



# 自适应动态规划最优控制方法

刘德荣

2016年11月6日

May 14–19, 2017

Anchorage, Alaska, USA

## Call for Papers

11月15日截稿

The annual International Joint Conference on Neural Networks (IJCNN) is the flagship conference of the IEEE Computational Intelligence Society and the International Neural Network Society. It covers a wide range of topics in the field of neural networks, from biological neural network modeling to artificial neural computation.

## Guidelines for Paper Submission

<http://www.ijcnn.org/paper-submission>

## Important Dates

- Special session & competition proposals submission: September 15, 2016
- Tutorial and workshop proposal submission: October 15, 2016 October 30, 2016
- Paper submission: November 15, 2016
- Paper decision notification: January 20, 2017
- Camera-ready submission: February 20, 2017

## Topics Covered

NEURAL NETWORK MODELS

### Follow Us!



IJCNN 2017



IJCNN 2017

### Sponsoring Organizations



INNS - International  
Neural Network Society

# 14<sup>th</sup> International Symposium on Neural Networks

June 21-23, 2017

Sapporo, Hokkaido, Japan

1月1日截稿



## Honorary Chair

Shun'ichi Amari, RIKEN Brain Science Institute, Tokyo, Japan

## General Chairs

Hidenori Kawamura, Hokkaido University, Sapporo, Japan

Jun Wang, City University of Hong Kong, Hong Kong

## Advisory Chairs

Kunihiro Fukushima and Takeshi Yamakawa, Fuzzy Logic Systems Institute, Fukuoka, Japan

## Steering Chairs

Haibo He, University of Rhode Island, Kingston, USA

Derong Liu, University of Illinois – Chicago, Chicago, USA

Jun Wang, City University of Hong Kong, Hong Kong

## Organizing Committee Chairs

Andrzej Cichocki, RIKEN Brain Science Institute, Tokyo, Japan

Min Han, Dalian University of Technology, Dalian, China

Bao-Liang Lu, Shanghai Jiao Tong University, Shanghai, China

Masahito Yamamoto, Hokkaido University, Sapporo, Japan

## Program Chairs

Fengyu Cong, Dalian University of Technology, Dalian, China

Andrew C.-S. Leung, City University of Hong Kong, Hong Kong

Qinglai Wei, CAS Institute of Automation, Beijing, China

## Special Sessions Chairs

Long Cheng, CAS Institute of Automation, Beijing, China

Satoshi Kurihara, Univ. of Electro-Communications, Tokyo, Japan

## Call for Papers

<https://conference.cs.cityu.edu.hk/isnn/>

**Sponsors/organizers:** [Hokkaido University](#) and [City University of Hong Kong](#)

**Technical co-sponsors:** Asia Pacific Neural Network Society, IEEE Computational Intelligence Society, International Neural Network Society, and Japanese Neural Network Society

Following the successes of previous events, the 14th International Symposium on Neural Networks (ISNN 2017) will be held in Sapporo, Hokkaido, Japan. Located in northern island of Hokkaido, Sapporo is the fourth largest Japanese city and a popular summer/winter tourist venue. ISNN 2017 aims to provide a high-level international forum for scientists, engineers, and educators to present the state of the art of neural network research and applications in related fields. The symposium will feature plenary speeches given by world renowned scholars, regular sessions with broad coverage, and special sessions focusing on popular topics.

## Call for Papers and Special Sessions

Prospective authors are invited to contribute high-quality papers to ISNN 2017. In addition, proposals for special sessions within the technical scopes of the symposium are solicited. Special sessions, to be organized by internationally recognized experts, aim to bring together researchers in special focused topics. Papers submitted for special sessions are to be peer-reviewed with the same criteria used for the contributed papers. Researchers interested in organizing special sessions are invited to submit formal proposals to ISNN 2017. A special session proposal should include the session title, a brief description of the scope and motivation, names, contact information and brief biographical information of the organizers.

## Topic Areas

Topics areas include, but not limited to, computational neuroscience, connectionist theory and cognitive science, mathematical modeling of neural systems, neurodynamic analysis, neurodynamic optimization and

# The 24th International Conference on Neural Information Processing (ICONIP 2017)

November 14–18, 2017, Guangzhou, China



## General Chair

Derong Liu, China

5月份截稿

大会总主席：  
刘德荣

## CALL FOR PAPERS

The 24th International Conference on Neural Information Processing will be held in Guangzhou, China, November 14–18, 2017. It is an annual event, organized since 1994 by the Asia Pacific Neural Network Society (APNNS, previously APNNA).

Guangzhou is the capital and largest city of Guangdong province, People's Republic of China. Located on the Pearl River, about 120 km (75 mi) north-northwest of Hong Kong and north-northeast of Macau, Guangzhou is a key national transportation hub and trading port. It is one of the five National Central Cities and holds sub-provincial administrative status. With a tropical climate, Guangzhou is warm all year long, and November is the best month over the year to visit Guangzhou with comfortable weather. There are so many snacks, delicious dim sum, great sea food and Cantonese cuisine. There is a great deal of tourist sites in Guangzhou, such as Beijing Road Walking Street, Haizhu Square, Sun Yat-sen's Memorial Hall, Yuexiu Park, Zhongxin Square, White Cloud Mountain, Canton Tower and so on.

ICONIP 2017 aims to provide a high-level international forum for scientists, researchers, educators, industrial professionals, and students worldwide to present state-of-the-art research results, address new challenges, and discuss trends in neural information processing and applications. ICONIP 2017 invites scholars in all areas of neural network theory and applications, computational neuroscience, machine learning and others. In addition to regular technical sessions with oral and poster presentations, the conference program will include special sessions and tutorials on topics of current interest.

ICONIP 2017 features plenary/keynote and panel discussion sessions by world leading researchers as

# Dynamic programming + Learning control



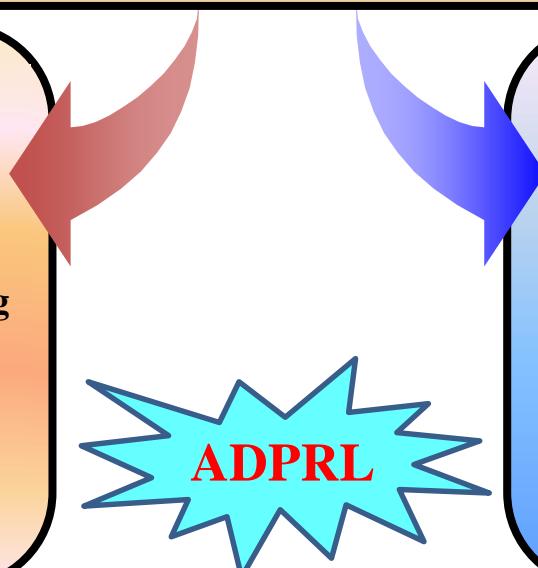
Curse-of-dimensionality - 1957

ADP - 1977

Adaptive dynamic programming  
Approximate dynamic programming  
Asymptotic dynamic programming  
Neural dynamic programming  
Adaptive critic designs

RL - 1984

Reinforcement learning  
Richard Sutton (1984) *Temporal Credit Assignment in Reinforcement Learning* (PhD thesis). University of Massachusetts.



ADP for self-learning control



# ADP for Self-Learning Control

Dynamic Programming

Adaptive Dynamic Programming

Iterative ADP

New Developments

Concluding Remarks



# Dynamic programming

## ● Consider

$$x_k = F(x_k, u_k)$$

- $x_k \in X \subset R^n$  is the state
- $u_k \in A \subset R^p$  is the control vector.

## ● Performance index (or cost to go)

$$J(x_k) = \sum_{i=k}^{\infty} \gamma^{i-k} U(x_i, u_i)$$

- $U$  is the given utility function
- $\gamma$  is the discount factor with  $0 < \gamma \leq 1$
- Infinite horizon problems.



# Examples of utility functions

Production output

A

Energy usage

B

Emissions

C

Tracking control errors

D

There are  
many similar  
examples in  
practice



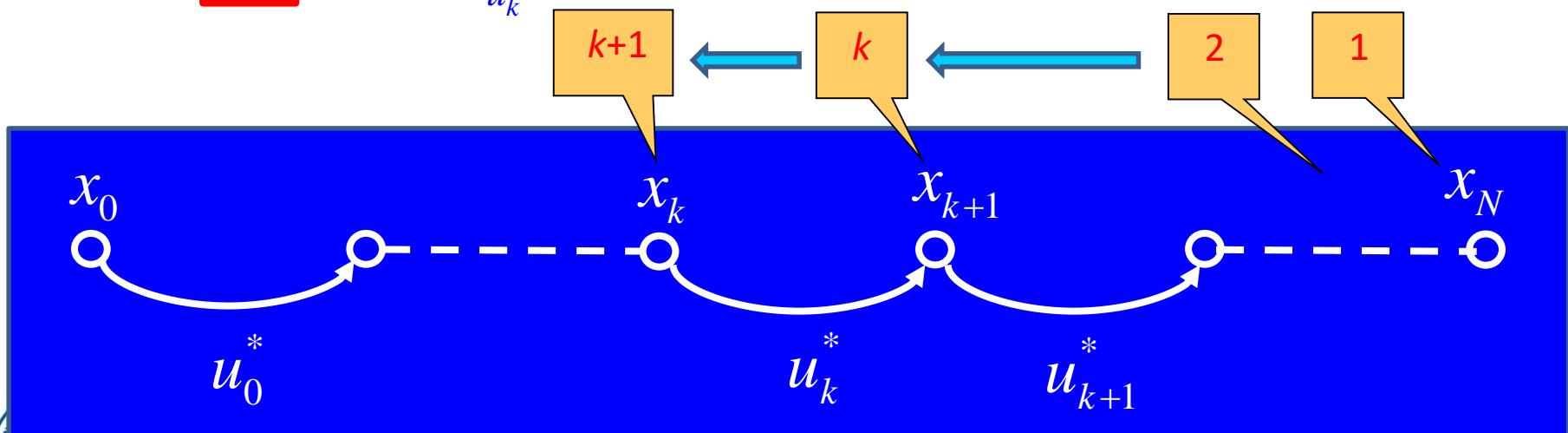
# Bellman Equation

- According to Bellman's principle of optimality,

$$J^*(x_k) = \min_{u_k} \{U(x_k, u_k) + \gamma J^*(x_{k+1})\}.$$

– The optimal control  $u_k^*$  at time  $k$  is the  $u_k$  that achieves this minimum, i.e.,

$$u_k^* = \arg \min_{u_k} \{U(x_k, u_k) + \gamma J^*(x_{k+1})\}.$$



# Issues with dynamic programming

- Dynamic programming is applicable to problems that minimize/maximize a cost

- Applicable to many engineering problems
- Model-based

- Backward numerical process

- Backward in time
- Unknown function  $J$

Curse of dimensionality

$$J(x_k) = \sum_{i=k}^{\infty} \gamma^{i-k} U(x_i, u_i)$$

- Computational complexity increases exponentially as the number of variables

- Only suitable for small problems in practice



# ADP for Self-Learning Control

**Dynamic Programming**

**Adaptive Dynamic Programming**

**Iterative ADP**

**New Developments**

**Concluding Remarks**



# Adaptive dynamic programming

Adaptive Dynamic Programming

ADP

Approximate Dynamic Programming

Asymptotic Dynamic Programming

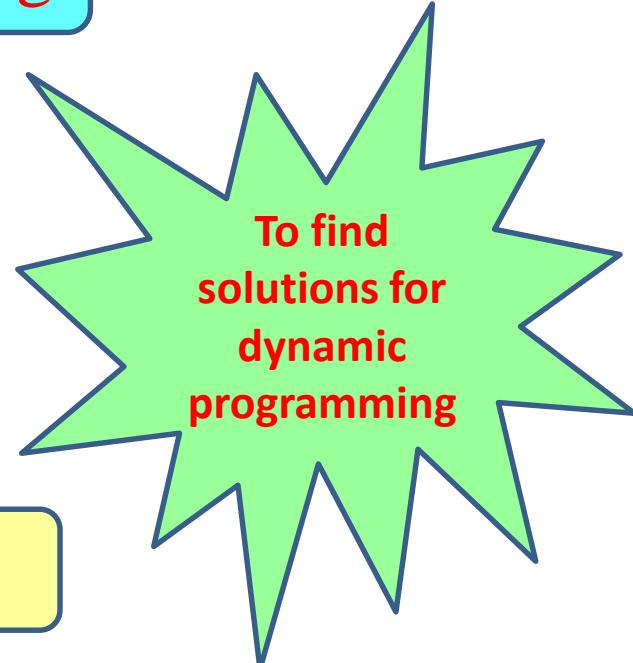
Adaptive Critic Designs

Neurodynamic Programming

Relaxed Dynamic Programming

Reinforcement Learning

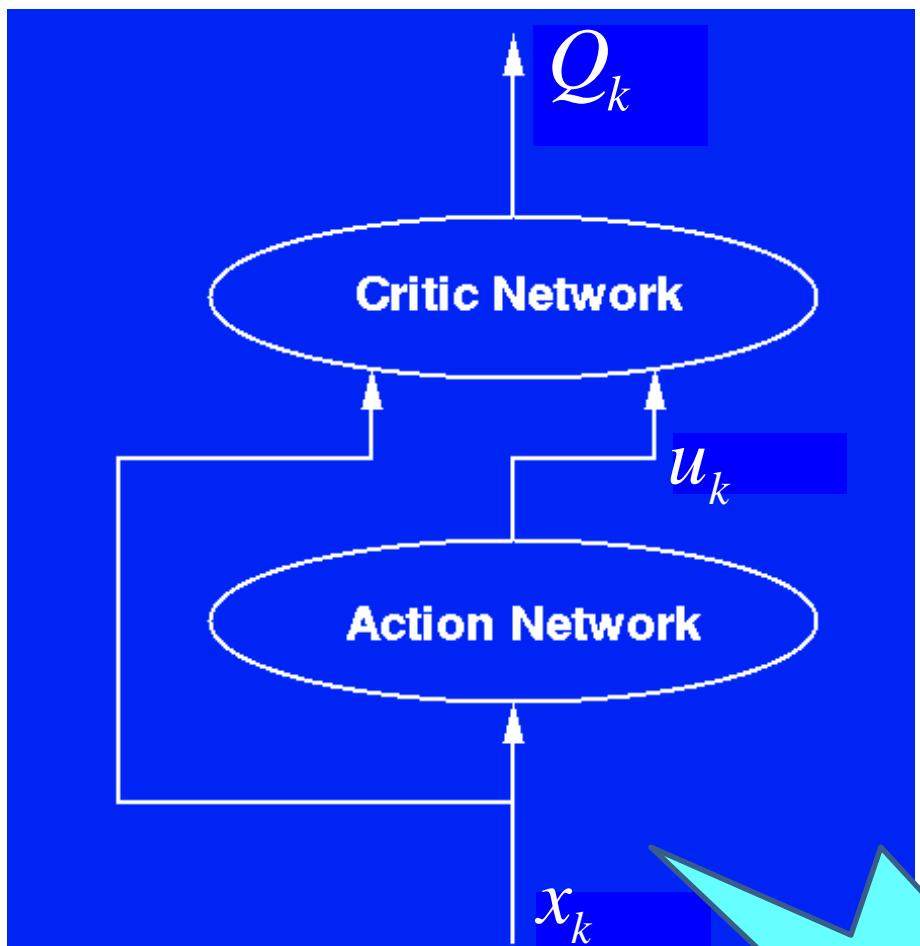
RL



To find  
solutions for  
dynamic  
programming



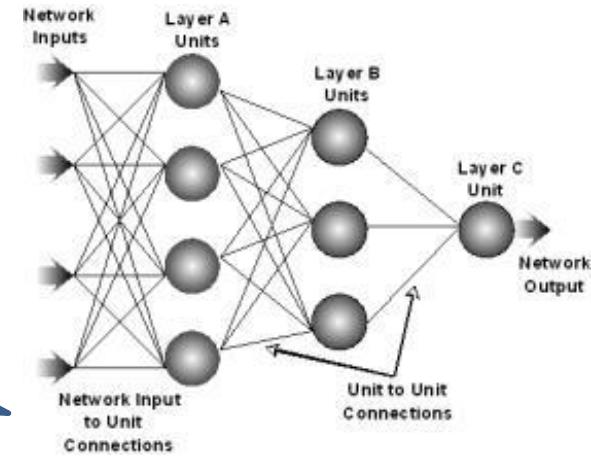
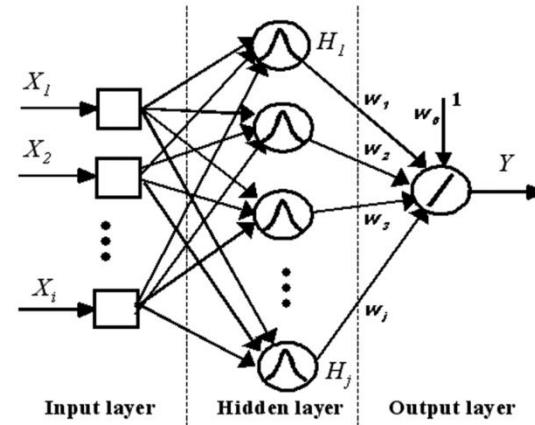
# Adaptive dynamic programming



Model-free HDP

Generate  $u_k$   
to minimize  
 $Q_k$

$$Q_k \approx J_{k+1}$$

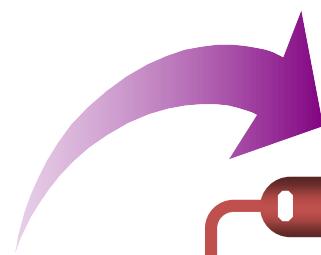


# Work done before 1997



Widrow et al.

**“Punish/Reward:  
Learning with a  
Critic in Adaptive  
Threshold  
Systems,” IEEE  
Transactions on  
Systems, Man,  
and Cybernetics  
(1973)**



Paul Werbos

**“Advanced  
Forecasting  
Methods for  
Global Crisis  
Warning and  
Models of  
Intelligence,”  
General Systems  
Yearbook (1977)**



Prokhorov/Wunsch

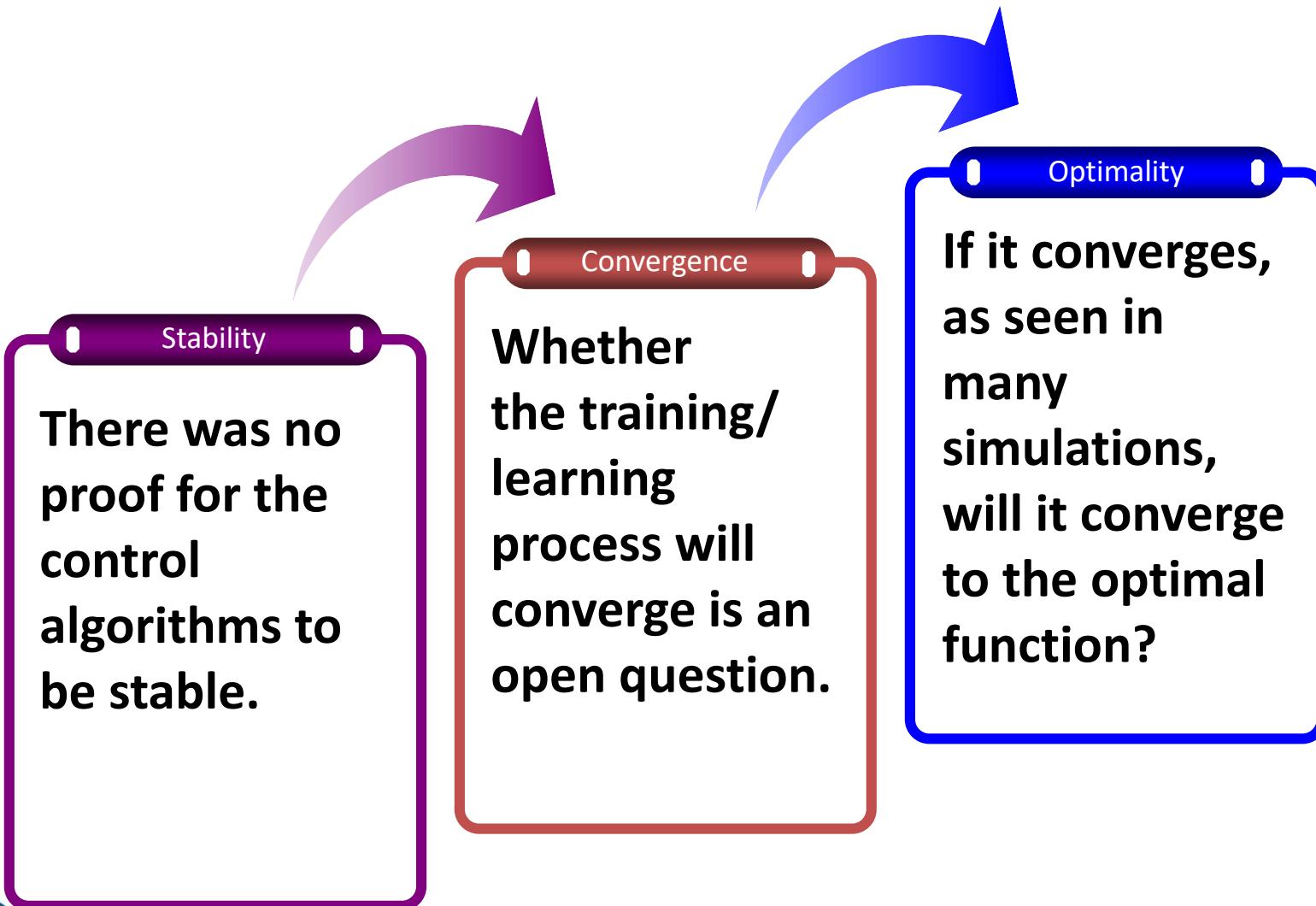
**“Adaptive Critic  
Designs,” IEEE  
Transactions on  
Neural Networks  
(1997)**



**Summarized  
all works  
before 1997.**



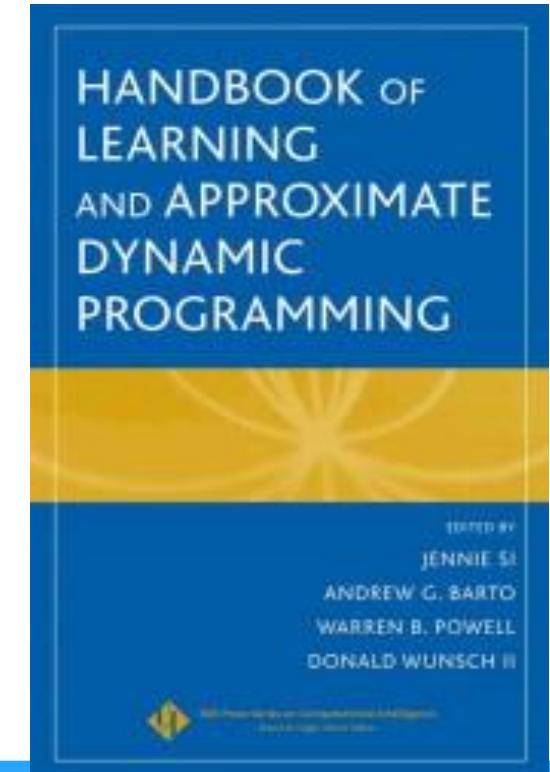
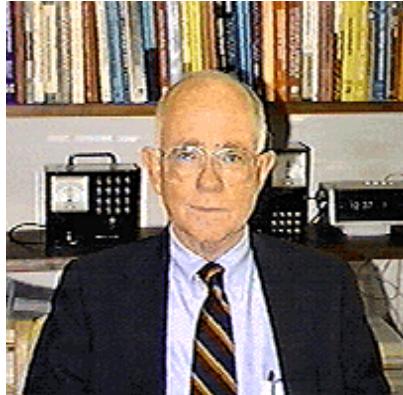
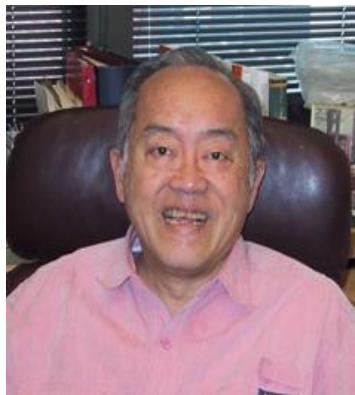
# Theoretical issues to be resolved



# Workshops

## ❖ Workshop on Learning and Approximate Dynamic Programming

- Sponsored by the National Science Foundation
- **April 8-10, 2002**, in Key West, Florida
- 29 researchers were invited
- Including **Larry Ho**, **Donald Wunsch II**, **Warren Powell**, **Andrew Barto**
- Published a book: Handbook of Learning and Approximate Dynamic Programming
- <http://www.fulton.asu.edu/~nsfadp/>

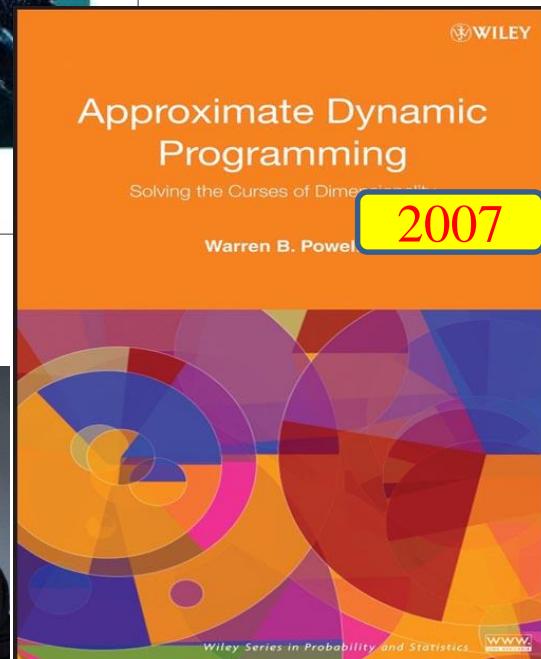
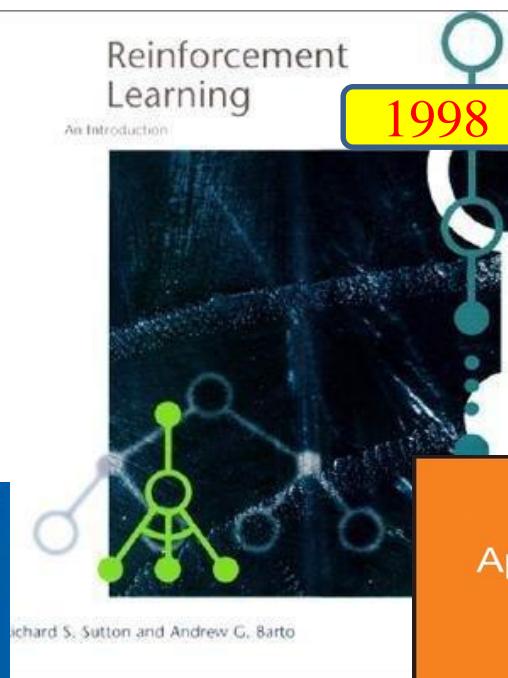
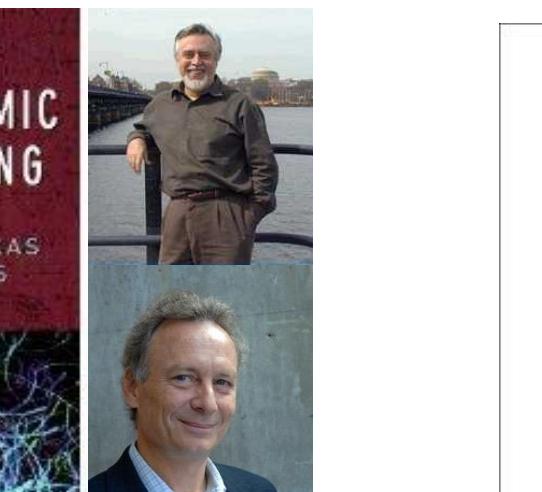
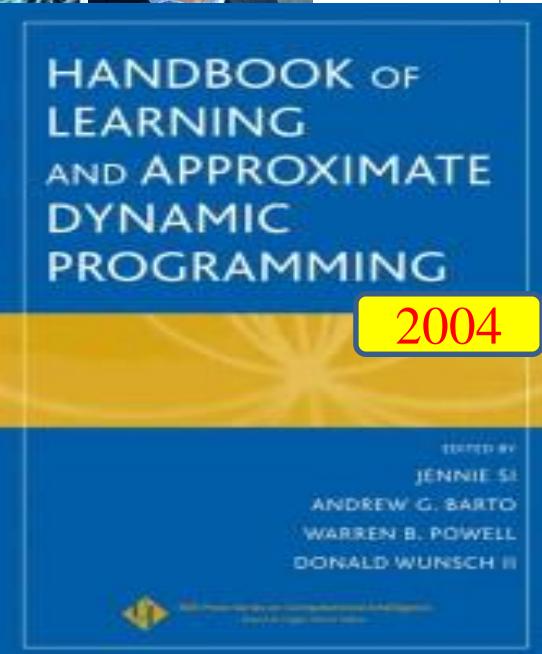
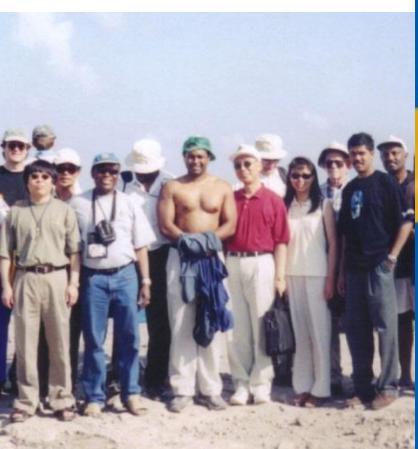
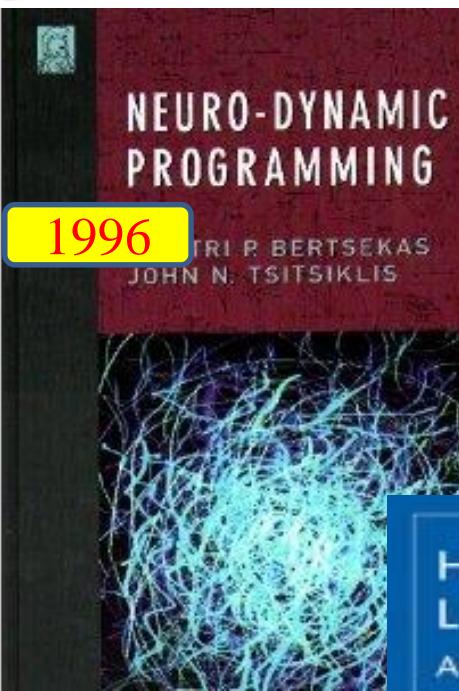


# Workshops

- ❖ Workshop on **Approximate** Dynamic Programming
  - Sponsored by the NSF
  - **April 3-6, 2006**, in Mexico
  - 42 researchers were invited
  - Including **Dimitri Bertsekas**, **Frank Lewis**, and **Paul Werbos**
  - Outreach to Mexican students and researchers
  - <http://www.fulton.asu.edu/~nsfadp/>



# Books by 2007



# ADP for Self-Learning Control

Dynamic Programming

Adaptive Dynamic Programming

Iterative ADP

New Developments

Concluding Remarks



# Iterative adaptive dynamic programming

- Bellman's principle of optimality

$$J^*(x_k) = \min_{u_k} \{U(x_k, u_k) + \gamma J^*(x_{k+1})\}.$$

$$\gamma = 1$$

We will use a function  $V(x_k)$  to approximate  $J^*(x_k)$ .

- We have

$$V(x_k) = \min_{u_k} \{U(x_k, u_k) + V(x_{k+1})\}.$$



# Iterative adaptive dynamic programming

- We need to solve for  $V(x_k)$  from

$$V(x_k) = \min_{u_k} \{U(x_k, u_k) + V(x_{k+1})\}.$$

- Considering an **algebraic equation**

$$x = f(x),$$

we can solve this equation by iterative method

$$x_{i+1} = f(x_i) \longrightarrow x_\infty = f(x_\infty)$$

starting from any initial value  $x_0$ , if the above iterative process is convergent.



# Iterative adaptive dynamic programming

- We can do the same for  $V(x_k)$  from

$$V(x_k) = \min_{u_k} \{U(x_k, u_k) + V(x_{k+1})\}$$

to use

$$V_{i+1}(x_k) = \min_{u_k} \{U(x_k, u_k) + V_i(x_{k+1})\}$$

starting from any initial function  $V_0(\cdot)$ , if the above iterative process is convergent. We can get

$$V_\infty(x_k) = \min_{u_k} \{U(x_k, u_k) + V_\infty(x_{k+1})\}$$



# Summary

- The solution of

$$J^*(x_k) = \min_{u_k} \{U(x_k, u_k) + J^*(x_{k+1})\},$$

can be obtained from

$$V_{i+1}(x_k) = \min_{u_k} \{U(x_k, u_k) + V_i(x_{k+1})\},$$

starting from any initial function  $V_0(\cdot)$ .

- If this iterative process is convergent. We can get

$$V_\infty(x_k) = \min_{u_k} \{U(x_k, u_k) + V_\infty(x_{k+1})\}.$$



# Relaxed dynamic programming

- Theorem due to Rantzer and his co-workers:

$$\left[ 1 + \frac{\alpha - 1}{(1 + \rho^{-1})^i} \right] J^* \leq V_i \leq \left[ 1 + \frac{\beta - 1}{(1 + \rho^{-1})^i} \right] J^*$$

where

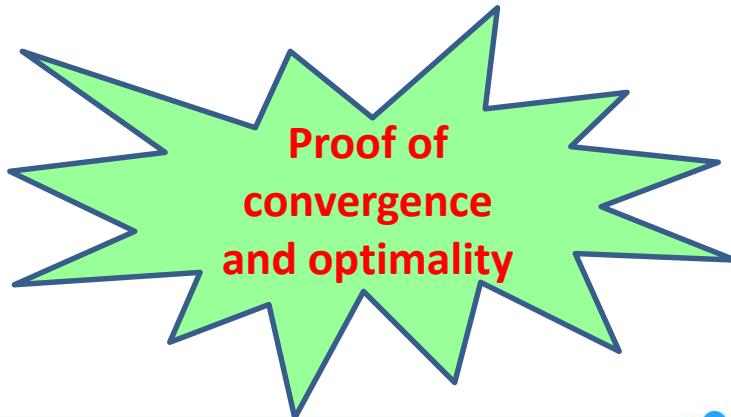
$$0 \leq \alpha J^* \leq V_0 \leq \beta J^*$$

$$J^*(F(x, u)) \leq \rho U(x, u)$$



Rantzer et al.

“Relaxing Dynamic Programming,”  
IEEE Transactions on Automatic  
Control (2006)  
and  
IEE Proceedings - Control Theory  
and Applications (2006)



# Relaxed dynamic programming

- Theorem due to Rantzer and his co-workers:

$$\left[ 1 + \frac{\alpha - 1}{(1 + \rho^{-1})^i} \right] J^* \leq V_i \leq \left[ 1 + \frac{\beta - 1}{(1 + \rho^{-1})^i} \right] J^*$$

where

$$0 \leq \alpha J^* \leq V_0 \leq \beta J^*$$

How to choose  $V_0$

Usually  $0 < \alpha < 1$  and  $\beta > 1$

$$J^*(F(x, u)) \leq \rho U(x, u)$$

Restrictions on  $J^*$  and  $U$

$U(x, u) \neq 0$ , and  $U(x, u) = 0$  only when  $J^*(F(x, u)) = 0$

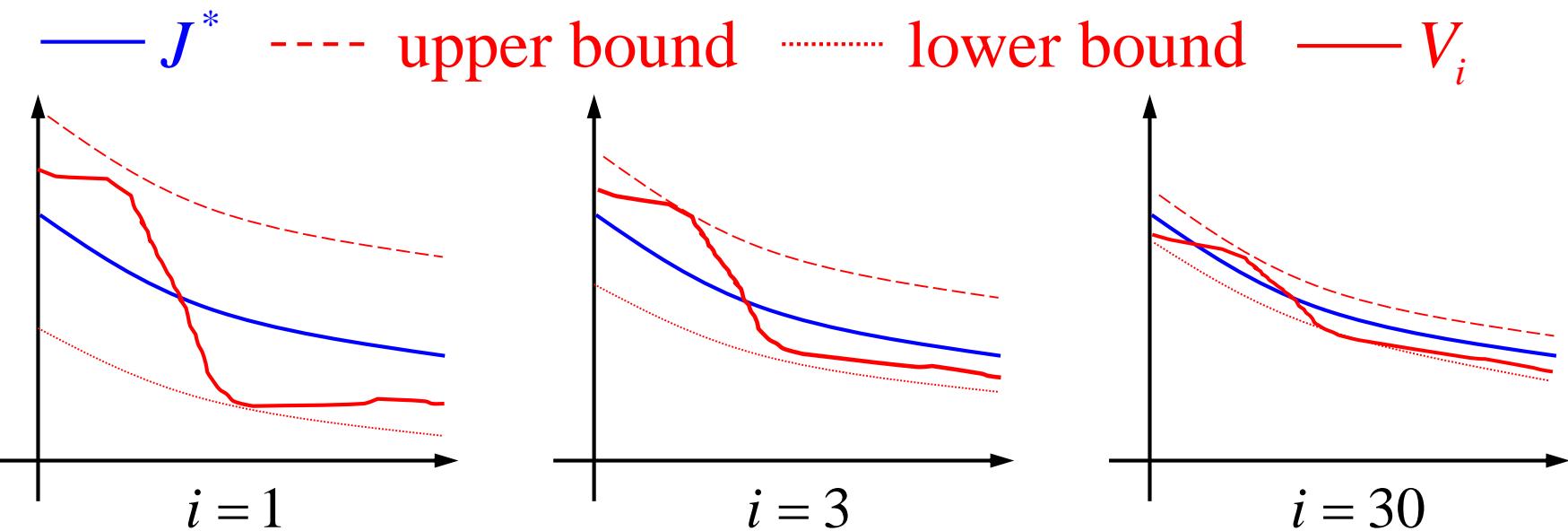
Large  $\rho \rightarrow$  slow convergence



# Relaxed dynamic programming

- Theorem due to Rantzer and his co-workers:

$$\left[ 1 + \frac{\alpha - 1}{(1 + \rho^{-1})^i} \right] J^* \leq V_i \leq \left[ 1 + \frac{\beta - 1}{(1 + \rho^{-1})^i} \right] J^* \quad \begin{array}{l} \alpha < 1 \\ \beta > 1 \end{array}$$



# Relaxed dynamic programming

- In the theorem due to Rantzer and his co-workers:

$$\left[ 1 + \frac{\alpha - 1}{(1 + \rho^{-1})^i} \right] J^* \leq V_i \leq \left[ 1 + \frac{\beta - 1}{(1 + \rho^{-1})^i} \right] J^*$$

where

$$0 \leq \alpha J^* \leq V_0 \leq \beta J^*$$

Rantzer suggested to choose  $V_0 \equiv 0$

$$\Rightarrow \alpha = 0$$



# Iterative adaptive dynamic programming

- Theorem due to Frank Lewis and his co-workers:
  - The limit  $V_\infty(x_k)$  of the value function sequence  $\{V_i\}$  exists;
  - $\lim_{i \rightarrow \infty} V_i(x_k) = V_\infty(x_k) = J^*(x_k);$
  - $\lim_{i \rightarrow \infty} v_i(x_k) = u^*(x_k).$



Al-Tamimi, Frank Lewis, Abu-Khalaf,  
“Discrete-Time Nonlinear HJB Solution  
Using Approximate Dynamic  
Programming: Convergence Proof,” IEEE  
Transactions on Systems, Man, and  
Cybernetics – Part B (2008)



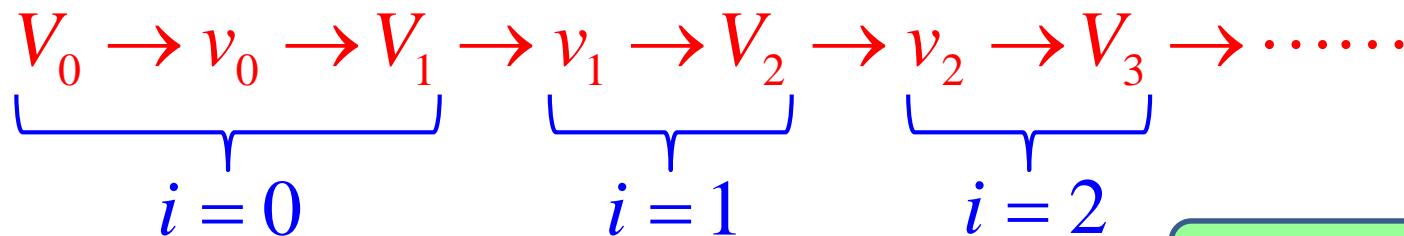
# Iterative adaptive dynamic programming

- Choose the **initial** value function as  $V_0(\cdot) \equiv 0$ .
- For  $i = 0, 1, 2, \dots$ , solve the control law as

$$v_i(x_k) = \arg \min_{u_k} \{U(x_k, u_k) + V_i(x_{k+1})\},$$

and update the value function as

$$V_{i+1}(x_k) = U(x_k, v_i(x_k)) + V_i(F(x_k, v_i(x_k))).$$

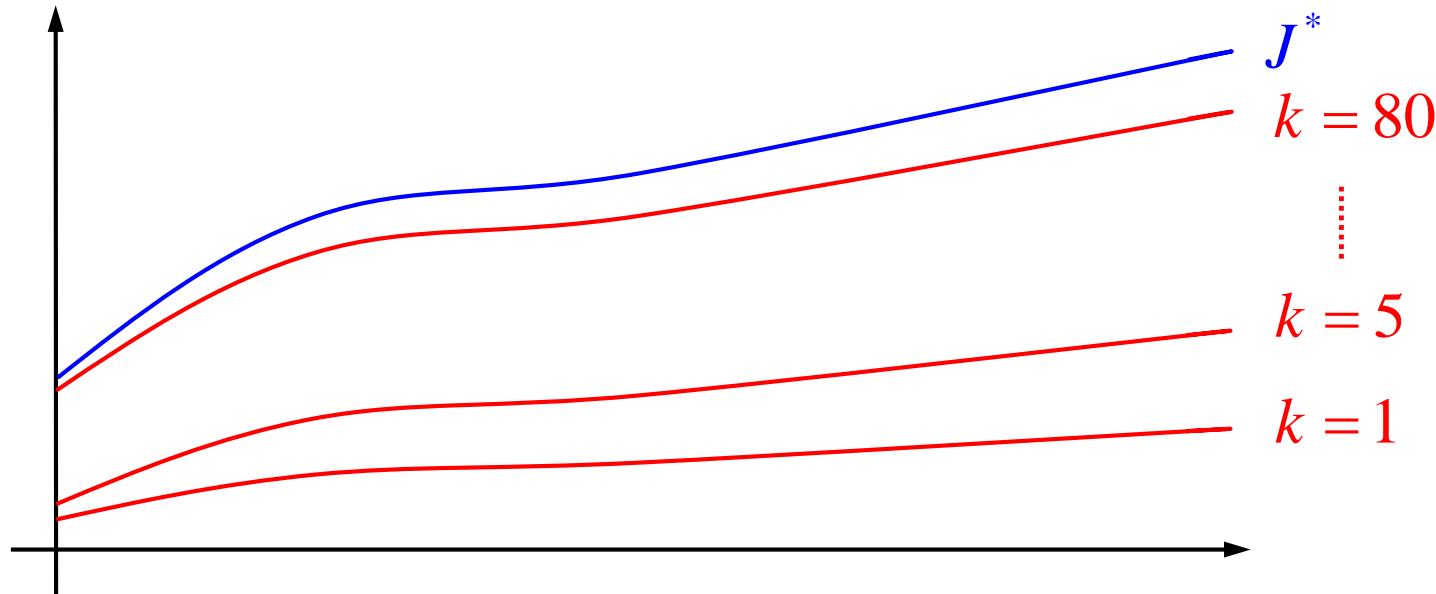


**Value iteration**



# Iterative adaptive dynamic programming

- Theorem due to Frank Lewis and his co-workers:
  - The sequence  $\{V_i\}$  is monotonically non-decreasing.

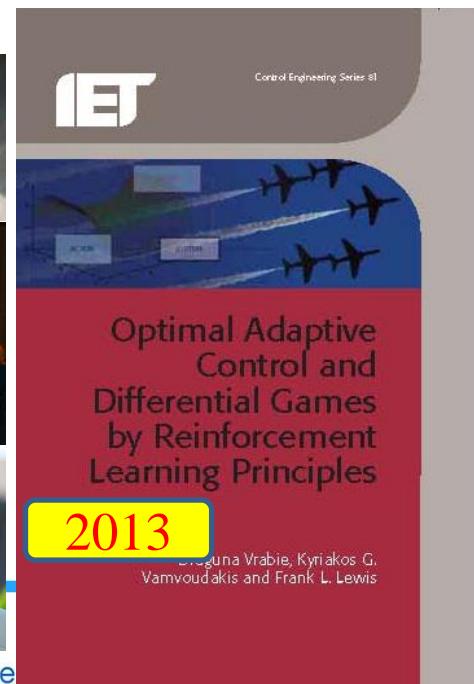
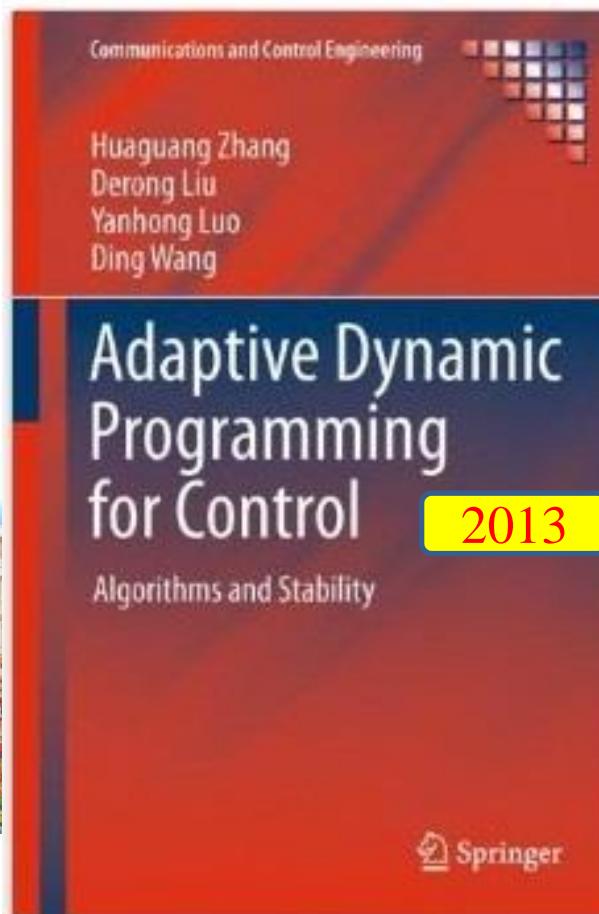
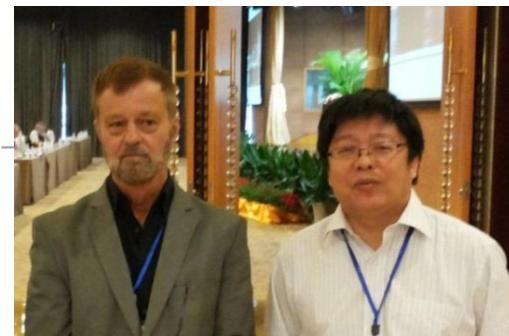
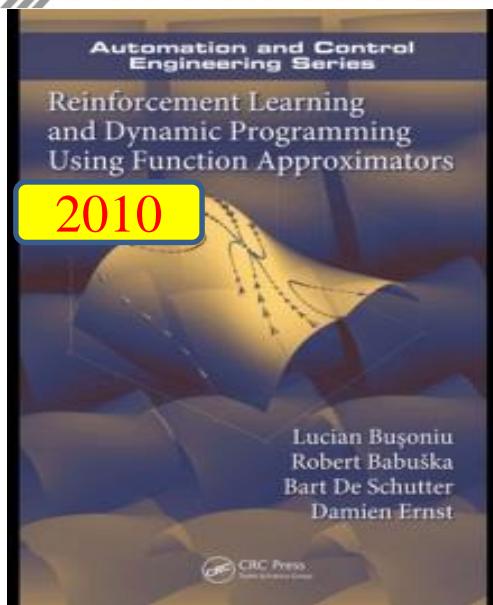


# Iterative adaptive dynamic programming

- In the theorem due to Frank Lewis and his co-workers:
  - The sequence  $\{V_i\}$  is monotonically non-decreasing because of the choice  $V_0 \equiv 0$ ;
  - From the  $i$ th iteration, one obtains  $v_i$  and  $V_{i+1}$ .



# More books by 2013



# A New Book



Derong Liu, Qinglai Wei, Ding Wang,  
Xiong Yang, Hongliang Li

2016

## Adaptive Dynamic Programming with Applications in Optimal Control

October 10, 2016

<b>1</b>	<b>Overview of Adaptive Dynamic Programming . . . . .</b>	1
1.1	Introduction . . . . .	1
1.2	Reinforcement Learning . . . . .	3
1.3	Adaptive Dynamic Programming . . . . .	7
1.3.1	Basic Forms of Adaptive Dynamic Programming . . . . .	10
1.3.2	Iterative Adaptive Dynamic Programming . . . . .	15
1.3.3	ADP for Continuous-Time Systems . . . . .	18
1.3.4	Remarks . . . . .	21
1.4	Related Books . . . . .	22
1.5	About This Book . . . . .	26
	References . . . . .	27

### Part I Discrete-Time Systems

<b>2</b>	<b>Value Iteration ADP for Discrete-Time Nonlinear Systems . . . . .</b>	37
2.1	Introduction . . . . .	37
2.2	Optimal Control of Nonlinear Systems Using General Value Iteration . . . . .	39
2.2.1	Convergence Analysis . . . . .	40
2.2.2	Neural Network Implementation . . . . .	49
2.2.3	Generalization to Optimal Tracking Control . . . . .	52
2.2.4	Optimal Control of Systems with Constrained Inputs . . . . .	56
2.2.5	Simulation Studies . . . . .	59
2.3	Iterative $\theta$ -Adaptive Dynamic Programming Algorithm for Nonlinear Systems . . . . .	67
2.3.1	Convergence Analysis . . . . .	68
2.3.2	Optimality Analysis . . . . .	75
2.3.3	Summary of Iterative $\theta$ -ADP Algorithm . . . . .	79
2.3.4	Simulation Studies . . . . .	81
2.4	Conclusions . . . . .	85
	References . . . . .	85

xv

Springer



北京科技大学 自动化学院

University of Science and Technology School of Automation and Electrical Engineering

# Iterative adaptive dynamic programming

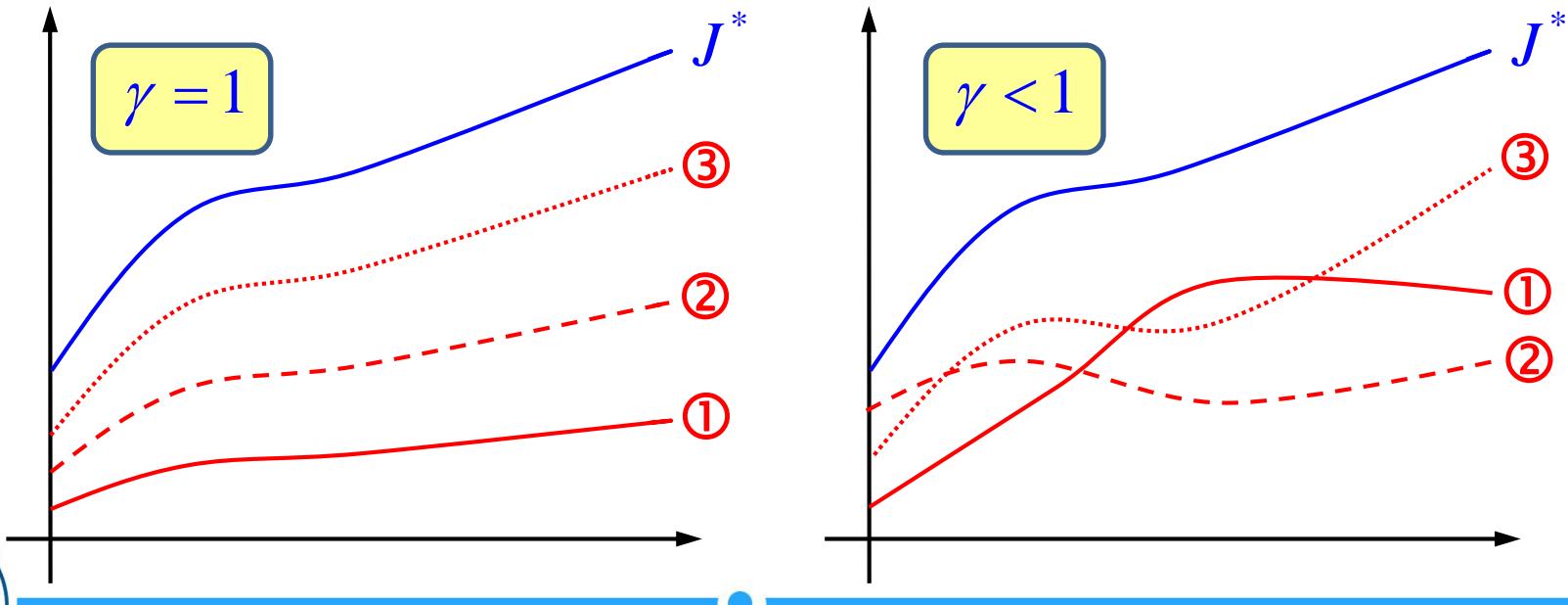
- Bellman Equation:

$$J^*(x_k) = \min_{u_k} \{U(x_k, u_k) + \gamma J^*(x_{k+1})\}$$

- Performance index (or cost to go)

$$J(x_k) = \sum_{i=k}^{\infty} \gamma^{i-k} U(x_i, u_i)$$

Most results  
are for  $\gamma = 1$ .

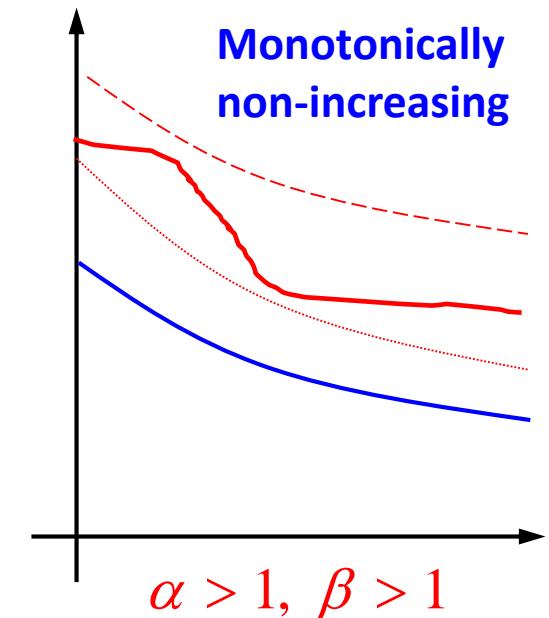
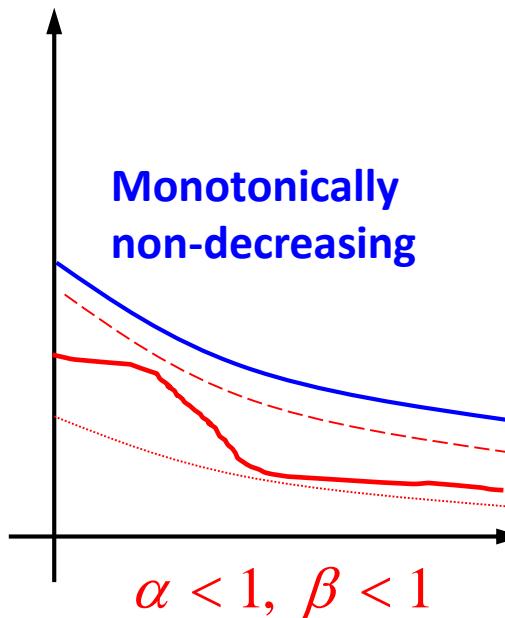
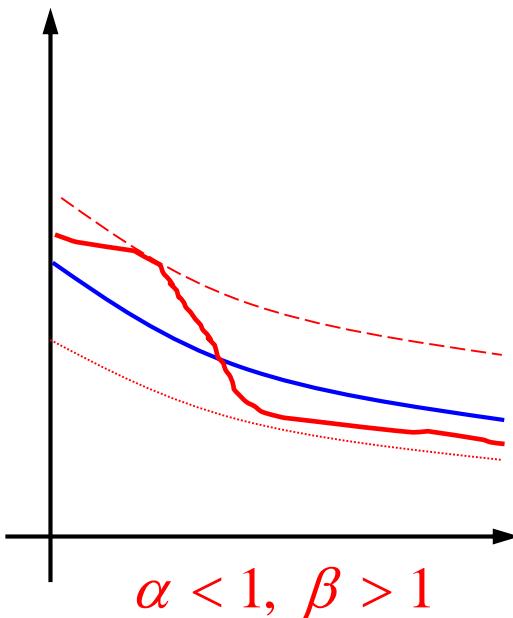


# Relaxed dynamic programming

- Theorem due to Rantzer and his co-workers:

$$\left[1 + \frac{\alpha - 1}{(1 + \rho^{-1})^i}\right] J^* \leq V_i \leq \left[1 + \frac{\beta - 1}{(1 + \rho^{-1})^i}\right] J^*$$

$$0 \leq \alpha J^* \leq V_0 \leq \beta J^*$$



# ADP for Self-Learning Control

**Dynamic Programming**

**Adaptive Dynamic Programming**

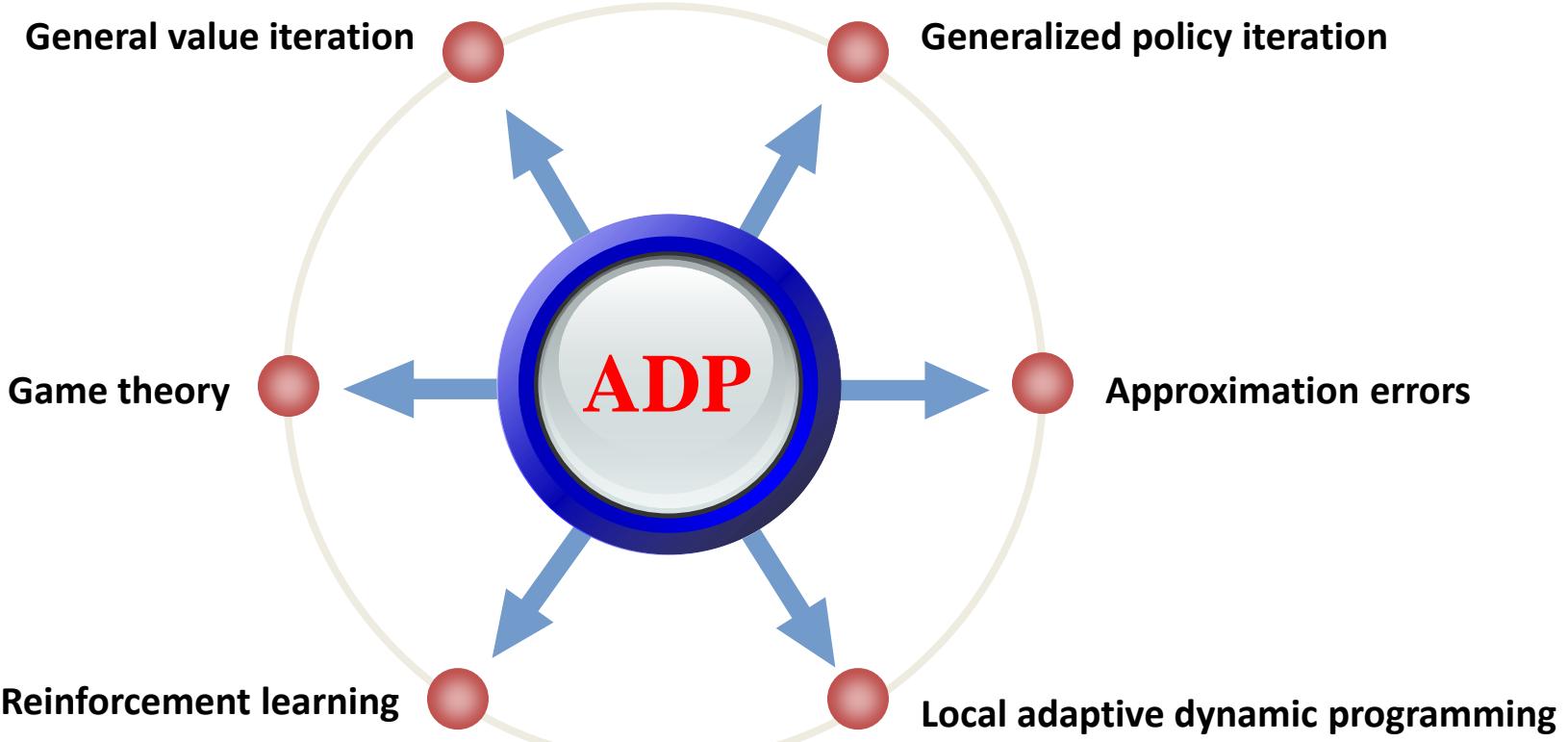
**Iterative ADP**

**New Developments**

**Concluding Remarks**



# New developments



# Value iteration

- Value iteration can be used to solve

$$J^*(x_k) = \min_{u_k} \{U(x_k, u_k) + J^*(x_{k+1})\}.$$

Choose the **initial** value function as  $V_0(\cdot) \equiv 0$ .

For  $i = 0, 1, 2, \dots$ , solve the control law as

$$v_i(x_k) = \arg \min_{u_k} \{U(x_k, u_k) + V_i(x_{k+1})\},$$

and update the value function as

$$V_{i+1}(x_k) = U(x_k, v_i(x_k)) + V_i(F(x_k, v_i(x_k))).$$

First, calculate

$$v_0(x_k) = \arg \min_{u_k} \{U(x_k, u_k) + V_0(x_{k+1})\}.$$

For  $i = 0, 1, 2, \dots$ , do value function update

$$\begin{aligned} V_i(x_k) &= \min_{u_k} \{U(x_k, u_k) + V_{i-1}(x_{k+1})\} \\ &= U(x_k, v_{i-1}(x_k)) + V_{i-1}(F(x_k, v_{i-1}(x_k))) \end{aligned}$$

and policy improvement

$$\begin{aligned} v_i(x_k) &= \arg \min_{u_k} \{U(x_k, u_k) + V_i(x_{k+1})\} \\ &= \arg \min_{u_k} \{U(x_k, u_k) + V_i(F(x_k, u_k))\}. \end{aligned}$$



# ① General value iteration

- The **initial** value function  $V_0$  can be chosen as:

$$V_0(x_k) = x_k^T P_0 x_k$$

–  $P_0$  is positive definite

$$V_0(x_k) = \Psi(x_k)$$

– positive semidefinite

First, calculate

$$v_0(x_k) = \arg \min_{u_k} \{U(x_k, u_k) + V_0(x_{k+1})\}.$$

For  $i = 0, 1, 2, \dots$ , do value function update

$$V_i(x_k) = \min_{u_k} \{U(x_k, u_k) + V_{i-1}(x_{k+1})\} = U(x_k, v_{i-1}(x_k)) + V_{i-1}(F(x_k, v_{i-1}(x_k)))$$

and policy improvement

$$v_i(x_k) = \arg \min_{u_k} \{U(x_k, u_k) + V_i(x_{k+1})\} = \arg \min_{u_k} \{U(x_k, u_k) + V_i(F(x_k, u_k))\}.$$



# ① General value iteration

- The **initial** value function  $V_0$  can be chosen as:

$$V_0(x_k) = x_k^T P_0 x_k$$

–  $P_0$  is positive definite

$$V_0(x_k) = \Psi(x_k)$$

– positive semidefinite

- ◎ If  $V_0 \leq V_1$ , then  $V_i \leq V_{i+1}$  for all  $i$ .
- ◎ If  $V_0 \geq V_1$ , then  $V_i \geq V_{i+1}$  for all  $i$ .

The monotonicity of  $\{V_i\}$  depends on the relationship between  $V_0$  and  $V_1$ .

$V_0 = 0 \Rightarrow$  monotonically non-decreasing;

$\theta$ -ADP:  $V_0 = \theta \Psi$ ,  $\theta > 0$  large and  $\Psi > 0 \Rightarrow$  monotonically non-decreasing.



# ① General value iteration

$\theta$ -ADP:  $V_0 = \theta \Psi$ ,  $\theta > 0$  large and  $\Psi > 0 \Rightarrow$   
monotonically non-decreasing.

$$\Psi \in \bar{\Psi}_x = \left\{ \Psi(x_k) : \Psi > 0 \text{ and } \exists v(x_k) \right.$$

s.t.  $\Psi(F(x_k, v(x_k))) < \Psi(x_k)$

## ◎ Value function update

$$V_i(x_k) = \min_{u_k} \{U(x_k, u_k) + V_{i-1}(x_{k+1})\}$$
$$= U(x_k, v_{i-1}(x_k)) + V_{i-1}(F(x_k, v_{i-1}(x_k)))$$

is more likely to obtain a **non-increasing** sequence.



# Policy iteration

- Policy iteration can be used to solve

$$J^*(x_k) = \min_{u_k} \{U(x_k, u_k) + J^*(x_{k+1})\}.$$

Choose an **initial** value function. Calculate

$$v_0(x_k) = \arg \min_{u_k} \{U(x_k, u_k) + V_0(x_{k+1})\}.$$

For  $i = 0, 1, 2, \dots$ , do value function update

$$\begin{aligned} V_i(x_k) &= \min_{u_k} \{U(x_k, u_k) + V_{i-1}(x_{k+1})\} \\ &= U(x_k, v_{i-1}(x_k)) + V_{i-1}(F(x_k, v_{i-1}(x_k))) \end{aligned}$$

and policy improvement

$$\begin{aligned} v_i(x_k) &= \arg \min_{u_k} \{U(x_k, u_k) + V_i(x_{k+1})\} \\ &= \arg \min_{u_k} \{U(x_k, u_k) + V_i(F(x_k, u_k))\}. \end{aligned}$$

**Value iteration**

Choose an **initial admissible control law**

$$v_0(x_k).$$

**Stable + cost function is finite.**

For  $i = 1, 2, \dots$ , do policy evaluation

$$V_i(x_k) = U(x_k, v_{i-1}(x_k)) + V_i(F(x_k, v_{i-1}(x_k)))$$

and policy improvement

$$\begin{aligned} v_i(x_k) &= \arg \min_{u_k} \{U(x_k, u_k) + V_i(x_{k+1})\} \\ &= \arg \min_{u_k} \{U(x_k, u_k) + V_i(F(x_k, u_k))\}. \end{aligned}$$

**Policy iteration**



# Policy iteration

- Policy iteration can be used to solve

$$J^*(x_k) = \min_{u_k} \{U(x_k, u_k) + J^*(x_{k+1})\}.$$

Choose an initial admissible control law

$$v_0(x_k).$$

For  $i = 1, 2, \dots$ , do policy evaluation

$$V_i(x_k) = U(x_k, v_{i-1}(x_k)) + V_i(F(x_k, v_{i-1}(x_k)))$$

and policy improvement

$$\begin{aligned} v_i(x_k) &= \arg \min_{u_k} \{U(x_k, u_k) + V_i(x_{k+1})\} \\ &= \arg \min_{u_k} \{U(x_k, u_k) + V_i(F(x_k, u_k))\}. \end{aligned}$$

Require infinite number  
of iterations to evaluate!



## ② Generalized policy iteration

Choose an initial admissible control law  $v_{(-1)}(x_k)$ .

Construct  $V_0$  from  $v_{(-1)}$  using

$$V_0(x_k) = U(x_k, v_{(-1)}(x_k)) + V_0(F(x_k, v_{(-1)}(x_k)))$$

and calculate  $v_0$  by

$$v_0(x_k) = \arg \min_{u_k} \{U(x_k, u_k) + V_0(x_{k+1})\} = \arg \min_{u_k} \{U(x_k, u_k) + V_0(F(x_k, u_k))\}.$$

For each  $i$ ,  $i = 1, 2, \dots$ , do  $N_i$ -step policy evaluation

$$V_{i,j_i}(x_k) = U(x_k, v_{i-1}(x_k)) + V_{i,j_i-1}(F(x_k, v_{i-1}(x_k))), \quad j_i = 1, 2, \dots, N_i,$$

to obtain  $V_i(x_k) = V_{i,N_i}(x_k)$  with  $V_{i,0}(x_{k+1}) = V_{i-1}(x_{k+1})$ .

For each  $i$ , do policy improvement

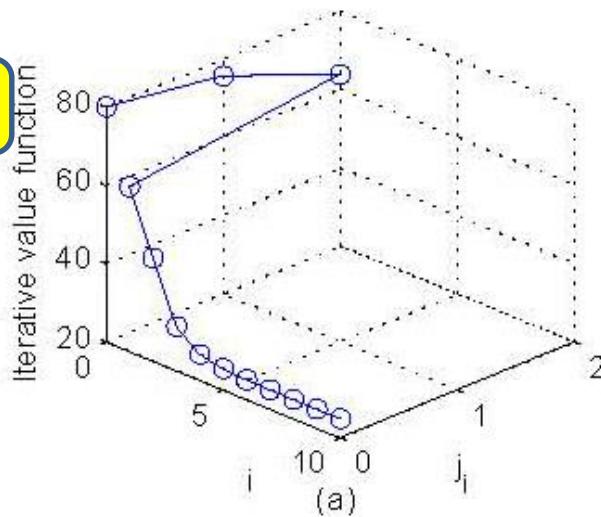
$$v_i(x_k) = \arg \min_{u_k} \{U(x_k, u_k) + V_i(x_{k+1})\} = \arg \min_{u_k} \{U(x_k, u_k) + V_i(F(x_k, u_k))\}.$$

$N_i = 1 \Rightarrow$  VI  
 $N_i = \infty \Rightarrow$  PI

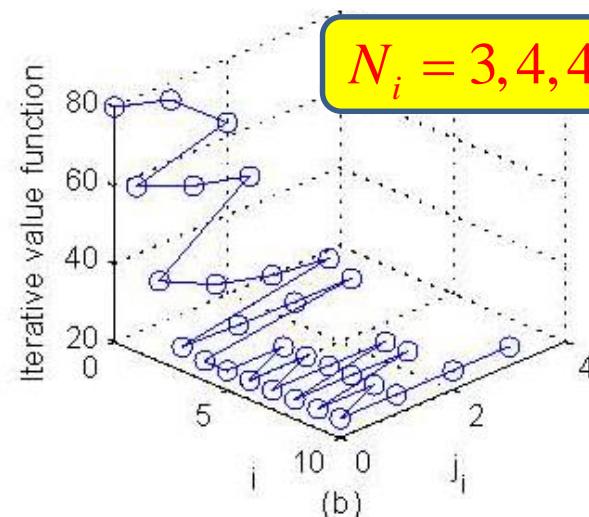


## ② Generalized policy iteration

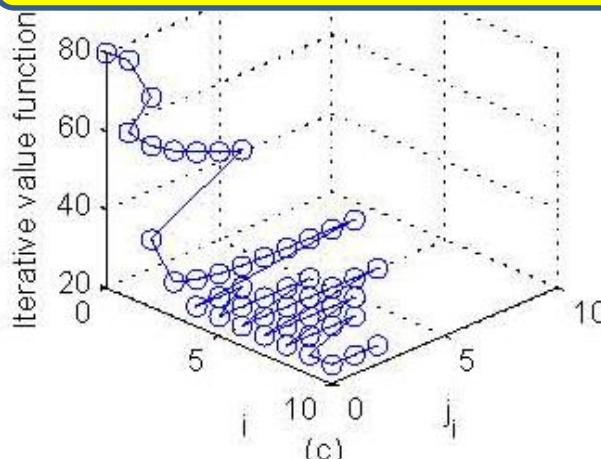
$$N_i = 1$$



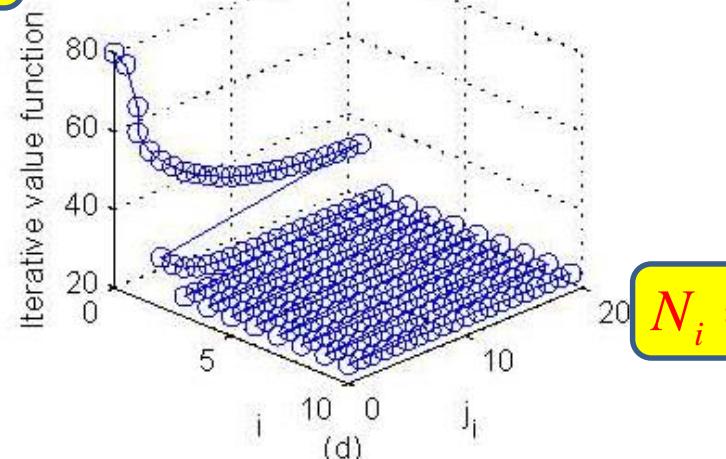
$$N_i = 3, 4, 4, 1, 2, 2, 3, 3, 2, 4.$$



$$N_i = 6, 1, 9, 3, 5, 7, 5, 4, 1, 3.$$



$$N_i = 20$$



# ③ ADP with approximation errors

- Considering NN approximation errors, realistic updates are

$$\hat{v}_i(x_k) = \arg \min_{u_k} \left\{ U(x_k, u_k) + \hat{V}_i(F(x_k, u_k)) \right\} + \eta_i(x_k)$$

$$\hat{V}_{i+1}(x_k) = U(x_k, \hat{v}_i(x_k)) + \hat{V}_i(F(x_k, \hat{v}_i(x_k))) + \pi_i(x_k)$$

- Convergence result now becomes

$$\hat{V}_i(x_k) \leq \sigma \left[ 1 + \sum_{j=1}^i \frac{\rho^j \sigma^{j-1} (\sigma - 1)}{(\rho + 1)^j} + \frac{\rho^i \sigma^i (\delta - 1)}{(\rho + 1)^i} \right] J^*(x_k)$$

$$V_0 \leq \delta J^*, \quad J^*(f(x, u)) \leq \rho U(x, u)$$



# ③ ADP with approximation errors

- The result becomes

$$\lim_{i \rightarrow \infty} \hat{V}_i(x_k) = \hat{V}_\infty(x_k) \leq \sigma \left[ 1 + \frac{\rho(\sigma - 1)}{1 - \rho(\sigma - 1)} \right] J^*(x_k)$$

if  $1 \leq \sigma \leq \frac{\rho + 1}{\rho}$

Finite neighborhood  
of the optimal solution

where

$$\hat{V}_i(x_k) \leq \sigma \Gamma_i(x_k)$$

$$\Gamma_i(x_k) = \min_{u_k} \left\{ U(x_k, u_k) + \hat{V}_{i-1}(x_{k+1}) \right\}$$



# ④ Local adaptive dynamic programming

- Every update is required to be done for the whole state space: VI, GVI, PI, GPI.
- Partition state space

$$\Omega_x = L_x^0 \cup L_x^1 \cup \dots \cup L_x^\zeta = \bigcup_{j=0}^{\zeta} L_x^j$$

- At each iteration step, only update in one of the subspaces  $L_x^j$ .

Asynchronous ADP



# ⑤ Connections to reinforcement learning

- TD, TD( $\lambda$ ), Q-learning, SARSA, Q( $\lambda$ ), SARSA( $\lambda$ ).
- ADP:
  - HDP  $\Leftrightarrow$  TD
  - ADHDHP  $\Leftrightarrow$  Q-learning
  - DHP, ADDHP
  - GDHP, ADGDHP
- SARSA, SARSA( $\lambda$ ): On-policy
$$Q_{i+1}(s_t, a_t) = Q_i(s_t, a_t) + \alpha [r_{t+1} + \gamma Q_i(s_{t+1}, a_{t+1}) - Q_i(s_t, a_t)].$$
Q-learning, Q( $\lambda$ ): Off-policy
$$Q_{i+1}(s_t, a_t) = Q_i(s_t, a_t) + \alpha \left[ r_{t+1} + \gamma \max_a \{Q_i(s_{t+1}, a)\} - Q_i(s_t, a_t) \right].$$
- Off-policy learning: Evaluating one policy while following another.



# ⑤ Connections to reinforcement learning

- HDP  $\Leftrightarrow$  TD
- ◎ In TD ( $\text{TD}(\lambda)$  as well):

$$V_{i+1}(s_t) = V_i(s_t) + \alpha[r_{t+1} + \gamma V_i(s_{t+1}) - V_i(s_t)]$$

[Target – OldValue]

- ◎ In model-free HDP, critic network training by minimizing:

$$E = \underbrace[U_{k+1} + \gamma V_i(x_{k+1}) - V_i(x_k)]_{{\text{Training target}}}$$

$\underbrace[-V_i(x_k)]_{{\text{To be trained}}}$



# ⑥ Game theory

- Consider

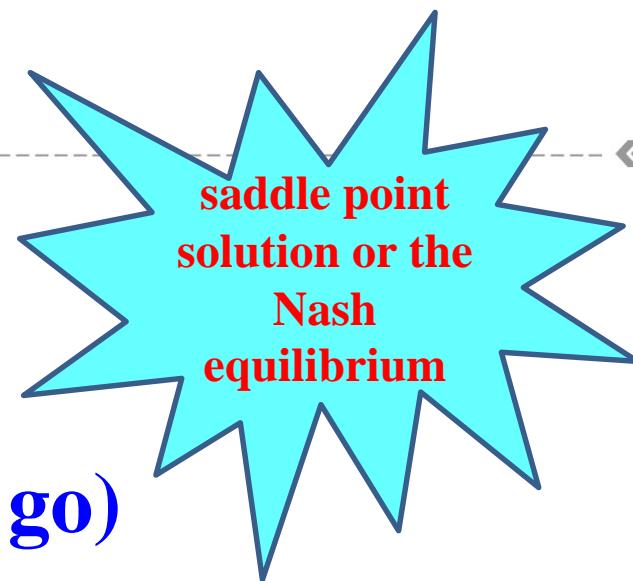
$$x_{k+1} = f(x_k) + g(x_k)u_k + h(x_k)w_k$$

- Performance index (or cost to go)

$$V(x_k, u_k, w_k) = \sum_{i=k}^{\infty} \left\{ x_i^T Q x_i + u_i^T R u_i - \gamma^2 w_i^T w_i \right\}$$

- Our goal is to find the control policy (player 1) and disturbance policy (player 2) so that

$$J^*(x_0) = \min_{u_k} \max_{w_k} \{V(x_0, u_k, w_k)\}$$



# ⑥ Game theory

- Iterative algorithm based on ADP

$$V_i(x_k) = \min_{u_k} \max_{w_k} \left\{ x_k^T Q x_k + u_k^T R u_k - \gamma^2 w_k^T w_k + V_{i-1} \left( f(x_k) + g(x_k) u_k + h(x_k) w_k \right) \right\}$$

$$u_i(x_k) = -\frac{1}{2} R^{-1} g^T(x_k) \frac{\partial V_i(x_{k+1})}{x_{k+1}}$$

$$w_i(x_k) = \frac{1}{2} \gamma^{-2} h^T(x_k) \frac{\partial V_i(x_{k+1})}{x_{k+1}}$$

- Convergence analysis + extensions to multiplayer games
- Linear systems and *nonlinear* systems



# ADP for Self-Learning Control

Dynamic Programming

Adaptive Dynamic Programming

Iterative ADP

New Developments

Concluding Remarks



# Concluding Remarks

- ① ADP has a long history. Theoretical development started in the 1970s.
- ② Adaptive dynamic programming is a robust learning control approach.
- ③ It has a close relationship to reinforcement learning.
- ④ Ultimate goal: To solve the curse of dimensionality – Need new ideas, to solve dynamic programming with less computation.



# For more information

- [derongliu@foxmail.com](mailto:derongliu@foxmail.com)
- <http://www.derongliu.org>
- 13810670526



# For more information

- [derongliu@foxmail.com](mailto:derongliu@foxmail.com)
- <http://www.derongliu.org>
- 13810670526

