

Dustin Van Tate Testa

1. Here's a partial solution using guesswork and what was given in the assignment because none of this makes sense. I'm gonna go kms now, thanks...

Handwritten notes on a piece of graph paper showing a partial solution to a problem. The notes include the formula $E = -\sum_{j=1}^K p_j \log_2(p_j)$ and a note "obtaining maximum value when $p_j = p_k = \frac{1}{K}$ ". There are also some scribbles and a note "I do not know how to take partial derivative of all infinite series of equations..."

Sent 5 days ago. Follow up?

Dismiss

[Ext] 13929.13930.202110 (Fall 2020 - Intro to Machine Learning (CS-484)): HW 3

Oct 10, 2020, 8:33 PM (5 days ago)

oabdulai@hawk.iit.edu

Hi all, I have received a number of emails about question 5(e) so I will address it for everyone. Once you obtain FPR and TPR from your predictions, use this to

Dustin Testa <dtesta1@hawk.iit.edu>

to Olu

Oct 10, 2020, 10:26 PM (5 days ago)

For question 1, in what context are the variables defined? It's been over two years since I've taken calculus and I'm not sure how to take the partial derivative of what seems to be an infinite series that isn't asymptotic?

Learning

Autumn 2020 Assignment 3

Question 1 (10 points)

Prove that $E = -\sum_{j=1}^K p_j \log_2(p_j)$ obtains its maximum value when $p_j = p_k = 1/K$.

Hint: (1) re-express $E = -\sum_{j=1}^K p_j \log_2(p_j) = -\sum_{j=1}^K p_j \log_2(p_j) - \sum_{j=1}^K p_j \log_2(p_j)$ (2) use this equality $\sum_{j=1}^K p_j = 1$ in calculating the partial derivatives $\partial E / \partial p_j, j=1, \dots, K-1$, and (3) solve the equations $\partial E / \partial p_j = 0, j=1, \dots, K-1$.

xxx

Reply Forward

If you tell me that you are too busy during office hours but don't respond to emails how am I going to pass this class?

2. The probability that any item in the dataset is misclassified is defined as:

$$P(\text{misclassification}) = \sum_{i=0}^N P(\text{put in wrong category}_i) \times P(\text{pulled from wrong category}_i)$$

Gini Impurity is defined as the probability that a randomly selected item in the dataset is put in wrong category

$$I_G = \sum_{i=0}^N p_i^2$$

Where p_i is the fraction of items having the same classification as the randomly selected one

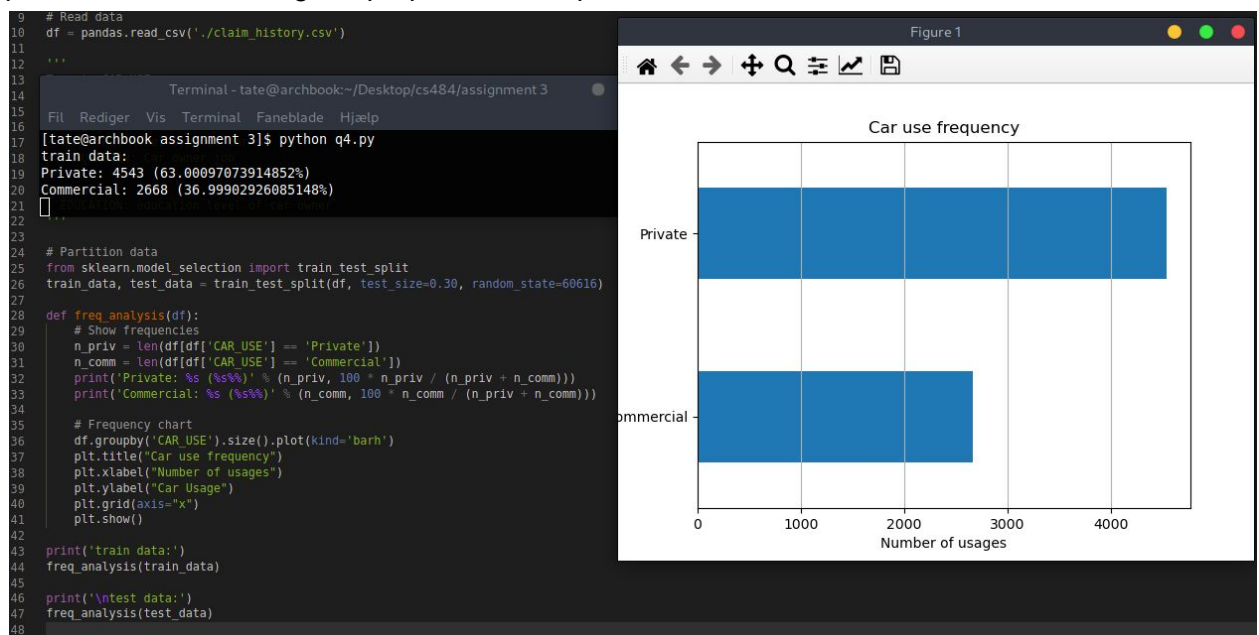
Thus because the set is entirely random $P(\text{put in wrong category})$ and $P(\text{pulled from wrong category})$ will be the same value for the entire dataset

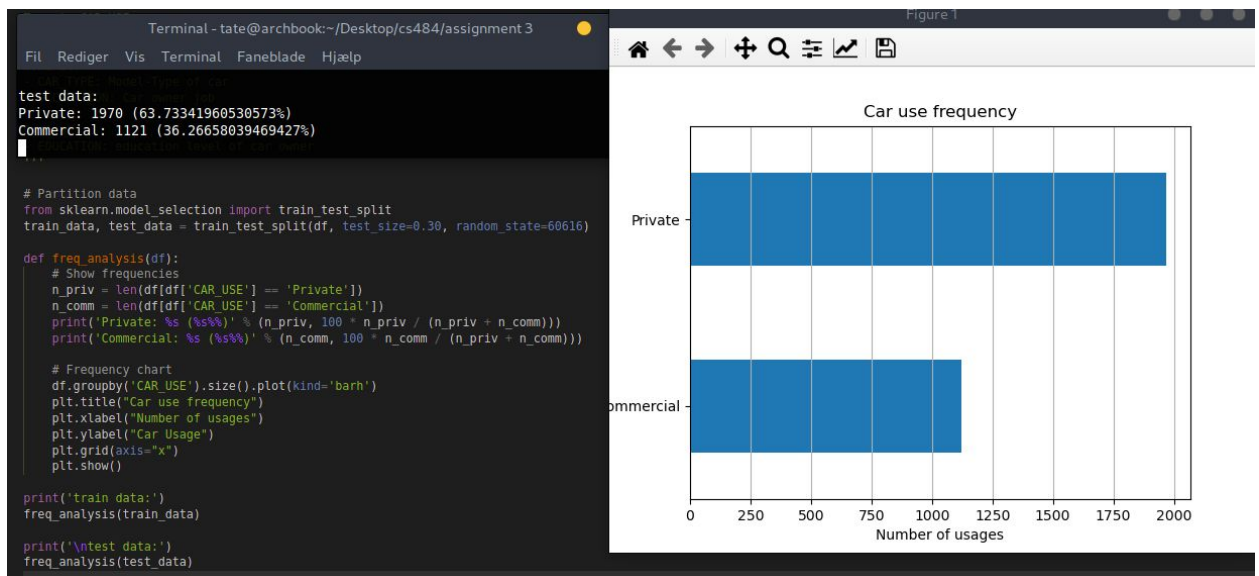
This makes the two equivalent.

3. If the set contains only a single misclassified node, the Gini Impurity would have to be 1 as the only item to select would be misclassified, giving it a 100% chance of being as such.

4.

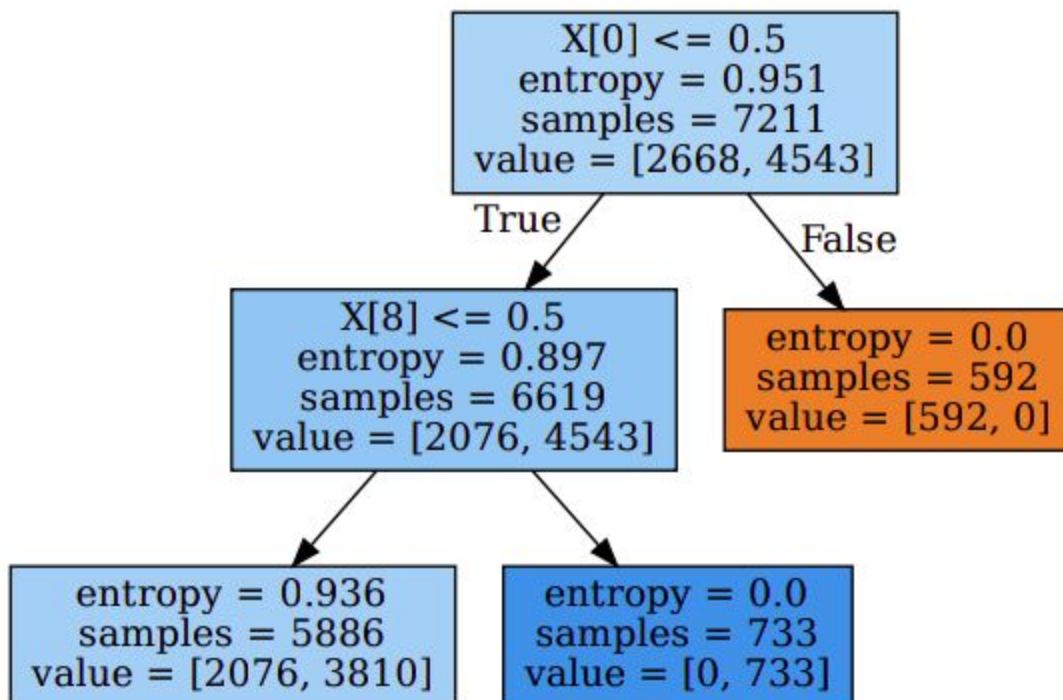
- a. The training data consisted of 63% private vehicles and the test data consisted of 63.7% private vehicles making the proportions comparable





- b. Because both groups have approximately the same proportions, the probability that a private vehicle will fall in the test group is just the ratio of the sizes. Therefore: $30/100 = 30\%$

5. This graph was generated using the the code in tree.py



- a. The entropy of first node is 0.951
- b. By enumerating over the column names for the one hot encoded training data we can see what it's checking for in each branch

- i. In the root node it's checking if the vehicle is a 'Panel Truck' and classifying the vehicle as commercial if so
- ii. Otherwise it goes to the second branch where it checks if the driver's Occupation is 'Lawyer'
- iii. These specific checks make sense as Panel Trucks are almost never used as private vehicles and Lawyers seldom drive commercial vehicles

```
0 print(graph)
1 graph.render('graph')
2
3 # print column names
4 for i, col in enumerate(one_hot_inputs):
5     print('X[%s] = %s' % (i, col))
6     if i > 10: break
```

```
[tate@archbook assignment 3]$ python tree.py
<graphviz.files.Source object at 0x7f4248dd2e50>
X[0] = CAR_TYPE_Panel Truck
X[1] = CAR_TYPE_Pickup
X[2] = CAR_TYPE_SUV
X[3] = CAR_TYPE_Sports Car
X[4] = CAR_TYPE_Van
X[5] = OCCUPATION_Clerical
X[6] = OCCUPATION_Doctor
X[7] = OCCUPATION_Home Maker
X[8] = OCCUPATION_Lawyer
X[9] = OCCUPATION_Manager
X[10] = OCCUPATION_Professional
X[11] = OCCUPATION_Student
[tate@archbook assignment 3]$
```

- c. 0.897 according to the graphviz plot at top
- d. Description of nodes:
 - i. From the left

Rule	Count	Entropy		
Occupation != Lawyer Car != panel truck	5886	0.936		
Occupation == Lawyer Car != panel truck	733	0.0		
Car == panel truck	592	0.0		

- e. https://en.wikipedia.org/wiki/Kolmogorov%E2%80%93Smirnov_test
<https://docs.scipy.org/doc/scipy/reference/generated/scipy.stats.kstest.html>
- f. Our tree only correctly classified

6.

- a. The misclassification rate for the test data was found to be about 27.82%

```
82
83 # Find misclassification rate for test data
84 one_hot_inputs = pandas.get_dummies(
85     test_data[['CAR_TYPE', 'OCCUPATION', 'EDUCATION']],
86     drop_first = True)
87 targets = test_data[['CAR_USE']]
88 print('misclassification rate: %s%%' % (100 - 100 * dt.sco
89
```

```
X[8] = OCCUPATION_Lawyer
X[9] = OCCUPATION_Manager
X[10] = OCCUPATION_Professional
X[11] = OCCUPATION_Student
misclassification rate: 27.822711096732448%
[tate@archbook assignment 3]$
```

Dustin Van Tate Testa

b.