

Name: - Prateek Mahajan

EE 511: Homework 4

Problem 1:-

I have read and understood the general instructions at the top of HW4 and I formally declare that all work I turn in for everything in this course will not contain or involve any cheating at all.

Problem 2:-

- a) Please refer to the end of this pdf for the RF, ridge & least square regressions.

While RF should typically perform better than ridge & least square regressions as the presence of various trees reduces both its bias & variance due to a lack of parameter tuning, I suspect my ~~new~~ RMSE results are higher than that of ridge & least squares regressions.

- b) Please refer

Since the provided dataset is high dimensional, has outliers & appears to be non-linear in nature, I would expect that the random forest model would perform better, as its various trees & structure make it more robust to outliers & overfitting, while also allowing it to represent non-linear relation.

- fits better.

However, due to a lack of parameter-tuning that is not the case here & so, the RMSE of RF is higher than ridge & least squares.

- b) Please refer to the end of this pdf for this question

Problem 3:-

- a) The computational complexity of optimizing a tree of depth d for dataset of size $n \times m$ with m features is difficult to generalise & depends on the method of optimisation. Therefore, assuming a binary tree, find some cases below.
- Worst Case Complexity : $O(n \log n \cdot mn \cdot 2^d)$
w/o optimisations
 $\xrightarrow{\text{sorting complexity}}$ $\xrightarrow{\text{traversing through split points}}$ $\xrightarrow{\text{checking all case w/o pruning}}$
 - Random Forest : $O(T \cdot n \cdot K \cdot \log n)$, where $T = \text{no. of trees}$
 $K = \sqrt{m}$ for classification & $K = m/3$ for regression
 - GBDT: $O(T m n d)$, where T, m, n & d have meanings as mentioned above
- b) The most expensive part of the GBDT computation is finding the split point at a given node. At this step, we typically need to parse the split nodes of all features to find the ideal one.

possible optimisation in large datasets is using sparse matrices (particularly for binary classification problems). We could also do pruning of some nodes/trees that contribute little to the entire model, thereby reducing computation. Another possibility is to use histograms-based methods to find the split points.

c) The "split-point" finding part of the computation for m features of a given node could be parallelised as there are no dependencies between the calculations b/w various features.

→ The simplifications of the left & right child of a node could be run on separate threads, as again these would be independent of one another / would have no dependencies on each other.

d) There are 2 major differences in how GBDT updates its model parameters vs. GD, SGD & CD:-

→ Each new tree/model in GBDT is given a negative impact of a loss function with respect to the predictions of previous trees in the GBDT. This is unlike the other model, which don't use info from an "ensemble".

→ SGD, GD & CD retain the same number of parameters throughout training. Depending on the data, GBDT may add more parameters during training as its leaf nodes increase.

- c) Once again GBDT should perform better than least squares regression & ridge regression. However, in this case, the RMSE I got was higher ('I'm suspecting due to a lack of hyperparameter tuning). Please refer to the end of the pdf for the data.
- f) Please refer to the end of the pdf for this qn.
- g) GBDT focuses on reducing ~~the~~ bias by focusing on correcting the output of the previous trees. Random forests, on the other hand, focuses on ~~on~~ reducing variance by using different subsets of data & features.

In the 3 experiments, GBDT & RF gave similar results for both classification problems, but RF outperformed GBDT in the regression problem. I would attribute this to the datasets. Since the boston dataset is noisy, reducing variance is more important & so, RF performs better. In the credit-g & breast cancer datasets GBDT & RF had similar performance, but GBDT edged out RF as the datasets were complex & reducing bias to prevent overfitting was more important than variance.

Problem: —

Please refer to the end of this pdf

gbdt

March 6, 2024

```
[1]: from multiprocessing import Pool
from functools import partial
import numpy as np
from numba import jit
```

```
[2]: #TODO: loss of least square regression and binary logistic regression
"""
pred() takes GBDT/RF outputs, i.e., the "score", as its inputs, and returns
predictions.

g() is the gradient/1st order derivative, which takes true values "true"
and scores as input, and returns gradient.

h() is the heassian/2nd order derivative, which takes true values "true"
and scores as input, and returns hessian.

"""

class loss(object):
    '''Loss class for mse. As for mse, pred function is pred(score).'''
    def pred(self,X_test,trees,log_threshold):
        return 0

    def g(self, y_true, pred):
        return 0

    def h(self, pred):
        return 0

class leastsquare(loss):
    '''Loss class for mse. As for mse, pred function is pred(score).'''
    def pred(self,X_test,trees,log_threshold):
        predictions = []
        #print("Tree Length : " + str(len(trees)))
        for test in X_test:
            #count = 0
            score = 0.0
            for tree in trees:
                curr_node = tree.root_node
                while curr_node.is_leaf is not True :
                    curr_node = curr_node.forward(test)
```

```

        #print(str(count) + " " + str(curr_node.weight))
        #count = count + 1
        score = score + curr_node.weight
    score = score/len(trees)
    predictions.append(score)
return predictions

def g(self, y_true, pred):
    return (-2*(y_true - pred))

def h(self, y_true, pred):
    return 2*np.ones(y_true.shape)

class logistic(loss):
    '''Loss class for log loss. As for log loss, pred function is logistic transformation.'''
    def pred(self,X_test,trees,log_threshold):
        predictions = []
        #print(len(trees))
        for test in X_test:
            temp = np.zeros(len(trees))
            for i in range(len(trees)):
                tree = trees[i]
                curr_node = tree.root_node
                while(curr_node.is_leaf == False):
                    curr_node = curr_node.forward(test)
                    temp[i] = curr_node.weight
            min = np.min(temp)
            max = np.max(temp)
            #print(min)
            #print(max)
            if max != min :
                temp = (temp - min)/(max - min)
            temp = np.rint(temp)
        unique, counts = np.unique(temp, return_counts=True)
        #print(unique)
        if(unique.size > 1):
            if counts[1] > counts[0] :
                predictions.append(1.0)
            else :
                predictions.append(0.0)
        else :
            predictions.append(unique[0])
    return predictions

def g(self, y_true, pred):
    return (((1-y_true)*np.exp(pred)) - y_true)/(np.exp(pred) + 1))

```

```

def h(self, y_true, pred):
    return (np.exp(pred)/((np.exp(pred) + 1)*(np.exp(pred) + 1)))*np.
    ↪ones(y_true.shape)

```

[3]: # TODO: class of a node on a tree

```
class TreeNode(object):
```

```
    '''
```

Data structure that are used for storing a node on a tree.

*A tree is presented by a set of nested TreeNodes,
with one TreeNode pointing two child TreeNodes,
until a tree leaf is reached.*

A node on a tree can be either a leaf node or a non-leaf node.

```
    '''
```

#TODO

```

def __init__(self, X1, X2, index_y, threshold, value_y, is_leaf):
    self.is_leaf = is_leaf
    self.left_child = X1
    self.right_child = X2
    self.feature_index = index_y
    self.threshold = threshold
    self.weight = value_y

def forward(self, x):
    #print(self.weight)
    #print(self.feature_index)
    #print(self.threshold)
    #print(self.is_leaf)
    #print(self.left_child.is_leaf)
    #print(self.right_child.is_leaf)
    if self.is_leaf == True :
        raise RuntimeError('This is a leaf node and you cannot forward from
    ↪it')
    elif x[self.feature_index] < self.threshold:
        #print(self.left_child.is_leaf)
        return self.left_child
    else :
        #print(self.right_child.is_leaf)
        return self.right_child

```

[4]: # TODO: class of single tree

```
class Tree(object):
```

```
    '''
```

Class of a single decision tree in GBDT

```

Parameters:
    n_threads: The number of threads used for fitting and predicting.
    max_depth: The maximum depth of the tree.
    min_sample_split: The minimum number of samples required to further_
        ↪split a node.
    lamda: The regularization coefficient for leaf prediction, also known_
        ↪as lambda.
    gamma: The regularization coefficient for number of TreeNode, also known_
        ↪as gamma.
    rf: rf*m is the size of random subset of features, from which we select_
        ↪the best decision rule,
    rf = 0 means we are training a GBDT.
    ...
    ...

def __init__(self, root_node = None, n_threads = None,
            max_depth = 3, min_sample_split = 10,
            lamda = 1, gamma = 0, pred = 0):
    self.n_threads = n_threads
    self.max_depth = max_depth
    self.min_sample_split = min_sample_split
    self.root_node = root_node
    self.lamda = lamda
    self.gamma = gamma
    self.pred = pred
    self.treenodes = []

def fit(self, X, y, loss):
    """
    train is the training data matrix, and must be numpy array (an n_train_
        ↪x m matrix).
    g and h are gradient and hessian respectively.
    """
    g = loss.g(y, self.pred)
    h = loss.h(y, self.pred)
    #print(g)
    #print(h)
    self.root_node = self.construct_tree(X, y, g, h, 0)
    #print(self.root_node.weight)
    return self

def split_dataset(self, X, y, feature_index, threshold):
    left_mask = X[:, feature_index] < threshold
    right_mask = X[:, feature_index] >= threshold
    return X[left_mask], X[right_mask], y[left_mask], y[right_mask]

def calculate_gain(self, g_l, h_l, g_r, h_r):

```

```

        G_l_tot = g_l.sum()
        G_r_tot = g_r.sum()
        H_l_tot = h_l.sum()
        H_r_tot = h_r.sum()
        gain = (0.5)*(((G_l_tot)*(G_l_tot)/(H_l_tot + self.lamda)) +_
        ((G_r_tot)*(G_r_tot)/(H_r_tot + self.lamda)) -_
        (((G_l_tot+G_r_tot)*(G_l_tot+G_r_tot))/(H_l_tot + H_r_tot+ self.lamda))) -_
        self.gamma
    return gain

def find_threshold(self, g, h, X, y, feature_index):
    """
    Given a particular feature  $p_j$ ,
    return the best split threshold  $\tau_j$  together with the gain that is
    achieved.
    """
    best_gain = 0
    best_threshold = 0
    thresholds = np.unique(X[:, feature_index])
    for threshold in thresholds:
        X_left, X_right, y_left, y_right = self.split_dataset(X, y, feature_index, threshold)
        if len(y_left) == 0 or len(y_right) == 0:
            continue
        left_mask = X[:, feature_index] < threshold
        right_mask = X[:, feature_index] >= threshold
        gain = self.calculate_gain(g[left_mask], h[left_mask],_
        g[right_mask], h[right_mask])
        #print("Gain : " +str(gain))
        if gain > best_gain:
            best_gain = gain
            best_threshold = threshold
    return [best_threshold, best_gain]

def construct_tree(self, X, y, g, h, current_depth):
    """
    Node Addition, which is recursively used to grow a tree.
    First we should check if we should stop further splitting.

    The stopping conditions include:
    1. tree reaches max_depth  $d_{max}$ 
    2. The number of sample points at current node is less than
    min_sample_split, i.e.,  $n_{min}$ 
    3. gain <= 0
    """
    if current_depth > self.max_depth or y.size < self.min_sample_split :
        return None

```

```

    best_feature, threshold, best_gain = self.find_best_decision_rule(X, y, u
↳g, h)
    #print("best_feature : " + str(best_feature))
    #print("threshold : " + str(threshold))
    #print("best_gain : " + str(best_gain))
    if best_gain == 0 :
        return None
    #print(X[:, best_feature])
    left_mask = X[:, best_feature] < threshold
    right_mask = X[:, best_feature] >= threshold
    node_left = self.construct_tree(X[left_mask], y[left_mask], u
↳g[left_mask], h[left_mask], current_depth + 1)
    node_right = self.construct_tree(X[right_mask], y[right_mask], u
↳g[right_mask], h[right_mask], current_depth + 1)
    h_val = h.sum()
    weight = -(g.sum()/(h_val + self.lamda))
    is_leaf = False
    if node_left is None and node_right is not None :
        node_left = TreeNode(None, None, 0, 0, weight, True)
    elif node_right is None and node_left is not None :
        node_right = TreeNode(None, None, 0, 0, weight, True)
    elif node_right is None and node_left is None :
        is_leaf = True
    node = TreeNode(node_left, node_right, best_feature, threshold, weight, u
↳is_leaf)
    self.treenodes.append(node)
    #print(node)
    return node

def find_best_decision_rule(self, X, y, g, h):
    """
    Return the best decision rule [feature, threshold], i.e., $(p_j, u
↳\tau_j)$ on a node j,
    train is the training data assigned to node j
    g and h are the corresponding 1st and 2nd derivatives for each data u
↳point in train
    g and h should be vectors of the same length as the number of data u
↳points in train

    for each feature, we find the best threshold by find_threshold(),
    a [threshold, best_gain] list is returned for each feature.
    Then we select the feature with the largest best_gain,
    and return the best decision rule [feature, threshold] together with its u
↳gain.
    """
    best_gain = 0

```

```

best_threshold = 0
best_feature = 0
for feature_index in range(X.shape[1]):
    #print("Feature " + str(feature_index))
    threshold, gain = self.find_threshold(g, h, X, y, feature_index)
    if gain > best_gain:
        best_gain = gain
        best_threshold = threshold
        best_feature = feature_index
return best_feature, threshold, best_gain

```

[5]: # TODO: class of Random Forest

```

class RF(object):
    """
    Class of Random Forest

    Parameters:
        n_threads: The number of threads used for fitting and predicting.
        loss: Loss function for gradient boosting.
            'mse' for regression task and 'log' for classification task.
            A child class of the loss class could be passed to implement
            ↵customized loss.
        max_depth: The maximum depth d_max of a tree.
        min_sample_split: The minimum number of samples required to further
            ↵split a node.
        lamda: The regularization coefficient for leaf score, also known as
            ↵lambda.
        gamma: The regularization coefficient for number of tree nodes, also
            ↵know as gamma.
        rf: rf*m is the size of random subset of features, from which we select
            ↵the best decision rule.
        num_trees: Number of trees.
    """

    def __init__(self, log_threshold = 0.5,
                 n_threads = 0, loss = 'mse',
                 max_depth = 3, min_sample_split = 10,
                 lamda = 1, gamma = 0,
                 rf = 0.99, num_trees = 100):
        self.log_threshold = log_threshold
        self.n_threads = n_threads
        if loss == 'mse':
            self.loss = leastsquare()
        else :
            self.loss = logistic()
        self.max_depth = max_depth
        self.min_sample_split = min_sample_split
        self.lamda = lamda

```

```

        self.gamma = gamma
        self.rf = rf
        self.num_trees = num_trees
        self.trees = []

    def fit(self, X, y):
        # X is n x m 2d numpy array
        # y is n-dim 1d array
        num_samples, num_features = X.shape
        for i in range(self.num_trees):
            #print ("Tree : " + str(i))
            m_rf = (int)(np.ceil(self.rf*num_features))
            m_rf_indices = np.random.choice(num_features, m_rf, replace=False)
            n_rf_indices = np.random.choice(num_samples, num_samples, □
            ↪replace=True)
            XTree = X[n_rf_indices] [:,m_rf_indices]
            #print(XTree.shape)
            yTree = y[n_rf_indices]
            prediction = 0
            tree = Tree(None, self.n_threads, self.max_depth, self.
            ↪min_sample_split, self.lamda, self.gamma, prediction)
            tree = tree.fit(XTree, yTree, self.loss)
            self.trees.append(tree)
        return self

    # test is the input to the model
    def predict(self, X_test):
        return self.loss.pred(X_test, self.trees, self.log_threshold)

```

[6]: # TODO: class of GBDT

```

class GBDT(object):
    ...
    Class of gradient boosting decision tree (GBDT)

    Parameters:
        n_threads: The number of threads used for fitting and predicting.
        loss: Loss function for gradient boosting.
            'mse' for regression task and 'log' for classification task.
            A child class of the loss class could be passed to implement
            ↪customized loss.
        max_depth: The maximum depth D_max of a tree.
        min_sample_split: The minimum number of samples required to further
            ↪split a node.
        lamda: The regularization coefficient for leaf score, also known as
            ↪lambda.
        gamma: The regularization coefficient for number of tree nodes, also
            ↪know as gamma.

```

```

learning_rate: The learning rate eta of GBDT.
num_trees: Number of trees.

def __init__(self, log_threshold = 0.5,
            learning_rate=0.5, n_threads = 0, loss = 'mse',
            max_depth = 3, min_sample_split = 10,
            lamda = 1, gamma = 0,
            rf = 0.99, num_trees = 100):
    self.loss_type = loss
    self.n_threads = n_threads
    if loss == 'mse':
        self.loss = leastsquare()
    else :
        self.loss = logistic()
    self.max_depth = max_depth
    self.log_threshold = log_threshold
    self.min_sample_split = min_sample_split
    self.lamda = lamda
    self.gamma = gamma
    self.rf = rf
    self.num_trees = num_trees
    self.learning_rate = learning_rate
    self.trees = []

def fit(self, X, y):
    # X is n x m 2d numpy array
    # y is n-dim 1d array
    num_samples, num_features = X.shape
    prediction = np.zeros(y.shape)
    for i in range(self.num_trees):
        #print ("Tree : " + str(i))
        tree = Tree(None, self.n_threads, self.max_depth, self.
        ↵min_sample_split, self.lamda, self.gamma, prediction)
        tree = tree.fit(X, y, self.loss)
        self.trees.append(tree)

        if self.loss_type == 'mse':
            new_pred = self.loss.pred(X, [tree], self.log_threshold)
        else :
            count0 = 0
            count1 = 1
            curr_node = tree.root_node
            new_pred = []
            for test in X:
                while(curr_node.is_leaf == False):
                    curr_node = curr_node.forward(test)
                if curr_node.weight > self.log_threshold:

```

```

        count1 = count1 + 1
    else :
        count0 = count0 + 1
    if count1 > count0 :
        new_pred.append(1.0)
    else :
        new_pred.append(0.0)
    prediction = self.learning_rate*prediction + np.array(new_pred)
return self

# test is the input to the model
def predict(self, X_test):
    return self.loss.pred(X_test, self.trees, self.log_threshold)

```

[7]: # TODO: Evaluation functions (you can use code from previous homeworks)

```

# RMSE
def root_mean_square_error(pred, y):
    diff_matrix = y - pred
    #print(y)
    #print(diff_matrix)
    rmse = diff_matrix**2
    rmse = rmse.sum()
    rmse = rmse/np.size(y)
    rmse = np.sqrt(rmse)
    return rmse

# precision
def accuracy(pred, y):
    dataset_size = y.size
    correct_predictions = 0
    #print(pred)
    #print(y)
    for i in range(dataset_size):
        if y[i] == (pred[i]):
            correct_predictions += 1
    precision = (correct_predictions/dataset_size)*100
    return precision

```

[8]: # Run Random Forest Code

```

from sklearn import datasets
boston = datasets.load_boston()
X = boston.data
y = boston.target

# train-test split
from sklearn.model_selection import train_test_split

```

```

X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.3,
    ↪random_state=8)
RF_regression = RF(0.5, n_threads = 0, loss = 'mse', max_depth = 8,
    ↪min_sample_split = 10, lamda = 1, gamma = 1, rf = 0.8, num_trees = 99)
RF_regression = RF_regression.fit(np.array(X_train), np.array(y_train))
pred_train = RF_regression.predict(X_train)
pred_test = RF_regression.predict(X_test)
rmse_test = root_mean_square_error(pred_test, y_test)
rmse_train = root_mean_square_error(pred_train, y_train)
print("Training RMSE of the Random Forest Method for the Boston Dataset is : ", ↪
    ↪rmse_train)
print("Testing RMSE of the Random Forest Method for the Boston Dataset is : ", ↪
    ↪rmse_test)
print("Training Set RMSE of Linear Regression for Boston Dataset (from HW2): 4. ↪
    ↪820626531838223")
print("Testing Set RMSE of Linear Regression for Boston Dataset (from HW2): 5. ↪
    ↪209217510530916")
print("Training Set RMSE of Ridge Regression for Boston Dataset (from HW2): 4. ↪
    ↪829777333975097")
print("Testing Set RMSE of Ridge Regression for Boston Dataset (from HW2): 5. ↪
    ↪189347305423606")

```

/Users/prateek/anaconda3/envs/AnacondaTest/lib/python3.11/site-packages/sklearn/utils/deprecation.py:87: FutureWarning: Function load_boston is deprecated; `load_boston` is deprecated in 1.0 and will be removed in 1.2.

The Boston housing prices dataset has an ethical problem. You can refer to the documentation of this function for further details.

The scikit-learn maintainers therefore strongly discourage the use of this dataset unless the purpose of the code is to study and educate about ethical issues in data science and machine learning.

In this special case, you can fetch the dataset from the original source::

```

import pandas as pd
import numpy as np

data_url = "http://lib.stat.cmu.edu/datasets/boston"
raw_df = pd.read_csv(data_url, sep="\s+", skiprows=22, header=None)
data = np.hstack([raw_df.values[:, :-1], raw_df.values[:, -1]])
target = raw_df.values[:, -1]

```

Alternative datasets include the California housing dataset (i.e. :func:`~sklearn.datasets.fetch_california_housing`) and the Ames housing dataset. You can load the datasets as follows::

```

from sklearn.datasets import fetch_california_housing
housing = fetch_california_housing()

for the California housing dataset and::

    from sklearn.datasets import fetch_openml
    housing = fetch_openml(name="house_prices", as_frame=True)

    for the Ames housing dataset.
    warnings.warn(msg, category=FutureWarning)

Training RMSE of the Random Forest Method for the Boston Dataset is :
9.229962989283528
Testing RMSE of the Random Forest Method for the Boston Dataset is :
8.975222658557659
Training Set RMSE of Linear Regression for Boston Dataset (from HW2):
4.820626531838223
Testing Set RMSE of Linear Regression for Boston Dataset (from HW2):
5.209217510530916
Training Set RMSE of Ridge Regression for Boston Dataset (from HW2):
4.829777333975097
Testing Set RMSE of Ridge Regression for Boston Dataset (from HW2):
5.189347305423606

```

```
[20]: # load data
from sklearn.datasets import fetch_openml
X, y = fetch_openml('credit-g', version=1, return_X_y=True, data_home='credit/'  

    ↪, as_frame=False)
y = np.array(list(map(lambda x: 1 if x == 'good' else 0, y)))

# train-test split
from sklearn.model_selection import train_test_split
X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.3,  

    ↪random_state=8)
#print(X_train)
RF_regression = RF(log_threshold = 0.5, n_threads = 0, loss = 'class', max_depth=  

    ↪= 8, min_sample_split = 10, lamda = 1, gamma = 1, rf = 0.5, num_trees = 99)
RF_regression = RF_regression.fit(np.array(X_train), np.array(y_train))
pred_train = RF_regression.predict(X_train)
pred_test = RF_regression.predict(X_test)
acc_test = accuracy(pred_test, y_test)
acc_train = accuracy(pred_train, y_train)
print("Training Accuracy of the Random Forest Method for the Credit-g Dataset  

    ↪is : ", acc_train)
print("Testing Accuracy of the Random Forest Method for the Credit-g Dataset is  

    ↪: ", acc_test)
```

```
Training Accuracy of the Random Forest Method for the Credit-g Dataset is :  
69.71428571428572  
Testing Accuracy of the Random Forest Method for the Credit-g Dataset is : 70.0
```

```
[10]: # load data  
from sklearn import datasets  
breast_cancer = datasets.load_breast_cancer()  
X = breast_cancer.data  
y = breast_cancer.target  
  
# train-test split  
from sklearn.model_selection import train_test_split  
X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.3, u  
↳random_state=8)  
RF_regression = RF(log_threshold = 0.5, n_threads = 0, loss = 'class', u  
↳max_depth = 10, min_sample_split = 10, lamda = 1, gamma = 1, rf = 0.5, u  
↳num_trees = 99)  
RF_regression = RF_regression.fit(np.array(X_train), np.array(y_train))  
pred_train = RF_regression.predict(X_train)  
pred_test = RF_regression.predict(X_test)  
acc_test = accuracy(pred_test, y_test)  
acc_train = accuracy(pred_train, y_train)  
print("Training Accuracy of the Random Forest Method for the Breast Cancer u  
↳Dataset is : ", acc_train)  
print("Testing Accuracy of the Random Forest Method for the Breast Cancer u  
↳Dataset is : ", acc_test)
```

```
Training Accuracy of the Random Forest Method for the Breast Cancer Dataset is :  
63.31658291457286  
Testing Accuracy of the Random Forest Method for the Breast Cancer Dataset is :  
61.40350877192983
```

```
[11]: # TODO: GBDT regression on boston house price dataset  
  
# load data  
from sklearn import datasets  
boston = datasets.load_boston()  
X = boston.data  
y = boston.target  
  
# train-test split  
from sklearn.model_selection import train_test_split  
X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.3, u  
↳random_state=8)  
learning_rate = 0.2
```

```

GBDT_regression = GBDT(0.5, learning_rate, n_threads = 0, loss = 'mse',  

    ↪max_depth = 8, min_sample_split = 10, lamda = 1, gamma = 1, rf = 0.5,  

    ↪num_trees = 99)  

GBDT_regression = GBDT_regression.fit(np.array(X_train), np.array(y_train))  

pred_train = GBDT_regression.predict(X_train)  

pred_test = GBDT_regression.predict(X_test)  

rmse_test = root_mean_square_error(pred_test, y_test)  

rmse_train = root_mean_square_error(pred_train, y_train)  

print("Training RMSE of the GBDT Method for the Boston Dataset is : ",  

    ↪rmse_train)  

print("Testing RMSE of the GBDT Method for the Boston Dataset is : ", rmse_test)

```

Training RMSE of the GBDT Method for the Boston Dataset is : 15.3444415615304
Testing RMSE of the GBDT Method for the Boston Dataset is : 15.524724610785293

[12]: # TODO: GBDT classification on credit-g dataset

```

# load data
from sklearn.datasets import fetch_openml
X, y = fetch_openml('credit-g', version=1, return_X_y=True, data_home='credit/  

    ↪', as_frame=False)
y = np.array(list(map(lambda x: 1 if x == 'good' else 0, y)))  

# train-test split
from sklearn.model_selection import train_test_split
X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.3,  

    ↪random_state=8)
learning_rate = 0.9
GBDT_regression = GBDT(0.5, learning_rate, n_threads = 0, loss = 'class',  

    ↪max_depth = 8, min_sample_split = 10, lamda = 5, gamma = 0.5, rf = 0.5,  

    ↪num_trees = 100)
GBDT_regression = GBDT_regression.fit(np.array(X_train), np.array(y_train))
pred_train = GBDT_regression.predict(X_train)
pred_test = GBDT_regression.predict(X_test)
acc_test = accuracy(pred_test, y_test)
acc_train = accuracy(pred_train, y_train)
print("Training Accuracy of the GBDT Method for the Credit-g Dataset is : ",  

    ↪acc_train)
print("Testing Accuracy of the GBDT Method for the Credit-g Dataset is : ",  

    ↪acc_test)

```

Training Accuracy of the GBDT Method for the Credit-g Dataset is :

68.85714285714286

Testing Accuracy of the GBDT Method for the Credit-g Dataset is :

71.33333333333334

```
[13]: # TODO: GBDT classification on breast cancer dataset

# load data
from sklearn import datasets
breast_cancer = datasets.load_breast_cancer()
X = breast_cancer.data
y = breast_cancer.target

# train-test split
from sklearn.model_selection import train_test_split
X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.3,
random_state=8)

# train-test split
from sklearn.model_selection import train_test_split
X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.3,
random_state=8)
learning_rate = 0.8
GBDT_regression = GBDT(0.5, learning_rate, n_threads = 0, loss = 'class',
max_depth = 4, min_sample_split = 10, lamda = 6, gamma = 0.8, rf = 0.5,
num_trees = 60)
GBDT_regression = GBDT_regression.fit(np.array(X_train), np.array(y_train))
pred_train = GBDT_regression.predict(X_train)
pred_test = GBDT_regression.predict(X_test)
acc_test = accuracy(pred_test, y_test)
acc_train = accuracy(pred_train, y_train)
print("Training Accuracy of the GBDT Method for the Breast Cancer Dataset is :",
acc_train)
print("Testing Accuracy of the GBDT Method for the Breast Cancer Dataset is :",
acc_test)
```

Training Accuracy of the GBDT Method for the Breast Cancer Dataset is :
63.31658291457286

Testing Accuracy of the GBDT Method for the Breast Cancer Dataset is :
61.40350877192983

[]:

March 9, 2024

1 Fruit Quality Classification Exercise using KNN classifiers

In this notebook, we'll go through and explore a fruit goodness dataset, build a few classifiers using standard pre-coded classifier code using python libraries, and then ultimately code up our own from-scratch KNN classifier that we have coded ourselves and compare how it works. Using our own code, we'll look at the k value to see how it effects the quality of the classifier on both the training and test sets. To address this problem, you need to write the code wherever it says **FILLIN** or ****FILLIN**** below.

The code you produce *must* be your own code, that you wrote yourself without help from anything other than your own brain.

1.0.1 Analysis of the histograms

From the above histograms, and from the perspective of designing our KNN classifier, the data is centered around zero for all features and has similar standard deviation in all of them. Therefore, this histogram tells us the data is normalised. In kNN, the data being normalised makes calculation easier, which is why this is important. Further, the histogram could help us in identifying a maximum value of “ k ” (number of nearest neighbours) to be used in our simulations. There is no point in setting a “ k ” value greater than the maximum number (or even second/third highest number depending on the data) of samples in a bin of any feature in the dataset. Hence, in this case, these histograms could help us limit k down to a value less than 1500

1.0.2 Analysis

From the above set of statistics, we see that all features are in the range (-8,+8), but seem to be centred around 0. Further, the features are concentrated near the mean (which is approximately zero), and so, the data appears normalised. Overall, does this dataset looks ok (i.e., free from errors or bugs) at least from what we can tell so far? Why or why not? Thus far, the dataset looks okay. The count of all features is the same, and the mean, standard deviation, min, max and 25/50/75% points seem reasonable. The data also seems normalised which is very useful in simplifying kNN calculations!

1.0.3 Analysis of balance

Regarding balance, we see that this data is fairly well balanced on a high level. The number of good and bad samples seems similar.

The above shows that looking at ‘Ripeness’ vs ‘Acidity’ shows that the data is difficult to separate linearly. That being said, there is a tiny amount of scope of a linear decision boundary, but it

would not be very clean. Lets try a few others.

The three above pairwise plots all show that, at least for the set of 2D axes we selected, regarding linear seperability, the data is not cleanly linearly seperable for any of the 2D axes combinations shown above.

1.0.4 Many scatter plots

In the above we plotted only 3 scatter plots, but we can plot all $m \choose 2$ of them at the same time as well as the feature histograms over the labels in one fell swoop. When you do not have too many features, this is a really useful plot to do to understand your data better, albeit realize that we're still only looking at all $m \choose 2$ scatter plots, there still could be higher order relationships amongst size > 2 groups of features that these plots will not reveal.

So the above plots tell three important but distinct properties about the data, namely that:

- (1) There are no redundant features in the data, as the correlation between features is low in general
- (2) Features such as size and weight, or sweetness and ripeness would intuitively (as per the physical world), be highly correlated. The data, however, shows otherwise
- (3) Data appears to be normalised, as it is centred around zero and does not seem very spread out.

1.0.5 Preparing data for prediction

What we next do is adjust the data to be ready to build a classifier. We'll first build several pre-built classifiers using sklearn, and then write code ourselves from scratch for a KNN classifier.

```
[26]: # We start with a random forest classifier with 100 trees and a random state of ↵42.
forest_model_1 = RandomForestClassifier(n_estimators=100, random_state=42)
# We call the `acc_score` function with the `forest_model_1` model as an ↵argument.
acc_score(forest_model_1)
```

ACC score: 0.90125

```
[27]: # We then try a random forest classifier with 1000 trees and a random state of ↵42.
# Does this boost the accuracy relative to the random forest classifier? Yes, ↵it does
forest_model_2 = RandomForestClassifier(n_estimators=1000, random_state=42)
# We call the `acc_score` function with the `forest_model_2` model as an ↵argument.
acc_score(forest_model_2)
```

ACC score: 0.905

```
[28]: import xgboost as xgb
# Next, we try an XGBoost classifier with 100 trees and a random state of 42.
# Does this celebrated model family do better than the random forest? No, it ↴
    ↴ gives the same accuracy as the random forest model
xgb_model_1 = xgb.XGBClassifier(objective='binary:logistic', n_estimators=100, ↴
    ↴ seed=42)
# We call the `acc_score` function with the `xgb_model_1` model as an argument.
acc_score(xgb_model_1)
```

ACC score: 0.90125

```
[29]: # Next we check if the XGBoost classifier with 100 trees is statistically ↴
    ↴ significantly better than the random forest classifier with 1000 trees.
# We can do this by using a t-test to compare the two models.
from scipy.stats import ttest_rel
y_preds_xgb = xgb_model_1.predict(X_test)
y_preds_forest = forest_model_2.predict(X_test)
ttest_rel(y_preds_xgb, y_preds_forest)
# The output is statistic=-1.480693188626695, pvalue=0.13908255678067352, ↴
    ↴ df=799. So there is no evidence that either classifier is statistically ↴
    ↴ significantly better than the other
# Thus, do we reject or do we fail to reject the null hypothesis that the two ↴
    ↴ models are statistically significantly different? Clearly, the pvalue is ↴
    ↴ greater than 0.05, so there is no evidence to reject the null hypothesis. ↴
    ↴ Hence, we fail to reject the null hypothesis
```

[29]: TtestResult(statistic=-1.480693188626695, pvalue=0.13908255678067352, df=799)

```
[30]: # Next, we try an XGBoost classifier with 1000 trees and a random state of 42. ↴
    ↴ Will this boost the accuracy more than the above classifiers?
xgb_model_2 = xgb.XGBClassifier(objective='binary:logistic', n_estimators=1000, ↴
    ↴ seed=42)
# We call the `acc_score` function with the `xgb_model_2` model as an argument.
acc_score(xgb_model_2)
```

ACC score: 0.8975

```
[31]: # Last but not least, we try a decision tree classifier with a random state of ↴
    ↴ 42.
decision_model = DecisionTreeClassifier(random_state=42)
# We call the `acc_score` function with the `decision_model` model as an ↴
    ↴ argument.
acc_score(decision_model)
```

ACC score: 0.80625

1.0.6 Conclusion on pre-built classifiers.

We should see some variability of the performance of each of the classifiers.

You are to record the accuracies of the sklearn models in the dataset and compare them with your own KNN classifier below.

Looking at the above results, we conclude that: - Random forest classifier seems to improve in performance as we increase the number of trees, while the same is not the case for the GBDT - Standard decision trees show worse performance than random forest and GBDT

1.0.7 On to building our own KNN classifier.

Having established the above baselines to compare against (which is *always* a very good idea to do), we next build our own KNN decision tree code and see if we can match or beat the above. We will use the variable `kvalue` to be the number of nearest neighbors that our KNN classifier uses.

```
[32]: from collections import Counter
from tqdm.autonotebook import tqdm

# This function will calculate the Euclidean distance between two rows of data.
def euclidean_distance(row1, row2):
    distance = np.linalg.norm(row1 - row2)
    return distance
```

```
[33]: # KNN Algorithm. This function takes in the training data, the test data and ↴the value of k.
# It runs the KNN algorithm using the kvalue `k` and returns the predicted ↴values for the given test data
# as a numpy array.
```

```
# Here you are to provide your own code for a KNN classifier. You may if you ↴like use the `euclidean_distance` function you provided above to help you.
def knn(X_train, y_train, X_test, kvalue):
    y_pred = []
    for x in X_test :
        distances = [euclidean_distance(x, x_train) for x_train in X_train]
        k_indices = np.argsort(distances)[:kvalue]
        k_nearest_labels = [y_train[i] for i in k_indices]
        most_common = np.argmax(np.bincount(k_nearest_labels))
        y_pred.append(most_common)
    return np.array(y_pred)
```

```
[34]: # This function computes the accuracy of a model by comparing the predicted ↴values with the actual values, both of which
# are given in the arguments. The function returns the accuracy as a float.
def accuracy(y_true, y_pred):
    accuracy = np.sum(y_true == y_pred) / len(y_true)
```

```

    return accuracy

[35]: # This sets up testing for different values of kvalue
k_values = range(1, 200, 1)

[36]: # We next ensure that the data is in NumPy array format since the code above
      ↴uses numpy matrices and arrays, not pandas DataFrames.
X_train_np = X_train.to_numpy()
X_test_np = X_test.to_numpy()
y_train_np = y_train.to_numpy().ravel()    # ravel in case y_train is a column
      ↴vector
y_test_np = y_test.to_numpy().ravel()      # ravel in case y_test is a column
      ↴vector

[37]: # Reset the accuracies arrays to empty lists
train_accuracies = []
test_accuracies = []
# We loop over the different values of k and calculate the accuracy for each
      ↴value of k.
for kvalue in tqdm(k_values, position=0, desc='k_values Progress'):
    # Predict labels for train set and test set
    y_pred_train = knn(X_train_np, y_train_np, X_train_np, kvalue)
    y_pred_test = knn(X_train_np, y_train_np, X_test_np, kvalue)

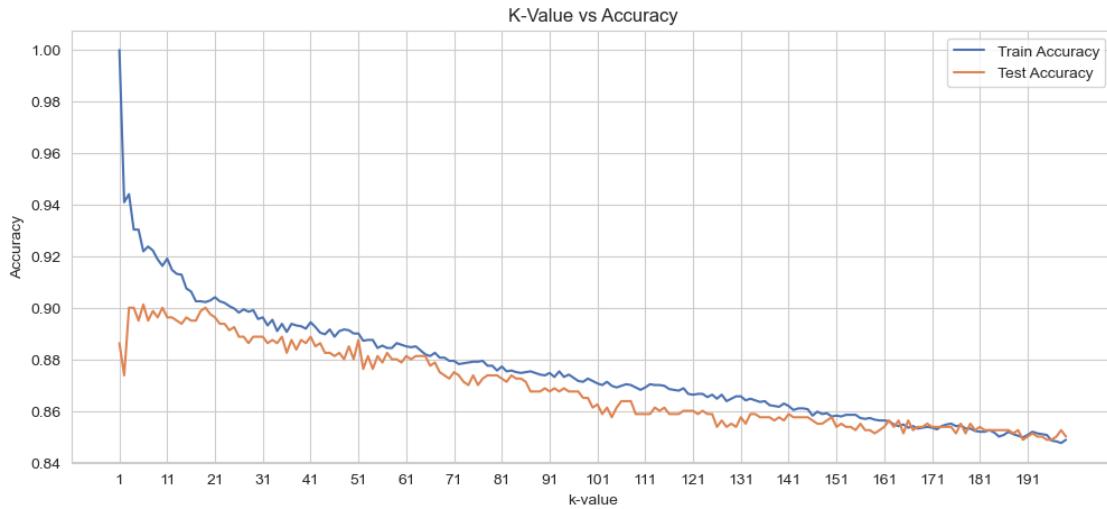
    # Calculate accuracy
    train_accuracies.append(accuracy(y_train_np, y_pred_train))
    test_accuracies.append(accuracy(y_test_np, y_pred_test))

k_values Progress:  0% | 0/199 [00:00<?, ?it/s]

[38]: # We plot the results of our KNN classifier to see how the accuracy changes
      ↴with the value of kvalue on both the training and test data,
      # seeing where the underfitting and overfitting regions are.
      # Plotting
plt.figure(figsize=(12, 5))
plt.plot(k_values, train_accuracies, label='Train Accuracy')
plt.plot(k_values, test_accuracies, label='Test Accuracy')
plt.legend()
plt.title('K-Value vs Accuracy')
plt.xlabel('k-value')
plt.ylabel('Accuracy')
# plt.xticks(list(k_values))  # Ensure that x-axis ticks match k values
tick_interval = 10
plt.xticks(range(1, len(k_values) + 1, tick_interval))
print('Best k-value:', k_values[np.argmax(test_accuracies)])
print('Max accuracy:', max(test_accuracies))
plt.show()

```

Best k-value: 6
Max accuracy: 0.90125



1.0.8 Questions

These are questions in the pdf file they are answered here.

Q1: What kvalue had the best training accuracy?

A1: $k = 1$

Q2: What kvalue had the best testing accuracy?

A2: $k = 6$

Q3: What range of kvalues lead to underfitting?

A3: $k > = 160$, as per the graph

Q4: What range of kvalues lead to overfitting?

A4: $k \in [1,5]$

Q5: Carefully explain and justify your answers by referring to specific regions of the plot and explain your knowledge of overfitting, underfitting, and fitting.

A5: Underfitting is the region where a model has high bias and low variance, and where test accuracy would outperform training accuracy. Overfitting is the region of this plot where a model has high variance and low bias, and where the training accuracy would outperform the test accuracy, as the model would be too specific to the training set. Finally, fitting is the region where the model has good bias and variance, such that the testing accuracy is typically at its highest, but is generally still slightly lower than training accuracy. Therefore, based on the above, in the graph, $k=6$ is the ideal fit of the model, $k<6$ is the region for overfitting, and underfitting occurs for $k > =160$.