

	cap- shape_b	cap- shape_c	cap- shape_e	cap- shape_s	cap- shape_x	cap- surface_f	cap- surface_g	cap- surface_s	cap- surface_y	...	population_s	population_v	population_p
0	0	0	0	0	0	1	0	0	1	0	...	1	0
1	0	0	0	0	0	1	0	0	1	0	...	0	0
2	1	0	0	0	0	0	0	0	1	0	...	0	0
3	0	0	0	0	0	1	0	0	0	1	...	1	0
4	0	0	0	0	0	1	0	0	1	0	...	0	0

5 rows × 17 columns

```
In [39]: #Encoding the target
from sklearn.preprocessing import LabelEncoder
encoder = LabelEncoder()
y = encoder.fit_transform(y)
print(y)

[1 0 0 ... 0 1 0]
```

```
In [40]: #We will start by finding the test set accuracy when training a
#Random classifier using 500 training samples:
#Splitting dataset to training and test and printing the sizes:
from sklearn.model_selection import train_test_split
X_train_val, X_test, y_train_val, y_test = train_test_split(X, y,
test_size=0.3845,
random_state=1)
print("X_train_val: {}, X_test: {}".format(X_train_val.shape, X_test.shape))
print("Y_train_val: {}, Y_test: {}".format(y_train_val.shape, y_test.shape))

X_train_val: (5000, 117), X_test: (3124, 117)
Y_train_val: (5000,), Y_test: (3124,)
```

```
In [41]: #Normalizing the features:
from sklearn.preprocessing import StandardScaler
normalizer = StandardScaler()
X_train_val = normalizer.fit_transform(X_train_val)
X_test = normalizer.transform(X_test)
```

```
In [42]: #We would like to use 10-fold cross validation:
folds = 10

#Getting the number of samples in the training and validation set
num_train_val = X_train_val.shape[0]

# shuffle the index of samples in the train_val set
index_of_samples = np.arange(num_train_val)
shuffle(index_of_samples)

# split the index of the train_val set into 10 folds
index_of_folds = index_of_samples.reshape(folds, -1)
print(index_of_folds)

best_acc = 0.0

#10-fold cross-validation:
sum_acc = 0.0

for fold in range(folds):
    index_of_folds_temp = index_of_folds.copy()

    valid_index = index_of_folds_temp[fold,:].reshape(-1) #get the index of the validation set
    train_index = np.delete(index_of_folds_temp, fold, 0).reshape(-1) #get the index of the training set

    # training set
    X_train = X_train_val[train_index]
    y_train = y_train_val[train_index]

    # validation set
    X_valid = X_train_val[valid_index]
    y_valid = y_train_val[valid_index]

    #Converting back into Dataframes to be able to be used by my algorithm
    X_train_df = pd.DataFrame(X_train,
y_train_df = pd.DataFrame(y_train, columns = ['class'])
X_valid_df = pd.DataFrame(X_valid)
Y_valid_df = pd.DataFrame(y_valid, columns = ['class'])

    appended = pd.concat([X_train_df,y_train_df], axis = 1)
    appendedValid = pd.concat([X_train_df,y_valid_df], axis = 1)

    treeshroom = Id3(appended,'class')

    #Determine the accuracy score
    accuracy = evaluate(treeshroom, appendedValid, 'class') #evaluating the test dataset
    print("Model Accuracy for Fold ", (fold + 1), ": ", accuracy)

    sum_acc += accuracy

averageAcc = sum_acc/folds
print("Average Accuracy: (0:0.4f)".format(averageAcc))

[(1560 3203 1584 ... 3631 4474 3178)
 (3375 570 3908 ... 1415 2342 692)
 [ 65 3217 2797 ... 3122 4579 3787]
 /
 (4312 1351 4045 ... 2580 964 787)
 (115 191 4855 ... 3890 3004 4198)
 [ 594 1128 2013 ... 1754 416 4051]]
Model Accuracy for Fold 1 : 1.0
Model Accuracy for Fold 2 : 1.0
Model Accuracy for Fold 4 : 1.0
Model Accuracy for Fold 5 : 1.0
Model Accuracy for Fold 6 : 1.0
Model Accuracy for Fold 7 : 1.0
Model Accuracy for Fold 9 : 0.998
Model Accuracy for Fold 10 : 1.0
Average Accuracy: 0.998
```

Thus, when using our model on 5000 training samples, we achieve a test score accuracy of .9998. Since the test score accuracy for 500 training samples was .9940, we see that the larger training dataset allowed the model to perform slightly better, as expected.

d) Report accuracy on training and test data when using a decision tree that has 10 vs 30 leaves.

I am a bit confused how we are supposed to complete part d. It was not specified in part a that our model should account for/allow us to input a certain number of leaves. Thus, the model I created cannot do that. I am not going to rewrite my entire model to account for this--which I already had to trouble doing since I never taught how to do this--so I need to use the sklearn's Decision Tree Classifier. If our model should be able to account for this, it should have been specified in part a, not in part d when a student has already spent the time building the model a certain way. I do not feel that this is really my fault, so I hope you do not take off points for me using the Decision Tree Classifier.

```
In [43]: #First, I will find the accuracy on the training and test data when using
# a decision tree that has 10 leaves:
#Separating Features and Target:
X = dataset.drop(['class'], axis = 1)
y = dataset['class']
```

```
In [44]: #Since the values in the dataset are all categorical, we can encode
#using pandas dummy variable and encode y using LabelEncoder.
X = pd.get_dummies(X)
X.head()
```

	cap- shape_b	cap- shape_c	cap- shape_e	cap- shape_k	cap- shape_s	cap- shape_x	cap- surface_f	cap- surface_g	cap- surface_s	cap- surface_y	...	population_s	population_v	population_p
0	0	0	0	0	0	0	1	0	0	1	0	...	1	0
1	0	0	0	0	0	0	1	0	0	1	0	...	0	0
2	1	0	0	0	0	0	0	0	1	0	0	0
3	0	0	0	0	0	0	1	0	0	0	1	...	1	0
4	0	0	0	0	0	0	1	0	0	1	0	...	0	0

5 rows × 17 columns

```
In [45]: #Encoding the target
from sklearn.preprocessing import LabelEncoder
encoder = LabelEncoder()
y = encoder.fit_transform(y)
print(y)

[1 0 0 ... 0 1 0]
```

```
In [46]: #Splitting dataset to training and test and printing the sizes:
from sklearn.model_selection import train_test_split
X_train_val, X_test, y_train_val, y_test = train_test_split(X, y,
test_size=0.3845,
random_state=1)
print("X_train_val: {}, X_test: {}".format(X_train_val.shape, X_test.shape))
print("Y_train_val: {}, Y_test: {}".format(y_train_val.shape, y_test.shape))

X_train_val: (5000, 117), X_test: (3124, 117)
Y_train_val: (5000,), Y_test: (3124,)
```

```
In [47]: #Normalizing the features:
from sklearn.preprocessing import StandardScaler
normalizer = StandardScaler()
X_train_val = normalizer.fit_transform(X_train_val)
X_test = normalizer.transform(X_test)
```

```
In [48]: #We would like to use 10-fold cross validation:
folds = 10

#Getting the number of samples in the training and validation set
num_train_val = X_train_val.shape[0]

# shuffle the index of samples in the train_val set
index_of_samples = np.arange(num_train_val)
shuffle(index_of_samples)

# split the index of the train_val set into 10 folds
index_of_folds = index_of_samples.reshape(folds, -1)
print(index_of_folds)

best_acc = 0.0

#10-fold cross-validation:
sum_acc = 0.0
sum_train_acc = 0.0

for fold in range(folds):
    index_of_folds_temp = index_of_folds.copy()

    valid_index = index_of_folds_temp[fold,:].reshape(-1) #get the index of the validation set
    train_index = np.delete(index_of_folds_temp, fold, 0).reshape(-1) #get the index of the training set

    # training set
    X_train = X_train_val[train_index]
    y_train = y_train_val[train_index]

    # validation set
    X_valid = X_train_val[valid_index]
    y_valid = y_train_val[valid_index]

    #WE NOTE THAT WE HAVE A MAX LEAF NODES OF 10 FOR 10 LEAVES
    clf_entropy = DecisionTreeClassifier(criterion='entropy', max_leaf_nodes=10, random_state=0)

    #Finally, fitting the model:
    clf_entropy.fit(X_train, y_train)

    y_pred_entropy = clf_entropy.predict(X_valid)
    y_pred_train_entropy = clf_entropy.predict(X_train) #for model comparison

    #Determine the accuracy score
    acc = accuracy_score(y_valid, y_pred_entropy)
    accTrain = accuracy_score(y_train, y_pred_train_entropy)

    print("Model accuracy score with criterion Entropy index: (0:0.4f) for fold".format(acc), (fold + 1))
    #Accuracy score for training set
    print("Training-set accuracy score with criterion Entropy index: (0:0.4f)".format(accTrain))

    sum_acc += acc
    sum_train_acc += accTrain

averageAcc = sum_acc/folds
averageAccTrain = sum_train_acc/folds

print("Average Accuracy Test: (0:0.4f)".format(averageAcc))
print("Average Accuracy Train: (0:0.4f)".format(averageAccTrain))

[(12400 1499 1116 ... 3837 3620 4222)
 (4232 2842 1865 ... 1968 2945 796)
 [ 958 5362 4487 ... 1206 3987 968]
 /
 [2245 2405 1944 ... 3769 4577 3335]
 (4926 2457 2722 ... 2320 264 1571)
 (1283 1424 4276 ... 1754 2748 4804)]
Model accuracy score with criterion Entropy index: 1.0000 for fold 1
Training-set accuracy score with criterion Entropy index: 0.9993
Model accuracy score with criterion Entropy index: 1.0000 for fold 2
Training-set accuracy score with criterion Entropy index: 0.9993
Model accuracy score with criterion Entropy index: 1.0000 for fold 3
Training-set accuracy score with criterion Entropy index: 0.9993
Model accuracy score with criterion Entropy index: 1.0000 for fold 4
Training-set accuracy score with criterion Entropy index: 0.9980 for fold 5
Model accuracy score with criterion Entropy index: 1.0000 for fold 6
Training-set accuracy score with criterion Entropy index: 0.9993
Model accuracy score with criterion Entropy index: 0.9993
Model accuracy score with criterion Entropy index: 1.0000 for fold 7
Training-set accuracy score with criterion Entropy index: 0.9993
Model accuracy score with criterion Entropy index: 1.0000 for fold 8
Training-set accuracy score with criterion Entropy index: 0.9980 for fold 9
Model accuracy score with criterion Entropy index: 0.9998
Training-set accuracy score with criterion Entropy index: 1.0000
Model accuracy score with criterion Entropy index: 0.9998
Average Accuracy Test: 0.9994
Average Accuracy Train: 0.9994
```

Thus, the Average accuracy on the training data is .9994 and on the testing data is .9994 when using 10 leaves.

We will now do the same procedure but with 30 leaves:

```
In [49]: #Separating Features and Target:
X = dataset.drop(['class'], axis = 1)
y = dataset['class']
```

```
In [50]: #Since the values in the dataset are all categorical, we can encode
#using pandas dummy variable and encode y using LabelEncoder.
X = pd.get_dummies(X)
X.head()
```

	cap- shape_b	cap- shape_c	cap- shape_e	cap- shape_k	cap- shape_s	cap- shape_x	cap- surface_f	cap- surface_g	cap- surface_s	cap- surface_y	...	population_s	population_v	population_p
0	0	0	0	0	0	0	1	0	0	1	0	...	1	0
1	0	0	0	0	0	0	1	0	0	1	0	...	0	0
2	1	0	0	0	0	0	0	0	1	0	0	0
3	0	0	0	0	0	0	1	0	0	0	1	...	1	0
4	0	0	0	0	0	0	1	0	0	1	0	...	0	0

5 rows × 17 columns

```
In [51]: #Encoding the target
from sklearn.preprocessing import LabelEncoder
encoder = LabelEncoder()
y = encoder.fit_transform(y)
print(y)

[1 0 0 ... 0 1 0]
```

```
In [52]: #Splitting dataset to training and test and printing the sizes:
from sklearn.model_selection import train_test_split
X_train_val, X_test, y_train_val, y_test = train_test_split(X, y,
test_size=0.3845,
random_state=1)
print("X_train_val: {}, X_test: {}".format(X_train_val.shape, X_test.shape))
print("Y_train_val: {}, Y_test: {}".format(y_train_val.shape, y_test.shape))

X_train_val: (5000, 117), X_test: (3124, 117)
Y_train_val: (5000,), Y_test: (3124,)
```

```
In [53]: #Normalizing the features:
from sklearn.preprocessing import StandardScaler
normalizer = StandardScaler()
X_train_val = normalizer.fit_transform(X_train_val)
X_test = normalizer.transform(X_test)
```

```
In [54]: #We would like to use 10-fold cross validation:
folds = 10

#Getting the number of samples in the training and validation set
num_train_val = X_train_val.shape[0]

# shuffle the index of samples in the train_val set
index_of_samples = np.arange(num_train_val)
shuffle(index_of_samples)

# split the index of the train_val set into 10 folds
index_of_folds = index_of_samples.reshape(folds, -1)
print(index_of_folds)

best_acc = 0.0

#10-fold cross-validation:
sum_acc = 0.0
sum_train_acc = 0.0

for fold in range(folds):
    index_of_folds_temp = index_of_folds.copy()

    valid_index = index_of_folds_temp[fold,:].reshape(-1) #get the index of the validation set
    train_index = np.delete(index_of_folds_temp, fold, 0).reshape(-1) #get the index of the training set

    # training set
    X_train = X_train_val[train_index]
    y_train = y_train_val[train_index]

    # validation set
    X_valid = X_train_val[valid_index]
    y_valid = y_train_val[valid_index]

    #WE NOTE THAT WE HAVE A MAX LEAF NODES OF 30 FOR 30 LEAVES
    clf_entropy = DecisionTreeClassifier(criterion='entropy', max_leaf_nodes=30, random_state=0)

    #Finally, fitting the model:
    clf_entropy.fit(X_train, y_train)

    y_pred_entropy = clf_entropy.predict(X_valid)
    y_pred_train_entropy = clf_entropy.predict(X_train) #for model comparison

    #Determine the accuracy score
    acc = accuracy_score(y_valid, y_pred_entropy)
    accTrain = accuracy_score(y_train, y_pred_train_entropy)

    print("Model accuracy score with criterion Entropy index: (0:0.4f) for fold".format(acc), (fold + 1))
    #Accuracy score for training set
    print("Training-set accuracy score with criterion Entropy index: (0:0.4f)".format(accTrain))

    sum_acc += acc
    sum_train_acc += accTrain

averageAcc = sum_acc/folds
averageAccTrain = sum_train_acc/folds

print("Average Accuracy Test: (0:0.4f)".format(averageAcc))
print("Average Accuracy Train: (0:0.4f)".format(averageAccTrain))

[(12754 657 2527 ... 3290 109 1583)
 (116 1243 4340 ... 1982 3437 3480)
 [ 932 711 2193 ... 3568 1460 3166]
 /
 [295 3917 793 ... 1740 1601 1971]
 (3886 3207 3710 ... 1422 2446 4010)
 (4759 4070 2613 ... 248 13 2988)]
Model accuracy score with criterion Entropy index: 0.9980 for fold 1
Training-set accuracy score with criterion Entropy index: 1.0000
Model accuracy score with criterion Entropy index: 1.0000 for fold 2
Training-set accuracy score with criterion Entropy index: 1.0000
Model accuracy score with criterion Entropy index: 1.0000 for fold 3
Training-set accuracy score with criterion Entropy index: 1.0000 for fold 4
Training-set accuracy score with criterion Entropy index: 1.0000 for fold 5
Model accuracy score with criterion Entropy index: 1.0000 for fold 6
Training-set accuracy score with criterion Entropy index: 1.0000 for fold 7
Training-set accuracy score with criterion Entropy index: 1.0000 for fold 8
Model accuracy score with criterion Entropy index: 1.0000 for fold 9
Training-set accuracy score with criterion Entropy index: 1.0000 for fold 10
Training-set accuracy score with criterion Entropy index: 1.0000
Model accuracy score with criterion Entropy index: 0.9998
Average Accuracy Test: 0.9998
Average Accuracy Train: 1.0000
```

With 30 leaves, we see the average accuracy for the testing data is .998 and 1.0 for the training data. These scores are only slightly better than with 10 leaves, suggesting that having a max depth of 10 (having 10 leaves) was sufficient enough. Still, a max_leaf_nodes=30 would be able to improve the accuracy ever so slightly, so we most likely would want just over 10 leaves.

e) Report attribute tests used in the decision tree with 10 leaves.

I am a bit confused what this question is asking. I have looked up the term "attribute tests" and have found no results. This terminology appears to be unique to the person who wrote this assignment, as I cannot find anything on it. The closest thing I could find online is "test attributes"; but this does not make sense in this context. If I were to guess, I would say that this question is asking us to discuss how the accuracy score was determined.

The accuracy score outputted was determined using sklearn's accuracy_score function, which takes as parameters the ground truth (correct) labels, and the predicted labels. By the numbers that match exactly, then there will be an accuracy score of 1.0. In the decision tree model written above, the predicted class function is run on the training data and the validation data to meet the desires of the previous question. The scores get printed for each fold, with their average taken at the end. The scores use criterion Entropy index, as specified when making the Decision Tree Classifier.

f) For a decision tree with 10 leaves report the number of positive and negative training examples at each internal node and at leaves.

To report the number of positive and negative training examples at each internal node and at leaves, one would have to examine these numbers 10 times since we put the model on 10 folds. Since this would be a lot of information to examine, I will report these numbers without doing the 10-fold cross validation.

```
In [55]: #Separating Features and Target:
X = dataset.drop(['class'], axis = 1)
y = dataset['class']
```

```
In [56]: #Since the values in the dataset are all categorical, we can encode
#using pandas dummy variable and encode y using LabelEncoder.
X = pd.get_dummies(X)
X.head()
```

	cap- shape_b	cap- shape_c	cap- shape_e	cap- shape_k	cap- shape_s	cap- shape_x	cap- surface_f	cap- surface_g	cap- surface_s	cap- surface_y	...	population_s	population_v	population_p
0	0	0	0	0	0	0	1	0	0	1	0	...	1	0
1	0	0	0	0	0	0	1	0	0	1	0	...	0	0
2	1	0	0	0	0	0	0	0	1	0	0	0
3	0	0	0	0	0	0	1	0	0	0	1	...	1	0
4	0	0	0	0	0	0	1	0	0	1	0	...	0	0

5 rows × 17 columns

```
In [57]: #Encoding the target
from sklearn.preprocessing import LabelEncoder
encoder = LabelEncoder()
y = encoder.fit_transform(y)
print(y)

[1 0 0 ... 0 1 0]
```

```
In [58]: #Splitting dataset to training and test and printing the sizes:
from sklearn.model_selection import train_test_split
X_train, X_test, y_train, y_test = train_test_split(X, y,
test_size=0.3845,
random_state=1)
print("X_train_val: {}, X_test: {}".format(X_train.shape, X_test.shape))
print("Y_train_val: {}, Y_test: {}".format(y_train.shape, y_test.shape))

X_train_val: (5000, 117), X_test: (3124, 117)
Y_train_val: (5000,), Y_test: (3124,)
```

```
In [59]: #WE NOTE THAT WE HAVE A MAX LEAF NODES OF 10 FOR 10 LEAVES
clf_entropy = DecisionTreeClassifier(criterion='entropy', max_leaf_nodes=10, random_state=0)

#Finally, fitting the model:
clf_entropy.fit(X_train, y_train)

y_pred_entropy = clf_entropy.predict(X_test)
y_pred_train_entropy = clf_entropy.predict(X_train) #for model comparison

#Determine the accuracy score
acc = accuracy_score(y_test, y_pred_entropy)
accTrain = accuracy_score(y_train, y_pred_train_entropy)

print("Model accuracy score with criterion Entropy index: (0:0.4f) for fold".format(acc))
#Accuracy score for training set
print("Training-set accuracy score with criterion Entropy index: (0:0.4f)".format(accTrain))

Model accuracy score with criterion Entropy index: 0.9984
Training-set accuracy score with criterion Entropy index: 0.9994
```

```
In [60]: #Importing matplotlib.pyplot as plt
from sklearn import tree
```

```
In [61]: plt.figure(figsize = (12,8))
tree.plot_tree(clf_entropy)
```

```
Out[61]: [Text(386.3076923076923, 398.64, 'X[27] <= 0.5\nentropy = 0.999\nsamples = 5000\nvalue = [2600, 2400]'),
Text(257.5384615384615, 326.15999999999997, 'X[20] <= 0.5\nentropy = 0.663\nsamples = 2810\nvalue = [485, 233 0]'),
Text(103.01538461538461, 253.67999999999998, 'X[53] <= 0.5\nentropy = 0.978\nsamples = 825\nvalue = [485, 34 0]'),
Text(154.52307692307693, 181.2, 'X[63] <= 0.5\nentropy = 0.936\nsamples = 523\nvalue = [183, 340]'),
Text(103.01538461538461, 108.71999999999997, 'X[110] <= 0.5\nentropy = 0.645\nsamples = 407\nvalue = [67, 34 0]'),
Text(11.50769230769231, 36.239999999999995, 'entropy = 0.0\nsamples = 340\nvalue = [0, 340]'),
Text(154.52307692307693, 36.239999999999995, 'entropy = 0.0\nsamples = 67\nvalue = [67, 0]'),
Text(106.03076923076924, 108.71999999999997, 'entropy = 0.0\nsamples = 116\nvalue = [116, 0]'),
Text(257.5384615384615, 181.2, 'entropy = 0.0\nsamples = 302\nvalue = [1985, 0]'),
Text(309.04615384615386, 253.67999999999998, 'entropy = 0.894\nsamples = 1993\nvalue = [0, 1993]'),
Text(15.076923076923073, 108.71999999999997, 'entropy = 0.0\nsamples = 20\nvalue = [0, 20]'),
Text(106.03076923076923, 108.71999999999997, 'entropy = 0.5\nentropy = 0.198\nsamples = 2182\nvalue = [2115, 6 7]'),
Text(63.5692307692308, 253.67999999999998, 'X[63] <= 0.5\nentropy = 0.086\nsamples = 2182\nvalue = [2115, 2 3]'),
Text(160.55384615384617, 181.2, 'X[35] <= 0.5\nentropy = 0.016\nsamples = 2109\nvalue = [2106, 3 3]'),
Text(109.04615384615385, 108.71999999999997, 'entropy = 0.164\nsamples = 124\nvalue = [121, 3]'),
Text(412.0615384615385, 108.71999999999997, 'entropy = 0.0\nsamples = 1985\nvalue = [1985, 0]'),
Text(566.5846153846154, 181.2, 'X[89] <= 0.5\nentropy = 0.894\nsamples = 29\nvalue = [9, 20]'),
Text(15.076923076923073, 108.71999999999997, 'entropy = 0.0\nsamples = 20\nvalue = [0, 20]'),
Text(618.0923076923077, 108.71999999999997, 'entropy = 0.0\nsamples = 9\nvalue = [9, 0]'),
Text(566.5846153846154, 253.67999999999998, 'entropy = 0.0\nsamples = 44\nvalue = [0, 44]')]
```

We now can see the number of positive and negative training examples at each internal node and at leaves for a decision tree with 10 leaves. We note that the column names are not the same as the original dataset due to the label encoding and dummy encoder we decided to use on the dataset. Nonetheless, we see that the tree has 10 leaves (as expected) with 6 levels and the number of positive and negative training examples reported in the 'value' = [] section of each node. For the most part, values are reported as either positive or negative (indicated in the value brackets by one of the values being 0 and the other being the total of the rest of the values given for the parent node), which is emphasized by our model's high accuracy score. Since there are 10 leaf nodes and 9 internal nodes, I will not explicitly write out the positive and negative training examples for each one (especially since we can see those values easily in the tree pictured above). However, for an example, we see that the very first internal node, the root, has a split of 2600 and 2400.

Problem 2:

Show that the entropy of a node in a decision tree never increases after splitting it into smaller successor nodes.

The professor's hypothesis that we do not need to complete this question if we are undergraduate students. Even though I am an undergraduate student, I will try to complete this proof.

Due to the large number of mathematical symbols required for this proof, I have written it in latex. However, I do not know how to paste latex into a Jupyter Notebook, so I am putting pictures of my work/steps below:

Let $Y = \{y_1, y_2, \dots, y_c\}$, where c represents the number of classes in attribute Y (the attribute we are hoping to predict). Also, let $X = \{x_1, x_2, \dots, x_k\}$ where k is the number of attribute values, a certain attribute X has. In the very beginning, when we are at the root of the tree, the entropy is given by the following equation:

$$E(Y) = -\sum_{j=1}^c P(y_j) \log_2 P(y_j) = \sum_{j=1}^c \sum_{i=1}^k P(x_i, y_j) \log_2 P(y_j),$$

which uses the following idea from the law of total probability: $P(y_j) = \sum_{i=1}^k P(x_i, y_j)$. After we split an attribute X , the entropy of each child node $X = x_i$ is

$$E(Y \mid x_i) = -\sum_{j=1}^c P(y_j \mid x_i) \log_2 P(y_j \mid x_i),$$

where the proportion of the sample where $X = x_i$ belong to class y_j is represented by $P(y_j \mid x_i)$. Then, after we split X , the weighted entropy of the subsequent, child nodes is as follows:

$$\begin{aligned} E(Y \mid X) &= \sum_{i=1}^k P(x_i) E(Y \mid x_i) \\ &= -\sum_{i=1}^k \sum_{j=1}^c P(x_i) P(y_j \mid x_i) \log_2 P(y_j \mid x_i) \\ &= -\sum_{i=1}^k \sum_{j=1}^c P(y_j, x_i) \log_2 P(y_j \mid x_i), \end{aligned}$$

from the law of total probability. To prove that the entropy never increases after splitting a node, we must show that $E(Y \mid X) \leq E(Y \mid X)$, which we can do by showing that $E(Y \mid X) - E(Y \mid X) \leq 0$.

$$\begin{aligned} E(Y \mid X) - E(Y) &= -\sum_{i=1}^k \sum_{j=1}^c P(y_j, x_i) \log_2 P(y_j \mid x_i) + \sum_{j=1}^c P(y_j) \log_2 P(y_j \mid x_i) \\ &= \sum_{i=1}^k \sum_{j=1}^c P(y_j, x_i) \log_2 P(y_j \mid x_i) - P(y_j, x_i) \log_2 P(y_j) \\ &= \sum_{i=1}^k \sum_{j=1}^c P(y_j, x_i) \log_2 \frac{P(y_j)}{P(y_j \mid x_i)} \\ &= \sum_{i=1}^k \sum_{j=1}^c P(y_j, x_i) \log_2 \frac{P(y_j) P(x_i)}{P(y_j \mid x_i)} \\ &\leq \log_2 \left(\sum_{i=1}^k \sum_{j=1}^c P(y_j, x_i) \frac{P(y_j) P(x_i)}{P(y_j \mid x_i)} \right) \text{ by Jensen's Inequality} \\ &= \log_2 \left[\sum_{i=1}^k P(x_i) \sum_{j=1}^c P(y_j) \right] \\ &= \log_2(1) \\ &= 0 \end{aligned}$$