

דו"ח מסכם תרגיל בית 4 – מבוא ללמידה ממוכנת:

מגישים: איתי גיא, ת.ז. – 305104184

רועי נתן, ת.ז. - 203842885

א. קריאת קובץ הנתונים *iris.data* – בוצע ע"י קריאת הקובץ שורה אחרי שורה והכנסת כל שורה לאובייקט *FeatureVector* שמתאר כל תצפית בנתונים אלו – פיצ'רים של עלי הכותרת והתיוג של השורה (המחרוזות של המחלקה שאליה היא מתאימה)

```
6 /**
7  * Handling FeatureVector operations
8  * @author Itay_Guy
9  */
10 public class FeatureVector {
11     private List<Double> _feature_vector; // numeric data
12     private String _name; //data labeling
13
14     public FeatureVector() {
15         this._feature_vector = new ArrayList<Double>();
16         this._name = new String();
17     }
18
19     public FeatureVector(String[] features) {
20         this();
21         for (String f : features) {
22             try {
23                 double d = Double.parseDouble(f);
24                 this._feature_vector.add(d);
25             } catch (Exception ex) {
26                 _name = new String(f.trim());
27             }
28         }
29     }
}
```

Training-set size = 120 , Test-set size = 30

ב. קידוד פונקציית הפיצול שמכינה את הנתונים בקובץ לאימון מודל – תוספת הפרמטרים הם בכדי למחזר את הפונקציה גם לטובת המחלקה של *DT* – באופן כללי כל הטיפול בנתונים מתבצע במחלקה בשם *Dataset*.

```

83  /**
84  * prepare data to the model usage
85  * @param split_perc    training-set size percentage to split
86  * @param uniform        selection distribution type - if false there is random type
87  * @param headers        usage for DT-ID3 only and it says if there are headers to this data file
88  * @throws IOException
89  */
90  public void prepare_data(int split_perc, boolean uniform, boolean headers) throws IOException {
91      List<String> lines = Files.readAllLines(this._data_name); // taking data from Path file you provided in the constructor
92      // handling headers:
93      if(headers) {
94          this._headers = new ArrayList<String>();
95          String[] variables = lines.get(0).split(",");
96          int i = 0;
97          for (String var : variables) {
98              this._headers.add(new String(var.trim()));
99              i++;
100          }
101          lines.remove(0); // remove the labeling from headers array
102      }
103      int lines_amount = lines.size();
104      int set_perc = (split_perc*lines.size())/100;
105      List<Integer> seq = new ArrayList<Integer>();
106      for (int i = 0; i < lines_amount; i++) {
107          seq.add(i);
108      }
109      List<Integer> permut = this.make_permut(seq, lines_amount); // make randomized data lines
110      if(!uniform) {
111          // this is not uniform distribution -> take the first split_perc data lines to training set and the remain to test set
112          for (int i = 0; i < lines_amount; i++) {
113              String[] data_line = lines.get(permut.get(i)).split(",");
114              if(i >= set_perc) {
115                  this._test_set.add(new FeatureVector(data_line));
116              } else {
117                  this._train_set.add(new FeatureVector(data_line));
118              }
119          }
120      } else {
121          // this is distributed uniformly -> handling the issue there is same amounts in the training set:
122          Map<String,Integer> types = new HashMap<String,Integer>();

```

```

123     for (String line : lines) {
124         String[] features = line.split(",");
125         String name = features[features.length - 1].trim();
126         if(types.containsKey(name)) {
127             types.put(name,types.get(name).intValue()+1);
128         } else {
129             types.put(name,1);
130         }
131     }
132     int uniform_chunk = lines_amount/types.keySet().size();
133     for (String key : types.keySet()) {
134         if(types.get(key).intValue() > uniform_chunk) {
135             types.put(key,types.get(key).intValue()-(types.get(key).intValue()-uniform_chunk));
136         }
137     }
138     boolean[] bucket = new boolean[lines_amount];
139     for (int i = 0; i < bucket.length; i++) {
140         bucket[i] = false;
141     }
142     for (int i = 0; i < lines_amount; i++) {
143         String[] features = lines.get(permut.get(i)).split(",");
144         String line_name = features[features.length - 1].trim();
145         if(this._train_set.size() < set_perc){
146             if(types.get(line_name).intValue() > 0) {
147                 types.put(line_name,types.get(line_name).intValue() - 1);
148                 this._train_set.add(new FeatureVector(features));
149                 bucket[permut.get(i)] = true;
150             }
151         } else {
152             break;
153         }
154     }
155     for (int i = 0; i < bucket.length; i++) {
156         if(bucket[i] == false) {
157             this._test_set.add(new FeatureVector(lines.get(i).split(",")));
158         }
159     }
160 }
161 this.prepare_class_mapping(this._train_set,this._test_set); // build class tag mapping
162 System.out.println("Training-set size = " + this._train_set.size() + " , Test-set size = " + this._test_set.size());
163 }

```

ג. פונקציית המרחק בין אירוסים נכתבה במחלקה *FeatureVector* כאשר היא מממשת באופן כללי את פונקציית מינקובסקי כאשר $p = 2$ זה מרחק אוקלידי לכן ניתן לראות אותו בפעולה עבור p -ים שונים שזה בעצם עבור סוגים שונים של פונקציות מרחק מאותה המשפחה.

```

39  /**
40   * Minkovski distance
41   * @param other FeatureVector to comparison with
42   * @param p if p=1 -> (Manhattan distance), if p=2 -> (euclidian distance),...
43   * @return the distance number
44   */
45  public double distance_p(FeatureVector other,int p) {
46      double sum = 0.0;
47      for (int i = 0; i < this._feature_vector.size(); i++) {
48          sum += Math.pow(Math.abs(this.getFeature(i) - other.getFeature(i)),p);
49      }
50      return Math.pow(sum,1.0/p);
51  }

```

7. פונקציית סיווג שמחזירה את k השכנים הקרובים ביותר לפרח המתקבל כפרמטר נכתבה במחלקה KNN שמממשת את המודל.

```
27  /**
28   * pick k neighbors
29   * @param data data-set object
30   * @param f sample from the test-set
31   * @param k for picking the k-neighbors
32   * @return k neighbors into array
33   */
34  public List<FeatureVector> make_subset(Dataset data, FeatureVector f, int k) {
35      List<FeatureVector> train_set = data.getTrainingset();
36      Collections.sort(train_set, new Comparator<FeatureVector>() { // sorting bottom up
37          @Override
38          public int compare(FeatureVector o1, FeatureVector o2) {
39              double x = f.distance_p(o1, P);
40              double y = f.distance_p(o2, P);
41              if(x < y) return -1;
42              else if(x > y) return 1;
43              else return 0;
44          }
45      });
46      List<FeatureVector> k_neighbors = new ArrayList<FeatureVector>();
47      for (int i = 0; i < k; i++) {
48          k_neighbors.add(new FeatureVector(train_set.get(i)));
49      }
50      return k_neighbors;
51  }
```

למשל עבור $k = 3$ נקבל:

```
For Test Example: Feature vector = [6.7,3.0,5.2,2.3,Label:Iris-virginica]
Neighbor 0: Feature vector = [6.9,3.1,5.1,2.3,Label:Iris-virginica]
Neighbor 1: Feature vector = [6.5,3.0,5.2,2.0,Label:Iris-virginica]
Neighbor 2: Feature vector = [6.9,3.1,5.4,2.1,Label:Iris-virginica]
```

קביעת סוג האירוס עפ"י המחלקה של הרוב מתוך k השכנים שהתקבלו:

```
53  /**
54   * finding the class that the sample from test-set is belong to
55   * @param best_k    array of best k neighbor from "make_subset"
56   * @param weighted  true if need weighted technique and false otherwise
57   * @return classified class for the test-set's sample
58   */
59  public String choose_best(List<FeatureVector> best_k,boolean weighted) { // best_k is already sorted
60      int k = best_k.size();
61      Map<String, HashMap<String,Double>> freq = new HashMap<String,HashMap<String,Double>>();
62      for (int i = 0; i < k; i++) {
63          double eps = 0.0;
64          if(weighted) {
65              eps = 1.0*k-i; // weight for each neighbor by their distance from f
66          }
67
68          HashMap<String,Double> h = new HashMap<String,Double>();
69          h.put("weight",1.0 + eps);
70          freq.put(best_k.get(i).getName() + "@" + i, h);
71      }
72      Map<String,Double> counter = new HashMap<String,Double>();
73      for (String key : freq.keySet()) {
74          HashMap<String,Double> h = freq.get(key);
75          double val = h.get("weight").doubleValue();
76          key = key.split("@")[0];
77          if(!counter.containsKey(key)) {
78              counter.put(key,val);
79          } else {
80              counter.put(key,counter.get(key) + val);
81          }
82      }
83      String max_label = "";
84      double max_appearances = 0.0;
85      for (String key : counter.keySet()) {
86          double val = counter.get(key).doubleValue();
87          if(val > max_appearances) {
88              max_appearances = val;
89              max_label = key;
90          }
91      }
92      return max_label;
93  }
```

- פרמטר *weighted* נכנס בכדי להשתמש בקוד עבור הדרישה המתקדמת בתרגיל.

ה. פונקציה הקובעת את אחוז הדיוק של קבוצת הבדיקה – הפונקציה כתובה במחלקה *KNN* שמממשת את המודל.

```
156 /**
157  * classification for same data examples
158  * @param data Data-set
159  * @param k number of neighbors
160  * @param weighted need to use weighted method or not
161  * @param need_confution_matrix_disp if we would like to display the confusion matrix
162  * @param validate_training_set true if we need to classify the training-set
163  * @return match accuracy of this model
164  */
165 public double classify_by_accuracy(Dataset data,int k,boolean weighted,boolean need_confution_matrix_disp,boolean validate_training_set) {
166     Map<String,HashMap<String,Integer>> confusion_matrix = this.classify(data, k, weighted,validate_training_set);
167     if(need_confution_matrix_disp) {
168         System.out.println(KNN.disp_confusion_matrix(confusion_matrix));
169     }
170     int total = 0;
171     int correct = 0;
172     //counting how many is correct predictions:
173     for (String key : confusion_matrix.keySet()) {
174         for(String name : confusion_matrix.get(key).keySet()){
175             if(key.equals(name)) {
176                 correct += confusion_matrix.get(key).get(name);
177             }
178             total += confusion_matrix.get(key).get(name);
179         }
180     }
181     if(!need_confution_matrix_disp) {
182         System.out.println("*****");
183     }
184     double accuracy = correct*100.0/total;
185     System.out.println("Accuracy = " + accuracy);
186     return accuracy;
187 }
```

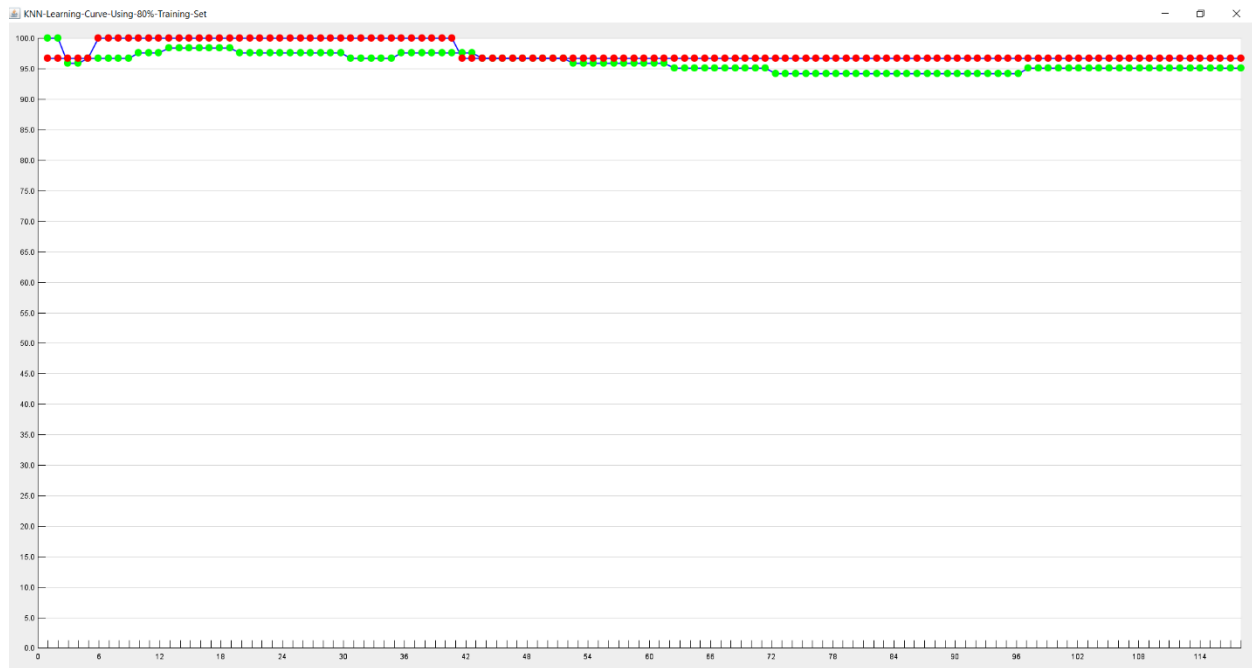
```

95  /**
96  * making classification for some data - training-set or test-set
97  * @param data Data-set
98  * @param k number of neighbors
99  * @param weighted need to use weighted method for finding the class between k neighbors
100 * @param validate_training_set if need to classify the training-set for learning curve for example
101 * @return confusion matrix
102 */
103 public Map<String,HashMap<String,Integer>> classify(Dataset data,int k,boolean weighted,boolean validate_training_set) {
104     List<FeatureVector> test = null;
105     if(validate_training_set) {
106         System.out.println("*****\nClassification: Training-Set");
107         test = data.getTrainingset();
108     } else {
109         System.out.println("*****\nClassification: Test-Set");
110         test = data.getTestSet();
111     }
112     Map<String,HashMap<String,Integer>> confusion_matrix = new HashMap<String,HashMap<String,Integer>>();
113     System.out.println("\nUsing k = " + k + " and Weighted-Sum = " + weighted + "\n*****");
114     for (FeatureVector f_test : test) {
115         System.out.println("For Test Example: " + f_test);
116         List<FeatureVector> best_k = this.make_subset(data,f_test, k);
117         int i = 0;
118         for (FeatureVector fv : best_k) {
119             System.out.println("\tNeighbor " + i + ": " + fv);
120             i++;
121         }
122         String res = this.choose_best(best_k, weighted); // weighted = true or false
123         if(res.equals(f_test.getName())) {
124             System.out.println("\t+ " + res + " Classified correctly");
125         } else {
126             System.out.println("\t- Truth value: " + f_test.getName() + " => but classified as: " + res);
127         }
128         // build confusion matrix
129         HashMap<String,Integer> h = new HashMap<String,Integer>();
130         if(confusion_matrix.keySet().contains(f_test.getName())){
131             if(confusion_matrix.get(f_test.getName()).keySet().contains(res)) {
132                 confusion_matrix.get(f_test.getName()).put(res,confusion_matrix.get(f_test.getName()).get(res) + 1);
133             } else {
134                 h.put(res,1);
135             }
136         }
137         confusion_matrix.put(f_test.getName(),h);
138     }
139 }
140
141 }
142 ArrayList<String> titles = new ArrayList<String>();
143 for (String name : confusion_matrix.keySet()) {
144     titles.add(name);
145 }
146 for (String name : confusion_matrix.keySet()) {
147     for (int i = 0; i < titles.size(); i++) {
148         if(!confusion_matrix.get(name).containsKey(titles.get(i))) {
149             confusion_matrix.get(name).put(titles.get(i),0);
150         }
151     }
152 }
153 return confusion_matrix;
154 }
155 ---

```

ז. ניתן לראות בגרף שצויר למטה את ההשוואה בין ה- $test - set$ ל- $training - set$ כאשר ב- $test - set$ (העקום האדום) גילינו שה- k האופטימלי מתקבל בטווח בין $k = 6$ ל- $k = 42$ עם דיוק של 100% ו-80% $train$ (אם צריך לתת k מסוים אז נבחר את $k = 7$).

ח. בנוסף, גילינו ע"י ציור 2 העקומים למטה שבאזור של $k = 5$ מתקבלת התאמה של הפונקציה. כאשר לפני כן ב- k נמוך נראה אזור של $underfitting$ מכיוון שאין דיוק של ה- $test - set$ כמו של ה- $training - set$ ובאזור שאחרי עד $k = 42$ מתקבל $overfitting$ מכיוון שדיוק ה- $training - set$ עולה על הדיוק של ה- $test - set$. בהמשך הגרף החל מהאזור של $k = 42$ נראה שהאלגוריתם מתייצב בעיקר באזור שבין $k = 42$ עד $k = 60$ (בערך) ומדייק באופן דומה ל- $training - set$.



דגימה של קבוצת הבדיקה – הדפסת תוצאת הסיווג לכל אירוס. הסיווג הנכון והאם הוא היה נכון וכן בסוף את הדיוק הכללי של

המסווג:

```
Distance Method: Minkovski Distance with P = 2
Training-set size = 120 , Test-set size = 30
Using k = 3 and Weighted-Sum = false
*****
For Test Example: Feature vector = [6.2,3.4,5.4,2.3,Label:Iris-virginica]
  Neighbor 0: Feature vector = [6.3,3.4,5.6,2.4,Label:Iris-virginica]
  Neighbor 1: Feature vector = [6.4,3.2,5.3,2.3,Label:Iris-virginica]
  Neighbor 2: Feature vector = [6.5,3.2,5.1,2.0,Label:Iris-virginica]
+ Iris-virginica Classified correctly
For Test Example: Feature vector = [5.2,4.1,1.5,0.1,Label:Iris-setosa]
  Neighbor 0: Feature vector = [5.5,4.2,1.4,0.2,Label:Iris-setosa]
  Neighbor 1: Feature vector = [5.1,3.8,1.5,0.3,Label:Iris-setosa]
  Neighbor 2: Feature vector = [5.4,3.7,1.5,0.2,Label:Iris-setosa]
+ Iris-setosa Classified correctly
For Test Example: Feature vector = [5.1,3.8,1.6,0.2,Label:Iris-setosa]
  Neighbor 0: Feature vector = [5.1,3.8,1.5,0.3,Label:Iris-setosa]
  Neighbor 1: Feature vector = [5.1,3.7,1.5,0.4,Label:Iris-setosa]
  Neighbor 2: Feature vector = [5.0,3.6,1.4,0.2,Label:Iris-setosa]
+ Iris-setosa Classified correctly
For Test Example: Feature vector = [4.9,3.1,1.5,0.1,Label:Iris-setosa]
  Neighbor 0: Feature vector = [4.9,3.1,1.5,0.1,Label:Iris-setosa]
  Neighbor 1: Feature vector = [4.9,3.1,1.5,0.1,Label:Iris-setosa]
  Neighbor 2: Feature vector = [4.9,3.0,1.4,0.2,Label:Iris-setosa]
+ Iris-setosa Classified correctly
For Test Example: Feature vector = [4.6,3.1,1.5,0.2,Label:Iris-setosa]
  Neighbor 0: Feature vector = [4.6,3.2,1.4,0.2,Label:Iris-setosa]
  Neighbor 1: Feature vector = [4.7,3.2,1.6,0.2,Label:Iris-setosa]
  Neighbor 2: Feature vector = [4.8,3.1,1.6,0.2,Label:Iris-setosa]
+ Iris-setosa Classified correctly
...
For Test Example: Feature vector = [4.9,2.5,4.5,1.7,Label:Iris-virginica]
  Neighbor 0: Feature vector = [5.4,3.0,4.5,1.5,Label:Iris-versicolor]
  Neighbor 1: Feature vector = [5.5,2.6,4.4,1.2,Label:Iris-versicolor]
  Neighbor 2: Feature vector = [5.6,2.7,4.2,1.3,Label:Iris-versicolor]
- Truth value: Iris-virginica => but classified as: Iris-versicolor
For Test Example: Feature vector = [7.3,2.9,6.3,1.8,Label:Iris-virginica]
  Neighbor 0: Feature vector = [7.4,2.8,6.1,1.9,Label:Iris-virginica]
  Neighbor 1: Feature vector = [7.1,3.0,5.9,2.1,Label:Iris-virginica]
  Neighbor 2: Feature vector = [7.7,2.8,6.7,2.0,Label:Iris-virginica]
+ Iris-virginica Classified correctly
For Test Example: Feature vector = [4.6,3.6,1.0,0.2,Label:Iris-setosa]
  Neighbor 0: Feature vector = [4.6,3.4,1.4,0.3,Label:Iris-setosa]
  Neighbor 1: Feature vector = [5.0,3.5,1.3,0.3,Label:Iris-setosa]
  Neighbor 2: Feature vector = [4.4,3.2,1.3,0.2,Label:Iris-setosa]
+ Iris-setosa Classified correctly
For Test Example: Feature vector = [6.4,3.2,4.5,1.5,Label:Iris-versicolor]
  Neighbor 0: Feature vector = [6.3,3.3,4.7,1.6,Label:Iris-versicolor]
  Neighbor 1: Feature vector = [6.7,3.1,4.4,1.4,Label:Iris-versicolor]
  Neighbor 2: Feature vector = [6.7,3.1,4.7,1.5,Label:Iris-versicolor]
+ Iris-versicolor Classified correctly
*****
Confusion Matrix:
      Iris-versicolor  Iris-virginica  Iris-setosa
Iris-versicolor      9              0              0
Iris-virginica       1             10              0
Iris-setosa          0              0             10
*****
```

```

Iris-versicolor Precision = 100.0
Iris-virginica Precision = 90.9090909090909
Iris-setosa Precision = 100.0
*****
Iris-versicolor Recall = 90.0
Iris-virginica Recall = 100.0
Iris-setosa Recall = 100.0
*****
F-[Iris-versicolor]-measure = 94.73684210526316
F-[Iris-virginica]-measure = 94.73684210526316
F-[Iris-setosa]-measure = 100.0
*****
Accuracy = 96.66666666666667

```

דגימה של קבוצת האימון – הדפסת תוצאת הסיווג לכל אירוס, הסיווג הנכון והאם הוא היה נכון וכן בסוף את הדיקת הכללי של המסווג:

```

Distance Method: Minkovski Distance with P = 2
Training-set size = 120 , Test-set size = 30
Classification: Training-Set

Using k = 3 and Weighted-Sum = false
*****
For Test Example: Feature vector = [4.8,3.4,1.6,0.2,Label:Iris-setosa]
  Neighbor 0: Feature vector = [4.8,3.4,1.6,0.2,Label:Iris-setosa]
  Neighbor 1: Feature vector = [5.0,3.4,1.5,0.2,Label:Iris-setosa]
  Neighbor 2: Feature vector = [5.0,3.4,1.6,0.4,Label:Iris-setosa]
  + Iris-setosa Classified correctly
For Test Example: Feature vector = [7.3,2.9,6.3,1.8,Label:Iris-virginica]
  Neighbor 0: Feature vector = [7.3,2.9,6.3,1.8,Label:Iris-virginica]
  Neighbor 1: Feature vector = [7.4,2.8,6.1,1.9,Label:Iris-virginica]
  Neighbor 2: Feature vector = [7.2,3.2,6.0,1.8,Label:Iris-virginica]
  + Iris-virginica Classified correctly
For Test Example: Feature vector = [7.7,2.6,6.9,2.3,Label:Iris-virginica]
  Neighbor 0: Feature vector = [7.7,2.6,6.9,2.3,Label:Iris-virginica]
  Neighbor 1: Feature vector = [7.7,2.8,6.7,2.0,Label:Iris-virginica]
  Neighbor 2: Feature vector = [7.7,3.0,6.1,2.3,Label:Iris-virginica]
  + Iris-virginica Classified correctly
For Test Example: Feature vector = [7.7,3.0,6.1,2.3,Label:Iris-virginica]
  Neighbor 0: Feature vector = [7.7,3.0,6.1,2.3,Label:Iris-virginica]
  Neighbor 1: Feature vector = [7.4,2.8,6.1,1.9,Label:Iris-virginica]
  Neighbor 2: Feature vector = [7.3,2.9,6.3,1.8,Label:Iris-virginica]
...

For Test Example: Feature vector = [6.7,3.1,4.7,1.5,Label:Iris-versicolor]
  Neighbor 0: Feature vector = [6.7,3.1,4.7,1.5,Label:Iris-versicolor]
  Neighbor 1: Feature vector = [6.9,3.1,4.9,1.5,Label:Iris-versicolor]
  Neighbor 2: Feature vector = [6.7,3.1,4.4,1.4,Label:Iris-versicolor]
  + Iris-versicolor Classified correctly
For Test Example: Feature vector = [5.4,3.4,1.5,0.4,Label:Iris-setosa]
  Neighbor 0: Feature vector = [5.4,3.4,1.5,0.4,Label:Iris-setosa]
  Neighbor 1: Feature vector = [5.4,3.4,1.7,0.2,Label:Iris-setosa]
  Neighbor 2: Feature vector = [5.2,3.4,1.4,0.2,Label:Iris-setosa]
  + Iris-setosa Classified correctly
*****

Confusion Matrix:

```

	Iris-versicolor	Iris-virginica	Iris-setosa
Iris-versicolor	39	0	0
Iris-virginica	3	36	0
Iris-setosa	0	0	42

```

*****

```

```

Iris-versicolor Precision = 100.0
Iris-virginica Precision = 92.3076923076923
Iris-setosa Precision = 100.0
*****
Iris-versicolor Recall = 92.85714285714286
Iris-virginica Recall = 100.0
Iris-setosa Recall = 100.0
*****
F-[Iris-versicolor]-measure = 96.2962962962963
F-[Iris-virginica]-measure = 96.2962962962963
F-[Iris-setosa]-measure = 100.0
*****
Accuracy = 97.5

```

2.

א. בפונקציית הדגימה (*prepare_data*) יש 2 מצבים – בהתאם לפרמטר *boolean uniform*.

הקוד לדגימה פרופורציונלית (מקרה בתוך הפונקציה *prepare_data*):

```

120     } else {
121         // this is distributed uniformly -> handling the issue there is same amounts in the training set:
122         Map<String,Integer> types = new HashMap<String,Integer>();
123         for (String line : lines) {
124             String[] features = line.split(",");
125             String name = features[features.length - 1].trim();
126             if(types.containsKey(name)) {
127                 types.put(name,types.get(name).intValue()+1);
128             } else {
129                 types.put(name,1);
130             }
131         }
132         int uniform_chunk = lines_amount/types.keySet().size();
133         for (String key : types.keySet()) {
134             if(types.get(key).intValue() > uniform_chunk) {
135                 types.put(key,types.get(key).intValue()-(types.get(key).intValue()-uniform_chunk));
136             }
137         }
138         boolean[] bucket = new boolean[lines_amount];
139         for (int i = 0;i < bucket.length;i++) {
140             bucket[i] = false;
141         }
142         for (int i = 0;i < lines_amount;i++) {
143             String[] features = lines.get(permut.get(i)).split(",");
144             String line_name = features[features.length - 1].trim();
145             if(this._train_set.size() < set_perc){
146                 if(types.get(line_name).intValue() > 0) {
147                     types.put(line_name,types.get(line_name).intValue() - 1);
148                     this._train_set.add(new FeatureVector(features));
149                     bucket[permut.get(i)] = true;
150                 }
151             } else {
152                 break;
153             }
154         }
155         for (int i = 0;i < bucket.length;i++) {
156             if(bucket[i] == false) {
157                 this._test_set.add(new FeatureVector(lines.get(i).split(",")));
158             }
159         }
160     }

```

ב. בפונקציית הסיווג יינתן משקל גבוהה יותר לאירוסים קרובים במידה ופרמטר *boolean weighted* יהיה שווה ל-*true*.

```
53  /**
54   * finding the class that the sample from test-set is belong to
55   * @param best_k    array of best k neighbor from "make_subset"
56   * @param weighted  true if need weighted technique and false otherwise
57   * @return classified class for the test-set's sample
58   */
59  public String choose_best(List<FeatureVector> best_k,boolean weighted) { // best_k is already sorted
60      int k = best_k.size();
61      Map<String, HashMap<String,Double>> freq = new HashMap<String,HashMap<String,Double>>();
62      for (int i = 0; i < k; i++) {
63          double eps = 0.0;
64          if(weighted) {
65              eps = 1.0*k-i; // weight for each neighbor by their distance from f
66          }
67
68          HashMap<String,Double> h = new HashMap<String,Double>();
69          h.put("weight",1.0 + eps);
70          freq.put(best_k.get(i).getName() + "@" + i, h);
71      }
72      Map<String,Double> counter = new HashMap<String,Double>();
73      for (String key : freq.keySet()) {
74          HashMap<String,Double> h = freq.get(key);
75          double val = h.get("weight").doubleValue();
76          key = key.split("@")[0];
77          if(!counter.containsKey(key)) {
78              counter.put(key,val);
79          } else {
80              counter.put(key,counter.get(key) + val);
81          }
82      }
83      String max_label = "";
84      double max_appearances = 0.0;
85      for (String key : counter.keySet()) {
86          double val = counter.get(key).doubleValue();
87          if(val > max_appearances) {
88              max_appearances = val;
89              max_label = key;
90          }
91      }
92      return max_label;
93  }
```

ג. ניתן לראות שע"י שימוש במרחק אוקלידי כפרמטר למרחק קיבלנו את התוצאות הבאות עבור $k = 3$:

	Iris-versicolor	Iris-virginica	Iris-setosa
Iris-versicolor	11	1	0
Iris-virginica	0	8	0
Iris-setosa	0	0	10

```

*****
Iris-versicolor Precision = 91.66666666666667
Iris-virginica Precision = 100.0
Iris-setosa Precision = 100.0
*****
Iris-versicolor Recall = 100.0
Iris-virginica Recall = 88.88888888888889
Iris-setosa Recall = 100.0
*****
F-[Iris-versicolor]-measure = 95.65217391304348
F-[Iris-virginica]-measure = 94.11764705882352
F-[Iris-setosa]-measure = 94.11764705882352
*****
Accuracy = 96.66666666666667

```

וע"י שימוש בפונקציית מינקובסקי אחרת - $p = 10$ קיבלנו את התוצאות הבאות עבור $k = 3$:

	Iris-versicolor	Iris-virginica	Iris-setosa
Iris-versicolor	10	0	0
Iris-virginica	0	9	0
Iris-setosa	0	0	11

```

*****
Iris-versicolor Precision = 100.0
Iris-virginica Precision = 100.0
Iris-setosa Precision = 100.0
*****
Iris-versicolor Recall = 100.0
Iris-virginica Recall = 100.0
Iris-setosa Recall = 100.0
*****
F-[Iris-versicolor]-measure = 100.0
F-[Iris-virginica]-measure = 100.0
F-[Iris-setosa]-measure = 100.0
*****
Accuracy = 100.0

```

כלומר, אכן ניתן לראות את השפעת השינוי על התוצאות באופן משמעותי מאוד בין מינקובסקי $p = 2$ מול מינקובסקי $p = 10$ (במאמר מוסגר ניתן לומר שעלות חישוב של מינקובסקי עם $p = 10$ לווקטורים גדולים יותר עלולה להיות עלות יקרה בזמן חישוב מודל KNN ולהאט משמעותית את זמן החישוב).

3.

א. מימוש עץ ההחלטה בוצע במחלקה בשם *DecisionTree* שנעזרת במחלקות *Node*, *Tree* ו-*FeatureVector* בכדי לייצג את הנתונים שבקובץ *votes.csv*.

* בפונקציה *fit* יש קריאה לפונקציית בנייה של העץ.

```
159 /**
160  * build new DT model
161  */
162 public void fit() {
163     List<FeatureVector> train = this._data.getTrainingset();
164     List<String> attributes = this._data.getAttributes(false); // false is for taking attribute without labels
165     String defc = this.majority_value(train); // pickup the majority label in the training-set samples
166     this._hypotheses_tree = this.build_decision_tree(train, attributes, defc);
167 }
168
169 private Tree build_decision_tree(List<FeatureVector> examples, List<String> attributes, String defc) {
170     if(examples == null || examples.isEmpty()) {
171         return new Tree(new Node(this._data.convert_class(defc), attributes));
172     } else if(this.has_homogenous_class(examples) || attributes.isEmpty()) {
173         return new Tree(new Node(this._data.convert_class(this.majority_value(examples)), attributes));
174     } else {
175         String best = this.choose_attribute(examples, attributes);
176         Tree tree = new Tree(new Node(best, attributes));
177         for(int vi : DecisionTree.BINARY_SET) { // build node options
178             List<FeatureVector> examplesi = this.choose_subtrain(examples, vi, best);
179             Tree subtree = this.build_decision_tree(examplesi, this.subtract_attr(attributes, best), this.majority_value(examples));
180             tree.add_branch(subtree, vi, best); // add new sub-tree to this vi side in this node
181         }
182         return tree;
183     }
184 }
```

ב. חלוקת הנתונים לקבוצת אימון וקבוצת בדיקה: שימוש בפונקציה שהוצגה הנ"ל בשם *prepare_data* כאשר הפרמטר *boolean header* שווה ל-*true* זה אומר שיש שמות לעמודות והן נאספות בפונקציה לאובייקט נפרד לצורך ביצוע מיפוי בין אינדקס לשם בהמשך.

Training-set size = 279 , Test-set size = 156

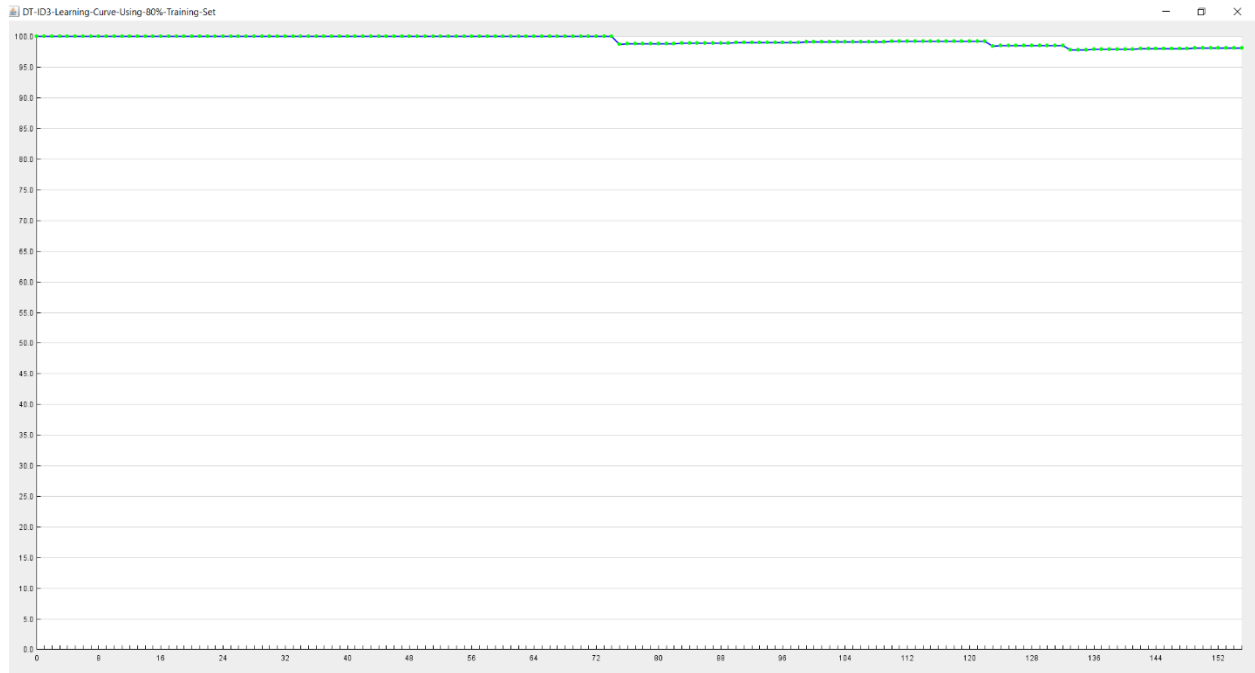
ג. פונקציית פיצול מקסימום *Gain* נכתבה במחלקה *DecisionTree* ולוקחת את הפיצ'ר בעל ה-*Gain* המקסימאלי באופן הפוך לאנטרופיה.

פונקציית ה-Gain:

```
296 /**
297  * computing Gain function
298  * @param train current training-set
299  * @param attributes current attributes at this node to compute the best from this
300  * @return best attribute by gain function
301  */
302 public String select_by_gain(List<FeatureVector> train, List<String> attributes) {
303     String target_name = train.get(0).getName();
304     double class_entropy = this.calc_entropy(train, target_name, true);
305     double max_cond_ent = 0.0;
306     String future_attr = attributes.get(0);
307     for (String cond_attr : attributes) {
308         double prob = this.cond_prob_by_class(train, target_name, cond_attr); // probability
309         double attr_entropy = this.calc_entropy(train, cond_attr, false); // attribute entropy value
310         double cond_entropy = class_entropy - prob*attr_entropy; //confitinal entropy computation
311         if(cond_entropy > max_cond_ent) { // comparison to pick the best
312             max_cond_ent = cond_entropy;
313             future_attr = cond_attr;
314         }
315     }
316     return future_attr;
317 }
```

השוואה בין עץ החלטה עם פונקציית פיצול מקסימום $Gain$ לבין KNN עם k אופטימלי:

הגרף של הרצת DT – דיוק כפונקציה של כמות תצפיות:



ניתן לראות בגרף הנ"ל שהאלגוריתם חוזה באופן מרשים את נתונים קבוצת הבדיקה וככל שהתצפיות גדלות כך **הטעות עולה** אבל לא בהרבה – הפרש של 10% לכל היותר בין הדגימות.

השוואה בין דיוקי המודלים:

השוואה	KNN with optimal k	$Decision Tree$ with $Gain$
זמן למידה עד שמגלה את h דיוק לפי שיטת השוואה	דורש זמן גבוה למציאת k - משווא ע"י "מרחקים" בין התצפיות כאשר הגדרת המרחק משפיע גם הוא על k - בכדי להיות אופטימלי ומשפר גם הוא את הדיוק כפי שראינו הנ"ל	מתייצב באופן מהיר יחסית פונקציית ה- $Gain$ כפי שהוסבר הנ"ל מוצאת באופן חכם יותר את הפיצ'רים שאיתם נפצל את העץ ולכן מכלילה את התצפיות ב- $test - set$ יותר ולכן דיוק המודל יכול לעבוד לאורך זמן ולהתאים ליותר נתונים כפי שהוצג בגרף
אחוזי דיוק	כאשר ה- k הוא אופטימלי אחוז הדיוק כאן גבוה מאוד בדומה ל- DT	אחוז הדיוק גבוה מאוד בדומה ל- KNN עם k אופטימלי

הדפסת ה-*confusion matrix, precision, recall, F1* ועקומת הלמידה ל-*KNN* ($k = 3, p = 2,80\%$ train):

Confusion Matrix:

	Iris-versicolor	Iris-virginica	Iris-setosa
Iris-versicolor	8	1	0
Iris-virginica	0	13	0
Iris-setosa	0	0	8

Iris-versicolor Precision = 88.88888888888889

Iris-virginica Precision = 100.0

Iris-setosa Precision = 100.0

Iris-versicolor Recall = 100.0

Iris-virginica Recall = 92.85714285714286

Iris-setosa Recall = 100.0

F-[Iris-versicolor]-measure = 94.11764705882352

F-[Iris-virginica]-measure = 96.2962962962963

F-[Iris-setosa]-measure = 96.2962962962963

הדפסת ה-*confusion matrix, precision, recall, F1* ועקומת הלמידה ל-*DecisionTree* (80% train):

Confusion Matrix:

=====

	democrat	republican
democrat	91	3
republican	2	60

democrat Precision: 0.9680851063829787

democrat Precision: 0.967741935483871

democrat Recall: 0.978494623655914

democrat Recall: 0.9523809523809523

F-[democrat]-measure: 0.9732620320855615

F-[democrat]-measure: 0.96

סיכום ההבדלים בין 2 אלגוריתמי הלמידה המונחית שהוצגו הנ"ל:

:KNN

יתרונות – מימוש פשוט לביצוע ושיטה מובנת ואף טריוויאלית, נותן ביצועים טובים על נתונים קטנים כאשר פונקציית המרחק בין התצפיות מוגדרת היטב.

חסרונות – דורש כמות אכסון גבוהה בכדי להריץ את המודל ביחס ל-*DT* שלא שומר את כל הנתונים אלא רק את הפונקציה שעץ ההחלטה מייצר וזה בעצם ייצוג קומפקטי יותר מ-*KNN* בהשוואה בין המודלים, מכמות אכסון גבוהה יוצא גם שבשלב הסיווג האלגוריתם רץ על כל התצפיות שניתנו לו בלמידה ולכן הוא בעצם לא בדיוק לומד אלא יותר אלגוריתם עצל ששומר נתונים בצורה נוחה ומביא באופן *broote force* את התוצאה ולכן זמן הריצה שהוא דורש בכמויות גדולות יותר של נתונים נחשבת משמעותית יותר ויותר.

:DT

יתרונות – ייצוג קומפקטי לפונקציה המקורבת h של התצפיות שנתנו לה בזמן הלמידה ולכן בניגוד ל- KNN תופס שטח אכסון קטן יותר משמעותית כפונקציה של הקלט, מכליל הרבה יותר טוב את הנתונים מאשר KNN בכמויות גדולות יותר של נתונים מכיוון שדורש פחות פרמטרי קלט מ- KNN כגון: k וגם עצם זה שמתייחס לאזורים במרחב ולא לכל תצפית. לכן, האלגוריתם הוא כמעט "עצמאי" למצוא את העץ שמתאר באופן הטוב ביותר את התצפיות שקיבל, מהר מאוד ניתן למצוא את כמות התצפיות שתיתן תוצאה טובה ומוכללת בניגוד ל- KNN שמקיימת את זה עבור k גדול יחסית (לפני כן זה *underfitting* ולאחר מכן זה *overfitting*), משתמש בפונקציית ה- $Gain$ שהיא בעצם פונקציית ה- $mutual information$ שמזהה תלות בין משתנים וע"י ביצוע מקסימום ניתן להשים פיצולים שמסבירים באופן כללי יותר את תופעת התצפיות מכאן נובע שעץ החלטה זה אלגוריתם לומד ומשתמש ביותר כלים חכמים ובעלי משמעות מאשר KNN ומכאן גם ההכללה הטובה יותר שלו.

חסרונות – ישנם פונקציות ידועות שנלמדו בהרצאה שאינן יודע לבצע בגלל צורת החלוקה של העולם שהוא מייצר, עלול להיכנס מהר מאוד ל-*overfitting* ולשם כך ישנם פעולות שמטפלות בתופעה זו אבל מייקרות את תהליך למידה בזמן הריצה, כאשר מספר המשתנים/הפיצ'רים גדל איתו גם כמות האפשרויות לפיצול בכל צומת ולכן אלגוריתם זה במקרים מסוימים יכול להיות יותר מורכב מ- KNN ופחות ברור לסביבה כך שמייצר בעיות כגון: *debugging*.