# CS145 Howework 2

**Important Note:** HW2 is due on **11:59 PM PT, Oct 30 (Friday, Week 4)**. Please submit through GradeScope.

## Print Out Your Name and UID

**Name: Yunong Ye, UID: 004757414**

## Before You Start

You need to first create HW2 conda environment by the given `cs145hw2.yml` file, which provides the name and necessary packages for this tasks. If you have `conda` properly installed, you may create, activate or deactivate by the following commands:

```
conda env create -f cs145hw2.yml
conda activate hw1
conda deactivate
```

OR

```
conda env create --name NAMEOFYOURCHOICE -f cs145hw2.yml
conda activate NAMEOFYOURCHOICE
conda deactivate
```

To view the list of your environments, use the following command:

```
conda env list
```

More useful information about managing environments can be found [here (https://docs.conda.io/projects/conda/en/latest/user-guide/tasks/manage-environments.html)](https://docs.conda.io/projects/conda/en/latest/user-guide/tasks/manage-environments.html).

You may also quickly review the usage of basic Python and Numpy package, if needed in coding for matrix operations.

In this notebook, you must not delete any code cells in this notebook. If you change any code outside the blocks (such as some important hyperparameters) that you are allowed to edit (between STRART/END YOUR CODE HERE), you need to highlight these changes. You may add some additional cells to help explain your results and observations.

In [66]:
```python
import numpy as np
import pandas as pd
import seaborn as sns
import sys
import random as rd
import matplotlib.pyplot as plt
%load_ext autoreload
%autoreload 2
```

```
The autoreload extension is already loaded. To reload it, use:
  %reload_ext autoreload
```

If you can successfully run the code above, there will be no problem for environment setting.

# 1. Decision trees

This workbook will walk you through a decision tree.

## 1.1 Attribute selection measures

For classification models, misclassification rate is usually used as the final performance measurement. However, for classification trees, when selecting which attribute to split, measurements people often use includes information gain, gain ratio, and Gini index. Let's investigate these different measurements through the following problem.

Note: below shows how to calculate the misclassification rate of a classification tree with $N$ total data points, $K$ classes of the value we want to predict, and $M$ leaf nodes.

In a node $m$, $m = 1, \ldots, M$, let's denote the number of data points using $N_m$, and the number of data points in class k as $N_{mk}$, so the class prediction under majority vote is $j = argmax_k N_{mk}$. The misclassification rate of this node m is $R_m = 1 - \frac{N_{mj}}{N_m}$. The total misclassification rate of the tree will be $R = \frac{\sum_{m=1}^{M} R_m * N_m}{N}$

### Questions

Note: this question is a pure "question answer" problem. You don't need to do any coding.

Suppose our dataset includes a total of 800 people with 400 males and 400 females, and our goal is to do gender classification. Consider two different possible attributes we can split on in a decision tree model. Split on the first attribute results in a node11 with 300 male and 100 female, and a node12 with 100 male and 300 female. Split on the second attribute results in in a node21 with 400 male and 200 female, and a node22 with 200 female only.

1. Which split do you prefer when the measurement is misclassifcation rate and why?
2. What is the entropy in each of these four node?
3. What is the information gain of each of the two splits?

4. Which split do you prefer if the measurement is information gain. Do you see why it is an uncertainty or impurity measurement?
5. What is the gain ratio (normalized information gain) of each of the two splits? Which split do you prefer under this measurement. Do you get the same conclusion as information gain?

**Your answer here:**

Note: you can use several code cells to help you compute the results and answer the questions. Again you don't need to do any coding.

```
In [67]: def entropy(n1,n2):
             total = n1 + n2
             p1 = n1 / total
             p2 = n2 / total
             res = -p1*np.log2(p1) - p2*np.log2(p2)
             return res

         print(entropy(300,100))
         print(entropy(400,200)*3/4)
         print(entropy(400,400))
         print(entropy(600,200))
```

```
0.8112781244591328
0.6887218755408672
1.0
0.8112781244591328
```

Please type your answer here!

answer 1:

The information gain of second attribute is larger than the first so I will prefer use second split. The fisrt split has a misclassification rate of 1/4.
The second split has a misclassification rate of $\frac{3}{4} * \frac{1}{3} = \frac{1}{4}$.
So I think these two are the same in misclassification rate.

answer 2:
The entropy of node11 is $-3/4 * log_2(3/4) - 1/4 * log_2(1/4) = 0.811278$
The entropy of node12 is $-1/4 * log_2(1/4) - 3/4 * log_2(3/4) = 0.811278$
The entropy of node21 is $-2/3 * log_2(2/3) - 1/3 * log_2(1/3) = 0.918296$
The entropy of node22 is $0 * log_2(0) - 1 * log_2(1) = 0$

answer 3:
The entropy is $-1/2 * log_2(1/2) - 1/2 * log_2(1/2) = 1$
The condtiional entropy of first attribute is
$1/2 * entropy(node11) + 1/2 * entropy(node12) = 0.811278$
The conditional entropy of second attribute is

$3/4 * enrtopy(node21) + 1/4 * entropy(node22) = 0.688722$
The information gain of first attribute is 0.188722
The information gain of second attribute is 0.311278

answer 4:
I will choose second attribute using information gain. Yes by calculating entropy it favors the maximum reduction of uncertainty or impurity when choosing attribute.

answer 5:
The splitinfo of first attribute is $-1/2 * log_2(1/2) - 1/2 * log_2(1/2) = 1$
The splitinfo of second attribute is $-3/4 * log_2(3/4) - 1/4 * log_2(1/4) = 0.811278$
The gain ratio of first attreibute is 0.188722.
The gain ratio of second attribute is 0.383688.
I still prefer the second split which is the same as using information gain as measurement.

# 1.2 Coding decision trees

In this section, we are going to use the decision tree model to predict the the animal `type` class of the `zoo` dataset. The dataset has been preprocessed and splited into `decision-tree-train.csv` and `decision-tree-test.csv` for you.

```
In [68]:  from hw2code.decision_tree import DecisionTree
          mytree = DecisionTree()
          mytree.load_data('./data/decision-tree-train.csv','./data/decision-tree-test.csv'
          # As a sanity check, we print out the size of the training data (80, 17) and test
          print('Training data shape: ', mytree.train_data.shape)
          print('Testing data shape:', mytree.test_data.shape)
```

```
Training data shape:  (80, 17)
Testing data shape: (21, 17)
```

## 1.2.1 Infomation gain

Complete the `make_tree` and `compute_info_gain` function in `decision_tree.py`.

Train you model using `info_gain` measure to classify `type` and print the test accuracy.

In [69]:
```python
from hw2code.decision_tree import DecisionTree
mytree = DecisionTree()
mytree.load_data('./data/decision-tree-train.csv','./data/decision-tree-test.csv'
test_acc= 0
#=======================#
# STRART YOUR CODE HERE  #
#=======================#
mytree.train('type', 'info_gain')
test_acc = mytree.test('type')


#=======================#
#   END YOUR CODE HERE   #
#=======================#
print('Test accuracy is: ', test_acc)
```

```
best_feature is:  legs
best_feature is:  fins
best_feature is:  toothed
best_feature is:  eggs
best_feature is:  hair
best_feature is:  hair
best_feature is:  toothed
best_feature is:  aquatic
Test accuracy is:  0.8571428571428571
```

## 1.2.2 Gain ratio

Complete the `compute_gain_ratio` function in `decision_tree.py` .

Train you model using `gain_ratio` measure to classify `type` and print the test accuracy.

In [70]:
```python
mytree = DecisionTree()
mytree.load_data('./data/decision-tree-train.csv','./data/decision-tree-test.csv'
test_acc = 0
#=======================#
# STRART YOUR CODE HERE  #
#=======================#
mytree.train('type', 'gain_ratio')
test_acc = mytree.test('type')
#=======================#
#   END YOUR CODE HERE   #
#=======================#
print('Test accuracy is: ', test_acc)
```

```
best_feature is:  feathers
best_feature is:  backbone
best_feature is:  airborne
best_feature is:  predator
best_feature is:  milk
best_feature is:  fins
best_feature is:  legs
Test accuracy is:  0.8095238095238095
```

**Question**

Which measure do you like the most and why?

**Your answer here:**

I prefer information gain as it has a higher training accuracy. I think this might be becuase that gain ratio prefers unbalanced split while the data here has a close to balanced split so the gain ratio is not needed to prevent huge-number splits getting preferred.

# 2. SVM

This workbook will walk you through a SVM.

## 2.1 Support vectors and decision boundary

Note: for this question you can work entirely in the Jupyter Notebook, no need to edit any .py files.

Consider classifying the following 20 data points in the 2-d plane with class label y

In [71]:
```python
ds = pd.read_csv('data/svm-2d-data.csv')
ds.head()
# This command above will print out the first five data points
# in the dataset with column names as "x1", "x2" and "y"
# You may use command "ds" to show the entire dataset, which contains 20 data poi
```

Out[71]:

|   | x1 | x2 | y |
|---|------|-------|---|
| 0 | 0.52 | -1.00 | 1 |
| 1 | 0.91 | 0.32 | 1 |
| 2 | -1.48 | 1.23 | 1 |
| 3 | 0.01 | 1.44 | 1 |
| 4 | -0.46 | -0.37 | 1 |

Suppose by solving the dual form of the quadratic programming of svm, we can derive the $\alpha_i$'s for each data point as follows: Among $j = 0, 1, \cdots, 19$ (note that the index starts from 0), $\alpha_1 = 0.5084$, $\alpha_5 = 0.4625$, $\alpha_{17} = 0.9709$, and $\alpha_j = 0$ for all other $j$.

**Questions**

1. Which vectors in the training points are support vectors?
2. What is the normal vector of the hyperplane $w$?
3. What is the bias $b$?
4. With the parameters $w$ and $b$, we can now use our SVM to do predictions. What is predicted label of $x_{new} = (2, -0.5)$? Write out your $f(x_{new})$.

5. A plot of the data points has been generated for you. Please change the `support_vec` variable such that only the support vectors are indicated by red circles. Please also fill in the code to draw the decision boundary. Does your prediction of part 4 seems right visually on the plot?

**Your answer here**

Note: you can use several code cells to help you compute the results and answer the questions. Again you don't need to edit any .py files.

Please type your answer here!

answer 1:
Vector 1, vector 5 and vector 17 are the support vectors.

answer 2:
x1 = (0.91,0.32), x5 = (0.41,2.04), x17 = (2.05,1.54)
$$w = \alpha_1 * x_1 + \alpha_5 * x_5 + \alpha_{17} * x_{17} = (2.642614, 2.601374)^T$$

answer 3:
$$b = 1/3 * (y_1 - w^T x_1) + 1/3 * (y_5 - w^T x_5) + 1/3 * (y_{17} - w^T x_{17}) = 1/3 * (-2.23721842) \,\text{-}$$
$$-6.01698926$$

answer 4:
$$f(x_{new}) = w * x_{new} + b = 2.642614 * 2 + 2.601374 * (-0.5) - 6.01698926 = -2.0324482($$
Predict label of $x_{new}$ is -1.

In [72]:
```python
# answer 5
# The prediction seems correct as (2,0,5) belongs in the -1 group.
x1_range = np.arange(-2, 5, 0.5)
x2_range = np.arange(-2, 4., 0.5)

fig, ax = plt.subplots(figsize=(10, 8))
ax = fig.gca()
ax.set_xticks(x1_range)
ax.set_yticks(x2_range)
ax.grid()
ax.scatter(ds['x1'], ds['x2'], c=ds['y'])

support_vec = ds
#=======================#
# STRART YOUR CODE HERE  #
#=======================#
support_vec = ds.loc[[1,5,17]]
x1_values = [i for i in np.arange(-1.5,4.5,0.01)]
w = [2.642614,2.601374]
b = -6.01698926
#w1*x1+w2*x2+b = 0 ---> x2 = (-b-w1x1)/w2
x2_values = []
for x in x1_values:
    x2_values.append((-b-w[0]*x)/w[1])
plt.plot(x1_values,x2_values)

#=======================#
#    END YOUR CODE HERE    #
#=======================#
ax.scatter(support_vec['x1'], support_vec['x2'], marker='o', facecolor='none', s=
sns.despine(ax=ax, left=True, bottom=True, offset=0)
plt.show()
```
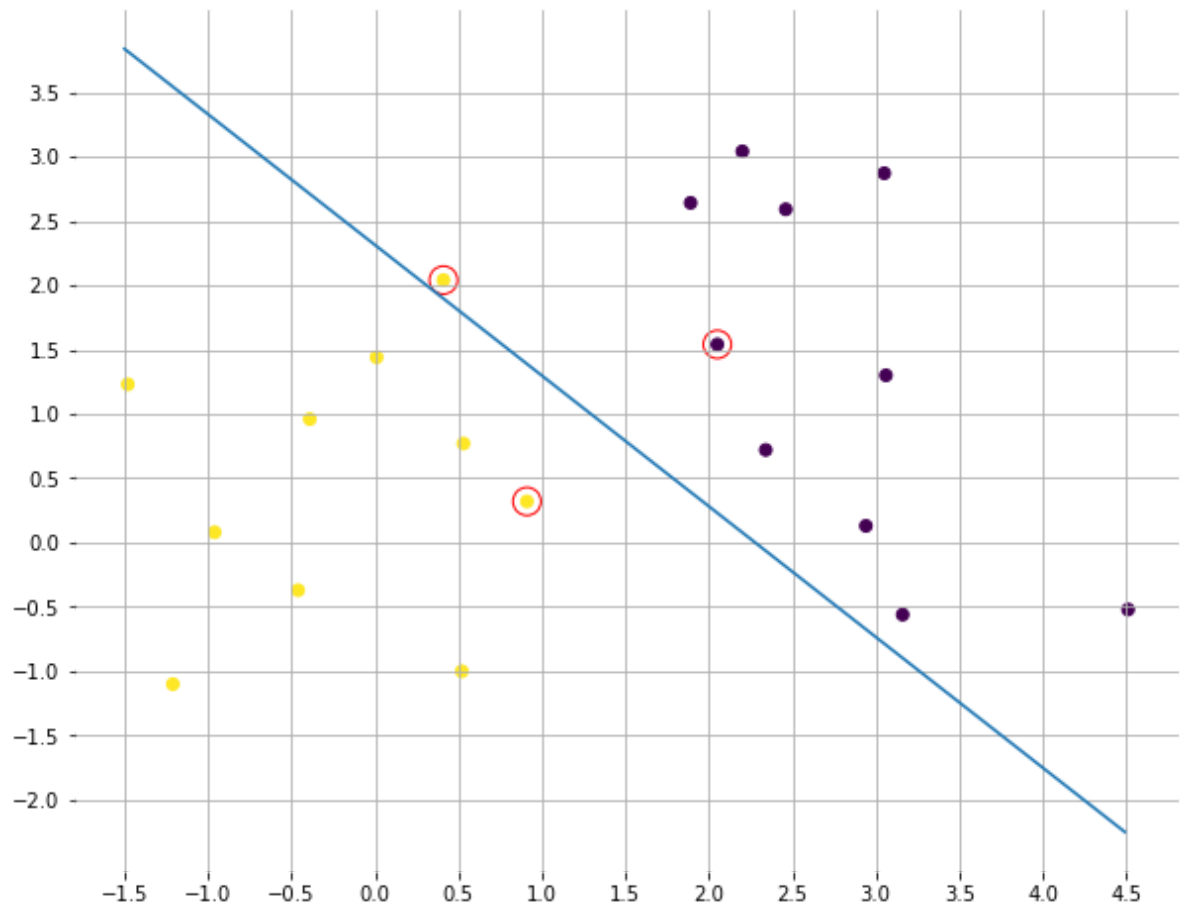
## 2.2 Coding SVM

In this section, we are going to use SVM for classifying the  y  value of 4-dimensional data points.
The dataset has been preprocessed and splited into  svm-train.csv  and  svm-test.csv  for
you.

For this question we are going to use the  cvxopt  package to help us solve the optimization
problem of SVM. You will see it in the .py files, but you don't need to any coding with it. For this
question, you only need to implement the right kernel function, and your kernel matrix  K  in
 svm.py  line 135 will be pluged in the cvxopt optimization problem solver.

For more information about  cvxopt  please refer to http://cvxopt.org/ (http://cvxopt.org/)

```
In [73]:  from hw2code.svm import SVM
          svm = SVM()
          svm.load_data('./data/svm-train.csv', './data/svm-test.csv')
          # As a sanity check, we print out the size of the training data (1098, 4) and (16
          print('Training data shape: ', svm.train_x.shape, svm.train_y.shape)
          print('Testing data shape:', svm.test_x.shape, svm.test_y.shape)
```

```
Training data shape:  (1098, 4) (1098,)
Testing data shape: (274, 4) (274,)
```

## 2.2.1 Linear kernel

## 2.2.1 Linear kernel

Complete the `SVM.predict` and `linear_kernel` function in `svm.py` . Train a hard margin SVM and a soft margin SVM with linear kernel. Print the test accuracy for both cases.

```
In [74]: from hw2code.svm import SVM
         svm_hard = SVM()
         svm_hard.load_data('./data/svm-train.csv', './data/svm-test.csv')
         hard_test_acc = 0
         #=======================#
         # STRART YOUR CODE HERE  #
         #=======================#
         svm_hard.train('linear_kernel')
         hard_test_acc = svm_hard.test()
         #=======================#
         #   END YOUR CODE HERE   #
         #=======================#

         svm_soft = SVM()
         svm_soft.load_data('./data/svm-train.csv', './data/svm-test.csv')
         soft_test_acc = 0
         #=======================#
         # STRART YOUR CODE HERE  #
         #=======================#
         svm_soft.train('linear_kernel',100)
         soft_test_acc = svm_soft.test()
         #=======================#
         #   END YOUR CODE HERE   #
         #=======================#
         print('Hard margin test accuracy is: ', hard_test_acc)
         print('Soft margin test accuracy is: ', soft_test_acc)
```

```
1098 support vectors out of 1098 points
30 support vectors out of 1098 points
Hard margin test accuracy is:  0.5547445255474452
Soft margin test accuracy is:  0.9890510948905109
```

**Questions**

Are these two results similar? Why or why not?

**Your Answer**

No the answers are different. For Hard Margin kernel, the cvxopt cannot find the solution as the test dataset is not linearly separable. So it aborted the job and made all vectors the support vectors which is an ill case. This results in bad test accuracy with only 55%, similar to random guessing. For Soft Margin, the model now tolerates some misclassified data and in this case had a optimal performance with 98% test accuracy.

## 2.2.2 Polynomial kernel

Complete the `polynomial_kernel` function in `svm.py`. Train a soft margin SVM with degree 3 polynomial kernel and parameter `C = 100` for the regularization term. Print the test accuracy.

```
In [75]: from hw2code.svm import SVM
         svm = SVM()
         svm.load_data('./data/svm-train.csv', './data/svm-test.csv')
         test_acc = 0
         #=======================#
         # STRART YOUR CODE HERE  #
         #=======================#
         svm.train("polynomial_kernel", 100)
         test_acc = svm.test()
         #=======================#
         #   END YOUR CODE HERE   #
         #=======================#
         print('Test accuracy is: ', test_acc)
```

```
19 support vectors out of 1098 points
Test accuracy is:  0.9562043795620438
```

**Questions**

Is the result better than linear kernel? Why or why not?

**Your Answer**

It does not seem better compared with the soft margin linear kernel. I am guessing the linearity is enough to classify the test dataset so adding higher-dimension attribute will not improve the accuracy.

## 2.2.3 Gaussian kernel

Complete the `gaussian_kernel` function using the `gaussian_kernel_point` in `svm.py`. Train a soft margin SVM with Gaussian kernel and parameter `C = 100` for the regularization term. Print the test accuracy.

In [76]:
```python
svm = SVM()
svm.load_data('./data/svm-train.csv', './data/svm-test.csv')
test_acc = 0
#=======================#
# STRART YOUR CODE HERE  #
#=======================#
svm.train("gaussian_kernel", 100)
test_acc = svm.test()
#=======================#
#   END YOUR CODE HERE   #
#=======================#
print('Test accuracy is: ', test_acc)
```

```
35 support vectors out of 1098 points
Test accuracy is:  1.0
```

## Questions

1. Is the result better than linear kernel and polynomial kernel? Why or why not?
2. Which one of these four models do you like the most and why?
3. (Bonus question, optional) Can you come up with a vectorized implementation of `gaussian_kernel` without calling `gaussian_kernel_point` ? Fill that in svm.py.

## Your Answer

Please write down your answers and/or observations here

answer 1: The result is better than both, probably because the test dataset has a normal distribution tendency so it fits the model well.

answer 2: I actually like the Soft Margin Linear Kernel as this reaches a high test accuracy with much less computation.

# End of Homework 2 :)

After you've finished the homework, please print out the entire `ipynb` notebook and two `py` files into one PDF file. Make sure you include the output of code cells and answers for questions. Prepare submit it to GradeScope. Also this time remember assign the pages to the questions on GradeScope

```python
1  import pandas as pd
2  import numpy as np
3  from pprint import pprint
4  import sys
5
6  # Reads the data from CSV files, each attribute column can be obtained via its name,
   e.g., y = data['y']
7  def getDataframe(filePath):
8      data = pd.read_csv(filePath)
9      return data
10
11 # predicted_y and y are the predicted and actual y values respectively as numpy
   arrays
12 # function prints the accuracy
13 def compute_accuracy(predicted_y, y):
14     acc = 100.0
15     acc = np.sum(predicted_y == y)/predicted_y.shape[0]
16     return acc
17
18 #Compute entropy according to y distribution
19 def compute_entropy(y):
20     entropy = 0.0
21     elements,counts = np.unique(y, return_counts = True)
22     n = y.shape[0]
23
24     for i in range(len(elements)):
25         prob = counts[i]/n
26         if prob!= 0:
27             entropy -= prob * np.log2(prob)
28     return entropy
29
30 #att_name: attribute name; y_name: the target attribute name for classification
31 def compute_info_gain(data, att_name, y_name):
32     info_gain = 0.0
33
34     #Calculate the values and the corresponding counts for the select attribute
35     vals, counts = np.unique(data[att_name], return_counts=True)
36     total_counts = np.sum(counts)
37
38     #Calculate the conditional entropy
39     #=======================#
40     # STRART YOUR CODE HERE  #
41     #=======================#
42     info = compute_entropy(data[y_name])
43     entropy = 0
44     for i in range(len(vals)):
45         probability = counts[i] / total_counts
46         val_data = data.loc[data[att_name] == vals[i]]
47         entropy += probability * compute_entropy(val_data[y_name])
48     info_gain = info - entropy
49     #=======================#
50     #   END YOUR CODE HERE   #
51     #=======================#
52
53     return info_gain
54
55
56 def comput_gain_ratio(data, att_name, y_name):
57     gain_ratio = 0.0
58     #Calculate the values and the corresponding counts for the select attribute
```

```python
59         vals, counts = np.unique(data[att_name], return_counts=True)
60         total_counts = np.sum(counts)
61
62         #Calculate the information for the selected attribute
63         att_info = 0.0
64         #=======================#
65         # STRART YOUR CODE HERE  #
66         #=======================#
67         for i in range(len(counts)):
68             probability = counts[i] / total_counts
69             att_info -= probability * np.log2(probability)
70         #=======================#
71         #   END YOUR CODE HERE    #
72         #=======================#
73         gain_ratio = 0.0 if np.abs(att_info) < 1e-9 else min(1, compute_info_gain(data,
    att_name, y_name) / att_info)
74         return gain_ratio
75
76 # Class of the decision tree model based on the ID3 algorithm
77 class DecisionTree(object):
78     def __init__(self):
79         self.train_data = pd.DataFrame()
80         self.test_data = pd.DataFrame()
81
82     def load_data(self, train_file, test_file):
83         self.train_data = getDataframe(train_file)
84         self.test_data = getDataframe(test_file)
85
86     def train(self, y_name, measure, parent_node_class= None):
87         self.y_name = y_name
88         self.measure = measure
89         self.tree = self.make_tree(self.train_data, parent_node_class)
90
91     def make_tree(self, train_data, parent_node_class = None):
92         data = train_data
93         features = data.drop(self.y_name, axis = 1).columns.values
94         measure =  self.measure
95         #Stopping condition 1: If all target_values have the same value, return this
    value
96         if len(np.unique(data[self.y_name])) <= 1:
97             leaf_value = -1
98             #=======================#
99             # STRART YOUR CODE HERE  #
100            #=======================#
101            leaf_value = data[self.y_name].values[0]
102            #=======================#
103            #   END YOUR CODE HERE    #
104            #=======================#
105            return leaf_value
106
107        #Stopping condition 2: If the dataset is empty, return the parent_node_class
108        elif len(data)== 0:
109            return parent_node_class
110
111        #Stopping condition 3: If the feature space is empty, return the majority
    class
112        elif len(features) == 0:
113            return np.unique(data[self.y_name])
    [np.argmax(np.unique(data[y_name],return_counts=True)[1])]
114
```

```python
115              # Not a leaf node, create an internal node
116              else:
117                  #Set the default value for this node --> The mode target feature value of
     the current node
118                  parent_node_class = np.unique(data[self.y_name])
     [np.argmax(np.unique(data[self.y_name],return_counts=True)[1])]
119
120                  #Select the feature which best splits the dataset
121                  if measure == 'info_gain':
122                      item_values = [compute_info_gain(data, feature, self.y_name) for
     feature in features] #Return the information gain values for the features in the
     dataset
123                  elif measure == 'gain_ratio':
124                      item_values = [comput_gain_ratio(data, feature, self.y_name) for
     feature in features] #Return the gain_ratio for the features in the dataset
125                  else:
126                      raise ValueError("kernel not recognized")
127
128                  best_feature_index = np.argmax(item_values)
129                  best_feature = features[best_feature_index]
130                  print('best_feature is: ', best_feature)
131
132                  #Create the tree structure. The root gets the name of the feature
     (best_feature)
133                  tree = {best_feature:{}}
134
135
136              #Grow a branch under the root node for each possible value of the root node
     feature
137
138              for value in np.unique(data[best_feature]):
139                  #Split the dataset along the value of the feature with the largest
     information gain and therwith create sub_datasets
140                  sub_data = data.where(data[best_feature] == value).dropna()
141
142                  #Remove the selected feature from the feature space
143                  sub_data = sub_data.drop(best_feature, axis = 1)
144
145                  #Call the ID3 algorithm for each of those sub_datasets with the new
     parameters --> Here the recursion comes in!
146                  subtree = self.make_tree(sub_data, parent_node_class)
147
148                  #Add the sub tree, grown from the sub_dataset to the tree under the root
     node
149                  tree[best_feature][value] = subtree
150
151              return tree
152
153
154      def test(self, y_name):
155          accuracy = self.classify(self.test_data, y_name)
156          return accuracy
157
158      def classify(self, test_data, y_name):
159          #Create new query instances by simply removing the target feature column from
     the test dataset and
160          #convert it to a dictionary
161          test_x = test_data.drop(y_name, axis=1)
162          test_y = test_data[y_name]
163
```

```
164         n = test_data.shape[0]
165         predicted_y = np.zeros(n)
166
167         #Calculate the prediction accuracy
168         for i in range(n):
169             predicted_y[i] = DecisionTree.predict(self.tree, test_x.iloc[i])
170
171         output = np.zeros((n,2))
172         output[:,0] = test_y
173         output[:,1] = predicted_y
174         accuracy = compute_accuracy(predicted_y, test_y.values)
175         return accuracy
176
177     def predict(tree, query):
178         # find the root attribute
179         default = -1
180         for root_name in list(tree.keys()):
181             try:
182                 subtree = tree[root_name][query[root_name]]
183             except:
184                 return default ## root_name does not appear in query attribute list
    (it is an error!)
185
186             ##if subtree is still a dictionary, recursively test next attribute
187             if isinstance(subtree,dict):
188                 return DecisionTree.predict(subtree, query)
189             else:
190                 leaf = subtree
191                 return leaf
192
193
```

```python
1  import numpy as np
2  from numpy import linalg
3  import cvxopt
4  import cvxopt.solvers
5  import sys
6  import pandas as pd
7  cvxopt.solvers.options['show_progress'] = False
8
9  # Reads the data from CSV files, converts it into Dataframe and returns x and y
   dataframes
10 def getDataframe(filePath):
11     dataframe = pd.read_csv(filePath)
12     y = dataframe['y']
13     x = dataframe.drop('y', axis=1)
14     y = y*2 -1.0
15     return x.to_numpy(), y.to_numpy()
16
17 def compute_accuracy(predicted_y, y):
18     acc = 100.0
19     acc = np.sum(predicted_y == y)/predicted_y.shape[0]
20     return acc
21
22 def gaussian_kernel_point(x, y, sigma=5.0):
23     return np.exp(-linalg.norm(x-y)**2 / (2 * (sigma ** 2)))
24
25 def linear_kernel(X, Y=None):
26     Y = X if Y is None else Y
27     m = X.shape[0]
28     n = Y.shape[0]
29     assert X.shape[1] == Y.shape[1]
30     kernel_matrix = np.zeros((m, n))
31     #======================#
32     # STRART YOUR CODE HERE  #
33     #======================#
34     for i in range(m):
35         for j in range(n):
36             kernel_matrix[i,j] = np.dot(X[i], Y[j])
37     #======================#
38     #   END YOUR CODE HERE   #
39     #======================#
40     return kernel_matrix
41
42 def polynomial_kernel(X, Y=None, degree=3):
43     Y = X if Y is None else Y
44     m = X.shape[0]
45     n = Y.shape[0]
46     assert X.shape[1] == Y.shape[1]
47     kernel_matrix = np.zeros((m, n))
48     #======================#
49     # STRART YOUR CODE HERE  #
50     #======================#
51     for i in range(m):
52         for j in range(n):
53             kernel_matrix[i,j] = (1+np.dot(X[i],Y[j])) ** degree
54     #======================#
55     #   END YOUR CODE HERE   #
56     #======================#
57     return kernel_matrix
58
59 def gaussian_kernel(X, Y=None, sigma=5.0):
```

```python
60      Y = X if Y is None else Y
61      m = X.shape[0]
62      n = Y.shape[0]
63      assert X.shape[1] == Y.shape[1]
64      kernel_matrix = np.zeros((m, n))
65      #=======================#
66      # STRART YOUR CODE HERE  #
67      #=======================#
68      for i in range(m):
69          for j in range(n):
70              kernel_matrix[i,j] = np.exp((-np.linalg.norm(X[i]-Y[j])**2)/ (2*
    (sigma**2)))
71      #=======================#
72      #   END YOUR CODE HERE   #
73      #=======================#
74      return kernel_matrix
75
76
77 # Bonus question: vectorized implementation of Gaussian kernel
78 # If you decide to do the bonus question, comment the gaussian_kernel function above,
79 # then implement and uncomment this one.
80 # def gaussian_kernel(X, Y=None, sigma=5.0):
81 #      return
82
83 class SVM(object):
84      def __init__(self):
85          self.train_x = pd.DataFrame()
86          self.train_y = pd.DataFrame()
87          self.test_x = pd.DataFrame()
88          self.test_y = pd.DataFrame()
89          self.kernel_name = None
90          self.kernel = None
91
92      def load_data(self, train_file, test_file):
93          self.train_x, self.train_y = getDataframe(train_file)
94          self.test_x, self.test_y = getDataframe(test_file)
95
96
97      def train(self, kernel_name='linear_kernel', C=None):
98          self.kernel_name = kernel_name
99          if(kernel_name == 'linear_kernel'):
100             self.kernel = linear_kernel
101         elif(kernel_name == 'polynomial_kernel'):
102             self.kernel = polynomial_kernel
103         elif(kernel_name == 'gaussian_kernel'):
104             self.kernel = gaussian_kernel
105         else:
106             raise ValueError("kernel not recognized")
107
108         self.C = C
109         if self.C is not None:
110             self.C = float(self.C)
111
112         self.fit(self.train_x, self.train_y)
113
114     # predict labels for test dataset
115     def predict(self, X):
116         if self.w is not None: ## linear case
117             n = X.shape[0]
118             predicted_y = np.zeros(n)
```

```python
119            #=======================#
120            # STRART YOUR CODE HERE   #
121            #=======================#
122            predicted_y = np.dot(self.w, np.transpose(X)) + self.b
123            #=======================#
124            #   END YOUR CODE HERE    #
125            #=======================#
126            return predicted_y
127
128        else: ## non-linear case
129            n = X.shape[0]
130            predicted_y = np.zeros(n)
131            #=======================#
132            # STRART YOUR CODE HERE   #
133            #=======================#
134            kernel_result = self.kernel(X,self.sv)
135            for i in range(n):
136                for j in range(self.sv.shape[0]):
137                    predicted_y[i] += self.a[j]*self.sv_y[j]*kernel_result[i,j]
138            #=======================#
139            #   END YOUR CODE HERE    #
140            #=======================#
141            return predicted_y
142
143    #===============================================#
144    # Please DON'T change any code below this line!  #
145    #===============================================#
146    def fit(self, X, y):
147        n_samples, n_features = X.shape
148        # Kernel matrix
149        K = self.kernel(X)
150
151        # dealing with dual form quadratic optimization
152        P = cvxopt.matrix(np.outer(y,y) * K)
153        q = cvxopt.matrix(np.ones(n_samples) * -1)
154        A = cvxopt.matrix(y, (1,n_samples),'d')
155        b = cvxopt.matrix(0.0)
156
157        if self.C is None:
158            G = cvxopt.matrix(np.diag(np.ones(n_samples) * -1))
159            h = cvxopt.matrix(np.zeros(n_samples))
160        else:
161            tmp1 = np.diag(np.ones(n_samples) * -1)
162            tmp2 = np.identity(n_samples)
163            G = cvxopt.matrix(np.vstack((tmp1, tmp2)))
164            tmp1 = np.zeros(n_samples)
165            tmp2 = np.ones(n_samples) * self.C
166            h = cvxopt.matrix(np.hstack((tmp1, tmp2)))
167
168        # solve QP problem
169        solution = cvxopt.solvers.qp(P, q, G, h, A, b)
170        # Lagrange multipliers
171        a = np.ravel(solution['x'])
172
173        # Support vectors have non zero lagrange multipliers
174        sv = a > 1e-5
175        ind = np.arange(len(a))[sv]
176        self.a = a[sv]
177        self.sv = X[sv]
178        self.sv_y = y[sv]
```

```python
179
180          print("%d support vectors out of %d points" % (len(self.a), n_samples))
181
182          # Intercept via average calculating b over support vectors
183          self.b = 0
184          for n in range(len(self.a)):
185              self.b += self.sv_y[n]
186              self.b -= np.sum(self.a * self.sv_y * K[ind[n],sv])
187          self.b /= len(self.a)
188
189          # Weight vector
190          if self.kernel_name == 'linear_kernel':
191              self.w = np.zeros(n_features)
192              for n in range(len(self.a)):
193                  self.w += self.a[n] * self.sv_y[n] * self.sv[n]
194          else:
195              self.w = None
196
197
198      def test(self):
199          accuracy = self.classify(self.test_x, self.test_y)
200          return accuracy
201
202      def classify(self, X, y):
203          predicted_y = np.sign(self.predict(X))
204          accuracy = compute_accuracy(predicted_y, y)
205          return accuracy
```