## Assignment 4

## Problem 1

**a)**

**Model Summary**

| Specifications | Growing Method | CRT |
|---|---|---|
| | Dependent Variable | V1 |
| | Independent Variables | V2, V3, V4, V5, V6, V7, V8, V9, V10, V11, V12, V13, V14, V15, V16, V17 |
| | Validation | Split Sample |
| | Maximum Tree Depth | 20 |
| | Minimum Cases in Parent Node | 10 |
| | Minimum Cases in Child Node | 5 |
| Results | Independent Variables Included | V12, V8, V11, V7, V10, V14, V13, V9, V16, V15, V17, V4, V3, V5, V2, V6 |
| | Number of Nodes | 973 |
| | Number of Terminal Nodes | 487 |
| | Depth | 20 |

The holdout method is used to split data, in which the training set contains 70% of total observations or 14,011 tuples, whereas a test set contains 5,989 tuples or 30% of the total observations. The decision tree is built using CHART method with the GINI index of 0.000. The tree contains 973 nodes, which 487 nodes are terminal nodes. The error rate for a training set is 12.4% and for a test set is 19.4%.

| Model (Parent, Child) | Depth | Partition | Accuracy | Error Rate |
|---|---|---|---|---|
| (10, 5) | 20 | Train | 87.6% | 12.4% |
| | | Test | 80.6% | 19.4% |
| (20, 10) | 20 | Train | 82.4% | 17.6% |
| | | Test | 78.0% | 22.0% |
| (50, 25) | 20 | Train | 73.3% | 26.7% |
| | | Test | 71.8% | 28.2% |
| (100, 50) | 10 | Train | 62.9% | 37.1% |
| | | Test | 62.9% | 37.1% |
| (200, 100) | 10 | Train | 58.0% | 42.0% |
| | | Test | 58.1% | 41.9% |

**b)** The table shows different accuracies of different models. The best model is the first model, containing 10 tuples in parent and 5 tuples in child node with the level of depth of 20. Among 5 models, this first model is the best model because it has the highest accuracy rate (80.6% for a test set) and the lowest error rate (19.4% for a test set).
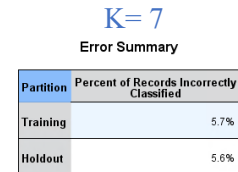
**c)** The misclassification matrix is shown on the left. It comes for the best model selected in part b). The overall accuracy rate is high. The rate for a training set is 87.6% and for a test set is 80.6%. The accuracy is a good measure for this dataset because each class of target variable has approximately equal frequency, so there is no problem of imbalanced class.

**Independent Variable Importance**

| Independent Variable | Importance | Normalized Importance |
|---|---|---|
| V14 | .344 | 100.0% |
| V11 | .291 | 84.6% |
| V15 | .275 | 80.0% |
| V12 | .270 | 78.6% |
| V10 | .264 | 77.0% |
| V9 | .250 | 72.8% |
| V16 | .244 | 71.2% |
| V13 | .229 | 66.6% |
| V7 | .223 | 64.9% |
| V8 | .215 | 62.5% |
| V17 | .205 | 59.8% |
| V2 | .113 | 32.8% |
| V4 | .111 | 32.2% |
| V6 | .109 | 31.9% |
| V3 | .104 | 30.4% |
| V5 | .079 | 22.9% |

Growing Method: CRT
Dependent Variable: V1

**d)** From the table above, the most important three attributes are ranked in order as following: V14 (x-ege mean edge count left to right), V11(xybar mean x y correlation), V15 (xegvy correlation of x-ege with y). The V14 attribute has the highest important rate and has the highest position in the tree, so the V14 is the most important attribute.

## Problem 2

a) The KNN is performed using Euclidean metric. All 16 attributes (predictors) are standardized using z-scores to make sure that there is no problem regarding a difference in scale. This prevents distance measures from being dominated by attributes that have larger ranges. I performed PCA to reduce dimension and ran the KNN with new components, but the accuracy rate was significantly lower than original attributes. So, I decide to use 16 attributes (normalized) in part b.

b) The misclassification matrices for (K=1, 3, 5, and 7) are shown below. The best value of K is equal to 1 because it is the smallest K, which has the lowest error rate of 4.7%. The K=1 model has the highest accuracy rate of 95.3%.



k Selection Error log

**K = 1**
Error Summary

| Partition | Percent of Records Incorrectly Classified |
|-----------|-------------------------------------------|
| Training  | 4.7% |
| Holdout   | 4.7% |

**K=3**
Error Summary

| Partition | Percent of Records Incorrectly Classified |
|-----------|-------------------------------------------|
| Training  | 5.0% |
| Holdout   | 5.0% |

**K= 5**
Error Summary

| Partition | Percent of Records Incorrectly Classified |
|-----------|-------------------------------------------|
| Training  | 5.4% |
| Holdout   | 5.6% |

**K= 7**
Error Summary

| Partition | Percent of Records Incorrectly Classified |
|-----------|-------------------------------------------|
| Training  | 5.7% |
| Holdout   | 5.6% |

| K | Partition | 1 | 3 | 5 | 7 |
|---|-----------|-----|-----|-----|-----|
| Accuracy | Train | 95.3% | 95.0% | 94.6% | 94.3% |
| | Test | 95.3% | 95.0% | 94.4% | 94.4% |
| Error Rate | Train | 4.7% | 5.0% | 5.4% | 5.7% |
| | Test | 4.7% | 5.0% | 5.6% | 5.6% |

K = 1



K=3

K = 5                                                      K=7





**c)** The model from KNN (k=1) has a better performance than the model from decision tree because the KNN model has a relatively higher accuracy rate of 95.3% (test set), compared to an accuracy rate of 80.6% (test set) from decision tree. The KNN model also has a lower error rate of 4.7%, compared to the error rate of 19.4% from decision tree.

## KNN Model

| K | Partition | 1 | 3 | 5 | 7 |
|---|---|---|---|---|---|
| Accuracy | Train | 95.3% | 95.0% | 94.6% | 94.3% |
| | Test | 95.3% | 95.0% | 94.4% | 94.4% |
| Error Rate | Train | 4.7% | 5.0% | 5.4% | 5.7% |
| | Test | 4.7% | 5.0% | 5.6% | 5.6% |

## Decision Tree

| Model (Parent, Child) | Depth | Partition | Accuracy | Error Rate |
|---|---|---|---|---|
| (10, 5) | 20 | Train | 87.6% | 12.4% |
| | | Test | 80.6% | 19.4% |

## Problem 3

**a)**

**1**. For K-means, the cluster centers or the mean points are calculated based on average of each of the attributes, which does not need to be an object in the data set.

**2**. First measure is a Jaccard coefficient which is similarity measure for data that has asymmetric binary variables. Second measure is cosine similarity which is used for the data that has Term-frequency vectors (sparse numeric data).

**3.**

## K= 3

**i  Final cluster centers**   **ii  number of elements in each cluster**   **iii class distribution**

**Final Cluster Centers**

|     | Cluster 1 | Cluster 2 | Cluster 3 |
|-----|-----------|-----------|-----------|
| V1  | 18.72     | 11.96     | 14.65     |
| V2  | 16.30     | 13.27     | 14.46     |
| V3  | .8851     | .8522     | .8792     |
| V4  | 6.2089    | 5.2293    | 5.5638    |
| V5  | 3.723     | 2.873     | 3.278     |
| V6  | 3.6036    | 4.7597    | 2.6489    |
| V7  | 6.066     | 5.089     | 5.192     |

**Number of Cases in each Cluster**

| Cluster | 1 | 61.000 |
|---------|---|--------|
|         | 2 | 77.000 |
|         | 3 | 72.000 |
| Valid   |   | 210.000 |
| Missing |   | .000   |

**V8 * Cluster Number of Case Crosstabulation**

Count

|     |       | Cluster Number of Case 1 | 2 | 3 | Total |
|-----|-------|---|---|----|-------|
| V8  | 1.000 | 1 | 9 | 60 | 70    |
|     | 2.000 | 60 | 0 | 10 | 70   |
|     | 3.000 | 0 | 68 | 2 | 70    |
| Total |     | 61 | 77 | 72 | 210  |

Actual Class

## K= 4

**i  Final cluster centers**   **ii  number of elements in each cluster**   **iii class distribution**

**Final Cluster Centers**

|     | Cluster 1 | Cluster 2 | Cluster 3 | Cluster 4 |
|-----|-----------|-----------|-----------|-----------|
| V1  | 11.94     | 14.42     | 17.75     | 19.52     |
| V2  | 13.27     | 14.35     | 15.88     | 16.65     |
| V3  | .8515     | .8795     | .8840     | .8844     |
| V4  | 5.2292    | 5.5239    | 6.0476    | 6.3501    |
| V5  | 2.867     | 3.253     | 3.614     | 3.812     |
| V6  | 4.8040    | 2.5904    | 3.1649    | 4.1641    |
| V7  | 5.095     | 5.127     | 5.921     | 6.184     |

**Number of Cases in each Cluster**

| Cluster | 1 | 75.000 |
|---------|---|--------|
|         | 2 | 67.000 |
|         | 3 | 40.000 |
|         | 4 | 28.000 |
| Valid   |   | 210.000 |
| Missing |   | .000   |

**V8 * Cluster Number of Case Crosstabulation**

Count

|     |       | Cluster Number of Case 1 | 2 | 3 | 4 | Total |
|-----|-------|---|----|----|----|-------|
| V8  | 1.000 | 8 | 58 | 4  | 0  | 70    |
|     | 2.000 | 0 | 6  | 36 | 28 | 70    |
|     | 3.000 | 67 | 3 | 0  | 0  | 70    |
| Total |     | 75 | 67 | 40 | 28 | 210  |

## K= 5

**i  Final cluster centers**   **ii  number of elements in each cluster**   **iii class distribution**

**Final Cluster Centers**

|     | Cluster 1 | Cluster 2 | Cluster 3 | Cluster 4 | Cluster 5 |
|-----|-----------|-----------|-----------|-----------|-----------|
| V1  | 16.56     | 14.69     | 19.15     | 12.09     | 11.98     |
| V2  | 15.39     | 14.47     | 16.47     | 13.31     | 13.29     |
| V3  | .8782     | .8809     | .8871     | .8571     | .8508     |
| V4  | 5.8882    | 5.5721    | 6.2689    | 5.2174    | 5.2414    |
| V5  | 3.481     | 3.286     | 3.773     | 2.901     | 2.880     |
| V6  | 4.1095    | 2.4079    | 3.4604    | 3.3438    | 5.6733    |
| V7  | 5.725     | 5.159     | 6.127     | 5.005     | 5.122     |

**Number of Cases in each Cluster**

| Cluster | 1 | 25.000 |
|---------|---|--------|
|         | 2 | 51.000 |
|         | 3 | 48.000 |
|         | 4 | 44.000 |
|         | 5 | 42.000 |
| Valid   |   | 210.000 |
| Missing |   | .000   |

**V8 * Cluster Number of Case Crosstabulation**

Count

|     |       | Cluster Number of Case 1 | 2 | 3 | 4 | 5 | Total |
|-----|-------|----|----|----|----|----|-------|
| V8  | 1.000 | 6  | 48 | 0  | 14 | 2  | 70    |
|     | 2.000 | 19 | 3  | 48 | 0  | 0  | 70    |
|     | 3.000 | 0  | 0  | 0  | 30 | 40 | 70    |
| Total |     | 25 | 51 | 48 | 44 | 42 | 210  |

## K= 6

**i  Final cluster centers**   **ii  number of elements in each cluster**   **iii class distribution**

**Final Cluster Centers**

|     | Cluster 1 | Cluster 2 | Cluster 3 | Cluster 4 | Cluster 5 | Cluster 6 |
|-----|-----------|-----------|-----------|-----------|-----------|-----------|
| V1  | 11.83     | 14.24     | 16.41     | 18.95     | 12.32     | 19.58     |
| V2  | 13.22     | 14.26     | 15.32     | 16.39     | 13.42     | 16.65     |
| V3  | .8500     | .8793     | .8783     | .8868     | .8580     | .8877     |
| V4  | 5.2156    | 5.4935    | 5.8640    | 6.2475    | 5.2659    | 6.3159    |
| V5  | 2.844     | 3.234     | 3.463     | 3.745     | 2.951     | 3.835     |
| V6  | 4.1684    | 2.3165    | 3.8501    | 2.7235    | 6.3367    | 5.0815    |
| V7  | 5.076     | 5.062     | 5.690     | 6.119     | 5.122     | 6.144     |

**Number of Cases in each Cluster**

| Cluster | 1 | 56.000 |
|---------|---|--------|
|         | 2 | 54.000 |
|         | 3 | 31.000 |
|         | 4 | 33.000 |
|         | 5 | 21.000 |
|         | 6 | 15.000 |
| Valid   |   | 210.000 |
| Missing |   | .000   |

**V8 * Cluster Number of Case Crosstabulation**

Count

|     |       | Cluster Number of Case 1 | 2 | 3 | 4 | 5 | 6 | Total |
|-----|-------|----|----|----|----|----|----|-------|
| V8  | 1.000 | 7  | 52 | 9  | 0  | 2  | 0  | 70    |
|     | 2.000 | 0  | 0  | 22 | 33 | 0  | 15 | 70    |
|     | 3.000 | 49 | 2  | 0  | 0  | 19 | 0  | 70    |
| Total |     | 56 | 54 | 31 | 33 | 21 | 15 | 210  |

**4.** The best value of K is 5 (elbow).

Ave_Distance

| K | Ave_Distance |
|---|---|
| 3 | 1.4915084 |
| 4 | 1.4089938 |
| 5 | 1.2235427 |
| 6 | 1.1620128 |

**Best K = 5**

**5.** From the class distribution in part iii, the best value of K would be 3 because the label is similar within the clusters, but it is different across clusters.

**6.** The normalization will influence the clustering results because the range of value changes. The attributes are normalized with z-scores before running K-mean (K=5). The number of elements in each cluster change after normalization. The cluster centers contain both positive and negative numbers. The distance between final clusters center changes. The frequency of data assigned to each cluster changes as well.

**After normalization:**

**Final Cluster Centers**

| | Cluster 1 | 2 | 3 | 4 | 5 |
|---|---|---|---|---|---|
| Zscore(V1) | 1.45626 | -1.08742 | -.75669 | .47893 | -.20755 |
| Zscore(V2) | 1.44408 | -1.04571 | -.78371 | .54612 | -.26019 |
| Zscore(V3) | .65695 | -1.18706 | -.18993 | .14168 | .55314 |
| Zscore(V4) | 1.43611 | -.90606 | -.77643 | .54258 | -.38641 |
| Zscore(V5) | 1.33994 | -1.19871 | -.63603 | .44951 | -.02707 |
| Zscore(V6) | -.15042 | .40134 | 1.61155 | .16569 | -.85882 |
| Zscore(V7) | 1.45420 | -.63393 | -.58677 | .55343 | -.72743 |

**Number of Cases in each Cluster**

| Cluster | 1 | 50.000 |
|---|---|---|
| | 2 | 55.000 |
| | 3 | 19.000 |
| | 4 | 28.000 |
| | 5 | 58.000 |
| Valid | | 210.000 |
| Missing | | .000 |

**Distances between Final Cluster Centers**

| Cluster | 1 | 2 | 3 | 4 | 5 |
|---|---|---|---|---|---|
| 1 | | 5.715 | 5.162 | 2.128 | 4.017 |
| 2 | 5.715 | | 1.724 | 3.610 | 2.767 |
| 3 | 5.162 | 1.724 | | 3.117 | 2.788 |
| 4 | 2.128 | 3.610 | 3.117 | | 2.252 |
| 5 | 4.017 | 2.767 | 2.788 | 2.252 | |

**b)**

**1. Single linkage algorithm**

### i  Dendogram



Dendrogram using Single Linkage
Rescaled Distance Cluster Combine

### ii Class distribution

**V8 * Single Linkage**
**Crosstabulation**

Count

|  |  | Single Linkage | | | |
|---|---|---|---|---|---|
|  |  | 1 | 2 | 3 | Total |
| V8 | 1.000 | 68 | 1 | 1 | 70 |
| **Class** | 2.000 | 70 | 0 | 0 | 70 |
|  | 3.000 | 68 | 2 | 0 | 70 |
| Total |  | 206 | 3 | 1 | 210 |

The dendogram is created from the nearest neighbors method. The class distribution is shown in the table above. The majority data falls in the first cluster.

## 2. Complete linkage and report

### i Dendogram



Dendrogram using Complete Linkage
Rescaled Distance Cluster Combine

### ii Class distribution

**V8 \* Complete Linkage Crosstabulation**

Count

| | | Complete Linkage | | | |
|---|---|---|---|---|---|
| | | 1 | 2 | 3 | Total |
| V8 | 1.000 | 48 | 2 | 20 | 70 |
| | 2.000 | 4 | 66 | 0 | 70 |
| | 3.000 | 0 | 0 | 70 | 70 |
| Total | | 52 | 68 | 90 | 210 |

The dendogram is created from the furthest neighbors method. The class distribution from complete linkage method is clearer than the single linkage shown in the table above.

C)                                    **executive summary**

## Overview

The seed dataset was retrieved from UCI Machine Learning Repository. It contains 210 random samples of kernels belonging to the three different varieties of wheat including Kama, Rosa and Canadian (class label). Each type contains 70 elements. The study combined a soft X-ray technique which used to detect the internal kernel structure with harvested wheat grain originating from experimental fields to measure kernels with seven geometric parameters (attributes) including area A, perimeter P, compactness $C = 4*pi*A/P^2$, length of kernel, width of kernel, asymmetry coefficient, and length of kernel groove. All seven attributes are numeric and continuously varied variables.

## Problem

The seed dataset contains sample kernels that needs to be clustered into groups and labeled. Kernels will be clustered based on their seven geometric parameters.

## Solutions

The raw dataset will be cleaned before doing the analysis in SPSS. We will implement a clustering analysis for this dataset by using two methods:

1. K-means clustering:  This method divides the observations into k clusters, in which each observation belongs to the cluster with the nearest mean (center of cluster). All seven attributes will be used in this clustering, and we will run the multiple numbers of clusters (k =3, 4, 5, 6) to indicate the best number of clusters for our seed datasets.

2. Hierarchical clustering: This method clusters observations by building a hierarchy of clusters. We will perform the "bottom-up" approach which will start from single cluster to merged clusters. We will perform both single linkage approach in which distance between two clusters is distance between two closest records and complete linkages approach in which distance between two clusters is distance between two farthest records

## Results and Recommendations:

The K-mean method suggests 5 clusters for 210 kernels in which each cluster contains kernels. Both single linkage and complete linkage approach suggest 3 clusters. However, the class distribution of complete linkage is clearer than single linkage. Among three methods, the class distribution is clearer and close to the actual class distribution when we use K-means with 3 clusters as it minimizes variability within the cluster, but maximizes variability between clusters. The graph shows the 7 attributes of each cluster. In recommendation, the 210 kernels should be clustered into 3 clusters using K-means method as shown in the left graph.