

Approximate Dynamic Programming via Iterated Bellman Inequalities

Yang Wang and Stephen Boyd

Information Systems Laboratory, Electrical Engineering, Stanford University

Stochastic Control

- discrete-time dynamics: $x_{t+1} = f(x_t, u_t, w_t)$, $t = 0, 1, \dots$,
 - $x_t \in \mathcal{X}$ is the state
 - $u_t \in \mathcal{U}$ is the input
 - $w_t \in \mathcal{W}$ is IID noise
 - x_0 is random, independent of w_t
- state feedback control policy: $u_t = \psi(x_t)$, $t = 0, 1, \dots$,
 - $\psi: \mathcal{X} \rightarrow \mathcal{U}$ is the state feedback function
- infinite horizon discounted cost:

$$J = \mathbf{E} \sum_{t=0}^{\infty} \gamma^t \ell(x_t, u_t)$$

- $\ell: \mathcal{X} \times \mathcal{U} \rightarrow \mathbf{R} \cup \{+\infty\}$ is stage cost function
- $\gamma \in (0, 1)$ is discount factor
- use infinite values of ℓ to encode constraints on state and input:
 - $(x_t, u_t) \in \mathcal{C} = \{(x, u) \mid \ell(x, u) < \infty\}$ a.s.
- stochastic control problem: choose feedback function ψ to minimize J
- infinite dimensional nonconvex problem: very hard to solve in general
- ψ^* and J^* denote optimal feedback function and objective value

Dynamic Programming

- can characterize optimal solution via dynamic programming
- find (unique) $V^*: \mathcal{X} \rightarrow \mathbf{R}$ that satisfies *Bellman equation*

$$V^*(z) = \min_{v \in \mathcal{U}} (\ell(z, v) + \gamma \mathbf{E} V^*(f(z, v, w_t))), \quad \forall z \in \mathcal{X}$$

(expectation is over w_t)

- optimal feedback function is then
 - $\psi^*(z) = \operatorname{argmin}_{v \in \mathcal{U}} (\ell(z, v) + \gamma \mathbf{E} V^*(f(z, v, w_t)))$
- optimal value of stochastic control problem is $J^* = \mathbf{E} V^*(x_0)$
- define *Bellman operator* \mathcal{T} , for $h: \mathcal{X} \rightarrow \mathbf{R}$

$$(\mathcal{T}h)(z) = \inf_{v \in \mathcal{U}} \{\ell(z, v) + \gamma \mathbf{E} h(f(z, v, w_t))\}$$

- Bellman equation is then $V^* = \mathcal{T}V^*$
- \mathcal{T} is monotone: $V_1 \leq V_2 \Rightarrow \mathcal{T}V_1 \leq \mathcal{T}V_2$
- \mathcal{T} is γ -contraction (w.r.t. sup-norm) so for any $V: \mathcal{X} \rightarrow \mathbf{R}$

$$V^* = \lim_{k \rightarrow \infty} \mathcal{T}^k V$$

('value iteration' converges)

Bounds and Approximate Policies

- we consider problems for which we cannot compute V^* , ψ^* , J^*
- **performance bound**:
 - find effectively computable lower bound J^{lb} on J^*
- **approximate policy**:
 - find effectively computable policy $\hat{\psi}$ with performance \hat{J}
- if $\hat{J} - J^{\text{lb}}$ small, we're done (for this problem), since $J^{\text{lb}} \leq J^* \leq \hat{J}$
- 'effectively computable' means, e.g., by solving tractable convex optimization problem(s)

Iterated Bellman inequality

- suppose \hat{V} satisfies *iterated Bellman inequality*
 - $\hat{V} \leq \mathcal{T}^M \hat{V}$
- then $\hat{V} \leq \mathcal{T}^M \hat{V} \leq \mathcal{T}^{2M} \hat{V}$ (by monotonicity of \mathcal{T})
- thus we get
 - $\hat{V}(x) \leq \lim_{k \rightarrow \infty} (\mathcal{T}^{kM} \hat{V})(x) = V^*(x), \quad \forall x \in \mathcal{X}$
 - (since \mathcal{T} is a γ -contraction)
- iterated Bellman inequality is a sufficient condition for $\hat{V} \leq V^*$
- \hat{V} is a value function *underestimator*
- then (by monotonicity of expectation) $\mathbf{E} \hat{V}(x_0) \leq \mathbf{E} V^*(x_0) = J^*$
- so value function underestimator \hat{V} yields performance bound $J^{\text{lb}} = \mathbf{E} \hat{V}(x_0)$
- can get better bounds with larger M

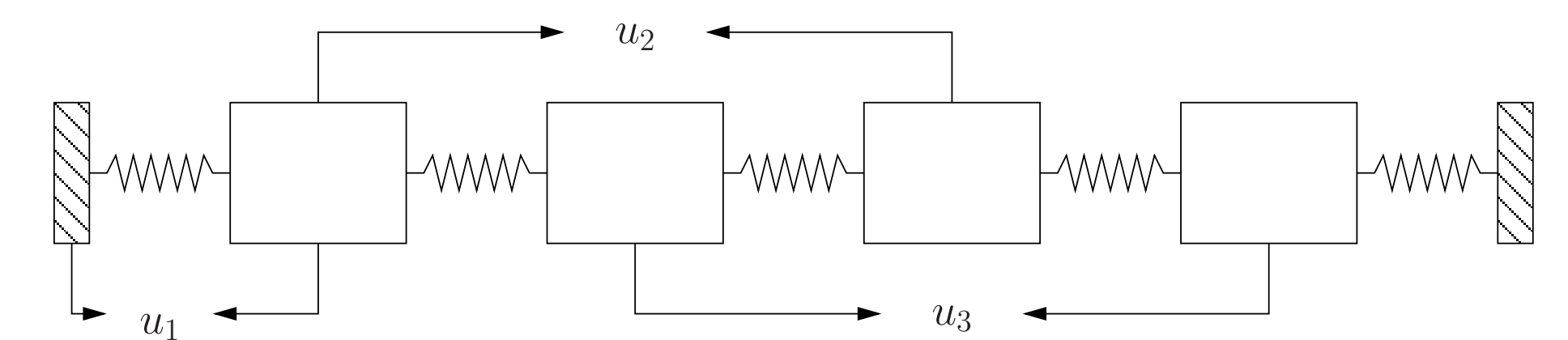
Bound Optimization

- a sufficient condition for the iterated Bellman inequality is
 - $\hat{V}_{i-1} \leq \mathcal{T} \hat{V}_i, \quad i = 1, \dots, M$
 - where we define $\hat{V}_0 = \hat{V}_M = \hat{V}$
- each \hat{V}_i is restricted to a finite-dimensional subspace
 - $\hat{V}_i = \sum_{j=1}^K \alpha_{ij} V^{(j)}$
 - $\alpha_{ij} \in \mathbf{R}$ are coefficients; $V^{(j)}: \mathcal{X} \rightarrow \mathbf{R}$ are (pre-selected) basis functions
- optimize performance bound by solving
 - maximize $\mathbf{E} \hat{V}(x_0) = \alpha_1 \mathbf{E} V^{(1)}(x_0) + \dots + \alpha_K \mathbf{E} V^{(K)}(x_0)$
 - subject to $\hat{V}_{i-1} \leq \mathcal{T} \hat{V}_i, \quad i = 1, \dots, M$
 - convex optimization problem in variables α_{ij}

ADP Policy

- approximate dynamic programming (ADP) policy is
 - $\psi^{\text{adp}}(z) = \operatorname{argmin}_{v \in \mathcal{U}} (\ell(z, v) + \gamma \mathbf{E} V^{\text{adp}}(f(z, v, w_t)))$
 - $V^{\text{adp}}: \mathcal{X} \rightarrow \mathbf{R}$ (which is to be chosen) is the *approximate value function*
 - when $V^{\text{adp}} = V^*$, this is optimal policy
- often performs well, even when V^{adp} is not close to V^*
- the underestimator \hat{V} is a natural choice for V^{adp}

Example: Discretized Mechanical System



- 4 masses connected by springs
- 3 input forces, with $|u_i| \leq 1$
- $w_t \sim \mathcal{N}(0, 0.1I)$
- $Q = 0.1I, R = 0.01I, \gamma = 0.95$

| Policy/Bound | Objective/Value |
|---|-----------------|
| ADP ($V^{\text{adp}} = \hat{V}$ from iterated bound) | 162.8 |
| J^{lb} (iterated Bellman with $M = 50$) | 154.8 |
| J^{lb} (iterated Bellman with $M = 1$) | 111.1 |
| Unconstrained LQR | 55.7 |

- iterated Bellman inequality gives better bounds
- performance of ADP policy very close to bound
- can close gap even further with a few simple methods ...
- ADP policy can be evaluated in *microseconds*

Summary

- applies to many other problem families with e.g., polynomial dynamics and stage costs, switching systems, multiplicative disturbances, ...
- in many cases, gives everything you want:
 - a provable lower bound on performance
 - a policy that comes close in performance, can be evaluated quickly