

goal: mathematical model for optimal control with random dynamics / cost

refs: *Neuro-Dynamic Programming*

Dimitri P. Bertsekas and John N. Tsitsiklis

chapter 2

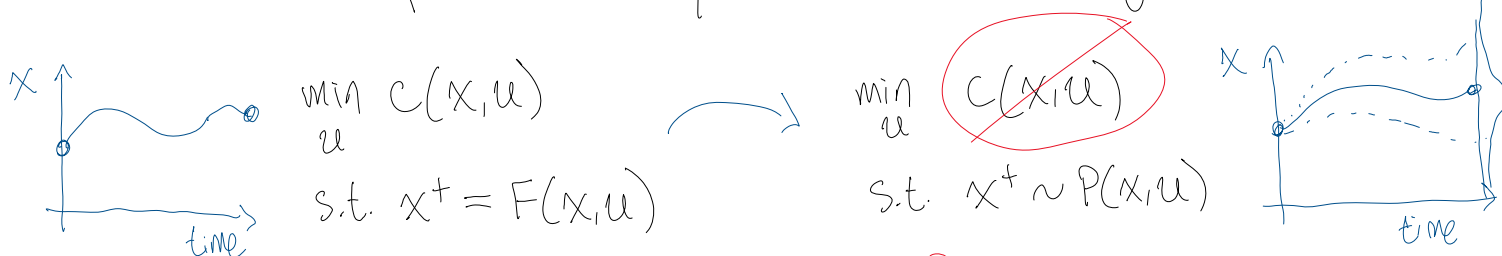
## Reinforcement Learning

An Introduction

Richard S. Sutton and Andrew G. Barto

chapter 3

- analogous to how MP/SDE naturally generalizes deterministic dynamics, we can consider optimal control problems with random dynamics.



→ what should we use for the cost function?  
i.e. what statistic should we choose? (e.g.  $E[\cdot]$ ,  $\text{Var}[\cdot]$ )

-  $E[C(x, u)]$  is the average / "expected" cost

↪ useful when system will run many times and/or unlikely outcomes are ok

- $\text{Var}[c(x,u)]$  is the "spread" in cost  
 $\leadsto$  useful when system will run few times and/or unlikely outcomes catastrophic
- also consider order statistics (median, interquartile, confidence intervals)

\* we'll focus on minimizing expected cost: 
$$\min_u E[c(x,u)]$$
  
s.t.  $x^+ \sim P(x,u)$

where  $c(x,u) = l(t, x_t) + \sum_{s=0}^{t-1} \mathcal{L}(s, x_s, u_s)$  — finite-horizon

or  $= \lim_{t \rightarrow \infty} \frac{1}{t} \sum_{s=0}^{t-1} \mathcal{L}(s, x_s, u_s)$  — infinite-horizon, time-averaged

or  $= \lim_{t \rightarrow \infty} \sum_{s=0}^t \gamma^s \cdot \mathcal{L}(s, x_s, u_s)$  — infinite-horizon, exponentially-discounted  
 $\uparrow \gamma \in (0,1)$   
termed a discount factor

def: Markov decision process (MDP) / stochastic optimal control prob. (SOP)

specified by  $(X, U, P, c)$  where  $X, U$  are sets,

$P: X \times U \rightarrow \Delta(X)$ ,  $c: X^T \times U^T \rightarrow \mathbb{R}$ ,  $T$  is a time interval, e.g.  
 $\uparrow \{x: T \rightarrow X\}$   $= [0, t]$  or  $= [0, \infty)$