

Dynamic Programming and Reinforcement Learning

Approximate Dynamic Programming (Week 3a)

Rainer Schlosser, Alexander Kastius

Hasso Plattner Institute (EPIC)

May 2, 2022

Outline

- Questions?
- Today: Approximate Dynamic Programming
 Problem Examples
 Forward Dynamic Programming
 Simulation-based Approaches

Recap: Last Week

- Markov Policies in Infinite Horizon MDPs
- Discounting for Future Rewards
- Bellman Equation & Recursive Problem Decomposition
- Value Iteration
- Policy Iteration

Solving MDP Problems

- Continuous Time Problems & Control Theory (not in focus)
- Discrete Time MDP Problems with **Recursive Solutions** (last week)
 - *Time-dependent* Framework, Backward induction (finite time)
 - *Time-independent* Framework, Value & Policy Iteration
 - Optimal numerical solutions via Bellman Equation (backwards)
- Discrete Time MDP Problems with **Approximate Solutions** (today)
 - Relaxation Concepts to Attack Larger Problem Sizes
 - **Simulation-based Heuristics** (today: forward dynamic programming)
 - Basis for Reinforcement Learning



MDP Problems with Different Complexities

Example	Objective	State	Action	Events	Rewards
----------------	------------------	--------------	---------------	---------------	----------------

Airline Tickets

Hotel/Rental/Rail

Apparel/Seasonal/Events

Perishable Products

Inventory Mgmt.

Durable Products

E-Commerce

Resource Allocations

Tetris/Chess/Go

Self-driving

Can We Solve All of Them?

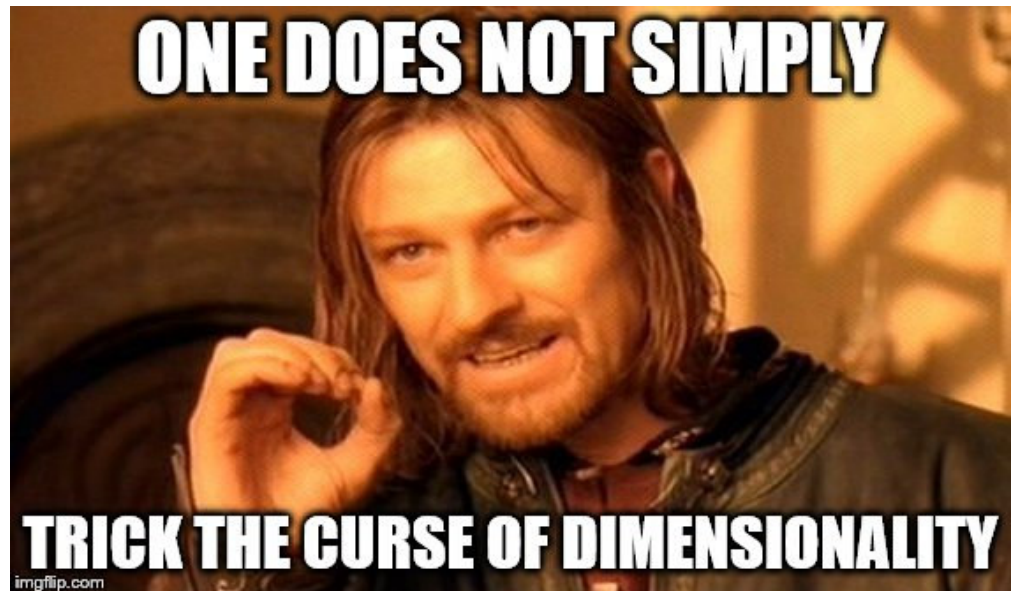
Example **Objective** **State** **Action** **Events** **Rewards**

Airline Tickets

Hotel/Rental/Rail
 Fashion/Seasonal
 Perishable Products

Inventory Mgmt.

Durable Products
 E-Commerce
 Resource Allocations
 Tetris/Chess/Go
 Self-driving . . .



Problem Sizes and Curse of Dimensionality

- State space: Compare $|S| = 10, 100, 1000, 10K, \dots$
- Action space: Compare $|A| = 10, 100, 1000, 10K, \dots$
- Event space: Compare $|I| = 10, 100, 1000, 10K, \dots$
- Time/Iterations Compare $T = 10, 100, 1000, 10K, \dots$ (cf. $\gamma \rightarrow 1$)
- Backward Induction, Policy & Value Iteration **become intractable**
- **Heuristic Options:** Clustering Approaches (not in focus)
 Simulation & Focus on relevant states
 Approximation of Value Functions/Policies

Curse of Dimensionality (Optimal DP Solution)

Example: Sell J types of products with N items each over T periods

Table 5. Optimal expected profits $V_0^*(\vec{N})$ and computation times of (4) - (5) for different $T = 10, 20, 50$ and $N = 5, 10, 20$ with $J = 3$, $\delta = 1$, $c = 0$, $L = 0.05$, $S := \{0, 1, \dots, N\}$, $I := \{0, 1, \dots, 4\}$, and $A := \{4, 8, \dots, 40\}$; Example 3.1.

T	N	$V_0^*(\vec{N})$	time
5	5	326.15	260s
10	5	498.61	641s
10	10	671.57	4 324s
20	10	1 044.49	8 332s
20	20	1 331.67	53 595s
50	10	1 138.12	26 110s
50	20	/	/
100	20	/	/

Schlosser, R. (2021). Scalable Relaxation Techniques to Solve Stochastic Dynamic Multi-Product Pricing Problems with Substitution Effects, *Journal of Revenue and Pricing Management* 20 (1), 54-65.

Approximate Dynamic Programming (ADP)

$$\pi(s) := \arg \max_{a \in A} \left\{ \sum_{i \in I} P(i, a, s) \cdot \left(r(i, a, s) + \gamma \cdot V_{(t+1)}(\Gamma(i, a, s)) \right) \right\}$$


- (1) Use **explicit function approximations** for $V_{(t+1)}(s')$ (offline)
 - **Aggregation**, enforced decomposition (use \tilde{V} of simpler problem)
 - Parametric approximation of $\tilde{V}(s', \theta)$ (NNs, QL, AC, LP, etc.)
- (2) Use **implicit value approximations** for $V_{(t+1)}(s' | a, s)$ (online)
 - **Forward DP** for a, s (via simulation, use full information)
 - Rollout of a heuristic base policy for a, s (via simulation, cp. Pol. It.)
 - Open-loop feedback control (cf. e.g., det. problem version)

(1) Example of Aggregation (Explicit Value Appr.)

Example: Sell J types of products with N items each over T periods


Table 9. Expect profits, cf. (6), and runtimes of a combined heuristic compared to the optimal solution for different $T = 10, 20, 50, 100$ and $N = 5, 10, 20$, cf. Table 5, $S_2 := \{0, 1, 5, N\}$, $I := \{0\}, \{1, 2\}, \{3, 4\}$, $A := \{10, 20, 30, 40\}$; Example 3.1.

relaxed!



T	N	$\bar{V}_0(\bar{N})$	\bar{V}_0/V_0^*	time	%time
5	5	308.90	94.7%	0.45s	0.17%
10	5	468.28	93.9%	0.78s	0.12%
10	10	637.32	94.9%	2.73s	0.06%
20	10	974.42	93.3%	5.64s	0.07%
20	20	1 255.51	94.3%	15.5s	0.03%
50	10	1 102.27	96.9%	14.3s	0.05%
50	20	2 005.91	/	34.0s	/
100	20	2 131.19	/	62.4s	/

relaxed!

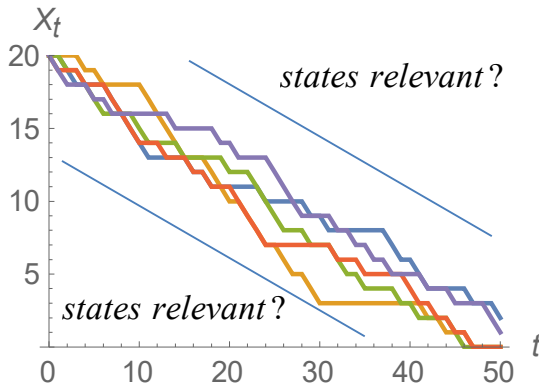


Schlosser, R. (2021). Scalable Relaxation Techniques to Solve Stochastic Dynamic Multi-Product Pricing Problems with Substitution Effects, *Journal of Revenue and Pricing Management* 20 (1), 54-65.

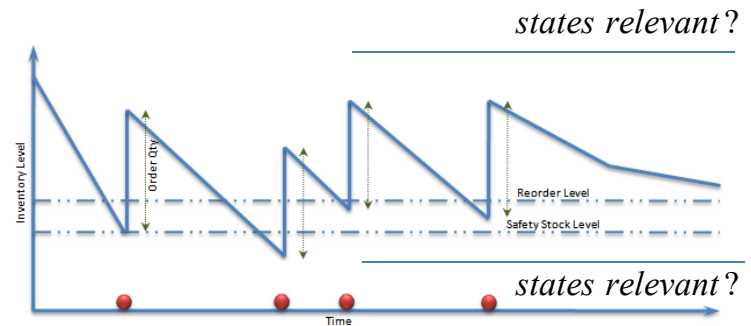
(2) Example of Simulation-Based Value Appr.

- In (1) we approximate the value function **for all** states
- But do we really need all states?

Airline Example



Inventory Example



Forward Dynamic Programming (Infinite Horizon)

- $$V^*(s) = \max_{a \in A} \left\{ \sum_{i \in I} P(i, a, s) \cdot (r(i, a, s) + \gamma \cdot V^*(\Gamma(i, a, s))) \right\} \quad (\text{Bellman equ.})$$
- **Forward Dynamic Programming:** Use a simulation-based approach:
 - (0) Start with $V(s) := 0, \forall s \in S$ and perform $k=0, \dots, K$ iterations with given s_0 :

Forward Dynamic Programming (Infinite Horizon)

- $V^*(s) = \max_{a \in A} \left\{ \sum_{i \in I} P(i, a, s) \cdot (r(i, a, s) + \gamma \cdot V^*(\Gamma(i, a, s))) \right\}$ (Bellman equ.)
- **Forward Dynamic Programming:** Use a simulation-based approach:
 - (0) Start with $V(s) := 0$, $\forall s \in S$ and perform $k=0, \dots, K$ iterations with given s_0 :
 - (1) Let $a_k(s_k) = \pi(s_k)$, i.e., apply the current action/policy based on current $V(s)$,
 where $\pi(s) := \arg \max_{a \in A} \left\{ \sum_{i \in I} P(i, a, s) \cdot (r(i, a, s) + \gamma \cdot V(\Gamma(i, a, s))) \right\}$, $\forall s \in S$

Forward Dynamic Programming (Infinite Horizon)

- $V^*(s) = \max_{a \in A} \left\{ \sum_{i \in I} P(i, a, s) \cdot (r(i, a, s) + \gamma \cdot V^*(\Gamma(i, a, s))) \right\}$ (Bellman equ.)
- **Forward Dynamic Programming:** Use a simulation-based approach:
 - (0) Start with $V(s) := 0, \forall s \in S$ and perform $k=0, \dots, K$ iterations with given s_0 :
 - (1) Let $a_k(s_k) = \pi(s_k)$, i.e., apply the current action/policy based on current $V(s)$,
 where $\pi(s) := \arg \max_{a \in A} \left\{ \sum_{i \in I} P(i, a, s) \cdot (r(i, a, s) + \gamma \cdot V(\Gamma(i, a, s))) \right\}, \forall s \in S$
 - (2) Improve $V(s_k) \leftarrow \sum_{i \in I} P(i, a_k(s_k), s_k) \cdot (r(i, a_k(s_k), s_k) + \gamma \cdot V(\Gamma(i, a_k(s_k), s_k)))$

Forward Dynamic Programming (Infinite Horizon)

- $V^*(s) = \max_{a \in A} \left\{ \sum_{i \in I} P(i, a, s) \cdot (r(i, a, s) + \gamma \cdot V^*(\Gamma(i, a, s))) \right\}$ (Bellman equ.)
- **Forward Dynamic Programming:** Use a simulation-based approach:
 - (0) Start with $V(s) := 0, \forall s \in S$ and perform $k=0, \dots, K$ iterations with given s_0 :
 - (1) Let $a_k(s_k) = \pi(s_k)$, i.e., apply the current action/policy based on current $V(s)$,
 where $\pi(s) := \arg \max_{a \in A} \left\{ \sum_{i \in I} P(i, a, s) \cdot (r(i, a, s) + \gamma \cdot V(\Gamma(i, a, s))) \right\}, \forall s \in S$
 - (2) Improve $V(s_k) \leftarrow \sum_{i \in I} P(i, a_k(s_k), s_k) \cdot (r(i, a_k(s_k), s_k) + \gamma \cdot V(\Gamma(i, a_k(s_k), s_k)))$
 - (3) Simulate next state $s_{k+1} \leftarrow \Gamma(i, a_k(s_k), s_k)$ according to $P(i, a_k(s_k), s_k), i \in I$

Discussion of Forward Dynamic Programming

- Visit/simulate relevant states (starting from an initial state s_0)
(not in a synchronous manner as in DP)
- Exploit full knowledge
 - Update the value function via expected rewards and state transitions
 - Simulate future states based on event/state transition probabilities P
- Apply a pure “greedy” policy based on current values $V(s)$
- Subsequently update V using the Bellman equation principle
- **Problem:** We may miss optimal paths (cf. loops)! **What can we do?**

Forward Dynamic Programming II (Infinite Horizon)

Forward Dynamic Programming (ε -greedy):

(0) Start with $V(s) := 0, \forall s \in S$ and perform $k=0, \dots, K$ iterations with given s_0 :

(1) Let $a_k(s_k) = \begin{cases} a \in A & \text{with prob. } \varepsilon \text{ play a random action} \\ \pi(s_k) & \text{with prob. } 1 - \varepsilon \end{cases}, \varepsilon \in (0, 1),$

i.e., apply a mixed *greedy/exploration* action/policy based on $V(s)$,

where $\pi(s) := \arg \max_{a \in A} \left\{ \sum_{i \in I} P(i, a, s) \cdot (r(i, a, s) + \gamma \cdot V(\Gamma(i, a, s))) \right\}, \forall s \in S$

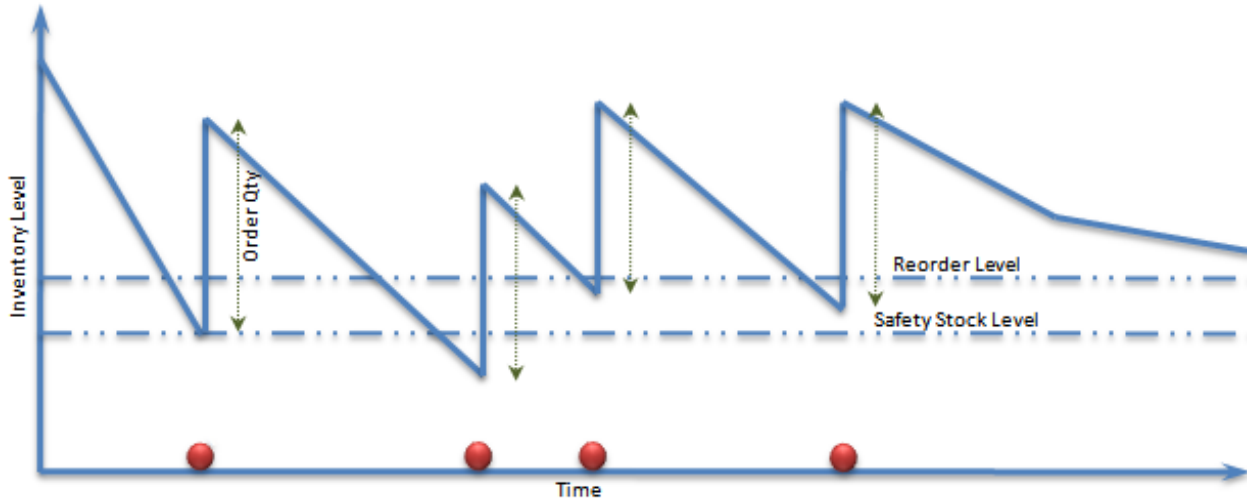
(2) Improve $V(s_k) \leftarrow \sum_{i \in I} P(i, a_k(s_k), s_k) \cdot (r(i, a_k(s_k), s_k) + \gamma \cdot V(\Gamma(i, a_k(s_k), s_k)))$ (check!)

(3) Simulate next state $s_{k+1} \leftarrow \Gamma(i, a_k(s_k), s_k)$ according to $P(i, a_k(s_k), s_k), i \in I$

Discussion Forward Dynamic Programming

- Visit/simulate relevant states (starting from an initial state s_0)
(not in a synchronous manner as in DP)
- Exploit full knowledge (vs. full knowledge is needed)
 - Improve the value function via expected rewards and state transitions
 - Simulate future states based on event/state transition probabilities P
- Apply a partly “*greedy*” policy based on current values $V(s)$
- Subsequently improve V using the Bellman equation principle
- The approach converges correctly (**Asymptotical optimal**)
- **Allows “good” heuristic solutions for larger problems in a feasible time**

ADP for the Inventory Management Example



Example Infinite Horizon MDP (Inventory Management)

- Framework: $t = 0, 1, 2, \dots, \infty$ Discrete time periods
- State: $s \in S$ Number of items left
- Actions: $a \in A$ Number of ordered items (replenish)
- Events: $i \in I, P(i, a, s)$ Demand i (e.g., 0,1,2,3 with prob. 1/4 each)
- Rewards: $r = r(i, a, s)$ Revenue – Order Cost – Holding Cost

$$:= p \cdot \min(i, s) - c \cdot a - h \cdot s - 1_{\{a>0\}} \cdot f$$
 e.g., for given price p , variable order cost c , holding h , and fixed order costs f
- New State: $s \rightarrow s' = \Gamma(i, a, s)$ Old – Sold + Replenish (end of period)
- Initial State: $s_0 \in S$ Initial items in $t=0$

ADP Results for the Inventory Management Example

- $\varepsilon = 0.05$ exploration probability
- $K \in \{0, \dots, 50\,000\}$ episodes/iterations (all < 10 sec)
- $\pi_{ADP}^{(K)} \approx \pi^*$ obtain good / near-optimal solutions based on V
- $V_{ADP}^{(K)} \approx V^*$ especially for (relevant) states with few inventory

<i>runs</i> K	500	1 000	2 000	10 000	50 000
$V_{ADP}^{(K)}(10) / V^*(10)$	0.63	0.77	0.93	0.97	1.00

- **At home:** play with K and ε as well as other parameters and study the quality of the ADP solution against the optimal one

ADP for Finite Horizon MDP Problems



Forward Dynamic Programming (Finite Horizon)

Forward Dynamic Programming (ε -greedy):

(0) Start with $V_t(s) := 0, \forall s \in S$. Use $k=0, \dots, K$ iterations **over $t=0, \dots, T-1$** from S_0 :

(1) Let $a_t^{(k)}(s_t^{(k)}) = \begin{cases} a \in A & \text{with prob. } \varepsilon_k \text{ play a random action} \\ \pi_t(s_t^{(k)}) & \text{with prob. } 1 - \varepsilon_k \end{cases}, \varepsilon_k \in (0, 1)$ ←

i.e., apply a mixed **exploration-exploitation** policy based on $V_t(s)$,

where $\pi_t(s) := \arg \max_{a \in A} \left\{ \sum_{i \in I} P_t(i, a, s) \cdot (r_t(i, a, s) + \gamma \cdot V_{t+1}(\Gamma_t(i, a, s))) \right\}, \forall s \in S$

(2) Improve $V_t(s_t^{(k)}) \leftarrow \sum_{i \in I} P_t(i, a_t^{(k)}(s_t^{(k)}), s_t^{(k)}) \cdot (r_t(i, a_t^{(k)}(s_t^{(k)}), s_t^{(k)}) + \gamma \cdot V_{t+1}(\Gamma_t(i, a_t^{(k)}(s_t^{(k)}), s_t^{(k)})))$

(3) Simulate state $s_{t+1}^{(k)} \leftarrow \Gamma_t(i, a_t^{(k)}(s_t^{(k)}), s_t^{(k)})$ according to $P_t(i, a_t^{(k)}(s_t^{(k)}), s_t^{(k)})$, $i \in I$

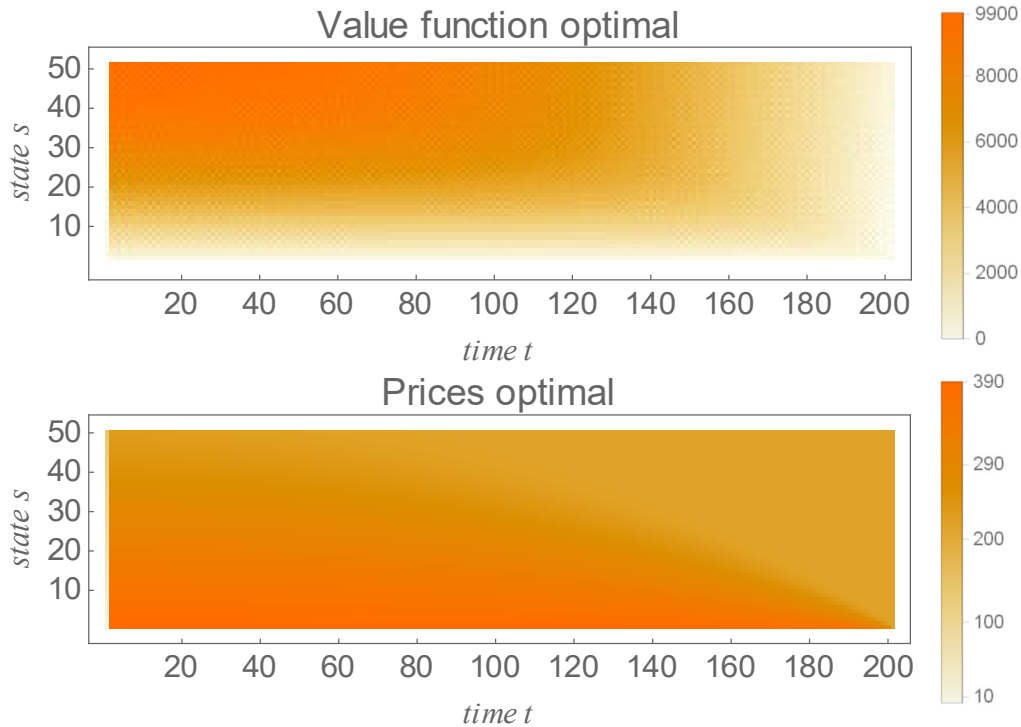
Example MDP (Selling Airline Tickets)

- Framework: $t = 0, 1, 2, \dots, T$ Time periods
- State: $s_t \in S := \{0, 1, \dots, N\}$ Items left
- Actions: $a_t \in A := \{5, 10, \dots, 400\}$ Price
- Events: $i_t \in I := \{0, 1\}$ with probabilities Demand
 $P_t(1, a, s) := (1 - a / 400) \cdot (1 + t) / T$ $P_t(0, a, s) = 1 - P_t(1, a, s)$
- Rewards: $r_t = r(i, a, s) := a \cdot \min(i, s)$ Revenue
- New State: $s_t \rightarrow s_{t+1} = \Gamma(i_t, a_t, s_t) := \max(0, s_t - i_t)$ Old – sold
- Initial State: $s_0 \in S, s_0 := N$ Initial items N
- Final Reward: $r_T(s) := f \cdot s$ with $f = 10$ Weight for freight

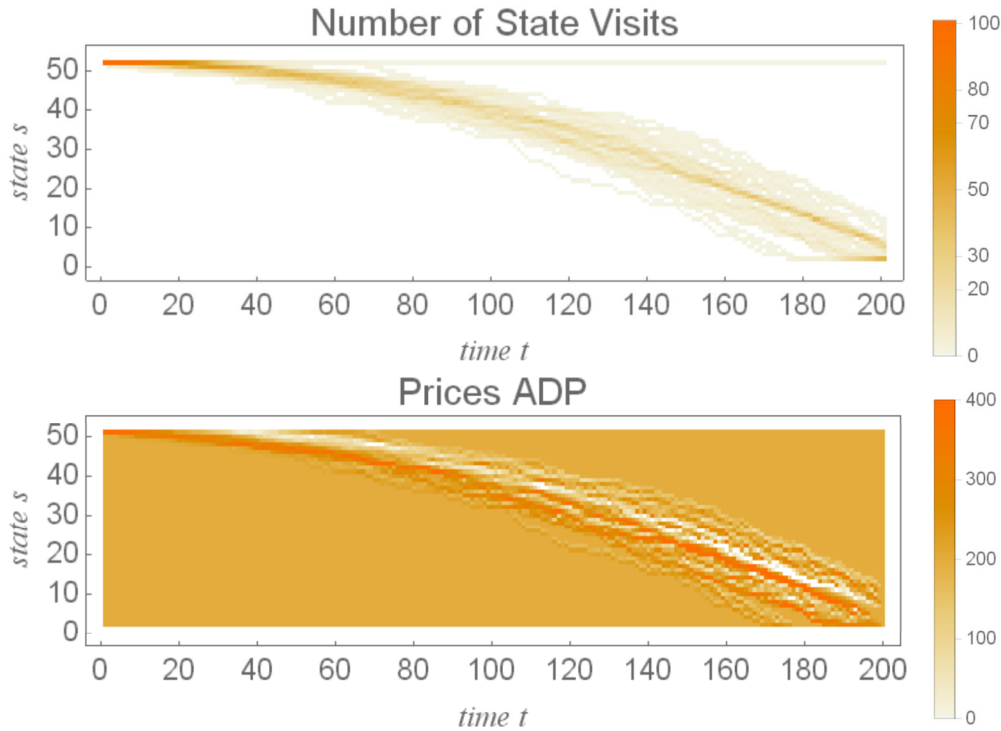
ADP for the Airline Example

- $\varepsilon_k = 0.1 + 0.4 \cdot (1 - k / K)$ exploration probability (for run k)
- $K \in \{0, \dots, 10\,000\}$ different numbers of episodes (T iterations)
- $\pi_{ADP}^{(K)} \approx \pi^*$ obtain good / (near-)optimal solutions (for s_0)
- $V_{ADP}^{(K)} \approx V^*$ especially for (relevant/achievable) states
- **At home:** play with K and ε_k as well as other model parameters to study the quality of the ADP solution against the optimal one

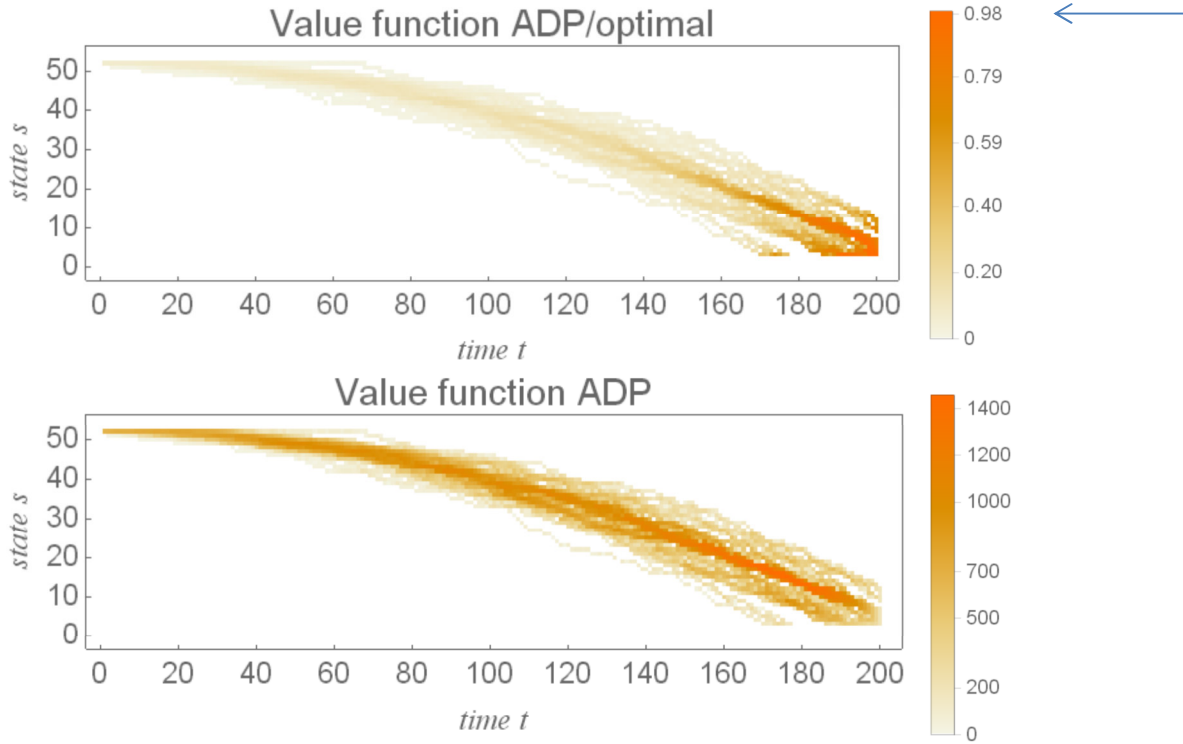
ADP Results for the Airline Example (*optimal*)



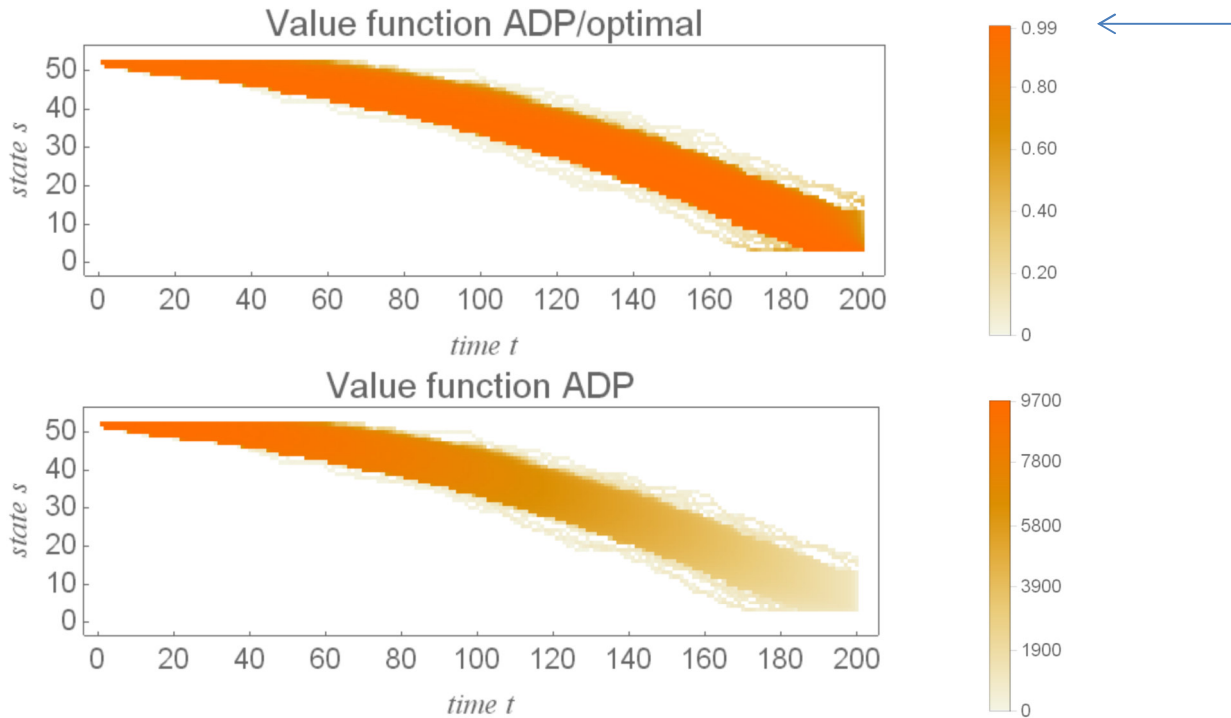
ADP Results for the Airline Example ($K=100$, 3 sec)



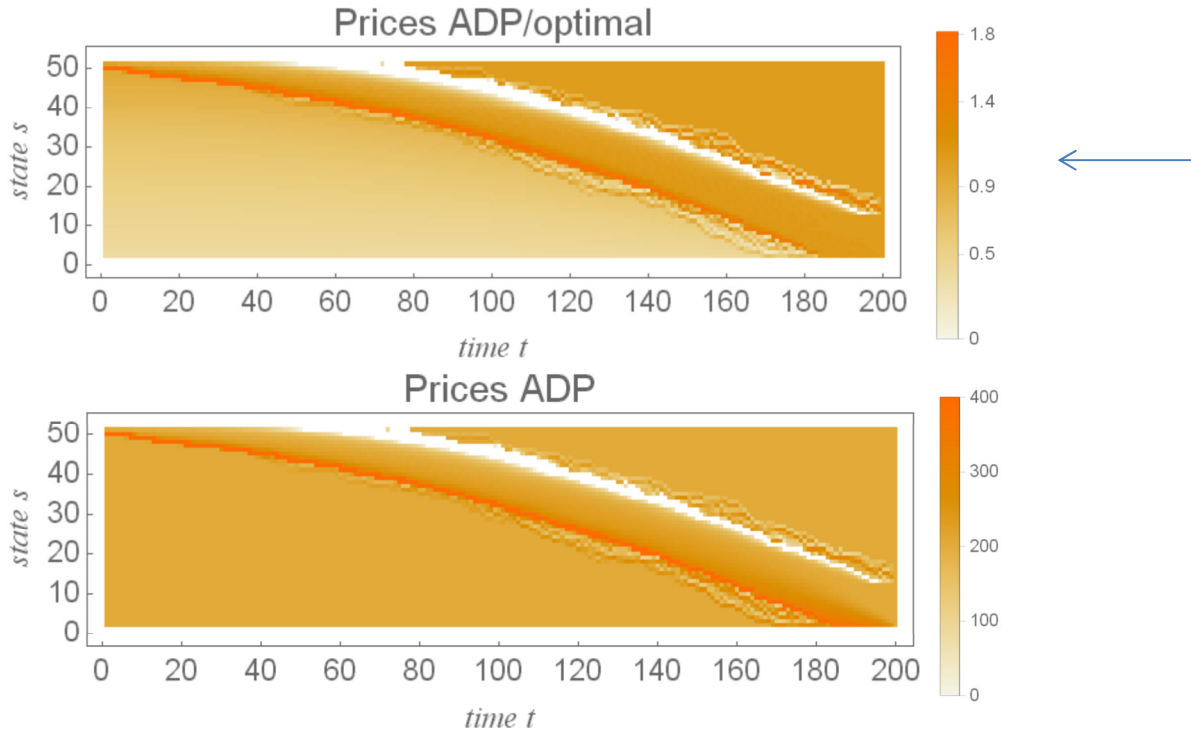
ADP Results for the Airline Example ($K=100, 3 \text{ sec}$)



ADP Results for the Airline Example ($K=10\,000$, 200 sec)



ADP Results for the Airline Example ($K=10\,000$, 200 sec)



Summary (Solving Discrete Time MDPs via ADP)

ADP (Forward Dynamic Programming)

- (+) provides near-optimal solutions for (in/)finite horizon MDPs
- (+) guaranteed convergence
- (+) numerically simple
- (+) general applicable
- (+) quickly obtain good heuristics

- (-) updates only for single “visited” states (cf. large state spaces)
- (-) results are stochastic (due to simulated next states)
- (-) hyperparameter tuning (e.g., exploration rate)
- (-) **full information required (cf. events & transitions)**

Next: QL, i.e., similar solution approaches requiring less information

Could You Solve Different Test Problems via ADP?

- Any Questions?
- Finite Horizon (*use ADP*)
 - Eating cake (deterministic utility)
 - Selling airline tickets (stochastic demand)
- Infinite Horizon (*use ADP*)
 - Car replacement problem (deterministic costs)
 - Inventory management (stochastic demand)

Overview

Week	Dates	Topic
1	April 21	Introduction
2	April 25/28	Finite + Infinite Time MDPs
3	May 2/5	Approximate Dynamic Programming (ADP) + DP Exercise
4	May 12	Q-Learning (QL) (not Mon May 9)
5	May 16/19	Deep Q-Networks (DQN)
6	May 23	DQN Extensions (Thu May 26 “Himmelfahrt”)
7	May 30/June 2	Policy Gradient Algorithms
8	June 9	Project Assignments (Mon June 6 “Pfingstmontag”)
9	June 13/16	Work on Projects: Input/Support
10	June 20/23	Work on Projects: Input/Support
11	June 27/30	Work on Projects: Input/Support
12	July 4/7	Work on Projects: Input/Support
13	July 11/14	Work on Projects: Input/Support
14	July 18/21	Final Presentations
	Sep 15	Finish Documentation