

CS145 Homework 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 NAMEOFOURCHOICE -f cs145hw2.yml
conda activate NAMEOFOURCHOICE
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 START/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]: 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 , $m = 1, \dots, M$, 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 $j = \operatorname{argmax}_k N_{mk}$.

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 misclassification 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 first 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 conditional entropy of first attribute is

$1/2 * \text{entropy}(\text{node11}) + 1/2 * \text{entropy}(\text{node12}) = 0.811278$

The conditional entropy of second attribute is

$$3/4 * \text{entropy}(\text{node21}) + 1/4 * \text{entropy}(\text{node22}) = 0.688722$$

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 attribute 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 splitted 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]: 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 your model using `gain_ratio` measure to classify `type` and print the test accuracy.

```
In [70]: 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 because 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]: 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 points
```

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 α_i 's for each data point as follows: Among $j = 0, 1, \dots, 19$ (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:

$x_1 = (0.91, 0.32)$, $x_5 = (0.41, 2.04)$, $x_{17} = (2.05, 1.54)$

$w = \alpha_1 * x_1 + \alpha_5 * x_5 + \alpha_{17} * x_{17} = (2.642614, 2.601374)^T$

answer 3:

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

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.

```

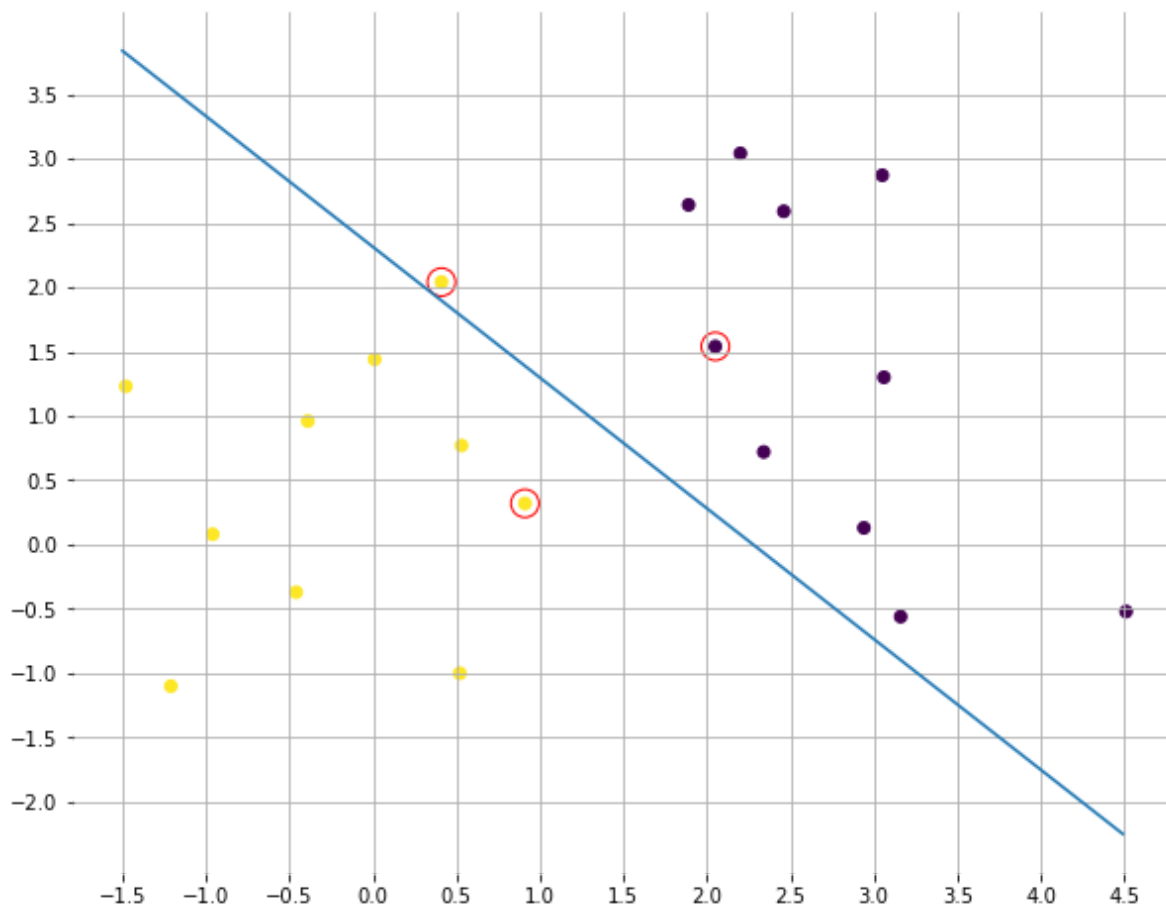
In [72]: # answer 5
# The prediction seems correct as (2,0,5) belongs in the -1 group.
x1_range = np.arange(-2, 5, 0.5)
x2_range = np.arange(-2, 4., 0.5)

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 splitted 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 plugged 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 (1098, 4)
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

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]: 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