

СОФИЙСКИ УНИВЕРСИТЕТ
"СВ. КЛИМЕНТ ОХРИДСКИ"
ФАКУЛТЕТ ПО МАТЕМАТИКА И ИНФОРМАТИКА
МАШИННО САМООБУЧЕНИЕ

спец. Изкуствен интелект, I курс, зимен семестър
учебна година 2024/2025

Изготвил:

Кристиян Симов
фак. номер 4MI3400288

Дата:

25. 10. 2024 г.
София

Домашна работа №2



Съдържание

1	Решение на задача №1	2
2	Решение на задача №2	4
3	Решение на задача №3	9
4	Решение на задача №4	12

1 Решение на задача №1

Нека имаме множество от обучаващи примери S дефинирано чрез таблицата:

Пример	Класификация	A_1	A_2
1	+	T	T
2	+	T	T
3	-	T	F
4	+	F	F
5	-	F	T
6	-	F	T

а) Формулата за изчисление на ентропия от информационната теория за произволно множество от примери S с булеви стойности на целевата функция (+ или -) , показваща неговата еднородност, е:

$$Entropy(S) \equiv -p_+ \log_2 p_+ - p_- \log_2 p_-,$$

където p_+ и p_- са съответно отношенията на броя на положителните и отрицателните примери към броят всички примери.

Прилагаме я към конкретното множество S и последователно получаваме:

$$\begin{aligned} Entropy([3_+, 3_-]) &= -\frac{3}{6} \log_2 \frac{3}{6} - \frac{3}{6} \log_2 \frac{3}{6} = -\frac{1}{2} \log_2 \frac{1}{2} - \frac{1}{2} \log_2 \frac{1}{2} = \\ &= -\frac{1}{2}(-1) - \frac{1}{2}(-1) = \frac{1}{2} + \frac{1}{2} = 1 \end{aligned}$$

Очаквано, получихме ентропия равна на 1, тъй като броят на положителните и отрицателните примери е еднакъв (в случая равен на 3).

б) Формулата за изчисление информационната печалба на атрибут A по отношение на произволно множество от примери S е:

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

където $Values(A)$ е множеството от възможни стойности на атрибута A , а множеството $S_v = \{s \in S | A(s) = v\}$.

Прилагаме я към атрибута A_2 по отношение на конкретното множество S и последователно получаваме:

$$\begin{aligned} Gain(S, A_2) &= Entropy(S) - \sum_{v \in Values(A_2)} \frac{|S_v|}{|S|} Entropy(S_v) = \\ &= 1 - \sum_{v \in \{T, F\}} \frac{|S_v|}{6} Entropy(S_v) = \\ &= 1 - \frac{|S_T|}{6} Entropy(S_T) - \frac{|S_F|}{6} Entropy(S_F) = 1 - \left(\frac{4}{6}\right)1 - \left(\frac{2}{6}\right)1 = 0 \end{aligned}$$

Очаквано, получихме печалба равна на 0, тъй като броят на положителните и отрицателните примери е еднакъв в подмножествата S_T и S_F .

2 Решение на задача №2

а) Нека имаме множество от обучаващи примери S дефинирано чрез таблицата:

Пример	Небе	Въздух	Влажност	Вятър	Вода	Прогноза	Харесва
1	Слънце	Топъл	Нормална	Силен	Топла	Същото	Да
2	Слънце	Топъл	Висока	Силен	Топла	Същото	Да
3	Дъжд	Студен	Висока	Силен	Топла	Промяна	Не
4	Слънце	Топъл	Висока	Силен	Студена	Промяна	Да

Тогава алгоритъмът ID3 ще премине през следните стъпки:

$$0) \quad S = \{x_1, x_2, x_3, x_4\}, A = \{A_{\text{Небе}}, A_{\text{Въздух}}, A_{\text{Влажност}}, A_{\text{Вятър}}, A_{\text{Вода}}, A_{\text{Прогноза}}\}$$

$$Entropy(S) = -p_+ \log_2 p_+ - p_- \log_2 p_- = -\frac{3}{4} \log_2 \frac{3}{4} - \frac{1}{4} \log_2 \frac{1}{4} \approx 0.811$$

$$Gain(A_{\text{Небе}}, S) \approx 0.811 - (0 + 0) = 0.811 \leftarrow best$$

$$Gain(A_{\text{Въздух}}, S) = Gain(A_{\text{Небе}}, S) \approx 0.811$$

$$Gain(A_{\text{Влажност}}, S) \approx 0.811 - (0 + \frac{3}{4} * 0.918) \approx 0.811 - 0.689 = 0.122$$

$$Gain(A_{\text{Вятър}}, S) \approx 0.811 - 0.811 = 0$$

$$Gain(A_{\text{Вода}}, S) = A_{\text{Влажност}}, S) \approx 0.122$$

$$Gain(A_{\text{Прогноза}}, S) \approx 0.811 - (0 + \frac{2}{4} * 1) = 0.811 - 0.5 = 0.311$$

$$1) \quad S_{\text{Слънце}} = \{x_1, x_2, x_4\}, A = \{A_{\text{Небе}}, A_{\text{Въздух}}, A_{\text{Влажност}}, A_{\text{Вятър}}, A_{\text{Вода}}, A_{\text{Прогноза}}\}$$

$$Entropy(S_{\text{Слънце}}) = -p_+ \log_2 p_+ - p_- \log_2 p_- = -\frac{3}{3}0 - \frac{0}{3}1 = 0 - 0 = 0$$

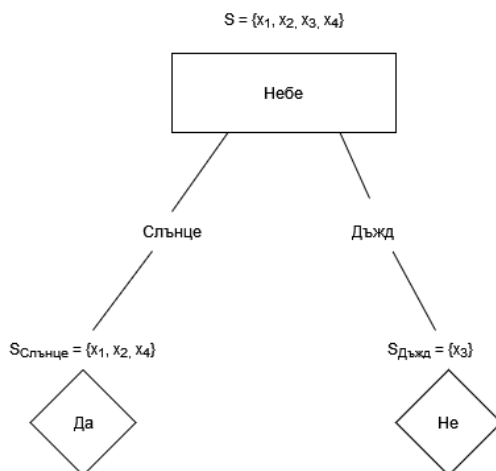
Множеството $S_{\text{Слънце}}$ е напълно еднородно - образуваме листо със знак "Да".

$$2) \quad S_{\text{Дъжд}} = \{x_3\}, A = \{A_{\text{Небе}}, A_{\text{Въздух}}, A_{\text{Влажност}}, A_{\text{Вятър}}, A_{\text{Вода}}, A_{\text{Прогноза}}\}$$

$$Entropy(S_{\text{Облаци}}) = -p_+ \log_2 p_+ - p_- \log_2 p_- = -\frac{0}{1}1 - \frac{1}{1}0 = 0 - 0 = 0$$

Множеството $S_{\text{Дъжд}}$ е напълно еднородно - образуваме листо със знак "Не".

Край - дървото е обучено и изглежда така:



Фигура 1: На изображението виждаме, че още след първото най-добро разделяне дървото е обучено успешно.

б) Нека към предходната таблица за S прибавим още един обучаващ пример:

Пример	Небе	Въздух	Влажност	Вятър	Вода	Прогноза	Харесва
1	Слънце	Топъл	Нормална	Силен	Топла	Същото	Да
2	Слънце	Топъл	Висока	Силен	Топла	Същото	Да
3	Дъжд	Студен	Висока	Силен	Топла	Промяна	Не
4	Слънце	Топъл	Висока	Силен	Студена	Промяна	Да
5	Слънце	Топъл	Нормална	Слаб	Топла	Същото	Не

Тогава алгоритъмът ID3 ще премине през следните стъпки:

$$0) \quad S = \{x_1, x_2, x_3, x_4, x_5\}, A = \{A_{\text{Небе}}, A_{\text{Въздух}}, A_{\text{Влажност}}, A_{\text{Вятър}}, A_{\text{Вода}}, A_{\text{Прогноза}}\}$$

$$Entropy(S) = -p_+ \log_2 p_+ - p_- \log_2 p_- = -\frac{3}{5} \log_2 \frac{3}{5} - \frac{2}{5} \log_2 \frac{2}{5} \approx 0.970$$

$$Gain(A_{\text{Небе}}, S) \approx 0.970 - (0 + \frac{4}{5} * 0.811) = 0.970 - 0.608 = 0.362 \leftarrow best$$

$$Gain(A_{\text{Въздух}}, S) = Gain(A_{\text{Небе}}, S) \approx 0.362$$

$$Gain(A_{\text{Влажност}}, S) \approx 0.970 - (\frac{2}{5} * 1 + \frac{3}{5} * 0.918) \approx 0.970 - 0.951 = 0.019$$

$$Gain(A_{\text{Вятър}}, S) \approx 0.970 - (\frac{1}{5} * 0 + \frac{4}{5} * 0.811) = 0.970 - 0.649 = 0.321$$

$$Gain(A_{\text{Вода}}, S) \approx 0.970 - (\frac{1}{5} * 0 + \frac{4}{5} * 1) = 0.970 - 0.8 = 0.170$$

$$Gain(A_{\text{Прогноза}}, S) = Gain(A_{\text{Влажност}}, S) \approx 0.019$$

$$1) \quad S_{\text{Слънце}} = \{x_1, x_2, x_4, x_5\}, A = \{A_{\text{Небе}}, A_{\text{Въздух}}, A_{\text{Влажност}}, A_{\text{Вятър}}, A_{\text{Вода}}, A_{\text{Прогноза}}\}$$

$$Entropy(S_{\text{Слънце}}) = -\frac{3}{4} \log_2 \frac{3}{4} - \frac{1}{4} \log_2 \frac{1}{4} \approx 0.811$$

$$Gain(A_{\text{Небе}}, S) \approx 0.811 - 0.811 = 0$$

$$Gain(A_{\text{Въздух}}, S) = Gain(A_{\text{Небе}}, S) = 0$$

$$Gain(A_{\text{Влажност}}, S) \approx 0.811 - (\frac{2}{4} * 0 + \frac{2}{4} * 1) \approx 0.811 - 0.5 = 0.311$$

$$Gain(A_{\text{Вятър}}, S) \approx 0.811 - (\frac{1}{4} * 1 + \frac{3}{4} * 0) = 0.811 - 0.25 = 0.561 \leftarrow best$$

$$Gain(A_{\text{Вода}}, S) \approx 0.811 - (\frac{1}{4} * 0 + \frac{3}{4} * 0.918) = 0.970 - 0.689 = 0.122$$

$$Gain(A_{\text{Прогноза}}, S) = Gain(A_{\text{Вода}}, S) \approx 0.122$$

$$2) \quad S_{\text{Силен}} = \{x_1, x_2, x_4\}, A = \{A_{\text{Небе}}, A_{\text{Въздух}}, A_{\text{Влажност}}, A_{\text{Вятър}}, A_{\text{Вода}}, A_{\text{Прогноза}}\}$$

$$Entropy(S_{\text{Силен}}) = -p_+ \log_2 p_+ - p_- \log_2 p_- = -\frac{3}{3}0 - \frac{0}{3}1 = 0 - 0 = 0$$

Множеството $S_{\text{Силен}}$ е напълно еднородно - образуваме листо със знак "Да".

$$3) \quad S_{\text{Слаб}} = \{x_5\}, A = \{A_{\text{Небе}}, A_{\text{Въздух}}, A_{\text{Влажност}}, A_{\text{Вятър}}, A_{\text{Вода}}, A_{\text{Прогноза}}\}$$

$$Entropy(S_{\text{Слаб}}) = -p_+ \log_2 p_+ - p_- \log_2 p_- = -\frac{0}{1}1 - \frac{1}{1}0 = 0 - 0 = 0$$

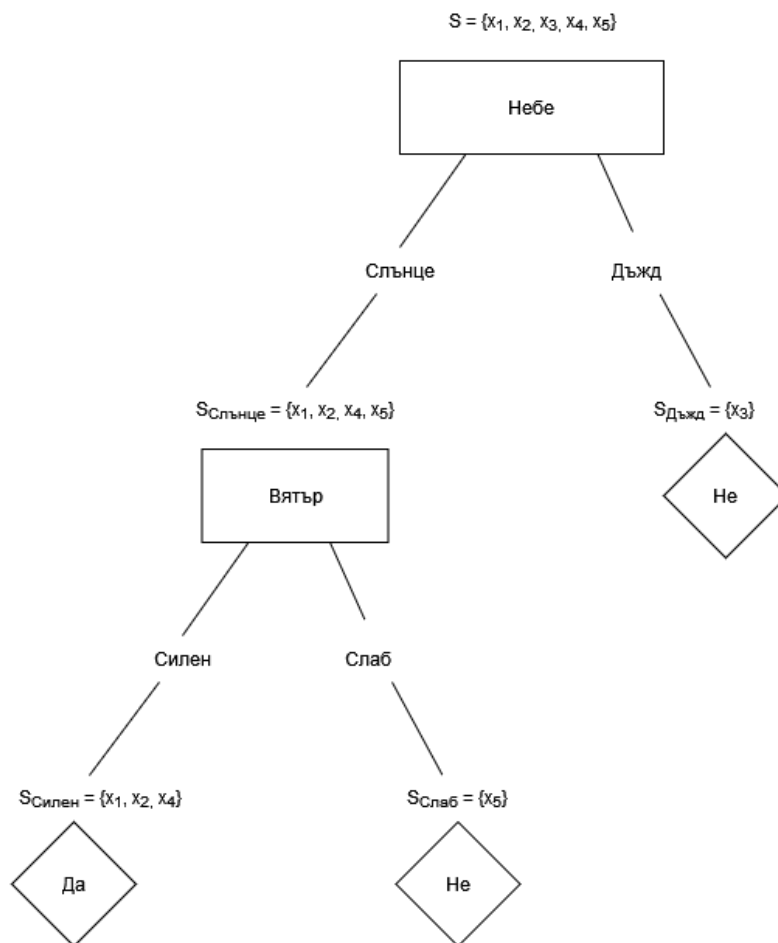
Множеството $S_{\text{Слаб}}$ е напълно еднородно - образуваме листо със знак "Не".

$$4) \quad S_{\text{Дъжд}} = \{x_3\}, A = \{A_{\text{Небе}}, A_{\text{Въздух}}, A_{\text{Влажност}}, A_{\text{Вятър}}, A_{\text{Вода}}, A_{\text{Прогноза}}\}$$

$$Entropy(S_{\text{Дъжд}}) = -p_+ \log_2 p_+ - p_- \log_2 p_- = -\frac{0}{1}1 - \frac{1}{1}0 = 0 - 0 = 0$$

Множеството $S_{\text{Дъжд}}$ е напълно еднородно - образуваме листо със знак "Не".

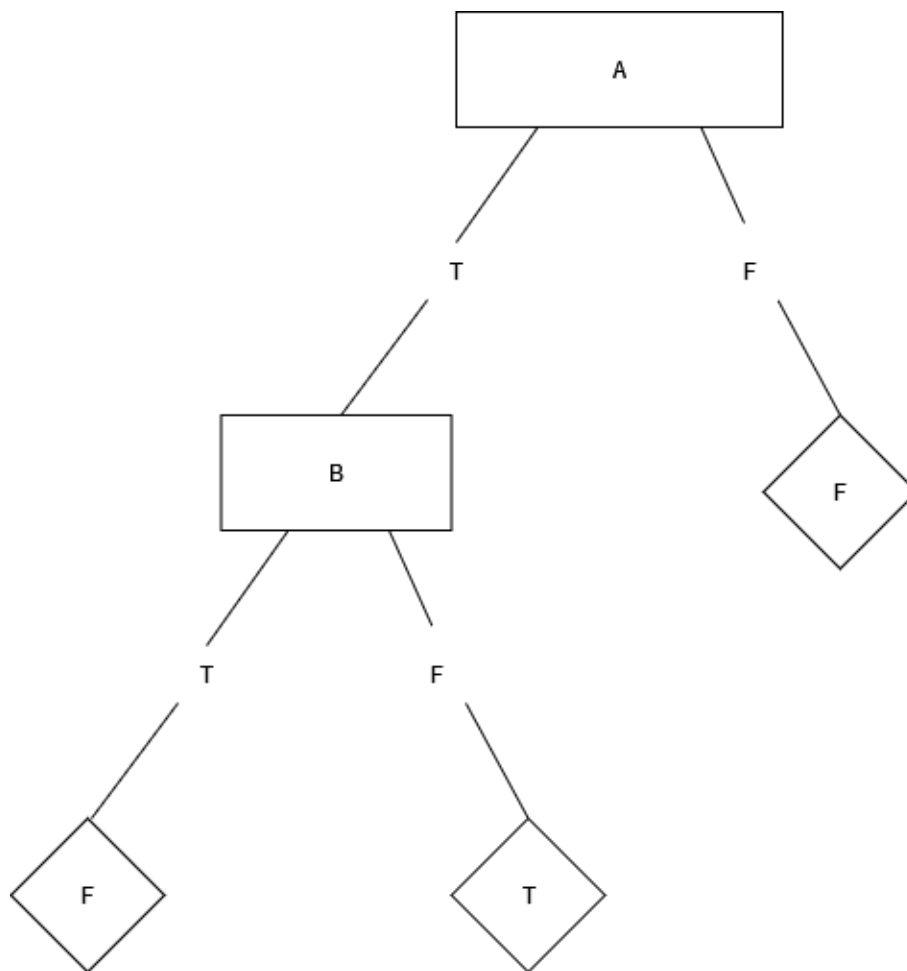
Край - дървото е обучено и изглежда така:



Фигура 2: На изображението виждаме, че вече след първото най-добро разделяне се налага да изберем още едно такова за множеството в лявото поддърво, след което дървото е обучено успешно.

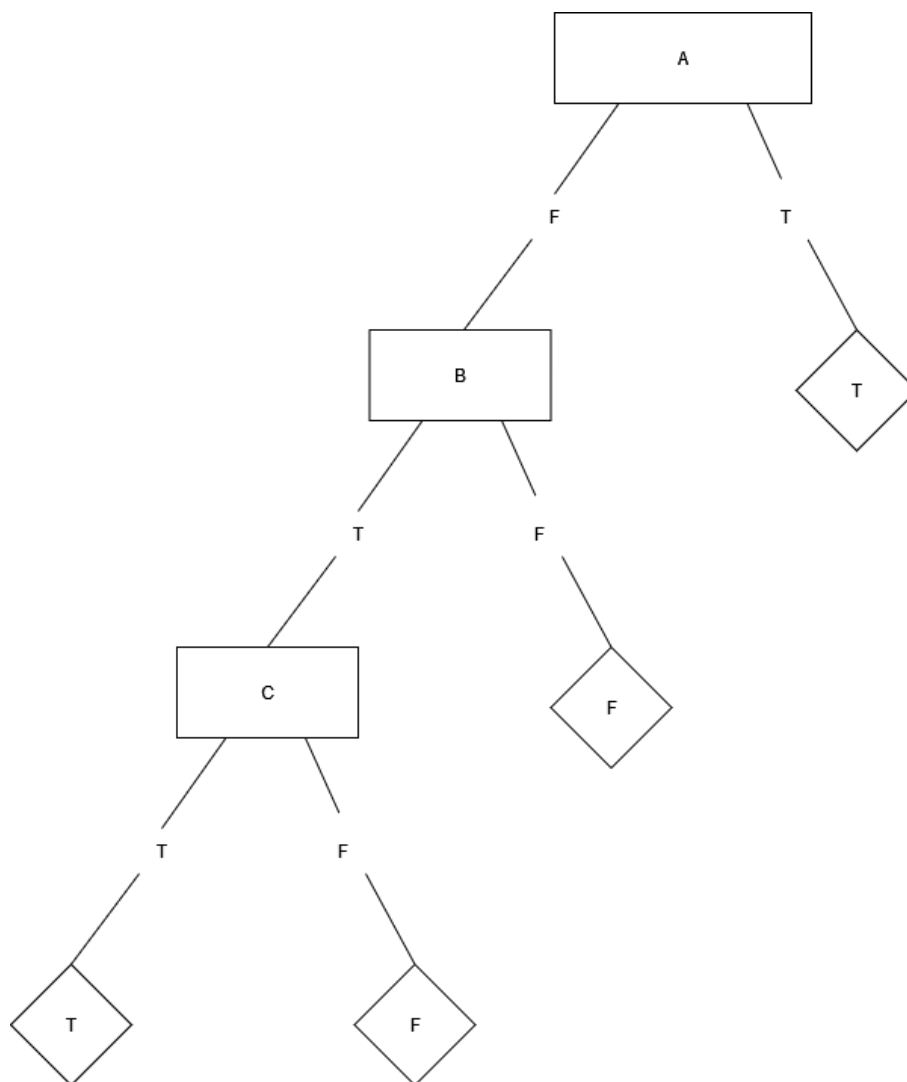
3 Решение на задача №3

a) $A \wedge \neg B$



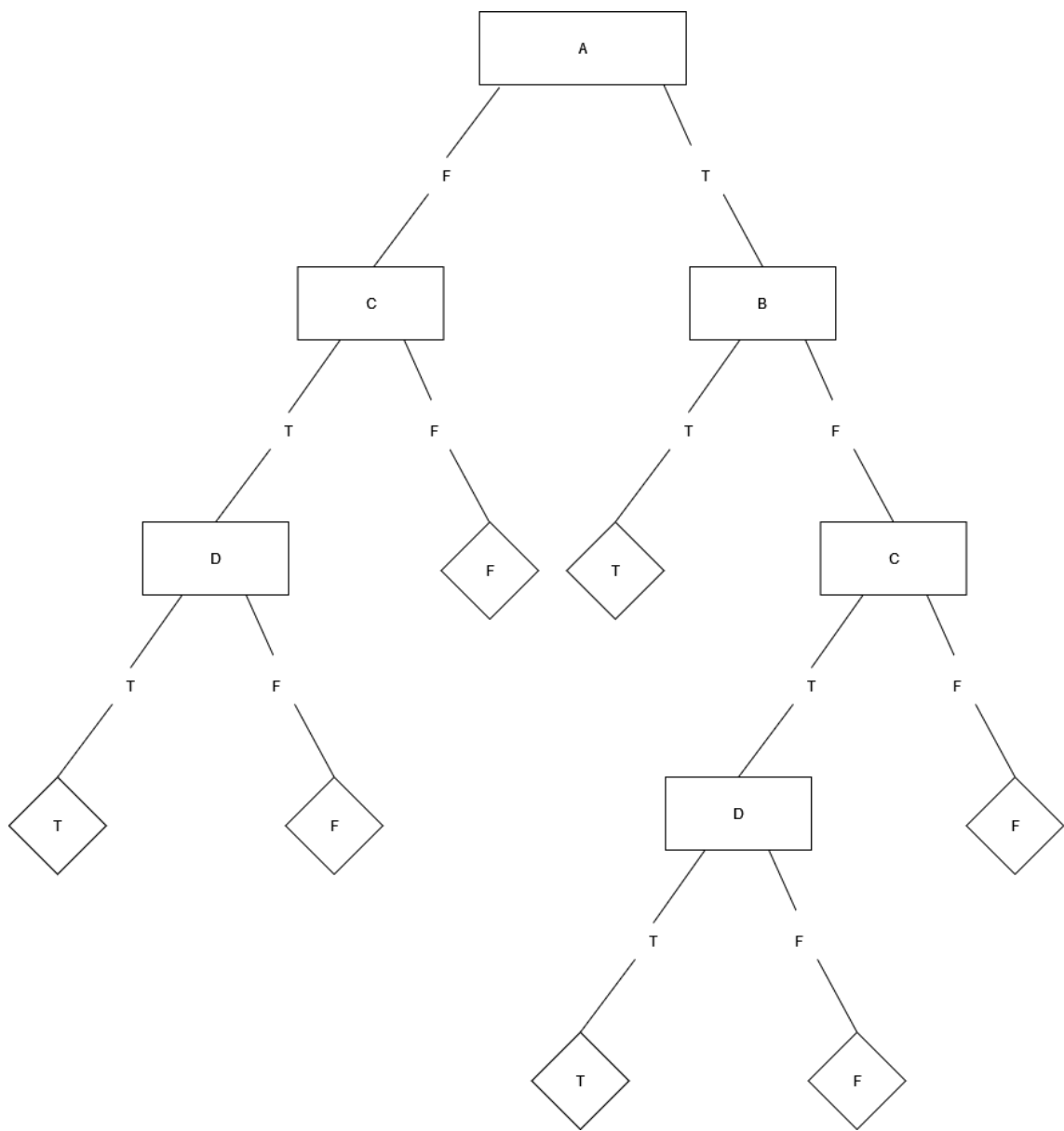
Фигура 3: а)

b) $A \vee (B \wedge C)$



Фигура 4: b)

с) $(A \wedge B) \vee (C \wedge D)$



Фигура 5: с)

4 Решение на задача №4

Нека D1 и D2 са класификационни дървета описващи булеви функции (като тези от Задача №3), такива че D2 е получено от D1 чрез заместване на листо (термален възел) в D1 с цяло поддърво T' .

Ще покажем, че твърдението:

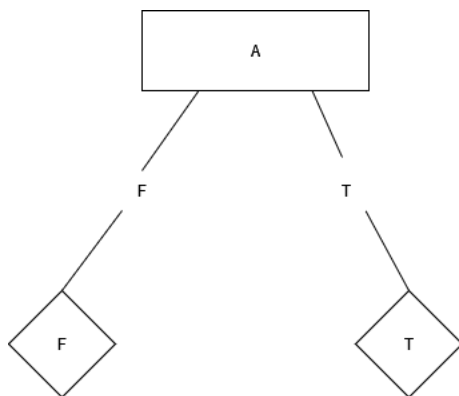
D1 е по-общо-или-равно-на D2

е невинаги в сила.

1) Нека за простота D1 се описва чрез израза:

A

Тогава D1 ще изглежда така:

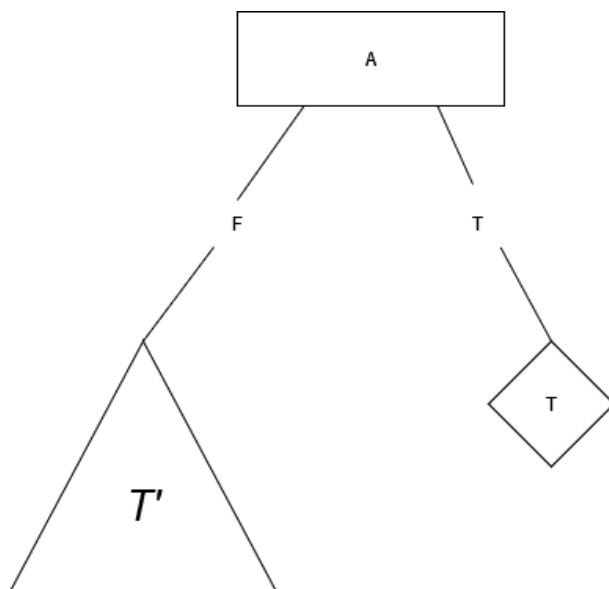


Фигура 6: D1

1) Нека получим израз за D2 от този на D1 чрез добавяне на непразния (състоящ се поне от променлива B) израз на поддърво T' посредством дизюнкция:

$A \vee T'$

Тогава D2 най-общо ще изглежда така:

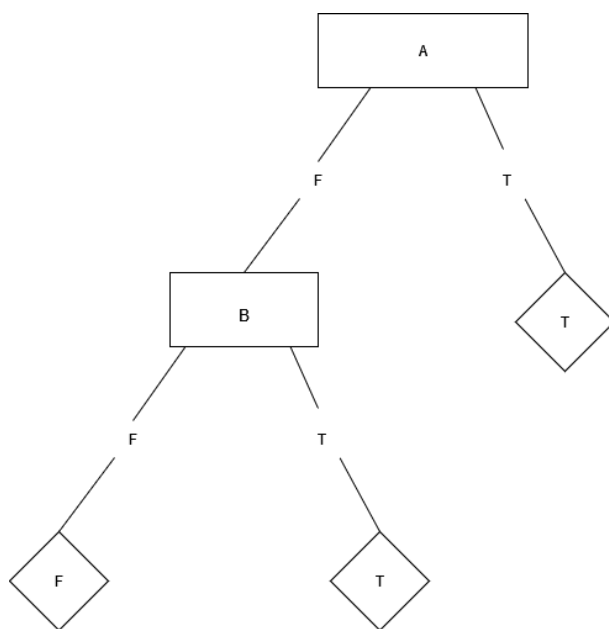


Фигура 7: D2

Забелязваме, че вече оценката на пример с $A = F$ не е директно F, а вече зависи от резултата от минаването по поддървото T' . Ако резултатът от това минаване е T, тогава общия резултат ще е T противно на резултатът F, получен при $A = F$ в D1. Това би означавало противоречие с твърдението, тъй като примерът $A = F \wedge T' = T$ се покрива от хипотезата D2, но не от хипотезата D1. Нека за простота изразът описващ поддървото T' е равен на B. Тоест изразът за D2 става:

$$A \vee B$$

Тогава D2 ще изглежда по следния начин:



Фигура 8: D2

Нека разгледаме примерът $x \equiv A = F \wedge B = T$. Той се покрива от хипотезата D2 ($D2(x) = T$), но не се покрива от хипотезата D1 ($D1(x) = F$). D1 го "изпуска". Достигнахме до противоречие с твърдението *D1 е по-общо-или-равно-на D2*. \square