Planning Under Uncertainty Computer Science cpsc322, Lecture 11

(Textbook Chpt 9.1-3)

June 12, 2012

Planning in Stochastic Environments



Planning Under Uncertainty: Intro

- **Planning** how to select and organize a sequence of actions/decisions to achieve a given goal.
- Deterministic Goal: A possible world in which some propositions are true

- Planning under Uncertainty: how to select and organize a sequence of actions/decisions to "<u>maximize the probability</u>" of "achieving a given goal"
 - <u>Goal under Uncertainty</u>: we'll move from all-ornothing goals to a richer notion: rating how *happy* the agent is in different possible worlds.

"Single" Action vs. Sequence of Actions

Set of primitive decisions that can be treated as a single macro decision to be made before acting one-off

- Agents makes observations
- Decides on an action
- Carries out the action

Sequential Decisions

Lecture Overview

- One-Off Decision
 - Example
 - Optimal Decision: Utilities / Preferences
 - Single stage Decision Networks
- Sequential Decisions
 - Representation
 - Policies
 - Finding Optimal Policies

One-off decision example

Delivery Robot Example not

Ner

- Robot needs to reach a certain room
- Atolse (Going through stairs may cause an accident.
- It can go the short way through long stairs, or the long way through short stairs (that reduces the chance of an accident but takes more time)
- Which Woy long (1) P(A=t | WW=long) < P(A=t | WW= short (1) P(A=t | WW=long) < P(A=t | WW= short (1) Short

• The Robot can choose to wear pads to protect itself or not (to protect itself in case of an accident) but pads slow it down

pods f i i + A=t i If there is an accident the Robot does not get to the room

Decision Tree for Delivery Robot

This scenario can be represented as the following decision tree



- The agent has a set of decisions to make (a macro-action it can perform)
- Decisions can influence random variables
- Decisions have probability distributions over outcomes

Decision Variables: Some general Considerations

- A possible world specifies a value for each random variable and each decision variable.
- For each assignment of values to all decision variables, the probabilities of the worlds satisfying that assignment sum to 1.



Lecture Overview

- One-Off Decision
 - Example
 - Optimal Decision: Utilities / Preferences
 - Single stage Decision Networks
- Sequential Decisions
 - Representation
 - Policies
 - Finding Optimal Policies

What are the optimal decisions for our Robot?

It all depends on how happy the agent is in different situations.

For sure getting to the room is better than not getting there..... but we need to consider other factors..



Utility / Preferences

Utility: a measure of desirability of possible worlds to an agent

 Let U be a real-valued function such that U(w) represents an agent's degree of preference for world w. [0, 100]

Would this be a reasonable utility function for our Robot?

-		V			
	Which way	Accident	Wear Pads	Utility	World
	short	true	true	35	w0, moderate damage
\rightarrow	short	false	true	95	w1, reaches room, quick, extra weight
-	long	true	true	30	w2, moderate damage, low energy
	long	false	true	75	w3, reaches room, slow, extra weight
	short	true	false	3	w4, severe damage
\rightarrow	short	false	false	100	w5, reaches room, quick
	long true	faise	false	0	w6, severe damage, low energy
コ	long	tene	false	80	w7, reaches room, slow
~1	t	3158			

Utility: Simple Goals

• Can simple (boolean) goals still be specified? avoid : ⁴ reaching the room⁴ Accident

must be

talse

V				1
	Which way	Accident	Wear Pads	Utility
	long	true	true	0
	long	true	false	0
$\left(\right)$	long	false	true	100
\succ	long	false	false	100
	short	true	true	0
\sim	short	true	false	\bigcirc
<pre></pre>	short	false	true	102
	short	false	false	100

Optimal decisions: How to combine Utility with Probability

What is the utility of achieving a certain probability distribution over possible worlds?



 It is its <u>expected utility/value i.e.</u>, its average utility, weighting possible worlds by their probability.

$$EU(wP=t, WW = short) = .2 \times 35 + .8 \times 75$$

Optimal decision in one-off decisions

 Given a set of *n* decision variables var_i (e.g., Wear Pads, Which Way), the agent can choose:

 $D = d_i$ for any $d_i \otimes \text{dom}(var_1) \times ... \times \text{dom}(var_n)$.



Optimal decision: Maximize Expected Utility

• The expected utility of decision $D = d_i$ is

$$\mathbb{E}(U \mid D = d_i) = \sum_{w \models D = d_i} P(w \mid D = d_i) U(w)$$

e.g.,
$$\mathbb{E}(U \mid D = \{WP = \text{isc}, WW = \text{short} \} =$$



 $P(w_4|D) \neq U(w_4) + P(w_5|D) \neq U(w_5)$

msx

short

 An optimal decision is the decision D = d_{max} whose expected utility is maximal:
 Wear Pads Which way EU

$$d_{\max} = \underset{d_i \in dom(D)}{\operatorname{arg\,max}} \mathbb{E}(U \mid D = d_i) \xrightarrow{}_{i \in dom(D)} \operatorname{true}_{\substack{\text{false} \\ i \in dom(D)}} \operatorname{true}_{\substack{\text{false} \\ \text{false} \\ \text{false} \\ i \in dom(D)}} \operatorname{true}_{\substack{\text{false} \\ i \in$$

Expected utility of a decision

• The expected utility of decision $D = d_i$ is

$$\mathbb{E}(U \mid D = d_i) = \sum_{w \models (D = d_i)} P(w) U(w)$$

 What is the expected utility of Wearpads=yes, Way=short ?



Lecture Overview

- One-Off Decision
 - Example
 - Optimal Decision: Utilities / Preferences
 - Single stage Decision Networks
- Sequential Decisions
 - Representation
 - Policies
 - Finding Optimal Policies

Single-stage decision networks

Extend belief networks with:

- **Decision nodes**, that the agent chooses the value for. Drawn as rectangle.
- Utility node, the parents are the variables on which the utility depends. Drawn as a diamond.
- Shows explicitly which decision nodes
 affect random variables



ds.	SHOL	laise	0.0	,	
des					
Which way	Accident	Wear Pads]	Uti	lity
long	true	true		30	
long	true	false		0	
long	false	true		75	
long	false	false		80	
short	true	true		35	
short	true	false		3	
short	false	true		95	

false

short

false

100

	Which	Accident	
	way		
_	long	true	0.01
	long	false	0.99
	short	true	0.2
	short	false	0.8



To find the optimal decision we can use VE:

- 1. Create a factor for each conditional probability and for the utility
- 2. Multiply factors and sum out all of the random variables (This creates a factor on \mathcal{D} that gives the expected utility for each d_{i})
- 3. Choose the ∂_{i} with the maximum value in the factor.

VE Example: Step 1, create initial factors



Which Way	Accider	nt	Utility		Ste f(A	p 2a: .,W,P	$compu$ $P) = f_1(A)$	te product (A,W) × $f_2(A)$	A,W,P)
	Wear Pad.	5	f(A=a,F	P=p,V	V=w) =	= f ₁ (A	=a,W=w)	× $f_2(A=a,W=$	=w,P=p)
Which way W	Accident A	f ₁ (A,W)							
long	true	0.01			Which	way W	Accident A	A Pads P	f(A,W,P)
long	false	0.99			lona		true	true	0.01 * 30
short	true	0.2			long		true	false	
short	false	0.8			long		false	true	
Which way W	Accident A	Pads P	f ₂ (A,W,P)]	long		false	false	???
long	true	true	30		short		true	false	
long	true	false	0		short		false	true	
long	false	true	75		short		false	false	
long	false	false	80						ļ
short	true	true	35						
short	true	false	3			0.99 *	30	0.01 * 80	
short	false	true	95						
short	false	false	100			0.99	* 80	0.8 * 30	23

Which Way	Accider	nt	Utility		Step 2a: f(A,W,F	$\mathbf{P} = f_1(\mathbf{A})$	e product ,W) × $f_2(A$	A,W,P)
	Wear Pad.	5	f(A=a,I	⊃=p,V	$V=w) = f_1(A)$	=a,W=w) >	< f ₂ (A=a,W=	=w,P=p)
Which way W	Accident A	f ₁ (A,W)						
long	true	0.01			Which way W	Accident A	Pads P	f(A,W,P)
long	false	0.99			long	true	true	0.01 * 30
short	true	0.2			long	true	false	0.01*0
short	false	0.8			long	false	true	0.99*75
Which way W	Accident A	Pads P	f ₂ (A,W,P)	1	long	false	false	0.99*80
			2	1	short	true	true	0.2*35
long	true	true	30		short	true	false	0.2*3
long	true	false	0		short	false	true	0.8*95
long	false	true	75		short	false	false	0.8*100
long	false	false	80					
short	true	true	35					
short	true	false	3					
short	false	true	95					
short	false	false	100					24



Step 2b: sum A out of the product f(A,W,P):

$$f_3(W,P) = \sum_A f(A,W,P)$$

Which way W	Pads P	f ₃ (W,P)
long long short short	true false true false	0.01*30+0.99*75=74.55 ??

0.2*35 + 0.2*0.3

0.2*35 + 0.8*95

0.99*80 + 0.8*95

0.8 * 95 + 0.8*100

Which way W	Accident A	Pads P	f(A,W,P)
long	true	true	0.01 * 30
long	true	false	0.01*0
long	false	true	0.99*75
long	false	false	0.99*80
short	true	true	0.2*35
short	true	false	0.2*3
short	false	true	0.8*95
short	false	false	0.8*100



Step 2b: sum A out of the product f(A,W,P):

$$f_3(W,P) = \sum_A f(A,W,P)$$

Which way W	Pads P	f ₃ (W,P)
long	true	0.01*30+0.99*75=74.55
long	false	
short	true	0.2*35+0.8*95=83
short	false	

Which way W	Accident A	Pads P	f(A,W,P)
long	true	true	0.01 * 30
long	true	false	0.01*0
long	false	true	0.99*75
long	false	false	0.99*80
short	true	true	0.2*35
short	true	false	0.2*3
short	false	true	0.8*95
short	false	false	0.8*100

VE example: step 3, choose decision with max E(U)



Step 2b: sum A out of the product f(A,W,P):

$$f_3(W,P) = \sum_A f(A,W,P)$$

Which way W	Pads P	f ₃ (W,P)
long	true	0.01*30+0.99*75=74.55
long	false	0.01*0+0.99*80=79.2
short	true	0.2*35+0.8*95=83
short	false	0.2*3+0.8*100=80.6

Which way W	Accident A	Pads P	f(A,W,P)
long	true	true	0.01 * 30
long	true	false	0.01*0
long	false	true	0.99*75
long	false	false	0.99*80
short	true	true	0.2*35
short	true	false	0.2*3
short	false	true	0.8*95
short	false	false	0.8*100

The final factor encodes the expected utility of each decision

- Alspace
- Thus, taking the short way but wearing pads is the best choice, with an expected utility of 83

Learning Goals for today's class – part 1 You can:

- Compare and contrast stochastic single-stage (one-off) decisions vs. multistage decisions
- Define a <u>Utility Function</u> on possible worlds
- Define and compute optimal one-off decision (max expected utility)
- Represent one-off decisions as single stage decision networks and compute optimal decisions by Variable Elimination

Lecture Overview

- One-Off Decision
 - Example
 - Optimal Decision: Utilities / Preferences
 - Single stage Decision Networks
- Sequential Decisions
 - Representation
 - Policies
 - Finding Optimal Policies

"Single" Action vs. Sequence of Actions

Set of primitive decisions that can be treated as a single macro decision to be made *before acting*

- Agent makes observations
- Decides on an action *L*
- Carries out the action ∠

Sequential decision problems

- A sequential decision problem consists of a sequence of decision variables D_1, \ldots, D_n .
- Each D_i has an information set of variables pD_i , whose value will be known at the time decision D_i is made. $PD_3 = \int D_2 \sqrt{3} \sqrt{4}$



Sequential decisions : Simplest possible

- Only one decision! (but different from one-off decisions)
- Early in the morning. Shall I take my umbrella today? (I'll have to go for a long walk at noon)
- Relevant Random Variables?



Policies for Sequential Decision Problem: Intro

• A **policy** specifies what an agent should do under each circumstance (for each decision, consider the parents of the decision node)

In the Umbrella "degenerate" case:



Sequential decision problems: "complete" Example

Utility

Call

 $PC = \{R\}$ $PC = \{R, CS, SS\}$

 $p(S \leq p)$

- A sequential decision problem consists of a sequence of decision variables D₁,....,D_n.
- Each D_i has an information set of variables pD_i, whose value will be known at the time decision D_i is made.

SeeSmoke



decisions are totally ordered

Alarm

Leaving

Tampering

• if a decision D_b comes before D_a , then

Fire

Smoke

Check Smoke

- D_b is a parent of D_a
- any parent of D_b is a parent of D_a

Lecture Overview

- One-Off Decision
 - Example
 - Optimal Decision: Utilities / Preferences
 - Single stage Decision Networks
- Sequential Decisions
 - Representation
 - Policies
 - Finding Optimal Policies

Policies for Sequential Decision Problems

- A policy is a sequence of $\delta_1, \dots, \delta_n$ decision functions $\delta_i : \operatorname{dom}(pD_i) \to \operatorname{dom}(D_i)$
- This policy means that when the agent has observed $O \in \text{dom}(pD_i)$, it will do $\delta_i(O) = \text{Example:} \quad \mathbf{agent}$



When does a possible world satisfy a policy?

- A possible world specifies a value for each random variable and each decision variable.
- **Possible world** *w* **satisfies policy** δ , written $w \models \delta$ if the value of each decision variable is the value selected by its decision function in the policy (when applied in *w*).



When does a possible world satisfy a policy?

• Possible world *w* satisfies policy δ , written $w \models \delta$ if the value of each decision variable is the value selected by its decision function in the policy (when applied in *w*).



Expected Value of a Policy

- Each possible world w has a probability P(w) and a utility U(w)
- The expected utility of policy δ is

$$\sum_{\substack{w \neq \delta}} P(w) * U(w)$$

The optimal policy is one with the Mag

expected utility.

Lecture Overview

- One-Off Decision
 - Example
 - Optimal Decision: Utilities / Preferences
 - Single stage Decision Networks
- Sequential Decisions
 - Representation
 - Policies
 - Finding Optimal Policies

Complexity of finding the optimal policy: how many policies? • How many assignments to parents? Utility Tampering Fire $C57 C 2^{3}$ Smoke Alarm How many decision functions? (binary decisions) 23 Leaving SeeSmoke Check Smoke How many policies? product 3 Report Call

If a decision *D* has *k* binary parents, how many assignments of values to the parents are there? k^2 2^k k Dk

 $b2^k$

- If there are *b* possible actions (possible values for D), how many ۲ different decision functions are there? $2^{(k+1)}$
- If there are d decisions, each with k binary parents and b possible ۲ actions, how many policies are there?

Finding the optimal policy more efficiently: VE

- 1. Create a factor for each conditional probability table and a \angle factor for the utility.
- 2. Sum out random variables that are not parents of a decision node.
- 3. Eliminate (aka sum out) the decision variables
- 4. Sum out the remaining random variables.
- 5. Multiply the factors: this is the expected utility of the optimal policy.





Eliminate the decision Variables: step3 details

- Select a variable *D* that corresponds to the latest decision to be made
 - this variable will appear in only one factor with its parents
- Eliminate *D* by maximizing. This returns:
 - The optimal decision function for D_{i} arg max_D f_{j}
 - A new factor to use in VE, max_D f
- Repeat till there are no more decision nodes.





Report	CheckSmoke
true	true &
false	talse

VE elimination reduces complexity of finding the optimal policy

- We have seen that, if a decision *D* has *k* binary parents, there are *b* possible actions, If there are d decisions,
- Then there are: $(b^{2^k})^d$ policies
 - Doing variable elimination lets us find the optimal policy after considering only d. b^{2k} policies (we eliminate one decision at a time)
 - VE is much more efficient than searching through policy space.
 - However, this complexity is <u>still doubly-exponential</u> we'll only be able to handle relatively small problems. + give up nonforgetting somp + give up nonforgetting somp



Big Picture: Planning under Uncertainty



Learning Goals for today's class – part 2

You can:

- Represent sequential decision problems as decision networks. And explain the non forgetting property
- Verify whether a possible world satisfies a policy and define the expected value of a policy
- Compute the <u>number of policies</u> for a decision problem
- Compute the optimal policy by Variable Elimination

Cpsc 322 Big Picture





Announcements

• Fill out Online Teaching Evaluations Survey.

• FINAL EXAM: Thur June14, 9:00-11:30 pm DMP 310 (NOT regular room) 67pts

Final will comprise: 10 -15 short questions + 3-4 problems

- Work on all practice exercises and sample problems
- While you revise the learning goals, work on review questions
 I may even reuse some verbatim ^(C)
 - Come to remaining Office hours! Today and Tomorrow (2-4 X150 (Learning Center))

STRIKE? http://vpstudents.ubc.ca/news/strike-action/