

EE4400 - Tutorial 1, Question 3

This question is on the Multi-Layer Perceptron (MLP) and using it to do classification. The aim is to find the best number of hidden nodes in the 3 hidden layers, assuming the same number of hidden nodes in each hidden layer. Cross-validation needs to be done on the training set. The MLP classifier with the best network size is then used for testing.

We shall use the Tensorflow Keras package to implement the MLP classifier. Obtain the data set “from sklearn.datasets import load_iris”. Import the necessary packages.

```
In [ ]: 1  ## Load data from scikit
        2  import numpy as np
        3  import pandas as pd
        4  print("pandas version: {}".format(pd.__version__))
        5  import sklearn
        6  from sklearn.datasets import load_iris
        7  from sklearn.model_selection import train_test_split
        8  from sklearn import metrics
        9
       10  from tensorflow import keras
       11  from keras.models import Sequential
       12  from keras.layers import Dense, Activation
       13  from keras.optimizers import SGD
```

pandas version: 1.1.5

(a) Load the data and split the database into two sets: 80% of samples for training, and 20% of samples for testing.

```
In [ ]: 1  ## Load data
        2  iris_dataset = load_iris()
        3  ## split dataset into training and test sets
        4  X_train, X_test, y_train, y_test = train_test_split(iris_dataset['data'],
        5  iris_dataset['target'],
        6  test_size=0.20,
        7  random_state=0)
```

Change y_train and y_test to categorical values (required for classification).

```
In [ ]: 1  y_train = keras.utils.to_categorical(y_train, num_classes = 3)
        2  y_test = keras.utils.to_categorical(y_test, num_classes = 3)
```

(b) Perform a 5-fold Cross-validation using only the training set to determine the best 3-layer MLP classifier with hidden_layer_sizes=(Nhidd,Nhidd,Nhidd) for Nhidd in range(1,11)^ for prediction. In other words, partition the training set into two sets, 4/5 for training and 1/5 for validation; and repeat this process until each of the 1/5 has been validated. ^ The assumption of hidden_layer_sizes=(Nhidd,Nhidd,Nhidd) is to reduce the search space in this exercise. In field applications, the search needs to consider different sizes for each hidden layer.

```
In [ ]: 1  def MLP_model(Nhidd):
        2  model = Sequential()
        3  model.add(Dense(Nhidd, activation='relu', input_shape=(4,)))
        4  model.add(Dense(Nhidd, activation='relu'))
        5  model.add(Dense(Nhidd, activation='relu'))
        6  model.add(Dense(3, activation='softmax'))
        7  model.compile(optimizer='adam',
        8  loss='categorical_crossentropy',
        9  metrics=['accuracy'])
       10  return model
```

```

In [ ]: 1 acc_train_array = []
2 acc_valid_array = []
3 for Nhidd in range(1,11):
4     acc_train_array_fold = []
5     acc_valid_array_fold = []
6     ## Random permutation of data
7     Idx = np.random.RandomState(seed=8).permutation(len(y_train))
8     ## Tuning: perform 5-fold cross-validation on the training set to determine the best network size
9     clf = MLP_model(Nhidd)
10    for k in range(0,5):
11        N = np.around((k+1)*len(y_train)/5)
12        N = N.astype(int)
13        Xvalid = X_train[Idx[N-24:N]] # validation features
14        Yvalid = y_train[Idx[N-24:N]] # validation targets
15        Idxtrn = np.setdiff1d(Idx, Idx[N-24:N])
16        Xtrain = X_train[Idxtrn] # training features in tuning loop
17        Ytrain = y_train[Idxtrn] # training targets in tuning loop
18        ## MLP Classification with same size for each hidden-layer (specified in question)
19        clf.fit(Xtrain, Ytrain, epochs = 100, verbose = 0)
20        ## trained output
21        y_est_p = clf.predict(Xtrain)
22        Ytrain_class = np.argmax(Ytrain, axis=1)
23        y_est_p_class = np.argmax(y_est_p, axis=1)
24        acc_train_array_fold += [metrics.accuracy_score(y_est_p_class, Ytrain_class)]
25        ## validation output
26        yt_est_p = clf.predict(Xvalid)
27        Yvalid_class = np.argmax(Yvalid, axis=1)
28        yt_est_p_class = np.argmax(yt_est_p, axis=1)
29        acc_valid_array_fold += [metrics.accuracy_score(yt_est_p_class, Yvalid_class)]
30    acc_train_array += [np.mean(acc_train_array_fold)]
31    acc_valid_array += [np.mean(acc_valid_array_fold)]
32    clf.summary()

```

Total params: 255
Trainable params: 255
Non-trainable params: 0

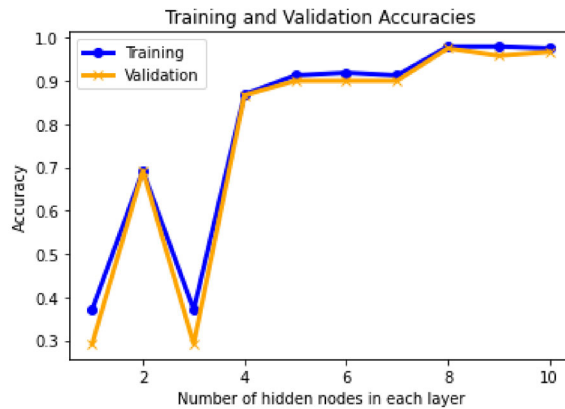
Model: "sequential_9"

Layer (type)	Output Shape	Param #
dense_36 (Dense)	(None, 10)	50
dense_37 (Dense)	(None, 10)	110
dense_38 (Dense)	(None, 10)	110
dense_39 (Dense)	(None, 3)	33

Total params: 303
Trainable params: 303
Non-trainable params: 0

(c) Provide a plot of the average 5-fold training and validation accuracies over the different network sizes, i.e. different number of nodes in the hidden layer. Determine the hidden layer size Nhidd that gives the best validation accuracy for the training set.

```
In [ ]: 1 ## plotting
2 import matplotlib.pyplot as plt
3 hiddensize = [x for x in range(1,11)]
4 plt.plot(hiddensize, acc_train_array, color='blue', marker='o', linewidth=3, label='Training')
5 plt.plot(hiddensize, acc_valid_array, color='orange', marker='x', linewidth=3, label='Validation')
6 plt.xlabel('Number of hidden nodes in each layer')
7 plt.ylabel('Accuracy')
8 plt.title('Training and Validation Accuracies')
9 plt.legend()
10 plt.show()
11 ## find the best hidden layer size that gives the best validation accuracy using only the training set
12 Nhidden = np.argmax(acc_valid_array,axis=0)+1
13 print('best hidden layer size =', Nhidden, 'based on 5-fold cross-validation on training set')
```



best hidden layer size = 8 based on 5-fold cross-validation on training set

```
In [ ]: 1 print(acc_train_array)
2 print(acc_valid_array)
```

[0.3708333333333333, 0.6916666666666667, 0.3708333333333333, 0.86875, 0.9125, 0.91875, 0.9125, 0.9791666666666667, 0.9791666666666666, 0.975]
[0.2916666666666663, 0.6916666666666667, 0.2916666666666663, 0.8666666666666668, 0.9, 0.9, 0.9, 0.975, 0.9583333333333333, 0.9666666666666668]

(d) Using the best hidden layer size Nhidd in the MLP classifier with hidden_layer_sizes=(Nhidd,Nhidd,Nhidd), evaluate the performance of the MLP by computing the prediction accuracy based on the 20% of samples for testing in part (a).

```
In [ ]: 1 ## perform evaluation
2 clf = MLP_model(Nhidd)
3 history=clf.fit(X_train, y_train, epochs = 100, batch_size = 32, verbose = 0)
4 ## trained output
5 y_test_predict = clf.predict(X_test)
6 y_test_class = np.argmax(y_test, axis=1)
7 y_test_predict_class = np.argmax(y_test_predict, axis=1)
8 test_accuracy = metrics.accuracy_score(y_test_predict_class,y_test_class)
9 print('test accuracy =', test_accuracy)
```

test accuracy = 0.8666666666666667

```
In [ ]: 1 plt.plot(history.history["loss"])
```

Out[9]: [matplotlib.lines.Line2D at 0x7f59fbc89908]

