Machine Learning

Final Project Report

**Team members:** Shilp Patel, Divy Patel, Smit Patel, Jay Patel


<div align="center">Machine Learning Final Project</div>

In this project, we had two primary tasks; one was to do classification with the given 5 train, test, label datasets and second one was to estimate the missing values in two additional datasets. After a little bit of research and execution, we have implemented all of the projects with the use of Java. To achieve the target for part-one, we used the KNN algorithm in order to predict the missing values in the data set. We utilized two dimensional arrays for our data sets to make our algorithm work smoothly and more essential. Our main goal to use the KNN algorithm was, it can help us very efficiently in handling the endless data and also the primary concern was in regarding the algorithm that it works very smoothly with multi-dimensions.

In the first part, the most important task was the classification part. In that, we had to check for the missing values in all of the datasets provided. Surprisingly, most of them had missing values, it was unpredicted. For this case, our KNN algorithm helped a lot in finding the nearest k-neighbors inside the missing values. We had many similar types of value, so we used Euclidean distance for our dataset. The KNN formula was used in order to get the results. After using that formula, we found the corresponding value for every nearest K neighbor. By calculating the mean value of the attribute between the closest K neighbors, except such missing values. After replacing the

present missed value with the equivalent mean value, the approach was used in this classification. To predict the labels for the test data, the same approach as above was used again. For instance, we used the KNN algorithm and the Euclidean distance to predict the values. The next thing we applied was that the decision set was calculated by dividing the distance by 1 and a request instance and obtaining for each possible mark the sum of all weights. The best prediction rating shall be the mark with the highest weight score.

In the second challenge, the Gene expression data are similar and reconsidered from the previous parts. The mean we found from the attribute value of the nearest neighbors was calculated using KNN. To get the absolute distances, we used Euclidean Distance Formula. The appropriate calculated mean was used to replace all the missing values in order to get the best results.

| Code | Comments |
|---|---|
| ```
String[] training = {"TrainData1.txt","TrainData2.txt", "TrainData3.txt", "TrainData4.txt","TrainData5.txt"};
    int [] trainSamples = {150, 100, 6300, 2547, 1119};
    int[] trainFeatures = {3312, 9182, 13, 112, 11};

String[] testing = {"TestData1.txt", "TestData2.txt", "TestData3.txt", "TestData4.txt", "TestData5.txt"};
    int[] testSamples = {53, 74, 2693, 1092, 480};
    int[] testFeatures = {3312, 9182, 13, 112, 11};

    fixed = false;
    for(int i = 0; i < training.length; i++) {

        getData(training[i], trainSamples[i], trainFeatures[i]);
        currentData = trainingData;
        fixData();
        writeData(training[i]);
        System.out.println();
    }
    for(int i = 0; i < testing.length; i++) {

        getData(testing[i], testSamples[i], testFeatures[i]);
        currentData = testingData;
        fixData();
        writeData(testing[i]);
        System.out.println();
    }
    fixed = true;
``` | This part of the code is used to find missing value estimates using KNN algorithm. Training Dataset are gathered from the file. This part is used to find estimated missing values. It outputs the complete data into a new file. This part is used to find estimated missing values. |
| ```
String[] newTraining = {"NewTrainData1.txt", "NewTrainData2.txt", "NewTrainData3.txt", "NewTrainData4.txt","NewTrainData
String[] newTesting = {"NewTestData1.txt", "NewTestData2.txt", "NewTestData3.txt", "NewTestData4.txt", "NewTestData5.txt
String[] labels = {"TrainLabel1.txt", "TrainLabel2.txt", "TrainLabel3.txt", "TrainLabel4.txt", "TrainLabel5.txt"};
for(int i = 0; i < newTraining.length; i++) {
    set = i+1;
    getData(newTraining[i], trainSamples[i], trainFeatures[i]);
    getData(newTesting[i], testSamples[i], testFeatures[i]);
    getData(labels[i],trainSamples[i],0);

    classify();
}
``` | This section concerns the classification of the test data. After performing the code above we now finally have 2d arrays named train_Data and test_Data & for the label we will use vectorized data named lab_data. Now we will find nearest k- neighbors for every case using classify(). |

```
68      String [] micros = {"MissingData1.txt", "MissingData2.txt"};
69      int [] microSamples = {14, 50};
70      int [] microFeatures = {242, 758};
71
72
73      fixed = false;
74      for(int i = 0; i < micros.length; i++) {
75
76          getData(micros[i], microSamples[i], microFeatures[i]);
77          currentData = microData;
78          fixData();
79          writeData(micros[i]);
80          System.out.println();
81      }
82      fixed = true;
83
84  }
```

**Question 2: Missing Value Estimation**
The following part is about Gene Expression Data. Datasets are gathered from the file.This part is used to find estimated missing values.It outputs the complete data into a new file.

```
89      System.out.println("Retrieving data from "+filename);
90      try {
91          FileReader fileReader = new FileReader("./input/"+filename);
92          BufferedReader bufferedReader = new BufferedReader(fileReader);
93
94          String line;
95
96
97          if(filename.contains("Label")) {
98
99              labelData = new String[samples];
100
101             for(int i = 0; i < labelData.length; i++) {
102                 line = bufferedReader.readLine().trim();
103                 labelData[i] = line;
104             }
105             System.out.println("Success!\n");
106
107         }
108
109         else if(filename.contains("Missing")) {
110             System.out.println(filename+" Samples: "+samples+" Features: "+features);
111             microData = new String[features][samples];
112
113             String[] lineSplit;
114             for(int i = 0; i < microData.length; i++) {
115                 line = bufferedReader.readLine().trim();
116                 lineSplit = line.split("\\s+|\\s*,\\s*");
117                 for (int j = 0; j < microData[i].length; j++) {
118                     microData[i][j] = lineSplit[j];
119                 }
120             }
```

Reads the training data collection and inserts a 2d string sequence called train_Data & test_Data. This codes shown are are used to for Labels and for Gene Data

```
123             double d = Math.sqrt(microData.length) / 2;
124             k = (int)d;
125             bufferedReader.close();
126             System.out.println("Success!\n");
127         }
128
129
130         else if(filename.contains("Train")) {
131
132             System.out.println(filename+" Samples: "+samples+" Features: "+features);
133             trainingData = new String[samples][features];
134
135             String[] lineSplit;
136             for(int i = 0; i < trainingData.length; i++) {
137                 line = bufferedReader.readLine().trim();
138                 lineSplit = line.split("\\s+|\\s*,\\s*");
139                 for (int j = 0; j < trainingData[i].length; j++) {
140                     trainingData[i][j] = lineSplit[j];
141                 }
142             }
```

The following formula is used to identify the specific K-value. Following code is used for training data.

```
144
145             double d = Math.sqrt(trainingData.length) / 2;
146             k = (int)d;
147             bufferedReader.close();
148             System.out.println("Success!\n");
149         }
150
151         else if(filename.contains("Test")) {
152
153
154             testingData = new String[samples][features];
155
156             String[] lineSplit;
157             for(int i = 0; i < testingData.length; i++) {
158                 line = bufferedReader.readLine().trim();
159                 lineSplit = line.split("\\s+|\\s*,\\s*");
160                 for (int j = 0; j < testingData[i].length; j++) {
161                     testingData[i][j] = lineSplit[j];
162                 }
163             }
164
165             double d = Math.sqrt(testingData.length) / 2;
166             k = (int)d;
167             bufferedReader.close();
168             System.out.println("Success!\n");
169         }
170     }
171     catch(IOException io) {
172         System.out.println("File: "+filename+" not found.");
173         System.exit(1);
174     }
175 }
176
```

The following formula is used to identify the specific K-value. For testing data. The following formula is used to identify the specific K-value

| Code | Description |
|---|---|
| ```<br>177  public static void fixData() {<br>178<br>179      System.out.println("Examining data for missing values..");<br>180      for(int i = 0; i < currentData.length; i++) {<br>181          int missing = 0;<br>182          System.out.print("Line "+(i+1)+"... ");<br>183          for(int j = 0; j < currentData[i].length; j++) {<br>184              if(currentData[i][j].contains("+99")) {<br>185                  missing++;<br>186                  String[] query = currentData[i];<br>187<br>188<br>189                  Distance[] kNearest = myKnn(i, query);<br>190<br>191<br>192                  double sum = 0;<br>193                  for(int n = 0; n < k; n++) {<br>194                      Double val = Double.parseDouble(currentData[kNearest[n].index][j]);<br>195                      sum += val;<br>196                  }<br>197<br>198                  String mean = Double.toString(sum / k);<br>199<br>200                  currentData[i][j] = mean;<br>201              }<br>202          }<br>203          System.out.println("Missing values fixed: "+missing);<br>204<br>205      }<br>206      System.out.println("All Missing Values fixed!!");<br>207  }<br>208<br>209  public static Distance[] myKnn(int index, String[] query) {<br>210<br>211      Distance[] distances;<br>212<br>213      if(fixed == false) {<br>214<br>``` | Queries for missing attributes, gets the nearest K-neighbors for data, and substitutes the missing value with the mean value for that attribute. Outputs the k nearest value. |
| ```<br>216          distances = new Distance[currentData.length-1];<br>217          int j = 0;<br>218<br>219          for(int i = 0; i < currentData.length; i++) {<br>220              String[] instance = currentData[i];<br>221              Distance d = new Distance(getDistance(instance, query), i);<br>222              if(i != index) {<br>223                  distances[j] = d;<br>224                  j++;<br>225              }<br>226          }<br>227      }<br>228      else {<br>229          distances = new Distance[trainingData.length];<br>230          for(int i = 0; i < trainingData.length ; i++) {<br>231              String[] instance = trainingData[i];<br>232              Distance d = new Distance(getDistance(instance, query), i);<br>233              distances[i] = d;<br>234          }<br>235      }<br>236      Arrays.sort(distances);<br>237      Distance[] kNearest = Arrays.copyOfRange(distances, 0, k);<br>238      return kNearest;<br>239  }<br>``` | Finds the closest k-neighbor to an object with the missing attribute and returns a collection of 'Distance' objects made up of the distance value and index of each of the closest k instances. |
| ```<br>241  public static double getDistance(String[] s1, String[] s2) {<br>242<br>243<br>244      double distance = 0;<br>245      for(int i = 0; i < s2.length ; i++) {<br>246          if(s1[i].contains("+99") || s2[i].contains("+99")) {<br>247              distance += 0;<br>248          }<br>249          else {<br>250              double d1 = Double.parseDouble(s1[i]);<br>251              double d2 = Double.parseDouble(s2[i]);<br>252              double d = Math.pow((d2 - d1),2);<br>253              distance += d;<br>254<br>255          }<br>256<br>257      }<br>258      return Math.sqrt(distance);<br>259  }<br>260<br>261  public static void writeData(String filename) {<br>262<br>263      String file;<br>264      if(filename.contains("Missing")) {<br>265          file = "./output/PatelNew"+filename;<br>266      }<br>267      else {<br>268          file = "./input/New"+filename;<br>269      }<br>270      try {<br>271          PrintWriter pw = new PrintWriter(file);<br>272          for(int k=0; k<currentData.length; k++) {<br>273              for(int l=0; l<currentData[k].length; l++) {<br>274                  pw.print(currentData[k][l]+"\t");<br>275              }<br>276              pw.print("\n");<br>277          }<br>278          pw.close();<br>279          System.out.println("Fixed data saved to: "+file);<br>280      }<br>``` | Calculates the difference between two vectors for Euclidean while overlooking any other incomplete data values. Writes the complete data into a new file. |

| Code | Description |
|---|---|
| ```<br>281        catch (FileNotFoundException e) {<br>282            System.out.println("File not found.");<br>283        }<br>284<br>285<br>286    }<br>287<br>288    public static void classify() {<br>289<br>290        String[] results = new String[testingData.length];<br>291<br>292        for(int i = 0; i < testingData.length; i++) {<br>293<br>294            String[] query = testingData[i];<br>295            Distance[] kNearest = myKnn(-1, query);<br>296<br>297            String[] labels = new String[kNearest.length];<br>298            for(int j = 0; j<kNearest.length; j++) {<br>299                labels[j] = labelData[kNearest[j].index];<br>300            }<br>301            int res = getWeighted(labels, kNearest);<br>302<br>303            results[i] = ""+res;<br>304        }<br>305<br>``` | This result is compared to every possible case in the testing data |
| ```<br>307    try {<br>308        PrintWriter pw = new PrintWriter("./output/PatelClassification"+set+".txt");<br>309        for(int k=0; k<results.length; k++) {<br>310            System.out.println(results[k]);<br>311                pw.print(results[k]+"\n");<br>312        }<br>313        pw.close();<br>314        System.out.println("Testing labels saved to: PatelClassification"+set+".txt\n");<br>315    }<br>316    catch (FileNotFoundException e) {<br>317        System.out.println("File not found.");<br>318    }<br>319    }<br>320<br>``` | This part creates a new file with all the results. |
| ```<br>320<br>321    public static int getWeighted(String[] labels, Distance[] kNearest){<br>322<br>323        int[] classes = new int[numClasses[set-1]];<br>324        double[] scores = new double[classes.length];<br>325        for(int i=0; i<classes.length; i++) {<br>326            classes[i] = i+1;<br>327        }<br>328<br>329        for(int j = 0; j<classes.length; j++) {<br>330            for(int k = 0; k<labels.length; k++) {<br>331                if(String.valueOf(classes[j]).equals(labels[k])) {<br>332                    scores[j] += 1/kNearest[k].distValue;<br>333                }<br>334            }<br>335        }<br>336        double[] sorted = scores.clone();<br>337        Arrays.sort(sorted);<br>338        for(int l=0; l<scores.length; l++) {<br>339            if(sorted[(sorted.length)-1] == scores[l]){<br>340                return classes[l];<br>341            }<br>342        }<br>343        return -1;<br>344    }<br>345 }<br>``` | This part gives out the label values regarding the weights(high) of the nearest k-neighbor. |
| ```<br>347 class Distance implements Comparable<Distance> {<br>348<br>349<br>350    double distValue;<br>351    int index;<br>352<br>353    Distance(double distValue, int index) {<br>354        this.distValue = distValue;<br>355        this.index = index;<br>356    }<br>357<br>358    public int compareTo(Distance d) {<br>359        if(distValue == d.distValue) {<br>360            return 0;<br>361        }<br>362        else if(distValue > d.distValue) {<br>363            return 1;<br>364        }<br>365        else {<br>366            return -1;<br>367        }<br>368    }<br>369<br>370    public String toString() {<br>371        return "Index: "+index+" Distance: "+distValue+"\n";<br>372    }<br>373<br>374 }<br>``` | The distance entity used in myKnn system includes the case index and the distance from the query-instance. This is made equivalent by the distance attribute for sorting. |

Screenshots taken throughout the process of the code and the output of code is shown below: