

PREPARED FOR SUBMISSION TO UNIVERSITY OF CAMBRIDGE

# **M1: Applied Data Science - Coursework Assignment**

**Adnan Siddiquei**

University of Cambridge

E-mail: [as3438@cam.ac.uk](mailto:as3438@cam.ac.uk)

---

## Contents

<b>1</b>	<b>Section A</b>	<b>1</b>
1.1	Q1 - Dataset A	1
1.1.1	Question 1a	1
1.1.2	Question 1b	1
1.1.3	Questions 1c and 1d	1
1.1.4	Questions 1e	3
1.2	Q2 - Dataset B	5
1.3	Q3 - Dataset C	6
<b>2</b>	<b>Section B</b>	<b>7</b>
<b>A</b>	<b>Q2: Duplicate Observations in <code>B_NoiseAdded.csv</code></b>	<b>7</b>

---

## 1 Section A

This section contains the answers to the questions in Section A of the coursework assignment. Wherever the phrase 'dataset' is used in this section, it refers to `A_NoiseAdded.csv`.

### 1.1 Q1 - Dataset A

#### 1.1.1 Question 1a

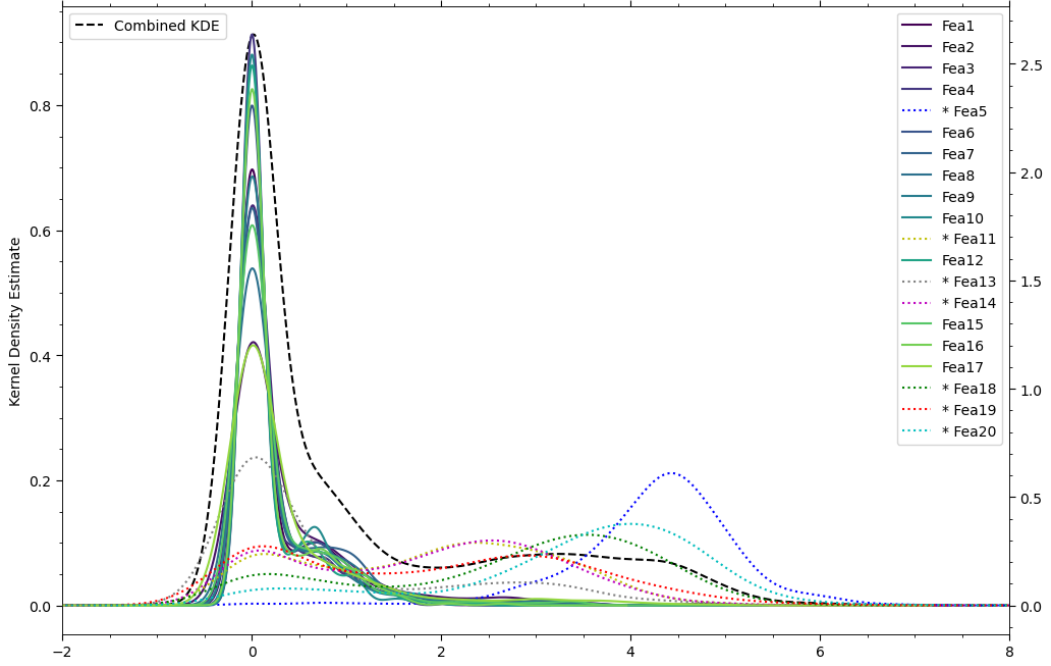
Fig(1) shows the combined and separated kernel density estimates of the first 20 features of the dataset. Several conclusions can be drawn from this plot. For most of the features, the KDE is concentrated around the zero value with very little variance, with the exception of 7 features: 5, 11, 13, 14, 18, 19 and 20. These features are likely to be more discriminative within classification algorithms, as they contain more variability.

#### 1.1.2 Question 1b

The biplot in Fig(2) shows a PCA of the entire dataset, after standardisation has been applied to each feature. The coloured arrows indicate the loadings of the features which contribute more than 2% to either principal component, which implies that none of the first 20 features are particularly discriminative. An interesting observation is that the most discriminative features are very discriminative with respect to the second principal component, whereas the less discriminative features (of which there are many, many more of) are more discriminative with respect to the first principal component. As a result, the scores on the biplot for the observations w.r.t. the first two principal components separate the observations more along the first principal component, rather than the second.

#### 1.1.3 Questions 1c and 1d

Fig(3a) shows the contingency table for two k-means clusterings of the dataset with  $k = 8$  and  $k = 3$ . The two clusterings (`kmeans_1` and `kmeans_2`) were formed by training two k-means models on half of the dataset, and then the other half were mapped onto the learned clusters. Standardisation was applied to each feature before clustering as k-means is sensitive to feature scaling, this is because k-means works by computing a centroid for each cluster such that the

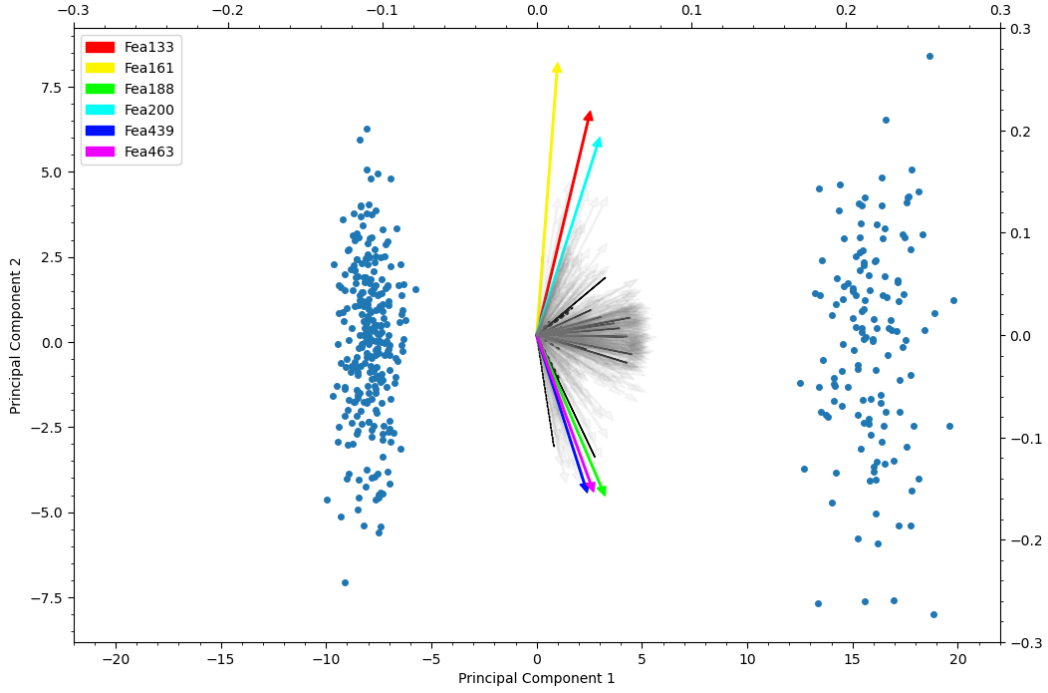


**Figure 1:** A kernel density estimate of the first 20 features of the A\_NoiseAdded.csv dataset, which has 20 plots in total, with the legend and corresponding y-axis on the right. A combined KDE has also been plotted, with its corresponding x-axis on the left. 7 of the 20 features have been highlighted in the legend with an asterisk, and have been coloured in slightly more contrasting colours and plotted with a dotted line. These features have a larger variance in their density, and are therefore more likely to be more discriminative.

sum of the squared Euclidean distances between each centroid and the observations in the cluster is minimised. Therefore, without standardisation, features with larger variances and scale will dominate the clustering. Predictions are then made by assigning new observations to the cluster whose centroid is closest to the observation [1]. Crucially, the labels assigned to each cluster from each k-means computation are arbitrary, therefore, the labels in `kmeans_2` were mapped back to the labels in `kmeans_1` by finding which centroids in each clustering were closest to each other. This step was crucial, otherwise the contingency table would be meaningless.

For  $k = 8$ , 59% of the observations lie on the leading diagonal, which indicates that the two clusterings were not very similar or stable, as a large proportion of the observations were assigned to different clusters in each clustering. `kmeans_1` clustered 85% of the observations into clusters 1, 2 and 3. `kmeans_2` clustered 86% of the observations into clusters 1 and 2. This indicates that both clusterings identified most of the data lives within a small set of clusters, which was expected as the dataset is labelled and so it is known that there are only 3 clusters - both clusterings identified at least 2 similar clusters (1 and 2). However, the presence of very small clusters indicates that there may be outliers present in the dataset or smaller clusters that the large  $k = 8$  is overfitting to.

Fig(3b) shows the contingency table for two k-means clusterings of the dataset with  $k = 3$ . A larger proportion of the observations lie on the leading diagonal compared to  $k = 8$ , which is expected because the number of clusters is both smaller and equal to the



**Figure 2:** A biplot of the first two principal components in the `A_NoiseAdded.csv` dataset. The coloured arrows indicate the first two principal component loading vectors for every feature which contributes to either principal component more than 2%, these are the most discriminative features. Every other loading vector has been plotted in very light grey so their general directions and magnitudes are visible. The loading vectors for the first 20 features have also been plotted in darker black lines, with the dotted black lines corresponding to the dotted KDEs in Fig(1). The loading vectors use the right and top axis of the plot. The blue dots indicate the scores for each observation in the dataset for the first two principal components. The scores use the left and bottom axis of the plot. All four axis are symmetric about the origin for ease of comparison.

actual number of clusters in the dataset. `kmeans_2` identified 3 distinct clusters whereas `kmeans_1` only identified 2 distinct clusters with a few remnant observations in assigned to cluster 2.

#### 1.1.4 Questions 1e

Fig(4) shows the k-means clusterings on the PCA plot shown in Fig(2). The PCA indicates that there are two clusters in the dataset when it is reduced to the first two principal components (which generally explain a large amount of the variance in a dataset). Fig(4) provides subtle evidence in favour of this as well, note that clusters 1 and 3 are stable, but cluster 2 jumps between the two groups in the PCA plots, whilst only capturing observations from either one of the groups but never both. The fact that cluster 2 never captures data from both groups in the PCA plot indicates that it is likely capturing a subset of the data from one of the PCA groups, based on random initialisation of the k-means algorithm. Its inability to capture data from both PCA groups at once, and the clear separation of clusters 1 and 3 indicate that there may only be two clusters in the dataset, as opposed to the three clusters

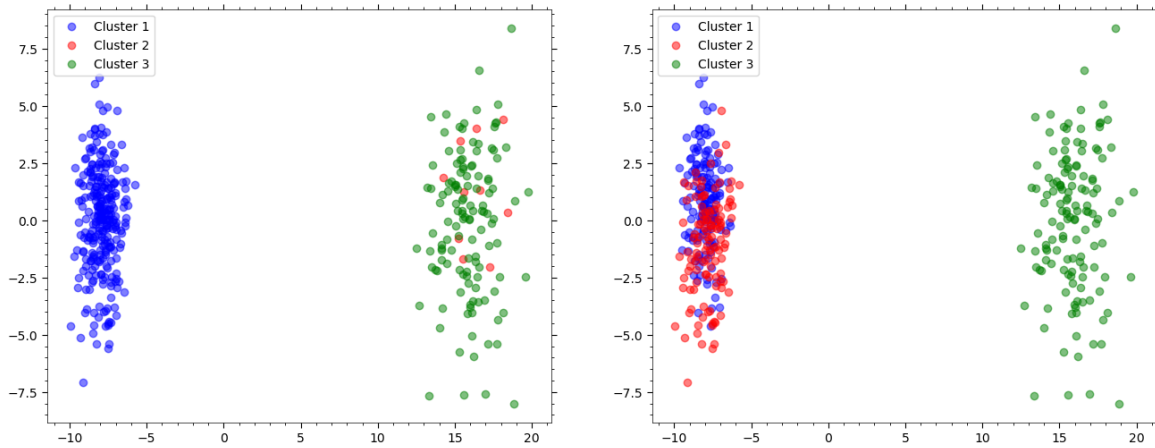
	Cluster 1	Cluster 2	Cluster 3	Cluster 4	Cluster 5	Cluster 6	Cluster 7	Cluster 8	Total
Cluster 1	<b>43</b>	0	9	8	0	4	14	1	<b>79</b>
Cluster 2	0	<b>185</b>	0	0	0	0	0	0	<b>185</b>
Cluster 3	0	87	<b>0</b>	0	0	0	0	0	<b>87</b>
Cluster 4	12	0	2	<b>5</b>	0	3	6	0	<b>28</b>
Cluster 5	1	0	0	0	<b>0</b>	0	0	0	<b>1</b>
Cluster 6	6	0	0	0	0	<b>0</b>	0	0	<b>6</b>
Cluster 7	12	0	0	0	1	1	<b>7</b>	0	<b>21</b>
Cluster 8	1	0	0	0	0	0	0	<b>0</b>	<b>1</b>
Total	<b>75</b>	<b>272</b>	<b>11</b>	<b>13</b>	<b>1</b>	<b>8</b>	<b>27</b>	<b>1</b>	

(a) A contingency table for two k-means clusterings of the `A_NoiseAdded.csv` dataset, with  $k = 8$ , the default `scikit-learn` value. 240 of the 408 observations lie on the leading diagonal.

	Cluster 1	Cluster 2	Cluster 3	Total
Cluster 1	<b>155</b>	117	0	<b>272</b>
Cluster 2	0	<b>0</b>	10	<b>10</b>
Cluster 3	0	0	<b>126</b>	<b>126</b>
Total	<b>155</b>	<b>117</b>	<b>136</b>	

(b) A contingency table for two k-means clusterings of the `A_NoiseAdded.csv` dataset, with  $k = 3$ . 281 of the 408 observations lie on the leading diagonal.

**Figure 3:** Contingency tables for two k-means clusterings of the `A_NoiseAdded.csv` dataset with number of clusters  $k = 3$  and  $k = 8$ . Each feature in the dataset was standardised before clustering. `kmeans_1` totals are on the right and `kmeans_2` totals are on the bottom.



**Figure 4:** The k-means clusterings performed with  $k = 3$ , shown on the contingency table in Fig(3b), plotted on the first two principal components of the dataset shown in Fig(2). The left plot is `kmeans_1` and the right plot is `kmeans_2`.

that the labels indicate.

Performing k-means before PCA has the advantage of being able to visualise the clusters separation in the original feature space, whilst performing PCA before k-means has the advantage of being able to reduce computational load and identify clusters along the first two principal components (which often explain the most variance in the dataset). Performing PCA first also provides a visual way to determine the number of clusters in the dataset if is

	Count
1	179
2	157
4	72
Missing	20
<b>Total</b>	<b>428</b>

(a) Raw `B_NoiseAdded.csv` dataset, before pre-processing. There are 428 observations with 20 missing labels and 20 duplicated observations (40 observations involved in the duplication).

	Count
1	177 (-3, 5, -12, 8)
2	163 (-3, 4, -5, 10)
4	68 (-4, 1, -3, 2)
Missing	0
<b>Total</b>	<b>408</b>

(b) `B_NoiseAdded.csv` after pre-processing, no missing labels or duplicated observations. The 4 numbers in the brackets indicate how the count in each classification changed due to, in order: 1 and 2 are counts exiting and entering (respectively) the class after correcting for mislabelling; 3 - count exiting after dropping duplicates; 4 - count entering after imputation of missing labels.

**Figure 5:** Summary of classifications for the `B_NoiseAdded.csv` dataset before and after pre-processing.

is unknown as PCA will separate data along the two most discriminative axis. If the number of clusters is unknown, performing PCA first would be a good idea. As to which is better, it can depend on the dataset and the number of features. It can often turn out that the first two principal components do not actually explain much of the variance, and as such, in these cases it can be better to perform k-means first.

## 1.2 Q2 - Dataset B

Fig.(5a) shows the summary of classifications for the `B_NoiseAdded.csv` dataset. It was identified that the dataset contained 20 duplicated observations (a total of 40 observations involved in the duplication), and 10 of these duplicates contained different labels across the two duplicates. See Appendix A for the full list of duplicates. The correct assignment for this mislabelling was determined using multinomial logistic regression on the labelled data to predict the labels on the mislabelled data. Multinomial logistic regression was then used again to predict the labels for the 20 missing observations, and Fig.(5b) shows the new summary of classifications.

Generally, there are multiple ways to handle missing labels. One such way is model based imputation, which was used here. This is where an appropriate model is trained on the labelled data to predict the missing labels. Another option could be to ignore the data with the missing labels, if the sample size is sufficiently large enough. Model based imputation has the advantage of using all the data, however it can introduce bias if the model is not a great fit for the data. Ignoring the data is advantageous in that it won't create bias if the sample size is large enough.

	Fea58	Fea142	Fea150	Fea233	Fea269	Fea299	Fea339	Fea355	Fea458	Fea466	Fea491
138	nan	nan	nan	nan	nan	nan	nan	nan	nan	nan	nan
143	nan	nan	nan	nan	nan	nan	nan	nan	nan	nan	nan
231	nan	nan	nan	nan	nan	nan	nan	nan	nan	nan	nan
263	nan	nan	nan	nan	nan	nan	nan	nan	nan	nan	nan
389	nan	nan	nan	nan	nan	nan	nan	nan	nan	nan	nan

(a) The samples and features with missing data.

	Fea58	Fea142	Fea150	Fea233	Fea269	Fea299	Fea339	Fea355	Fea458	Fea466	Fea491
138	0.0	4.0	0.36	1.08	0.0	0.0	2.31	0.0	1.42	0.0	0.0
143	0.0	4.25	0.0	0.64	0.32	0.0	2.55	0.0	1.09	0.72	0.0
231	0.0	4.14	0.14	1.0	0.0	0.0	1.65	0.0	0.5	0.0	0.0
263	0.0	4.33	0.0	0.58	0.0	0.0	1.91	0.0	0.92	0.14	0.0
389	0.0	4.86	0.0	0.14	0.14	0.0	1.87	0.0	0.0	1.04	0.0

(b) The imputed data.

**Figure 6:** Samples and features with missing data, for the `C_MissingData.csv` dataset.

Missing at random (MAR) is when the likelihood that a label is missing is independent of the label itself, and missing not at random (MNAR) is when the likelihood of the label missing is in some way correlated with the value of the label. Looking at the 4th number in the brackets of Fig.(5b), the missing labels in each class as percentages of the total count in the respective classes are 4.5%, 6.1% and 2.9% (for 1, 2, and 4 respectively). This gives no significant indication of MNAR.

### 1.3 Q3 - Dataset C

Fig(6a) shows all the features that have missing values, and all the samples which are missing those features. There are a few ways to handle missing data. A straightforward method is static imputation where every sample missing a given feature, is imputed with the same value for that feature, such as the mean of that feature. This is computationally inexpensive but can reduce the variance in the dataset and introduce bias if there are many missing values. Another method is model based imputation where an applicable model is chosen to estimate the value of the missing feature for each sample. An example is the K Nearest Neighbours approach, where a sample missing a feature is imputed with the mean of the K nearest neighbours. This can reduce the bias introduced and leave the variance unaffected, but is heavily dependent on the correct choice of model. Multiple imputation is where the missing data is imputed using a probabilistic model multiple times, to create multiple datasets. The multiple imputed values can then either be averaged or the multiple datasets can then be individually analysed, Multiple imputation helps capture the uncertainty in what the missing data is, which is useful in the case that there is a large amount of missing data.

Fig(6b) shows the imputed data, with aforementioned KNN approach using  $k = 5$ , and also by first splitting samples with missing data into their 3 separate classes before applying KNN. This model was chosen because the data is known to be a classification dataset and as such, would lead to the best imputation as nearby data would most likely resemble what the missing data would be.

The imputed data does not change the distributions of the features by any statistically significant amount. Fig(7) shows how the variances of the KNN imputed data compares to the original dataset and a mean imputed version of the dataset. This indicates that the KNN imputation retains more of the original variance in the dataset compared to a static

	Var(KNN) / Var(orig) %	Var(mean) / Var(orig) %	Var(KNN) / Var(mean) %	KS(orig, KNN)
Fea58	-1.21	-1.22	0.02	1.0
Fea142	-1.0	-1.12	0.12	1.0
Fea150	-0.37	-1.22	0.87	1.0
Fea233	-1.12	-1.22	0.11	1.0
Fea269	-1.15	-1.23	0.08	1.0
Fea299	-1.22	-1.22	0.0	1.0
Fea339	-1.1	-1.2	0.1	1.0
Fea355	-1.22	-1.22	-0.0	1.0
Fea458	-1.03	-1.2	0.17	1.0
Fea466	-0.29	-1.02	0.74	1.0
Fea491	-1.22	-1.21	-0.0	1.0

**Figure 7:** A comparison of the variances of each feature after imputation. Column 1 shows the variance of each feature after KNN impuration, as a percentage of the original dataset. Column 2 shows this for mean imputation. Column 3 compares column 1 against column 2. Column 4 shows the p-value of the Kolmogorov-Smirnov test comparing the original dataset to the KNN imputed dataset.

mean imputation. Likewise, the Kolmogorov-Smirnov (KS) test was used to compare the distributions of the original dataset and the KNN imputed dataset, and the p-value of this test for each feature indicates that KNN imputation did not affect the original distributions of the features by any significant amount.

## 2 Section B

### A Q2: Duplicate Observations in `B_NoiseAdded.csv`

The duplicated pairs in `B_NoiseAdded.csv`, the number is the sample number: [[ 74 146], [220 291], [147 409], [ 44 193], [ 66 253], [ 28 260], [384 396], [188 249], [ 83 198], [175 311], [166 424] [344 389], [117 297], [351 352], [120 359], [119 382], [100 173], [ 30 101], [ 46 107], [210 305]]

## References

- [1] scikit-learn developers, *scikit-learn k-means documentation*. Available at: <https://scikit-learn.org/stable/modules/clustering.html#k-means> [Accessed: 13-Dec-2023].