

Robust Optimization and Learning in Approximate Dynamic Programming Using Kernels

Theodore B. Trafalis

School of Industrial and Systems Engineering
University of Oklahoma

Talk

2nd Winter School on Data Analytics (WSDA 2017)
Nizhny Novgorod, Russia
November 03-04, 2017

Decision making is ubiquitous in science and engineering. Dynamic Programming (DP) provides a structured approach, but the **CURSE OF DIMENSIONALITY** makes finding an optimal policy:

- challenging and sometimes even impossible to calculate
- lacking computation efficiency
- requiring off-line calculations
- ignoring uncertainty within the problem parameters

Decision making is ubiquitous in science and engineering. Dynamic Programming (DP) provides a structured approach, but the **CURSE OF DIMENSIONALITY** makes finding an optimal policy:

- challenging and sometimes even impossible to calculate
- lacking computation efficiency
- requiring off-line calculations
- ignoring uncertainty within the problem parameters

Decision making is ubiquitous in science and engineering. Dynamic Programming (DP) provides a structured approach, but the **CURSE OF DIMENSIONALITY** makes finding an optimal policy:

- challenging and sometimes even impossible to calculate
- lacking computation efficiency
- requiring off-line calculations
- ignoring uncertainty within the problem parameters

Decision making is ubiquitous in science and engineering. Dynamic Programming (DP) provides a structured approach, but the **CURSE OF DIMENSIONALITY** makes finding an optimal policy:

- challenging and sometimes even impossible to calculate
- lacking computation efficiency
- requiring off-line calculations
- ignoring uncertainty within the problem parameters

Since similar situations require similar actions we propose to use **CLUSTERING** (grouping). Therefore:

- find optimal and near optimal policies
- make computation more efficient
- search for policies on-line
- Incorporate robustness to get stable policies

Since similar situations require similar actions we propose to use **CLUSTERING** (grouping). Therefore:

- find optimal and near optimal policies
- make computation more efficient
- search for policies on-line
- Incorporate robustness to get stable policies

Since similar situations require similar actions we propose to use **CLUSTERING** (grouping). Therefore:

- find optimal and near optimal policies
- make computation more efficient
- search for policies on-line
- Incorporate robustness to get stable policies

Since similar situations require similar actions we propose to use **CLUSTERING** (grouping). Therefore:

- find optimal and near optimal policies
- make computation more efficient
- search for policies on-line
- Incorporate robustness to get stable policies

- The agent needs to take an action given the actual state and the reward from taking that action.
- The action taken does not need to increase the immediate reward but the expected cumulative reward.
- The rewards are used to evaluate how well the action taken will help achieve the already set goal.

- The agent needs to take an action given the actual state and the reward from taking that action.
- The action taken does not need to increase the immediate reward but the expected cumulative reward.
- The rewards are used to evaluate how well the action taken will help achieve the already set goal.

- The agent needs to take an action given the actual state and the reward from taking that action.
- The action taken does not need to increase the immediate reward but the expected cumulative reward.
- The rewards are used to evaluate how well the action taken will help achieve the already set goal.

An MDP is defined by:

- S , a state space.
- A , an action space.
- \mathcal{P} , a transition probability distribution, where:
$$\mathcal{P}_{ss'}^a = Pr \{s_{t+1} = s' | s_t = s, a_t = a\}.$$
- $g_{ss'}^a$ immediate reward/cost.
- π , the policy that the agent needs to find, is a mapping from S to A .

An MDP is defined by:

- S , a state space.
- A , an action space.
- \mathcal{P} , a transition probability distribution, where:
$$\mathcal{P}_{ss'}^a = Pr \{s_{t+1} = s' | s_t = s, a_t = a\}.$$
- $g_{ss'}^a$ immediate reward/cost.
- π , the policy that the agent needs to find, is a mapping from S to A .

An MDP is defined by:

- S , a state space.
- A , an action space.
- \mathcal{P} , a transition probability distribution, where:
$$\mathcal{P}_{ss'}^a = Pr \{s_{t+1} = s' | s_t = s, a_t = a\}.$$
- $g_{ss'}^a$, immediate reward/cost.
- π , the policy that the agent needs to find, is a mapping from S to A .

An MDP is defined by:

- S , a state space.
- A , an action space.
- \mathcal{P} , a transition probability distribution, where:
$$\mathcal{P}_{ss'}^a = Pr \{s_{t+1} = s' | s_t = s, a_t = a\}.$$
- $g_{ss'}^a$ immediate reward/cost.
- π , the policy that the agent needs to find, is a mapping from S to A .

An MDP is defined by:

- S , a state space.
- A , an action space.
- \mathcal{P} , a transition probability distribution, where:
$$\mathcal{P}_{ss'}^a = Pr \{s_{t+1} = s' | s_t = s, a_t = a\}.$$
- $g_{ss'}^a$ immediate reward/cost.
- π , the policy that the agent needs to find, is a mapping from S to A .

Value Functions

To evaluate how good a particular policy is, the agent needs to evaluate $\mathcal{J}_\pi(\mathbf{s})$ and $Q^\pi(\mathbf{s}, a)$.

- The state-value function:

$$\mathcal{J}_\pi(\mathbf{s}) = E_\pi \left\{ \sum_{k=t}^{\infty} \gamma^{k-t} g_{s_k s_{k+1}}^\pi | s_t = \mathbf{s} \right\}$$

Where ($0 < \gamma < 1$).

- the action-value function:

$$Q_\pi(\mathbf{s}, a) = E_\pi \left\{ \sum_{k=t}^{\infty} \gamma^{k-t} g_{s_k s_{k+1}}^\pi | s_t = \mathbf{s}, a_t = a \right\}$$

$E_\pi \{ \}$ is the expected value achieved following the policy π .

Value Functions

To evaluate how good a particular policy is, the agent needs to evaluate $\mathcal{J}_\pi(\mathbf{s})$ and $Q^\pi(\mathbf{s}, a)$.

- The state-value function:

$$\mathcal{J}_\pi(\mathbf{s}) = E_\pi \left\{ \sum_{k=t}^{\infty} \gamma^{k-t} g_{s_k s_{k+1}}^\pi | s_t = \mathbf{s} \right\}$$

Where ($0 < \gamma < 1$).

- the action-value function:

$$Q_\pi(\mathbf{s}, a) = E_\pi \left\{ \sum_{k=t}^{\infty} \gamma^{k-t} g_{s_k s_{k+1}}^\pi | s_t = \mathbf{s}, a_t = a \right\}$$

$E_\pi \{ \}$ is the expected value achieved following the policy π .

Value Functions -continued-

- The value functions \mathcal{J}_π satisfies a recursive relationship:

$$\mathcal{J}_\pi(\mathbf{s}) = \sum_{a \in A(\mathbf{s})} \pi(\mathbf{s}, a) \sum_{\mathbf{s}'} P_{\mathbf{s}\mathbf{s}'}^a [g_{\mathbf{s}\mathbf{s}'}^a + \gamma \mathcal{J}_\pi(\mathbf{s}')]]$$

- The optimal policy π^* is the one that maximizes \mathcal{J}^π

$$\mathcal{J}^*(\mathbf{s}) = \max_{\pi} \mathcal{J}^\pi(\mathbf{s})$$

- The optimal action-value function Q^* is defined by:

$$Q^*(\mathbf{s}, a) = \max_{\pi} Q^\pi(\mathbf{s}, a)$$

Value Functions -continued-

- The value functions \mathcal{J}_π satisfies a recursive relationship:

$$\mathcal{J}_\pi(\mathbf{s}) = \sum_{a \in A(\mathbf{s})} \pi(\mathbf{s}, a) \sum_{s'} P_{ss'}^a [g_{ss'}^a + \gamma \mathcal{J}_\pi(\mathbf{s}')]]$$

- The optimal policy π^* is the one that maximizes \mathcal{J}^π

$$\mathcal{J}^*(\mathbf{s}) = \max_{\pi} \mathcal{J}^\pi(\mathbf{s})$$

- The optimal action-value function Q^* is defined by:

$$Q^*(\mathbf{s}, a) = \max_{\pi} Q^\pi(\mathbf{s}, a)$$

Value Functions -continued-

- The value functions \mathcal{J}_π satisfies a recursive relationship:

$$\mathcal{J}_\pi(\mathbf{s}) = \sum_{a \in A(\mathbf{s})} \pi(\mathbf{s}, a) \sum_{s'} P_{ss'}^a [g_{ss'}^a + \gamma \mathcal{J}_\pi(\mathbf{s}')]]$$

- The optimal policy π^* is the one that maximizes \mathcal{J}^π

$$\mathcal{J}^*(\mathbf{s}) = \max_{\pi} \mathcal{J}^\pi(\mathbf{s})$$

- The optimal action-value function Q^* is defined by:

$$Q^*(\mathbf{s}, a) = \max_{\pi} Q^\pi(\mathbf{s}, a)$$

Value Functions -continued-

Bellman's equation

The value function \mathcal{J}_π has the following properties:

- \mathcal{J}_π is the unique solution to:

$$\mathcal{J}_\pi = \mathcal{T}_\pi \mathcal{J}_\pi$$

- The operator \mathcal{T}_π is defined by:

$$\mathcal{T}_\pi \mathcal{J} = g_\pi + \gamma \mathcal{P}_\pi \mathcal{J}$$

- The optimal value function $\mathcal{J}^* = \min_{\pi} \mathcal{J}_\pi$ is the unique solution to Bellman's equation:

$$\mathcal{J}^* = \mathcal{T} \mathcal{J}^*$$

- where the operator \mathcal{T} is:

$$\mathcal{T} \mathcal{J} = \min_{\pi} \mathcal{T}_\pi \mathcal{J}$$

Value Functions -continued-

Bellman's equation

The value function \mathcal{J}_π has the following properties:

- \mathcal{J}_π is the unique solution to:

$$\mathcal{J}_\pi = \mathcal{T}_\pi \mathcal{J}_\pi$$

- The operator \mathcal{T}_π is defined by:

$$\mathcal{T}_\pi \mathcal{J} = \mathbf{g}_\pi + \gamma \mathcal{P}_\pi \mathcal{J}$$

- The optimal value function $\mathcal{J}^* = \min_{\pi} \mathcal{J}_\pi$ is the unique solution to Bellman's equation:

$$\mathcal{J}^* = \mathcal{T} \mathcal{J}^*$$

- where the operator \mathcal{T} is:

$$\mathcal{T} \mathcal{J} = \min_{\pi} \mathcal{T}_\pi \mathcal{J}$$

Value Functions -continued-

Bellman's equation

The value function \mathcal{J}_π has the following properties:

- \mathcal{J}_π is the unique solution to:

$$\mathcal{J}_\pi = \mathcal{T}_\pi \mathcal{J}_\pi$$

- The operator \mathcal{T}_π is defined by:

$$\mathcal{T}_\pi \mathcal{J} = \mathbf{g}_\pi + \gamma \mathcal{P}_\pi \mathcal{J}$$

- The optimal value function $\mathcal{J}^* = \min_{\pi} \mathcal{J}_\pi$ is the unique solution to Bellman's equation:

$$\mathcal{J}^* = \mathcal{T} \mathcal{J}^*$$

- where the operator \mathcal{T} is:

$$\mathcal{T} \mathcal{J} = \min_{\pi} \mathcal{T}_\pi \mathcal{J}$$

Value Functions -continued-

Bellman's equation

The value function \mathcal{J}_π has the following properties:

- \mathcal{J}_π is the unique solution to:

$$\mathcal{J}_\pi = \mathcal{T}_\pi \mathcal{J}_\pi$$

- The operator \mathcal{T}_π is defined by:

$$\mathcal{T}_\pi \mathcal{J} = \mathbf{g}_\pi + \gamma \mathcal{P}_\pi \mathcal{J}$$

- The optimal value function $\mathcal{J}^* = \min_{\pi} \mathcal{J}_\pi$ is the unique solution to Bellman's equation:

$$\mathcal{J}^* = \mathcal{T} \mathcal{J}^*$$

- where the operator \mathcal{T} is:

$$\mathcal{T} \mathcal{J} = \min_{\pi} \mathcal{T}_\pi \mathcal{J}$$

Reduced States' Set

Introduction

- This approach aims to decrease the complexity of the problem by reducing the number of states
- Many real life problems have states that are similar or closed to each other
- This similarity yields the same or similar value function
- \Rightarrow calculate the value function just for one of the similar states
- To determine similar states we will use clustering

Reduced States' Set

Introduction

- This approach aims to decrease the complexity of the problem by reducing the number of states
- Many real life problems have states that are similar or closed to each other
- This similarity yields the same or similar value function
- \Rightarrow calculate the value function just for one of the similar states
- To determine similar states we will use clustering

Reduced States' Set

Introduction

- This approach aims to decrease the complexity of the problem by reducing the number of states
- Many real life problems have states that are similar or closed to each other
- This similarity yields the same or similar value function
- ⇒ calculate the value function just for one of the similar states
- To determine similar states we will use clustering

Reduced States' Set

Introduction

- This approach aims to decrease the complexity of the problem by reducing the number of states
- Many real life problems have states that are similar or closed to each other
- This similarity yields the same or similar value function
- \Rightarrow calculate the value function just for one of the similar states
- To determine similar states we will use clustering

Reduced States' Set

Introduction

- This approach aims to decrease the complexity of the problem by reducing the number of states
- Many real life problems have states that are similar or closed to each other
- This similarity yields the same or similar value function
- \Rightarrow calculate the value function just for one of the similar states
- To determine similar states we will use clustering

Reduced States' Set

Notations

- We need to construct a subset \bar{S} using clustering. The value function will be calculated for the elements in \bar{S} only
- $\{S_1, S_2, \dots, S_I\}$ be a partition of S
- \bar{s}_i , an element of S_i , be the **state representative** for the cluster S_i
- $\bar{S} = \{\bar{s}_1, \bar{s}_2, \dots, \bar{s}_I\}$

Reduced States' Set

Notations

- We need to construct a subset \bar{S} using clustering. The value function will be calculated for the elements in \bar{S} only
- $\{S_1, S_2, \dots, S_I\}$ be a partition of S
- \bar{s}_i , an element of S_i , be the **state representative** for the cluster S_i
- $\bar{S} = \{\bar{s}_1, \bar{s}_2, \dots, \bar{s}_I\}$

Reduced States' Set

Notations

- We need to construct a subset \bar{S} using clustering. The value function will be calculated for the elements in \bar{S} only
- $\{S_1, S_2, \dots, S_I\}$ be a partition of S
- \bar{s}_i , an element of S_i , be the **state representative** for the cluster S_i
- $\bar{S} = \{\bar{s}_1, \bar{s}_2, \dots, \bar{s}_I\}$

Reduced States' Set

Notations

- We need to construct a subset \bar{S} using clustering. The value function will be calculated for the elements in \bar{S} only
- $\{S_1, S_2, \dots, S_I\}$ be a partition of S
- \bar{s}_i , an element of S_i , be the **state representative** for the cluster S_i
- $\bar{S} = \{\bar{s}_1, \bar{s}_2, \dots, \bar{s}_I\}$

Reduced States' Set

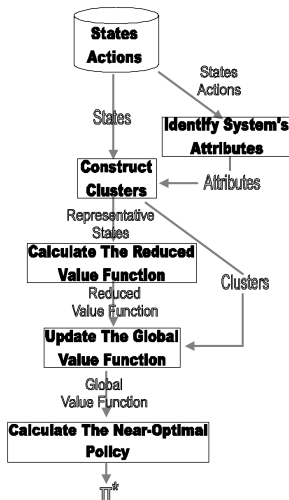


Figure: Reduced States' set diagram

Reduced States' Set

Experiments and results

- Consider an example of a city evacuation plan.
- There are few possible exits from the city.
- The Objective:
Draft a policy that will guide the evacuees to get out of the city.

Reduced States' Set

Experiments and results

- Consider an example of a city evacuation plan.
- There are few possible exits from the city.
- The Objective:
Draft a policy that will guide the evacuees to get out of the city.

Reduced States' Set

Experiments and results

- Consider an example of a city evacuation plan.
- There are few possible exits from the city.
- The Objective:
Draft a policy that will guide the evacuees to get out of the city.

Reduced States' Set

Experiments and results -City representation-

- To simplify the modeling, we suppose that there are 256 blocks (16 columns and 16 rows).
- At each intersection we can go on all 4 directions.
- The only available exits from the city are the north-west and south-east blocks.

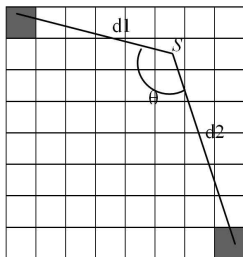


Figure: Simplified city representation with two exits

Reduced States' Set

Experiments and results -City representation-

- To simplify the modeling, we suppose that there are 256 blocks (16 columns and 16 rows).
- At each intersection we can go on all 4 directions.
- The only available exits from the city are the north-west and south-east blocks.

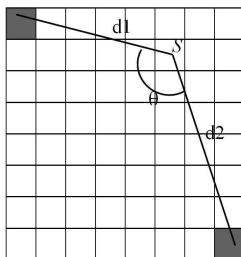


Figure: Simplified city representation with two exits

Reduced States' Set

Experiments and results -City representation-

- To simplify the modeling, we suppose that there are 256 blocks (16 columns and 16 rows).
- At each intersection we can go on all 4 directions.
- The only available exits from the city are the north-west and south-east blocks.

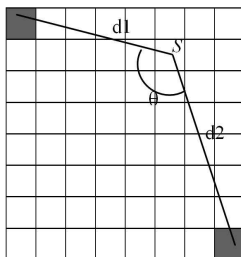


Figure: Simplified city representation with two exits

Reduced States' Set

Experiments and results -Attributes and rewards-

- 254 states and 2 terminal ones.
- The reward to go from one block to another is -1.
- The reward to get to one of the terminal states is 0.
- 3 attributes: d_1 , d_2 , and θ .

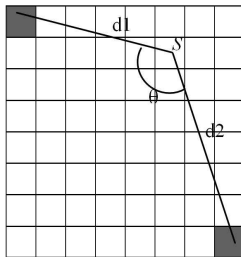


Figure: Attributes

Reduced States' Set

Experiments and results -Attributes and rewards-

- 254 states and 2 terminal ones.
- The reward to go from one block to another is -1.
- The reward to get to one of the terminal states is 0.
- 3 attributes: d_1 , d_2 , and θ .

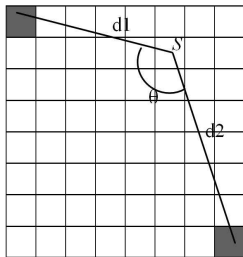


Figure: Attributes

Reduced States' Set

Experiments and results -Attributes and rewards-

- 254 states and 2 terminal ones.
- The reward to go from one block to another is -1.
- The reward to get to one of the terminal states is 0.
- 3 attributes: d_1 , d_2 , and θ .

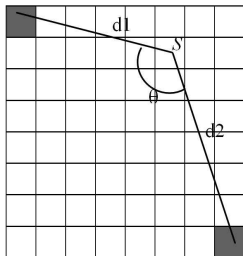


Figure: Attributes

Reduced States' Set

Experiments and results -Attributes and rewards-

- 254 states and 2 terminal ones.
- The reward to go from one block to another is -1.
- The reward to get to one of the terminal states is 0.
- 3 attributes: d_1 , d_2 , and θ .

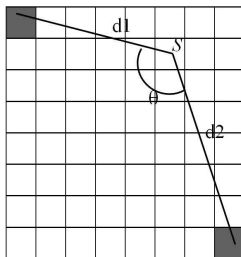


Figure: Attributes

Reduced States' Set

Experiments and results -Results-

- The accuracy rate is introduced to compare the different policies obtained:

$$A_r = \frac{\text{number of accurate states}}{\text{number of states (256)}}$$

	Optimal	Reduced	Reduced
# of clusters	256	73	56
Time (sec)	211	70	33
# of accurate states	256	256	196
A_r	1	1	.77

Table: Comparing Computation Time and Accuracy

Curse of dimensionality

- Consider finite sets with fixed cardinality for the states and actions.
- Need huge amount of resources to find an optimal policy.
- Need off-line algorithms.
- However
In real world problems, the environment might change and include new states or actions.
- A policy that was optimal for a previous environment might not be optimal anymore.
- Use the reduced states' set method in an on-line framework.

Curse of dimensionality

- Consider finite sets with fixed cardinality for the states and actions.
- Need huge amount of resources to find an optimal policy.
- Need off-line algorithms.
- However
In real world problems, the environment might change and include new states or actions.
- A policy that was optimal for a previous environment might not be optimal anymore.
- Use the reduced states' set method in an on-line framework.

Curse of dimensionality

- Consider finite sets with fixed cardinality for the states and actions.
- Need huge amount of resources to find an optimal policy.
- Need off-line algorithms.
- However
 - In real world problems, the environment might change and include new states or actions.
 - A policy that was optimal for a previous environment might not be optimal anymore.
 - Use the reduced states' set method in an on-line framework.

Curse of dimensionality

- Consider finite sets with fixed cardinality for the states and actions.
- Need huge amount of resources to find an optimal policy.
- Need off-line algorithms.
- However
In real world problems, the environment might change and include new states or actions.
- A policy that was optimal for a previous environment might not be optimal anymore.
- Use the reduced states' set method in an on-line framework.

Curse of dimensionality

- Consider finite sets with fixed cardinality for the states and actions.
- Need huge amount of resources to find an optimal policy.
- Need off-line algorithms.
- However
In real world problems, the environment might change and include new states or actions.
- A policy that was optimal for a previous environment might not be optimal anymore.
- Use the reduced states' set method in an on-line framework.

Curse of dimensionality

- Consider finite sets with fixed cardinality for the states and actions.
- Need huge amount of resources to find an optimal policy.
- Need off-line algorithms.
- However
In real world problems, the environment might change and include new states or actions.
- A policy that was optimal for a previous environment might not be optimal anymore.
- Use the reduced states' set method in an on-line framework.

On-line Reinforcement Learning using Kernels

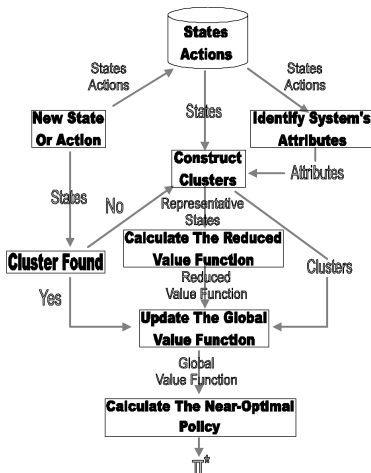


Figure: On-line Reduced States' set diagram

On-line Reinforcement Learning using Kernels

Experiments and results -City representation-

Consider the same city example with the following changes:

- Some blocks in the city might be inaccessible or some of the actions might become impossible to take.
- Due to the congestion, one of the emergency exits might be temporarily disabled.
- The emergency services might succeed in increasing the number of emergency exits.
- To make the problem feasible, we assume that during the whole horizon, there is at least one terminal state.
- The reward to get to an inaccessible block is $-\infty$.

On-line Reinforcement Learning using Kernels

Experiments and results -City representation-

Consider the same city example with the following changes:

- Some blocks in the city might be inaccessible or some of the actions might become impossible to take.
- Due to the congestion, one of the emergency exits might be temporarily disabled.
- The emergency services might succeed in increasing the number of emergency exits.
- To make the problem feasible, we assume that during the whole horizon, there is at least one terminal state.
- The reward to get to an inaccessible block is $-\infty$.

On-line Reinforcement Learning using Kernels

Experiments and results -City representation-

Consider the same city example with the following changes:

- Some blocks in the city might be inaccessible or some of the actions might become impossible to take.
- Due to the congestion, one of the emergency exits might be temporarily disabled.
- The emergency services might succeed in increasing the number of emergency exits.
- To make the problem feasible, we assume that during the whole horizon, there is at least one terminal state.
- The reward to get to an inaccessible block is $-\infty$.

On-line Reinforcement Learning using Kernels

Experiments and results -City representation-

Consider the same city example with the following changes:

- Some blocks in the city might be inaccessible or some of the actions might become impossible to take.
- Due to the congestion, one of the emergency exits might be temporarily disabled.
- The emergency services might succeed in increasing the number of emergency exits.
- To make the problem feasible, we assume that during the whole horizon, there is at least one terminal state.
- The reward to get to an inaccessible block is $-\infty$.

On-line Reinforcement Learning using Kernels

Experiments and results -City representation-

Consider the same city example with the following changes:

- Some blocks in the city might be inaccessible or some of the actions might become impossible to take.
- Due to the congestion, one of the emergency exits might be temporarily disabled.
- The emergency services might succeed in increasing the number of emergency exits.
- To make the problem feasible, we assume that during the whole horizon, there is at least one terminal state.
- The reward to get to an inaccessible block is $-\infty$.

On-line Reinforcement Learning using Kernels

Experiments and results

Events	Optimal	Reduced			Reduced		
	time (sec)	time (sec)	clusters	A_r	time (sec)	clusters	A_r
- 1 block	220	75	74	1	36	57	.769
- 1 block	220	80	75	1	1	57	.762
- 1 exit	350	170	75	1	62	57	.773
- 1 block	345	1	75	1	1	57	.754
+ 1 block	347	1	75	1	1	57	.754
+ 1 exit	219	76	75	1	35	57	.773
+ 1 block	222	72	74	1	33	56	.773
+ 1 block	221	73	74	1	1	56	.766
Total	2144	548	-	-	170	-	-

Table: Comparing Computation Time and Accuracy for On line Reduced states' set

Theoretical Proof

Modified Asynchronous Policy Iteration

- We used a Modified Asynchronous Policy Iteration:
- The value iterations are executed according to the following updates:

$$\mathcal{J}_{k+1}(s) = \begin{cases} (\mathcal{T}_{\pi_k} \mathcal{J}_k)(s), & \text{if } s \in \bar{S}, \\ \mathcal{J}_{k+1}(s_i), & \text{s.t. } s \in S_i, \text{ otherwise,} \end{cases}$$

the policy is kept unchanged by setting $\pi_{k+1} = \pi_k$

- The policy iterations are executed according to the following updates:

$$\pi_{k+1}(s) = \begin{cases} \operatorname{argmin}_{a \in A(s)} \sum_{s'=0}^{|S|} \mathcal{P}_{ss'}^a (g_{ss'}^a + \mathcal{J}_k(s')), & \text{if } s \in \bar{S}, \\ \pi_{k+1}(s_i), & \text{s.t. } s \in S_i, \text{ otherwise,} \end{cases}$$

the value function estimates are kept unchanged by setting $\mathcal{J}_{k+1} = \mathcal{J}_k$

Theoretical Proof

Modified Asynchronous Policy Iteration

- We used a Modified Asynchronous Policy Iteration:
- The value iterations are executed according to the following updates:

$$\mathcal{J}_{k+1}(s) = \begin{cases} (\mathcal{T}_{\pi_k} \mathcal{J}_k)(s), & \text{if } s \in \bar{S}, \\ \mathcal{J}_{k+1}(s_i), & \text{s.t. } s \in S_i, \text{ otherwise,} \end{cases}$$

the policy is kept unchanged by setting $\pi_{k+1} = \pi_k$

- The policy iterations are executed according to the following updates:

$$\pi_{k+1}(s) = \begin{cases} \operatorname{argmin}_{a \in A(s)} \sum_{s'=0}^{|S|} \mathcal{P}_{ss'}^a (g_{ss'}^a + \mathcal{J}_k(s')), & \text{if } s \in \bar{S}, \\ \pi_{k+1}(s_i), & \text{s.t. } s \in S_i, \text{ otherwise,} \end{cases}$$

the value function estimates are kept unchanged by setting $\mathcal{J}_{k+1} = \mathcal{J}_k$

Theoretical Proof

Modified Asynchronous Policy Iteration

- We used a Modified Asynchronous Policy Iteration:
- The value iterations are executed according to the following updates:

$$\mathcal{J}_{k+1}(\mathbf{s}) = \begin{cases} (\mathcal{T}_{\pi_k} \mathcal{J}_k)(\mathbf{s}), & \text{if } \mathbf{s} \in \bar{\mathcal{S}}, \\ \mathcal{J}_{k+1}(\mathbf{s}_i), & \text{s.t. } \mathbf{s} \in \mathcal{S}_i, \text{ otherwise,} \end{cases}$$

the policy is kept unchanged by setting $\pi_{k+1} = \pi_k$

- The policy iterations are executed according to the following updates:

$$\pi_{k+1}(\mathbf{s}) = \begin{cases} \operatorname{argmin}_{a \in A(\mathbf{s})} \sum_{s'=0}^{|\mathcal{S}|} \mathcal{P}_{ss'}^a (g_{ss'}^a + \mathcal{J}_k(s')), & \text{if } \mathbf{s} \in \bar{\mathcal{S}}, \\ \pi_{k+1}(\mathbf{s}_i), & \text{s.t. } \mathbf{s} \in \mathcal{S}_i, \text{ otherwise,} \end{cases}$$

the value function estimates are kept unchanged by setting $\mathcal{J}_{k+1} = \mathcal{J}_k$

Theoretical Proof

Theorem

- Let $s \in \bar{S}_i$ and $s \notin \bar{S}$, let $s' \in \bar{S}_j$.
- Let $\mathcal{P}_{ss'}^{\pi_k(s)} > 0$ and $s'' \in \bar{S}_j$ such that:
 1. $\mathcal{P}_{\bar{s}_i s''}^{\pi_k(\bar{s}_i)} > 0$ and $\left| \mathcal{P}_{ss'}^{\pi_k(s)} - \mathcal{P}_{\bar{s}_i s''}^{\pi_k(\bar{s}_i)} \right| \leq \epsilon_1$
 2. $\left| g_{ss'}^{\pi_k(s)} - g_{\bar{s}_i s''}^{\pi_k(\bar{s}_i)} \right| \leq \epsilon_2$
- then $\exists \epsilon \geq 0$ such that:
$$|(\mathcal{T}_{\pi_k} \mathcal{J}_k)(s) - \mathcal{J}_{k+1}(s)| \leq \epsilon \forall s \notin \bar{S}.$$

Theoretical Proof

Theorem

- Let $s \in \bar{S}_i$ and $s \notin \bar{S}$, let $s' \in \bar{S}_j$.
- Let $\mathcal{P}_{ss'}^{\pi_k(s)} > 0$ and $s'' \in \bar{S}_j$ such that:
 1. $\mathcal{P}_{\bar{s}_i s''}^{\pi_k(\bar{s}_i)} > 0$ and $\left| \mathcal{P}_{ss'}^{\pi_k(s)} - \mathcal{P}_{\bar{s}_i s''}^{\pi_k(\bar{s}_i)} \right| \leq \epsilon_1$
 2. $\left| g_{ss'}^{\pi_k(s)} - g_{\bar{s}_i s''}^{\pi_k(\bar{s}_i)} \right| \leq \epsilon_2$
- then $\exists \epsilon \geq 0$ such that:
 $|\mathcal{T}_{\pi_k} \mathcal{J}_k(s) - \mathcal{J}_{k+1}(s)| \leq \epsilon \forall s \notin \bar{S}$.

Theoretical Proof

Theorem

- Let $s \in \bar{S}_i$ and $s \notin \bar{S}$, let $s' \in \bar{S}_j$.
- Let $\mathcal{P}_{ss'}^{\pi_k(s)} > 0$ and $s'' \in \bar{S}_j$ such that:
 1. $\mathcal{P}_{\bar{s}_i s''}^{\pi_k(\bar{s}_i)} > 0$ and $\left| \mathcal{P}_{ss'}^{\pi_k(s)} - \mathcal{P}_{\bar{s}_i s''}^{\pi_k(\bar{s}_i)} \right| \leq \epsilon_1$
 2. $\left| g_{ss'}^{\pi_k(s)} - g_{\bar{s}_i s''}^{\pi_k(\bar{s}_i)} \right| \leq \epsilon_2$
- then $\exists \epsilon \geq 0$ such that:
$$|(\mathcal{T}_{\pi_k} \mathcal{J}_k)(s) - \mathcal{J}_{k+1}(s)| \leq \epsilon \forall s \notin \bar{S}.$$

- Conduct a set of experiments on the class of shortest path problems.
- Compare the time it takes to find an optimal policy using our algorithm and without reducing the states' set cardinality.
- Define the following two percentages:

$$A = \frac{\textit{Clustering Time} + \textit{Reduced Time}}{\textit{Optimal Time}} \times 100,$$

$$B = \frac{\textit{Reduced Time}}{\textit{Optimal Time}} \times 100.$$

- Conduct a set of experiments on the class of shortest path problems.
- Compare the time it takes to find an optimal policy using our algorithm and without reducing the states' set cardinality.
- Define the following two percentages:

$$A = \frac{\textit{Clustering Time} + \textit{Reduced Time}}{\textit{Optimal Time}} \times 100,$$

$$B = \frac{\textit{Reduced Time}}{\textit{Optimal Time}} \times 100.$$

- Conduct a set of experiments on the class of shortest path problems.
- Compare the time it takes to find an optimal policy using our algorithm and without reducing the states' set cardinality.
- Define the following two percentages:

$$A = \frac{\textit{Clustering Time} + \textit{Reduced Time}}{\textit{Optimal Time}} \times 100,$$

$$B = \frac{\textit{Reduced Time}}{\textit{Optimal Time}} \times 100.$$

Theoretical Proof

Experiments and Results

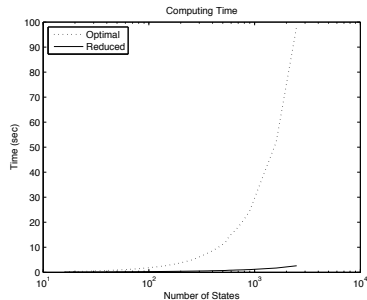
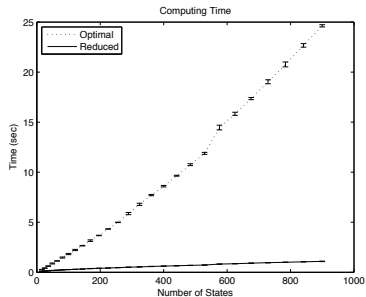


Figure: Optimal Vs Reduced Computation Times

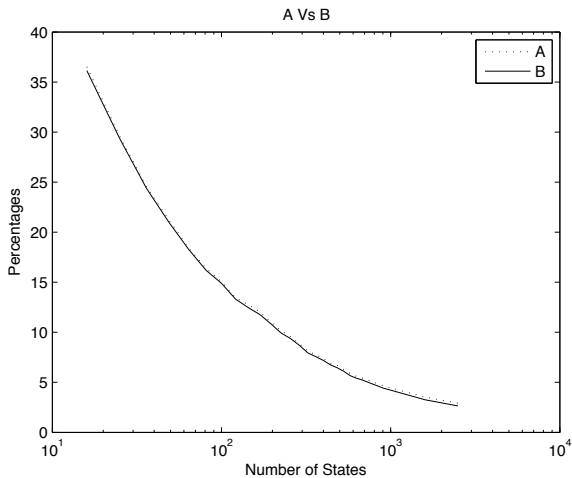


Figure: The Impact of the Clustering Time

Kernel Based Reinforcement Learning

Kernel Methods

- Kernel methods play a major role in Machine Learning.
- They provide a simple framework for manipulating nonlinear relationships.
- Instantaneous adaptation of former linear algorithms.
- Require modest computational resources.

Kernel Based Reinforcement Learning

The Mercer Kernels

- A kernel is a continuous symmetric real-valued function defined on compact subsets of \mathbb{R}^n , $k : (x, y) \mapsto k(x, y)$.
- A Mercer kernel is a nonnegative definite kernel.
- The domain of a Mercer kernel is called the input space.
- The quantity $k(x, y)$ can be used to represent measures of angle and measures of distance.
- Angles and distances are between inputs mapped in a higher dimensional Hilbert space.
- The Hilbert space is called the feature space.

Kernel Based Reinforcement Learning

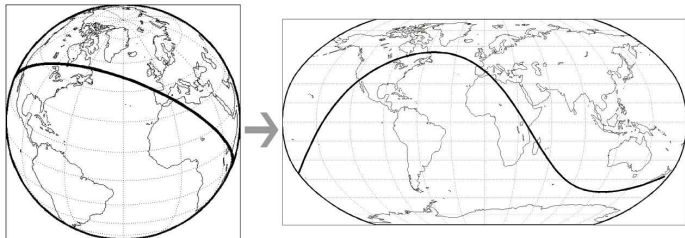
The Kernel Trick

- Mercer's theorem suggests a particular decomposition of Mercer kernels.
- A Mercer kernel can be expressed as a dot product between two inputs mapped in the feature space, $k(x, y) = \langle \phi(x) \cdot \phi(y) \rangle$.
- Explicit knowledge of the map ϕ and the feature space is not required. The only thing of importance is the kernel itself.
- With k we can calculate angles and distances between elements in the feature space without knowing the map ϕ .

Kernel Based Reinforcement Learning

Purpose of Kernel Methods

- Kernel methods simplify the representation of nonlinear patterns in the input space.
- The intersection between hyperplanes and the manifold has a non-trivial reciprocal image in the input space.
- Instead of searching complex patterns in the input space, we use kernel methods.



- Use Linear approximate dynamic programming to speed up the algorithm.
- Incorporate robustness to get stable policies.
- Use kernels to define the concept of the neighborhood of a state.
- Combine all the above techniques to generate an on-line robust RL algorithm.

RORLK: Algorithm Diagram

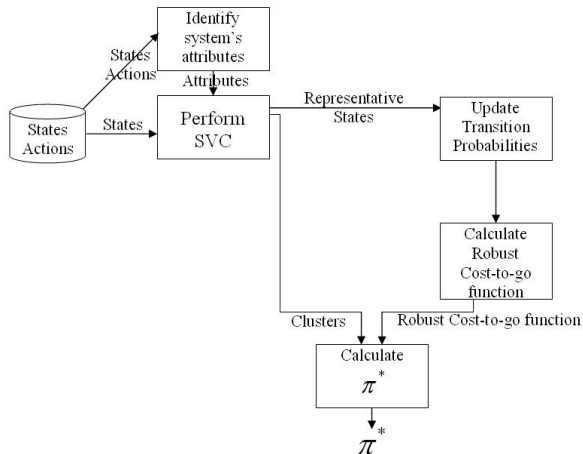


Figure: Robust On-line RL using Kernels diagram

RORLK: Algorithm Diagram

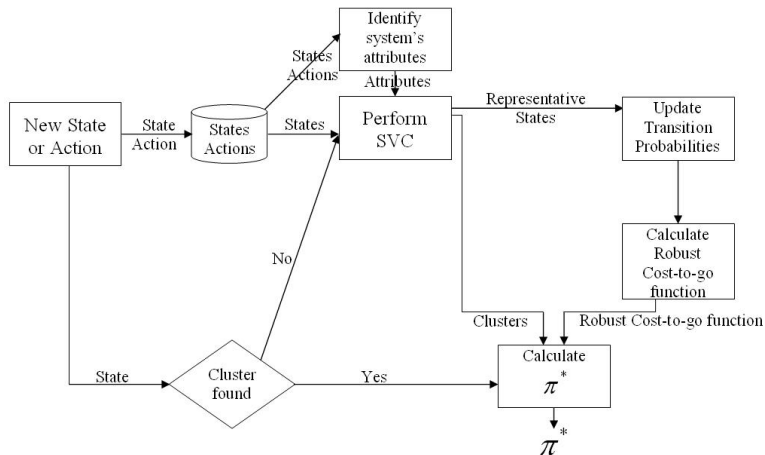


Figure: Robust On-line RL using Kernels diagram

Reduced States' Set

Support Vector Clustering -Introduction-

- This approach aims to decrease the complexity of the problem by reducing the number of states.
- Many real life problems have states that are similar or closed to each other.
- This similarity yields the same or similar value function.
- \Rightarrow calculate the value function just for one of the similar states.
- To determine similar states we will use clustering.

Reduced States' Set

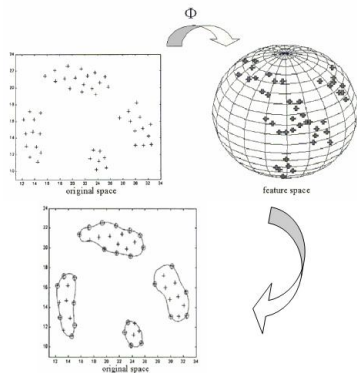
Support Vector Clustering -Introduction-

- Clustering is a part of data mining that consists of grouping a set of data according to various attributes.
- Many methods have been developed.
- The effectiveness of each approach depends on the nature of the data.
- Ben-Hur (2001) derived SVC from SVM.

Reduced States' Set

Support Vector Clustering -Main steps-

- Map data points to a feature space.
- Search for the smallest sphere enclosing all data points.
- Map the sphere back and generate clusters.



Reduced States' Set

Support Vector Clustering -Formulation-

- $\{X_i\}_{1 \leq i \leq n}$ a data set of n points in the input space.
- $X_i \in \mathbb{R}^d$ and d is the number of attributes.
- We need to find the smallest sphere with radius R .
- The optimization formulation is the following:

$$\min_{R,a,\xi_i} \quad R^2 + C \sum_i \xi_i$$

Subject to:

$$\begin{aligned} \|\Phi(X_i) - a\|^2 &\leq R^2 + \xi_i, \quad i = 1 \dots n \\ \xi_i &\geq 0 \end{aligned}$$

Reduced States' Set

Support Vector Clustering -Dual problem-

- To solve this optimization problem, we write the Lagrangian:

$$L(R, a, \mu, \beta, \xi) = R^2 - \sum_i \left(R^2 + \xi_i - \|\Phi(X_i) - a\|^2 \right) \beta_i - \sum_i \xi_i \mu_i + C \sum_i \xi_i$$

Where $\beta_i \geq 0$ and $\mu_i \geq 0$ are the Lagrange multipliers.

- Using the Karush-Kuhn-Tucker complementary slackness conditions, the Lagrangian becomes:

$$L(\beta) = \sum_i \Phi(X_i)^2 \beta_i - \sum_{ij} \beta_i \beta_j \Phi(X_i) \cdot \Phi(X_j)$$

With the constraints $0 \leq \beta_i \leq C$.

Reduced States' Set

Support Vector Clustering -Dual problem-

- Replacing the dot product with the Gaussian kernel function:

$$K(X_i, X_j) = e^{-q\|X_i - X_j\|^2}$$

Where q is the width parameter.

- we get:

$$\max_{\beta} \quad \sum_i K(X_i, X_i)\beta_i - \sum_{ij} \beta_i\beta_j K(X_i, X_j)$$

Subject to:

$$\begin{aligned} 0 \leq \beta_i \leq C, \quad i = 1 \dots n \\ \sum_i \beta_i = 1 \end{aligned}$$

Reduced States' Set

Support Vector Clustering -Distance from the center and Radius-

- The distance of each image in the feature space from the sphere center a :

$$\begin{aligned}R^2(X) &= \|\Phi(X_i) - a\|^2 \\ &= K(X, X) - 2 \sum_i \beta_i K(X_i, X) + \sum_{ij} \beta_i \beta_j K(X_i, X_j)\end{aligned}$$

- The radius of the sphere is as follows:

$$R_{\text{Sphere}} = \frac{\sum_{X_i \text{ is SV}} R(X_i)}{\text{Number of SVs}}$$

- clusters are defined as the connected components of the graph induced by A , where A is the adjacency matrix defined by:

$$A_{ij} = \begin{cases} 1, & \text{if for all } y \in (X_i, X_j), R(y) \leq R_{sphere} \\ 0, & \text{Otherwise} \end{cases}$$

Reduced States' Set

Reduced States' set Diagram

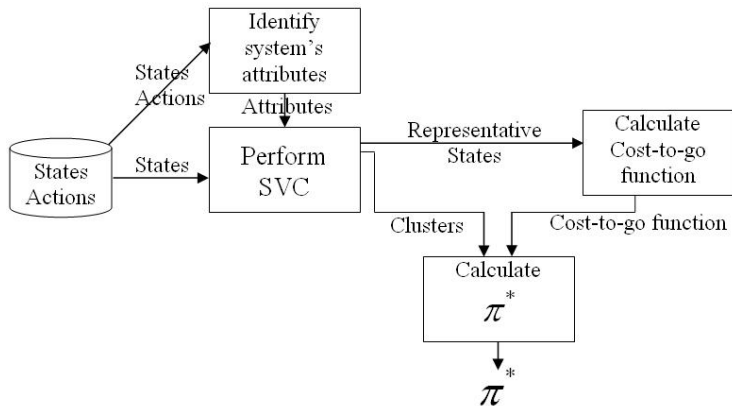


Figure: Reduced States' set diagram

Reduced Linear Dynamic Programming

Linear Dynamic Programming

Bellman's equation can be solved using the following Linear Dynamic Programming (LDP):

$$\max c^T \mathcal{J}$$

$$\text{subject to : } g_s^a + \gamma \sum_{s' \in \mathcal{S}} \mathcal{P}_{ss'}^a \mathcal{J}(s') \geq \mathcal{J}(s), \forall s \in \mathcal{S}, \forall a \in A(s)$$

Reduced Linear Dynamic Programming

Reduced Linear Dynamic Programming

Using the partitioning $S = \{S_1, S_2, \dots, S_l\}$ we can reduced the LDP to:

$$\max \sum_{i=1}^l \bar{c}_i \mathcal{J}(\bar{s}_i)$$

$$\text{subject to : } g_{\bar{s}_j}^a + \gamma \sum_{i=1}^l |S_i| \mathcal{P}_{S\bar{s}_i}^a \mathcal{J}(\bar{s}_i) \geq \mathcal{J}(\bar{s}_j), \forall j = 1 \dots l, \forall a \in A(\bar{s}_j)$$

$$\text{where } \bar{c}_i = \sum_{s \in S_i} c_s$$

Reduced Linear Approximate Dynamic Programming

Linear Approximate Dynamic Programming

\mathcal{J}^* is approximated using a linear combination of preselected basis functions:

$$\phi_k : \mathcal{S} \mapsto \mathbb{R}, k = 1, \dots, K$$

\Rightarrow Generate a weight vector $\tilde{r} \in \mathbb{R}^K$, such that:

$$\mathcal{J}^*(s) \approx \sum_{k=0}^K \phi_k(s) \tilde{r}$$

Using the vector \tilde{r} we get the Linear Approximate DP (LADP):

$$\max_r C^T \Phi r$$

subject to

$$g_{ss''}^a + \gamma \sum_{s' \in \mathcal{S}} \mathcal{P}_{ss'}^a(\Phi r)(s') \geq (\Phi r)(s), \forall s \in \mathcal{S}, \forall a \in A(s)$$

where $\Phi = \begin{bmatrix} \phi_1(\mathbf{1}) & \cdots & \phi_K(\mathbf{1}) \\ \vdots & & \vdots \\ \phi_1(|\mathcal{S}|) & \cdots & \phi_K(|\mathcal{S}|) \end{bmatrix}$

Reduced Linear Approximate Dynamic Programming

Reduced Linear Approximate Dynamic Programming

Using the partitioning $S = \{S_1, S_2, \dots, S_l\}$ we can reduced the

LADP to:

$$\max \bar{c}^T \bar{\Phi} r$$

$$\text{subject to } g_{\bar{s}_j}^a + \gamma \sum_{i=1}^l |S_i| \mathcal{P}_{\bar{s}_j \bar{s}_i}^a (\bar{\Phi} r) (\bar{s}_i) \geq (\bar{\Phi} r) (\bar{s}_j), \forall j = 1 \dots l, \forall a \in A(\bar{s}_j),$$

$$\text{where } \bar{\Phi} = \begin{bmatrix} \phi_1(\bar{s}_1) & \cdots & \phi_K(\bar{s}_1) \\ \vdots & & \vdots \\ \phi_1(\bar{s}_l) & \cdots & \phi_K(\bar{s}_l) \end{bmatrix}.$$

- The LADP yields good approximation to the value function
- How to choose the basis function and K
- We propose to use the partitioning $\mathcal{S} = \{S_1, S_2, \dots, S_l\}$ and set:

$$\phi_k(\mathbf{s}) = \kappa(\mathbf{s}, \bar{\mathbf{s}}_k), \quad \forall k = 1 \dots l$$

- The approximation will be done using the system's parameters
- $K = l$ is the number of clusters

- The LADP yields good approximation to the value function
- How to choose the basis function and K
- We propose to use the partitioning $\mathcal{S} = \{S_1, S_2, \dots, S_l\}$ and set:

$$\phi_k(s) = \kappa(s, \bar{s}_k), \forall k = 1 \dots l$$

- The approximation will be done using the system's parameters
- $K = l$ is the number of clusters

- The LADP yields good approximation to the value function
- How to choose the basis function and K
- We propose to use the partitioning $\mathcal{S} = \{\mathcal{S}_1, \mathcal{S}_2, \dots, \mathcal{S}_l\}$ and set:

$$\phi_k(\mathbf{s}) = \kappa(\mathbf{s}, \bar{\mathbf{s}}_k), \quad \forall k = 1 \dots l$$

- The approximation will be done using the system's parameters
- $K = l$ is the number of clusters

- The LADP yields good approximation to the value function
- How to choose the basis function and K
- We propose to use the partitioning $\mathcal{S} = \{\mathcal{S}_1, \mathcal{S}_2, \dots, \mathcal{S}_l\}$ and set:

$$\phi_k(\mathbf{s}) = \kappa(\mathbf{s}, \bar{\mathbf{s}}_k), \quad \forall k = 1 \dots l$$

- The approximation will be done using the system's parameters
- $K = l$ is the number of clusters

- The LADP yields good approximation to the value function
- How to choose the basis function and K
- We propose to use the partitioning $\mathcal{S} = \{\mathcal{S}_1, \mathcal{S}_2, \dots, \mathcal{S}_l\}$ and set:

$$\phi_k(\mathbf{s}) = \kappa(\mathbf{s}, \bar{\mathbf{s}}_k), \quad \forall k = 1 \dots l$$

- The approximation will be done using the system's parameters
- $K = l$ is the number of clusters

The LADP becomes Kernelized LADP:

$$\max c^T \mathcal{K}r$$

$$\text{subject to } g_s^a + \gamma \sum_{s' \in \mathcal{S}} \mathcal{P}_{ss'}^a(\mathcal{K}r)(s') \geq (\mathcal{K}r)(s), \forall s \in \mathcal{S}, \forall a \in A(s),$$

where

$$\mathcal{K} = \begin{bmatrix} \kappa(\mathbf{s}_1, \bar{\mathbf{s}}_1) & \cdots & \kappa(\mathbf{s}_1, \bar{\mathbf{s}}_l) \\ \vdots & & \vdots \\ \kappa(\mathbf{s}_{|\mathcal{S}|}, \bar{\mathbf{s}}_1) & \cdots & \kappa(\mathbf{s}_{|\mathcal{S}|}, \bar{\mathbf{s}}_l) \end{bmatrix}$$

The RLADP becomes Kernelized RLADP:

$$\max \bar{c}^T \bar{\mathcal{K}} r$$

$$\text{subject to } g_{\bar{s}_j}^a + \gamma \sum_{i=1}^l |\mathcal{S}_i| \mathcal{P}_{\bar{s}_j \bar{s}_i}^a (\bar{\mathcal{K}} r) (\bar{s}_i) \geq (\bar{\mathcal{K}} r) (\bar{s}_j), \forall j = 1 \cdots l, \forall a \in \mathbf{A}(\bar{s}_j)$$

where

$$\bar{\mathcal{K}} = \begin{bmatrix} \kappa(\bar{\mathbf{s}}_1, \bar{\mathbf{s}}_1) & \cdots & \kappa(\bar{\mathbf{s}}_1, \bar{\mathbf{s}}_l) \\ \vdots & & \vdots \\ \kappa(\bar{\mathbf{s}}_l, \bar{\mathbf{s}}_1) & \cdots & \kappa(\bar{\mathbf{s}}_l, \bar{\mathbf{s}}_l) \end{bmatrix}$$

Experiments and Results

- We used the same shortest path problems used in section
- The reduce linear dynamic programming did not find the optimal policies all the time
- We introduce the reduction in this experiments:

$$R_r = \frac{|S| - |\bar{S}|}{|S|}$$

Experiments and Results

- We used the same shortest path problems used in section
- The reduce linear dynamic programming did not find the optimal policies all the time
- We introduce the reduction in this experiments:

$$R_r = \frac{|S| - |\bar{S}|}{|S|}$$

Experiments and Results

- We used the same shortest path problems used in section
- The reduce linear dynamic programming did not find the optimal policies all the time
- We introduce the reduction in this experiments:

$$R_r = \frac{|S| - |\bar{S}|}{|S|}$$

Experiments and Results

Computation Time

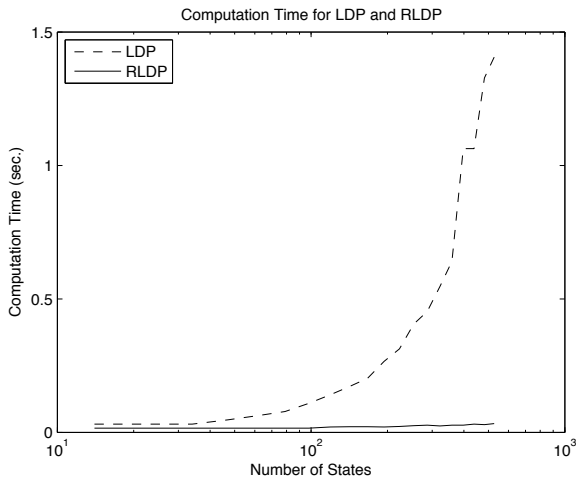


Figure: Computation Time -LDP Vs RLDP-

Experiments and Results

Accuracy and Reduction Rates

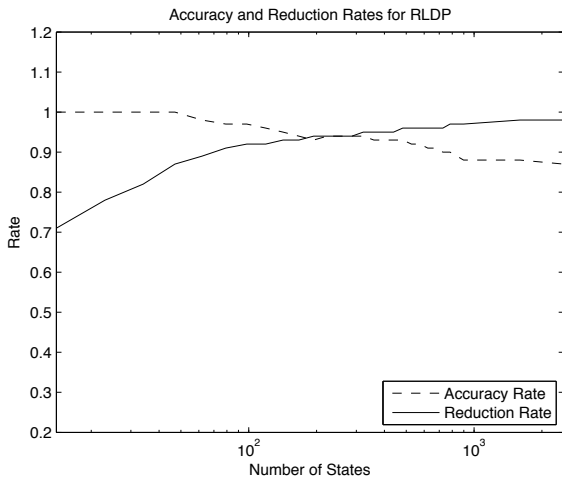


Figure: Accuracy and Reduction Rates -LDP Vs RLDP-

Robust Linear Dynamic Programming

Introduction

- Most developed algorithms ignore the uncertainty within the transition probability matrices.
- **However**, the DP optimal solution is sensitive to perturbation in transition matrices and in the cost/reward.
- **Therefore**, we need to account for the uncertainty.
- Different approaches have been developed.

Robust Linear Dynamic Programming

Introduction

- Most developed algorithms ignore the uncertainty within the transition probability matrices.
- **However**, the DP optimal solution is sensitive to perturbation in transition matrices and in the cost/reward.
- **Therefore**, we need to account for the uncertainty.
- Different approaches have been developed.

Robust Linear Dynamic Programming

Introduction

- Most developed algorithms ignore the uncertainty within the transition probability matrices.
- **However**, the DP optimal solution is sensitive to perturbation in transition matrices and in the cost/reward.
- **Therefore**, we need to account for the uncertainty.
- Different approaches have been developed.

Robust Linear Dynamic Programming

Introduction

- Most developed algorithms ignore the uncertainty within the transition probability matrices.
- **However**, the DP optimal solution is sensitive to perturbation in transition matrices and in the cost/reward.
- **Therefore**, we need to account for the uncertainty.
- Different approaches have been developed.

Robust Linear Dynamic Programming

Introduction

- Uncertainty within the transition probability matrices has been studied extensively
- But not the uncertainty within the reward/cost values
- We have augmented the LDP to make the policies robust with respect to the reward/cost values
- We have used Bertsimas and Sim's robust formulation to mitigate the effect of uncertainty within the transition probability matrices

Robust Linear Dynamic Programming

Introduction

- Uncertainty within the transition probability matrices has been studied extensively
- But not the uncertainty within the reward/cost values
- We have augmented the LDP to make the policies robust with respect to the reward/cost values
- We have used Bertsimas and Sim's robust formulation to mitigate the effect of uncertainty within the transition probability matrices

Robust Linear Dynamic Programming

Introduction

- Uncertainty within the transition probability matrices has been studied extensively
- But not the uncertainty within the reward/cost values
- We have augmented the LDP to make the policies robust with respect to the reward/cost values
- We have used Bertsimas and Sim's robust formulation to mitigate the effect of uncertainty within the transition probability matrices

Robust Linear Dynamic Programming

Introduction

- Uncertainty within the transition probability matrices has been studied extensively
- But not the uncertainty within the reward/cost values
- We have augmented the LDP to make the policies robust with respect to the reward/cost values
- We have used Bertsimas and Sim's robust formulation to mitigate the effect of uncertainty within the transition probability matrices

Consider the following linear optimization problem:

$$\max c^T x$$

$$\text{subject to : } Ax \leq b$$

- Each element b_i of b can take a finite number of values $\{b_i^1, \dots, b_i^q\}$
- In practice, a convex combination of the $\{b_i^1, \dots, b_i^q\}$ is used as the expected value for b_i
- We assume that $\text{probability}(b_i = b_i^w) = \frac{1}{q}, \forall w = 1 \dots q$

Consider the following linear optimization problem:

$$\max c^T x$$

$$\text{subject to : } Ax \leq b$$

- Each element b_i of b can take a finite number of values $\{b_i^1, \dots, b_i^q\}$
- In practice, a convex combination of the $\{b_i^1, \dots, b_i^q\}$ is used as the expected value for b_i
- We assume that $\text{probability}(b_i = b_i^w) = \frac{1}{q}, \forall w = 1 \dots q$

Consider the following linear optimization problem:

$$\max c^T x$$

$$\text{subject to : } Ax \leq b$$

- Each element b_i of b can take a finite number of values $\{b_i^1, \dots, b_i^q\}$
- In practice, a convex combination of the $\{b_i^1, \dots, b_i^q\}$ is used as the expected value for b_i
- We assume that $\text{probability}(b_i = b_i^w) = \frac{1}{q}, \forall w = 1 \dots q$

Robust Linear Dynamic Programming

Augmented Linear Dynamic Programming

- For each realization b_i^w we associate the decision variable x_i^w .
- Hence, $E \{ \max c^T x^w | Ax^w \leq b^w \}$ is equivalent to:

$$\max \sum_i \frac{c_i}{q} \sum_{v=1}^q x_i^v$$

$$\text{subject to: } \sum_{ij, j \neq i} \frac{a_{ij}}{q} \sum_{v=1}^q x_j^v + a_{ii} x_i^w \leq b_i^w, \forall i, w$$

- The linear dynamic formulation becomes:

$$\max \sum_s \frac{c_s}{q} \sum_{v=1}^q \mathcal{J}(s^v)$$

$$\text{subject to: } g_{s^w}^a + \gamma \sum_{s' \in S, s' \neq s^w} \frac{\mathcal{P}_{s^w s'}^a}{q} \sum_{v=1}^q \mathcal{J}(s'^v) \geq (1 - \gamma \mathcal{P}_{s^w s^w}^a) \mathcal{J}(s^w), \forall s^w, a, w$$

Robust Linear Dynamic Programming

Augmented Linear Dynamic Programming

- For each realization b_i^w we associate the decision variable x_i^w .
- Hence, $E \{ \max c^T x^w | Ax^w \leq b^w \}$ is equivalent to:

$$\max \sum_i \frac{c_i}{q} \sum_{v=1}^q x_i^v$$

$$\text{subject to : } \sum_{ij, j \neq i} \frac{a_{ij}}{q} \sum_{v=1}^q x_j^v + a_{ii} x_i^w \leq b_i^w, \forall i, w$$

- The linear dynamic formulation becomes:

$$\max \sum_s \frac{c_s}{q} \sum_{v=1}^q \mathcal{J}(s^v)$$

$$\text{subject to : } g_{s^w}^a + \gamma \sum_{s' \in S, s' \neq s^w} \frac{\mathcal{P}_{s^w s'}^a}{q} \sum_{v=1}^q \mathcal{J}(s'^v) \geq (1 - \gamma \mathcal{P}_{s^w s^w}^a) \mathcal{J}(s^w), \forall s^w, a, w$$

Robust Linear Dynamic Programming

Augmented Linear Dynamic Programming

- For each realization b_i^w we associate the decision variable x_i^w .
- Hence, $E \{ \max c^T x^w | Ax^w \leq b^w \}$ is equivalent to:

$$\max \sum_i \frac{c_i}{q} \sum_{v=1}^q x_i^v$$

$$\text{subject to : } \sum_{ij, j \neq i} \frac{a_{ij}}{q} \sum_{v=1}^q x_j^v + a_{ii} x_i^w \leq b_i^w, \forall i, w$$

- The linear dynamic formulation becomes:

$$\max \sum_s \frac{c_s}{q} \sum_{v=1}^q \mathcal{J}(s^v)$$

$$\text{subject to : } g_{s^w}^a + \gamma \sum_{s' \in \mathcal{S}, s' \neq s^w} \frac{\mathcal{P}_{s^w s'}^a}{q} \sum_{v=1}^q \mathcal{J}(s'^v) \geq (1 - \gamma \mathcal{P}_{s^w s^w}^a) \mathcal{J}(s^w), \forall s^w, a, w$$

Robust Linear Dynamic Programming

- Use Bertsimas and Sim's formulation to account for uncertainty within the transition probability matrices
- We assume that each entry $\mathcal{P}_{ss'}^a$ is modeled as a symmetric and bounded random variable $\tilde{\mathcal{P}}_{ss'}^a$, such that:

$$\tilde{\mathcal{P}}_{ss'}^a \in \left[\mathcal{P}_{ss'}^a - \hat{\mathcal{P}}_{ss'}^a, \mathcal{P}_{ss'}^a + \hat{\mathcal{P}}_{ij}^a \right]$$

- We set $\alpha_{s^w s'^w}^a = (1 - \delta_{s^w s'^w}) \frac{\delta_{s^w s'^w} - \gamma \mathcal{P}_{s^w s'^w}^a}{q + (1-q)\delta_{s^w s'^w}}$ and

$$\hat{\alpha}_{s^w s'^w}^a = (1 - \delta_{s^w s'^w}) \frac{\gamma \hat{\mathcal{P}}_{s^w s'^w}^a}{q + (1-q)\delta_{s^w s'^w}}$$

Robust Linear Dynamic Programming

- Use Bertsimas and Sim's formulation to account for uncertainty within the transition probability matrices
- We assume that each entry $\mathcal{P}_{ss'}^a$ is modeled as a symmetric and bounded random variable $\tilde{\mathcal{P}}_{ss'}^a$, such that:

$$\tilde{\mathcal{P}}_{ss'}^a \in \left[\mathcal{P}_{ss'}^a - \hat{\mathcal{P}}_{ss'}^a, \mathcal{P}_{ss'}^a + \hat{\mathcal{P}}_{ij}^a \right]$$

- We set $\alpha_{sws'w'}^a = (1 - \delta_{sws'w'}) \frac{\delta_{sws'w'} - \gamma \mathcal{P}_{sws'w'}^a}{q + (1-q)\delta_{sws'w'}}$ and

$$\hat{\alpha}_{sws'w'}^a = (1 - \delta_{sws'w'}) \frac{\gamma \hat{\mathcal{P}}_{sws'w'}^a}{q + (1-q)\delta_{sws'w'}}$$

Robust Linear Dynamic Programming

- Use Bertsimas and Sim's formulation to account for uncertainty within the transition probability matrices
- We assume that each entry $\mathcal{P}_{ss'}^a$ is modeled as a symmetric and bounded random variable $\tilde{\mathcal{P}}_{ss'}^a$, such that:

$$\tilde{\mathcal{P}}_{ss'}^a \in \left[\mathcal{P}_{ss'}^a - \hat{\mathcal{P}}_{ss'}^a, \mathcal{P}_{ss'}^a + \hat{\mathcal{P}}_{ij}^a \right]$$

- We set $\alpha_{s^w s'^w}^a = (1 - \delta_{s^w s'^w}) \frac{\delta_{s^w s'^w} - \gamma \mathcal{P}_{s^w s'^w}^a}{q + (1-q)\delta_{s^w s'^w}}$ and

$$\hat{\alpha}_{s^w s'^w}^a = (1 - \delta_{s^w s'^w}) \frac{\gamma \hat{\mathcal{P}}_{s^w s'^w}^a}{q + (1-q)\delta_{s^w s'^w}}$$

The Robust Linear Dynamic Programming is the following:

$$\max c_\omega^T \mathcal{J}_\omega \quad (1)$$

$$\text{subject to: } \sum_{j=1}^{|\mathcal{S}|} \sum_{v=1}^q \alpha_{s^w s_j^v}^a \mathcal{J}(s_j^v) + z_s^{aw} \Gamma_s^{aw} + \sum_{j=1}^{|\mathcal{S}|} \sum_{v=1}^q \rho_{s^w s_j^v} \leq g_{s^w}^a, \forall s^w, a \quad (2)$$

$$z_s^{aw} + \rho_{s^w s_j^v} \geq \hat{\alpha}_{s^w s_j^v}^a y_j^v, \forall s^w, a, j, v \quad (3)$$

$$-y_j^v \leq \mathcal{J}(s_j^v) \leq y_j^v, \forall j, v \quad (4)$$

$$\rho_{s^w s_j^v} \geq 0, \forall s^w, j, v \quad (5)$$

$$z \geq 0 \quad (6)$$

$$y \geq 0. \quad (7)$$

Robust Linear Dynamic Programming

The Reduced Robust Linear Dynamic Programming is the following:

$$\max \sum_{i=1}^l \frac{\bar{c}_i}{q} \sum_{v=1}^q \mathcal{J}(\bar{s}_i) \quad (8)$$

$$\text{subject to: } \sum_{j=1}^l \sum_{v=1}^q |\bar{S}_j| \alpha_{\bar{s}_j^w \bar{s}_j^v}^a \mathcal{J}(\bar{s}_j^v) + z_{\bar{s}_i}^{aw} \Gamma_{\bar{s}_i}^{aw} + \sum_{j=1}^l \sum_{v=1}^q p_{s^w s^v} \leq g_{\bar{s}_i^w}^a, \quad \forall \bar{s}_i^w, a \quad (9)$$

$$z_{\bar{s}_i}^{aw} + p_{\bar{s}_i^w \bar{s}_j^v} \geq \hat{\alpha}_{\bar{s}_i^w \bar{s}_j^v}^a y_j^v, \quad \forall \bar{s}_i^w, a, j, v \quad (10)$$

$$-y_j^v \leq \mathcal{J}(\bar{s}_j^v) \leq y_j^v, \quad \forall j, v \quad (11)$$

$$p_{\bar{s}_i^w \bar{s}_j^v} \geq 0, \quad \forall \bar{s}_i^w, j, v \quad (12)$$

$$z \geq 0 \quad (13)$$

$$y \geq 0. \quad (14)$$

Example: Strategy for American Options

Definitions

- An option is the right to engage in a future transaction on some underlying security for a certain prescribed price known as the exercise or strike price
- There are two basic types of options: American and European options
- There are two kinds of American options: call options and put options

Example: Strategy for American Options

Definitions

- An option is the right to engage in a future transaction on some underlying security for a certain prescribed price known as the exercise or strike price
- There are two basic types of options: American and European options
- There are two kinds of American options: call options and put options

Example: Strategy for American Options

Definitions

- An option is the right to engage in a future transaction on some underlying security for a certain prescribed price known as the exercise or strike price
- There are two basic types of options: American and European options
- There are two kinds of American options: call options and put options

Example: Strategy for American Options

Basic concepts

- K is the strike price
- T is the expiration date
- x_t is the price of the underlying asset at t
- r is the risk free interest rate
- σ is the volatility
- The intrinsic value of a put option for the holder is:

$$f(x_t) = \begin{cases} K - x_t, & \text{if } x_t \leq K \text{ and } t \leq T \\ 0, & \text{otherwise} \end{cases}$$

Example: Strategy for American Options

Basic concepts

- K is the strike price
- T is the expiration date
- x_t is the price of the underlying asset at t
- r is the risk free interest rate
- σ is the volatility
- The intrinsic value of a put option for the holder is:

$$f(x_t) = \begin{cases} K - x_t, & \text{if } x_t \leq K \text{ and } t \leq T \\ 0, & \text{otherwise} \end{cases}$$

Example: Strategy for American Options

Basic concepts

- K is the strike price
- T is the expiration date
- x_t is the price of the underlying asset at t
- r is the risk free interest rate
- σ is the volatility
- The intrinsic value of a put option for the holder is:

$$f(x_t) = \begin{cases} K - x_t, & \text{if } x_t \leq K \text{ and } t \leq T \\ 0, & \text{otherwise} \end{cases}$$

Example: Strategy for American Options

Basic concepts

- K is the strike price
- T is the expiration date
- x_t is the price of the underlying asset at t
- r is the risk free interest rate
- σ is the volatility
- The intrinsic value of a put option for the holder is:

$$f(x_t) = \begin{cases} K - x_t, & \text{if } x_t \leq K \text{ and } t \leq T \\ 0, & \text{otherwise} \end{cases}$$

Example: Strategy for American Options

Basic concepts

- K is the strike price
- T is the expiration date
- x_t is the price of the underlying asset at t
- r is the risk free interest rate
- σ is the volatility
- The intrinsic value of a put option for the holder is:

$$f(x_t) = \begin{cases} K - x_t, & \text{if } x_t \leq K \text{ and } t \leq T \\ 0, & \text{otherwise} \end{cases}$$

Example: Strategy for American Options

Basic concepts

- K is the strike price
- T is the expiration date
- x_t is the price of the underlying asset at t
- r is the risk free interest rate
- σ is the volatility
- The intrinsic value of a put option for the holder is:

$$f(x_t) = \begin{cases} K - x_t, & \text{if } x_t \leq K \text{ and } t \leq T \\ 0, & \text{otherwise} \end{cases}$$

Example: Strategy for American Options

American Options and DP

- American options can be exercised at any time before maturity
- It is important to set exercising strategies that maximize the profit
- American options can be implemented as optimal stopping problems (OSPs)
- OSPs can be written in the form of a Bellman equation (DP)

Example: Strategy for American Options

American Options and DP

- American options can be exercised at any time before maturity
- It is important to set exercising strategies that maximize the profit
- American options can be implemented as optimal stopping problems (OSPs)
- OSPs can be written in the form of a Bellman equation (DP)

Example: Strategy for American Options

American Options and DP

- American options can be exercised at any time before maturity
- It is important to set exercising strategies that maximize the profit
- American options can be implemented as optimal stopping problems (OSPs)
- OSPs can be written in the form of a Bellman equation (DP)

Example: Strategy for American Options

American Options and DP

- American options can be exercised at any time before maturity
- It is important to set exercising strategies that maximize the profit
- American options can be implemented as optimal stopping problems (OSPs)
- OSPs can be written in the form of a Bellman equation (DP)

Example: Strategy for American Options

American Options and DP

- The value function is the expected value of the intrinsic value under the risk neutral assumption:

$$\mathcal{J}(x_t) = \max_{t \in [0, T]} E_{\pi} \{ e^{-rt} f(x_t) \}$$

- The evolution of x_t is simulated using the Binomial Options Pricing Model (BOPM)

$$x_{t+1} = \begin{cases} ux_t, & \text{with probability } p \\ dx_t, & \text{with probability } 1 - p \end{cases} .$$

Example: Strategy for American Options

American Options and DP

- The value function is the expected value of the intrinsic value under the risk neutral assumption:

$$\mathcal{J}(x_t) = \max_{t \in [0, T]} E_{\pi} \{ e^{-rt} f(x_t) \}$$

- The evolution of x_t is simulated using the Binomial Options Pricing Model (BOPM)

$$x_{t+1} = \begin{cases} ux_t, & \text{with probability } p \\ dx_t, & \text{with probability } 1 - p \end{cases} .$$

Example: Strategy for American Options

The Binomial Options Pricing Model

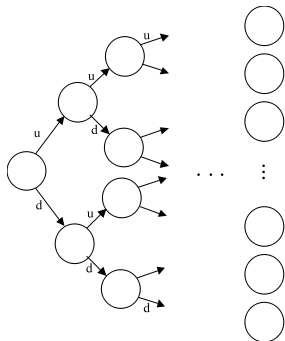


Figure: The Binomial Options Pricing Model

Example: Strategy for American Options

Robust Binomial Options Pricing Model

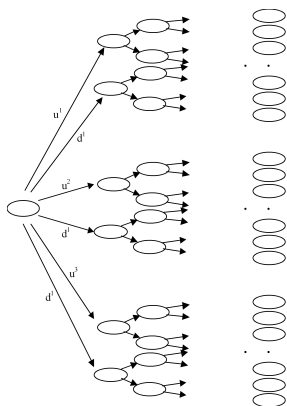


Figure: The Robust Binomial Options Pricing Model

Example: Strategy for American Options

Numerical Example

The parameters used are:

	$K = 1$	$x_0 = .5$	$e^{-r} = .99$	
	Low Volatility	Moderate Volatility	High Volatility	Average
u	$\frac{13}{12}$	$\frac{8}{6}$	$\frac{9}{5}$	1.406
d	$\frac{12}{13}$	$\frac{6}{8}$	$\frac{5}{9}$.711

Table: Parameters Used for the Exercising Strategy

Example: Strategy for American Options

Intrinsic Value

The intrinsic value for this example is:

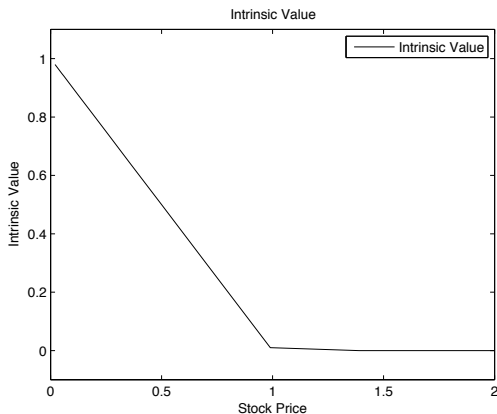


Figure: The Intrinsic Value for the Put Option

Example: Strategy for American Options

Results

Time	Low Volatility	Moderate Volatility	High Volatility	Average
1	0.393	0.281	0.154	0.253
2	0.393	0.281	0.154	0.253
3	0.393	0.281	.278	0.253
4	0.393	0.281	.278	0.356
5	0.426	0.375	.278	0.356
6	0.426	0.375	.278	0.356
7	0.461	0.375	.278	0.356
8	0.5	0.5	0.5	0.5
9	0.7	0.7	0.7	0.7
10	1	1	1	1

Table: Exercising Strategies -Robust Vs Average-

Example: Strategy for American Options

Exercising Strategies

The intrinsic value for this example is:

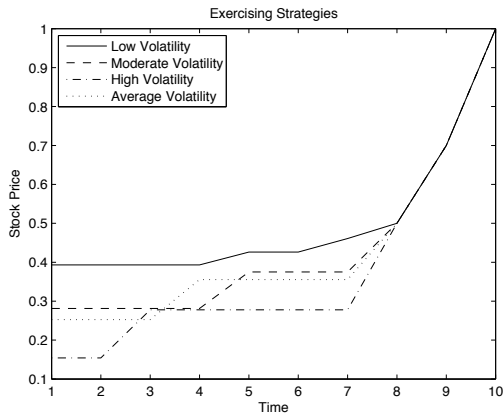


Figure: Exercising Strategies for the Robust Formulation and the Average

Conclusions and Future Research

Conclusions

- This work was motivated by the “curse of dimensionality” in reinforcement learning and dynamic programming
- Coped with this problem using clustering
- Conducted a mathematical analysis to support the results obtained
- Developed an on line dynamic programming procedure
- Presented a novel LP algorithm to mitigate uncertainties within the transition probability matrices and the cost/reward

Conclusions and Future Research

Conclusions

- This work was motivated by the “curse of dimensionality” in reinforcement learning and dynamic programming
- Coped with this problem using clustering
- Conducted a mathematical analysis to support the results obtained
- Developed an on line dynamic programming procedure
- Presented a novel LP algorithm to mitigate uncertainties within the transition probability matrices and the cost/reward

Conclusions and Future Research

Conclusions

- This work was motivated by the “curse of dimensionality” in reinforcement learning and dynamic programming
- Coped with this problem using clustering
- Conducted a mathematical analysis to support the results obtained
- Developed an on line dynamic programming procedure
- Presented a novel LP algorithm to mitigate uncertainties within the transition probability matrices and the cost/reward

Conclusions and Future Research

Conclusions

- This work was motivated by the “curse of dimensionality” in reinforcement learning and dynamic programming
- Coped with this problem using clustering
- Conducted a mathematical analysis to support the results obtained
- Developed an on line dynamic programming procedure
- Presented a novel LP algorithm to mitigate uncertainties within the transition probability matrices and the cost/reward

Conclusions and Future Research

Conclusions

- This work was motivated by the “curse of dimensionality” in reinforcement learning and dynamic programming
- Coped with this problem using clustering
- Conducted a mathematical analysis to support the results obtained
- Developed an on line dynamic programming procedure
- Presented a novel LP algorithm to mitigate uncertainties within the transition probability matrices and the cost/reward

Conclusions and Future Research

Future Research

- Develop a linear dynamic programming that uses multicriteria optimization (state relevant weight)
- Find a good trade off between the reduction and the accuracy rates
- Develop the concept presented in this talk for an infinite horizon and for continuous states' space

Conclusions and Future Research

Future Research

- Develop a linear dynamic programming that uses multicriteria optimization (state relevant weight)
- Find a good trade off between the reduction and the accuracy rates
- Develop the concept presented in this talk for an infinite horizon and for continuous states' space

Conclusions and Future Research

Future Research

- Develop a linear dynamic programming that uses multicriteria optimization (state relevant weight)
- Find a good trade off between the reduction and the accuracy rates
- Develop the concept presented in this talk for an infinite horizon and for continuous states' space

Thank you