

## 0816170 Homework 3

### Part 1

#### 1. Gini and Entropy

```
[26] 1 print("Gini of data is ", gini(data))
```

```
Gini of data is 0.4628099173553719
```

```
[27] 1 print("Entropy of data is ", entropy(data))
```

```
Entropy of data is 0.9456603046006402
```

#### 2. Decision Tree

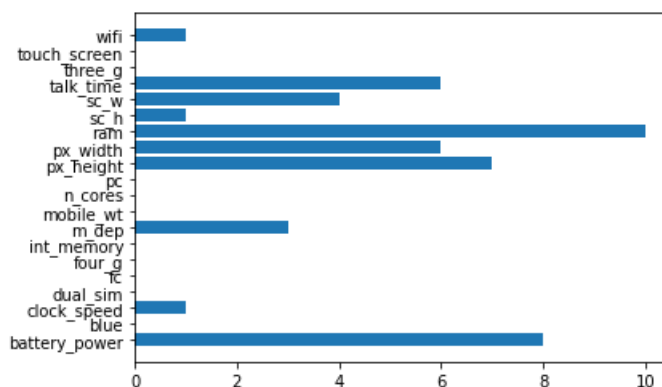
##### 2.1 Max\_depth = 3 and Max\_depth = 10

```
Decision Tree  
max_depth = 3, accuracy score: 0.92  
max_depth = 10, accuracy score: 0.93
```

##### 2.2 Criterion = 'gini' and Criterion = 'entropy'

```
Decision Tree  
Criterion= 'gini', accuracy score: 0.92  
Criterion= 'entropy', accuracy score: 0.9333333333333333
```

#### 3. Feature importance



## 4. AdaBoost

### 4.1

```
AdaBoost
n_estimators=10, accuracy score: 0.9366666666666666
n_estimators=100, accuracy score: 0.9733333333333334
```

## 5. Random Forest

### 5.1

```
Random Forest
n_estimators=10, accuracy score: 0.93
n_estimators=100, accuracy score: 0.96
```

### 5.2

```
Random Forest
random features, accuracy score: 0.9333333333333333
all features, accuracy score: 0.9666666666666667
```

## 6. My\_model

使用 adaboost · n\_estimators = 150 (經實驗發現最高可以到達 0.98 的準確率)

將 train\_df 跟 val\_df 一起送進去訓練

得到最終的 my\_model

大約在 google colab 上跑 5 分鐘

```
[63] 1 from sklearn.metrics import accuracy_score
      2
      3 def train_your_model(data):
      4     ## Define your model and training
      5     x_train = data.drop(labels=["price_range"], axis="columns")
      6     x_train = x_train.values
      7     y_train = data["price_range"].values
      8
      9     ada = 150
     10     adatest = AdaBoost(ada)
     11     adatest.fit(x_train, y_train)
     12     predtest = adatest.predict(x_val)
     13     return adatest

[64] 1 trainval_df = train_df.append(val_df)
      2 my_model = train_your_model(trainval_df)

1 test_df = pd.read_csv('x_test.csv')
2 x_test = test_df.values
3 y_pred = adal00.predict(x_test)

1 assert y_pred.shape == (500, )
```

## Final result

☞ \*\*\* We will check your result for Question 3 manually \*\*\* (5 points)  
\*\*\* We will check your result for Question 6 manually \*\*\* (20 points)  
Approximate score range: 45.0 ~ 70.0  
\*\*\* This score is only for reference \*\*\*

## Part 2

1. ① 因為 decision tree 會想盡辦法把 train data 中的兩種 class 分開來, 因此只要還有 branch 資料不純, decision tree 會為了把不純的地方挑出來而選用一些不太有代表性的 feature 和 threshold, 如此就可能導致 overfit

② 可能達到 100%, 若 training set 的 X 資料皆獨一無二, decision tree 總會找到 feature 和 threshold 把不同 class 的 data 隔開 (code 實作中也有發生)

③

(1) pre-pruning: 例如調整 max-depth 控制 tree 的深度  
讓 tree 不會過度 split

(2) post-pruning: 先生成一棵完整的 tree, 再移除部分 branch, 讓 tree 不要那麼複雜

(3) random forest: 應用 bootstrap sampling 和 data aggregation  
使 tree 減少 overfitting 的機率

2

a True, 根據 update equation 知, 若分類錯誤, 該 data 的 weight 變大

b True, 在訓練過程中, weak classifier 被迫嘗試分類更困難的 examples 若某 example 一直被分類錯, 該 example 的 weight 會增加, 而  $\epsilon_t$  of the  $t^{\text{th}}$  weak classifier 也因此傾向增加

c false, 在部分情況中, 若 training set 中有 data 無法被我們使用的 weak classifier 分割, 無論迭代多少次都無法達到 zero training error

3

① misclassification rates

$$\text{model A : } \frac{200+0}{800} = \frac{1}{4}$$

$$\text{model B : } \frac{100+100}{800} = \frac{1}{4}$$

$\Rightarrow$  their misclassification rates are equal

② model A

$(200, 400)$

$$\text{gini} = 1 - \left(\frac{1}{3}\right)^2 - \left(\frac{2}{3}\right)^2 = \frac{4}{9}$$

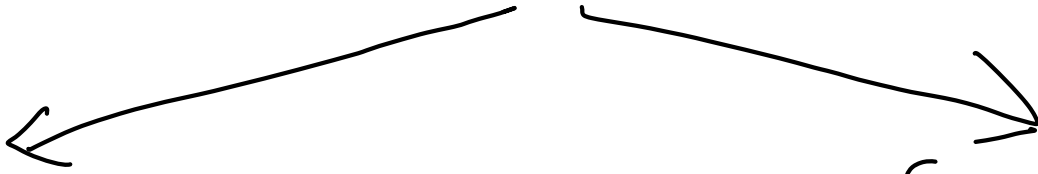
$$\text{entropy} = -\frac{1}{3} \log_2 \frac{1}{3} - \frac{2}{3} \log_2 \frac{2}{3} \approx 0.918$$

$(200, 0)$

$$\text{Gini} = 1 - 1^2 = 0$$

$$\text{entropy} = -1 \log_2 1 = 0$$

model B



$(300, 100)$

$$\text{gini} = 1 - \left(\frac{3}{4}\right)^2 - \left(\frac{1}{4}\right)^2 = \frac{3}{8}$$

$$\text{entropy} = -\frac{3}{4} \log_2 \frac{3}{4} - \frac{1}{4} \log_2 \frac{1}{4} \doteq 0.811$$

$(100, 300)$

$$\text{gini} = 1 - \left(\frac{1}{4}\right)^2 - \left(\frac{3}{4}\right)^2 = \frac{3}{8}$$

$$\text{entropy} = -\frac{1}{4} \log_2 \frac{1}{4} - \frac{3}{4} \log_2 \frac{3}{4} \doteq 0.811$$