

Part 1: Search

Value agent: based on current state; planning agent: consequence of actions; choose best
A search problem is consisted of ① A state space, ② A set of actions

③ a transition model ④ An action cost ⑤ Start state ⑥ Goal States

	property	Optimal tree Graph	Complete tree Graph	Time	Space
BFS	FIFO queue Cost=1	✓	✓	bs	bs
DFS	LIFO stack each time search depth = 1, 2, ... uniform cost search	X	X	dm	bm
A* search	fun = g(n) + h(n)	✓	X	b ^{costs}	b ^{h(n)}

- ① admissibility: $\forall n \quad 0 \leq h(n) \leq h^*(n)$
- ② consistency: $\forall A, C, H(n) \leq cost(A, C)$.
- ③ dominate: $\forall n \quad h(n) \geq h(b(n))$.

Theorem 1: if h is admissible, then A* yields an optimal solution in tree search.

Theorem 2: if h is consistent, then A* yields an optimal solution in Graph search.

greedy search: select the one with lowest heuristic value to expand.
not complete & optimal.

Local Search:

→ not L&O

① Hill climbing search: best neighbor

② Simulated Annealing Search:

function SIMULATED_ANNEALING(problem, schedule) returns a 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

$\Delta E \leftarrow next.value - current.value$

 if $\Delta E > 0$ then current ← next $P = 1$

 else current ← next only with probability $e^{\Delta E / T}$



③ Local Beam search: best K result each time

④ Genetic Algorithm:

function GENETIC-ALGORITHM(population, FITNESS-FN) returns an individual
inputs: population, a set of individuals
 FITNESS-FN, a function that measures the fitness of an individual

repeat

 new_population ← empty set

 for i = 1 to SIZE(population) do

 x ← RANDOM-SELECTION(population, FITNESS-FN)

 y ← RANDOM-SELECTION(population, FITNESS-FN)

 child ← REPRODUCE(x, y)

 if (small random probability) then child ← MUTATE(child)

 add child to new_population

 population ← new_population

until some individual is fit enough, or enough time has elapsed

return the best individual in population, according to FITNESS-FN

function REPRODUCE(x, y) returns an individual

inputs: x, y, parent individuals

$n \leftarrow LENGTH(x); c \leftarrow$ random number from 1 to n

return APPEND(SUBSTRING(x, 1, c), SUBSTRING(y, c + 1, n))

2. Game → zero-sum game. adversarial search problems

① Min-Max Tree

② Alpha-Beta Pruning: $O(CB^{\frac{M}{2}})$ [last adversarial game can use]
α: MAX's best option on path to root
β: MIN's best option on path to root

```
def max-value(state, α, β):
    initialize v = -∞
    for each successor of state:
        v = max(v, value(successor, α, β))
        if v ≥ β return v
        α = max(α, v)
    return v
```

```
def min-value(state, α, β):
    initialize v = +∞
    for each successor of state:
        v = min(v, value(successor, α, β))
        if v ≤ α return v
        β = min(β, v)
    return v
```

Max Pruning leaf = $(b-1)(M-1)$; remain the left branch of each branch from root
number leaves ≠ possible values = # leaf node - max pruning node.

③ Exploit Max: cut from V, cut K largest cut from A: cut K smallest

④ Monte Carlo Tree Search.

while

$$S1 \text{ Tree traversal}(UCB_1(c_n)) = \frac{u(c_n)}{n(c_n)} + C \times \sqrt{\frac{\log N(PARENT(c_n))}{n(c_n)}}$$

S2 Node Expansion.

S3 Simulation (Rollout)

S4 Backpropagation

3 Logic

valid: true for all models

satisfiable: \exists model. true

unsatisfiable: not true for all models

conjunctive normal form: $(P_1 \vee \dots \vee P_i) \wedge \dots \wedge (Q_j \vee \dots \vee Q_k)$

entail: $A \models B$.

① $A \models B$ iff $A \Rightarrow B$ is valid ② $A \models B$ iff $A \wedge \neg B$ is unsatisfiable.

Model checking

DPLL → is a variant of DFS

function DPLL-SATISFIABLE(s) returns true or false sound & correct & complete

inputs: s , a sentence in propositional logic

clauses ← the set of clauses in the CNF representation of s

symbols ← a list of the proposition symbols in s

return DPLL(clauses, symbols, {})

function DPLL(clauses, symbols, model) returns true or false

if every clause in clauses is true in model then return true Early Termination. (one symbol T)
if some clause in clauses is false in model then return false
 $P, value \leftarrow \text{FIND-PURE-SYMBOL}(symbols, clauses, model)$ Pure Symbol Heuristic (all A or $\neg A$ and assign as Y)
if P is non-null then return DPLL(clauses, symbols - P , model $\cup \{P=value\}$) and assign as Y
 $P, value \leftarrow \text{FIND-UNIT-CLAUSE}(clauses, model)$ Unit Clause Heuristic Eg. $B \vee F \vee \neg F$ only visit in clause
if P is non-null then return DPLL(clauses, symbols - P , model $\cup \{P=value\}$) only visit in clause
 $P \leftarrow \text{FIRST}(symbols); rest \leftarrow \text{REST}(symbols)$
return DPLL(clauses, rest, model $\cup \{P=true\}$) or
DPLL(clauses, rest, model $\cup \{P=false\}$)

three rules of inference

① Modus Ponens: if knowledge base contain $A \wedge A \Rightarrow B$, we can infer B

② And-Elimination: if knowledge base contain $A \wedge B$. we can infer A and B .

③ Resolution: if knowledge base contain A and B . we can infer $A \wedge B$

Forward chaining KB?

function PL-FC-ENTAILS?(KB, q) returns true or false

inputs: KB, the knowledge base, a set of propositional definite clauses

q, the query, a proposition symbol

count ← a table, where $count[c]$ is the number of symbols in c 's premise

inferred ← a table, where $inferred[s]$ is initially false for all symbols

agenda ← a queue of symbols, initially symbols known to be true in KB

while agenda is not empty do

$p \leftarrow \text{POP}(agenda)$

 if $p = q$ then return true

 if $inferred[p] = \text{false}$ then

$inferred[p] \leftarrow \text{true}$

 for each clause c in KB where p is in $c.\text{PREMISE}$ do

 decrement $count[c]$

 if $count[c] = 0$ then add $c.\text{CONCLUSION}$ to agenda

return false

4 Bayes Net (tree-like net usually not holds ??))

mutually independent $A \perp \!\!\! \perp B$: $P(A, B) = P(A) \cdot P(B)$

conditional independent $B \perp \!\!\! \perp A | C$: $P(A, B | C) = P(A | C) \cdot P(B | C)$; $P(C | A, B) = P(C | B)$

$P(Q_1, \dots, Q_k | E_1, \dots, E_l)$: Q_i : Query variable E_i : Evidence variable. C_i : Hidden variable.

④ Infer by Enumeration.

1. Collect all the rows consistent with the observed evidence variables.

2. Sum out (marginalize) all the hidden variables.

3. Normalize the table so that it is a probability distribution (i.e. values sum to 1)

Zero-sum game

if one takes suboptimal other than opt. the opponent's utility ↑

Bayes Network Representation. CDAG /

② Variable Elimination. O(e^n)

```

function ELIMINATION-ASK( $X, e, bn$ ) returns a distribution over  $X$ 
  inputs:  $X$ , the query variable
     $e$ , observed values for variables  $E$ 
     $bn$ , a Bayesian network specifying joint distribution  $P(X_1, \dots, X_n)$ 
  unnormalized  $P$ 
  factors  $\leftarrow []$ 
  for each var in ORDER( $bn.VARS$ ) do
    factors  $\leftarrow [MAKE-FACTOR(var, e)|factors]$ 
    if var is a hidden variable then factors  $\leftarrow \text{SUM-OUT}(var, factors)$ 
  return NORMALIZE(POINTWISE-PRODUCT(factors))

```

order of elimination: best \rightarrow NP-hard problem:

- ① Maybe good? ② best out-neighbour ③ min-weight ④ minfill: min size of factor

Approximate.

$$\text{Prior Sampling} \quad \text{random. } P(\text{sample}) = \prod_i P(\text{initial Parent}(x_i))$$

④ Rejection Sampling: consistent with given info? reject it!

⑤ Likelihood weight (evidence $\rightarrow P$, Query \rightarrow Sampling) | consistent evidence

```

function LIKELIHOOD-WEIGHTING( $X, e, bn, N$ ) returns an estimate of  $P(X|e)$ 
  inputs:  $X$ , the query variable
     $e$ , observed values for variables  $E$ 
     $bn$ , a Bayesian network specifying joint distribution  $P(X_1, \dots, X_n)$ 
     $N$ , the total number of samples to be generated
  local variables:  $W$ , a vector of weighted counts for each value of  $X$ , initially zero
  for  $j = 1$  to  $N$  do
     $x, w \leftarrow \text{WEIGHTED-SAMPLE}(bn, e)$ 
     $W[x] \leftarrow W[x] + w$  where  $x$  is the value of  $X$  in  $x$ 
  return NORMALIZE( $W$ )

```

```

function WEIGHTED-SAMPLE( $bn, e$ ) returns an event and a weight
   $w \leftarrow 1; x \leftarrow$  an event with  $n$  elements initialized from  $e$ 
  foreach variable  $X_i$  in  $X_1, \dots, X_n$  do
    if  $X_i$  is an evidence variable with value  $x_i$  in  $e$ 
      then  $w \leftarrow w \times P(X_i = x_i | \text{parents}(X_i))$ 
    else  $x[i] \leftarrow$  a random sample from  $P(X_i | \text{parents}(X_i))$ 
  return  $x, w$ 

```

Gibbs Sampling. (consistent) \leadsto evidence

```

function GIBBS-ASK( $X, e, bn, N$ ) returns an estimate of  $P(X|e)$ 
  local variables:  $N$ , a vector of counts for each value of  $X$ , initially zero
     $Z$ , the non-evidence variables in  $bn$ 
     $x$ , the current state of the network, initially copied from  $e$ 
  initialize  $x$  with random values for the variables in  $Z$ 
  for  $j = 1$  to  $N$  do
    for each  $Z_i$  in  $Z$  do
      set the value of  $Z_i$  in  $x$  by sampling from  $P(Z_i | mb(Z_i))$ 
       $N[x] \leftarrow N[x] + 1$  where  $x$  is the value of  $X$  in  $x$ 
    return NORMALIZE( $N$ )

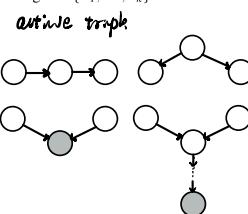
```

Type of Environment In AI

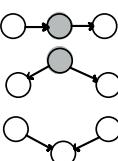
- ① Fully observable / Partially Observable: sense the complete state at each point time
- ② Static / dynamic: the environment keep change?
- ③ Stochastic / deterministic: a uniqueness in the agent's current state determine next state
- ④ know physics: agent know the transition model.

Conditional independence.

1. Shade all observed nodes $\{Z_1, \dots, Z_k\}$ in the graph.
2. Enumerate all undirected paths from X to Y .
3. For each path:
 - (a) Decompose the path into triples (segments of 3 nodes).
 - (b) If all triples are active, this path is active and d-connects X to Y .
4. If no path d-connects X and Y , then X and Y are d-separated, so they are conditionally independent given $\{Z_1, \dots, Z_k\}$



active triple



inactive triple

Markov Models

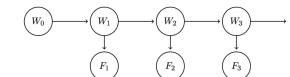
Memoryless / Markov Property: $W_0 \perp\!\!\!\perp \{w_1, \dots, w_{t-1}\} | w_t$

Stationary Assumption: for all values of i : $P(W_{i+1}|W_i)$ is identical.

The Mini-Forwarding Algorithm: $P(W_{i+1}) = \sum_{w_i} P(W_{i+1}|w_i) \cdot P(w_i)$

we must compute the stationary distribution: $P(W_{t+1}) = P(w_t)$

Hidden Markov Models:



Markov Property: ① $F_t \perp\!\!\!\perp W_0 | W_1 : F_t$ only depends on W_t

② $\forall i=2, \dots, n, W_i \perp\!\!\!\perp \{w_0, \dots, w_{i-1}, F_1, \dots, F_{i-1}\} | W_i$

Assumptions: ① $P(w_i | w_i)$ is stationary ② sensor model $P(F_i | W_i)$ is stationary.

Forwarding Algorithm: $B'(W_{i+1}) = P(W_{i+1}|f_1, \dots, f_i) = \sum_{w_i} P(W_{i+1}, w_i | f_1, \dots, f_i)$

$$B(W_i) = P(W_i | f_1, \dots, f_i)$$

① time elapse update: $B'(W_{i+1}) = \sum_{w_i} P(W_{i+1}|w_i)B(w_i)$

② observation update: $B(W_{i+1}) \propto P(f_{i+1} | W_{i+1})B'(W_{i+1}) = P(f_{i+1} | W_{i+1}) \sum_{w_i} P(W_{i+1}|w_i)B(w_i)$

Viterbi Algorithm: $m_t(x_t) = \max_{x_{t+1}} P(x_{t+1}, e_{1:t})$

$$= \max_{x_{t+1}} P(e_{t+1}|x_t) P(x_t|x_{t-1}) P(x_{t-1}, e_{1:t-1})$$

$$= P(e_{t+1}|x_t) \max_{x_{t+1}} P(x_t|x_{t-1}) \max_{x_{t-1}} P(x_{t-1}, e_{1:t-1})$$

$$= P(e_{t+1}|x_t) \max_{x_{t+1}} P(x_t|x_{t-1}) \cdot m_{t-1}(x_{t-1})$$

Result: Most likely sequence of hidden states $x_{1:N}^*$

/* Forward pass

for $t = 1$ to N do

 for $x_t \in \mathcal{X}$ do

 if $t = 1$ then

$$m_t[x_t] = P(x_t)P(e_0|x_t)$$

 else

$$a_t[x_t] = \arg \max_{x_{t-1}} P(x_t|x_{t-1})m_{t-1}[x_{t-1}]; \rightarrow \text{best } x_{t-1},$$

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

 end

end

/* Find the most likely path's ending point

$$x_N^* = \arg \max_{x_N} m_N[x_N]$$

/* Work backwards through our most likely path and find the hidden states

for $t = N$ to 2 do

$$x_{t-1}^* = a_t[x_t^*]$$

end

Particle Filtering?

• Time Elapse Update - Update the value of each particle according to the transition model. For a particle in state t_i , sample the updated value from the probability distribution given by $P(T_{i+1}|t_i)$.

• Observation Update - During the observation update for particle filtering, we use the sensor model $P(F_i | T_i)$ to weight each particle according to the probability dictated by the observed evidence and the particle's state. Specifically, for a particle in state t_i with sensor reading f_i , assign a weight of $P(f_i | t_i)$. The algorithm for the observation update is as follows:

1. Calculate the weights of all particles as described above.
2. Calculate the total weight for each state.
3. If the sum of all weights across all states is 0, reinitialize all particles.
4. Else, normalize the distribution of total weights over states and resample your list of particles from this distribution.

Utility & Rationality, decision.

Axioms of Rationality:

• Orderability: $(A \succ B) \vee (B \succ A) \vee (A \sim B)$

A rational agent must either prefer one of A or B , or be indifferent between the two.

• Transitivity: $(A \succ B) \wedge (B \succ C) \Rightarrow (A \succ C)$

If a rational agent prefers A to B and B to C , then it prefers A to C .

• Continuity: $A \succ B \succ C \Rightarrow \exists p [p, A; (1-p), C] \sim B$

If a rational agent prefers A to B but B to C , then it's possible to construct a lottery L between A and C such that the agent is indifferent between L and B with an appropriate selection of p .

• Substitutability: $A \sim B \Rightarrow [p, A; (1-p), C] \sim [p, B; (1-p), C]$

A rational agent indifferent between two prizes A and B is also indifferent between any two lotteries which only differ in substitutions of A for B or B for A .

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

If a rational agent prefers A over B , then given a choice between lotteries involving only A and B the agent prefers the lottery assigning the highest probability to A .

Value of perfect information:

$$VPI(E'|e) = MEU(e, E') - MEU(e)$$

$$MEU(e) = \max_a \sum_s P(s|e)U(s, a)$$

$$MEU(e, e') = \max_a \sum_s P(s|e, e')U(s, a)$$

Property:

① Nonnegativity. $\forall E', e \ VPI(E'|e) \geq 0$

② Nonadditivity. $VPI(E_j, E_k|e) \neq VPI(E_j|e) + VPI(E_k|e)$ in general.

③ Order-independence. $VPI(E_j, E_k|e) = VPI(E_j|e) + VPI(E_k|e, E_j) = VPI(E_k|e) + VPI(E_j|e, E_k)$

Markov Decision Processes.

Bellman equation: $\tau \rightarrow$ focus on long term,

$$U^*(s) = \max_a \sum_{s'} T(s, a, s')[R(s, a, s') + \gamma U^*(s')] \quad U^*(s) = \max_a Q^*(s, a)$$

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

① Value Iteration

1. $\forall s \in S$, initialize $U_0(s) = 0$. This should be intuitive, since setting a time limit of 0 timesteps means no actions can be taken before termination, and so no rewards can be acquired.

2. Repeat the following update rule until convergence: $\rightarrow U_{k+1}(s) = U_k(s)$

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

② Policy Extraction

$$\forall s \in S, \pi(s) = \operatorname{argmax}_a Q^*(s, a) = \operatorname{argmax}_a \sum_{s'} T(s, a, s')[R(s, a, s') + \gamma U^*(s')]$$

③ Q-Value Iteration:

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

④ Policy Iteration:

1. Define an *initial policy*. This can be arbitrary, but policy iteration will converge faster the closer the initial policy is to the eventual optimal policy.

2. Repeat the following until convergence:

- Evaluate the current policy with **policy evaluation**. For a policy π , policy evaluation means computing $U^\pi(s)$ for all states s , where $U^\pi(s)$ is expected utility of starting in state s when following π :

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

Define the policy at iteration i of policy iteration as π_i . Since we are fixing a single action for each state, we no longer need the max operator which effectively leaves us with a system of $|S|$ equations generated by the above rule. Each $U^{\pi_i}(s)$ can then be computed by simply solving this system. Alternatively, we can also compute $U^{\pi_i}(s)$ by using the following update rule until convergence, just like in value iteration:

$$U_{i+1}^{\pi_i}(s) \leftarrow \sum_{s'} T(s, \pi_i(s), s')[R(s, \pi_i(s), s') + \gamma U_i^{\pi_i}(s')]$$

However, this second method is typically slower in practice.

- Once we've evaluated the current policy, use **policy improvement** to generate a better policy. Policy improvement uses policy extraction on the values of states generated by policy evaluation to generate this new and improved policy:

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

If $\pi_{i+1} = \pi_i$, the algorithm has converged, and we can conclude that $\pi_{i+1} = \pi_i = \pi^*$.

Reinforcement Learning: - Active RL \rightarrow policy evaluation \wedge Direct Evaluation
- Passive RL \rightarrow learn policy \wedge model-based RL
- model-free RL \rightarrow Q-Learning

① Direct Learning Reward(A) = $\frac{1}{\text{# times A appears}} \sum R(A, a, s)$

② Temporal Difference Learning

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

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

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

③ Q-Learning:

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

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

Approximating Q-Learning \rightarrow feature-based representation.

$$V(s) = w_1 \cdot f_1(s) + w_2 \cdot f_2(s) + \dots + w_n \cdot f_n(s) = \vec{w} \cdot \vec{f}(s)$$

$$Q(s, a) = w_1 \cdot f_1(s, a) + w_2 \cdot f_2(s, a) + \dots + w_n \cdot f_n(s, a) = \vec{w} \cdot \vec{f}(s, a)$$

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

$$w_i \leftarrow w_i + \alpha \cdot \text{difference} \cdot f_i(s, a)$$

$$Q(s, a) \leftarrow Q(s, a) + \alpha \cdot \text{difference}$$

④ Model-Based Learning

$$f(s, a, s') = \frac{R(s, a, s')}{\# C(s, a)}$$

Generate T-examples by value / policy iteration

ϵ -Greedy Policy \wedge explore with ϵ

ϵ too large \rightarrow behavior randomly \rightarrow init ϵ large, then tune ϵ

ϵ too small \rightarrow less exploration \rightarrow slow to reach optimal policy

To avoid manual tuning? exploration function: deterministic way to explore.

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

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

model-based RL: estimate the T & R function, with samples obtained during exploration. Then use these estimates to solve MDP with value/policy iteration.

model-free RL: estimate V or Q -value of state directly.

not estimate T / R .

③ can be offline / online RL.

Both of them ① use exploration, to learn about environment (R & T).

② require online interaction with environment.

TD learning \rightarrow offline.

Q learning \rightarrow online.

Machine Learning \wedge supervised learning

procedure \wedge train validation. \rightarrow can modify hyperparameters

\ test

Naive Bayes \rightarrow classification problem

$$\begin{aligned} \text{prediction } (f_1, \dots, f_n) &= \operatorname{argmax}_{y \in \{1, \dots, k\}} P(y = y | f_1, \dots, f_n) \\ &= \operatorname{argmax}_{y \in \{1, \dots, k\}} P(Y = y, F_1 = f_1, \dots, F_n = f_n) \\ &= \operatorname{argmax}_{y \in \{1, \dots, k\}} \prod_{i=1}^n P(F_i = f_i | Y = y) \end{aligned}$$

Parameter Estimation \rightarrow MLE \rightarrow special case of MAP (max a priori)

assumptions: ① X_i is identically distributed.

iid \rightarrow ② Each sample x_i is conditionally independent of the others

③ All possible θ are equally likely before we seen any data (uniform prior)

$$\frac{\partial}{\partial \theta} \log L(\theta) = 0$$

Maximum Likelihood for Naive Bayes (Bernoulli dists)

$$L(\theta) = \prod_{i=1}^n P(F_i = f_i | Y = y) = \prod_{i=1}^n \theta^{f_i} (1 - \theta)^{1 - f_i}$$

$$\log L(\theta) = \log \theta \sum_{i=1}^n f_i + \log (1 - \theta) \sum_{i=1}^n (1 - f_i)$$

$$\frac{\partial}{\partial \theta} \log L(\theta) = 0 \Rightarrow \theta = \frac{1}{n} \sum_{i=1}^n f_i$$

Smoothing: $\text{PALE}(A) = \frac{\text{Count}(A) + k}{N + k|X|}$

$$P_{\text{APR}}(x) = \frac{\text{Count}(x) + k}{N + k|X|}, P_{\text{SAPR}}(x) = \frac{\text{Count}(x, y) + k}{N + k|X|}, P_{\text{SAPR}, \theta} = \frac{1}{|X|} \quad \theta \text{ can reduce overfitting}$$

Linear Regression

$$\text{Loss}(h_w) = \frac{1}{2} \sum_{i=1}^n \|x_i(t_i) \cdot h_w(x_i)\|^2 = \frac{1}{2} \|y - Xw\|^2$$

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

Perceptron:

$$\text{Linear Classifier: classification } (x) = h_w(x) = w^T f(x) \quad \begin{cases} + & \text{if } h_w(x) > 0 \\ - & \text{if } h_w(x) \leq 0. \end{cases}$$

Binary Perceptron:

1. Initialize all weights to 0: $w = 0$

2. For each training sample, with features $f(x)$ and true class label $y^* \in \{-1, +1\}$, do:

(a) Classify the sample using the current weights, let y be the class predicted by your current w :

$$y = \text{classify}(x) = \begin{cases} +1 & \text{if } h_w(x) = w^T f(x) > 0 \\ -1 & \text{if } h_w(x) = w^T f(x) \leq 0 \end{cases}$$

(b) Compare the predicted label y to the true label y^* :

- If $y = y^*$, do nothing
- Otherwise, if $y \neq y^*$, then update your weights: $w \leftarrow w + y^* f(x)$

3. If you went through every training sample without having to update your weights (all samples predicted correctly), then terminate. Else, repeat step 2

bias: $w^T f(x) + b = 0$

Optimization:

① gradient based method

Gradient ascent: $\text{obj} = \text{loss}(w) = \frac{1}{2} \|y - Xw\|^2$

Gradient descent: $\text{minimize } w \leftarrow w - \eta \nabla_w \text{loss}(w)$

Stochastic gradient descent: each iteration, use one sample

Mini-batch gradient descent: use size m samples

Least squares: gradient descent: $\text{loss}(h_w) = \frac{1}{2} \|y - Xw\|^2$

$w \leftarrow w - \eta (X^T y - X^T X w)$

$\sum_{i=1}^n t_{i,k} = 1$

$$\begin{aligned} \ell(w_1, \dots, w_K) &= \prod_{i=1}^n P(y_i | f(x_i); w) = \prod_{i=1}^n \prod_{k=1}^K \left(\frac{e^{w_k^T f(x_i)}}{\sum_{l=1}^K e^{w_l^T f(x_i)}} \right)^{t_{i,k}} \\ \log \ell(w_1, \dots, w_K) &= \sum_{i=1}^n \sum_{k=1}^K t_{i,k} \log \left(\frac{e^{w_k^T f(x_i)}}{\sum_{l=1}^K e^{w_l^T f(x_i)}} \right) \nabla_{w_j} \log \ell(w) = \sum_{i=1}^n \nabla_{w_j} \sum_{k=1}^K t_{i,k} \log \left(\frac{e^{w_k^T f(x_i)}}{\sum_{l=1}^K e^{w_l^T f(x_i)}} \right) = \sum_{i=1}^n \left(t_{i,j} - \frac{e^{w_j^T f(x_i)}}{\sum_{l=1}^K e^{w_l^T f(x_i)}} \right) f(x_i). \end{aligned}$$

Neural Network:

Theorem: (Universal Function Approximators) A two-layer neural network with a sufficient number of neurons can approximate any continuous function to any desired accuracy.

Measuring Accuracy:

$$\text{binary: } f^{(0)}(w) = \frac{1}{n} \sum_{i=1}^n \operatorname{sgn}(w^T f(x_i)) = y_i$$

$$\text{P of class } i: \sigma(x)_i = \frac{e^{w_i^T f(x)}}{\sum_{j=1}^J e^{w_j^T f(x)}} = P(y_i = j | f(x_i); w) \quad \ell(w) = \prod_{i=1}^n P(y_i | f(x_i); w)$$

$$\text{Sigmoid: } \delta(x) = \frac{1}{1 + e^{-x}}$$

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

How to choose activation func?
if we want to fit a non-linear func.
use non-linear activation func.

1 Dynamic Bayes Nets

Idea: Repeat fixed Bayes net structure at each time

Particle Filtering: simulating motion of set of particles through state graph to approx prob dist

1) Prediction Step: Sample new state from transition model

$$\text{particle } i \text{ state } x_i^{(t)} \quad x_{i+1}^{(t)} = P(x_{i+1} | x_i^{(t)})$$

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

$$w_i^{(t+1)} = P(e_{t+1} | x_i^{(t)})$$

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:

$A > B$: 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) Orderliness: 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 & 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 $A > B \Rightarrow (p \geq q \Leftrightarrow [p, A; (1-p), B] \geq [q, A; (1-q), B])$

3 Decision Networks

used to model state, actions, utilities of subsequent states

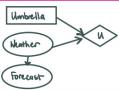
Chance Nodes: outcome has associated probability



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

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

4) Markov Decision Process (MDP)

Solving nondeterministic search problems

Markov Decision Process properties:

- State Space S - terminal state s_t^*
- actions A - start state s_0
- Transition Function $T(s, a, s')$ probability taking action a at state s results in s'
- Reward Function $R(s, a, s')$ positive reward for each step

Policy π : $S \rightarrow A$, 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 $U_0(s) = 0$

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

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

$$U^*(s) = \sum_a T(s, \pi(s), s') [R(s, \pi(s), s') + \gamma U^*(s')]$$

3) Use policy improvement to find better policy

$$\pi_{k+1}(s) = \arg \max_a \sum_{s'} T(s, a, s') [R(s, a, s') + \gamma U^*(s')]$$

5 Machine Learning

learning parameters of specified model given data

Decision Trees

simple nonlinear classifier but susceptible to overfitting

Entropy: uncertainty about variable \uparrow entropy \uparrow uncertainty

$$H(V) = -\sum_n P(V_n) \log_2 P(V_n)$$

entropy of binary variable

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

Information Gain: info gained by splitting on that feature

$$\text{Gain}(\text{Type}) = B\left(\frac{p}{p+n}\right) - \sum_{i=1}^n \frac{p_i+n_i}{p+n} B\left(\frac{p_i}{p_i+n_i}\right)$$

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:

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



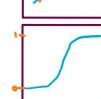
Logistic Regression

linear classifier for categorical variables

Logistic function:

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

Find optimal using gradients



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

Perceptron

Linear classifier: classification using linear combination of features

Activation: $activation(x) = h_w(x) = w^T f(x)$

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

Perceptron Algorithm:

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

$$w = w + y^* f(x)$$

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

Prediction for class label becomes $\text{prediction}(f) = \arg \max_p P(Y=f) \prod_i P(F_i=f|Y)$, normalize

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

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

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

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

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 sign threshold, softmax func to classify

Log Likelihood:

$$\log L(w) = \prod_i \log P(y_i | f(x_i), w)$$

Activation functions

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

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



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, 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 gives policy to follow and learns values of state exploring

1) Direct Evaluation: fix policy π and have agent experience, can compute by averages

2) Temporal Difference (TD) learning: learning from every experience

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

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

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

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

$$Q\text{-Value Iteration: } Q_{k+1}(s, a) \leftarrow \sum_s T(s, a, s') [R(s, a, s') + \gamma \max_a Q_k(s', a)]$$

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

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

Exploration functions: modified update $Q(s, a) = (1-\alpha) Q(s, a) + \alpha [R(s, a, s') + \gamma \max_a Q(s', a)]$

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

times a state visited

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