

Lab Course Machine Learning

Exercise Sheet 5

December 5th, 2021

Syed Wasif Murtaza Jaffri-311226

Exercise 1: Backward search for variable selection

```
In [23]: import numpy as np
import math
import pandas as pd
import matplotlib.pyplot as plt
from matplotlib import cm
from matplotlib.ticker import LinearLocator, FormatStrFormatter
from sklearn.metrics import r2_score, mean_squared_error
import warnings
import itertools
warnings.filterwarnings('ignore')

In [24]: df = pd.read_csv('bank.csv', sep = ';')
df.columns

Out [24]: Index(['age', 'job', 'marital', 'education', 'default', 'balance', 'housing', 'loan', 'contact', 'day', 'month', 'duration', 'campaign', 'pdays', 'previous', 'outcome', 'y'],
      dtype='object')
```

If required drop out the rows with missing values or NA. (hope) in next lectures we will handle sparse data, which will allow us to use records with missing values.

```
In [25]: is_nan = df.isnull()
row_has_NaN = is_nan.any(axis=1)
rows_with_NaN = df[row_has_NaN]
rows_with_NaN

Out [25]:
```

age	job	marital	education	default	balance	housing	loan	contact	day	month	duration	campaign	pdays	previous	outcome	y
30	unemployed	married	primary	no	1787	no	no	cellular	19	oct	79	1	-1	-1	0	0
3	blue-collar	married	tertiary	no	1476	yes	yes	unknown	3	jun	199	4	-1	-1	0	0
4	59	manager	married	secondary	no	0	yes	no	unknown	5	may	226	1	-1	-1	0
7	39	technician	married	secondary	no	147	yes	no	cellular	6	may	151	2	-1	-1	0
8	41	entrepreneur	married	tertiary	no	221	yes	no	unknown	14	may	57	2	-1	-1	0
...
4513	49	blue-collar	married	secondary	no	322	no	no	cellular	14	aug	356	2	-1	-1	0
4515	32	services	single	secondary	no	473	yes	no	cellular	7	jul	324	5	-1	-1	0
4516	33	services	married	secondary	no	-333	yes	no	cellular	30	jul	629	5	-1	-1	0
4517	57	self-employed	married	tertiary	yes	-3313	yes	yes	unknown	9	may	153	1	-1	-1	0
4518	57	technician	married	secondary	no	295	no	no	cellular	19	aug	151	11	-1	-1	0

357 rows x 17 columns

1 Convert any non-numeric values to numeric values. For example you can replace a country name with an integer value or more appropriately use hot-one encoding.

transformData() takes a dataframe and for each column gets unique values in that column and assign a value from zero to number of possible values in that column. Finally, applying this dictionary for one column to next columns. In this way it creates a hashmap for dataframe.

```
In [27]: def transformData(df):
    cols = df.columns
    num_cols = df.get_numeric_data().columns
    colList = list(set(cols) - set(num_cols))
    dictMap = dict()
    for c in colList:
        values = df[c].unique()
        keys = list(range(0, len(keys)))
        dictionary = dict(zip(keys, values))
        dictMap.update(dictionary)
    df = df.replace(dictMap)
    return df, dictMap

In [28]: ds, dictMap = transformData(df)

In [29]: dictMap

Out [29]: {'married': 0,
'single': 1,
'divorced': 2,
'married': 0,
'yes': 1,
'cellular': 0,
'unknown': 3,
'telephone': 2,
'failure': 1,
'other': 2,
'success': 3,
'unemployed': 0,
'services': 1,
'management': 2,
'blue-collar': 3,
'self-employed': 4,
'entrepreneur': 5,
'admin.': 7,
'student': 8,
'housemaid': 9,
'retired': 10,
'primary': 6,
'secondary': 1,
'tertiary': 2,
'oct': 0,
'may': 1,
'apr': 2,
'jun': 3,
'feb': 4,
'aug': 5,
'jan': 6,
'jul': 7,
'nov': 8,
'sep': 9,
'mar': 10,
'dec': 11}
```

4 Normalize (Standardize) the data

```
In [30]: day = ds.loc[:, ds.columns != 'y']
nds = (day-day.mean())/day.std()
nds

Out [30]:
```

age	job	marital	education	default	balance	housing	loan	contact	day	month	duration	campaign	
0	1.056153	-1.595307	-0.716234	-1.644573	-0.130744	0.121058	-1.141925	-0.424709	-0.738102	0.374011	-1.444582	-0.711782	-0.5767
1	-0.772489	-1.199305	-0.716234	-0.309004	-0.130744	1.118521	2.354032	-0.738102	-0.595961	-1.076641	-0.169175	-0.5767	
2	0.583394	-0.803303	0.721640	1.026566	-0.130744	-0.024142	0.875521	2.354032	-0.738102	0.1010271	-0.708700	-0.303865	-0.5767
3	1.056153	-0.803303	-0.716234	1.026566	-0.130744	0.017724	0.875521	2.354032	1.450546	-1.565932	-0.340759	-0.249989	0.3879
4	1.685850	-0.407301	-0.716234	-0.309004	-0.130744	-0.472701	0.875521	-0.424709	1.450546	-1.323439	-1.076641	-0.146086	-0.5767
...
4516	-0.772489	-1.199305	-0.716234	-0.309004	-0.130744	-0.583345	0.875521	-0.424709	-0.738102	1.707721	1.131006	0.250287	0.7094
4517	1.496746	-0.011299	-0.716234	1.026566	0.7646823	-1.573497	0.875521	2.354032	1.450546	-0.838453	-1.076641	-0.427010	-0.5767
4518	1.496746	0.384703	-0.716234	-0.309004	-0.130744	-0.374682	-1.141925	-0.424709	-0.738102	0.374011	0.395124	-0.434706	2.6388
4519	-1.248256	-0.407301	-0.716234	-0.309004	-0.130744	-0.094914	-1.141925	-0.424709	-0.738102	-1.202193	0.027183	-0.519368	-0.3879
4520	0.267573	0.780704	0.721640	1.026566	-0.130744	-0.095247	0.875521	2.354032	-0.738102	-1.565932	-0.708700	0.311859	-0.2552

4521 rows x 16 columns

3 Split the data into a train/test splits according to the ratios 80%:20%.

```
In [33]: nds_train = nds.loc[0:math.floor(len(nds)*0.8)]
nds_test = nds.loc[math.floor(len(nds)*0.8)+1:]
len(nds) == len(nds_train) + len(nds_test)

Out [33]:
```

age	job	marital	education	default	balance	housing	loan	contact	day	month	duration	campaign	
0	1.056153	-1.595307	-0.716234	-1.644573	-0.130744	0.121058	-1.141925	-0.424709	-0.738102	0.374011	-1.444582	-0.711782	-0.5767
1	-0.772489	-1.199305	-0.716234	-0.309004	-0.130744	1.118521	2.354032	-0.738102	-0.595961	-1.076641	-0.169175	-0.5767	
2	0.583394	-0.803303	0.721640	1.026566	-0.130744	-0.024142	0.875521	2.354032	-0.738102	0.1010271	-0.708700	-0.303865	-0.5767
3	1.056153	-0.803303	-0.716234	1.026566	-0.130744	0.017724	0.875521	2.354032	1.450546	-1.565932	-0.340759	-0.249989	0.3879
4	1.685850	-0.407301	-0.716234	-0.309004	-0.130744	-0.472701	0.875521	-0.424709	1.450546	-1.323439	-1.076641	-0.146086	-0.5767
...
4516	-0.772489	-1.199305	-0.716234	-0.309004	-0.130744	-0.583345	0.875521	-0.424709	-0.738102	1.707721	1.131006	0.250287	0.7094
4517	1.496746	-0.011299	-0.716234	1.026566	0.7646823	-1.573497	0.875521	2.354032	1.450546	-0.838453	-1.076641	-0.427010	-0.5767
4518	1.496746	0.384703	-0.716234	-0.309004	-0.130744	-0.374682	-1.141925	-0.424709	-0.738102	0.374011	0.395124	-0.434706	2.6388
4519	-1.248256	-0.407301	-0.716234	-0.309004	-0.130744	-0.094914	-1.141925	-0.424709	-0.738102	-1.202193	0.027183	-0.519368	-0.3879
4520	0.267573	0.780704	0.721640	1.026566	-0.130744	-0.095247	0.875521	2.354032	-0.738102	-1.565932	-0.708700	0.311859	-0.2552

904 rows x 17 columns

5 Implement logistic regression and mini-batch Gradient Ascent.

The implementation of Gradient Ascent and logistic regression is same as last lab, for mini batches I am splitting data in N number of batches of size 50, and for each batch calling Gradient Ascent Func(). Beta returned from last batch will be used as initial Beta for the next batch.

```
In [34]: def sig(X, B):
    return 1 / (1 + np.exp(-(1)*np.dot(X, B)))

def logLike(X, Y, B):
    return np.sum(Y*sig(X, B)) - np.log(1* np.exp(sig(X, B)))
# For logloss: https://ml-cheatsheet.readthedocs.io/en/latest/loss_functions.html
# ref given in exercise was not running
def logloss(yHat, y):
    sum0 = 0
    for i in range(len(yHat)):
        if y[i] == 1:
            sum = sum + np.log(yHat[i])
        else:
            sum = sum + np.log(1 - yHat[i])
    return sum

def miniBatch(X, Y, batchSize):
    totalRecords = len(X)
    numberOfBatches = np.ceil(totalRecords/ batchSize)
    X_mini = np.array_split(X, numberOfBatches)
    Y_mini = np.array_split(Y, numberOfBatches)
    B = np.zeros(shape=(len(X[0]), 1))
    for i in range (len(X_mini)):
        B = GradAscent(X_mini[i], Y_mini[i])
    return B

def aicScore(X, Y, B):
    return (-2)*(logLike(X, Y, B)) + 2*(len(B))

def mapProbability(Y, prob):
    y = np.zeros(shape=(len(Y),))
    for i in range(len(Y)):
        if Y[i] > 0.5:
            y[i] = 1
        else:
            y[i] = 0
    return y

def errorEval(Y_pred, Y_test):
    error = 0
    for i in range(len(Y_pred)):
        if Y_pred[i] != Y_test[i]:
            error += 1
    return error/len(Y_pred)

def GradAscent(X, Y):
    B_old = np.zeros(shape=(len(X[0]), 1))
    fold_min_fnew_list = []
    numberIterations = 100
    log_loss = []
    for i in range (numberIterations):
        y_hat = sig(X, B_old)
        B_hat = B_old + (1/n)* np.dot(X.T, (Y-y_hat))
        L_old = L
        L = logLike(X, Y, B_hat)
        B_old = B_hat
        return B_hat

def kFoldCrossValidation(X, Y, k):
    X_splits = np.array_split(X, k)
    Y_splits = np.array_split(Y, k)
    for i in range (len(Y_splits)):
        print(len(Y_splits[i]))

In [35]: # For AIC when V column is dropped
X_without_v = Xtrain.loc[:, Xtrain.columns != 'V'].to_numpy()
bias_column = np.ones(shape=(len(X_without_v), 1))
X_used = np.append(bias_column, X_used, axis=1)
B_used = miniBatch(X_used, Ytrain, to_numpy(), batchSize)
aic_used = aicScore(X_used, Ytrain, to_numpy(), B_used)
for V in Vused:
    # For AIC when V column is dropped
    X_without_v = Xtrain.loc[:, Xtrain.columns != 'V'].to_numpy()
    bias_column = np.ones(shape=(len(X_without_v), 1))
    X_without_v = np.append(bias_column, X_without_v, axis=1)
    B_without_v = miniBatch(X_without_v, Ytrain, to_numpy(), batchSize)
    aic_without_v = aicScore(X_without_v, Ytrain, to_numpy(), B_without_v)
    gain = aic_used - aic_without_v
    if (gain > gainBest):
        gainBest = gain
        v_best = V
    if improvement:
        Vused.remove(v_best)
    return Vused

In [36]: X_train = nds_train.loc[:, nds_train.columns != 'y']
Y_train = nds_train['y']
params = backwardSearch(X_train, Y_train)
X_train = nds_train[params]
Y_train = nds_train['y']
X_test = (nds_train[params]).to_numpy()
bias_column = np.ones(shape=(len(X_test), 1))
Y_test = np.append(bias_column, X_test, axis=1)
Y_test = (nds_train['y']).to_numpy()
Y = X_train.to_numpy()
bias_column = np.ones(shape=(len(X_train), 1))
X = np.append(bias_column, X, axis=1)
B = np.zeros(len(X[0]))
B_hat = miniBatch(X, Y, 50)

Report the final error on Test set.
```

```
In [37]: Y_hat = sig(X_test, B_hat)
Y_pred = mapProbability(Y_hat)
error = errorEval(Y_pred, Y_test)
print (error)
11.169477467514515
```

Exercise 2: Regularization for Logistic Regression

1. Pick a range of α_0 and λ defined on grid. You can choose fixed batchSize

```
In [162]: mu = [10**(i-4), 10**(i-7), 10**(i-10), 10**(i-14), 10**(i-17)]
lambda = [np.exp(-4), np.exp(-7), np.exp(-10), np.exp(-14), np.exp(-17)]

Out [162]: [0.0001, 1e-07, 1e-10, 1e-14, 1e-17]

In [163]: lambda = [np.exp(-4), np.exp(-7), np.exp(-10), np.exp(-14), np.exp(-17)]

Out [163]: [0.0181563888873418, 0.002478752716663585, 0.00033946262790251185, 4.539992976248485e-05]

In [164]: grid = list(itertools.product(mu, lambda))

In [165]: def accuracyValue(Y_pred, Y_test):
    acc = 0
    for i in range(len(Y_pred)):
        if Y_pred[i] == Y_test[i]:
            acc += 1
    return (acc/len(Y_pred))*100

def logLike(X, Y, B, lambda):
    return np.sum(Y*sig(X, B)) - np.log(1* np.exp(sig(X, B))) - (lambda*(np.linalg.norm(B))**2)

def miniBatch(X, Y, batchSize, mu, lambda, numberIterations):
    totalRecords = len(X)
    numberOfBatches = np.ceil(totalRecords/ batchSize)
    X_mini = np.array_split(X, numberOfBatches)
    Y_mini = np.array_split(Y, numberOfBatches)
    B = np.zeros(shape=(len(X[0]), 1))
    for i in range (len(X_mini)):
        B = GradAscent(X_mini[i], Y_mini[i], mu, lambda, numberIterations)
    return B

def logloss(yHat, y):
    sum0 = 0
    for i in range(len(yHat)):
        if y[i] == 1:
            sum = sum + np.log(yHat[i])
        else:
            sum = sum + np.log(1 - yHat[i])
    return sum

def GradAscent(X, Y, mu, lambda, numberIterations):
    B_old = np.zeros(shape=(len(X[0]), 1))
    fold_min_fnew_list = []
    numberIterations = 100
    log_loss = []
    for i in range (numberIterations):
        y_hat = sig(X, B_old)
        B_hat = B_old + (1/n)* np.dot(X.T, (Y-y_hat)) - (2*lambda*B_old)
        L_old = L
        L = logLike(X, Y, B_hat)
        B_old = B_hat
        return B_hat
```

2. Implement k-fold cross-validation protocol for search.

For k-folds dividing whole dataset into equal 5 parts and in one iteration taking one part as test and other 4 parts as training. Finding Accuracy for each iteration.

```
In [223]: def kFoldCrossValidation(X, Y, k, mu, lambda, numberIterations, batchSize):
    chunkSize = int(np.floor(len(X)/k))
    meanAcc = 0
    accLogList = []
    for i in range (k):
        X_test = X[i*chunkSize:(i+1)*chunkSize].to_numpy()
        Y_test = Y[i*chunkSize:(i+1)*chunkSize].to_numpy()
        index = list(range(i*chunkSize, (i+1)*chunkSize))
        X_train = np.delete(X, index, 0)
        Y_train = np.delete(Y, index, 0)
        B = miniBatch(X_train, Y_train, batchSize, mu, lambda, numberIterations)
        # For training data
        Y_hat_train = sig(X_train, B)
        Y_pred_train = mapProbability(Y_hat_train)
        train_accuracy = accuracyValue(Y_pred_train, Y_train)
        logLossTrain = logloss(Y_hat_train, Y_train)
        # For test data
        Y_hat = sig(X_test, B)
        Y_pred = mapProbability(Y_hat)
        meanAcc = accuracyValue(Y_pred, Y_test)
        test_accuracy = accuracyValue(Y_pred, Y_test)
        logLossTest = logloss(Y_hat, Y_test)
        accLogList.append([meanError, train_accuracy, test_accuracy, logLossTrain[0], logLossTest[0]])
    return (meanAcc(k), np.array(accLogList))

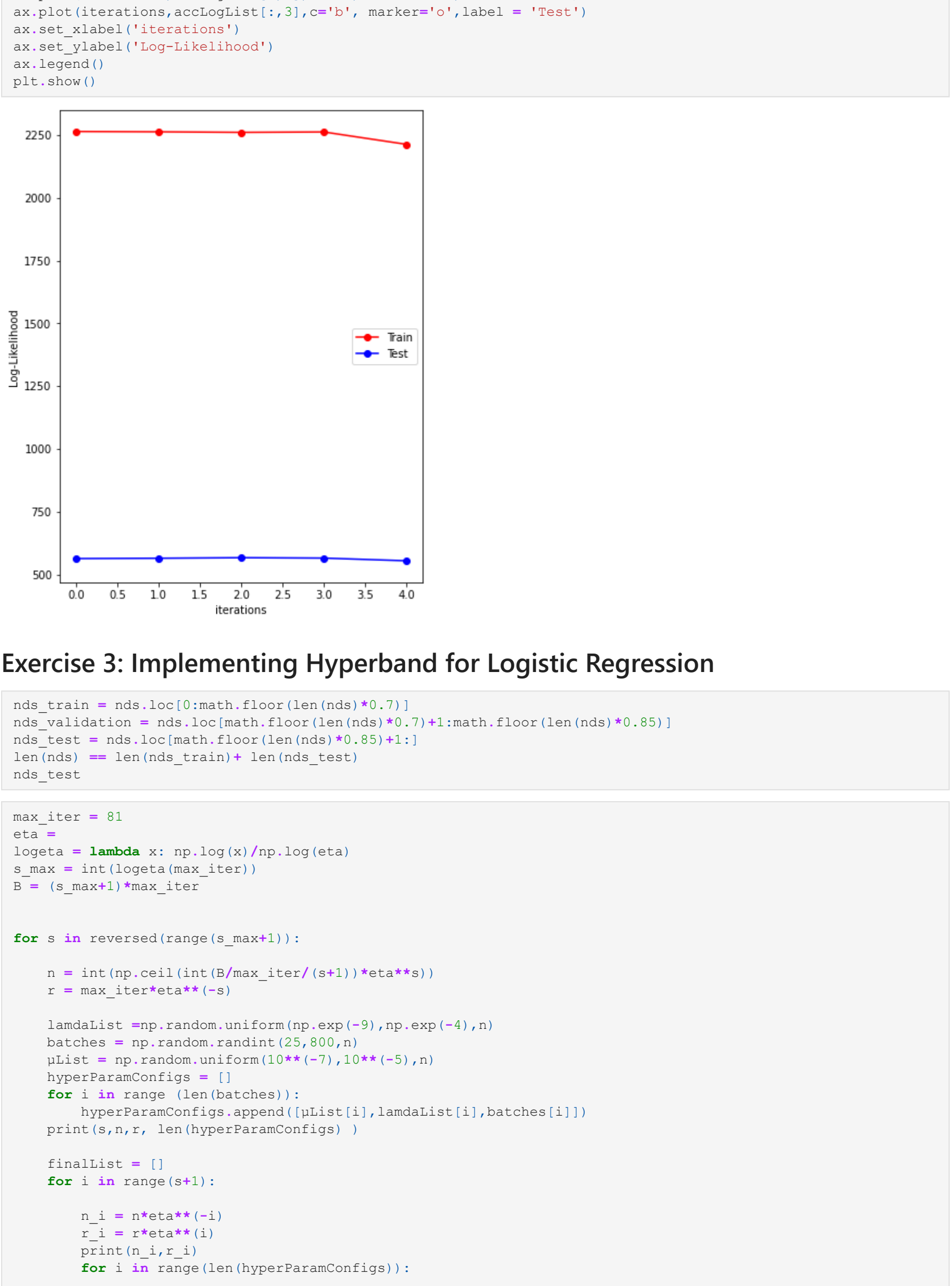
In [168]: X = nds.loc[:, nds.columns != 'y'].to_numpy()
Y = nds['y'].to_numpy()
bias_column = np.ones(shape=(len(X), 1))
X = np.append(bias_column, X, axis=1)
listMeanError = []
i=0
numberIterations=100
meanList = []
for mu, lambda in grid:
    meanError, accLogList = kFoldCrossValidation(X, Y, 5, mu, lambda, numberIterations)
    meanList.append(list([meanError, mu, lambda]))

In [169]: print(max(meanList))
[86.06194690265487, 0.0001, 0.0181563888873418]

In [170]: print(np.array(meanList))
[[ 86.0619469e+01 1.00000000e-04 1.83156389e-02]
 [ 86.0619469e+01 1.00000000e-04 2.47875218e-03]
 [ 86.0619469e+01 1.00000000e-04 3.35462628e-04]
 [ 86.0619469e+01 1.00000000e-04 4.53999298e-05]
 [ 86.0619469e+01 1.00000000e-04 3.35462628e-02]
 [ 86.0619469e+01 1.00000000e-07 2.47875218e-03]
 [ 86.0619469e+01 1.00000000e-07 3.35462628e-04]
 [ 86.0619469e+01 1.00000000e-07 4.53999298e-05]
 [ 86.0619469e+01 1.00000000e-10 2.47875218e-03]
 [ 86.0619469e+01 1.00000000e-10 3.35462628e-04]
 [ 86.0619469e+01 1.00000000e-10 4.53999298e-05]
 [ 86.0619469e+01 1.00000000e-14 2.47875218e-03]
 [ 86.0619469e+01 1.00000000e-14 3.35462628e-04]
 [ 86.0619469e+01 1.00000000e-14 4.53999298e-05]
 [ 86.0619469e+01 1.00000000e-17 2.47875218e-03]
 [ 86.0619469e+01 1.00000000e-17 3.35462628e-04]
 [ 86.0619469e+01 1.00000000e-17 4.53999298e-05]]

3. Keep track of mean performance (i.e. Classification Accuracy value) across k – folds for each set of hyperparameters.
```

```
In [171]: fig = plt.figure(figsize=(14,10))
ax = plt.axes(projection='3d')
ax.view_init(elev=10., azim=80.)
ax.scatter(np.array(meanList[:,1]), np.array(meanList[:,2]), np.array(meanList[:,0]), c='r', marker='o')
ax.set_xlabel('mu')
ax.set_ylabel('lambda')
ax.set_zlabel('Accuracy Perc')
plt.show()
```



4. Finally, for the optimal value of alpha0 and λ , train your model on complete training data and evaluate on Test data.

```
In [172]: muBest = max(meanList[:,1])
lambdaBest = max(meanList[:,2])

In [173]: print(muBest)
0.0001

In [174]: print(lambdaBest)
0.0181563888873418

In [175]: X_train = nds_train.loc[:, nds_train.columns != 'y']
Y_train = nds_train['y']
X_test = nds_train[params]
Y_test = nds_train['y']
X = X_train.to_numpy()
Y = Y_train.to_numpy()
bias_column = np.ones(shape=(len(X_train), 1))
X = np.append(bias_column, X, axis=1)
B = np.zeros(len(X[0]))
B_hat = miniBatch(X, Y, 50, muBest, lambdaBest, numberIterations)
Y_hat = sig(X_test, B_hat)
Y_pred = mapProbability(Y_hat)
logLoss = logloss(Y_hat, Y_test)
acc = accuracyValue(Y_pred, Y_test)
print('Accuracy:', acc)
print('Log-likelihood:', logLoss)

Accuracy: 89.13464196848217
Log-likelihood: [2204.04901852]
```

5. Plot Train and Validation Accuracy and Log-likelihood metrics per k – fold iteration.

```
In [177]: X = nds.loc[:, nds.columns != 'y'].to_numpy()
Y = nds['y'].to_numpy()
bias_column = np.ones(shape=(len(X), 1))
X = np.append(bias_column, X, axis=1)
meanError, accLogList = kFoldCrossValidation(X, Y, 5, muBest, lambdaBest, numberIterations)
iterations = list(range(0, 5))

In [178]: # (train accuracy, test accuracy, logLossTrain[0], logLossTest[0])
accLogList

Out [178]: array([[ 85.42991429, 86.61504425, 2264.03756485, 563.37793643],
 [ 85.70638651, 85.0884956, 2262.76155206, 564.65400643],
 [ 85.7932817, 85.1699115, 2260.50101069, 566.91454783],
 [ 85.65100207, 85.77000885, 2260.13845967, 565.27702881],
 [ 87.97345867, 87.27876106, 2212.87838349, 554.53023419]])

In [179]: iterations

Out [179]: [0, 1, 2, 3, 4]
```

```
In [180]: fig = plt.figure(figsize=(6,8))
ax = plt.axes()
ax.plot(iterations, accLogList[:,0], c='r', marker='o', label='Train')
ax.plot(iterations, accLogList[:,1], c='b', marker='o', label='Test')
ax.set_xlabel('iterations')
ax.set_ylabel('Accuracy')
ax.legend()
plt.show()
```



```
In [181]: fig = plt.figure(figsize=(6,8))
ax = plt.axes()
ax.plot(iterations, accLogList[:,2], c='r', marker='o', label='Train')
ax.plot(iterations, accLogList[:,3], c='b', marker='o', label='Test')
ax.set_xlabel('Log-likelihood')
ax.legend()
plt.show()
```


Exercise 3: Implementing Hyperband for Logistic Regression

```
In [ ]: nds_train = nds.loc[0:math.floor(len(nds)*0.7)]
nds_validation = nds.loc[math.floor(len(nds)*0.7)+1:math.floor(len(nds)*0.85)]
nds_test = nds.loc[math.floor(len(nds)*0.85)+1:]
len(nds) == len(nds_train) + len(nds_test)

In [227]: max_iter = 81
eta = 0.1
logeta = lambda x: np.log(x)/np.log(eta)
s_max = int(logeta(max_iter))
B = (s_max**2)*max_iter

for s in reversed(range(s_max+1)):
    n = int(np.ceil(int(B*max_iter/(s**1))*eta**s))
    lambdaList = np.random.uniform(np.exp(-9), np.exp(-4), n)
    batches = np.random.randint(25, 800, n)
    pList = np.random.uniform(10**(s-7), 10**(s-5), n)
    hyperParamConfigs = []
    for i in range (len(batches)):
        hyperParamConfigs.append([pList[i], lambdaList[i], batches[i]])
    pList = s, n, z, len(hyperParamConfigs)

    finalList
```