

Artificial Intelligence  
Assignment - 2

Samip Subedi  
070BCT 536

1. Constraint satisfaction problems (CSPs) are mathematical problems defined as a set of objects whose state must satisfy a number of constraints or limitations. A CSP mainly contains three components:  $X, D \& C$ , where
- $X$  is a set of variables  $\{x_1, x_2, \dots, x_n\}$
  - $D$  is a set of Domains  $\{D_1, D_2, \dots, D_n\}$
  - $C$  is a set of constraints that specify allowable combinations of values.

a) FORTY

TEN

+ TEN

SIXTY

Here, ~~STORY~~

$$X = \{F, O, R, T, Y, E, N, S, I, X\}$$

$$D = \{\text{Integers, Alphabets}\}$$

$$C = \{0 \leq X \leq 9; X_i \neq X_j, 0 \leq c_i \leq 2\}$$

Now,

$$\begin{array}{r} c_4 & c_3 & c_2 & c_1 \\ F & O & R & T & Y \\ + & & & & \\ \hline S & I & X & T & Y \end{array}$$

Hence,

$$Y + 2N = Y + 10c_1$$

$$c_1 + T + 2E = T + 10c_2$$

$$c_2 + R + 2T = X + 10c_3$$

$$c_3 + O = I + 10c_4$$

$$c_4 + F = S$$

$$C_4 = 1 \quad (\because F \neq S, C_4 \neq 2 \text{ because } C_4 \text{ comes from } C_3 + O)$$

$$\Rightarrow S = F + 1$$

$$\text{Also, } C_1 = 0$$

$$N = 0 \text{ or } 5$$

$$E = 0 \text{ or } 5$$

$$N \neq 5 \quad (\because C_1 = 1 \text{ doesn't balance } T + C_1 + 0 \times 2 = T)$$

$$\therefore N = 0, E = 5$$

The problem becomes

$$\begin{array}{r}
 1 \quad C_3 \quad C_2 \quad 0 \\
 F \quad 0 \quad R \quad T \quad Y \\
 \quad \quad T \quad 5 \quad \textcircled{0} \\
 + \quad \quad T \quad 5 \quad \textcircled{0} \\
 \hline
 S \quad I \quad X \quad T \quad Y
 \end{array}$$

$$\text{Here, } 0 = 9 \text{ and } C_3 = 2 \quad (C_3 \neq 0 \text{ or } 1 \text{ as } 0 \neq I) \Rightarrow I = 1$$

$$\text{Let } T = 8 \Rightarrow C_2 = 1 :$$

$$\begin{array}{r}
 1 \quad 2 \quad 1 \quad 0 \\
 F \quad 9 \quad R \quad 8 \quad Y \\
 \quad \quad 8 \quad 5 \quad \textcircled{0} \\
 + \quad \quad 8 \quad 5 \quad \textcircled{0} \\
 \hline
 S \quad I \quad X \quad 8 \quad Y
 \end{array}$$

Now, Remaining integers are  $\{2, 3, 4, 6, 7\}$

$$\therefore F = 2, S = 3, \text{ etc.}$$

$$\text{or } F = 6, S = 7$$

$$\therefore F = 2, S = 3, R = 7, T = 8, Y = 6.$$

lets take the first option, and  $R = 7 \Rightarrow X = 4$ . Finally  $Y = 6$ .

$\therefore$  The solution is:

$$\begin{array}{r}
 1 \quad 2 \quad 1 \quad 0 \\
 2 \quad 9 \quad 7 \quad 8 \quad 6 \\
 \quad \quad 8 \quad 5 \quad \textcircled{0} \\
 + \quad \quad 8 \quad 5 \quad \textcircled{0} \\
 \hline
 3 \quad 1 \quad 4 \quad 8 \quad 6
 \end{array}$$

$$\begin{array}{lll}
 \therefore N = 0 & E = 5 \\
 I = 1 & Y = 6 \\
 F = 2 & R = 7 \\
 S = 3 & T = 8 \\
 X = 4 & O = 9
 \end{array}$$

$$\begin{array}{r} \text{b) LOGIC} \\ + \text{LOGIC} \\ \hline \text{PROLOG} \end{array}$$

Here,

$$X = \text{set of variables} = \{L, O, G, I, C, P, R\}$$

$$D = \text{set of Domains} = \{\text{Integers, Alphabets}\}$$

$$C = \text{set of constraints}$$

$$= \{0 \leq x_i \leq 9, x_i \neq x_j \text{ if } i \neq j, 0 \leq c_i \leq 1\} \text{ where } c_i \text{ is the carry involved in the problem.}$$

So, the problem becomes

$$\begin{array}{r} c_5 \ c_4 \ c_3 \ c_2 \ c_1 \\ L \ O \ G \ I \ C \\ + \text{LOGIC} \\ \hline \text{PROLOG} \end{array}$$

Here,

$$2C = G_1 + 10C_1$$

$$2I + C_1 = O + 10C_2$$

$$2G + C_2 = L + 10C_3$$

$$2O + C_3 = 2O + 10C_4$$

$$2L + C_4 = R + 10C_5$$

$$c_5 = P$$

We can see that:

$$2O + C_3 = O + 10C_4$$

$$\Rightarrow C_3 = O, O = 0$$

$$\text{And, } P = c_5 = 1 \text{ } (P \neq 0)$$

i.e.

$$\begin{array}{r} 1 \ c_4 \ c_3 \ c_2 \ c_1 \\ L \ O \ G \ I \ C \\ + \text{LOGIC} \\ \hline 1 \ R \ O \ L \ O \ G_1 \end{array}$$

Let  $C = 2 \Rightarrow G = 4$  and  $C_1 = 0$

$$\Rightarrow 2I = 0 + 10C_2$$

$$\text{So, } I = 5 \Rightarrow C_2 = 1 \text{ and } L = 4+4+1 = 9$$

The problem is now:

$$\begin{array}{r} 9 \ 0 \ 4 \ 5 \ 2 \\ + 9 \ 0 \ 4 \ 5 \ 2 \\ \hline 1 \ 8 \ 0 \ 9 \ 0 \ 4 \end{array}$$

$$\therefore L = 9, O = 0, G = 4, I = 5, C = 2, P = 1 = R = 8.$$

$\Rightarrow$  WRONG

$$\begin{array}{r} \text{WRONG} \\ + \text{WRONG} \\ \hline \text{RIGHT} \end{array}$$

Here,

$X = \text{Set of variables} = \{W, R, O, N, G, I, H, T\}$

$D = \text{Set of Domains} = \{\text{Integers, Alphabets}\}$

$C = \text{Set of constraints}$   
 $= \{0 \leq x_i \leq 9, x_i \neq x_j \forall i, j, 0 \leq c_i \leq 1\}$  where  $c_i$  is the carry involved in the problem.

So, the problem becomes:

$$\begin{array}{r} c_4 \ c_3 \ c_2 \ c_1 \\ \text{WRONG} \\ + \text{WRONG} \\ \hline \text{RIGHT} \end{array}$$

$$\Rightarrow 2G = T + 10C_1$$

$$2N + C_1 = H + 10C_2$$

$$2O + C_2 = G + 10C_3$$

$$2R + C_3 = I + 10C_4$$

$$2W + C_4 = R$$

Assume  $G = 6 \Rightarrow T = 2, C_1 = 1$ , the problem becomes

$$\begin{array}{r} c_4 \ c_3 \ c_2 \ t \\ \text{WRONG} \\ + \text{WRONG} \\ \hline \text{RIGHT} \end{array}$$

Assume:

$$G = 5 \Rightarrow T = 0, C_1 = 1$$

$$N = 6 \Rightarrow H = 3, C_2 = 1$$

$$O = 2 \Rightarrow G = 5, C_3 = 0$$

$$R = 9 \Rightarrow I = 8, C_4 = 1$$

$$W = 4 \Rightarrow L = 9$$

So, the final solution is

$$\begin{array}{r} & & 1 \\ 4. & 9 & 2 & 6 & 5 \\ + & 4 & 9 & 2 & 6 & 5 \\ \hline 9 & 8 & 5 & 3 & 0 \end{array}$$

$$\therefore W = 4$$

$$R = 9$$

$$O = 2$$

$$N = 6$$

$$G = 5$$

$$I = 8$$

$$H = \cancel{3}$$

$$T = 0$$

2

"All married employees earning Rs 300,000 or more per year in Nepal pay taxes. All unmarried employees earning Rs 250,000 or more per year in Nepal pay taxes. The president of Nepal earns Rs. 2,500,000 and has to pay maximum taxes. No other employee earns more than the president. Some of the Nepalese citizens earn less than Rs 100 per day and they don't have to pay any taxes."

$\Rightarrow$  Let's assume, for FOL, the following predicates:

- i)  $\text{married}(x)$ ;  $x$  is married
- ii)  $\text{employed}(x, y)$ ;  $x$  is employed as  $y$  (profession)
- iii) ~~earnpayyear~~  $\text{earns}(x, z)$ ;  $x$  earns  $z$
- iv)  $\text{greaterthan}(z)$ ;  $z$  is greater than  $z$ ;  $\text{lt}(z)$ ;  $z$  is less than  $z$
- v)  $\text{equalto}(z)$ ;  $z$  is equal to  $z$
- vi)  $\text{paytax}(a)$ ; paytax of amount  $a$
- vii)  $\text{max}(x)$ ;  $x$  is maximum of all
- viii)  $\text{Nepalesecitizen}(x)$ ;  $x$  is a Nepalese citizen
- ix)  $\text{earnpayday}(x, z)$ ;  $x$  earns  $z$  per day
- x)  $\text{lessthan}(a, b)$ ;  $a$  is less than  $b$
- xi)  $\text{is}(x, y)$ ;  $x$  is  $y$

The first sentence is transformed to:

a)  $\forall x \exists t (\text{married}(x) \wedge \text{employed}(x, y) \wedge \text{earnpayyear}(x, \text{greaterthan}(300,000)))$   
 $\rightarrow \text{paytax}(t)$

i.e. for all person 'x', ~~for all employment 'y'~~, where  $x$  is married, is employed as  $y$  and earns greater than Rs 300,000 per year, there exists some amount of tax 't' to be paid by  $x$ .

b)  $\forall x \exists t (\neg \text{married}(x) \wedge \text{employed}(x, y) \wedge \text{earnpayyear}(x, \text{greaterthan}(250,000)))$   
 $\rightarrow \text{paytax}(x)$

i.e. for all person 'x', ~~for all employment 'y'~~, if  $x$  is unmarried, is employed as  $y$  and earns greater than Rs 250,000 per year, there exists some amount of tax 't' to be paid by  $x$ .

c)  $\exists x \exists t \text{employed}(x, \text{President}) \wedge \text{earnpayyear}(x, \text{equalto}(2,500,000)) \rightarrow \text{paytax}($   
 $\max(t))$   
i.e. There exists a person  $x$  who is employed as the President of Nepal, earns Rs 2,500,000 per year and pays the maximum tax  $t$ .

The facts are transformed to FOL as:

- a)  $\forall x \exists t (\text{married}(x) \wedge \text{employed}(x) \wedge \text{earnperyear}(x, \text{greaterthan}(300,000))$   
→  $\text{paytax}(t)$   
i.e. For all married, employed person  $x$  who earns more than Rs 300,000 per year has to pay a tax ' $t$ '.
- b)  $\forall x \exists t (\neg \text{married}(x) \wedge \text{employed}(x) \wedge \text{earnperyear}(x, \text{greaterthan}(250,000))$   
→  $\text{paytax}(t)$   
i.e. for all unmarried, employed person  $x$  who earns more than Rs 250,000 per year has to pay a tax ' $t$ '.
- c)  $\exists t [\text{employed}(\text{President}) \wedge \text{earnperyear}(\text{President}, \text{equalto}(2,500,000)) \wedge \text{paytax}(\max(t))]$   
i.e. The employed President of Nepal earns Rs 2,500,000 per year and there's a tax ' $t$ ', which is the maximum of all.
- d)  $\forall x \text{ employed}(x) \wedge \neg \text{is}(x, \text{President}) \rightarrow \text{less than}(\text{earnperyear}(x, a), \text{earnperyear}(\text{President}, b))$   
i.e. For all employed person  $x$  who is not a president, the earning of  $x$  is less than that of the President.
- e)  $\exists x \text{ Nepalesecitizen}(x) \wedge \text{earnperday}(\text{less than } \text{lt}(100)) \rightarrow \text{paytax}(0)$   
i.e. There exists some Nepalese citizen  $x$ , who earns less than Rs 100 per day and ~~he~~ doesn't have to pay ~~to~~ any tax, or a tax of Rs 0.

3

Given facts are:

- John likes all kinds of food
- Apples are food
- Chicken is food
- Anything anyone eats and isn't killed by is food
- Bill eats peanuts and is still alive
- Sue eats everything Bill eats.

Representing ~~the~~<sup>in</sup> predicates, the above facts: ~~are~~:

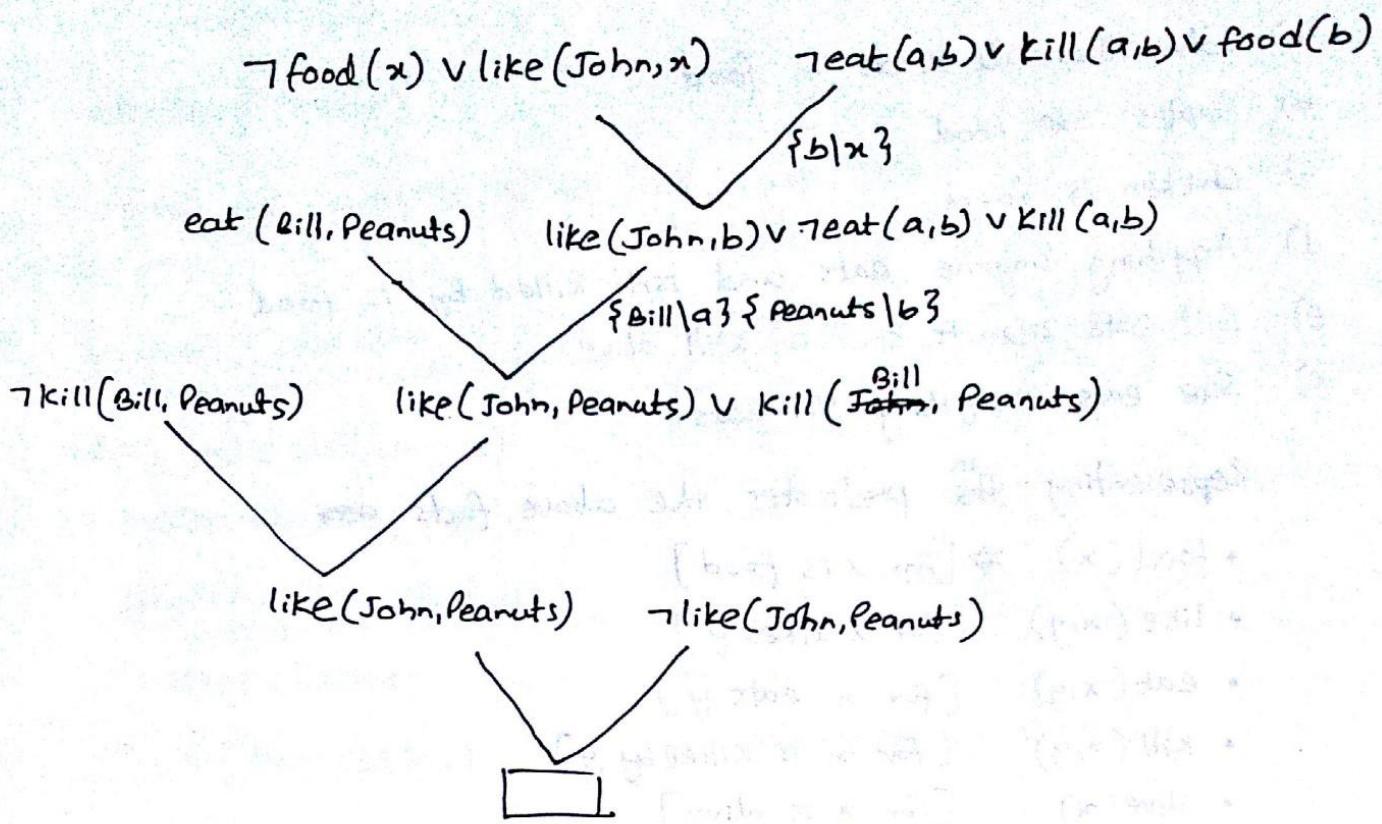
- $\text{food}(x) \Leftrightarrow [\text{for } x \text{ is food}]$
- $\text{like}(x,y) \quad [\text{for } x \text{ likes } y]$
- $\text{eat}(x,y) \quad [\text{for } x \text{ eats } y]$
- $\text{kill}(x,y) \quad [\text{for } x \text{ is killed by } y]$

The FOPC representation of the facts are:

- $\forall x \text{ food}(x) \rightarrow \text{like}(\text{John},x)$   
 $\equiv \forall x [\neg \text{food}(x) \vee \text{like}(\text{John},x)]$   
 $\equiv \neg \text{food}(x) \vee \text{like}(\text{John},x)$
- $\text{food}(\text{Apples})$
- $\text{food}(\text{chicken})$
- $\forall x \forall y (\text{eat}(x,y) \wedge \neg \text{kill}(x,y)) \rightarrow \text{food}(y)$   
 $\equiv \forall x \forall y \neg (\text{eat}(x,y) \wedge \neg \text{kill}(x,y)) \vee \text{food}(y)$   
 $\equiv \forall x \forall y (\neg \text{eat}(x,y) \vee \text{kill}(x,y) \vee \text{food}(y)) \equiv \neg \text{eat}(\text{a},\text{b}) \vee \text{kill}(\text{a},\text{b}) \vee \text{food}(\text{b})$
- $\text{eat}(\text{Bill}, \text{Peanuts}) \wedge \neg \text{kill}(\text{Bill}, \text{Peanuts})$
- $\forall x \text{ eat}(\text{Bill},x) \rightarrow \text{eat}(\text{Sue},x) \equiv \neg \text{eat}(\text{Bill},x) \vee \text{eat}(\text{Sue},x)$

The conclusion is: "John likes peanuts" which is negated  
 $\neg \text{like}(\text{John}, \text{peanuts})$

## Using Resolution Refutation:



Hence, John likes peanuts.

4. Given facts are:

- Steve only likes easy courses
- Science courses are hard
- All the courses in the basket weaving department are easy
- BK301 is a basket weaving course

The predicates involved in the above facts are

- $\text{easy}(x)$  ;  $x$  course  $x$  is easy
- $\text{like}(x, y)$  ;  $x$  likes  $y$
- $\text{Science}(x)$  ;  $x$  is a Science course
- $\text{basket}(x)$  ;  $x$  is a basket weaving course

∴ The facts in FOPC form are:

- a)  $\forall x \text{ easy}(x) \rightarrow \text{like}(\text{Steve}, x)$   
 $\equiv \forall x \neg \text{easy}(x) \vee \text{like}(\text{Steve}, x)$   
 $\equiv \neg \text{easy}(x) \vee \text{like}(\text{Steve}, x)$
- b)  $\forall x \text{ Science}(x) \rightarrow \neg \text{easy}(x)$   
 $\equiv \neg \text{Science}(x) \vee \neg \text{easy}(x)$

$$c) \forall x \text{ basket}(x) \rightarrow \text{easy}(x)$$

$$\equiv \neg \text{basket}(x) \vee \text{easy}(x)$$

$$d) \text{basket(BK301)}$$

To find out what'd Steve like; we construct a resolution table as:

### Symbolic Sentences

- i)  $\neg \text{easy}(x) \vee \text{like}(\text{steve}, x)$
- ii)  $\neg \text{basket}(x) \vee \text{easy}(x)$
- iii)  $\text{like}(\text{steve}, x) \vee \neg \text{basket}(x)$
- iv)  $\text{basket}(\text{BK301})$
- v)  $\text{like}(\text{steve}, \text{BK301})$

### Rules

- i) Given
- ii) Given
- iii) Resolution of i) & ii)
- iv) Given
- v) Resolution of iii) & iv)

∴ Steve likes BK301.

5

The steps to obtain CNF are as follows:

- i) Eliminate  $\leftrightarrow$  replacing  $\alpha \leftrightarrow \beta$  as  $(\alpha \rightarrow \beta) \wedge (\beta \rightarrow \alpha)$
- ii) Eliminate  $\rightarrow$  replacing  $\alpha \rightarrow \beta$  as  $\neg \alpha \vee \beta$
- iii) Move ' $\neg$ ' inwards as

$$\neg(\neg \alpha) \equiv \alpha$$

$$\neg(\alpha \vee \beta) \equiv \neg \alpha \wedge \neg \beta$$

$$\neg(\alpha \wedge \beta) \equiv \neg \alpha \vee \neg \beta$$

- iv) Apply distributive law: distribute  $\vee$  over  $\wedge$ , i.e.

$$\alpha \vee (\beta \wedge \gamma) \equiv (\alpha \vee \beta) \wedge (\alpha \vee \gamma)$$

The resolution rule can be applied only to disjunctions of literals, so it would seem to be relevant only to knowledge bases and queries consisting of such disjunctions. Every sentence of propositional logic is logically equivalent to a conjunction of disjunctions of literals. Hence, after conversion to CNF, the clauses can be used as input to a resolution procedure.

### Conversion to CNF:

- a)  $\neg(P \wedge Q) \vee (P \vee Q) \rightarrow R$
- $$\begin{aligned} &\equiv \neg[\neg(P \wedge Q) \wedge (\neg P \vee Q)] \vee R \\ &\equiv [(\neg P \wedge Q) \vee \neg(\neg P \vee Q)] \vee R \\ &\equiv [(\neg P \wedge Q) \vee (P \wedge \neg Q)] \vee R \\ &\equiv (((P \wedge Q) \vee \neg P) \wedge (P \wedge Q) \vee \neg Q)) \vee R \\ &\equiv ((P \vee \neg P) \wedge (Q \vee \neg P)) \wedge ((P \vee \neg Q) \wedge (Q \vee \neg Q)) \vee R \\ &\equiv ((Q \vee \neg P) \wedge (P \vee \neg Q)) \vee R \\ &\equiv (Q \vee \neg P \vee R) \wedge (P \vee \neg Q \vee R) // \end{aligned}$$
- b)  $\neg(P \vee \neg Q) \wedge (R \rightarrow S)$
- $$\begin{aligned} &\equiv \neg(P \vee \neg Q) \wedge (\neg R \vee S) \\ &\equiv (\neg P \wedge Q) \wedge (\neg R \vee S) \\ &\equiv \neg P \wedge (\neg R \vee S) \wedge Q \wedge (\neg R \vee S) \\ &\equiv (\neg P \wedge \neg R) \vee (\neg P \wedge S) \wedge (Q \wedge \neg R) \vee (Q \wedge S) \\ &\equiv ((\neg P \wedge \neg R) \vee \neg P) \wedge ((\neg P \wedge \neg R) \vee S) \wedge ((Q \wedge \neg R) \vee Q) \wedge ((Q \wedge \neg R) \vee S) \\ &\equiv (\neg P \vee \neg R) \wedge (\neg P \vee S) \wedge (\neg R \vee Q) \wedge (\neg R \vee S) \wedge (Q \vee Q) \wedge (\neg R \vee Q) \wedge (Q \vee S) \\ &\equiv \neg P \wedge (\neg R \vee \neg P) \wedge (\neg P \vee S) \wedge (\neg R \vee S) \wedge Q \wedge (\neg R \vee Q) \wedge (Q \vee S) \wedge (\neg R \vee S) // \end{aligned}$$

$$\begin{aligned}
 c) \quad & P \rightarrow ((Q \wedge R) \leftrightarrow S) \\
 \equiv & \neg P \vee [ (Q \wedge R) \leftrightarrow S ] \\
 \equiv & \neg P \vee [ ( (Q \wedge R) \rightarrow S ) \wedge ( S \rightarrow (Q \wedge R) ) ] \\
 \equiv & \neg P \vee [ \neg (Q \wedge R) \vee S \wedge \neg S \vee (Q \wedge R) ] \\
 \equiv & ( \neg P \vee \neg (Q \wedge R) \vee S ) \wedge ( \neg P \vee \neg S \vee (Q \wedge R) ) \\
 \equiv & ( \neg P \vee ( \neg Q \vee \neg R ) \vee S ) \wedge \neg P \vee ( \neg S \vee Q \wedge R ) \\
 \equiv & \neg P \vee ( \neg Q \vee \neg R \vee S ) \wedge ( \neg P \vee ( \neg S \vee Q ) \wedge \neg P \vee ( \neg S \vee R ) ) \\
 \equiv & ( \neg P \vee \neg Q \vee \neg R \vee S ) \wedge ( \neg P \vee \neg S \vee Q ) \wedge ( \neg P \vee \neg S \vee R ) \\
 = & \underline{(\neg P \vee \neg Q \vee \neg R \vee S)} \wedge \underline{(\neg P \vee \neg S \vee Q)} \wedge \underline{(\neg P \vee \neg S \vee R)} \wedge \underline{(\neg P \vee R)} \\
 \equiv & (\neg P \vee \neg Q \vee \neg R \vee S) \wedge (\neg P \vee \neg S \vee Q) \wedge (\neg P \vee \neg S \vee R) //
 \end{aligned}$$

6

## Difference between inference and reasoning:

### Inference

- i) Inference is a general term representing the derivation of new knowledge from existing knowledge and axioms (~~as~~ rules of derivation)
- ii) It is done within a single step.
- iii) It can be one of many kinds, such as induction, deduction & abduction. For eg: modus ponens.

### Reasoning

- i) Reasoning is in context of a goal, i.e. derivation of chains of conclusions that lead to the desired goal.
- ii) It is done using multiple inferences.
- iii) Resolution is a particular kind of reasoning involving the resolution rule.

### Why probabilistic reasoning in AI?

=> To make a good decision, an agent cannot simply assume what the world is like and act according to those assumptions. It must consider multiple possible contingencies and their likelihood. Agents in uncertain environments must be able to keep track of the current state of the environment. Hence, probabilistic reasoning is important in AI as it is the way of modelling dynamic situations.

For eg: When repairing a car, we ~~can~~ assume that whatever is broken, remains broken during the process of diagnosis. Now, let's consider treating a diabetic patient. As in the case of car repair, we have evidences such as recent insulin doses, food intake, blood sugar measurements and other physical signs. The task is to assess ~~a~~ the ~~current~~ ~~state~~ state of the patient. Given ~~this~~ information, the doctor makes a decision about the patient's food intake and insulin dose. Blood sugar levels and measurements thereof can change rapidly over time depending on one's recent food intake and insulin doses, one's metabolic activity, the time of day, and so on. To assess the current state from the history of evidences and to predict the outcomes of treatment actions, these dynamic changes must be modelled. This is done in the AI equivalent of computing by probabilistic reasoning.

7

Given:

In the university,

Probability of a student being a <sup>man</sup> male,  $P(m) = \frac{2}{5} = 0.4$

" " " " female, <sup>woman</sup>  $P(f) = \frac{3}{5} = 0.6$

Probability of a man being over 6 ft tall,  $P(t|m) = 4\% = 0.04$

" " woman " " ",  $P(t|f) = 1\% = 0.01$

Probability of an over 6 feet tall being a woman  $P(w|t) = ?$

$$\therefore P(w|t) = \frac{P(t|w) \times P(w)}{P(t)}$$

$$= \frac{P(t|w) \times P(w)}{P(t|w) \times P(w) + P(t|m) \times P(m)}$$

$$= \frac{0.01 \times 0.6}{0.01 \times 0.6 + 0.04 \times 0.4}$$

$$\therefore P(w|t) = 0.2727 //$$

8

Given:

In the country,

i) Probability of an adult being a male,  $P(m) = 51\% = 0.51$

" " " " female,  $P(f) = 49\% = 0.49$

ii) Probability of a male smoking cigars,  $P(c|m) = 9.5\% = 0.095$

" " " " female " ",  $P(c|f) = 1.7\% = 0.017$

Probability of an <sup>a selected</sup> adult subject being a male,  $P(m|c) = ?$

$$\therefore P(m|c) = \frac{P(c|m) \times P(m)}{P(c)}$$

$$= \frac{P(c|m) \times P(m)}{P(c|m) \times P(m) + P(c|f) \times P(f)}$$

$$= \frac{0.095 \times 0.51}{0.095 \times 0.51 + 0.017 \times 0.49}$$

$$\therefore P(m|c) = 0.8533 //$$

## 9 Selection of best attribute in a decision tree:

The scheme used in decision tree learning for selecting attributes is designed to minimize the depth of the tree. The measure should have its maximum value when the attribute is perfect and its minimum value when it is of no use at all. One suitable measure is the expected amount of information provided by the attribute, which is given by:

$$\text{Information } (A) = \sum_{i=1}^v \frac{P_i + n_i}{P+n} \times \text{Entropy} \left( \frac{P_i}{P_i + n_i}, \frac{n_i}{P_i + n_i} \right)$$

where,  $P$  is the number of positive examples,  $n$  is the negative example for the distinct values of  $A$ .

Now, the heuristic used for choosing attribute is information gain, i.e. choose attribute with the largest gain.

$$\begin{aligned}\text{Gain}(S, A) &= \text{Entropy}(S) - I(S, A) \\ &= E(S) - \frac{|S|}{|S|} \times E(S_i)\end{aligned}$$

where,  $E(S)$  is the overall entropy.

$$\begin{aligned}E(S) &= \frac{-P}{P+n} \log_2 \left( \frac{P}{P+n} \right) - \frac{n}{P+n} \log_2 \left( \frac{n}{P+n} \right) \\ &= -P^+ \log_2 P^+ - P^- \log_2 P^-\end{aligned}$$

where,

$P^+$  and  $P^-$  are the positive and negative probabilities for the target variable.

~~Given~~ Given Sample Data is

<u>Outlook</u>	<u>Temperature</u>	<u>Humidity</u>	<u>Windy</u>	<u>Play golf (target variable)</u>
Rainy	Hot	High	False	No
Rainy	Hot	High	True	No
Overcast	Hot	High	False	Yes
Sunny	Mild	High	False	Yes
Sunny	Cool	Normal	False	Yes
Sunny	Cool	Normal	True	No
overcast	Cool	Normal	True	Yes
Rainy	Mild	High	False	No
Rainy	Cool	Normal	False	Yes
Sunny	Mild	Normal	False	Yes
Rainy	Mild	Normal	True	Yes
Overcast	Mild	High	True	Yes
Overcast	Hot	Normal	False	Yes
Sunny	Mild	High	True	No

For the root attribute of decision tree from the above data:

a) ~~Starting with Outlook~~:

We have:

Play Golf : Yes: 9 , No: 5

Attributes:

Temperature : Hot : Yes: 2 , No: 2

Cold : Yes: 3 , No: 1

Mild : Yes: 4 , No: 2

Humidity : High : Yes: 3 , No: 4

Normal : Yes: 6 , No: 1

Outlook : Rainy : Yes: 2 , No: 3

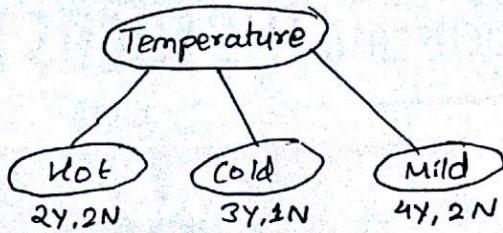
Overcast: Yes: 4 , No: 0

Sunny : Yes: 3 , No: 2

Windy : False: Yes: 6 , No: 2

True : Yes: 3 , No: 3

a) For Temperature:



$$\text{Overall Entropy before splitting} = \frac{-9}{14} \log_2 \left( \frac{9}{14} \right)$$

~~(S)~~ E(S)

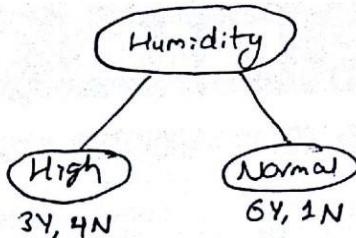
$$= \frac{-5}{14} \log_2 \left( \frac{5}{14} \right)$$

$$= 0.9403$$

$$\begin{aligned} \text{Entropy after splitting, } E_S(T) &= \frac{4}{14} \underbrace{\left( -p^+ \log_2 p^+ - p^- \log_2 p^- \right)}_{\text{Hot}} + \frac{5}{14} \underbrace{\left( -p^+ \log_2 p^+ - p^- \log_2 p^- \right)}_{\text{Cold}} \\ &\quad + \frac{6}{14} \underbrace{\left( -p^+ \log_2 p^+ - p^- \log_2 p^- \right)}_{\text{Mild}} \\ &= \frac{4}{14} \left( -\frac{2}{4} \log_2 \frac{2}{4} - \frac{2}{4} \log_2 \frac{2}{4} \right) + \frac{5}{14} \left( -\frac{3}{4} \log_2 \frac{3}{4} - \frac{1}{4} \log_2 \frac{1}{4} \right) \\ &\quad + \frac{6}{14} \left( -\frac{4}{6} \log_2 \frac{4}{6} - \frac{2}{6} \log_2 \frac{2}{6} \right) \\ &= \frac{4}{14} \times 1 + \frac{5}{14} \times 0.8113 + \frac{6}{14} \times 0.9183 \\ &= 0.911 \end{aligned}$$

$$\text{Information Gain, } G(S, \text{Temperature}) = \frac{0.9403 - 0.911}{\text{before split - after split}} = 0.0293.$$

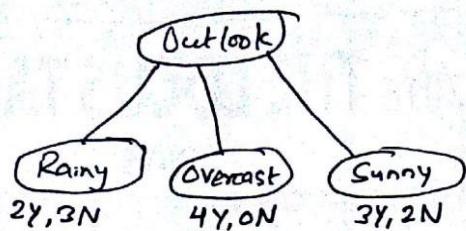
b) For Humidity:



$$\begin{aligned} \text{Entropy after splitting } E_S(H) &= \frac{7}{14} \left( -\frac{3}{7} \log_2 \frac{3}{7} - \frac{4}{7} \log_2 \frac{4}{7} \right) + \frac{7}{14} \left( -\frac{6}{7} \log_2 \frac{6}{7} - \frac{1}{7} \log_2 \frac{1}{7} \right) \\ &= \frac{7}{14} \times 0.9852 + \frac{7}{14} \times 0.5917 \\ &= 0.7885 \end{aligned}$$

$$\begin{aligned} \therefore \text{Information Gain, } G(S, H) &= 0.9403 - 0.7885 = 0.15185. \end{aligned}$$

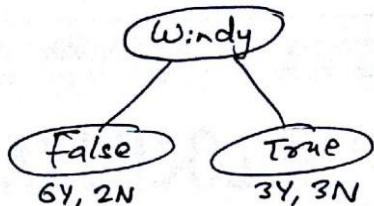
c) For Outlook:



$$\begin{aligned}
 \text{Entropy after splitting, } E_S(O) &= \frac{5}{14} \left( -\frac{2}{5} \log_2 \frac{2}{5} - \frac{3}{5} \log_2 \frac{3}{5} \right) + \frac{4}{14} \left( -\frac{4}{4} \log_2 \frac{4}{4} - \frac{0}{4} \log_2 \frac{0}{4} \right) \\
 &\quad + \frac{5}{14} \left( -\frac{3}{5} \log_2 \frac{3}{5} - \frac{2}{5} \log_2 \frac{2}{5} \right) \\
 &= \frac{5}{14} \times 0.971 + \frac{4}{14} \times 0 + \frac{5}{14} \times 0.971 \\
 &= 0.6935
 \end{aligned}$$

$$\text{Information Gain, } G_I(S, O) = 0.9403 - 0.6935 = 0.2467$$

d) For Windy:



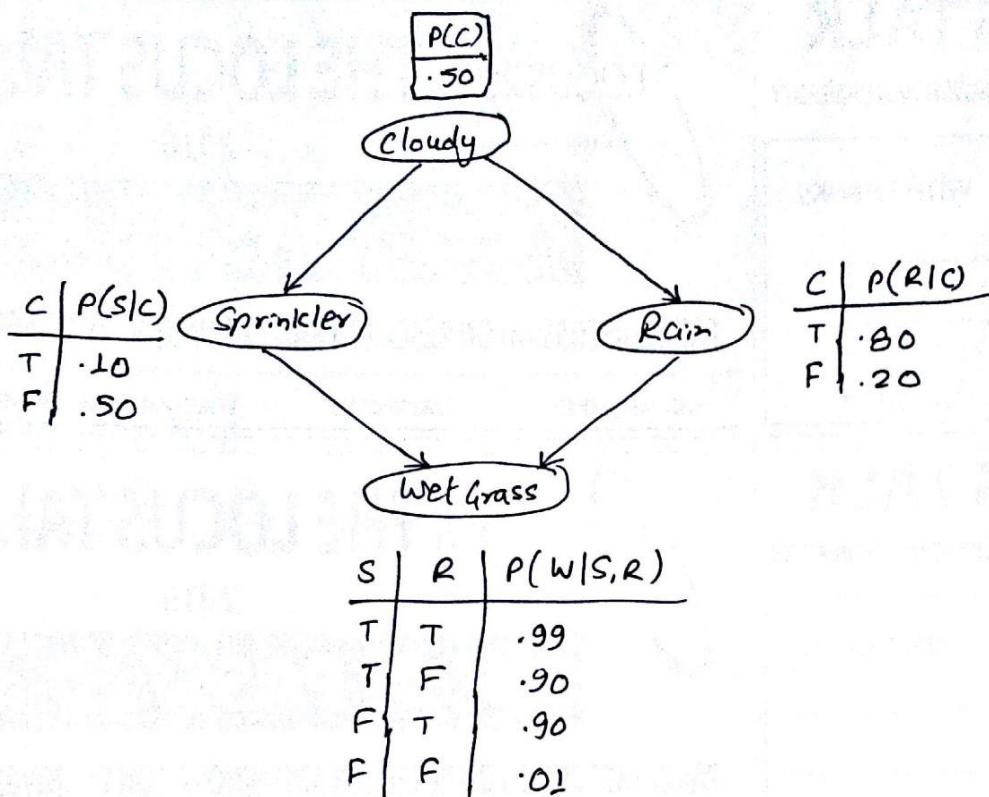
$$\begin{aligned}
 \text{Entropy after splitting } E_S(W) &= \frac{8}{14} \left( -\frac{6}{8} \log_2 \frac{6}{8} - \frac{2}{8} \log_2 \frac{2}{8} \right) + \frac{6}{14} \left( -\frac{3}{6} \log_2 \frac{3}{6} - \frac{3}{6} \log_2 \frac{3}{6} \right) \\
 &= \frac{8}{14} \times 0.8113 + \frac{6}{14} \times 1 \\
 &= 0.8922
 \end{aligned}$$

$$\begin{aligned}
 \text{Information Gain, } G_I(S, W) &= 0.9403 - 0.8922 \\
 &= 0.04813.
 \end{aligned}$$

Since the highest information gain is for outlook, so, the root node is outlook.

11

Given network is



We are required to find the

No. of samples (in percentage) of type  $C = T, S = T, R = F, W = T$  using direct sampling method.

$$\text{i.e. } P(C \wedge S \wedge R \wedge W)$$

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

$$= 0.50 \times 0.10 \times 0.20 \times 0.90$$

$$= 0.9 \times 10^{-3}$$

$$= 0.9\% //$$