

Machine Learning

02524541

April 2024

Question 1

Part 1

(1)

Figure 1 shows the mean absorption over all the samples at each wave length. It can be seen that as the wavelength increases, the mean level of absorption increases from roughly -1 to 1.6. This shows that as the wavelength increases from 800 to 999, the level of absorption increases. Between the wavelengths of 900 and 925, the mean absorption increases and then decreases dramatically.

Feature	Correlation
Temperature	0.164507
Water	0.005536
Ethanol	0.005536
2-Propanol	0.005536

Table 1: Correlations between Various Features and Absorption

In the table above, it is clear that each component has a positive correlation with the mean absorption over the wavelengths. Although the data is not linear, the correlation coefficients can show a general increase in the mean absorption over each feature.

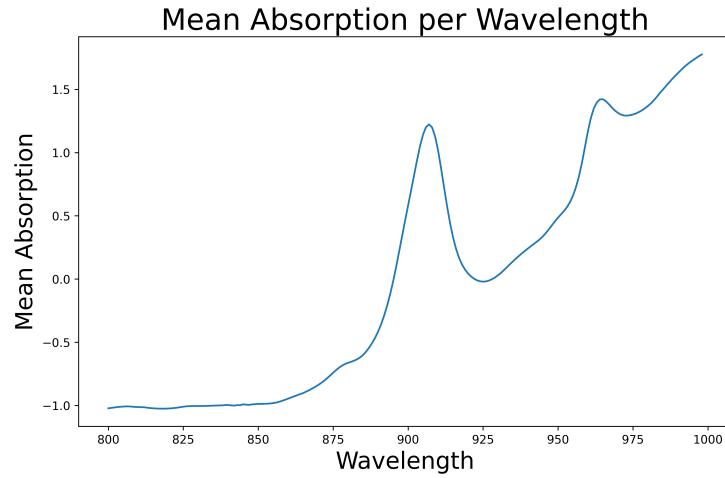


Figure 1: Mean Wavelength Absorption

The temperatures take values between 30 to 70. The concentrations of water, Ethanol and 2-Propanol are take values between 0 and 1.

(2)

Hierarchical clustering was performed,

The 2 metrics used to decide how the hierarchical clustering was performed were "euclidean" and "cityblock". Each of these metrics has different ways to conclude the distance between a new value.

The 3 methods used to decide how the hierarchical clustering is done were "complete", "single" and "average". Each of these methods has different ways to conclude the distance between a new value.

Below is the graphs of each hierarchical cluster using different metrics and methods, this was performed using the `linkage` and `dendrogram` function as part of the `scipy.cluster.hierarchy` library.

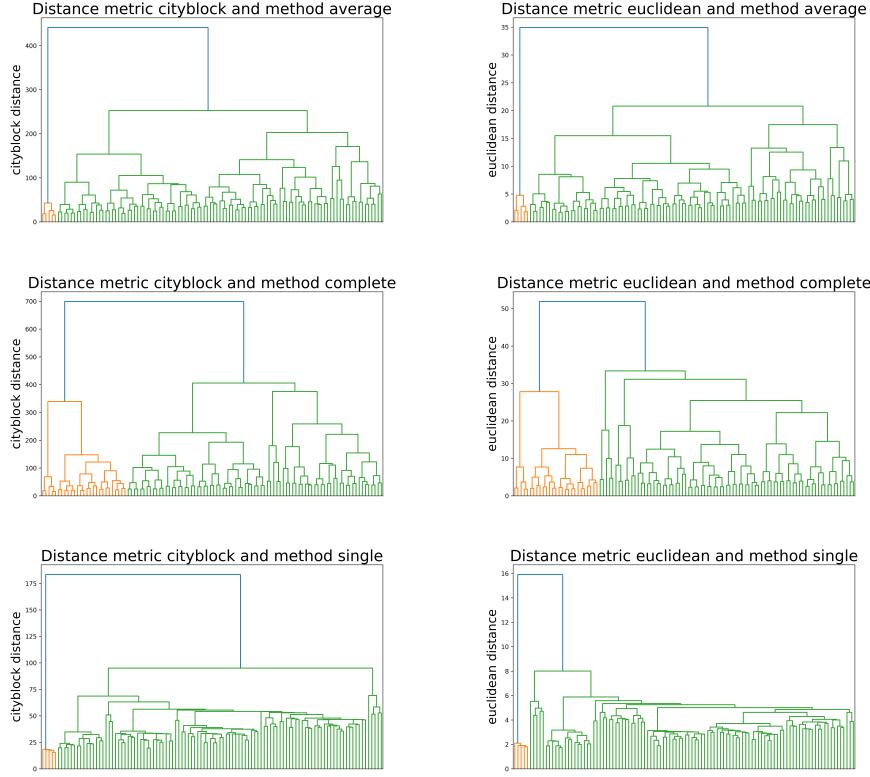


Figure 2: Hierarchical Clustering for Different Metric and Methods

As can be seen in the graphs above, different distance metrics produce different types of data splits. For example, the single and average distance metrics show less clearly distinct clusters or groups compared to the euclidean distance metrics.

Using the `fcluster` function, as part of the `scipy.cluster.hierarchy` library, the data was split into 4 distinct groups using 2 different combinations of methods and metrics.

Group	Water	Ethanol	2-Propanol	Temperature
1	0.000	1.000	0.000	49.966
2	0.231641	0.511095	0.257264	53.063913
3	0.494467	0.186763	0.318770	46.5115
4	0.324083	0.000	0.675917	48.242353

Table 2: Mean Values of Features for each group using "average" linkage method and "cityblock" distance metric

Group	Water	Ethanol	2-Propanol	Temperature
1	0.154217	0.685316	0.160468	54.1
2	0.0	0.0	1.0	50.018
3	0.738271	0.000000	0.261729	35.714286
4	0.385679	0.253295	0.361026	49.428406

Table 3: Mean Values of Features for each Group using "complete" linkage method and "euclidean" distance metric

When separating the data into different clusters, similar components and temperatures tend to cluster together. This is why there are clusters with unnaturally high or unnaturally low mean values for the components and the temperature values.

(3)

The principal component analysis of the spectra was performed on the space using the PCA function as part of the `sklearn` library. Plotted below is the cumulative variance for each number of principal components.

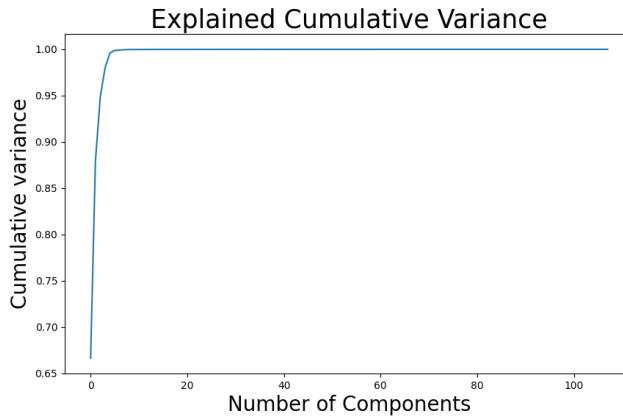


Figure 3: Cumulative Variance for each Number of Principal Components

The smallest amount of principal components to account for 99.99% of the variance is 14.

Dimension reduction was performed by taking the standardised spectra data points that accounted for 99.99% of the variance. The mean of the overall data set and the reduction data sets were charted below to compare the data sets. As can be seen the reduced data set has a smaller mean for temperature and 2-propanol, and a larger mean for ethanol.

Features	14 Principle Features Mean	All Features Mean
Temperature	42.873571	49.734722
Water	0.344916	0.332811
Ethanol	0.392625	0.333156
2-Propanol	0.262459	0.334033

Table 4: Mean Values for 14 Principle Component Features and Overall Data

Hierarchical clustering was performed on the reduced data set, using the same metrics and methods used in the previous part of this question.

The clusters are plotted below,

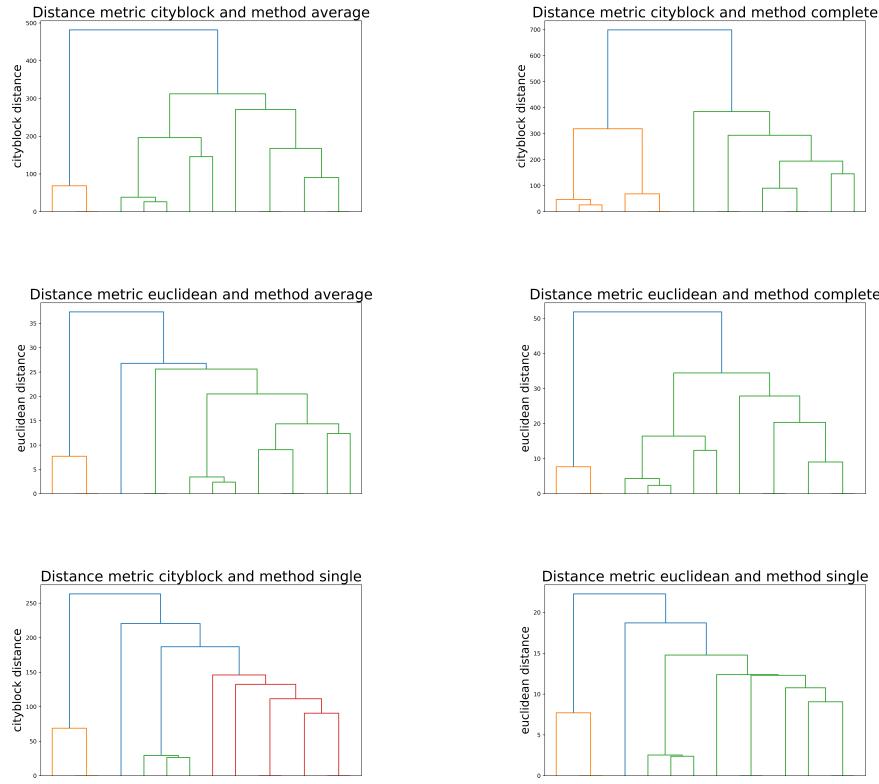


Figure 4: Hierarchical Clustering for Different Metric and Methods with Reduced Data Set

Group	Water	Ethanol	2-Propanol	Temperature
1	0.0	1.0	0.0	43.346667
2	0.232409	0.499349	0.268242	42.01
3	0.733357	0.0	0.266643	46.028
4	0.0	0.0	1.0	30.0

Table 5: Mean Values of Features for each group using "complete" linkage method and "euclidean" distance metric, for the Reduced Data

Group	Water	Ethanol	2-Propanol	Temperature
1	0.0	1.0	0.0	43.346667
2	0.232409	0.499349	0.268242	42.01
3	0.733357	0.0	0.266643	46.028
4	0.0	0.0	1.0	30.0

Table 6: Mean Values of Features for each group using "complete" linkage method and "euclidean" distance metric, for the Reduced Data

Just like the groups from the overall data set, different groups have different means for the component and the temperature values. This is seen in the fact that, the single and average distance metrics show less clearly distinct clusters/groups compared to the euclidean distance metrics.

Part 2

(1)

- A sample of 40 chemical samples was randomly selected from the 108 original chemical samples using the `.sample()` function from a NumPy data set,
- Then random noise was added using the parameters specified in the coursework to each element
- 80% of the 40×199 measurements were replaced by `Nan` values

This was all done with the author's student number (2524541) set as the random seed value.

(2)

From the 40 components, the component with the smallest amount of missing wavelength data points was chosen. The reason for this is because with more data points the accuracy of the error is likely to decreases and a more nuanced understanding of the data can be found.

A combination of the following Gaussian processes was tested using k-fold validation with 3 folds and with the `model_selection.KFold` function as part of the `sklearn` package.

The Exponential Sine Squared (ESS) kernel is denoted as $k(x, x')_{\text{ESS}}$, the Radical basis function (RBF) is denoted as $k(x, x')_{\text{RBF}}$, the White Noise kernel is denoted as $k(x, x')_{\text{WN}}$ and finally the Matern kernel is denoted as $k(x, x')_M$. All these kernels were implemented using the `sklearn.model_selection.kernels` library with the default hyperparameters set for each kernel.

As explained in Question 1 Part 1 (a), the absorption does not follow any simple pattern like a linear relationship or a quadratic one. For this reason, these kernels were picked because they can model complex relationships like the one seen between the wavelength and the absorption.

The kernel combinations below were tested using k-fold validation.

$$\begin{aligned} & k(x, x')_M + k(x, x')_{\text{RBF}} + k(x, x')_{\text{WN}} \\ & k(x, x')_M + k(x, x')_{\text{RBF}} \\ & k(x, x')_M \\ & k(x, x')_{\text{ESS}} + k(x, x')_M + k(x, x')_{\text{RBF}} \\ & k(x, x')_{\text{RBF}} \end{aligned}$$

A different combinations of all of these was tested using 3-fold validation to test which combination worked best to predict the absorption at any given the wavelength. A box plot of the results are plotted below with each different combination tested.

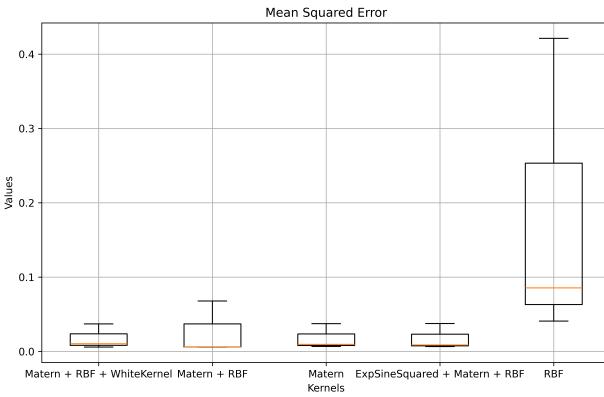


Figure 5: Boxplot of MSE for Different Kernels

Using the combinations of the Matern, RBF and white noise kernels, the mean squared error was the smallest at 0.015.

Using this kernel combination, the Gaussian regressor was trained using all the values in sample i^* . The predictions for each wavelength compared to the true absorption was plotted below.

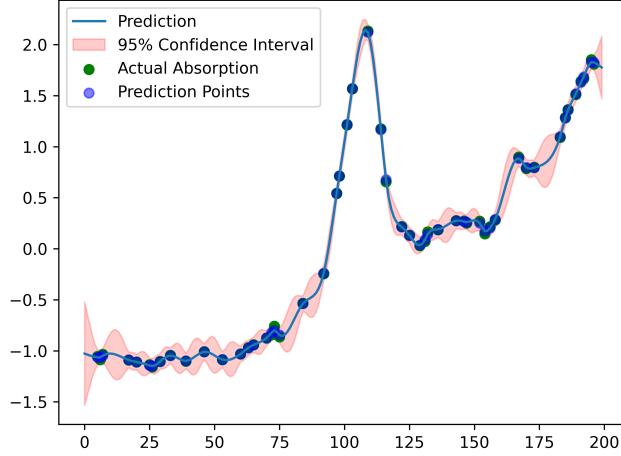


Figure 6: Predictions Vs Actual Values

(3)

To perform this Gaussian process considering the data set \mathcal{D} excluding the i^* -th chemical sample, a data set of the covariants and of the response were created.

The data set of the covariants contained each of the 39 chemical and temperature components repeated 199 times with each of the wavelengths. The response data set contained the corresponding absorption at each of the covariants data sets points.

The kernel was chosen by testing how different kernels fitted the data after being put thought the Gaussian process.

This process was very similar to the process found in Part 2 section 2 of this question, were k-fold cross validation below.

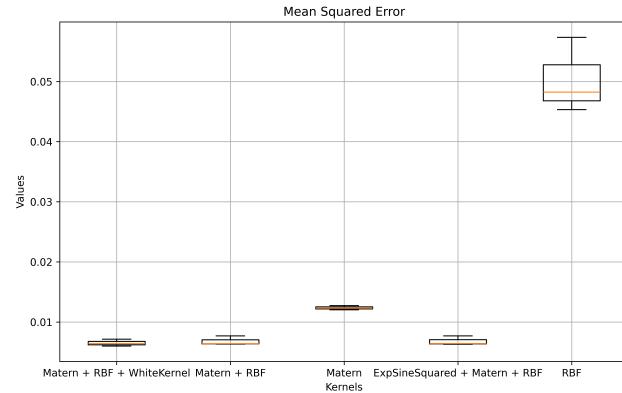


Figure 7: Boxplot of MSE for Different Kernels

Just like the previous section, the combinations of the Matern, RBF and white noise kernel is the smallest.

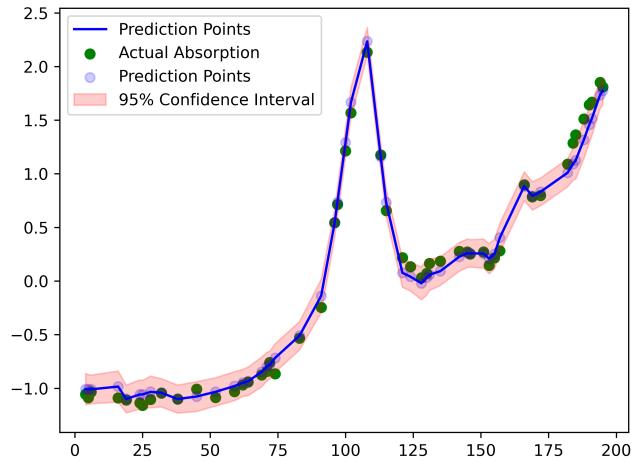


Figure 8: Predictions Vs Actual Values

The mean squared error between the predicted values and the actual values is 0.0079. This is a better predictor than the Gaussian process cross validations in the previous part of this question.

(4)

(a)

The model found to be the best from testing in Part 2 section 3 was used to train over the entire data set \mathcal{D} .

A list of 100 evenly spaced values between the minimum and maximum temperature value in \mathcal{D} were created. Each temperature on this list was used with the test components and the 199 wavelengths to find 199 predictions of the absorption's at each wavelength and the σ (standard deviation) value associated with each prediction.

Then the Negative Log Predictive Density (NLPD) is calculated using the following equations,

$$\text{Wavelength}_i \text{ NLPD} = \left(\log(\sigma_i \sqrt{2\pi}) + \frac{1}{2} \left(\frac{(y_{\text{pred},i} - y_{\text{test},i})^2}{\sigma_i^2} \right) \right)$$

$$\text{Mean NLPD} = \frac{1}{n} \sum_{i=1}^n \left(\log(\sigma_i \sqrt{2\pi}) + \frac{1}{2} \left(\frac{(y_{\text{pred},i} - y_{\text{test},i})^2}{\sigma_i^2} \right) \right)$$

This was calculated for every temperature value in the list of possible temperature values and is plotted below.

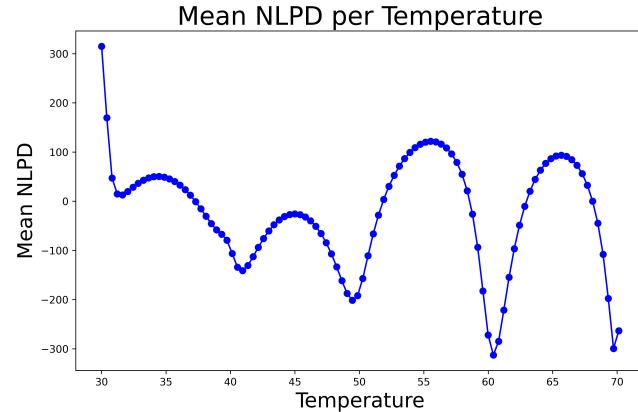


Figure 9: Histogram of correlation coefficients.

The Negative Log Predictive Density (NLPD) has the smallest mean when the temperature is 60.40. From this Gaussian process, which has been shown to be very accurate, one can conclude that the temperature of this chemical sample is 60.40 or a value very close.

(b)

A similar process to the one described above was preformed using the test temperature 69.99 and the second set of observations of the absorption at each of the 199 wavelengths.

The important difference is that there are 3 components in the covariants while only one covariant was used in (a). This means that to create a list of components to test, 3 different factors must be considered

To do this a list of 15 evenly spaced numbers between 0 and 1 was created, then a $15^3 \times 3$ dimensions data frame was created which contains every single possible combination of this original list.

Using the test temperature provided and the 199 wavelengths, predictions and their associated value of sigma was calculated for each combination of 3 components.

Using these prediction values the mean NLPD was calculated for each combination of components.

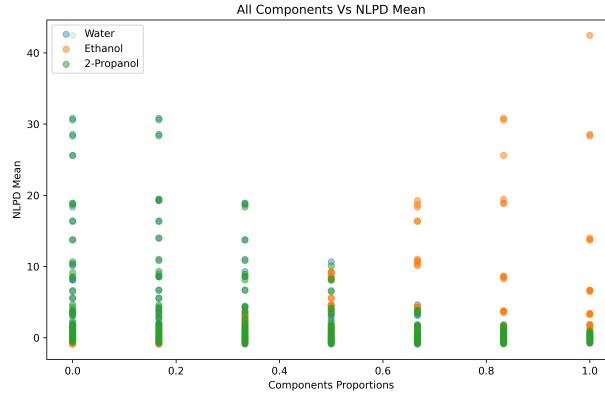


Figure 10: Histogram of correlation coefficients.

The mixture of the 3 components that lead to the smallest NLPD are tabled below, the exact component values are likely not these exact values but something close.

Component	Best Value
Water	0.11111111
Ethanol	0.0
2-Propanol	0.55555556

Table 7: Components for Minimum Value of NLPD mean

Question 2

(1)

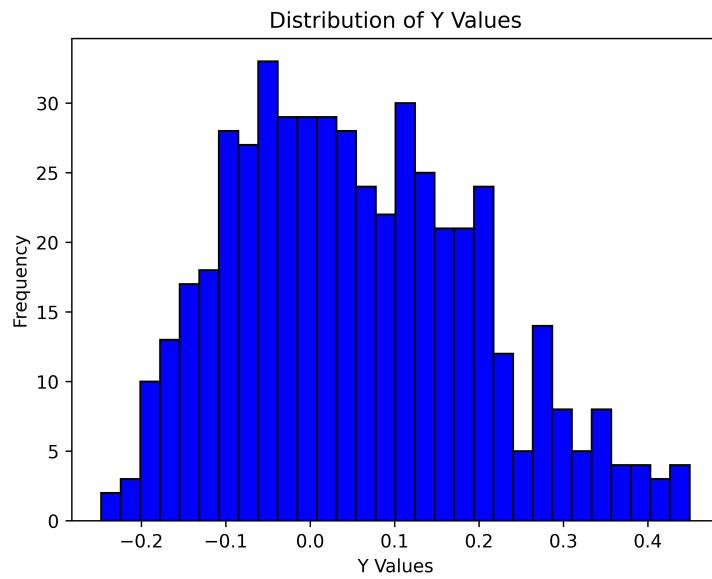


Figure 11: Histogram of 'y' Values

The histogram above shows the minimum value on the x-axis is roughly -0.3 and the maximum is roughly 0.5. This histogram is relatively well rounded between these values, and the mean is 0.056.

There are 49 covariants and they are very diverse; the below graph shows the distribution of the 49 covariants. Most of them have a mean less than one and a few outliers have a mean significantly greater than 1.

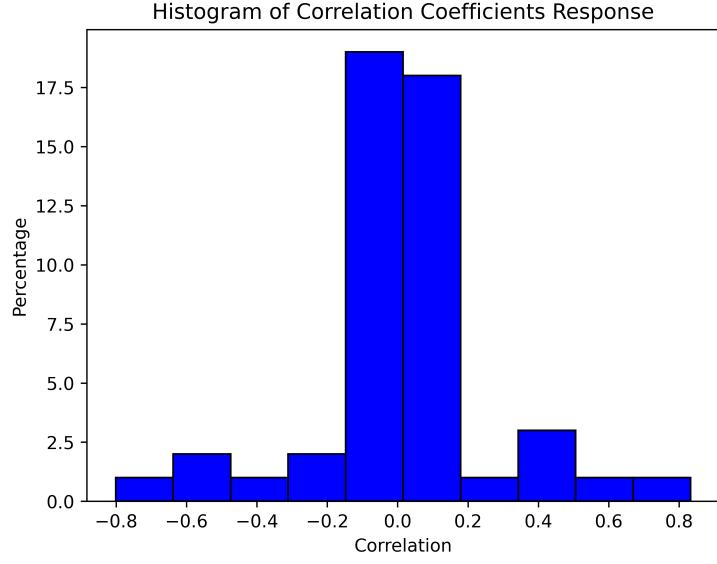


Figure 12: Histogram of Correlations

As can be seen in the graph above, many of the covariants are not strongly correlated with the response variable. Only 10 of the covariants have a absolute correlation value greater than 0.2.

(2)

The 4 machine learning models to predict the output of interest Y given the other biological features were, Decision Trees Regression, Decision Tree Regression, a k-nearest neighbour regression model and a custom Neural network.

The data was split into a training data set and a test data set using a 70%-30% split. All testing of hyperparameters was tested on the training dataset.

Throughout this section, the k-fold cross validation procedure is used to test hyperparameters on the training dataset.

The k-fold cross validation procedure used throughout uses 5 folds. This number of folds is chosen to reduce bias and variance while also not requiring excessive computational power.

The k-fold cross validation used involves splitting the data into 5 subsets. For each iteration 4 of the 5 subsets is used to train the model. Then, the model is tested against the subset not used for training. The step is repeated 5 times.

The loss function used to evaluate the performance of all the machine learning models was the Mean Squared Error (MSE) loss function.

Decision Trees

Since the output of the data is continuous and not discrete, a Regression Decision Tree is used instead of a Classifier Decision Tree. The models were made using the `DecisionTreeRegressor` function as part of the library `sklearn.tree`.

Different Loss criteria were set as hyperparameters in the model and then these hyperparameters were tested using k-fold validation with a MSE loss function. These were the 'absolute error', 'squared error' and 'friedman mse' loss functions.

Fitting Loss function	Mean MSE
Squared Error	8.83×10^{-3}
Absolute Error	8.13×10^{-3}
Friedman Mean Squared Error	9.04×10^{-3}

Table 8: Training Loss Function

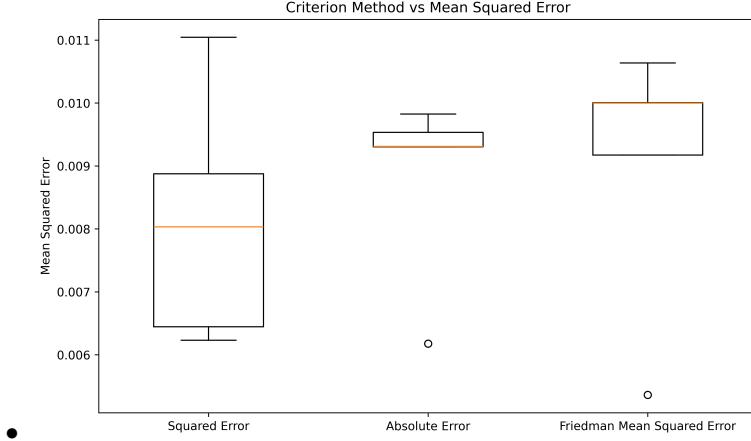


Figure 13: Boxplot of MSE form each Loss Function

The performance of each hyperparameter was tested by using the k-fold cross validation procedure explained above on the training data.

It was tested whether a random training splitter or a best training splitter produced a smaller MSE error. This was tested by performing k-fold validation on the training data.

Splitter	Mean MSE
Best	$8.125722277687000 \times 10^{-3}$
Random	$11.758628541753999 \times 10^{-3}$

Table 9: MSE Comparison for Different Data Splitters

As can be seen in the table above the best training split produces a lower MSE error. The hyperparameters were chosen that minimised the MSE through the cross validations.

Random Forest Regression

As can be seen in the histogram above in question 2 part 1, most of the covariants have a very small correlation with the response variable. In the Decision Trees Regression model it was both tested using the entire data set and using the 10 covariants in which the correlation with the predictor has a absolute value greater than 0.2. The reduced data frame of covariants is referred to as `X_filtered` in the code in the appendix below.

Since the output of the data is continuous and not discrete, a Random Forest Regression is used instead of a Random Forest Classifier. The models were made using the `RandomForestRegressor` function as part of the library `sklearn.ensemble`.

As can be seen in the table below, the reduced dataset has a noticeably lower mean MSE under k-fold cross validation than when training is performed on the entire dataset.

n Estimators	MSE Full	MSE Filtered	Difference
100	4.815×10^{-3}	4.033×10^{-3}	0.786×10^{-3}
350	4.742×10^{-3}	4.018×10^{-3}	0.724×10^{-3}
500	4.714×10^{-3}	4.019×10^{-3}	0.695×10^{-3}
1000	4.706×10^{-3}	4.008×10^{-3}	0.698×10^{-3}

Table 10: MSE with full and reduced data sets against different number of estimators

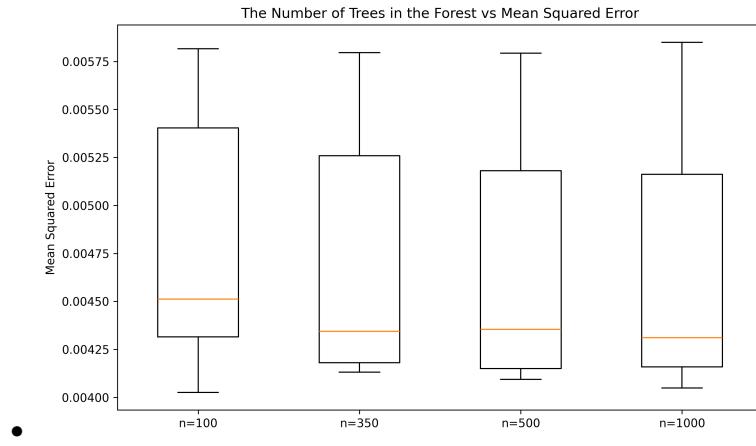


Figure 14: Boxplot of MSE Loss on complete data set against Different Numbers of Estimators

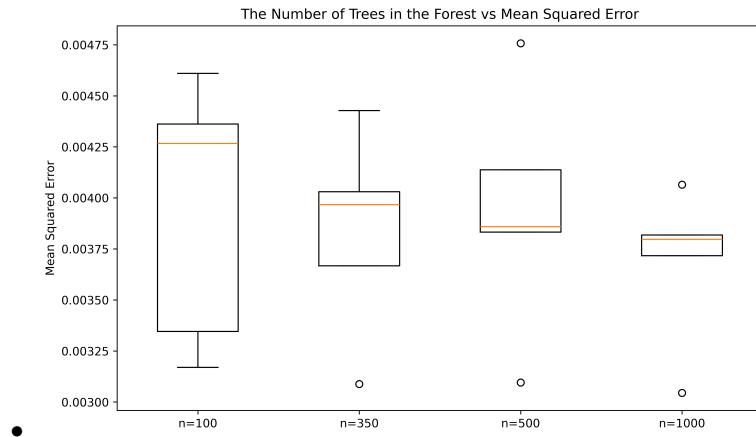


Figure 15: Boxplot of MSE Loss on reduced data set against Different Numbers of Estimators

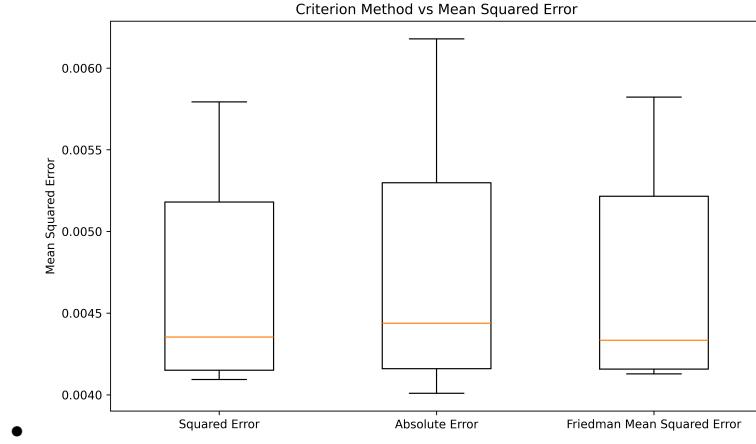


Figure 16: Boxplot of MSE Loss on complete data set against Different Loss Criteria

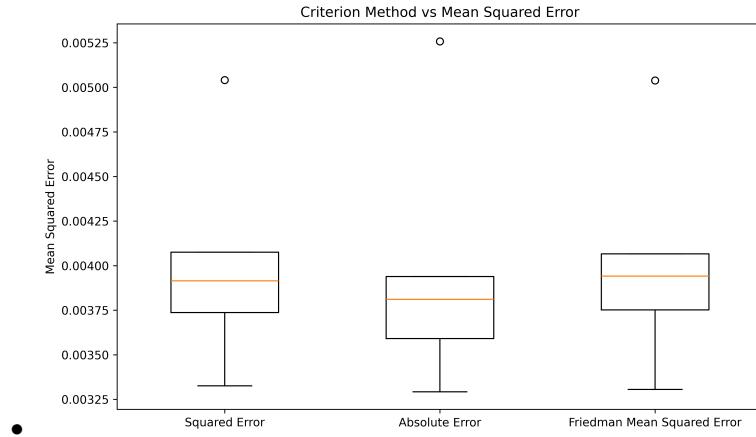


Figure 17: Boxplot of MSE Loss on reduced data set against Different Loss Criteria

As can be seen in boxplots and table above as the number of estimators increases the smaller the mean MSE under k-fold cross validation becomes. However as the number of estimators increases so does the computational costs increase.

Different Loss criteria were set as hyperparameters for the Random Forest Regression model. Then, the hyperparameters were tested using k-fold validation with a MSE loss function. These were the 'absolute error', 'squared error' and 'friedman mse' loss functions.

From these results, 500 estimators were chosen. The mean squared error function was used as the loss function.

K-Nearest neighbour

A k-Nearest neighbour algorithms were implemented using the `KNeighborsRegressor` functions as part of the `sklearn.neighbors` library.

To find the best hyperparameters different values for k between 3 and 8 were tested in k-nearest neighbour algorithm. The k-nearest neighbour algorith was also tested using the 2 different distance metrics 'chebyshev' and 'manhattan'.

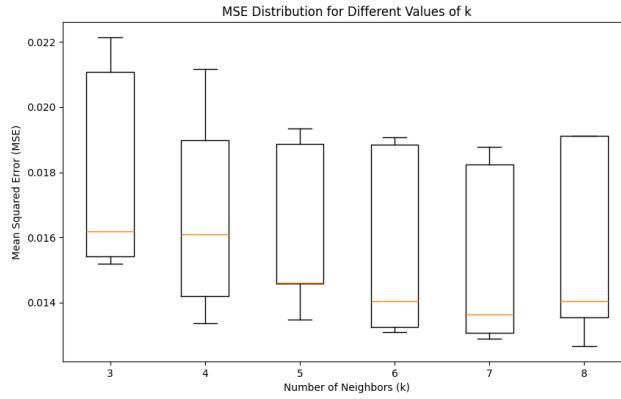


Figure 18: MSE for different values of k with chebyshev distance metric

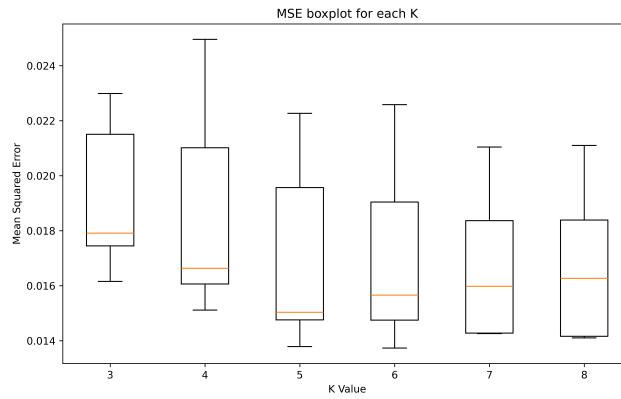


Figure 19: MSE for different values of k with manhattan distance metric

k Value	Mean MSE (Chebyshev)	Mean MSE (Manhattan)
3	0.018	0.0192
4	0.0168	0.0188
5	0.0162	0.0171
6	0.0157	0.0172
7	0.0153	0.0168
8	0.0157	0.0168

Table 11: Mean MSE values for different k values using Chebyshev and Manhattan Distance

The lowest MSE is found when the distance metric is 'Chebyshev' and with $k = 7$. So this was used on the test data.

Neural network

A neural network was used to fit to the data. The neural network is a combination of the Rectified Linear Unit (ReLU) function and linear functions. This was implemented using functions and optimiser from the `torch` library. The exact architecture of the neural network can be found in the code appendix.

(3)

Method	Mean Squared Error
Random Forest Regression	0.0029
Decision Tree Regression	0.006
Neural Network	0.01922
K-nearest neighbour	0.015

Table 12: Summary of Mean Squared Error (MSE) for each Machine Learning model

(4)

From each of the models, it is clear that the Decision Tree Regression and Random Forest Regression preformed better than the Neural Network and k nearest neighbour with given hyperparameters. The Random Forest Regression had the smallest MSE on the Test data set.

Because of the fact that the Random Forest Regression on the reduced data set consistently preforms better under k-fold cross validation

Code Appendix

```
1 import pandas as pd
2 import numpy as np
3 import matplotlib.pyplot as plt
4
5 from sklearn.model_selection import train_test_split
6 from sklearn.metrics import mean_squared_error
7 from sklearn.svm import SVR
8 from sklearn.model_selection import cross_val_score, KFold
9
10 import seaborn as sns
11
12 from scipy.cluster.hierarchy import linkage, dendrogram
13 from scipy.cluster.hierarchy import cut_tree
14
15 seed=2524541
16
17 np.random.seed(seed)
```

```
1 # Import Dataset
2 components = pd.read_csv('components.csv', header = None, sep=' ')
3 spectra = pd.read_csv('spectra.csv', header=None, sep=' ')
4 temperatures = pd.read_csv('temperatures.csv', header=None)
```

Question 1

Part 1

(a)

```
1 plt.plot(np.mean(spectra, axis=0))
2
3 plt.title('Mean Absorption per Wavelength', size = 25)
4
5 plt.ylabel('Mean Absorption', size =20)
6
7 plt.xlabel('Wavelength', size = 20)
8
9 step = 25
10 original_indices = np.arange(0, 201, step)
11 mapped_indices = original_indices + 800
12
13 plt.xticks(ticks=original_indices, labels=mapped_indices)
14
15 plt.gcf().set_size_inches(10, 6)
16
```

```

17 plt.savefig('Q1_analysis_1.png', dpi=600, bbox_inches='tight')
18
19 plt.show()

```

```

1 names = ["water", "ethanol", "2-propanol"]
2 for i in range(3):
3     print(f"The Correlation Coefficient between components
        ↪ {names[i]} and absorption is",
4         np.corrcoef(np.mean(spectra, axis=1), components[1])[0][1])
5 print(f"The Correlation Coefficient between temperature {names[i]}
        ↪ and absorption is",
6         np.corrcoef(np.mean(spectra, axis=1), temperatures[0])[0][1])

```

(b)

```

1 # Import hierarchy cluster packages need
2 from sklearn.decomposition import PCA
3 from scipy.cluster.hierarchy import linkage, fcluster, dendrogram

```

```

1 # standardized the spectra values
2 standardized_spectra = (spectra - spectra.mean()) / spectra.std()
3
4 Metrics = ["euclidean", "cityblock"]
5 Methods = ["complete", "single", "average"]
6
7 for method in Methods:
8     for metric in Metrics:
9         fig, axes = plt.subplots(figsize=(10, 6))
10        linkage_data = linkage(standardized_spectra, method=method,
11                               ↪ metric=metric)
12        dendrogram(linkage_data, ax=axes)
13        # remove x values
14        axes.set_xticklabels([])
15        axes.set_title(f"Distance metric {metric} and method
16                      ↪ {method}", size=25)
17        axes.set_ylabel(f'{metric} distance', size=20)
18        fig.savefig(f"Q1b_{method}_{metric}.png", dpi=600)

```

```

1 # Linkage data with method='complete' and metric='euclidean'
2 linkage_data = linkage(standardized_spectra, method='complete',
3                         ↪ metric='euclidean')
4 cluster_labels = fcluster(linkage_data, t=4, criterion='maxclust')
5
6 for i in range(1,5):
7     components_values = components.iloc[cluster_labels==i].mean()
8     print(f"Group {i}")

```

```

8     print("Water Mean", components_values[0])
9     print("Ethanol Mean", components_values[1])
10    print("2-propanol Mean", components_values[2])
11    print("Temperatures
12      ↪  Mean",temperatures.iloc[cluster_labels==i].mean())
13    print('-----')
1
# Linkage data with method='average' and metric='cityblock'
2 linkage_data = linkage(standardized_spectra, method='average',
3   ↪  metric='cityblock')
4 cluster_labels = fcluster(linkage_data, t=4, criterion='maxclust')
5
# Print Values
6 for i in range(1,5):
7     components_values = components.iloc[cluster_labels==i].mean()
8     print(f"Group {i}")
9     print("Water Mean", components_values[0])
10    print("Ethanol Mean", components_values[1])
11    print("2-propanol Mean", components_values[2])
12    print("Temperatures
13      ↪  Mean",temperatures.iloc[cluster_labels==i].mean())
14    print('-----')

```

(c)

```

1 # Define principal component analysis
2 pca_calculator = PCA()
3
4 pca_calculator.fit(standardized_spectra)
5
6 # Get variance explained on every value
7 var_explained = pca_calculator.explained_variance_ratio_
8 cumvar_explained = np.cumsum(var_explained)
9
10
11 # Plot Explained Cumulative Variance
12 plt.figure(figsize=(10,6))
13
14 plt.plot(cumvar_explained, linestyle='--')
15 plt.xlabel('Number of Components', size = 20)
16 plt.ylabel('Cumulative variance', size = 20)
17 plt.title('Explained Cumulative Variance', size = 25)
18 plt.savefig('Q1_c_cumvar.png')
19
20 plt.show()
21
22

```

```

23
24 # get the number of values that explain 99.99% of the variance
25 min_cumvar_explained = len(cumvar_explained) -
→ sum((cumvar_explained>=0.9999)*1) + 1
26
27 print("Number of components needed to explain 99.99% is",
→ min_cumvar_explained)

```

```

1 n=4
2 linkage_data = linkage(Simplified_data, method='complete',
→ metric='euclidean')
3 cluster_labels = fcluster(linkage_data, t=n, criterion='maxclust')
4
5 for i in range(1,n+1):
6     components_values =
→ components.iloc[Simplified_data.index[cluster_labels==i]].mean()
7     print(f"Group {i}")
8     print("Water Mean", components_values[0])
9     print("Ethanol Mean", components_values[1])
10    print("2-propanol Mean", components_values[2])
11    print("Temperatures
→ Mean",temperatures.iloc[Simplified_data.index[cluster_labels==i]].mean())
12    print('-----')

```

```

1 # Print mean values for different data sets
2 print(temperatures.mean(),components.mean())
3 print(temperatures.loc[Simplified_data.index].mean(),
4 ,components.loc[Simplified_data.index].mean())

```

```

1 linkage_data = linkage(Simplified_data, method='average',
→ metric='cityblock')
2 cluster_labels = fcluster(linkage_data, t=n, criterion='maxclust')
3
4 for i in range(1,n+1):
5     components_values =
→ components.iloc[Simplified_data.index[cluster_labels==i]].mean()
6     print(f"Group {i}")
7     print("Water Mean", components_values[0])
8     print("Ethanol Mean", components_values[1])
9     print("2-propanol Mean", components_values[2])
10    print("Temperatures Mean"
11
→ ,temperatures.iloc[Simplified_data.index[cluster_labels==i]].mean())
12    print('-----')

```

Part 2

(a)

```
1 seed_value = 2524541
2 n=40
3 components_sampled = components.sample(n=n,
   ↪ random_state=seed_value)
4 index = components_sampled.index
5
6 spectra_noisy_sampled = spectra.iloc[index]
7 temperatures_sampled = temperatures.iloc[index]
8 components_sampled = components.iloc[index]
9
10 np.random.seed(seed)
11
12 num_rows = 40
13 num_columns = 199
14 num_data_points = num_rows * num_columns
15 random_data = np.random.normal(scale=np.sqrt(0.0025),
   ↪ size=spectra_noisy_sampled.shape)
16 spectra_noisy_sampled += random_data
17
18 np.random.seed(seed)
19 NaN_index = np.array([True]*int(num_data_points * .8) +
   ↪ [False]*int(num_data_points * .2))
20 np.random.shuffle(NaN_index)
21 NaN_index = NaN_index.reshape(spectra_noisy_sampled.shape)
22
23 spectra_noisy_sampled[NaN_index] = np.nan
```

(b)

```
1 # import kernel packages
2 from sklearn.gaussian_process import GaussianProcessRegressor
3 from sklearn.gaussian_process.kernels import RBF, WhiteKernel,
   ↪ Matern, ExpSineSquared
4
5 # Get the row with the least amount of na values
6 nan_per_row = spectra_noisy_sampled.isna().mean(axis=1)
7 nan_min_row = nan_per_row.idxmin()
8 min_na_value_sample = spectra_noisy_sampled.loc[nan_min_row]
9 X = np.arange(1,200)[min_na_value_sample.notna()].reshape(-1, 1)
10 y = min_na_value_sample.to_numpy()[min_na_value_sample.notna()]
11
12 # Define the k-fold cross validation
13 def k_fold_cross_validation(kernel, X=X, y=y, k=3, seed=seed):
```

```

3     gp = GaussianProcessRegressor(kernel=kernel,
4         ↪ n_restarts_optimizer=10)
5     K_Fold = KFold(n_splits=k, shuffle=True, random_state=seed)
6     mse = cross_val_score(gp, X, y, cv=K_Fold,
7         ↪ scoring='neg_mean_squared_error')
8     return -mse

```

```

1 # Use k-fold validations to see which kernel creates the best
2 ↪ scores
3 mses = []
4
5 kernel = Matern() + RBF() + WhiteKernel()
6 mses.append(k_fold_cross_validation(kernel=kernel, X=X, y=y, k = 3,
7   ↪ seed = seed))
8
9 kernel = Matern() + RBF()
10 mses.append(k_fold_cross_validation(kernel=kernel, X=X, y=y, k = 3,
11   ↪ seed = seed))
12
13 kernel = ExpSineSquared() + Matern() + RBF()
14 mses.append(k_fold_cross_validation(kernel=kernel, X=X, y=y, k = 3,
15   ↪ seed = seed))
16
17 kernel = RBF()
18 mses.append(k_fold_cross_validation(kernel=kernel, X=X, y=y, k = 3,
19   ↪ seed = seed))

```

```

1 # Name of each combinations
2 knames = ["Matern + RBF + WhiteKernel", "Matern + RBF",
3           "Matern", "ExpSineSquared + Matern + RBF", "RBF"]
4 # Plot boxplots of the mses
5 plt.figure(figsize=(10, 6))
6 plt.boxplot(mses, labels=knames)
7 plt.title("Mean Squared Error")
8 plt.xlabel("Kernels")
9 plt.ylabel("Values")
10 plt.savefig("Question_1_part_2a_boxplot.png", dpi=600)
11 plt.show()

```

```

1 # Find the smallest MSE and the kernel combinations
2 mse = [np.mean(mse_group) for mse_group in mses]
3

```

```

4 print("The smallest MSE is with kernel
      ↳ combination",knames[np.argmin(mse)])
5
6 print("The smallest MSE is",np.min(mse))
7

```

```

1 # All wave lengths
2 X_long = np.linspace(0,199,200).reshape(-1,1)
3 # Kernel with lowest k-fold cross validation
4 kernel = Matern() + RBF() + WhiteKernel()
5 # fit GaussianProcess
6 gp = GaussianProcessRegressor(kernel=kernel, alpha=5e-10)
7 gp.fit(X, y)
8 #Calculate the prediction for i*
9 y_pred_train = gp.predict(X, return_std=False)
10 #Calculate the prediction for every wavelength
11 y_pred_long, sigma = gp.predict(X_long, return_std=True)

```

```

1 # Plot the results with the confidence intervals
2 plt.plot(X_long, y_pred_long, label='Prediction')
3 plt.fill_between(X_long.squeeze(),y_pred_long - sigma *
      ↳ 1.96,y_pred_long + sigma * 1.96,color='red',
4                  alpha=0.2,label='95% Confidence Interval')
5 plt.scatter(X, y, color='green', label='Actual Absorption')
6 plt.scatter(X, y_pred_train, color='blue', label='Prediction
      ↳ Points',alpha = 0.5)
7 plt.legend()
8 # Adding labels and legend
9 plt.savefig("Question_1_part_2a_predictions.png", dpi=600)
10 plt.show()

```

(c)

```

1 # Consider entire index except the
2 filtered_spectra = spectra_noisy_sampled.drop(index=nan_min_row)
3 filtered_components = components_sampled.drop(index=nan_min_row)
4 filtered_temperatures =
5      ↳ temperatures_sampled.drop(index=nan_min_row)
6
7 spectra_star = spectra_noisy_sampled.loc[nan_min_row]
8 components_star = components_sampled.iloc[nan_min_row].to_frame().T
9 temperatures_star = temperatures_sampled.loc[nan_min_row]
10
11 # Create data frame fit for gaussian regression
12 spectra_longer = pd.DataFrame({'Column': np.repeat(np.arange(0,
      ↳ 199), 39)})

```

```

12 components_longer = pd.concat([filtered_components] * 199,
13   ↪ ignore_index=True)
14 temperatures_longer = pd.concat([filtered_temperatures] * 199,
15   ↪ ignore_index=True)
16 y = filtered_spectra.values.flatten(order='F')
17 X =
18   ↪ pd.concat([temperatures_longer,components_longer,spectra_longer],axis=1)
19 # Create data frame fit for gaussian regression
20 spectra_longer = pd.DataFrame({'1': np.repeat(np.arange(0, 199),
21   ↪ 1)})
22 components_star = pd.concat([components_star] * 199,
23   ↪ ignore_index=True)
24 temperatures_star = pd.concat([temperatures_star] * 199,
25   ↪ ignore_index=True)
26 y_star = spectra_star.values.flatten(order='F')
27 X_star =
28   ↪ pd.concat([temperatures_star,components_star,spectra_longer],axis=1)
29 # remove na values
30 indices = np.where(~np.isnan(y))[0]
31 y = y[indices]
32 X = X.iloc[indices,:].to_numpy()
33 # remove na values
34 indices = np.where(~np.isnan(y_star))[0]
35 y_star = y_star[indices]
36 X_star = X_star.iloc[indices,:].to_numpy()

```

```

1 mses = []
2
3 kernel = Matern() + RBF() + WhiteKernel()
4 mses.append(k_fold_cross_validation(kernel=kernel, X=X, y=y, k = 3,
5   ↪ seed = seed))
6 kernel = Matern() + RBF()
7 mses.append(k_fold_cross_validation(kernel=kernel, X=X, y=y, k = 3,
8   ↪ seed = seed))
9 kernel = Matern()
10 mses.append(k_fold_cross_validation(kernel=kernel, X=X, y=y, k = 3,
11   ↪ seed = seed))
12 kernel = ExpSineSquared() + Matern() + RBF()
13 mses.append(k_fold_cross_validation(kernel=kernel, X=X, y=y, k = 3,
14   ↪ seed = seed))

```

```

15 kernel = RBF()
16 mses.append(k_fold_cross_validation(kernel=kernel, X=X, y=y, k = 3,
→ seed = seed))

```

```

1 # Find the smallest MSE and the kernel combinations
2 mse = [np.mean(mse) for mse in mses]
3
4 print("The smallest MSE is with kernel
→ combination",knames[np.argmin(mse)])
5
6 print("The smallest MSE is",np.min(mse))

```

```

1 knames = ["Matern + RBF + WhiteKernel","Matern + RBF",
2      "Matern","ExpSineSquared + Matern + RBF","RBF"]
3 plt.figure(figsize=(10, 6))
4 plt.boxplot(mses, labels=knames)
5 plt.title("Mean Squared Error")
6 plt.xlabel("Kernels")
7 plt.ylabel("Values")
8 plt.grid(True)
9 plt.savefig("Question_1_part_2b_boxplot.png", dpi=600)
10 plt.show()

```

```

1 # Kernel
2 kernel = Matern() + RBF() + WhiteKernel()
3 # fit GP
4 gp_regressor = GaussianProcessRegressor(kernel=kernel, alpha=5e-10)
5
6 gp_regressor.fit(X, y)
7
8 y_pred_star, sigma = gp_regressor.predict(X_star, return_std=True)
9

```

```

1 plt.plot(X_star[:,4], y_pred_star, color='blue', label='Prediction
→ Points')
2 plt.scatter(X_star[:,4], y_star, color='green', label='Actual
→ Absorption')
3 plt.scatter(X_star[:,4], y_pred_star, color='blue',
→ label='Prediction Points',alpha = 0.2)
4 plt.fill_between(X_star[:, 4], y_pred_star - sigma * 1.96,
→ y_pred_star + sigma * 1.96, color='red',
→ alpha=0.2,
→ label='95%\% Confidence Interval')
5 plt.legend()
6 plt.savefig("Question_1_part_2b_predictions.png", dpi=600)
9 plt.show()

```

```

1 mse = mean_squared_error(y_star, y_pred_star)
2 print("The mean squared error is",mse)

```

4

(a)

```

1 # Create data frame fit for gaussian regression
2 spectra_longer = pd.DataFrame({'Column': np.repeat(np.arange(0,
   ↪ 199), 40)})
3 components_longer = pd.concat([components_sampled] * 199,
   ↪ ignore_index=True)
4 temperatures_longer = pd.concat([temperatures_sampled] * 199,
   ↪ ignore_index=True)
5 y = spectra_noisy_sampled.values.flatten(order='F')
6 X =
   ↪ pd.concat([temperatures_longer,components_longer,spectra_longer]
   ↪ ,axis=1)
7
8
9 # remove na values
10 indices = np.where(~np.isnan(y))[0]
11 y = y[indices]
12 X = X.iloc[indices,:].to_numpy()

```

```

1 kernel = RBF() + WhiteKernel() + Matern()
2 gp_regressor = GaussianProcessRegressor(kernel=kernel)
3 gp_regressor.fit(X, y)

```

```

1 components_test = pd.read_csv('components_test.csv',header = None,
   ↪ sep = ' ')
2 spectra_test = pd.read_csv('spectra_test.csv',header=None, sep =
   ↪ ' ')
3 temperatures_test =
   ↪ pd.read_csv('temperatures_test.csv',header=None)
4
5 n = 100
6 min_temperatures = np.min(temperatures)
7 max_temperatures = np.max(temperatures)
8 temperature_grid = np.linspace(min_temperatures, max_temperatures,
   ↪ n)

```

```

1 spectra_longer_test = pd.DataFrame({'Column': np.arange(0, 199)})
2 components_longer_test =
   ↪ pd.concat([components_test.loc[0].to_frame().transpose()] * 199
   ↪ , ignore_index=True)
3
4
5 one_d_array = spectra_test.iloc[0].values.flatten(order='F')

```

```

6 y_test_1 = pd.DataFrame({'y': one_d_array}).to_numpy().squeeze()
7 X_test_1 =
    ↳ pd.concat([components_longer_test,spectra_longer_test],axis=1).to_numpy()

```

```

1 nlpd_mean = []
2 for i in range(n):
3
4     constant_column = np.full((X_test_1.shape[0], 1),
    ↳ temperature_grid[i])
5
6     X_test_np = np.hstack((constant_column, X_test_1))
7
8     y_pred, sigma = gp_regressor.predict(X_test_np,
    ↳ return_std=True)
9     nlpd = -np.log(sigma * np.sqrt(2 * np.pi)) - 0.5 * (((y_pred -
    ↳ y_test_1) ** 2) / (sigma**2))
10    nlpd_mean.append(-nlpd.mean())
11

```

```

1 plt.figure(figsize=(10, 6))
2 plt.plot(temperature_grid, nlpd_mean, marker='o', linestyle='--',
    ↳ color='blue')
3 plt.title('Mean NLPD per Temperature', fontsize=25)
4 plt.xlabel('Temperature', fontsize=20)
5 plt.ylabel('Mean NLPD', fontsize=20)
6 plt.savefig("Question_1_part_2c_tempatures.png", dpi=600)
7 plt.grid(True)
8
9 plt.show()
10
11 print("Minimum value for nlpd mean is when the temperature is",
    ↳ temperature_grid[np.argmin(nlpd_mean)])

```

(b)

```

1 import itertools
2
3 n = 10
4 numbers = np.linspace(0, 1, n)
5
6 components_combinations = itertools.product(numbers, repeat=3)
7
8 components_combinations = np.array(list(components_combinations))

```

```

1 spectra_longer_test = pd.DataFrame({'Column': np.arange(0,
    ↳ 199)}).to_numpy()

```

```

2 temperatures_longer_test =
3     ↪ np.repeat([temperatures_test[0][1]],199).reshape(199,1)
4 one_d_array = spectra_test.iloc[1].values.flatten(order='F')
5 y_test_1 = pd.DataFrame({'y': one_d_array}).to_numpy().squeeze()
6 X_test_1 =
7     ↪ np.hstack([temperatures_longer_test,spectra_longer_test])

```

```

1 nlpd_mean = []
2 for i in range(components_combinations.shape[0]):
3     constant_column = np.full((X_test_1.shape[0], 3),
4         ↪ components_combinations[i])
4     X_test_np = np.hstack((X_test_1[:,[0]],constant_column,
5         ↪ X_test_1[:,[1]]))
5     y_pred, sigma = gp_regressor.predict(X_test_np,
6         ↪ return_std=True)
6     nlpd=np.log(sigma * np.sqrt(2 * np.pi)) + 0.5 * ((y_pred -
7         ↪ y_test_1) ** 2) / sigma**2
7     nlpd_mean.append(np.mean(nlpd))

```

```

1 plt.figure(figsize=(10, 6))
2 plt.scatter(components_combinations[:, 0], nlpd_mean,
3     ↪ label='Water', alpha=.4)
3 plt.scatter(components_combinations[:, 1], nlpd_mean,
4     ↪ label='Ethanol', alpha=0.5)
4 plt.scatter(components_combinations[:, 2], nlpd_mean,
5     ↪ label='2-Propanol', alpha=.5)
5
6 plt.title('All Components Vs NLPD Mean')
7 plt.xlabel('Components Proportions')
8 plt.ylabel('NLPD Mean')
9 plt.legend(loc = 'upper left')
10 plt.savefig("Question_1_part_2_4b_NLPD.png", dpi=600)
11 plt.show()
12
13 print("Minimum value for nlpd mean is when the componts are",
    ↪ components_combinations[np.argmin(nlpd_mean)])

```

Question 2

1

```

1 data = pd.read_csv('02524541.csv')
2 y = data.y.to_numpy()
3 X = data.iloc[:, 1: ].to_numpy()

```

```

1 plt.hist(y, bins=30, color='blue', edgecolor='black')
2 plt.title("Distribution of Y Values")
3 plt.xlabel("Y Values")
4 plt.ylabel("Frequency")
5
6
7 plt.savefig('Q2_y_hist.png', dpi=600, format='png',
   ↪ bbox_inches='tight')
8
9 plt.show()

```

```

1 low_corr_columns = correlations[abs(correlations) < 0.2].index
2 filtered_data = data.drop(columns=low_corr_columns)
3 X_filtered = filtered_data.iloc[:, 1: ].to_numpy()

```

2

```

1 X_train, X_test, y_train, y_test= train_test_split(X,y,
   ↪ random_state=seed, test_size=0.3,
                           shuffle=True)
2
3
4 correlations = data.corr()['y'].drop('y')
5 low_corr_columns = correlations[abs(correlations) < 0.2].index
6 filtered_data = data.drop(columns=low_corr_columns)
7 X_filtered = filtered_data.iloc[:, 1: ].to_numpy()
8 X_train_filtered, X_test_filtered, y_train_filtered,
   ↪ y_test_filtered= train_test_split(X_filtered,y,
9
   ↪ random_state=seed,
10
   ↪ test_size=0.3,
11
   ↪ shuffle=True)
12 kf = KFold(n_splits=5, shuffle=True, random_state=seed)

```

Random Forest Regressor

```

1 from sklearn.ensemble import RandomForestRegressor
2
3 ##Random Forest output 1
4 model_comp = []
5 model = RandomForestRegressor(n_estimators=100, random_state=seed)
6 cross_val_scores = cross_val_score(model, X_train, y_train, cv=kf,
   ↪ scoring='neg_mean_squared_error')
7 model_comp.append(-cross_val_scores)
8 print("n=100",np.mean(-cross_val_scores))
9 model = RandomForestRegressor(n_estimators=350, random_state=seed)

```

```

10 cross_val_scores = cross_val_score(model, X_train, y_train, cv=kf,
11     ↪ scoring='neg_mean_squared_error')
12 print("n=350",np.mean(-cross_val_scores))
13 model_comp.append(-cross_val_scores)
14
15 model = RandomForestRegressor(n_estimators=500, random_state=seed)
16 cross_val_scores = cross_val_score(model, X_train, y_train, cv=kf,
17     ↪ scoring='neg_mean_squared_error')
18 print("n=500",np.mean(-cross_val_scores))
19 model_comp.append(-cross_val_scores)
20
21 model = RandomForestRegressor(n_estimators=1000, random_state=seed)
22 cross_val_scores = cross_val_score(model, X_train, y_train, cv=kf,
23     ↪ scoring='neg_mean_squared_error')
24 print("n=1000",np.mean(-cross_val_scores))
25 model_comp.append(-cross_val_scores)

```

```

1  ##Random Forest output 1
2  model_comp_filtered = []
3  model = RandomForestRegressor(n_estimators=100, random_state=seed)
4  cross_val_scores = cross_val_score(model, X_train_filtered,
5      ↪ y_train, cv=kf, scoring='neg_mean_squared_error')
6  model_comp_filtered.append(-cross_val_scores)
7  print("n=100",np.mean(-cross_val_scores))
8
9  model = RandomForestRegressor(n_estimators=350, random_state=seed)
10 cross_val_scores = cross_val_score(model, X_train_filtered,
11      ↪ y_train, cv=kf, scoring='neg_mean_squared_error')
12 print("n=350",np.mean(-cross_val_scores))
13 model_comp_filtered.append(-cross_val_scores)
14
15 model = RandomForestRegressor(n_estimators=500, random_state=seed)
16 cross_val_scores = cross_val_score(model, X_train_filtered,
17      ↪ y_train, cv=kf, scoring='neg_mean_squared_error')
18 print("n=500",np.mean(-cross_val_scores))
19 model_comp_filtered.append(-cross_val_scores)
20
21 model = RandomForestRegressor(n_estimators=1000, random_state=seed)
22 cross_val_scores = cross_val_score(model, X_train_filtered,
23      ↪ y_train, cv=kf, scoring='neg_mean_squared_error')
24 print("n=1000",np.mean(-cross_val_scores))
25 model_comp_filtered.append(-cross_val_scores)

```

```

1  labels = ['n=100', 'n=350', 'n=500', 'n=1000']
2
3  plt.boxplot(model_comp, labels = labels, widths=0.5)
4
5  plt.ylabel('Mean Squared Error')
6
7  plt.title('The Number of Trees in the Forest vs Mean Squared
   ↪   Error')
8
9  plt.savefig('Random_Forest_output_1.png', dpi=600)
10
11 plt.show()

```

```

1  labels = ['n=100', 'n=350', 'n=500', 'n=1000']
2
3  plt.boxplot(model_comp_filtered, labels = labels, widths=0.5)
4
5  plt.ylabel('Mean Squared Error')
6
7  plt.title('The Number of Trees in the Forest vs Mean Squared
   ↪   Error')
8
9
10 plt.savefig('Random_Forest_filtered_output_1.png', dpi=600)
11
12 plt.show()
13

```

```

1  ##Random Forest output 2
2  model_comp_filtered = []
3  n=500
4  model = RandomForestRegressor(n_estimators=n, random_state=seed,
   ↪   criterion='squared_error')
5  cross_val_scores = cross_val_score(model, X_train_filtered,
   ↪   y_train, cv=kf, scoring='neg_mean_squared_error')
6  model_comp_filtered.append(-cross_val_scores)
7  print("criterion='squared_error'",np.mean(-cross_val_scores))
8  model = RandomForestRegressor(n_estimators=n, random_state=seed,
   ↪   criterion='absolute_error')
9  cross_val_scores = cross_val_score(model, X_train_filtered,
   ↪   y_train, cv=kf, scoring='neg_mean_squared_error')
10 print("criterion='absolute_error'",np.mean(-cross_val_scores))
11 model_comp_filtered.append(-cross_val_scores)
12 model = RandomForestRegressor(n_estimators=n, random_state=seed,
   ↪   criterion='friedman_mse')
13 cross_val_scores = cross_val_score(model, X_train_filtered,
   ↪   y_train, cv=kf, scoring='neg_mean_squared_error')
14 print("criterion='friedman_mse'",np.mean(-cross_val_scores))

```

```

15 model_comp_filtered.append(-cross_val_scores)

1 ##Random Forest output 2
2 model_comp = []
3 n=500
4 model = RandomForestRegressor(n_estimators=n, random_state=seed,
   ↪ criterion='squared_error')
5 cross_val_scores = cross_val_score(model, X_train, y_train, cv=kf,
   ↪ scoring='neg_mean_squared_error')
6 model_comp.append(-cross_val_scores)
7 print("criterion='squared_error'",np.mean(-cross_val_scores))
8 model = RandomForestRegressor(n_estimators=n, random_state=seed,
   ↪ criterion='absolute_error')
9 cross_val_scores = cross_val_score(model, X_train, y_train, cv=kf,
   ↪ scoring='neg_mean_squared_error')
10 print("criterion='absolute_error'",np.mean(-cross_val_scores))
11 model_comp.append(-cross_val_scores)
12 model = RandomForestRegressor(n_estimators=n, random_state=seed,
   ↪ criterion='friedman_mse')
13 cross_val_scores = cross_val_score(model, X_train, y_train, cv=kf,
   ↪ scoring='neg_mean_squared_error')
14 print("criterion='friedman_mse'",np.mean(-cross_val_scores))
15 model_comp.append(-cross_val_scores)

1 labels = ['Squared Error', 'Absolute Error', 'Friedman Mean Squared
   ↪ Error']
2
3 plt.boxplot(model_comp_filtered, labels = labels, widths=0.5)
4
5 plt.ylabel('Mean Squared Error')
6
7 plt.title('Criterion Method vs Mean Squared Error for Filtered
   ↪ Data')
8
9 plt.savefig('Random_Forest_filtered_output_2.png', dpi=600)
10
11 plt.show()

1 plt.boxplot(model_comp, widths=0.5, labels = labels)
2
3 plt.ylabel('Mean Squared Error')
4
5 plt.title('Criterion Method vs Mean Squared Error')
6
7 plt.savefig('Random_Forest_output_2.png', dpi=600)
8
9 plt.show()

```

K Nearest Neighbor Regressor

```
1 from sklearn.neighbors import KNeighborsRegressor
2
3 k_values = [3, 4, 5, 6, 7, 8]
4
5 mses = []
6 for k in k_values:
7     KNN = KNeighborsRegressor(n_neighbors=k, metric='manhattan')
8     mse = -cross_val_score(KNN, X_train, y_train, cv=kf,
9         scoring='neg_mean_squared_error')
10    mses.append(mse)
11 mses2 = [np.mean(mse) for mse in mses]
12
13 for i in range(len(k_values)):
14     print(k_values[i], mses2[i])
15
16 print("Min K value", k_values[np.argmin(mses2)])
17
```

```
1 plt.figure(figsize=(10, 6))
2 plt.boxplot(mses, labels = k_values)
3 plt.xlabel('K values')
4 plt.ylabel('Mean Squared Error')
5 plt.title('MSE Distribution for Different Values of k')
6 plt.savefig('Knn_output1.png')
7 plt.show()
8
```

```
1 k_values = [3, 4, 5, 6, 7, 8]
2
3 mses = []
4 for k in k_values:
5     KNN = KNeighborsRegressor(n_neighbors=k, metric='chebyshev')
6     mse = -cross_val_score(KNN, X_train, y_train, cv=kf,
7         scoring='neg_mean_squared_error')
8     mses.append(np.mean(mse))
9 mses2 = [np.mean(mse) for mse in mses]
10
11 for i in range(len(k_values)):
12     print(k_values[i], mses2[i])
13
14 print("Min K value", k_values[np.argmin(mses2)])
```

```

1 # Plot the box plots for MSEs
2 plt.figure(figsize=(10, 6))
3 plt.boxplot([mse_results[k] for k in k_values], labels=k_values)
4 plt.xlabel('K values')
5 plt.ylabel('Mean Squared Error')
6 plt.title('MSE Distribution for Different Values of k')
7 plt.savefig('Knn_output2.png', dpi = 600)
8 plt.show()
9

```

Decision Trees

```

1 from sklearn.tree import DecisionTreeRegressor

1 ##Decision Tree output 1
2 mse_comp = []
3 model =
4     ↪ DecisionTreeRegressor(random_state=seed,criterion='absolute_error'
5                             ,splitter='best')
6 cross_val_scores = cross_val_score(model, X_train
7                                     , y_train, cv=kf,
8                                     ↪ scoring='neg_mean_squared_error')
9 print("absolute_error",np.mean(-cross_val_scores))
10 mse_comp.append(-cross_val_scores)
11 model =
12     ↪ DecisionTreeRegressor(random_state=seed,criterion='squared_error'
13                             ,splitter='best')
14 cross_val_scores = cross_val_score(model, X_train
15                                     , y_train, cv=kf,
16                                     ↪ scoring='neg_mean_squared_error')
17 print("squared_error",np.mean(-cross_val_scores))
18 mse_comp.append(-cross_val_scores)
19 model =
20     ↪ DecisionTreeRegressor(random_state=seed,criterion='friedman_mse'
21                             ,splitter='best')
22 cross_val_scores = cross_val_score(model, X_train
23                                     , y_train, cv=kf,
24                                     ↪ scoring='neg_mean_squared_error')
25 print("friedman_mse",np.mean(-cross_val_scores))
26 mse_comp.append(-cross_val_scores)

```

```

1 labels = ['Squared Error', 'Absolute Error', 'Friedman Mean Squared
2           ↪ Error']
3 plt.boxplot(mse_comp, labels = labels)
4
5 plt.ylabel('Mean Squared Error')

```

```

6
7 plt.title('Criterion Method vs Mean Squared Error')
8
9 plt.savefig('Decision_Tree_output_1.png', dpi=600,
10    ↪ bbox_inches='tight')
11 plt.show()

```

```

1 ##Decision Tree output 2
2 mse_comp = []
3 model =
4    ↪ DecisionTreeRegressor(random_state=seed,criterion='absolute_error'
5 ,splitter='best')
6 cross_val_scores = cross_val_score(model, X_train, y_train, cv=kf,
7    ↪ scoring='neg_mean_squared_error')
8 print("best",np.mean(-cross_val_scores))
9 mse_comp.append(-cross_val_scores)
10 model =
11    ↪ DecisionTreeRegressor(random_state=seed,criterion='absolute_error'
12 ,splitter='random')
12 cross_val_scores = cross_val_score(model, X_train, y_train, cv=kf,
13    ↪ scoring='neg_mean_squared_error')
14 print("random",np.mean(-cross_val_scores))
15 mse_comp.append(-cross_val_scores)

```

Neural Network

```

1 X_train_nn = torch.tensor(X_train, dtype=torch.float32)
2 y_train_nn = torch.tensor(y_train, dtype=torch.float32)
3 train_data = TensorDataset(X_train_nn, y_train_nn)
4 train_data = DataLoader(dataset=train_data, batch_size=10,
5    ↪ shuffle=True)
6
7 class Neural_Network(nn.Module):
8     def __init__(self, dim):
9         nn.Module.__init__(self)
10        self.layer1 = nn.Linear(dim, 128)
11        self.layer2 = nn.Linear(128, 64)
12        self.layer3 = nn.Linear(64, 1)
13
14        def forward(self, x):
15            x = F.relu(self.layer1(x))
16            x = F.relu(self.layer2(x))
17            x = F.relu(self.layer3(x))
18            return x
19
20 model = Neural_Network(49)

```

```

20
21 criterion = nn.MSELoss()
22
23 Adam = optim.Adam(model.parameters(), lr=0.001)
24
25 epochs = 1000
26
27 for epoch in range(epochs):
28     for X_vals, y_vals in train_data:
29         Adam.zero_grad()
30         y_pred = model(X_vals)
31         loss = criterion(y_pred, y_vals)
32         loss.backward()
33         Adam.step()
34     print(loss)

```

3

```

1 rf = RandomForestRegressor(n_estimators=350, random_state=seed)
2
3 rf.fit(X_train_filtered, y_train_filtered)
4
5 y_pred = rf.predict(X_test_filtered)
6
7 mse = mean_squared_error(y_test_filtered, y_pred)
8
9 print("Mean Squared Error is",mse)

```

```

1 dt = DecisionTreeRegressor(random_state=seed)
2
3 dt.fit(X_train, y_train)
4
5 y_pred = dt.predict(X_test)
6
7 mse = mean_squared_error(y_test, y_pred)
8
9 print("Mean Squared Error is",mse)

```

```

1 knn = KNeighborsRegressor(n_neighbors=7, metric='manhattan')
2 knn.fit(X_train, y_train)
3
4 y_pred = knn.predict(X_test)
5
6 mse = mean_squared_error(y_test, y_pred)
7 print("Mean Squared Error is", mse)

```

```
1 X_test_nn = torch.tensor(X_test, dtype=torch.float32)
2 y_test_nn = torch.tensor(y_test, dtype=torch.float32)
3 test_data = TensorDataset(X_test_nn, y_test_nn)
4 test_data = DataLoader(dataset=test_data, batch_size=10,
5 → shuffle=False)
6 loss_count = 0
7 count = 0
8 with torch.no_grad():
9     for X_vals, y_vals in test_data:
10         y_pred = model(X_vals)
11         loss = criterion(y_pred, y_vals)
12         loss_count += loss.item()
13         count +=1
14 mse = loss_count / count
15 print("Mean Squared Error is",mse)
```