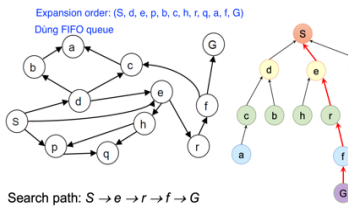
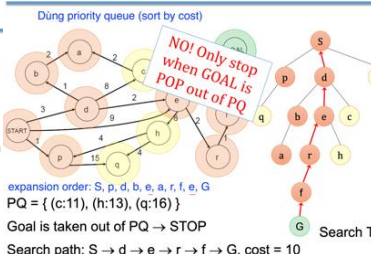


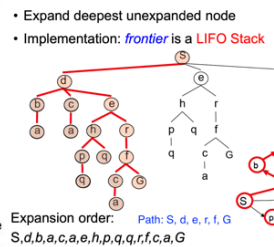
## Breadth-first search: An example



## Uniform-cost search: An example



## Depth-first search (DFS)



## Comparison of BFS and DFS

	DFS	BFS
Space complexity	Linear space	Maybe the whole search space
Time complexity	Same, better on the average (many goals, no loops, and no infinite paths)	Same, better in worst-cases
In general	better if many goals, not many loops, and much better in terms of memory.	better if goal is not deep, infinite paths, many loops, or small search space

DFS in use

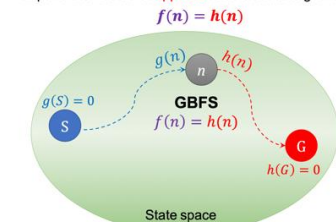
- The goal test is applied to each node when it is **generated** rather than when it is selected for expansion.
- Avoid repeated states by checking new states against those on the path from the root to the current node.

## Iterative deepening search (IDS)

- General strategy, often used in combination with **depth-first tree search** to find the best depth limit
- function ITERATIVE-DEEPENING-SEARCH(*problem*) returns a solution, or failure for depth = 0 to ∞ do  
  result ← DEPTH-LIMITED-SEARCH(*problem*, depth)  
  if result = cutoff then return result
- Gradually increase the limit until a goal is found.
  - The depth limit reaches the depth *d* of the **shallowest goal node**.

## Greedy best-first search

- Expand the node that **appears** to be closest to goal using

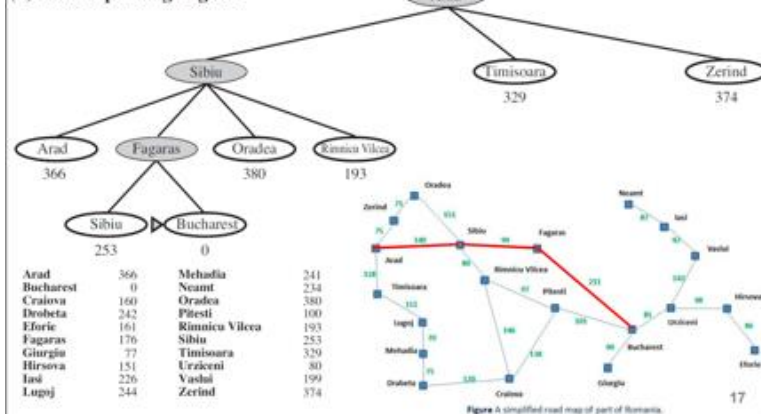


## A\* search

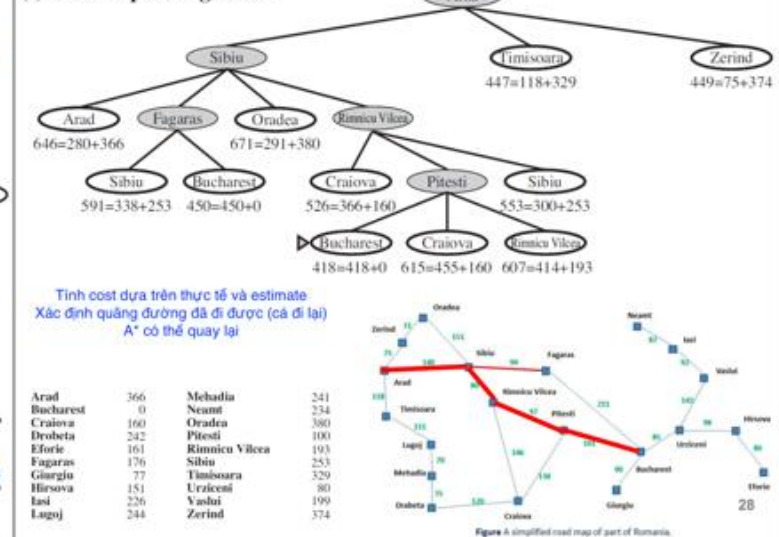
- The **most widely known** form of best-first search
  - Use **heuristic** to guide search, but not only
  - Avoid expanding paths that are already expensive
  - Ensure to compute a path with minimum cost
- Evaluate nodes by  $f(n) = g(n) + h(n)$
- where  $g(n)$  is the cost to reach the node *n* and  $h(n)$  is the cost to get from *n* to the goal
  - $f(n)$  = estimated cost of the cheapest solution through *n*

## Greedy best-first search: An example

### (d) After expanding Fagaras



### (f) After expanding Pitesti



## Conditions for optimality: Admissibility

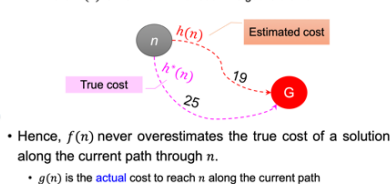
- $h(n)$  must be an **admissible heuristic**
- Never overestimate the cost to reach the goal → **optimistic**
- E.g., the straight-line distance  $h_{SLD}$

## Conditions for optimality: Consistency

- Admissibility is insufficient for graph search.
  - The optimal path to a repeated state could be discarded if it is not the first one selected.
  - $h(n)$  is **consistent** if for every node *n*, every successor *n'* of *n* generated by any action *a*,  
 $h(n) \leq c(n, a, n') + h(n')$
  - Every consistent heuristic is also admissible.
- Triangle inequality

## Conditions for optimality: Admissibility

- $h(n)$  is **admissible** if for every node *n*,  $h(n) \leq h^*(n)$
- where  $h^*(n)$  is the **true cost** to reach the goal state from *n*



## Conditions for optimality: Admissibility

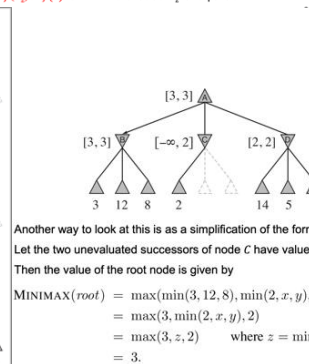
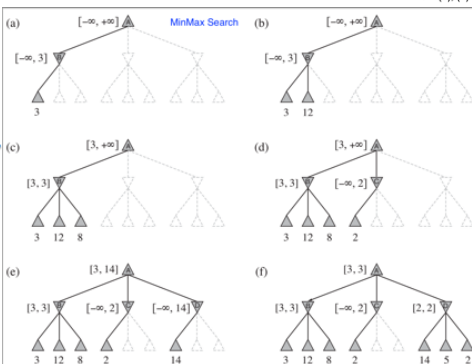
- If  $h(n)$  is **admissible**, A\* using **TREE-SEARCH** is optimal

- Suppose some suboptimal goal  $G_2$  has been generated and is in the frontier.
- Let *n* be an unexpanded node in the frontier such that *n* is on a shortest path to an optimal goal  $G_1$ .
- since  $h(G_2) = 0$  since  $G_2$  is suboptimal since  $h(G_2) = 0$   
 $f(G_2) = g(G_2)$   
 $f(G_2) > f(n)$   
 $h(n) \leq h^*(n)$  since *n* is admissible  
 $g(n) + h(n) \leq g(n) + h^*(n)$   
 $f(n) \leq f(G_2)$   
From (1), (2):  $f(G_2) > f(n) \rightarrow A^*$  will never select  $G_2$  for expansion

## Conditions for optimality: Consistency

- If  $h(n)$  is **consistent**, A\* using **GRAPH-SEARCH** is optimal

- If  $h(n)$  is consistent, the values of  $f(n)$  along any path are non-decreasing.
- Suppose  $n'$  is a successor of *n* →  $g(n') = g(n) + c(n, a, n')$
- $f(n') = g(n') + h(n') = g(n) + c(n, a, n') + h(n') \geq g(n) + h(n) = f(n)$
- Whenever A\* selects a node *n* for expansion, the **optimal path** to that node has been found.
- Proof by contradiction: There would have to be another frontier node  $n'$  on the optimal path from the start node to *n* (by the graph separation property)
- $f$  is nondecreasing along any path →  $f(n') < f(n) \rightarrow n'$  would have been selected first

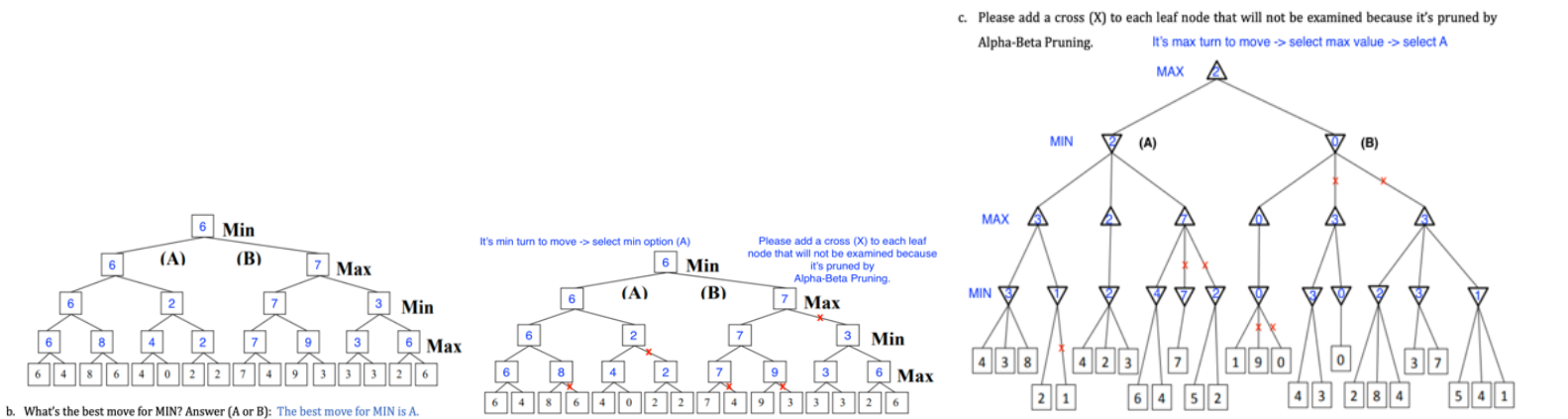


Another way to look at this is as a simplification of the formula for MINIMAX.

Let the two unevaluated successors of node *C* have values *x* and *y*.

Then the value of the root node is given by

MINIMAX(*root*) = max(min(3, 12, 8), min(2, *x*, *y*), min(14, 5, 2))  
= max(3, min(2, *x*, *y*), 2)  
= max(3, *z*, 2) where  $z = \min(2, x, y) \leq 2$   
= 3.



## Propositional logic: Semantics

Each model specifies true/false for each proposition symbol.

E.g.,  $m_1 = \{P_{1,2} = \text{false}, P_{2,2} = \text{false}, P_{3,1} = \text{true}\}$ , 8 possible models

Rules for evaluating truth with respect to a model  $m$

$P$	$Q$	$\neg P$	$P \wedge Q$	$P \vee Q$	$P \Rightarrow Q$	$P \Leftrightarrow Q$
true	true	false	true	true	true	true
true	false	true	false	true	false	false
false	true	false	false	true	true	false
false	false	true	false	false	true	true

Simple recursive process evaluates an arbitrary sentence.

E.g.,  $\neg P_{1,2} \wedge (P_{2,2} \vee P_{3,1}) = \text{true} \wedge (\text{true} \vee \text{false}) = \text{true} \wedge \text{true} = \text{true}$

## Logical equivalence

Two sentences,  $\alpha$  and  $\beta$ , are **logically equivalent** if they are true in the same set of models.

$$\alpha \equiv \beta \text{ iff } \alpha \models \beta \text{ and } \beta \models \alpha$$

$(\alpha \wedge \beta) \equiv (\beta \wedge \alpha)$	commutativity of $\wedge$
$(\alpha \vee \beta) \equiv (\beta \vee \alpha)$	commutativity of $\vee$
$((\alpha \wedge \beta) \wedge \gamma) \equiv (\alpha \wedge (\beta \wedge \gamma))$	associativity of $\wedge$
$((\alpha \vee \beta) \vee \gamma) \equiv (\alpha \vee (\beta \vee \gamma))$	associativity of $\vee$
$\neg(\neg\alpha) \equiv \alpha$	double-negation elimination
$(\alpha \Rightarrow \beta) \equiv (\neg\beta \Rightarrow \neg\alpha)$	contraposition
$(\alpha \Rightarrow \beta) \equiv (\neg\alpha \vee \beta)$	implication elimination
$(\alpha \Leftrightarrow \beta) \equiv ((\alpha \Rightarrow \beta) \wedge (\beta \Rightarrow \alpha))$	biconditional elimination
$\neg(\alpha \wedge \beta) \equiv (\neg\alpha \vee \neg\beta)$	De Morgan
$\neg(\alpha \vee \beta) \equiv (\neg\alpha \wedge \neg\beta)$	De Morgan
$(\alpha \wedge (\beta \vee \gamma)) \equiv ((\alpha \wedge \beta) \vee (\alpha \wedge \gamma))$	distributivity of $\wedge$ over $\vee$
$(\alpha \vee (\beta \wedge \gamma)) \equiv ((\alpha \vee \beta) \wedge (\alpha \vee \gamma))$	distributivity of $\vee$ over $\wedge$

## Conversion to CNF

- Eliminate  $\Leftrightarrow$ :  $\alpha \Leftrightarrow \beta \equiv (\alpha \Rightarrow \beta) \wedge (\beta \Rightarrow \alpha)$
- Eliminate  $\Rightarrow$ :  $\alpha \Rightarrow \beta \equiv \neg\alpha \vee \beta$
- The operator  $\neg$  appears only in literals: "move  $\neg$  inwards"
  - $\neg\neg\alpha \equiv \alpha$  (double-negation elimination)
  - $\neg(\alpha \wedge \beta) \equiv \neg\alpha \vee \neg\beta$  (De Morgan)
  - $\neg(\alpha \vee \beta) \equiv \neg\alpha \wedge \neg\beta$  (De Morgan)
- Apply the distributivity law to distribute  $\vee$  over  $\wedge$ 
  - $(\alpha \wedge \beta) \vee \gamma \equiv (\alpha \vee \gamma) \wedge (\beta \vee \gamma)$

## Horn clauses and Definite clauses

- Definite clause:** a disjunction of literals of which **exactly one is positive**.
  - E.g.,  $\neg P \vee \neg Q \vee R$  is a definite clause, whereas  $\neg P \vee Q \vee R$  is not.
- Horn clause:** a disjunction of literals of which **at most one is positive**.
  - All definite clauses are Horn clauses
- Goal clause:** clauses with **no positive literals**
- Horn clauses are closed under resolution
  - Resolving two Horn clauses will get back a Horn clause.

## Backus normal form (BNF)

$CNF Sentence$	$\rightarrow Clause_1 \wedge \dots \wedge Clause_n$
$Clause$	$\rightarrow Literal_1 \vee \dots \vee Literal_m$
$Literal$	$\rightarrow Symbol \mid \neg Symbol$
$Symbol$	$\rightarrow P \mid Q \mid R \mid \dots$
$Horn Clause Form$	$\rightarrow Definite Clause Form \mid Goal Clause Form$
$Definite Clause Form$	$\rightarrow (Symbol_1 \wedge \dots \wedge Symbol_i) \Rightarrow Symbol$
$Goal Clause Form$	$\rightarrow (Symbol_1 \wedge \dots \wedge Symbol_i) \Rightarrow \text{False}$

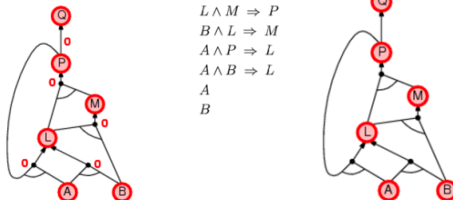
## The DPLL procedure

**function DPLL-SATISFIABLE( $s$ )** returns true or false  
**inputs:**  $s$ , a sentence in propositional logic  
 $clauses \leftarrow$  the set of clauses in the CNF representation of  $s$   
 $symbols \leftarrow$  a list of the proposition symbols in  $s$   
**return** DPLL( $clauses, symbols, \{\}$ )

**function DPLL( $clauses, symbols, model$ )** returns true or false  
If every clause in  $clauses$  is true in  $model$  then **return** true  
If some clause in  $clauses$  is false in  $model$  then **return** false  
 $P, value \leftarrow$  FIND-PURE-SYMBOL( $symbols, clauses, model$ )  
If  $P$  is non-null then **return** DPLL( $clauses, symbols - P, model \cup \{P=value\}$ )  
 $P, value \leftarrow$  FIND-UNIT-CLAUSE( $clauses, model$ )  
If  $P$  is non-null then **return** DPLL( $clauses, symbols - P, model \cup \{P=value\}$ )  
 $P \leftarrow$  FIRST( $symbols$ );  $rest \leftarrow$  REST( $symbols$ )  
**return** DPLL( $clauses, rest, model \cup \{P=true\}$ ) or DPLL( $clauses, rest, model \cup \{P=false\}$ )

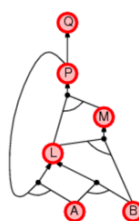
## Forward chaining: An example

$P \Rightarrow Q$   
 $L \wedge M \Rightarrow P$   
 $B \wedge L \Rightarrow M$   
 $A \wedge P \Rightarrow L$   
 $A \wedge B \Rightarrow L$   
 $A$   
 $B$



## Backward chaining: An example

$P \Rightarrow Q$   
 $L \wedge M \Rightarrow P$   
 $B \wedge L \Rightarrow M$   
 $A \wedge P \Rightarrow L$   
 $A \wedge B \Rightarrow L$   
 $A$   
 $B$



Q? ☒  
P? ☒  
L? ☒  
A? ☒  
B? ☒  
M? ☒  
L? ☒  
B? ☒

## Improvements in DPLL

- Early termination:** A clause is true if **any literal is true**, and a sentence is false if **any clause is false**.
  - Avoid examination of entire subtrees in the search space
  - E.g.,  $(A \vee B) \wedge (A \vee C)$  is true if  $A$  is true, regardless of  $B$  and  $C$
- Pure symbol heuristic:** a **pure symbol** always appears with the same "sign" in all clauses.
  - E.g.,  $(A \vee \neg B), (\neg B \vee \neg C), (A \vee C)$ ,  $A$  and  $B$  are pure,  $C$  is impure.
  - Make a pure symbol true  $\rightarrow$  Doing so never make a clause false
- Unit clause heuristic:** there is **only one literal in the clause** and thus this literal must be true
  - Unit propagation:** if the model contains  $B = \text{true}$  then  $(\neg B \vee \neg C)$  simplifies to a unit clause  $\neg C \rightarrow C$  must be false (so that  $\neg C$  is true)  $\rightarrow A$  must be true (so that  $A \vee C$  is true)

## Quantifiers: Universal quantification

Expressions of general rules  $\forall \langle \text{variables} \rangle \langle \text{sentence} \rangle$

- E.g., "All kings are persons."  $\forall x \text{ King}(x) \Rightarrow \text{Person}(x)$
- E.g., "Students of FIT are smart."  $\forall x \text{ Student}(x, \text{FIT}) \Rightarrow \text{Smart}(x)$

$\forall x P$  is true in a model  $m$  iff  $P$  is true with  $x$  being each possible object in the model.

It is equivalent to the **conjunction of instantiations** of  $P$ .

$\text{Student}(\text{Lan}, \text{FIT}) \Rightarrow \text{Smart}(\text{Lan})$   
 $\wedge \text{Student}(\text{Tuan}, \text{FIT}) \Rightarrow \text{Smart}(\text{Tuan})$   
 $\wedge \text{Student}(\text{Long}, \text{FIT}) \Rightarrow \text{Smart}(\text{Long})$   
 $\wedge \dots$

## Quantifiers: Existential quantification

Expressions of "some cases"  $\exists \langle \text{variables} \rangle \langle \text{sentence} \rangle$

- E.g., "Some students of FIT are smart."  $\exists x \text{ Student}(x, \text{FIT}) \wedge \text{Smart}(x)$

$\exists x P$  is true in a model  $m$  iff  $P$  is true with  $x$  being some possible object in the model.

It is equivalent to the **disjunction of instantiations** of  $P$ .

$\text{Student}(\text{Lan}, \text{FIT}) \wedge \text{Smart}(\text{Lan})$   
 $\vee \text{Student}(\text{Tuan}, \text{FIT}) \wedge \text{Smart}(\text{Tuan})$   
 $\vee \text{Student}(\text{Long}, \text{FIT}) \wedge \text{Smart}(\text{Long})$   
 $\vee \dots$

## FOL definite clause: An example

Consider the following problem

The law says that it is a crime for an American to sell weapons to hostile nations. The country **Nono**, an enemy of **America**, has some missiles, and all of its missiles were sold to it by Colonel **West**, who is **American**.

$\text{American}(x) \wedge \text{Weapon}(y) \wedge \text{Sells}(x, y, z) \wedge \text{Hostile}(z) \Rightarrow \text{Criminal}(x)$

$\exists x \text{ Owns}(\text{Nono}, x) \wedge \text{Missile}(x)$

$\text{Missile}(x) \Rightarrow \text{Owns}(\text{Nono}, x) \Rightarrow \text{Sells}(\text{West}, x, \text{Nono})$

$\text{Missile}(x) \Rightarrow \text{Weapon}(x)$

$\text{Enemy}(x, \text{America}) \Rightarrow \text{Hostile}(x)$

$\text{Owns}(\text{Nono}, M_1)$

$\text{Missile}(M_1)$

$\text{American}(\text{West})$

$\text{Enemy}(\text{Nono}, \text{America})$

## Quantifiers: A common mistake to avoid

Typically,  $\wedge$  is the main connective with  $\exists$

**Common mistake:** using  $\Rightarrow$  as the main connective with  $\exists$

- $\exists x \text{ Student}(x, \text{FIT}) \Rightarrow \text{Smart}(x)$
- It is true even with anyone who is not at FIT.

## CNF for First-order logic

**First-order resolution** requires that sentences be in CNF.

For example, the sentence

$\forall x \text{ American}(x) \wedge \text{Weapon}(y) \wedge \text{Sells}(x, y, z) \wedge \text{Hostile}(z) \Rightarrow \text{Criminal}(x)$

becomes, in CNF,

$\neg \text{American}(x) \vee \neg \text{Weapon}(y) \vee \neg \text{Sells}(x, y, z) \vee \neg \text{Hostile}(z) \vee \text{Criminal}(x)$

**Every sentence of first-order logic can be converted into an inferentially equivalent CNF sentence.**

The CNF sentence will be unsatisfiable just when the original sentence is unsatisfiable  $\rightarrow$  perform **proofs by contradiction**.

## Conversion to CNF

*Everyone who loves all animals is loved by someone.*

$\forall x [\forall y \text{ Animal}(y) \Rightarrow \text{Loves}(x, y)] \Rightarrow [\exists y \text{ Loves}(y, x)]$

1. **Eliminate implications**

$\forall x [\neg \forall y \neg \text{Animal}(y) \vee \text{Loves}(x, y)] \vee [\exists y \text{ Loves}(y, x)]$

$\forall x [\neg \neg \neg \text{Animal}(y) \vee \text{Loves}(x, y)] \vee [\exists y \text{ Loves}(y, x)]$

$\forall x [\neg \neg \neg \text{Animal}(y) \vee \neg \text{Loves}(x, y)] \vee [\exists y \text{ Loves}(y, x)]$

$\forall x [\neg \neg \neg \text{Animal}(y) \wedge \neg \text{Loves}(x, y)] \vee [\exists y \text{ Loves}(y, x)]$

$\forall x [\neg \neg \neg \text{Animal}(y) \wedge \neg \text{Loves}(x, y)] \vee [\exists y \text{ Loves}(y, x)]$

$\forall x [\neg \neg \neg \text{Animal}(y) \wedge \neg \text{Loves}(x, y)] \vee [\exists y \text{ Loves}(y, x)]$



Conversion to CNF

- Everyone who loves all animals is loved by someone.
- $\forall x [\forall y \text{ Animal}(y) \Rightarrow \text{Loves}(x,y)] \Rightarrow [\exists y \text{ Loves}(y,x)]$
3. **Standardize variables:** each quantifier uses a different one
- $\forall x [\exists y \text{ Animal}(y) \wedge \neg \text{Loves}(x,y)] \vee [\exists z \text{ Loves}(z,x)]$
4. **Skolemize:** remove existential quantifiers by elimination
- Simple case: translate  $\exists x P(x)$  into  $P(A)$ , where  $A$  is a new constant.
  - However,  $\forall x [\text{Animal}(A) \wedge \neg \text{Loves}(x,A)] \vee [\text{Loves}(B,x)]$  has an entirely different meaning.
  - The arguments of the **Skolem function** are all universally quantified variables in whose scope the existential quantifier appears.
- $\forall x [\text{Animal}(F(x)) \wedge \neg \text{Loves}(x,F(x))] \vee [\text{Loves}(G(x),x)]$

Conversion to CNF

- Everyone who loves all animals is loved by someone.
- $\forall x [\forall y \text{ Animal}(y) \Rightarrow \text{Loves}(x,y)] \Rightarrow [\exists y \text{ Loves}(y,x)]$
5. **Drop universal quantifiers**
- $[\text{Animal}(F(x)) \wedge \neg \text{Loves}(x,F(x))] \vee [\text{Loves}(G(x),x)]$
6. **Distribute  $\vee$  over  $\wedge$**
- $[\text{Animal}(F(x)) \vee \text{Loves}(G(x),x)] \wedge [\neg \text{Loves}(x,F(x)) \vee \text{Loves}(G(x),x)]$

Inference by enumeration

- Consider the following query
  - The hidden variables are *Earthquake* and *Alarm*.
  - Using initial letters for the variables, we have
- $P(B | j, m) = \alpha P(B, j, m) = \alpha \sum_e P(B, j, m, e, a)$
- For simplicity, we do this for *Burglary* = true.
- $P(b | j, m) = \alpha \sum_a \sum_e P(b) P(e) P(a | b, e) P(j | a) P(m | a)$
- Complexity:**  $O(n^{2^n})$  for a network of  $n$  Boolean variables

ID3 Decision tree algorithm

- The remaining examples are **all positive** (or **all negative**),  $\rightarrow$  DONE, it is possible to **answer Yes or No**.
- There are **some positive** and **some negative** examples  $\rightarrow$  **choose the best attribute** to split them

ID3 Decision tree: An example

Alternate?

True: 3 T, 3 F

False: 3 T, 3 F

Example	Alternate	Heavy	NotHeavy	Smelly	Spotted	Smooth	Edible
A	T	F	F	F	F	F	F
B	T	F	F	F	F	F	F
C	T	F	F	F	F	F	F
D	T	F	F	F	F	F	F
E	T	F	F	F	F	F	F
F	T	F	F	F	F	F	F
G	T	F	F	F	F	F	F
H	T	F	F	F	F	F	F
I	T	F	F	F	F	F	F
J	T	F	F	F	F	F	F
K	T	F	F	F	F	F	F
L	T	F	F	F	F	F	F
M	T	F	F	F	F	F	F
N	T	F	F	F	F	F	F
O	T	F	F	F	F	F	F
P	T	F	F	F	F	F	F
Q	T	F	F	F	F	F	F
R	T	F	F	F	F	F	F
S	T	F	F	F	F	F	F
T	T	F	F	F	F	F	F
U	T	F	F	F	F	F	F
V	T	F	F	F	F	F	F
W	T	F	F	F	F	F	F
X	T	F	F	F	F	F	F
Y	T	F	F	F	F	F	F
Z	T	F	F	F	F	F	F
A	F	T	T	T	T	T	T
B	F	T	T	T	T	T	T
C	F	T	T	T	T	T	T
D	F	T	T	T	T	T	T
E	F	T	T	T	T	T	T
F	F	T	T	T	T	T	T
G	F	T	T	T	T	T	T
H	F	T	T	T	T	T	T
I	F	T	T	T	T	T	T
J	F	T	T	T	T	T	T
K	F	T	T	T	T	T	T
L	F	T	T	T	T	T	T
M	F	T	T	T	T	T	T
N	F	T	T	T	T	T	T
O	F	T	T	T	T	T	T
P	F	T	T	T	T	T	T
Q	F	T	T	T	T	T	T
R	F	T	T	T	T	T	T
S	F	T	T	T	T	T	T
T	F	T	T	T	T	T	T
U	F	T	T	T	T	T	T
V	F	T	T	T	T	T	T
W	F	T	T	T	T	T	T
X	F	T	T	T	T	T	T
Y	F	T	T	T	T	T	T
Z	F	T	T	T	T	T	T

- Calculate **Average Entropy** of attribute Alternate

$AE_{Alternate} = P(Alt = T) \times H(Alt = T) + P(Alt = F) \times H(Alt = F)$

$AE_{Alternate} = \frac{6}{12} \left[ -\left(\frac{3}{6} \log_2 \frac{3}{6}\right) - \left(\frac{3}{6} \log_2 \frac{3}{6}\right) \right] + \frac{6}{12} \left[ -\left(\frac{3}{6} \log_2 \frac{3}{6}\right) - \left(\frac{3}{6} \log_2 \frac{3}{6}\right) \right] = 1$

a.  $H_{Edible} = H[3+, 5-] \stackrel{def}{=} -\frac{3}{8} \log_2 \frac{3}{8} - \frac{5}{8} \log_2 \frac{5}{8} = \frac{3}{8} \log_2 \frac{8}{3} + \frac{5}{8} \log_2 \frac{8}{5} = \frac{3}{8} \cdot 3 - \frac{3}{8} \log_2 3 + \frac{5}{8} \cdot 3 - \frac{5}{8} \log_2 5 = 3 - \frac{3}{8} \log_2 3 - \frac{5}{8} \log_2 5 = 3 - \frac{3}{8} \cdot 1.58 - \frac{5}{8} \cdot 2.32 = 3 - 0.59 - 1.44 = 0.9544$

b.  $H_{NotHeavy} = H[3+, 5-] \stackrel{def}{=} -\frac{3}{8} \log_2 \frac{3}{8} - \frac{5}{8} \log_2 \frac{5}{8} = \frac{3}{8} \log_2 \frac{8}{3} + \frac{5}{8} \log_2 \frac{8}{5} = \frac{3}{8} \cdot 3 - \frac{3}{8} \log_2 3 + \frac{5}{8} \cdot 3 - \frac{5}{8} \log_2 5 = 3 - \frac{3}{8} \log_2 3 - \frac{5}{8} \log_2 5 = 3 - 0.59 - 1.44 = 0.9544$

c.  $H_{Smooth} = H[3+, 5-] \stackrel{def}{=} -\frac{3}{8} \log_2 \frac{3}{8} - \frac{5}{8} \log_2 \frac{5}{8} = \frac{3}{8} \log_2 \frac{8}{3} + \frac{5}{8} \log_2 \frac{8}{5} = \frac{3}{8} \cdot 3 - \frac{3}{8} \log_2 3 + \frac{5}{8} \cdot 3 - \frac{5}{8} \log_2 5 = 3 - \frac{3}{8} \log_2 3 - \frac{5}{8} \log_2 5 = 3 - 0.59 - 1.44 = 0.9544$

d.  $H_{NotHeavy} = H[3+, 5-] \stackrel{def}{=} -\frac{3}{8} \log_2 \frac{3}{8} - \frac{5}{8} \log_2 \frac{5}{8} = \frac{3}{8} \log_2 \frac{8}{3} + \frac{5}{8} \log_2 \frac{8}{5} = \frac{3}{8} \cdot 3 - \frac{3}{8} \log_2 3 + \frac{5}{8} \cdot 3 - \frac{5}{8} \log_2 5 = 3 - \frac{3}{8} \log_2 3 - \frac{5}{8} \log_2 5 = 3 - 0.59 - 1.44 = 0.9544$

ID3 Decision tree: An example

Alternate?

True: 3 T, 3 F

False: 3 T, 3 F

**Information Gain** is the difference in entropy from before to after the set  $S$  is split on the selected attribute.

$IG(Alternate, S) = H(S) - AE_{Alternate} = 1 - 1 = 0$

Node 1: Smooth = 0

Node 2: Smooth = 1

Node 3: Smooth = 0

Everyone who loves all animals is loved by someone.

Anyone who kills an animal is loved by no one.

Jack loves all animals.

Either Jack or Curiosity killed the cat, who is named Tuna.

Did Curiosity kill the cat?

- A.  $\text{Animal}(F(x) \vee \text{Loves}(G(x), x) \wedge \neg \text{Loves}(x, F(x)) \vee \text{Loves}(G(x), x)$
- B.  $\neg \text{Loves}(y, x) \vee \neg \text{Animal}(x) \vee \neg \text{Kills}(x, x)$
- C.  $\neg \text{Animal}(x) \vee \neg \text{Loves}(Jack, x)$
- D.  $\text{Kills}(Jack, Tuna) \vee \text{Kills}(Curiosity, Tuna)$
- E.  $\text{Cat}(Tuna)$
- F.  $\neg \text{Cat}(x) \vee \text{Animal}(x)$
- G.  $\neg \text{Kills}(Curiosity, Tuna)$

The wet grass example

- What is the probability that it is raining, given the grass is wet?
- Using the expansion for the joint probability function  $P(G, S, R)$  and the conditional probabilities from the CPTs stated in the diagram
- The numerical results (subscripted by the associated variable values) are

Variable elimination: Factorization

- First, we sum out  $A$  from the product of  $f_3, f_4$ , and  $f_5$ .
- Now we are left with  $P(B | j, m) = \alpha f_1(B) \times \sum_e f_2(E) \times f_6(B, E)$
- Next, we sum out  $E$  from the product of  $f_2$  and  $f_6$
- The following expression can be evaluated by taking pointwise product and normalizing the result.

A purity measure with entropy

- Entropy** is a measure of the uncertainty of a random variable  $V$  with values  $v_k$ .
- $H(V) = -\sum_k P(v_k) \log_2 \frac{1}{P(v_k)} = -\sum_k P(v_k) \log_2 P(v_k)$
- $v_k$  is a class in  $V$  (e.g., yes/no in binary classification)
- $P(v_k)$  is the proportion of the number of elements in class  $v_k$  to the number of elements in  $V$

Full joint distribution with BN

- An entry in the joint distribution is the probability of a variable assignment, such as  $P(X_1 = x_1 \wedge \dots \wedge X_n = x_n)$ .
- where  $\text{parent}(X_i)$  denotes the values of  $\text{Parent}(X_i)$  that appear in  $x_1, \dots, x_n$ .
- Thus, it is the product of the appropriate elements of the CPTs in the Bayesian network.
- A Bayesian network can be used to answer any query, by summing all the relevant joint entries.

Inference by enumeration

- A query can be answered by computing sums of products of conditional probabilities from the Bayesian network.
- $P(X | e) = \alpha P(X, e) = \alpha \sum_y P(X, e, y)$
- where  $\alpha$  stands for the constant denominator term, which is usually simplified during calculation.

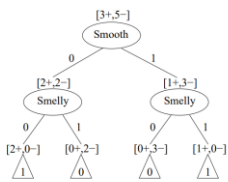
Variable elimination: An example

- Consider the following network. Calculate  $P(B | \neg c)$ .
- Irrelevant variable:  $D$ . Observed variable:  $C = \neg c$ .
- Sum out  $A$  to have  $f_4(B) = \sum_a P(B | a) \times P(a)$
- Join  $f_1$  and  $f_4$ :  $f_5(B, \neg c) = f_1(B) \times f_4(B)$
- Finally, we have  $P(B | \neg c) = \alpha f_5(B, \neg c)$
- Assume that  $B = b$ . Normalize for  $B$ :

$P(b | \neg c) = \frac{f_5(b, \neg c)}{f_5(b, c) + f_5(b, \neg c)}$

**Example**

	NotHeavy	Smelly	Spotted	Smooth	Edible
A	1	0	0	0	1
B	1	0	1	0	1
C	0	1	0	1	1
D	0	0	0	1	0
E	1	1	1	0	0
F	1	0	1	1	0
G	1	0	0	1	0
H	0	1	0	0	0
I	0	1	1	1	?
J	1	1	0	1	?
K	1	1	0	0	?



IF (Smooth = 0 AND Smelly = 0) OR (Smooth = 1 AND Smelly = 1)  
THEN Edible;  
ELSE ~Edible;

Classification of test instances:

U	Smooth = 1, Smelly = 1 ⇒ Edible = 1
V	Smooth = 1, Smelly = 1 ⇒ Edible = 1
W	Smooth = 0, Smelly = 1 ⇒ Edible = 0

- $P(H)$  (prior probability): the initial probability
  - E.g., X will buy computer, regardless of age, income, ...
- $P(X)$ : the probability that sample data is observed
  - E.g., X is 31.40 and has a medium income, regardless of the buying
- $P(X|H)$  (likelihood): the probability of observing the sample X, given that the hypothesis holds
  - E.g., given that X will buy computer, the probability that X is 31.40 and has a medium income
- $P(H|X) = \frac{P(X|H)P(H)}{P(X)}$  (posterior probability)
  - E.g., given that X is 31.40 and has a medium income, the probability that X will buy computer

## Naïve Bayesian classification

- **Class-conditional independence:** There are no dependence relationships among the attributes
- The **naïve Bayesian classification** formula is written as

$$P(X|C_i) = \prod_{k=1}^n P(x_k|C_i) = P(x_1|C_i) \times P(x_2|C_i) \times \dots \times P(x_n|C_i)$$

- $A_k$  is categorical:  $P(x_k|C_i)$  is the number of tuples in  $C_i$  having value  $x_k$  for  $A_k$  divided by  $|C_i|$  (# of tuples of  $C_i$  in  $D$ )
- $A_k$  is continuous:  $P(x_k|C_i) = g(x_k, \mu_{C_i}, \sigma_{C_i})$  with the Gaussian distribution  $g(x, \mu, \sigma) = \frac{1}{\sigma\sqrt{2\pi}} e^{-\frac{(x-\mu)^2}{2\sigma^2}}$

- Count class distributions only → computation cost reduced

## Naïve Bayesian classification: An example

P(buys_computer = "yes")	10/16
P(buys_computer = "no")	6/16

	buys_computer = "yes"	buys_computer = "no"
age = "<=30"	3/12	4/8
age = "31...40"	5/12	1/8
age = ">40"	4/12	3/8
income = "low"	4/12	2/8
income = "medium"	5/12	3/8
income = "high"	3/12	3/8
student = "yes"	7/11	2/7
student = "no"	4/11	5/7
credit_rating = "fair"	7/11	3/7
credit_rating = "excellent"	4/11	4/7

## Back-propagation learning rule

- **Step 1 – Initialization:** Initial weights and thresholds are assigned to random numbers.
  - The numbers may be uniformly distributed in the range  $\left(-\frac{2.4}{F_i}, \frac{2.4}{F_i}\right)$  (Haykin, 1999), where  $F_i$  is the total number of inputs of neuron
  - The weight initialization is done on a neuron-by-neuron basis
- **Step 2 – Activation:** At iteration  $p$ , apply the  $p^{\text{th}}$  example, which has inputs  $x_1(p), x_2(p), \dots, x_n(p)$  and desired outputs  $y_{d1}(p), y_{d2}(p), \dots, y_{dL}(p)$ .
  - (a) Calculate the actual output, from  $n$  inputs, of neuron  $j$  in the hidden layer.
 
$$y_j(p) = \sigma \left( \sum_{i=1}^n x_i(p)w_{ij}(p) - \theta_j \right)$$

$$\sigma(x) = \frac{1}{1 + e^{-x}}$$
  - (b) Calculate the actual output, from  $k$  inputs, of neuron  $m$  in the hidden layer.
 
$$y_k(p) = \sigma \left( \sum_{j=1}^m y_j(p)w_{jk}(p) - \theta_k \right)$$

## Foundation: Based on Bayes' Theorem

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

Likelihood
Class Prior Probability

Posterior Probability
Predictor Prior Probability

$$P(c|X) = P(x_1|c) \times P(x_2|c) \times \dots \times P(x_n|c) \times P(c)$$

## Classification with Bayes' Theorem

- Let  $D$  be a training set of tuples and associated class labels
- Each tuple is represented by a  $n$ -attribute  $X = (x_1, x_2, \dots, x_n)$
- Suppose there are  $m$  classes  $C_1, C_2, \dots, C_m$
- Classification is to derive the **maximum posteriori**  $P(C_i|X)$  from **Bayes' theorem**

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

- $P(X)$  is constant for all classes, only  $P(X|C_i)P(C_i)$  needs to be maximized.

age	income	student	credit_rating	buys_computer
<=30	medium	yes	fair	?

- $P(X|C_i)$ 
  - $P(X|buys\_computer = "yes") = 2/9 \times 4/9 \times 6/9 \times 6/9 = 0.044$
  - $P(X|buys\_computer = "no") = 3/5 \times 2/5 \times 1/5 \times 2/5 = 0.019$
- $P(X|C_i) \times P(C_i)$ 
  - $P(X|buys\_computer = "yes") \times P(buys\_computer = "yes") = 0.028$
  - $P(X|buys\_computer = "no") \times P(buys\_computer = "no") = 0.007$
- $P(C_i|X)$ 
  - $P(buys\_computer = "yes" | X) = 0.8$
  - $P(buys\_computer = "no" | X) = 0.2$

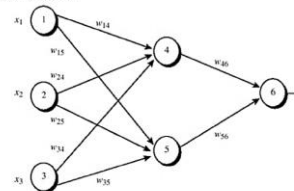
Therefore, X belongs to class ("buys\_computer = yes")

## Perceptron learning rule

- **Step 1 – Initialization:** Initial weights  $w_1, w_2, \dots, w_n$  and threshold  $\theta$  are randomly assigned to small numbers (usually in  $[-0.5, 0.5]$ , but not restricted to).
- **Step 2 – Activation:** At iteration  $p$ , apply the  $p^{\text{th}}$  example, which has inputs  $x_1(p), x_2(p), \dots, x_n(p)$  and desired output  $y_d(p)$ , and calculate the actual output
 
$$Y(p) = \sigma \left( \sum_{i=1}^n x_i(p)w_i(p) - \theta \right)$$

$$\sigma(x) = \begin{cases} 1 & \text{if } x \geq 0 \\ 0 & \text{if } x < 0 \end{cases}$$
 where  $n$  is the number of perceptron inputs and  $step$  is the activation function
- **Step 3 – Weight training**
  - Update the weights  $w_i$ :  $w_i(p+1) = w_i(p) + \Delta w_i(p)$  where  $\Delta w_i(p)$  is the weight correction at iteration  $p$
  - The **delta rule** determines how to adjust the weights:  $\Delta w_i(p) = \alpha \times x_i(p) \times e(p)$  where  $\alpha$  is the learning rate ( $0 < \alpha < 1$ ) and  $e(p) = Y_d(p) - Y(p)$
- **Step 4 – Iteration:** Increase iteration  $p$  by one, go back to Step 2 and repeat the process until convergence.

**Problem.** Consider the following neuron network, which includes 3 input neurons, 2 hidden neurons and 1 output neuron.



Initial input, weight and bias values are

X1	X2	X3	W14	W15	W24	W25	W34	W35	W46	W56	θ4	θ5	θ6
1	0	1	0.2	-0.3	0.4	0.1	-0.5	0.2	-0.3	-0.2	-0.4	0.2	0.1

The expected output value is 1. The learning rate is 0.9  
Knowing that the actual output at some neuron  $j$  is calculated as follows:

$$y_j(p) = \text{sigmoid} \left[ \sum_{i=1}^n x_i(p) \times w_{ij}(p) + \theta_j \right]$$

## Back-propagation learning rule

- **Step 3 – Weight training:** Update the weights in the back-propagation network and propagate backward the errors associated with output neurons.
  - (a) Calculate the **error gradient** for neuron  $k$  in the output layer
 
$$\delta_k(p) = y_k(p) \times [1 - y_k(p)] \times [y_{dk}(p) - y_k(p)]$$

$$e_k(p)$$
 Calculate the weight corrections:  $\Delta w_{jk}(p) = \alpha \times y_j(p) \times \delta_k(p)$   
Update the weights at the output neurons:  $w_{jk}(p+1) = w_{jk}(p) + \Delta w_{jk}(p)$
  - (b) Calculate the error gradient for neuron  $j$  in the hidden layer
 
$$\delta_j(p) = y_j(p) \times [1 - y_j(p)] \times \sum_{k=1}^l \delta_k(p)w_{jk}(p)$$
 Calculate the weight corrections:  $\Delta w_{ij}(p) = \alpha \times x_i(p) \times \delta_j(p)$   
Update the weights at the hidden neurons:  $w_{ij}(p+1) = w_{ij}(p) + \Delta w_{ij}(p)$
- **Step 4: Iteration:** Increase iteration  $p$  by one, go back to Step 2 and repeat the process until the selected error criterion is satisfied.

a) Ignore all biases (*precision to 3 decimal places*).

*Ignore all biases – Forward*

Output at neuron 4	0.426
Output at neuron 5	0.475
Output at neuron 6	0.527

*Ignore all biases – Backward*

Error gradient at neuron 6	0.118
Error gradient at neuron 5	-0.006

Error gradient at neuron 4	-0.009
Update w46	-0.255
Update w56	-0.150
Update w14	0.192
Update w15	-0.305
Update w24	0.400
Update w25	0.100
Update w34	-0.508
Update w35	0.195

b) Consider all biases such that each bias is treated as a neuron and thus it will be also updated (*precision to 3 decimal places*).

*(Consider all biases – Forward*

Output at neuron 4	0.332
Output at neuron 5	0.525
Output at neuron 6	0.552

$$y(p) = \text{sigmoid}[\sum_{i=1}^n x_i(p) \times w_{ij}(p) + \theta_j]$$

where  $n$  is the number of inputs of neuron  $j$ ,  $w_{ij}$  is the corresponding link from a neuron  $i$  in the previous layer to neuron  $j$ , and  $\theta_j$  is the bias at neuron  $j$ .

Present all calculations required to perform the backpropagation once (I.e., one forward pass and one backward pass) on the given neural network in the following cases

*Consider all biases – Backward*

Error gradient at neuron 6	0.111
Error gradient at neuron 5	-0.006
Error gradient at neuron 4	-0.007
Update w46	-0.267
Update w56	-0.148
Update w14	0.193
Update w15	-0.305
Update w24	0.400
Update w25	0.100
Update w34	-0.507
Update w35	0.195
Update bias 6	-0.407
Update bias 5	0.195
Update bias 4	0.200