



תרגיל 8 – Classification

הוראות הגשה:

1. בתרגיל הבא יש לענות על השאלות באמצעות שימוש בקוד פייתון ושימוש ב-Scikit-Learn.
2. יש להגיש את העבודה בזוגות בלבד.
3. התרגיל יוגש כמחברת colab, כאשר לתיבת הגשה יש להגיש קובץ המכיל לינק למחברת עם הרשאות קריאה למייל yaelhoc@post.bgu.ac.il + tamarin@post.bgu.ac.il
4. שם הקובץ יהיה מספרי הזהות של המגישים בצורה הבאה: זהות1_זהות2
במחברת הפתרון, יש לציין את מספר השאלה עליה עניתם עבור כל חלק בפתרון

ענו על השאלות הבאות באמצעות הנתונים על מחירי יהלומים:

1. חלקו את המידע ל `test` ו `train`. באמצעות אלגוריתם KNN ($K=1$) בנו מודל שחזה `clarity` של היהלום באמצעות נתוני `carat`, `depth`, `price`, `table`, `x`, `y`. לאחר מכן, חשבו את מדדי `f1_score` ו `accuracy` עבור המודל שיצרתם.
 2. חזרו על בניית המודל בסעיף 1 עבור ערכי k שונים. באמצעות `plotly express` ציירו גרפים של ביצועי המודלים עבור ערכי k שונים. כלומר, יש לצייר גרפים שבהם ציר ה- X הוא ערך ה- K של המודלים ואילו ערכי ה- Y הם מדדי `accuracy` ו-`f1_score` שמתאימים לכל מודל.
 3. חלקו את המידע ל `test` ו `train`. לאחר מכן, באמצעות אלגוריתמי Decision ו KNN ($K=1,3,5$) בנו מודלים שחוזים את `clarity` של היהלום באמצעות עמודות `carat`, `depth`, `price`, `table`, `x`, `y`. לפי מדד `accuracy`, מבין המודלים שיצרתם, איזה מודל חזה את `clarity` טוב יותר?
 4. חזרו על בניית המודלים בדומה לסעיף הקודם, רק הפעם בנו את המודל על ידי הוספת מידע מעמודות `color` ו `clarity` (כלומר, בנו את המודלים בתוספת שתי העמודות הנוספות). האם `accuracy` של המודלים השתפרו?
רמז: יש להשתמש ב-`LabelEncoder` בסעיף זה
 5. בדומה לסעיף 4, בנו מודלי KNN ($K=3$) החוזים את `clarity` של היהלום, רק שהפעם על מנת לבנות את המודלים, השתמשו בגדלים שונים של נתונים:
5%,10%,20%, 50%,70%,80%,90%,
כלומר, יש לבנות את המודל רק על ידי שימוש ב-5% מהמידע, 10% מהמידע, וכו',
צרו גרף שבו ציר ה- X הוא גודל ה-`trainset` באחוזים ואילו ציר ה- Y הוא `accuracy` של המסווג.
- הערה חשובה:** חשוב להשתמש ב-`test` באותם נתונים בדיוק עבור כל המודלים