



# Font Classification

kNN Application

## Table of Contents

<b>Scope of the Report .....</b>	<b>2</b>
<b>Data Cleaning and Treatment.....</b>	<b>2</b>
<b>1 Preliminary Analysis of the Data Set .....</b>	<b>3</b>
<b>1.1 Statistical Parameter Analysis of Feature X210 .....</b>	<b>3</b>
1.1.1 Mean of Feature X210 for each class .....	3
1.1.2 Student's t-test .....	3
1.1.3 Histogram of X210 Feature for each class.....	4
1.1.4 Kolmogorov-Smirnov test.....	5
<b>1.2 Correlation between each feature .....</b>	<b>5</b>
1.2.1 Correlation Value for 400 features and top 10 pairs of correlated pixels .....	5
1.2.2 Correlation Value Comparison .....	6
<b>1.3 Standardize the Data Set .....</b>	<b>6</b>
<b>1.4 Generate the complete Data Set with Response Variable and Features .....</b>	<b>7</b>
<b>2 Data Automatic Classification .....</b>	<b>8</b>
<b>2.1 Split the Standardize Data Set into Training Set and Test Set .....</b>	<b>8</b>
<b>2.2 Preliminary k-Nearest Neighbor (kNN) Automatic Classifier.....</b>	<b>8</b>
<b>2.3 Secondary k-Nearest Neighbor (kNN) Automatic Classifier .....</b>	<b>9</b>
2.3.1 Repeat the kNN process .....	9
2.3.2 Determine the best k .....	10
<b>2.4 Confusion Matrix .....</b>	<b>10</b>
<b>2.5 Misclassified Cases .....</b>	<b>11</b>
<b>3 Principal Component Analysis (PCA).....</b>	<b>14</b>
<b>3.1 Eigenvalue and Eigenvector .....</b>	<b>14</b>
<b>3.2 Percentage of Explained Variance (PEV).....</b>	<b>14</b>
<b>3.3 Apply Principal Component Analysis.....</b>	<b>15</b>
<b>3.4 Apply kNN Automatic Classifier to Data Set after Principal Component Analysis.....</b>	<b>15</b>
<b>3.5 Visualize the First Two Component.....</b>	<b>16</b>
<b>3.6 Visualize misclassified cases .....</b>	<b>18</b>

## Scope of the Report

This report is going to classify the font by its features which are the gray scale value extracted from a 20 x 20 pixels image. Thus, each case will have 400 features. The interested data is downloaded from “University of California Irvine Repository of Machine Learning Datasets”. There are 153 .csv files and each file contain a bunch of cases with some information about the case such as if the font is bolded or presented in italic style, etc. The most important things are the 400 features came from the gray scale value of each pixel. We are going to select 4 out of 153 files. Data will be cleaned and treated to apply the k-Nearest Neighbor (kNN) Automatic Classifier. We would like to know if it is possible to use gray scale value of each pixel to predict the correct font. If so, what is the accuracy of prediction.

Firstly, the data will be split into 20% of test set and 80% of training set randomly. Then, use these two sets to find out the best hyperparameter,  $k$ , in kNN model and identify the misclassified data and explain them.

Second, Principal Component Analysis method will be introduced in Part 3 to reduce the number of features. This can help to reduce the computation time but keep the prediction accuracy as high as possible.

## Data Cleaning and Treatment

The four font files listed below were chosen for this study:

1. GILL.csv
2. LEELAWADEE.csv
3. ROMAN.csv
4. TECHNIC.csv

The informative features of these datasets are

- font: font tyle of each case
- italics: 1 if the character is italic and 0 otherwise
- strength: 0.4 for normal characters and 0.7 for bold characters
- as well as 400 others that correspond to a specific pixel.

A single case in the font dataset describes a digitized image of a specific character typed in a specific font. Each image has size  $20 \times 20 = 400$  pixels. Each pixel is assigned an integer value between 0 and 255, known as a “gray level”.

These 9 columns are discarded: fontVariant, m\_label, orientation, m\_top, m\_left, originalH, originalW, h, w while 3 columns, font, strength, italic, were kept along with the 400 columns named r0c0, r0c1, r0c2, ..., r19c18, r19c19. Any row containing missing numerical data were also discarded.

We then defined GILL as CL1, LEELAWADIE as CL2, ROMAN as CL3, and TECHNIC as CL4. The size of each class and total size of all the classes are shown in Table 1 below.

Table 1 Selected Files and Size of each class and total

	File	Size
CL1	GILL.csv	1459
CL2	LEELAWADEE.csv	1378
CL3	ROMAN.csv	1194
CL4	TECHNIC.csv	1380
Total		5411

Since each case has about 22 ~ 27% of the total data set, it can be considered as balanced data. Thus, further treatment such as cloning, perturbation, or SMOTE which are used to deal with imbalanced data are not required. These four classes combined make a 5411 x 400 data frame.

## 1 Preliminary Analysis of the Data Set

### 1.1 Statistical Parameter Analysis of Feature X210

Feature X210 is the pixel located at row 10 and column 9. This example should give some insight about the information for the entire data set. Since this data is the gray level image of a font, it can be imaged that most of the values in each feature are close to either 0 (white) or 255 (black). There are very few numbers between 1 and 254.

#### 1.1.1 Mean of Feature X210 for each class

Firstly, the mean of this feature among all the cases is calculated. Under this circumstance, the means can be interpreted as the ratio of 0 and 255. If the mean is higher, more 255 are contained in this feature which means that most of the cases have black color at this pixel and vice versa. From Table 2, it can be preliminarily observed that the means for the four classes are noticeably distinct. The difference in mean will be further checked by using Student's t-test.

Table 2 Mean of Feature X210

	Mean
CL1	146.7
CL2	68.3
CL3	59.1
CL4	82.7

#### 1.1.2 Student's t-test

All the two-samples Student's t-test has Null Hypothesis as below:

$$H_0: \mu_{class A} = \mu_{class B}$$

$H_0$  can be rejected if the p-value is smaller than the significance level  $\alpha$ , which is set as 0.1 for the Student's t-tests. Table 3 shows the results from the two-samples Student's t-tests for each pair of four classes.

Table 3 Student's t-test between two classes and compare P-value with  $\alpha = 0.1$

1st Class	2nd Class	p-value	p-value < 0.1
CL1	CL2	8.07E-74	TRUE
CL1	CL3	3.51E-92	TRUE
CL1	CL4	2.92E-48	TRUE
CL2	CL3	1.88E-02	TRUE
CL2	CL4	4.79E-04	TRUE
CL3	CL4	6.60E-09	TRUE

Since all the p-values are essentially close to 0 and significantly smaller than  $\alpha$ , all the Null Hypothesis,  $H_0$ , can be rejected. The t-test results support the observation in Section 1.1.1 that the class means are significantly different.

### 1.1.3 Histogram of X210 Feature for each class

The histograms are plotted with density instead of frequency since the size are different for each class. From Figure 1, it can be observed that Class 1 (Upper Left) has a higher density of pixels with a grey level of 255, and a lower grey value of 0. The other three classes have a high density of grey levels at the 0 value and a lower density near the 255 values.

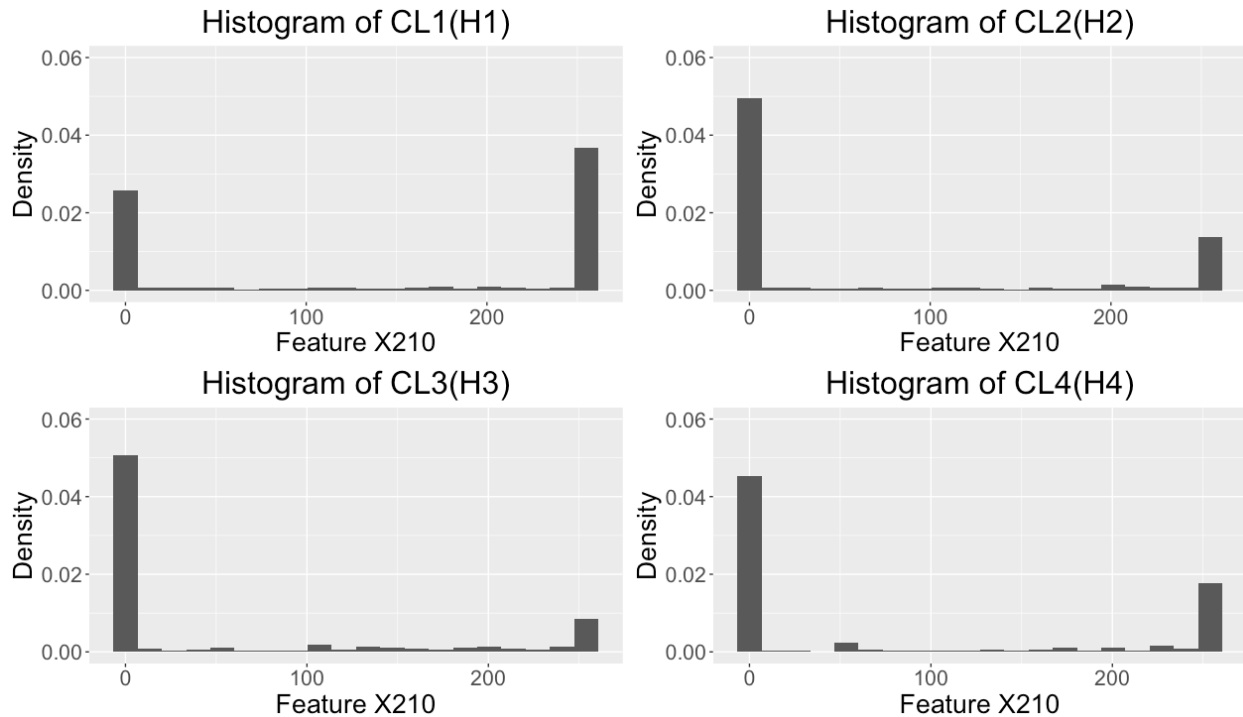


Figure 1 Histogram of X210 feature of each class. Up Left: GILL (Class 1), Up Right: LEELAWADEE (Class 2), Down Left: ROMAN (Class 3), Down Right: TECHNIC (Class 4)

However, visual comparison is not enough. Thus, the two-samples Kolmogorov-Smirnov tests which are used to compare the cumulative distribution between two classes need to be conducted to support the interpretation from histograms.

#### 1.1.4 Kolmogorov-Smirnov test

All the two-samples ks-test has Null Hypothesis as below:

$H_0$ : Class A and Class B have same distributions

$H_0$  can be rejected if the p-value is smaller than  $\alpha$ . Table 4 shows the results from the two-samples ks-tests and the p-value for each pair of four classes.

Table 4 4 Kolmogorov-Smirnov Tests between two classes

1st Class	2nd Class	p-value	p-value < 0.1
CL1	CL2	0.00E+00	TRUE
CL1	CL3	0.00E+00	TRUE
CL1	CL4	0.00E+00	TRUE
CL2	CL3	7.42E-04	TRUE
CL2	CL4	8.71E-04	TRUE
CL3	CL4	2.85E-09	TRUE

Since all p-values are essentially close to 0 and significantly smaller than  $\alpha$ , all the Null Hypothesis,  $H_0$ , can be rejected which means that all the classes have different distribution at Feature X210.

Overall, we can conclude that Feature X210 has a good discriminate power between the four classes.

### 1.2 Correlation between each feature

#### 1.2.1 Correlation Value for 400 features and top 10 pairs of correlated pixels

Table 5 shows the top 10 highest correlation values among 400 features. It can be obviously observed that 10 out of 10 of these paired pixels are located adjacent to each other.

It is reasonable because any kind of character in any font will have thickness and to represent the thickness of a font, the adjacent pixels will have to be both in black (value = 255). Meanwhile, any character will have large white space. To represent this big white space, several adjacent pixels will have to be all in white (value = 0).

Table 5 Top 10 highest correlated pixels

	Pixel 1	Pixel 2	Correlation
1	(7,19)	(8,19)	0.93
2	(9,19)	(10,19)	0.93
3	(9,18)	(10,18)	0.93
4	(19,9)	(19,10)	0.93
5	(11,0)	(12,0)	0.93
6	(9,0)	(10,0)	0.93
7	(5,9)	(5,10)	0.93
8	(18,9)	(18,10)	0.93
9	(12,1)	(13,1)	0.93
10	(10,0)	(11,0)	0.93

### 1.2.2 Correlation Value Comparison

In contrast, two pixels which are far away from each other will have less chance to have same color since different character will have different pattern. Because the data set is about the “font”, it is reasonable to take letter “L” and “H” as examples.

If the letters are cut into 20 x 20 pixels as shown in Figure 2, the left middle pixel which is highlighted by yellow box is black but the right middle pixel which is highlighted in red is white for “L” but for “H”, both are black. There are some letters fall between either the first category as letter “L” or second category as letter “H”. Thus, it is difficult to predict the color of a pixel by using a pixel far away from it.

Refer to Section 1.2.1, it can be observed that if the two pixels are adjacent, it is highly possible that they will be both black (highlighted in green in Figure 2) or both white (highlighted in blue in Figure 2) or one white and one black if they are at the transition zone of the letter and the white space.

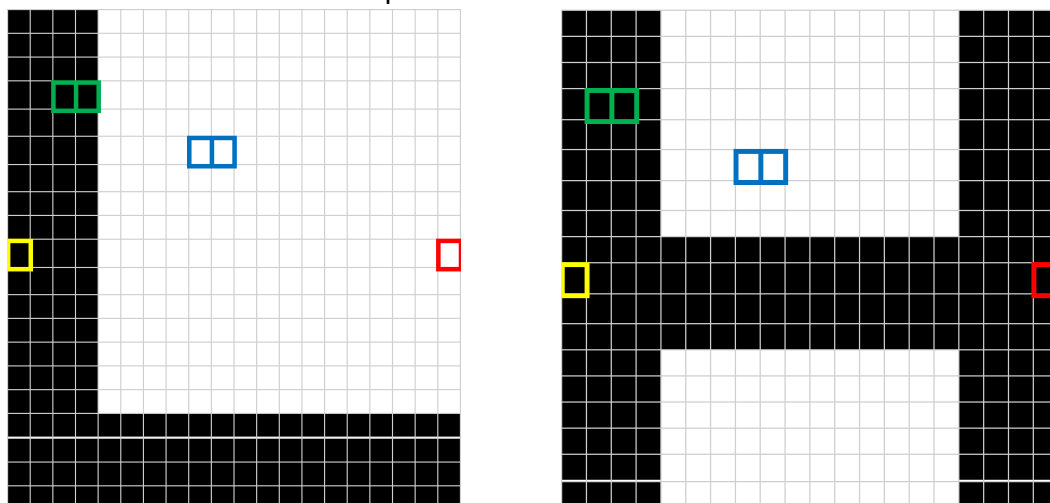


Figure 2 Letter “L” and “H” as examples to explain the correlation value

Table 6 shows a good example on this comparison. While comparing the correlation value of Pixel(6, 13) and Pixel(7, 13) which are next to each other and correlation value of Pixel(6, 1) and Pixel(6, 18) which are apart, higher correlation appears to the former one and lower correlation happens to the latter one.

Table 6 Correlation Value between Pixel(6, 13) and Pixel(7, 13) and between Pixel(6, 1) and Pixel(6, 18)

Pixel 1	Pixel 2	Correlation
(6,13)	(7,13)	0.82
(6,1)	(6,18)	0.58

### 1.3 Standardize the Data Set

To have more meaningful distance (eliminate the unit) between two cases while apply kNN algorithm in next section, all the features will be standardized as below.

$$SX(n) = \frac{[X(n) - \text{Mean}(X)]}{\text{Standard Deviation}(X)}$$

The data after standardizing will have size of 5411 x 400. Table 7 only shows the first 6 cases with first 7 features.

*Table 7 First 6 cases with First 7 features of Standardized Data*

	r0c0	r0c1	r0c2	r0c3	r0c4	r0c5	r0c6
1	-0.6364723	-0.691387	-0.7248442	-0.778314	-0.8469891	-0.7320565	-0.2619737
2	-0.6364723	-0.691387	-0.7248442	-0.778314	-0.8469891	-0.7320565	-0.2619737
3	-0.6364723	-0.691387	-0.7248442	-0.778314	-0.8469891	-0.7320565	-0.2619737
4	-0.6364723	-0.691387	-0.7248442	-0.778314	-0.8469891	-0.7320565	-0.2619737
5	-0.6364723	-0.691387	-0.7248442	-0.778314	-0.8469891	-0.9485551	-1.0578108
6	1.62125672	0.36714278	-0.1556462	-0.742095	-0.8469891	-0.9485551	-1.0578108

#### 1.4 Generate the complete Data Set with Response Variable and Features

Since kNN algorithm is a supervised classifier, it is important to have the original response variable for the data set to calculate the accuracy of prediction. Thus, each case will be assigned a number to represent its class number as the response variable. This column will be named as “TRUC” and combined with the standardized data set to get the complete data set to be ready for applying kNN in next step.

The standardized data together with true class column i.e., TRUC will have dimension of 5411 x 401. Table 8 only shows the first 6 cases with True Class and first 7 features.

*Table 8 First 6 cases with True Class and First 7 features of Standardized Data*

	TRUC	r0c0	r0c1	r0c2	r0c3	r0c4	r0c5	r0c6
1	1	-0.6364723	-0.691387	-0.7248442	-0.778314	-0.8469891	-0.7320565	-0.2619737
2	1	-0.6364723	-0.691387	-0.7248442	-0.778314	-0.8469891	-0.7320565	-0.2619737
3	1	-0.6364723	-0.691387	-0.7248442	-0.778314	-0.8469891	-0.7320565	-0.2619737
4	1	-0.6364723	-0.691387	-0.7248442	-0.778314	-0.8469891	-0.7320565	-0.2619737
5	1	-0.6364723	-0.691387	-0.7248442	-0.778314	-0.8469891	-0.9485551	-1.0578108
6	1	1.62125672	0.36714278	-0.1556462	-0.742095	-0.8469891	-0.9485551	-1.0578108



## 2 Data Automatic Classification

### 2.1 Split the Standardize Data Set into Training Set and Test Set

To reduce the bias toward any class for either training set or test set, evenly and randomly sampling the data shall be performed. The complete test set will be randomly selected from 20% of all the cases in each class and then combine as a complete test set. The remaining 80% cases in each class are combined as complete training set. Table 9 shows that there are 1080 cases of test set and 4331 cases of training set.

*Table 9 Size of Test Set and Training Set*

	Test	Train
CL1	291	1168
CL2	275	1103
CL3	238	956
CL4	276	1104
Total	1080	4331

### 2.2 Preliminary k-Nearest Neighbor (kNN) Automatic Classifier

Both training and test sets are used to determine the best value of  $k$ . Since there are no pre-defined statistical methods to find the most favorable value of  $k$ , this section and next section will repeat the kNN Automatic Classifier to find the best  $k$  which will end up with higher prediction accuracy i.e., performance in both training and test sets for the model without overfitting and underfitting. However, there are some empirical choices for the  $k$  as below:

1.  $k_{Max} \leq 50$
2.  $k_{Max} \leq \sqrt{\text{size of training set}}$

Since the complete training data set has 4331 cases, the first empirical choice is going to be applied. Thus, the preliminary kNN Automatic Classifier will be conducted for  $k = 5, 10, 15, 20, 30, 40, 50$ .

From Figure 3 below, both training and test performance curves decrease as  $k$  goes up with peak performance at  $k = 5$ . There is no overfitting or underfitting for any  $k$ . This is a good sign. Thus, it is wise to conduct a second kNN Automatic Classifier with  $k$  between 3 to 10 ( $k = 1$  or  $2$  are not good choices for any situation) to check if there is any value of  $k$  that can get a higher performance than  $k = 5$ .

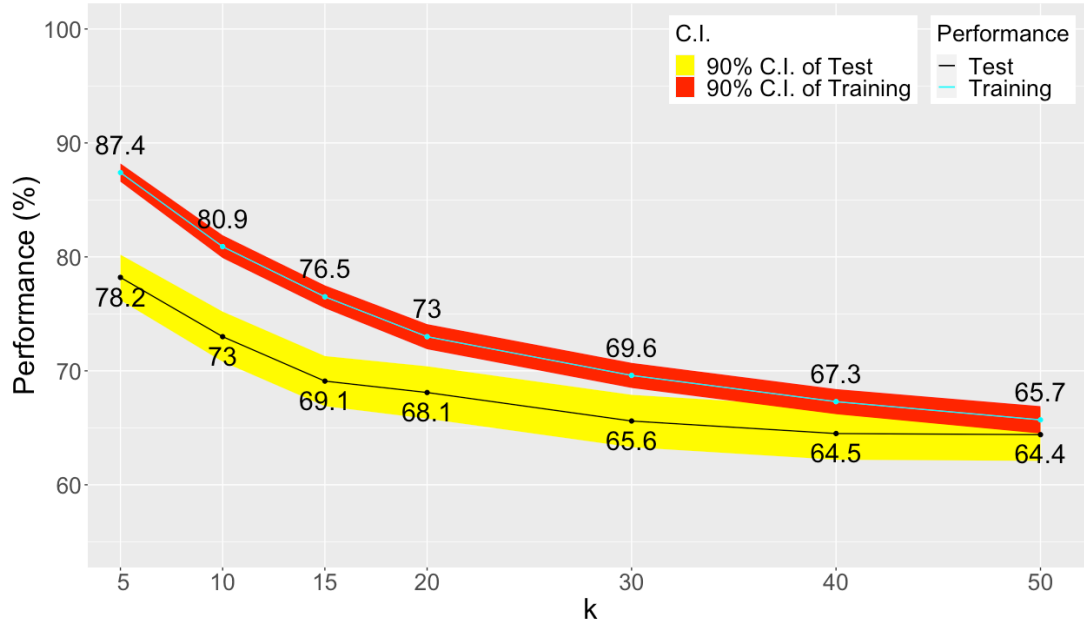


Figure 3 Training and Test Performance with  $k = 5, 10, 15, 20, 30, 40, 50$

## 2.3 Secondary k-Nearest Neighbor (kNN) Automatic Classifier

### 2.3.1 Repeat the kNN process

From Figure 4 below, it can be observed that test performance reaches around 80% while  $k = 3, 4, 5$ . There is no overfitting and underfitting for any  $k$ . Besides, 90% confidence interval of the test performance reach the highest point when  $k = 3$ . Thus, we choose  $k = 3$  as the best  $k$  for these data set.

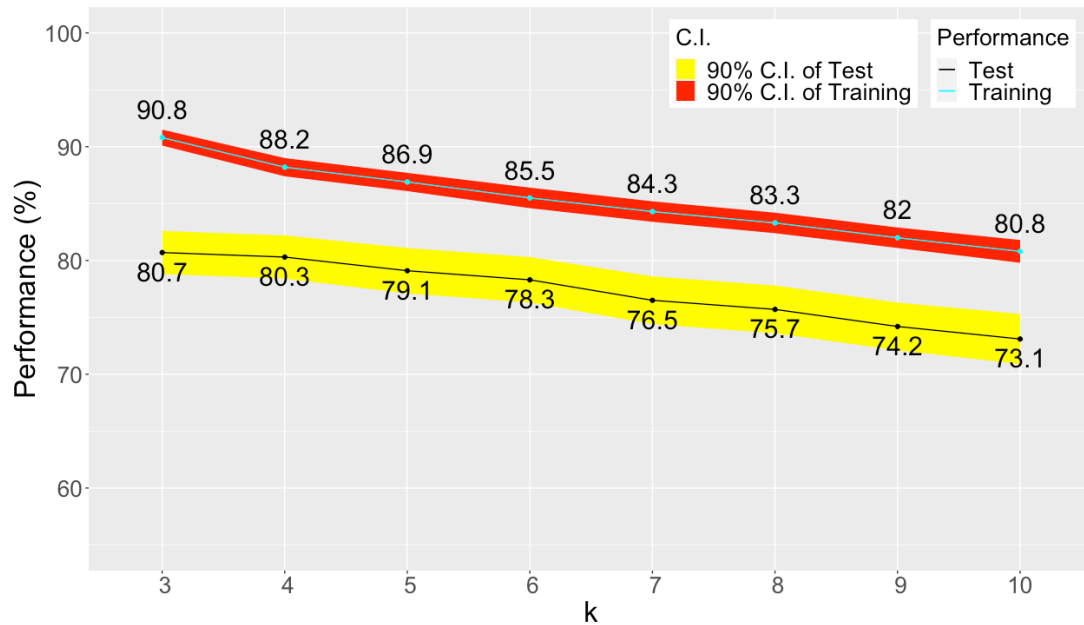


Figure 4 Training and Test Performance with  $k = 3, 4, 5, 6, 7, 8, 9, 10$

### 2.3.2 Determine the best $k$

Redo the kNN Automatic Classifier again with  $k = 3$  and compute the test performance, 90% confidence interval of the test performance, and the computing time for future comparison.

90% confidence interval = Test Performance  $\pm 1.6 \times$  Margin on Performance  
where

$$\bullet \text{ Margin on Performance} = \sqrt{\frac{\text{Test Performance} \times (1 - \text{Test Performance})}{\text{data size}}}$$

Margin on Performance will decrease as the data set size increase. Table 10 is the Performance Table for  $k = 3$ . The performing computer is Apple M1 chip (3.2 GHz) with 16 GB of RAM.

Table 10 Performance Table for Test Set when  $k = 3$

	Computing Time (s)	Test Performance (%)	90% C.I. (%)
Best k	2.8	80.8	[ 78.9 , 82.7 ]

### 2.4 Confusion Matrix

The confusion Matrix are built to have a better understanding about the correct classification and misclassification rate of each class. Table 11 shows that Class 4 has the best performance among all the classes. Meanwhile, Class 2 has the lowest performance.

Table 11 Confusion Table in Percentage when  $k = 3$ .

	True CL1	True CL2	True CL3	True CL4
Pred CL1	82.8	10.5	4.2	6.2
Pred CL2	3.4	65.5	3.8	1.8
Pred CL3	2.4	8.4	83.6	0.4
Pred CL4	11.3	15.6	8.4	91.7

Figure 5 below shows the 90% confidence interval for all the 4 classes. In conclusion, we have 90% confidence that the performance will be in below ranges for each class:

1. Class 1: Between 79.3% to 86.3%
2. Class 2: Between 60.9% to 70.1%
3. Class 3: Between 79.8% to 87.4%
4. Class 4: Between 89.0% to 94.4%

It is obvious that Class 1 and Class 3 have big portion of overlap for the confidence interval which means that even though the test performance of Class 3 is higher than Class 2, there will be some other data set cause the test performance of Class 2 higher than Class 3.

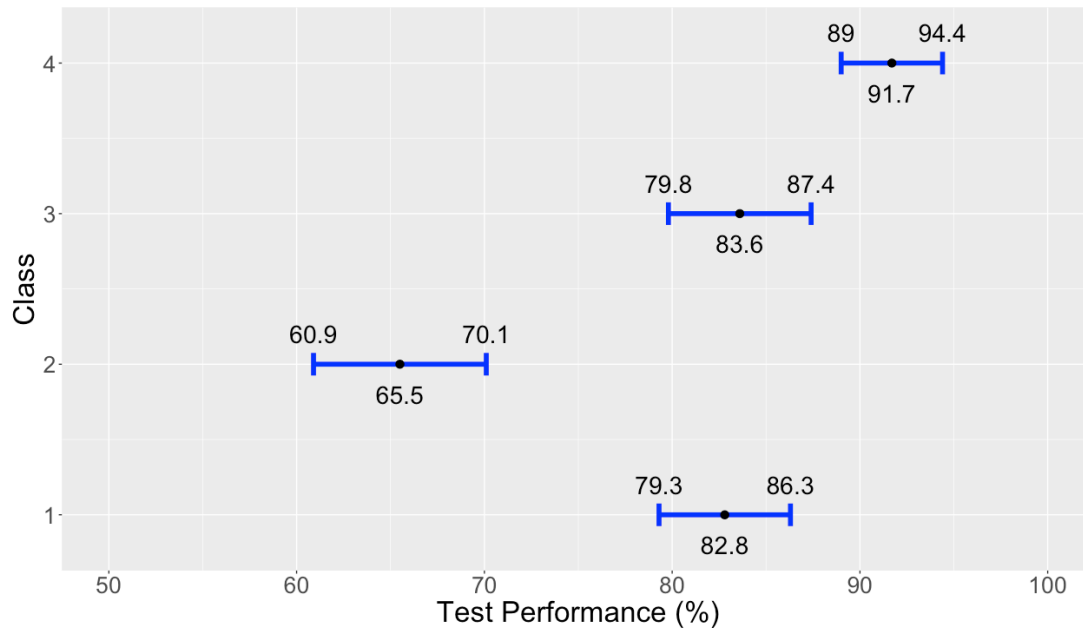


Figure 5 90% Confidence Interval of Each Class Test Performance

## 2.5 Misclassified Cases

Since the best value of  $k$  is 3, the kNN algorithm is to find out the 3 nearest neighbor of each case and figure out which classes are these neighbors in. Then, the case will then be classified as the class which the majority of its 3 nearest neighbors from. Thus, to identify why these misclassified cases are misclassified, it is wise to list all the  $k$  nearest neighbors for these cases. Table 12 shows all the cases which are under Class 2 but misclassified as Class 1, 3 or 4. The first case of each ERR list are picked and computed the distance between it and all cases.

Table 12 Cases which have True Class = 2 but misclassified as other classes. Left: Misclassified as Class 1. Middle: Misclassified as Class 3. Right: Misclassified as Class 4. The highlighted cases are selected to interpret the misclassification situation.

Test Result	Case Number	TRUC	Test Result	Case Number	TRUC	Test Result	Case Number	TRUC
1	2161	2	3	1886	2	4	1566	2
1	1468	2	3	2033	2	4	1992	2
1	2182	2	3	1951	2	4	2582	2
1	2737	2	3	1540	2	4	2815	2
1	2717	2	3	1466	2	4	2681	2
1	1928	2	3	2484	2	4	2059	2
1	2645	2	3	1591	2	4	2829	2
1	2355	2	3	1599	2	4	2836	2
1	2298	2	3	1777	2	4	2054	2
1	2531	2	3	2185	2	4	2411	2
1	2771	2	3	1552	2	4	2557	2
1	2171	2	3	1955	2	4	2685	2
1	2173	2	3	2793	2	4	2255	2
1	1666	2	3	2289	2	4	2833	2
1	2790	2	3	1845	2	4	1970	2
1	2175	2	3	1784	2	4	2786	2
1	1606	2	3	2091	2	4	2101	2
1	2048	2	3	2253	2	4	2488	2
1	2158	2	3	2117	2	4	1683	2
1	1494	2	3	2575	2	4	2713	2
1	2134	2	3	1908	2	4	2102	2
1	2763	2	3	2288	2	4	1976	2
1	2628	2	3	2591	2	4	1512	2
1	2824	2				4	2741	2
1	2142	2				4	2016	2
1	2743	2				4	1799	2
1	2295	2				4	2581	2
1	2775	2				4	2804	2
1	1486	2				4	2118	2
						4	2113	2
						4	2143	2
						4	1562	2
						4	2801	2
						4	2658	2
						4	2624	2
						4	1746	2
						4	2688	2
						4	2807	2
						4	2286	2
						4	1624	2
						4	2802	2
						4	2041	2
						4	1556	2

As per the description of knn function in R, *If there are ties for the  $k$ th nearest vector, all candidates are included in the vote.* Thus, the  $k$  nearest neighbor may include more than number of  $k$ .

Table 13 to Table 15 shows that  $k$  nearest neighbor for each misclassified case in ERR21, ERR23, and ERR24.

*Table 13 3 nearest neighbors to the case#2161 whose true class is 2 but predicted as class 1*

	Case #	Distance	TRUC
Neighbor 1	1472	6.60166414	2
Neighbor 2	884	10.0779043	1
Neighbor 3	1176	15.7395261	1

*Table 14 3 nearest neighbors to the case#1886 whose true class is 2 but predicted as class 3*

	Case #	Distance	TRUC
Neighbor 1	2575	7.92002605	2
Neighbor 2	3782	19.0384011	3
Neighbor 3	3057	19.4566534	3

*Table 15 3 nearest neighbors to the case#1566 whose true class is 2 but predicted as class 4*

	Case #	Distance	TRUC
Neighbor 1	2255	8.15466677	2
Neighbor 2	1562	11.577114	2
Neighbor 3	4054	12.976193	4
Neighbor 4	4298	12.976193	4

In Table 15, noted that Case# 1566 has 4 nearest neighbors even though the chosen  $k$  value is 3 due to the same distances for 3<sup>rd</sup> and 4<sup>th</sup> neighbor. Since there is a tie for the vote, kNN model will randomly pick from Class 2 or Class 4 which means that if we redo the kNN again, this case may not be misclassified.

In conclusion, most of the nearest neighbors for all the 3 misclassified cases don't belong to Class 2. That's the reason of them being misclassified.

### 3 Principal Component Analysis (PCA)

Dimension reduction is often a key step in automatic classification. One would like to have a new number of features much less than the number of features in the original data, but the new features should still provide as much information as the original ones. Principal Component Analysis constructs new features as linear combination of the original features. Also, the new features are standardized and decorrelated to avoid redundancy.

#### 3.1 Eigenvalue and Eigenvector

To perform the PCA, eigenvalue and eigenvector are calculated. Since there will be information loss between the data compression and decompression, it is important to keep as many information as possible. Thus, the Mean Squared Error of Decompression (MSE) i.e., information loss shall be as small as possible. Percentage of Explain Variance (PEV) which is a function of MSE is set to be 90% in this report.

$$PEV(r) = 1 - \frac{MSE(r)}{m} = \frac{\sum_{i=1}^r L_i}{m}$$

where

- $m$  is the number of original features
- $L_i$  is the sorted eigenvalue

In our case, there are 400 features in the original data set. A  $400 \times 400$  eigenvector will be computed together with 400 eigenvalue. Figure 6 shows the eigenvalue ( $L$ ) versus  $m$ .

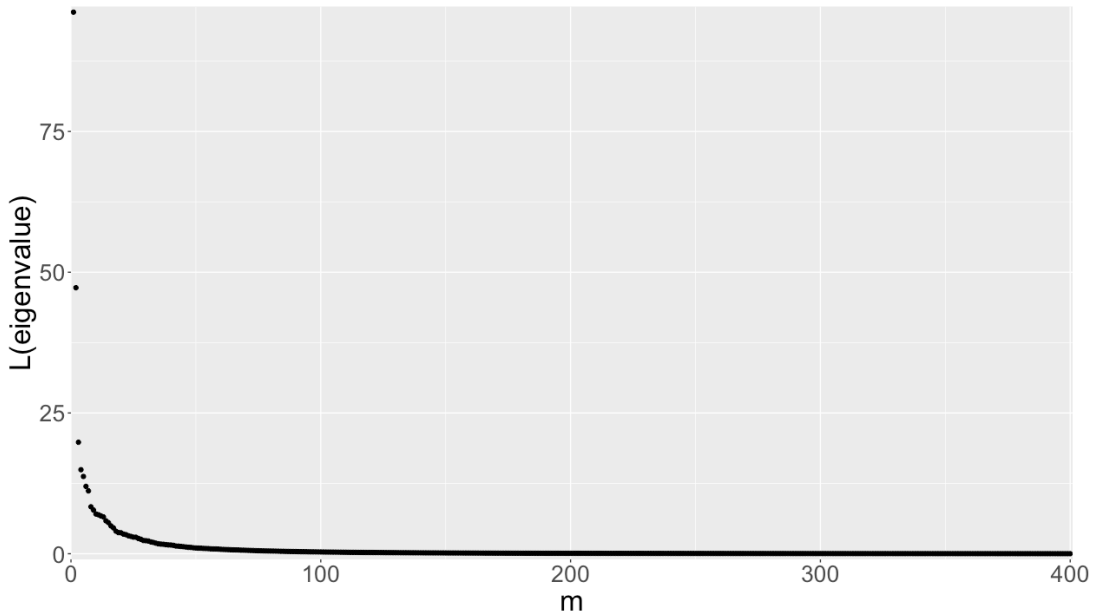


Figure 6 Eigenvalues vs m

#### 3.2 Percentage of Explained Variance (PEV)

In this report, PEV is set to be equal to or larger than 90%. Thus, the smallest integer “ $r$ ” such that  $PEV(r) \geq 90\%$  is 62. Figure 7 shows the PEV versus  $m$ .

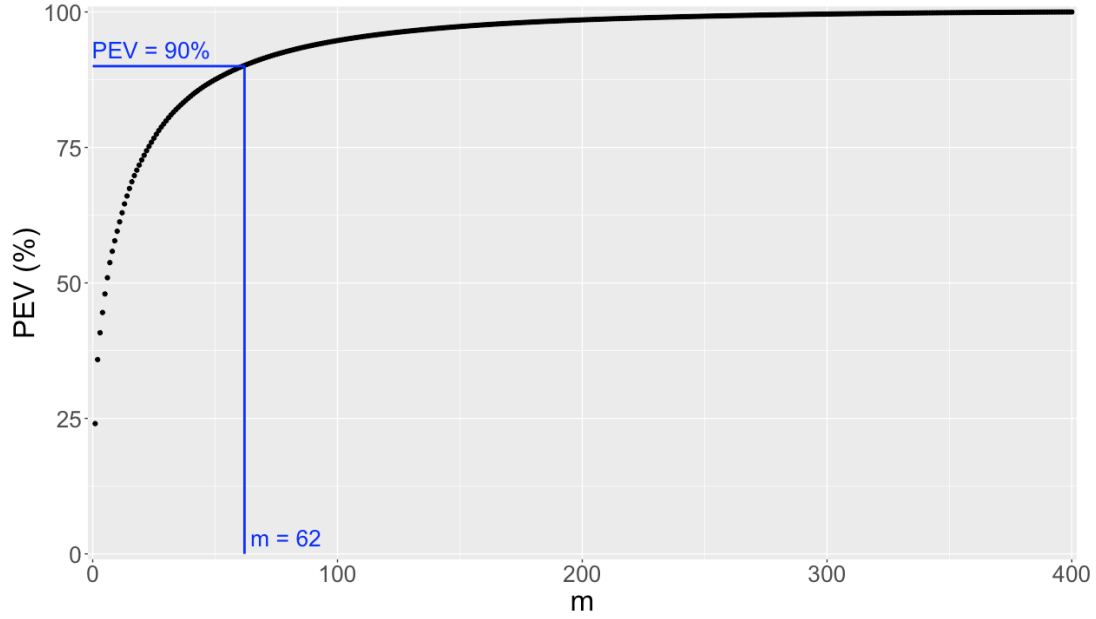


Figure 7 Percentage of Explained Variance (PEV) vs  $m$

### 3.3 Apply Principal Component Analysis

The new features after Principal Component Analysis (PCA) are produced by the multiplication of the Standardize Data (SDATA) and the Eigenvector ( $W$ ).

$$ZDATA_{N \times 400} = SDATA_{N \times 400} \times W_{400 \times 400}$$

The final dimension of the data after PCA, ZDATA, is  $5411 \times 400$ .

However, since we have  $r = 62$  to get  $PEV \geq 90\%$ , only the first 62 principal components will be taken from ZDATA to have a data set with size of  $5411 \times 62$ . Together with the TRUC column, the final data set with size of  $5411 \times 63$  is generated.

Since the data after PCA has dimension of  $5411 \times 63$ , Table 16 only shows the first 6 cases with TRUC and first 7 features in below.

Table 16 First 6 cases with First 7 features of Data after PCA

	TRUC	V1	V2	V3	V4	V5	V6	V7
1	1	20.3309051	-1.2448531	-2.770295	2.10850216	2.09203383	-1.4044534	2.92570844
2	1	17.8302992	-0.657016	-4.1377791	3.397947	-0.1339552	0.72800757	2.77131667
3	1	20.3309051	-1.2448531	-2.770295	2.10850216	2.09203383	-1.4044534	2.92570844
4	1	17.8302992	-0.657016	-4.1377791	3.397947	-0.1339552	0.72800757	2.77131667
5	1	-5.7027004	3.04603417	-5.9035089	-1.2941225	4.68990459	0.80037051	-2.4853168
6	1	7.2542036	1.27501693	1.64668112	1.47610997	-5.8213176	-1.2358181	1.10017072

### 3.4 Apply kNN Automatic Classifier to Data Set after Principal Component Analysis

To know the power of PCA and the performance against the original data, the newly obtained PCA data set will be applied to kNN automatic classifier at best  $k$  as concluded in Section 2.3. Figure 8 shows the comparison of test performance between original data set and data set after PCA. It is observed that two kinds of data sets have a big portion of



overlapping in their confidence intervals which means doing PCA will not compromise too much on the performance for kNN.

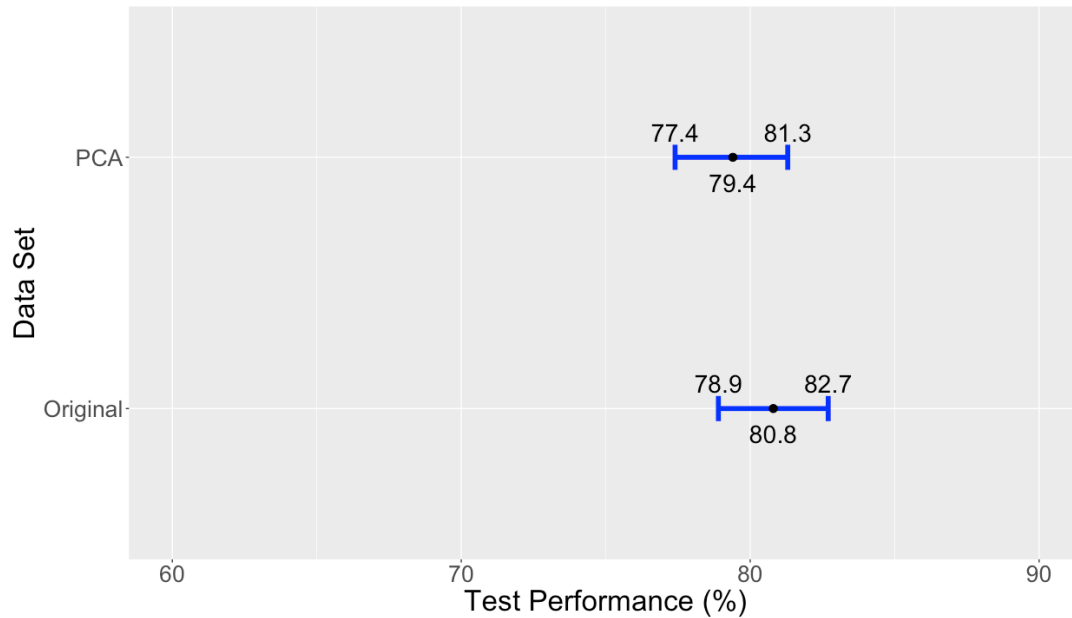


Figure 8 90% Confidence Interval of Test Performance of Original Data and PCA Data

Table 17 shows that the Test Performance decreased by 1.4% while it can save up to 72.3% of the computing time which is a good trade off.

Table 17 Comparison between Standardize Data Set and Data after Principal Component Analysis

	Computing Time (s)	Test Performance (%)
Original	2.8	80.8
PCA	0.9	79.4

### 3.5 Visualize the First Two Component

From  $PEV$  versus  $m$  in Figure 7, we can know that  $PEV(2) \cong 36\%$ . Therefore, first 2 principal components can only explain about 36% of information provided by the features, which is not powerful enough to accurately separate the classes on a 2D diagram. However, we may still get some information or surprise from these plots.

In Figure 9, we can observe a good separation between Class 1 and Class 2. Most of the Class 1 spread at right side and Class 2 spread at left side of the figure. However, there are still some overlapping around  $V1 = 0$ . This may cause the difficulty of classification.

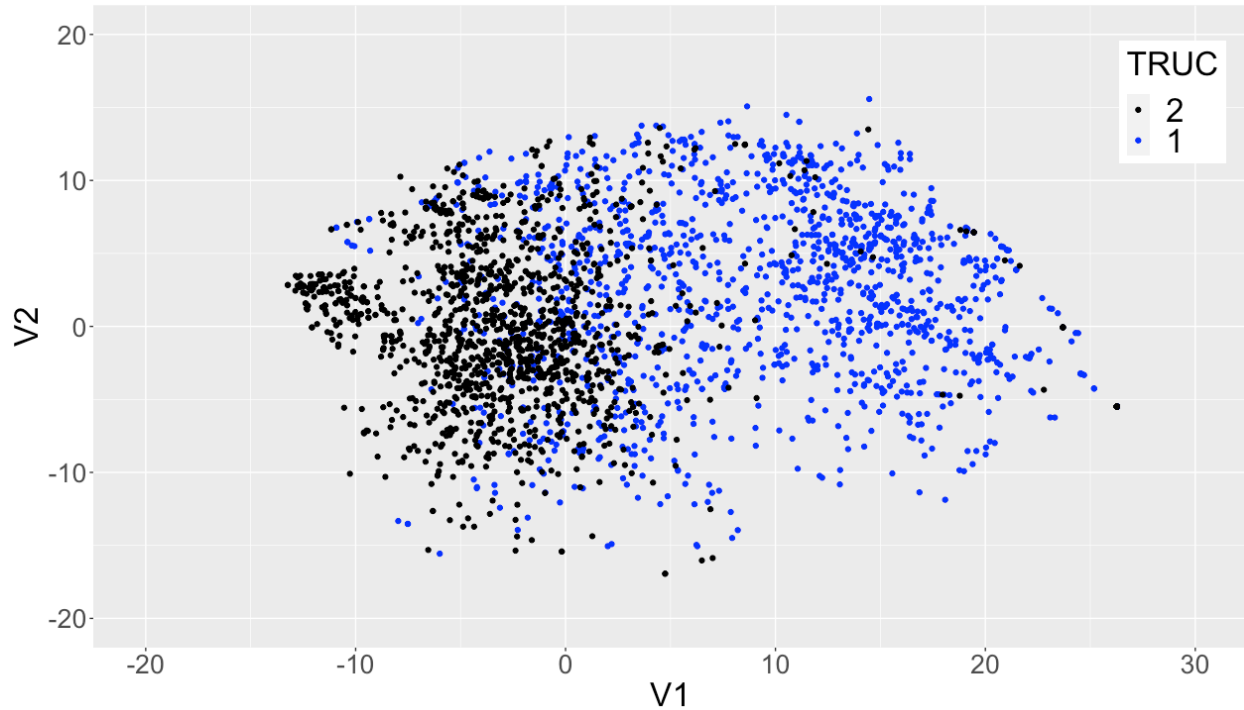


Figure 9 Class 2 and Class 1 Distribution on a planar graph based on first 2 principal components, V1 and V2.

In Figure 10, we can observe Class 3 and Class 2 have a slightly worse separation compared to Class 1 and Class 2, since there are some points of Class 2 which are covered by the points for Class 3 around V1 from -15 to 0 and V2 from -5 to 5.

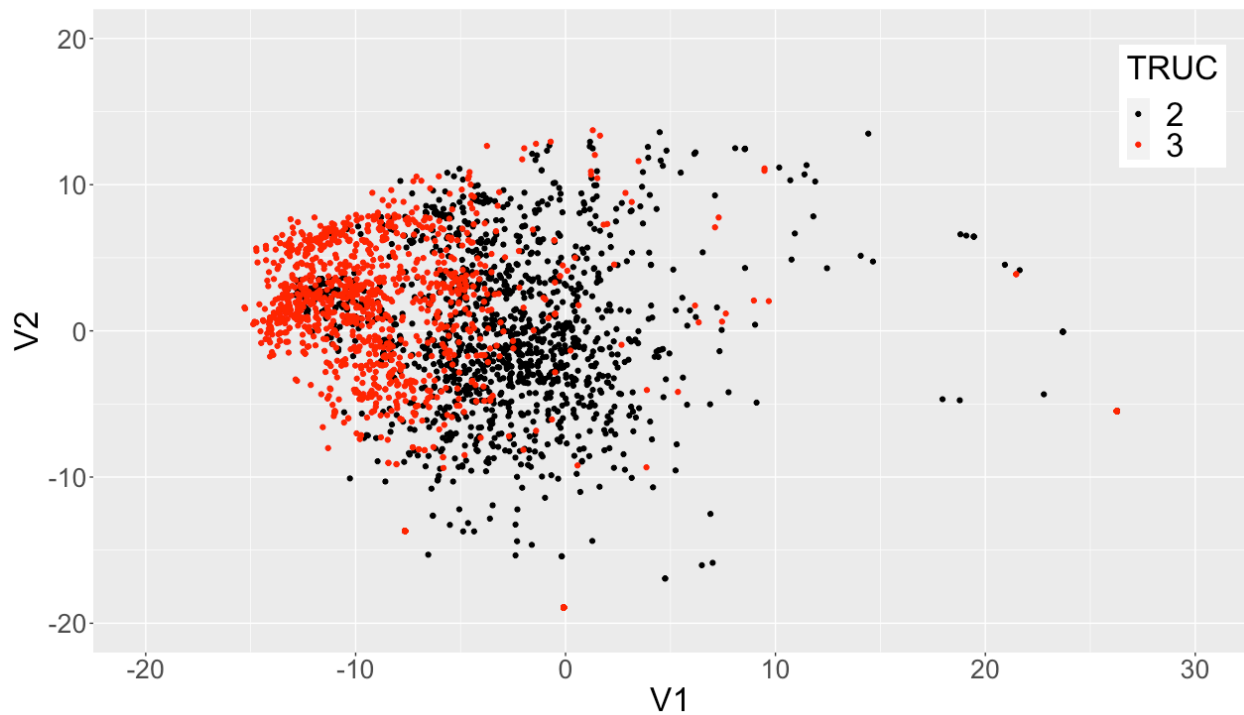


Figure 10 Class 2 and Class 3 Distribution on a planar graph based on first 2 principal components, V1 and V2.

In Figure 11, we can observe a poor separation between Class 2 and Class 4. It is difficult to distinguish the two classes from the plot. Most of the Class 2 and Class 4 are spread between V1 from -10 to 5 and V2 from -10 to 10. This is not a surprise that among all the misclassified situations, there are 15.6% of the Class 2 cases being classified as Class 4 which is the highest error rate in the entire confusion matrix in Section 2.4.

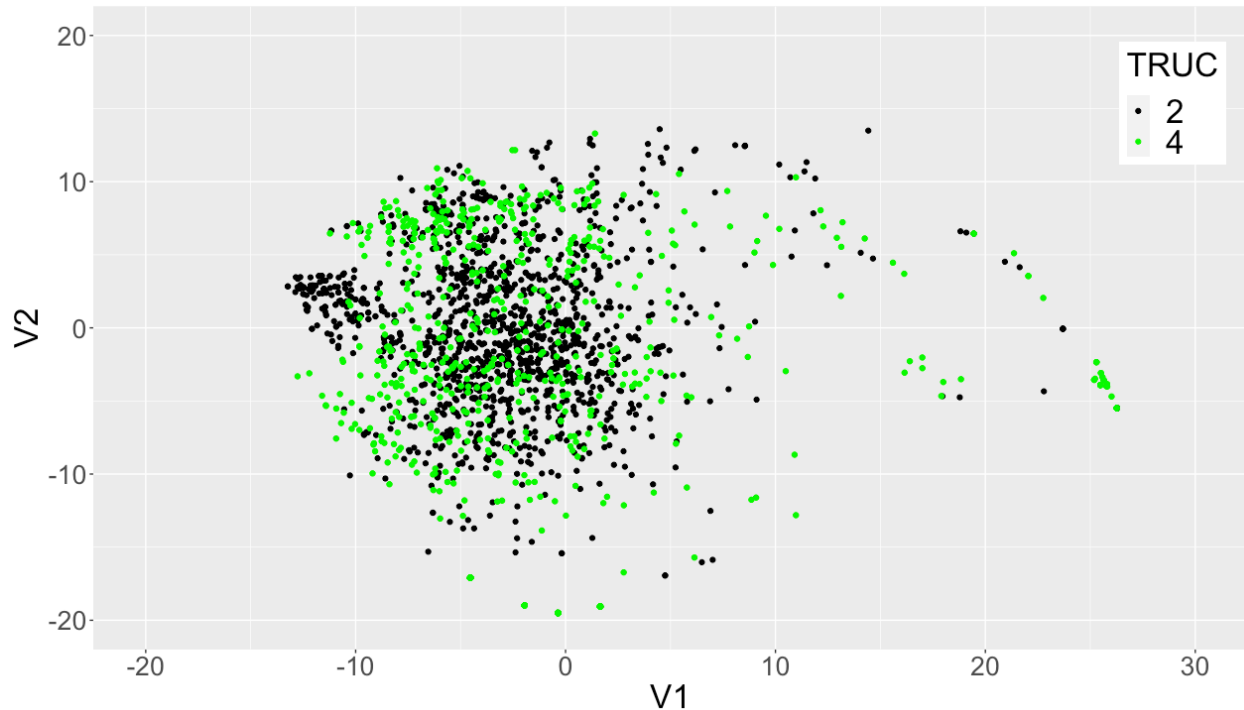


Figure 11 Class 2 and Class 4 Distribution on a planar graph based on first 2 principal components, V1 and V2.

### 3.6 Visualize misclassified cases

Figure 12 demonstrates that the misclassified Class 2 points seem to correspond with the range and density of the points in the other classes. For example, the Class 2 points that were misclassified as Class 3 are in the net area of V1 from -15 to 0 and V2 from -10 to 10, which is where the Class 3 points are mostly concentrated. Points misclassified as Class 1 and Class 4 also are in areas with high concentrations of their incorrect classes.

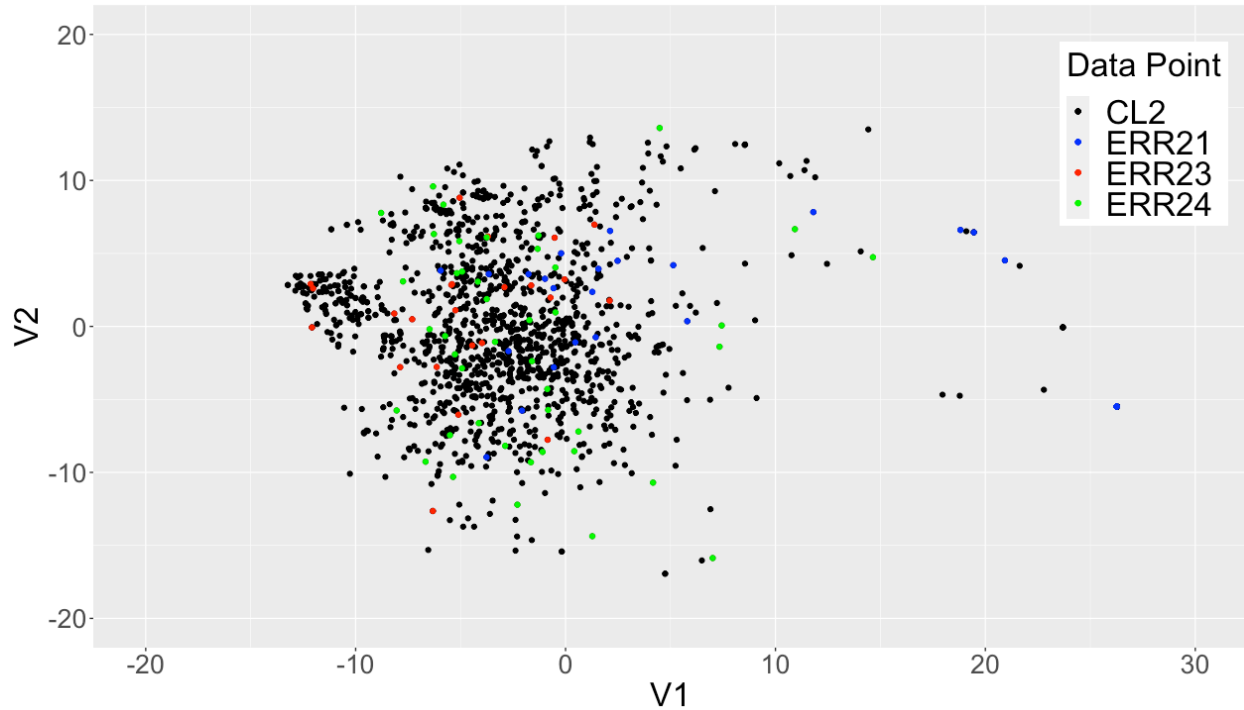


Figure 12 All the Class 2 cases on the planar projection and highlighted cases which shall be belongs to Class 2 but misclassified as Class 1, 3, or 4

As mentioned before, since  $PEV(2) \cong 36\%$ , there are a lot of information is missing while we try to classify the point visually on a 2-dimensional plot. Even if we plot them into 3-dimensional, i.e.,  $PEV(3)$  we still have only around 40% of the information. It is hard to visually identify the which class shall each case belongs to. So, the best way to classify these data is using the kNN algorithm which means to calculate the distance between cases and figure out which are the k-nearest neighbor so that it can vote for the classification and interpret mathematically and easily.