

# Optimal Control, Lecture 3: Value Iteration and Policy Iteration

Anders Hansson

Division of Automatic Control  
Linköping University

# Contents

- ▶ Bellman Equation
- ▶ Value Iteration (VI)
- ▶ Policy Iteration (PI)
- ▶ Approximate PI
- ▶ Linear Quadratic Control

# Infinite Time Horizon Optimization

$$\begin{aligned} J^*(x_0) = \min & \quad \sum_{k=0}^{\infty} \gamma^k f(x_k, u_k) \\ \text{subject to} & \quad x_{k+1} = F(x_k, u_k) \\ & \quad x_0 \text{ given} \\ & \quad u_k \in U(x_k) \end{aligned} \tag{1}$$

with  $0 < \gamma \leq 1$  discount factor.

# Bellman Equation

Assume  $0 \in U(0)$ ,  $f(0, 0) = 0$ ,  $F(0, 0) = 0$  and that  $f$  is strictly positive definite. If there exists a strictly positive definite  $V$  such that the Bellman equation

$$V(x) = \min_{u \in U(x)} \{f(x, u) + \gamma V(F(x, u))\}$$

holds, then

- ▶ (a)  $J^*(x) = V(x)$
- ▶ (b) The minimizing argument in the Bellman equation is an optimal feedback for (3) that results in a globally convergent closed loop system if  $\gamma$  is sufficiently close to one.

## Value Iterations (VI)

Change iteration index in the dynamic programming recursion so that we iterate forward instead:

$$V_{k+1}(x) = \min_{u \in U(x)} \{f_k(x, u) + \gamma V_k(F_k(x, u))\} \quad (2)$$

with initial value  $V_0(x) = 0$ .

If one has a clever guess of an approximate solution to the Bellman equation, this can be used as initial value instead.

## VI for LQ Control

$$\begin{aligned} \min & \quad \sum_{k=0}^{\infty} \gamma^k (x_k^T S x_k + u_k^T R u_k) \\ \text{subject to} & \quad x_{k+1} = A x_k + B u_k \\ & \quad x_0 \text{ given} \end{aligned} \tag{3}$$

Let  $V_k(x) = x^T P_k x$  with  $P_0 = 0$ . Similarly as in LQ Example in previous lecture

$$P_{k+1} = S + \gamma A^T P_k A - \gamma^2 A^T P_k B (R + B^T P_k B)^{-1} B^T P_k A$$

# Proof of Convergence

Based on Bellman operator:

$$T(V)(x) = \min_{u \in U(x)} \{f(x, u) + \gamma V(F(x, u))\}. \quad (4)$$

Details on white board.

# Policy Iterations (PI)

Bellman policy operator:

$$T_{\mu}(V)(x) = f(x, \mu(x)) + \gamma V(F(x, \mu(x))) \quad (5)$$

for a given function  $\mu$ .

Iterate starting with initial  $\mu_0$ :

1. Solve (policy evaluation step)

$$V_k(x) = T_{\mu_k}(V_k)(x), \quad (6)$$

2. Solve (policy improvement step)

$$\mu_{k+1}(x) = \operatorname{argmin}_{u \in U(x)} \{f(x, u) + \gamma V_k(F(x, u))\}. \quad (7)$$

Proof of convergence on white board.



## LQ Control

Guess that  $V_k(x) = x^T P_k x$  and that  $\mu_k(x) = -L_k x$ .

**Policy evaluation step:**

$$x^T P_k x = x^T S x + x^T L_k^T R L_k x + \gamma x^T (A - B L_k)^T P_k (A - B L_k) x$$

for given  $L_k$ . Solution from Lyapunov equation

$$P_k - \gamma (A - B L_k)^T P_k (A - B L_k) = S + L_k^T R L_k,$$

**Policy improvement step:**

$$\mu_{k+1}(x) = \underset{u}{\operatorname{argmin}} \left\{ x^T S x + u^T R u + \gamma (A x - B u)^T P_k (A x - B u) \right\}.$$

with solution  $\mu_{k+1}(x) = -L_{k+1} x$ , where

$$L_{k+1} = \gamma (R + \gamma B^T P_k B)^{-1} B^T P_k A.$$

## Approximate Evaluation of $V_k$

It holds that (6) implies

$$\begin{aligned} V_k(x_0) &= f(x_0, \mu_k(x_0)) + \gamma V_k(F(x_0, \mu_k(x_0))) \\ &= f(x_0, \mu_k(x_0)) + \gamma V_k(x_1) \\ &= f(x_0, \mu_k(x_0)) + \gamma f(x_1, \mu_k(x_1)) + \gamma^2 V_k(x_2) \\ &\vdots \\ &= \sum_{i=0}^{N-1} \gamma^i f(x_i, \mu_k(x_i)) + \gamma^N V_k(x_N), \end{aligned} \tag{8}$$

where  $x_{i+1} = F(x_i, \mu_k(x_i))$ .

In case  $N$  is large and  $\mu_k$  is stabilizing we have that  $x_N$  is close to zero and that also  $V_k(x_N)$  is close to zero.

Approximate evaluation of  $V_k(x_0)$  obtained by simulating the dynamical system and add up the incremental costs.

## Approximation of $V$

Let

$$\beta_k^s = \sum_{i=0}^{N-1} \gamma^i f(x_i, \mu_k(x_i)),$$

where  $x_{i+1} = F(x_i, \mu_k(x_i))$ ,  $x_0 = x^s$ ,  $1 \leq s \leq r$ .

Let  $\tilde{V}(x, a)$  be a linear regression or an Artificial Neural Network (ANN) with parameter  $a$  that should approximate  $V_k$  in (6).

Find the approximation of  $V_k$  by solving

$$\text{minimize } \frac{1}{2} \sum_{s=1}^r \left( \tilde{V}(x^s, a) - \beta_k^s \right)^2$$

with variable  $a$ . The solution is denoted  $a_k$ . Then perform exact policy improvement step

$$\mu_{k+1}(x) = \operatorname{argmin}_{u \in U(x)} \left\{ f(x, u) + \gamma \tilde{V}(F(x, u), a_k) \right\}. \quad (9)$$

## Approximate LQ Control

Assume one input and two states and let  $\varphi(x) = (x_1^2, x_2^2, 2x_1x_2)$ .<sup>1</sup> Let

$$\tilde{V}(x, a) = a^T \varphi(x),$$

With

$$\tilde{P} = \begin{bmatrix} a_1 & a_3 \\ a_3 & a_2 \end{bmatrix}$$

we have

$$\tilde{V}(x, a) = x^T \tilde{P} x.$$

Hence true value function  $V(x) = x^T P x$  and approximate value function  $\tilde{V}(x, a)$  agree if  $\tilde{P} = P$ .

---

<sup>1</sup>Notice that the indices refer to components of the vector and not to time.

## Approximate LQ Control

With an abuse of notation  $a_k \in \mathbf{R}^3$ , which defines  $\tilde{P}_k$ , is obtained as solution to Least Squares (LS) problem

$$\text{minimize } \frac{1}{2} \sum_{s=1}^r (\varphi^T(x^s)a - \beta_k^s)^2$$

with variable  $a$ . The solution  $a_k$  satisfies normal equations

$$\Phi_k^T \Phi_k a_k = \Phi_k^T \beta_k,$$

where

$$\Phi_k = \begin{bmatrix} \varphi^T(x^1) \\ \vdots \\ \varphi^T(x^r) \end{bmatrix}, \quad \beta_k = \begin{bmatrix} \beta_k^1 \\ \vdots \\ \beta_k^r \end{bmatrix},$$

and where

$$\beta_k^s = \sum_{i=0}^{N-1} \gamma^i (x_i^T S x_i + \mu_k(x_i)^T R \mu_k(x_i)),$$

and where  $x_{i+1} = Ax_i + B\mu_k(x_i)$  with initial values  $x^s$ .

## Approximate LQ Control

It is crucial to choose  $x^s$  such that  $\Phi_k^T \Phi_k$  is invertible, which holds if  $r \geq 3$ .

We define

$$\begin{aligned}\tilde{Q}_k(x, u, a) &= f(x, u) + \gamma \tilde{V}(Ax + Bu, a) \\ &= x^T Sx + u^T Ru + \gamma(Ax + Bu)^T \tilde{P}(Ax + Bu) \\ &= \begin{bmatrix} x \\ u \end{bmatrix}^T \begin{bmatrix} S + \gamma A^T \tilde{P} A & \gamma A^T \tilde{P} B \\ \gamma B^T \tilde{P} A & R + \gamma B^T \tilde{P} B \end{bmatrix} \begin{bmatrix} x \\ u \end{bmatrix}.\end{aligned}\quad (10)$$

The solution to (9) is given by

$$\mu_{k+1}(x) = \underset{u}{\operatorname{argmin}} \tilde{Q}_k(x, u, a_k) = -\gamma(R + \gamma B^T \tilde{P}_k B)^{-1} B^T \tilde{P}_k A x$$

assuming  $R + \gamma B^T \tilde{P}_k B$  positive definite. Hence

$$\mu_{k+1}(x) = -L_{k+1}x,$$

where  $L_{k+1} = \gamma(R + \gamma B^T \tilde{P}_k B)^{-1} B^T \tilde{P}_k A$ .