

# AA203 Optimal and Learning-based Control

## Lecture 5

### Dynamic Programming

Autonomous Systems Laboratory  
Daniele Gammelli

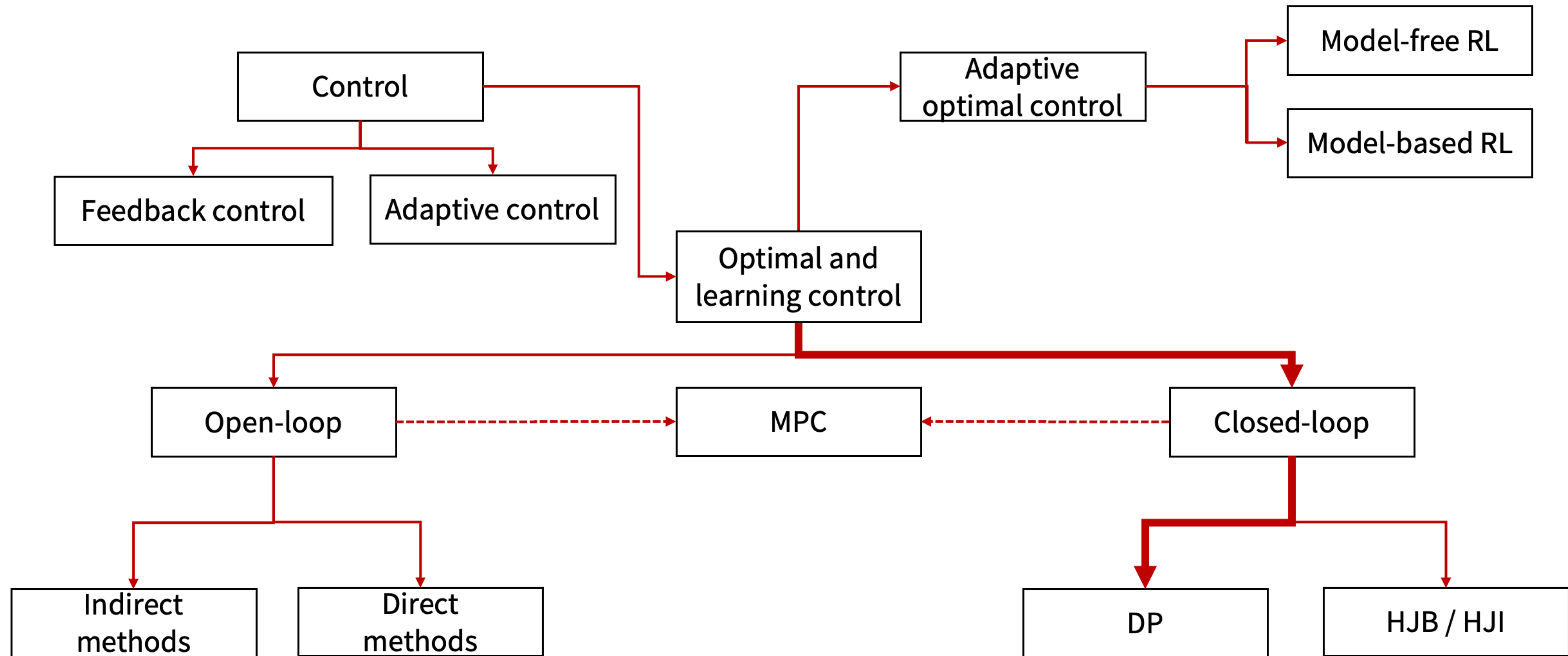


**Stanford University**



**Autonomous Systems Laboratory  
Stanford Aeronautics & Astronautics**

# Roadmap



# Outline of the next two lectures

Intro to dynamic programming (DP) and *principle of optimality*

The dynamic programming algorithm

Dynamic programming in control:

- Discrete LQR
- Stochastic Optimal Control Problem / Markov Decision Process (MDP): Stochastic LQR
- Policy Iteration and Value Iteration

This lecture

# Dynamic Programming

A method for solving complex problems, by:

- Breaking them down into subproblems
- Combining solutions to subproblems

Dynamic Programming is a very general solution method for problems which have two properties:

- Optimal substructure (*Principle of optimality* applies)
  - Optimal solution can be decomposed into subproblems, e.g., shortest path
- Overlapping subproblems
  - Subproblems recur many times
  - Solutions can be cached and reused

# Other applications of Dynamic Programming

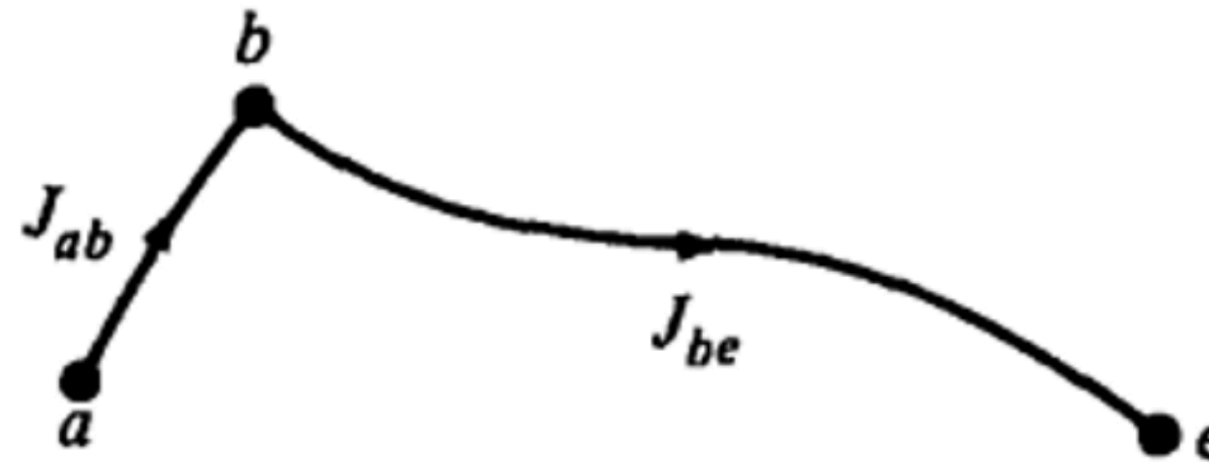
Dynamic programming is used across a wide variety of domains, e.g.

- Scheduling algorithms
- Graph algorithms (e.g., shortest path algorithms)
- Graphical models in ML (e.g., Viterbi algorithm)
- Etc.

# Principle of optimality

The *key concept* behind the dynamic programming approach is the *principle of optimality*

Suppose the optimal path for a multi-stage decision-making problem with additive cost structure is



Multi-stage:

- First decision yields segment a-b with cost  $J_{ab}$
- Remaining decisions yield segments b-e with cost  $J_{be}$

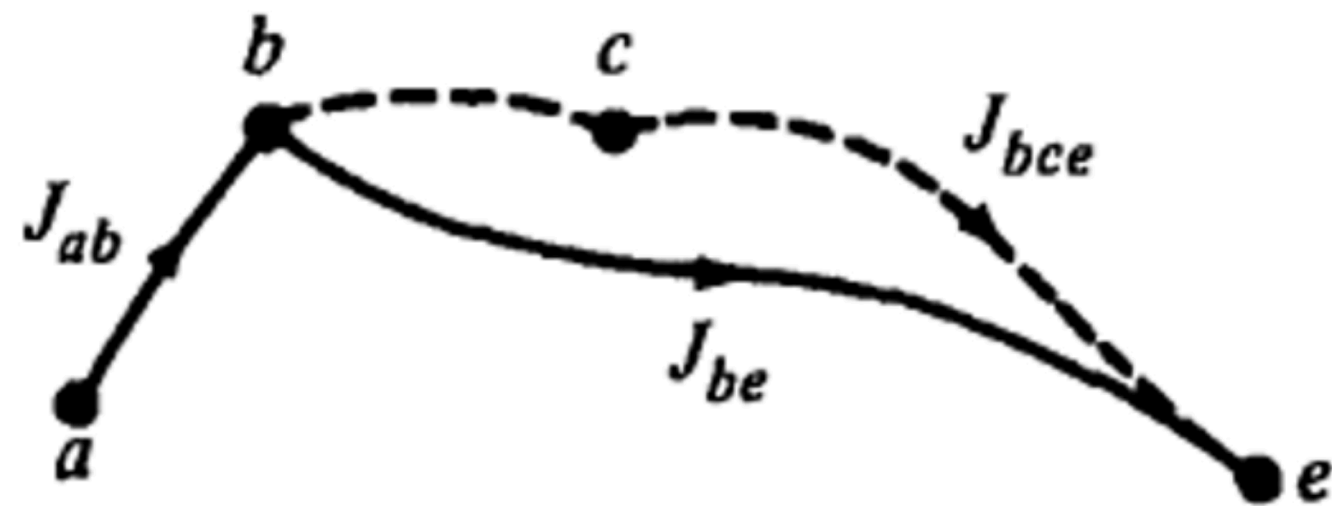
Additive cost:

- The optimal cost is then  $J_{ab}^* = J_{ab} + J_{be}$

# Principle of optimality

**Claim:** if  $a - b - e$  is the optimal path from  $a$  to  $e$ , then  $b - e$  is the optimal path from  $b$  to  $e$

**Proof:** suppose  $b-c-e$  is the optimal path from  $b$  to  $e$ . Then



$$J_{bce} < J_{be}$$

$$J_{ab} + J_{bce} < J_{ab} + J_{be} = J_{ae}^*$$

**Contradiction!**

# Principle of optimality

## Principle of optimality (for discrete-time systems):

Let  $\pi^* := \{ \pi_0^*, \pi_1^*, \dots, \pi_{N-1}^* \}$  be an optimal policy.

Assume state  $x_k$  is reachable.

Consider the subproblem whereby we are at  $x_k$  at time  $k$  and we wish to minimize the cost-to-go from time  $k$  to time  $N$ .

Then the truncated policy  $\{ \pi_k^*, \pi_{k+1}^*, \dots, \pi_{N-1}^* \}$  is optimal for the subproblem.

**Tail** policies are optimal for **tail** subproblems

Notation: for brevity  $\pi_k^* (\mathbf{x}_k) = \pi^* (\mathbf{x}_k, k)$



# Applying the principle of optimality

Consider the case where we want to find the optimal path from b to f, and that we know the cost of the optimal path from {c, d, e} to f.

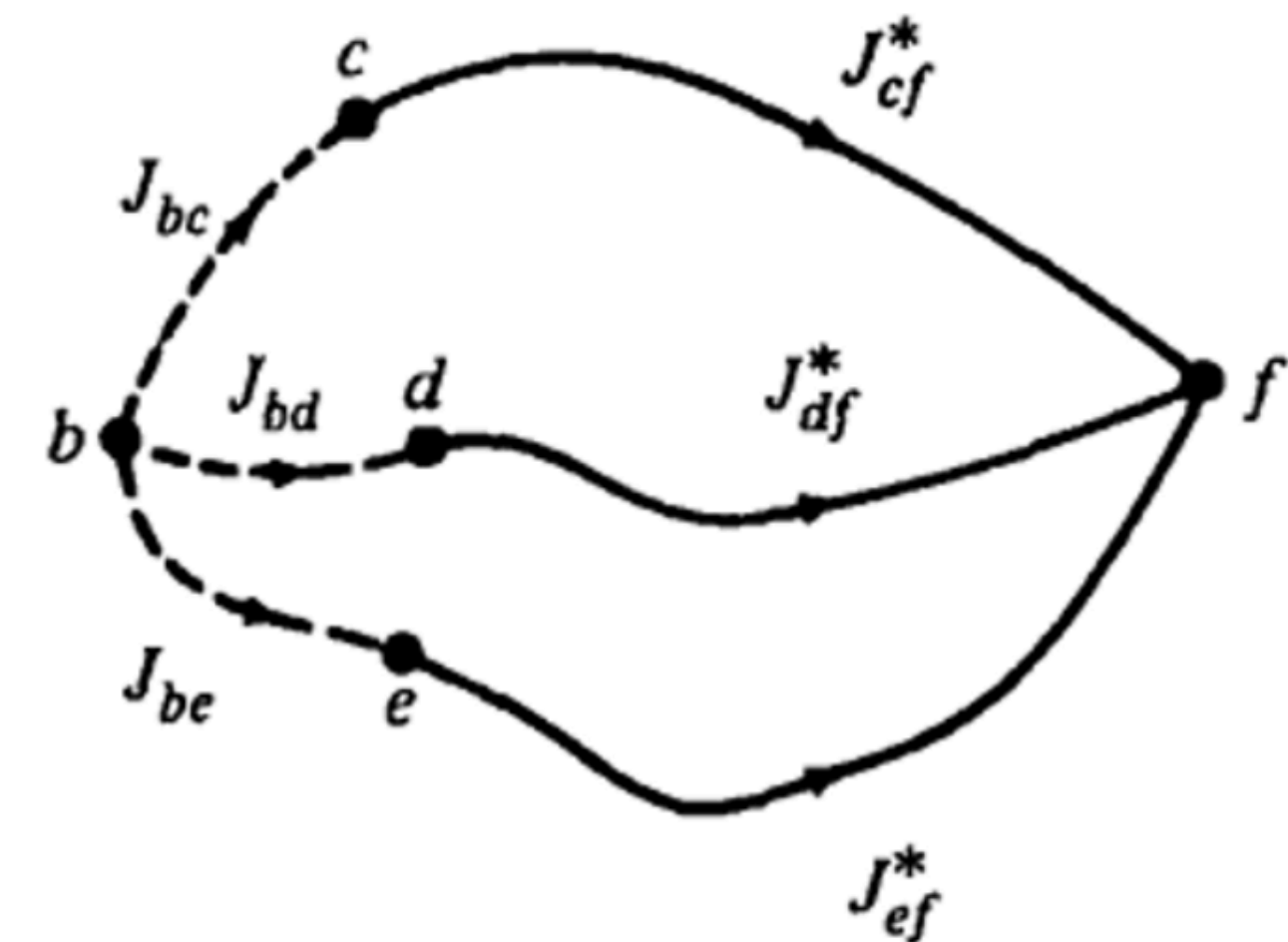
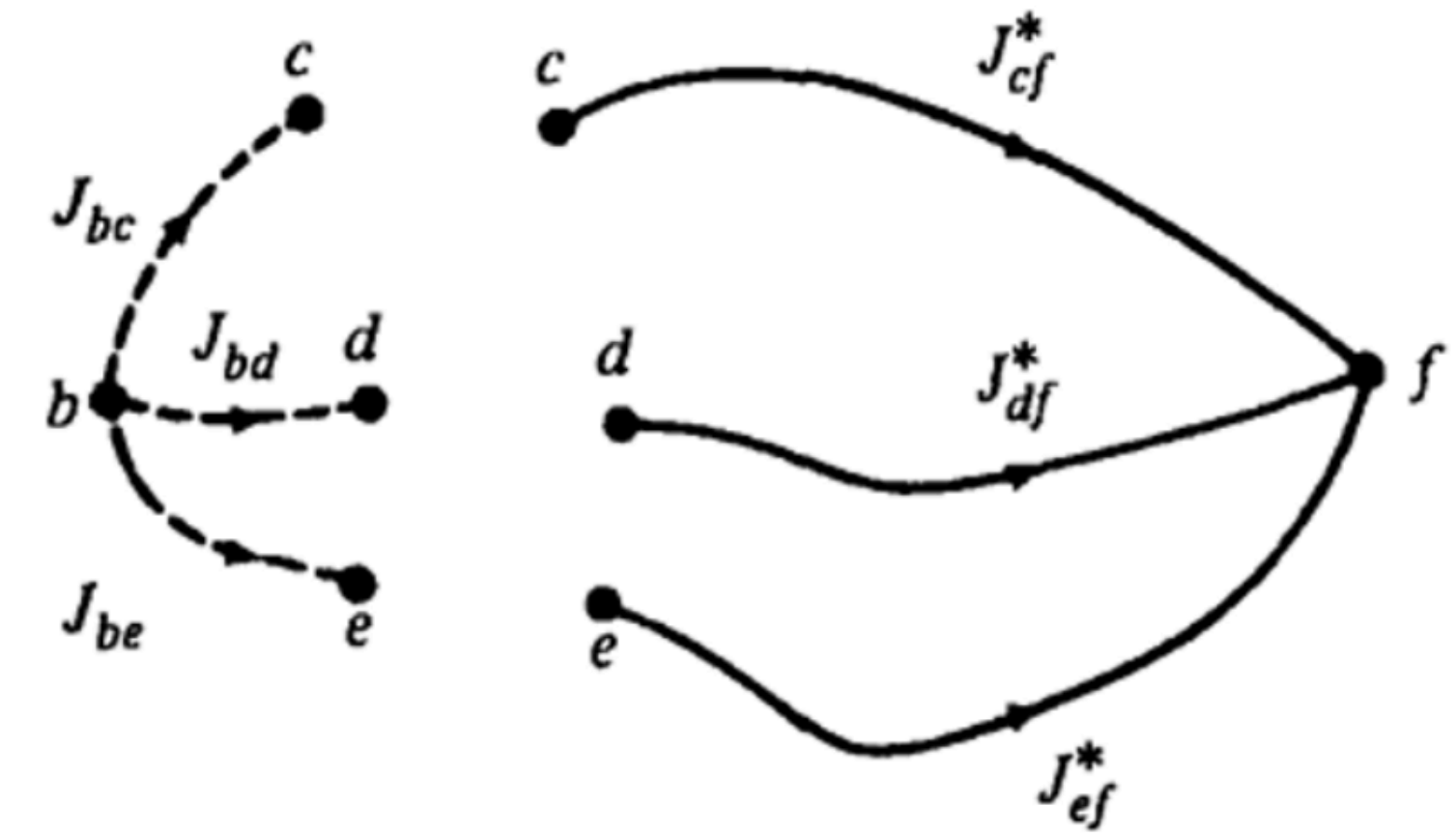
The principle of optimality tells us that the optimal policy is comprised of optimal sub-policies

Hence, the optimal trajectory is found by comparing:

$$C_{bcf} = J_{bc} + J_{cf}^*$$

$$C_{bdf} = J_{bd} + J_{df}^*$$

$$C_{bef} = J_{be} + J_{ef}^*$$

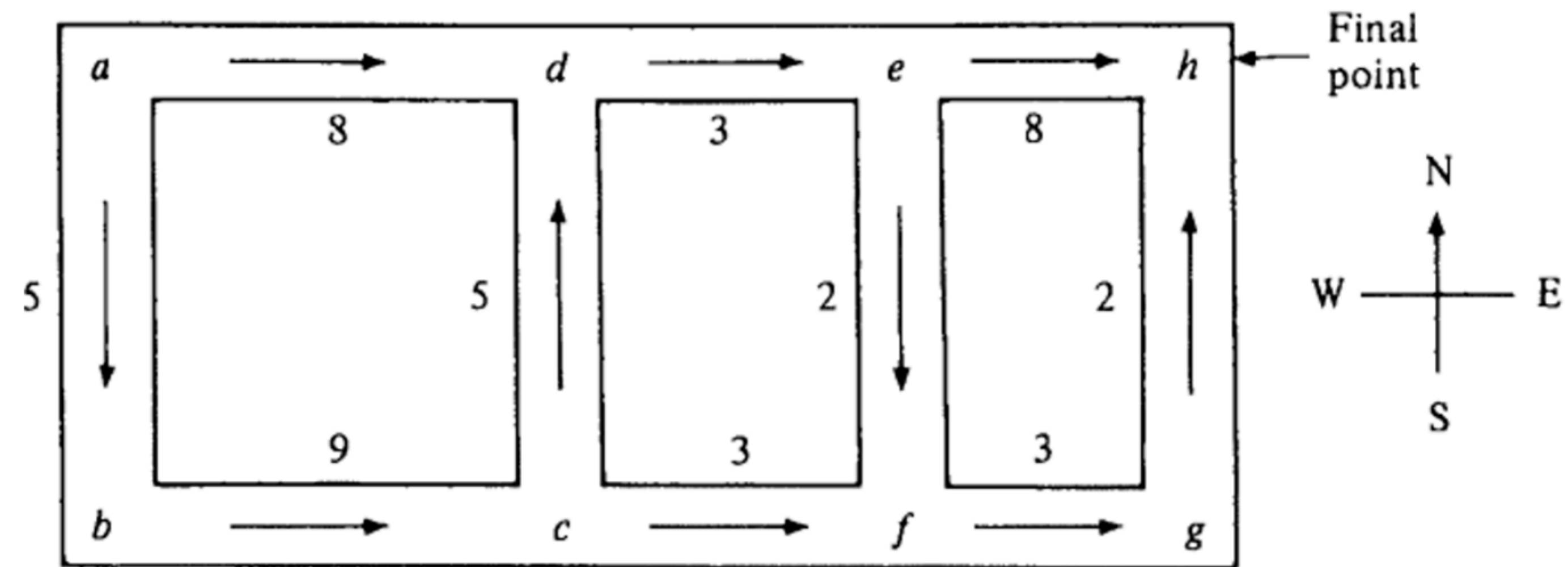


The “cost-to-go” allows us to only compute one-step look-ahead

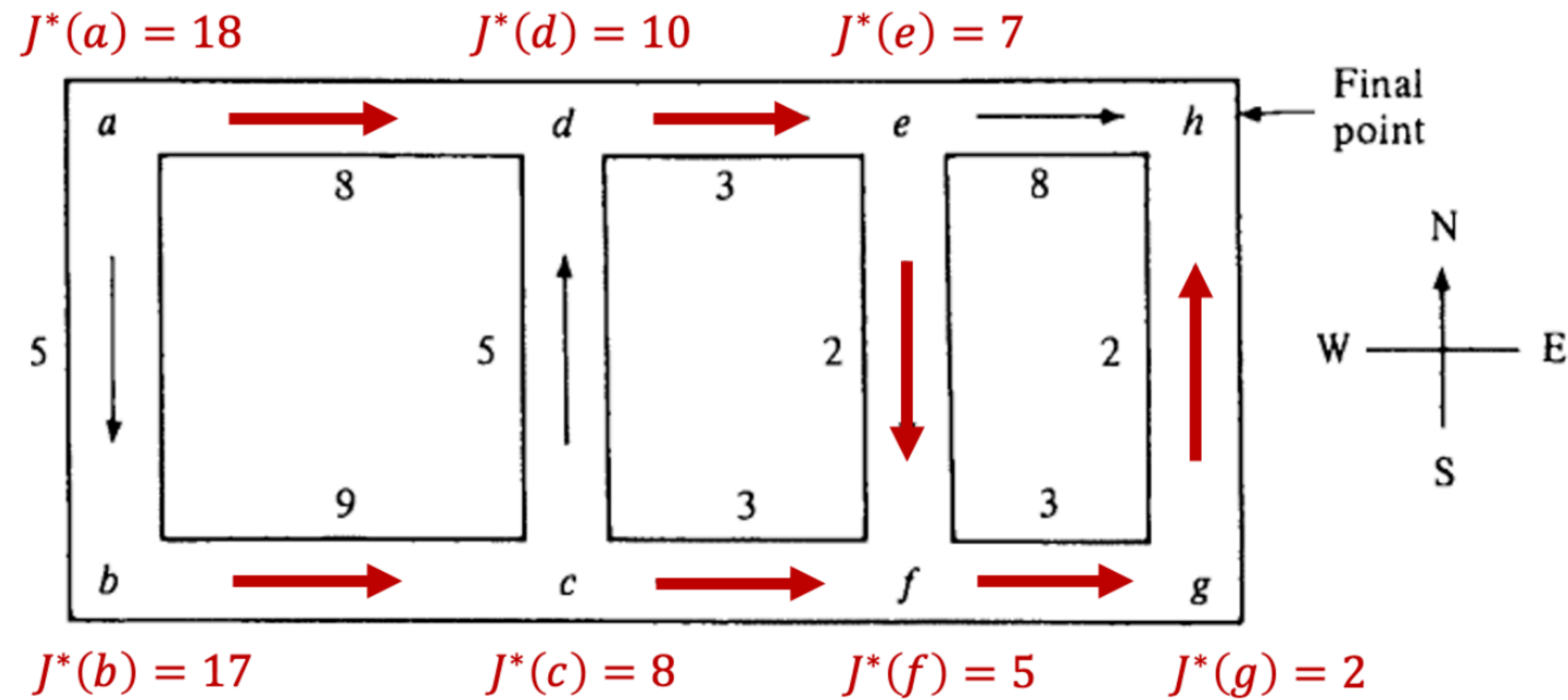
# Applying the principle of optimality

- Need only to compare the concatenations of immediate decisions and optimal decisions → significant decrease in computation/possibilities
- In practice: carry out this procedure **backward** in time

# Example



# Example



# DP Algorithm

$$\text{Model: } \mathbf{x}_{k+1} = f(\mathbf{x}_k, \mathbf{u}_k, k), \quad \mathbf{u}_k \in U(\mathbf{x}_k)$$

$$\text{Cost: } J(\mathbf{x}_0) = h_N(\mathbf{x}_N) + \sum_{k=0}^{N-1} g(\mathbf{x}_k, \pi_k(\mathbf{x}_k), k)$$

## DP Algorithm:

For every initial state  $\mathbf{x}_0$ , the optimal cost  $J^*(\mathbf{x}_0)$  is equal to  $J_0^*(\mathbf{x}_0)$ , given by the last step of the following algorithm, which proceeds backward in time from stage  $N - 1$  to stage  $0$ :

$$J_N^*(\mathbf{x}_N) = h_N(\mathbf{x}_N)$$

$$J_k^*(\mathbf{x}_k) = \min_{\mathbf{u}_k \in U(\mathbf{x}_k)} g(\mathbf{x}_k, \mathbf{u}_k, k) + J_{k+1}^*(f(\mathbf{x}_k, \mathbf{u}_k, k)), \quad k = 0, \dots, N - 1$$

Furthermore, if  $\mathbf{u}_k^* = \pi_k^*(\mathbf{x}_k)$  minimizes the right hand side of the above equation for each  $\mathbf{x}_k$  and  $k$ , the policy  $\{\pi_0^*, \pi_1^*, \dots, \pi_{N-1}^*\}$  is optimal

# Comments

- Discretization (from differential equations to difference equations)
- Quantization (from continuous to discrete state variables / controls)
- Guaranteed to converge to a global minimum
- Constraints, in general, simplify the numerical procedure
- Optimal control in **closed-loop** form
- Curse of dimensionality (both computationally and w.r.t. memory)
- Typically involves
  - *offline* computation of optimal costs (backward)
  - *online* planning through (forward) construction of solution

# Outline

Intro to dynamic programming (DP) and *principle of optimality*

The dynamic programming algorithm

Dynamic programming in control:

- Discrete LQR
- Stochastic Optimal Control Problem / Markov Decision Process (MDP): Stochastic LQR
- Policy Iteration and Value Iteration

# Discrete LQR

- Canonical application of dynamic programming for control
- One case where DP can be solved analytically (in general, DP algorithm must be performed numerically)

**Discrete (Deterministic) LQR** : Select control inputs to minimize

$$J_0(\mathbf{x}_0) = \frac{1}{2} \mathbf{x}_N^T \mathbf{Q}_N \mathbf{x}_N + \frac{1}{2} \sum_{k=0}^{N-1} (\mathbf{x}_k^T \mathbf{Q}_k \mathbf{x}_k + \mathbf{u}_k^T \mathbf{R}_k \mathbf{u}_k + 2\mathbf{x}_k^T \mathbf{S}_k \mathbf{u}_k)$$

Subject to dynamics

$$\mathbf{x}_{k+1} = \mathbf{A}_k \mathbf{x}_k + \mathbf{B}_k \mathbf{u}_k, \quad k \in \{0, 1, \dots, N-1\}$$

Assuming

$$\mathbf{Q}_k = \mathbf{Q}_k^T \succeq 0, \quad \mathbf{R}_k = \mathbf{R}_k^T \succ 0, \quad \begin{bmatrix} \mathbf{Q}_k & \mathbf{S}_k \\ \mathbf{S}_k^T & \mathbf{R}_k \end{bmatrix} \succeq 0 \quad \forall k$$



# Extensions

Many important extensions, some of which we'll cover later in this class

- Cost with linear terms, affine dynamics: can consider today's analysis with augmented dynamics

$$\mathbf{y}_{k+1} = \begin{bmatrix} \mathbf{x}_{k+1} \\ 1 \end{bmatrix} = \begin{bmatrix} A_k & c_k \\ 0 & 1 \end{bmatrix} \begin{bmatrix} \mathbf{x}_k \\ 1 \end{bmatrix} + \begin{bmatrix} B \\ 0 \end{bmatrix} \mathbf{u}_k = \tilde{A} \mathbf{y}_k + \tilde{B} \mathbf{u}_k$$

- Tracking LQR:  $\mathbf{x}_k$ ,  $\mathbf{u}_k$  represent small deviations (“errors”) from a nominal trajectory (possibly with nonlinear dynamics)

- Stochastic systems

$$\mathbf{x}_{k+1} = A_k \mathbf{x}_k + B_k \mathbf{u}_k + \mathbf{w}_k, \quad \mathbf{w}_k \sim \mathcal{N}(0, \Sigma_{\mathbf{w}_k})$$



# Discrete LQR - Trajectory Optimization

We can then solve this problem by applying the NOC (which, due to the problem's convexity, are also SOC)

$$\begin{aligned} \min_{\mathbf{z}} \quad & \frac{1}{2} \mathbf{z}^T \mathbf{W} \mathbf{z} \\ \text{s.t.} \quad & \mathbf{C} \mathbf{z} + \mathbf{d} = \mathbf{0} \end{aligned}$$

Specifically:

$$L(\mathbf{z}, \boldsymbol{\lambda}) = \frac{1}{2} \mathbf{z}^T \mathbf{W} \mathbf{z} + \boldsymbol{\lambda}^T (\mathbf{C} \mathbf{z} + \mathbf{d})$$

$$\nabla_{\mathbf{z}} L = \frac{1}{2} \mathbf{W} \mathbf{z} + \frac{1}{2} \mathbf{W}^T \mathbf{z} + \mathbf{C}^T \boldsymbol{\lambda} = \mathbf{W} \mathbf{z} + \mathbf{C}^T \boldsymbol{\lambda} = \mathbf{0}$$

$$\nabla_{\boldsymbol{\lambda}} L = \mathbf{C} \mathbf{z} + \mathbf{d} = \mathbf{0}$$

Compactly,  $\begin{bmatrix} \mathbf{z}^* \\ \boldsymbol{\lambda}^* \end{bmatrix} = \begin{bmatrix} \mathbf{W} & \mathbf{C}^T \\ \mathbf{C} & \mathbf{0} \end{bmatrix}^{-1} \begin{bmatrix} \mathbf{0} \\ -\mathbf{d} \end{bmatrix}$  Solving this requires  $\mathcal{O}[(N(m+n))^3]$

# Discrete LQR - Dynamic programming

Solving through DP allows us to

(1) Solve in  $\mathcal{O}[N(m+n)^3]$  vs  $\mathcal{O}[(N(m+n))^3]$

(2) Obtain a *closed-loop* policy  $\pi(\mathbf{x}_t)$

First step: 
$$J_N^*(\mathbf{x}_N) = \frac{1}{2} \mathbf{x}_N^T Q_N \mathbf{x}_N = \frac{1}{2} \mathbf{x}_N^T P_N \mathbf{x}_N$$

Proceeding backward in time:

$$\begin{aligned} J_{N-1}^*(\mathbf{x}_{N-1}) &= \min_{\mathbf{u}_{N-1}} \frac{1}{2} \left( \begin{bmatrix} \mathbf{x}_{N-1} \\ \mathbf{u}_{N-1} \end{bmatrix}^T \begin{bmatrix} Q_{N-1} & S_{N-1} \\ S_{N-1}^T & R_{N-1} \end{bmatrix} \begin{bmatrix} \mathbf{x}_{N-1} \\ \mathbf{u}_{N-1} \end{bmatrix} + \mathbf{x}_N^T P_N \mathbf{x}_N \right) \\ &= \min_{\mathbf{u}_{N-1}} \frac{1}{2} \left( \begin{bmatrix} \mathbf{x}_{N-1} \\ \mathbf{u}_{N-1} \end{bmatrix}^T \begin{bmatrix} Q_{N-1} & S_{N-1} \\ S_{N-1}^T & R_{N-1} \end{bmatrix} \begin{bmatrix} \mathbf{x}_{N-1} \\ \mathbf{u}_{N-1} \end{bmatrix} + \right. \\ &\quad \left. (A_{N-1} \mathbf{x}_{N-1} + B_{N-1} \mathbf{u}_{N-1})^T P_N (A_{N-1} \mathbf{x}_{N-1} + B_{N-1} \mathbf{u}_{N-1}) \right) \end{aligned}$$

# Discrete LQR - Dynamic programming

Unconstrained NOC:

$$\nabla_{u_{N-1}} J_{N-1}(\mathbf{x}_{N-1}) = R_{N-1} \mathbf{u}_{N-1} + S_{N-1}^T \mathbf{x}_{N-1} + B_{N-1}^T P_N (A_{N-1} \mathbf{x}_{N-1} + B_{N-1} \mathbf{u}_{N-1}) = \mathbf{0}$$

$$\begin{aligned} \implies \mathbf{u}_{N-1}^* &= -(R_{N-1} + \underbrace{B_{N-1}^T P_N B_{N-1}}_{\text{does not depend on time}})^{-1} (B_{N-1}^T P_N A_{N-1} + S_{N-1}^T) \mathbf{x}_{N-1} \\ &:= F_{N-1} \mathbf{x}_{N-1} \end{aligned}$$

Note also that SOC hold:  $\nabla_{u_{N-1}}^2 J_{N-1}(\mathbf{x}_{N-1}) = R_{N-1} + B_{N-1}^T P_N B_{N-1} \succ 0$

To obtain the optimal cost-to-go, we plug in the optimal policy to obtain:

$$\begin{aligned} J_{N-1}^*(\mathbf{x}_{N-1}) &= \frac{1}{2} \mathbf{x}_{N-1}^T (Q_{N-1} + A_{N-1}^T P_N A_{N-1} - \\ &\quad (A_{N-1}^T P_N B_{N-1} + S_{N-1})(R_{N-1} + B_{N-1}^T P_N B_{N-1})^{-1} (B_{N-1}^T P_N A_{N-1} + S_{N-1}^T)) \mathbf{x}_{N-1} \\ &:= \frac{1}{2} \mathbf{x}_{N-1}^T P_{N-1} \mathbf{x}_{N-1} \end{aligned}$$

Notice that:

- The optimal policy is a time-varying linear feedback policy (i.e., we can just store the matrices  $F_k$ )
- The cost-to-go is a quadratic function of the state at each step (!)

Additionally:

- In the infinite horizon case, this is guaranteed to converge to the optimal policy (as long as there exist a policy that can drive the system to zero)
- Often most convenient to use steady state  $F_\infty$

# Discrete LQR - Dynamic programming

Proceeding by induction, we derive the Riccati recursion:

1.  $P_N = Q_N$
2.  $F_k = -(R_k + B_k^T P_{k+1} B_k)^{-1} (B_k^T P_{k+1} A_k + S_k^T)$
3.  $P_k = Q_k + A_k^T P_{k+1} A_k -$   
 $(A_k^T P_{k+1} B_k + S_k)(R_k + B_k^T P_{k+1} B_k)^{-1} (B_k^T P_{k+1} A_k + S_k^T)$
4.  $\pi_k^*(\mathbf{x}_k) = F_k \mathbf{x}_k$
5.  $J_k^*(\mathbf{x}_k) = \frac{1}{2} \mathbf{x}_k^T P_k \mathbf{x}_k$

Which enables us to

- Compute the policy backwards in time (and store it)
- Apply the policy forward in time

# Outline

Intro to dynamic programming (DP) and *principle of optimality*

The dynamic programming algorithm

Dynamic programming in control:

- Discrete LQR
- Stochastic Optimal Control Problem / Markov Decision Process (MDP): Stochastic LQR
- Policy Iteration and Value Iteration

# Stochastic Optimal Control Problem: Markov Decision Problem (MDP)

- **System:**  $\mathbf{x}_{k+1} = f_k(\mathbf{x}_k, \mathbf{u}_k, \mathbf{w}_k), k = 0, \dots, N-1$
- **Probability distribution:**  $\mathbf{w}_k \sim P_k(\cdot | \mathbf{x}_k, \mathbf{u}_k)$
- **Control constraints:**  $\mathbf{u}_k \in U(\mathbf{x}_k)$
- **Policies:**  $\pi = \{\pi_0, \dots, \pi_{N-1}\}$ , where  $\mathbf{u}_k = \pi_k(\mathbf{x}_k)$
- **Expected Cost:**

$$J_\pi(\mathbf{x}_0) = \mathbb{E}_{\mathbf{w}_k, k=0, \dots, N-1} \left[ g_N(\mathbf{x}_N) + \sum_{k=0}^{N-1} g_k(\mathbf{x}_k, \pi_k(\mathbf{x}_k), \mathbf{w}_k) \right]$$

**Stochastic Optimal Control Problem:**

$$J^*(x_0) = \min_{\pi} J_\pi(x_0)$$



# Key points

- Discrete-time model
- Markovian model
- Objective: find optimal **closed-loop** policy
- Additive cost (central assumption in DP)
- Risk-neutral formulation

Other communities use different notation:

[Powell, W. B. *AI, OR and control theory: A Rosetta Stone for stochastic optimization*. Princeton University, 2012.]

# The DP algorithm (stochastic case)

## Principle of optimality:

- Let  $\pi^* := \left\{ \pi_0^*, \pi_1^*, \dots, \pi_{N-1}^* \right\}$  be an optimal policy
- Consider the tail subproblem

$$\mathbb{E}_{w_k} \left[ g_N(\mathbf{x}_N) + \sum_{k=i}^{N-1} g_k(\mathbf{x}_k, \pi_k(\mathbf{x}_k), \mathbf{w}_k) \right]$$

the tail policy  $\left\{ \pi_i^*, \dots, \pi_{N-1}^* \right\}$  is optimal for the tail subproblem

## Intuition:

- DP first solves ALL tail subproblems at the final stage
- At the generic step, it solves ALL tail subproblems of a given time length, using solution of tail subproblems of shorter length

# DP Algorithm (stochastic case)

Like in the deterministic case, start with:

$$J_N(x_N) = g_N(x_N)$$

and iterate backwards in time using

$$J_k(x_k) = \min_{u_k \in U(x_k)} \mathbb{E}_{w_k} \left[ g_k(x_k, u_k, w_k) + J_{k+1}(f(x_k, u_k, w_k)) \right], \quad k = 0, \dots, N-1$$

for which the optimal cost  $J^*(\mathbf{x}_0)$  is equal to  $J_0(\mathbf{x}_0)$  and the optimal policy is constructed by setting

$$\pi_k^*(x_k) = \operatorname{argmin}_{u_k \in U(x_k)} \mathbb{E}_{w_k} \left[ g_k(x_k, u_k, w_k) + J_{k+1}(f(x_k, u_k, w_k)) \right]$$

# Next time

- Stochastic Dynamic Programming
- Infinite-Horizon MDPs
- Value Iteration
- Policy Iteration