

#### Policy Learning and Markov Decision Processes

Thursday 24 October 2002

William H. Hsu

#### **Department of Computing and Information Sciences, KSU**

http://www.kddresearch.org http://www.cis.ksu.edu/~bhsu

Readings: Chapter 17, Russell and Norvig Sections 13.1-13.2, Mitchell



CIS 732: Machine Learning and Pattern Recognition

#### **Lecture Outline**

- Readings: Chapter 17, Russell and Norvig; Sections 13.1-13.2, Mitchell
- Suggested Exercises: 17.2, Russell and Norvig; 13.1, Mitchell
- This Week's Paper Review: Temporal Differences [Sutton 1988]
- Making Decisions in Uncertain Environments
  - Problem definition and framework (MDPs)
  - Performance element: computing optimal policies given *stepwise reward* 
    - Value iteration
    - Policy iteration
  - Decision-theoretic agent design
    - Decision cycle
    - Kalman filtering
    - Sensor fusion aka data fusion
  - <u>Dynamic Bayesian networks (DBNs) and dynamic decision networks (DDNs)</u>
- Learning Problem: Acquiring Decision Models from Rewards
- Next Lecture: Reinforcement Learning



# In-Class Exercise: Elicitation of Numerical Estimates [1]

- Almanac Game [Heckerman and Geiger, 1994; Russell and Norvig, 1995]
  - Used by decision analysts to calibrate numerical estimates
  - Numerical estimates: include <u>subjective probabilities</u>, other forms of knowledge
- Question Set 1 (Read Out Your Answers)
  - Number of passengers who flew between NYC and LA in 1989 3 million Population of Warsaw in 1992 1.6 million Year in which Coronado discovered the Mississippi River 1541 Number of votes received by Carter in the 1976 presidential election 41 million Number of newspapers in the U.S. in 1990 1611 Height of Hoover Dam in feet 221 Number of eggs produced in Oregon in 1985 649 million Number of Buddhists in the world in 1992 295 million Number of deaths due to AIDS in the U.S. in 1981 132 Number of U.S. patents granted in 1901 25546



**CIS 732: Machine Learning and Pattern Recognition** 

# In-Class Exercise: Elicitation of Numerical Estimates [2]

- Calibration of Numerical Estimates
  - Try to revise your bounds based on results from first question set
  - Assess your own penalty for having too wide a CI versus guessing low, high
- Question Set 2 (Write Down Your Answers)

<ul> <li>Year of birth of Zsa Zsa Gabor</li> </ul>	1917
<ul> <li>Maximum distance from Mars to the sun in miles</li> </ul>	155 million
<ul> <li>Value in dollars of exports of wheat from the U.S. in 1992</li> </ul>	4.5 billion
<ul> <li>Tons handled by the port of Honolulu in 1991</li> </ul>	11 million
<ul> <li>Annual salary in dollars of the governor of California in 1993</li> </ul>	120000
<ul> <li>Population of San Diego in 1990</li> </ul>	1.1 million
<ul> <li>Year in which Roger Williams founded Providence, RI</li> </ul>	1636
<ul> <li>Height of Mt. Kilimanjaro in feet</li> </ul>	19340
<ul> <li>Length of the Brooklyn Bridge in feet</li> </ul>	1595
<ul> <li>Number of deaths due to auto accidents in the U.S. in 1992</li> </ul>	41710



**CIS 732: Machine Learning and Pattern Recognition** 

# In-Class Exercise: Elicitation of Numerical Estimates [3]

- Descriptive Statistics
  - 50%, 25%, 75% guesses (median, first-second quartiles, third-fourth quartiles)
  - Box plots [Tukey, 1977]: actual frequency of data within 25-75% bounds
  - What kind of descriptive statistics do you think might be <u>informative</u>?
  - What kind of descriptive graphics do you think might be informative?
- Common Effects
  - Typically about half (50%) in first set
  - Usually, see some improvement in second set
  - Bounds also widen from first to second set (second system effect [Brooks, 1975])
  - Why do you think this is?
  - What do you think the ramifications are for interactive elicitation?
  - What do you think the ramifications are for learning?
- Prescriptive (Normative) Conclusions
  - Order-of-magnitude ("back of the envelope") calculations [Bentley, 1985]
  - <u>Value-of-information (VOI)</u>: framework for selecting questions, precision



#### **Overview:**

# **Making Decisions in Uncertain Environments**

- Problem Definition
  - Given: stochastic <u>environment</u>, outcome *P*(*Result* (*action*) | *Do*(*action*), *state*)
  - Return: a policy f : state  $\rightarrow$  action
- Foundations of <u>Sequential Decision Problems</u> and Policy Learning
  - <u>Utility</u> function:  $U : state \rightarrow value$
  - U(State): analogy with P(State) = agent's belief as distributed over event space
  - Expresses desirability of state according to decision-making agent
- Constraints and Rational Preferences
  - Definition: a <u>lottery</u> is defined by the set of outcomes of a random scenario and a probability distribution over them (e.g., denoted [*p*, *A*; 1 - *p*, *B*] for outcomes *A*, *B*)
  - Properties of <u>rational preference</u> (ordering on utility values)
    - <u>Total ordering</u>: antisymmetric, <u>transitive</u>, and  $\forall A, B.(A \succ B) \lor (B \succ A) \lor (A \sim B)$
    - <u>Continuity</u>:  $A \succ B \succ C \Rightarrow \exists p . [p, A; 1-p, C] \sim [1, B] \equiv B$
    - Substitutability:  $A \sim B \Rightarrow [p, A; 1-p, C] \sim [p, B; 1-p, C]$
    - <u>Monotonicity</u>:  $A \succ B \Rightarrow (p > q \Leftrightarrow [p, A; 1-p, B] \succ [q, A; 1-q, B])$
    - <u>Decomposability</u>:  $[p, A; 1-p, [q, B; 1-q, C]] \sim [p, A; (1-p)q, B; (1-p)(1-q), C]$

# <u>Markov Decision Processes</u> and <u>Markov Decision Problems</u>

- <u>Maximum Expected Utility (MEU)</u>
  - $E[U(action | D)] = \sum_i P(Result_i(action) | Do(action), D) \cdot U(Result_i(action))$
  - D denotes agent's available evidence about world
  - <u>Principle</u>: rational agent should choose actions to maximize expected utility
- <u>Markov Decision Processes (MDPs</u>)
  - Model: probabilistic state transition diagram, associated actions A: state  $\rightarrow$  state
  - <u>Markov property</u>: transition probabilities from any given state depend only on the state (not previous history)
  - Observability
    - Totally observable (MDP, <u>TOMDP</u>), aka <u>accessible</u>
    - Partially observable (POMDP), aka inaccessible, hidden
- <u>Markov Decision Problems</u>
  - Also called MDPs
  - Given: a stochastic environment (process model, utility function, and D)
  - Return: an optimal policy f : state  $\rightarrow$  action



# **Value Iteration**

- Value Iteration: Computing Optimal Policies by Dynamic Programming
  - Given: transition model *M*, reward function *R*: state  $\rightarrow$  value
  - $M_{ij}(a)$  denotes probability of moving from state *i* to state *j* via action *a*
  - Additive utility function on state sequences:  $U[s_0, s_1, ..., s_n] = R(s_0) + U[s_1, ..., s_n]$
- Function Value-Iteration (M, R)
  - Local variables U, U': "current" and "new" utility functions, initially identical to R
  - REPEAT
    - *U* ← *U*'
    - FOR each state *i* DO
       // dynamic programming update

 $U'[i] \leftarrow R[i] + max_a \sum_j M_{ij}(a) \cdot U[j]$ 

UNTIL Close-Enough (U, U')

– RETURN U

// approximate utility function on all states

- Result: Provably Optimal Policy [Bellman and Dreyfus, 1962]
  - Use computed U by maximizing utility  $U(next action | s_i)$
  - Evaluation: RMS error of U or expected difference  $U^*$  U (policy loss)



# **Policy Iteration**

- Policy Iteration: Another Algorithm for Calculating Optimal Policies
  - Given: transition model *M*, reward function *R*: *state*  $\rightarrow$  *value*
  - <u>Value determination</u> function: estimates current *U* (e.g., by solving linear system)
- Function Policy-Iteration (M, R)
  - Local variables U: initially identical to R; P: policy, initially optimal under U
  - REPEAT
    - *U* ← *Value-Determination* (*P*, *U*, *M*, *R*); *unchanged*? ← true
    - FOR each state *i* DO // dynamic programming update

IF  $max_a \sum_j M_{ij}(a) \cdot U[j] > \sum_j M_{ij}(P[i]) \cdot U[j]$  THEN

 $P[i] \leftarrow R[i] + arg max_a \sum_j M_{ij}(a) \cdot U[j]; unchanged? \leftarrow false$ 

UNTIL unchanged?

– RETURN P

// optimized policy

- Guiding Principle: Value Determination Simpler than Value Iteration
  - Reason: action in each state is fixed by the policy
  - Solutions: use value iteration without *max*; solve linear system



# Applying Policies: Decision Support, Planning, and Automation

- Decision Support
  - Learn an <u>action-value</u> function (to be discussed soon)
  - Calculate MEU action in current state
  - Open loop mode: recommend MEU action to agent (e.g., user)
- Planning
  - Problem specification
    - <u>Initial</u> state  $s_0$ , <u>goal</u> state  $s_G$
    - Operators (actions, preconditions = applicable states, effects = transitions)
  - Process: computing policy to achieve goal state
  - Traditional: <u>symbolic</u>; <u>first-order logic (FOL</u>), subsets thereof
  - "Modern": <u>abstraction</u>, <u>conditionals</u>, <u>temporal constraints</u>, <u>uncertainty</u>, etc.
- Automation
  - Direct application of policy
  - Caveats: partially observable state, uncertainty (measurement error, etc.)



# **Decision-Theoretic Agents**

- Function *Decision-Theoretic-Agent* (*Percept*)
  - <u>Percept</u>: agent's input; collected evidence about world (from sensors)
  - COMPUTE updated <u>probabilities for current state</u> based on available evidence, including current percept and previous action
  - COMPUTE <u>outcome probabilities</u> for actions, given action descriptions and probabilities of current state
  - SELECT <u>action</u> with highest expected utility, given probabilities of outcomes and utility functions
  - **RETURN** action
- Decision Cycle
  - Processing done by rational agent at each step of action
  - Decomposable into prediction and estimation phases
- Prediction and Estimation
  - <u>Prediction</u>: compute pdf over expected states, given knowledge of previous state, effects of actions
  - Estimation: revise belief over current state, given prediction, new percept



# **Kalman Filtering**

- Intuitive Idea
  - Infer "where we are" in order to compute outcome probabilities, select action
  - Inference problem: estimate Bel(X(t))
- Problem Definition
  - Given: action history, new percept
  - Return: estimate of probability distribution over current state
- Assumptions
  - State variables: real-valued, normal (Gaussian) distribution
  - Sensors: unbiased (mean = 0), normally distributed (Gaussian) noise
  - Actions: can be described as vector of real values, one for each state variable
  - New state: linear function of previous state, action
- Interpretation as Bayesian Parameter Estimation
  - Technique from classical control theory [Kalman, 1960]
  - Good success even when not all assumptions are satisfied
  - Prediction:  $B\hat{e}I(X(t)) = \sum_{x \in Y} P(X(t) | X(t-1) = x(t-1), action(t-1)) \cdot BeI(X(t-1) = x(t-1))$
  - <u>Estimation</u>:  $Bel(X(t)) = \stackrel{X(t-1)}{\alpha} \cdot P(percept(t) | X(t)) \cdot B\hat{e}l(X(t))$







# **Sensor and Data Fusion**

- Intuitive Idea
  - Sensing in uncertain worlds
  - Compute estimates of <u>conditional probability tables (CPTs</u>)
    - Sensor model (how environment generates sensor data): P(percept(t) | X(t))
    - Action model (how actuators affect environment): P(X(t) | X(t-1), action(t-1))
  - Use estimates to implement Decision-Theoretic-Agent : percept  $\rightarrow$  action
- Assumption: <u>Stationary</u> Sensor Model
  - Stationary sensor model:  $\forall t . P(percept(t) | X(t)) = P(percept(t) | X)$ 
    - <u>Circumscribe</u> (exhaustively describe) percept <u>influents</u> (variables that affect sensor performance)
    - *NB*: this does *not* mean sensors are immutable or unbreakable
  - <u>Conditional independence</u> of sensors given true value
- Problem Definition
  - Given: multiple sensor values for same state variables
  - Return: combined sensor value
  - Inferential process: sensor fusion, aka sensor integration, aka data fusion

Sensor Model

**S(t)** 

Sensor Model

# **Dynamic Bayesian Networks (DBNs)**

- Intuitive Idea
  - State of environment evolves over time
    - Evolution modeled by conditional pdf: P(X(t) | X(t-1), action(i-1))
    - Describes how state depends on previous state, action of agent
  - <u>Monitoring</u> scenario
    - Agent can only observe (and predict): P(X(t) | X(t-1))
    - State evolution model, aka Markov chain
  - Probabilistic projection
    - Predicting <u>continuation</u> of observed X(t) values (see last lecture)
    - Goal: use results of prediction and monitoring to make decisions, take action
- Dynamic Bayesian Network (aka Dynamic Belief Network)
  - Bayesian network *unfolded through time* (one note for each state and sensor variable, at each step)
  - Decomposable into prediction, rollup, and estimation phases
  - Prediction: as before; <u>rollup</u>: compute  $B\hat{e}I(X(t))$ ; estimation: <u>unroll</u> X(t + 1)

# **Dynamic Decision Networks (DDNs)**

- Augmented Bayesian Network [Howard and Matheson, 1984]
  - <u>Chance nodes</u> (ovals): denote random variables as in BBNs
  - <u>Decision nodes</u> (rectangles): denote points where agent has choice of actions
  - <u>Utility nodes</u> (diamonds): denote agent's utility function (e.g., in chance of death)
- Properties
  - Chance nodes: related as in BBNs (CI assumed among nodes not connected)
  - Decision nodes: choices can influence chance nodes, utility nodes (directly)
  - Utility nodes: conditionally dependent on joint pdf of parent chance nodes and decision values at parent decision nodes \_\_\_\_\_
  - See Section 16.5, Russell and Norvig
- Dynamic Decision Network
  - aka dynamic influence diagram
  - DDN : DBN :: DN : BBN



- Inference: over predicted (unfolded) sensor, decision variables



# Learning to Make Decisions in Uncertain Environments

- Learning Problem
  - Given: interactive environment
    - No notion of *examples* as assumed in supervised, unsupervised learning
    - Feedback from environment in form of rewards, penalties (reinforcements)
  - Return: utility function for decision-theoretic inference and planning
    - Design 1: utility function on states,  $U : state \rightarrow value$
    - Design 2: <u>action-value function</u>, *Q* : *state* × *action* → *value* (expected utility)
  - Process
    - Build predictive model of the environment
    - Assign credit to components of decisions based on (current) predictive model
- Issues
  - How to <u>explore environment</u> to acquire feedback?
  - <u>Credit assignment</u>: how to propagate positive credit and negative credit (<u>blame</u>)
     back through decision model *in proportion to importance*?



# Terminology

- Making Decisions in Uncertain Environments
  - Policy learning
    - Performance element: decision support system, planner, automated system
    - Performance *criterion*: <u>utility</u> function
    - Training signal: <u>reward</u> function
  - MDPs
    - <u>Markov Decision Process (MDP): model for decision-theoretic planning (DTP)</u>
    - <u>Markov Decision Problem (MDP): problem specification for DTP</u>
    - <u>Value iteration</u>: iteration over actions; decomposition of utilities into rewards
    - Policy iteration: iteration over policy steps; value determination at each step
  - <u>Decision cycle</u>: processing (inference) done by a rational agent at each step
  - Kalman filtering: estimate belief function (pdf) over state by iterative refinement
  - <u>Sensor and data fusion</u>: combining multiple sensors for same state variables
  - <u>Dynamic Bayesian network (DBN): temporal</u> BBN (*unfolded through time*)
  - <u>Dynamic decision network (DDN)</u>: temporal <u>decision network</u>
- Learning Problem: Based upon <u>Reinforcements</u> (Rewards, Penalties)



### **Summary Points**

- Making Decisions in Uncertain Environments
  - Framework: <u>Markov Decision Processes</u>, <u>Markov Decision Problems</u> (MDPs)
  - Computing policies
    - Solving MDPs by <u>dynamic programming</u> given a *stepwise reward*
    - Methods: value iteration, policy iteration
  - Decision-theoretic agents
    - Decision cycle, Kalman filtering
    - Sensor fusion (*aka* data fusion)
  - <u>Dynamic Bayesian networks</u> (DBNs) and <u>dynamic decision networks</u> (DDNs)
- Learning Problem
  - Mapping from observed actions and rewards to decision models
  - Rewards/penalties: reinforcements
- Next Lecture: Reinforcement Learning
  - Basic model: passive learning in a known environment
  - Q learning: policy learning by <u>a</u>daptive <u>dynamic programming</u> (ADP)

