

Coursework 2 Solution

Team: *18019783* and *21049846*

Department of Computer Science

University College London

December 2021

Abstract

This coursework includes three parts. Each section refers to answer of one question in the coursework. Appendix A, B, C includes corresponding codes for each part. All codes are written in Python.

Contents

1	Part 1: Kernel perceptron (Handwritten Digit Classification)	1
1.1	Introduction	1
1.2	Theory	1
1.2.1	Perceptron	1
1.2.2	Two-class polynomial kernel perceptron	2
1.2.3	Generalisation to k -class polynomial kernel perceptron	2
1.2.4	Gaussian kernel perceptron	3
1.2.5	One-versus-one algorithm	3
1.3	Experimental methods and Results	3
1.3.1	Basic results	3
1.3.2	Cross-validation	4
1.3.3	Confusion matrix	5
1.3.4	Hardest-to-predict images	7
1.4	Discussion	7
2	Part 2: Spectral Clustering	8
2.1	Experiments	8
2.2	Questions	10
3	Part 3: Sparse learning	11
	Appendices	15
A	Part 1 code	15
B	Part 2 code	34
C	Part 3 code	37

Label	0	1	2	3	4	5	6	7	8	9
Occurrence	1553	1269	929	824	852	716	834	792	708	821

Table 1

1 Part 1: Kernel perceptron (Handwritten Digit Classification)

1.1 Introduction

In this part of the coursework, several multi-class classification methods are investigated and compared. The data set used is *zipcombo.dat*, which is also known as MNIST handwritten digit data set [1]. This data set includes in total 9298 handwritten digits from 0 to 9. Each sample contains a label and 16×16 pixel values as attributes. Examples of digits in the data set is shown in Fig.(1). The number of occurrence of each label is shown in Table (1).

We convert single label of each digit into one-hot encoded label vector of size 1×10 . If the label is 5 then the corresponding element in the label vector is +1 and the rests are -1. This method is also called "one-versus-rest" classification method. More details are discussed in Section 1.2.2.



Figure 1: Digit labels in MNIST data set.

1.2 Theory

1.2.1 Perceptron

Perceptron algorithm works for binary class classification problem. If we have input $\mathbf{x} \in \mathbb{R}^n$, where n is the number of attributes, and true label $t \in \{-1, 1\}$, then the prediction $y(\mathbf{x}) = f(\mathbf{w}^T \phi(\mathbf{x}))$, where $f(\cdot)$ is an activation function and $\phi(\mathbf{x})$ is a basis function. For the case where $f(a) = 1$ if $a \geq 0$ and $f(a) = -1$ elsewhere, we may make class decisions based on the value of $\mathbf{w}^T \phi(\mathbf{x})$. If the classifier makes a correct decision, $\mathbf{w}^T \phi_n t_n \geq 0$. The perceptron criterion is given by

$$E_p(\mathbf{w}) = - \sum_{n \in M} \mathbf{w}^T \phi(\mathbf{x}_n) t_n,$$

where $M = \{n : \mathbf{w}^T \phi_n t_n < 0\}$.

The weight vector \mathbf{w} may be updated using stochastic gradient descent (SGD) algorithm in

online learning problem. The new weight vector $\mathbf{w}^{(t+1)}$ is given by

$$\mathbf{w}^{(t+1)} = \mathbf{w}^{(t)} - \eta \nabla E_n(\mathbf{w}) = \mathbf{w}^{(t)} - \eta(\phi_n t_n),$$

where η is the learning rate, and $E_n(\mathbf{w}) = -\mathbf{w}^T \phi(\mathbf{x}_n) t_n$. If the data points in the data set is linearly separable, then perceptron SGD will converge until a optimal decision boundary is reached. In usual online perceptron, the update rule is

$$\mathbf{w}^{(t+1)} = \mathbf{w}^{(t)} - \phi_n t_n.$$

1.2.2 Two-class polynomial kernel perceptron

In this problem, basis function $\phi(\mathbf{x})$ is equal to \mathbf{x} , and $f(\cdot)$ is a sign function, $sgn(\cdot)$, that returns the sign of the value inside the brackets. Thus, if we initialise the weight vector by zeros, and by observation we may see that at some time t , the weight vector is given by

$$\mathbf{w}^{(t)} = \sum_{i=0}^{t-1} \alpha_i \mathbf{x}_i,$$

where in a single epoch, α_i is the true label value, and in many epochs it represents the number of times where some sample \mathbf{x}_i is mis-classified. The predicted label is then given by

$$\begin{aligned} \hat{y}_t &= sgn\left(\sum_{i=0}^{t-1} \alpha_i \mathbf{x}_i \cdot \mathbf{x}_t\right) \\ &= sgn\left(\sum_{i=0}^{t-1} \alpha_i K(\mathbf{x}_i, \mathbf{x}_t)\right), \end{aligned}$$

where $K(\mathbf{x}_i, \mathbf{x}_t)$ is the kernel function, which stores a subset of seen training samples. We can see that as written in kernels, the primal form is replaced by dual form of prediction. The kernel chosen is the inner product of two sample vectors, which makes the feature space to higher dimensions by using higher degree polynomial kernels. In this question we choose a polynomial kernel function of the form $K_d(\mathbf{p}, \mathbf{q}) = (\mathbf{p} \cdot \mathbf{q})^d$, where $d \in \{1, \dots, 7\}$ is the degree of the polynomial, and \mathbf{p}, \mathbf{q} are two sample vectors. For the case where $d = 1$, the complexity of computing the prediction is $O(m^3 + mn)$ for the dual form, and $O(n^3 + n)$ for the primal form, where m is the number of samples and n is the number of experts. As d increases, the complexity of making prediction in primal form increases exponentially, which reveals that it is more tractable to compute in dual form when in high-dimensional feature space. The kernel method also allows the classifier to learn a non-linear decision boundary.

1.2.3 Generalisation to k -class polynomial kernel perceptron

The method we implemented for multi-class classification is called "One-versus-rest" (OvR). This method uses binary classification algorithm to solve multi-class classification problem. For example, for samples have label 9, the problem becomes a binary classification, that is, 9 is one class and the rest is another class. For each example in MNIST dataset, its label is predicted by comparing the 'confidence' belonging to different labels, where the label with the greatest

confidence will belong to class (+1) and the rest will belong to class (-1). The confidence $C^{(k)}$ of some sample \mathbf{x}_t belongs to class k is calculated by

$$C^{(k)} = \sum_{i=1}^{t-1} \alpha_i^{(k)} K(\mathbf{x}_i, \mathbf{x}_t),$$

where $\alpha_i^{(k)}$ is the number of times a sample \mathbf{x}_i is mis-classified into class k , and $K(\mathbf{x}_i, \mathbf{x}_t)$ is the inner product of two sample vectors to the power of d . This kernel is a measure of similarity of two samples, the more similar the two samples, the larger their inner product. The value of d helps to amplify such similarity, thus to improve the corresponding confidence. In MNIST data set, $k \in \{0, 1, \dots, 9\}$, thus for each sample we may obtain a confidence vector, each element represents how confident this sample belongs to a corresponding label. In training process, all α s are initialised to zero, and $\alpha_i^{(k)}$ is increases by one while $\alpha_i^{(\hat{k})}$ is reduced by one if a sample \mathbf{x}_i has true label k but mis-predicted as \hat{k} . If the label is predicted correctly, the value of the corresponding α does not change.

1.2.4 Gaussian kernel perceptron

Gaussian kernel is generally expressed as

$$K(\mathbf{p}, \mathbf{q}) = \theta_0 \exp(-\theta_1 * \|\mathbf{p} - \mathbf{q}\|^2), \quad (1)$$

where \mathbf{p} and \mathbf{q} are two sample vectors. In this question, we choose the amplification factor, θ_0 , to be 1. θ_1 is called the kernel width, it scales the distance between \mathbf{p} and \mathbf{q} . If θ_1 is large, the exponential decays faster as two sample vectors get further apart. As θ_1 turns to infinity, the Gram matrix becomes an identity matrix. In order to preserve more features and do not overfit, we choose $\theta_1 \in \{0.005, 0.010, 0.015, 0.020, 0.025, 0.030, 0.035\}$ after some trial values of θ_1 from 0.001 to 4. Within this range, we makes the best c of each run distributed similarly to the distribution of the best d .

1.2.5 One-versus-one algorithm

One-versus-one algorithm is also a binary classification method. Unlike one-versus-rest method, it has $k(k-1)/2$ classifiers, where k is the number of classes. Denoting $f_{ij}(\mathbf{x})$ as a binary classifier to classify the label of some sample \mathbf{x} , where i and j are two different classes in k classes, and here i denotes positive class and j denotes a negative class. Then we have $f_{ij} = -f_{ji}$, and the label is predicted by

$$f(\mathbf{x}) = \arg \max_{i \in k} \left(\sum_{j \in k, j \neq i} f_{ij}(\mathbf{x}) \right),$$

which means the prediction is based on the maximum number of vote on label i .

1.3 Experimental methods and Results

1.3.1 Basic results

We perform 20 runs for $d = \{1, \dots, 7\}$ and $c = \{0.005, \dots, 0.035\}$, each run randomly splits *zipcombo.dat* into 80% train and 20% test. For each run, the training set undergoes 20 epochs.

degree of polynomial kernel	Training set error rate(%)	Test set error rate(%)
1	5.2635 ± 0.3069	8.6452 ± 1.0026
2	0.1876 ± 0.0865	3.6371 ± 0.6118
3	0.1096 ± 0.0683	3.2554 ± 0.4260
4	0.0652 ± 0.0448	3.1882 ± 0.4586
5	0.0578 ± 0.0276	3.1183 ± 0.4057
6	0.0437 ± 0.0289	3.1317 ± 0.3797
7	0.0450 ± 0.0309	3.0887 ± 0.4721

Table 2: One-versus-rest polynomial kernel perceptron algorithm.

Gaussian kernel width	Training set error rate(%)	Test set error rate(%)
0.005	0.1123 ± 0.0710	3.2688 ± 0.4537
0.010	0.0511 ± 0.0301	2.9731 ± 3.6133
0.015	0.0410 ± 0.0263	3.0054 ± 0.3900
0.020	0.0370 ± 0.0204	3.1559 ± 0.4244
0.025	0.0202 ± 0.0274	3.2043 ± 0.4721
0.030	0.0215 ± 0.0220	3.2070 ± 0.3858
0.035	0.0101 ± 0.0130	3.4220 ± 0.4816

Table 3: One-versus-rest Gaussian kernel perceptron algorithm.

The basic results of using polynomial kernel perceptron algorithm and Gaussian kernel perceptron algorithm are shown in Table 2 and Table 3 respectively. The result for one-versus-one polynomial kernel perceptron is shown in Table 4.

1.3.2 Cross-validation

We perform 20 runs, and for each run we split the dataset, *zipcombo.bat*, into 80% training set and 20% test set. For each training set we perform 5-fold cross-validation and run 20 epochs and select the best degree of polynomial kernel d^* , or the best kernel width for the Gaussian

degree of polynomial kernel	Training set error rate(%)	Test set error rate(%)
1	1.6671 ± 0.1976	5.9247 ± 0.6547
2	0.0578 ± 0.0231	3.3226 ± 0.4798
3	0.0417 ± 0.0230	3.1747 ± 0.5306
4	0.0350 ± 0.0202	3.2419 ± 0.3541
5	0.0329 ± 0.0182	3.2769 ± 0.3723
6	0.0343 ± 0.0177	3.2527 ± 0.4342
7	0.0316 ± 0.0159	3.4620 ± 0.5090

Table 4: One-versus-one polynomial kernel perceptron algorithm.

Best degree	5	5	4	5	3
Test set error rate(%)	3.2796	3.4409	3.5484	3.2258	2.6882
Best degree	4	4	4	4	7
Test set error rate(%)	3.4946	3.6022	2.6882	3.2258	2.8495
Best degree	5	3	5	4	4
Test set error rate(%)	3.1720	2.7419	3.3871	2.3656	2.7419
Best degree	3	4	7	4	5
Test set error rate(%)	3.0108	3.7634	2.6344	3.2796	2.6344

Table 5

Best width	0.02	0.02	0.025	0.02	0.02
Test set error rate(%)	3.0108	2.7957	3.4946	3.2258	2.8495
Best width	0.025	0.015	0.035	0.02	0.02
Test set error rate(%)	3.4946	3.1183	3.0108	2.9570	3.0108
Best width	0.03	0.02	0.02	0.015	0.02
Test set error rate(%)	3.6559	3.0645	3.0645	2.3118	3.1183
Best width	0.025	0.02	0.025	0.015	0.02
Test set error rate(%)	3.0108	4.2473	3.3871	3.1720	2.9570

Table 6

kernel c^* . Then we obtain 20 best degrees and widths, and use each of them to re-train on the corresponding full training set, and record the test error on the test set. The same procedure is applied to one-versus-one polynomial kernel perceptron as well.

The results of using one-versus-rest polynomial kernel perceptron algorithm and Gaussian kernel perceptron algorithm are shown in Table 5 and Table 6 respectively. The mean and standard deviation for d^* is 4.45 ± 1.10 , and for c^* is 0.0205 ± 0.0039 . Both of them are approximately the median of their range, while the mean c^* is a bit large. This reveals that the value range of c we cross-validated over is reasonable. The mean and standard deviation for the test set error rate is $(3.0887 \pm 0.3964) \%$ for the polynomial kernel and $(3.1479 \pm 0.3871) \%$ for the Gaussian kernel. The result for one-versus-one algorithm is shown in Table 7. The mean and standard deviation for the corresponding d^* is 3.75 ± 0.97 , and for the test set error rate is $(3.1290 \pm 0.4895) \%$.

1.3.3 Confusion matrix

The confusion matrix is of size (10×10) , each element represents the number of times some true label (matrix row index) is mis-classified as some false label (matrix column index). The final confusion matrix is normalized over 20 runs, each run with the corresponding training/test sets split and best degree of polynomial kernel. The result is divided into two parts and shown in Table 8 and Table 9.

Best degree	3	5	5	4	3
Test set error rate(%)	2.6344	3.1183	3.4946	3.4409	2.8495
Best degree	4	3	6	4	3
Test set error rate(%)	3.3871	3.4409	2.7419	2.9570	3.1720
Best degree	4	3	5	3	4
Test set error rate(%)	3.2796	2.6882	3.6022	2.1505	2.9570
Best degree	3	4	4	2	3
Test set error rate(%)	2.1505	4.0323	3.4946	3.6559	3.3333

Table 7

	0	1	2	3	4
0	0.0 ± 0.0	0.0 ± 0.0	0.0089 ± 0.0126	0.0041 ± 0.0074	0.0046 ± 0.0105
1	0.0 ± 0.0	0.0 ± 0.0	0.0 ± 0.0	0.0 ± 0.0	0.0203 ± 0.0282
2	0.0164 ± 0.0190	0.0016 ± 0.0049	0.0 ± 0.0	0.0207 ± 0.0203	0.0355 ± 0.0226
3	0.0073 ± 0.0133	0.0007 ± 0.0032	0.0224 ± 0.0185	0.0 ± 0.0	0.0033 ± 0.0086
4	0.0017 ± 0.0054	0.0170 ± 0.0190	0.0161 ± 0.0193	0.0025 ± 0.0061	0.0 ± 0.0
5	0.0234 ± 0.0188	0.0007 ± 0.0032	0.0060 ± 0.0098	0.0265 ± 0.0192	0.0158 ± 0.0178
6	0.0226 ± 0.0179	0.0067 ± 0.0111	0.0078 ± 0.0104	0.0 ± 0.0	0.0132 ± 0.0127
7	0.0010 ± 0.0046	0.0053 ± 0.0107	0.0138 ± 0.0138	0.0023 ± 0.0057	0.0156 ± 0.0154
8	0.0220 ± 0.0176	0.0122 ± 0.0132	0.0158 ± 0.0189	0.0410 ± 0.0315	0.0165 ± 0.0150
9	0.0078 ± 0.0120	0.0040 ± 0.0089	0.0027 ± 0.0067	0.0092 ± 0.0119	0.0432 ± 0.0280

Table 8: Confusion matrix part I.

	5	6	7	8	9
0	0.0168 ± 0.0220	0.0121 ± 0.0151	0.0008 ± 0.0034	0.0085 ± 0.0146	0.0049 ± 0.0077
1	0.0040 ± 0.0086	0.0073 ± 0.0124	0.0074 ± 0.0094	0.0008 ± 0.0037	0.0015 ± 0.0047
2	0.0093 ± 0.0129	0.0070 ± 0.0108	0.0279 ± 0.0220	0.0049 ± 0.0078	0.0027 ± 0.0065
3	0.0699 ± 0.0337	0.0 ± 0.0	0.0095 ± 0.0103	0.0250 ± 0.0206	0.0085 ± 0.0114
4	0.0042 ± 0.0075	0.0247 ± 0.0222	0.0067 ± 0.0112	0.0 ± 0.0	0.0234 ± 0.0170
5	0.0 ± 0.0	0.0174 ± 0.0193	0.0017 ± 0.0051	0.0189 ± 0.0203	0.0122 ± 0.0157
6	0.0097 ± 0.0104	0.0 ± 0.0	0.0 ± 0.0	0.0082 ± 0.0158	0.0010 ± 0.0044
7	0.0039 ± 0.0107	0.0 ± 0.0	0.0 ± 0.0	0.0064 ± 0.0113	0.0279 ± 0.0226
8	0.0266 ± 0.0238	0.0050 ± 0.0079	0.0076 ± 0.0114	0.0 ± 0.0	0.0081 ± 0.0108
9	0.0034 ± 0.0070	0.0 ± 0.0	0.0292 ± 0.0322	0.0065 ± 0.0151	0.0 ± 0.0

Table 9: Confusion matrix part II.

1.3.4 Hardest-to-predict images

There are several criterion we may choose to determine the hardest-predict samples. For instance, we may perform 20 runs and find the confidence of each sample in different labels and sum them, then compare the confidence of the final predicted label of each sample to find the smallest five. We may also compare the difference between the two largest confidences of each sample to determine the hardest-to-predict samples. In this question, we run 20 times, for each run we split the data into 80% training set and 20% test set. We perform 5-fold cross-validation on the training set and obtain the best degree of polynomial kernel for the corresponding split. We then train the α matrix on the 80% training set using the best degree. Then we perform the test on the entire data set and record mis-classified samples. By repeating these steps we obtain the hardest-to-predict samples as shown in Fig.(2).

It is not surprising these five samples are the hardest to predict. They are 'distinctive' hand-written digits, and even the human eyes can hardly distinguish the correct labels.

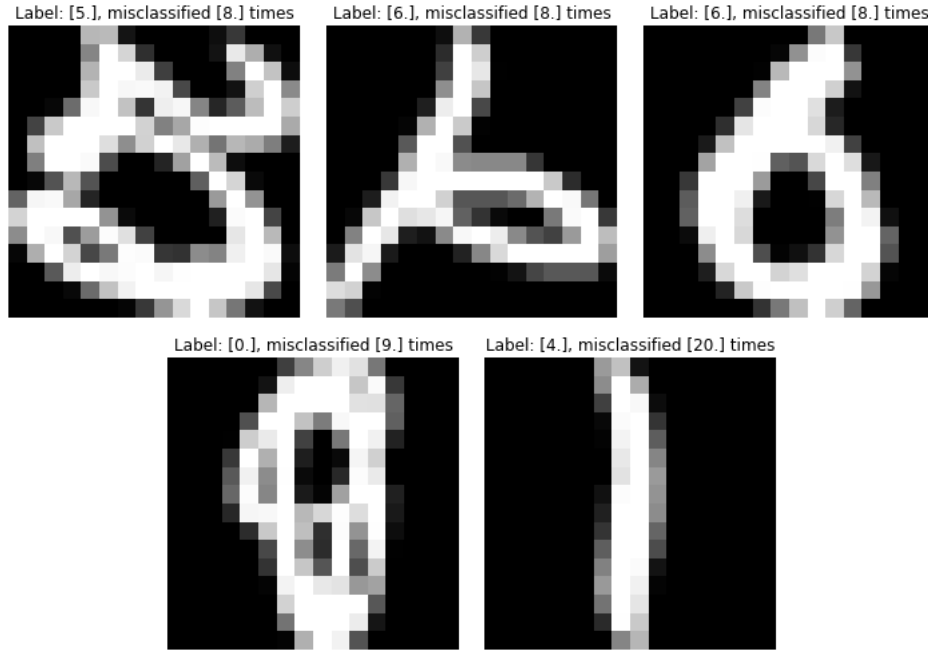


Figure 2: The five hardest-to-predict images in the data set.

1.4 Discussion

In cross validation, the parameter θ_0 , as shown in Eq.(1) in Section 1.2.4 is not cross-validated over. As stated earlier, θ_0 is an amplification factor. By cross-validating over this parameter, the difference between classes is enlarged, which may help the classifier to distinguish different labels more easily. The same idea may also apply on the polynomial kernel, where an amplification factor may amplify the similarity between two samples. However, as in the polynomial kernel we take the inner product, multiplying by the same value has no effect on the classifier's ability

to distinguish different labels. A more general expression of polynomial kernel is

$$K(\mathbf{p}, \mathbf{q}) = (\mathbf{p} \cdot \mathbf{q} + \theta)^d,$$

where θ is an offset term which shifts the similarity measure by a certain amount, but it may have very limited effect on the discrimination ability of the classifier.

Although one-versus-one(OvO) algorithm requires $O(k^2)$ classifiers while one-versus-rest(OvR) requires $O(k)$ classifiers, each classifier in OvO is smaller. This makes OvO algorithm seem more efficient. If the time taken for building each classifier is approximately the same, then the efficiency of OvR may exceed that of OvO. OvR algorithm can be very effective if regularised least-square(RLS) classification algorithm applied. In RLS method, OvR algorithm can be computed very effectively by adapting matrix factorization [2]. On the other hand, OvO on support vector machine (SVM) algorithm may have faster training and obtain more accurate result if we have a large-scale data set and large number of classes [3].

The results shown in Table 2 and 3 indicate that the Gaussian kernel has relatively lower training set and test set error rate compared with the polynomial kernel. However, this depends on the chosen degree or kernel width. Different choices may result in different classifier performance. The same idea also applies to the test error obtained by the best degree or width after cross-validation.

2 Part 2: Spectral Clustering

2.1 Experiments

1. When $c = 19.698$, which is $2^{4.3}$, the two-moons data is clustered correctly with an accuracy of 100%.

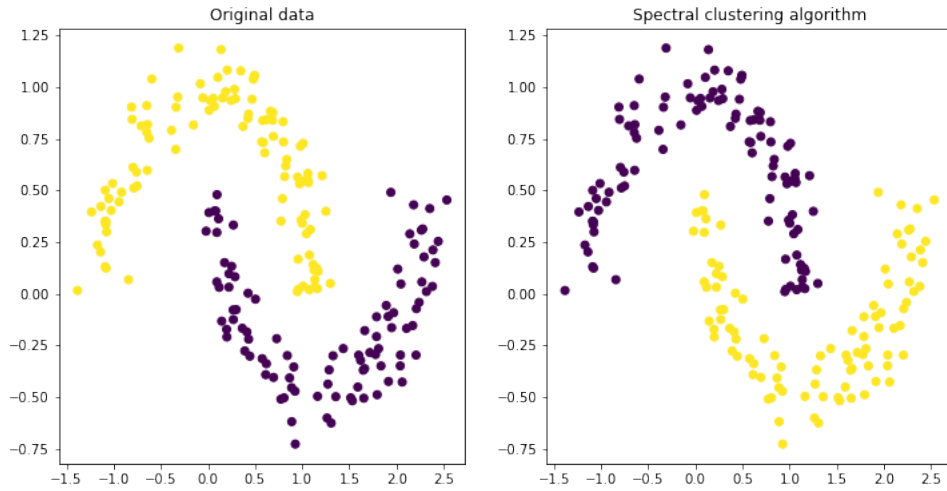


Figure 3: Two moons dataset

2. Since the gaussian data is randomly generated, the best c changes with every different generated data. For the data in this figure, when $c = 235.0$, the data is clustered correctly with an accuracy of 100%.

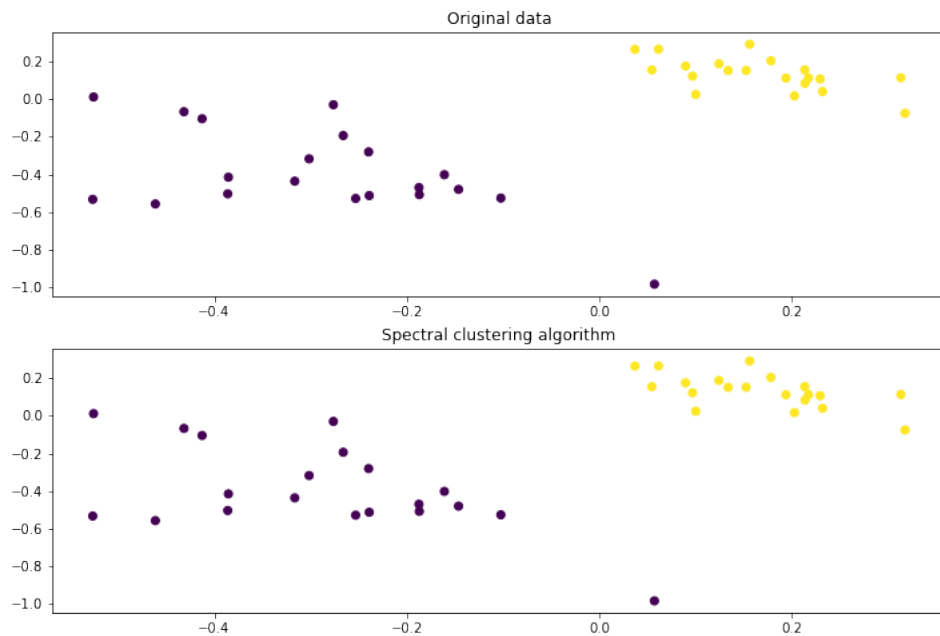


Figure 4: Gaussian dataset

3. For the data 'dtrain123', the best c is 0.01, and this c gives a correct cluster percentage (CP) of 0.91053.

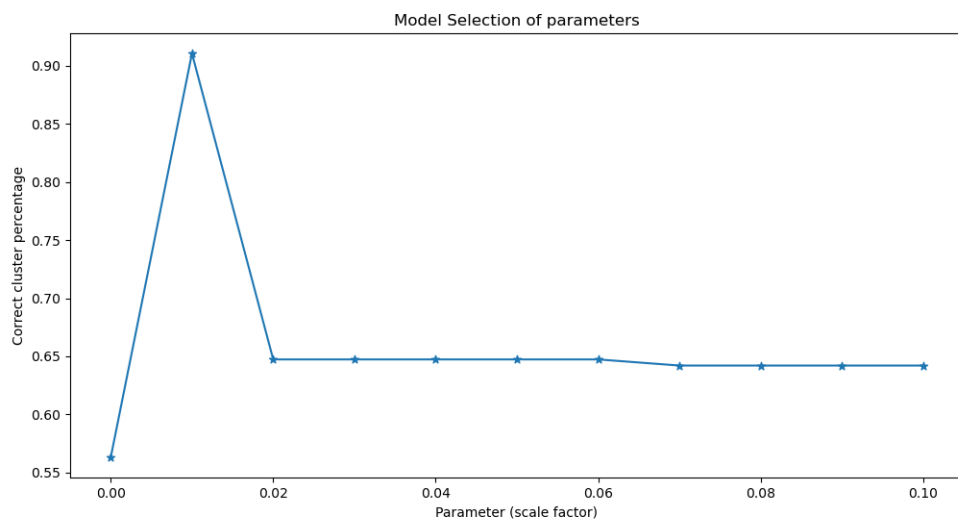


Figure 5: dtrain123 dataset

2.2 Questions

1. Explain why $CP(c)$ is a reasonable measure of cluster correctness.

Since we know the true labels of the data, a reasonable measure of cluster correctness is the percentage of the number of correctly predicted labels. Since we only have two clusters, if we have more wrongly predicted labels than correctly predicted labels, this means we have labelled the clusters in opposite labels of the true labels, which is also correct since we only need to cluster the data into two sections. So we need to calculate the percentage of the maximum of the number of correctly predicted labels and the number of wrongly predicted labels.

2. Explain why the first eigenvalue of the Laplacian is zero and the corresponding eigenvector is the constant vector.

Since $L = \text{diag}(\text{sum}(W)) - W$, the sum of all rows of L is 0. Then the rows are not linearly independent. Hence, it is not full rank. Hence, the nullity is not zero. Then there exist $Lv = 0 = 0v$. Again, the sum of all rows of L is 0, if we have $v = (1, \dots, 1)$, the elements of Lv is actually the sum of all rows of L , which is 0. Hence, $Lv = 0 = 0v$, L has eigenvalue 0 with corresponding eigenvector $(1, \dots, 1)$. Since Laplacian matrix is positive semi-definite, all eigenvalue of L is greater or equal to 0. So L has the smallest eigenvalue (the first eigenvalue) as 0 and the corresponding eigenvector is the constant vector $(1, \dots, 1)$.

3. In your own words (please explicitly cite any reference), provide a tentative explanation why spectral clustering “works” (1 paragraph only).

From [4], we can see that the spectral clustering we are doing here is actually the construction of a partition A_1, A_2 of the graph. We solve this by solving the optimization problem: $\min_{A \subset V} \text{RatioCut}(A, \bar{A})$. Then we can rewrite this using the graph's Laplacian matrix as [5]

$$f'Lf = |V| \cdot \text{RatioCut}(A, \bar{A})$$

So the optimization problem becomes $\min_{A \subset V} f'Lf$, where the vector f is orthogonal to the constant one vector. Then by the Rayleigh-Ritz theorem, the solution of f is the eigenvector corresponding to the second smallest eigenvalue of L . Now, to obtain the partition, we need to transform the value of f_i by a discrete indicator function. The simplest and easiest choice is the sign function. So we use $\text{sign}(f_i)$ as the classification label. In this way, we are able to classify data into two parts with label $(-1, +1)$.

4. Recall the parameter c in the definition (see equation (1)) of the weight matrix. Explain why and how c influences the quality of the clustering.

Here we use the Gaussian kernel as weight matrix, so c is defined as $\frac{1}{2\sigma^2}$. If c is small, the variance is large and the differences among those distance would be small and this weight matrix won't

reflect the relationships of those points well(the edges of the graph). If c is big, the variance is small and c works as an amplifier to amplify the influence of distance. However, for too large c , the influence of distance would be over amplified and give less accurate results.

3 Part 3: Sparse learning

(a) The graph below contains four plots of the sample complexity of the perceptron, winnow, least squares, and 1-nearest neighbours algorithms under dimension n . All those four plots have clear trends: linear relationship for perceptron, logarithmic relationship for winnow, linear relationship for least squares, exponential relationship for 1-nearest neighbours.

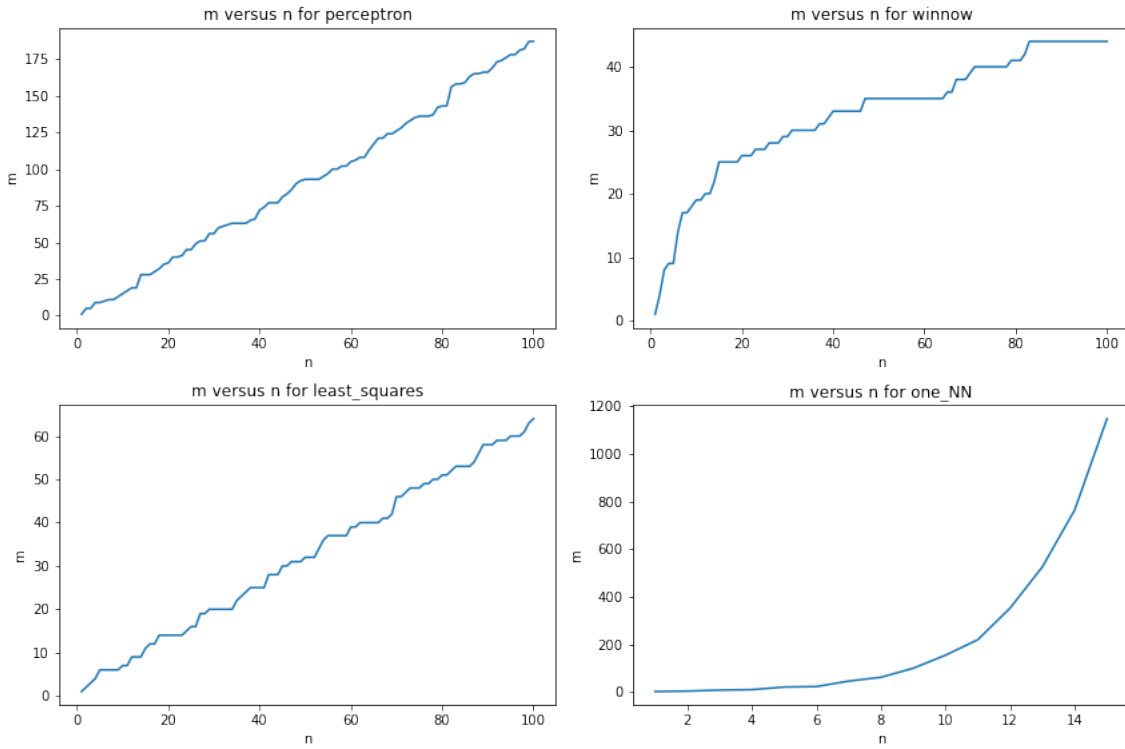


Figure 6: Estimated number of samples (m) to obtain 10% generalisation error versus dimension (n) for all four algorithms.

(b) i) My algorithm:

1. Run the algorithms for $1 \leq m \leq 10000$ and every n such that $1 \leq n \leq \max_n$, where \max_n is 15 for 1-nn and 100 for other three algorithms.
2. Choose the initial value of m for the i -th n as the final value of m for the $(i - 1)$ -th n . For $n = 1$, we run the algorithm from $m = 1$.
3. Generate m training samples and 5000 test samples with n features.
4. Run the algorithm for 10 times and calculate the mean generalisation error. If it is no more than 10%, move on to next n . Otherwise, try $m + 1$ and repeat step 3 and 4.

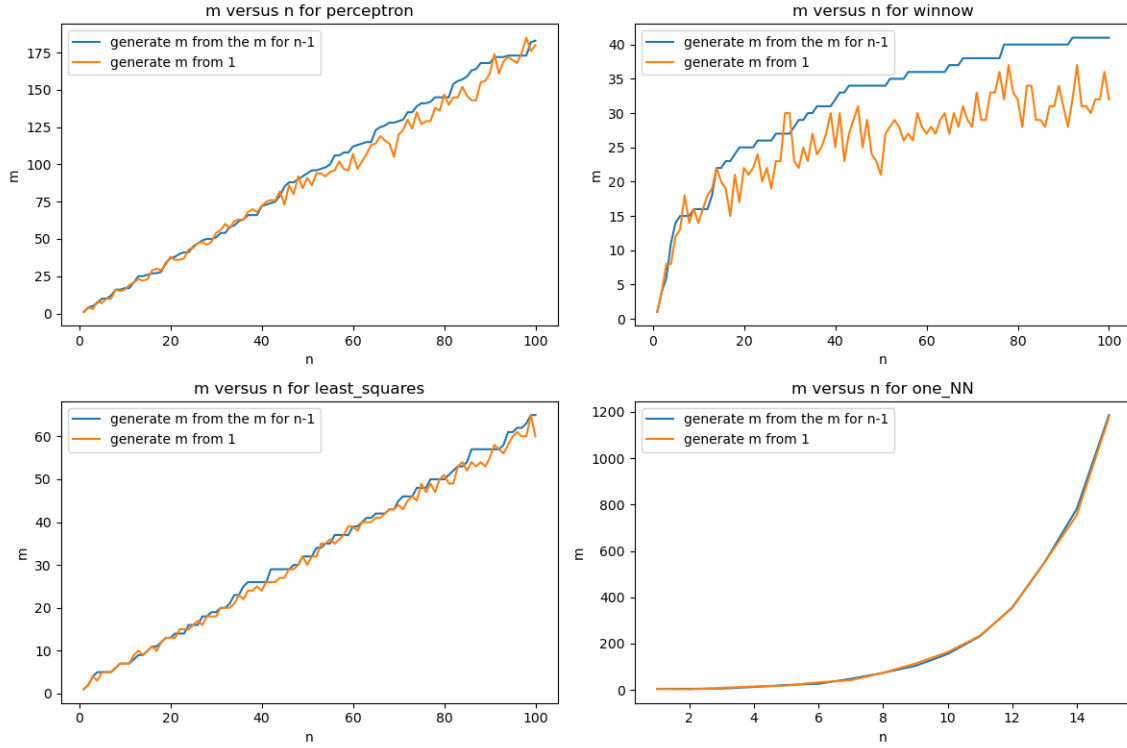


Figure 7: The comparison between my algorithm (blue line) and the normal algorithm (orange line).

ii) Trade offs and biases:

1. Sample bias: Since the test size does not cover all the combinations that can be generated, we may have biased data that will result in higher or lower generalisation error.

2. Algorithm bias: Since my algorithm assumes m always increases as n increases, if we have biased data that produces higher generalisation error, and results in a larger m for n , the final m for $n + 1$ will likely be higher than the value in our normal algorithm which tests all m starting from 1. And this is why we can observe some flat line segments on the graph of my algorithm.

3. My algorithm trade-off accuracy and computation time: Since not all m has been tested for each n and there could have smaller m that fits the criteria, my accuracy is relatively lower. But the differences are not too big for our chosen n , since the intuition is m will grow as n grows. And computation time is largely reduced especially for 1-nn. For example, in the 1-nn graph, m is around 800 and 1200 for $n = 14$ and 15. In my method, we only try m from 800 until we find 1200 and save the time that we needed for $m = 1, \dots, 800$. Also we will have longer computation time if we have more runs or larger test size or max_n . We know larger max_n can give us a clearer trend since we can see the relationship for more n , while bigger test size and more runs can give us a more stable trend since bias is reduced. I can't afford large max_n for 1-nn, so I chose relatively large test size and more runs to make the trend stable enough to be revealed. However, as they are not large enough, my accuracy will be lower than running through all combinations.

(c) As we stated in part(a), graphs of perceptron and least squares show linear relationships, graph of winnow shows a logarithm relationship, graph of 1-nn shows an exponential relationship. So we fit the data as those functions to estimate how m grows as a function of n and this tells us which of $O(\cdot)$, $\Omega(\cdot)$, $\Theta(\cdot)$ the algorithms belong to. And here is the graph of m versus n for all four algorithms and their fitted lines with the corresponding function form, where a, b in the graph are the parameters of the lines.

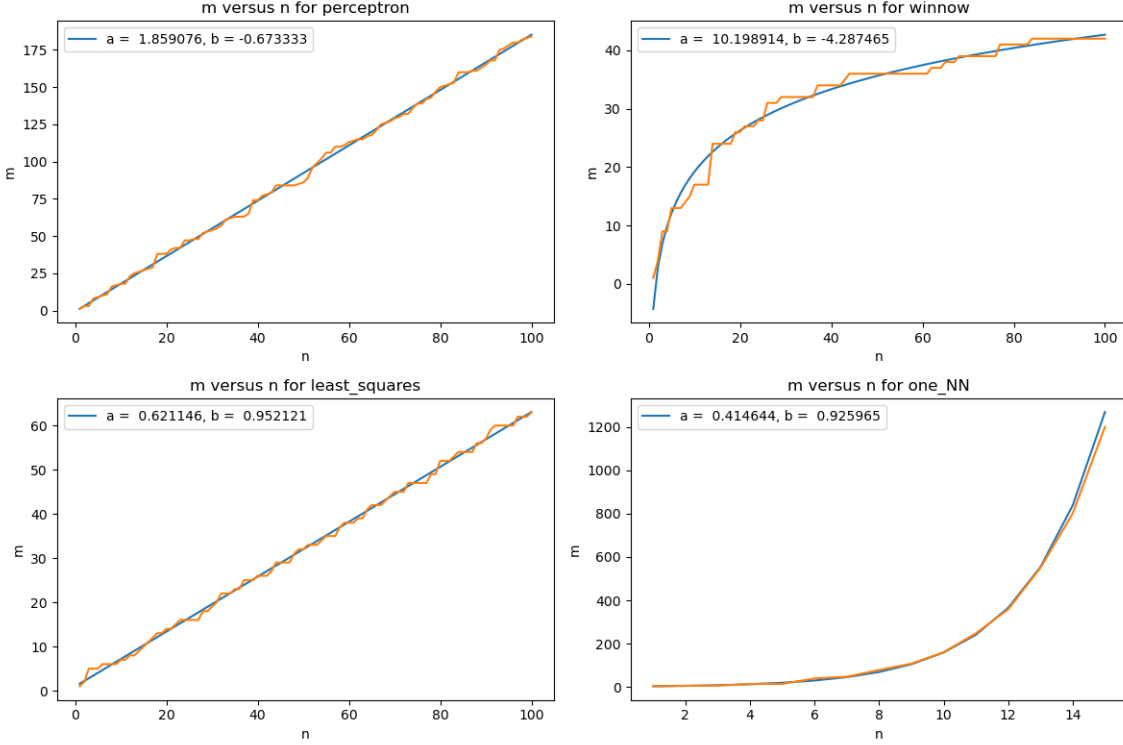


Figure 8: The plot of m and n (orange line) and the fitted line (blue line).

Here are the results we get from the graph and the actual computation time that each algorithm used when generating m and n in the graph:

Algorithm	function form	fitted function	$O(\cdot)$, $\Omega(\cdot)$, $\Theta(\cdot)$	time(seconds)
perceptron	$a * n + b$	$1.86n - 0.67$	$\Theta(n)$	8
winnow	$a * \log(n) + b$	$10.20 * \log(n) - 4.29$	$\Theta(\log(n))$	3
least squares	$a * n + b$	$0.62 * n + 0.95$	$\Theta(n)$	5
1-nn	$exp(b) * exp(a)^n$	$2.53 * 1.51^n$	$\Theta(1.51^n)$	2434

From the computation time, we can see while 1-nn requires a lot more time, other three algorithms only need seconds. For computational complexity (m), winnow requires the least m which is about 45 when $n = 100$ and the growth of m will be even slower as it has a logarithm relationship. Least squares are the second-best and perceptron are the third-best. They both have linear relationships but least squares requires smaller m for the same n . The most demanding one is again 1-nn, which not only requires high m , also the growth of m will be faster and faster.

It is easy to understand as adding 1 to n means adding 1 feature to each point of the data, and when calculating the distance for every pair of test and train points, we need to calculate the distance in one more dimension.

(d) Since s and the examples are sampled uniformly at random, the probability of making a mistake on example (x_s, y_s) is $\frac{M}{m}$, where M is the mistake bound for algorithm perceptron. By the Novikoff Theorem for Perceptron Bound, We have

$$M \leq \left(\frac{R}{\gamma}\right)^2$$

where $R := \max_t \|x_t\|$ and $\gamma \leq (\mathbf{v} \cdot \mathbf{x}_t)y_t$ and \mathbf{v} is a vector with $\|\mathbf{v}\| = 1$.

Since all entries of data are sampled from -1 and $+1$, the maximized γ is 1 and $R = \max_t \|x_t\| = \sqrt{n}$. As a result, the tightest bound we can have for M is

$$M \leq \left(\frac{\sqrt{n}}{1}\right)^2 = \frac{n}{1} = n$$

$$\text{So } \hat{p}_{m,n} = \frac{M}{m} = \frac{n}{m}$$

References

- [1] Y. Lecun, L. Bottou, Y. Bengio, and P. Haffner, “Gradient-based learning applied to document recognition,” *Proceedings of the IEEE*, vol. 86, no. 11, p. 2278–2324, 1998.
- [2] L. Rosasco, “Lecture 8- regularized least squares classification,” Spring 2014.
- [3] J. Milgram, M. Cheriet, and R. Sabourin, *Tenth International Workshop on Frontiers in Handwriting Recognition*.
- [4] U. Von Luxburg, “A tutorial on spectral clustering,” *Statistics and computing*, vol. 17, no. 4, pp. 395–416, 2007.
- [5] Wikipedia, “Eigenvalues and eigenvectors — Wikipedia, the free encyclopedia.” <http://en.wikipedia.org/w/index.php?title=Eigenvalues%20and%20eigenvectors&oldid=1057253625>, 2021. [Online; accessed 07-December-2021].

Appendices

A Part 1 code

```
1 # -*- coding: utf-8 -*-
2 """
3 Created on Thu Dec  2 03:45:39 2021
4
5 @author: 21049846
6
7 0078 SL CW2
8 """
9 import numpy as np
10 from matplotlib import pyplot as plt
11 from sklearn.model_selection import train_test_split
12 import seaborn as sns
13 import math
14 import statistics as stat
15 import time
16 from scipy.spatial.distance import cdist
17 from itertools import combinations
18
19 dataset = np.loadtxt('zipcombo.dat')
20 t = dataset[:,0].reshape((-1,1))
21 X = dataset[:,1:]
22
23 class_dict = {}
24 for i in range(int(min(t)),int(max(t)+1)):
25     class_dict[str(i)] = np.count_nonzero(t==i)
26 num_classes = len(class_dict)
27
28 def shuffle(train, label, rs):
29     """
30     split zipcombo.dat into training and test set.
31     """
32     full_data_list = []
33     full_data_list.append(train_test_split(train, label, test_size=0.2,
34     random_state=rs))
35     return full_data_list[0]
36
37 def polynomial_gram_matrix(X1, X2, d):
38     """
39     :input:
40         X1: input matrix.
41         X2: input matrix.
42         d: degree of polynomial kernel function.
43     :output:
44         K: gram matrix
45     """
46     K = np.matmul(X1,X2.T)**d
```



```

46     return K
47
48 def gaussian_gram_matrix(X1, X2, c):
49     """
50     :input:
51         X1: input matrix.
52         X2: input matrix.
53         c: Gaussian kernel width.
54     :output:
55         K: gram matrix
56     """
57     K = np.exp(-c * cdist(X1,X2)**2) # problem
58
59     return K
60
61 def test_kernel_perceptron(train_set, label, d):
62     """
63     :input:
64         train_set: training set matrix.
65         label: label vector.
66         d: degree of polynomial kernel function.
67     :output:
68         mean and standard deviation of training set error and test set error for
69         20 runs.
70     """
71     run = 20
72     epoches = 20
73     k = 10
74     train_error = []
75     confusion = []
76     test_error = []
77     confusion_test = []
78
79     for r in range(run):
80         start = time.time()
81
82         data = shuffle(train_set, label, r)
83
84         m = data[0].shape[0]
85         alpha = np.zeros((k,m))
86         K = polynomial_gram_matrix(data[0], data[0], d)
87         confusion_matrix = np.zeros((k,k))
88         for e in range(epoches):
89             num_false = 0
90             for t in range(m):
91                 true_label = data[2][t]
92                 pred_y_value = np.dot(alpha, K[:,t].reshape((m,-1)))
93                 confidence = pred_y_value
94                 pred_label = np.argmax(pred_y_value)
95                 if pred_label != true_label:
96                     confusion_matrix[int(data[2][t]), int(pred_label)] += 1

```

```

96         alpha[int(pred_label),t] -= 1
97         alpha[int(true_label),t] += 1
98         num_false += 1
99         train_error_rate = num_false / m
100
101     train_error.append(train_error_rate)
102
103     confusion.append(confusion_matrix)
104     end = time.time()
105
106     time_taken = end - start
107
108     print("current run: {}, current train error rate: {}, time taken: {}".
109           format(r+1, train_error[-1], time_taken))
110
111     start_test = time.time()
112     m_test = data[1].shape[0]
113     K_test = polynomial_gram_matrix(data[0],data[1],d)
114     confusion_matrix_test = np.zeros((k,k))
115     num_false_test = 0
116
117     for t in range(m_test):
118         true_label = data[3][t]
119         pred_y_value_test = np.dot(alpha, K_test[:,t].reshape((m,-1)))
120         confidence_test = pred_y_value_test
121         pred_label = np.argmax(confidence_test)
122         if pred_label != true_label:
123             num_false_test += 1
124             confusion_matrix_test[int(true_label), int(pred_label)] += 1
125     test_error_rate = num_false_test / m_test
126
127     test_error.append(test_error_rate)
128     confusion_test.append(confusion_matrix_test)
129     end_test = time.time()
130
131     time_taken_test = end_test - start_test
132     print("current run: {}, current test error rate: {}, time taken: {}".
133           format(r+1, test_error[-1], time_taken_test))
134
135     mean_train_error = stat.mean(train_error)
136     std_train_error = stat.stdev(train_error)
137     avg_confusion_mat = confusion[-1] / sum(confusion)
138
139     mean_test_error = stat.mean(test_error)
140     std_test_error = stat.stdev(test_error)
141     avg_confusion_matrix_test = confusion_test[-1] / sum(confusion_test)
142
143     return [mean_train_error, std_train_error], [mean_test_error, std_test_error],
144           [avg_confusion_mat, avg_confusion_matrix_test]
145
146 def k_fold_cross_val(k, data_set):

```

```

144     """
145     :input:
146         k: number of folds
147         data_set: dataset to be split into k folds
148     :output:
149         partitioned_set: a list contains k sub-lists, each sub-list is one fold
of the dataset
150     """
151     if len(data_set) % k == 0:
152         partitioned_set = np.reshape(data_set, (k, int(len(data_set)/k), -1))
153     else:
154         new_data_set = data_set.tolist()
155         new_length = len(data_set) - len(data_set) % k
156         left = [new_data_set[i] for i in range(new_length, len(data_set))]
157         partitioned_set = [new_data_set[i] for i in range(new_length)]
158         partitioned_set = np.reshape(partitioned_set, (k, int(new_length / k),
-1))
159         partitioned_set = partitioned_set.tolist()
160         for i in range(len(left)):
161             partitioned_set[i].append(left[i])
162
163     return partitioned_set
164
165 def cv_kernel_perceptron(train_set, label):
166     """
167     :input:
168         train_set: training set matrix.
169         label: label vector.
170     :output:
171         a list contains best degree of polynomial kernel and a list of
corresponding test set error.
172     """
173     run = 20
174     epoches = 20
175     k = 10
176     best_d = []
177     test_error_best_d = []
178
179     for r in range(run):
180         start = time.time()
181         data = shuffle(train_set, label, r)
182
183         cv_data_train = k_fold_cross_val(5, data[0])
184         cv_data_label = k_fold_cross_val(5, data[2])
185
186         test_error_cv = []
187         confusion_test_cv = []
188         for d in range(1, 8):
189             test_error = []
190             confusion_test = []
191             for i in range(5):

```

```

192         cv_test = np.array(cv_data_train[i])
193         cv_test_label = np.array(cv_data_label[i])
194         cv_train = []
195         cv_train_label = []
196         for j in range(5):
197             if j != i:
198                 cv_train = cv_train + cv_data_train[j]
199                 cv_train_label = cv_train_label + cv_data_label[j]
200         cv_train = np.array(cv_train)
201         cv_train_label = np.array(cv_train_label)
202
203         cv_m = cv_train.shape[0]
204         cv_alpha = np.zeros((k,cv_m))
205         cv_K = polynomial_gram_matrix(cv_train,cv_train,d)
206
207         for e in range(epochs):
208             for t in range(cv_m):
209                 true_label = cv_train_label[t]
210                 pred_y_value = np.dot(cv_alpha,cv_K[:,t].reshape((cv_m
, -1)))
211
212                 confidence = pred_y_value
213                 pred_label = np.argmax(confidence)
214
215                 if pred_label != true_label:
216                     cv_alpha[int(pred_label),t] -= 1
217                     cv_alpha[int(true_label),t] += 1
218
219             end = time.time()
220
221             time_taken = end - start
222
223             print("current run: {}, time taken: {}".format(r+1, time_taken))
224
225             start_test = time.time()
226             m_test_cv = cv_test.shape[0]
227             K_test_cv = polynomial_gram_matrix(cv_train,cv_test,d)
228             confusion_matrix_test_cv = np.zeros((k,k))
229             num_false_test = 0
230
231             for t in range(m_test_cv):
232                 true_label = cv_test_label[t]
233                 pred_y_value_test = np.dot(cv_alpha, K_test_cv[:,t].reshape
((cv_m, -1)))
234
235                 confidence_test = pred_y_value_test
236                 pred_label = np.argmax(confidence_test)
237                 if pred_label != true_label:
238                     num_false_test += 1
239                     confusion_matrix_test_cv[int(true_label), int(pred_label
)] += 1
240
241             test_error_rate = num_false_test / m_test_cv
242             test_error.append(test_error_rate)

```

```

240         confusion_test.append(confusion_matrix_test_cv)
241         end_test = time.time()
242         time_taken_test = end_test - start_test
243         print("current run: {}, d: {}, current test error rate: {}, time
taken: {}".format(r+1, d, test_error[-1], time_taken_test))
244
245         mean_test_error_cv = stat.mean(test_error)
246         test_error_cv.append(mean_test_error_cv)
247         avg_confusion_matrix_test = confusion_test[-1] / sum(confusion_test)
248         confusion_test_cv.append(avg_confusion_matrix_test)
249
250         d_star = np.argmin(test_error_cv) + 1
251         best_d.append(d_star)
252         start = time.time()
253         data = shuffle(train_set, label, r)
254         m = data[0].shape[0]
255         alpha = np.zeros((k, m))
256         K = polynomial_gram_matrix(data[0], data[0], d_star)
257         for e in range(epochs):
258             for t in range(m):
259                 true_label = data[2][t]
260                 pred_y_value = np.dot(alpha, K[:, t].reshape((m, -1)))
261                 confidence = pred_y_value
262                 pred_label = np.argmax(pred_y_value)
263
264                 if pred_label != true_label:
265                     alpha[int(pred_label), t] -= 1
266                     alpha[int(true_label), t] += 1
267
268             end = time.time()
269             time_taken = end - start
270             print("current run: {}, best d: {}, time taken: {}".format(r+1, d_star,
time_taken))
271
272             start_test = time.time()
273             m_test = data[1].shape[0]
274             K_test = polynomial_gram_matrix(data[0], data[1], d_star)
275             confusion_matrix_test = np.zeros((k, k))
276             num_false_test = 0
277
278             for t in range(m_test):
279                 true_label = data[3][t]
280                 pred_y_value_test = np.dot(alpha, K_test[:, t].reshape((m, -1)))
281                 confidence_test = pred_y_value_test
282                 pred_label = np.argmax(confidence_test)
283                 if pred_label != true_label:
284                     num_false_test += 1
285                     confusion_matrix_test[int(true_label), int(pred_label)] += 1
286             test_error_rate = num_false_test / m_test
287
288             test_error_best_d.append(test_error_rate)

```

```

289         confusion_test.append(confusion_matrix_test)
290         end_test = time.time()
291
292         time_taken_test = end_test - start_test
293         print("current run: {}, current test error rate: {}, time taken: {}".
format(r+1, test_error_best_d[-1], time_taken_test))
294
295     return best_d, test_error_best_d
296
297 def confusion_kernel_perceptron(train_set, label, d, r):
298     """
299     :input:
300         train_set: training set matrix.
301         label: label vector.
302         d: degree of polynomial kernel function.
303         r: random state, corresponding to a specific split of data set.
304     :output:
305         a confusion matrix on the test set using polynomial kernel perceptron.
306     """
307     epoches = 20
308     k = 10
309     train_error = []
310     test_error = []
311     confusion_test = []
312     start = time.time()
313     data = shuffle(train_set, label, r)
314     m = data[0].shape[0]
315     alpha = np.zeros((k, m))
316     K = polynomial_gram_matrix(data[0], data[0], d)
317     for e in range(epoches):
318         num_false = 0
319         for t in range(m):
320             true_label = data[2][t]
321             pred_y_value = np.dot(alpha, K[:, t].reshape((m, -1)))
322             confidence = pred_y_value
323             pred_label = np.argmax(confidence)
324             #print("test")
325             if pred_label != true_label:
326                 alpha[int(pred_label), t] -= 1
327                 alpha[int(true_label), t] += 1
328                 num_false += 1
329
330         train_error_rate = num_false / m
331
332     train_error.append(train_error_rate)
333     end = time.time()
334     time_taken = end - start
335     print("train error rate: {}, time taken: {}".format(train_error[-1],
time_taken))
336
337     start_test = time.time()

```

```

338 m_test = data[1].shape[0]
339 K_test = polynomial_gram_matrix(data[0], data[1], d)
340 confusion_matrix_test = np.zeros((k, k))
341 num_false_test = 0
342
343 for t in range(m_test):
344     true_label = data[3][t]
345     pred_y_value_test = np.dot(alpha, K_test[:, t].reshape((m, -1)))
346     confidence_test = pred_y_value_test
347     pred_label = np.argmax(confidence_test)
348     if pred_label != true_label:
349         num_false_test += 1
350         confusion_matrix_test[int(true_label), int(pred_label)] += 1
351 test_error_rate = num_false_test / m_test
352 test_error.append(test_error_rate)
353 end_test = time.time()
354 time_taken_test = end_test - start_test
355 print("current test error rate: {}, time taken: {}".format(test_error[-1],
356 time_taken_test))
357
358 return confusion_matrix_test
359
360 def test_gaussian_perceptron(train_set, label, c):
361     """
362     :input:
363         train_set: training set matrix.
364         label: label vector.
365         c: Gaussian kernel width.
366     :output:
367         mean and standard deviation of training set error and test set error for
368         20 runs.
369     """
370     run = 20
371     epoches = 20
372     k = 10
373     train_error = []
374     confusion = []
375     test_error = []
376     confusion_test = []
377
378     for r in range(run):
379         start = time.time()
380         data = shuffle(train_set, label, r)
381         m = data[0].shape[0]
382         alpha = np.zeros((k, m))
383         K = gaussian_gram_matrix(data[0], data[0], c)
384         confusion_matrix = np.zeros((k, k))
385         for e in range(epoches):
386             num_false = 0
387             for t in range(m):
388                 true_label = data[2][t]

```

```

387         pred_y_value = np.dot(alpha, K[:,t].reshape((m,-1)))
388         confidence = pred_y_value
389         pred_label = np.argmax(confidence)
390
391         if pred_label != true_label:
392             confusion_matrix[int(data[2][t]),int(pred_label)] += 1
393             alpha[int(pred_label),t] -= 1
394             alpha[int(true_label),t] += 1
395             num_false += 1
396
397         train_error_rate = num_false / m
398
399         train_error.append(train_error_rate)
400         confusion.append(confusion_matrix)
401         end = time.time()
402         time_taken = end - start
403         print("current run: {}, current train error rate: {}, time taken: {}".
404               format(r+1, train_error[-1], time_taken))
405         #print("confidence vector is {}".format(confidence))
406
407         start_test = time.time()
408         m_test = data[1].shape[0]
409         K_test = gaussian_gram_matrix(data[0], data[1], c)
410         confusion_matrix_test = np.zeros((k,k))
411         num_false_test = 0
412
413         for t in range(m_test):
414             true_label = data[3][t]
415             pred_y_value_test = np.dot(alpha, K_test[:,t].reshape((m,-1)))
416             confidence_test = pred_y_value_test
417             pred_label = np.argmax(confidence_test)
418             if pred_label != true_label:
419                 num_false_test += 1
420                 confusion_matrix_test[int(true_label), int(pred_label)] += 1
421             test_error_rate = num_false_test / m_test
422             test_error.append(test_error_rate)
423             confusion_test.append(confusion_matrix_test)
424             end_test = time.time()
425             time_taken_test = end_test - start_test
426             print("current run: {}, current test error rate: {}, time taken: {}".
427                   format(r+1, test_error[-1], time_taken_test))
428
429         mean_train_error = stat.mean(train_error)
430         std_train_error = stat.stdev(train_error)
431         mean_test_error = stat.mean(test_error)
432         std_test_error = stat.stdev(test_error)
433
434         return [mean_train_error, std_train_error], [mean_test_error, std_test_error]
435
436 def cv_gaussian_perceptron(train_set, label):
437     """

```



```

436     :input:
437         train_set: training set matrix.
438         label: label vector.
439     :output:
440         a list contains best Gaussian kernel width and a list of
441     corresponding test set error.
442     """
443     run = 20
444     epoches = 20
445     k = 10
446     best_c = []
447     test_error_best_c = []
448     c = [0.005,0.010,0.015,0.020,0.025,0.030,0.035]
449
450     for r in range(run):
451
452         data = shuffle(train_set,label,r)
453         cv_data_train = k_fold_cross_val(5,data[0])
454         cv_data_label = k_fold_cross_val(5,data[2])
455         test_error_cv = []
456         confusion_test_cv = []
457         for d in range(0,7):
458             test_error = []
459             confusion_test = []
460             for i in range(5):
461                 cv_test = np.array(cv_data_train[i])
462                 cv_test_label = np.array(cv_data_label[i])
463                 cv_train = []
464                 cv_train_label = []
465                 for j in range(5):
466                     if j != i:
467                         cv_train = cv_train + cv_data_train[j]
468                         cv_train_label = cv_train_label + cv_data_label[j]
469                 cv_train = np.array(cv_train)
470                 cv_train_label = np.array(cv_train_label)
471                 cv_m = cv_train.shape[0]
472                 cv_alpha = np.zeros((k,cv_m))
473                 cv_K = gaussian_gram_matrix(cv_train,cv_train,c[d])
474                 start = time.time()
475
476                 for e in range(epoches):
477                     for t in range(cv_m):
478                         true_label = cv_train_label[t]
479                         pred_y_value = np.dot(cv_alpha,cv_K[:,t].reshape((cv_m
480 , -1)))
481
482                         confidence = pred_y_value
483                         pred_label = np.argmax(confidence)
484
485                         if pred_label != true_label:
486                             cv_alpha[int(pred_label),t] -= 1
487                             cv_alpha[int(true_label),t] += 1

```

```

486
487         end = time.time()
488         time_taken = end - start
489         print("current run: {}, time taken: {}".format(r+1, time_taken))
490         start_test = time.time()
491         m_test_cv = cv_test.shape[0]
492         K_test_cv = gaussian_gram_matrix(cv_train, cv_test, c[d])
493         confusion_matrix_test_cv = np.zeros((k,k))
494         num_false_test = 0
495
496         for t in range(m_test_cv):
497             true_label = cv_test_label[t]
498             pred_y_value_test = np.dot(cv_alpha, K_test_cv[:,t].reshape
((cv_m,-1)))
499             confidence_test = pred_y_value_test
500             pred_label = np.argmax(confidence_test)
501             if pred_label != true_label:
502                 num_false_test += 1
503                 confusion_matrix_test_cv[int(true_label), int(pred_label
)] += 1
504             test_error_rate = num_false_test / m_test_cv
505             test_error.append(test_error_rate)
506             confusion_test.append(confusion_matrix_test_cv)
507             end_test = time.time()
508             time_taken_test = end_test - start_test
509             print("current run: {}, c: {}, current test error rate: {}, time
taken: {}".format(r+1, c[d], test_error[-1], time_taken_test))
510
511             mean_test_error_cv = stat.mean(test_error)
512             test_error_cv.append(mean_test_error_cv)
513             avg_confusion_matrix_test = confusion_test[-1] / sum(confusion_test)
514             confusion_test_cv.append(avg_confusion_matrix_test)
515
516             d_star = np.argmin(test_error_cv) + 1
517             c_star = c[d_star]
518             best_c.append(c_star)
519             start = time.time()
520             data = shuffle(train_set, label, r)
521             m = data[0].shape[0]
522             alpha = np.zeros((k,m))
523             K = gaussian_gram_matrix(data[0], data[0], c_star)
524
525             for e in range(epochs):
526                 for t in range(m):
527                     true_label = data[2][t]
528                     pred_y_value = np.dot(alpha, K[:,t].reshape((m,-1)))
529                     confidence = pred_y_value
530                     pred_label = np.argmax(pred_y_value)
531                     if pred_label != true_label:
532                         alpha[int(pred_label),t] -= 1
533                         alpha[int(true_label),t] += 1

```

```

534
535     end = time.time()
536     time_taken = end - start
537     print("current run: {}, best c: {}, time taken: {}".format(r+1, c_star,
time_taken))
538     start_test = time.time()
539     m_test = data[1].shape[0]
540     K_test = gaussian_gram_matrix(data[0], data[1], c_star)
541     confusion_matrix_test = np.zeros((k,k))
542     num_false_test = 0
543
544     for t in range(m_test):
545         true_label = data[3][t]
546         pred_y_value_test = np.dot(alpha, K_test[:,t].reshape((m,-1)))
547         confidence_test = pred_y_value_test
548         pred_label = np.argmax(confidence_test)
549         if pred_label != true_label:
550             num_false_test += 1
551             confusion_matrix_test[int(true_label), int(pred_label)] += 1
552     test_error_rate = num_false_test / m_test
553
554     test_error_best_c.append(test_error_rate)
555     confusion_test.append(confusion_matrix_test)
556     end_test = time.time()
557     time_taken_test = end_test - start_test
558     print("current run: {}, current test error rate: {}, time taken: {}".
format(r+1, test_error_best_c[-1], time_taken_test))
559
560     return best_c, test_error_best_c
561
562 def hardest_to_predict(train_set, label, d, r):
563     """
564     :input:
565         train_set: training set matrix.
566         label: label vector.
567         d: degree of polynomial kernel function.
568         r: random state , corresponding to a specific split of data set.
569     :output:
570         a list contains the number of times each sample is mis-classified.
571     """
572     epoches = 20
573     k = 10
574     start = time.time()
575     data = shuffle(train_set, label, r)
576     m = data[0].shape[0]
577     alpha = np.zeros((k,m))
578     K = polynomial_gram_matrix(data[0], data[0], d)
579
580     for e in range(epoches):
581
582         for t in range(m):

```

```

583         true_label = data[2][t]
584         pred_y_value = np.dot(alpha, K[:,t].reshape((m,-1)))
585         confidence = pred_y_value
586         pred_label = np.argmax(confidence)
587
588         if pred_label != true_label:
589             alpha[int(pred_label),t] -= 1
590             alpha[int(true_label),t] += 1
591
592     end = time.time()
593     time_taken = end - start
594     print("current r: {}, d: {}, time taken: {}".format(r+1, d, time_taken))
595     start_test = time.time()
596     m_whole = train_set.shape[0]
597     misprediction = np.zeros((m_whole,1))
598     K_test = polynomial_gram_matrix(data[0], train_set, d)
599
600     for t in range(m_whole):
601         true_label = label[t]
602         pred_y_value_test = np.dot(alpha, K_test[:,t].reshape((m,-1)))
603         confidence_test = pred_y_value_test
604         pred_label = np.argmax(confidence_test)
605         if pred_label != true_label:
606             misprediction[t] += 1
607
608     end_test = time.time()
609     time_taken_test = end_test - start_test
610     print("current r: {}, time taken: {}".format(r+1, time_taken_test))
611
612     return misprediction
613
614 def OvO_kernel_test(train_set, label, d):
615     """
616     :input:
617         train_set: training set matrix.
618         label: label vector.
619         d: degree of polynomial kernel function.
620     :output:
621         mean and standard deviation of training set error and test set error for
622         20 runs.
623     """
624     run = 20
625     epoches = 20
626     k = 10
627     classifier = list(combinations(range(k), 2))
628     # Hashmap
629     alpha_converter = dict(zip(tuple(classifier), tuple(range(len(classifier)))))
630
631     train_error = []
632     test_error = []

```

```

632     for r in range(run):
633         start = time.time()
634         data = shuffle(train_set, label, r)
635         m = data[0].shape[0]
636         alpha = np.zeros((int(k*(k-1)/2), m))
637         K = polynomial_gram_matrix(data[0], data[0], d)
638
639         for e in range(epoches):
640             num_false = 0
641             for t in range(m):
642                 true_label = data[2][t]
643                 pred_y_value = np.dot(alpha, K[:, t].reshape((m, -1)))
644                 confidence = pred_y_value
645                 binary_classifier = np.sign(confidence).clip(0, None)
646                 classifiers = np.zeros((len(classifier)))
647                 for c_index, b_index in enumerate(binary_classifier):
648                     classifiers[c_index] = classifier[int(c_index)][int(b_index)]
649
650             pred_label = np.bincount(classifiers.astype(int)).argmax()
651             if pred_label != true_label:
652                 num_false += 1
653
654             for loc, (a, b) in enumerate(alpha_converter.keys()):
655                 pred_y = classifiers[loc]
656
657                 if true_label == a and pred_y != true_label:
658                     alpha[loc, t] -= 1
659                 if true_label == b and pred_y != true_label:
660                     alpha[loc, t] += 1
661
662             train_error_rate = num_false / m
663
664             train_error.append(train_error_rate)
665             end = time.time()
666             time_taken = end - start
667             print("current run: {}, current train error rate: {}, time taken: {}".
668                   format(r+1, train_error[-1], time_taken))
669             start_test = time.time()
670             m_test = data[1].shape[0]
671             K_test = polynomial_gram_matrix(data[0], data[1], d)
672             num_false_test = 0
673
674             for t in range(m_test):
675                 true_label = data[3][t]
676                 pred_y_value_test = np.dot(alpha, K_test[:, t].reshape((m, -1)))
677                 confidence_test = pred_y_value_test
678                 binary_classifier_test = np.sign(confidence_test).clip(0, None)
679                 classifiers_test = np.zeros((len(classifier)))
680                 for c_index, b_index in enumerate(binary_classifier_test):

```

```

    )]

681
682         pred_label_test = np.bincount(classifiers_test.astype(int)).argmax()
683
684         if pred_label_test != true_label:
685             num_false_test += 1
686             test_error_rate = num_false_test / m_test
687             test_error.append(test_error_rate)
688             end_test = time.time()
689             time_taken_test = end_test - start_test
690             print("current run: {}, current test error rate: {}, time taken: {}".
format(r+1, test_error[-1], time_taken_test))
691
692         mean_train_error = stat.mean(train_error)
693         std_train_error = stat.stdev(train_error)
694         mean_test_error = stat.mean(test_error)
695         std_test_error = stat.stdev(test_error)
696
697         return [mean_train_error, std_train_error], [mean_test_error, std_test_error]
698
699 def cv_0v0_perceptron(train_set, label):
700     """
701     :input:
702         train_set: training set matrix.
703         label: label vector.
704     :output:
705         a list contains best Gaussian kernel width and a list of corresponding
test set error.
706     """
707     run = 20
708     epoches = 20
709     k = 10
710     classifier = list(combinations(range(k), 2))
711     alpha_converter = dict(zip(tuple(classifier), tuple(range(len(classifier)))))
712
713     best_d = []
714     test_error_best_d = []
715
716     for r in range(run):
717         start = time.time()
718         data = shuffle(train_set, label, r)
719         cv_data_train = k_fold_cross_val(5, data[0])
720         cv_data_label = k_fold_cross_val(5, data[2])
721         test_error_cv = []
722         for d in range(1, 8):
723             test_error = []
724             for i in range(5):
725                 cv_test = np.array(cv_data_train[i])
726                 cv_test_label = np.array(cv_data_label[i])
727                 cv_train = []

```

```

728         for j in range(5):
729             if j != i:
730                 cv_train = cv_train + cv_data_train[j]
731                 cv_train_label = cv_train_label + cv_data_label[j]
732         cv_train = np.array(cv_train)
733         cv_train_label = np.array(cv_train_label)
734
735         cv_m = cv_train.shape[0]
736         cv_alpha = np.zeros((int(k*(k-1)/2), cv_m))
737         cv_K = polynomial_gram_matrix(cv_train, cv_train, d)
738
739         for e in range(epochs):
740             for t in range(cv_m):
741                 true_label = cv_train_label[t]
742                 pred_y_value = np.dot(cv_alpha, cv_K[:,t]).reshape((cv_m
, -1)))
743                 confidence = pred_y_value
744                 binary_classifier = np.sign(confidence).clip(0, None)
745
746                 classifiers = np.zeros((len(classifier)))
747                 for c_index, b_index in enumerate(binary_classifier):
748                     classifiers[c_index] = classifier[int(c_index)][int(
b_index)]
749
750                 pred_label = np.bincount(classifiers.astype(int)).argmax
()
751
752                 for loc, (a,b) in enumerate(alpha_converter.keys()):
753                     pred_y = classifiers[loc]
754
755                     if true_label == a and pred_y != true_label:
756                         cv_alpha[loc,t] -= 1
757                     if true_label == b and pred_y != true_label:
758                         cv_alpha[loc,t] += 1
759
760             end = time.time()
761             time_taken = end - start
762             print("current run: {}, time taken: {}".format(r+1, time_taken))
763             start_test = time.time()
764             m_test_cv = cv_test.shape[0]
765             K_test_cv = polynomial_gram_matrix(cv_train, cv_test, d)
766             num_false_test = 0
767
768             for t in range(m_test_cv):
769                 true_label = cv_test_label[t]
770                 pred_y_value_test = np.dot(cv_alpha, K_test_cv[:,t]).reshape
((cv_m, -1)))
771                 confidence_test = pred_y_value_test
772                 binary_classifier_test = np.sign(confidence_test).clip(0,
None)
773                 classifiers_test = np.zeros((len(classifier)))

```

```

774         for c_index, b_index in enumerate(binary_classifier_test):
775             classifiers_test[c_index] = classifier[int(c_index)][int
(b_index)]
776
777         pred_label_test = np.bincount(classifiers_test.astype(int)).
argmax()
778
779         if pred_label_test != true_label:
780             num_false_test += 1
781             test_error_rate = num_false_test / m_test_cv
782             test_error.append(test_error_rate)
783             end_test = time.time()
784
785             time_taken_test = end_test - start_test
786             print("current run: {}, d: {}, current test error rate: {}, time
taken: {}".format(r+1, d, test_error[-1], time_taken_test))
787
788             mean_test_error_cv = stat.mean(test_error)
789             test_error_cv.append(mean_test_error_cv)
790
791             d_star = np.argmin(test_error_cv) + 1
792             best_d.append(d_star)
793             start = time.time()
794             data = shuffle(train_set, label, r)
795             m = data[0].shape[0]
796             alpha = np.zeros((int(k*(k-1)/2), m))
797             K = polynomial_gram_matrix(data[0], data[0], d_star)
798             for e in range(epochs):
799                 num_false = 0
800                 for t in range(m):
801                     true_label = data[2][t]
802                     pred_y_value = np.dot(alpha, K[:, t].reshape((m, -1)))
803                     confidence = pred_y_value
804                     binary_classifier = np.sign(confidence).clip(0, None)
805
806                     classifiers = np.zeros((len(classifier)))
807                     for c_index, b_index in enumerate(binary_classifier):
808                         classifiers[c_index] = classifier[int(c_index)][int(b_index)
]
809
810                     pred_label = np.bincount(classifiers.astype(int)).argmax()
811                     if pred_label != true_label:
812                         num_false += 1
813
814                     for loc, (a, b) in enumerate(alpha_converter.keys()):
815                         pred_y = classifiers[loc]
816
817                         if true_label == a and pred_y != true_label:
818                             alpha[loc, t] -= 1
819                         if true_label == b and pred_y != true_label:
820                             alpha[loc, t] += 1

```



```

821         end = time.time()
822         time_taken = end - start
823         print("current run: {}, best d: {}, time taken: {}".format(r+1, d_star,
time_taken))
824         start_test = time.time()
825         m_test = data[1].shape[0]
826         K_test = polynomial_gram_matrix(data[0], data[1], d_star)
827         num_false_test = 0
828
829         for t in range(m_test):
830             true_label = data[3][t]
831             pred_y_value_test = np.dot(alpha, K_test[:,t].reshape((m,-1)))
832             confidence_test = pred_y_value_test
833             binary_classifier_test = np.sign(confidence_test).clip(0, None)
834             classifiers_test = np.zeros((len(classifier)))
835             for c_index, b_index in enumerate(binary_classifier_test):
836                 classifiers_test[c_index] = classifier[int(c_index)][int(b_index
)]
837             pred_label_test = np.bincount(classifiers_test.astype(int)).argmax()
838
839             if pred_label_test != true_label:
840                 num_false_test += 1
841
842             test_error_rate = num_false_test / m_test
843             test_error_best_d.append(test_error_rate)
844             end_test = time.time()
845             time_taken_test = end_test - start_test
846             print("current run: {}, current test error rate: {}, time taken: {}".
format(r+1, test_error_best_d[-1], time_taken_test))
847
848         return best_d, test_error_best_d
849
850 if __name__ == "__main__":
851     print("start")
852     print(class_dict)
853     fig, axes = plt.subplots(2,5)
854     axes = axes.flatten()
855     for i in range(num_classes):
856         index = np.where(t == i)
857         image = np.reshape(X[index], (16, 16))
858         axes[i].imshow(image, cmap='gray')
859         axes[i].axis('off')
860     fig.subplots_adjust(wspace=0, hspace=0)
861     plt.tight_layout()
862     plt.show()
863
864     for i in range(1,8):
865         print("test_kernel_perceptron(X,t,{}".format(i))
866         test_kernel_perceptron(X,t,i)
867
868     c = [0.005, 0.01, 0.015, 0.02, 0.025, 0.03, 0.035]

```

```

869     for j in range(8):
870         print("test_gaussian_perceptron(X,t,{})".format(c[j]))
871         test_gaussian_perceptron(X,t,c[j])
872
873
874     print("cv_kernel_perceptron(X,t)")
875     best_d, test_error_d = cv_kernel_perceptron(X,t)
876     best_c, test_error_c = cv_gaussian_perceptron(X,t)
877
878     mean_d, std_d = stat.mean(best_d), stat.stdev(best_d)
879     mean_error_d, std_error_d = stat.mean(test_error_d), stat.stdev(test_error_d
880 )
881     mean_c, std_c = stat.mean(best_c), stat.stdev(best_c)
882     mean_error_c, std_error_c = stat.mean(test_error_c), stat.stdev(test_error_c
883 )
884
885     # best_d = [5,5,4,5,3,4,4,4,4,7,5,3,5,4,4,3,4,7,4,5]
886     # best_c = [0.02,0.02,
887 0.025,0.02,0.02,0.025,0.015,0.015,0.02,0.02,0.03,0.02,0.02,0.015,0.02,0.025,0.02,
888 0.025,0.015,0.02]
889
890     confusion_matrix_list = []
891     misclassified = []
892     for i in range(len(best_d)):
893         misclassified.append(hardest_to_predict(X,t,best_d[i],i))
894         cmt = confusion_kernel_perceptron(X,t,best_d[i],i)
895         norm_cmt = cmt / np.sum(cmt)
896         confusion_matrix_list.append(norm_cmt)
897
898
899     mean_cm = np.array(sum(confusion_matrix_list)) / len(confusion_matrix_list)
900     std_cm_list = []
901     for i in range(10):
902         for j in range(10):
903             element_list = []
904             for r in range(len(best_d)):
905                 element_list.append(confusion_matrix_list[r][i,j])
906             std = stat.stdev(element_list)
907             std_cm_list.append(std)
908     std_cm = np.array(std_cm_list).reshape((10,10))
909
910     Hp_index = np.argmax(sum(misclassified).reshape((len(sum(
911 misclassified)),),-5)[-5:]
912
913     for i in range(len(Hp_index)):
914         index = int(Hp_index[i])
915
916         pics = np.reshape(X[index],(16,16))
917
918         plt.imshow(pics, cmap='gray')

```

```

915     plt.axis('off')
916     plt.title("Label: {}, misclassified {} times".format(t[index], sum(
misclassified)[index]))
917     plt.show()
918
919     for i in range(1,8):
920         print("0v0_kernel_test(X,t,{}).format(i))
921         0v0_kernel_test(X,y,i)
922
923     ovo_d. ovo_test_error = cv_0v0_perceptron(X,t)
924     mean_ovo_d, std_ovo_d = stat.mean(ovo_d), stat.stdev(ovo_d)
925     mean_ovo_error, std_ovo_error = stat.mean(ovo_test_error), stat.stdev(
ovo_test_error)

```

Listing 1: Part I

B Part 2 code

```

1  # import libraries
2  import numpy as np
3  from matplotlib import pyplot as plt
4
5  # Spectral Clustering Algorithm
6  def weight(X, c):
7      # Calculate Adjacency weight matrix W.
8      l,n = X.shape
9      W = np.ones((l,l))
10     for i in range(l):
11         for j in range(l):
12             W[i,j] = np.exp(-c * (np.linalg.norm(X[i]-X[j]))**2)
13     return W
14
15  def laplacian(W):
16     # Calculate Graph Laplacian L.
17     l,l = W.shape
18     D = np.zeros((l,l))
19     for i in range(l):
20         D[i,i] = np.sum(W[i])
21     L = D - W
22     return L
23
24  def sec_eig(L):
25     # Find the second smallest eigenvector of L.
26     val, vec = np.linalg.eig(L)
27     a = np.argsort(val)[1]
28     return vec[:,a]
29
30  def sign(x):
31     # Return the sign of x.
32     if x >= 0:

```

```

33     y = 1
34     else:
35         y = -1
36     return y
37
38 def spectral_cluster(X, c):
39     # Use data X and c to give weight W, laplacian L,
40     # cluster vector v2, and use them to do the clustering.
41     W = weight(X, c)
42     L = laplacian(W)
43     v2 = sec_eig(L)
44     y = np.ones(v2.shape)
45     for i in range(len(v2)):
46         y[i] = sign(v2[i])
47     return y
48
49 def best_cluster(X, Y, c):
50     # Implement spectral_cluster function with different c
51     # to determine the best c which gives the best cluster.
52     # The return value is the best c, corresponding correctness and prediction.
53     l,n = X.shape
54     final_y = []
55     final_c = 0
56     final_cor = 0
57     for i in c:
58         y = spectral_cluster(X, i)
59         cor = max((y == Y).sum(), (-y == Y).sum())
60         if final_cor < cor:
61             final_cor = cor
62             final_c = i
63             final_y = y
64     final_cor_rate = final_cor / l
65
66     return final_c, final_cor_rate, final_y

```

Listing 2: Spectral Clustering Algorithm

```

1 # Load data
2 twomoons = np.loadtxt('C:/Users/james007/Desktop/twomoons.dat.txt')
3 X_moons = twomoons[:, 1:]
4 Y_moons = twomoons[:, 0]
5
6 # Give a sequence of c and find the best cluster of data
7 c1 = np.linspace(0, 100, 101)[1:]
8 moons_c, moons_cor_rate, moons_y = best_cluster(X_moons, Y_moons, c1)
9 print(f'The best c is {moons_c} and it gives a correctness of {moons_cor_rate:.2%}')
10
11 # Plot of the comparision
12 plt.figure(figsize=(12,6))
13
14 # Plot of original data

```

```

15 plt.subplot(121)
16 plt.scatter(X_moons[:, 0], X_moons[:, 1], c = Y_moons)
17 plt.title('Original data')
18
19 # Plot of spectral clustering
20 plt.subplot(122)
21 plt.scatter(X_moons[:, 0], X_moons[:, 1], c = moons_y)
22 plt.title('Spectral clustering algorithm')
23
24 plt.savefig('P2E1')

```

Listing 3: Experiment 1

```

1 I = np.eye(2)
2
3 # Generate random class '-1'
4 x_1 = np.random.multivariate_normal((-0.3,-0.3), 0.04*I, 20)
5 y_1 = -1 * np.ones(20)
6
7 # Generate random class '+1'
8 x_2 = np.random.multivariate_normal((0.15,0.15), 0.01*I, 20)
9 y_2 = np.ones(20)
10
11 # Combine two classes to give the data and label
12 X_random = np.array(list(x_1) + list(x_2))
13 Y_random = np.array(list(y_1) + list(y_2))
14
15 # Give a sequence of c and find the best cluster of data
16 c2 = np.linspace(0, 2000, 2001)[1:]
17 random_c, random_cor_rate, random_y = best_cluster(X_random, Y_random, c2)
18 print(f'The best c is {random_c} and it gives a correctness of{random_cor_rate:
    .2%}')
19
20 # Plot of the comparision
21 plt.figure(figsize=(12,8))
22
23 # Plot of original data
24 plt.subplot(211)
25 plt.scatter(X_random[:, 0], X_random[:, 1], c = Y_random)
26 plt.title('Original data')
27
28 # Plot of spectral clustering
29 plt.subplot(212)
30 plt.scatter(X_random[:, 0], X_random[:, 1], c = random_y)
31 plt.title('Spectral clustering algorithm')
32
33 plt.savefig('P2E2')

```

Listing 4: Experiment 2

```

1 # Load data
2 dtrain123 = np.loadtxt('C:/Users/james007/Desktop/dtrain123.dat.txt')

```

```

3
4 # Define dtrain13 to only contain data with label (1,3)
5 dtrain13 = []
6 for i in range(dtrain123.shape[0]):
7     if dtrain123[i, 0] == 1 or dtrain123[i, 0] == 3:
8         dtrain13.append(dtrain123[i])
9 dtrain13 = np.array(dtrain13)
10
11 X_dt = dtrain13[:, 1:]
12 Y_dt = dtrain13[:, 0]
13
14 #Relabel the data, change label (1,3) to label (1,-1)
15 classy = np.ones(Y_dt.shape)
16 for i in range(len(Y_dt)):
17     classy[i] = 2 * (Y_dt[i] == 1) -1
18
19 # Give a sequence of c and find the best cluster of data
20 # Give the result of best c and best CP
21 c3 = np.linspace(0, 0.1, 11)
22 CP = np.ones(len(c3))
23 for i in range(len(c3)):
24     y = spectral_cluster(X_dt, c3[i])
25     cor = max((y == classy).sum(), (-y == classy).sum())
26     CP[i] = cor / len(classy)
27 best_CP = np.max(CP)
28 best_c = c3[np.argmax(CP)]
29 print(f'The best c is {best_c} and the corresponding CP is{best_CP: .5f}')
30
31 # Plot of correctness vs c
32 plt.figure(figsize=(12,6))
33 plt.plot(c3, CP)
34 plt.xlabel('Parameter (scale factor)')
35 plt.ylabel('Correct cluster percentage')
36 plt.title('Model Selection of parameters')
37 plt.scatter(c3, CP, marker='*')
38
39 plt.savefig('P2E3.png')

```

Listing 5: Experiment 3

C Part 3 code

```

1 # Import libraries
2
3 import numpy as np
4 from matplotlib import pyplot as plt
5 from matplotlib import gridspec as gridspec
6 from datetime import datetime
7
8 # Data Generation function

```

```

9
10 def sample_01(m,n):
11     # Generate samples from 0 and 1
12     X = np.random.choice([0,1], (m,n))
13     Y = X[:, 0]
14     return X, Y
15
16 def sample_11(m,n):
17     # Generate samples from -1 and 1
18     X = np.random.choice([-1,1], (m,n))
19     Y = X[:, 0]
20     return X, Y
21
22 # Four algorithms
23
24 def perceptron(m, n):
25     # Perceptron algorithm: We generate the data and predict by w@x.
26     # Then check if the result is correct and update correspondingly.
27     X_train, Y_train = sample_11(m, n)
28     X_test, Y_test = sample_11(test_size, n)
29     W = np.zeros(n)
30     Y_pred = np.zeros(m)
31     for t in range(m):
32         Y_pred[t] = np.sign(W @ X_train[t])
33         if Y_pred[t] * Y_train[t] <= 0:
34             W += Y_train[t] * X_train[t]
35
36     # Use the final W to predict and calculate the error.
37     Y_res = np.sign(X_test @ W)
38     error = np.count_nonzero(Y_res != Y_test)
39     return Y_res, error
40
41 def winnow(m, n):
42     # Perceptron algorithm: We generate the data and predict by checking whether
43     # w@x is less than n.
44     # Then check if the result is correct and update correspondingly.
45     X_train, Y_train = sample_01(m, n)
46     X_test, Y_test = sample_01(test_size, n)
47     W = np.ones(n)
48     Y_pred = np.ones(m)
49     for t in range(m):
50         if W @ X_train[t] < n:
51             Y_pred[t] = 0
52         if Y_pred[t] != Y_train[t]:
53             W *= np.float_power(2, (Y_train[t] - Y_pred[t]) * X_train[t])
54
55     # Use the final W to predict and calculate the error.
56     Y_res = np.where(X_test @ W < n, 0, 1)
57     error = np.count_nonzero(Y_res != Y_test)
58     return Y_res, error

```

```

59 def least_squares(m, n):
60     # Least squares algorithm: We generate the data and predict by sign(x@w).
61     # Then calculate the error.
62     X_train, Y_train = sample_11(m, n)
63     X_test, Y_test = sample_11(test_size, n)
64     W = np.linalg.pinv(X_train) @ Y_train
65     Y_res = np.sign(X_test @ W)
66     error = np.count_nonzero(Y_res != Y_test)
67     return Y_res, error
68
69 def one_NN(m, n):
70     # One nearest neighbor algorithm: We generate the data and calculate the
71     # distance
72     # between every pair of train and test point, classify the test point as the
73     # nearest train point label.
74     X_train, Y_train = sample_11(m, n)
75     X_test, Y_test = sample_11(test_size, n)
76     Y_res = np.zeros(test_size)
77     for i in range(test_size):
78         distance = np.linalg.norm(X_train - X_test[i], ord = 2, axis = 1)
79         Y_res[i] = Y_train[np.argmin(distance)]
80     error = np.count_nonzero(Y_res != Y_test)
81     return Y_res, error
82
83 # My method of finding sample complexity.
84 def complexity(algorithm, max_n):
85     ''' My algorithm: try different m for each n no more than max_n for 10 runs
86     and calculate the mean error to see whether it is <= 0.1.
87     Use the final m for n as the starting m for n+1.
88
89     Input:
90         algorithm(string): name of the algorithm.
91         max_n(int): the maximum of n.
92
93     Return:
94         res(list): the final m for each n no more than max_n.
95     '''
96     res = np.zeros(max_n)
97     m = 1
98     for n in range(max_n):
99         while m < 10000:
100             error = []
101             for i in range(runs):
102                 error.append(algorithm(m, n+1)[1])
103             if np.mean(error) / test_size <= 0.1:
104                 res[n] = m
105                 break
106             m += 1
107     return res
108
109 # The normal method of finding sample complexity.

```



```

108 def complex_m(algorithm, max_n):
109     ''' Normal algorithm: try different m for each n no more than max_n for 10
        runs
110         and calculate the mean error to see whether it is <= 0.1.
111         Use 1 as the starting m for each n.
112
113         Input:
114             algorithm(string): name of the algorithm.
115             max_n(int): the maximum of n.
116
117         Return:
118             res(list): the final m for each n no more than max_n.
119     '''
120     res = np.zeros(max_n)
121     for n in range(max_n):
122         m = 1
123         while m < 10000:
124             error = []
125             for i in range(runs):
126                 error.append(algorithm(m, n+1)[1])
127             if np.mean(error) / test_size <= 0.1:
128                 res[n] = m
129                 break
130             m += 1
131     return res

```

Listing 6: Four Algorithms and two method functions

```

1 # Plot functions
2
3 def plot_complex(algorithm, max_n):
4     ''' Plot the graph of the called algorithm.
5
6         Input:
7             algorithm(string): name of the algorithm.
8             max_n(int): the maximum of n.
9     '''
10     sequence_n = [i for i in range(1, max_n+1) ]
11     res = complexity(algorithm, max_n)
12     plt.plot(sequence_n, res)
13     plt.xlabel('Dimension (n)')
14     plt.ylabel('Estimated number of samples (m)')
15     plt.title(f'm to obtain 10% generalisation error versus n for {algorithm.
        __name__}.')
16     plt.show()
17
18 def plot_all(algorithm_seq, max_n_seq):
19     ''' Plot the graphs of the all algorithms contained in the algorithm_seq.
20         And use the corresponding max_n in the max_n_seq.
21
22         Input:
23             algorithm_seq(list): list of names of the algorithms.

```

```

24         max_n_seq(list): list of all maximum value of n for each algorithm.
25     '''
26     fig = plt.figure(figsize=(12,8))
27     gs = gridspec.GridSpec(2, 2)
28     ax = [algorithm_seq[k].__name__ for k in range(4)]
29     ax[0] = plt.subplot(gs[0,0])
30     ax[1] = plt.subplot(gs[0,1])
31     ax[2] = plt.subplot(gs[1,0])
32     ax[3] = plt.subplot(gs[1,1])
33
34     for k in range(4):
35         algorithm = algorithm_seq[k]
36         res = complexity(algorithm, max_n_seq[k])
37         sequence_n = np.array([i for i in range(1, max_n_seq[k]+1)])
38
39         ax[k].plot(sequence_n, res)
40         ax[k].set_xlabel = 'n', ylabel = 'm'
41         ax[k].set_title(f'm versus n for {algorithm.__name__}')
42
43     plt.tight_layout()
44     fig.savefig('algorithm15.png')
45
46 def plot_fit(algorithm_seq, max_n_seq):
47     ''' Plot the graphs of the all algorithms contained in the algorithm_seq.
48         And use the corresponding max_n in the max_n_seq.
49         Also plot the fitted line for each graph using the specified function
50         form.
51
52         Input:
53             algorithm_seq(list): list of names of the algorithms.
54             max_n_seq(list): list of all maximum value of n for each algorithm.
55
56         Return:
57             Duration(list): the computation time for each algorithm.
58     '''
59     fig = plt.figure(figsize=(12,8))
60     gs = gridspec.GridSpec(2, 2)
61     ax = [algorithm_seq[k].__name__ for k in range(4)]
62     ax[0] = plt.subplot(gs[0,0])
63     ax[1] = plt.subplot(gs[0,1])
64     ax[2] = plt.subplot(gs[1,0])
65     ax[3] = plt.subplot(gs[1,1])
66
67     Duration = []
68     for k in range(4):
69         algorithm = algorithm_seq[k]
70         sequence_n = np.array([i for i in range(1, max_n_seq[k]+1)])
71
72         start_time = datetime.now()
73         res = complexity(algorithm, max_n_seq[k])
74         end_time = datetime.now()

```

```

74     Duration.append(end_time - start_time)
75
76     if k == 1:
77         a, b = np.polyfit(np.log(sequence_n), res, 1)
78         y = a * np.log(sequence_n) + b
79     elif k == 3:
80         a, b = np.polyfit(sequence_n, np.log(res), 1)
81         y = np.exp(a * sequence_n + b)
82     else:
83         a, b = np.polyfit(sequence_n, res, 1)
84         y = a * sequence_n + b
85
86     ax[k].plot(sequence_n, y, label = f'a = {a: 2f}, b = {b: 2f}')
87     ax[k].plot(sequence_n, res)
88     ax[k].set_xlabel = 'n', ylabel = 'm'
89     ax[k].set_title(f'm versus n for {algorithm.__name__}')
90     ax[k].legend()
91
92     plt.tight_layout()
93     fig.savefig('algorithm15 fit.png')
94     return Duration
95
96 def plot_comparison(algorithm_seq, max_n_seq):
97     ''' Plot the graphs of the all algorithms contained in the algorithm_seq.
98         And use the corresponding max_n in the max_n_seq.
99         Use both my method and the normal method to produce m and see the
100        comparison in graphs.
101
102        Input:
103            algorithm_seq(list): list of names of the algorithms.
104            max_n_seq(list): list of all maximum value of n for each algorithm.
105
106        Return:
107            Duration1(list): the computation time for each algorithm using my
108            method.
109            Duration2(list): the computation time for each algorithm using the
110            normal method.
111        '''
112     fig = plt.figure(figsize=(12,8))
113     gs = gridspec.GridSpec(2, 2)
114     ax = [algorithm_seq[k].__name__ for k in range(4)]
115     ax[0] = plt.subplot(gs[0,0])
116     ax[1] = plt.subplot(gs[0,1])
117     ax[2] = plt.subplot(gs[1,0])
118     ax[3] = plt.subplot(gs[1,1])
119
120     Duration1 = []
121     Duration2 = []
122     for k in range(4):
123         algorithm = algorithm_seq[k]
124         sequence_n = np.array([i for i in range(1, max_n_seq[k]+1)])

```

```

122
123     start_time = datetime.now()
124     res1 = complexity(algorithm, max_n_seq[k])
125     end_time = datetime.now()
126     Duration1.append(end_time - start_time)
127
128     start_time = datetime.now()
129     res2 = complex_m(algorithm, max_n_seq[k])
130     end_time = datetime.now()
131     Duration2.append(end_time - start_time)
132
133
134     ax[k].plot(sequence_n, res1, label = 'generate m from the m for n-1')
135     ax[k].plot(sequence_n, res2, label = 'generate m from 1')
136     ax[k].set(xlabel = 'n', ylabel = 'm')
137     ax[k].set_title(f'm versus n for {algorithm.__name__}')
138     ax[k].legend()
139
140     plt.tight_layout()
141     fig.savefig('algorithm15 comparison.png')
142     return Duration1, Duration2
143
144 # Specify the parameters not included in the algorithm.
145 runs = 10
146 test_size = 5000
147 max_n_seq = [100, 100, 100, 15]
148 algorithm_seq = [perceptron, winnow, least_squares, one_NN]
149
150 # Run the algorithms and plot
151 plot_all(algorithm_seq, max_n_seq)
152 plot_fit(algorithm_seq, max_n_seq)
153 plot_comparison(algorithm_seq, max_n_seq)

```

Listing 7: Plot functions and implementation