

1. בתרגיל הבא יש לענות על השאלות באמצעות שימוש בקוד פייתון ושימוש ב-Scikit-Learn.

ענו על השאלות הבאות באמצעות הנתונים על מחירי יהלומים:

1. חלקו את המידע ל train ו test. באמצעות אלגוריתם KNN ( $K=3$ ) 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 seaborn 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 train ו test. לאחר מכן, באמצעות אלגוריתמי KNN ( $K=1,3,5,7$ ) ו Decision Tree בנו מודלים שחזים את cut של היהלום באמצעות עמודות carat, depth, price, table, x, and y. לפי מדד accuracy, מבין המודלים שיצרתם, איזה מודל חזה את cut 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 ( $K=5$ ) החזים את cut של היהלום, רק שהפעם על מנת לבנות את המודלים, השתמשו בגדלים שונים של נתונים: 5%, 10%, 20%, 50%, 75%, 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 trainset באחוזים ואילו ציר ה-Y הוא accuracy של המסווג.  
**הערה חשובה:** חשוב להשתמש כ-test באותם נתונים בדיוק עבור כל המודלים