, j) + \alpha J_{N-1}^*(j)) \quad \forall i$

VERIFICATION: $J_0^*(i) = G(i)$; $J_1^*(i) = \dots + \alpha G(j)$

FOR AN N-STEP POLICY π_N WITH INITIAL STATE $i_0 = i$, WE SPLIT

$$\pi_N := \{ \pi_{N-1}, u \}$$

$$\begin{aligned} J_N^*(i) &\equiv \min_{\pi} J_N^\pi(i) = \min_{\{ \pi_{N-1}, u \}} \left\{ \sum_j P_{ij}(u) g(i, u, j) + \alpha J_{N-1}^{\pi_{N-1}}(j) \right\} \\ &= \min_u \sum_j P_{ij}(u) g(i, u, j) + \alpha \min_{\pi_{N-1}} J_{N-1}^{\pi_{N-1}}(j) \\ &= \min_u \sum_j P_{ij}(u) g(i, u, j) + \alpha J_{N-1}^*(j) \end{aligned}$$

II. DP FOR ∞ -HORIZON PROBLEMS

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

TAKING THE LIMIT TO THE FINITE HORIZON $N \rightarrow \infty$

$$(*) \quad J^*(i) = \min_{u \in \mathcal{U}(i)} \sum_{j=1}^n p_{ij}(u) (g(i, u, j) + \alpha J^*(j))$$

↳ THE MINIMIZER $u^* \rightarrow \mu(i)$ OPTIMAL FEEDBACK

NOTATION: $T(J)(i) \equiv (TJ)(i) := \min_u \sum_j p_{ij}(u) [g(i, u, j) + \alpha J(j)]$

BELLMAN OPERATOR $\stackrel{(*)}{\Rightarrow} J^* = TJ^*$
FIXED POINT!

MORE NOTATION: $T_{\mu} J(i) := \sum_j P_{ij}(\mu(i)) \dots$

REPLACING $\mu(i)$ by $\mu(i)$
 $\mu \in \mathcal{U}(i)$

OPTIMALITY: $J^* = T J^* \leq T_{\mu} J^*$, FOR ANY POLICY μ

MORE NOTATION II: $P_{\mu} \in \mathbb{R}^{m \times m}$, $P_{\mu}^{ij} = P_{ij}(\mu(i))$

$$\Rightarrow \boxed{T_{\mu} J} = \underbrace{\bar{g}_{\mu}}_{\in \mathbb{R}^m} + \underbrace{P_{\mu}}_{\in \mathbb{R}^m} J$$

WHERE $\bar{g}_{\mu}(i) = \sum_{j=0}^n P_{ij}(\mu(i)) g(i, \mu(i), j)$

MORE NOTATION III: $(T^k J)(i) = T(T^{k-1} J)(i) \dots, T^0 J(i) = J(i)$
(SAME FOR $(T_{\mu_1} T_{\mu_2} \dots T_{\mu_n} J)(i)$)

PROPERTIES OF T AND T_μ

MONOTONICITY: $\mathbb{J} \preceq \bar{\mathbb{J}} \implies J(i) \leq \bar{J}(i), i=1, \dots, n$

FOR ANY n -DIMENSIONAL VECTORS \mathbb{J} AND $\bar{\mathbb{J}}$, FOR ANY STATIONARY POLICY μ WE HAVE

$$(T^k \mathbb{J})(i) \leq (T^k \bar{\mathbb{J}})(i)$$

$$(T_\mu^k \mathbb{J})(i) \leq (T_\mu^k \bar{\mathbb{J}})(i)$$

PROOF: USE THE DEFINITION OF $T^k \mathbb{J}$ AND INDUCTION

$$T^0 \mathbb{J} = \mathbb{J} \leq \bar{\mathbb{J}} = T^0 \bar{\mathbb{J}}$$

$$T^k \mathbb{J} = T(T^{k-1} \mathbb{J}) = \min_{\mu \in \mathcal{U}_j} \sum_j p_{ij}(\mu) \{ g(i, \mu, j) + \alpha T^{k-1} \mathbb{J} \} \leq \min_{\mu} \dots + \alpha T^{k-1} \bar{\mathbb{J}}$$
$$T(T^{k-1} \bar{\mathbb{J}}) = T^k \bar{\mathbb{J}}$$

FIXED POINT PROPERTIES

ASSUME THERE EXISTS A TERMINATION STATE x_0 SUCH THAT $p_{00}(n) = 1$, $g(a, \mu, 0) = 0 \forall \mu \in \mathcal{U}(a)$ AND AT LEAST ONE POLICY μ SUCH THAT THE PROBABILITY OF REACHING x_0 AFTER n STAGES IS > 0 .

THEN: a) THE OPTIMAL COST-TO-GO $J^* \in \mathbb{R}^M$ IS THE
UNIQUE SOLUTION TO $J^* = T J^*$

b) $\lim_{K \rightarrow \infty} T^{(K)} J = J^* \quad \forall J \in \mathbb{R}^M$

c) A STATIONARY POLICY IS OPTIMAL IFF

$$T_\mu J^* = T J^*$$

Obs: T IS A CONTRACTION MAPPING

$$\|TJ - T\bar{J}\|_{\infty} \leq \beta \|J - \bar{J}\|_{\infty}$$

$$0 \leq \beta < 1, \quad \forall J, \bar{J}.$$

$$\text{WHERE } \|J\|_{\infty} = \max_{i=1, \dots, m} \frac{|J(i)|}{\xi(i)}$$

BANACH FIXED POINT SUGGESTS A SOLUTION METHOD:

VALUE ITERATION

GIVEN ANY $J^0(i), \forall i \Leftrightarrow J^0 \in \mathbb{R}^m$, ITERATE

$$J^{k+1} = TJ^k, \quad \text{UNTIL CONVERGENCE}$$

+ : Simple, Robust

- : Slow depending on β
if $\beta \rightarrow 1$

ALTERNATIVE FOR SOLVING $J^* = (TJ^*)$:

- NONLINEAR SYSTEM OF EQNS. (VI \Leftrightarrow FPI)

→ USE **NEWTON'S METHOD** $x_{k+1} = x_k - \frac{f(x_k)}{f'(x_k)}$

POLICY ITERATION:

GIVEN μ^0 : I) POLICY EVALUATION STEP:

$$\text{Solve for } J^{k+1} \quad J^{k+1} = T_{\mu^k} J^{k+1}$$

II) POLICY UPDATE STEP:

$$\mu^{k+1}(i) = \underset{\mu}{\operatorname{argmin}} \sum_j p_{ij}(\mu) [g(i, \mu, j) + \alpha J^{k+1}(j)]$$

Eventually $J^k \rightarrow J^*$, $\mu^k \rightarrow \mu^*$

+ : NEWTON METHOD, FASTER CONVERGENCE

- : NEWTON METHOD, NEED GOOD INITIAL GUESS μ_0

SUMMARY:

- * DYNAMIC PROGRAMMING
- * FINITE HORIZON CONTROL
- * INFINITE HORIZON CONTROL
- * STATIONARY POLICIES
- * BELLMAN OPERATOR
- * VALUE ITERATION
- * POLICY ITERATION