

# Topics in Reinforcement Learning: Lessons from AlphaZero for (Sub)Optimal Control and Discrete Optimization

Arizona State University  
Course CSE 691, Spring 2023

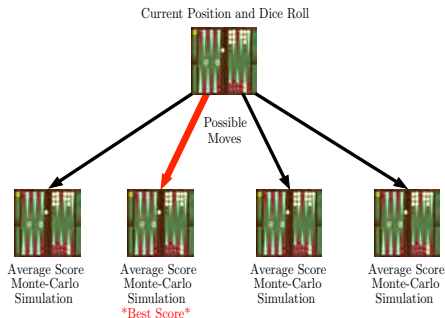
Links to Class Notes, Videolectures, and Slides at  
<http://web.mit.edu/dimitrib/www/RLbook.html>

Dimitri P. Bertsekas  
[dimitrib@mit.edu](mailto:dimitrib@mit.edu), [dbertsek@asu.edu](mailto:dbertsek@asu.edu)

Lecture 1  
Course Introduction and Overview

- 1 AlphaZero - Off-Line Training and On-Line Play
- 2 History, General Concepts
- 3 About the Course and its Connections to Various Fields
- 4 Dynamic Programming - Deterministic Problems
- 5 Examples: Finite-State/Discrete/Combinatorial DP Problems
- 6 Organizational Issues

# Chess and Backgammon - Off-Line Training and On-Line Play



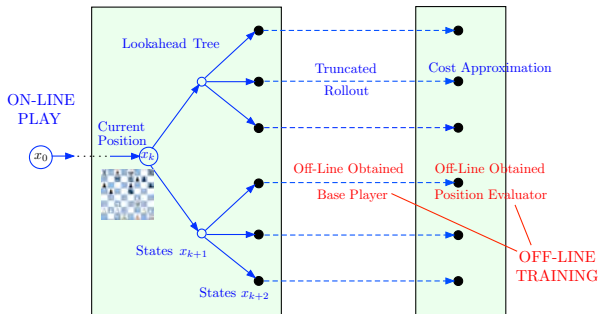
Both AlphaZero (2017) and TD-Gammon (1996) involve **two algorithms**:

- **Off-line training** of value and/or policy neural network approximations
- **On-line play** by multistep lookahead, rollout, and cost function approximation

Strong connections to DP, policy iteration, and RL-type methodology

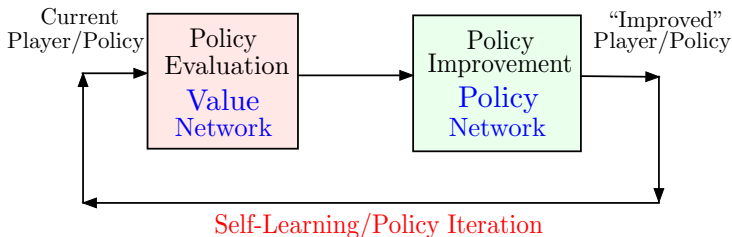
- We aim to understand this methodology, so it applies far more generally
- For example, in **control system design** (model predictive and adaptive control), and discrete optimization by rollout

# On-Line Play in AlphaZero/AlphaGo/TD-Gammon: Approximation in Value Space (Also Called “On-Line Tree Search”)



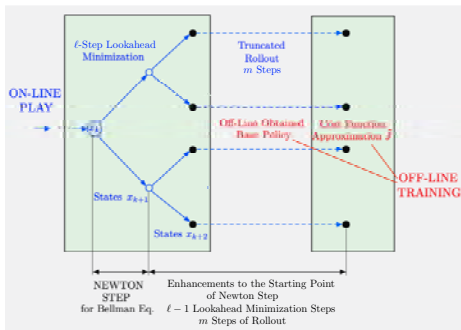
- On-line play uses the results of off-line training, which are: A position evaluator and a base player
- It aims to improve the base player by:
  - ▶ Searching forward for several moves through the lookahead tree
  - ▶ Simulating the base player for some more moves at the tree leaves
  - ▶ Approximating the effect of future moves by using the terminal position evaluation
  - ▶ Calculating the “values” of the available moves at the root and playing the best move
- Similarities with Model Predictive Control (MPC) (which involves continuous spaces) and discrete optimization by rollout (which uses a heuristic as base player)

# Off-Line Training in AlphaZero: Approximate Policy Iteration (PI) Using Self-Generated Data



- The current player is used to train an improved player, and the process is repeated
- The current player is “evaluated” by playing many games
- Its evaluation function is represented by a value neural net through training
- The current player is “improved” by using a form of approximate multistep lookahead minimization, called Monte-Carlo Tree Search (MCTS)
- The “improved player” is represented by a policy neural net through training
- TD-Gammon uses similar PI algorithm for off-line training of a value network (does not use MCTS and does not use a policy network)
- MPC and discrete optimization by rollout often use rudimentary forms of off-line training

# Some Major Empirical Observations



The AlphaZero on-line player plays much better than the off-line-trained player

TD-Gammon plays much better with truncated rollout than without rollout (Tesauro, 1996)

We will aim for explanations, insights, and generalizations through abstract Bellman operators, visualization, and the central role of Newton's method

# Aims of the Course

Provide a unifying framework for several areas of large scale computation:

- **Reinforcement learning** (RL) as practiced by the AI community
- **Approximate dynamic programming** (DP) as practiced by parts of the optimization/OR community
- **Model predictive and adaptive control** as practiced by the control systems community
- Parts of **discrete optimization** as practiced by the algorithms/CS community

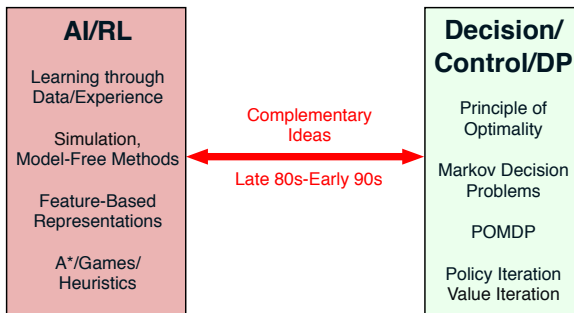
We rely on:

- The algorithmic theory of **exact, approximate, and abstract DP**
- **The paradigm of AlphaZero/TD-Gammon** and similar design architectures
- Intuitive **visualization** based on a Bellman operator formalism and Newton's method

We aim, through unification and abstraction, to:

- **Bridge the gap between cultures** of different communities
- Bring to bear the power of RL to a **very broad range of applications**

# Evolution of Approximate DP/RL: A Fruitful Synergy



## Historical highlights

- Exact DP, optimal control (Bellman, Shannon, and others 1950s ...)
- **AI/RL and Decision/Control/DP ideas meet** (late 80s-early 90s)
- First major successes: Backgammon programs (Tesauro, 1992, 1996)
- Algorithmic progress, analysis, applications, first books (mid 90s ...)
- Machine Learning, BIG Data, Robotics, Deep Neural Networks (mid 2000s ...)
- AlphaGo and AlphaZero (DeepMind, 2016, 2017)



# Approximate DP/RL Methodology is now Ambitious and Universal

Exact DP applies (in principle) to a very broad range of optimization problems

- Deterministic  $\longleftrightarrow$  Stochastic
- Combinatorial optimization  $\longleftrightarrow$  Optimal control w/ infinite state/control spaces
- One decision maker  $\longleftrightarrow$  Two player games
- ... BUT is plagued by the **curse of dimensionality** and **need for a math model**

Approximate DP/RL overcomes the difficulties of exact DP by:

- **Approximation** (use neural nets and other architectures to reduce dimension)
- **Simulation** (use a computer model in place of a math model)

State of the art:

- **Broadly applicable methodology**: Can address a very broad range of challenging problems. Deterministic-stochastic-dynamic, discrete-continuous, games, etc
- There are **no methods that are guaranteed to work** for all or even most problems
- There are **enough methods to try with a reasonable chance of success** for most types of optimization problems
- **Role of the theory**: Structure mathematically the methodology, guide the art, delineate the sound ideas (from the crazy and unhinged ideas)

## A Relevant Quotation from 25 Years Ago

From preface of Neuro-Dynamic Programming, Bertsekas and Tsitsiklis, 1996

A few years ago our curiosity was aroused by reports on new methods in reinforcement learning, a field that was developed primarily within the artificial intelligence community, starting a few decades ago. **These methods were aiming to provide effective suboptimal solutions to complex problems of planning and sequential decision making under uncertainty, that for a long time were thought to be intractable.**

**Our first impression was that the new methods were ambitious, overly optimistic, and lacked firm foundation.** Yet there were claims of impressive successes and indications of a solid core to the modern developments in reinforcement learning, suggesting that the correct approach to their understanding was through dynamic programming.

Three years later, **after a lot of study, analysis, and experimentation, we believe that our initial impressions were largely correct.** This is indeed an ambitious, often ad hoc, methodology, but for reasons that we now understand much better, **it does have the potential of success with important and challenging problems.**

**This assessment still holds true!**

# References of this Course

This course is research-oriented. It aims:

- To explore the state of the art of approximate DP/RL at a graduate level
- To explore in depth some special research topics (rollout, policy iteration)
- To provide the opportunity for you to explore research in the area

Main references:

- Bertsekas, Reinforcement Learning and Optimal Control, 2019
- Bertsekas, Rollout, Policy Iteration, and Distributed Reinforcement Learning, 2020
- Bertsekas, Lessons from AlphaZero for Optimal, Model Predictive, and Adaptive Control, 2022 (on-line)
- Bertsekas: Class notes based on the above, and focused on our special RL topics.
- Slides, papers, and videos from the 2019-2022 ASU courses; check my web site

Supplementary references

- Exact DP: Bertsekas, Dynamic Programming and Optimal Control, Vol. I (2017), Vol. II (2012) (also contains approximate DP material), Abstract DP (2022)
- Bertsekas and Tsitsiklis, Neuro-Dynamic Programming, 1996
- Sutton and Barto, 1998, Reinforcement Learning (new edition 2018, on-line)

## RL uses Max/Value, DP uses Min/Cost

- **Reward of a stage** = (Opposite of) Cost of a stage.
- **State value** = (Opposite of) State cost.
- **Value (or state-value) function** = (Opposite of) Cost function.

## Controlled system terminology

- **Agent** = Decision maker or controller.
- **Action** = Decision or control.
- **Environment** = Dynamic system.

## Methods terminology

- **Learning** = Solving a DP-related problem using simulation.
- **Self-learning (or self-play in the context of games)** = Solving a DP problem using simulation-based policy iteration.
- **Planning vs Learning distinction** = Solving a DP problem with model-based vs model-free simulation.

- Reinforcement learning uses **transition probability notation**

$$p(s, a, s')$$

( $s, s'$  are states,  $a$  is action), which is standard in finite-state problems (MDP)

- Control theory uses **discrete-time system equation**

$$x_{k+1} = f(x_k, u_k, w_k)$$

which is standard in continuous spaces problems

- Operations research uses both notations [typically  $p_{ij}(u)$  for transition probabilities]

**These two notational systems are mathematically equivalent but:**

- Transition probabilities are cumbersome for deterministic problems and continuous spaces problems
- System equations are cumbersome for finite-state problems

**We will use both notational systems, depending on the context**

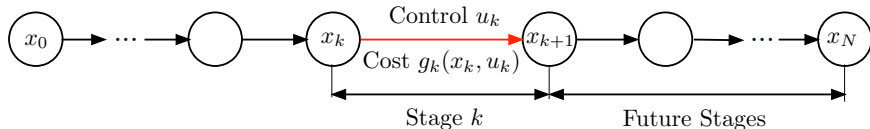
## A Fifteen-Minute Break

All our lectures will have a 15-minute break, somewhere in the middle

Catch our breath and think about issues relating to the first half of the lecture.

A short discussion/questions/answers period will follow each break.

# Finite Horizon Deterministic Optimal Control Model



- System

$$x_{k+1} = f_k(x_k, u_k), \quad k = 0, 1, \dots, N-1$$

where  $x_k$ : State,  $u_k$ : Control chosen from some set  $U_k(x_k)$

- Cost function:

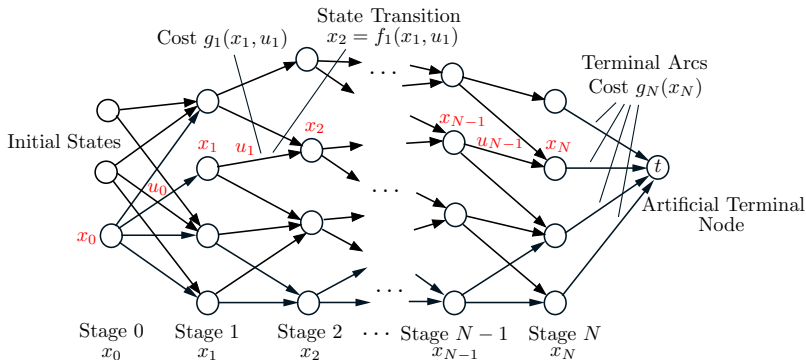
$$g_N(x_N) + \sum_{k=0}^{N-1} g_k(x_k, u_k)$$

- For given initial state  $x_0$ , minimize over control sequences  $\{u_0, \dots, u_{N-1}\}$

$$J(x_0; u_0, \dots, u_{N-1}) = g_N(x_N) + \sum_{k=0}^{N-1} g_k(x_k, u_k)$$

- Optimal cost function  $J^*(x_0) = \min_{\substack{u_k \in U_k(x_k) \\ k=0, \dots, N-1}} J(x_0; u_0, \dots, u_{N-1})$

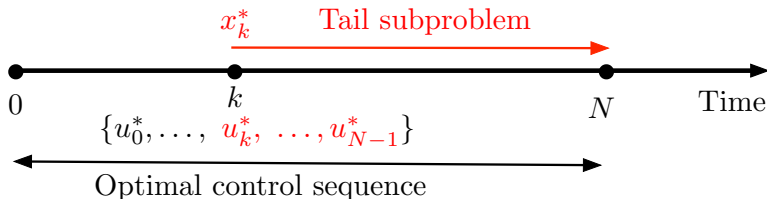
# A Special Case: Finite Number of States and Controls



- Nodes correspond to states  $x_k$
- Each arc corresponds to a state-control pair  $(x_k, u_k)$  (start node is  $x_k$ ; end node is  $x_{k+1} = f_k(x_k, u_k)$ )
- An arc corresponding to  $(x_k, u_k)$  has a cost  $g_k(x_k, u_k)$ .
- The cost to optimize is the sum of the arc costs from the initial node/state  $x_0$  to a terminal node  $t$ .
- The problem is equivalent to finding a minimum cost/shortest path from  $x_0$  to  $t$ .



## Principle of Optimality: A Very Simple Idea



### Principle of Optimality

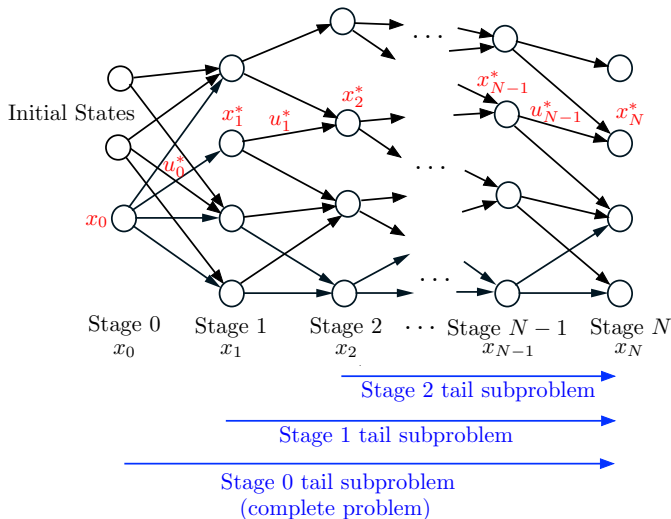
**THE TAIL OF AN OPTIMAL SEQUENCE IS OPTIMAL FOR THE TAIL SUBPROBLEM**

Let  $\{u_0^*, \dots, u_{N-1}^*\}$  be an optimal control sequence with corresponding state sequence  $\{x_1^*, \dots, x_N^*\}$ . Consider the **tail subproblem** that starts at  $x_k^*$  at time  $k$  and minimizes over  $\{u_k, \dots, u_{N-1}\}$  the “cost-to-go” from  $k$  to  $N$ ,

$$g_k(x_k^*, u_k) + \sum_{m=k+1}^{N-1} g_m(x_m, u_m) + g_N(x_N).$$

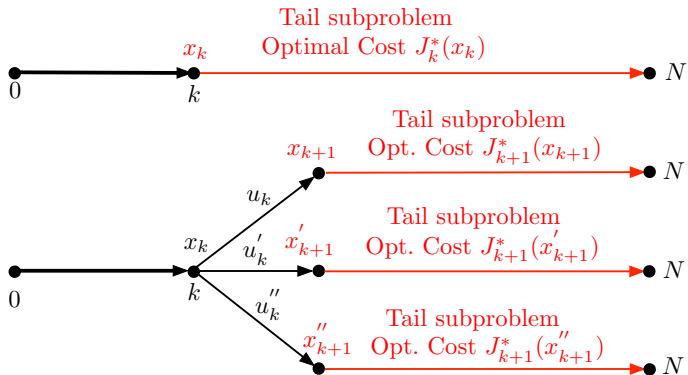
Then the tail optimal control sequence  $\{u_k^*, \dots, u_{N-1}^*\}$  is optimal for the tail subproblem.

# Principle of Optimality for the Generic Finite-State Problem



THE TAIL OF AN OPTIMAL SEQUENCE IS OPTIMAL FOR THE TAIL SUBPROBLEM

## From One Tail Subproblem to the Next



By the principle of optimality: To solve the tail subproblem that starts at  $x_k$

- Consider every possible  $u_k$  and solve the tail subproblem that starts at next state  $x_{k+1} = f_k(x_k, u_k)$ . This gives the "cost of  $u_k$ "
- Optimize over all possible  $u_k$

# DP Algorithm: Solves All Tail Subproblems Efficiently by Using the Principle of Optimality

## Idea of the DP algorithm

Solve **all** the tail subproblems of a given time length using the solution of **all the tail subproblems of shorter time length**

DP Algorithm: Produces the optimal costs  $J_k^*(x_k)$  of the  $x_k$ -tail subproblems

Start with

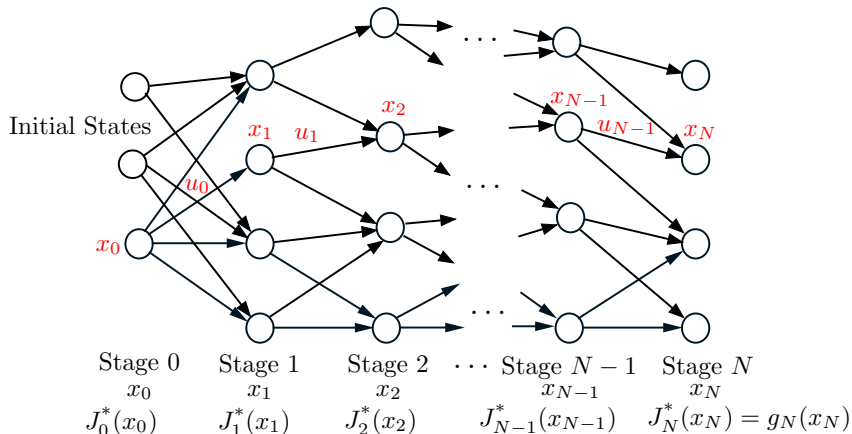
$$J_N^*(x_N) = g_N(x_N), \quad \text{for all } x_N,$$

and for  $k = 0, \dots, N - 1$ , let

$$J_k^*(x_k) = \min_{u_k \in U_k(x_k)} \left[ g_k(x_k, u_k) + J_{k+1}^*(f_k(x_k, u_k)) \right], \quad \text{for all } x_k.$$

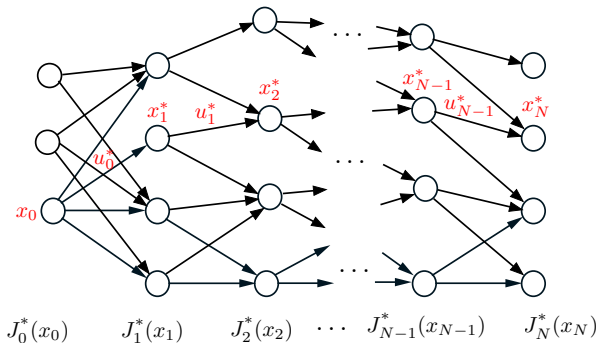
The optimal cost  $J^*(x_0)$  is obtained at the last step

# DP Algorithm for Generic Finite-State Problem



$$J_k^*(x_k) = \min_{u_k \in U_k(x_k)} \left[ g_k(x_k, u_k) + J_{k+1}^*(f_k(x_k, u_k)) \right], \quad \text{for all } x_k.$$

# Construction of Optimal Control Sequence $\{u_0^*, \dots, u_{N-1}^*\}$



Start with

$$u_0^* \in \arg \min_{u_0 \in U_0(x_0)} \left[ g_0(x_0, u_0) + J_1^*(f_0(x_0, u_0)) \right]$$

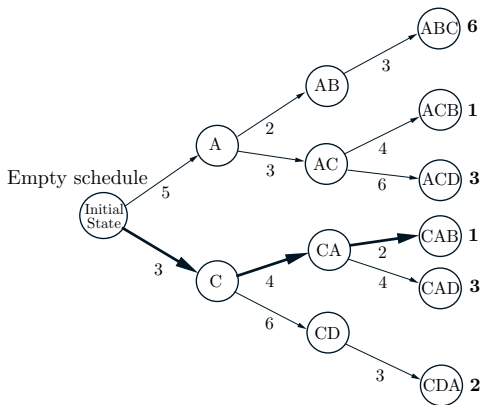
This takes you to

$$x_1^* = f_0(x_0, u_0^*).$$

Sequentially, going forward, for  $k = 1, 2, \dots, N-1$ , set

$$u_k^* \in \arg \min_{u_k \in U_k(x_k^*)} \left[ g_k(x_k^*, u_k) + J_{k+1}^*(f_k(x_k^*, u_k)) \right], \quad x_{k+1}^* = f_k(x_k^*, u_k^*).$$

# Discrete-State Deterministic Scheduling Example

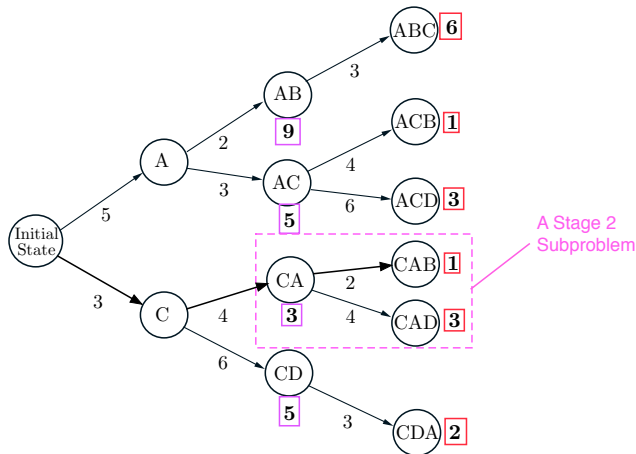


Find optimal sequence of operations A, B, C, D (A must precede B and C must precede D)

## DP Problem Formulation

- States: Partial schedules; Controls: Stage 0, 1, and 2 decisions; Cost data shown along the arcs
- Recall the DP idea: Break down the problem into smaller pieces (tail subproblems)
- Start from the last decision and go backwards

## DP Algorithm: Stage 2 Tail Subproblems

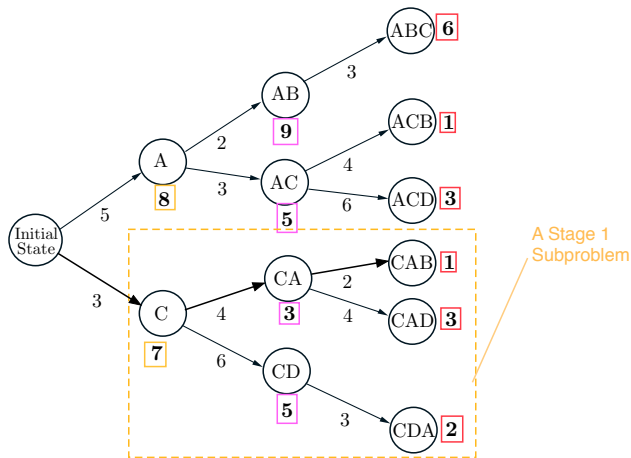


Solve the stage 2 subproblems (using the terminal costs - in red)

At each state of stage 2, we record the optimal cost-to-go and the optimal decision



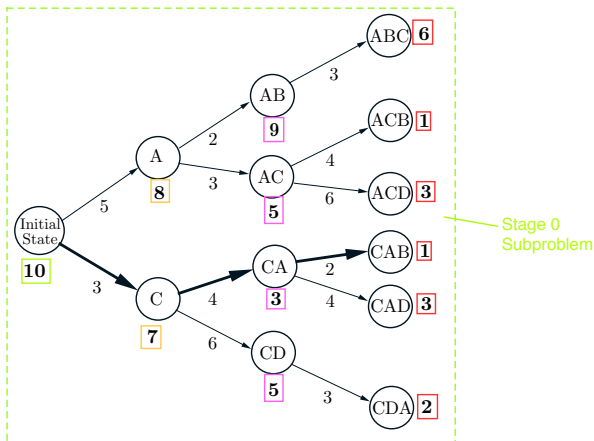
## DP Algorithm: Stage 1 Tail Subproblems



Solve the stage 1 subproblems (using the optimal costs of stage 2 subproblems - in purple)

At each state of stage 1, we record the optimal cost-to-go and the optimal decision

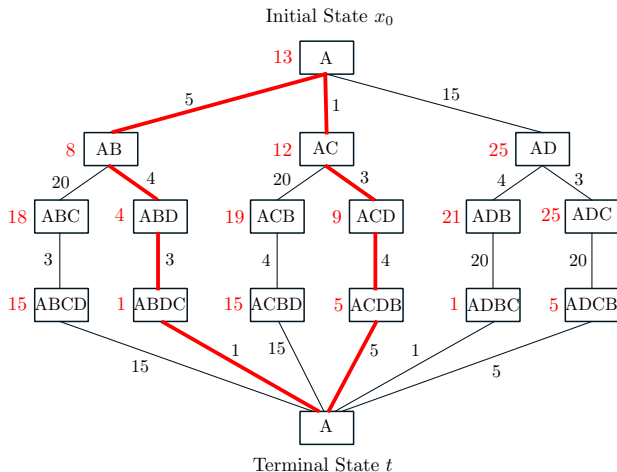
## DP Algorithm: Stage 0 Tail Subproblems



Solve the stage 0 subproblem (using the optimal costs of stage 1 subproblems - in orange)

- The stage 0 subproblem is the entire problem
- The optimal value of the stage 0 subproblem is the optimal cost  $J^*$ (initial state)
- Construct the optimal sequence going forward

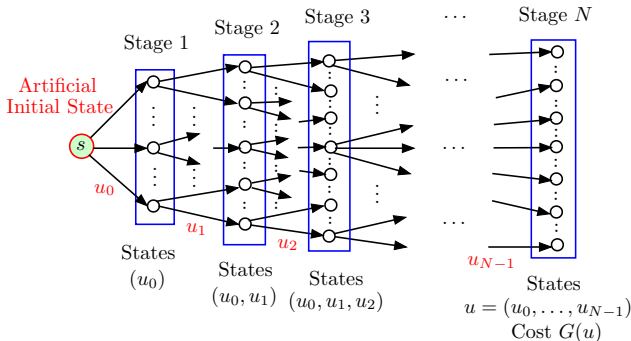
# Discrete Optimization: Traveling Salesman Example; Cities A,B,C,D



Matrix of Intercity  
Travel Costs

	5	1	15
5		20	4
1	20		3
15	4	3	

# General Discrete Optimization



Minimize  $G(u)$  subject to  $u \in U$

- Assume that **each solution  $u$  has  $N$  components:  $u = (u_0, \dots, u_{N-1})$**
- View the components as the controls of  $N$  stages
- Define  $x_k = (u_0, \dots, u_{k-1})$ ,  $k = 1, \dots, N$ , and introduce artificial start state  $x_0 = s$
- Define just terminal cost as  $G(u)$ ; all other costs are 0

**This formulation typically makes little sense for exact DP, but often makes a lot of sense for approximate DP/approximation in value space**

# Approximate DP Algorithm - Connection to Reinforcement Learning

**Exact DP algorithm:** Start with

$$u_0^* \in \arg \min_{u_0 \in U_0(x_0)} \left[ g_0(x_0, u_0) + J_1^*(f_0(x_0, u_0)) \right]$$

This takes you to

$$x_1^* = f_0(x_0, u_0^*).$$

Sequentially, going forward, for  $k = 1, 2, \dots, N - 1$ , set

$$u_k^* \in \arg \min_{u_k \in U_k(x_k^*)} \left[ g_k(x_k^*, u_k) + J_{k+1}^*(f_k(x_k^*, u_k)) \right], \quad x_{k+1}^* = f_k(x_k^*, u_k^*).$$

**Approximation in Value Space - Use Some  $\tilde{J}_k$  in Place of  $J_k^*$  (off-line training)**

Start with

$$\tilde{u}_0 \in \arg \min_{u_0 \in U_0(x_0)} \left[ g_0(x_0, u_0) + \tilde{J}_1(f_0(x_0, u_0)) \right]$$

This takes you to

$$\tilde{x}_1 = f_0(x_0, \tilde{u}_0).$$

Sequentially, going forward, for  $k = 1, 2, \dots, N - 1$ , set (**on-line play**)

$$\tilde{u}_k \in \arg \min_{u_k \in U_k(\tilde{x}_k)} \left[ g_k(\tilde{x}_k, u_k) + \tilde{J}_{k+1}(f_k(\tilde{x}_k, u_k)) \right], \quad \tilde{x}_{k+1} = f_k(\tilde{x}_k, \tilde{u}_k).$$

## Stochastic finite horizon problems

The next state  $x_{k+1}$  is also affected by a **random parameter** (in addition to  $x_k$  and  $u_k$ ).  
**More difficult than deterministic** (not equivalent to a shortest path problem).

## Infinite horizon problems

The exact DP theory is **mathematically more complex**, but also more elegant.

## Stochastic partial state information problems

We will **convert them to problems of perfect state information**, and then apply DP. Very hard to solve even approximately ... but offer great promise for applications.

## Minimax/game problems

The exact DP theory is **substantially more complex** ... but the most spectacular successes of RL involve games. We will treat lightly.

# Course Aims and Requirements

**Our principal aim:** To help you to think about how RL applies to your research interests

## Requirements:

- Homework (25%): A total of 3-4
- Research-oriented term paper (75%). A choice of:
  - ▶ A **mini-research project**. You may work in teams of 1-3 persons. You are encouraged to try. Selected class presentations at the end.
  - ▶ A **read-and-report term paper** based on 2-3 research publications (chosen by you in consultation with me)

**Notation:** People in AI/RL, Control Theory, and Operations Research focus on different problems and use different notations

- **AI/RL and OR focus on discrete/finite-state problems which are stochastic [Markovian Decision Problems (MDP)]**. Use transition probabilities  $p_{ij}(u)$  to describe the uncertainty.
- **Control theorists use system equation notation  $x_{k+1} = f_k(x_k, u_k, w_k)$** . This notation is well-suited for continuous-state problems and deterministic problems.
- **You are strongly encouraged to use the notation and terminology of the course.**

## Topics to be covered

- Introduction to exact and approximate dynamic programming
- Approximation in value and policy space
- Off-line training, on-line play, and Newton's method
- Rollout and approximate policy iteration
- Model predictive and adaptive control
- Multiagent and multiprocessor reinforcement learning
- Training of feature-based approximation architectures and neural networks
- Policy networks and approximation in policy space

The first four lectures will aim to provide an introduction and overview of the subject, which will facilitate selecting and focusing on some research area. The remaining lectures will develop the topics listed above in greater depth.



## Math requirements for this course are modest

Calculus, elementary probability, minimal use of vector-matrix algebra. Our objective is to use math to the extent needed to develop **insight** into the mechanism of various methods, and to be able to start research.

## However a math framework is critically important

Human insight can only develop within some structure of human thought ... math reasoning is most suitable for this purpose

## On machine learning (from NY Times Article, Dec. 2018)

"What is frustrating about machine learning is that the algorithms can't articulate what they're thinking. **We don't know why they work, so we don't know if they can be trusted** ... As human beings, we want more than answers. **We want insight**. This is going to be a source of tension in our interactions with computers from now on."

## About the Next Lecture

### We will cover:

- DP algorithm for stochastic problems
- Approximation in value space

PLEASE READ AS MUCH OF THE CLASS NOTES AS YOU CAN

Watch the video of Lecture 2 of the 2022 or 2021 offering of the class at my web site  
<http://web.mit.edu/dimitrib/www/RLbook.html>