

K-NN אלגוריתם

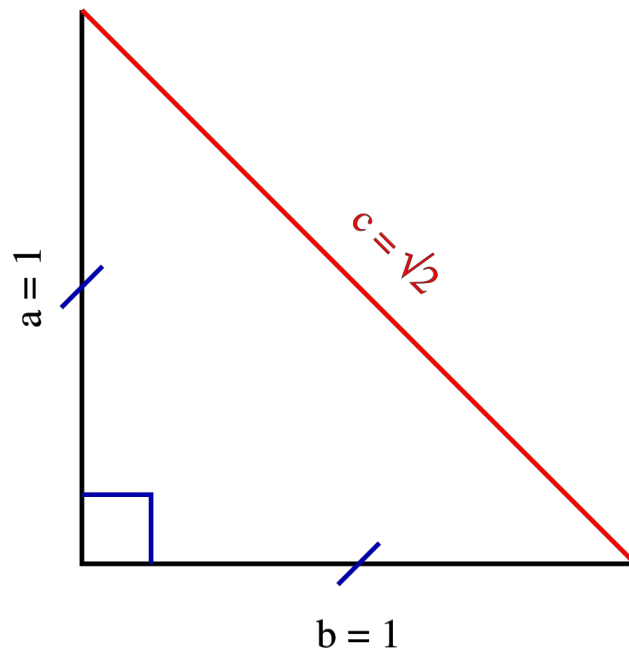
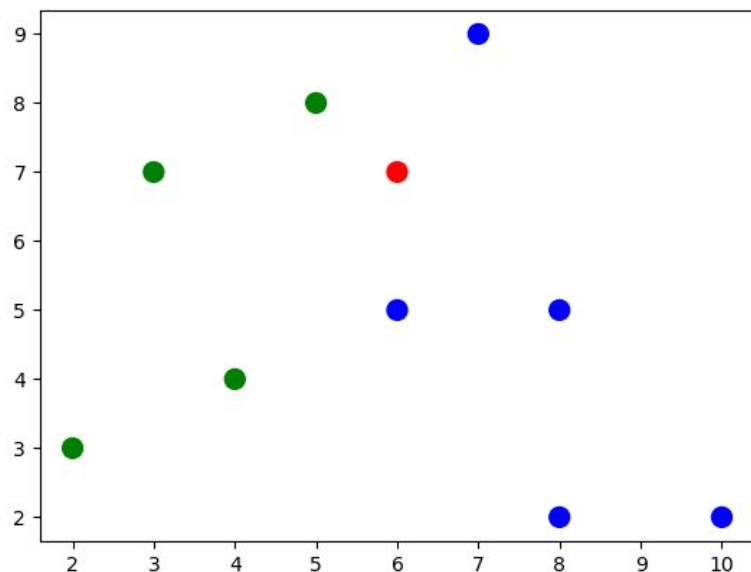
לפיתוח מכונה לומדת

K-Nearest Neighbors algorithm

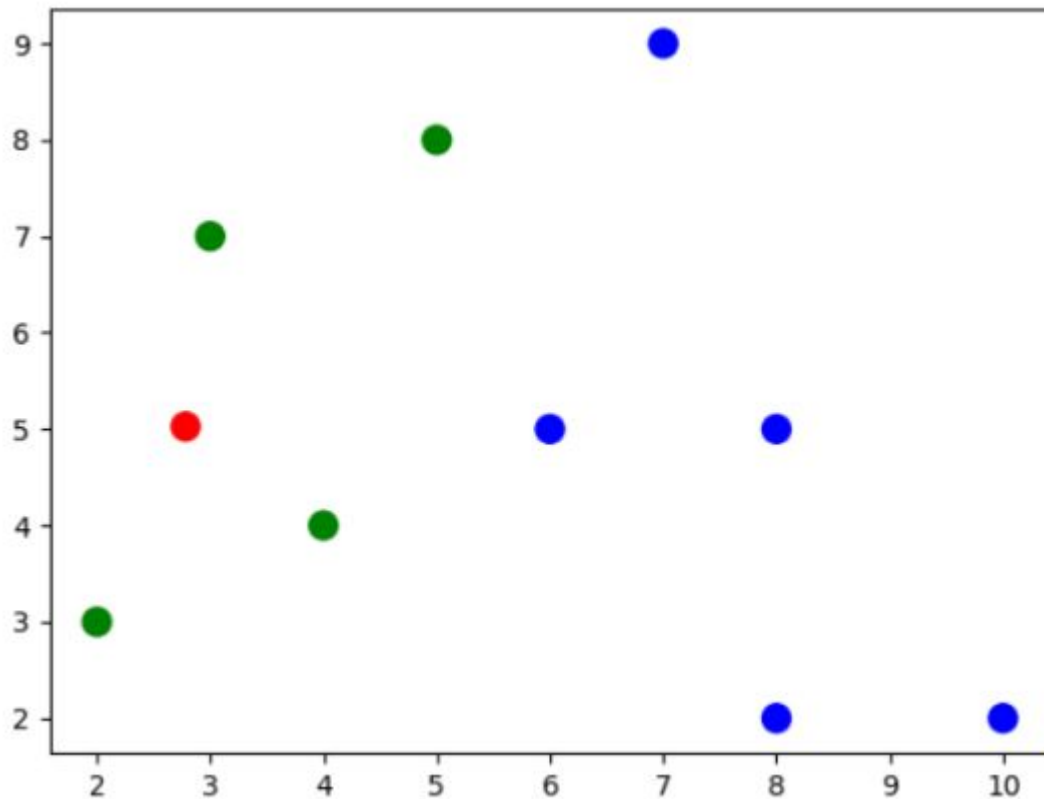
גדי הרמן

אלגוריתם K-Nearest Neighbors algorithm או בקיצור K-NN

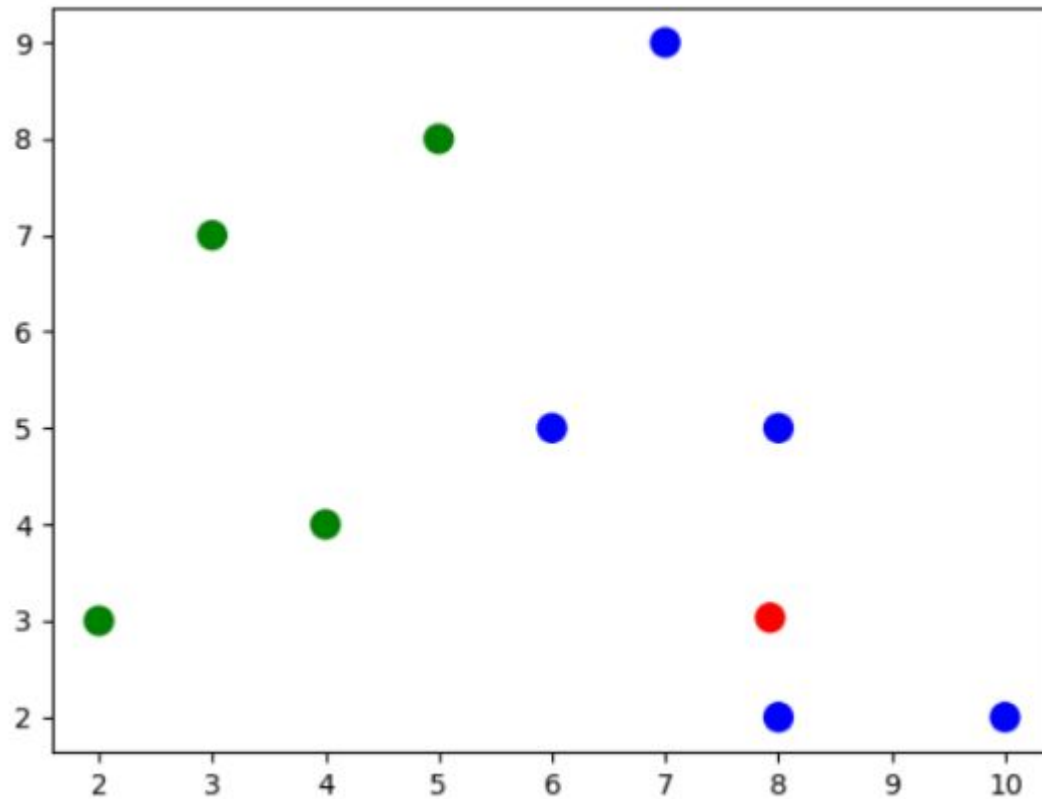
K-NN מבוסס בין היתר על מדידת מרחק אוקלידי בין נקודות במרחב או בשם הפשוט יותר שלו, שימוש בפיתגורס כדי לחשב את היתר במשולש ישר זווית.



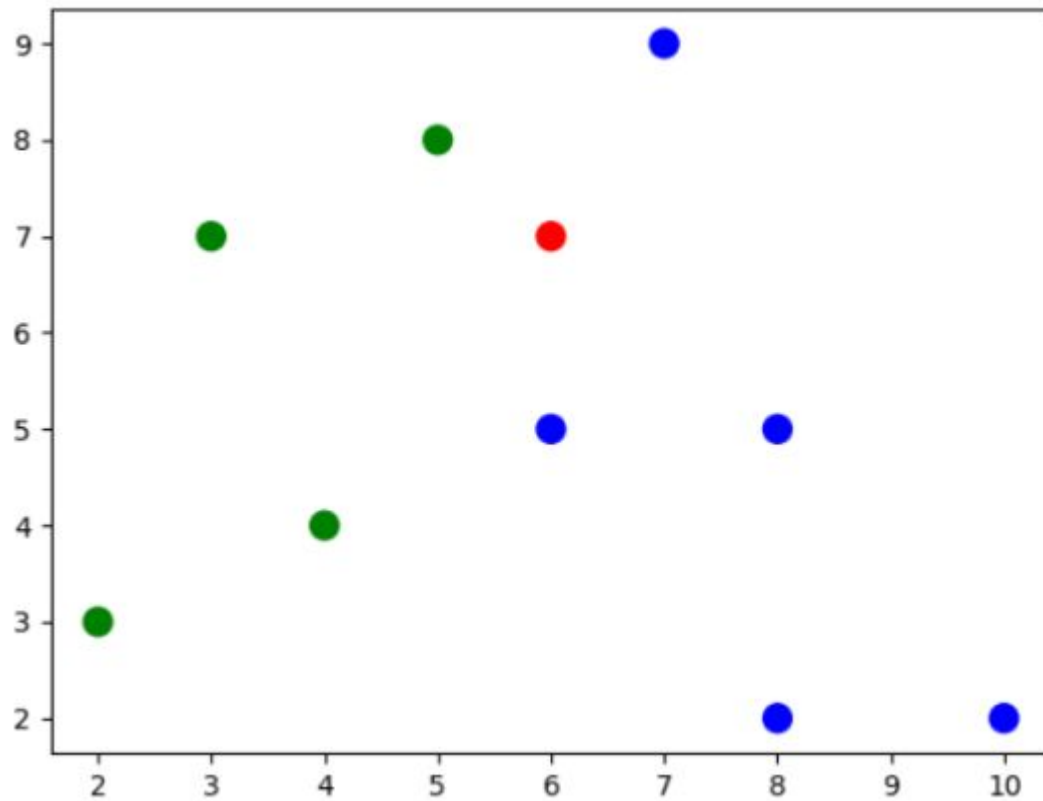
מי החברים שלי?



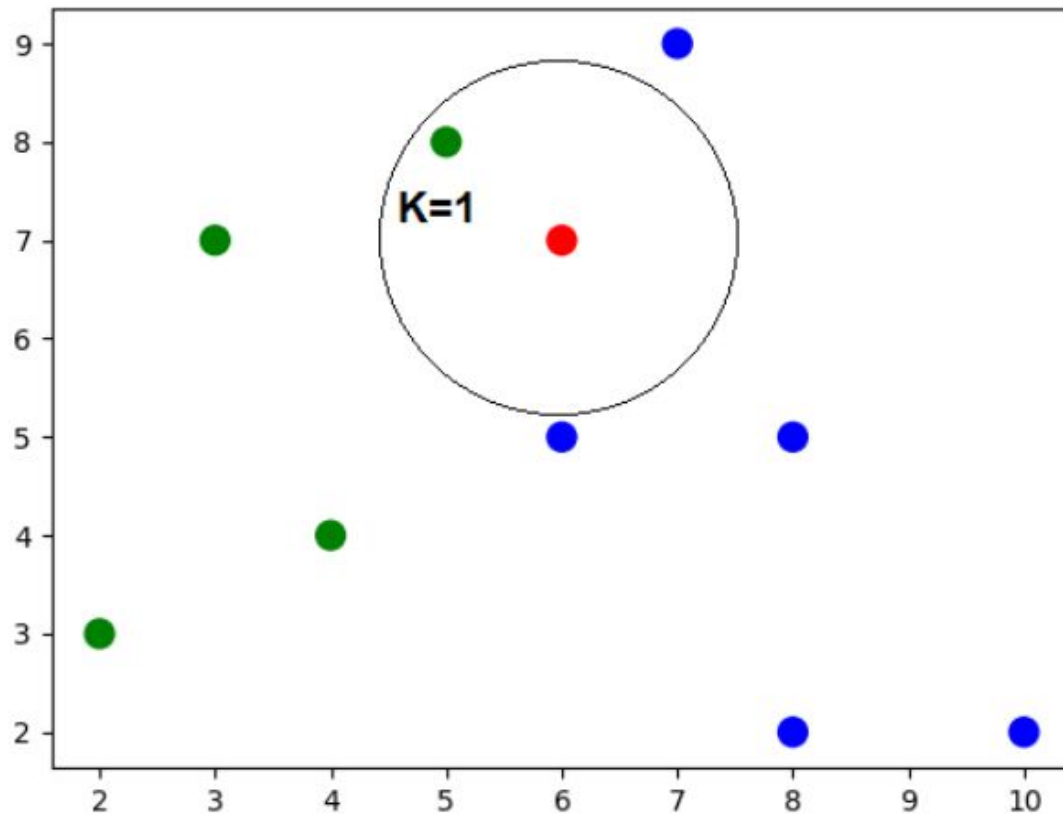
מי החברים שלי?



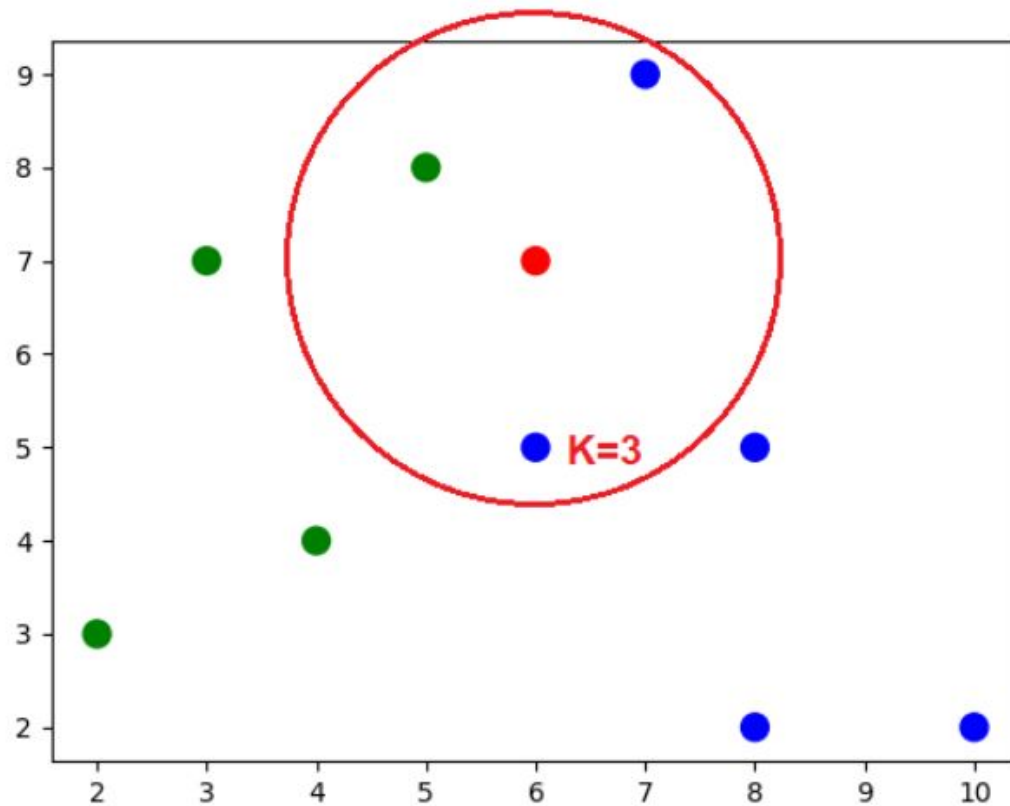
מי החברים שלי?



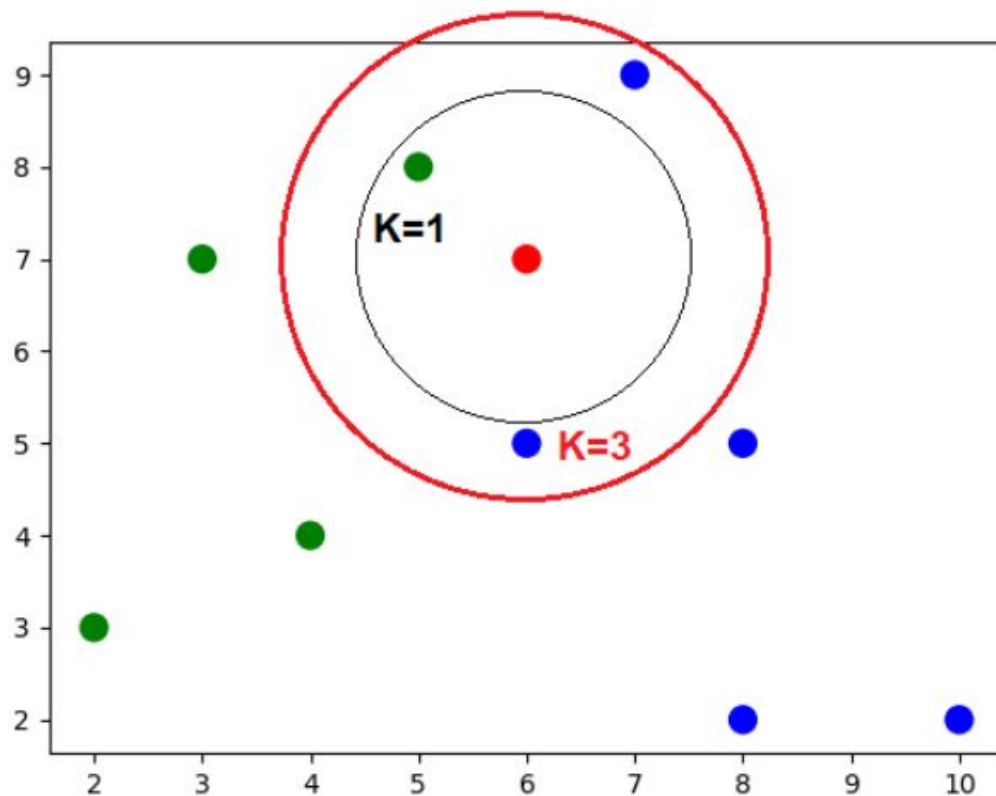
מי החברים שלי? $K=1$



מי החברים שלי? $K=3$

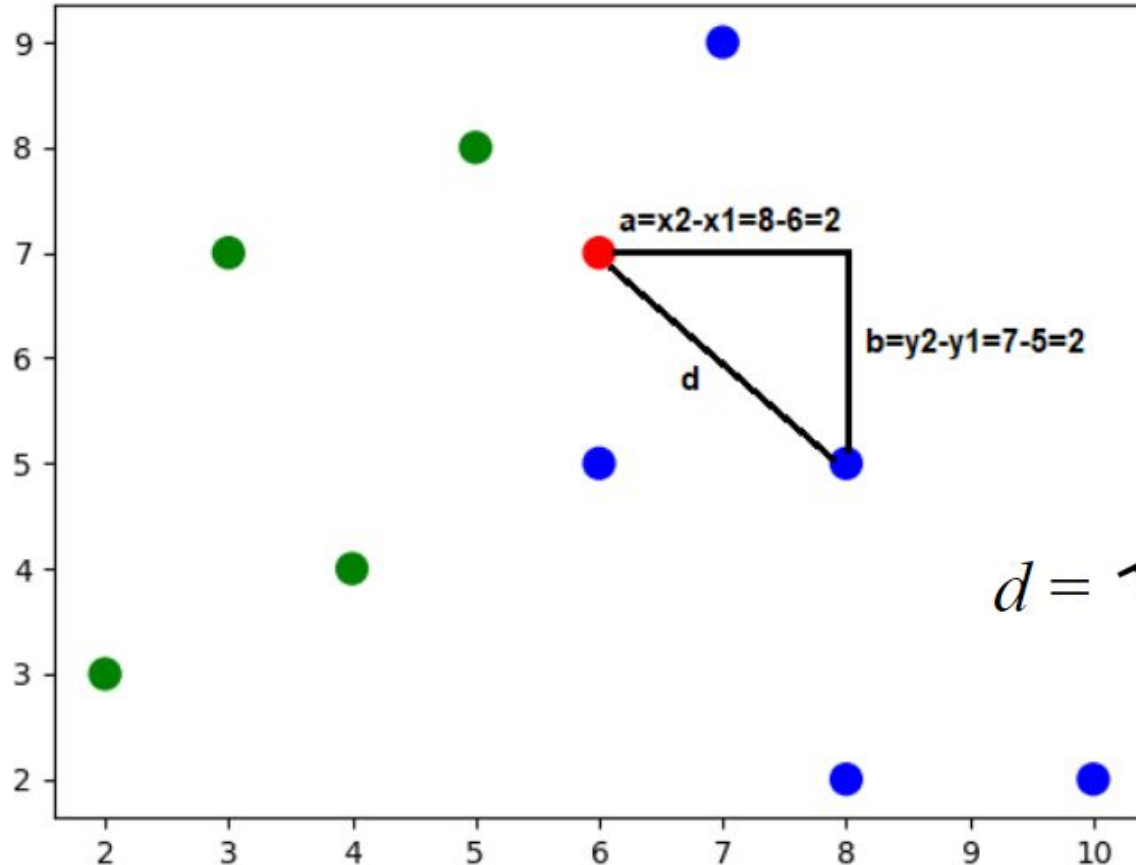


מי החברים שלי? $K=1$ או $K=3$



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.



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

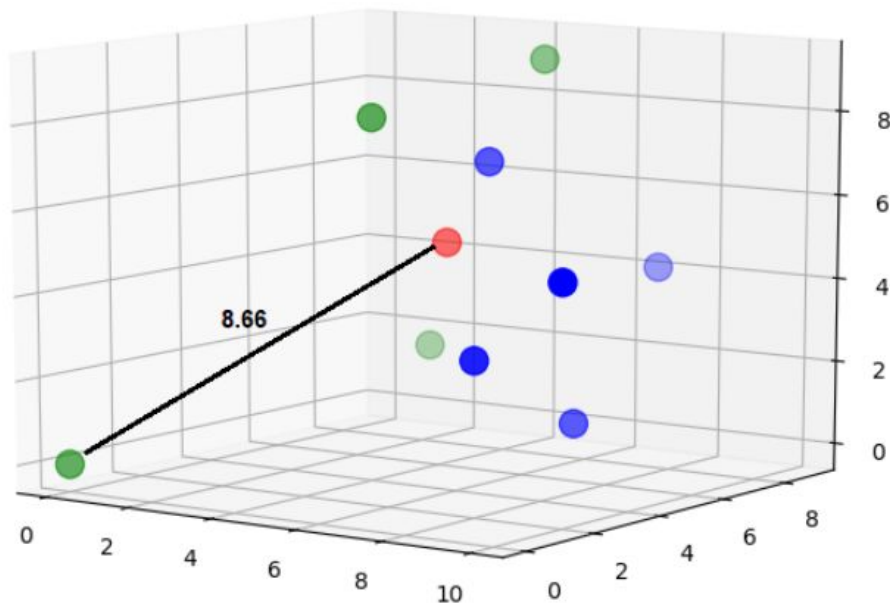
מימוש מרחק אוקלידי בין נקודות בקוד

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

```
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)
```

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.

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, התגית ואת ערכי הנקודה שנבדקה.
- לאחר סיום מדידת כל הנקודות נמיין את המערך לפי המרחקים.
- לבסוף נעבור על הלולאה הממוינת ונחלץ ממנה את K האיברים שאותם אנו רוצים לקבל.
- הפעולה תחזיר פרמטר אחד שהוא מספר הקבוצה.

מימוש הפעולה predict בתוך KNN

```
def takeSecond(elem):  
    return elem[1]  
  
def predict(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=takeSecond)  
    print(distances)  
    neighbors = []  
    for i in range(K):  
        neighbors.append(distances[i])  
    out = [row[-1] for row in neighbors]  
    return max(out, key=out.count)
```

מימוש הפעולה predict בתוך KNN

```
def takeSecond(elem):  
    return elem[1]  
  
def predict(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=takeSecond)  
    print(distances)  
    neighbors = []  
    for i in range(K):  
        neighbors.append(distances[i])  
    out = [row[-1] for row in neighbors]  
    return max(out, key=out.count)
```

```
f takeSecond(elem):  
    return [[array([5., 8.]), 1.4142135623730951, 1], [array.  
special variables  
function variables  
0: [array([5., 8.]), 1.4142135623730951, 1]  
1: [array([6., 5.]), 2.0, 2]  
2: [array([7., 9.]), 2.23606797749979, 2]  
3: [array([8., 5.]), 2.8284271247461903, 2]  
4: [array([3., 7.]), 3.0, 1]  
5: [array([4., 4.]), 3.605551275463989, 1]  
6: [array([8., 2.]), 5.385164807134504, 2]  
7: [array([2., 3.]), 5.656854249492381, 1]  
8: [array([10., 2.]), 6.4031242374328485, 2]  
len(): 9  
f predict  
dista  
for t  
d  
d  
distances.sort(key=takeSecond)  
print(distances)
```


מימוש הפעולה predict בתוך KNN

```
def takeSecond(elem):  
    return elem[1]  
  
def predict(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=takeSecond)  
    print(distances)  
    neighbors = []  
    for i in range(K):  
        neighbors.append(distances[i])  
    out = [row[-1] for row in neighbors]  
    return max(out, key=out.count)
```

```
distances.append([t, dist, l[0]])  
distance [[array([5., 8.]), 1.4142135623730951, 1], [array...  
print(di special variables  
neighbor function variables  
0: [array([5., 8.]), 1.4142135623730951, 1]  
1: [array([6., 5.]), 2.0, 2]  
2: [array([7., 9.]), 2.23606797749979, 2]  
for i in len(): 3  
    neighbors.append(distances[i])  
out = [row[-1] for row in neighbors]
```

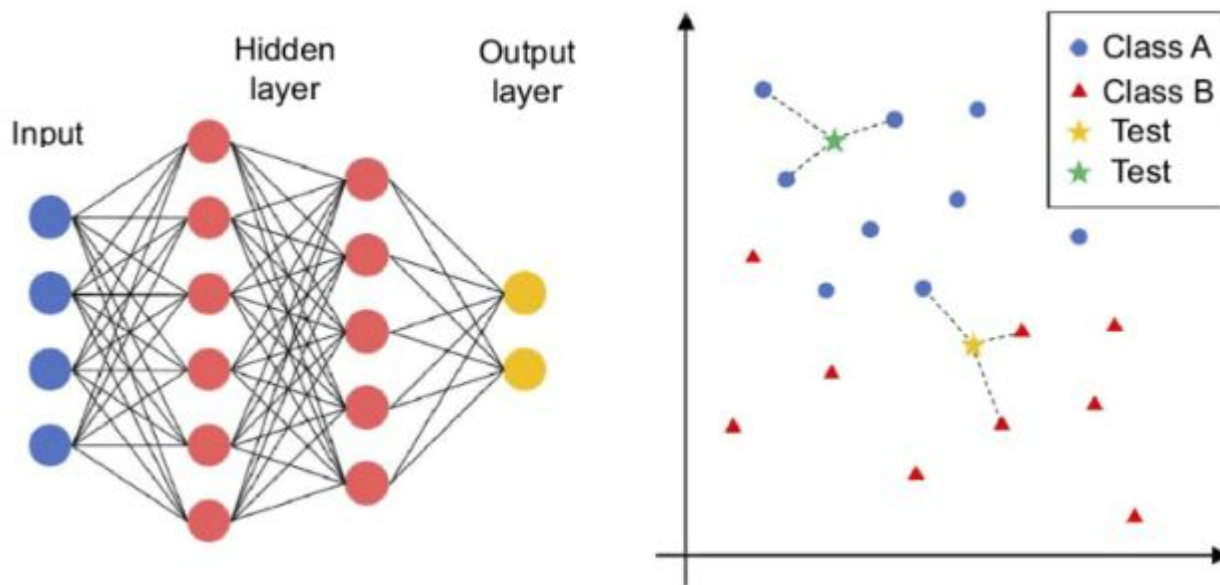

מימוש הפעולה predict בתוך KNN

```
def takeSecond(elem):  
    return elem[1]  
  
def predict(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=takeSecond)  
    print(distances)  
    neighbors = []  
    for i in range(K):  
        neighbors.append(distances[i])  
    out = [row[-1] for row in neighbors]  
    return max(out, key=out.count)
```

```
p  
n  
f  
[1, 2, 2]  
> special variables  
> function variables  
> 0: 1  
> 1: 2  
> 2: 2  
len(): 3  
distances[i])  
out = [row[-1] for row in neighbors]
```

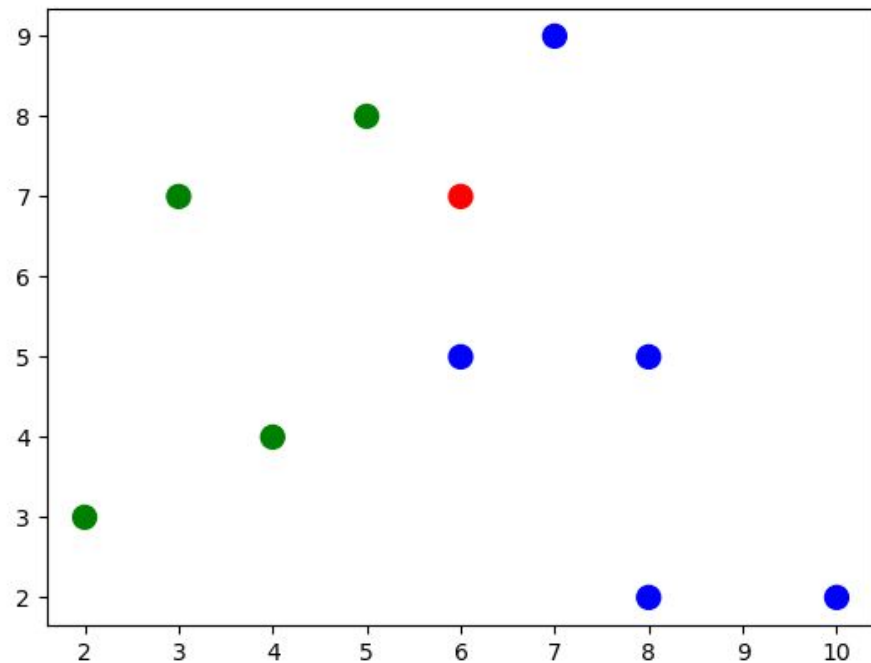
סיכום

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.



נבחן 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()
```



TRAINING



תרגול 8:
מימוש קוד KNN