# CT4101 MACHINE LEARNING - ASSIGNMENT 2
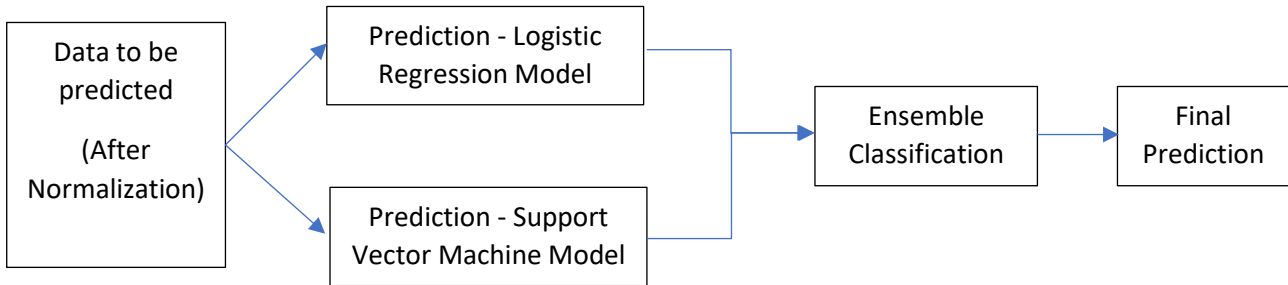
Prakhar Gurawa (20231064)
Yashitha Agarwal (20230091)
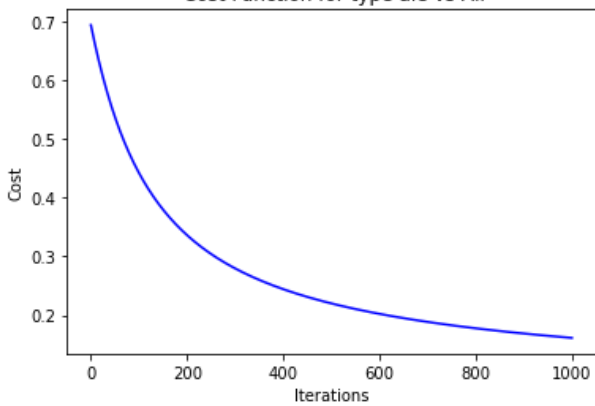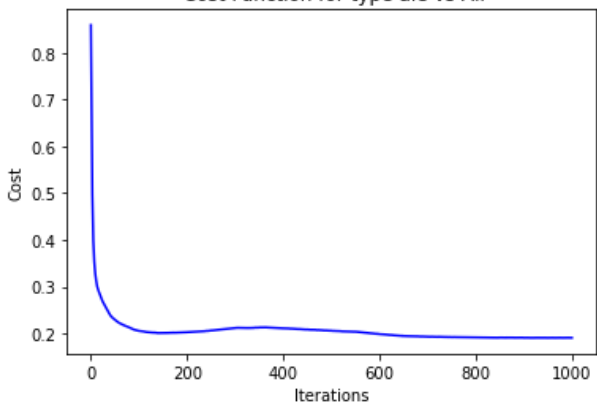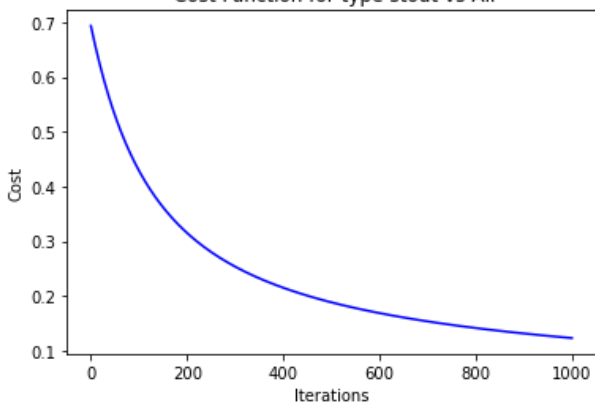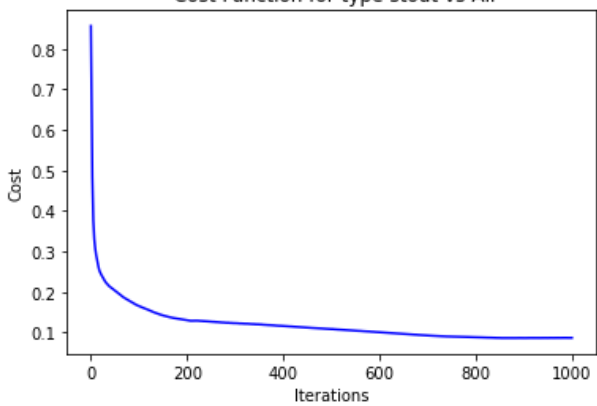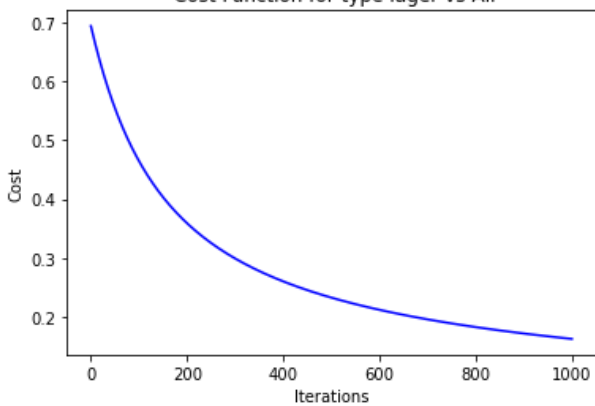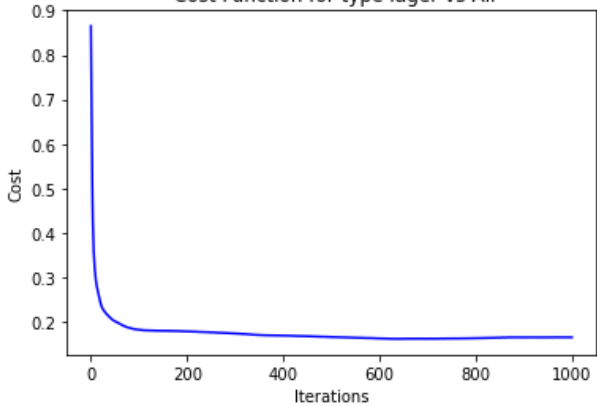
**Algorithm description and design decisions:** For this assignment, we have considered two machine learning models, namely Logistic Regression and Support Vector Machine classifiers, which have been implemented from scratch. Further, we have integrated both the models, based on their separate accuracy scores, to create an ensemble classification. We have also implemented standardization, which a feature scaling technique, from scratch.



1. **Logistic Regression:** In logistic regression model, the primary aim is to maximise the probability of minimising the loss function. To achieve this, first the features with weights are passed to the sigmoid function, which returns a result between 0 and 1. The result from the sigmoid function is evaluated using the loss function and then using gradient descent, weight of the feature is modified, to the optimum value, based on the derivatives of the loss function.

2. **Support Vector Machine:** In support vector machine model, the goal is to maximise the margin, which is the distance between the optimal hyperplane and the support vectors. The best line (hyperplane) for any given dataset separates the different classes with maximum margin. The algorithm does this by first finding the points that are closest to the line for each of the classes, called support vectors, and then calculates the distance between the line and support vectors. The line which has the maximum margin is chosen as the hyperplane.

3. **Ensemble Classification:** The ensemble classification integrates the predictions from the Logistic Regression model and the Support Vector Machine model, which will help to get a better prediction and make the algorithm more robust. Data to be predicted is passed to the ensemble classification, where the individual prediction from Logistic Regression model and Support Vector Machine model is taken. If the individual prediction of both algorithms is the same, the final prediction is taken to be the same, otherwise if both models predict different results, the prediction form the model with higher individual score is considered as the final prediction.

**Reference Implementation:** The Logistic Regression model and Support Vector Machine model from scikit learn have been used as a reference implementation to assess the performance of our implemented algorithms and the observations are reported below. Additionally, the ensemble classifier has also been implemented using scikit learn and the performance has been compared with the own implementation.

| Logistic Regression Model | Support Vector Machine Model |
|---|---|
| The sigmoid function has a s-shaped curve and exists between 0 and 1. Since the model needs to predict a probability as the output and probability values exist only between 0 to 1, the sigmoid function is used. | Hyperplanes are used as a decision boundary for the data points. Based on which side of the hyplerplane a given data point is located, the class to which the data point belongs to is determined. |
| Cost Function: The loss function used by Logistic Regression model is the below Binary cross entropy loss function. $$L(y, \hat{y}) = -\frac{1}{m} \sum_{i=1}^{m} [y_i \cdot \log(\hat{y}_i) + (1 - y_i) \cdot \log(1 - \hat{y}_i)]$$ | Cost Function: SVM uses Hinge loss function below. Regularization parameter is used to balance the cost function, which is set to (1 / Number of iterations). $$min_w \|w\|^2 + \sum_{i=1}^{n} (1 - y_i \langle x_i, w \rangle)_+$$ |

| Gradient descent in Logistic Regression is used to minimise the loss function. This is done by moving in the direction where the steepest (optimal) point is located on the given search landscape, using the below equation. $$w = w - x_i \cdot (\hat{y} - y)$$ | Gradient update in SVM depends on whether the prediction from our model is correct or incorrect. In case of missclassification : $$w = w + \alpha \cdot (y_i \cdot x_i - 2\lambda w)$$ In case of no missclassification : $$w = w - \alpha \cdot (2\lambda w)$$ |
|---|---|
| Number of iterations is set to 1000 for LR model. | Number of iterations is set to 1000 for SVM model. |



Cost Function for type ale vs All



Cost Function for type ale vs All



Cost Function for type stout vs All



Cost Function for type stout vs All



Cost Function for type lager vs All



Cost Function for type lager vs All

**Implementation Details:** Since the given problem is a multiclass classification problem, we have used the technique of one vs all in individual algorithms.

**Tests:** To ensure comprehensive tests are done, first the given data is imported into a dataset using pandas. The 'beer_id' column is removed from the dataset as it does not contribute towards identifying the type of beer. With the remaining dataset, we create dependent and independent (style) features, which is then split into training (two-thirds) and test (one-third) sets, shuffled over each iteration. Feature scaling is implemented for the dependent features, after which the data is fitted and scored over ten iterations, for each algorithm – Linear

Regression, Support Vector Machine and Ensemble Classifier of our own implementation and using scikit learn. The accuracy for each model over each iteration is shown in the images below.

**Results:**

1. Logistic Regression Classifier – Own implementation vs. Scikit Learn:

```
Logistic Regression Classifier Learning

Accuracy  0  =  0.9807692307692307
Accuracy  1  =  0.9807692307692307
Accuracy  2  =  0.9423076923076923
Accuracy  3  =  0.9423076923076923
Accuracy  4  =  0.9615384615384616
Accuracy  5  =  1.0
Accuracy  6  =  0.9807692307692307
Accuracy  7  =  0.9807692307692307
Accuracy  8  =  0.9615384615384616
Accuracy  9  =  0.9615384615384616

Mean accuracy =  0.9692307692307691
```

```
Logistic Regression Classifier Learning - Scikit

Accuracy  0  =  0.9615384615384616
Accuracy  1  =  0.9807692307692307
Accuracy  2  =  0.9615384615384616
Accuracy  3  =  0.9807692307692307
Accuracy  4  =  1.0
Accuracy  5  =  0.9423076923076923
Accuracy  6  =  0.9615384615384616
Accuracy  7  =  0.9615384615384616
Accuracy  8  =  0.9807692307692307
Accuracy  9  =  1.0

Mean accuracy =  0.973076923076923
```

2. Support Vector Machine Classifier – Own implementation vs. Scikit Learn:

```
SVM Classifier Learning

Accuracy  0  =  0.9038461538461539
Accuracy  1  =  0.9423076923076923
Accuracy  2  =  1.0
Accuracy  3  =  0.9807692307692307
Accuracy  4  =  0.9615384615384616
Accuracy  5  =  0.9615384615384616
Accuracy  6  =  0.9615384615384616
Accuracy  7  =  0.9615384615384616
Accuracy  8  =  0.9230769230769231
Accuracy  9  =  0.9230769230769231

Mean accuracy =  0.951923076923077
```

```
SVM Classifier Learning - Scikit

Accuracy  0  =  0.9807692307692307
Accuracy  1  =  0.9807692307692307
Accuracy  2  =  0.9615384615384616
Accuracy  3  =  0.9615384615384616
Accuracy  4  =  0.9615384615384616
Accuracy  5  =  0.9423076923076923
Accuracy  6  =  0.9423076923076923
Accuracy  7  =  0.9615384615384616
Accuracy  8  =  0.9230769230769231
Accuracy  9  =  0.9423076923076923

Mean accuracy =  0.9557692307692307
```

3. Ensemble Classifier – Own implementation vs. Scikit Learn

```
Ensemble Classifier Learning

Accuracy  0  =  0.9615384615384616
Accuracy  1  =  0.9615384615384616
Accuracy  2  =  0.9615384615384616
Accuracy  3  =  1.0
Accuracy  4  =  0.9807692307692307
Accuracy  5  =  0.9807692307692307
Accuracy  6  =  0.9807692307692307
Accuracy  7  =  1.0
Accuracy  8  =  0.9423076923076923
Accuracy  9  =  1.0

Mean accuracy =  0.976923076923077
```

```
Ensemble Classifier Learning - Scikit

Accuracy  0  =  0.9807692307692307
Accuracy  1  =  0.9807692307692307
Accuracy  2  =  0.9230769230769231
Accuracy  3  =  0.9807692307692307
Accuracy  4  =  0.9807692307692307
Accuracy  5  =  0.9615384615384616
Accuracy  6  =  0.9807692307692307
Accuracy  7  =  0.9423076923076923
Accuracy  8  =  0.9615384615384616
Accuracy  9  =  0.9807692307692307

Mean accuracy =  0.9673076923076922
```

**Observations:** From the results above, starting with the Logistic Regression classifier, it can be observed that the mean accuracy of the own implementation 96.9%, is just 0.4% lower than the scikit learn version of 97.3%. Similarly, comparing the mean accuracy for SVM model, again the own implementation, 95.2% performs just 0.4% lower than the scikit implementation value of 95.6%. However, for the ensemble classifier, there is an improvement by 1% in the own implementation, which scored 97.7% when compared to the 96.7%.

On analysing the three classifiers from the own implementation alone, it can be seen that the Logistic Regression (96.9%) model performs better than the Support Vector Machine (95.2%) model by 1.7%, while the ensemble classifier (97.7%) model gives the best performance among the three.

**Conclusion:** Overall, it can be seen that the Logistic Regression and Support Vector Machine models that have been implemented from scratch performs almost as good as the implementations using scikit learn. Also, the results above show that the performance of the ensemble classifier from our own implementation is slightly better when compared to the scikit learn implementation.

**Future Work:** The ensemble model implemented here, uses just two algorithms to make a final prediction. For further improvements, this model can be extended to include more than two algorithms, in which case it will become a voting classifier. For instance, if we decide to include five different algorithms, and three algorithms predict 'Class A' and the remaining algorithms predict 'Class B', then the final prediction would be 'Class A' as it has the highest vote among the five models.

**Contribution of each team member:** This assignment has been done as a group of two members, with the overview of the contribution of each member, towards the coding part, given below. Additionally, more specific code contributions are also mentioned in comments on the source code files.

Prakhar Gurawa (20231064):
- Own implementation: *SupportVectorMachine.py*, *EnsembleClassifier.py* (partially)
- Scikit implementation: *SupportVectorMchine_Scikit.py*, *EnsembleClassifier_Scikit.py* (partially)
- Report: Reference Implementation, Comparison Table, Implementation Details, Observation, Future work, References

Yashitha Agarwal (20230091):
- Own implementation: *LogisticRegression.py*, *EnsembleClassifier.py* (partially)
- Scikit implementation: *LogisticRegression_Scikit.py*, *EnsembleClassifier_Scikit.py* (partially)
- Report: Algorithm description and design decisions, Tests, Results, Conclusion, Contribution of each team member, Code execution

**Code execution:** All the python files required to compile the solutions and the dataset are zipped into a single folder. In order to execute the ensemble classification model (implementation from scratch), run the *EnsembleClassifier.py*, after ensuring that *SupportVectorMachine.py* and *LogisticRegression.py* are present in the same folder, as it imports those python files along with the *beer.txt* dataset. To execute our own implementation of individual classification models, Logistic Regression or Support Vector Machine, run the *LogisticRegression.py* or *SupportVectorMachine.py* file respectively.

On executing the *LogisticRegression.py*, *SupportVectorMachine.py* and *EnsembleClassifier.py* individually, corresponding **output files** showing the predicted and actual value for each model over the 10 iterations are created with the file names **LR_Results.csv**, **SVM_Results.csv** and **Ensemble_Results.csv** respectively. Scikit implementation of both the algorithms, and the ensemble classification are also present in the same folder.

**References:**
1. https://towardsdatascience.com/logistic-regression-from-scratch-69db4f587e17
2. https://towardsdatascience.com/logistic-regression-detailed-overview-46c4da4303bc
3. https://towardsdatascience.com/activation-functions-neural-networks-1cbd9f8d91d6
4. https://stackoverflow.com/questions/47966728/how-to-fix-float-object-has-no-attribute-exp?noredirect=1&lq=1
5. https://stackoverflow.com/questions/56594598/change-1s-to-0-and-0s-to-1-in-numpy-array-without-looping/56594688
6. https://datascience.stackexchange.com/questions/22470/python-implementation-of-cost-function-in-logistic-regression-why-dot-multiplic
7. https://stackoverflow.com/questions/56594598/change-1s-to-0-and-0s-to-1-in-numpy-array-without-looping/56594688
8. https://towardsdatascience.com/svm-implementation-from-scratch-python-2db2fc52e5c2
9. https://towardsdatascience.com/https-medium-com-pupalerushikesh-svm-f4b42800e989
10. https://ml-cheatsheet.readthedocs.io/en/latest/gradient_descent.html
11. https://towardsdatascience.com/multi-class-classification-one-vs-all-one-vs-one-94daed32a87b

```python
"""Final Source Code - Appendix"""

'''
LOGISTIC REGRESSION - OWN IMPLEMENTATION
'''

# -*- coding: utf-8 -*-
'''
CT4101 MACHINE LEARNING - ASSIGNMENT 2
Prakhar Gurawa (20231064)
Yashitha Agarwal (20230091)

Code by: Yashitha Agarwal (20230091)
'''

# References:
# https://towardsdatascience.com/logistic-regression-from-scratch-69db4f587e17 (A nice
# explaination of logistic regression and its mathematics)

# We are importing all necessary libraries to implement our model
import matplotlib.pyplot as plt
import numpy as np
import pandas as pd
from sklearn.model_selection import train_test_split


class LogisticRegression:

    def __init__(self,alpha=0.01,iterations=1000): # constructor function to
    intialize learning rate alpha and number of iterations
        self.alpha = alpha
        self.iterations = iterations

    def sigmoid(self, z): # utility function to find sigmoid values of input
        # Reference :
        # https://stackoverflow.com/questions/47966728/how-to-fix-float-object-has-no-at
        # tribute-exp?noredirect=1&lq=1
        z = np.array(z,dtype=float)
        return 1 / (1 + np.exp(-z))

    def transformMultiClass(self,y,c): # utility function to convert given
    mutilclass vector to 2 class 0/1 vector
        # Reference :
        # https://stackoverflow.com/questions/56594598/change-1s-to-0-and-0s-to-1-in-num
        # py-array-without-looping/56594688
        y_copy = y.copy()
        indices_C = y==c            # all indices with class 'c'
        indices_notC = y!=c         # all indices with classes other than 'c'
        y_copy[indices_C]=1         # convert values to 1 for all class 'c'
        y_copy[indices_notC]=0      # convert values to 0 for all class other than 'c'
        return y_copy

    def gradientDescent(self,X,y,weight,h): # function to calculate gradient descent
    for logistic regression
        # By calculus the derivative of cost function wrt weights comes out to be
        xi*(y_pred - y)
        """
        # Reference :
        # https://datascience.stackexchange.com/questions/22470/python-implementation-of
        # -cost-function-in-logistic-regression-why-dot-multiplic
        Cost function for Logistic Regression  = -1/m * np.sum(np.dot(Y,np.log(A)) +
        np.dot(1-Y, np.log(1-A)))
        dw = 1/m * np.dot(X, dz.T)
        """
        costGradient = np.dot(X.T, (h - y)) / len(y) # gradient of cost function for
        logistic regression (calculated using chain rule of calculus)
        weight = weight - self.alpha * costGradient  # gradient descent process
        (updating weight on basis of cost gradients)
        return weight

    def costFunction(self,y,weight,h): # defines cost function for logic regression
        costValue = (1 / len(y)) * (np.sum(-y.T.dot(np.log(h)) - (1 -
```

```python
           y).T.dot(np.log(1 - h)))) # Binary cross entropy loss function
58            return costValue
59
60        def fit(self,X,y): # starts the logistic regression model by fitting our given
          dataset
61            self.costs = []
62            self.weights = []
63            X = np.insert(X, 0, 1, axis=1) # adding 1 for bias term
64            classes = set(y)                # storing unique classes (our predicted
              output will be one of them)
65            # Using concept of one-vs-all where a particular class is treated as 1 and
              all other 0 and this process is repeated for all classes
66            for c in classes:
67                # Gradient descent for class 'c'
68                y_onevall = self.transformMultiClass(y,c)
69                weight = np.zeros(X.shape[1]) # initializing weights with 0 at beginning
                  of logistic regression
70                cost = list()
71                for itr in range(self.iterations):
72                    z = X.dot(weight)
73                    h = self.sigmoid(z)
74                    weight = self.gradientDescent(X,y_onevall,weight,h)
75                    cst = self.costFunction(y_onevall,weight,h)
76                    cost.append(cst)
77                self.weights.append((weight,c))
78                self.costs.append((cost,c))
79            return self
80
81        def predict(self,X): # predict class values for given independent features
82            X = np.insert(X, 0, 1, axis=1)
83            X_prediction = list() # storing predicted classes
84            for x in X:
85                class_predictions = [(self.sigmoid(x.dot(weight)),c) for weight,c in
                  self.weights] # This loop runs n times for n classes (multi class
                  logistic regression one vs all)
86
                  X_prediction.append(max(class_predictions)[1])
                          # append the class with maximum prediction value (probablity)
87            return X_prediction
88
89        def score(self,X,y): # function to calculate number of matches between actual
          classes and predicted classes by our model
90            size = len(y)
91            return sum(self.predict(X)==y)/size # number of matches divided by total
              inputs
92
93        def plotCost(self,costs): # utility function to plot cost value per class
94             for cost,c in costs   :
95                    plt.plot(range(len(cost)),cost,'blue')
96                    plt.title(" Cost Function for type " + str(c) +" vs All")
97                    plt.xlabel("Iterations")
98                    plt.ylabel("Cost")
99                    plt.show()
100
101
102   if __name__=="__main__":
103       # importing dataset using pandas
104       filename = 'beer.txt'
105       header_list =
          ['caloric_value','nitrogen','turbidity','style','alcohol','sugars','biterness','be
          er_id','colour','degree_of_fermentation']
106       data = pd.read_csv(filename,sep='\t', header=None,dtype=str,names=header_list)
107
108       # creating dependent and independent features
109       X = data.drop(['style','beer_id'],axis=1).values
110       y = data['style'].values
111
112       # creating a pandas dataframe for storing results
113       predictions_final = pd.DataFrame(columns=['Iteration','Predicted Value','Actual
          Value'])
114
115       # data stardaization pre-processing
```

```python
116        def feature_scaling(X):
117            X = X.astype(np.float)
118            mean = np.mean(X, axis=0)
119            sd = np.std(X, axis=0)
120            X_scaled= (X -  mean) / sd
121            return X_scaled
122
123        X = feature_scaling(X)
124
125        scores=list()
126        print("Logistic Regression Classifier Learning\n")
127        for i in range(10):
128            X_train,X_test,y_train,y_test  = train_test_split(X, y, train_size = 2/3,
                   shuffle = True) # split data
129            model = LogisticRegression(alpha=0.01,iterations=1000)
130            model.fit(X_train, y_train)
131            prediction = model.predict(X_test)
132            score = model.score(X_test,y_test)
133            # storing all the predictions and actual values in the dataframe
134            for (p , a) in zip(prediction, y_test):
135                predictions_final = predictions_final.append({'Iteration': i, 'Predicted
                       Value': p, 'Actual Value': a}, ignore_index=True)
136            print("Accuracy ",i," = ",score)
137            scores.append(score)
138
139        print("\nMean accuracy = ",np.mean(scores))
140        model.plotCost(model.costs)
141
142        # output the results to csv file
143        predictions_final.to_csv('LR_Results.csv', index=False)
144
145    '''
146    ----------------------------------------------------------------------------------------
       ----------------------------------------------------------------------------------------
147    '''
148
149    '''
150    LOGISTIC REGRESSION - SCIKIT LEARN
151    '''
152
153    # -*- coding: utf-8 -*-
154    '''
155    CT4101 MACHINE LEARNING - ASSIGNMENT 2
156    Prakhar Gurawa (20231064)
157    Yashitha Agarwal (20230091)
158
159    Code by: Yashitha Agarwal (20230091)
160    '''
161
162    # We are importing all necessary libraries to implement our model
163    import pandas as pd
164    import numpy as np
165    from sklearn.linear_model import LogisticRegression
166    from sklearn.preprocessing import StandardScaler
167    from sklearn.model_selection import train_test_split
168
169    # importing dataset using pandas
170    filename = 'beer.txt'
171    header_list =
       ['caloric_value','nitrogen','turbidity','style','alcohol','sugars','biterness','beer_i
       d','colour','degree_of_fermentation']
172    data = pd.read_csv(filename,sep='\t', header=None,dtype=str,names=header_list)
173
174    # creating dependent and independent features
175    X = data.drop(['style','beer_id'],axis=1).values
176    y = data['style'].values
177
178    # data stardaization pre-processing
179    scaler = StandardScaler()
180    X = scaler.fit_transform(X)
181
182    scores=list()
```

```
183    print("Logistic Regression Classifier Learning - Scikit\n")
184    for i in range(10):
185        X_train,X_test,y_train,y_test = train_test_split(X, y, train_size = 2/3,
           shuffle = True) # split data
186        model = LogisticRegression()
187        model.fit(X_train, y_train)
188        prediction = model.predict(X_test)
189        score = model.score(X_test,y_test)
190        print("Accuracy ",i," = ",score)
191        scores.append(score)
192
193    print("\nMean accuracy = ",np.mean(scores))
194
195    '''
196    --------------------------------------------------------------------------------
       --------------------------------------------------------------------------------
197    '''
198
199    '''
200    SUPPORT VECTOR MACHINE - OWN IMPLEMENTATION
201    '''
202
203    # -*- coding: utf-8 -*-
204    '''
205    CT4101 MACHINE LEARNING - ASSIGNMENT 2
206    Prakhar Gurawa (20231064)
207    Yashitha Agarwal (20230091)
208
209    Code by: Prakhar Gurawa (20231064)
210    '''
211
212    #We are importing all necessary libraries to implement our model
213    import matplotlib.pyplot as plt
214    import numpy as np
215    import pandas as pd
216    from sklearn.model_selection import train_test_split
217
218    class SVM:
219
220        def __init__(self, alpha=0.001,iterations=1000): # constructor function to
           intialize learning rate alpha, lambda param and number of iterations
221            self.alpha = alpha
222            self.iterations = iterations
223            # self.lambda_param = lambda_param (if in future user wants to give some
               well defined lambda param say for example : lambda_param=0.01 , currently
               taking as 1/epoch)
224
225        def transformMultiClass(self,y,c): # utility function to convert given
           mutilclass vector to 2 class 0/1 vector
226            # Reference :
               https://stackoverflow.com/questions/56594598/change-1s-to-0-and-0s-to-1-in-num
               py-array-without-looping/56594688
227            y_copy = y.copy()
228            indices_C = y==c            # all indices with class 'c'
229            indices_notC = y!=c         # all indices with classes other than 'c'
230            y_copy[indices_C]=1         # convert values to 1 for all class 'c'
231            y_copy[indices_notC]=-1     # convert values to 0 for all class other than 'c'
232            return y_copy
233
234        def compute_cost(self,W, X, Y):
235            # Reference :
               https://towardsdatascience.com/svm-implementation-from-scratch-python-2db2fc52
               e5c2
236            # calculate hinge loss - SVM uses hinge loss function
237            N = X.shape[0]
238            distances = 1 - Y * (np.dot(X, W))
239            distances[distances < 0] = 0  # equivalent to max(0, distance)
240            hinge_loss = 10000 * (np.sum(distances) / N) # 10000 is regularization
               parameter (fixed for now)
241            # calculate cost
242            cost = 1 / 2 * np.dot(W, W) + hinge_loss
243            return cost
```

```python
244
245
246         def fit(self,X,y):
247             n_samples, n_features = X.shape
248             self.costs = []
249             self.weights = []
250             classes = set(y)                    # storing unique classes (our predicted
                    output will be one of them)
251             # Using concept of one-vs-all where a particular class is treated as 1 and
                all other 0 and this process is repeated for all classes
252             for c in classes:
253                 # weight updation for class 'c'
254                 y_onevall = self.transformMultiClass(y,c)
255                 weight = np.zeros(n_features) # initializing weights with 0 at beginning
                    of SVM Classifier
256                 bias = 0                        # initailizing bias as 0 at beginning of
                    SVM Classifier
257                 cost = list()
258                 for itr in range(self.iterations):
259                     for idx, x_i in enumerate(X):    # iterating through full dataset
260                         lambda_param = 1/(itr+1)     # regularization param ( will
                            decrease with iteration)
261                         condition = y_onevall[idx] * (np.dot(x_i, weight) - bias) >= 1 #
                            standard condition of SVM Classifier
262                         if condition:
263                             weight = weight - self.alpha * (2 * lambda_param * weight) #
                                gradient update when no misclassification
264                         else:
265                             weight = weight - self.alpha * (2 * lambda_param * weight -
                                np.dot(x_i, y_onevall[idx])) # gradient update when
                                misclassification
266                             bias = bias - self.alpha * y_onevall[idx]
267
268                     cst=self.compute_cost(weight,X,y_onevall)
269                     cost.append(cst)
270                 self.weights.append((weight,bias,c))
271                 self.costs.append((cost,c))
272             return self
273
274         def predict(self,X):
275             X_prediction = list() # storing predicted classes
276             for x in X:
277                 class_predictions = [(np.dot(x, weight) - bias ,c) for weight,bias,c in
                    self.weights] # This loop runs n times for n classes (multi class
                    logistic regression one vs all)
278                 X_prediction.append(max(class_predictions)[1])
279             return X_prediction
280
281         def score(self,X,y): # function to calculate number of matches between actual
            classes and predicted classes by our model
282             size = len(y)
283             return sum(self.predict(X)==y)/size # number of matches divided by total
                inputs
284
285         def plotCost(self,costs): # utility function to plot cost value per class
286             for cost,c in costs:
287                 scaledCost = [c/10000 for c in cost ] # Dividing by 10000 for scaled
                    output
288                 plt.plot(range(len(cost)),scaledCost,'blue')
289                 plt.title(" Cost Function for type " + str(c) +" vs All")
290                 plt.xlabel("Iterations")
291                 plt.ylabel("Cost")
292                 plt.show()
293
294     if __name__=="__main__":
295         # importing dataset using pandas
296         filename = 'beer.txt'
297         header_list =
            ['caloric_value','nitrogen','turbidity','style','alcohol','sugars','biterness','be
            er_id','colour','degree_of_fermentation']
298         data = pd.read_csv(filename,sep='\t', header=None,dtype=str,names=header_list)
299         data.dtypes
```

```python
300
301       # creating dependent and independent features
302       y = data['style'].values
303       X = data.drop(['style','beer_id'],axis=1).values
304
305       # creating a pandas dataframe for storing results
306       predictions_final = pd.DataFrame(columns=['Iteration','Predicted Value','Actual
          Value'])
307
308       # data stardaization pre-processing
309       def feature_scaling(X):
310           X = X.astype(np.float)
311           mean = np.mean(X, axis=0)
312           sd = np.std(X, axis=0)
313           X_scaled= (X -  mean) / sd
314           return X_scaled
315
316       X = feature_scaling(X)
317
318       scores=list()
319       print("SVM Classifier Learning\n")
320       for i in range(10):
321           X_train,X_test,y_train,y_test = train_test_split(X, y, train_size = 2/3,
              shuffle = True) # split data
322           model = SVM(alpha=0.001,iterations=1000)
323           model.fit(X_train, y_train)
324           prediction = model.predict(X_test)
325           score = model.score(X_test,y_test)
326           # storing all the predictions and actual values in the dataframe
327           for (p , a) in zip(prediction, y_test):
328               predictions_final = predictions_final.append({'Iteration': i, 'Predicted
                  Value': p, 'Actual Value': a}, ignore_index=True)
329           print("Accuracy ",i," = ",score)
330           scores.append(score)
331
332
333       print("\nMean accuracy = ",np.mean(scores))
334       model.plotCost(model.costs)
335
336       # output the results to csv file
337       predictions_final.to_csv('SVM_Results.csv', index=False)
338
339   '''
340   -------------------------------------------------------------------------------------
      -------------------------------------------------------------------------------------
341   '''
342
343   '''
344   SUPPORT VECTOR MACHINE - SCIKIT LEARN
345   '''
346
347   # -*- coding: utf-8 -*-
348   '''
349   CT4101 MACHINE LEARNING - ASSIGNMENT 2
350   Prakhar Gurawa (20231064)
351   Yashitha Agarwal (20230091)
352
353   Code by: Prakhar Gurawa (20231064)
354   '''
355
356   # We are importing all necessary libraries to implement our model
357   import pandas as pd
358   import numpy as np
359   from sklearn.svm import SVC
360   from sklearn.preprocessing import StandardScaler
361   from sklearn.model_selection import train_test_split
362
363   # importing dataset using pandas
364   filename = 'beer.txt'
365   header_list =
      ['caloric_value','nitrogen','turbidity','style','alcohol','sugars','biterness','beer_i
      d','colour','degree_of_fermentation']
```

```python
366    data = pd.read_csv(filename,sep='\t', header=None,dtype=str,names=header_list)
367
368    # creating dependent and independent features
369    X = data.drop(['style','beer_id'],axis=1).values
370    y = data['style'].values
371
372    # data stardaization pre-processing
373    scaler = StandardScaler()
374    X = scaler.fit_transform(X)
375
376    scores=list()
377    print("SVM Classifier Learning - Scikit\n")
378    for i in range(10):
379        X_train,X_test,y_train,y_test = train_test_split(X, y, train_size = 2/3, shuffle
           = True) # split data
380        model = SVC()
381        model.fit(X_train, y_train)
382        prediction = model.predict(X_test)
383        score = model.score(X_test,y_test)
384        print("Accuracy ",i," = ",score)
385        scores.append(score)
386
387
388    print("\nMean accuracy = ",np.mean(scores))
389
390    '''
391    ------------------------------------------------------------------------------------
       ------------------------------------------------------------------------------------
392    '''
393
394    '''
395    ENSEMBLE CLASSIFIER - OWN IMPLEMENTATION
396    '''
397
398    # -*- coding: utf-8 -*-
399    '''
400    CT4101 MACHINE LEARNING - ASSIGNMENT 2
401    Prakhar Gurawa (20231064)
402    Yashitha Agarwal (20230091)
403
404    Code by: Combined effort (specific parts mentioned in comments)
405    '''
406
407    # We are importing all necessary libraries to implement our model
408    import numpy as np
409    import pandas as pd
410    from sklearn.model_selection import train_test_split
411    from LogisticRegression import LogisticRegression
412    from SupportVectorMachine import SVM
413
414    class EnsembleClassifier:
415
416        # SVM and LR classifers are considered to make final model more robust
417        # Scratch implementation of SVM : SupportVectorMachine.py
418        # Scratch implementation of Logistic Regression : LogisticRegression.py
419
420        # Code by: Yashitha Agarwal (20230091)
421        def __init__(self,lrAlpha=0.01,svmAlpha=0.01,iterations=1000): # Constructor
           function to initalize individual hyperparameters for LR and SVM
422            self.lrAlpha = lrAlpha
423            self.svmAlpha =svmAlpha
424            self.iterations = iterations
425            self.lrModel = None
426            self.svmModel = None
427
428        # Code by: Yashitha Agarwal (20230091)
429        def fit(self,X,y):
430            self.lrModel = LogisticRegression(self.lrAlpha,self.iterations)
431            self.svmModel = SVM(self.svmAlpha,self.iterations)
432            self.lrModel.fit(X,y)    # Fitting independent and dependent feature in
               Logistic Regression
433            self.svmModel.fit(X,y)   # Fitting independent and dependent feature in
```

```python
                    Support Vector Machine Classifier

        # Code by: Prakhar Gurawa (20231064)
        def predict(self,X,y):
            lrScore = self.lrModel.score(X,y)
            svmScore = self.svmModel.score(X,y)
            lrPrediction = self.lrModel.predict(X)      # Prediction of Logistic
                Regression model
            svmPrediction = self.svmModel.predict(X)    # Prediction of Support Vector
                Machine model
            # Currently we are considering only two algorithms and considering
                prediction of that higher score classifier in case of disagreement
            finalPrediction = list() # Storing predicted classes
            for i in range(len(lrPrediction)):
                if lrPrediction[i] == svmPrediction[i]:
                    finalPrediction.append(lrPrediction[i]) # Case 1: Both LR ans SVM
                        predict to same class
                else:                                        # Case 2: Disagreement
                    between LR and SVM classifiers
                    if lrScore > svmScore:
                        finalPrediction.append(lrPrediction[i])
                    else:
                        finalPrediction.append(svmPrediction[i])
            # Future work : If we have more than two algorithms we will make this as a
                voting classifier.
            # Mutiple classifer are considered and majority of prediction is taken as
                final prediction.
            return finalPrediction # Final Predictions using ensemble

        # Code by: Prakhar Gurawa (20231064)
        def score(self,X,y): # Function to calculate number of matches between actual
            classes and predicted classes by our model
            size = len(y)
            return sum(self.predict(X,y)==y)/size # Number of matches divided by total
                inputs

    # Code by: Yashitha Agarwal (20230091)
    # importing dataset using pandas
    filename = 'beer.txt'
    header_list =
    ['caloric_value','nitrogen','turbidity','style','alcohol','sugars','biterness','beer_i
    d','colour','degree_of_fermentation']
    data = pd.read_csv(filename,sep='\t', header=None,dtype=str,names=header_list)

    # creating dependent and independent features
    X = data.drop(['style','beer_id'],axis=1).values
    y = data['style'].values

    # creating a pandas dataframe for storing results
    predictions_final = pd.DataFrame(columns=['Iteration','Predicted Value','Actual
    Value'])


    # Code by: Prakhar Gurawa (20231064)
    # data stardaization pre-processing
    def feature_scaling(X):
        X = X.astype(np.float)
        mean = np.mean(X, axis=0)
        sd = np.std(X, axis=0)
        X_scaled= (X -  mean) / sd
        return X_scaled

    X = feature_scaling(X)

    scores=list()
    print("Ensemble Classifier Learning\n")
    for i in range(10):
        X_train,X_test,y_train,y_test  = train_test_split(X, y, train_size = 2/3,
            shuffle = True) # split data
        model = EnsembleClassifier()
        model.fit(X_train, y_train)
        # passing y_test to predict since it is needed to get the score of LR and SVM
```

```python
        model
492    prediction = model.predict(X_test, y_test)
493    score = model.score(X_test,y_test)
494    # storing all the predictions and actual values in the dataframe
495    for (p , a) in zip(prediction, y_test):
496        predictions_final = predictions_final.append({'Iteration': i, 'Predicted
               Value': p, 'Actual Value': a}, ignore_index=True)
497    print("Accuracy ",i," = ",score)
498    scores.append(score)
499
500 print("\nMean accuracy = ",np.mean(scores))
501
502 # output the results to csv file
503 predictions_final.to_csv('Ensemble_Results.csv', index=False)
504
505 '''
506 --------------------------------------------------------------------------------
     --------------------------------------------------------------------------------
507 '''
508
509 '''
510 ENSEMBLE CLASSIFIER - SCIKIT LEARN
511 '''
512
513 # -*- coding: utf-8 -*-
514 '''
515 CT4101 MACHINE LEARNING - ASSIGNMENT 2
516 Prakhar Gurawa (20231064)
517 Yashitha Agarwal (20230091)
518
519 Code by: Combined effort (specific parts mentioned in comments)
520 '''
521
522 # We are importing all necessary libraries to implement our model
523 import numpy as np
524 import pandas as pd
525 from sklearn.model_selection import train_test_split
526 from sklearn.preprocessing import StandardScaler
527 from sklearn.linear_model import LogisticRegression
528 from sklearn.svm import SVC
529
530 class EnsembleClassifier:
531
532    # SVM and LR classifers are considered to make final model more robust
533    # Code by: Yashitha Agarwal (20230091)
534    def fit(self,X,y):
535        self.lrModel = LogisticRegression()
536        self.svmModel = SVC()
537        self.lrModel.fit(X,y)   # Fitting independent and dependent feature in
               Logistic Regression
538        self.svmModel.fit(X,y)  # Fitting independent and dependent feature in
               Support Vector Machine Classifier
539
540    # Code by: Prakhar Gurawa (20231064)
541    def predict(self,X,y):
542        lrScore = self.lrModel.score(X,y)
543        svmScore = self.svmModel.score(X,y)
544        lrPrediction = self.lrModel.predict(X)       # Prediction of Logistic
               Regression model
545        svmPrediction = self.svmModel.predict(X)    # Prediction of Support Vector
               Machine model
546        # Currently we are considering only two algorithms and considering
               prediction of that higher score classifier in case of disagreement
547        finalPrediction = list() # Storing predicted classes
548        for i in range(len(lrPrediction)):
549            if lrPrediction[i] == svmPrediction[i]:
550                finalPrediction.append(lrPrediction[i]) # Case 1: Both LR ans SVM
                       predict to same class
551            else:                                        # Case 2: Disagreement
                   between LR and SVM classifiers
552                if lrScore > svmScore:
553                    finalPrediction.append(lrPrediction[i])
```

```python
554                else:
555                    finalPrediction.append(svmPrediction[i])
556        # Future work : If we have more than two algorithms we will make this as a
           voting classifier.
557        # Mutiple classifer are considered and majority of prediction is taken as
           final prediction.
558        return finalPrediction # Final Predictions using ensemble
559
560    # Code by: Prakhar Gurawa (20231064)
561    def score(self,X,y): # Function to calculate number of matches between actual
       classes and predicted classes by our model
562        size = len(y)
563        return sum(self.predict(X,y)==y)/size # Number of matches divided by total
           inputs
564
565
566 # Code by: Yashitha Agarwal (20230091)
567 # importing dataset using pandas
568 filename = 'beer.txt'
569 header_list =
    ['caloric_value','nitrogen','turbidity','style','alcohol','sugars','biterness','beer_i
    d','colour','degree_of_fermentation']
570 data = pd.read_csv(filename,sep='\t', header=None,dtype=str,names=header_list)
571
572 # creating dependent and independent features
573 X = data.drop(['style','beer_id'],axis=1).values
574 y = data['style'].values
575
576 scaler = StandardScaler()
577 X = scaler.fit_transform(X)
578
579 # Code by: Prakhar Gurawa (20231064)
580 scores=list()
581 print("Ensemble Classifier Learning - Scikit\n")
582 for i in range(10):
583     X_train,X_test,y_train,y_test  = train_test_split(X, y, train_size = 2/3,
        shuffle = True) # split data
584     model = EnsembleClassifier()
585     model.fit(X_train, y_train)
586     # passing y_test to predict since it is needed to get the score of LR and SVM
        model
587     prediction = model.predict(X_test, y_test)
588     score = model.score(X_test,y_test)
589     print("Accuracy ",i," = ",score)
590     scores.append(score)
591
592 print("\nMean accuracy = ",np.mean(scores))
```