

CS 540

1.1 2.1 2.2 2.3

deterministic

stochastic

episodic

sequential

static

semi-dynamic

discrete

multi-agent

why AI is hard:-

- Large data set (high res. images)
- environment
- many solns
- Inaccurate as a data ("noise")
- Domain ~~is~~ knowledge

Uninformed SEARCH 3.1-3.4

HW1

* 5 components

- ① States → initial states, Goal states
- ② Action set
- ③ Transition model
- ④ goal test
- ⑤ PATH COST

$F_c \rightarrow$

$\leftarrow F$

$F_c \rightarrow$

$\leftarrow F_s$

$F_D \rightarrow$

$\leftarrow F \quad F_s \rightarrow$

Basic Search Tech

What knowledge does the Agent need?

WATER JUG PROBLEM

(0,0) // START STATE

(4,0) (0,3)

(3,0)

(3,3)

(4,2)

(0,2)

(2,0)

node

(0,0)

edge (0,3)

Expanding

action ((0,0))

The generated, but not yet expanded, states define the Frontier (aka Open or Fringe) list.

Directed graph to formalize the search

$$G = (V, E)$$

Frontier

TREE SEARCH ALG:

Frontier = {s}, where s is the start state.

Uninformed Search on Trees

BFS - Queue used for the Frontier list

DFS - stack

O-ABC

DFS

VCS

Priority Queue

Evaluating Search Strategies

Completeness

Optimality / Admissibility

Time Complexity

SPACE Complexity

Breadth is finite

Informed SEARCH

domain knowledge to guide selection of the best path

All domain knowledge used in the search is encoded in the heuristic function, h .

weak method.

Best First Search.

Greedy Best First Search

Beam Search

★
$$h^*(n) \leq h(n) \leq k \cdot h^*(n)$$

$$h^*(n)$$

local search & informed search

solution.

we don't care about path.

① Random start

② Stochastic HC

→ Simulated annealing

$$P = \frac{e^{(f(\text{newnode}) - f(\text{currentnode})) / T}}{f(s) > f(t)}$$

Boltzman's equation

Genetic Algorithms

1011

0111

2 operations { crossover
mutate

Selection: Finding the Fittest

Game PLAYING

Types of Games

~~game~~ Utility Function

Minimax Principle

Assume both players

The computer assumes after it moves

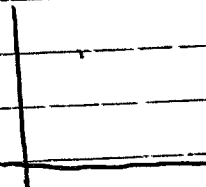
MINIMAX ALGORITHM

STATIC Board Evaluator

① KNN

Store all training example in memory
- Take in test example

② SVM



Decision TREES

OFFICE HOUR: T/R 2-300 pm

Exam

Multiple choice matching 40

T/F

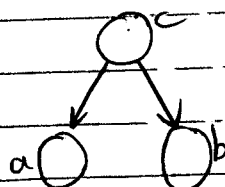
- Problem solving : 60
Open ended answer

X

Neural Network

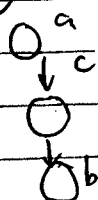
				w_1	w_2	w_0
1	1	0	< 0			
1	0	1	> 0	1	0	
0	1	1	> 0	0	1	
0	0	0	< 0			

①

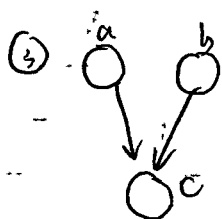


a & b are cond. indep. given c

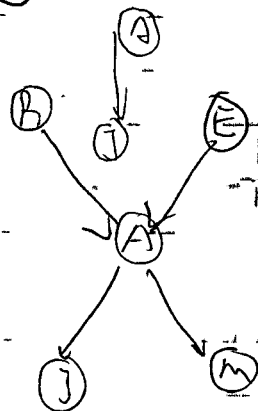
②



a & b are cond. indep. given c



$a \nVdash b$ and dependent given c



a. state space actions transition model goal test path cost determine if it is a goal node. cost function

b. BFS Queue
DFS Stack
UCS Priority Queue
IDS ~~map~~ DFS to depth. } has advantage of BFS: complete/optimal
DFS: limited space

	BFS	DFS	UCS	IDS
TIME complexity	$O(b^d)$	$O(b^d)$	$O(b^{(b^d)})$	$O(b^d)$
Space complexity	$O(b^d)$	$O(b^d)$	$O(b^{(b^d)})$	$O(b^d)$
Complete	✓	✓	✓	✓
Optimal	✓	✗	✓	✓

c. Informed Search. uses domain-specific information in some way
a. heuristic: is computable from the current state description
it estimates the "goodness" of node n
b. Greedy Best-First search
not optimal search. $f(n) = g(n) + h(n)$ → the estimate cost from node n to goal
cost from start to current node n cost path from node n to goal
never over-estimated $\underline{h(n) \leq h^*(n)} \Rightarrow$ h is an admissible heuristic function
actual

c. TIME complexity
space.
complete / optimal / Admissible

d. admissible / consistent \Rightarrow A heuristic h is consistent.
e. $h(n) \leq c(n, n') + h(n')$ consistency is a stronger condition than admissibility
f. local searching:

for every node n ,
the estimated cost of reaching the goal from n
step cost of getting to n'
estimated cost of reaching the goal from n'

$h(n) = h^*(n) \Rightarrow$ no necessary work is performed
 $h(n) = 0$ same as UCS

h is close to $h^* \Rightarrow$ the fewer extra nodes that will be expanded
 $h_1(n) \leq h_2(n) \leq h^*(n)$ A_2^* is better informed

f. Local searching: every node is a solution
action go from one solution to another
can stop at any time and have a valid solution
goal: to find a better/best solution
evaluation function

hill-climbing: Pick initials
Pick t in neigh (s) with the largest $f(t)$
 $f(t) \leq f(s)$ then stop and return $s = t$
 \Rightarrow stop at a local maximum
stochastic restart

- Pick initial state, s
 Randomly pick state t from neighbors of s
 if $f(t)$ better than $f(s)$
 then $s = t$ else with small probability ϵ
 $p = e^{\Delta E / T}$ can performed multiple backward steps in a row to escape a local max different from local search.

h. genetic algorithms
 proportional selection

fitness = $\frac{\text{Fitness of individual}}{\text{Sum Fitness for all individuals}}$

zero-sum

Discrete Finite:
 Stochastic.

Backgammon
 Monopoly

a. Deterministic
 Fully Observable (perfect info)
 Checkers
 Chess, Go
 Othello

Partially observed (imperfect info) strategy
 Battleship

Bridge, Poker
 Scrabble

b. utility function, to map each terminal state of the board to a score indicating the value of that outcome to the computer.

Greedy search: Expand the search tree to the terminal states on each branch.
 Evaluate the utility of each terminal board configuration. make the initial move that results in the board config with the maximum value.

c. Minimax Complexity Space: PFs $O(b^d)$ time BFS $O(b^d)$

e: because utility function is defined only for terminal states
 static Evaluation function use heuristics to estimate the value of non-terminal states
 is used to estimate how good the current board configuration is for the computer
 should agree with utility function when calculated at terminal node.

f.g. Pruning can be used to ignore some branches.

Max value $\alpha \rightarrow$ come from children

$\beta \rightarrow$ come from parents

min value $\alpha \rightarrow$ parents
 $\beta \rightarrow$ children

h. how to deal with limited time
 IDS iterative deepening search

$O(b \cdot d)$

Horizon Effect: the computer has a limited horizon, it cannot see that this significant event could happen

sol = quiescence search
 secondary search

4. Unsupervised machine learning

a. solve classification problems / learn models of data (data fitting) / understand that are unknown to humans ("data mining")
 and improve efficiency of human learning / discover new things or structures

b. Inductive learning: generalize from a given set of (training) examples so that accurate predictions can be made about future examples. \odot learn unknown function $f(x) = y$

c. $x \rightarrow$ represent a specific object. $y \rightarrow$ label

h (hypothesis) function is learned that approximates f

difference between unsupervised/supervised: gives a set of (x, y) pairs and only the x 's are given

f. clustering: group training sample into clusters such that examples in the same cluster are similar, and examples in different clustering are different.

g h i. HAC k-means

mean shift clustering

j. distortion: determine how good our clusters are.

5) Supervised Machine Learning

b) Concept learning: Determine from a given set of examples if a given example is or is not an instance of the concept/class/category. If it is \rightarrow positive / negative

c. xmn d. Inductive bias. \Rightarrow learning can be viewed as searching the Hypothesis Space H of possible functions. Inductive Bias is used when one h is chosen over another / - it needed to generalize beyond the specific training examples beyond the specific training examples

e. SVM

slack variables: $\min w, b, \xi \text{ s.t. } w^T x + b \leq \xi$

trade off parameter \rightarrow select variables.

map data to high dimensional space
 kernels.

1) Decision trees a. measure response usually discretized but can be continuous classification
 In supervised learning { Training set: n pairs of example, label
 A predictor $f: X \rightarrow Y$
 Hypothesis space: space of predictors
 Find the "best" function in the hypothesis space that general well
 performance measure: MSE for regression

c. Accuracy $acc = \frac{1}{n} \sum_{i=1}^n 1_{y_i = \hat{y}_i}$
 ID3 build tree (examples, questions, defaults)
 IF empty (examples) then return the attribute of majority class
 IF (examples have same label y) then return y
 IF empty (questions) then return majority vote in examples
 $q = \text{best-question}(\text{examples}, \text{questions})$ create and return an internal node with n children
 let there be n answers to q
 The i-th child is built by calling

i. info gain $I(Y;X) = H(Y) - H(Y|X)$
 $H(Y) = -\sum_{i=1}^n P(Y=y_i) \log_2 P(Y=y_i)$
 $H(Y|X) = -\sum_{i=1}^n P(Y=y_i|X=v) \log_2 P(Y=y_i|X=v)$
 ii. see purity node iii. root node is also internal node

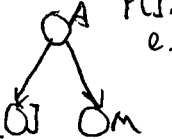
f. overfitting: why too many features
 how to overcome: 1. Randomly split data into train and tune 2. build a full tree using only TRAIN
 3. Prune the tree based on the TUNE set
 Algo: Prune (tree, T, TUNE set)
 1. Compute T's accuracy on TUNE, call it ACT
 2. For every internal node N
 a) New tree $T_N = \text{copy of } T$, but prune (delete) the subtree under N
 b) N becomes a leaf node in T_N . The label is the majority vote of TRAIN-examples reaching N
 c) $ACT_N = T_N$'s accuracy on TUNE
 3. let T_X be the tree (among the T_N 's and T) with the largest ACT
 - ACT: set $T \leftarrow T_X$ / prune

j. ensemble: Aggregation of predictions of multiple classifiers with the goal of improving accuracy by reducing the variance of an estimated prediction function
 k. bagging: bootstrap sample, Construct a classifier \Rightarrow calculate error \Rightarrow repeat
 l. random forests: bagging, Randomized choose attributes
 m. boosting: sequential production of classifiers, each is dependent on the previous one, make examples misclassified by current classifier more important
 n. linear perceptron: unique global minimum $Wd \leftarrow Wd - \alpha \sum_{i=1}^n (a_i - y_i) x_i d$
 o. step function $g(h) = 0$ if $h < 0$; $g(h) = 1$ if $h \geq 0$
 p. sigmoid function $g(h) = 1/(1 + \exp(-h))$ $g'(h) = g(h)(1-g(h))$
 q. linear perceptron: $a = w_0 * x_0 + w_1 * x_1 + \dots + w_n * x_n$ b. biased term

r. step function is discontinuous, cannot use gradient descent
 s. sigmoid function $E(W) = \frac{1}{2} \sum_{i=1}^n (a_i - y_i)^2$ $\frac{\partial E}{\partial w_d} = \sum_{i=1}^n (a_i - y_i) a_i (1 - a_i) x_{id}$
 t. Even with a non-linear sigmoid function, the decision boundary a perceptron can be produced
 u. Calculate output h_1 out(h_1) = sigmoid(h_1)
 v. Calculate delta terms
 w. outer layers: $E_{total} = \frac{1}{2} \sum (out - y_i)^2$
 $\frac{\partial E_{total}}{\partial w} = (0 - target) * (1 - 0) * out(h_1) = W$
 x. hidden layers: $\frac{\partial E_{total}}{\partial w} = z_i * out(h_1) * (1 - out(h_1)) * W_5 * (0 - target) * 0 * (1 - 0) =$
 y. Deep learning: multiple layer NNs

3) Uncertainty
 a. Event space / Mutually exclusive events / Random variables
 b. Probability Distributions: $P(A)$ the set of values $\{P(a_1), P(a_2), \dots, P(a_n)\}$ sum=1
 c. The axioms of Probability: ① $0 \leq P(A) \leq 1$ ② $P(\text{true}) = 1$ $P(\text{false}) = 0$ ③ $P(A \cup B) = P(A) + P(B) - P(A \cap B)$
 d. joint probability: $P(A, B)$ full joint probability distribution table n random variables, each taking k values
 e. marginalization: marginal probabilities
 f. conditional: $0 \leq P(a|e) \leq 1$ $P(a|e) + P(a|\neg e) = 1$ $P(\neg a|e) = 1 - P(a|e)$
 g. Probability

Summary
 Conditional Probability: $P(A|B) = P(A, B) / P(B)$
 Product Rule: $P(A, B) = P(A|B)P(B)$
 Chain Rule: $P(A, B, C, D) = P(A|B, C, D)P(B|C, D)P(C|D)P(D)$
 Conditionalized version of chain Rule:
 $P(A, B|C) = P(A|B, C)P(B|C)$
 Bayes' Rule: $P(A|B) = P(B|A)P(A) / P(B)$
 Conditionalized version of Bayes' Rule:
 $P(A|B, C) = \frac{P(B|A, C)P(A|C)}{P(B|C)}$
 Addition/conditioning rule: $P(A) = P(A, B) + P(A, \neg B)$
 $P(A) = P(A|B)P(B) + P(A|\neg B)P(\neg B)$

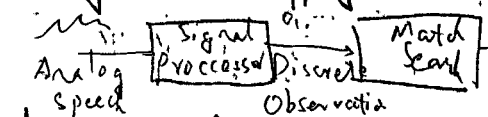
independence: $P(A, B) = P(A) \times P(B)$ $P(A|B) = P(A)$
 conditional independence:
 $P(A|B, C) = P(A|C)$ $P(B|A, C) = P(B|C)$ $P(A, B|C) = P(A|C)P(B|C)$
 4. Bayesian Networks table $k^n \rightarrow$ combination
 CPT: Conditional Probability Table
 $P(I, M|A) = P(I|A)P(M|A)$
 e. d-separation:


A and B are independent given nothing else, but are dependent given C

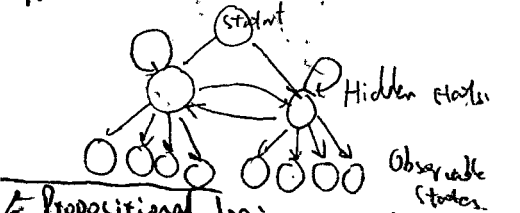
d. Bayes Net must be acyclic $O(Nk^m)$!
 $P(x_1, \dots, x_N) = \prod_i P(x_i | \text{parents}(x_i))$ \rightarrow Inference by Enumeration
 g. variable elimination - bookkeeping algo to be more efficiency
 h. Parameter learning - problem: unseen event, smoothing
 1. add-one smoothing: $P(A|B, E) = \frac{1}{\#(A) + 1} \frac{\#(A) + 1}{\#(A) + 1 + \#(\neg A) + 1}$
 2. add-delta smoothing: $P = \frac{(\text{delta} + \text{true})}{(\text{delta} + \text{true} + \text{smooth})}$

$P(d) = \sum \sum \sum \sum P(v, s, t, l, a, b, x, d)$
 $P(d) = \sum_v P(v) \sum_s P(s) \sum_t P(t|s) \sum_l P(l|s) P(a|t, l) P(x|a, l, d)$
 $P(t) = \sum_l P(t|l)$

j. Naive Bayes Net. Maximum Likelihood Estimate: $V = \text{argmax}_v P(Y=v) \prod_{i=1}^n P(X_i = x_i | Y=v)$ classifier.
 Assumption: all evidence variables are conditionally independent of each other given the class variable.
 Input (evidence): C, Z, H $P(I|C, Z, H) = P(I, C, Z, H) / P(C, Z, H)$ $P(I|C, Z, H) = P(I, C, Z, H) / P(C, Z, H)$
 Output (query): I $P(I, C, Z, H) = P(I)P(C|I)P(Z|I)P(H|I)$ $P(\neg I, C, Z, H) = P(\neg I)P(C|\neg I)P(Z|\neg I)P(H|\neg I)$

b. Speech Recognition. a. mapping an acoustic signal into a string of words
 Task:  a. Model: $P(\text{words} | \text{signal}) = \frac{P(\text{signal} | \text{words}) P(\text{words})}{P(\text{signal})}$
 Language: $P(\text{words})$

c. first-order markov Assumption: $P(w_i | w_1, \dots, w_{i-1}) \approx P(w_i | w_{i-1})$
 e. LM: $P(w_1, w_2, \dots, w_n) = P(w_1)P(w_2|w_1) \dots P(w_n|w_{n-1})$ bigram model
 Acoustic Model: $P(\text{signal} | \text{words}) \Rightarrow P(\text{phones} | \text{word})$ (mm)
 HMM model: State q_{t+1} is conditionally independent of $\{q_{t+2}, q_{t+3}, \dots, q_i\}$ given q_t
 trigram model: $P(w_i | w_1, \dots, w_{i-1}) \approx P(w_i | w_{i-1}, w_{i-2})$
 HMM: $P(\text{signal} | \text{phones})$



$P(q_1, q_2) = P(q_2 | q_1)P(q_1)$
 $P(q_2) = P(q_2 = H | q_1 = H)P(q_1 = H) + P(q_2 = H | q_1 = S)P(q_1 = S)$
 $P(q_3 = H) = P(q_3 = H | q_1 = H, q_2 = H)P(q_1 = H, q_2 = H) + P(q_3 = H | q_1 = H, q_2 = S)P(q_1 = H, q_2 = S) + \dots$

6. Propositional logic
 Syntax Precedence: highest to lowest $\neg, \wedge, \vee, \Rightarrow, \Leftrightarrow$
 The semantics (meaning) of a sentence is the set of interpretations in which the sentence evaluates to T
 the semantics of the sentence $P \vee Q$ is the set of 6 interpretations:
 $P=T, Q=T, R=T$ or F out
 $P=T, Q=F, R=T$ or F out
 $P=F, Q=T, R=T$ or F

P	Q	$\neg P$	$P \wedge Q$	$P \vee Q$	$P \Rightarrow Q$	$P \Leftrightarrow Q$
F	F	T	F	F	T	T
F	T	T	F	T	T	F
T	F	F	F	F	F	F
T	T	F	T	T	T	T

a. Interpretation: is a complete T/F assignment to all proposition symbols
 $P, Q, R \Rightarrow 8$ interpretation
 b. knowledge base (KB): a set of sentences.
 A model of a KB is an interpretation in which all sentences in KB are true
 d. Entailment: is the relation of a sentence β logically following from other sentences α , $\alpha \models \beta$ iff α is true, β is true
 if $\alpha \models \beta$ is valid / $\alpha \wedge \neg \beta$ is not satisfiable.
 Soundness & Completeness.

Soundness: any wff that follows deductively from a set of axioms, KB, is valid (i.e., true in all models).
 Completeness: can all valid sentences (i.e., true in all models of KB) can be proved from KB and hence are theorems.
 f. Inference by Enumeration: complete. slow: takes exponential time.
 g. Resolution Sound
 Complete