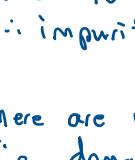


# Assignment 1.4

Thursday, February 1, 2024

6:35 PM

4. I a)



Using code in q-4.py, 81 out of 314 women did not survive and 109 out of 577 men survived  
 $\therefore \text{impurity} = \frac{81}{314+577} + \frac{109}{314+577}$

$$= 21.324\% \rightarrow \text{percentage of not correct classification}$$

b) There are multiple ways to split by a numerical value (i.e., domain of  $x \in [0, 100]$ ), can choose any value between 0 and 100 to split, impurity becomes an important factor when selecting the threshold to split a numerical variable (Want to minimize impurity with chosen threshold)

c) Code in q-4.py

Splitting on age  $\leq 25$ : 255 young died, 197 old survived

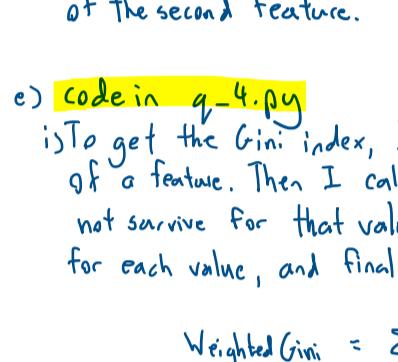
$$\text{impurity: } \frac{255+197}{314+577} = 0.4738 = 47.38\%$$

Splitting on age  $\geq 65$ : 539 young died, 1 old survived  
 $\text{impurity: } \frac{539+1}{314+577} = 0.606060 = 60.60\%$

Since splitting age at  $\leq 25$  leads to less impurity, it results in a better decision tree

d) Split: Gender  $\rightarrow \leq 25$

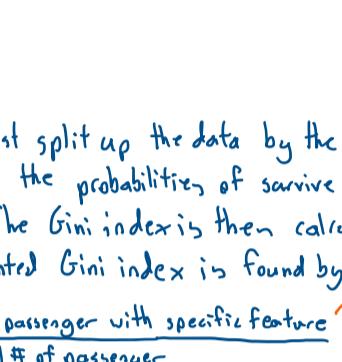
Code in q-4.py



$$\text{impurities: } \frac{422}{891} = 47.36\%$$

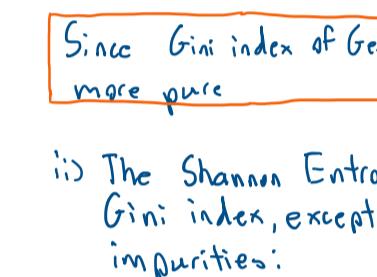
Split: Gender  $\rightarrow \geq 65$

Code in q-4.py



$$\text{impurities: } \frac{549}{891} = 60.61\%$$

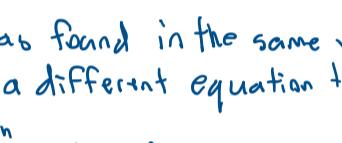
Split:  $\leq 25 \rightarrow \text{Gender}$



$$\text{impurities: } \frac{190}{891} = 21.324\%$$

Split:  $\geq 65 \rightarrow \text{Gender}$

Code in q-4.py



$$\text{impurities: } \frac{190}{891} = 21.324\%$$

The impurities are the same as the impurities calculated in parts a) and c). This is because the decision of whether someone survives is solely based on the second feature. Thus, it does not matter what the split is with the first feature, the impurities of the model is solely dependent on the split of the second feature.

e) code in q-4.py

i) To get the Gini index, I first split up the data by the values of a feature. Then I calculated the probabilities of survive and not survive for that value. The Gini index is then calculated for each value, and final weighted Gini index is found by:

$$\text{Weighted Gini} = \sum \frac{\#\text{of passenger with specific feature}}{\text{total \# of passenger}} \times \text{Gini index of specific feature}$$

$$\text{ie: Male/Female}$$

$$\text{Gini index of specific feature}$$

$$\text{ie: } \frac{\#\text{of male}}{\#\text{of passenger}} \times \text{Gini male} + \frac{\#\text{of Female}}{\#\text{of passenger}} \times \text{Gini female}$$

$$\downarrow 1 - (\text{P(male survive})^2 + \text{P(male not survive})^2)$$

Result: Weighted Gini of Gender: 0.333  
 Weighted Gini of Age: 0.417

Since Gini index of Gender is lower, the split is better and more pure

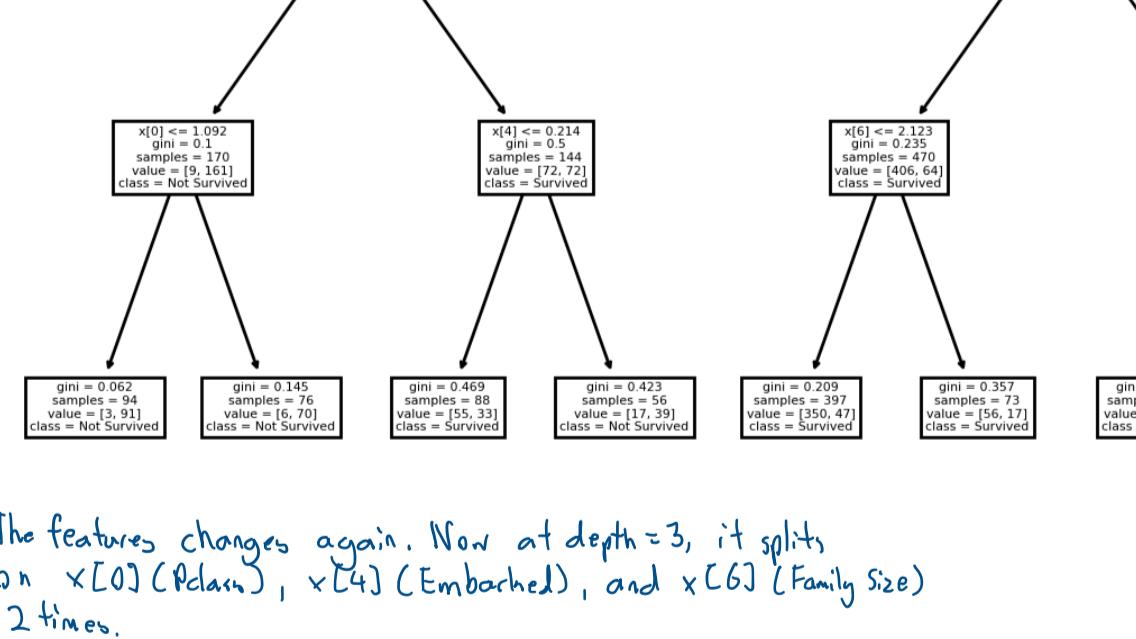
ii) The Shannon Entropy was found in the same way as the Gini index, except with a different equation to measure impurities:

$$H(x) = - \sum_{i=1}^n p_i \log_2 p_i$$

Results: Weighted Shannon Entropy of Gender: 0.7405  
 Weighted Shannon Entropy of Age: 0.8639

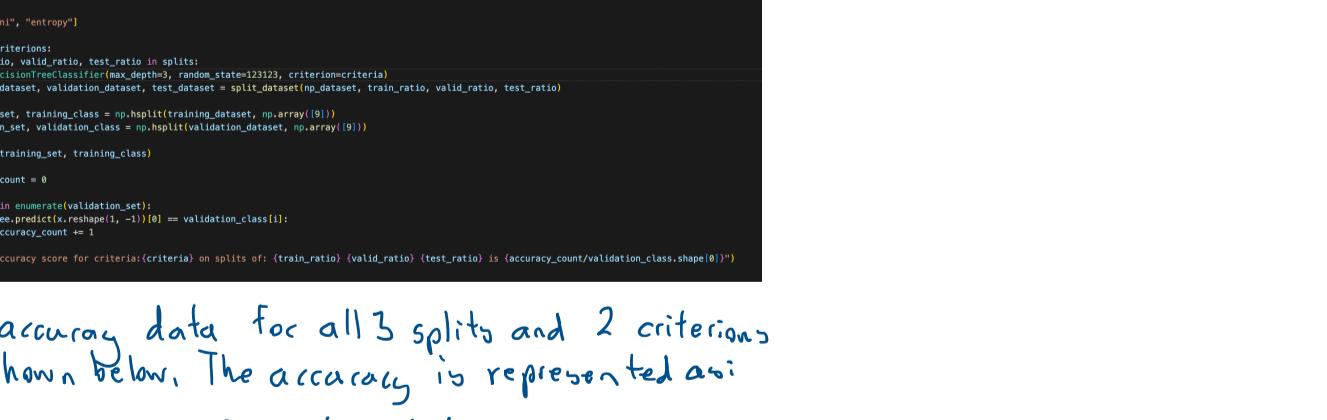
Shannon Entropy shows that Gender feature is more pure and will get a better splits which matches with the results of the Gini Index

II. a) Code in q-4.py



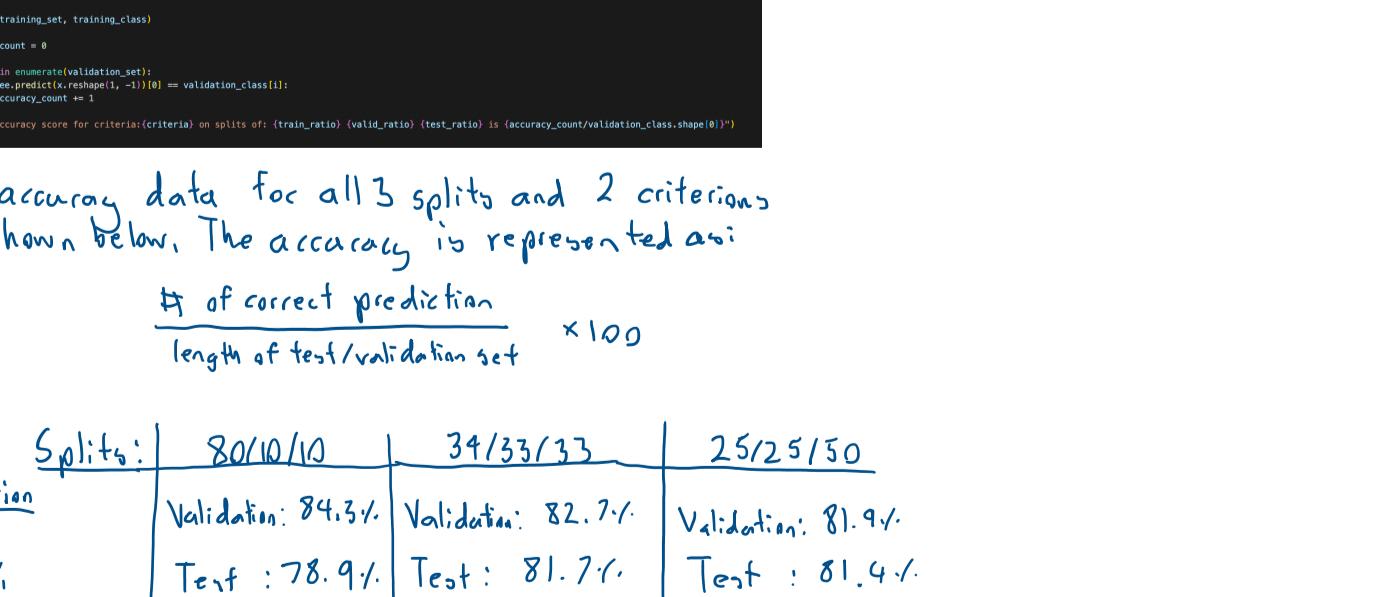
Features  $x[0]$  (Pclass/social class) and  $x[5]$  (has cabin) used to split at depth=2

b)



The features do change before it splits on  $x[1]$  (Age) 3 times and  $x[3]$  (Fare) once. Now it splits on Age 2 times and Fare both 2 times

c)



The features changes again. Now at depth=3, it splits on  $x[0]$  (Pclass/social class),  $x[4]$  (Embarked), and  $x[6]$  (Family size) 2 times.

III. Code in q-4.py

# part 3

if \_\_name\_\_ == "\_\_main\_\_":
 np.random.seed(42)
 train\_index = math.floor(validation\_ratio \* dataset.shape[0])
 validation\_index = dataset.shape[0] - train\_index
 validation\_dataset = dataset[train\_index:]
 test\_index = dataset[:train\_index]
 validation\_set = dataset[:validation\_index]
 test\_set = dataset[validation\_index:]

train\_set, validation\_set, test\_set = train\_test\_split(dataset, validation\_ratio=0.2, random\_state=42)

train\_set, validation\_set, test\_set = shuffle(train\_set, validation\_set, test\_set)

train\_set, validation\_set, test\_set = train\_set.reset\_index(drop=True), validation\_set.reset\_index(drop=True), test\_set.reset\_index(drop=True)

train\_set, validation\_set, test\_set = train\_set.drop(['Name'], axis=1), validation\_set.drop(['Name'], axis=1), test\_set.drop(['Name'], axis=1)

train\_set, validation\_set, test\_set = train\_set.fillna('Unknown'), validation\_set.fillna('Unknown'), test\_set.fillna('Unknown')

train\_set, validation\_set, test\_set = train\_set.drop(['Cabin'], axis=1), validation\_set.drop(['Cabin'], axis=1), test\_set.drop(['Cabin'], axis=1)

train\_set, validation\_set, test\_set = train\_set.drop(['Ticket'], axis=1), validation\_set.drop(['Ticket'], axis=1), test\_set.drop(['Ticket'], axis=1)

train\_set, validation\_set, test\_set = train\_set.drop(['PassengerId'], axis=1), validation\_set.drop(['PassengerId'], axis=1), test\_set.drop(['PassengerId'], axis=1)

train\_set, validation\_set, test\_set = train\_set.drop(['SibSp'], axis=1), validation\_set.drop(['SibSp'], axis=1), test\_set.drop(['SibSp'], axis=1)

train\_set, validation\_set, test\_set = train\_set.drop(['Parch'], axis=1), validation\_set.drop(['Parch'], axis=1), test\_set.drop(['Parch'], axis=1)

train\_set, validation\_set, test\_set = train\_set.drop(['Fare'], axis=1), validation\_set.drop(['Fare'], axis=1), test\_set.drop(['Fare'], axis=1)

train\_set, validation\_set, test\_set = train\_set.drop(['Age'], axis=1), validation\_set.drop(['Age'], axis=1), test\_set.drop(['Age'], axis=1)

train\_set, validation\_set, test\_set = train\_set.drop(['Sex'], axis=1), validation\_set.drop(['Sex'], axis=1), test\_set.drop(['Sex'], axis=1)

train\_set, validation\_set, test\_set = train\_set.drop(['Pclass'], axis=1), validation\_set.drop(['Pclass'], axis=1), test\_set.drop(['Pclass'], axis=1)

train\_set, validation\_set, test\_set = train\_set.drop(['Embarked'], axis=1), validation\_set.drop(['Embarked'], axis=1), test\_set.drop(['Embarked'], axis=1)

train\_set, validation\_set, test\_set = train\_set.drop(['Name'], axis=1), validation\_set.drop(['Name'], axis=1), test\_set.drop(['Name'], axis=1)

train\_set, validation\_set, test\_set = train\_set.drop(['Title'], axis=1), validation\_set.drop(['Title'], axis=1), test\_set.drop(['Title'], axis=1)

train\_set, validation\_set, test\_set = train\_set.drop(['Deck'], axis=1), validation\_set.drop(['Deck'], axis=1), test\_set.drop(['Deck'], axis=1)

train\_set, validation\_set, test\_set = train\_set.drop(['Fare'], axis=1), validation\_set.drop(['Fare'], axis=1), test\_set.drop(['Fare'], axis=1)

train\_set, validation\_set, test\_set = train\_set.drop(['SibSp'], axis=1), validation\_set.drop(['SibSp'], axis=1), test\_set.drop(['SibSp'], axis=1)

train\_set, validation\_set, test\_set = train\_set.drop(['Parch'], axis=1), validation\_set.drop(['Parch'], axis=1), test\_set.drop(['Parch'], axis=1)

train\_set, validation\_set, test\_set =