```python
 1  import pandas as pd
 2  import numpy as np
 3  import sys
 4  import random as rd
 5
 6  #insert an all-one column as the first column
 7  def addAllOneColumn(matrix):
 8      n = matrix.shape[0] #total of data points
 9      p = matrix.shape[1] #total number of attributes
10
11      newMatrix = np.zeros((n,p+1))
12      newMatrix[:,1:] = matrix
13      newMatrix[:,0] = np.ones(n)
14
15      return newMatrix
16
17  # Reads the data from CSV files, converts it into Dataframe and returns x and
    y dataframes
18  def getDataframe(filePath):
19      dataframe = pd.read_csv(filePath)
20      y = dataframe['y']
21      x = dataframe.drop('y', axis=1)
22      return x, y
23
24  # train_x and train_y are numpy arrays
25  # function returns value of beta calculated using (0) the formula beta =
    (X^T*X)^ -1)*(X^T*Y)
26  def getBeta(train_x, train_y):
27      n = train_x.shape[0] #total of data points
28      p = train_x.shape[1] #total number of attributes
29
30      beta = np.zeros(p)
31      #=========================#
32      # STRART YOUR CODE HERE   #
33      #=========================#
34      # 1. Calculate the transpose
35      # 2. first_term = (X^T*X)^ -1)
36      # 3. second_term = (X^T*Y)
37      # 4. calculate beta = first_term * second_term
38      train_x_transpose = np.transpose(train_x)
39      first_term = np.linalg.inv(np.matmul(train_x_transpose, train_x))
40      second_term = np.matmul(train_x_transpose, train_y)
41      beta = np.matmul(first_term, second_term)
42      #=========================#
43      #   END YOUR CODE HERE    #
44      #=========================#
45      return beta
46
47  # train_x and train_y are numpy arrays
48  # lr (learning rate) is a scalar
49  # function returns value of beta calculated using (1) batch gradient descent
50  def getBetaBatchGradient(train_x, train_y, lr, num_iter):
51      beta = np.random.rand(train_x.shape[1])
52
53      n = train_x.shape[0] #total of data points
54      p = train_x.shape[1] #total number of attributes
55
56
57      beta = np.random.rand(p)
```

```python
59      for iter in range(0, num_iter):
60          deriv = np.zeros(p)
61          for i in range(n):
62              #=======================#
63              # STRART YOUR CODE HERE  #
64              #=======================#
65               x_i = train_x[i]
66               y_i = train_y[i]
67               x_i_transpose = np.transpose(x_i)
68               A = np.matmul(x_i_transpose, beta)
69               deriv += x_i * (A - y_i)
70              #=======================#
71              #   END YOUR CODE HERE   #
72              #=======================#
73          deriv = deriv / n
74          beta = beta - deriv.dot(lr)
75      return beta
76
77 # train_x and train_y are numpy arrays
78 # lr (learning rate) is a scalar
79 # function returns value of beta calculated using (2) stochastic gradient
   descent
80 def getBetaStochasticGradient(train_x, train_y, lr):
81      n = train_x.shape[0] #total of data points
82      p = train_x.shape[1] #total number of attributes
83
84      beta = np.random.rand(p)
85
86      epoch = 100;
87      for iter in range(epoch):
88          indices = list(range(n))
89          rd.shuffle(indices)
90          for i in range(n):
91              idx = indices[i]
92              #=======================#
93              # STRART YOUR CODE HERE  #
94              #=======================#
95               y_i = train_y[idx]
96               x_i = train_x[idx]
97               x_i_transpose = np.transpose(x_i)
98               A = np.matmul(x_i_transpose, beta)
99               beta = beta + (lr * (y_i - A) * x_i)
100             #=======================#
101             #   END YOUR CODE HERE   #
102             #=======================#
103     return beta
104
105
106 # Linear Regression implementation
107 class LinearRegression(object):
108     # Initializes by reading data, setting hyper-parameters, and forming
   linear model
109     # Forms a linear model (learns the parameter) according to type of beta
   (0 - closed form, 1 - batch gradient, 2 - stochastic gradient)
110     # Performs z-score normalization if z_score is 1
111     def __init__(self,lr=0.005, num_iter=1000):
112         self.lr = lr
113         self.num_iter = num_iter
114         self.train_x = pd.DataFrame()
```

```python
116             self.test_x = pd.DataFrame()
117             self.test_y = pd.DataFrame()
118             self.algType = 0
119             self.isNormalized = 0
120
121     def load_data(self, train_file, test_file):
122             self.train_x, self.train_y = getDataframe(train_file)
123             self.test_x, self.test_y = getDataframe(test_file)
124
125     def normalize(self):
126             # Applies z-score normalization to the dataframe and returns a
     normalized dataframe
127             self.isNormalized = 1
128             means = self.train_x.mean(0)
129             std = self.train_x.std(0)
130             self.train_x = (self.train_x - means).div(std)
131             self.test_x = (self.test_x - means).div(std)
132
133     # Gets the beta according to input
134     def train(self, algType):
135             self.algType = algType
136             newTrain_x = addAllOneColumn(self.train_x.values) #insert an all-one
     column as the first column
137             print('Learning Algorithm Type: ', algType)
138             if(algType == '0'):
139                 beta = getBeta(newTrain_x, self.train_y.values)
140                 #print('Beta: ', beta)
141
142             elif(algType == '1'):
143                 beta = getBetaBatchGradient(newTrain_x, self.train_y.values,
     self.lr, self.num_iter)
144                 #print('Beta: ', beta)
145             elif(algType == '2'):
146                 self.lr = 0.0005  #lr to 0.0005 so that the beta does converge
147                 beta = getBetaStochasticGradient(newTrain_x, self.train_y.values,
     self.lr)
148                 #print('Beta: ', beta)
149             else:
150                 print('Incorrect beta_type! Usage: 0 - closed form solution, 1 -
     batch gradient descent, 2 - stochastic gradient descent')
151
152
153             return beta
154
155     # Predicts the y values of all test points
156     # Outputs the predicted y values to the text file named "logistic-
     regression-output_algType_isNormalized" inside "output" folder
157     # Computes MSE
158     def predict(self,x, beta):
159             newTest_x = addAllOneColumn(x)
160             self.predicted_y = newTest_x.dot(beta)
161             return self.predicted_y
162
163
164     # predicted_y and test_y are the predicted and actual y values
     respectively as numpy arrays
165     # function prints the mean squared error value for the test dataset
166     def compute_mse(self,predicted_y, y):
167             mse = np.sum((predicted_y - y)**2)/predicted_y.shape[0]
```

```
169
170
171
```