



---

# Different Machine Learning Approaches to Quantify the Relationship of Spectral Data on Virus Load Estimation and Virus Classification

---

Module Report

submitted towards differentiation for the degree of

Msc Artificial Intelligence

by

**Joseph McInerney**

EE ECS

Queen's University Belfast

**Lecturer**

Professor Hui Wang

November 2024

# Contents

|          |   |           |
|----------|---|-----------|
| <b>1</b> | <b>Introduction</b>                             | <b>2</b>  |
| <b>2</b> | <b>Methodology</b>                              | <b>3</b>  |
| 2.1      | Regression . . . . .                            | 3         |
| 2.1.1    | Data Preprocessing . . . . .                    | 3         |
| 2.1.2    | Feature Engineering . . . . .                   | 3         |
| 2.1.3    | Model Selection . . . . .                       | 4         |
| 2.2      | Classification . . . . .                        | 4         |
| 2.2.1    | Preprocessing . . . . .                         | 4         |
| 2.2.2    | Feature Engineering . . . . .                   | 4         |
| 2.2.3    | Model Selection . . . . .                       | 5         |
| 2.3      | Clustering . . . . .                            | 5         |
| 2.3.1    | Preprocessing . . . . .                         | 5         |
| 2.3.2    | Feature Engineering . . . . .                   | 6         |
| 2.3.3    | Model Selection . . . . .                       | 6         |
| <b>3</b> | <b>Results</b>                                  | <b>7</b>  |
| 3.1      | Regression . . . . .                            | 7         |
| 3.1.1    | Multicollinearity . . . . .                     | 7         |
| 3.1.2    | Dimensionality Reduction Comparison . . . . .   | 7         |
| 3.1.3    | Bias Variance Model Comparison on MSE . . . . . | 8         |
| 3.2      | Classification . . . . .                        | 8         |
| 3.2.1    | Model Comparison . . . . .                      | 8         |
| 3.3      | Clustering . . . . .                            | 8         |
| 3.3.1    | K-means . . . . .                               | 8         |
| 3.3.2    | AGC . . . . .                                   | 10        |
| 3.3.3    | DBSCAN . . . . .                                | 11        |
| <b>4</b> | <b>Discussion</b>                               | <b>12</b> |
| 4.1      | Regression . . . . .                            | 12        |
| 4.1.1    | Multicollinearity . . . . .                     | 12        |
| 4.1.2    | Dimensionality Reduction Comparison . . . . .   | 12        |
| 4.1.3    | Bias Variance Model Comparison on MSE . . . . . | 12        |
| 4.2      | Classification . . . . .                        | 12        |
| 4.3      | Clustering . . . . .                            | 12        |
| 4.3.1    | K-means . . . . .                               | 12        |
| 4.3.2    | AGC . . . . .                                   | 13        |
| 4.3.3    | DBSCAN . . . . .                                | 13        |

|          |                          |           |
|----------|--------------------------|-----------|
| <b>5</b> | <b>Conclusion</b>        | <b>14</b> |
| 5.1      | Regression . . . . .     | 14        |
| 5.2      | Classification . . . . . | 14        |
| 5.3      | Clustering . . . . .     | 14        |
|          | <b>Bibliography</b>      | <b>15</b> |
| <b>A</b> | <b>Appendix</b>          | <b>16</b> |
| A.1      | Code . . . . .           | 16        |

# 1 | Introduction

The link between viral load and virus severity was established with Covid-19 [1]. This motivates further the importance of potentially being able to predict viral load using machine learning techniques on spectral data. This would permit, not only virus classification [4] but virus severity predictions enabling quick and proportional treatment. In attempting to cluster the viral load levels, one could understand further the underlying patterns in spectral data that correlate to virus severity.

## 2 | Methodology

### 2.1 Regression

#### 2.1.1 Data Preprocessing

Data was read from an Excel spreadsheet with access provided by Queen's University, Belfast into a Pandas DataFrame in order to manipulate the data. As the data had previously been used for a study, the data was relatively clean and minimal preprocessing steps had to be taken. It was established that there were no missing values in the data set. The data set was of size  $N = 3801$  and  $p = 512$ . The data was split based on method being either '*dmem*' or '*pbs*'. *Dmem* had 1831 entries compared to *pbs* which had 1970 entries. A correlation matrix for both sub data sets was generated in order to assess co-linearity.

The design matrix  $X$  was standardized using Sklearn's StandardScalar so that each of the variables has mean 0 in order to prepare the data for Principal Components Analysis (PCA) [2, p.512].

Outliers were removed using Sklearn's supervised algorithm IsolationForest(). Isolation forest leverages properties of outliers being uncommon and unusual to isolate them [3]. It has linear time complexity, so is ideal for the task at hand in removing outliers with a model reducing the bias seen in other techniques such as z-scores, and interquartile range where thresholds must be defined *prior* to outlier detection.

#### 2.1.2 Feature Engineering

Spectrum data often contains wave lengths that are strongly correlated with each other. As such, particular attention to collinearity was taken. A correlation matrix of the features was generated for both data sets to demonstrate this common phenomenon when working with spectral data.

PCA was the first method to reduce dimensionality in order to eliminate collinearity. PCA projects the features to a lower dimensional space, while aiming to maintain the important information contained between the features. This can be useful for data visualisation, time complexity and coefficient estimation.

A closely related method for dimensionality reduction is Partial Least Squares. It is a supervised approach, as it identifies principal components relative to their covariance with the response variable. In the context of regression, PLS can be a more appropriate selection compared to PCA due to this guarantee that the features chosen will be the most effective at predicting the response variable[2, p.260,261], which is the viral load in this case, the target.

Another way to determine which features are of most importance to the model is using shrinkage methods. These methods shrink coefficient estimates which can significantly reduce their variance [2]. Lasso was included along with PCA and PLS as to serve the goal of inference. The Lasso algorithm shrinks the features to a subset of the original features, it can be directly interpreted as to which wave lengths are the most important when it comes to viral load estimation.

### 2.1.3 Model Selection

The ideal model in regression is usually determined by the model that minimizes Mean Squared Error (MSE). MSE is the average distance of error terms from the signal. It can be represented as a combination of 3 components [2, p.32], as shown below in fig 2.1.

$$E \left( y_0 - \hat{f}(x_0) \right)^2 = \text{Var}(\hat{f}(x_0)) + [\text{Bias}(\hat{f}(x_0))]^2 + \text{Var}(\epsilon).$$

**Figure 2.1:** Expected Test MSE

Here we can see MSE represented as a function of the bias and variance of the model and the variance of error terms. Bias is the price one pays for modelling real world data using constrained models that have a priori assumptions. Variance is the extent to which the model would change if trained on different data. To minimize MSE, variance and bias should be low to prevent under-fitting (high bias) or over-fitting (high variance).

The approach taken for this regression task was to compare the MSE between models that differ in bias and variance in order to select a model that balances the two and minimizes MSE. The following are the list of models to be evaluated in order of increasing variance, and decreasing bias.

*[Linear Regression, Polynomial Regression, Neural Network]*

Linear regression will be used as a baseline to evaluate the various dimension reduction techniques. The design matrix of the most effective technique will then be used for both polynomial regression and a basic neural network (NN).

## 2.2 Classification

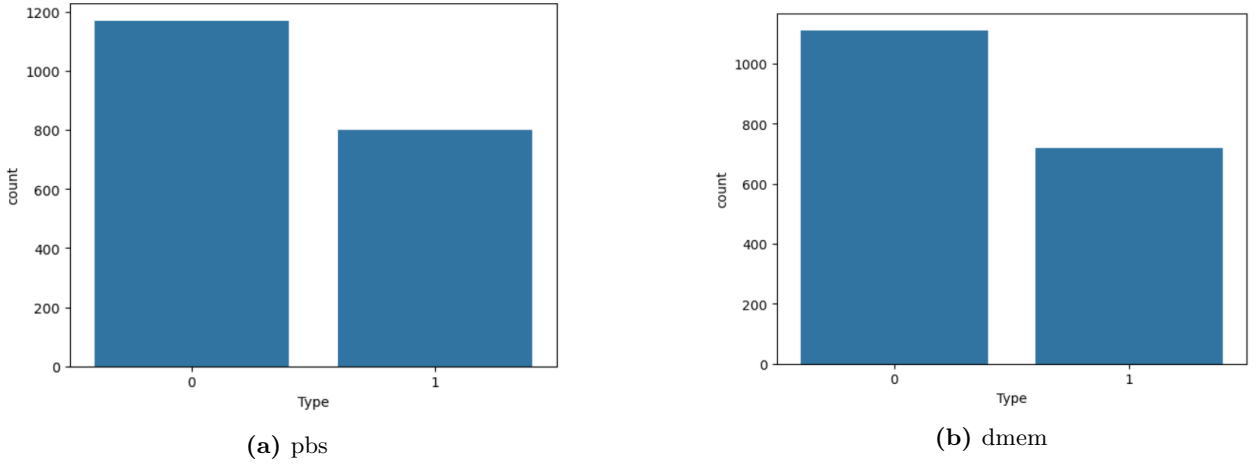
The aim of the classification method of this project was to reproduce the accuracy displayed previously in using spectral data to correctly identify and discriminate between virus types [4]. As such a PLS-DA was implemented. This model was then compared against other classic approaches to classification.

### 2.2.1 Preprocessing

The virus is either of type 'X' or of type 'Y'. This was binary encoded, where 'X' : 1 and 'Y': 0. This step was crucial as the data needed to be numeric for PLS. Then outliers were removed using Isolation Forest in for more robust classification modelling. The design matrix  $X$  remained the same from regression, being the standardized spectral data.

### 2.2.2 Feature Engineering

Cross validation was used in order to select the optimal number of components for PLS-DA to prevent multicollinearity present in spectral data. As shown below in figure 2.2, the two virus classes are relatively unbalanced.



**Figure 2.2:** A Consistent Moderate Virus Class Imbalance Across Datasets

PCA was used to generate the design matrix for the other classification approaches. To be cautious with parameter estimation, the *f1* metric was considered for model evaluation to combine both precision and recall, taking into consideration the minority class accuracy. However, *accuracy* was selected instead, to facilitate comparison to the original paper [4]. Classification accuracy is measured as follows in equation 2.1:

$$\text{Accuracy} = \frac{\text{TN} + \text{TP}}{\text{TN} + \text{TP} + \text{FN} + \text{FP}} \quad (2.1)$$

### 2.2.3 Model Selection

The PLS-DA model was selected on the basis of its previous predictive power with similar data [4], then other classical models were selected in order to serve as comparison. These models were Gaussian Naive Bayes (NB), Logistic Regression, K Nearest Neighbours (KNN), Linear Discriminate Analysis (LDA) and Quadratic Discriminate Analysis (QDA). None of these models. The performance of these models can also be revealing of properties underpinning the data. The model selection also compares the inclusion of the target variable in dimension reduction, as PLS-DA includes the target variable, whereas the other approaches used PCA.

## 2.3 Clustering

### 2.3.1 Preprocessing

Generally, as clustering is an unsupervised learning algorithm, there is no target variable. However, in this case, we do have a target variable. The motivation for clustering even when the target variable is present is that often this is not the case. So, it is useful to implement unsupervised learning algorithms and evaluate them, even when the target variable is present. In this paper, the target variable is the viral load level. So, the viral load  $y$ , is split from the spectral data  $X$ .

The absence of missing values is verified, then outliers were removed. As the target variable is not to be known until evaluation, an unsupervised approach to outlier removal is selected. Local Outlier Factor (LOF) was selected. This is a density based algorithm, that identifies outliers as being in a sparser region compared to their neighbours, where, unlike z-scores or interquartile range, normality is not assumed, the relationships between wavelengths are considered and finally no a priori thresholds

are defined. All this makes LOF an appropriate outlier detection algorithm for spectral data prior to clustering.

### 2.3.2 Feature Engineering

In clustering, there are some important decisions to be made. These decisions may appear small, but they can have big impacts on clustering [2, p.532,533]. The correct decision is not always obvious, as such, it can often be helpful to compare variations of clustering method specifications. The choice whether to standardize the data before clustering is one such decision. For model specification, previously with classification and regression tasks, prediction can be deemed as being more important than inference. However, with clustering, while it can be used for predictions, often it can also make up part of the exploratory stage of research. So while standardization, can improve precision and robustness of statistical learning algorithms, it can also sometimes hide something that would not otherwise be seen. Therefore it was decided to perform clustering on two forms of the spectral data. One that was standardized and dimensionally reduced using PCA, and the other being the original spectral data with just outliers removed.

### 2.3.3 Model Selection

3 different models were selected. The K-means algorithm was selected as a classic approach. Scatter plot visualisation is only applicable to the PCA reduced data. However, K-means does require the number of clusters to be predefined. Different values of K can be compared, but ultimately must be previously specified.

To avoid this, agglomerative (bottom-up) hierarchical clustering (AGC) is implemented [2, p.525], which also provides an interpretable dendrogram that remains functional in high dimensions. The measures of dissimilarity that link leaves and branches can vary. The choice of measure can significantly impact the AGC clustering. As such all measures available using the SKlearn library are used. This includes: complete, average, single, and Ward. Complete linkage computes the largest between cluster dissimilarity. Average linkage computes the average between cluster dissimilarity. Single computes the minimum between cluster similarity. So these dissimilarity measures all measure dissimilarity taking different perspectives on what best quantifies the distance between two clusters. Ward, on the other hand takes a different, variance based approach, which merges clusters that minimize between cluster variance. All of these dissimilarity measures are to be compared.

Finally the last model used in the clustering approach is DBSCAN. Where both K-means and AGC assign all data points a cluster, DBSCAN can detect noise based on density and assign outliers. This advantage may be mitigated by the fact that the data was preprocessed in identifying outliers using LOF, another density based algorithm.

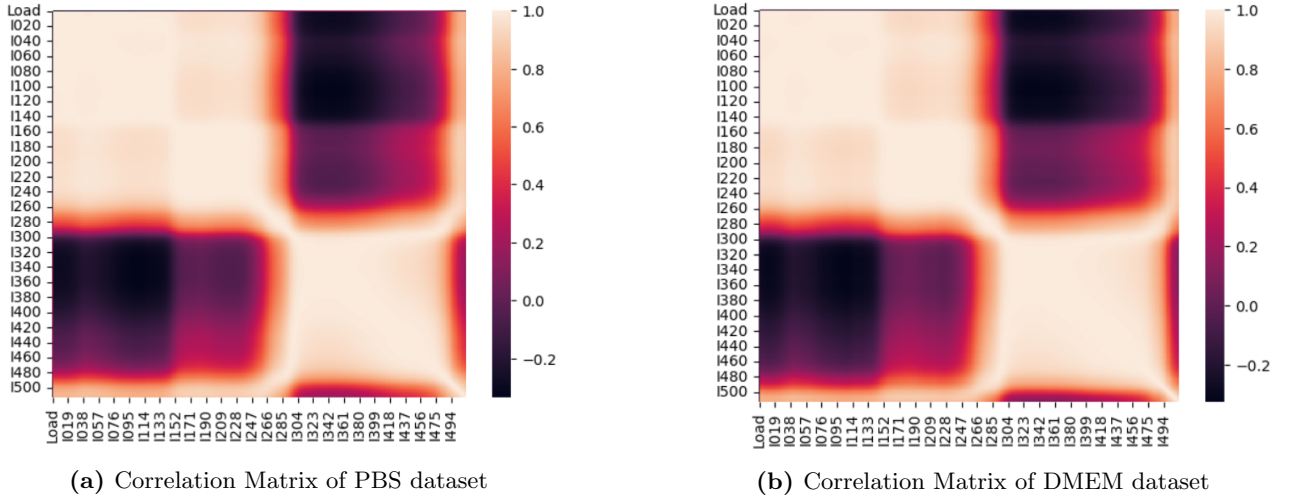


# 3 | Results

## 3.1 Regression

### 3.1.1 Multicollinearity

As displayed in figure 3.1 below, we can see lots of multicollinearity of the wave lengths. This motivated the preprocessing decision to reduce dimensionality before regression.



**Figure 3.1:** Correlation Matrices of Features and Target Variable Load

### 3.1.2 Dimensionality Reduction Comparison

The following table 3.1 is a comparison of linear model performance comparing between PCA, PLS and Lasso.

|      | PCA   | PLS  | Lasso |
|------|-------|------|-------|
| pbs  | 3.25  | 3.48 | 5.27  |
| dmem | 15.26 | 5.58 | 5.14  |

**Table 3.1:** A comparison of **MSE** scores between linear regression models using features obtained from different dimensionality reduction techniques.

Here we can see variance of MSE between data sets and between dimension reduction methods. To gain an understanding of the underlying structure of our feature space, the number of features remaining after applying the various methods was also recorded as shown below in table 3.2

|      | 1    | 2    | 3      | 4     | NN   |
|------|------|------|--------|-------|------|
| pbs  | 2.01 | 4.67 | 16.019 | 91.95 | 3.89 |
| dmem | 5.44 | 8.62 | 8.80   | 10.60 | 6.50 |

**Table 3.3:** MSE Comparison in Relation to Model Complexity

|      | PCA | PLS | Lasso |
|------|-----|-----|-------|
| pbs  | 3   | 5   | 0     |
| dmem | 3   | 1   | 0     |

**Table 3.2:** Number of Features After Applying Dimension Reduction Methods

### 3.1.3 Bias Variance Model Comparison on MSE

The PLS reduced data was selected as the design matrix based on both the MSE values and the number of features across the various methods. PCA was ruled out due to the significantly poor MSE performance with the *dmem* data set. Lasso was rejected as its optimal number of components was 0. This means that systematic shrinking of features actually improved the model every single time until no features remained. 5-fold cross validation was used for all models to ensure robustness of results and appropriateness of comparison between models.

## 3.2 Classification

### 3.2.1 Model Comparison

The following table 3.4 compares the classification accuracy across classification approaches and data sets. It is important to note that all the models except for the PLS-DA model had their design matrix generated using PCA.

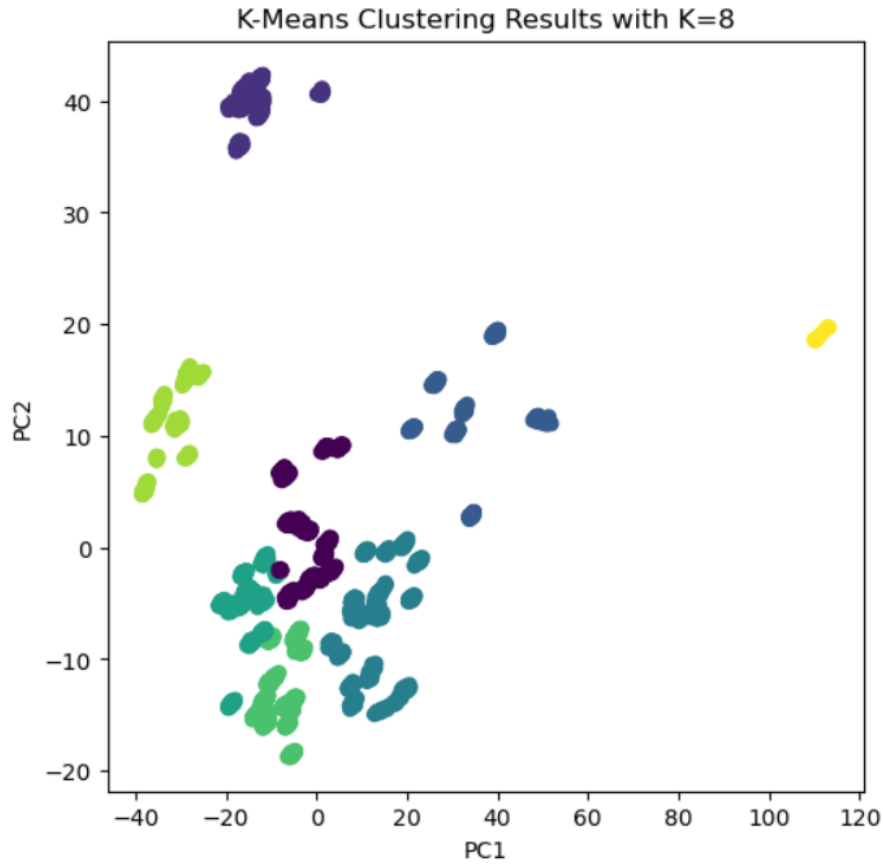
|      | PLS-DA | NB   | LR   | KNN  | LDA  | QDA  |
|------|--------|------|------|------|------|------|
| pbs  | 0.98   | 0.60 | 0.55 | 0.60 | 0.55 | 0.57 |
| dmem | 0.99   | 0.63 | 0.63 | 0.74 | 0.63 | 0.62 |

**Table 3.4:** Virus Type Classification **Accuracy** Comparison by Model

## 3.3 Clustering

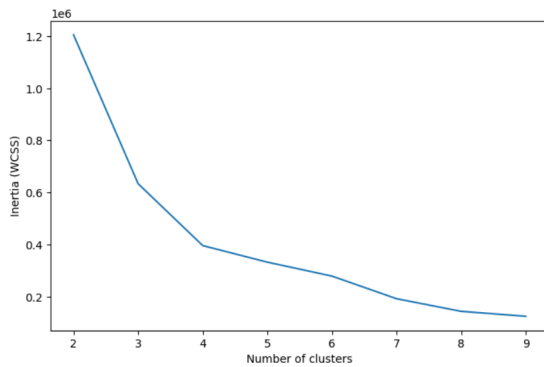
### 3.3.1 K-means

As the goal is to see whether the clustering algorithm will group the data by viral load levels. The visual representation of K-means where K is set to 8 is the most pertinent as there are 8 different load levels. The following figure 3.2.

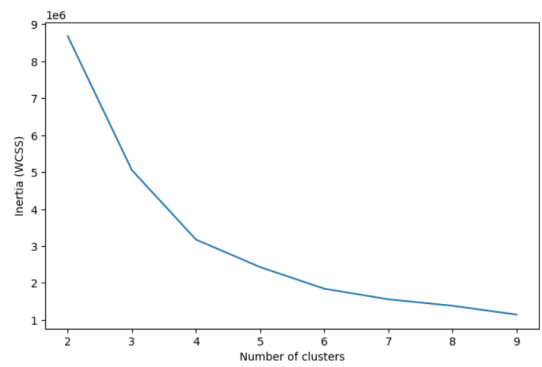


**Figure 3.2:** K-mans Clustering of PCA Reduced Data Using the First Two Prinicple Components where K is 8

Figure 3.3 represents the inertia scores relative to the number of clusters between 1 through to 10. This inertia is essentially the average distance between points in a cluster and that clusters centroid. This error term appears to be decreasing as K increases.



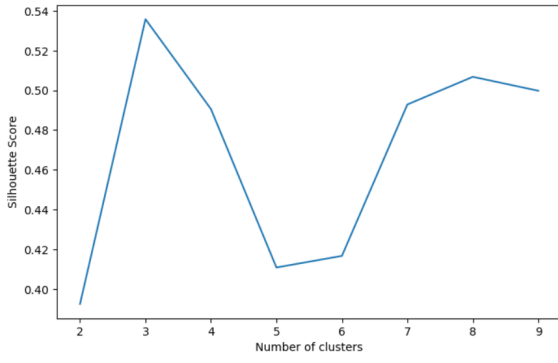
(a) Standardized and PCA Reduced Design Matrix



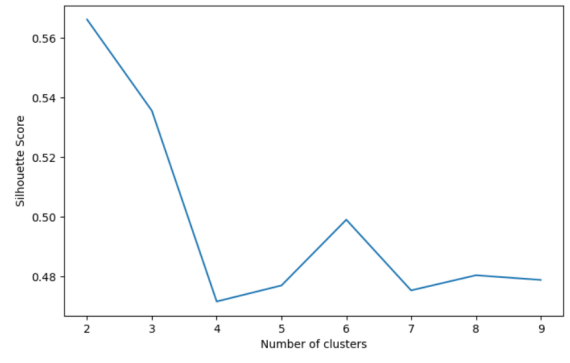
(b) Original Design Matrix

**Figure 3.3:** A Comparison of Inertia Between PCA Reduced  $X$  and Original  $X$  Relative to K

Figure 3.4 compares the silhouette scores between the two data sets. The silhouette score ranges from -1 to 1, and measures whether a data point is in the correct cluster. It does so by calculating the pairwise difference of a data point to its nearest centroid versus its assigned centroid. This is then normalized, where positive scores correspond to accurate cluster assignment.



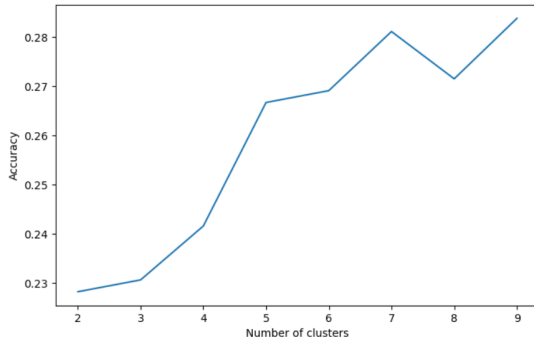
(a) Standardized and PCA Reduced Design Matrix



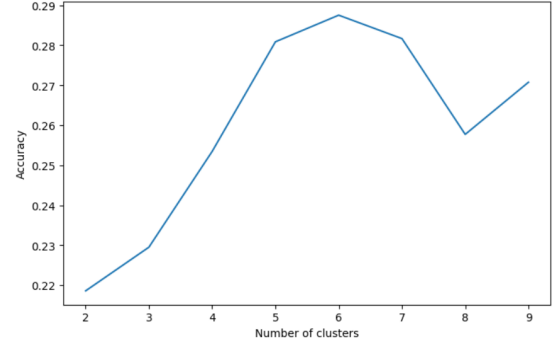
(b) Original Design Matrix

**Figure 3.4:** A Comparison of Silhouette Score Between PCA Reduced  $X$  and Original  $X$  Relative to  $K$

The accuracy measure corresponds to the proportion of data points that are assigned to the correct cluster. This evaluation metric differs to the previous two, as it requires true values. Another comparison between the PCA dimensionally reduced, and original data is shown below in figure 3.5. Here we can see an upwards trend of accuracy as  $K$  increases with accuracy slightly higher in the PCA reduced design matrix.



(a) Standardized and PCA Reduced Design Matrix

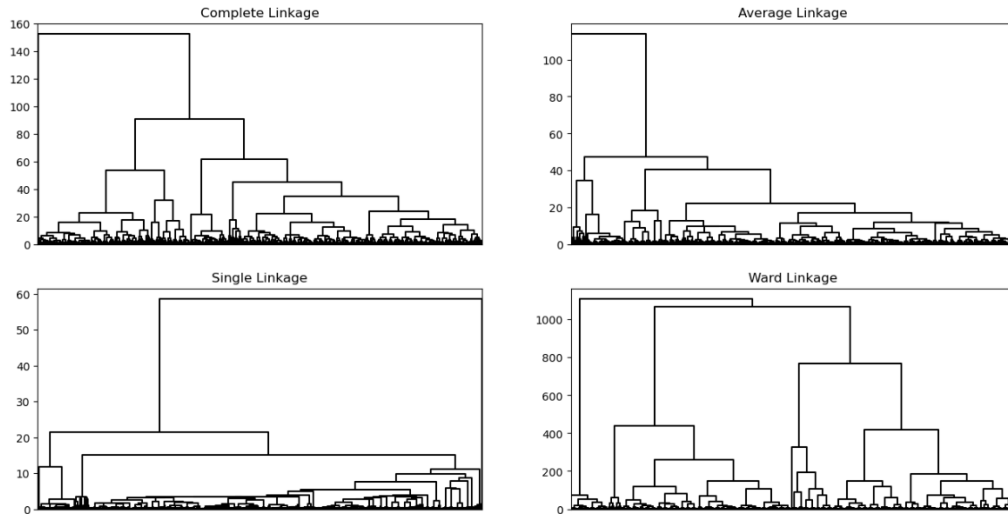


(b) Original Design Matrix

**Figure 3.5:** A Comparison of Accuracy Scores Between PCA Reduced  $X$  and Original  $X$  Relative to  $K$

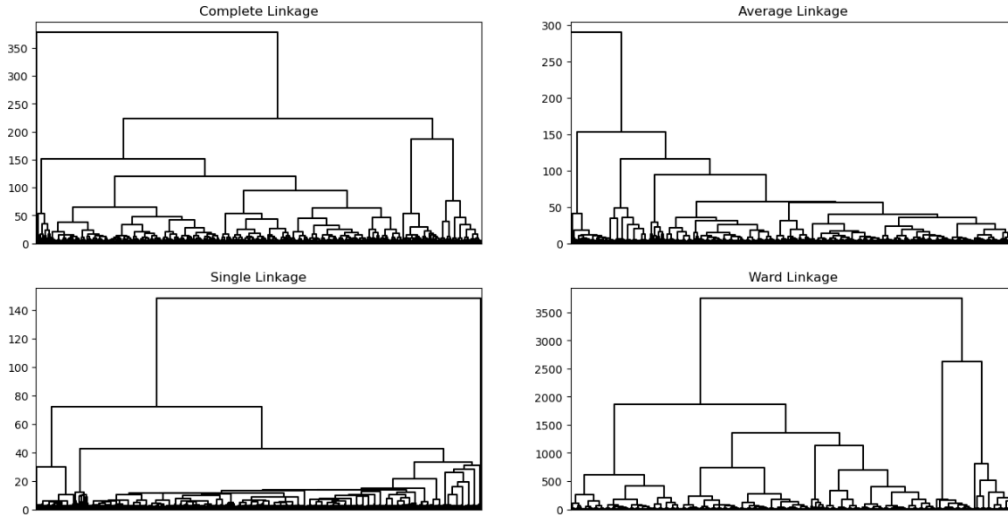
### 3.3.2 AGC

No cutting of the dendrogram was predefined. The following figure 3.6 compares the clustering relative to the different linkage methods for the PCA dimensionally reduced data.



**Figure 3.6:** A Comparison of Dendrograms by Dissimilarity Metric Using AGC on PCA Dimensionally Reduced Data

As done with K-means, we will also compare this to AGC with the original design matrix. This can be seen below in figure 3.7



**Figure 3.7:** A Comparison of Dendrograms by Dissimilarity Metric Using AGC Original Design Matrix

### 3.3.3 DBSCAN

The following table 3.5 displays the results of the DBSCAN clustering.

|                  | PCA | UNSTD |
|------------------|-----|-------|
| No. Clusters     | 30  | 124   |
| No. Noise Points | 0   | 176   |

**Table 3.5:** DBSCAN Clustering Results Displaying the Number of Cluster and Number of Noise Points Identified Relative to Whether  $X$  was PCA Dimensionally Reduced ( $PCA$ ) or Not ( $unstd$ )

## 4 | Discussion

### 4.1 Regression

#### 4.1.1 Multicollinearity

The overview of the correlation matrix resembles much that of a 2 by 2 identity matrix. This leads to the hypothesis that there are only 2 uncorrelated features. Upon dimension reduction, there was no consensus among methods as to how many independent features there were in the data. PCA was the only method with consensus between data sets, it was also the only unsupervised method. This could indicate a significant difference between the two data sets in how the spectral data relates to the viral load.

#### 4.1.2 Dimensionality Reduction Comparison

This importance of the target variable in reducing the data was echoed in dimension reduction, where PLS outperformed PCA.

#### 4.1.3 Bias Variance Model Comparison on MSE

Ultimately, the linear model and the NN did not differ significantly. With the flexibility of the NN and its similitude to the linear model, it is possible that the NN just learned a linear relationship between the features and the viral load. This is further motivated by the increase in MSE as complexity increases, as seen in the polynomial models with degrees of freedom between 2 and 4.

### 4.2 Classification

The significantly higher accuracy of the PLS-DA model compared to the other models that used an unsupervised approach to generate the design matrix demonstrates the power of adapting PLS Regression, which classically is used on continuous data, for a classification approach. It shows the importance of the relationship between the target variable and the features in dimension reduction. The KNN model outperformed the other models.

KNN tends to outperform these other algorithms when there are proportionally few features and many observations [2, p.164]. The PCA reduced data had 3 features and over 1000 observations. The NB model performs relatively well, this could be due to the effective meeting of its assumption that features are independent. This is achieved as the design matrix after PCA has orthogonal features. The seemingly small sample size might also favour NB.

### 4.3 Clustering

#### 4.3.1 K-means

The decrease in inertia as  $k$  increases is expected as smaller clusters naturally lead to less distance between points and centroids. This is consistent regardless of dimension reduction. No obvious elbow

is present, most noticeably, not at  $k = 8$ . So we can see that the data has not been sub-grouped by load level by the K-means algorithm. The magnitude of the distance between the two data sets is significantly different. The higher dimensional original data set has approximately 10 times the inertia. This is understandable as clusters become further apart in high dimensional space.

Accuracy and Silhouette scores on the other hand remain consistent regardless of dimension reduction. Silhouette scores seem to vary with no obvious pattern, this weak assigning of data points to their correct labels is demonstrated with accuracy scores of 25%. While 25% might not seem very precise, considering that there are 8 load levels, random guessing and balanced class would yield 12.5% accuracy. So, the clusters can actually be seen to double the probability of correctly assigning load level to a data point compared to random guessing.

### 4.3.2 AGC

Regardless of dissimilarity measure, all dendrograms except for Ward's linkage with the unstandardized data identify outliers in the data. This could indicate poor outlier detection with LOF. The number of neighbours used in LOF was not cross validated, this could have led to more robust outlier identification and removal. The Ward linkage appears to generate the dendrogram with the most well-defined clusters. Potentially, Ward's focus on variance and not outliers has made it less susceptible to the outliers in the data. Overall, there is no evidence of the measures generating a tree with 8 clear and distinct clusters representing the different load levels.

### 4.3.3 DBSCAN

DBSCAN identified 176 noise points of the unstandardized and reduced data and 0 with the PCA reduced data. This was even after identifying outlier using LOF, and subsequently removing them. This highlights potential over identification of points as noise. This could potentially be from the trouble of using distance metric in high dimensional space, or maybe the violation of the assumption of independence that DBSCAN carries.

## 5 | Conclusion

### 5.1 Regression

For the estimation of viral loads, *pbs* method seem to provide data that is more effective at predicting viral load than *dmem* method. The spectral data does not predict the viral load to a great accuracy using the models tested. This could motivate exploring other potential predictors or modelling strategies.

### 5.2 Classification

The application of machine learning techniques such as PLS-DA can be highly accurate classifying viruses using spectral data. This could increase diagnostic capacity of disease, which is especially important in times of outbreak.

### 5.3 Clustering

The clustering approach demonstrated how important the small decisions, such as whether to standardize data and the number of clusters can have a large impact on evaluation metrics and visualization. Its exploration of the patterns underpinning the spectral data and how they might relate to load levels revealed similarly to the regression approach, that there was no obvious inference to be made using spectral data and load levels.



# Bibliography

- [1] Fajnzylber, J., Regan, J., Coxen, K., et al. 2020, Nature Communications, 11, 5493, doi: [10.1038/s41467-020-19057-5](https://doi.org/10.1038/s41467-020-19057-5)
- [2] James, G., Witten, D., Hastie, T., Tibshirani, R., & Taylor, J. 2023, An Introduction to Statistical Learning: with Applications in Python (Springer)
- [3] Liu, F. T., Ting, K. M., & Zhou, Z.-H. 2008, 2008 Eighth IEEE International Conference on Data Mining, 413. <https://api.semanticscholar.org/CorpusID:6505449>
- [4] Song, W., Wang, H., Power, U. F., et al. 2023, IEEE Sensors Journal, 23, 9981, doi: [10.1109/JSEN.2022.3207222](https://doi.org/10.1109/JSEN.2022.3207222)

# A | Appendix

## A.1 Code

```
!/usr/bin/env python coding: utf-8
```

```
In[1]:
```

```
Standard Imports import numpy as np import pandas as pd import seaborn as sns import matplotlib.pyplot as plt
```

```
Data Preprocessing from sklearn.preprocessing import StandardScaler from sklearn.decomposition import PCA from mpl_toolkits.mplot3d import Axes3D from sklearn.cross_decomposition import PLSRegression from
```

```
Regression from sklearn.model_selection import cross_val_idate, cross_val_score, GridSearchCV, train_test_split, KFold
```

```
Deep Learning import tensorflow as tf from tensorflow import keras from tensorflow.keras import layers from keras.models import Sequential from scikeras.wrappers import KerasRegressor from keras.layers import Dense, Input
```

```
Classification from sklearn.linear_model import LogisticRegression from sklearn.model_selection import cross_val_idate, cross_val_score
```

```
CLustering from sklearn.cluster import KMeans, DBSCAN from sklearn.metrics import silhouette_score, accuracy_score
```

```
First Look at Dataset
```

```
In[3]:
```

```
virus_df = pd.read_excel('2022QUB_virus_data.xlsx', sheet_name = 'Data')
```

```
In[4]:
```

```
virus_df.info()
```

```
In[5]:
```

```
virus_df.head()
```

```
In[6]:
```

```
Function for graphing null values in data frame def graph_nulls(data): weather_df = dataplt.figure(figsize = (10,6))sns.heatmap(weather_df.isnull(), cmap = 'viridis', cbar = False)plt.title('Null Values in Each Column')plt.show()
```

```
1. Regression
```

```
In[8]:
```

```
Subset data set by matrix, being careful of preceding whitespace ! dmem_virus_df = virus_df.loc[virus_df['Matrix'] == 'pbs']
```

```
print(f'N — dmem: {dmem_virus_df.shape[0]}, pbs : {pbs_virus_df.shape[0]}')
```

```
In[9]:
```

```
sns.heatmap(pbs_virus_df.corr(numeric_only = True))
```

```
In[10]:
```

```
sns.heatmap(dmem_virus_df.corr(numeric_only = True))
```

```
Relationship Between Viral Load and Spectrum Values
```

```
Split Target Variable and Design Matrix
```

```
In[13]:
```

```
Split predictor and response variables y_dmem_reg = dmem_virus_df['Load']X_dmem = dmem_virus_df.drop(['Load'], axis = 1)
```

```
y_pbs_reg = pbs_virus_df['Load']X_pbs = pbs_virus_df.drop(['Load', 'SID', 'Type', 'Matrix'], axis = 1)
```

```
Target Variable Distribution
```

```
In[15]:
```

```

sns.countplot(data=dmem_virus_df, x = y_dmem_reg)
In[16]:
sns.countplot(data=pbs_virus_df, x = y_pbs_reg)
Standardize
In[18]:
Use standard scaler to standardize features, to prepare for feature engineering. scaler = Standard-
Scaler() X_dmem_std = scaler.fit_transform(X_dmem) X_pbs_std = scaler.fit_transform(X_pbs)
Remove Outliers
In[20]:
Use isolation forest to remove outliers. def remove_outliers_isoforest(X, y) : clf = IsolationForest(random_state=
42) outlier_pred = clf.fit_predict(X) RemoveOutliersIdentified X_clean_standardized = X[outlier_pred ==
1] y_clean = y[outlier_pred == 1]
return X_clean_standardized, y_clean, outlier_pred
X_dmem_std_clean, y_dmem_clean, outlier_pred_dmem = remove_outliers_isoforest(X_dmem_std, y_dmem_reg) X_pbs_std_clean,
remove_outliers_isoforest(X_pbs_std, y_pbs_reg)
In[21]:
dmem_lr = LinearRegression() lr_pca_dmem_results = cross_validate(lr, X_dmem_std_clean, y_dmem_clean, cv =
5, scoring = 'neg_mean_squared_error') print(lr_pca_dmem_results['test_score']) print(f'MSE : -np.mean(lr_pca_dmem_r
In[22]:
pbs_lr = LinearRegression() lr_pca_pbs_results = cross_validate(lr, X_pbs_std_clean, y_pbs_clean, cv =
5, scoring = 'neg_mean_squared_error') print(lr_pca_pbs_results['test_score']) print(f'MSE : -np.mean(lr_pca_pbs_r
In [ ]:
Dimension Reduction
PCA
In[25]:
def pca_scee(X, range) : pca = PCA(n_components = range) scores = pca.fit_transform(X)
fig, axes = plt.subplots(1, 2, figsize=(15, 6)) ticks = np.arange(pca.n_components)+1 ax = axes[0] ax.plot(ticks,
o') ax.set_xlabel('Principal Component'); ax.set_ylabel('Proportion of Variance Explained') ax.set_ylim([0, 1]) ax.set_
ax = axes[1] ax.plot(ticks, pca.explained_variance_ratio_.cumsum(), marker = 'o') ax.set_xlabel('Principal Comp
In[26]:
Scree plots for PCA in order to determine number of components, no titles as report in in Latex
and easier to add titles then.
In[27]:
pca_scee(X_dmem_std_clean, 11)
In[28]:
pca_scee(X_pbs_std_clean, 11)
In[29]:
Set PCA to 3 components pca = PCA(n_components = 3) X_pca_dmem = pca.fit_transform(X_dmem_std_clean) X_p
pca.fit_transform(X_pbs_std_clean)
Conovert back to df X_pca_dmem_df = pd.DataFrame(X_pca_dmem, columns = ['PC1', 'PC2', 'PC3']) X_pca_pbs_df =
pd.DataFrame(X_pca_pbs, columns = ['PC1', 'PC2', 'PC3'])
In[30]:
X_pca_dmem_and_y = X_pca_dmem_df.copy() X_pca_dmem_and_y['Load'] = y_dmem_clean sns.heatmap(X_pca_dmem_and_y.c
True), cmap = "viridis")
In[31]:
X_pca_pbs_and_y = X_pca_pbs_df.copy() X_pca_pbs_and_y['Load'] = y_pbs_clean sns.heatmap(X_pca_pbs_and_y.corr(numeric_o
True), cmap = "viridis")
In[32]:
Function to visualise the response variable in relation to 3 principle components def pca_visual_three_dim(X_pca_df,
target_and_PC = X_pca_df.copy() target_and_PC['Load'] = y
3D scatter plot fig = plt.figure(figsize=(10, 8)) ax = fig.add_subplot(111, projection = '3d')
scatter = ax.scatter( target_and_PC['PC1'], target_and_PC['PC2'], target_and_PC['PC3'], c = target_and_PC['Load'],
coolwarm', edgecolor = 'k')

```

```

Add color bar for 'Load' values cbar = plt.colorbar(scatter, ax=ax, pad=0.1) cbar.set_label('Load')
Label axes and set title ax.set_xlabel('PC1')ax.set_ylabel('PC2')ax.set_zlabel('PC3')plt.title("3D Scatter Plot of Load")
plt.show()
In[33]:
pca_visual_three_dim(X_pca_dmem_df, y_dmem_clean)
In[34]:
pca_visual_three_dim(X_pca_pbs_df, y_pbs_clean)
In[35]:
lr = LinearRegression() lr_pca_dmem_results = cross_validate(lr, X_pca_dmem_df, y_dmem_clean, cv =
5, scoring = 'r2')print(f'R^2 Linear Regression with PC Admem matrix : np.mean(lr_pca_dmem_results['test_score'])')
lr = LinearRegression() lr_pca_pbs_results = cross_validate(lr, X_pca_pbs_df, y_pbs_clean, cv = 5, scoring = '
r2')print(f'R^2 Linear Regression with PC Apbs matrix : np.mean(lr_pca_pbs_results['test_score'])')
PLS
In[37]:
def fit_PLS(X, y, max_components): Specify number of components to test component range = range(1, max_componen
1)mean_cv_scores = []
for n_components in component_range : pls = PLSRegression(n_components = n_components)cv_scores =
cross_validate(pls, X, y, cv = 5, scoring = 'neg_mean_squared_error')mean_cv_score = -cv_scores.mean()mean_cv_score
n_components, MSE : mean_cv_score")
optimal_components = component_range[np.argmin(mean_cv_scores)]print(f"Optimal number of components :
optimal_components")print(f"Min cross - validated MSE : min(mean_cv_scores)")
optimal_pls = PLSRegression(n_components = optimal_components)X_optimal = optimal_pls.fit_transform(X,
return X_optimal, optimal_pls
y_dmem_std_clean = scaler.fit_transform(y_dmem_clean.to_frame())y_pbs_std_clean = scaler.fit_transform(y_pbs_clean.to_frame())
In[38]:
X_dmem_PLS, dmem_PLS_model = fit_PLS(X_dmem_std_clean, y_dmem_clean, 10)X_pbs_PLS, pbs_PLS_model =
fit_PLS(X_pbs_std_clean, y_pbs_clean, 10)
Shrinkage Methods
Lasso
In[41]:
Use Lasso instead def lasso(X, y): lasso = Lasso() parameters='alpha': [0.001, 0.01, 0.1, 0.5,
1.0, 100.0, 1000.0] lasso_ regressor = GridSearchCV(lasso, parameters, scoring = 'r2')lasso_model =
lasso_ regressor.fit(X, y)best_lasso_model = lasso_model.best_estimator_.print(f'Optimal Lasso : best_lasso_model')nonzero_coef =
np.sum(best_lasso_model.coef_ != 0)print(f'Number of predictors with non-zero coefficients : nonzero_predictors')
cross_validate(lasso_model, X, y, cv = 5, scoring = 'neg_mean_squared_error')print(-cv_results['test_score'])cv_err =
np.mean(-cv_results['test_score'])print(f'MSE : cv_err')
In[42]:
lasso(X_dmem_std_clean, y_dmem_clean)
In[43]:
lasso(X_pbs_std_clean, y_pbs_clean)
Move Beyond Linearity
Polynomials
In[46]:
def poly(X, y, degree): Set up a range of polynomial degrees to test degrees = range(1, degree+1)
for degree in degrees: print(degree) Generate polynomial features poly = PolynomialFeatures(degree)
X_poly = poly.fit_transform(X)
model = LinearRegression()
MSE scores cv_scores = cross_validate(model, X_poly, y, cv = 5, scoring = 'neg_mean_squared_error')
print(f'Degree degree - Average MSE score: -np.mean(cv_scores) : .4f')print(f'MSE scores for each fold :
cv_scores')
In[47]:
poly(X_dmem_PLS, y_dmem_clean, 4)
In[48]:

```

```

poly(XpbsPLS, ypbsclean, 4)
Neural Network
In[50]:
Must set default paramaters for Grid Search to work def createmodel(optimizer = 'adam', inputdim =
8, activation = 'relu', numhiddenilayers = 1) : InitialiseNeuralNetorknn = Sequential()tf.keras.Sequential, w
... - > output.Inputnn.add(Input(shape = (inputdim,)))Hiddenlayersforiinrange(numhiddenilayers) :
nn.add(Dense(64, activation = activation))Outputlayer, loadestimate.nn.add(Dense(1))
.compile just sets configuration nn.compile(optimizer=optimizer, loss='meansquarederror', metrics =
[keras.metrics.RootMeanSquaredError()])
return nn
Function to perform regression given design matrix and target variable def NNregression(X, y) :
UseKerasRegressorwrapperinordertousesklearntoolswithKeras.model = KerasRegressor(model =
createmodel, inputdim = X.shape[1], numhiddenilayers = 1, epochs = 1, verbose = 0)
paramgrid = 'epochs' : [1, 2, 5, 10, 20], 'numhiddenilayers' : [1, 2]
grid = GridSearchCV(estimator=model, paramgrid = paramgrid, scoring = 'negmeansquarederror', cv =
5)gridresult = grid.fit(X, y)
return gridresult.bestparams, - gridresult.bestscore
In[51]:
dmemnnparamspls, dmemnnMSEpls = NNregression(XdmemPLS, ydmemclean)pbsnnparamspls, pbsnnMSEpls =
NNregression(XpbsPLS, ypbsclean)
In[52]:
print('PLS') print(f'dmem — optimal parameters: dmemnnparamsplsMSE : dmemnnMSEpls')print(f'pbs|op
pbsnnparamsplsMSE : pbsnnMSEpls')
2. Classification
Preprocessing
In[55]:
new target variable target variable yclfdmem = dmemvirusdf['Type']yclfpbs = pbsvirusdf['Type']makeybinar
(1, 0), handlespacesybindmem = yclfdmem.str.strip().map('X' : 1, 'Y' : 0)ybinpbs = yclfpbs.str.strip().map('X' :
Outlier Removal
In[57]:
def removeoutliers;soforest(X, y) : clf = IsolationForest(randomstate = 42)outlierpred =
clf.fitpredict(X)RemoveOutliersidentifieXcleansstandardized = X[outlierpred == 1]yclean = y[outlierpred ==
1]
return Xcleansstandardized, ycclean, outlierpred
Already standardized with standard scaler Xdmemstdcleanclf, ydmemcleanclf, outlierpreddmem =
removeoutliers;soforest(Xdmemstd, ybindmem)Xpbsstdcleanclf, ypbscleanclf, outlierpredpbs = removeoutliers;
In[58]:
sns.countplot(x=ybindmem)
In[59]:
sns.countplot(x=ybinpbs)
PCA
Scree Plot
In[62]:
pcascree(Xdmemstdcleanclf, 11)
In[63]:
pcascree(Xpbsstdcleanclf, 11)
Correlation with Target
In[65]:
Set PCA to 3 components pca = PCA(ncomponents = 3)Xpcadmemclf = pca.fittransform(Xdmemstdcleanclf)
pca.fittransform(Xpbsstdcleanclf)
In[66]:
Conovert back to df Xpcadmemclfdf = pd.DataFrame(Xpcadmemclf, columns = ['PC1', 'PC2', 'PC3'])Xpcapbsclf =
pd.DataFrame(Xpcapbsclf, columns = ['PC1', 'PC2', 'PC3'])

```

```

In[67]:
X_pca_dmem_and_y_clf = X_pca_dmem_clf.copy()X_pca_dmem_and_y_clf['Load'] = y_dmem_lean_clf sns.heatmap(X_pca_dmem_and_y_clf, cmap = "viridis")
In[68]:
X_pca_pbs_and_y_clf = X_pca_pbs_clf.copy()X_pca_pbs_and_y_clf['Load'] = y_pbs_lean_clf sns.heatmap(X_pca_pbs_and_y_clf, cmap = "viridis")
In[69]:
pca_visualize_3d(X_pca_pbs_clf, y_dmem_lean_clf)
In[70]:
pca_visualize_3d(X_pca_pbs_clf, y_pbs_lean_clf)
PLS
In[72]:
from sklearn.cross_decomposition import PLSRegression from sklearn.linear_model import LogisticRegression
def fit_PLS_for_classification(X, y, max_components) : component_range = range(1, max_components+
1) mean_cv_score_relative_to_number_of_components = []
for n_components in component_range : Fit_PLS = PLSRegression(n_components = n_components) X_transformed = pls.fit_transform(X, y)[0]
- DA_clf = LogisticRegression(max_iter = 1000) cv_scores = cross_val_score(clf, X_transformed, y, cv =
5, scoring = 'f1') mean_cv_score = np.mean(cv_scores) mean_cv_scores.append(mean_cv_score)
optimal number of components based on cross-validated score optimal_components = component_range[np.argmax
np.max(mean_cv_scores)] transformation_using_the_optimal_number_of_components optimal_pls = PLSRegression(n_components = optimal_components) X_optimal = optimal_pls.fit_transform(X, y)[0]
return X_optimal, optimal_pls, optimal_components, best_acc
In[73]:
X_pbs_PLS_clf, pbs_PLS_model_clf, num_pbs, score_pbs = fit_PLS_for_classification(X_pbs_std_lean_clf, y_pbs_lean_clf, 20, num_pbs, accuracy : score_pbs')
X_dmem_PLS_clf, dmem_PLS_model_clf, num_dmem, score_dmem = fit_PLS_for_classification(X_dmem_std_lean_clf, y_dmem_lean_clf, 20, num_dmem, accuracy : score_dmem')
Normality of Features
In[75]:
def graph_data_hist(data) : Plots histograms for numerical features plt.rc('font', size = 14) plt.rc('axes', labelsiz
14, titlesize = 14) plt.rc('legend', fontsize = 14) plt.rc('xtick', labelsiz = 10) plt.rc('ytick', labelsiz =
10)
data.hist(bins=50, figsize=(12, 8)) plt.show()
In[76]:
graph_data_hist(pd.DataFrame(X_pca_dmem))
In[77]:
graph_data_hist(pd.DataFrame(X_pca_pbs))
Ratio of Features to Observations
In[79]:
Function for visualising confusion matrix for evaluating classification models. def CM_visualise(X, y, model) :
ConfusionMatrixGenerator cross-validated predictions y_pred_log = cross_val_predict(model, X, y, cv =
10) Create and plot the confusion matrix conf_matrix = confusion_matrix(y, y_pred_log) disp = ConfusionMatrixDis
conf_matrix) disp.plot(cmap = 'Blues') plt.title(f"Confusion Matrix for model") plt.show()
Naive Bayes
In[81]:
gnb = GaussianNB()
In[82]:
cv = np.mean(cross_val_score(gnb, X_pca_pbs, y_pbs_lean_clf, cv = 5, scoring = 'accuracy')) print(cv)
In[83]:
CM_visualise(X_pca_pbs, y_pbs_lean_clf, gnb)
In[84]:
cv = cross_val_score(gnb, X_pca_dmem, y_dmem_lean_clf, cv = 5, scoring = 'accuracy') print(np.mean(cv))

```

```

In[85]:
CM_visualise(X_pca_dmem, y_dmem_leanclf, gnb)
Logistic Regression
In[87]:
log_model = LogisticRegression(max_iter = 1000)
In[88]:
cv = cross_val_score(log_model, X_pca_pbs, y_pbs_leanclf, cv = 5, scoring = 'accuracy')print(np.mean(cv))
In[89]:
CM_visualise(X_pca_pbs, y_pbs_leanclf, log_model)
In[90]:
cv = cross_val_score(log_model, X_pca_dmem, y_dmem_leanclf, cv = 5, scoring = 'accuracy')print(np.mean(cv))
In[91]:
CM_visualise(X_pca_dmem, y_dmem_leanclf, log_model)
KNN Classifier
In[93]:
knnclf = KNeighborsClassifier(n_neighbors = 5)
In[94]:
cv = cross_val_score(knnclf, X_pca_pbs, y_pbs_leanclf, cv = 5, scoring = 'accuracy')print(np.mean(cv))
In[95]:
CM_visualise(X_pca_pbs, y_pbs_leanclf, knnclf)
In[96]:
cv = cross_val_score(knnclf, X_pca_dmem, y_dmem_leanclf, cv = 5, scoring = 'accuracy')print(np.mean(cv))
In[97]:
CM_visualise(X_pca_dmem, y_dmem_leanclf, knnclf)
LDA
In[99]:
lda = LDA()
In[100]:
cv = cross_val_score(lda, X_pca_pbs, y_pbs_leanclf, cv = 5, scoring = 'accuracy')print(np.mean(cv))
In[101]:
CM_visualise(X_pca_pbs, y_pbs_leanclf, lda)
In[102]:
cv = cross_val_score(lda, X_pca_dmem, y_dmem_leanclf, cv = 5, scoring = 'accuracy')print(np.mean(cv))
In[103]:
CM_visualise(X_pca_dmem, y_dmem_leanclf, lda)
QDA
In[105]:
qda = QDA()
In[106]:
cv = cross_val_score(qda, X_pca_pbs, y_pbs_leanclf, cv = 5, scoring = 'accuracy')print(np.mean(cv))
In[107]:
CM_visualise(X_pca_pbs, y_pbs_leanclf, qda)
In[108]:
cv = cross_val_score(qda, X_pca_dmem, y_dmem_leanclf, cv = 5, scoring = 'accuracy')print(np.mean(cv))
In[109]:
CM_visualise(X_pca_dmem, y_dmem_leanclf, qda)
3. Clustering
Preprocessing
In[112]:
Use full data set cluster_data = virus_df
Split X and y
In[114]:

```

```

Spectrum data as design matrix  $X_{cluster} = cluster\_data[cluster\_data.columns[4 : 517]]$  Forevaluatingclusterin.
cluster_data['Load']
Standardization
In[116]:
Standardize, so as to compare performance with and without. scaler = StandardScaler()  $X_{cluster\_std} =$ 
scaler.fit_transform( $X_{cluster}$ )
Outlier Removal (Unsupervised)
In[118]:
Unsupervised outlier removal. def outlier_removal_of( $X, y$ ) :  $LOF = LocalOutlierFactor(n\_neighbors =$ 
2)outlier_pred =  $LOF.fit\_predict(X)$ 
 $X_{clean} = X[outlier\_pred == 1]$   $y_{clean} = y[outlier\_pred == 1]$ 
return  $X_{clean}, y_{clean}$ 
In[119]:
 $X_{cluster\_std\_clean}, y_{cluster\_std\_clean} = outlier\_removal\_of(X_{cluster\_std}, y_{cluster\_load})$   $X_{cluster\_clean}, y_{cluster\_clean} =$ 
outlier_removal_of( $X_{cluster}, y_{cluster\_load}$ )
PCA
In[121]:
Use PCA as it is an unsupervised dimension reduction (colinearity removal) method.  $pca\_score(X_{cluster\_std\_clean})$ 
In[122]:
 $pca = PCA(n\_components = 3)$   $X_{cluster\_pca} = pca.fit\_transform(X_{cluster\_std\_clean})$ 
Kmeans
In[124]:
Define range for k values (number of centroids) def plot_k_means( $X, y, max\_k = 10$ ) :  $K = range(2, max\_k)$ 
fits = [] inertia = [] silhouette_scores = [] accuracies = [] Iterate over different k values for in  $K$  :
 $kmeans = KMeans(n\_clusters = k, init = 'k-means++', random\_state = 42)$   $kmeans.fit(X)$  fits.append( $kmeans.inertia_$ )
inertia.append( $kmeans.inertia_$ )
distance of  $x_i$  to centroid => how tightly packed (defined) clusters are. inertia.append( $kmeans.inertia_$ )
is point in the right cluster - 1 - > 1 score = silhouette_score( $X, kmeans.labels_$ ) silhouette_scores.append(score)
labels =  $kmeans.labels_$  label_mapping = np.zeros_like(labels)
for i in range(k): mask = (labels == i) use mode() to find most frequent label. label_mapping[mask] =
mode( $y[mask]$ )[0]
accuracy = accuracy_score( $y, label\_mapping$ ) accuracies.append(accuracy)
if  $X.shape[1] > 4$ : Scatter plot for the clusters based on PC1 and PC2 plt.figure(figsize=(6, 6))
plt.scatter( $X[:, 0], X[:, 1], c=kmeans.labels_$ ) plt.xlabel('PC1') plt.ylabel('PC2') plt.title(f'K - Means Clustering Results for k')
plt.show()
return inertia, silhouette_scores, accuracies
In[125]:
inertia_load, silhouette_load, accuracies_load = plot_k_means( $X_{cluster\_pca}, y_{cluster\_std\_clean}, 10$ )
In[126]:
Plot Elbow Method graph plt.figure(figsize=(8, 5)) plt.plot(range(2, len(inertia_load)+2), inertia_load) plt.title('Elbow Method')
In[127]:
Plot Silhouette Analysis graph plt.figure(figsize=(8, 5)) plt.plot(range(2, len(inertia_load)+2), silhouette_load) plt.title('Silhouette Analysis')
In[128]:
Accuracy by K graph plt.figure(figsize=(8, 5)) plt.plot(range(2, len(accuracies_load)+2), accuracies_load) plt.title('Accuracy by K')
In[129]:
inertia_load_u_nstd, silhouette_load_u_nstd, accuracies_load_u_nstd = plot_k_means( $X_{cluster\_clean}, y_{cluster\_clean}, 10$ )
In[130]:
Plot Elbow Method graph plt.figure(figsize=(8, 5)) plt.plot(range(2, len(inertia_load_u_nstd)+2), inertia_load_u_nstd)
In[131]:
Plot Silhouette Analysis graph plt.figure(figsize=(8, 5)) plt.plot(range(2, len(inertia_load_u_nstd) +
2), silhouette_load_u_nstd) plt.title('Silhouette Analysis') plt.xlabel('Number of clusters') plt.ylabel('Silhouette Score')
In[132]:
Accuracy by K graph plt.figure(figsize=(8, 5)) plt.plot(range(2, len(accuracies_load_u_nstd)+2), accuracies_load_u_nstd)
Hierarchical

```



```

In[134]:
Complete Linkage hc_complete = HClust(distance_threshold = 0, n_clusters = None, linkage = 'complete')
AverageLinkagehc_ave = HClust(distance_threshold = 0, n_clusters = None, linkage = 'average')
Singlelinkagehc_sing = HClust(distance_threshold = 0, n_clusters = None, linkage = 'single')

Centroid linkage hc_ward = HClust(distance_threshold = 0, n_clusters = None, linkage = 'ward')
PCA
In[136]:
No preconceived chopping cargs = 'color_threshold' : -np.inf, 'above_threshold_color' : 'black'
In[137]:
def plot_bylinkage(X) : BylinkageTypecomplete, average, single, centroidfig, ax = plt.subplots(2, 2, figsize =
(16, 8))
complete hc_complete.fit(X)linkage_complete = compute_linkage(hc_complete)dendrogram(linkage_complete, ax
ax[0, 0], **cargs)ax[0, 0].set_title("CompleteLinkage")ax[0, 0].set_xticks([])averagehc_ave.fit(X)linkage_ave =
compute_linkage(hc_ave)dendrogram(linkage_ave, ax = ax[0, 1], **cargs)ax[0, 1].set_title("AverageLinkage")ax[0, 1]
single hc_sing.fit(X)linkage_sing = compute_linkage(hc_sing)dendrogram(linkage_sing, ax = ax[1, 0], **
cargs)ax[1, 0].set_title('SingleLinkage')ax[1, 0].set_xticks([])
centroid hc_ward.fit(X)linkage_ward = compute_linkage(hc_ward)dendrogram(linkage_ward, ax =
ax[1, 1], ** cargs)ax[1, 1].set_title('WardLinkage')ax[1, 1].set_xticks([])
plt.show()
In[138]:
plot_bylinkage(X_cluster_pca)
In[139]:
plot_bylinkage(X_cluster_clean)
DBSCAN
In[141]:
dbscan = DBSCAN(eps=3, min_samples = 2).fit(X_cluster_pca)n_clusters = len(set(dbscan.labels_)) -
(1 if -1 in dbscan.labels_ else 0)n_noise = list(dbscan.labels_).count(-1)
print('pca') print(f'Number of clusters: n_clusters")print(f"Number of noise points : n_noise")
In[247]:
dbscan = DBSCAN(eps=3, min_samples = 2).fit(X_cluster_clean)n_clusters = len(set(dbscan.labels_)) -
(1 if -1 in dbscan.labels_ else 0)n_noise = list(dbscan.labels_).count(-1)
print('unstd') print(f'Number of clusters: n_clusters")print(f"Number of noise points : n_noise")

```