

# Reinforcement Learning and Optimal Control

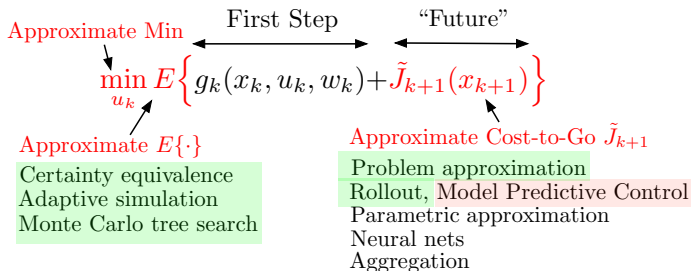
ASU, CSE 691, Winter 2019

Dimitri P. Bertsekas  
dimitrib@mit.edu

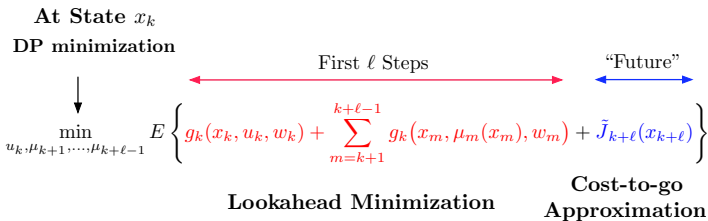
Lecture 5

- 1 Review of Approximation in Value Space and Rollout
- 2 On-Line Rollout for Deterministic Infinite Spaces Problems
- 3 Model Predictive Control
- 4 Parametric Approximation Architectures

# Recall Approximation in Value Space



## ONE-STEP LOOKAHEAD



## MULTISTEP LOOKAHEAD

# The Pure Form of Rollout - A Review

At State  $x_k$   
DP minimization

$$\min_{u_k, \mu_{k+1}, \dots, \mu_{k+\ell-1}} E \left\{ g_k(x_k, u_k, w_k) + \sum_{m=k+1}^{k+\ell-1} g_k(x_m, \mu_m(x_m), w_m) + \tilde{J}_{k+\ell}(x_{k+\ell}) \right\}$$

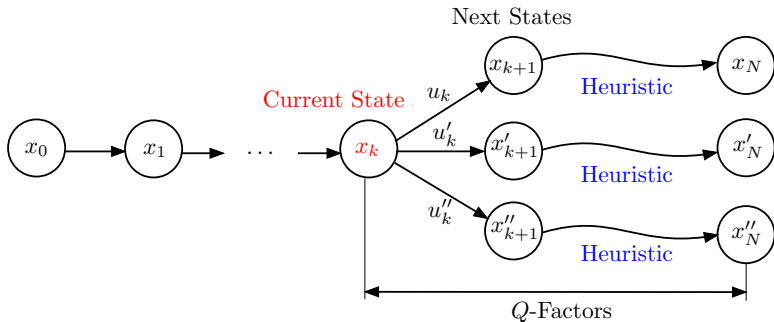
First  $\ell$  Steps      “Future”

Lookahead Minimization      Heuristic Cost  
Run the Base Policy

Use a suboptimal/heuristic policy at the end of limited lookahead

- The heuristic is called **base policy** (or default policy).
- The lookahead policy is called **rollout policy**.
- **Policy improvement**; connection with **policy iteration**.
- Involves **simulation and on-line implementation**; suitable for **on-line replanning**.
- **Deterministic rollout** lends itself to on-line implementation.

# General Structure of Deterministic Rollout



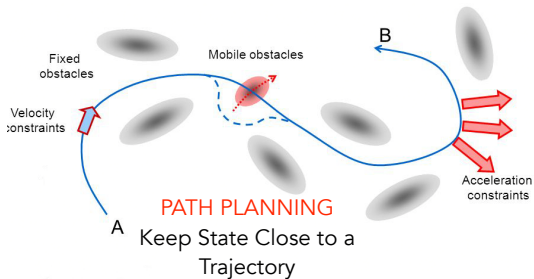
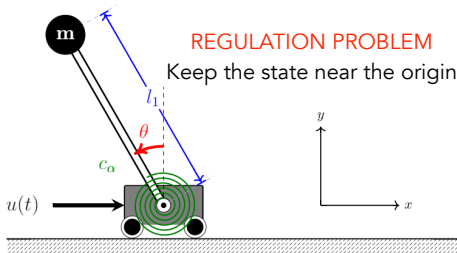
- At state  $x_k$ , for every pair  $(x_k, u_k)$ ,  $u_k \in U_k(x_k)$ , we generate a Q-factor

$$\tilde{Q}_k(x_k, u_k) = g_k(x_k, u_k) + H_{k+1}(f_k(x_k, u_k))$$

using the base heuristic [ $H_{k+1}(x_{k+1})$  is the heuristic cost starting from  $x_{k+1}$ ].

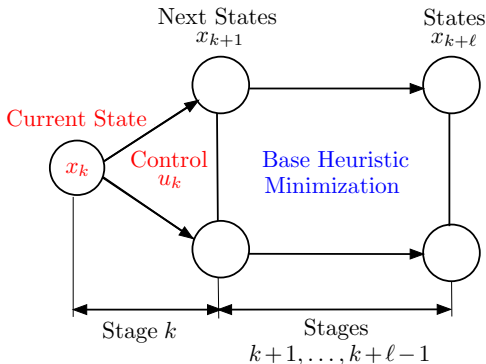
- We select the control  $u_k$  with minimal Q-factor.
- We move to next state  $x_{k+1}$ , and continue.
- A key question for today's lecture: **What if we have a continuous/infinite control set?**

# Classical Control Problems - Infinite Control Spaces



Need to deal with state and control constraints; linear-quadratic is often inadequate

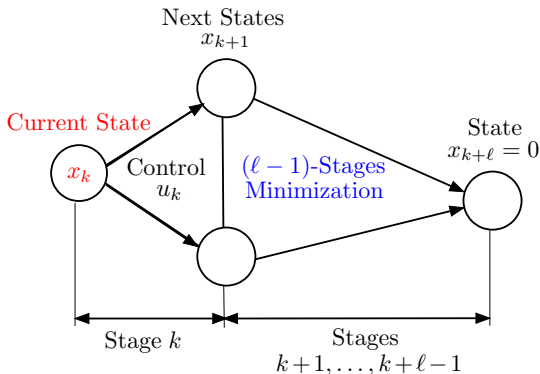
# On-Line Rollout for Deterministic Infinite-Spaces Problems



Suppose the control space is infinite

- One possibility is discretization of  $U_k(x_k)$ ; but **excessive number of Q-factors**.
- Another possibility is to use **optimization heuristics** that look  $(\ell - 1)$  steps ahead.
- Seamlessly combine the  $k$ th stage minimization and the optimization heuristic into **a single  $\ell$ -stage deterministic optimization**.
- Can solve it by **nonlinear programming/optimal control methods** (e.g., quadratic programming, gradient-based).

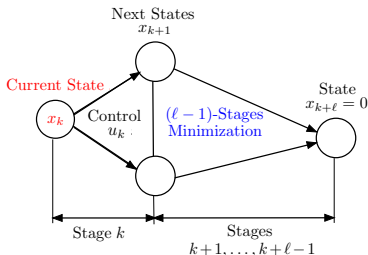
# Model Predictive Control for Regulation Problems



- System:  $x_{k+1} = f_k(x_k, u_k)$
- Cost per stage:  $g_k(x_k, u_k) \geq 0$ , **the origin 0 is cost-free and absorbing.**
- State and control constraints:  $x_k \in X_k$ ,  $u_k \in U_k(x_k)$  for all  $k$
- **At  $x_k$  solve an  $\ell$ -step lookahead version of the problem**, requiring  $x_{k+l} = 0$  while satisfying the state and control constraints.
- If  $\{\tilde{u}_k, \dots, \tilde{u}_{k+l-1}\}$  is the control sequence so obtained, apply  $\tilde{u}_k$ .



# Relation to Rollout

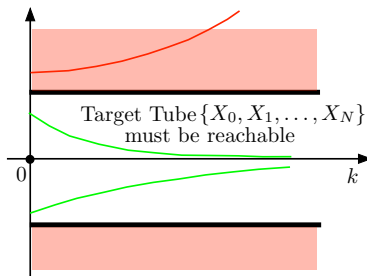


- It is rollout with **base heuristic** the **( $\ell - 1$ )-step min** (0 is cost-free and absorbing).
- This heuristic is **sequentially improving** (not sequentially consistent)

$$\min_{u_k \in U_k(x_k)} [g_k(x_k, u_k) + H_{k+1}(f_k(x_k, u_k))] \leq H_k(x_k)$$

where  $H_k(x_k)$ ,  $H_{k+1}(x_{k+1})$ : optimal heuristic costs starting at  $x_k$  and  $x_{k+1}$ .

- **Sequential improvement implies "stability"**:  $\sum_{k=0}^{\infty} g_k(x_k, u_k) \leq H_0(x_0) < \infty$ , where  $\{x_0, u_0, x_1, u_1, \dots\}$  is the state and control sequence generated by MPC.
- Major issue: **How do we know that the optimization of the base heuristic is solvable** (e.g., there exists  $\ell$  such that we can drive  $x_{k+l}$  to 0 for all  $x_k \in X_k$  while observing the state and control constraints).

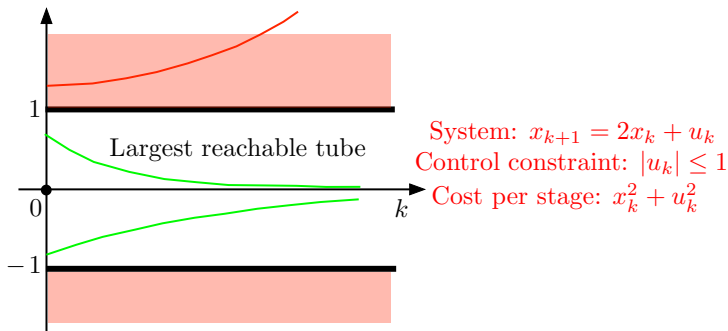


- The tube of state constraint sets  $\{X_0, X_1, \dots, X_N\}$  is **reachable** if the state  $x_k$  can be kept within it for all  $k$  when the initial state  $x_0$  belongs to  $X_0$ .
- If  $\{X_0, X_1, \dots, X_N\}$  is not reachable, MPC will not work; if it is reachable MPC will “typically” work. We may try to **extract a reachable subset**  $\{\bar{X}_0, \bar{X}_1, \dots, \bar{X}_N\}$ , with  $\bar{X}_k \subset X_k$ , for all  $k$ . Then **use  $\bar{X}_k$  in place of  $X_k$** .
- **Reachability algorithm**: Start with  $\bar{X}_N = X_N$ , and proceed backwards

$$\bar{X}_k = \{x_k \in X_k \mid \text{for some } u_k \in U_k(x_k) \text{ we have } f_k(x_k, u_k) \in \bar{X}_{k+1}\}.$$

- Generally, it is difficult to compute the sets  $\bar{X}_k$  of the target tube, but **algorithms that produce inner approximations have been constructed**.

## A Working Break: Challenge Question



- Is it true that  $\{(-1, 1), (1, 1), \dots, (-1, 1)\}$  is the largest reachable tube?
- Is the tube

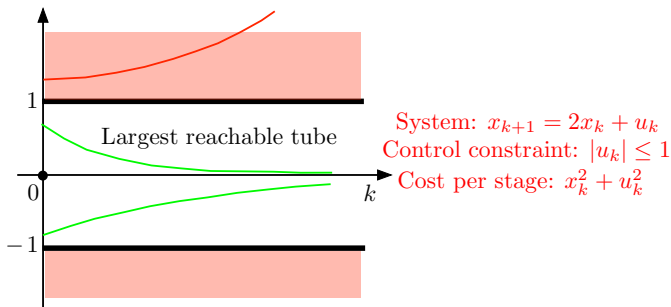
$$\{[-2, 2], [-2, 2], \dots, [-2, 2]\}$$

reachable? How about the tube

$$\{[-1/2, 1/2], [-1/2, 1/2], \dots, [-1/2, 1/2]\}$$

- How will MPC with  $\ell = 2$  work starting from  $x_0 = 1/2$  and from  $x_0 = 2$ ?

## Some Answers (see the textbook for details)



- If  $|x_k| > 1$  the state cannot be brought back towards 0; if  $|x_k| < 1$  it can.
- If  $|x_k| \leq 1/2$  the state can be driven to 0 in one step; if  $1/2 < |x_k| < 1$  the state can be driven to 0 in finitely many steps (the number increases as  $|x_k|$  is closer to 1).
- If  $|x_k| = 1$  the state can at best be kept where it is.
- $\{[-1, 1], [-1, 1], \dots, [-1, 1]\}$  is the largest reachable tube.
- $\{(-1, 1), (1, 1), \dots, (-1, 1)\}$  is the largest tube from within which the state can be driven to 0 in a finite number of steps.
- For  $\ell = 2$ , MPC must start from  $|x_0| \leq 1/2$ . It has the form  $\tilde{u}_k = -(5/3)x_k$ . It is stable and drives the state to 0 asymptotically (not in a finite number of steps).

## Approximation Architectures

- A class of functions  $\tilde{J}(x, r)$  that depend on  $x$  and a **vector**  $r = (r_1, \dots, r_m)$  of  $m$  **"tunable" scalar parameters** (or weights).
- We adjust  $r$  to change  $\tilde{J}$  and "match" the cost function approximated.
- **Training the architecture**: The algorithm to choose  $r$  (typically use **data/regression**).
- Architectures are **linear or nonlinear**, depending on whether  $\tilde{J}(x, r)$  is linear or nonlinear in  $r$ .
- Architectures are **feature-based** if they depend on  $x$  via a feature vector  $\phi(x)$ ,

$$\tilde{J}(x, r) = \hat{J}(\phi(x), r),$$

where  $\hat{J}$  is some function. Idea: **Features capture dominant nonlinearities**.

- A **linear feature-based architecture**:

$$\tilde{J}(x, r) = \sum_{\ell=1}^m r_{\ell} \phi_{\ell}(x),$$

where  $r_{\ell}$  and  $\phi_{\ell}(x)$  are the  $\ell$ th components of  $r$  and  $\phi(x)$ .

- **Local vs global**: Change in a single weight affects  $\tilde{J}$  locally vs globally.

# Generic Example Architectures

- **Piecewise constant approximation** (local): Partition the state space into subsets  $S_1, \dots, S_m$ . Let the  $\ell$ th feature be defined by membership in the set  $S_\ell$ , i.e.,  $\phi_\ell(x) = 1$  if  $s \in S_\ell$  and  $\phi_\ell(x) = 0$  if  $s \notin S_\ell$ . The architecture

$$\tilde{J}(x, r) = \sum_{\ell=1}^m r_\ell \phi_\ell(x),$$

is piecewise constant with value  $r_\ell$  for all  $x$  within the set  $S_\ell$ .

- **Quadratic polynomial approximation** (global):  $\tilde{J}(x, r)$  is quadratic in the components  $x^i$  of  $x$ . Consider features

$$\phi_0(x) = 1, \quad \phi_i(x) = x^i, \quad \phi_{ij}(x) = x^i x^j, \quad i, j = 1, \dots, n.$$

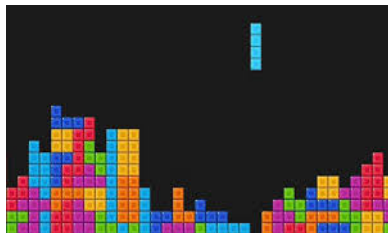
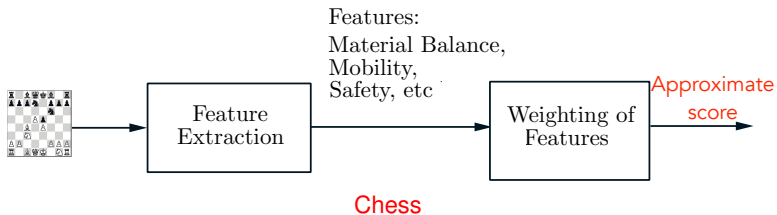
A linear feature-based approximation architecture:

$$\tilde{J}(x, r) = r_0 + \sum_{i=1}^n r_i x^i + \sum_{i=1}^n \sum_{j=i}^n r_{ij} x^i x^j$$

The parameter vector  $r$  has components  $r_0$ ,  $r_i$ , and  $r_{ij}$ .

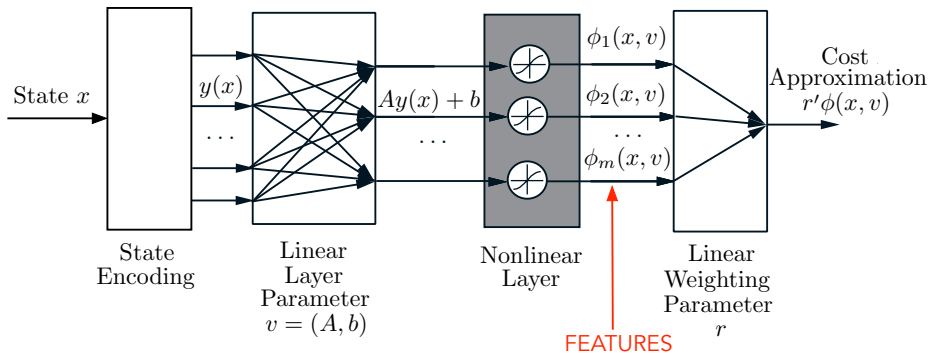
- **General polynomial architectures**: Polynomials in the components  $x^1, \dots, x^n$ . Another possibility: **Polynomials of features**.

# Examples of Domain-Specific Feature-Based Architectures



Tetris

# Neural Nets: An Architecture that does not Require Knowledge of Features





### We will cover:

- Training of parametric approximation architectures
- Neural networks; how do we use
- Sequential Dynamic Programming Approximation
- Q-factor Parametric Approximation

**PLEASE READ AS MUCH OF SECTIONS 3.1.3, 3.2-3.4 AS YOU CAN**  
**PLEASE DOWNLOAD THE LATEST VERSIONS FROM MY WEBSITE**