

עבודה 2- ניתוח מאגרי מידע

בעיה 3 (עץ החלטות):

א. על מנת לבנות את עץ ההחלטות נחשב Entropy עבור כל אופציה בקטגוריה ונחשב עבור כל קטגוריה

את ה Information Gain. נבחר את השאלה לפי הקטגוריה עם הרווח מידע הגבוה ביותר:

תחילה נחשב Entropy עבור כל הדאטא סט, כאשר 5 משחקים לא משחקים ו-10 כן (15 סה"כ):

$$entropy = -p_+ * \log_2 p_+ - p_- * \log_2 p_- = -\frac{10}{15} * \log_2 \frac{10}{15} - \frac{5}{15} * \log_2 \frac{5}{15} = 0.918$$

על מנת להבין איזה שאלה כדאי לשאול ראשונה, נחשב לכל קטגוריה entropy לפי הנוסחא למעלה,

והנוסחא בה נשתמש לחישוב ה Information Gain:

$$information\ gain = entropy(s) - \sum_v \frac{|s_v|}{|s|} * entropy(s_v)$$

outlook מזג

$$\text{outlook, sunny: } \left. \begin{array}{l} p_+ = 2 \\ p_- = 3 \end{array} \right\} Entropy = 0.971$$

$$\text{outlook, rainy: } \left. \begin{array}{l} p_+ = 4 \\ p_- = 2 \end{array} \right\} Entropy = 0.918$$

$$\text{outlook, overcast: } \left. \begin{array}{l} p_+ = 4 \\ p_- = 0 \end{array} \right\} Entropy = 0$$

$$Information\ Gain\ (outlook) = Entropy(s) - \sum_v \frac{|s_v|}{|s|} * Entropy(s_v)$$

$$= 0.918 - \left[\frac{5}{15} * 0.971 + \frac{6}{15} * 0.918 + \frac{4}{15} * 0 \right] = 0.227$$

temp טמפר

$$\text{temp, hot: } \left. \begin{array}{l} p_+ = 2 \\ p_- = 2 \end{array} \right\} Entropy = 1$$

$$\text{temp, mild: } \left. \begin{array}{l} p_+ = 4 \\ p_- = 2 \end{array} \right\} Entropy = 0.918$$

$$\text{temp, cool: } \left. \begin{array}{l} p_+ = 4 \\ p_- = 1 \end{array} \right\} Entropy = 0.722$$

$$Info\ gain\ (temp) = 0.918 - \left[\frac{4}{15} * 1 + \frac{6}{15} * 0.918 + \frac{5}{15} * 0.722 \right] = 0.043$$

: humidity גבוה

$$\text{humidity, high: } \left. \begin{array}{l} P+ = 3 \\ P- = 4 \end{array} \right\} \text{Entropy} = 0.985$$

$$\text{humidity, normal: } \left. \begin{array}{l} P+ = 7 \\ P- = 1 \end{array} \right\} \text{Entropy} = 0.544$$

$$\text{info gain} = 0.918 - \left[\frac{7}{15} \cdot 0.985 + \frac{8}{15} \cdot 0.544 \right] = 0.168$$

: Windy גבוה

$$\text{Windy, weak: } \left. \begin{array}{l} P+ = 6 \\ P- = 1 \end{array} \right\} \text{Entropy} = 0.592$$

$$\text{Windy, strong: } \left. \begin{array}{l} P+ = 4 \\ P- = 4 \end{array} \right\} \text{Entropy} = 1$$

$$\text{info gain} = 0.918 - \left[\frac{7}{15} \cdot 0.592 + \frac{8}{15} \cdot 1 \right] = 0.108$$

סיכום info gain עבור כלל הדאטא (לשאלה הראשונה):

Outlook = 0.227

Temp = 0.043

Humidity = 0.168

Windy = 0.108

כיוון ש- **outlook** בעל הרווח הכי גבוה הוא יהיה השאלה הראשונה.

עבור השאלה השניה:

:Outlook = overcast

כל המשחקים שוחקו ולכן $\text{entropy} = 0$, לכן לא נשאל שאלה נוספת וישר נשחק.

:Outlook = sunny

(חושב מקודם) $\text{Entropy(s)} = 0.971$

קיימים 5 משחקים: 3-לא, 2-כן

:temp 7/28

$$\text{sunny} \rightarrow \text{temp, hot: } \left. \begin{array}{l} p^+ = 0 \\ p^- = 2 \end{array} \right\} \text{Entropy} = 0$$

$$\text{sunny} \rightarrow \text{temp, mild: } \left. \begin{array}{l} p^+ = 1 \\ p^- = 1 \end{array} \right\} \text{Entropy} = 1$$

$$\text{sunny} \rightarrow \text{temp, cool: } \left. \begin{array}{l} p^+ = 1 \\ p^- = 0 \end{array} \right\} \text{Entropy} = 0$$

$$\text{info gain (temp)} = 0.971 - \left[\frac{2}{5} \cdot 0 + \frac{2}{5} \cdot 1 + \frac{1}{5} \cdot 0 \right] = 0.571$$

:humidity 7/28

$$\text{sunny} \rightarrow \text{humidity, high: } \left. \begin{array}{l} p^+ = 0 \\ p^- = 3 \end{array} \right\} \text{Entropy} = 0$$

$$\text{sunny} \rightarrow \text{humidity, normal: } \left. \begin{array}{l} p^+ = 2 \\ p^- = 0 \end{array} \right\} \text{Entropy} = 0$$

$$\text{info gain (humidity)} = 0.971 - 0 = 0.971$$

:Windy 7/28

$$\text{sunny} \rightarrow \text{Windy, Weak: } \left. \begin{array}{l} p^+ = 1 \\ p^- = 1 \end{array} \right\} \text{Entropy} = 1$$

$$\text{sunny} \rightarrow \text{Windy, strong: } \left. \begin{array}{l} p^+ = 1 \\ p^- = 2 \end{array} \right\} \text{Entropy} = 0.918$$

$$\text{info gain (windy)} = 0.971 - \left[\frac{2}{5} \cdot 1 + \frac{3}{5} \cdot 0.918 \right] = 0.020$$

סיכום info gain עבור sunny:

Temp = 0.571

Humidity = 0.971

Windy = 0.02

הרווח הכי גבוה הוא מ- **humidity** ולכן נבחר בו לשאלה השנייה.

:Outlook = rainy

(חושב קודם) Entropy = 0.918

קיימים 6 משחקים: 4-כן, 2-לא

temp נמדד

$$\text{rainy} \rightarrow \text{temp, hot: } \left. \begin{array}{l} p^+ = 0 \\ p^- = 0 \end{array} \right\} \text{Entropy} = 0$$

$$\text{rainy} \rightarrow \text{temp, mild: } \left. \begin{array}{l} p^+ = 2 \\ p^- = 1 \end{array} \right\} \text{Entropy} = 0.918$$

$$\text{rainy} \rightarrow \text{temp, cool: } \left. \begin{array}{l} p^+ = 2 \\ p^- = 1 \end{array} \right\} \text{Entropy} = 0.918$$

$$\text{info gain (temp)} = 0.918 - \left[\frac{3}{6} \cdot 0.918 + \frac{3}{6} \cdot 0.918 \right] = 0$$

humidity נמדד

$$\text{rainy} \rightarrow \text{humidity, high: } \left. \begin{array}{l} p^+ = 1 \\ p^- = 1 \end{array} \right\} \text{Entropy} = 1$$

$$\text{rainy} \rightarrow \text{humidity, normal: } \left. \begin{array}{l} p^+ = 3 \\ p^- = 1 \end{array} \right\} \text{Entropy} = 0.811$$

$$\text{info gain (humidity)} = 0.918 - \left[\frac{2}{6} \cdot 1 + \frac{4}{6} \cdot 0.811 \right] = 0.044$$

Windy נמדד

$$\text{rainy} \rightarrow \text{Windy, weak: } \left. \begin{array}{l} p^+ = 3 \\ p^- = 0 \end{array} \right\} \text{Entropy} = 0$$

$$\text{rainy} \rightarrow \text{Windy, strong: } \left. \begin{array}{l} p^+ = 1 \\ p^- = 2 \end{array} \right\} \text{Entropy} = 0.918$$

$$\text{info gain (windy)} = 0.918 - \left[\frac{3}{6} \cdot 0 + \frac{3}{6} \cdot 0.918 \right] = 0.459$$

סיכום info gain עבור sunny:

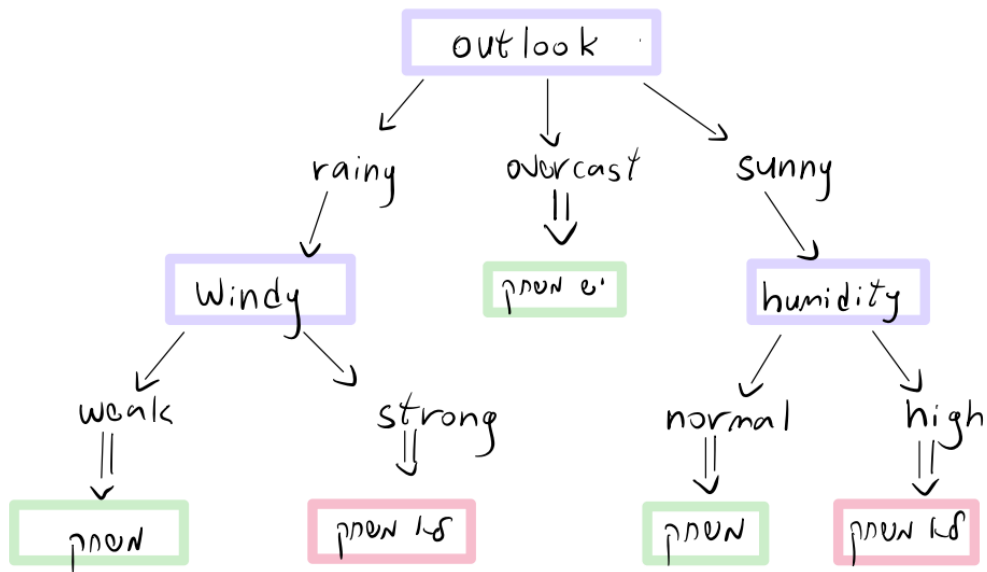
Temp = 0

Humidity = 0.044

Windy = 0.459

הרווח הכי גבוה הוא מ- windy ולכן נבחר בו לשאלה השניה.

מקרא: שאלות, תשובה חיובית למשחק, תשובה שלילית למשחק



- ב. בעץ החלטות זה השגיאה עומדת על 1/15 טעויות כלומר: 6.667% סיכוי לטעות.
- ג. לפי עץ ההחלטות שלי, הימים בהם הייתי משחקת הם 1,2 ו-3 (משחקת ב-D16, D17, D18 ולא הייתי משחקת ב-D19) כלומר משחקת בימים בהם מעונן או גשום ורוח חלשה.