

1. Jifan Zhang Sub86474

(a)

$$C_1 : \begin{cases} +1 & x > a \\ -1 & \text{other} \end{cases} \quad C_2 : \begin{cases} +1 & x < b \\ -1 & \text{other} \end{cases} \quad C_3 : \begin{cases} +1, & x < +\infty \end{cases}$$

$a < b$, we need to judge the value of $0.5C_3(x) - C_1(x) - C_2(x)$.

$$\textcircled{1} \quad x \leq a$$

$$\textcircled{2} \quad a \leq x \leq b$$

$$\textcircled{3} \quad x \geq b$$

$$0.5C_3(x) - C_1(x) - C_2(x)$$

$$0.5C_3(x) - C_1(x) - C_2(x)$$

$$= 0.5 - (-1) - 1$$

$$= 0.5 - 1 - (-1)$$

$$= 0.5$$

$$= 0.5 - 1 - 1$$

$$= 0.5$$

when $x \leq a$, or $x \geq b$, the $f(x)$ is positive.

$f(x)$ is a threshold classifier, because the $f(x)$ can return $+1$ or -1 by the value of x #

(b)

This is because OOB error is a part of bagging method, so we don't need extra codes. And OOB error uses every data to do model training and performance evaluation, the OOB error method doesn't need to split training data and testing data #

2.

(a)

Question 2

```
In 10 1 # (a)
2 from sklearn.datasets import fetch_california_housing
3 from sklearn.model_selection import train_test_split
4
5 data = fetch_california_housing()
6 print(data.DESCR)
7
8 X_train, X_test, y_train, y_test = train_test_split(data.data, data.target,
9 test_size=0.3,
10 random_state=50)
```

Executed at 2023.10.14 12:26:28 in 73ms

.. _california_housing_dataset:

California Housing dataset

Data Set Characteristics:

:Number of Instances: 20640

:Number of Attributes: 8 numeric, predictive attributes and the target

:Attribute Information:

The value of d is 8 in this dataset

(b)

```
In 11  1 #(b)
2 from sklearn.ensemble import RandomForestRegressor
3 from sklearn.metrics import mean_squared_error
4
5 rf = RandomForestRegressor(n_estimators=100, random_state=50)
6
7 rf.fit(X_train, y_train)
8
9 y_train_pred = rf.predict(X_train)
10 y_test_pred = rf.predict(X_test)
11
12 train_mse = mean_squared_error(y_train, y_train_pred)
13 test_mse = mean_squared_error(y_test, y_test_pred)
14
15 print(f'Training MSE: {train_mse}')
16 print(f'Test MSE: {test_mse}')

```

Executed at 2023.10.14 12:26:49 in 1ds 704ms

```
Training MSE: 0.03561237258828922
Test MSE: 0.2643704261029829
```

I hope this won't cause you any inconvenience. Looking forward to your response. The value of m is 8, because we use default hyperparameters, so the max_features is 8.

#

(c)

```
In 12  1 #(c)
2 trees_mse = []
3 for tree in rf.estimators_:
4     y_pred = tree.predict(X_test)
5     trees_mse.append(mean_squared_error(y_test, y_pred))
6
7 print("MSE for Each Tree in the Test Set: ", trees_mse)
8 print("MSE for the Entire Forest on the Test Set: ", test_mse)

```

Executed at 2023.10.14 12:27:04 in 210ms

```
MSE for Each Tree in the Test Set: [0.5492299401658591, 0.5302397605532784, 0.6049794481668281, 0.602028885702487, 0.5947437957668927, 0.6430340320716892, 0.5653728307504038, 0.6013802910773901, 0.5450129434605135, 0.5501374440352228, 0.5672762809702519, 0.5622726671523094, 0.5258212810607558, 0.5702052210381137, 0.5208931283826711, 0.625478728014664, 0.5581362527012274, 0.5805645501952842, 0.594155706095898, 0.5871157760478197, 0.604312105305168, 0.5761633583569928, 0.607930533690003, 0.6093420094099483, 0.5430173128630813, 0.5709032669198805, 0.5813048339838987, 0.5401538015579296, 0.5882778607182977, 0.6104711977498384, 0.539875976260788, 0.5816350887257105, 0.5848523531307331, 0.5571811815448159, 0.5573483312449774, 0.557944039910998, 0.5419462107537951, 0.5462771541001453, 0.570012139390859, 0.5548241841713663, 0.526002003130943, 0.5752917458474484, 0.5670126634930394, 0.5813917092937984, 0.5760583168619186, 0.5499357693941375, 0.5736766030568636, 0.5451472552500322, 0.5780992297623546, 0.561666576464438, 0.5977621740194767, 0.6007376482569928, 0.5894707276436046, 0.5816620031989663, 0.554349996038049, 0.6262230803051033, 0.6056559472432815, 0.5929205906766473, 0.5886772080406977, 0.5703550322495154, 0.5453610392779231, 0.5971695462964631, 0.5952409924689761, 0.6086908845751453, 0.5603233383765504, 0.5742959016875646, 0.547251770634496, 0.5740553610058623, 0.5438541189630524, 0.5814659174075983, 0.5316846526984496, 0.5608903980372738, 0.5985906824239018, 0.5415347226368702, 0.5482515852564599, 0.5868951662474967, 0.6271884934541828, 0.5486584532405685, 0.5849934046366116, 0.5417149131185238, 0.5925992621331717, 0.5523712344748384, 0.5873022618941053, 0.5326184679602551, 0.5503006222430555, 0.572200906195091, 0.568743442099968, 0.5528136558831396, 0.6009660923685884, 0.6060363620942021, 0.563557882474968, 0.54802112889677, 0.5847901830987886, 0.5716677477661014, 0.6093378005148093, 0.5558999161929263, 0.5702055045131945, 0.5834549520683623, 0.5692988724718991, 0.5786533481697351]
MSE for the Entire Forest on the Test Set: 0.2643704261029829
```

#

(d)

```
In 13  1 #(d)
2 import numpy as np
3 from scipy.stats import pearsonr
4
5 predictions = []
6 for tree in rf.estimators_:
7     predictions.append(tree.predict(X_test))
8
9 correlations = []
10 for i in range(len(predictions)):
11     for j in range(i+1, len(predictions)):
12         corr, _ = pearsonr(predictions[i], predictions[j])
13         correlations.append(corr)
14
15 mean_corr = np.mean(correlations)
16 print("the average of all pairwise correlations: ", mean_corr)

```

Executed at 2023.10.14 12:27:18 in 5s 85ms

```
the average of all pairwise correlations: 0.7659047740859255
```

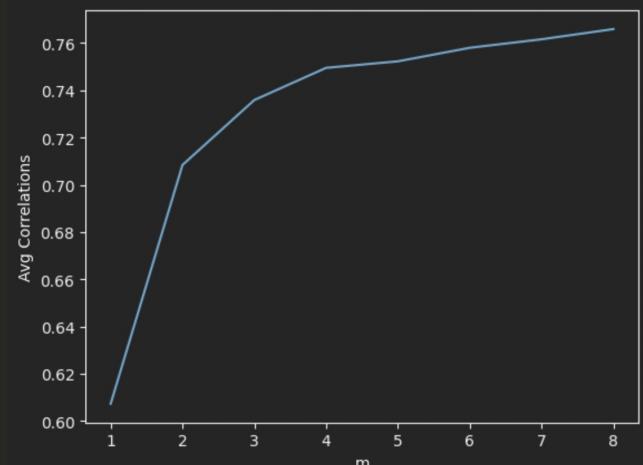
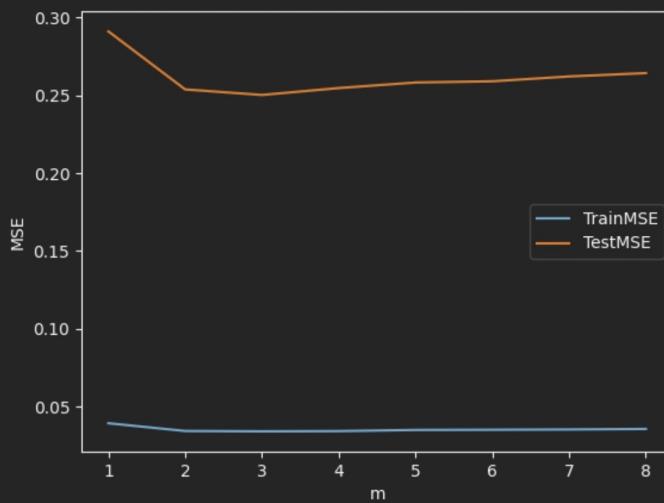
#

(2)

```
In 14  1 #(@)
2 import matplotlib.pyplot as plt
3 d = 8
4 mses_train = []
5 mses_test = []
6 mean_corrs = []
7
8 for m in range(1, d+1):
9     rf = RandomForestRegressor(n_estimators=100, max_features=m,
10                             random_state=50)
11    rf.fit(X_train, y_train)
12    y_train_pred = rf.predict(X_train)
13    y_test_pred = rf.predict(X_test)
14    mse_train = mean_squared_error(y_train, y_train_pred)
15    mse_test = mean_squared_error(y_test, y_test_pred)
16    mses_train.append(mse_train)
17    mses_test.append(mse_test)
18    predictions = []
19    for tree in rf.estimators_:
20        predictions.append(tree.predict(X_test))
21    correlations = []
22    for i in range(len(predictions)):
23        for j in range(i+1, len(predictions)):
24            corr, _ = pearsonr(predictions[i], predictions[j])
25            correlations.append(corr)
26    mean_corr = np.mean(correlations)
27    mean_corrs.append(mean_corr)
28
29 print("m\tTrain MSE\tTest MSE\tAvg Correlations")
30 for m, mse_train, mse_test, avg_corr in zip(range(1, d+1), mses_train,
31                                              mses_test, mean_corrs):
32     print(f"{m}\t{mse_train:.4f}\t{mse_test:.4f}\t{avg_corr:.4f}")
33
34
35 plt.figure()
36 plt.plot(range(1, d+1), mses_train, label='TrainMSE')
37 plt.plot(range(1, d+1), mses_test, label='TestMSE')
38 plt.xlabel('m')
39 plt.ylabel('MSE')
40 plt.legend()
41 plt.show()
42
43
44 plt.figure()
45 plt.plot(range(1, d+1), mean_corrs)
46 plt.xlabel('m')
47 plt.ylabel('Avg Correlations')
48 plt.show()
```

Executed at 2023.10.14 12:28:57 in 1m 33s 998ms

m	Train MSE	Test MSE	Avg Correlations
1	0.0393	0.2912	0.6073
2	0.0343	0.2539	0.7084
3	0.0341	0.2503	0.7359
4	0.0342	0.2548	0.7495
5	0.0350	0.2583	0.7523
6	0.0351	0.2591	0.7580
7	0.0353	0.2622	0.7616
8	0.0356	0.2644	0.7659



(f)

(f)

The average correlation increases as m increases, this is because in a random forest, the m determines number of features in each splitting node. When m increases the correlation becomes bigger, it means their predictions increases.

(g)

(g)

False, When we choose m , we need to consider the m that corresponds to the smallest MSE value, rather than choosing the smallest m value directly. Only choosing the m value with the smallest MSE will indicate the best fit.

#

3.

We have $L_\lambda(w) = L(w) + \frac{1}{2}\lambda\|w\|^2$. we need to calculate gradient of $L_\lambda(w)$

$$\begin{aligned} \nabla L_\lambda(w) &= \nabla L(w) + \nabla \frac{1}{2}\lambda\|w\|^2 \\ &= \nabla L(w) + \lambda w \end{aligned}$$

we can get $\nabla L_\lambda(w) - \nabla L(w) = \lambda w$

$$\text{by } w_{\text{new}} = (I - \eta \lambda)w - \eta \nabla L(w)$$

we can proof it.

#

$$w_{\text{new}} = w - \eta \nabla L_\lambda(w)$$

$$= w - \eta (\nabla L(w) + \lambda w)$$

$$= w - \eta \nabla L(w) - \eta \lambda w$$

$$= (I - \eta \lambda)w - \eta \nabla L(w)$$

4.

(b)

```
def predict_proba(self, X):
    """
    Predict the class distributions for given input examples.

    Parameters
    -----
    X: input examples, represented as an input array of shape (n_sample,
        n_features).

    Returns
    -----
    y: predicted class distributions, represented as an array of shape (n_sample,
        n_classes)
    """

    # Replace pass with your code
    proba = self.predict_proba(X)
    return self.classes_[np.argmax(proba, axis=1)]
```

Modify the code below to implement the idea for avoiding numerical flow as described in Q4 (b)

```
scores = X @ self.coef_.T + self.intercept_
scores -= np.max(scores, axis=1).reshape(-1, 1)
scores = np.exp(scores)
return scores / scores.sum(axis=1).reshape(-1, 1)
```

#

(C)

```

14 def fit(self, X, y, lr=0.1, momentum=0, niter=100):
15     """
16     Train a multiclass logistic regression model on the given training set.
17
18     Parameters
19     -----
20     X: training examples, represented as an input array of shape (n_sample,
21     |   n_features).
22     y: labels of training examples, represented as an array of shape
23     |   (n_sample,) containing the classes for the input examples
24     lr: learning rate for gradient descent
25     niter: number of gradient descent updates
26     momentum: the momentum constant (see assignment task sheet for an explanation)
27
28     Returns
29     -----
30     self: fitted model
31
32     self.classes_ = np.unique(y)
33     self.class2int = dict((c, i) for i, c in enumerate(self.classes_))
34     y = np.array([self.class2int[c] for c in y])
35
36     n_features = X.shape[1]
37     n_classes = len(self.classes_)
38
39     self.intercept_ = np.zeros(n_classes)
40     self.coef_ = np.zeros((n_classes, n_features))
41
42     # Implement your gradient descent training code here; uncomment the code below to do "random training"
43     # self.intercept_ = np.random.randn(*self.intercept_.shape)
44     # self.coef_ = np.random.randn(*self.coef_.shape)
45     for iteration in range(niter):
46         softmax = self.predict_proba(X)
47
48         # Compute the gradient of the log-loss
49         error = softmax - np.eye(n_classes)[y]
50         grad_coef = error.T @ X
51         grad_intercept = np.sum(error, axis=0)
52
53         # Update model parameters using gradient descent
54         self.coef_ -= lr * grad_coef
55         self.intercept_ -= lr * grad_intercept
56
57     return self
58
59
60 if __name__ == '__main__':
61     X, y = fetch_covtype(return_X_y=True)
62     X_tr, X_ts, y_tr, y_ts = train_test_split(X, y, test_size=0.3, random_state=42)
63
64     clf = LogisticRegression()
65     clf.fit(X_tr, y_tr)
66     print(accuracy_score(y_tr, clf.predict(X_tr)))
67     print(accuracy_score(y_ts, clf.predict(X_ts)))

```

Executed at 2023.10.14 12:29:33 in 15s 99ms

```

0.479938924240487
0.4800463557921792

```

#

Best learning rate: 0.0001
Best number of iterations: 800
Best average log-loss: 20.7053
Training accuracy: 0.4753
Test accuracy: 0.4745

(d)

```

102 if __name__ == '__main__':
103     X, y = fetch_covtype(return_X_y=True)
104     X_tr, X_ts, y_tr, y_ts = train_test_split(X, y, train_size=0.3,
105                                              test_size=0.5, random_state=42)
106
107     learning_rates = [0.0001, 0.001, 0.01, 0.1]
108     iterations = [50, 100, 500, 800]
109
110     best_log_loss = float('inf')
111     best_lr = None
112     best_niter = None
113     best_model = None
114
115     for lr in learning_rates:
116         for niter in iterations:
117             clf = LogisticRegression()
118             clf, log_losses = clf.fit(X_tr, y_tr, lr=lr, niter=niter)
119             avg_log_loss = np.mean(log_losses)
120             if avg_log_loss < best_log_loss:
121                 best_log_loss = avg_log_loss
122                 best_lr = lr
123                 best_niter = niter
124                 best_model = clf
125
126     print(f"Best learning rate: {best_lr}")
127     print(f"Best number of iterations: {best_niter}")
128     print(f"Best average log-loss: {best_log_loss:.4f}")
129     print(f"Training accuracy: {accuracy_score(y_tr, best_model.predict(X_tr)):.4f}")
130     print(f"Test accuracy: {accuracy_score(y_ts, best_model.predict(X_ts)):.4f}")

```

Executed at 2023.10.14 12:39:05 in 9m 26s 919ms

In order to find the optimal pair of hyperparameters, I have experimented with four different learning rates and four different numbers of iterations. This approach involves varying both the learning rate and the number of iterations to identify the best combination.

BC- ηg_i)

(e)

↑

$$\text{when } t=2, \quad W_{2+1} = W_3 = W_2 - \eta g_2 + \beta(W_2 - W_1)$$

$$= W_2 - \eta g_2 - \beta \eta g_1$$

$$= W_2 - \eta(g_2 + \beta g_1)$$

$$\text{we assume } W_{t+1} = W_t - \eta(g_t + \beta g_{t-1} + \dots + \beta^{t-1} g_1)$$

$$W_{t+2} = W_{t+1} - \eta g_{t+1} + \beta(W_{t+1} - W_t)$$

$$= [W_t - \eta(g_t + \beta g_{t-1} + \dots + \beta^{t-1} g_1)] - \eta g_{t+1} + \beta([W_t - \eta(g_t + \beta g_{t-1} + \dots + \beta^{t-1} g_1)] - W_t)$$
$$= W_t - \eta(g_t + \beta_{t-1} + \dots + \beta^{t-1} g_1) - \eta g_{t+1} - \beta \eta(g_t + \beta g_{t-1} + \dots + \beta^{t-1} g_1)$$

$$= W_{t+1} - \eta(g_{t+1} + \beta(g_t + \dots + \beta^{t-1} g_1))$$

$$= W_{t+1} - \eta(g_{t+1} + \beta g_t + \dots + \beta^t g_1)$$

by mathematical induction we can proof $W_{t+1} = W_t - \eta(g_t + \beta g_{t-1} + \dots + \beta^{t-1} g_1)$ #.

(f)

```

12 def fit(self, X, y, lr=0.0001, momentum=0.9, niter=800):
13     best_log_loss = float('inf')
14
15     log_loss_list = []
16
17     self.classes_ = np.unique(y)
18     self.class2int = dict((c, i) for i, c in enumerate(self.classes_))
19     y = np.array([self.class2int[c] for c in y])
20
21     n_features = X.shape[1]
22     n_classes = len(self.classes_)
23
24     self.intercept_ = np.zeros(n_classes)
25     self.coef_ = np.zeros((n_classes, n_features))
26
27     X_tensor = torch.tensor(X, dtype=torch.float32)
28     y_tensor = torch.tensor(y, dtype=torch.long)
29
30     self.model = torch.nn.Linear(n_features, n_classes)
31
32     optimizer = optim.SGD(self.model.parameters(), lr=lr, momentum=momentum)
33
34     loss_fn = torch.nn.CrossEntropyLoss()
35
36     for iteration in range(niter):
37         optimizer.zero_grad()
38         output = self.model(X_tensor)
39         loss = loss_fn(output, y_tensor)
40         loss.backward()
41         optimizer.step()
42
43         with torch.no_grad():
44             y_pred_proba = torch.softmax(output, dim=1)
45             training_log_loss = log_loss(y_tensor.numpy(), y_pred_proba.numpy(), labels=np.arange(n_classes))
46             log_loss_list.append(training_log_loss)
47
48             if training_log_loss < best_log_loss:
49                 best_log_loss = training_log_loss
50
51     avg_log_loss = sum(log_loss_list) / len(log_loss_list)
52
53     with torch.no_grad():
54         self.intercept_ = self.model.bias.numpy()
55         self.coef_ = self.model.weight.numpy()
56
57         self.training_log_loss = avg_log_loss |
58         return momentum, avg_log_loss

```

```

102 if __name__ == '__main__':
103     X, y = fetch_covtype(return_X_y=True)
104     X_tr, X_ts, y_tr, y_ts = train_test_split(X, y, test_size=0.3, random_state=42)
105
106     best_overall_momentum = None
107     best_overall_log_loss = float('inf')
108     best_overall_train_acc = 0.0
109     best_overall_test_acc = 0.0
110
111     clf = LogisticRegression()
112
113     # Loop through different momentum values
114     for momentum in [0.1, 0.5, 0.9, 0.99]: # Example momentum values
115         best_momentum, best_log_loss = clf.fit(X_tr, y_tr, lr=0.0001, momentum=momentum, niter=800)
116
117         if best_log_loss < best_overall_log_loss:
118             best_overall_log_loss = best_log_loss
119             best_overall_momentum = best_momentum
120
121     # Compute train and test accuracies
122     train_preds = clf.predict(X_tr)
123     test_preds = clf.predict(X_ts)
124     train_acc = accuracy_score(y_tr, train_preds)
125     test_acc = accuracy_score(y_ts, test_preds)
126
127     if train_acc > best_overall_train_acc:
128         best_overall_train_acc = train_acc
129     if test_acc > best_overall_test_acc:
130         best_overall_test_acc = test_acc
131
132     print(f"Best Momentum: {best_overall_momentum}, Minimum Training Log-Loss: {best_overall_log_loss}")
133     print(f"Train Accuracy: {train_acc}, Test Accuracy: {test_acc}")
134
135     print(f"Overall Best Train Accuracy: {best_overall_train_acc}, Overall Best Test Accuracy: {best_overall_test_acc}")
136

```

Executed at 2023.10.14 10:48:11 in 8m 23s 280ms

Momentum: 0.1, Minimum Training Log-Loss: 9.490934452604215
 Train Accuracy: 0.4761524238520019, Test Accuracy: 0.47596727556453094

Momentum: 0.5, Minimum Training Log-Loss: 9.454841827561578
 Train Accuracy: 0.36372286751182664, Test Accuracy: 0.36407655590233157

Momentum: 0.9, Minimum Training Log-Loss: 9.333311641419847
 Train Accuracy: 0.33188675905071946, Test Accuracy: 0.33361253901230037

✓ Momentum: 0.99, Minimum Training Log-Loss: 8.70195239409016
 Train Accuracy: 0.5938978333349725, Test Accuracy: 0.5937557371029925
 Best Momentum: 0.99, Overall Minimum Training Log-Loss: 8.70195239409016
 Overall Best Train Accuracy: 0.5938978333349725, Overall Best Test Accuracy: 0.5937557371029925

By (d) when got the best learning rate and iterator number
 So we just use different momentum to generate result.

(g)

```

15 def fit(self, X, y, lr=0.1, momentum=0, niter=100):
16     """
17         Train a multiclass logistic regression model on the given training set.
18
19     Parameters
20     -----
21     X: training examples, represented as an input array of shape (n_sample,
22     |   n_features).
23     y: labels of training examples, represented as an array of shape
24     |   (n_sample,) containing the classes for the input examples
25     lr: learning rate for gradient descent
26     niter: number of gradient descent updates
27     momentum: the momentum constant (see assignment task sheet for an explanation)
28
29     Returns
30     -----
31     self: fitted model
32     """
33     np.random.seed(50)
34     self.classes_ = np.unique(y)
35     self.class2int = dict((c, i) for i, c in enumerate(self.classes_))
36     y = np.array([self.class2int[c] for c in y])
37
38     n_features = X.shape[1]
39     n_classes = len(self.classes_)
40
41     self.intercept_ = np.zeros(n_classes)
42     self.coef_ = np.zeros((n_classes, n_features))
43
44     log_losses = []
45     # Implement your gradient descent training code here; uncomment the code below to do "random training"
46     self.intercept_ = np.random.randn(*self.intercept_.shape)
47     self.coef_ = np.random.randn(*self.coef_.shape)
48     for iteration in range(niter):
49         softmax = self.predict_proba(X)
50
51
52         error = softmax - np.eye(n_classes)[y]
53         grad_coef = error.T @ X
54         grad_intercept = np.sum(error, axis=0)
55
56         adjusted_softmax = softmax + 1e-15
57         log_loss = -np.log(adjusted_softmax[range(len(y)), y]).mean()
58         log_losses.append(log_loss)
59
60         self.coef_ -= lr * grad_coef
61         self.intercept_ -= lr * grad_intercept
62     self.log_losses = log_losses
63     return self, self.log_losses
64
65 if __name__ == '__main__':
66     X, y = fetch_covtype(return_X_y=True)
67
68     scaler = StandardScaler()
69     X = scaler.fit_transform(X)
70
71     X_tr, X_ts, y_tr, y_ts = train_test_split(X, y, train_size=0.3,
72                                             test_size=0.5, random_state=42)
73
74     learning_rates = [0.001, 0.01, 0.1, 1]
75     iterations = [50, 100, 500, 1000]
76
77     best_log_loss = float('inf')
78     best_lr = None
79     best_niter = None
80     best_model = None
81
82     for lr in learning_rates:
83         for niter in iterations:
84             clf = LogisticRegression()
85             clf, log_losses = clf.fit(X_tr, y_tr, lr=lr, niter=niter)
86             avg_log_loss = np.mean(log_losses)
87             if avg_log_loss < best_log_loss:
88                 best_log_loss = avg_log_loss
89                 best_lr = lr
90                 best_niter = niter
91                 best_model = clf
92
93     print(f"Best learning rate: {best_lr}")
94     print(f"Best number of iterations: {best_niter}")
95     print(f"Best average log-loss: {best_log_loss:.4f}")
96     print(f"Training accuracy: {accuracy_score(y_tr, best_model.predict(X_tr)):.4f}")
97     print(f"Test accuracy: {accuracy_score(y_ts, best_model.predict(X_ts)):.4f}")

```

Executed at 2023.10.14 12:49:21 from 10.5.83.715

Best learning rate: 0.001
Best number of iterations: 1000
Best average log-loss: 12.3414
Training accuracy: 0.5955
Test accuracy: 0.5921

It is evident that by standardizing the features, i.e., scaling them to a mean of 0 and a standard deviation of 1, the accuracy of training and testing is significantly improved. The best models trained on data that have been normalized for features exhibit significantly lower mean log-loss compared to models trained on un-normalized data.

#