Approximate Dynamic Programming for High-Dimensional Problems

2007 IEEE Symposium on Approximate Dynamic Programming and Reinforcement Learning April, 2007



Warren Powell *CASTLE Laboratory* **Princeton University** http://www.castlelab.princeton.edu

© 2007 Warren B. Powell, Princeton University





The fractional jet ownership industry



NetJets Inc.

XEE

F

and in case of







Planning for a risky world

Weather

- Robust design of emergency response networks.
- Design of financial instruments to hedge against weather emergencies to protect individuals, companies and municipalities.
- Design of sensor networks and communication systems to manage responses to major weather events.



Disease

- Models of disease propagation for response planning.
- Management of medical personnel, equipment and vaccines to respond to a disease outbreak.
- Robust design of supply chains to mitigate the disruption of transportation systems.



© 2007 Warren B. Powell

Blood management



Managing financial portfolios

■ Money can be invested and then reinvested....



© 2007 Warren B. Powell

Energy management

Applications

- Jet fuel hedging Designing strategies to hedge against fluctuations in jet fuel (and other commodities).
- Valuing energy contracts
- Planning the use of future technologies
- R&D portfolio management





Research in ADP

- Convergence proofs
- Rate of convergence research
- Design and evaluation of approximation strategies
- Design of advanced approximation strategies

NRG Energy

2

21個

14 J.

Challenges

Real-time control

- » Scheduling aircraft, pilots, generators, tankers
- » Pricing stocks, options

■ Near-term tactical planning

- » Can I accept a customer request?
- » Should I lease equipment?
- » Do I have to purchase extra energy?

Strategic planning

- » What is the right equipment mix?
- » What is the value of this contract?
- » What is the value of more reliable aircraft?

Outline

- The languages of dynamic programming
- A resource allocation model
- The post-decision state variable
- Example: A discrete resource: the nomadic trucker
- The states of our system
- Example: A continuous resource: blood inventory management
- Approximation methods
 - » Lookup tables and aggregation
 - » Basis functions
- Stepsizes
- Exploration vs. exploitation
- Applications

Languages

■ The languages of "optimization over time"

	Engineering	OR/AI/Probability	OR/Math programming
	Optimal	Markov decision	Stochastic
Discipline	control	processes	programming
Decision (English)	Control	Action	Decision
Decision (math)	u	а	Х
"Value function" (English)	Cost-to-go	Value function	Recourse function
"Value function" (Math)	J	V	Q
State variable	Х	S	Huh? (oh, "tenders")
Optimality equations	Hamilton-Jacobi	Bellman	Huh?

Languages

- "Approximate dynamic programming" has been discovered independently by different communities under different names:
 - » Neuro-dynamic programming
 - » Reinforcement learning
 - » Forward dynamic programming
 - » Adaptive dynamic programming
 - » Heuristic dynamic programming
 - » Iterative dynamic programming



How to land a plane:



- » Control: angle, velocity, acceleration, pitch, yaw...
- » Noise: wind, measurement

$$V_{t}(x_{t}) = \max_{u} \left(C(x_{t}, u_{t}) + EV_{t+1}(x_{t+1}) \right)$$

© 2007 Warren B. Powell

Slide 18



■ Where to send a plane:



- » Control: Where to send the plane to accomplish a goal.
- » Noise: demands on the system, equipment failures.

$$V_t(S_t) = \max_{a} \left(C(S_t, a_t) + EV_{t+1}(S_{t+1}) \right)$$

© 2007 Warren B. Powell

Slide 19



■ How to manage a fleet of aircraft:



- » Control: Which plane to assign to each customer.
- » Noise: demands on the system, equipment failures. $V_t(S_t) = \max_x \left(C(S_t, x_t) + EV_{t+1}(S_{t+1}) \right)$

© 2007 Warren B. Powell

Slide 20

A progression of models

Major problem classes

	Simple attributes	Complex attributes
Single entity	Textbook Markov decision process	Classical Al applications
Multiple entities	Classical OR applications	Opportunity for combining AI/OR

Sample applications

- Single entity problems
 - » Playing a board game
 - » Routing a truck around the country
 - » Planning a set of courses through college
- Storage problems (single resource class)
 - » Maintaining product inventories
 - » Purchasing commodity futures (oil, orange juice, ...)
- Managing multiple resource classes
- » Blood inventories
- » Fleet management (with different equipment types)
- Managing multiple, discrete resources
- » Locomotives, jets, people











Asset acquisition problems

Storage problems



$$R_{t+1} = \left[R_t + x_t - \hat{D}_{t+1} \right]^+$$
$$P_{t+1} = P_t + \hat{P}_t$$

Multiple inventory types

Managing blood inventories



Multiple inventory types

Managing blood inventories over time





Multiple discrete assets



Outline

- The languages of dynamic programming
- A resource allocation model
- The post-decision state variable
- Example: A discrete resource: the nomadic trucker
- The states of our system
- Example: A continuous resource: blood inventory management
- Approximation methods
 - » Lookup tables and aggregation
 - » Basis functions
- Stepsizes
- Exploration vs. exploitation
- Applications

A resource allocation model

Modeling resources:

- » The state of a single resource:
 - a = The attributes of a single resource

 $a \in A$ The attribute space

» The state of multiple resources:

 R_{ta} = The number of resources with attribute *a*

$$R_t = (R_{ta})_{a \in A}$$
 The resource state vector

» The information process:

 \hat{R}_{ta} = The change in the number of resources with attribute *a*.

A resource allocation model

Modeling demands:

- » The attributes of a single demand:
 - b = The attributes of a demand to be served.
 - $b \in \mathbf{B}$ The attribute space
- » The demand state vector:

 D_{tb} = The number of demands with attribute b

$$D_t = (D_{tb})_{b \in \mathbf{B}}$$
 The demand state vector

» The information process:

 \hat{D}_{tb} = The change in the number of demands with attribute *b*.

A resource allocation model

■ The system:

» The state vector

$$S_t = \left(R_t, D_t\right)$$

» The information process:

 $W_t = \text{Exogenous changes to resources and demands}$ = $\left(\overrightarrow{R_t}, D_t \right)$
■ The three states of our system

» The state of a single resource/entity

$$a_t = \begin{bmatrix} a_{t1} \\ a_{t2} \\ a_{t3} \end{bmatrix}$$

» The resource state vector

$$R_{t} = \begin{bmatrix} R_{ta_{1}} \\ R_{ta_{2}} \\ R_{ta_{3}} \end{bmatrix}$$

» The system state vector

$$S_t = \left(R_t, D_t\right)$$







Outline

- The languages of dynamic programming
- A resource allocation model
- The post-decision state variable
- Example: A discrete resource: the nomadic trucker
- The states of our system
- Example: A continuous resource: blood inventory management
- Approximation methods
 - » Lookup tables and aggregation
 - » Basis functions
- Stepsizes
- Exploration vs. exploitation
- Applications



Laying the foundation

Dynamic programming review:

» Let:

 S_t = "State" of our "system" at time t.

 x_t = "Action" that we take to change the system.

 $C(S_t, x_t)$ = Contribution earned when we take action x from state S_t .

» We model system dynamics using:

 $p(S_{t+1} | S_t, x_t) =$ Probability that action x_t takes us from state S_t to state S_{t+1}

Laying the foundation

Bellman's equation:

» Standard form:

$$V_t(S_t) = \max_x \left(C_t(S_t, x_t) + \sum_{s'} p(s' | S_t, x_t) V_{t+1}(S_{t+1} = s') \right)$$

» Expectation form:

$$V_t(S_t) = \max_{x} \left(C_t(S_t, x_t) + E\left\{ V_{t+1}(S_{t+1}(S_t, x_t)) | S_t \right\} \right)$$









Bellman's equation

■ We just solved Bellman's equation:

$$V_t(S_t) = \max_{x \in X} C_t(S_t, x_t) + E\{V_{t+1}(S_{t+1}) \mid S_t\}$$

» We found the value of being in each state by stepping backward through the tree.

Bellman's equation

■ The challenge of dynamic programming:

$$V_t(S_t) = \max_{x \in X} \left(C_t(S_t, x_t) + E\left\{ V_{t+1}(S_{t+1}) \mid S_t \right\} \right)$$

The curses of dimensionality

What happens if we apply this idea to our blood problem?

- » State variable is:
 - The supply of each type of blood, along with its age
 - 8 blood types
 - 6 ages
 - = 48 "blood types"
 - The demand for each type of blood
 - 8 blood types
- » Decision variable is how much of 48 blood types to supply to 8 demand types.
 - 216- dimensional decision vector
- » Random information
 - Blood donations by week (8 types)
 - New demands for blood (8 types)



The curses of dimensionality

■ The challenge of dynamic programming:

$$V_t(S_t) = \max_{x \in \mathcal{X}} \left(C_t(S_t, x_t) + E\left\{ V_{t+1}(S_{t+1}) \,|\, S_t \right\} \right)$$

Three curses

Problem: Curse of dimensionality State space

Outcome space

Action space (feasible region)

The curses of dimensionality

■ The computational challenge:





 $V_t(S_t) = \max_{x \in X} \left(C_t(S_t, x_t) + E\left\{ V_{t+1}(S_{t+1}) \mid S_t \right\} \right)$



■ New concept:

- » The "pre-decision" state variable:
 - S_t = The information required to make a decision x_t
 - Same as a "decision node" in a decision tree.
- » The "post-decision" state variable:
 - S_t^x = The state of what we know immediately after we make a decision.
 - Same as an "outcome node" in a decision tree.

■ Pre-decision, state-action, and post-decision



 3^9 states $3^9 \times 9$ state-action pairs 3^9 states

A single, complex entity

Pre- and post-decision attributes for our nomadic truck driver:



■ Pre-decision: resources and demands



$$S_t^x = S^{M,x}(S_t, x_t)$$



Slide 60



© 2007 Warren B. Powell

Slide 61



System dynamics

It is traditional to assume you are given the one-step transition matrix:

 $p(S_{t+1} | S_t, x_t)$ = Probability that action x_t takes us from state S_t to state S_{t+1}

- » Computing the transition matrix is impossible for the vast majority of problems.
- We are going to assume that we are given a *transition function*:

 $S_{t+1} = S^{M}(S_{t}, x_{t}, W_{t+1})$

- » This is at the heart of any simulation model.
- » Often rule-based. Very easy to compute, even for large-scale problems.

Working with pre- and post-decision states
» The "usual" transition function:

 $S_{t+1} = S^{M}(S_{t}, x_{t}, W_{t+1})$ From S_{t} to S_{t+1} .

» The transition function broken into two steps:

 $S_{t}^{x} = S^{M,x}(S_{t}, x_{t})$ The pure effect of a decision $S_{t+1} = S^{M,W}(S_{t}^{x}, W_{t+1})$ The effect of the exogenous information

The transition function

- Actually, we have three transition functions:
 - » The attribute transition function:

 $a_t^x = a^{M,x}(a_t, x_t)$ The pure effect of a decision $a_{t+1} = a^{M,W}(a_t^x, W_{t+1})$ The effect of the exogenous information

» The resource transition function

 $R_{t}^{x} = R^{M,x} \left(R_{t}, x_{t} \right)$ The pure effect of a decision $R_{t+1} = R^{M,W} \left(R_{t}^{x}, W_{t+1} \right)$ The effect of the exogenous information » The general transition function:

 $S_{t}^{x} = S^{M,x}(S_{t}, x_{t})$ The pure effect of a decision $S_{t+1} = S^{M,W}(S_{t}^{x}, W_{t+1})$ The effect of the exogenous information

Bellman's equations with the post-decision state

- Bellman's equations broken into stages:
 - » Optimization problem (making the decision): $V_t(S_t) = \max_x \left(C_t(S_t, x_t) + V_t^x \left(S_t^{M, x}(S_t, x_t) \right) \right)$
 - Note: this problem is deterministic!
 - » Simulation problem (the effect of exogenous information):

$$V_t^x(S_t^x) = E\left\{V_{t+1}(S^{M,W}(S_t^x, W_{t+1})) \mid S_t^x\right\}$$

Bellman's equations with the post-decision state

■ Challenges

» For most practical problems, we are not going to be able to compute $V_t^x(S_t^x)$.

$$V_t(S_t) = \max_x \left(C_t(S_t, x_t) + V_t^x(S_t^x) \right)$$

» Concept: replace it with an approximation $\overline{V_t}(S_t^x)$ and solve

$$V_t(S_t) = \max_x \left(C_t(S_t, x_t) + \overline{V_t}(S_t^x) \right)$$

- » So now we face:
 - What should the approximation look like?
 - How do we estimate it?

Approximating the value function

- For "resource allocation" problems, we have been using:
 - » Linear (in the resource state):

$$\overline{V}_t(R_t^x) = \sum_{a \in \mathcal{A}} \overline{v}_{ta} \cdot R_{ta}^x$$

Best when assets are complex, which means that R_{ta} is small (typically 0 or 1).

» Piecewise linear, separable:

$$\overline{V_t}(R_t^x) = \sum_{a \in \mathcal{A}} \overline{V_{ta}}(R_{ta}^x)$$

Best when assets are simple, which means that R_{ta} may be larger.

Our general algorithm

Step 1: Start with a post-decision state $S_{t-1}^{x,n}$ Step 2: Obtain Monte Carlo sample of $W_t(\omega^n)$ and compute the next pre-decision state: $S_t^n = S^{M,W}(S_{t-1}^{x,n}, W_t(\omega^n))$ Simulation Step 3: Solve the deterministic optimization using an approximate value function: $\hat{v}_t^n = \max_x \left(C_t(S_t^n, x_t) + \overline{V}_t^{n-1}(S^{M,x}(S_t^n, x_t)) \right)$ Optimization

to obtain x_t^n .

Step 4: Update the value function approximation $\overline{V}_{t-1}^{n}(S_{t-1}^{x,n}) = (1 - \alpha_{n-1})\overline{V}_{t-1}^{n-1}(S_{t-1}^{x,n}) + \alpha_{n-1}\hat{v}_{t}^{n}$ Statistics Step 5: Find the next post-decision state: $S_{t}^{x,n} = S^{M,x}(S_{t}^{n}, x_{t}^{n})$

© 2007 Warren B. Powell

Slide 69

Competing updating methods

- Comparison to other methods:
 - » Classical MDP (value iteration)

$$V^{n}(S) = \max_{x} \left(C(S, x) + EV^{n-1}(S_{t+1}) \right)$$

- S Classical ADP (pre-decision state): $\hat{v}_t^n = \max_x \left(C_t(S_t^n, x_t) + \sum_{s'} p(s' | S_t^n, x_t) \overline{V_{t+1}}(s') \right)$ $\overline{V_t^n}(S_t^n) = (1 \alpha_{n-1}) \overline{V_t^{n-1}}(S_t^n) + \alpha_{n-1} \hat{v}_t^n \qquad \hat{v}_t \text{ updates } \overline{V_t}(S_t)$
- » Our method (update $\overline{V}_t^{x,n-1}$ around post-decision state): $\hat{v}_t^n = \max_x \left(C_t(S_t^n, x_t) + \overline{V}_t^{x,n-1}(S^{M,x}(S_t^n, x_t)) \right)$ $\overline{V}_{t-1}^n(S_{t-1}^{x,n}) = (1 - \alpha_{n-1}) \left(\overline{V}_{t-1}^{n-1}(S_{t-1}^{x,n}) + \alpha_{n-1} \hat{v}_t^n \right)$ \hat{v}_t updates $\overline{V}_{t-1}(S_{t-1}^x)$

© 2007 Warren B. Powell

APPROXIMATE DYNAMIC PROGRAMMING

Solving the curses of dimensionality

Warren B. Powell



A JOHN WILEY & SONS, INC., PUBLICATION

Outline

- The languages of dynamic programming
- A resource allocation model
- The post-decision state variable
- Example: A discrete resource: the nomadic trucker
- The states of our system
- Example: A continuous resource: blood inventory management
- Approximation methods
 - » Lookup tables and aggregation
 - » Basis functions
- Stepsizes
- Exploration vs. exploitation
- Applications
■ The previous post-decision state: trucker in Texas



■ Pre-decision state: we see the demands



■ We use initial value function approximations...



• ... and make our first choice: x^1



■ Update the value of being in Texas.



Now move to the next state, sample new demands and make a new decision



■ Update value of being in NY



Move to California.



Make decision to return to TX and update value of being in CA



Back in TX, we repeat the process, observing a different set of demands.



We get a different decision and a new estimate of the value of being in TX



Updating the value function:

Old value:

 $\overline{V}^1(TX) = \$450$

New estimate:

 $\hat{v}^2(TX) = \$800$

How do we merge old with new?

$$\overline{V}^{2}(TX) = (1 - \alpha)\overline{V}^{1}(TX) + (\alpha)\hat{v}^{2}(TX)$$
$$= (0.90)\$450 + (0.10)\$800$$
$$= \$485$$

■ An updated value of being in TX



Outline

- The languages of dynamic programming
- A resource allocation model
- The post-decision state variable
- Example: A discrete resource: the nomadic trucker
- The states of our system
- Example: A continuous resource: blood inventory management
- Approximation methods
 - » Lookup tables and aggregation
 - » Basis functions
- Stepsizes
- Exploration vs. exploitation
- Applications

■ Now let's take a look at what we just did:

- a_t = Attribute of our nomadic trucker at time *t* before the decision is made.
- \hat{D}_t = Vector of demands that are revealed at time *t*.
- $S_t = (a_t, \hat{D}_t)$ Pre-decision state variable.
- a_t^x = Attribute of our nomadic trucker at time *t* after the decision is made.

 $S_t^x = a_t^x$ Post-decision state variable.

A single, complex entity

Pre- and post-decision attributes for our nomadic truck driver:



Multiple, complex entities

■ Notation for multiple entities:

- » The truck state vector:
 - a = The attributes of the truck
 - $a \in A$ The attribute space

 $R_{ta}^{truck} = \text{The number of trucks with attribute } a$ $R_{t}^{truck} = \left(R_{ta}^{truck}\right)_{a \in A} \text{The truck state vector}$

» The information process:

 \hat{R}_{ta}^{truck} = The change in the number of trucks with attribute *a*.

Multiple, complex entities

- Modeling the fleet management problem:
 - » The load state vector:
 - b = The attributes of a load to be moved.
 - $b \in \mathbf{B}$ The attribute space
 - R_{tb}^{load} = The number of tasks with attribute *b* $R_{t}^{load} = \left(R_{tb}^{load}\right)_{b \in B}$ The load state vector
 - » The information process:

 \hat{R}_{tb}^{load} = The change in the number of loads with attribute *b*.

Multiple, complex entities

- Modeling the fleet management problem:
 - » The resource state vector (a.k.a. "physical state")

$$R_t = \left(R_t^{truck}, R_t^{load}\right)$$

- » The information process:
 - $\hat{R}_{t} = \text{The number of new arrivals (of drivers and loads)}$ during time interval t. $= \left(\vec{R}_{t}^{\text{insuck}}, R_{t}^{\text{load}} \right)$ $= W_{t}$

■ The state of a single, simple entity:

$$a = [Location]$$



$$a \in \mathbf{A}$$
 $|\mathbf{A}| \approx 100 - 10,000$

■ The state of a single, complex entity:

 $a = \begin{bmatrix} Time \\ Location \\ Equipment type \\ Home base \\ Operator attributes \\ Time in service \\ Maintenance status \end{bmatrix}$



 $a \in \mathbf{A}$ $|\mathbf{A}| \approx 10^{10} - 10^{100}$ The curse of dimensionality!

■ Multiple, complex entities



The number of *dimensions* of our state variable is equal to the size of the *state space* for a single entity problem.

The curse of curses.

■ What is a state variable?

» A minimally dimensioned function of history that necessary and sufficient to compute the decision function, the transition function and the contribution function.

■ The three states of our system

» The state of a single resource/entity

$$a_t = \begin{bmatrix} a_{t1} \\ a_{t2} \\ a_{t3} \end{bmatrix}$$

» The state of all our resources

$$R_{t} = \begin{bmatrix} R_{ta_{1}} \\ R_{ta_{2}} \\ R_{ta_{3}} \end{bmatrix}$$

» The state of knowledge

$$S_t = (R_t, \overline{\theta}_t)$$
 $\overline{\theta}_t = \text{Estimates of "other parameters"}$

■ The state variable



■ The state variable



What is missing from our state variable?



Outline

- The languages of dynamic programming
- A resource allocation model
- The post-decision state variable
- Example: A discrete resource: the nomadic trucker
- The states of our system
- Example: A continuous resource: blood inventory management
- Approximation methods
 - » Lookup tables and aggregation
 - » Basis functions
- Stepsizes
- Exploration vs. exploitation
- Applications

Blood management

Managing blood inventories



Blood management

Managing blood inventories over time











Updating the value function approximation

Estimate the gradient at R_t^n



Updating the value function approximation

• Update the value function at $R_{t-1}^{x,n}$


Updating the value function approximation

• Update the value function at $R_{t-1}^{x,n}$



Updating the value function approximation

Update the value function at $R_{t-1}^{x,n}$



Updating the value function approximation

■ The updated function may not be concave:











A projection algorithm (SPAR)



A projection algorithm (SPAR)



A projection algorithm (SPAR)







© 2007 Warren B. Powell



© 2007 Warren B. Powell



© 2007 Warren B. Powell



Outline

- The languages of dynamic programming
- A resource allocation model
- The post-decision state variable
- Example: A discrete resource: the nomadic trucker
- The states of our system
- Example: A continuous resource: blood inventory management
- Approximation methods
 - » Lookup tables and aggregation
 - » Basis functions
- Stepsizes
- Exploration vs. exploitation
- Applications

Multiattribute resources

Assets can have a number of attributes:









	Location	[Location]	[Location]	Location	
	Equipment type	ETA	ETA	ETA	
<i>a</i> =		Equipment type	Bus. segment	A/C type	
		Train priority	Single/team	Fuel level	
		Pool	Domicile	Home shop	
		Due for maint	Drive hours	Crew	
		Home shop	Duty hours	Eqpt1	
			8 day history	М	
			Days from home	Eqpt100	

Boxcars



Multiattribute resources

A

■ The evolution of attributes:

$$a = \begin{bmatrix} \text{Time} \\ \text{Location} \end{bmatrix} \begin{bmatrix} \text{Time} \\ \text{Location} \\ \text{Boxcar type} \end{bmatrix} \begin{bmatrix} \text{Time} \\ \text{Location} \\ \text{Boxcar type} \\ \text{Time to dest.} \end{bmatrix} \begin{bmatrix} \text{Time} \\ \text{Location} \\ \text{Boxcar type} \\ \text{Time to dest.} \\ \text{Repair status} \end{bmatrix} \begin{bmatrix} \text{Time} \\ \text{Location} \\ \text{Boxcar type} \\ \text{Time to dest.} \\ \text{Repair status} \\ \text{Shipper pool} \end{bmatrix}$$

$$\approx 4,000 \quad 40,000 \quad 1,680,000 \quad 5,040,000 \quad 50,400,000$$

The states of our system

■ Multiple, complex entities



The number of *dimensions* of our state variable is equal to the size of the *state space* for a single entity problem.

The curse of curses.





■ Different levels of aggregation



 $|\mathbf{A}| \approx 12,000,000 \quad 600,000 \quad 6,000 \quad 2,000$

Aggregation

Aggregation:

- » Exact methods
 - We have to use the same level of aggregation throughout (in particular, the transition matrix and value function).



Aggregation

Approximate DP

- » We only need to discretize the value function. We can capture the full state variable in the transition function:
 - Decision function:

$$x_t = \arg\max_x \left(C_t(S_t, x_t) + \overline{V_t}(S_t^x) \right)$$

Value function using aggregated state

• Transition functions

$$S_t = S^{M,W}(S_{t-1}^x, W_t(\omega))$$
$$S_t^x = S^{M,x}(S_t, x_t)$$

No aggregation

Updating the value of a driver:



Estimating value functions

» Most disaggregate level

$$\overline{v}^{n}\begin{pmatrix} Location \\ Fleet \\ Domicile \end{pmatrix} = (1-\alpha)\overline{v}^{n-1}\begin{pmatrix} Location \\ Fleet \\ Domicile \end{pmatrix} + \alpha\hat{v}\begin{pmatrix} Location \\ Fleet \\ Domicile \end{pmatrix}) + \alpha\hat{v}(\begin{array}{c} Domicile \\ Domicile \\ DOThrs \\ DaysFromHome \end{pmatrix})$$

Estimating value functions

» Middle level of aggregation

$$\overline{v}^{n} \begin{pmatrix} Location \\ Fleet \end{pmatrix} = (1-\alpha)\overline{v}^{n-1} \begin{pmatrix} Location \\ Fleet \end{pmatrix} + \alpha \hat{v} \begin{pmatrix} Location \\ Fleet \end{pmatrix} \\ Domicile \\ DOThrs \\ DaysFromHome \end{bmatrix}$$

Estimating value functions

» Most aggregate level

$$\overline{v}^{n}([Location]) = (1-\alpha)\overline{v}^{n-1}([Location]) + \alpha \hat{v}(\begin{bmatrix} Location \\ Fleet \\ Domicile \\ DOThrs \\ DaysFromHome \end{bmatrix}$$





Using different levels of aggregation:

- » Pick the (single) level of aggregation that produces the best overall results.
- » Pick the level of aggregation that produces the lowest variance for each state.
- » Use a weighted sum of estimates at each level of aggregation (weight depends only on the level of aggregation):

$$\overline{v}_a^n = \sum_g w^{(g,n)} \overline{v}_a^{(g,n)}$$

» Use a weighted sum, but where the weights depend on the state (attribute):

$$\overline{v}_a^n = \sum_g w_a^{(g,n)} \overline{v}_a^{(g,n)}$$



State-dependent weighted aggregation:

» There may be hundreds of thousands of weights, so these have to be easy to compute.



Both can be computed using simple recursive formulas.

© 2007 Warren B. Powell





© 2007 Warren B. Powell
Aggregation for table lookup



© 2007 Warren B. Powell

Outline

- The languages of dynamic programming
- A resource allocation model
- The post-decision state variable
- Example: A discrete resource: the nomadic trucker
- The states of our system
- Example: A continuous resource: blood inventory management
- Approximation methods
 - » Lookup tables and aggregation
 - » Basis functions
- Stepsizes
- Exploration vs. exploitation
- Applications

Basis functions

- Aggregation works well when we have a state (attribute) space with very little structure.
- But what if we have some structure? Consider our inventory problem:



Basis functions

- Approximating the value function:
 - » We have to exploit the structure of the value function (e.g. concavity).
 - » We might approximate the value function using a simple polynomial

$$\overline{V_t} (R_t \mid \theta) = \theta_0 + \theta_1 R_t + \theta_2 R_t^2$$

» .. or a complicated one:

$$\overline{V_t} (R_t | \theta) = \theta_0 + \theta_1 R_t + \theta_2 R_t^2 + \theta_3 \ln(R_t) + \theta_4 \sin(R_t)$$

» Sometimes, they get really messy:

$$\overline{V}_{t}(R \mid \theta) = \theta^{(0)} + \sum_{s} \sum_{t'=t}^{t+2} \theta^{(1)}_{t,st'} R_{t,st'} + \sum_{w} \sum_{t'=t}^{t+3} \theta^{(1)}_{t,wt'} R_{t,wt'}$$

 $+\sum_{w}\sum_{t'}\theta_{t,wt'}^{(2)}R_{t,wt'}^{2} + \sum_{s}\sum_{t'}\theta_{t,st'}^{(2)}R_{t,st'}^{2}$ 8-14

1,2,3 4,5,6,7

$$+\sum_{s} \theta_{ts}^{(3)} \left(R_{t,st} - \frac{1}{S} \sum_{s'} R_{t,s't} \right)^{2}$$
 15

$$+\sum_{s} \theta_{ts}^{(4)} \left(\left(\sum_{t'=t}^{t+1} R_{t,st'} \right) - \frac{1}{2S} \sum_{s'} \sum_{t'=t}^{t+1} R_{t,s't'} \right)^{2}$$
 16

$$+\sum_{s} \theta_{ts}^{(5)} \left(\left(\sum_{t'=t}^{t+2} R_{t,st'} \right) - \frac{1}{3S} \sum_{s'} \sum_{t'=t}^{t+2} R_{t,s't'} \right)^{2}$$
 17

$$+\sum_{w}\sum_{s}\theta_{t,ws}^{(ws,1)}R_{t,wt}R_{t,st}$$
18

$$+\sum_{w}\sum_{s}\theta_{t,ws}^{(ws,2)}R_{t,wt}\left(\sum_{t'=t}^{t+2}R_{t,st'}\right)$$
19

$$+\sum_{w}\sum_{s} \theta_{t,ws}^{(ws,3)} \left(\sum_{t'=t}^{t+3} R_{t,wt'}\right) \left(\sum_{t'=t}^{t+2} R_{t,st'}\right)$$

$$20$$

$$+\sum_{w}\sum_{s} \theta_{t,ws}^{(ws,4)} \left(\sum_{t'=t} R_{t,wt'}\right) \left(\sum_{t'=t} R_{t,st'}\right)$$
²¹

$$+\sum_{w}\sum_{s}\theta_{t,sw}^{(ws,5)}R_{t,s,t+2}R_{t,wt}R_{t,w,t+2}$$
22

We can write a model of the observed value of being in a state as:

$$\hat{v} = \theta_0 + \theta_1 R_t + \theta_2 R_t^2 + \theta_3 \ln(R_t) + \theta_4 \sin(R_t) + \varepsilon$$

■ This is often written as a generic regression model:

$$Y = \theta_0 + \theta_1 X_1 + \theta_2 X_2 + \theta_3 X_3 \qquad + \theta_4 X_4$$

The ADP community refers to the independent variables as *basis functions*:

$$Y = \theta_0 \varphi_0(R) + \theta_1 \varphi_1(R) + \theta_2 \varphi_2(R) + \theta_3 \varphi_3(R) + \theta_4 \varphi_4(R)$$

 $= \sum_{f \in \mathcal{F}} \theta_f \varphi_f(R) \qquad \varphi_f(R) \text{ are also known as } features.$

© 2007 Warren B. Powell

Basis functions

• Methods for estimating θ

- » Generate observations $\mathfrak{E}_{v^2}, ..., v^N$, and use traditional regression methods to fit θ .
- » Use recursive statistics update θ^n after each iteration:

Basis functions

Notes:

- » When using basis functions, we are basically drawing on the entire field of statistics.
- » Designing basis functions (independent variables) is mostly art.
- » In special cases, the resulting algorithm can produce optimal solutions.
- » Most of the time, we are hoping for "good" solutions.
- » In some cases, it can work terribly.
- » As a general rule you have to use problem structure. Value function approximations have to capture the right structure. Blind use of polynomials will rarely be successful.

Outline

- The languages of dynamic programming
- A resource allocation model
- The post-decision state variable
- Example: A discrete resource: the nomadic trucker
- The states of our system
- Example: A continuous resource: blood inventory management
- Approximation methods
 - » Lookup tables and aggregation
 - » Basis functions
- Stepsizes
- Exploration vs. exploitation
- Applications

Stepsizes:

» Fundamental to ADP is an updating equation that looks



■ Theory:

» Many convergence results require:

$$\sum_{n=1}^{\infty} \alpha_{n-1} = \infty$$
$$\sum_{n=1}^{\infty} (\alpha_{n-1})^2 < \infty$$

» For example:

n=1

$$\alpha_{n-1} = \frac{1}{n}$$

Practice

- » 1/n "doesn't work"
- » Constant stepsizes
- » Various stepsize rules
 - Deterministic
 - Stochastic

Rate of convergence



Rate of convergence

The challenge of stepsizes: » When have we converged?



We need to improve our understanding of adaptive stepsizes.

© 2007 Warren B. Powell

Deterministic stepsize rules:

$$\alpha_n = \frac{1}{n} \text{ or } \frac{1}{n^\beta}$$
$$\alpha_n = \frac{a}{a+n-1}$$
$$\frac{b}{n+a} + a$$
$$\alpha_n = \frac{\frac{b}{n}}{\frac{b}{n} + a + n^\beta}$$

Basic averaging

Slows the rate of descent

"Search then converge"

$$\alpha_n = \frac{\alpha_n}{1 + \alpha_n - \overline{\alpha}}$$

McClain's formula

- The right stepsize rule depends on the rate of change in the value function.
- This varies widely for different parameters in the same problem:



© 2007 Warren B. Powell

- The right stepsize rule depends on the rate of change in the value function.
- This varies widely for different parameters in the same problem:



© 2007 Warren B. Powell

- The right stepsize rule depends on the rate of change in the value function.
- This varies widely for different parameters in the same problem:



© 2007 Warren B. Powell

Stochastic stepsize rules:

Let: ε^n = New observation - old estimate

Tunable parameters

$$\alpha_n = \alpha^0 \frac{a}{a+K^n} \qquad K^n = K^{n-1} + \mathbb{1}_{\left\{\varepsilon^n \varepsilon^{n-1} < 0\right\}} \qquad a$$

Stochastic gradient adaptive stepsize (Benveniste et al.)

$$\alpha_n = \left[\alpha^{n-1} + a\varepsilon^n \ \varphi^n \right]_{\alpha^{\min}}^{\alpha^{\max}} \qquad a, \ \alpha^{\max}, \ \alpha^{\min}$$

where:

$$\varphi^n = \left(1 - \alpha^{n-1}\right)\varphi^n + \varepsilon^n$$



■ Bias-adjusted Kalman filter



As σ^2 increases, stepsize decreases As β^n increases, stepsize increases.



Bias-adjusted Kalman filter

» Properties:

$$\alpha_n \rightarrow 1 \text{ as } \sigma^2 \rightarrow 0$$

$$\alpha_n \to 1/n \text{ as } \beta \to 0 \text{ or } \sigma^2 \to \infty$$



Deterministic data: predictions and stepsizes



■ Low noise stochastic data: predictions



1 11 21 31 41 51 61 71 81 91 101 111 121 131 141 151 161 171 181 191 201 211 221 231 241 251 261 271 281 291

© 2007 Warren B. Powell







■ High noise stochastic data: predictions



1 10 19 28 37 46 55 64 73 82 91 100 109 118 127 136 145 154 163 172 181 190 199 208 217 226 235 244 253 262 271 280 289 298

© 2007 Warren B. Powell





■ Bias-adjusted Kalman filter learning rate vs. other



Outline

- The languages of dynamic programming
- A resource allocation model
- The post-decision state variable
- Example: A discrete resource: the nomadic trucker
- The states of our system
- Example: A continuous resource: blood inventory management
- Approximation methods
 - » Lookup tables and aggregation
 - » Basis functions
- Stepsizes
- Exploration vs. exploitation
- Applications



- What decision do we make?
 - » The one we think is best?
 - Exploitation

- » Or do we make a decision just to try something and learn more about the result?
 - Exploration

- Exploration vs. exploitation with the nomadic trucker
 - » Pure exploitation

Information collection



Pure exploitation

© 2007 Warren B. Powell





Resource allocation



Strategies for overcoming the exploitation trap

- » Generalization:
 - Visit one state, learn something about other states
 - Exploitation with more general learning

Information collection



Pure exploration

© 2007 Warren B. Powell

Outline

- The languages of dynamic programming
- A resource allocation model
- The post-decision state variable
- Example: A discrete resource: the nomadic trucker
- The states of our system
- Example: A continuous resource: blood inventory management
- Approximation methods
 - » Lookup tables and aggregation
 - » Basis functions
- Stepsizes
- Exploration vs. exploitation
- Applications


Schneider National



		0522	8
1.0	dr_29812_Sys_6		
1.0	dr_29137_Sys_6		
1.0	dr_29901_Sys_6	1.0 0360918	
1.0	dr_29985_Sys_6	1.0 0320349	
1.0	dr_30156_Sys_6	1.0 0624671	
1.0	dr_30197_Sys_6	1.0 0622613	
0.1	dr_30293_Sys_6	1.0 0102029	
1.0	dr_27387_Sys_6	1.0 0624671	
0,1	dr_27461_Sys_6	1.0 0500451	
1.0	dr_27917_Sys_6	1.0 0504475	
1.0	dr_27970_Sys_6	1.0 0102029	
1.0	dr_28466_Sys_6	1.0 0303311	
0.1	dr_28535_Sys_6	1.0 0303311	
0.1	dr_28875_Sys_6	1.0 0523526	
0.1	dr_29130_Sys_6	1.0 0523526	
0.1	dr_29220_Sys_6	1.0 0442432	
0.1	dr_29383_Sys_6	1.0 0102029	
1.0	dr_34741_Sys_7	1.0 0622613	
1.0	dr_34643_Sys_7	-M	
1.0	dr_34696_Sys_7		















© 2007 Warren B. Powell

Slide 189





© 2007 Warren B. Powell















Implementation metrics

Results from the real world:





The planning process



The flow of information

In practice, there are a number of parallel information processes taking place:



Time









Engineering practice



A car distribution problem



Engineering practice



Engineering practice



Assignments to booked orders.

Repositioning movements based on forecasts

A car distribution problem



A car distribution problem










































Tanker study

