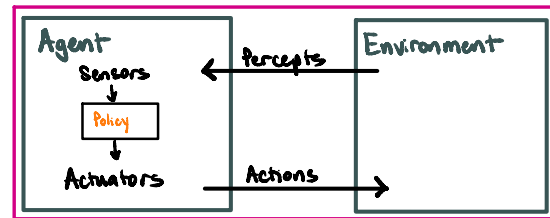


## 1 Introduction

- An agent perceives its environment through sensors and acts upon it through actuators
- A rational agent chooses action to max expected value



## 2 Search

Search Problems have:

- State Space  $S$
- Initial State  $S_0$
- Goal test  $G(s)$
- Action cost  $C(s, a, s')$
- Transition model  $T(s, a)$
- Actions of state  $A(s)$

Solutions: action sequence that reaches goal state

↳ Optimal Solution: least cost among solutions

### Depth First Search (DFS)

Strat: expand deepest node first

Use: LIFO stack

Optimal?: No, finds "leftmost" sol

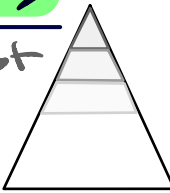


### Breadth First Search (BFS)

Strat: expand shallowest node first

Use: FIFO Queue

Optimal: IF costs are equal

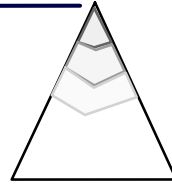


### Uniform Cost Search (UCS)

Strat: expand lowest cost from root  $g(n)$

Use: Priority Queue by cost

Optimal?: Yes

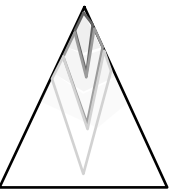


### A\* Search

Strat: expand node most likely on optimal path using heuristic

Use: Priority Queue by  $g(n) + h(n)$

Optimal: Yes, if heuristic admissible & consistent



Admissible:  $0 \leq h(n) \leq h^*(n)$

heuristic  $\leq$  actual cost  $\uparrow$  true cost to goal

Consistent:  $h(A) - h(C) \leq c(A, C)$

heuristic "arc" cost  $\leq$  actual arc cost

Graph Search: tracks explored nodes on list  $\uparrow$  prevent revisiting nodes

Tree Search: map paths w/ tree, node appear multiple times  $\uparrow$  less memory

### Greedy Search

Strategy: Expand node w/ lowest heuristic value, believed closest to goal

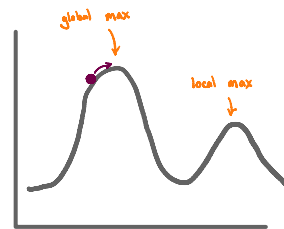
Use: priority queue with heuristics, estimated forward cost

Optimal: No, not guaranteed to find goal state, unpredictable

# Local Search

## Hill Climbing

Idea: Start wherever, repeat: move to the best neighboring state. If no neighbors better than current, quit



```
function HILL-CLIMBING(problem) returns state
  current ← make-node(problem.initial-state)
  loop do
    neighbor ← highest valued successor of current
    if neighbor.value ≤ current.value then
      return current.state
    current ← neighbor
```

Random-restart hill climbing  
trivially complete  
- restarts from random initial state

## Simulated Annealing

Idea: random walk + hill climbing, choose randomly, choose worse w/ some probability according to temperature, temp starts high w/ more "bad" moves allowed and decreases

```
function SIMULATED-ANNEALING(problem, schedule) returns state
  current ← problem.initial-state
  for t=1 to ∞ do
    T ← schedule(t)
    if T=0 then return current
    next ← a randomly selected successor of current
    ΔE ← next.value - current.value
    if ΔE > 0 then current ← next
    else current ← next w/ prob  $e^{-\Delta E/T}$ 
```

If T decreased slowly enough, will converge to optimal state

## Local Beam Search

multiple interaction searches

Idea: track k states at each iteration, k threads share information good threads attract other threads, chooses k best successors

## Genetic Algorithms

Break and recombine states

Idea: beam search k states in population, each state evaluated w/ evaluation function (fitness func), offspring produced by crossing parent strings at crossover point

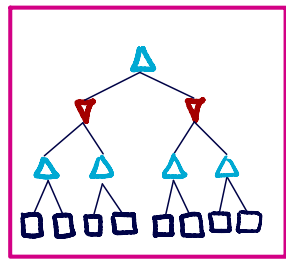
### 3 Games

Zero sum games: our gain is directly equivalent to opponent's loss

#### Minimax algorithm to counter optimal opponent moves

Idea: Zero sum also assuming opponent plays optimally

terminal utilities/state value: optimal score attainable



function minimax-decision(s) returns action

return action a in Actions(s) w/ highest minimax-value (Result(s,a))

function minimax\_value(s) returns a value

if Terminal-Test(s) then return Utility(s)

if Player(s) = MAX then return  $\max_{a \in \text{Actions}(s)} \text{minimax\_value}(\text{Result}(s,a))$

if Player(s) = MIN then return  $\min_{a \in \text{Actions}(s)} \text{minimax\_value}(\text{Result}(s,a))$

#### Alpha-Beta Pruning

optimization to minimax:  $O(b^m) \Rightarrow O(b^{m/2})$

Idea: Stop looking as soon as you know n's value can at best equal optimal value of n's parent

$\alpha$ : MAX's best option on path to root

$\beta$ : MIN's best option on path to root

```
def max-value(state,  $\alpha$ ,  $\beta$ )
  initialize  $v = -\infty$ 
  for each successor of state:
     $v = \max(v, \text{value}(\text{successor}, \alpha, \beta))$ 
    if  $v \geq \beta$  return  $v$ 
     $\alpha = \max(\alpha, v)$ 
  return  $v$ 
```

```
def min-value(state,  $\alpha$ ,  $\beta$ )
  initialize  $v = +\infty$ 
  for each successor of state:
     $v = \min(v, \text{value}(\text{successor}, \alpha, \beta))$ 
    if  $v \leq \alpha$  return  $v$ 
     $\beta = \min(\beta, v)$ 
  return  $v$ 
```

Evaluation functions are used in depth-limited minimax to output an estimate of the true minimax value of the node

$$\text{Eval}(s) = w_1 f_1(s) + w_2 f_2(s) + \dots + w_n f_n(s)$$

weight associated w/ feature      features of state

# Expectimax

algorithm with randomness using Expected value

Idea: Chance nodes calculated by taking expected utility of children

function decision(s) returns action

return action a in Actions(s) w/ highest value (Result(s,a))

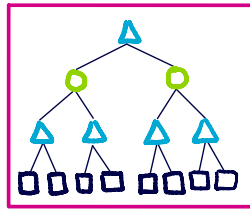
function value(s) returns a value

if Terminal-Test(s) then return Utility(s)

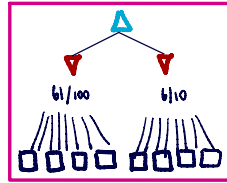
if Player(s) = MAX then return  $\max_{a \in \text{Actions}(s)} \text{value}(\text{Result}(s,a))$

if Player(s) = MIN then return  $\min_{a \in \text{Actions}(s)} \text{value}(\text{Result}(s,a))$

if Player(s) = CHANCE then return  $\sum_{a \in \text{Actions}(s)} \text{Pr}(a) * \text{value}(\text{Result}(s,a))$



# Monte Carlo Tree Search (MCTS)



Idea:

1) Evaluation by rollouts: From state s play many times using policy and count wins and losses

2) Selective search: explore parts of the tree, without constraints on horizon, that will improve decision at root

# UCB Algorithm

$$UCB1(n) = \frac{U(n)}{N(n)} + C \sqrt{\frac{\log N(\text{PARENT}(n))}{N(n)}}$$

↖ total wins for Player(Parent(n))  
↗ user specific param to balance exploration vs. exploitation  
↖ total rollouts from n

# 4 Logic

Maintain a knowledge base w/ logical sentences that can make logical inferences

# Logic Syntax

Symbol	Meaning	Description
$\neg$	not	$\neg P$ true iff P is false
$\wedge$	and (conjunction)	$A \wedge B$ true iff both A true and B true
$\vee$	or (disjunction)	$A \vee B$ true iff either A true or B true
$\Rightarrow$	implication	$A \Rightarrow B$ true unless A true and B false
$\Leftrightarrow, \equiv$	biconditional	$A \Leftrightarrow B (A \equiv B)$ true iff either both A, B true or false

**Conjunctive normal form (CNF):** conjunction of clauses that are disjunction of literals  
 $(P_1 \vee \dots \vee P_i) \wedge \dots \wedge (P_j \vee \dots \vee P_n)$

Converting to CNF:

- 1) Eliminate  $\Leftrightarrow$
- 2) Eliminate  $\Rightarrow$
- 3) Not's ( $\neg$ ) only on literals
- 4) Reduce and distribute

$$\begin{aligned}
 (A \Leftrightarrow B) &\equiv ((A \Rightarrow B) \wedge (B \Rightarrow A)) \\
 A \Rightarrow B &\equiv \neg A \vee B \\
 \neg(A \wedge B) &\equiv (\neg A \vee \neg B), \neg(A \vee B) \equiv (\neg A \wedge \neg B) \\
 (A \wedge (B \vee C)) &\equiv ((A \wedge B) \vee (A \wedge C)) \\
 (A \vee (B \wedge C)) &\equiv ((A \vee B) \wedge (A \vee C))
 \end{aligned}$$

**model:** assignment of true/false to all proposition symbols  
**valid:** true in all models  
**satisfiable:** at least one model where true  
**unsatisfiable:** not true in any models

First Order Logic (FOL)

uses objects as base components

Use function symbols to name objects in terms

- quantifiers:
- 1) universal quantifier  $\forall$ : "for all"
  - 2) existential quantifier  $\exists$ : "there exists"

Operator Precedence  
 $\neg, =, \wedge, \vee, \Rightarrow, \Leftrightarrow$

Propositional Logical Inference

**Entails ( $\models$ ):** in all models where A is true, B is as well,  $M(A) \subseteq M(B)$

Show by:

- 1)  $A \models B$  iff  $A \Rightarrow B$  valid direct proof
- 2)  $A \not\models B$  iff  $A \wedge \neg B$  unsatisfiable proof by contradiction

Model Checking

DPLL Algorithm

depth-first, back-tracking search w/ tricks SAT solver

**Idea:** Solve satisfiability problem given in CNF, continue assigning truth values until found or cannot be found, backtrack w/ conditions:

- 1) **Early Termination:** clause true if any symbol true, sentence false if any single clause is false
- 2) **Pure Symbol Heuristic:** symbol only shows up in pos/neg form assigned immediately  $A = \text{true}$   
 $(A \vee B)(\neg B \vee C) \wedge (\neg C \vee A)$
- 3) **Unit Clause Heuristic:** clause w/ one literal or literal w/ many falses

# Theorem Proving

## Resolution Algorithm

Idea: iteratively apply to knowledge base either  $\perp$  inferred or nothing left to infer

## Forward Chaining Algorithm

Idea: data driven reasoning, iterating through every implication statement where LHS known, adding RHS to facts using Generalized Modus Ponens

Solve inference in FOL by Generalized Modus Ponens or propositionalization translating problem into propositional logic and using SAT solver

## 5 Probability

$\perp$ : Conditional independence

Marginal distribution:

$$P(A, B) = \sum_c P(A, B, C)$$

Bayes' rule:

$$P(A|B, C) = \frac{P(A, B, C)}{P(B, C)} = \frac{P(A, B|C)}{P(B|C)}$$

$$P(A|B) = \frac{P(B|A)P(A)}{P(B)}$$

Chain Rule:

$$P(A, B, C) = P(A, B|C)P(C) = P(A|B, C)P(B|C)P(C)$$

A conditionally independent given C:

$$P(A, B|C) = P(A|C)P(B|C)$$

A independent of B given C:

$$P(A|B, C) = P(A|C)$$

A, B independent:

$$P(A, B) = P(A)P(B)$$

Query Variables

unknown, trying to compute

$$P(Q_1, \dots, Q_k | e_1, \dots, e_k)$$

Evidence variables observed, values known

Hidden Variables

values present but not in distribution we are computing

## Bayes Nets

Idea: Joint distribution distributed across smaller probability tables with Directed Acyclic Graph (DAG)

- 1) DAG w/ node per variable  $X$
- 2) conditional distribution for each node  $P(X | A_1, \dots, A_n)$  where  $A_i$  is  $i^{\text{th}}$  parent of  $X$ , stored as conditional probability table (CPT)

Each node is conditionally independent of all ancestor nodes in graph, given all of parents

Bayes Net Inference (calculating joint probability)

Eliminate vars by:

1. Join (multiply together) all factors involving  $X$
2. Sum out  $X$

$$P(X=1 | A=+, B=+) = P(X=1) P(A=+ | X=1) P(B=+ | X=1)$$

Prior Sampling: Randomly generate samples

Rejection Sampling: early reject any sample inconsistent w/ evidence

Likelihood weighting: Ensure we never generate a bad sample

- 1) manually set variables to equal evidence
- 2) weight each sample by probability of the evidence variables given the sampled variables

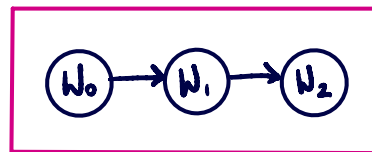
Gibbs Sampling

set all variables to random value, repeatedly pick one variable at a time, clear value, resample

## 6 Markov Models

Chain like infinite-length Bayes Net

Need to know initial state and transition model



Mini-Forward Algorithm

$$Pr(W_{i+1}) = \sum_{w_i} Pr(W_{i+1} | w_i) Pr(w_i)$$

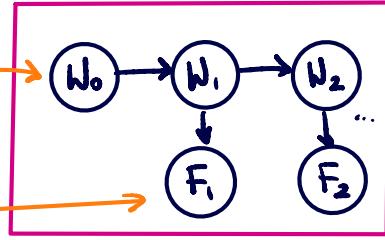
Stationary Distribution

$$Pr(W_{t+1}) = Pr(W_t) = \sum_{w_t} Pr(W_{t+1} | w_t) Pr(w_t)$$

# Hidden Markov Models

allows to observe some evidence at timestep and affects belief distribution

State variable encodes belief evidence var



## Forward Algorithm

$$B'(W_{t+1}) = \sum_{w_i} \Pr(W_{t+1} | w_i) B(w_i)$$

Time elapse Update

$$P(x_t | o_{0:t-1}) = \sum_{x_{t-1}} P(x_t | x_{t-1}) P(x_{t-1} | o_{0:t-1})$$

$o$ : observation

$$B(W_{t+1}) \propto \Pr(f_{t+1} | W_{t+1}) B'(W_{t+1})$$

Observation update

$$B(W_{t+1}) \propto \Pr(f_{t+1} | W_{t+1}) \sum_{w_i} \Pr(W_{t+1} | w_i) B(w_i)$$

$$P(x_t | o_{0:t}) \propto P(x_t, o_t | o_{0:t-1}) = P(o_t | x_t) P(x_t | o_{0:t-1})$$

## Viterbi Algorithm

Dynamic Programming algorithm to solve

$$\operatorname{argmax}_{x_{1:N}} P(x_{1:N} | e_{1:N}) = \operatorname{argmax}_{x_{1:N}} P(x_{1:N}, e_{1:N})$$

$$m_t[x_t] = P(e_t | x_t) m_{t-1} P(x_t | x_{t-1}) m_{t+1}[x_{t-1}]$$



# 1 Dynamic Bayes Nets track multiple variables over time, using multiple source evidence

Idea: Repeat fixed Bayes net structure at each time

**Particle Filtering** simulating motion of set of particles through state graph to approx prob distr

1) Prediction Step: sample new state from transition model

particle  $j$  state  $x_t^{(j)}$   $x_{t+1}^{(j)} \sim P(x_{t+1} | x_t^{(j)})$

2) Update Step: weight each sample on evidence, then normalize

$w_t^{(j)} = P(e_{t+1} | x_t^{(j)})$

**Exact Inference**: Calc exact probability, more accurate but hard w/ large belief distribution

# 2 Rational Decisions

**Principle of Maximum Expected Utility (MEU)**: rational agents always select action to max utility

Preferences:

$A \succ B$  : prefers A over B       $A \sim B$  : indifferent between A and B

**Lottery**: A received w/ probability  $p$ , B w/ probability  $1-p$

$L = [p, A; (1-p), B]$

**Axioms of Rationality**:

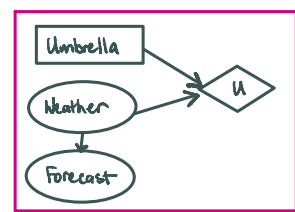
- 1) **Orderability**: must prefer A or B, or be indifferent
- 2) **Transitivity**: if prefers A to B and B to C, prefers A to C
- 3) **Continuity**: if prefers A to B, B to C, lottery L w/ A, C possible st. indifferent between L and B for some  $p$
- 4) **Substitutability**: indifferent between A and B also indifferent for lottery substitutions
- 5) **Monotonicity**: prefers A over B, chooses lottery w/ highest A prob

$(A \succ B) \vee (B \succ A) \vee (A \sim B)$   
 $(A \succ B) \wedge (B \succ C) \Rightarrow (A \succ C)$   
 $A \succ B \succ C \Rightarrow \exists p [p, A; (1-p), C] \sim B$   
 $A \sim B \Rightarrow [p, A; (1-p), C] \sim [p, B; (1-p), C]$

$A \succ B \Rightarrow (p \geq q \Leftrightarrow [p, A; (1-p), B] \succeq [q, A; (1-q), B])$

# 3 Decision Networks used to model state, actions, utilities of subsequent states

- Chance Nodes**: outcome has associated probability
- Action Nodes**: choices between actions we can choose
- Utility Nodes**: children of chance & action representing utility



**Value of Perfect Information (VPI)** expected improvement in decision quality from observing value

$VPI(E'|e)$  is value of observing new evidence  $E'$  given current evidence  $e$

$$VPI(E'|e) = MEU(e, E') - MEU(e) = \left[ \sum_{e'} P(e'|e) \max_a \sum_s P(s|e, e') U(s, a) \right] - \max_a \sum_s P(s|e) U(s, a)$$

**VPI Properties**

- 1) **Nonnegativity**: more informed decision
- 2) **Nonadditivity**: can change how much we care about  $E_k$
- 3) **Order-independence**: order of observation doesn't matter

$\forall E', e, VPI(E'|e) \geq 0$   
 $VPI(E_j, E_k | e) \neq VPI(E_j | e) + VPI(E_k | e)$   
 $VPI(E_i, E_j | e) = VPI(E_j, E_i | e)$

## 4) Markov Decision Process (MDP)

solving nondeterministic search problems

Markov Decision Process properties:

- State space  $S$
- actions  $A$
- Transition function  $T(s, a, s')$  probability taking action  $a$  at state  $s$  results in  $s'$
- Reward function  $R(s, a, s')$  pos/neg reward for each step
- terminal state  $(s)$
- Start state  $S_0$
- discount factor  $\gamma$  multiply to reward for exponential decay in reward "time constraint"
- $P(s'|s, a)$

Policy  $\pi: S \rightarrow A$ , mapping each state to an action

State Optimal Value  $U^*(s)$ : expected utility of optimal agent from state  $s$

Optimal Q Value  $Q^*(s, a)$ : expected utility of optimal agent from state  $s$  taking action  $a$

Policy Extraction: used to determine policy given some state value function

### Value Iteration

DP algorithm to compute MDP until convergence of values

1)  $\forall s \in S$ , set  $U_0(s) = 0$

2) Repeat until convergence:

$$\forall s \in S, U_{k+1}(s) \leftarrow \max_a \sum_{s'} T(s, a, s') [R(s, a, s') + \gamma U_k(s')]$$

Runtime:  $O(|S|^2|A|)$

### Policy Iteration

more efficient method to converge to optimal policy

1) Define an initial policy

2) Use policy evaluation for current policy

3) Use policy improvement to find better policy

$$U^\pi(s) = \sum_{s'} T(s, \pi(s), s') [R(s, \pi(s), s') + \gamma U^\pi(s')]$$

$$\pi_{k+1}(s) = \operatorname{argmax}_a \sum_{s'} T(s, a, s') [R(s, a, s') + \gamma U^k(s')]$$

## 5) Machine Learning

learning parameters of specified model given data

### Decision Trees

simple non-linear classifier but susceptible to over-fitting

Entropy: uncertainty about variable  $\uparrow$  entropy  $\uparrow$  uncertainty

entropy of binary variable

$$H(V) = - \sum_k P(v_k) \log_2 P(v_k)$$

$$B(q) = -q \log_2 q - (1-q) \log_2 (1-q)$$

Information Gain: info gained by splitting on that feature

$$\text{Gain}(\text{Type}) = B\left(\frac{P}{p+n}\right) - \sum_{k=1}^K \frac{P_k + n_k}{p+n} B\left(\frac{P_k}{P_k + n_k}\right)$$

### Linear Regression

linear classifier for continuous var output

L2 Loss function: metric to measure our model using squared difference  $(y - h(x))^2$

Least squares:

$$\hat{W} = (X^T X)^{-1} X^T y$$

### Logistic Regression

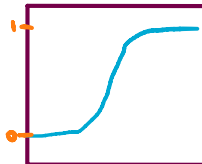
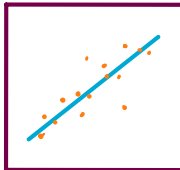
linear classifier for categorical variables

Logistic function:

$$h_w(x) = \frac{1}{1 + e^{-wx}}$$

Find optimal using gradients

Models  $P(c|x)$  where  $c$ : class,  $x$ : evidence



# Perceptron

Linear classifier: classification using linear combination of features

Activation:  $activation_w(x) = h_w(x) = w \cdot f(x)$

$$classify(x) = \begin{cases} + & \text{if } h_w(x) > 0 \\ - & \text{if } h_w(x) < 0 \end{cases}$$

## Perceptron Algorithm:

- 1) Initialize all weights to 0:  $w = 0$
- 2) Classify  $y$  and compare to true label  $y^*$ 
  - if  $y = y^*$ , do nothing
  - if  $y \neq y^*$ , update  $w \leftarrow w + y^* f(x)$

# Naive Bayes

model features as Bayes Net, assumes each feature is independent from others

Prediction for class label becomes  $prediction(F) = \underset{y_i}{\operatorname{argmax}} P(Y=y_i) \prod_j P(F_j=f_j | Y=y_i)$ , normalize

Maximum Likelihood Estimation (MLE): given i.i.d., method to learn probability, count

Likelihood:  $\mathcal{L}(\theta) = \prod_{i=1}^n P_{\theta}(x_i)$  take gradient  $\frac{\partial}{\partial \theta} \mathcal{L}(\theta) = 0$  to get max, log likelihood  $\log \mathcal{L}(\theta)$

Laplace smoothing: mitigate overfitting, assume seen  $k$  extra each of each outcome

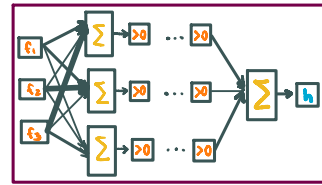
$$P_{lap,k}(x) = \frac{\text{count}(x) + k}{N + k |X|}$$

# Neural Networks

Multi-layer perceptron map data to higher dimension then classify

Universal func approximator: two layer NN w/ sufficient neurons can approx any cont. func  
can use indicator function of sgn threshold, softmax func to classify

Log Likelihood:  $\log \mathcal{L}(w) = \prod_{i=1}^n \log P(y_i | f(x_i); w)$



## Activation functions

sigmoid:  $\sigma(x) = \frac{1}{1+e^{-x}}$



ReLU:  $f(x) = \begin{cases} 0 & \text{if } x < 0 \\ x & \text{if } x \geq 0 \end{cases}$



## Backpropagation

determine gradient of output w/ respect to each of inputs

- 1) Forward Pass: compute values through computation graphs
- 2) Backwards Pass: compute gradients taking advantage of chain rule

# 6 Reinforcement Learning

method for solving MDP w/o transition & reward functions

Episode: collection of samples which are  $(s, a, s', r)$  tuples

Model-based learning: estimate transition, reward functions w/ samples before using policy/value iter  
- count times arrived in state and normalize counts, Law of Large Numbers will converge on optimal

Model-free learning: estimate Q values directly w/o constructing model

↳ Passive reinforcement learning: agent given policy to follow and learns values of state exploring

- 1) Direct Evaluation: fix policy  $\pi$  and have agent experience, can compute by averages
- 2) Temporal Difference (TD) learning: learning from every experience

$$\text{sample} = R(s, \pi(s), s') + \gamma V^{\pi}(s')$$

$$V_k^{\pi}(s) \leftarrow (1-\alpha) V_{k-1}^{\pi}(s) + \alpha \cdot \text{sample}$$

↳ Active Reinforcement learning: learning agent can use feedback to iteratively update policy

3) Q-Learning: Bypass model by directly learning Q values

Q-Value Iteration:

$$Q_{k+1}(s,a) \leftarrow \sum_{s'} T(s,a,s') [R(s,a,s') + \gamma \max_{a'} Q_k(s',a')]$$

Q value samples:

$$\text{sample} = R(s,a,s') + \gamma \max_{a'} Q(s',a')$$

$$Q(s,a) \leftarrow (1-\alpha) Q(s,a) + \alpha \cdot \text{sample}$$

Feature Based Representation: allow model to generalize learning experiences, store as linear func

Weight update

$$w_i \leftarrow w_i - \alpha \cdot \text{sample} \frac{\partial Q_w}{\partial w_i}$$

negative for gradient descent

Exploration vs. Exploitation

$\epsilon$ -greedy:  $0 \leq \epsilon \leq 1$ , act randomly w/ prob  $\epsilon$ , should be lowered over time

Exploration functions:

modified update

$$Q(s,a) = (1-\alpha) Q(s,a) + \alpha [R(s,a,s') + \gamma \max_{a'} Q(s',a')]$$

$$f(s,a) = Q(s,a) + \frac{k}{N(s,a)}$$

← predetermined

← # times Q state visited

Regret: metric to measure model, difference between total reward for optimal and reward in our model