

שאלה 6: MDP

0	0	100
S	0	-50

נתון ה MDP בציור: agent נמצא ב S. אם יגיע לאחד ממשבצות שבשורה הימנית המשחק יסתיים. עבור כל נסיון של מהלך ממשבצת למשבצת סמוכה העלות היא 5. ההנחה היא שעבור כל תזוזה יש 75% סיכוי להצליח ו 10% סיכוי להזרק לשני שני הצדדים ו 5% סיכוי להזרק אחורה. אם ה agent נתקע בקיר, הוא נשאר באותה המשבצת.

עליכם לבצע שתי איטרציות של Value Iteration.

יש לרשום את ה utility אחרי האיטרציה הראשונה של כל משבצת ריקה

יש לרשום את ה utility אחרי האיטרציה השנייה של כל משבצת ריקה וכן לרשום מהי ה policy האופטימאלי (בצורת חץ) לאחר שתי האיטרציות בכל משבצת ריקה.

שאלה 6: MDP

0	0	100
S	0	-50

נתון ה MDP בציור: agent נמצא ב S. אם יגיע לאחד ממשבצות שבשורה הימנית המשחק יסתיים. עבור כל נסיון של מהלך ממשבצת למשבצת סמוכה העלות היא 5. ההנחה היא שעבור כל תזוזה יש 75% סיכוי להצליח ו 10% סיכוי להזרק לשני שני הצדדים ו 5% סיכוי להזרק אחורה. אם ה agent נתקע בקיר, הוא נשאר באותה המשבצת.

עליכם לבצע שתי איטרציות של Value Iteration.

יש לרשום את ה utility אחרי האיטרציה הראשונה של כל משבצת ריקה

-5	70	100
-5	-7.5	-50

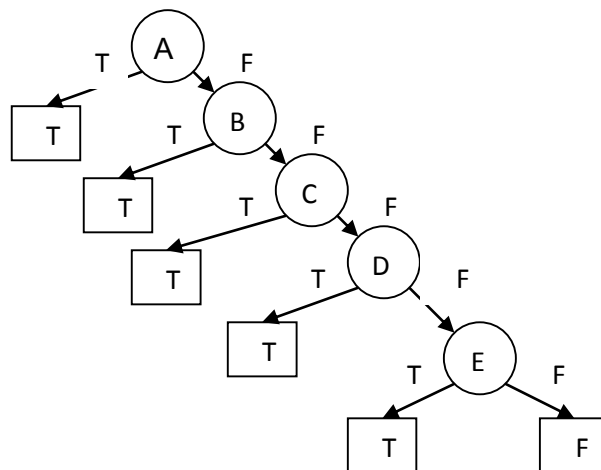
יש לרשום את ה utility אחרי האיטרציה השנייה של כל משבצת ריקה וכן לרשום מהי ה policy האופטימאלי (בצורת חץ) לאחר שתי האיטרציות בכל משבצת ריקה.

46.25	76	100
- 10.125	41.625	-50

עץ החלטה

א. צייר עץ החלטה (decision tree) המתאר את הנוסחה הבוליאנית הבאה: (6)

$$A \vee B \vee C \vee D \vee E$$



נתונה טבלת האימון הבאה:

A	B	C	Target
T	T	F	F
F	T	F	T
F	T	T	T
F	F	T	F
T	F	F	F

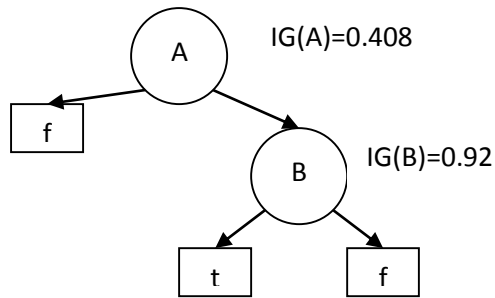
א. צייר עץ החלטה לקבוצת הדוגמאות. אם קיימות מספר תכונות עם אותו gain יש לבחור לפי סדר הא"ב. יש לציין בקודקודי העץ את ערכי ה-gain הרלוונטיים.

ב. מהי הפונקציה הבוליאנית המיוצגת בעץ שמצאת
תשובה:

א.

$$IG(A) = I(2/5, 3/5) - (2/5 * I(0/2, 2/2) + 3/5 * I(2/3, 1/3)) = 0.96 - (2/5 * 0 + 3/5 * 0.92) = 0.96 - 0.552 = 0.408$$

$$IG(B) = I(2/3, 1/3) - (2/3 * I(2/2, 0/2) + 1/3 * I(0/1, 1/1)) = 0.92$$



ב. $(-A^B)$

שאלה 2: עצי החלטה

הנח כי נתונות לך הדוגמאות הבאות. $F1, F2, F3$ יכולים לקבל את הערכים a, b, c, d.

	<u>$F1$</u>	<u>$F2$</u>	<u>$F3$</u>	<u>Output</u>
<i>ex1</i>	<i>b</i>	<i>d</i>	<i>c</i>	+
<i>ex2</i>	<i>c</i>	<i>d</i>	<i>b</i>	-
<i>ex3</i>	<i>c</i>	<i>a</i>	<i>c</i>	+
<i>ex4</i>	<i>b</i>	<i>a</i>	<i>b</i>	-
<i>ex5</i>	<i>a</i>	<i>b</i>	<i>a</i>	+
<i>ex6</i>	<i>a</i>	<i>b</i>	<i>a</i>	-
<i>ex7</i>	<i>d</i>	<i>c</i>	<i>c</i>	+

א) מה ה-Gain ש-information gain יחשב לכל אחד מה-features $F1, F2$ ו- $F3$?

מה ה-feature שייבחר?

תשובה:

נחשב את ה-Gain הכללי - $I(4/7, 3/7) = 0.9852$

$$IG(F1) = I(4/7, 3/7) - (2/7 * I(1/2, 1/2) + 2/7 * I(1/2, 1/2) + 2/7 * I(1/2, 1/2) + 1/7 * I(1, 0)) = 0.9852 - 6/7 = 0.1280$$

$$IG(F2) = I(4/7, 3/7) - (2/7 * I(1/2, 1/2) + 2/7 * I(1/2, 1/2) + 1/7 * I(1, 0) + 2/7 * I(1/2, 1/2)) = 0.1280$$

$$IG(F3) = I(4/7, 3/7) - (2/7 * I(1/2, 1/2) + 2/7 * I(0, 1) + 3/7 * I(1, 0) + 0) = 0.6994$$

ולכן F3 ייבחר.

(ב) בלי קשר לתשובתך ב- א), הנח כי F1 נבחר להיות בקודקוד השורש.

הראה כיצד יראו הקריאות הרקורסיביות של אלגוריתם ה- Decision Tree [רמה אחת]. הראה מה הפרמטרים שישלחו (אבל אל תבצע את החישובים).

תשובה:

עבור $v1=a$ $DTL(\{aba+, aba-\}, \{F2, F3\}, +)$

עבור $v2=b$ $DTL(\{bdc+, bab-\}, \{F2, F3\}, +)$

עבור $v3=c$ $DTL(\{cdb-, cac+\}, \{F2, F3\}, +)$

עבור $v4=d$ $DTL(\{dcc+\}, \{F2, F3\}, +)$

(ג) בלי קשר לתשובתך ב- א) ו ב), הנח כי F3 נבחר להיות בקודקוד השורש.

הראה כיצד יראו הקריאות הרקורסיביות של אלגוריתם ה- Decision Tree [רמה אחת]. הראה מה הפרמטרים שישלחו (אבל אל תבצע את החישובים).

תשובה:

עבור $v1=a$ $DTL(\{aba+, aba-\}, \{F1, F2\}, +)$

עבור $v2=b$ $DTL(\{cdb-, bab-\}, \{F1, F2\}, +)$

עבור $v3=c$ $DTL(\{bdc+, cac+, dcc+\}, \{F1, F2\}, +)$

עבור $v4=d$ $DTL(\{\}, \{F1, F2\}, +)$

נתונה טבלה ובה נתוני אימון. מוגדרים שלושה מאפיינים F1, F2 ו-F3. לכל מאפיין ישנם שלושה ערכים אפשריים: a, b או c.

- א. עליכם לבנות עץ החלטה המייצג את מאפיין הפלט output על סמך האלגוריתם שגלמם בשיעור. אם קיימות מספר תכונות עם אותו f1 gain עדיף על f2 ו-f2 עדיף על f3. ערך ה-default כאשר אין דוגמא הוא false.
- ב. ציינו בקודקודי העץ הרלוונטיים את ה-gain של התכונה שבחרתם.
- ג. כתבו את הנוסחה הלוגית המתארת את העץ שקבלתם.

טבלת הנתונים:

#	f1	f2	f3	output
1	b	b	a	t
2	a	a	a	f
3	b	a	a	f
4	b	c	a	t
5	c	b	a	f
6	a	c	b	f
7	b	a	c	t
8	a	a	c	f
9	c	b	c	f

$$I(\frac{p}{p+n}, \frac{n}{p+n}) = -\frac{p}{p+n} \log_2 \frac{p}{p+n} - \frac{n}{p+n} \log_2 \frac{n}{p+n}$$

$$remainder(A) = \sum_{i=1}^v \frac{p_i + n_i}{p+n} I(\frac{p_i}{p_i + n_i}, \frac{n_i}{p_i + n_i})$$

$$IG(A) = I(\frac{p}{p+n}, \frac{n}{p+n}) - remainder(A)$$

$$I(3/9,6/9)=0.92$$

$$\text{Reminder}(f1)=3/9*0+4/9*I(3/4,1/4)+2/9*0=4/9*0.81=\underline{\mathbf{0.36}}$$

$$\begin{aligned}\text{Reminder}(f2)&=4/9*I(1/4,3/4)+3/9*I(1/3,2/3)+2/9*I(1/2,1/2)=4/9*0.81+3/9*0.92+2/9*1 \\ &=\underline{\mathbf{0.889}}\end{aligned}$$

$$\text{Reminder}(f3)=5/9*I(2/5,3/5)+1/9*I(0,1)+3/9*I(1/3,2/3)=5/9*0.97+0+3/9*0.92=\underline{\mathbf{0.845}}$$

$$\mathbf{IG(f1)}=0.92-0.36=\underline{\mathbf{0.56}}$$

#	f1	f2	f3	output
1	B	b	a	t
3	B	a	a	f
4	B	c	a	t
7	B	a	c	t

$$I(3/4,1/4)=0.81$$

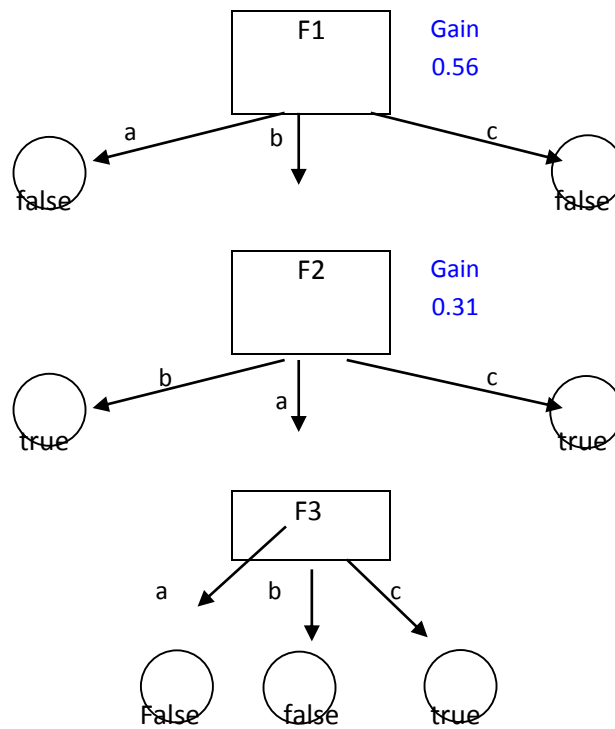
$$\text{Remainder}(f2)=2/4*I(1/2,1/2)+1/4*I(1,0)+1/4*I(1,0)=0.5+0+0=\underline{\mathbf{0.5}}$$

$$\text{Remainder}(f3)=3/4*I(2/3,1/3)+1/4*I(1,0)=3/4*0.92=\underline{\mathbf{0.69}}$$

$$\mathbf{IG(f2)}=0.81-0.5=\mathbf{0.31}$$

#	f1	f2	f3	output
3	B	a	a	f
7	B	a	c	t

--	--	--	--	--



Boolean Function:

$$f1 = b \wedge \{ (f2 = b) \vee (f2 = c) \vee [(f2 = a) \wedge (f3 = c)] \}$$

2. (תשס"ג ב') נתונה סוכנות הכריות עם database גדול של אנשים. על מנת לספק שירות יעיל אתם מעוניינים לבנות מודל לכל לקוח חדש בהתבסס על קבוצה מייצגת קטנה, אותה הלקוח כבר קטלג (ה-training set). ברגע שהמודל נלמד ניתן לסנן עבורו את ה-database ולשלוח רק את המועמדים הרלוונטיים.
נתונה קבוצת הדוגמאות הבאות ומטרתכם ללמוד את הפונקציה Candidate המחזירה yes / no, המציין את פוטנציאל המועמד ללקוח.

<i>i</i>	<i>Sex</i>	<i>Status</i>	<i>Age</i>	<i>Education</i>	<i>Location</i>	<i>Looks</i>	<i>Candidate</i>
1	M	Single	<30	HS	North	Good	No
2	F	Single	30-40	BA	Center	Good	Yes
3	F	Single	>40	MA	South	Great	No
4	M	Divorced	30-40	MA	South	Great	No
5	F	Single	<30	MA	North	Great	No
6	F	Single	30-40	HS	Center	Good	No
7	F	Single	30-40	BA	South	Good	Yes
8	F	Single	>40	HS	Center	Great	No
9	F	Divorced	30-40	MA	South	Great	Yes
10	F	Divorced	<30	HS	South	Good	No
11	F	Single	>40	BA	Center	Good	No
12	F	Divorced	<30	MA	Center	Great	No

ציירו עץ החלטה לקבוצת הדוגמאות. אם קיימות מספר תכונות עם אותו gain יש לבחור לפי סדר הא"ב. יש לציין בקודקודי העץ את ערכי ה-gain הרלוונטיים.

פתרון:

$$I(3/12, 9/12) = -(1/4)\log(1/4) - (3/4)\log(3/4) = 0.5 + 0.311 = 0.811$$

$$\text{Remainder}(\text{Age}) = (4/12) * I(0/4, 4/4) + (5/12) * I(3/5, 2/5) + (3/12) * I(0/3, 3/3) \\ = 0 + (5/12) * [-0.6\log(0.6) - 0.4\log(0.4)] + 0 = 0.405$$

$$IG(\text{Age}) = I(3/12, 9/12) - \text{Remainder}(\text{Age}) = 0.811 - 0.405 = 0.41$$

בצורה דומה נבדוק את שאר המאפיינים. המאפיין AGE נבחר כמאפיין בעל ה-IG הגדול ביותר.

נחלק את טבלת הדוגמאות על פי ערכי AGE האפשריים. אם גיל המתמודד קטן מ-30 או גדול מ-40 לכל הדוגמאות יש את אותו הערך ולכן נגדיר אותם כעלים. עתה נותרנו עם כל הדוגמאות בהן גיל המתמודד הוא בין 30 ל-40:

<i>I</i>	<i>Sex</i>	<i>Status</i>	<i>Age</i>	<i>Education</i>	<i>Location</i>	<i>Looks</i>	<i>Candidate</i>
2	F	Single	30-40	BA	Center	Good	Yes
4	M	Divorced	30-40	MA	South	Great	No
6	F	Single	30-40	HS	Center	Good	No
7	F	Single	30-40	BA	South	Good	Yes
9	F	Divorced	30-40	MA	South	Great	Yes

על מנת לבחור את המאפיין לתת העץ נחשב עבור כל מאפיין את IG:

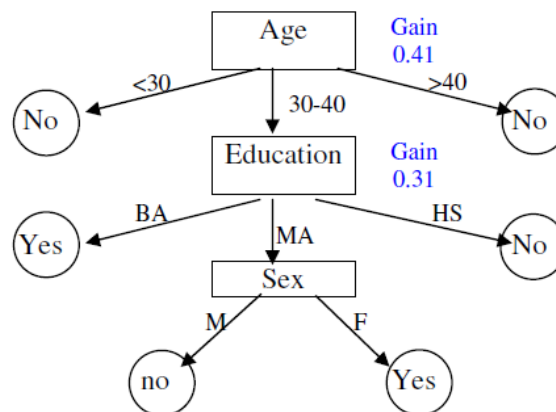
$$IG(\text{Education}) = I(3/5, 2/5) - \text{Remainder}(\text{Education})$$

$$I(3/5, 2/5) = -0.6\log(0.6) - 0.4\log(0.4) = 0.442 + 0.529 = 0.971$$

$$\text{Remainder}(\text{Education}) = (1/5)I(0/1, 1/1) + (2/5)I(2/2, 0/2) + (2/5)I(1/2, 1/2) = 0 + 0 + 0.4 = 0.4$$

$$\text{IG}(\text{Education}) = 0.971 - 0.4 = 0.571$$

העץ המתקבל הוא :



3. (תשס"א א') נתונות הדוגמאות הבאות לעץ החלטה בנושא של האם להצביע בבחירות או לא. הפרמטרים הם מזג אוויר, מידת תמיכה, וסקרים.

החלטה	סקרים	מידת תמיכה	מזג אוויר
לא	תיקו	מועטה	סוער
כן	שלי מוביל	גדולה	בהיר
כן	שלי מוביל	גדולה	מעונן
לא	יריב מוביל	גדולה	סוער
כן	תיקו	מועטה	מעונן
לא	תיקו	גדולה	סוער
לא	שלי מוביל	מועטה	בהיר
כן	יריב מוביל	מועטה	בהיר
לא	שלי מוביל	מועטה	מעונן
כן	יריב מוביל	מועטה	מעונן

א. בנו עץ החלטה מתאים, במקרה שה-gain שווה בחרו עפ"י סדר הא"ב.

ב. ציינו בקודקודי העץ הרלוונטיים את ה-gain של התכונה שבחרתם.

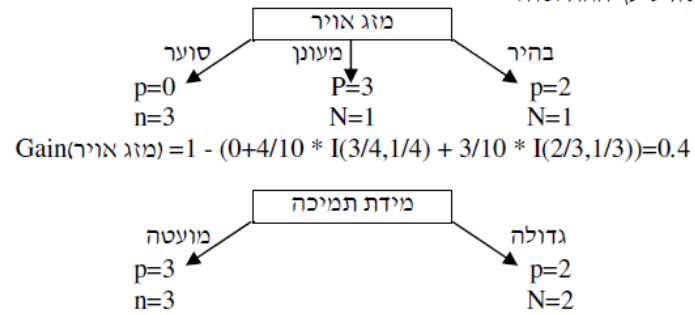
ג. תארו נוסחה לוגית שיכולה לתאר את העץ שקבלתם.

פתרון:

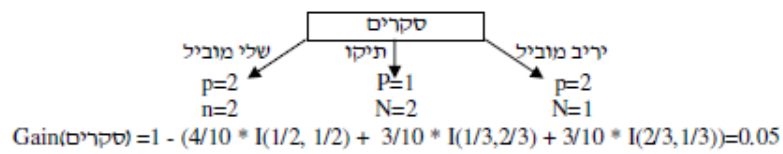
א.ב. בניית עץ החלטה לשאלה האם להצביע בבחירות:
נחשב את האנטרופיה ההתחלתית –

$$p=5, n=5 \rightarrow I(1/2, 1/2) = 1$$

נבחר תכונה לשורש עץ ההחלטה:



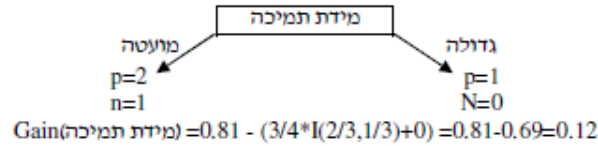
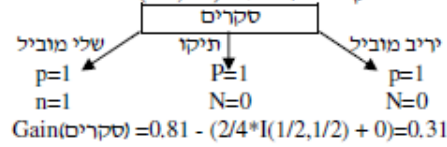
Gain(מידת תמיכה)=0



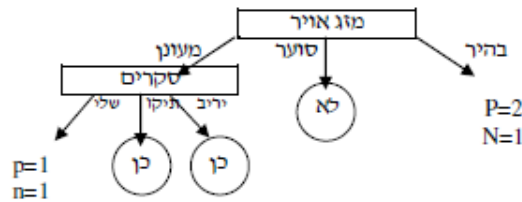
ה-Gain הכי גדול הוא של התכונה "מוזג אוויר" ולכן היא תהיה שורש העץ.



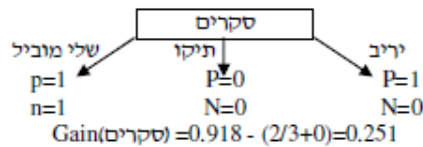
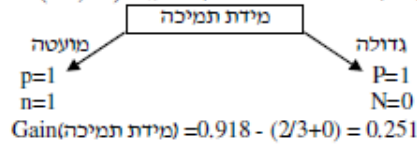
נבחר תכונה עבור תת העץ השמאלי בו מוזג האוויר הוא מעונן :
 $I(3/4, 1/4) = 0.81$ היא עץ זה היא



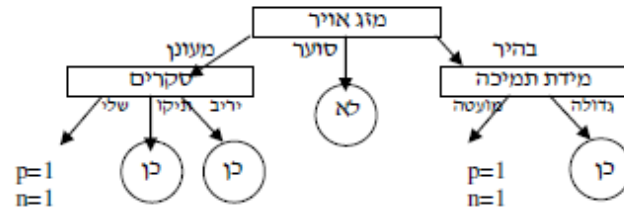
שורש תת העץ יהיה אם כן "סקרים"



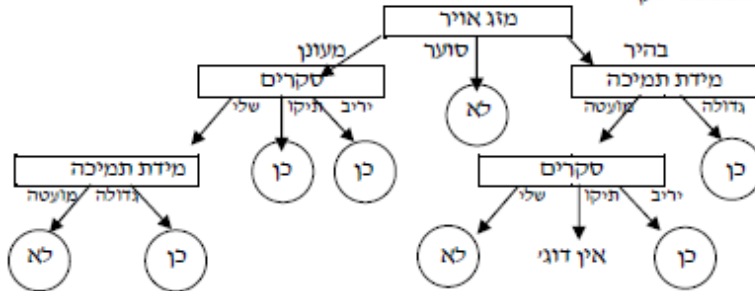
נבחר תכונה לתת העץ הימני. האנטרופיה ההתחלתית היא $I(1/3, 2/3) = 0.918$



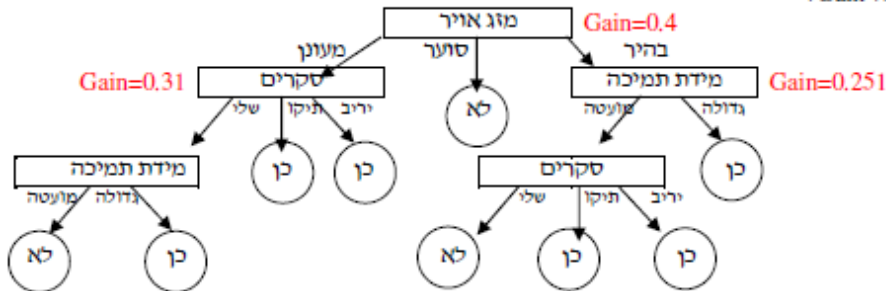
שורש תת העץ יהיה אם כן "מידת תמיכה", שכן במקרה שה-Gain שווה יש לבחור לפי סדר הא"ב. עד כה העץ הוא:



נציב את התכונות הנוותרות בתתי העץ:



היות ואין דוגמאות עבור סקרים=תיקן ולא ניתן לנו ערך ברירת מחדל, נבחר עפ"י הרוב "כן", והעץ המתקבל עם ערכי ה-Gain:



ג. נסמן ב-W את מזג האוויר, T את מידת התמיכה ו-S את הסקרים. הנוסחא שמתארת את העץ היא:

$$((W = \text{בהיר}) \wedge (T = \text{גדולה})) \vee ((W = \text{בהיר}) \wedge (T = \text{מועטה}) \wedge \neg(S = \text{שלי})) \vee ((W = \text{מועון}) \wedge (T = \text{גדולה})) \vee ((W = \text{מועון}) \wedge \neg(S = \text{שלי}))$$

$$((n-1)2^{1/2}) - (5-20)) + ((n-1)2^{1/2}) - (1-112+114))$$

ב. אם לא פצלנו את עץ ההחלטה עפ"י התכונה המפרידה ביותר בכל שלב, נקבל פתרון שלא יתאים לכל הדוגמאות של הקלט.

בטבלה שלהלן מופיעות רשומות של 12 פציינטים. לכל אחד מהם יש נתונים עבור המאפיינים הבאים: מין, גיל (האם מעל 60), האם הוא סוכרתי, האם יש לו לחץ דם גבוה, האם ה-EKG שלו חורג מהנורמה, וסיווג (classification) - האם הפציינט סובל מהפרעה בקצב הלב?

פציינט	מין	מעל 60	סכרתי	לחץ דם גבוה	EKG	הפרעה בקצב הלב
1	ז	+	+	-	-	-
2	ז	-	-	+	+	+
3	ז	-	+	+	-	-
4	ז	+	-	-	+	+
5	ז	+	+	+	-	+
6	ז	-	+	+	-	+
7	נ	-	-	+	-	-
8	נ	+	+	+	+	+
9	נ	-	+	-	+	+
10	נ	+	-	-	-	-
11	נ	+	+	-	-	-
12	נ	+	-	+	+	+

א. חשבו את האנטרופיה: $H(\text{HasArrhythmia}|\text{Sex}=\text{Female})$
 (השתמשו בנוסחה המופיעה בראש עמוד 704 בספר הלימוד).

$$H(\text{HasArrhythmia}|\text{Sex} = \text{Female}) = -(0.5 \log_2 0.5 + 0.5 \log_2 0.5) = 1$$

ב. מה היא התכונה שתיבחר להיות בשורש עץ ההחלטה?
 הניחו כי:

$$\text{Cost}(\text{Sex})=\text{Cost}(\text{Over60})=1$$

$$\text{Cost}(\text{Diabetic})=3$$

$$\text{Cost}(\text{HighPulse})=2$$

$$\text{Cost}(\text{AbnormalEKG})=5$$

נשתמש בביטוי $A = \frac{\text{Gain}^2(S,A)}{\text{Cost}(A)}$ כדי לבחור מי יהיה בשורש עץ ההחלטה. בדוגמא עגלתי לשתי ספרות אחרונות.

$$\text{gain}(\text{Sex})=0.98 - 0.5 * E(3, 3) - 0.5 * E(4, 2) = 0.02$$

$$A = 0.02^2 / 1.0 = 0$$

$$\text{gain}(\text{over 60})=0.98 - 0.583 * E(4, 3) - 0.417 * E(3, 2) = 0$$

$$A = 0^2 / 1.0 = 0$$

$$\text{gain}(\text{Diabetic})=0.98 - 0.583 * E(4, 3) - 0.417 * E(3, 2) = 0$$

$$A = 0^2 / 3.0 = 0$$

$$\text{gain}(\text{HighPulse})=0.98 - 0.5 * E(5, 2) - 0.5 * E(2, 5) = 0.12$$

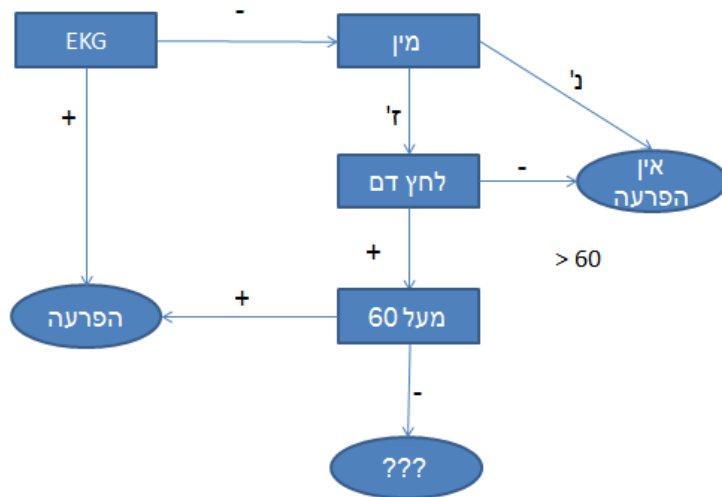
$$A = 0.12^2 / 2.0 = 0.01$$

$$\text{gain}(\text{AbnormalEKG})=0.98 - 0.417 * E(5, 0) - 0.583 * E(2, 5) = 0.48$$

$$A = 0.48^2 / 5.0 = 0.05$$

EKG קיבל את התוצאות הטובות ביותר ולכן הוא יהיה בשורש עץ ההחלטה.

ג. בנו עץ החלטה לחיזוי קבלה לאוניברסיטה. פרטו את כל שלבי הבניה.



ד. נניח כי עבור קבוצה אחרת של פציינטים ידוע גילם המדויק.

גילאי הדוגמאות החיוביות (positive examples) הם: {40,60,62,64,70,74,75,82}

וגילאי הדוגמאות השליליות (negative examples) הם {33,35,42,45,49,52,58,59,80}.

נניח שכל יתר התכונות ב-data set פחות "טובות" מתכונת הגיל, כך שנרצה לפצל את העץ

ב-Age=k על-ידי חלוקת קבוצת הדוגמאות לשתי תת-קבוצות: Age<k ו-Age≥k.

בהתבסס על תוספת אינפורמציה (information gain), איזו חלוקה נרצה לבחור ?

הדוגמאות השליליות מסומנות באדום.

33, 35, 40, 42, 45, 49, 52, 58, 59, 60, 62, 64, 70, 74, 75, 80, 82

מהשרטוט קל לראות שאם נשים את הגבול בין 59-60 נהיה זקוקים למספר הנמוך ביותר של ביטים כדי לסווג דוגמה חדשה, לכן remindern יהיה הנמוך ביותר ולכן הgain יהיה מקסימלי.