

CRYPTArithmetic Problem

↳ TYPE of constraint satisfaction Problem

↳ Constraints

- ↳ digits can be assigned a word / alphabet ($0-a$)^{range}
- ↳ no 2 letters have same value
- ↳ sum of digits must be shown in problem
- ↳ only 1 carry forward

↳ Start from left most digit

↳ If it is carry

↳ Assign it 1

$\frac{1}{10}$ always 1

Q)

$$\begin{array}{r} T_2 \ 0_1 \\ + \ G \ 0_1 \\ \hline 0 \ U \ T_2 \end{array}$$

$2+G = 7 \text{ or } 10$
 as prob is carrying
 will be 10 or >
 $U=1 X$
 as 0:1
 can't be same

Letter	Digit
T	—
O	—
G	—
U	—

Q)

$$\begin{array}{r} C_5 \ C_4 \ C_3 \ C_2 \ C_1 \rightarrow \text{carry} \\ S_1 \ E_5 \ N_6 \ D_7 \\ + M_1 \ O_0 \ R_8 \ E_5 \\ \hline M_1 \ D_0 \ N_6 \ E_5 \ Y_2 \end{array}$$

L	D
S	9
E	5
N	6
D	7
M	1
O	0
R	8
Y	2

1. $S+M > 10$
 $S+1 > 10$
 \downarrow
 $q+1 = 0$

2. $E+O = N$
 $E+O+C_1 = N$
 $\downarrow \quad \downarrow$
 $5+0+1 = 6$

If $E+O=N$
 $E=N$
 which means
 there should be
 a carry
 $C_1=1$

3. $N+R=E$

$6+R=5$

\downarrow
 q
 8

$6+q=15$
 \downarrow
 carry
 $\{ \}$

$6+8+C_1=15$
 $6+8+1=15$

$S=q$

4. $D+5=Y+10$

$\downarrow \quad \downarrow \quad \downarrow$
 $4 \quad 5 \quad 6$
 $y+5=9 \quad E=5 \quad x \times$
 $x \quad x$
 $5=9$

$S=9$

$7+5=12$
 $y=2$

$$0) \begin{array}{r} \text{C}_5 \text{ C}_4 \text{ C}_3 \text{ C}_2 \text{ C}_1 \text{ C}_0 \\ \text{E}_8 \text{ A}_1 \text{ T}_9 \\ + \text{T}_9 \text{ H}_L \text{ A}_1 \text{ T}_9 \\ \hline \text{A}_1 \text{ P}_0 \text{ P}_0 \text{ L}_3 \text{ E}_8 \end{array}$$

$$\begin{aligned} 2. \quad T &\neq P & 8+H = 0+10 \\ T+C_3 &= C_4+P \\ 9+1 &= 10 \end{aligned}$$

$$\begin{array}{r} L \quad D \\ A \quad 1 \\ P \quad 0 \\ L \quad 3 \\ E \quad 8 \\ T \quad 9 \\ H \quad 2 \end{array}$$

$$0) \begin{array}{r} \text{C}_5 \text{ C}_4 \text{ C}_3 \text{ C}_2 \text{ C}_1 \\ \text{S}_1 \text{ O}_9 \text{ M}_3 \text{ E}_4 \\ + \text{T}_8 \text{ I}_5 \text{ M}_3 \text{ E}_4 \\ \hline \text{S}_1 \text{ P}_0 \text{ E}_4 \text{ N}_6 \text{ T}_8 \end{array}$$

$$\begin{array}{l} MM=N \\ EE=T \end{array} \rightarrow \text{EVEN} \quad \begin{array}{c} 0 \\ 2 \\ 4 \\ 6 \\ 8 \end{array} \quad \begin{array}{l} T=P \rightarrow E=4 \\ N=6 \rightarrow M=3 \end{array}$$

$$\begin{array}{l} 1+T+C_3=10 \\ \downarrow \\ 8 \end{array} \quad \begin{array}{l} 0+I=E+10 \\ =4+10 \end{array}$$

$$\begin{array}{l} 0+I=14 \\ \begin{array}{c} 1 \\ 3 \\ 7 \\ 9 \\ \times \end{array} \quad \begin{array}{c} 1 \\ 8 \\ 9 \\ 5 \\ \times \end{array} \end{array}$$

$$0) \begin{array}{r} \text{B}_1 \text{ A}_9 \text{ S}_8 \text{ E}_3 \\ + \text{B}_1 \text{ A}_6 \text{ L}_5 \text{ L}_5 \\ \hline \text{G}_1 \text{ A}_4 \text{ M}_9 \text{ E}_3 \text{ S}_8 \end{array}$$

$$\begin{aligned} E+L &= S \rightarrow C \\ E+L &= S+10 \rightarrow C \\ \hookrightarrow E &= S-L+10 \end{aligned}$$

$$\begin{aligned} 1. \quad B+B &= A+C \\ \begin{array}{c} 1 \\ 6 \\ 7 \\ 8 \\ 9 \\ \times \end{array} & \quad \begin{array}{c} 9 \\ 6 \\ 7 \\ 8 \\ 9 \\ \times \end{array} \quad \begin{array}{l} \text{Taking} \\ 9 \text{ } E(7,2) \quad B(6) \end{array} \\ 2. \quad S-E &= S \\ \begin{array}{c} 5 \\ 6 \\ 7 \\ 8 \\ 9 \\ \times \end{array} & \quad \begin{array}{c} 5 \\ 6 \\ 7 \\ 8 \\ 9 \\ \times \end{array} \quad \begin{array}{l} 6+6=R \\ C \quad A \end{array} \\ S-E &= S, E(8,3) \quad B(7) \\ 7+7 &= 14 \end{aligned}$$

$$\begin{aligned} S+L &= E \\ S-L+10 &= S+L \\ 2L &= 10 \\ L &= 5 \end{aligned}$$

Adversarial Search

↳ competition b/w 2 people

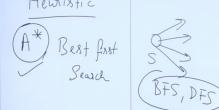
↳ 1 person's loss, is the other's gain

↳ both players have acts to complete into

↳ no chance → dice

e.g. TIC TAK TOE
CHECKERS
CHESS

Environments
↳ deterministic
↳ fully observable

Informed Search	Uninformed Search	Adversarial Search
Heuristic  A* Best first Search	 BFS, DFS	Game  Minimax Planning

Evaluation function

- Evaluation function or static evaluator is used to evaluate the "goodness" of a game position
 - Contract: heuristic search where evaluation function was a non-negative estimate of the cost from the start node to a goal and passing through the given node ↗
- Zero-sum assumption lets us use a single evaluation function to describe goodness of a board wrt both players
 - $f(n) \gg 0$: position n good for me and bad for you
 - $f(n) \ll 0$: position n bad for me and good for you
 - $f(n)$ near 0: position n is a neutral position
 - $f(n) = +\infty$: win for me
 - $f(n) = -\infty$: win for you

Evaluation function examples

- Example of an evaluation function for Tic-Tac-Toe
 $f(n) = [\# \text{ of 3-lengths open for me}] - [\# \text{ of 3-lengths open for you}]$ where a 3-length is a complete row, column, or diagonal
- Alan Turing's function for chess
 - $f(n) = w(n)/b(n)$ where $w(n) = \text{sum of the point value of white's pieces}$ and $b(n) = \text{sum of black's}$
- Most evaluation functions specified as a weighted sum of position features
 $f(n) = w_1 * \text{feat}_1(n) + w_2 * \text{feat}_2(n) + \dots + w_n * \text{feat}_n(n)$
- Example features for chess are piece count, piece placement, squares controlled, etc.
- Deep Blue had >8K features in its evaluation function

Introduction to Game Playing

↳ Minimax AI/ID

↳ Alpha beta ($\alpha-\beta$) Pruning

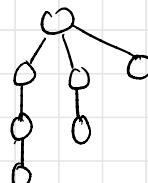
→ use Depth First Search

Max: My turn

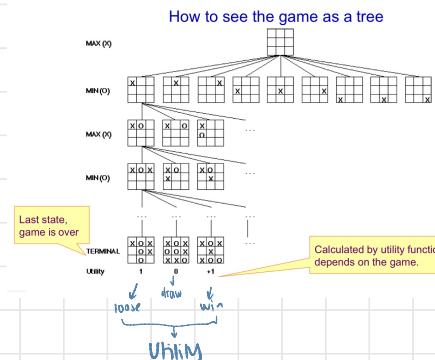
Min: Opponent's turn

Utility: Winner Prize

$b \rightarrow$ branching factor =
 $m \rightarrow$ max Path length =



Game Tree (2-player, deterministic, turns)



Minimax Algo

↳ Backtracking Algo

↳ Best move strategy used

↳ MAX will try to maximize its utility \rightarrow best move

↳ MIN will try to minimize its utility \rightarrow worst move

Time Complexity

$O(b^d)$

depth

branching factor

moves by MAX \rightarrow always at top

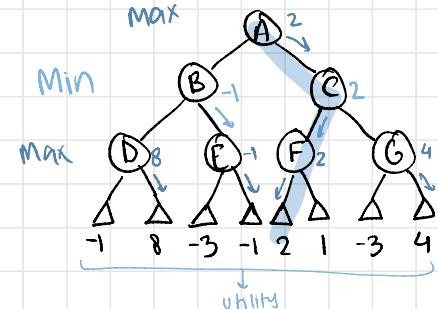
counter moves by MIN

Minimax Tree Evaluation

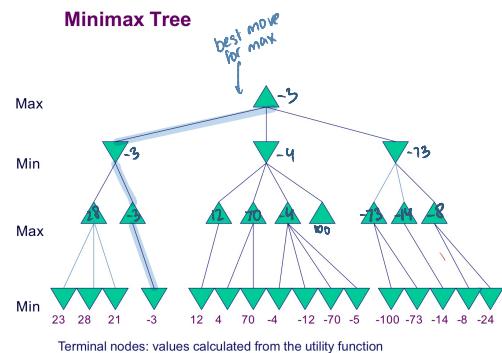
For the MAX player

1. Generate the game as deep as time permits
2. Apply the evaluation function to the leaf states
3. Back-up values
 - At MIN assign minimum payoff move
 - At MAX assign maximum payoff move
4. At root, MAX chooses the operator that led to the highest payoff

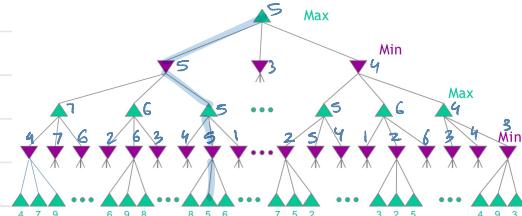
Game Tree



Minimax Tree



MiniMax Example-2



Properties

Complete

Yes, if tree is finite

Optimal

Yes, against an optimal opponent

Time complexity

b^m

Space complexity

b^m

For chess, $b \approx 35$, $m \approx 100$ for "reasonable" games

\rightarrow exact solution completely infeasible, since it's too big

Instead, we limit the depth based on various factors, including time available.

Part A) Under the Flag of FAST NUCES- IEEE an AI-Enabled Tech Exhibition 2022 project exhibition has been organized. This exhibition aims to provide students' work with broader visibility and make their work accessible and available for the interested people. With the policy of having 3 members in group and zero tolerance for plagiarism. If found it will lead to Disciplinary Action. So, you and two of your friends decided to participate (Consider yourself=X, Friend1=Y, Friend2=Z). X and Y sincerely did hard work, and both cooperated with each other to analyze and solve the occurring challenges. But Z always considered himself a mastermind and never participated in group activity. But suddenly Z got stuck and copy pasted the code from internet without bringing this thing in his friends' knowledge.

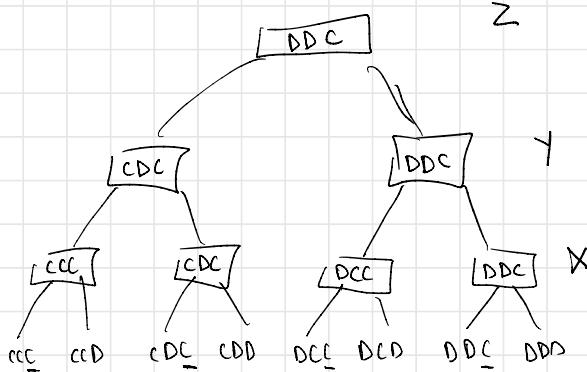
During evaluation it was found that a particular chunk of code has been copied and now three of them will be called for DC. Where either each participant will Confess or Deny. As X and Y always worked as a team so they decided to give priority to co-operate (i.e. Confess or Deny) opposite decision(X confess and Y Deny or vice versa). But always aiming for Z to Confess because he did wrong so only he will be liable for breaking rule. In case if in any circumstances both must take opposite decision still their aim would be same for Z to confess. Meanwhile Z realized that he did wrong and try to make sure his friends won't be suffering furthermore due to him.

So, using Min max draw a Game tree and Identify what decision should be taken by Z?

(Note: Tuple utilities to be assigned as X, Y, Z)

Part B) A game is played between max and min. Draw a tree considering MIN as a first player and branching factor is 2. Given terminal values below, show backed up values, best decision available at the root and branches that will get pruned (explain why they are getting pruned if any)?

5,6,4,10,7,8,4,5

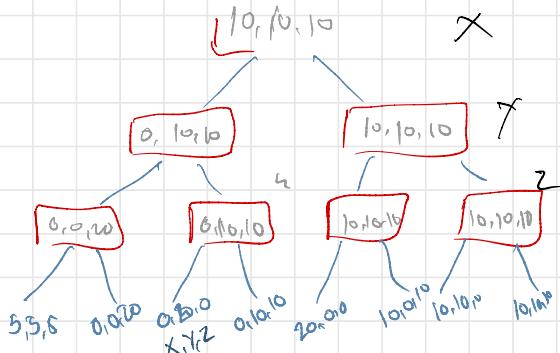
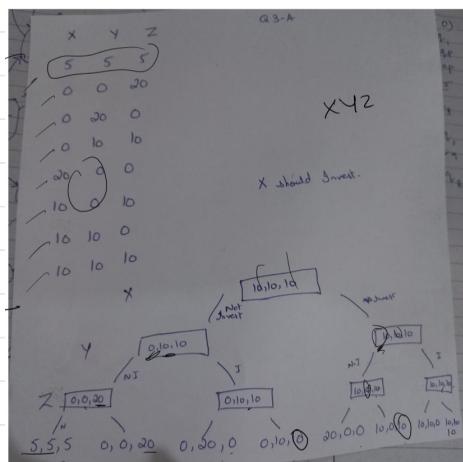


P1 - Z cooperates
P2 - Y, X oppose

A: Suppose that three companies, Company X, Company Y, and Company Z, are competing in the market for a new product. Each company can choose to invest in research and development (R&D) or not invest in R&D. The offer is:

- A. If all companies invest in R&D, they will all develop a high-quality product and receive a benefit of \$10 million each. ↗ ↗
- B. If two companies invest in R&D, they will all develop a high-quality product and receive a benefit of \$10 million each and the other companies receive nothing. ↗ ↗
- C. If only one company invests in R&D, that company will develop a superior product and receive a benefit of \$20 million while the other companies receive nothing. ↗ ↗
- D. If no company invests in R&D, they will all develop a mediocre product and receive a benefit of \$5 million each. ↗ ↗

Now choose the best strategy for Company X using Min-Max concept and provide reasoning for your choice. Define properly how you generated utility values based on which the result/ strategy has been figured.



X	Y	Z
10	10	10
10	10	0
0	10	10

ALPHA-BETA PRUNING (α - β)

cut off

improves performance
of MinMax algo

↳ cut off search by exploring less no. of nodes

↳ if curr max > successor min
then don't explore that min subtree

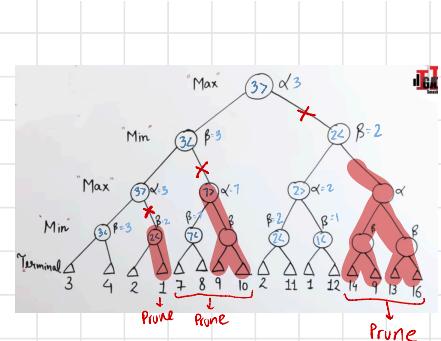
↳ saves time

↳ avoids unnecessary computation

$\alpha \rightarrow$ Max node

$\beta \rightarrow$ Min node

Pruning: ways of determining that certain branches will not be useful.

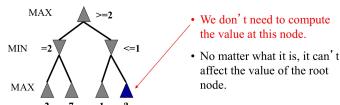


Time Complexity

↳ $O(b^{d/2}) \rightarrow$ best/avg case

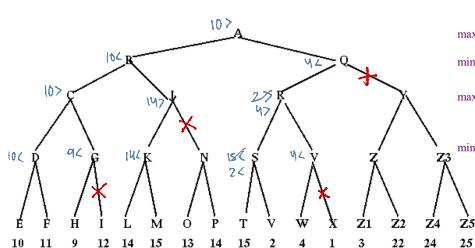
↳ $O(b^d) \rightarrow$ worst case

• Basic idea: "If you have an idea that is surely bad, don't take the time to see how truly awful it is." -- Pat Winston

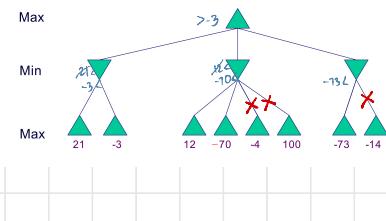


- We don't need to compute the value at this node.
- No matter what it is, it can't affect the value of the root node.

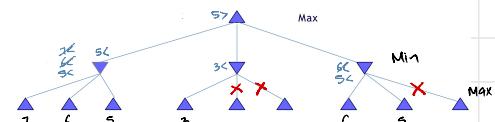
Exercise



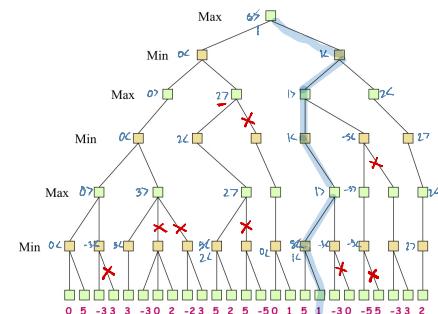
β Cut Example



Alpha-Beta Example 2



Alpha-Beta Tic-Tac-Toe Example 2

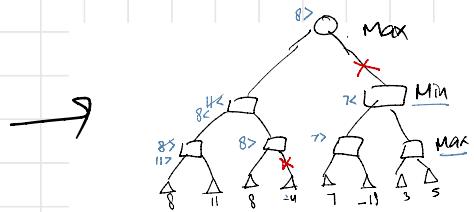


Question 1: Alpha-Beta pruning [10 points]

Draw a tree with 2 moves for Max and 1 move for Min and branching factor=2. Label the nodes alphabetically. That is, label the root A, the next B, and so on in a left to right breadth first fashion. Given terminal values below, show backed up values, best decision available at the root, branch(es) that will get pruned and explain why they will get pruned.

6, 11, 4, -4, 7, -18, 3, 5

Question 2: CD [10 points]

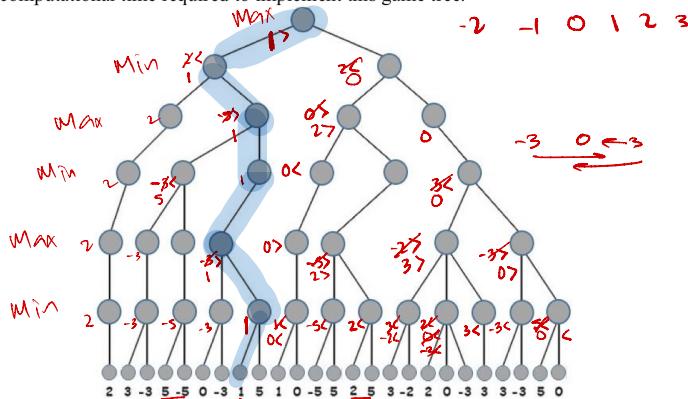


Yes, BFS can be used instead of minmax in a game play to find all the game states that can be obtained in upto ∞ moves from some starting game state using a queue -

- Using BFS, all possible sequence of N moves can be easily figured out from starting position and situation after each sequence of N moves.
- In a BFS all reachable N moves can be easily found. The game states are reachable in N moves and will be at Nth-level of game state tree rooted at 's(state)'.
- & BFS can be used but there are some drawbacks i.e:
 - we have to store all the nodes (level wise) and their children until you reach terminating search depth - consuming more memory as compared to minmax/DFS which require linear amount of memory.

A) Can Breadth First Search (BFS) be used instead of Min-Max Algorithm in game play? Justify your answer.

B) Use Min-Max algorithm on the following game tree to find the best possible path for Max player that ensures to maximize its score. Also, apply Alpha-Beta pruning to reduce the computational time required to implement this game tree.



-3 -2 0 1
→

α - β Pruning

Can store information along an entire *path*, not just at most recent levels!

Keep along the path:

- α : best MAX value found on this path
(initialize to most negative utility value)
- β : best MIN value found on this path
(initialize to most positive utility value)

Why is it called α - β ?

- α is the value of the best (i.e., highest-value) choice found so far at any choice point along the path for *max*
- If v is worse than α , *max* will avoid it
→ prune that branch
- Define β similarly for *min*



Pruning at MAX node

- α is possibly updated by the MAX of successors evaluated so far
- If the value that would be returned is ever $> \beta$, then stop work on this branch
- If all children are evaluated without pruning, return the MAX of their values

Pruning at MIN node

- β is possibly updated by the MIN of successors evaluated so far
- If the value that would be returned is ever $< \alpha$, then stop work on this branch
- If all children are evaluated without pruning, return the MIN of their values

Imperfect Decisions

Complete search is impractical for most games

Alternative: search the tree only to a certain depth

- Requires a cutoff-test to determine where to stop
- Replaces the terminal test
- The nodes at that level effectively become terminal leave nodes
- Uses a heuristics-based evaluation function to estimate the expected utility of the game from those leave nodes.

Utility Evaluation Function

Very game-specific

Take into account knowledge about game

"Stupid" utility

- 1 if player 1 wins
- -1 if player 0 wins
- 0 if tie (or unknown)
- Only works if we can evaluate complete tree
- But, should form a basis for other evaluations

Utility Evaluation

Need to assign a numerical value to the state

- Could assign a more complex utility value, but then the min/max determination becomes trickier.

Typically assign numerical values to lots of individual factors:

- $a = \# \text{player 1's pieces} - \# \text{player 2's pieces}$
- $b = 1$ if player 1 has queen and player 2 does not, -1 if the opposite, or 0 if the same
- $c = 2$ if player 1 has 2-rook advantage, 1 if a 1-rook advantage, etc.

Utility Evaluation

The individual factors are combined by some function

Usually a linear weighted combination is used:

- $u = \alpha a + \beta b + \gamma c$
- Different ways to combine are also possible

Notice: quality of utility function is based on:

- What features are evaluated
- How those features are scored
- How the scores are weighted/combined

Absolute utility value doesn't matter – relative value does.

Evaluation Functions

If you had a perfect utility evaluation function, what would it mean about the minimax tree?

You would never have to evaluate more than one level deep!

Typically, you can't create such perfect utility evaluations, though.

Evaluation Functions for Ordering

As mentioned earlier, order of branch evaluation can make a big difference in how well you can prune

A good evaluation function might help you order your available moves:

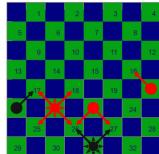
- Perform one move only
- Evaluate board at that level
- Recursively evaluate branches in order from best first move to worst first move (or vice-versa if at a MIN node)

Checkers Case Study

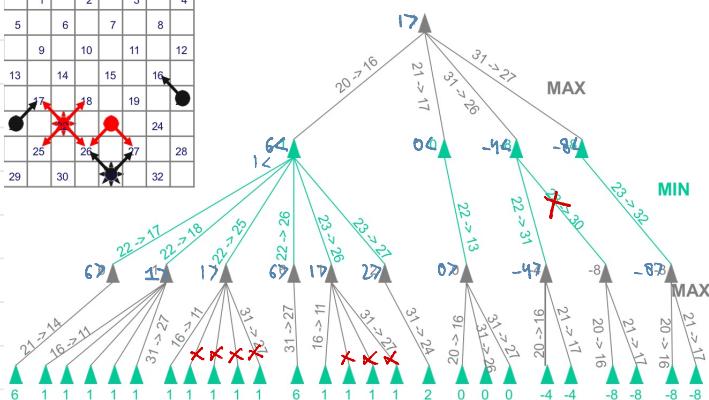
- Initial board configuration

- Black single on 20
 single on 21
 king on 31
- Red single on 23
 king on 22
- Evaluation function

$$E(s) = (5x_1 + x_2) - (5r_1 + r_2)$$
where
 x_1 = black king advantage,
 x_2 = black single advantage,
 r_1 = red king advantage,
 r_2 = red single advantage



Checkers MiniMax Example



-6 -5 -1 -3 -2 -1 12 3 15 6 7 8

S ← → B

ishma hafeez notes

repst
tree

Summary

- CSPs are a special kind of problem: states defined by values of a fixed set of variables, goal test defined by constraints on variable values
- Backtracking=depth-first search with one variable assigned per node
- Variable ordering and value selection heuristics help significantly
- Forward checking prevents assignments that lead to failure.
- Constraint propagation does additional work to constrain values and detect inconsistencies.
- The CSP representation allows analysis of problem structure. Tree structured CSPs can be solved in linear time.
- Iterative min-conflicts is usually effective in practice.

A) Explain the following in one or two statements.

- 1) **Zero-Sum game:** Agents have opposite utilities (values on outcomes), single value that one maximizes and the other minimizes
- 2) **Successor function** in min-max game: It defines what the legal moves a player can make are.
- 3) **Significance of game tree:** Game tree shows the best possible move of max and min player at each level.
- 4) **Types of constraints:** Unary Constraint and binary constraint.
- 5) **backtracking search** is used for a depth-first search that chooses values for one variable at a time and backtracks when a variable has no legal values left to assign.
- 6) **Binary CSP:** each constraint relates (at most) two variables
- 7) **Logic:** Logic is a formal system for manipulating facts so that true conclusion may be drawn
- 8) **Propositional logic:** concrete statements that are either true or false and **Predicate logic:** allows statements to contain variables, functions, and quantifiers
- 9) **Knowledge Base:** A set of statements that encode assertions about the world in a formal knowledge representation language.
- 10) Game theory is powerful tool for modeling cooperative behavior in many wireless networking applications such as cognitive radio networks, wireless system, physical layer security, virtual MIMO, among others.

ishma hafeez

notes

repsit
veet

Bayesian Network

↳ used when representing uncertainty

↳ uses a conditional Probability table (CPT)

↳ User Directed Acyclic Graph (DAG)

$$\text{eg } P(b|e, m) = \frac{P(b, e, m)}{P(e, m)}$$

ignored

= P sab kay parents

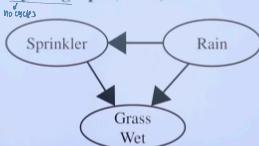
~~steps~~

↳ Make CBT of each

↳ simple of independent variable

↳ for dependent variable

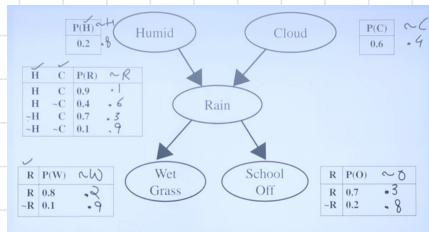
It is probabilistic graphical model that represents a set of variables and their conditional dependencies via a directed acyclic graph (DAG).



↳ depends on Rain

↳ independent

Sprinkler	Rain	Grass Wet
$\begin{array}{c cc} & S & \bar{S} \\ \hline F & 0.4 & 0.6 \\ T & 0.01 & 0.99 \end{array}$	$\begin{array}{c cc} & R & \bar{R} \\ \hline R & 0.2 & 0.8 \\ \bar{R} & 0.1 & 0.9 \end{array}$	$\begin{array}{c cc} & W & \bar{W} \\ \hline S, R & 0.0 & 1 \\ \sim S, R & 0.8 & 0.2 \\ S, \bar{R} & 0.9 & 0.1 \\ \sim S, \bar{R} & 0.99 & 0.01 \end{array}$



$$P(D, R, \bar{H}, C)$$

$$= P(D|R) \times P(R|\bar{H}, C) \times P(\bar{H}) \times P(C)$$

$$= 0.1 \times 0.1 \times 0.8 \times 0.6$$

$$P(W, S, R)$$

$$= P(W \cap S \cap R)$$

$$= P(W|S, R) \times P(S|R) \times P(R)$$

$$= 0.99 \times 0.01 \times 0.2$$

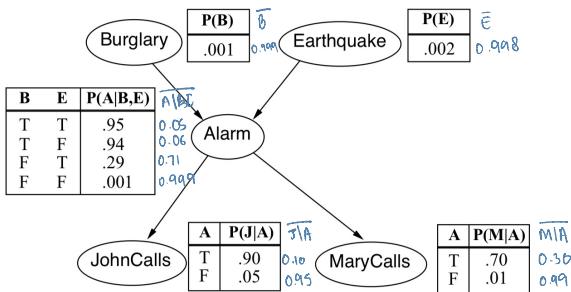
=

$$P(W, \bar{S}, \bar{R})$$

$$= P(W|\bar{S}, \bar{R}) \times P(\bar{S}|\bar{R}) \times P(\bar{R})$$

$$= 0.0 \times 0.6 \times 0.8$$

$$= 0$$



$$\varnothing) P(J, M, A | \bar{B}, \bar{E})$$

$$\begin{aligned} & \cdot P(J|A) \times P(M|A) \times P(A|\bar{B}, \bar{E}) \times P(\bar{B}) \times P(\bar{E}) \\ & \rightarrow 0.9 \times 0.7 \times 0.001 \times 0.999 \times 0.998 \end{aligned}$$

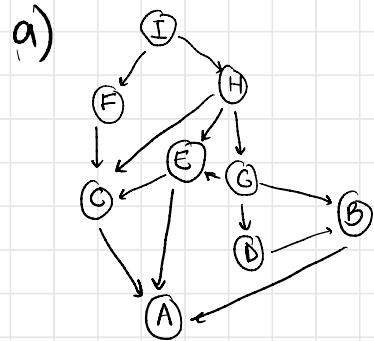
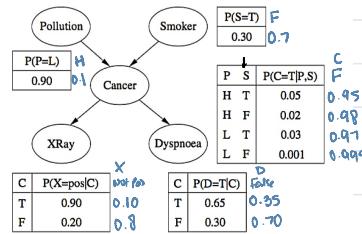
a) Draw the Bayesian Network that corresponds to the following conditional probability:

$$P(A | B, C, E) P(B | D, G) P(C | E, F, H) P(D | G) P(E | G, H) P(F | I) P(G | H) P(H | I) P(I)$$

b) Suppose we have a robot on a grid world with a noisy sensor and we know the distribution for P(Observation|Location) of sensor readings given the robot's grid square. We also have knowledge about the locations in the form of a prior, P(Location). However, we are really interested in the probability of being in a grid square given the sensor reading, P(Location|Observation). How could you compute this probability? Short answer.

c) Consider the Bayesian Network given below.

1. Calculate $P(\text{Cancer})$
2. Calculate $P(\text{Pollution} | \text{Cancer})$
3. Calculate $P(\text{Xray} | \text{Cancer})$
4. Calculate $P(\text{Dyspnoea} | \text{Cancer})$



b) $P(L|O) = \frac{P(O|L) \times P(L)}{P(O)}$ given

c i) $P(\text{Cancer}) = P(\text{Cancer} | S, P) \times P(S) \times P(P)$

C	S	P	$C=T S, P$	$S \times P$		
0	0	0	0.001	0.1	0.9	0.00063
0	0	1	0.02	0.1	0.1	0.0014
1	0	0	0.03	0.3	0.9	0.0081
1	1	0.05	0.3	0.1		0.0015
					\downarrow	0.01163

c ii) $P(P|C) = \frac{P(C|P) \times P(P)}{P(C)}$ $\therefore P(P=L) = 0.90$
 $\therefore P(C=T) = 0.01163$

$$\begin{aligned} P(c|P) &= \frac{P(C, P)}{P(P)} \\ &\rightarrow \frac{P(C|P, S) \times P(P) \times P(S)}{P(P)} \\ &= P(C|P, S) \times P(S) \end{aligned}$$

c iii) $P(X|C) = 0.9$

c iv) $P(D|C) = 0.65$

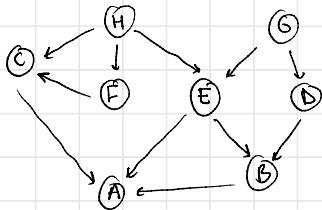
S	$C=T P=L, S$	S	
0	0.001	0.7	$= 0.0007$
1	0.03	0.3	$= \frac{0.009}{0.0097}$
			\downarrow

$$\begin{aligned} P(P|C) &= \frac{0.0097 \times 0.9}{0.01163} \\ &= 0.7506 \end{aligned}$$

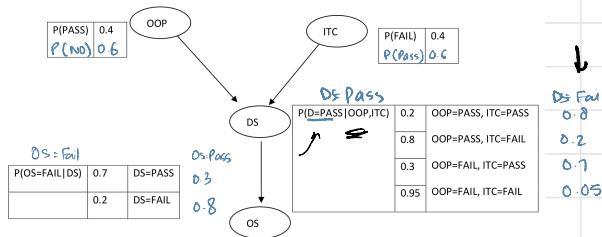
$P(X|C) = \frac{P(C|X) P(X)}{P(C)}$

$$\begin{aligned} P(C|X) &= \frac{P(C, X)}{P(X)} \\ &= \frac{P(C=T | P, S) \times P(P) \times P(S) \times P(X|C)}{P(X)} \end{aligned}$$

- a. Draw the Bayesian Network that corresponds to this conditional probability: $P(A | B, C, E)$
 $P(B | D, E) P(C | F, H) P(D | G) P(E | G, H) P(F | H) P(G) P(H)$.



- b. Calculate $P(OS=PASS | ITC=PASS)$ using the network given below:



$$\delta) P(OS=Pass | ITC=Pass)$$

$$P(OS=Pass, ITC=Pass)$$

$$P(ITC=Pass)$$

$$P(OOP|DS) \times P(DS|ITC, OOP) \times P(ITC) \times P(OOP)$$

$$P(OS)$$

DS	OOP	$(O=Pass DS)$	$DS ITC=Pass, OOP$	$P(OOP)$	
1	1	0.3	0.2	0.4	$= 0.024$
1	0	0.3	0.3	0.6	$= 0.054$
0	1	0.8	0.8	0.4	$= 0.256$
0	0	0.8	0.7	0.6	$= 0.336$
					$\underline{\underline{0.67}}$

Semantics of Bayesian Networks

- ↳ Representing full joint probability distribution
 - ↳ method for constructing Bayesian network
 - ↳ compactness and node ordering
- ↳ Conditional Independence relation in Bayesian network

↳ construction of Bayesian Networks from Probabilities.

1. No Parent

$$P(A) \quad \textcircled{A}$$

2. 1 Parent

$$P(A|B) \quad \textcircled{A} \leftarrow \textcircled{B}$$

3. 2 Parents

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

```
graph TD; B((B)) --> A((A)); C((C)) --> A
```

4. 3 Parents

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

```
graph TD; B((B)) --> A((A)); C((C)) --> A; D((D)) --> A
```

Global semantics

"Global" semantics defines the full joint distribution as the product of the local conditional distributions:

$$P(x_1, \dots, x_n) = \prod_{i=1}^n P(x_i | \text{parents}(X_i))$$

e.g., $P(j \wedge m \wedge a \wedge \neg b \wedge \neg e)$

$$\begin{aligned} &= P(j|a)P(m|a)P(a|\neg b, \neg e)P(\neg b)P(\neg e) \\ &= 0.9 \times 0.7 \times 0.001 \times 0.999 \times 0.998 \\ &\approx 0.00063 \end{aligned}$$



↓
Now we solve it

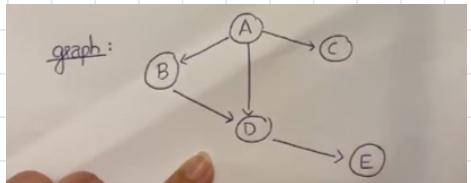
$$P(A=\text{on}) = 0.6$$

$$P(B=\text{on} | A) = \begin{cases} 0.1 & A=\text{off} \\ 0.9 & A=\text{on} \end{cases}$$

$$P(C=\text{on} | A) = \begin{cases} 0.8 & A=\text{off} \\ 0.5 & A=\text{on} \end{cases}$$

$$P(D=\text{on} | A, B) = \begin{cases} 0.1 & A=\text{off}, B=\text{off} \\ 0.9 & A=\text{on}, B=\text{off} \\ 0.3 & A=\text{off}, B=\text{on} \\ 0.95 & A=\text{on}, B=\text{on} \end{cases}$$

$$P(E=\text{on} | D) = \begin{cases} 0.8 & D=\text{off} \\ 0.1 & D=\text{on} \end{cases}$$



$$\begin{aligned} P(E=\text{on} | A=\text{on}) &= \frac{P(E=\text{on}, A=\text{on})}{P(A=\text{on})} \\ &= \frac{P(E=\text{on}|D) \times P(D|A=\text{on}, B) \times P(A=\text{on}) \times P(B|A=\text{on})}{P(A=\text{on})} \\ &= P(E=\text{on}|D) \times P(D|A=\text{on}, B) \times P(B|A=\text{on}) \end{aligned}$$

as we don't know whether B is on or off

D	B	P(E=on D)	P(D A=on, B)	P(B A=on)	
off	on	0.8	0.05	0.95	_____
on	off	0.1	0.9	0.05	_____
on	on	0.1	0.95	0.95	_____
off	off	0.8	0.1	0.05	_____
$(0.8 \times 0.05 + 0.95) + (0.1 \times 0.9 \times 0.05) + (0.1 \times 0.95 \times 0.95) + (0.8 \times 0.1 \times 0.05) = [0.13675]$					SUM

solved in previous question

$$(3) P(A=\text{on} | E=\text{on}) = \frac{P(E=\text{on}|A=\text{on}), P(A=\text{on})}{P(E=\text{on})}$$

$$P(A) = \frac{P(E=\text{on}|D) \times P(D|A,B) \times P(A) \times P(B|A)}{P(B)}$$

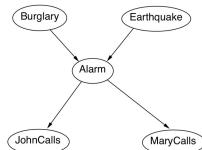
$$\begin{aligned} \text{Denominator: } P(E=\text{on}) &= P(E=\text{on}|D) \times P(D|A,B) \times P(A) \times P(B|A) \\ &= \frac{0.13675 \times 0.6}{P(E=\text{on})} \end{aligned}$$

A	B	D	P(E=on D)	P(D A,B)	P(A)	P(B A)	
on	on	off	0.8		0.6		
on	off	on		0.1		0.4	
on	on	on		0.9		0.6	
on	off	off		0.95		0.4	
off	on	off	0.8	0.7	0.4	0.1	0.13675
off	off	on	0.1	0.1	0.4	0.9	
off	on	on	0.1	0.3	0.7	0.1	
off	off	off	0.8	0.9	0.4	0.9	

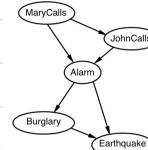
↳ COMPACTNESS

↳ A general property of locally structured systems

aka sparse



A	A	B, E		
M	J	A	B	E
2 ¹	2 ¹	2 ²	2 ⁰	2 ⁰
2	2	4	1	1
$= 10$				



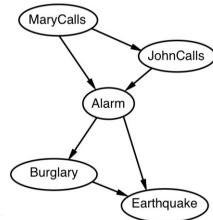
No

M	J	M, J	A	B	B, A
2 ⁰	2 ¹	2 ²	2 ¹	2 ¹	2 ²
= 1	2	4	2	4	= 13
$\text{Network is less compact, less complex}$					

↳ NODE ORDERING

Q) M, J, A, B, E
I
root

- Suppose we decide to add the nodes in the order MaryCalls, JohnCalls, Alarm, Burglary, Earthquake.
- Adding Mary Calls:** No parents. M
- Adding JohnCalls:** If Mary calls, that probably means the alarm has gone off, which of course would make it more likely that John calls. Therefore, JohnCalls needs Mary Calls as a parent
- Adding Alarm:** Clearly, if both call, it is more likely that the alarm has gone off than if just one or neither call, so we need both MaryCalls and JohnCalls as parents.
- Adding Burglary:** If we know the alarm state, then the call from John or Mary might give us information about our phone ringing or Mary's music, but not about burglary:
- P(Burglary|Alarm, JohnCalls, MaryCalls) = P(Burglary|Alarm)**
- Hence we need just Alarm as parent.
- Adding Earthquake:** if the alarm is on, it is more likely that there has been an E|B, A earthquake. But if we know that there has been a burglary, then that explains the alarm, and the probability of an earthquake would be only slightly above normal. Hence, we need both Alarm and Burglary as parents.



$$P(J|M) = P(J)? \text{ No}$$

$$P(A|J, M) = P(A|J)? \quad P(A|J, M) = P(A)? \text{ No}$$

$$P(B|A, J, M) = P(B|A)? \text{ Yes}$$

$$P(B|A, J, M) = P(B)? \text{ No}$$

$$P(E|B, A, J, M) = P(E|A)? \text{ No}$$

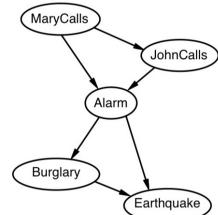
$$P(E|B, A, J, M) = P(E|A, B)? \text{ Yes}$$

↳ Conditional independence

- ↳ A node is conditionally independent of its non descendants, given its parents

$$X_i \perp\!\!\!\perp \text{non descendants}(X_i) \mid \text{Parents}(X_i)$$

↓ independent ↓ given



1. J independent of B,E given A

- ↳ J conditionally independent of B,E

- ↳ J dependent on A

2. B independent of J,M given A

- ↳ B conditionally independent of J,M

- ↳ B dependent on A,E

J

?

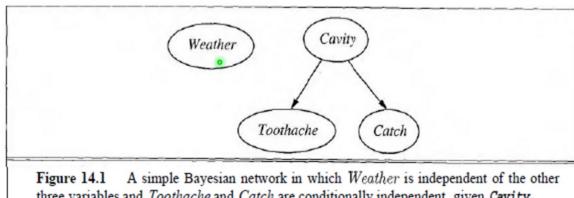


Figure 14.1 A simple Bayesian network in which Weather is independent of the other three variables and Toothache and Catch are conditionally independent, given Cavity.

W independent of other 3 variables

Toothache conditionally independent of Catch, given cavity

Toothache and catch independent to each other

but both depend on cavity

D Separation

X, Y = independent variable

E = set of observed variables

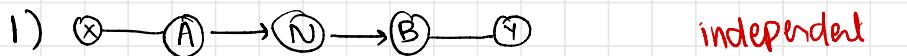
↳ E d-separates X and Y

If E blocks every un-directed Path b/w X and Y

any arrow
doesn't matter direction

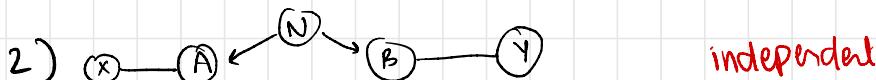
↳ If E d-separates X and Y

then X and Y are independent, given E



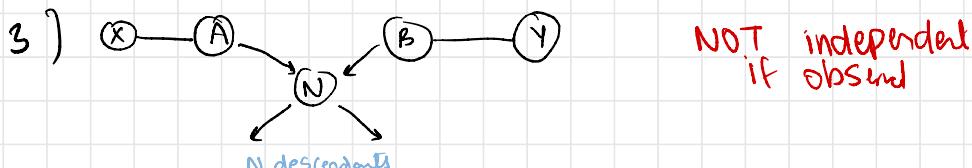
if N is Observed

it blocks the Path b/w X and Y



if N is Observed

it blocks the Path b/w X and Y

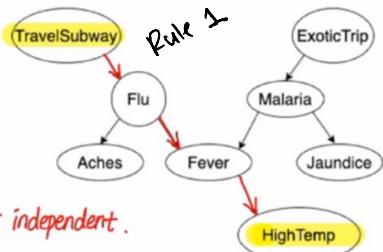


If N and N 's descendants are **NOT** observed

then they block the Path b/w X and Y

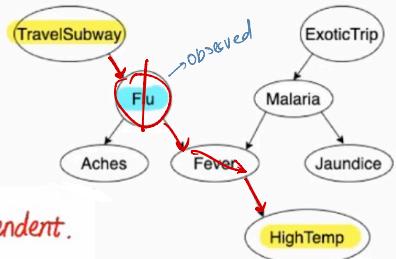
CQ: Applying D-separation

CQ 1a: Are **TravelSubway** and **HighTemp** independent?



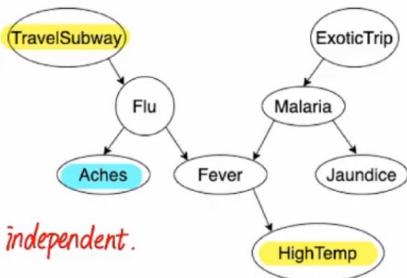
Not independent.

CQ 1b: Are **TravelSubway** and **HighTemp** conditionally independent given **Flu**?



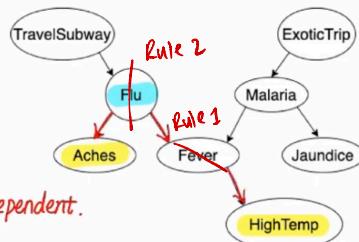
Independent.

CQ 1c: Are **TravelSubway** and **HighTemp** conditionally independent given **Aches**?



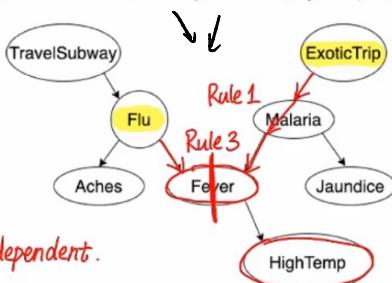
Not independent.

CQ 2b: Are **Aches** and **HighTemp** conditionally independent given **Flu**?



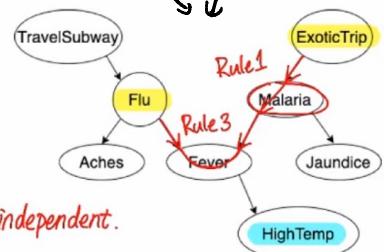
Independent.

CQ 3a: Are **Flu** and **ExoticTrip** conditionally independent?



Independent.

CQ 3b: Are **Flu** and **ExoticTrip** conditionally independent given **HighTemp**?



Not independent.

grey = dependent area

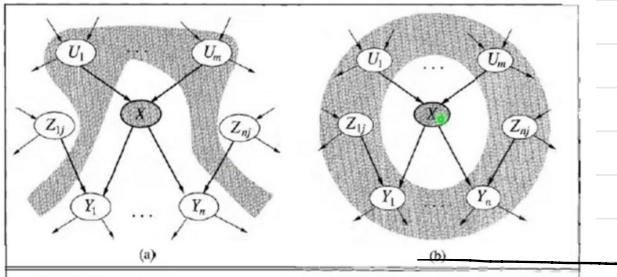
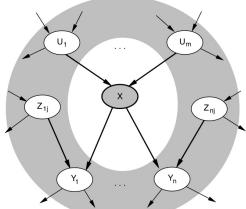


Figure 14.4 (a) A node X is conditionally independent of its non-descendants (e.g., the Z_{ij} 's) given its parents (the U_i 's shown in the gray area). (b) A node X is conditionally independent of all other nodes in the network given its Markov blanket (the gray area).

- a) X independent to Z_{ij} given U_i
 ↳ X conditionally independent of Z_{ij}
 ↳ X dependent on U_i
- b) X independent to all nodes
 ↳ X conditionally independent of all nodes

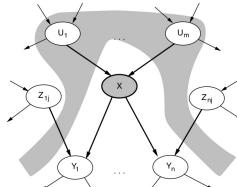
Markov blanket

Each node is conditionally independent of all others given its
 Markov blanket: parents + children + children's parents



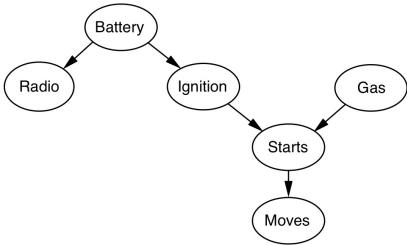
Local semantics

Local semantics: each node is conditionally independent of its nondescendants given its parents



Theorem: Local semantics \Leftrightarrow global semantics

Example



Are Gas and Radio independent? Given Battery? Ignition? Starts? Moves?

Compact conditional distributions

- CPT grows exponentially with number of parents
- CPT becomes infinite with continuous-valued parent or child
- Solution: canonical distributions that are defined compactly
- Deterministic nodes are the simplest case:
 $X = f(\text{Parents}(X))$ for some function f
- E.g., Boolean functions
 $\text{NorthAmerican} \Leftrightarrow \text{Canadian} \vee \text{US} \vee \text{Mexican}$
- E.g., numerical relationships among continuous variables
 $\frac{\partial \text{Level}}{\partial t} = \text{inflow} + \text{precipitation} - \text{outflow} - \text{evaporation}$

Compact conditional distributions contd.

Noisy-OR distributions model multiple noninteracting causes

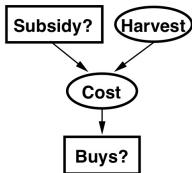
- Parents $U_1 \dots U_k$ include all causes (can add leak node)
- Independent failure probability q_i for each cause alone
 $\Rightarrow P(X|U_1 \dots U_j, \neg U_{j+1} \dots \neg U_k) = 1 - \prod_{i=1}^j q_i$

<i>Cold</i>	<i>Flu</i>	<i>Malaria</i>	$P(\text{Fever})$	$P(\neg \text{Fever})$
F	F	F	0.0	1.0
F	F	T	0.9	0.1
F	T	F	0.8	0.2
F	T	T	0.98	$0.02 = 0.2 \times 0.1$
T	F	F	0.4	0.6
T	F	T	0.94	$0.06 = 0.6 \times 0.1$
T	T	F	0.88	$0.12 = 0.6 \times 0.2$
T	T	T	0.988	$0.012 = 0.6 \times 0.2 \times 0.1$

Number of parameters linear in number of parents

Hybrid (discrete+continuous) networks

Discrete (*Subsidy?* and *Buys?*); continuous (*Harvest* and *Cost*)



Option 1: discretization—possibly large errors, large CPTs
 Option 2: finitely parameterized canonical families

- 1) Continuous variable, discrete+continuous parents (e.g., *Cost*)
- 2) Discrete variable, continuous parents (e.g., *Buys?*)

Continuous child variables

Need one conditional density function for child variable given continuous parents, for each possible assignment to discrete parents

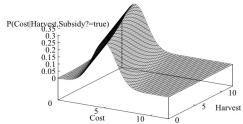
Most common is the linear Gaussian model, e.g.,:

$$\begin{aligned}
 P(\text{Cost} = c | \text{Harvest} = h, \text{Subsidy?} = \text{true}) \\
 &= N(a_t h + b_t, \sigma_t)(c) \\
 &= \frac{1}{\sigma_t \sqrt{2\pi}} \exp\left(-\frac{1}{2} \left(\frac{c - (a_t h + b_t)}{\sigma_t} \right)^2\right)
 \end{aligned}$$

Mean *Cost* varies linearly with *Harvest*, variance is fixed

Linear variation is unreasonable over the full range
 but works OK if the likely range of *Harvest* is narrow

Continuous child variables

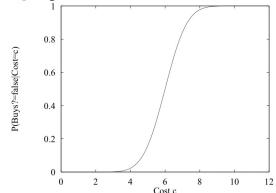


All-continuous network with LG distributions
 \Rightarrow full joint distribution is a multivariate Gaussian

Discrete+continuous LG network is a **conditional Gaussian** network i.e., a multivariate Gaussian over all continuous variables for each combination of discrete variable values

Discrete variable w/ continuous parents

Probability of *Buy?* given *Cost* should be a "soft" threshold:



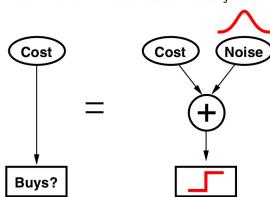
Probit distribution uses integral of Gaussian:

$$\Phi(x) = \int_{-\infty}^x N(0, 1)(x) dx$$

$$P(\text{Buy?} = \text{true} \mid \text{Cost} = c) = \Phi((-c + \mu)/\sigma)$$

Why the probit?

1. It's sort of the right shape
2. Can view as hard threshold whose location is subject to noise

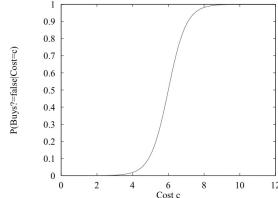


Discrete variable contd.

Sigmoid (or logit) distribution also used in neural networks:

$$P(\text{Buy?} = \text{true} \mid \text{Cost} = c) = \frac{1}{1 + \exp(-2\frac{c+\mu}{\sigma})}$$

Sigmoid has similar shape to probit but much longer tails:



Summary

Bayes nets provide a natural representation for (causally induced) conditional independence

Topology + CPTs = compact representation of joint distribution

Generally easy for (non)experts to construct

Canonical distributions (e.g., noisy-OR) = compact representation of CPTs

Continuous variables \Rightarrow parameterized distributions (e.g., linear Gaussian)

BAYES THEOREM

conditional probability

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

A will occur
Given that B already

BAYES RULE

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

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

$$\hookrightarrow P(A|B) P(B) = P(B|A) P(A)$$

Examples

Suppose there are three bowls B1, B2, B3 and bowl B1 has 2 red and 4 blue coins; bowl B2 has 1 red and 2 blue coins; bowl B3 contains 5 red and 4 blue coins. Suppose the probabilities for selecting the bowls is not the same but are:-

- $P(B1) = 1/3$
- $P(B2) = 1/6$
- $P(B3) = 1/2$

Now, let us compute, assuming that a red coin was drawn what will be the probability that it came from bowl B1.



$$\textcircled{1} \quad P(R) = P(B1 \cap R) + P(B2 \cap R) + P(B3 \cap R)$$

$$= P(\text{selecting } B1) * P(\text{Number of Red coins / total number of coins in } B1) + \\ P(\text{selecting } B2) * P(\text{Number of Red coins in } B2 / \text{total number of coins in } B2) \\ + P(\text{selecting } B3) * P(\text{Number of Red coins in } B3 / \text{total number of coins in } B3)$$

$$= \frac{1}{3} \times \frac{2}{6} + \frac{1}{6} \times \frac{1}{3} + \frac{1}{2} \times \frac{5}{9}$$

$$= \frac{4}{9}$$

$$\textcircled{2} \quad P(R|B1) = \frac{P(B1 \cap R)}{P(B1)} = \frac{\left(\frac{1}{3} \times \frac{2}{6} \right)}{4/9}$$

$$= \frac{2}{8}$$

NAIVE BAYES CLASSIFIER

↳ supervised learning

↳ based on applying Bayes Theorem

↳ assumes no dependencies b/w features

BAYES RULE

$$P(Y|X) = \frac{P(X|Y) \times P(Y)}{P(X)}$$

$$P(Y|X_1, X_2, \dots, X_n) = P(X_1|Y) \times P(X_2|Y) \dots P(X_n|Y) \times P(Y)$$

Person	X_1 (Yes/No)	X_2 (Yes/No)	Fever (Yes/No)	
1	Yes	No	Yes	Step-1: Prior Probability: $P(\text{fever} = \text{yes}) = 7/10$ $P(\text{fever} = \text{no}) = 3/10$
2	No	Yes	Yes	
3	Yes	Yes	Yes	
4	No	No	No	
5	Yes	No	Yes	Step-2: Conditional Probability:
6	No	No	Yes	
7	Yes	No	Yes	
8	Yes	No	No	
9	No	Yes	Yes	
10	No	Yes	No	

Given Person(flu, Covid)
 $P(\text{yes} | \text{flu, Covid}) = P(\text{flu} | \text{yes}) * P(\text{covid} | \text{yes}) * P(\text{yes})$

$P(\text{NO} | \text{flu, Covid}) = P(\text{flu} | \text{NO}) * P(\text{covid} | \text{NO}) * P(\text{NO})$

Q) Given Person(flu, Covid), is fever Yes or No

$$P(\text{yes} | \text{flu, Covid})$$

$$\cdot P(\text{flu} | \text{yes}) \times P(\text{covid} | \text{yes}) \times P(\text{yes})$$

$$= \frac{3}{7} \times \frac{4}{7} \times \frac{7}{10}$$

$$= 0.17$$

↓

$$P(\text{NO} | \text{flu, Covid})$$

$$\cdot P(\text{flu} | \text{NO}) \times P(\text{covid} | \text{NO}) \times P(\text{NO})$$

$$= \frac{2}{3} \times \frac{2}{3} \times \frac{3}{10}$$

$$= 0.13$$

Higer prob so prediction is Yes fever

fruit = {Yellow, Sweet, long}

Fruit	Yellow	Sweet	Long	Total
Orange	350	450	0	650
Banana	400	300	350	400
Others	50	100	50	150
Total	800	850	400	1200

$$P(\text{Yellow}|\text{orange}) = \frac{P(\text{orange}|\text{yellow}) P(\text{yellow})}{P(\text{orange})}$$

$$= \frac{350}{800} \times \frac{800}{1200} = 0.53$$

$$\frac{650}{1200}$$

$$P(\text{Sweet}|\text{orange}) = \frac{P(\text{o|s}) P(s)}{P(o)} = 0.69$$

$$P(\text{fruit}|\cancel{\text{orange}}) = \frac{P(y|o) \times P(s|o) \times P(l|o)}{P(o)} = 0.53 \times 0.69 \times 0 = 0$$

$$P(\text{fruit}|\text{Banana}) = \frac{P(y|B) \times P(s|B) \times P(l|B)}{P(B)} = 1 \times 0.75 \times 0.89 = 0.65$$

$$P(\text{fruit}|others) = \frac{P(y|o) \times P(s|o) \times P(l|o)}{P(o)} = 0.33 \times 0.66 \times 0.33 = 0.072$$

TYPES OF NAIVE BAYES ALGO

1. Gaussian

- ↳ used in classification
- ↳ assumes features follow a normal distribution

2. Bernoulli

- ↳ binomial model
- ↳ used when features are binary

e.g. text classification with a bag of words

3. Multinomial Naive Bayes

- ↳ used on multinomially distributed data

For example, let's say, we have a text classification problem. Here we can consider bernoulli trials which is one step further and instead of "word occurring in the document", we have "count how often word occurs in the document", you can think of it as "number of times outcome number x_{-i} is observed over the n trials".

What to Do About Errors

- Problem: there's still spam in your inbox
- Need more **features** – words aren't enough!
 - Have you emailed the sender before?
 - Have 1M other people just gotten the same email?
 - Is the sending information consistent?
 - Is the email in ALL CAPS?
 - Do inline URLs point where they say they point?
 - Does the email address you by (your) name?
- Naïve Bayes models can incorporate a variety of features, but tend to do best in homogeneous cases (e.g. all features are word occurrences)

DECISION TREES

Decision Tree

- ↳ It is a tree structure classifier
- ↳ A graphical representation of all possible decisions related to a problem
- ↳ Used in
 - ↳ Classification → Using decision tree classifiers
 - ↳ Regression

Example of a Decision Tree

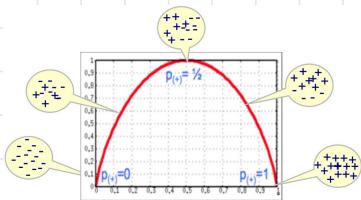


How to learn such a tree from past experience?

Entropy

- ↳ tells how pure/impure a data set is

$$E\{+x, -y\} = -\left(\frac{N_{++}}{N} \log_2 \frac{N_{++}}{N}\right) - \left(\frac{N_{+-}}{N} \log_2 \frac{N_{+-}}{N}\right)$$



Information Gain

Entropy measures purity at each node, information gain looks at all nodes together and the expected drop in entropy after split.

Gain(S, A) = expected reduction in entropy due to sorting on A

$$\text{Gain}(S, A) = \text{Entropy}(S) - \sum_{v \in \text{Values}(A)} \frac{|S_v|}{|S|} \cdot \text{Entropy}(S_v)$$

Maximum Gain(S, A) is selected!

Example	$\begin{matrix} 9 \text{ Yes} & 5 \text{ No} \\ \text{Wind} & \end{matrix}$ $E=0.94$ $\text{Gain}(S, \text{Wind}) = 0.94 - (8/14) \cdot 0.81 - (6/14) \cdot 1 = 0.048$
	$\begin{matrix} 6/2 & 3/3 \\ \text{Weak} & \text{Strong} \\ E=0.81 & E=1 \end{matrix}$

16

PROS

↳ Interpretable

- ↳ humans can understand decisions
- ↳ Easily handles irrelevant attributes e.g. $G=0$
- ↳ very compact
- ↳ no. of nodes $\ll D$ after pruning
- ↳ fast at testing time
- ↳ $O(D \cdot \text{Depth})$

CONS

↳ Greedy Approach

- ↳ may not find best tree
- ↳ instances are represented by attribute pairs
- ↳ what if we have discrete input values
- ↳ target function has discrete output values
- ↳ hence can't have continuous no. output values
- ↳ training data may contain errors / miss attributes
- ↳ uncertainty in the data

Uncertainty in the data (e.g., suppose we have two exact days/features, one with "yes" and one with "no". ➔ no classifier can help in such totally uncertain data.)

Outlook: sunny

Wind: strong

Yes: no

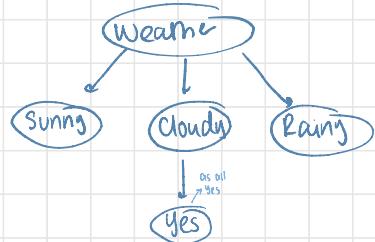
IG of entire dataset (S)

$(S, \text{Weather}) = 0.246 \rightarrow \text{MAX} \rightarrow \text{Hence ROOT NODE}$

$(S, \text{Temp}) = 0.029$

$(S, \text{Humid}) = 0.15$

$(S, \text{Wind}) = 0.0478$



FOR SUNNY

Day	Weather	Temperature	Humidity	Wind	Play Football?
Day 1	Sunny	Hot	High	Weak	No
Day 2	Sunny	Hot	High	Strong	No
Day 8	Sunny	Mild	High	Weak	No
Day 9	Sunny	Cool	Normal	Weak	Yes
Day 11	Sunny	Mild	Normal	Strong	Yes

Calculate IG of Temperature

Step1: Entropy of Sunny {+2,-3} = $-\frac{2}{5} \log_2 \frac{2}{5} - \frac{3}{5} \log_2 \frac{3}{5} = 0.97$

Step2: Entropy of all attributes:

Entropy of Hot {+0,+2} = $-\frac{0}{2} \log_2 \frac{0}{2} - \frac{2}{2} \log_2 \frac{2}{2} = 0$

Entropy of Mild {+1,-1} = $-\frac{1}{2} \log_2 \frac{1}{2} - \frac{1}{2} \log_2 \frac{1}{2} = 1.0$

Entropy of Cool {+1,-0} = $-\frac{1}{1} \log_2 \frac{1}{1} - \frac{0}{1} \log_2 \frac{0}{1} = 0$

Information Gain = Entropy(Sunny) - $\frac{2}{5} \text{Ent(H)} - \frac{2}{5} \text{Ent(M)} - \frac{1}{5} \text{Ent(C)}$
 $= 0.57$

Calculate IG of Humidity

Step1: Entropy of Sunny {+2,-3} = $-\frac{2}{5} \log_2 \frac{2}{5} - \frac{3}{5} \log_2 \frac{3}{5} = 0.97$

Step2: Entropy of all attributes:

Entropy of High {+0,-3} = $-\frac{0}{3} \log_2 \frac{0}{3} - \frac{3}{3} \log_2 \frac{3}{3} = 0$

Entropy of Normal {+2,-0} = $-\frac{2}{2} \log_2 \frac{2}{2} - \frac{0}{2} \log_2 \frac{0}{2} = 0$

Information Gain = Entropy(Sunny) - $\frac{3}{5} \text{Ent(H)} - \frac{2}{5} \text{Ent(M)}$
 $= 0.97$

Calculate IG of Wind

Step1: Entropy of Sunny {+2,-3} = $-\frac{2}{5} \log_2 \frac{2}{5} - \frac{3}{5} \log_2 \frac{3}{5} = 0.97$

[Step2: Entropy of all attributes:

Entropy of Strong {+1,-1} = $-\frac{1}{2} \log_2 \frac{1}{2} - \frac{1}{2} \log_2 \frac{1}{2} = 1$

Entropy of Weak {+1,-2} = $-\frac{1}{3} \log_2 \frac{1}{3} - \frac{2}{3} \log_2 \frac{2}{3} = 0.918$

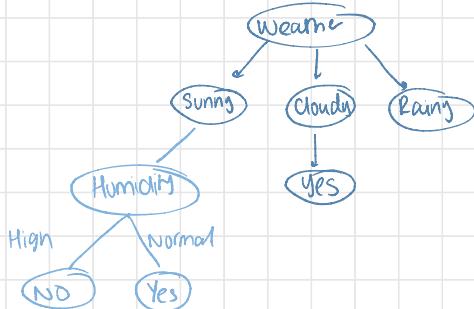
Information Gain = Entropy(Sunny) - $\frac{2}{5} \text{Ent(S)} - \frac{3}{5} \text{Ent(W)}$
 $= 0.019$

IG Sunny

$(\text{Sunny}, \text{Temp}) = 0.51$

$(\text{Sunny}, \text{Humid}) = 0.91 \rightarrow \text{MAX}$

$(\text{Sunny}, \text{Wind}) = 0.019$



FDR RAIN

Day	Weather	Temperature	Humidity	Wind	Play Football?
Day 4	Rain	Mild	High	Weak	Yes
	Rain	Cool	Normal	Weak	Yes
Day 6	Rain	Cool	Normal	Strong	No
Day 10	Rain	Mild	Normal	Weak	Yes
	Rain	Mild	High	Strong	No

IG Rain

$$(Rain, Temp) = 0.019$$

$$(Rain, Humid) = 0.019$$

$$(Rain, Wind) = 0.97 \rightarrow \text{Max}$$

Calculate IG of Temperature

Step1: Entropy of Rain $\{+3,-2\} = -\frac{3}{5} \log_2 \frac{3}{5} - \frac{2}{5} \log_2 \frac{2}{5} = 0.97$

Step2: Entropy of all attributes:

Entropy of Hot $\{+0,-2\} = -\frac{0}{2} \log_2 \frac{0}{2} - \frac{0}{2} \log_2 \frac{0}{2} = 0$

Entropy of Mild $\{+2,-1\} = -\frac{2}{3} \log_2 \frac{2}{3} - \frac{1}{3} \log_2 \frac{1}{3} = 0.918$

Entropy of Cool $\{+1,-1\} = -\frac{1}{2} \log_2 \frac{1}{2} - \frac{1}{2} \log_2 \frac{1}{2} = 1.0$

Information Gain = Entropy(Rain) $- \frac{0}{5} \text{Ent}(H) - \frac{3}{5} \text{Ent}(M) - \frac{2}{5} \text{Ent}(C)$
 $= 0.019$

Calculate IG of Humidity

Step1: Entropy of Rain $\{+3,-2\} = -\frac{3}{5} \log_2 \frac{3}{5} - \frac{2}{5} \log_2 \frac{2}{5} = 0.97$

Step2: Entropy of all attributes:

Entropy of High $\{+1,-1\} = -\frac{1}{2} \log_2 \frac{1}{2} - \frac{1}{2} \log_2 \frac{1}{2} = 1$

Entropy of Normal $\{+2,-1\} = -\frac{2}{3} \log_2 \frac{2}{3} - \frac{1}{3} \log_2 \frac{1}{3} = 0.918$

Information Gain = Entropy(Rain) $- \frac{2}{5} \text{Ent}(S) - \frac{3}{5} \text{Ent}(W)$
 $= 0.019$

Calculate IG of Wind

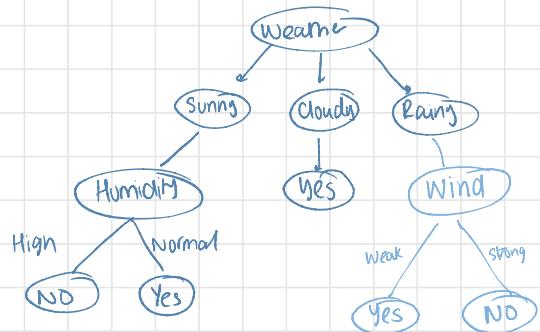
Step1: Entropy of Rain $\{+3,-2\} = -\frac{3}{5} \log_2 \frac{3}{5} - \frac{2}{5} \log_2 \frac{2}{5} = 0.97$

Step2: Entropy of all attributes:

Entropy of Strong $\{+0,-2\} = -\frac{0}{2} \log_2 \frac{0}{2} - \frac{2}{2} \log_2 \frac{2}{2} = 0$

Entropy of Weak $\{+3,-0\} = -\frac{3}{3} \log_2 \frac{3}{3} - \frac{0}{3} \log_2 \frac{0}{3} = 0$

Information Gain = Entropy(Rain) $- \frac{2}{5} \text{Ent}(S) - \frac{3}{5} \text{Ent}(W)$
 $= 0.97$



In this problem, the following dataset has been used to learn a decision tree which predicts that if students pass **Artificial Intelligence** course (Yes or No), based on their Study (Yes, No) and previous Grade (High, Medium, or Low).

Study	Grade	Pass
Yes	Low	Yes
No	Medium	No
No	Low	No
Yes	Medium	Yes
No	High	No
Yes	High	Yes

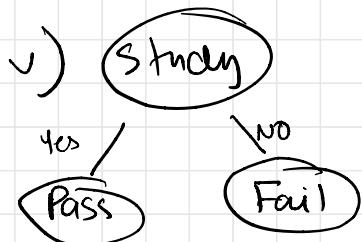
i. What is the entropy $H(\text{Pass})?$ ≈ 1

ii. What is the entropy $H(\text{Pass} | \text{Grade})?$ ≈ 1

iii. What is the entropy $H(\text{Pass} | \text{Study})?$ ≈ 0

iv. Which attribute would you consider as the root node? What was the information gain of the attribute you chose as the root node? **Study**

v. Draw the decision tree.



Page 7 of 8

↳ entire dataset entropy

$$S = \{-3, +3\} = -\left(\frac{3}{6} \log_2 \frac{3}{6}\right) - \left(\frac{3}{6} \log_2 \frac{3}{6}\right) = 1$$

IG Grade

↳ entropy of all variables

$$\text{High: } \{+1, -1\} = -\left(\frac{1}{2} \log_2 \frac{1}{2}\right) - \left(\frac{1}{2} \log_2 \frac{1}{2}\right) = 1$$

$$\text{Medium: } \{+1, -1\} = -\left(\frac{1}{2} \log_2 \frac{1}{2}\right) - \left(\frac{1}{2} \log_2 \frac{1}{2}\right) = 1$$

$$\text{Low: } \{+1, -1\} = -\left(\frac{1}{2} \log_2 \frac{1}{2}\right) - \left(\frac{1}{2} \log_2 \frac{1}{2}\right) = 1$$

↳ Weighted entropy

$$\frac{2}{6}(1) + \frac{2}{6}(1) + \frac{2}{6}(1) = 1$$

$$\hookrightarrow \text{IG: } 1 - \frac{2}{6}(1) - \frac{2}{6}(1) - \frac{2}{6}(1) = 0$$

IG Study

↳ entropy of each variable

$$\text{Yes: } \{+3, -3\} = 0$$

$$\text{No: } \{0, -3\} = 0$$

↳ Weighted entropy

$$0 \cdot 0 = 0$$

$$\hookrightarrow \text{IG} = 1 - 0 = 1$$

↓
Root node

A) What is Machine Learning? Name its types and discuss the difference between classification and clustering.

B) Following data is being provided from 10 patients, having some symptoms similar to that of Covid-19. Use this data to design and draw a decision tree for an automated system

which can decide if a patient should be tested for covid-19 or not. Also, show the calculations of Entropy and Gain.

Patient ID	Headache	Cough	Fever	Flu	Pain in body	Diagnosis
1	Yes	Yes	Yes	Yes	Yes	Covid-19
2	No	No	No	Yes	Yes	Allergy
3	Yes	Yes	No	Yes	No	Cold
4	Yes	No	Yes	No	No	Covid-19
5	No	Yes	No	Yes	No	Cold
6	No	No	No	Yes	No	Allergy
7	No	No	Yes	No	No	Covid-19
8	Yes	No	No	Yes	Yes	Allergy
9	No	Yes	No	Yes	Yes	Cold
10	Yes	Yes	No	Yes	Yes	Cold

$$\text{Entropy (all values)} = \{1, 3, -1\} : -\left(\frac{3}{10} \log_2 \frac{3}{10}\right) - \left(\frac{7}{10} \log_2 \frac{7}{10}\right) = 0.88$$

16 Headache

$$\text{Yes } \{2, -3\} : -\left(\frac{2}{5} \log_2 \frac{2}{5}\right) - \left(\frac{3}{5} \log_2 \frac{3}{5}\right) = 0.91$$

$$\text{No } \{1, 4\} : -\left(\frac{1}{5} \log_2 \frac{1}{5}\right) - \left(\frac{4}{5} \log_2 \frac{4}{5}\right) = 0.72$$

$$IG = 0.88 - \frac{5}{10} (0.91) - \frac{5}{10} (0.72) = 0.035$$

16 Cough

$$\text{Yes } \{1, 4\} : -\left(\frac{1}{5} \log_2 \frac{1}{5}\right) - \left(\frac{4}{5} \log_2 \frac{4}{5}\right) = 0.72$$

$$\text{No } \{2, -3\} : -\left(\frac{2}{5} \log_2 \frac{2}{5}\right) - \left(\frac{3}{5} \log_2 \frac{3}{5}\right) = 0.91$$

$$IG = 0.88 - \frac{5}{10} (0.72) - \frac{5}{10} (0.91) = 0.035$$

16 Fever

DO Rest

MACHINE LEARNING ALGOS

1. SUPERVISED LEARNING

- ↳ learns mapping b/w input and target variables

L. Classification

- ↳ predicts a class label

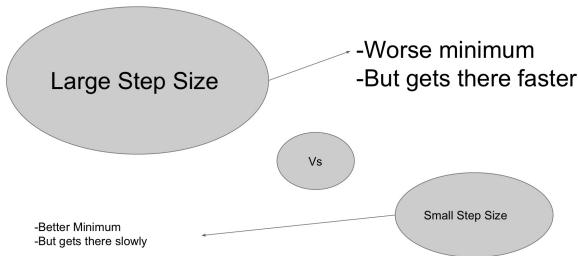
R. Regression

- ↳ predicts a numerical label

2. UNSUPERVISED LEARNING

- ↳ it learns on its own
- ↳ no input given

Optimizers are our friends



LINEAR REGRESSION

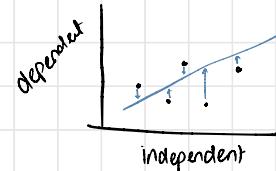
Regression

↳ Predicts a continuous outcome variable

- Predicting the value of a company's future stock price using its past and existing financial info
- Predicting the amount of rainfall
- Predicting ...
- Key difference from classification**
- We measure prediction errors differently
- This leads us to quite different learning models and algorithms

Linear Regression

↳ Predicts by plotting a straight line that fits well to the dataset



$$y = mx + c$$

Annotations for the equation:

- m : slope
- c : intercept → value of y when $x=0$
- dependent Variable
- independent Variable

$$1. \bar{X} = \frac{\text{Sum } X}{N}$$

Annotations for the mean calculation:

- mean x
- mean y

$$2. \bar{Y} = \frac{\text{Sum } Y}{N}$$

Annotations for the mean calculation:

- $D(x)$ Deviation (x)
- $D(y)$ Deviation (y)

	X	Y	$X - \bar{X}$	$Y - \bar{Y}$	Product of Deviations	$(D(x))^2$	y_{predict}	$(Y - \text{Predict})$
1.								
2.								
3.								

$$4. m = \frac{\text{Product } D}{\text{Sum } D}$$

Annotations for the slope calculation:

- gradient

$$5. c = \bar{Y} - m \bar{X}$$

Product D

Sum D

Sum
↓
error

Diameter(X) In Inches	Price(Y) In Dollars	Mean(X)	Mean(Y)	Deviations(X)	Deviations(Y)	Product of Deviations	Sum of Product of Deviations	Square of Deviations for X
8	10	10	13	-2	-3	6	12	4
10	13			0	0	0	0	
12	16			2	3	6	12	4

Calculate $m = \text{Sum of product of deviations} / \text{Sum of square of deviation for } X$

Calculate $b = \text{Mean of } Y - (m * \text{Mean of } X)$

$m = \frac{12}{8} = 1.5$

$b = 13 - (1.5 \times 10) = -2$

$c = 13 - (1.5 \times 8) = -2$

A scatter plot with a horizontal axis labeled from 0 to 14 and a vertical axis labeled from 0 to 18. Three data points are plotted at (8, 10), (10, 13), and (12, 16). A straight line is drawn through these points, representing the linear regression model.

TABLE. B.

[Use $x = r \cos \theta$ and $y = r \sin \theta$ to transform angles into the x and y coordinates, and radius = 10].

Table.B Channels and angle values

Channels	Angles (in degree)
A	34
B	45
C	55
D	48
E	49

X	Y	\bar{x}	\bar{y}	$D(x)$	$D(y)$	$D(x) \times D(y)$	$(D(x))^2$
8.29	5.59	1.43	-1.57	-2.2451	2.0449		
7.07	7.07	0.21	-0.094	-0.01974	0.0441		
5.73	8.19	-1.13	1.026	-1.159	1.276		
6.69	7.43	-0.11	0.266	-0.045	0.0289		
6.56	7.54	-0.3	0.376	-0.1128	0.09		
				-3.58164	3.4839		

$$\bar{x} = \frac{8.29 + 7.07 + 5.73 + 6.69 + 6.56}{5} = 6.86$$

$$\bar{y} = \frac{5.59 + 7.07 + 8.19 + 7.43 + 7.54}{5} = 7.164$$

$$m = \frac{-3.58164}{3.4839} = -1.028$$

$$C = \bar{y} - m\bar{x}$$

$$C = 7.164 - (-1.028)(6.86) = 14.216$$

$$y = -1.028x + 14.216$$

21K4638

Date:

???

no of weeks	add cars	\bar{x}	y	$(x - \bar{x})$	$(y - \bar{y})$	y_{Predict}	$(y - y_{\text{Predict}})^2$
168	212	-29.33	-61.1	42.6	212.272	3567.43	
428	300	230.67	198.6	70.6	363.592	4013.41	
299	311	98.67	66.6	81.6	289.514	505.6	
392	365	194.67	162.6	135.6	342.64	492.9	
80	167	-117.33	-149.4	-62.4	161.056	35.3	
56	149	-141.33	-173.4	-80.4	141.088	3.655	
352	366	154.67	122.6	-136.6	319.36	2175.28	
444	310	246.67	214.6	80.6	372.904	3956.91	
168	192	-29.33	-61.1	-37.4	212.272	410.9	
200	229	2.67	-29.4	-0.4	230.896	3.59	
4	88	-193.33	-255.9	-141.4	116.824	830.8	
52	118	-145.33	-175.9	-111.4	144.76	716.09	
20	62	-177.33	-201.9	-167.4	126.136	4112.6	
228	319	30.67	-1.4	89.6	247.192	5156.38	
72	193	-125.33	-157.9	-36.4	156.4	1339.56	

sum = 27587

↓
error

$$\bar{x} = 197.33 \quad \bar{y} = 229.4$$

$$x = \frac{2960}{45} \quad y = \frac{3441}{45}$$

$$y = a + bx$$

$$a = 114.496$$

$$b = 0.582 \quad y = 114.496 + 0.582x$$

$$r = 0.895$$

Select

1. Gradient Descent Algo

↳ optimises y-intercept and m → for Linear Regression

Repeat until convergence {

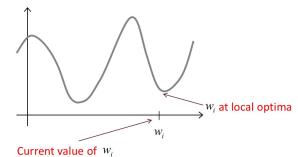
$$w_i := w_i - \alpha \frac{\partial}{\partial w_i} J(w_0, w_1) \quad \text{for } i = 0 \text{ and } i = 1$$

}

Learning rate

Partial derivative

What if you initialize w_i at a minimum



↳ If α too small

↳ gradient descent can be slow

↳ If α too large

↳ gradient descent can overshoot the minimum

↳ It may fail to converge/diverge

Gradient descent algorithm for Linear regression

Repeat until convergence {

$$w_0 := w_0 - \alpha \sum_n [y_n - (w_0 + w_1 x_n)]$$

$$w_1 := w_1 - \alpha 2 \sum_n [y_n - (w_0 + w_1 x_n)] x_n$$

}

This particular version is called "Batch" gradient descent

<https://youtu.be/LzXakYdtAOo?si=-58VyX5Ko310W8Bd>

2. Least Squared Method

↳ to find optimal

$$1. \text{mean}(x) = \bar{x}$$

$$2. \text{mean}(y) = \bar{y}$$

$$3. \text{mean}(xy) = \bar{xy}$$

$$4. \text{mean}(x^2) = \bar{x^2}$$

$$5. m = \frac{(\bar{x} \times \bar{y}) - \bar{xy}}{(\bar{x})^2 - \bar{x^2}} \rightarrow \text{gradient}$$

$$6. c = \bar{y} - (m \times \bar{x}) \rightarrow y\text{-intercept}$$

$$y = mx + c$$

\swarrow find \searrow

x	y
1	2
2	1
4	3

$$\bar{x} = \frac{4+2+1}{3} =$$

$$\bar{y} = \frac{2+1+3}{3} =$$

$$\bar{xy} = \frac{(4 \times 2) + (2 \times 1) + (3 \times 1)}{3},$$

$$\bar{x^2} = \frac{1^2 + 2^2 + 4^2}{3},$$

$$m = 0.428$$

$$c = 1$$

$$y = \underline{\underline{0.428n + 1}}$$

Are there other methods to find optimal w

► Closed form solution exists using Linear Algebra

$$y = w^T X - b$$

► The method is called Least squares method

↳ Co-relation co-efficient $\rightarrow R$

↳ -1 very correlated \square

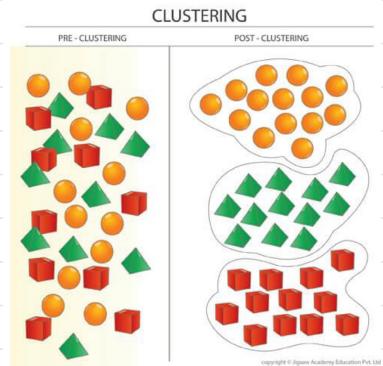
↳ 0 not correlated at all \square

↳ 1 +very correlated $\square \rightarrow$ linear

CLUSTERING

Clustering

- ↳ Unsupervised ML technique
- ↳ Clusters are inferred from the data without human input
- ↳ Groups similar instances
 - ↳ Euclidean distance



Unsupervised

- ↳ It learns on its own how to classify the data by implicitly learned features
- ↳ Requires data but no labels

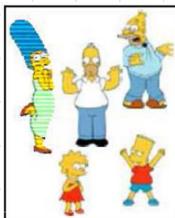
- Detect patterns e.g. in

- Group emails or search results
 - Customer shopping patterns
 - Regions of images
- Useful when don't know what you're looking for
 - But: can get gibberish

CLUSTERING ALGOS

Partition Algos

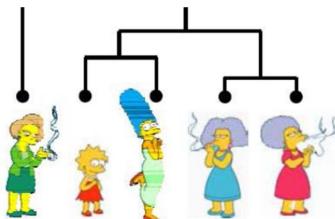
- ↳ K-means
- ↳ Mixture of Gaussian
- ↳ Spectral Clustering



↳ FLAT

Hierarchical Algos

- ↳ Agglomerative → bottoms up
- ↳ Divisive → top down



K-Means / KNN

↳ An iterative clustering algo

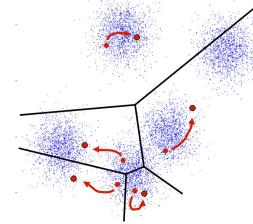
↳ Guarantees convergence in a finite no of iterations

- Initialize: Pick K random points as cluster centers

- Alternate:

1. Assign data points to closest cluster center
2. Change the cluster center to the average of its assigned points

- Stop when no points' assignments change



$$\text{Euclidean Distance} = \sqrt{(x_2 - x_1)^2 + (y_2 - y_1)^2}$$

↓ observed value ↓ centroid value

Steps

taking $K=2$

↳ Choose any 2 points for 2 clusters C_1 , C_2

↳ Find euclidean distance b/w

next point, C_1 and next point, C_2

↳ Which ever cluster gives smallest distance

↳ Get Avg of that cluster and Point

↳ Update cluster

REPEAT

Q)

sr_no	age	amount
C1	20	500
C2	40	1000
C3	30	800
C4	18	300
C5	28	1200
C6	35	1400
C7	45	1800

$$K_1 = (C_1) = (20, 500)$$

$$d(K_1, C_3) = \sqrt{(30-20)^2 + (800-500)^2} = 300.0$$

$$K_2 = (C_2) = (40, 1000)$$

$$d(K_2, C_3) = \sqrt{(30-40)^2 + (800-1000)^2} = 200.$$

smaller 300
update min cluster

$K_1 + C_3$

$$\left(\frac{40+20}{2}, \frac{1000+500}{2} \right) = 35, 900$$

$$K_1 = (C_1) = (20, 500)$$

$$d(K_1, C_4) = \sqrt{(18-20)^2 + (300-500)^2} =$$

↓
Smaller

$K_1 + C_4$

$$\left(\frac{18+20}{2}, \frac{300+500}{2} \right) =$$

!

Repeat till all Points done

Assume that the Central Unit (CU) makes a decision about a user X received signal values from the BS in terms of distance and power as shown in **Table. C** between **11:00 to 3:00 PM**. The CU is intelligent and can classify the signals based on distance and power levels. Use **KNN algorithm** with $k = 3$ to classify the signal 5.

Table.C Mobile signal with respect to power and distance

Signals	Distance	Power	Class (Normal and Abnormal)
1	8	20	Normal
2	4	18	Normal
3	5	16	Abnormal
4	10	15	Normal
5	32	10	??

a)

$$d(1,5) = \sqrt{(32-8)^2 + (10-20)^2} = 26 \rightarrow N \rightarrow \text{NORMAL}$$

$$d(2,5) = \sqrt{(32-4)^2 + (10-18)^2} = 29.19 \rightarrow N$$

$$d(3,5) = \sqrt{(32-5)^2 + (10-16)^2} = 27.65 \rightarrow A$$

$$d(4,5) = \sqrt{(32-10)^2 + (10-15)^2} = 22.56$$

b)

Table. D Patient dataset

Patient Name	Address	Age	Body temperature	Height	Blood pressure	Cough	Respiratory Problem	Weight	Admitted to the hospital
Aslam	DHA	46	102	6	109	Yes	Severe	60	Yes
Umair	Malir	44	97	6.11	98	No	Less severe	55	No
Zain	Johar	36	99	5.9	85	No	Less severe	56	No
Umer	Gulshan	37	100	5.8	83	No	Severe	59	Yes
Ali	Saddar	39	97	5.3	95	No	Less severe	76	No
Ziauddin	Johar	40	98	5.6	120	Yes	Severe	78	Yes
Waqar	DHA	42	99	5.5	180	Yes	Less severe	80	No
Saleem	Bahria town	44	102	5.8	110	Yes	Severe	??	Yes

$$\frac{P_1 + P_2 + P_7}{3} = \frac{60 + 55 + 80}{3} = 65$$

check euclidean distance

A) The outbreak of the coronavirus pandemic has completely affected the educational system worldwide, leading to the near-total closures of academia impacting about 98.5 percent of the world's student population. So, as an ingenuity action HEC allowed universities and degree awarding institutes to continue online education and conduct examinations by making a new policy without compromising the health and safety factors. For which FAST NUACES has decided a policy of SU and Letter grade for final examination. But in an online meeting with classmates you noticed that it was difficult for everyone to choose which policy will be more effective.

So, you are required to make a classifier using K-Means Algorithm to divide the observation into TWO clusters (**C1=S/U Policy**) and (**C2=Letter Grade**) using the following dataset rows according to your roll number digits.

For Example as shown in Figure 1 below, For 19k-2858 choose all relevant rows even if a digit is repeated twice or thrice (choose the 1st, 9th, 2nd, 8th, 5th and 8th row again) and make a new table.

07th July 2020, 9:00 am – 12:30 pm

Digit	A	B
0	3	3
1	8	5
2	4	4
3	2	4
4	7	7
5	5	8
6	3	5
7	4	8
8	6	9
9	9	6

Digit	A	B
1	8	5
9	9	6
2	4	4
8	6	9
5	5	8
8	6	9

Figure 1: Example data selection for Roll No. 19k-2858

173730

→

Digit	A	B
1	8	5
7	4	8
3	2	4
7	4	8
3	2	4
0	3	3

Taking 3 and 7 as random seeds

$$C_1 = (3)(2,4)$$

$$1(8,5) (C_1, 1) = \sqrt{(4-2)^2 + (5-4)^2} = 6.08$$

$$C_2 = (7)(4,8)$$

$$(C_2, 1) = \sqrt{(8-4)^2 + (5-8)^2} = 5$$

$$2,4$$

$$7(4,8) = \sqrt{(4-2)^2 + (8-4)^2} = 4.47$$

$$\frac{8+4}{2} = \frac{12}{2}$$

$$C_2(7,1) = (6,6.5)$$

$$\sqrt{(4-6)^2 + (8-6.5)^2} = 2.5$$

$$C_1(3) = (2,4)$$

$$C_2(7,1,7) = (5,7.25)$$

$$3(2,9) = \sqrt{(2-2)^2 + (4-0)^2} = 0$$

$$C_1(3,3) = (2,4)$$

$$C_2(7,1,7) = (5,7.25)$$

$$0(3,3) = \sqrt{(3-2)^2 + (3-4)^2} = 1.414$$

$$\sqrt{(3-5)^2 + (3-7.25)^2} = 4.64$$

$$C_1(3,3,0) = (2.5, 3.5)$$

$$C_2(7,1,7) = (5,7.25)$$

Figure 1. Example data selection for k-means

- B) Consider a fictitious document collection. There are 6 documents in this collection and these are represented as 6 points in two-dimensional vector space as follow:

D1	(2,1)
D2	(1,1)
D3	(4,1)
D4	(1,2)
D5	(2,2)
D6	(4,2)

Suppose that the distance between a pair of documents is measured by the Euclidean distance between their corresponding points. Show that how the k-means algorithm (with k=2) clusters these documents, using D2 and D3 as seeds.

$c_1(D_2)$ (1,1)	$c_2(D_3)$ (4,1)
$D_1(2,1)$ $\sqrt{(2-1)^2 + (1-1)^2} = 1 \checkmark$	$\sqrt{(2-4)^2 + (1-1)^2} = 2$
$c_1(D_2, D_1) = (1.5, 1)$	
$D_4(1,2)$ $\sqrt{(1-1.5)^2 + (2-1)^2} = 1.1 \checkmark$	$\sqrt{(1-4)^2 + (2-1)^2} = 3.16$
$c_1(D_2, D_1, D_4) = (1.25, 1.5)$	
$D_5(2,2)$ $\sqrt{(2-1.25)^2 + (2-1.5)^2} = 0.901 \checkmark$	$\sqrt{(2-4)^2 + (2-1)^2} = 2.23$
$c_1(D_2, D_1, D_4, D_5) = (1.625, 1.75)$	
$D_6(4,2)$ $\sqrt{(4-1.625)^2 + (2-1.75)^2} = 2.38$	$\sqrt{(4-4)^2 + (2-1)^2} = 1 \checkmark$
$c_1(D_2, D_1, D_4, D_5)$	$c_2(D_3, D_6)$

NEURAL NETWORKS

$$y_{in} = \sum_{i=1}^n w_i x_i + b \xrightarrow{\text{biasness}}$$

$$y_{in} = w_1 x_1 + w_2 x_2 + \dots + w_n x_n + b$$

Activation Functions

↳ Sigmoidal $f(y) = \frac{1}{1+e^{-y}}$

↳ RELU $= \max(y, 0)$

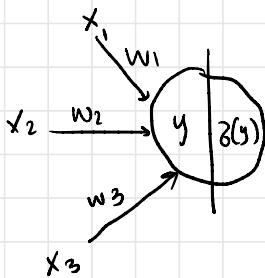
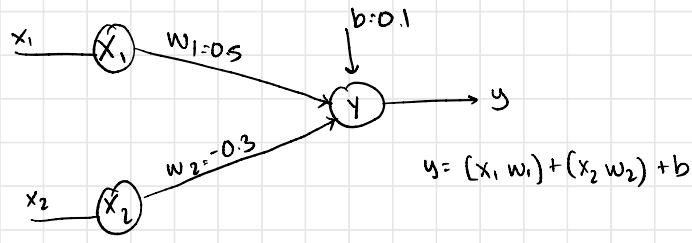
↳ Step function $y = \begin{cases} 1 & \text{if } n > 0 \\ 0 & \text{if } n \leq 0 \end{cases}$

so ours b/w
(0,1)

$$\begin{aligned} w_1 &= 0.5 \\ w_2 &= -0.3 \\ b_1 &= 0.1 \\ &\downarrow \\ &\text{biasness} \end{aligned}$$

Hows
Study

Hows
sleep



$$\begin{aligned} y &= w_1 x_1 + w_2 x_2 + w_3 x_3 + b \\ f(y) &= \frac{1}{1+e^{-y}} \end{aligned}$$

Weight Update

$$w_{new} = w_{old} + L(t-o)x_i$$

↑ learning rate ↑ target ↑ observed

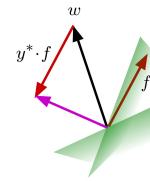
$L=0.5$ Threshold :=

Problem: Why can't a perceptron represent XOR?

Solution: Intuitively, a perceptron is a linear classifier (given that we use the step function as the activation function), but XOR is not linearly separable. In the drawing below, you can see that AND and OR are linearly separable, but there is no way to linearly separate the shaded outputs of XOR.

Binary Perceptron

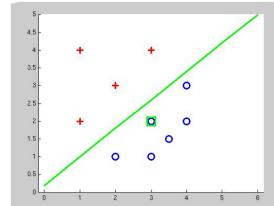
- ↳ inputs → feature values
- ↳ each feature has a weight
- ↳ sum is the activation
- ↳ if activation is
 - ↳ +ve → output +1
 - ↳ -ve → output -1



▪ Separable Case

- ↳ Start with $w=0$
- ↳ Classify with current weights

$$y = \begin{cases} +1 & \text{if } w \cdot f(x) \geq 0 \\ -1 & \text{if } w \cdot f(x) \leq 0 \end{cases}$$



- ↳ e.g. $y=y^*$
↳ if correct → no change
- ↳ if wrong → adjust the weight vector
by adding/subtracting the feature vector

$$w = w + y^* \cdot f$$

Multiclass Perceptron

↳ If we have multiple classes

↳ A weight vector for each class w_y

↳ Score activation of class y $w_y \cdot f(x)$

↳ Prediction highest score wins $y = \arg \max_y w_y \cdot f(x)$

↳ Start with all $w=0$

↳ Pick up training examples 1 by 1

↳ Predict with current weights

$$y = \arg \max_y w_y \cdot f(x)$$

↳ If correct \rightarrow no change

↳ If wrong \rightarrow lower score of wrong ans

$$w_y = w_y - f(n)$$

raise score of right ans

$$w_y = w_y + f(n)$$

"win the vote"

"win the election"

"win the game"

w_{SPORTS}

BIAS : 1
win : 0
game : 0
vote : 0
the : 0
...

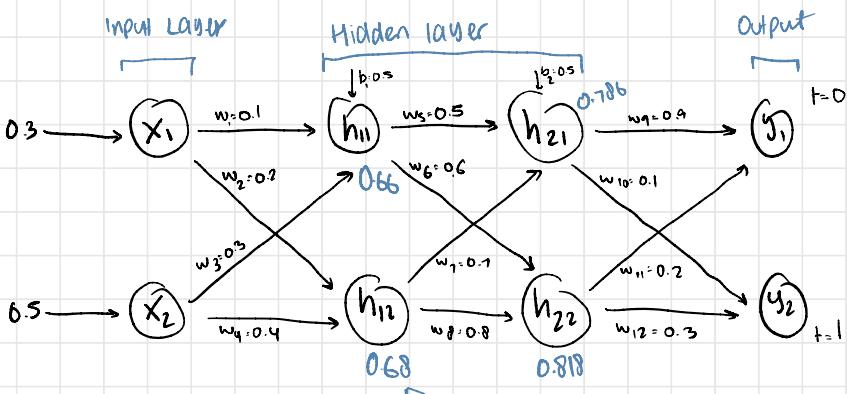
$w_{POLITICS}$

BIAS : 0
win : 0
game : 0
vote : 0
the : 0
...

w_{TECH}

BIAS : 0
win : 0
game : 0
vote : 0
the : 0
...

Feed Forward Neural Network



$$h_{11} = w_1 x_1 + w_3 x_2 + b_1$$

$$(0.1 \times 0.3) + (0.3 \times 0.5) + 0.5 = 0.68$$

$$\delta(h_{11}) = \frac{1}{1 + e^{-0.68}} = 0.66$$

$$\begin{aligned} h_{21} &= \delta(h_{11}) w_5 + \delta(h_{12}) w_7 + b_2 \\ &= (0.66 \times 0.5) + (0.68 \times 0.7) + 0.5 \\ &= 1.306 \end{aligned}$$

$$\delta(h_{21}) = \frac{1}{1 + e^{-1.306}} = 0.786$$

$$\begin{aligned} y_1 &= w_9 \delta(h_{21}) + w_{11} \delta(h_{22}) \\ &= 0.1636 \end{aligned}$$

$$\delta(y_1) = 0.54$$

$$h_{12} = x_2 w_4 + x_1 w_2 + b_1$$

$$(0.5 \times 0.4) + (0.3 \times 0.2) + 0.5 = 0.76$$

$$\delta(h_{12}) = \frac{1}{1 + e^{-0.76}} = 0.68$$

$$\begin{aligned} h_{22} &= \delta(h_{11}) w_6 + \delta(h_{12}) w_8 + b_2 \\ &= (0.66 \times 0.6) + (0.68 \times 0.8) + 0.5 \\ &= 1.504 \end{aligned}$$

$$\delta(h_{22}) = \frac{1}{1 + e^{-1.504}} = 0.818$$

$$\begin{aligned} y_2 &= w_{10} \delta(h_{21}) + w_{12} \delta(h_{22}) \\ &= 0.324 \end{aligned}$$

$$\delta(y_2) = 0.58$$

$$\begin{aligned} \text{Error} &= \text{mean squared error} = \frac{1}{2} \left[(\underline{y_0 - y_T})^2 + (\underline{y_0 - y_T})^2 \right] \\ &= \frac{1}{2} \left[(0.54 - 0)^2 + (0.58 - 1)^2 \right] = 0.234 \end{aligned}$$

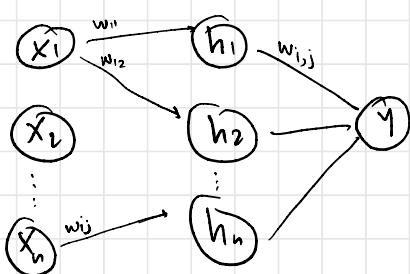
Get corrected weight for net iteration

NOW back PROPOgation

↳ change weights and bias

NOT DONE

- A) Suppose you had a neural network with linear activation functions. That is, for each unit the output is some constant c times the weighted sum of the inputs. Assume that the network has one hidden layer. For a given assignment to the weights w , write down equations for the value of the units in the output layer as a function of w and the input layer x , without any explicit mention of the output of the hidden layer.



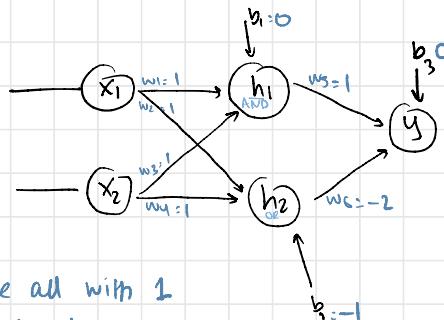
$$H = c \left(\sum_k^n w_{ij} x_k \right) + b_{in}$$

$$O = (\sum w_{ij} \times H_i) = c \left[\sum_i w_{ij} \times c \left(\sum_k^n w_{ik} x_k \right) + b_m \right] + b_o$$

B). The spread spectrum (SS) has been a very popular technique to encode messages in military communication. In SS, a message signal is combined with the secret chip sequence using **XOR** operation. The resultant message is unable to decode by the unknown receiver and transmitted through wireless channel. Draw a simple Neural Network Architecture which can achieve the **XOR** operation of the following given data shown in Table. E.

Table.E Input and output values

Message signal: x_1	Secret chip sequence: x_2	Resultant message: y
0	0	0
1	1	0
1	0	1
0	1	1



Initialize all with 1

Do trial and error

$$h_1 = x_1 w_1 + x_2 w_3 + b_1 = 1$$

1 1 0 1 0 0

$$h_2 = x_1 w_2 + x_2 w_4 + b_2 = 0$$

1 1 0 1 1 -1

$$y = h_1 w_5 + h_2 w_6 + b_3 = 0$$

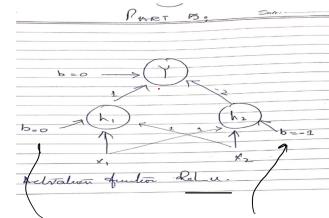
1 -2 0 1 0 0

Problem: Why can't a perceptron represent XOR?

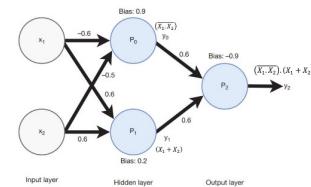
Solution: Intuitively, a perceptron is a linear classifier (given that we use the step function as the activation function), but XOR is not linearly separable. In the drawing below, you can see that AND and OR are linearly separable, but there is no way to linearly separate the shaded outputs of XOR.



Activation function
 $\text{ReLU} = \max(y, 0)$



Combining Multiple Perceptrons Cont'd



Weight Update

$$w_{\text{new}} = w_{\text{old}} + L(t - o) x_i$$

learning rate L target t desired o

$L=0.5$ Threshold: 1

AND Gate - PERCEPTRON TRAINING RULE

w1 = 1.2, w2 = 0.6 Threshold = 1 and Learning Rate n = 0.5

3. A=1, B=0 and Target = 0

- $wi \cdot xi = 1 * 1.2 + 0 * 0.6 = 1.2$
- This is greater than the threshold of 1, so the output = 1

A	B	$A \wedge B$
0	0	0
0	1	0
1	0	0
1	1	1

$$wi = wi + n(t - o)xi$$

$$w1 = 1.2 + 0.5(0 - 1)1 = 0.7$$

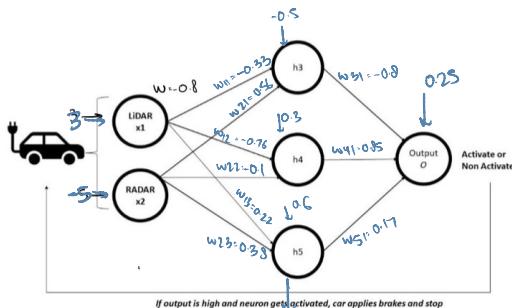
$$w2 = 0.6 + 0.5(0 - 1)0 = 0.6$$

[Subscribe to Mahesh Huddar](#)

Visit: vtupulse.com

of output neuron O based on the given data as follows using only feedforward propagation and for 1 epoch. Use Sigmoid function = $f(x) = \frac{1}{1-e^{-x}}$ for all the activation functions.

LIDAR x1 = 3	RADAR x2 = -5	w11 = -0.33	w12 = -0.76	w13 = 0.22	w21 = 0.56	w22 = -0.1	w23 = 0.38
W31 = -0.8	W41 = 0.85	w51 = 0.17	b3 = -0.5	b4 = 0.3	b5 = 0.6	b6 = 0.25	b7 = 0.1



hidden layer

$$h_3 = (3 \times -0.33) + (-5 \times 0.56) + (-0.5) \\ = -4.29$$

$$\delta(h_3) = -0.0139$$

$$\delta(h_4) = -0.295$$

$$\delta(h_5) = -1.155$$

$$\delta(h_3) = \frac{1}{1 - e^{(-4.29)}} = -0.0139$$

$$y = (-0.0139 \times -0.8) + (-0.295 \times 0.85) + (-1.155 \times 0.7) + 0.15$$

$$= 0.32232$$

$$\delta(y) = 3.629$$

$$O = \underline{3.629}$$

$$h_3 = (3 \times 0.22) + (-5 \times 0.38) + 0.6 \\ = -0.64$$

$$\delta(h_3) = -1.155$$

Properties of Perceptron

↳ Separability

- ↳ true if some parameter get the training set perfectly correct

↳ Convergence

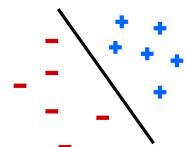
- ↳ If the training set is separable
Perception will eventually converge → binary case

↳ Mistake bound

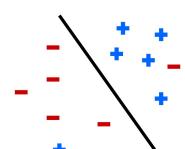
- ↳ the max no. of mistakes related to the → binary case
margin/degree of separability

$$\text{mistakes} < \frac{K}{\delta^2}$$

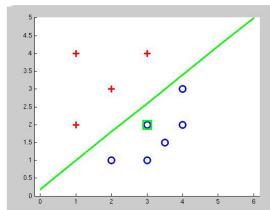
Separable



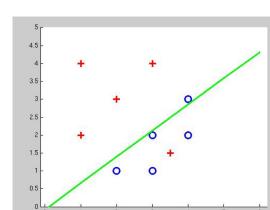
Non-Separable



▪ Separable Case



▪ Non-Separable Case

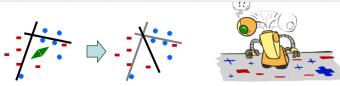


Problems with Perceptrons

1. Noise

↳ If the data isn't separable, weights might thrash

SOL → Avg weight vectors over time can help average perceptron



2. Mediocre generalization

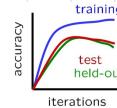
↳ finds a 'barely' separating solution



3. Overtraining

↳ test accuracy usually rises, then falls

↳ Overtraining is a kind of overfitting



Fixing me Perceptrons

- Idea: adjust the weight update to mitigate these effects

- MIRA*: choose an update size that fixes the current mistake...

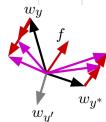
- ... but, minimizes the change to w

$$\min_w \frac{1}{2} \sum_y ||w_y - w'_y||^2$$

$$w_{y^*} \cdot f(x) \geq w_y \cdot f(x) + 1$$

- The +1 helps to generalize

- *Margin Infused Relaxed Algorithm



Guessed y instead of y^* on example x with features $f(x)$

$$w_y = w'_y - \tau f(x)$$

$$w_{y^*} = w'_y + \tau f(x)$$

Minimum Correcting Update

$$\min_w \frac{1}{2} \sum_y ||w_y - w'_y||^2$$

$$w_{y^*} \cdot f \geq w_y \cdot f + 1$$

$$\min_{\tau} ||\tau f||^2$$

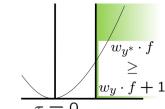
$$w_{y^*} \cdot f \geq w_y \cdot f + 1$$

$$(w'_y + \tau f) \cdot f = (w'_y - \tau f) \cdot f + 1$$

$$\tau = \frac{(w'_y - w'_{y^*}) \cdot f + 1}{2f \cdot f}$$

$$w_y = w'_y - \tau f(x)$$

$$w_{y^*} = w'_y + \tau f(x)$$



min not $\tau=0$, or would not have made an error, so min will be where equality holds

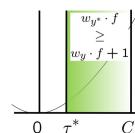
Maximum Step Size

- In practice, it's also bad to make updates that are too large

- Example may be labeled incorrectly
- You may not have enough features
- Solution: cap the maximum possible value of τ with some constant C

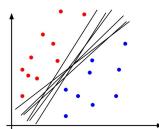
$$\tau^* = \min \left(\frac{(w'_y - w'_{y^*}) \cdot f + 1}{2f \cdot f}, C \right)$$

- Corresponds to an optimization that assumes non-separable data
- Usually converges faster than perceptron
- Usually better, especially on noisy data



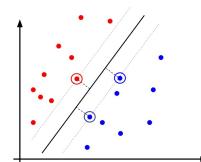
Linear Separators

- Which of these linear separators is optimal?



Support Vector Machines

- Maximizing the margin: good according to intuition, theory, practice
- Only support vectors matter; other training examples are ignorable
- Support vector machines (SVMs) find the separator with max margin
- Basically, SVMs are MIRA where you optimize over all examples at once



$$\begin{aligned} \min_w & \frac{1}{2} ||w - w'||^2 \\ w_{y^*} \cdot f(x_i) & \geq w_y \cdot f(x_i) + 1 \end{aligned}$$

$$\begin{aligned} \text{SVM} \quad \min_w & \frac{1}{2} ||w||^2 \\ \forall i, y_i w_{y^*} \cdot f(x_i) & \geq w_y \cdot f(x_i) + 1 \end{aligned}$$

Classification: Comparison

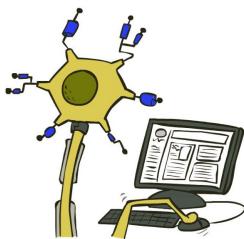
Naïve Bayes

- Builds a model training data
- Gives prediction probabilities
- Strong assumptions about feature independence
- One pass through data (counting)

Perceptrons / MIRA:

- Makes less assumptions about data
- Mistake-driven learning
- Multiple passes through data (prediction)
- Often more accurate

Web Search



Extension: Web Search

- Information retrieval:
 - Given information needs, produce information
 - Includes, e.g. web search, question answering, and classic IR
- Web search: not exactly classification, but rather ranking



Feature-Based Ranking

$x = \text{"Apple Computer"}$

$$f(x, \cdot) = [0.3 \ 5 \ 0 \ 0 \ \dots]$$

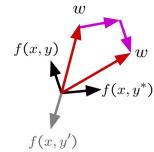
Screenshot of a search results page for "Apple Computer" showing various links.

$$f(x, \cdot) = [0.8 \ 4 \ 2 \ 1 \ \dots]$$

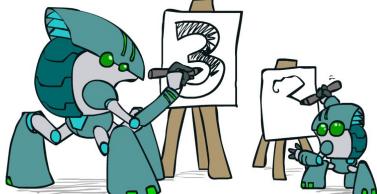
Screenshot of a search results page for "Apple Computer" showing various links.

Perceptron for Ranking

- Inputs x
- Candidates y
- Many feature vectors: $f(x, y)$
- One weight vector: w
 - Prediction: $y = \arg \max_y w \cdot f(x, y)$
 - Update (if wrong): $w = w + f(x, y^*) - f(x, y)$



Apprenticeship



ishma hafeez
notes
repost free!

Pacman Apprenticeship!

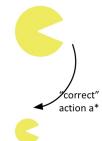
- Examples are states s



- Candidates are pairs (s, a)
- "Correct" actions: those taken by expert
- Features defined over (s, a) pairs: $f(s, a)$
- Score of a q-state (s, a) given by:

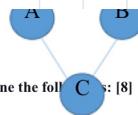
$$w \cdot f(s, a)$$

- How is this VERY different from reinforcement learning?



$$\forall a \neq a^*, w \cdot f(a^*) > w \cdot f(a)$$

[Demo: Pacman Apprentice (L22)]



A) Define the foll : [8]

1- Define AI agent.

An agent perceives its environment through sensors and acts upon it through actuators (or effectors, depending on whom you ask). The Agent takes its action(s), which changes the environment, and process repeats.

2- What is back propagation and why do we use it?

Algorithm for training neural networks with hidden layers. Calculate error for output layer and propagate error back layer wise and update weights.

3- How supervised and reinforcement learning are related?

Both Supervised and RL get feedback:
Supervise Learning get feedback in form of Labeled Class
RL get feedback in form Reward.

4- Define the role of activation function and give the name of popular activation function.

Activation function is non-linear function which is used activate/ trigger particular neuron. Sigmoid, tanh, RELU are mostly commonly used activation functions.

5- Differentiate between Markov decision process and reinforcement learning.

In MDP transition probability and rewards are known. In RL transitions and rewards are unknown and it is estimated through interacting with Environment

6- Given the figure below, determine the followings:

- a) Whether A and B are conditionally independent given C is observable (NO)
No. If C is observed than probability A can effect the probability of B so they are not independent.
- b) Whether A and B are conditionally independent given C is unobservable (YES)
Yes. They are conditional independent if C is unobserved. Probability of A does not affect probability of B.

not independent

independent

7- What is the concept of gradient descent in neural network?

Gradient descent is used to minimize the loss function of neural network. We choose any starting point in weight space and then move to a neighboring point that is downhill, repeating until we converge on the minimum possible.

Reasoning Under Uncertainty in AI

↳ It refers to situations where there is a lack of complete info/knowledge about a particular aspect leading to ambiguity and unpredictability



Reasoning strategy to handle uncertainty

1. Probabilistic Models
2. Bayesian Networks
3. Monte Carlo Methods
4. Decision Theory
5. Fuzzy Logic
6. Qualitative Reasoning

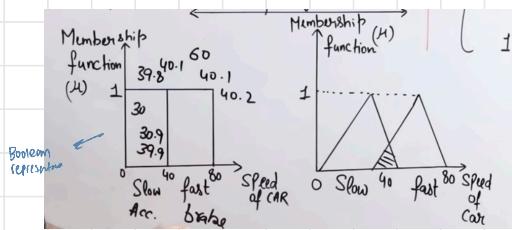
FUZZY LOGIC

↳ represents uncertainty

↳ represent with degree

↳ represent the belongings of a member $\rightarrow \mu$

of a crisp set to fuzzy set



If we use

if speed becomes 40

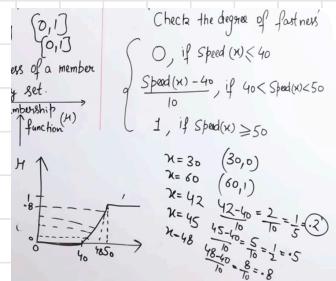
it will break

then with speed = 30

it will accelerate

so it will keep oscillating b/w the two

no flexibility



Operations in fuzzy logic

→ Union: $\max\{\mu_A(x), \mu_B(x)\}, x \in U$

→ Intersection: $\min\{\mu_A(x), \mu_B(x)\}, x \in U$

→ Complement: $\mu_{\bar{A}}(x) = [1 - \mu_A(x)], x \in U$

→ Bold Union: $\mu_{A+B} = \min[1, \mu_A(x) + \mu_B(x)]$

→ Bold Intersection: $\mu_{A \cdot B} = \max[0, \mu_A(x) + \mu_B(x) - 1]$

→ Equality $A=B$ if $\mu_A(x) = \mu_B(x) \forall x \in S$

$$U = \{5, 10, 20, 25, 30, 40\}$$

$$A = \{(10, 0.2), (20, 0.4), (25, 0.7), (30, 0.9), (40, 1)\}$$

$$B = \{(10, 0.4), (20, 0.1), (25, 0.9), (30, 0.2), (40, 0.6)\}$$

Probabilistic Inference

Samples:

- $s_1: (\neg a, \neg b, c, d);$
- $s_2: (a, \neg b, \neg c, d);$
- $s_3: (\neg a, \neg b, c, \neg d);$
- $s_4: (\neg a, \neg b, \neg c, \neg d);$
- $\cancel{s_5: (\neg a, b, c, d);}$
- $s_6: (a, \neg b, \neg c, \neg d);$
- $\cancel{s_7: (\neg a, b, \neg c, d);}$
- $\cancel{s_8: (a, b, \neg c, \neg d);}$
- $s_9: (\neg a, \neg b, \neg c, d);$
- $s_{10}: (a, \neg b, c, \neg d);$

Queries:

$$\begin{aligned} P(\neg c) &\approx \frac{6}{10} = \frac{\text{fav.}}{\text{Total}} = \frac{1}{2} \\ P(\neg a \wedge b) &\approx \frac{2}{10} \\ P(\neg b \mid \neg d) &\approx \frac{4}{5} = P(A \mid B) \\ P(a \mid b, \neg c) &= \frac{1}{2} \quad \text{find } \downarrow \quad \text{given } \downarrow \\ P(A \mid B) &= \frac{\text{case } A, B}{\text{case } B} \quad \text{Q: } \downarrow \quad \text{Evid.} \\ Q: \downarrow & \quad \text{Evid.} \end{aligned}$$

Samples:

- $s_1: (\neg a, \neg b, c, d); \quad 0.4$
- $s_2: (a, \neg b, \neg c, d); \quad 0.2$
- $s_3: (\neg a, \cancel{\neg b}, c, \cancel{\neg d}); \quad 0.1$
- $s_4: (\neg a, \cancel{\neg b}, \neg c, \neg d); \quad 0.5$
- $s_5: (\neg a, b, c, d); \quad 0.3$
- $s_6: (a, \cancel{\neg b}, \neg c, \neg d); \quad 0.4$
- $s_7: (\neg a, b, \neg c, d); \quad 0.6$
- $s_8: (\neg a, b, \neg c, \neg d); \quad 0.2$
- $s_9: (a, \neg b, \neg c, d); \quad 0.5$
- $s_{10}: (a, \cancel{\neg b}, c, \neg d); \quad 0.1$

Queries:

$$\begin{aligned} P(\neg c) &= \frac{6}{10} = \frac{24}{33} \\ P(\neg a, b) & \\ P(\neg b \mid \neg d) &= \frac{1.1}{1.3} \\ P(a \mid b, \neg c) & \end{aligned}$$



Fahad Sherwani

20 May

Dear Students,

The pattern for final exam will be around 8

- 10 questions. 80 • 85% content to be included is studied after MID 2 while 15 • 20% before MID 2.

The Syllabus includes:

- 1) Content before MID 2.
 - 2) Bayesian Networks
 - 3) Naive Bayes
 - 4) Decision Trees
 - 5) Linear Regression
 - 6) Neural Networks
 - 7) K-means/KNN
- ~~✗~~ Performance matrices (Confusion matrix, I will share a tutorial about it tomorrow)
~~✗~~ Some part of Reinforcement learning/Q-Learning (Not confirmed yet, but I will share a tutorial about it too).

Thank you.

Conditional Probability

Conditional probability is the probability of an event A occurring given that event B has occurred. It is denoted as $P(A|B)$, where P represents probability. The

$$P(A|B) = P(A \cap B) / P(B)$$

Probability Theorems

bn

HMM

} Lab 7

} Lab 8

Lab 9, 10

Supervised

- ↳ Regression → output real value
 - ↳ Linear
 - ↳ Polynomial
 - ↳ SVM
 - ↳ Decision Tree
 - ↳ Random Forest

- ↳ Classification → output categories
 - ↳ Logistic
 - ↳ SVM
 - ↳ Decision Tree
 - ↳ Random Forest
 - ↳ Naive Bayes

Unsupervised

- ↳ Clustering → group similar data together
 - ↳ Exclusive Partitioning
 - ↳ Agglomerative
 - ↳ Overlapping
 - ↳ Probabilistic

- ↳ Hierarchical

- ↳ K-means
- ↳ Principal Component Analysis
- ↳ Gaussian Mixture Models
- ↳ DBSCAN

- ↳ Association → identify patterns
 - ↳ Apriori
 - ↳ FP Growth

Reinforcement

- ↳ Q Learning, MDP

SUPERVISED

Regression & Classification

Code Implementations

Linear Regression
Definition- Linear regression is a statistical method used to model the relationship between a dependent variable (target) and one or more independent variables (features) by fitting a linear equation to the observed data.

Use Cases

- Predicting house prices based on features like area, number of bedrooms, etc.
- Forecasting sales revenue based on advertising spend.
- Estimating the demand for a product based on historical sales data.

```
from sklearn.linear_model import LinearRegression
import numpy as np

# Sample data (house areas in sq. ft. and prices in USD)
X = np.array([[1500], [1800], [2200], [2500]]) # Input features
y = np.array([300000, 350000, 400000, 450000]) # Target variable (house prices)

# Create and fit the linear regression model
model = LinearRegression()
model.fit(X, y)

# Predict price for a new house with area 2000 sq. ft.
new_area = np.array([[2000]])
predicted_price = model.predict(new_area)
print("Predicted Price: ", predicted_price) # Output: [383333.333333]
```

Polynomial Regression

Definition- Polynomial regression extends linear regression by adding polynomial terms to the model, allowing it to capture non-linear relationships between features and target

Use Cases

- Modelling physical phenomena where the relationship between variables is not linear.
- Predicting the birth rate of organisms based on time and environmental factors.
- Capturing the effect of interactions between variables in the model.

```
from sklearn.preprocessing import PolynomialFeatures
from sklearn.linear_model import LinearRegression
import numpy as np

# Sample data (input features and target variable)
X = np.array([[1], [2], [3], [4]]) # Quadratic relationship (y = x^2)
y = np.array([1, 4, 9, 16]) # Quadratic relationship (y = x^2)

# Transform features to polynomial features
poly_features = PolynomialFeatures(degree=2)
X_poly = poly.fit_transform(X)

# Create and fit the polynomial regression model
model = LinearRegression()
model.fit(X_poly, y)

# Predict for new values
X_new = np.array([5])
X_new_poly = np.poly1d(poly.fit_transform(X_new))
predicted_value = model.predict(X_new_poly)
print("Predicted Value (Polynomial Regression): ", predicted_value)
# Output: [25.]
```

Support Vector Machine (SVM) Regression

Definition- Support Vector Machine (SVM) regression finds a hyperplane that best fits the data by maximizing the margin between support vectors and the hyperplane.

Use Cases

- Predicting stock prices based on historical data.
- Forecasting future demand for products or services.
- Estimating the price of real estate based on property features.

```
from sklearn.svm import SVR
import numpy as np

# Sample data (input features and target variable)
X = np.array([[1], [2], [3], [4]]) # Linear relationship (y = x + 1)
y = np.array([2, 3, 4, 5]) # Linear relationship (y = x + 1)

# Create and fit the SVM regression model
svm_regressor = SVR(kernel='linear')
svm_regressor.fit(X, y)

# Predict for new values
X_new = np.array([10])
predicted_value = svm_regressor.predict(X_new)
print("Predicted Value (SVM Regression): ", predicted_value)
# Output: [6.]
```

Decision Tree Regression

Definition- Decision tree regression divides the feature space into regions and predicts the average target value within each region.

Use Cases

- Predicting the price of a used car based on its age, mileage, and other features.
- Estimating the energy consumption of a building based on various parameters.
- Forecasting the sales volume of a product based on marketing campaigns.

```
from sklearn.tree import DecisionTreeRegressor
import numpy as np

# Sample data (input features and target variable)
X = np.array([[1], [2], [3], [4]]) # Linear relationship (y = x + 1)
y = np.array([2, 3, 4, 5]) # Linear relationship (y = x + 1)

# Create and fit the decision tree regression model
tree_regressor = DecisionTreeRegressor()
tree_regressor.fit(X, y)

# Predict for new values
X_new = np.array([[5]])
predicted_value = tree_regressor.predict(X_new)
print("Predicted Value (Decision Tree Regression): ", predicted_value)
# Output: [6.]
```

Random Forest Regression

Definition- Random forest regression is an ensemble learning method that combines multiple decision trees and averages their predictions to improve accuracy and reduce overfitting.

Use Cases

- Predicting the yield of a crop based on weather conditions and properties.
- Estimating the price of a product based on customer demographic and market data.
- Forecasting the sales revenue of a retail store based on historical data.

```
from sklearn.ensemble import RandomForestRegressor
import numpy as np

# Sample data (input features and target variable)
X = np.array([[1], [2], [3], [4]]) # Linear relationship (y = x + 1)
y = np.array([2, 3, 4, 5]) # Linear relationship (y = x + 1)

# Create and fit the random forest regression model
forest_regressor = RandomForestRegressor()
forest_regressor.fit(X, y)

# Predict for new values
X_new = np.array([[5]])
predicted_value = forest_regressor.predict(X_new)
print("Predicted Value (Random Forest Regression): ", predicted_value)
# Output: [6.]
```

Code Implementations

Logistic Regression
Definition- Logistic regression is a statistical method used for binary or multiclass classification by modeling the probability of a categorical outcome using a logistic (sigmoid) function.

Use Cases

- Predicting whether an email is spam or not.
- Classifying customer churn (yes/no).
- Identifying the type of flower based on petal and sepal measurements.

```
from sklearn.linear_model import LogisticRegression
from sklearn.datasets import load_iris
from sklearn.model_selection import train_test_split
from sklearn.metrics import accuracy_score

# Load the Iris dataset
iris = load_iris()
X = iris.data
y = iris.target

# Split data into train and test sets
X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2,
random_state=42)

# Create and fit the logistic regression model
model = LogisticRegression()
model.fit(X_train, y_train)

# Predict on the test set
y_pred = model.predict(X_test)

# Evaluate accuracy
accuracy = accuracy_score(y_test, y_pred)
print("Accuracy (Logistic Regression): ", accuracy)
```

Support Vector Machines (SVM)
Definition- Support Vector Machines (SVM) are supervised learning models used for classification and regression tasks by finding the optimal hyperplane that best separates classes in a high-dimensional space.

Use Cases

- Text classification (e.g., sentiment analysis).
- Image classification.
- Anomaly detection in network traffic.

```
from sklearn.svm import SVC
from sklearn.datasets import load_iris
from sklearn.model_selection import train_test_split
from sklearn.metrics import accuracy_score

# Load the Iris dataset
iris = load_iris()
X = iris.data
y = iris.target

# Split data into train and test sets
X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2,
random_state=42)

# Create and fit the SVM classifier
model = SVC()
model.fit(X_train, y_train)

# Predict on the test set
y_pred = model.predict(X_test)

# Evaluate accuracy
accuracy = accuracy_score(y_test, y_pred)
print("Accuracy (SVM): ", accuracy)
```

Decision Trees

Definition- Decision trees are hierarchical models that partition the feature space into regions and assign class labels to each region based on majority voting.

Use Cases

- Customer segmentation based on purchasing behavior.
- Medical diagnosis based on patient symptoms.
- Credit scoring for loan approval.

```
from sklearn.tree import DecisionTreeClassifier
from sklearn.datasets import load_iris
from sklearn.model_selection import train_test_split
from sklearn.metrics import accuracy_score

# Load the Iris dataset
iris = load_iris()
X = iris.data
y = iris.target

# Split data into train and test sets
X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2,
random_state=42)

# Create and fit the decision tree classifier
model = DecisionTreeClassifier()
model.fit(X_train, y_train)

# Predict on the test set
y_pred = model.predict(X_test)

# Evaluate accuracy
accuracy = accuracy_score(y_test, y_pred)
print("Accuracy (Decision Tree): ", accuracy)
```

Random Forests

Definition- Random forests are ensemble learning methods that build multiple decision trees and aggregate their predictions to improve accuracy and reduce overfitting.

Use Cases

- Predicting the species of plants based on botanical features.
- Detecting fraudulent transactions in financial transactions.
- Recommender systems based on user preferences.

```
from sklearn.ensemble import RandomForestClassifier
from sklearn.datasets import load_iris
from sklearn.model_selection import train_test_split
from sklearn.metrics import accuracy_score

# Load the Iris dataset
iris = load_iris()
X = iris.data
y = iris.target

# Split data into train and test sets
X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2,
random_state=42)

# Create and fit the random forest classifier
model = RandomForestClassifier()
model.fit(X_train, y_train)

# Predict on the test set
y_pred = model.predict(X_test)

# Evaluate accuracy
accuracy = accuracy_score(y_test, y_pred)
print("Accuracy (Random Forest): ", accuracy)
```

Naive Bayes

Definition- Naive Bayes classifiers apply Bayes' theorem with strong (naive) independence assumptions between features.

Use Cases

- Text classification tasks such as spam detection and sentiment analysis.
- Document categorization based on content.
- Medical diagnosis based on symptoms and test results.

```
from sklearn.naive_bayes import GaussianNB
from sklearn.datasets import load_iris
from sklearn.model_selection import train_test_split
from sklearn.metrics import accuracy_score

# Load the Iris dataset
iris = load_iris()
X = iris.data
y = iris.target

# Split data into train and test sets
X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2,
random_state=42)

# Create and fit the Naive Bayes classifier
model = GaussianNB()
model.fit(X_train, y_train)

# Predict on the test set
y_pred = model.predict(X_test)

# Evaluate accuracy
accuracy = accuracy_score(y_test, y_pred)
print("Accuracy (Naive Bayes): ", accuracy)
```

VNSUP ERG SUD

CLUSTERING

Hierarchical Clustering

Definition: Hierarchical clustering is a method of cluster analysis that builds a hierarchy of clusters by either recursively merging or splitting clusters based on distance criteria.

Use Cases:

- Groupping similar documents or articles in text mining.
- Classifying taxonomies or dendograms in biological sciences.
- Customer segmentation based on purchasing behavior.

K-means Clustering

Definition: K-means clustering partitions data into k clusters by iteratively assigning data points to the nearest cluster centroid and updating centroids based on the points assigned to each cluster.

Use Cases:

- Market segmentation based on customer demographics.
- Image segmentation for object recognition.
- Analyzing clickstream data for user behavior pattern.

Principal Component Analysis (PCA)

Definition: Principal Component Analysis (PCA) is a dimensionality reduction technique that projects data onto a lower-dimensional space while preserving the maximum variance.

Use Cases:

- Feature extraction and data visualization.
- Noise reduction in signal processing.
- Enhancing the performance of machine learning models by reducing multicollinearity.

Gaussian Mixture Model (GMM)

Definition: Gaussian Mixture Models (GMMs) assume that the data is generated from a mixture of several Gaussian distributions and estimate the parameters of those distributions to perform probabilistic clustering.

Use Cases:

- Image segmentation in computer vision.
- Density estimation for anomaly detection.
- Modeling complex data distributions in unsupervised learning.

Density-Based Spatial Clustering of Applications with Noise (DBSCAN)

Definition: DBSCAN is a density-based clustering algorithm that groups together closely packed points based on a density threshold.

Use Cases:

- Identifying spatial hotspots in geographic data.
- Anomaly detection in cybersecurity.
- Clustering large datasets with irregular shapes and noise.

ASSOCIATE LEARNING

FP-Growth Algorithm

Apriori Algorithm

Definition: The Apriori algorithm is an association rule mining algorithm that discovers frequent items in a transaction database. It uses an iterative approach to generate candidate items and prunes the search space by leveraging the "apriori" property.

Use Cases:

- Market basket analysis to identify frequently co-occurring items in customer transactions.
- Recommendation systems to suggest related products based on purchase history.
- Web usage mining to understand sequential patterns in user behavior.

FP-Growth Algorithm

Definition: The FP-growth (Frequent Pattern Growth) algorithm is a frequent itemset mining algorithm that constructs a compact data structure called the FP-tree to efficiently mine frequent itemsets without generating candidate itemsets.

Use Cases:

- Market basket analysis for identifying association rules.
- Bioinformatics for identifying frequent sequence patterns in DNA sequences.
- Telecommunications for analyzing call patterns and customer behaviors.

from mlxtend.frequent_patterns import apriori
from mlxtend.preprocessing import TransactionEncoder
import pandas as pd

Sample transaction dataset (list of lists)
dataset = [[1, 2, 3], [2, 3, 4], [1, 3, 4],
[1, 2, 4], [1, 2, 3, 4], [2, 3, 5],
[1, 2, 5], [1, 3, 5], [2, 4, 5],
[3, 4, 5]]

```
# One-hot encode the dataset
t2 = TransactionEncoder()
t2_ary = t2.transform(dataset)
df_t2 = pd.DataFrame(t2_ary, columns=t2.columns_)

# Apply Apriori algorithm to find frequent itemsets
frequent_items = apriori(df_t2, min_support=0.4, use_colnames=True)
print(frequent_items)
```

from mlxtend.frequent_patterns import fprowth
from mlxtend.preprocessing import TransactionEncoder
import pandas as pd

Sample transaction dataset (list of lists)
dataset = [[1, 2, 3], [2, 3, 4], [1, 3, 4],
[1, 2, 4], [1, 2, 3, 4], [2, 3, 5],
[1, 2, 5], [1, 3, 5], [2, 4, 5],
[3, 4, 5]]

```
# One-hot encode the dataset
t2 = TransactionEncoder()
t2_ary = t2.transform(dataset)
df_t2 = pd.DataFrame(t2_ary, columns=t2.columns_)

# Apply FP-Growth algorithm to find frequent itemsets
frequent_items = fprowth(df_t2, min_support=0.4, use_colnames=True)
print(frequent_items)
```

Evaluating Non-Supervised Learning Models

Evaluating non-supervised learning models is an important step in ensuring that the model is effective and useful. However, it can be more