

K-NN אלגוריתם

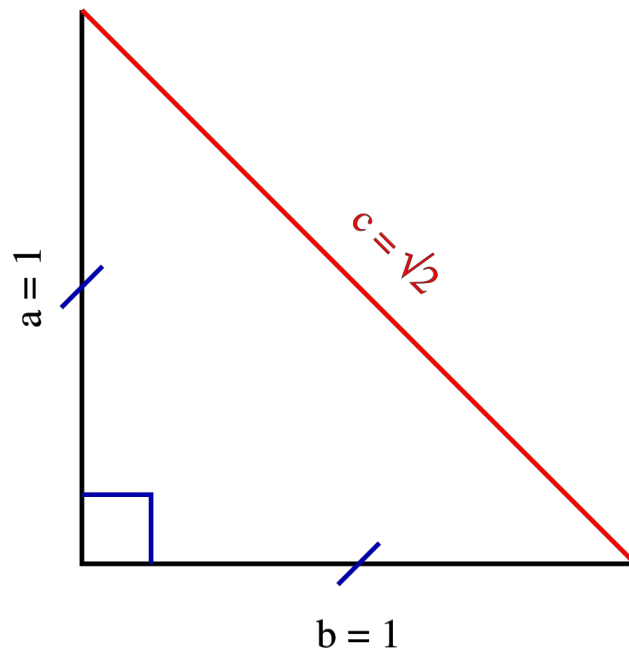
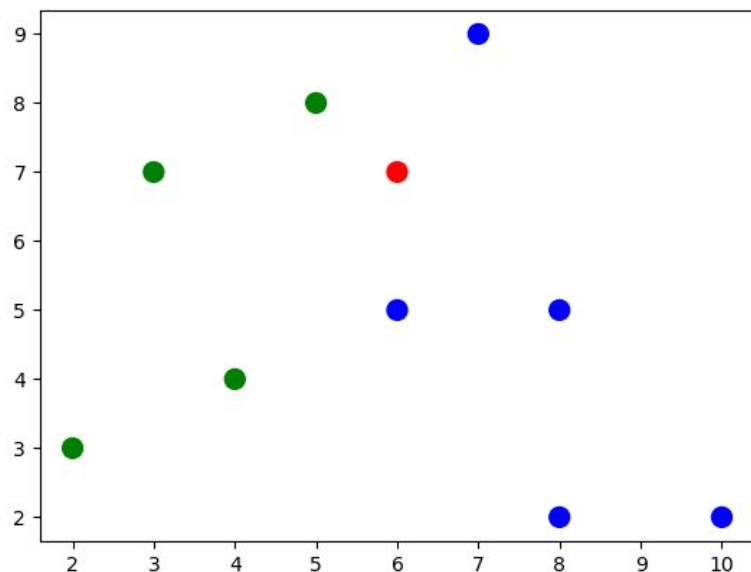
לפיתוח מכונה לומדת

K-Nearest Neighbors algorithm

גדי הרמן

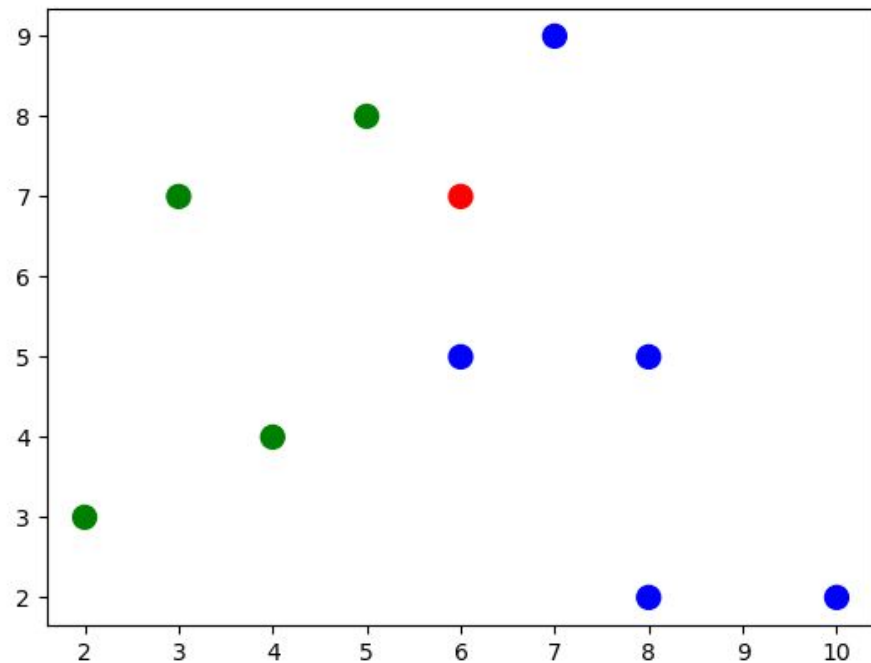
אלגוריתם K-Nearest Neighbors algorithm או בקיצור K-NN

K-NN מבוסס בין היתר על מדידת מרחק אוקלידי בין נקודות במרחב או בשם הפשוט יותר שלו, שימוש בפיתגורס כדי לחשב את היתר במשולש ישר זווית.

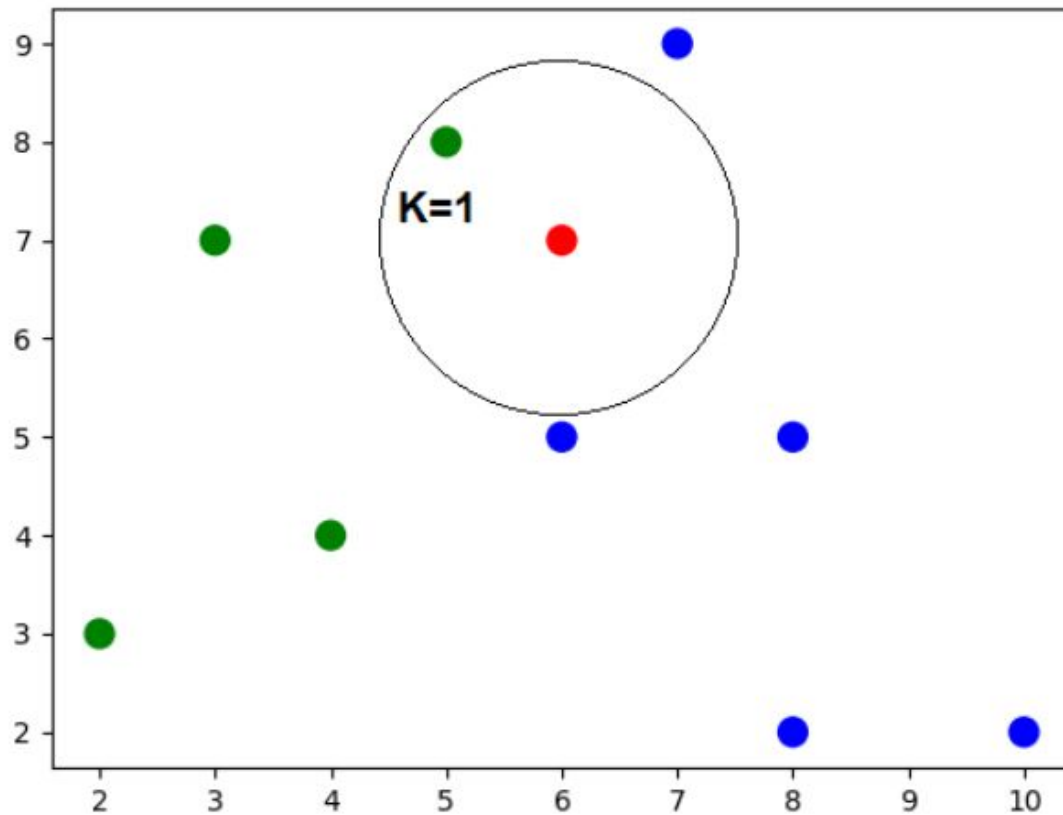


נבחן תחילה את הקוד הבא

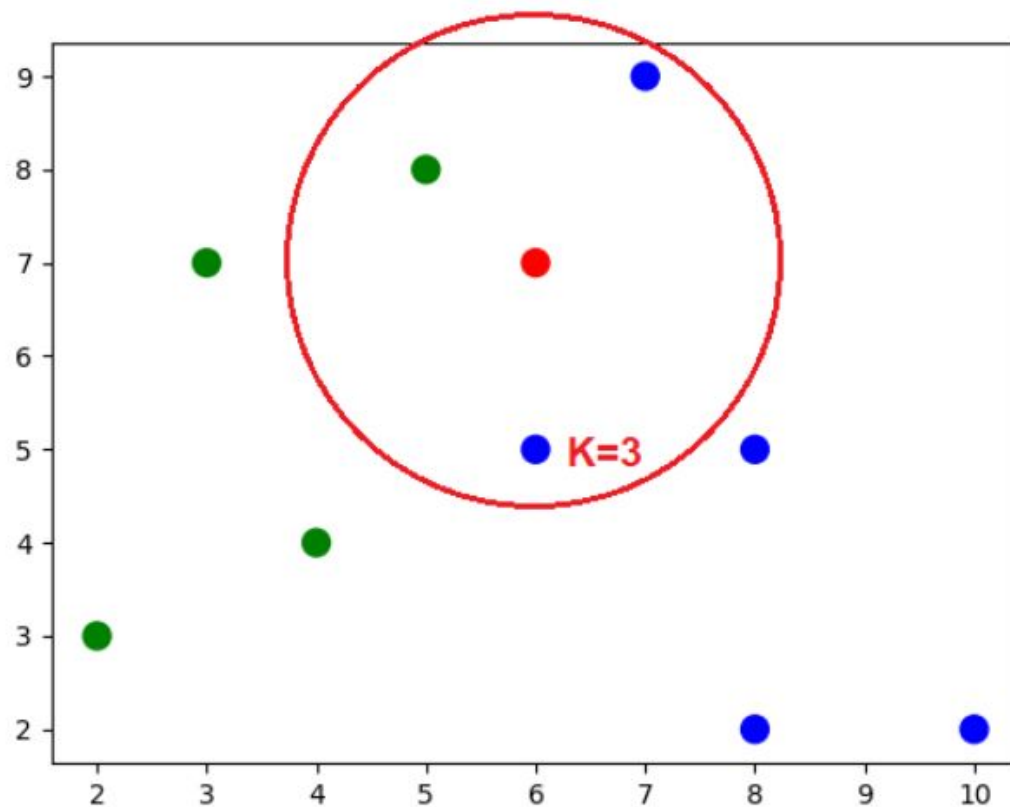
```
1 import matplotlib.pyplot as plt
2 import numpy as np
3 data = np.array( [ [ 6 , 7],
4                   [ 2 , 3],
5                   [ 3 , 7],
6                   [ 4 , 4],
7                   [ 5 , 8],
8                   [ 6 , 5],
9                   [ 7 , 9],
10                  [ 8 , 5],
11                  [ 8 , 2],
12                  [10 , 2]  ])
13 categories = np.array([0,1,1,1,1,2,2,2,2,2])
14 colormap = np.array(['r', 'g', 'b'])
15 plt.scatter(data[:,0], data[:,1], s=100, c=colormap[categories])
16 plt.show()
```



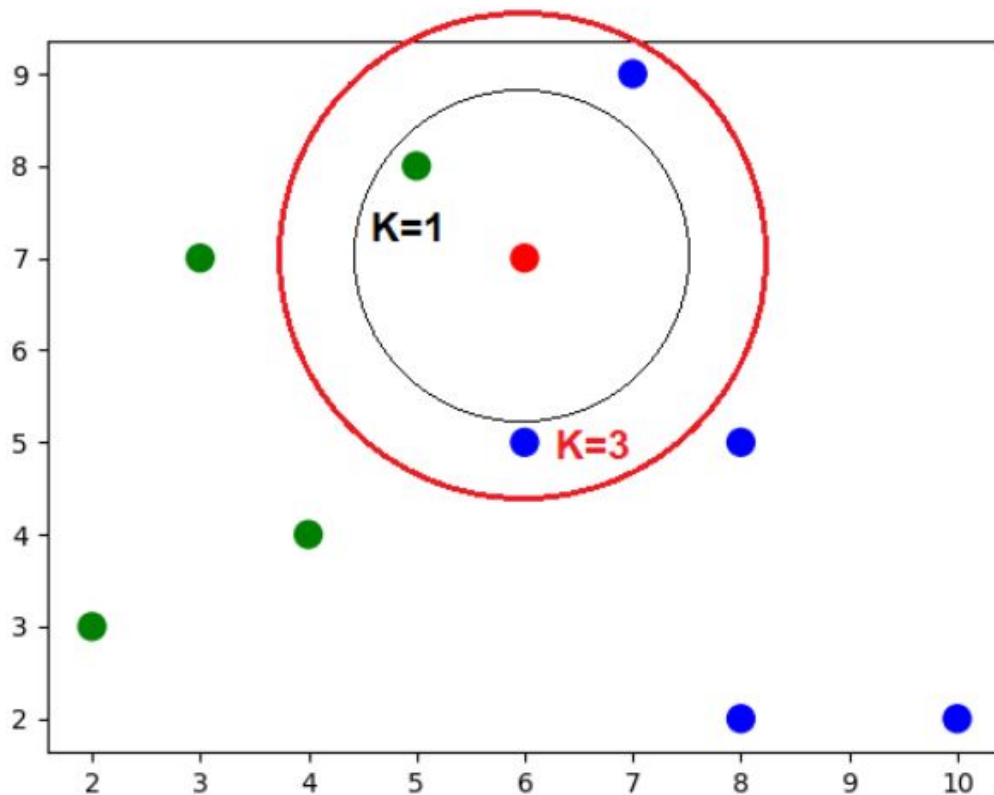
מי החברים שלי? $K=1$



מי החברים שלי? $K=3$

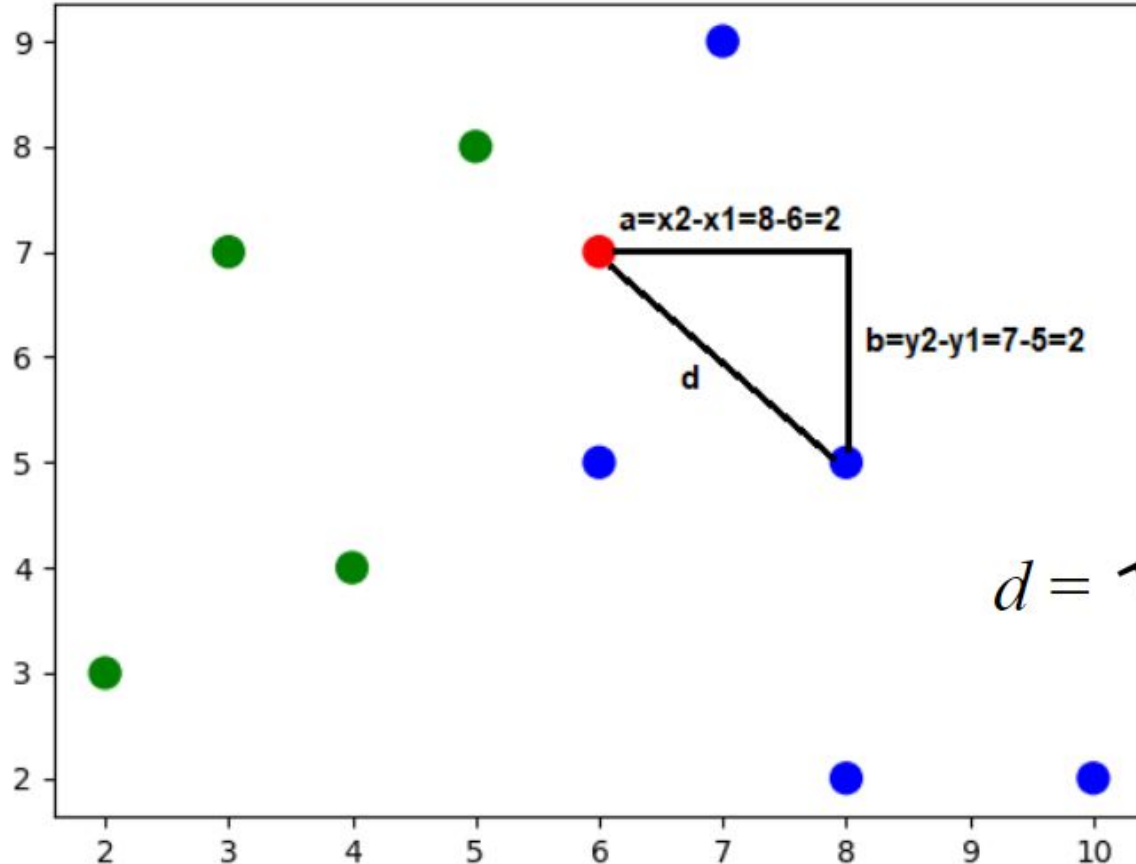


מי החברים שלי? $K=1$ או $K=3$



מרחק אוקלידי בין נקודות

כדי לממש את אלגוריתם KNN אנו זקוקים לנוסחה לחישוב מרחק בין נקודות. נדגים זאת על ידי מציאת המרחק בין 2 נקודות על פני מישור דו-מימדי.



$$d = \sqrt{(x_2 - x_1)^2 + (y_2 - y_1)^2}$$

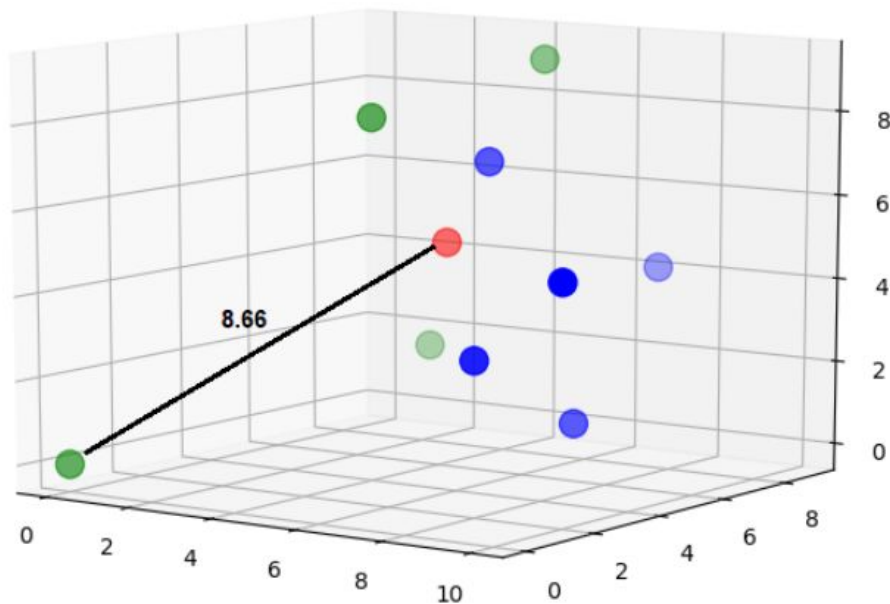
מימוש מרחק אוקלידי בין נקודות בקוד

$$d = \sqrt{(x_2 - x_1)^2 + (y_2 - y_1)^2}$$

```
def euclidean_distance(p1, p2):  
    dx = float(p1[0]-p2[0])  
    dy = float(p1[1]-p2[1])  
    d = np.power(dx,2) + np.power(dy,2)  
    return np.sqrt(d)
```


מרחק אוקלידי על גבי מערכת מרובת צירים

```
data = np.array( [ [ 5.0 , 5.0, 5.0],  
[ 0.0 , 0.0, 0.0],  
[ 3.0 , 7.0, 2.0],  
[ 4.0 , 4.0, 8.0],  
[ 5.0 , 8.0, 9.0],  
[ 6.0 , 5.0, 7.0],  
[ 7.0 , 9.0, 4.0],  
[ 8.0 , 5.0, 1.0],  
[ 8.0 , 2.0, 3.0],  
[10.0 , 2.0, 5.0] ])
```



$$d = \sqrt{(x_2 - x_1)^2 + (y_2 - y_1)^2 + (z_2 - z_1)^2}$$

$$d = \sqrt{(5 - 0)^2 + (5 - 0)^2 + (5 - 0)^2} = 8.66$$

מרחק אוקלידי על גבי מערכת מרובת צירים

נכתוב מחדש את הפעולה euclidean_distance כדי שתתאים לחישוב מרחק אוקלידי ביותר מ-2 מימדים.

חישוב מרחק ב-N מימדים

```
def euclidean_distance(p1, p2):  
    d = 0.0  
    for i in range(len(p1)):  
        a = float(p1[i])  
        b = float(p2[i])  
        d += np.power((a-b),2)  
    d = np.sqrt(d)  
    return d
```

חישוב מרחק ב-2 מימדים

```
def euclidean_distance(p1, p2):  
    dx = float(p1[0]-p2[0])  
    dy = float(p1[1]-p2[1])  
    d = np.power(dx,2) + np.power(dy,2)  
    return np.sqrt(d)
```

יישום מרחק אוקלידי מרובה צירים באלגוריתם KNN

נממש את הפעולה בשם KNN שתקבל:

- מערך נקודות כפרמטר בשם train.
- נקודה שאליה אנו רוצים לחפש את הנקודות הקרובות ביותר בשם test
- פרמטר נוסף בשם K שקובע את מספר הנקודות אותה הפעולה תחזיר.

הפעולה תבצע את האלגוריתם הבא:

- תעבור בלולאה על כל הערכים במערך train ובכל פעם תזמן את הפעולה euclidean_distance
- הערך המוחזר מפעולה euclidean_distance נכנס למערך פנימי הכולל את המרחק ואת ערכי הנקודה שנבדקה.
- לאחר סיום מדידת כל הנקודות נמיין את המערך לפי המרחקים.
- נמחק את האיבר הראשון במערך בגלל שהוא מודד מרחק לנקודה test כך שברור שהערך תמיד יהיה אפס.
- לבסוף נעבור על הלולאה הממוינת ונחלץ ממנה את K האיברים שאותם אנו רוצים לקבל.

יישום מרחק אוקלידי מרובה צירים באלגוריתם KNN

```
def KNN(train, test, K):  
    distances = []  
    NN = []  
    for p in train:  
        dist = euclidean_distance(test, p)  
        distances.append((p, dist))  
    distances.sort(key=lambda dist: dist[1])  
    distances = np.delete(distances, [0], axis=0)  
    for i in range(K):  
        NN.append(distances[i][0])  
    return NN
```

מימוש הפעולה predict בתוך KNN

הפעולה predict

תקבל

- מערך נקודות כפרמטר בשם `train`.
- נקודה בודדת כפרמטר בשם `test` שאליה אנו רוצים לחפש את הנקודות הקרובות ביותר.
- מערך תגיות בשם `lbl` המציין לאיזו קבוצה שייכת כל נקודה במערך `train`.
- פרמטר בשם `K` שקובע את מספר הנקודות שעל פיהם נקבע לאיזו קבוצה שייכת הנקודה `test`.

תחזיר

- פרמטר אחד שהוא מספר הקבוצה.

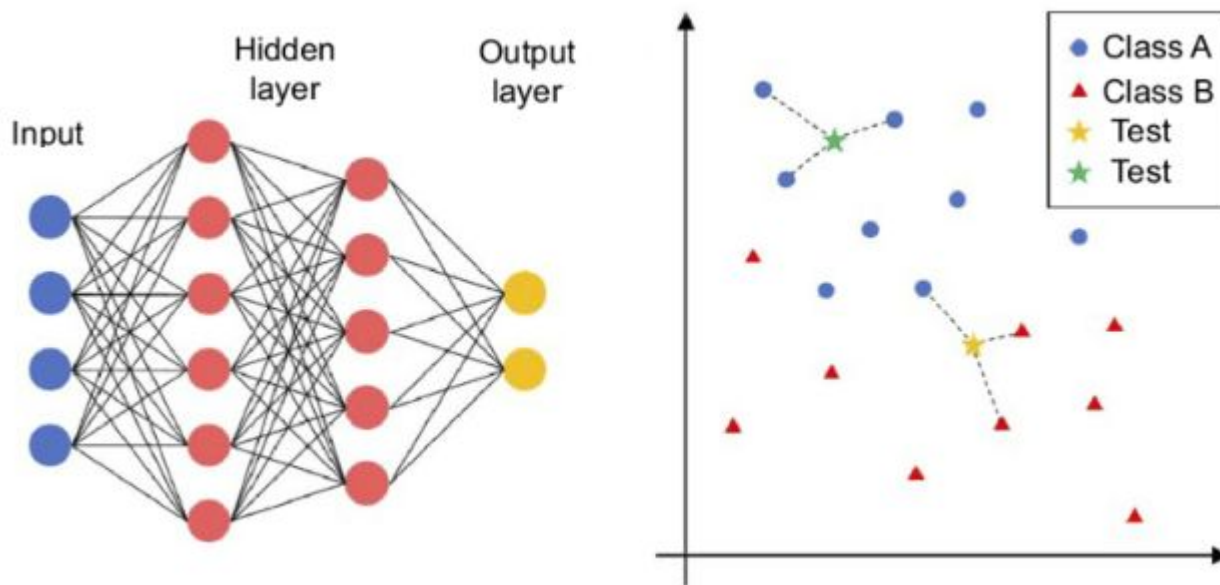
מימוש הפעולה predict בתוך KNN

```
def predict(train, test, lbl, K):  
    neighbors = KNN(train, test, lbl, K)  
    out = [row[-1] for row in neighbors]  
    return max(set(out), key=out.count)
```

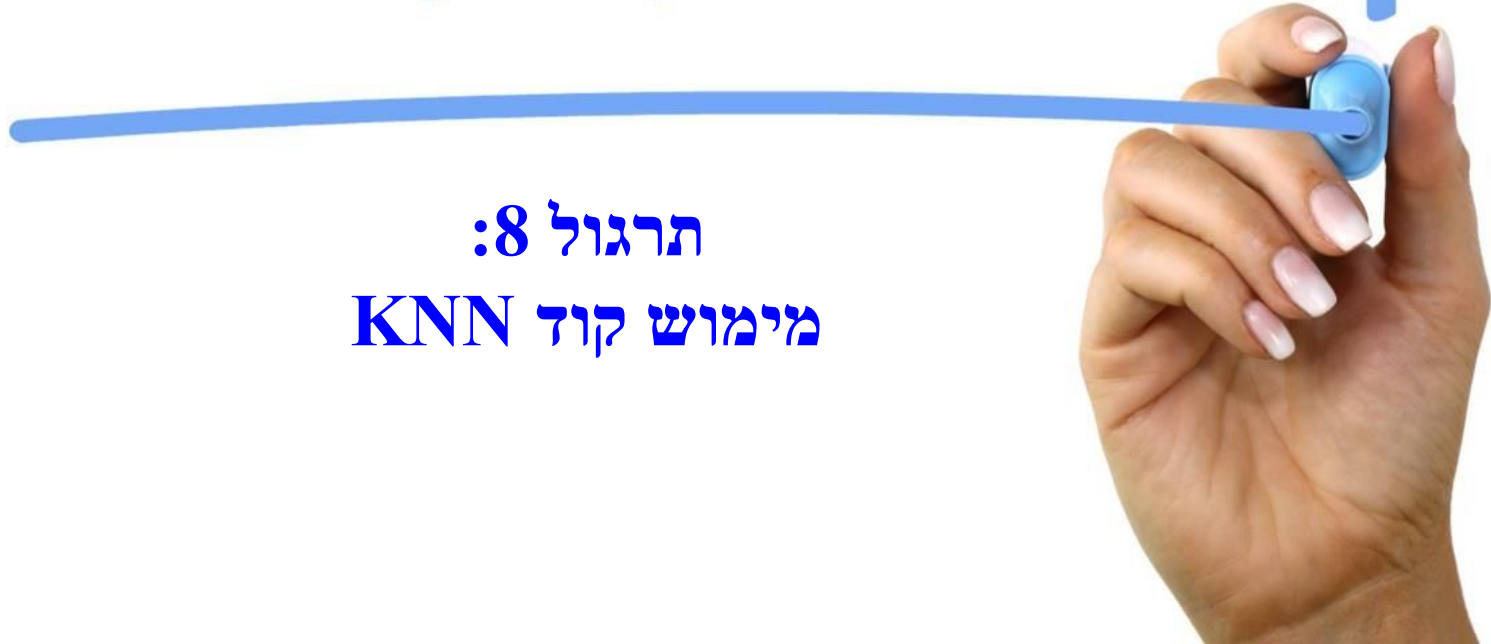
```
def KNN(train, test, lbl, K):  
    distances = []  
    for t, l in zip(train, lbl):  
        dist = euclidean_distance(test, t)  
        distances.append((t, dist, l[0]))  
    distances.sort(key=lambda dist: dist[1])  
    NN = []  
    for i in range(K):  
        NN.append(distances[i])  
    return NN
```

סיכום

KNN הוא אלגוריתם עצלן לימודית כי בניגוד ל ANN שקודם לומד ולאחר מכן יכול לסווג. אלגוריתם KNN לומד ומסווג רק כאשר מתקבלת בקשה. מכאן שזה גם אחד החסרונות של KNN. אלגוריתם זה יתקשה לתפקד כאשר מדובר על כמויות גדולות של אימון (למידה). בקיצור אלגוריתם עצלן! כאילו התלמיד לומד תוך כדי הבחינה.



TRAINING



תרגול 8:
מימוש קוד KNN