

## תרגיל בית 5 – Decision Trees, Ensemble Methods

יש להגיש שני קבצים נפרדים: קובץ PDF ובו פתרון התרגיל כולל הפלטים של החלק המעשי וקובץ נוסף ובו הקוד שכתבתם. יש להקפיד על תשובות ברורות ומסודרות ועל קוד מסודר ומתועד היטב. רק אחד מבין חברי הזוג צריך להגיש את הפתרון.

שאלות על התרגיל יש לכתוב בפורום תרגילי הבית באתר הקורס. התרגיל מנוסח בלשון נקבה אך מתייחס לשני המינים.

### שאלה 1

נתון אוסף התצפיות הבא:

<u>ID</u>	<u>Family heart attacks</u>	<u>Gender</u>	<u>Cholesterol</u>	<u>Blood pressure</u>	<u>Heart attack</u>
<u>1</u>	<u>Yes</u>	<u>Male</u>	<u>160</u>	<u>High</u>	<u>No</u>
<u>2</u>	<u>Yes</u>	<u>Male</u>	<u>260</u>	<u>Normal</u>	<u>Yes</u>
<u>3</u>	<u>No</u>	<u>Female</u>	<u>245</u>	<u>Normal</u>	<u>No</u>
<u>4</u>	<u>No</u>	<u>Male</u>	<u>170</u>	<u>Normal</u>	<u>No</u>
<u>5</u>	<u>Yes</u>	<u>Female</u>	<u>230</u>	<u>High</u>	<u>Yes</u>
<u>6</u>	<u>No</u>	<u>Male</u>	<u>215</u>	<u>Normal</u>	<u>No</u>
<u>7</u>	<u>Yes</u>	<u>Female</u>	<u>240</u>	<u>Normal</u>	<u>No</u>
<u>8</u>	<u>No</u>	<u>Male</u>	<u>235</u>	<u>High</u>	<u>Yes</u>

- נרצה לנבא האם אדם יקבל התקף לב או לא. על פי הנלמד בהרצאה ובתרגול, ציירו את עץ ההחלטה המתאים לנתונים. על מנת לחשב השתמשו ב Entropy.
- כעת השתמשו בGini. האם העץ הינו אותו עץ?
- בהסתמך על הסעיפים הקודמים, האם השימוש בGini וEntropy יכול ליצור את אותו העץ? אם כן - האם העץ תמיד יהיה אותו עץ? במידה וכן - הסבירו. במידה ולא - הראו דוגמה נגדית.

## שאלה 2 – שאלה ממבחן

אלגוריתם ה AdaBoost מבצע שלושה צעדים:

$$h_t = WL(D^{(t)}, S) \quad (i)$$

$$w_t = \log\left(\frac{1}{\epsilon_t} - 1\right)/2 \text{ וגם } \epsilon_t = \sum_{i=1}^m D_i^{(t)} 1_{[h_t(x_i) \neq y_i]} \quad (ii)$$

$$D_i^{(t+1)} \propto D_i^{(t)} e^{-w_t y_i h_t(x_i)} \quad (iii)$$

א. הסבירו את ההנחה על  $WL()$  ואת כל אחד מהצעדים באלגוריתם.

$$ב. \text{ הראו ש } 2\sqrt{\epsilon_t(1 - \epsilon_t)} = \sum_{i=1}^m D_i^{(t)} e^{-w_t y_i h_t(x_i)}.$$

$$ג. \text{ הראו ש } 1/2 = \sum_{i=1}^m D_i^{(t+1)} 1_{[h_t(x_i) \neq y_i]}.$$

ד. מהו המספר המינימלי של צעדים הנדרש באלגוריתם ה AdaBoost בכדי שפונקציית ההפסד שלו תהיה 0 על מדגם האימון?

### שאלה 3 – חלק רטוב

בחלק זה של התרגיל נממש עץ החלטה ע"פ האלגוריתם CART שנלמד בהרצאה ובתרגול. נעבוד עם מדד Entropy.

1. Data עליו נעבוד הינו Data אודות אבחון של סרטן השד.
  - a. על מנת לטעון את Data אפשר להשתמש בפונקציה load\_breast\_cancer של sklearn. קישור - [https://scikit-learn.org/stable/modules/generated/sklearn.datasets.load\\_breast\\_cancer.html#sklearn.datasets.load\\_breast\\_cancer](https://scikit-learn.org/stable/modules/generated/sklearn.datasets.load_breast_cancer.html#sklearn.datasets.load_breast_cancer)
  - b. לפני מתחילים לעבוד, מומלץ לקרוא את תיאור הדאטה - [https://archive.ics.uci.edu/ml/datasets/Breast+Cancer+Wisconsin+\(Diagnostic\)](https://archive.ics.uci.edu/ml/datasets/Breast+Cancer+Wisconsin+(Diagnostic))
2. ממשו את מדד Entropy
  - a. קלט הפונקציה – רשימה של נתונים
  - b. פלט הפונקציה – ערך Entropy
3. ממשו את בניית העץ על פי האלגוריתם שראינו בתרגול ע"פ אלגוריתם CART
  - a. מומלץ להיעזר בחבילות מובנות המממשות עצים (כדוגמת anytree) אך אין חובה
4. פצלו את Data למדגם אימון (80%) ומדגם מבחן (20%), את הפיצול יש לבצע באמצעות הפונקציה train\_test\_split עם random\_state=3.
5. בנו את העץ על פי מדגם האימון ובדקו את התוצאות על מדגם המבחן
6. על מנת לבדוק את עצמכם, בדקו מול הרצה של האלגוריתם הממומש בsklearn עם אותם הפרמטרים (הציון בשאלה זו יקבע לפי דיוק הפלט אל מול הפלט מsklearn)
7. ציירו את העץ שקיבלתם (ניתן לצייר אותו בכל דרך שהיא ולצרף לפתרון)
8. הציגו את התוצאות בConfusion Matrix