## IC 272: DATA SCIENCE - III
## LAB ASSIGNMENT – IV
## Data classification using K-nearest neighbor classifier and Bayes classifier with unimodal Gaussian density

**Student's Name: Yash Sharma**                     **Mobile No: 8802131138**

**Roll Number:  B20241**                              **Branch:          CSE**

**1      a.**

| | Prediction Outcome | |
|---|---|---|
| **True  Label** | 93 | 25 |
| | 19 | 200 |

Figure 1 KNN Confusion Matrix for K = 1

| | Prediction Outcome | |
|---|---|---|
| **True  Label** | 92 | 26 |
| | 9 | 210 |

Figure 2 KNN Confusion Matrix for K = 3

|  | Prediction Outcome | |
|---|---|---|
| **True Label** | 92 | 26 |
|  | 10 | 209 |

**Figure 3 KNN Confusion Matrix for K = 5**

**b.**

**Table 1 KNN Classification Accuracy for K = 1, 3 and 5**

| K | Classification Accuracy (in %) |
|---|---|
| 1 | 86.944 |
| 3 | 89.614 |
| 4 | 89.318 |

**Inferences:**

1. The highest classification accuracy is obtained with K =2.
2. Increasing the value of k increases the accuracy for some iterations but only up to some limit i.e., k=3 here and then their onwards it seems to be decreasing.
3. When k=1 we estimate our probability based on a single sample: the closest neighbor. This is very sensitive to all sort of distortions like noise, outliers, mislabeling of data, and so on. By using a higher value for k, we tend to be more robust against those distortions. Thus, increasing k seems to be increasing the classification accuracy.
4. As the classification accuracy increases with the increase in value of K i.e., k from 1 to 3, the number of diagonal elements increases.
5. As we increase k, the areas predicting each class will be more "smoothed", since it's the majority of the k-nearest neighbors which decide the class of any point. Thus, it makes our model predict more test samples correctly with their respective positive or negative class. Thus, we see increase in diagonal elements.
6. As the classification accuracy increases with the increase in value of K i.e., k from 1 to 3, the number of off-diagonal elements decreases.
7. As we see the diagonal elements decrease as k increases as accuracy increases thus proportionally off-diagonal elements decreases (as total test samples are fixed) as now our model makes more accurate and correct predictions (diagonal elements) thereby decreasing off-diagonal elements.

**2    a.**

| | Prediction Outcome | |
|---|---|---|
| **True Label** | 111 | 7 |
| | 6 | 213 |

**Figure 4 KNN Confusion Matrix for K = 1 post data normalization**

| | Prediction Outcome | |
|---|---|---|
| **True Label** | 112 | 6 |
| | 5 | 214 |

**Figure 5 KNN Confusion Matrix for K = 3 post data normalization**

| | Prediction Outcome | |
|---|---|---|
| **True Label** | 111 | 7 |
| | 3 | 216 |

**Figure 6 KNN Confusion Matrix for K = 5 post data normalization**

**b.**

**Table 2 KNN Classification Accuracy for K = 1, 3 and 5 post data normalization**

| K | Classification Accuracy (in %) |
|---|---|
| 1 | 96.1424 |
| 3 | 96.7359 |
| 5 | 97.0326 |

**Inferences:**

1. Data normalization increases classification accuracy.
2. As we know KNN method is based on finding the Euclidean distance between a given test tuple with all other input tuples. Euclidean distance calculated on non-normalized data does not give real distance between those two tuples as any attribute in data set with big range can caused a biased result while calculating the Euclidean distance, Thus, normalized data gives the accurate Euclidean distance thereby increasing the accuracy as seen.
3. The highest classification accuracy is obtained with K =5.
4. Increasing the value of k increases the accuracy.
5. When k=1 we estimate our probability based on a single sample: the closest neighbor. This is very sensitive to all sort of distortions like noise, outliers, mislabeling of data, and so on. By using a higher value for k, we tend to be more robust against those distortions. Thus, increasing k seems to be increasing the classification accuracy.
6. As the classification accuracy increases with the increase in value of K i, the number of diagonal elements increases.
7. As we increase k, the areas predicting each class will be more "smoothed", since it's the majority of the k-nearest neighbors which decide the class of any point. Thus, it makes our model predict more test samples correctly with their respective positive or negative class. Thus, we see increase in diagonal elements.
8. As the classification accuracy increases with the increase in value of K, the number of off-diagonal elements decreases.
9. As we see the diagonal elements decrease with increase in k and thereby accuracy increases. Thus, proportionally off-diagonal elements decrease (as total test samples are fixed) and now our model makes more accurate and correct predictions (diagonal elements) thereby decreasing off-diagonal elements.

**3**

|  | Prediction Outcome | |
|---|---|---|
| **True Label** | 102 | 16 |
| | 3 | 216 |

**Figure 7 Confusion Matrix obtained from Bayes Classifier**

The classification accuracy obtained from Bayes Classifier is    94.362 %.

**Table 3 Mean for class 0 and class 1**

| S. No. | Attribute Name | Mean | |
|---|---|---|---|
| | | **Class 0** | **Class 1** |
| 1. | X_Minimum | Removed | |
| 2. | X_Maximum | 273.418 | 723.656 |
| 3. | Y_Minimum | Removed | |
| 4. | Y_Maximum | 1583169.659 | 1431588.69 |
| 5. | Pixels_Areas | 7779.663 | 585.967 |
| 6. | X_Perimeter | 393.835 | 54.491 |
| 7. | Y_Perimeter | 273.183 | 45.658 |
| 8. | Sum_of_Luminosity | 843350.275 | 62191.126 |
| 9. | Minimum_of_Luminosity | 53.326 | 96.236 |
| 10. | Maximum_of_Luminosity | 135.762 | 130.452 |
| 11. | Length_of_Conveyer | 1382.762 | 1480.018 |
| 12. | TypeOfSteel_A300 | Removed | |
| 13. | TypeOfSteel_A400 | Removed | |
| 14. | Steel_Plate_Thickness | 40.073 | 104.214 |
| 15. | Edges_Index | 0.123 | 0.385 |
| 16. | Empty_Index | 0.459 | 0.427 |
| 17. | Square_Index | 0.592 | 0.513 |
| 18. | Outside_X_Index | 0.108 | 0.02 |
| 19. | Edges_X_Index | 0.55 | 0.608 |
| 20. | Edges_Y_Index | 0.523 | 0.831 |
| 21. | Outside_Global_Index | 0.288 | 0.608 |
| 22. | LogOfAreas | 3.623 | 2.287 |

| 23. | Log_X_Index | 2.057 | 1.227 |
|-----|-------------|-------|-------|
| 24. | Log_Y_Index | 1.848 | 1.318 |
| 25. | Orientation_Index | -0.314 | 0.136 |
| 26. | Luminosity_Index | -0.115 | -0.116 |
| 27. | SigmoidOfAreas | 0.925 | 0.543 |

In Fig. 8 and 9 representing covariance matrices for class 0 and class 1 respectively the column numbers and row numbers correspond to attribute with serial number as in Table 3.

Data classification using K-nearest neighbor classifier and Bayes classifier with unimodal Gaussian density

| | 23 | 22 | 21 | 20 | 19 | 18 | 17 | 16 | 15 | 14 | 13 | 12 | 11 | 10 | 9 | 8 | 7 | 6 | 5 | 4 | 3 | 2 | 1 |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| 1 | -30.09 | 18.08 | 27.679 | -31.15 | -47.8 | -76.32 | 30.839 | 22.505 | 16.654 | -1.52 | 4.696 | -6.929 | 25.36 | 16.734 | 1237.6 | 2040.9 | 3886.07 | -3E+07 | -12944 | -15751 | -3E+05 | -6E+07 | **4.675E+04** |
| 2 | 73812 | -50711 | -82047 | 73014 | 1E+05 | 2E+05 | -86404 | -35306 | -19166 | 4295 | -59251 | 21948 | -47711 | -1E+05 | -8E+06 | -6E+06 | -6E+06 | 4.9E+10 | 2E+08 | 8E+07 | 1E+09 | **2E+12** | -6E+07 |
| 3 | 575 | -300.2 | 980.33 | 2840.7 | 1427 | 3456.9 | 556.08 | -354.6 | -1121 | 223.06 | 200.2 | 585.23 | -492.1 | 547.01 | 10070 | 6294.5 | -2E+05 | 1E+07 | 1E+07 | 7E+06 | **1E+08** | 1E+09 | -3.207E+05 |
| 4 | 28.521 | -15.7 | 72.436 | 169.13 | 68.412 | 183.06 | 45.342 | -13.28 | -67.82 | 10.994 | 10.596 | 38.161 | -24.09 | 31.924 | 771.6 | 769.59 | -7764 | 6E+08 | 706257 | **442771** | 7E+06 | 8E+07 | -1.575E+04 |
| 5 | 19.506 | -21.06 | 105.12 | 207.79 | 44.055 | 176.64 | 63.25 | 13.411 | -65.42 | 6.496 | -16.55 | 44.182 | -17.57 | 10.207 | 1492.1 | 1492.1 | -6894 | 8E+08 | **1E+06** | 706257 | 1E+07 | 2E+08 | -1.294E+04 |
| 6 | 62063.3 | -22291 | 96509.5 | 278177 | 157341 | 361545 | 60033 | -50985 | -123181 | 25471 | 44602 | 58475 | -53267 | 49760 | 2E+06 | 777671 | -2E+07 | **8.2E+11** | 8E+08 | 6E+08 | 1E+07 | 5E+10 | -3.261E+07 |
| 7 | -6.557 | 4.448 | 3.817 | -10.747 | -12.861 | -22.187 | 4.759 | 4.623 | 3.739 | -1.46 | 1.078 | -1.75 | 3.932 | -1.973 | -153.83 | 439.24 | **1458.21** | -2E+07 | -6894 | -7764 | -2E+05 | -6E+06 | 3686.07 |
| 8 | -2.737 | 2.716 | 4.136 | -1.529 | -4.358 | -5.859 | 4.207 | 1.575 | -0.142 | -0.35 | 2.058 | -0.222 | 1.769 | -0.791 | 2.285 | **333.38** | 439.24 | 777671 | 1492.1 | 769.59 | 6294.5 | -6E+06 | 2040.905 |
| 9 | 0.211 | -0.485 | 4.37 | 2.645 | -0 | 2.03 | 4.536 | -0.534 | -2.697 | -0.19 | 3.926 | 0.806 | 1.322 | -1.821 | **2521.6** | 2.285 | -153.83 | 2E+06 | 1492.1 | 771.6 | 10070 | -8E+06 | 1237.644 |
| 10 | 0.005 | -0.008 | -0.022 | 0.019 | 0.041 | 0.041 | -0.021 | -0.015 | 0.003 | 0.019 | -0.015 | 0.015 | -0.009 | **0.73** | -1.821 | -0.791 | -1.973 | 49760 | 10.207 | 31.924 | 547.01 | -114611 | 16.734 |
| 11 | -0.028 | 0.016 | 0.024 | -0.038 | -0.05 | -0.084 | 0.026 | 0.022 | 0.015 | -0.01 | 0.007 | -0.009 | **0.029** | -0.009 | 1.322 | 1.769 | 3.932 | -53267 | -17.57 | -24.09 | -492.1 | -47711 | 25.36 |
| 12 | 0.015 | -0.003 | 0.005 | 0.036 | 0.03 | 0.052 | 0.003 | -0.012 | -0.018 | 0.005 | 0.005 | **0.015** | -0.009 | 0.015 | 0.806 | -0.222 | -1.75 | 58475 | 44.182 | 38.161 | 585.23 | 21948 | -6.929 |
| 13 | -0.01 | 0.016 | 0.069 | 0.023 | 0.03 | 0.001 | 0.07 | -0.001 | -0.036 | -0 | **0.064** | 0.005 | 0.007 | -0.015 | 3.926 | 2.058 | 1.078 | 44602 | -16.55 | 10.596 | 200.2 | -59251 | 4.696 |
| 14 | 0.007 | -0 | -0.01 | 0.014 | 0.021 | 0.029 | -0.01 | -0.007 | -0.002 | **0.005** | -0.004 | 0.005 | -0.006 | 0.019 | -0.192 | -0.353 | -1.46 | 25471 | 6.496 | 10.994 | 223.06 | 4294.7 | -1.516 |
| 15 | -0.026 | 0.003 | -0.045 | -0.073 | -0.04 | -0.099 | -0.039 | 0.023 | **0.057** | -0 | -0.036 | -0.018 | 0.015 | -0.015 | -2.697 | -0.142 | 3.739 | -123181 | -65.42 | -67.82 | -1121 | -19166 | 16.654 |
| 16 | -0.031 | 0.014 | 0.023 | -0.045 | -0.06 | -0.099 | 0.025 | **0.031** | 0.023 | -0 | -0.001 | -0.012 | 0.022 | -0.015 | -0.534 | 1.575 | 4.623 | -50985 | 13.411 | -13.28 | -354.6 | -35306 | 22.505 |
| 17 | -0.033 | 0.033 | 0.138 | 0.019 | -0.07 | -0.058 | **0.203** | 0.025 | -0.039 | -0.01 | 0.07 | 0.003 | 0.026 | -0.021 | 4.536 | 4.207 | 4.759 | 60033 | 63.25 | 45.342 | 556.08 | -86404 | 30.839 |
| 18 | 0.135 | -0.067 | -0.044 | 0.247 | 0.267 | **0.471** | -0.058 | -0.099 | -0.098 | 0.029 | 0.001 | 0.052 | -0.084 | 0.041 | 2.03 | -5.859 | -22.187 | 361545 | 176.6 | 183.06 | 3456.9 | 168070 | -76.32 |
| 19 | 0.082 | -0.044 | -0.066 | 0.124 | **0.168** | 0.267 | -0.073 | -0.063 | -0.039 | 0.021 | -0.02 | 0.03 | -0.054 | 0.041 | -0 | -4.358 | -12.861 | 157341 | 44.055 | 68.412 | 1427 | 111448 | -47.782 |
| 20 | 0.065 | -0.025 | 0.029 | **0.157** | 0.124 | 0.247 | 0.019 | -0.045 | -0.073 | 0.014 | 0.023 | 0.036 | -0.038 | 0.019 | 2.645 | -1.529 | -10.747 | 278177 | 207.79 | 169.13 | 2840.7 | 73014 | -31.147 |
| 21 | -0.028 | 0.031 | **0.133** | 0.029 | 0.124 | 0.247 | 0.138 | 0.023 | -0.045 | -0.01 | 0.069 | 0.005 | 0.024 | -0.022 | 4.37 | 4.136 | 3.817 | 96509.5 | 105.12 | 72.436 | 980.33 | -82047 | 27.679 |
| 22 | -0.026 | **0.027** | 0.031 | -0.025 | -0.04 | -0.067 | 0.033 | 0.014 | 0.003 | -0 | 0.016 | -0.003 | 0.016 | -0.008 | -0.485 | 2.716 | 4.448 | -22291 | -21.06 | -15.7 | -300.2 | -50711 | 18.083 |
| 23 | **0.049** | -0.026 | -0.028 | 0.065 | 0.082 | 0.135 | -0.033 | -0.031 | -0.026 | 0.007 | -0.01 | 0.015 | -0.028 | 0.005 | 0.211 | -2.737 | -6.557 | 62063.3 | 19.506 | 28.521 | 575.04 | 73812 | -30.093 |

Figure 8: Covariance matrix for class 0

| | 23 | 22 | 21 | 20 | 19 | 18 | 17 | 16 | 15 | 14 | 13 | 12 | 11 | 10 | 9 | 8 | 7 | 6 | 5 | 4 | 3 | 2 | 1 |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| 1 | -14.22 | -8.452 | -25.9 | -21.2 | 5.306 | -13.78 | -16.56 | -5.018 | 6.695 | 1.504 | 10.89 | -3.806 | 8.914 | -1933 | 13220 | -744 | -1224.8 | -2E+06 | -1974 | 1101.1 | -22255 | 1E+08 | **256526.31** |
| 2 | -37675 | -14718 | -1E+05 | -63427 | 64300 | 15298 | -74705 | -22199 | 64533 | 13457 | -38010 | -19251 | 23556 | -4E+07 | 4E+06 | -43296 | -4E+06 | 3E+10 | 5E+06 | 2E+07 | 3E+08 | **3E+12** | 1E+08 |
| 3 | 218.95 | -32.38 | 65.419 | 355.12 | 330.78 | 653.05 | 55.178 | -96.06 | -101.6 | 52.909 | -90.63 | 35.619 | -47.65 | 4262.2 | -23835 | -300.3 | -15632 | 5E+08 | 129451 | 178492 | **5E+06** | 3E+08 | -2254.624 |
| 4 | 15.51 | -1.119 | -3.758 | 16.86 | 23.56 | 36.62 | -2.152 | -9.176 | -4.85 | 3.972 | -7.318 | 4.156 | -1.332 | 282.11 | -1447 | 30.15 | -570.12 | 2E+07 | 55546.9 | **9807.2** | 178492 | 2E+07 | 1101.079 |
| 5 | 13.01 | -1.556 | 11.045 | 21.025 | 10.681 | 29.028 | 7.11 | -2.367 | -8.612 | 1.204 | -6.496 | 2.952 | -2.244 | 438.56 | -1139 | -79.15 | -557.42 | 1E+07 | **5000.6** | 55546.9 | 129451 | 5E+06 | -1973.565 |
| 6 | 22865 | -2282 | 6364.1 | 36735 | 34740 | 67783 | 5462.3 | -10272 | -10535 | 5578 | -9653 | 3985.1 | -4689 | 343512 | -3E+06 | 84723 | -1E+06 | **5.1E+10** | 1E+07 | 2E+07 | 5E+08 | 3E+10 | -2334975.6 |
| 7 | -1.984 | 3.684 | -2.503 | -3.287 | -1.3 | -5.043 | -2.224 | -0.833 | 0.427 | -0.15 | 0.775 | 0.591 | 1.066 | -204.8 | -993.3 | 348.05 | **733.909** | -1E+06 | -557.4 | -570.1 | -15632 | -4E+06 | -1224.809 |
| 8 | -0.96 | 2.786 | -2.874 | -2.165 | 0.678 | -1.504 | -2.018 | -1.09 | 0.878 | 0.044 | -0.267 | -0.025 | 0.429 | -205.4 | -381.1 | **406.46** | 348.045 | 84723 | -79.15 | 30.15 | -300.3 | -43296 | -744.043 |
| 9 | -5.967 | -4.547 | -7.431 | -10.57 | -1.44 | -7.953 | -3.138 | 1.971 | 6.591 | -0.7 | 2.468 | -5.16 | -0.09 | 1243.4 | **23101** | -381.1 | -993.31 | -3E+06 | -1139 | -1447 | -23835 | 4E+06 | 13220.079 |
| 10 | 2.39 | -1.662 | 7.846 | 5.403 | -1.38 | 3.627 | 6.623 | 2.058 | -3.443 | -0.17 | -1.134 | 0.699 | -1.331 | **5645.3** | 1243.4 | -205.4 | -204.84 | 343512 | 438.56 | 282.11 | 4262.2 | -4E+07 | -1932.619 |
| 11 | -0.004 | 0.005 | -0.024 | -0.017 | 0.005 | -0.012 | -0.017 | -0.003 | 0.008 | 0 | 0.011 | -0.001 | **0.09** | -1.331 | -0.09 | 0.429 | 1.066 | -4689 | -2.244 | -1.332 | -47.65 | 23556 | 8.914 |
| 12 | 0.024 | 0.002 | -0.004 | 0.022 | 0.022 | 0.026 | -0.008 | -0.011 | -0.012 | 0.001 | -0.002 | **0.02** | -0.001 | 0.699 | -5.16 | -0.025 | 0.591 | 3985.08 | 2.952 | 4.156 | 35.619 | -19251 | -3.806 |
| 13 | -0.028 | 0.001 | -0.021 | -0.033 | -0.02 | -0.053 | -0.016 | 0.015 | 0.02 | -0 | **0.082** | -0.002 | 0.011 | -1.134 | 2.468 | -0.267 | 0.775 | -9652.6 | -6.496 | -7.318 | -90.63 | -38010 | 10.893 |
| 14 | 0.005 | 0 | -0.01 | 0.001 | 0.012 | 0.012 | -0.01 | -0.01 | -0.014 | **0.002** | -0 | 0.002 | 0 | -0.165 | -0.698 | 0.044 | -0.151 | 5577.97 | 1.204 | 3.972 | 52.909 | 13457 | 1.504 |
| 15 | -0.045 | 0.004 | -0.103 | -0.086 | -0.103 | -0.066 | -0.068 | -0.014 | **0.065** | -0.014 | 0.02 | -0.012 | 0.008 | -3.443 | 6.591 | 0.878 | 0.427 | -10535 | -8.612 | -4.85 | -101.6 | 64533 | 6.695 |
| 16 | -0.017 | -0.007 | 0.086 | 0.024 | -0.06 | -0.025 | 0.064 | **0.049** | -0.014 | -0.01 | 0.015 | -0.011 | -0.003 | 2.058 | 1.971 | -1.09 | -0.833 | -10272 | -2.367 | -9.176 | -96.06 | -22199 | -5.018 |
| 17 | 0.022 | -0.015 | 0.229 | 0.113 | -0.07 | 0.048 | **0.227** | 0.064 | -0.068 | -0.01 | -0.016 | -0.008 | -0.017 | 6.623 | -3.138 | -2.018 | -2.224 | 5462.3 | 7.11 | -2.152 | 55.178 | -74705 | -16.564 |
| 18 | 0.147 | -0.019 | 0.073 | 0.177 | 0.116 | **0.271** | 0.048 | -0.025 | -0.066 | 0.012 | -0.053 | 0.026 | -0.012 | 3.627 | -7.953 | -1.504 | -5.043 | 67782.7 | 29.028 | 36.62 | 653.05 | 15298 | -13.781 |
| 19 | 0.065 | 0 | -0.101 | 0.017 | **0.119** | 0.116 | -0.073 | -0.058 | 0.011 | 0.012 | -0.02 | 0.022 | 0.005 | -1.376 | -1.44 | 0.678 | -1.299 | 34740.3 | 10.681 | 23.557 | 330.78 | 64300 | 5.306 |
| 20 | 0.103 | -0.017 | 0.169 | **0.178** | 0.017 | 0.177 | 0.113 | 0.024 | -0.086 | 0.001 | -0.033 | 0.022 | -0.017 | 5.403 | -10.57 | -2.165 | -3.287 | 36734.8 | 21.025 | 16.864 | 355.12 | -63427 | -21.204 |
| 21 | 0.041 | 0.025 | **0.302** | 0.169 | -0.101 | 0.073 | 0.229 | 0.086 | -0.103 | -0.01 | -0.021 | -0.004 | -0.024 | 7.846 | -7.431 | -2.874 | -2.503 | 6364.12 | 11.045 | -3.758 | 65.419 | -119870 | -25.896 |
| 22 | -0.009 | **0.025** | 0.025 | -0.017 | 0 | -0.019 | -0.015 | -0.007 | 0.004 | 0 | 0.001 | 0.002 | 0.005 | -1.662 | -4.547 | 2.786 | 3.684 | -2282.4 | -1.556 | -1.119 | -32.38 | -14718 | -8.452 |
| 23 | **0.102** | -0.009 | 0.041 | 0.103 | 0.065 | 0.147 | 0.022 | -0.017 | -0.045 | 0.005 | -0.028 | 0.024 | -0.004 | 2.39 | -5.987 | -0.96 | -1.984 | 22864.8 | 13.014 | 15.508 | 218.95 | -37675 | -14.221 |

Figure 9: Covariance matrix for class 1

**Inferences:**

1. Accuracy of Bayes Classifier is 94.362 %. The reason that it is lesser than the accuracy obtained via "KNN on normalized data" (97.003 %) can be because K-NN does better because of its inherent nature to optimize locally and make predictions based on maximum class of the k-nearest neighbors around it. But still we see Bayes Classifier still isn't very far from the accuracy 97.003 % thus it may not be regarded as a bad approach just based upon this data prediction results.

2. As seen from covariance matrix the diagonal elements are all positive and the initial values also are very big in magnitude. They are positive because they represent the variance of the respective attribute.

3. The off-diagonal elements in the covariance matrix are both negative and positive are of varying magnitude (some of big magnitude and some of less). They simply represent the covariance between any two given attributes. The attributes "Y_Maximum" and "Sum_of_Luminosity" have highest covariance while attributes "X_Maximum" and "Y_Maximum" have least covariance for class 0 and The attributes "Y_Maximum" and "Sum_of_Luminosity" have highest covariance while attributes "Steel_Plate_Thickness" and "Y_Maximum" have least covariance for class 1.

**4**

Table 4 Comparison between classifiers based upon classification accuracy

| S. No. | Classifier | Accuracy (in %) |
|--------|------------|-----------------|
| 1. | KNN | 89.614 |
| 2. | KNN on normalized data | 97.033 |
| 3. | Bayes | 94.362 |

**Inferences:**

1. K-NN classifier on normalized data has the highest accuracy while K-NN classifier on not-normalized data has least accuracy.

2. KNN Classifiers < Bayes Classifier < K-NN Classifier on normalized data.

3. The reason behind the fact that K-NN performed on normalized data has higher accuracy than K-NN performed on normal data is as follow. As we know KNN method is based on finding the Euclidean distance between a given test tuple with all other input tuples. Euclidean distance calculated on non-normalized data does not give real distance between those two tuples as any attribute in data set with big range can caused a biased result while calculating the Euclidean distance, Thus, normalized data gives the accurate Euclidean distance thereby increasing the accuracy as seen.

4. The reason for the Bayes classifier having low accuracy than K-NN (on normalized data) is as follow. K-NN does better because of its inherent nature to optimize locally and make predictions based on maximum class of the k-nearest neighbors around it. But still we see Bayes Classifier still isn't very far from the accuracy 97.003 % thus it may not be regarded as a bad approach just based upon this data prediction results.

5. Generally, Bayes Classifier is more accurate and beneficial on complex data sets as their K-NN method would fail due to its high time complexity and its inherent nature to optimize locally. Here, data set is not that much big and complex thus K-NN shows best accuracy among all other classifiers.