Lecture 17

Introduction to Reinforcement Learning: Q Learning

Thursday 29 October 2002

William H. Hsu

Department of Computing and Information Sciences, KSU

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

Readings: Sections 13.3-13.4, Mitchell Sections 20.1-20.2, Russell and Norvig

KSU

CIS 732: Machine Learning and Pattern Recognition

Lecture Outline

- Readings: Chapter 13, Mitchell; Sections 20.1-20.2, Russell and Norvig
 - Today: Sections 13.1-13.4, Mitchell
 - Review: "Learning to Predict by the Method of Temporal Differences", Sutton
- Suggested Exercises: 13.2, Mitchell; 20.5, 20.6, Russell and Norvig
- Control Learning
 - Control policies that choose optimal actions
 - MDP framework, continued
 - Issues
 - Delayed reward
 - Active learning opportunities
 - Partial observability
 - Reuse requirement
- Q Learning
 - Dynamic programming algorithm
 - Deterministic and nondeterministic cases; convergence properties



Control Learning

- Learning to Choose Actions
 - Performance element
 - Applying policy in uncertain environment (last time)
 - Control, optimization objectives: belong to intelligent agent
 - Applications: automation (including mobile robotics), information retrieval
- Examples
 - Robot learning to dock on battery charger
 - Learning to choose actions to optimize factory output
 - Learning to play Backgammon
- Problem Characteristics
 - <u>Delayed reward</u>: loss signal may be <u>episodic</u> (e.g., win-loss at end of game)
 - **Opportunity for active exploration:** <u>situated learning</u>
 - Possible partially observability of states
 - Possible need to learn multiple tasks with same sensors, <u>effectors</u> (e.g., <u>actuators</u>)



CIS 732: Machine Learning and Pattern Recognition

Example: TD-Gammon

- Learns to Play Backgammon [Tesauro, 1995]
 - Predecessor: *NeuroGammon* [Tesauro and Sejnowski, 1989]
 - Learned from examples of labelled moves (very tedious for human expert)
 - Result: strong computer player, but not grandmaster-level
 - TD-Gammon: first version, 1992 used reinforcement learning
- Immediate Reward
 - +100 if win
 - -100 if loss
 - 0 for all other states
- Learning in *TD-Gammon*
 - Algorithm: temporal differences [Sutton, 1988] next time
 - Training: playing 200000 1.5 million games against itself (several weeks)
 - Learning curve: improves until ~1.5 million games
 - Result: now approximately equal to best human player (won World Cup of Backgammon in 1992; among top 3 since 1995)



Reinforcement Learning: Problem Definition



- Interactive Model
 - State (may be partially observable), incremental reward presented to agent
 - Agent selects actions based upon (current) policy
 - Taking action puts agent into new state in environment
 - New reward: <u>reinforcement</u> (feedback)
 - Agent uses decision cycle to estimate new state, compute outcome distributions, select new actions
- Reinforcement Learning Problem
 - Given
 - Observation sequence $\mathbf{S}_0 \xrightarrow{\mathbf{a}_0: \mathbf{r}_0} \mathbf{S}_1 \xrightarrow{\mathbf{a}_1: \mathbf{r}_1} \mathbf{S}_2 \xrightarrow{\mathbf{a}_2: \mathbf{r}_2} \cdots$
 - **Discount factor** $\gamma \in [0, 1)$
 - Learn to: choose actions that maximize $r(t) + \gamma r(t+1) + \gamma^2 r(t+2) + \dots$



CIS 732: Machine Learning and Pattern Recognition

Quick Review: Markov Decision Processes

- <u>Markov Decision Processes (MDPs</u>)
 - Components
 - Finite set of states S
 - Set of actions A
 - At each time, agent
 - observes state $s(t) \in S$ and chooses action $a(t) \in A$;
 - then receives reward r(t),
 - and state changes to s(t + 1)
 - Markov property, aka <u>Markov assumption</u>: $s(t + 1) = \delta(t + 1)$ and r(t) = r(s(t), a(t))
 - i.e., r(t) and s(t + 1) depend only on *current* state and action
 - Previous history s(0), s(1), ..., s(t 1): <u>irrelevant</u>
 - i.e., s(t + 1) conditionally independent of s(0), s(1), ..., s(t 1) given s(t)
 - δ, *r*: may be nondeterministic; not necessarily known to agent
 - Variants: totally observable (accessible), partially observable (inaccessible)
- Criterion for *a*(*t*): Total Reward <u>Maximum Expected Utility (MEU)</u>



Agent's Learning Task

- Performance Element
 - Execute actions in environment, observe results
 - Learn action policy π : *state* \rightarrow *action* that maximizes <u>expected discounted reward</u> $E[r(t) + \gamma r(t + 1) + \gamma^2 r(t + 2) + ...]$ from any starting state in *S*
 - γ ∈ [0, 1)
 - Discount factor on future rewards
 - Expresses preference for rewards sooner rather than later
- Note: Something New!
 - Target function is π : *state* \rightarrow *action*
 - However...
 - We have no training examples of form < state, action>
 - Training examples are of form << state, action>, reward>



Department of Computing and Information Sciences

Value Function

- First Learning Scenario: Deterministic Worlds
 - Agent considers adopting policy π from <u>policy space</u> Π
 - For each possible policy $\pi \in \Pi$, can define an evaluation function over states: $V^{\pi}(s) \equiv r(t) + \gamma r(t+1) + \gamma^2 r(t+1) + \dots$

$$\equiv \sum_{i=0}^{\infty} \boldsymbol{\gamma}^{i} \boldsymbol{r}(\boldsymbol{t}+\boldsymbol{i})$$

where r(t), r(t + 1), r(t + 2), ... are generated by following policy π starting at state s

- Restated, task is to learn optimal policy π^*

$$\pi^* \equiv \arg \max_{\pi} V^{\pi}(s), \forall s$$

Finding Optimal Policy





CIS 732: Machine Learning and Pattern Recognition

What to Learn

- Idea
 - Might have agent try to learn evaluation function V^{π^*} (abbreviated V^*)
 - Could then perform lookahead search to choose best action from any state s, because:

$$\pi^*(s) \equiv \arg\max_a [r(s, a) + V^*(\delta(s, a))]$$

- Problem with Idea
 - Works well if agent knows
 - δ : state × action \rightarrow state
 - $r: state \times action \rightarrow R$
 - When agent doesn't know δ and *r*, cannot choose actions this way



Department of Computing and Information Sciences

Q Function

- Solution Approach
 - Define new function very similar to V^*

$$Q(s, a) \equiv r(s, a) + \gamma V * (\delta(s, a))$$

- If agent learns Q, it can choose optimal action even without knowing δ !
- Using Learned Q
 - Q: evaluation function to be learned by agent
 - Apply Q to select action
 - Idealized, computed policy (this is your brain without *Q*-learning):

$$\boldsymbol{\pi}^{\star}(\boldsymbol{s}) \equiv \arg \max_{\boldsymbol{a}} [r(\boldsymbol{s}, \boldsymbol{a}) + \boldsymbol{V}^{\star}(\boldsymbol{\delta}(\boldsymbol{s}, \boldsymbol{a}))]$$

• Approximated policy (this is your brain with *Q*-learning):

$$\pi^*(s) \equiv \arg \max_a Q(s, a)$$



Department of Computing and Information Sciences

Training Rule to Learn Q

- Developing Recurrence Equation for Q
 - Note: Q and V* closely related

$$V^{*}(s) \equiv arg \max_{a'} Q(s, a')$$

- Allows us to write **Q** recursively as

$$\begin{aligned} \boldsymbol{Q}(\boldsymbol{s}(t), \boldsymbol{a}(t)) &= \boldsymbol{r}(\boldsymbol{s}(t), \boldsymbol{a}(t)) + \boldsymbol{\gamma} \boldsymbol{V} \left(\boldsymbol{\delta}(\boldsymbol{s}(t), \boldsymbol{a}(t)) \right) \\ &= \boldsymbol{r}(\boldsymbol{s}(t), \boldsymbol{a}(t)) + \boldsymbol{\gamma} \max_{\boldsymbol{a}'} \boldsymbol{Q}(\boldsymbol{s}(t+1), \boldsymbol{a}') \end{aligned}$$

- Nice! Let \hat{Q} denote learner's current approximation to Q
- Training Rule

$$\hat{Q}(s, a) \leftarrow r(s, a) + \gamma \max_{a'} \hat{Q}(s', a')$$

- s': state resulting from applying action a in state s
- (Deterministic) transition function δ made implicit
- <u>Dynamic programming</u>: iterate over table of possible *a*' values



Q Learning for **Deterministic Worlds**

- (Nonterminating) Procedure for Situated Agent
- Procedure *Q-Learning-Deterministic* (*Reinforcement-Stream*)
 - *Reinforcement-Stream*: consists of <<<u>s</u>tate, <u>a</u>ction>, <u>r</u>eward> tuples
 - FOR each <s, a> DO
 - Initialize table entry $\hat{Q}(s, a) \leftarrow 0$
 - Observe current state s
 - WHILE (true) DO
 - Select action a and execute it
 - <u>Receive</u> immediate reward *r*
 - Observe new state s'
 - Update table entry for $\hat{Q}(s, a)$ as follows

$$\hat{Q}(s, a) \leftarrow r(s, a) + \gamma \max_{a'} \hat{Q}(s', a')$$

• <u>Move</u>: record transition from *s* to *s*'



Updating the **Q** Estimate



- Example: Propagating Credit (Q) for Candidate Action
 - Ø glyph: denotes mobile robot
 - Initial state: s₁ (upper left)
 - Let discount factor γ be 0.9
 - Qestimate

$$\hat{Q}(s_1, a_{right}) \leftarrow r(s, a) + \gamma \max_{a'} \hat{Q}(s_2, a')$$
$$\leftarrow 0 + 0.9 \cdot \max \{63, 81, 100\}$$
$$\leftarrow 90$$

- Property
 - If rewards nonnegative, then \hat{Q} increases monotonically between 0 and true Q
 - $r(s, a) \ge 0 \implies \forall s, a, n . \hat{Q}_{n+1}(s, a) \ge \hat{Q}_n(s, a) \land \forall s, a, n . 0 \le \hat{Q}_n(s, a) \le Q(s, a)$

Convergence

- Claim •
 - $-\hat{Q}$ converges to Q
 - Scenario: deterministic world where <s, a> are observed (visited) infinitely often
- Proof •
 - Define full interval: interval during which each < s, a> is visited
 - During each full interval, largest error in \hat{Q} table is reduced by factor of γ
 - Let \hat{Q}_n be table after *n* updates and Δ_n be the maximum error in \hat{Q}_n ; that is,

$$\boldsymbol{\Delta}_{n} = \max_{\mathbf{s}, \mathbf{a}} \left| \hat{\boldsymbol{Q}}_{n}(\mathbf{s}, \mathbf{a}) - \boldsymbol{Q}(\mathbf{s}, \mathbf{a}) \right|$$

- For any table entry $\hat{Q}_n(s, a)$, updated error in revised estimate $\hat{Q}_{n+1}(s, a)$ is $\left|\hat{Q}_{n+1}(s,a)-Q(s,a)\right| = \left|r+\gamma \max \hat{Q}_{n}(s',a')-(r+\gamma \max Q_{n}(s',a'))\right|$ $= \gamma \max \hat{Q}_n(s', a') - \max Q(s', a')$ $\leq \gamma \max \left| \hat{Q}_n(s',a') - Q(s',a') \right| \leq \gamma \max_{s'',s'} \left| \hat{Q}_n(s'',a') - Q(s'',a') \right|$ $\left|\hat{\boldsymbol{Q}}_{n+1}(\boldsymbol{s},\boldsymbol{a})-\boldsymbol{Q}(\boldsymbol{s},\boldsymbol{a})\right|\leq \boldsymbol{\gamma}\boldsymbol{\Delta}_{\boldsymbol{\mu}}$ $- \underline{\text{Note}}: \text{ used general fact, } \left| \max_{a} f_1(a) - \max_{a} f_2(a) \right| \le \max_{a} |f_1(a) - f_2(a)|$

CIS 732: Machine Learning and Pattern Recognition

Nondeterministic Case

- Second Learning Scenario: Nondeterministic World (Nondeterministic MDP)
 - What if reward and next state are <u>nondeterministically selected</u>?
 - i.e., reward function and transition function are nondeterministic
 - Nondeterminism may express many kinds of uncertainty
 - Inherent uncertainty in dynamics of world
 - Effector exceptions (qualifications), side effects (ramifications)
- Solution Approach
 - Redefine V, Q in terms of expected values

$$V^{\pi}(s) \equiv E[r(t) + \gamma r(t+1) + \gamma^{2} r(t+1) + \dots]$$
$$\equiv E\left[\sum_{i=0}^{\infty} \gamma^{i} r(t+i)\right]$$
$$Q(s, a) \equiv E[r(s, a) + \gamma V * (\delta(s, a))]$$

- Introduce decay factor; retain some of previous Q value
- Compare: momentum term in ANN learning



Nondeterministic Case: Generalizing *Q* Learning

- Q Learning Generalizes to Nondeterministic Worlds
 - Alter training rule to

$$\hat{\boldsymbol{Q}}_{n}(\boldsymbol{s},\boldsymbol{a}) \leftarrow (1 - \boldsymbol{\alpha}_{n}) \hat{\boldsymbol{Q}}_{n-1}(\boldsymbol{s},\boldsymbol{a}) + \boldsymbol{\alpha}_{n} \left[\boldsymbol{r} + \boldsymbol{\gamma} \max_{\boldsymbol{a}'} \hat{\boldsymbol{Q}}_{n-1}(\boldsymbol{s}',\boldsymbol{a}') \right]$$

- Decaying weighted average

$$\alpha_n = \frac{1}{1 + visits_n(s, a)}$$

- *visits_n*(*s*, *a*): total number of times $\langle s, a \rangle$ has been visited by iteration *n*, inclusive
- r: observed reward (may also be stochastically determined)
- Can Still Prove Convergence of \hat{Q} to Q [Watkins and Dayan, 1992]
- Intuitive Idea
 - $\alpha \in [0, 1]$: discounts estimate by number of visits, *prevents oscillation*
 - Tradeoff
 - More gradual revisions to Q
 Q
 - Able to deal with stochastic environment: P(s' | s, a), P(r | s, s', a)



Terminology

- <u>Reinforcement Learning (RL</u>)
 - RL: learning to choose optimal actions from <*state*, *reward*, *action*> observations
 - Scenarios
 - <u>Delayed reward</u>: reinforcement is deferred until end of <u>episode</u>
 - <u>Active learning</u>: agent can control collection of experience
 - <u>Partial observability</u>: may only be able to observe rewards (*must infer state*)
 - Reuse requirement: sensors, effectors may be required for multiple tasks
- <u>Markov Decision Processes (MDPs</u>)
 - Markovity (*aka* Markov property, Markov assumption): Cl assumption on states over time
 - <u>Maximum expected utility (MEU)</u>: maximum expected *total reward* (under additive decomposition assumption)
- Q Learning
 - Action-value function Q: state \times action \rightarrow value
 - Q learning: training rule and dynamic programming algorithm for RL



Summary Points

- Control Learning
 - Learning policies from < state, reward, action> observations
 - Objective: choose optimal actions given new percepts and incremental rewards
 - Issues
 - Delayed reward
 - Active learning opportunities
 - Partial observability
 - Reuse of sensors, effectors
- Q Learning
 - Action-value function Q: *state* \times *action* \rightarrow *value* (expected utility)
 - Training rule
 - Dynamic programming algorithm
 - Q learning for deterministic worlds
 - Convergence to true Q
 - Generalizing *Q* learning to nondeterministic worlds
- Next Week: More Reinforcement Learning (Temporal Differences)

