# ELG5255[EG]: Applied Machine Learning

# Assignment 4

**Decision Tree and Ensemble Methods**

Group: 25

The report will go as follows: #) for steps to perform the intended tasks including screenshots of the code we used to tackle this question along with any description needed or relevant figures.

---

# Part #1

**# Numerical Questions**

**Q1)**

We calculate the Gini for each feature to choose the best split. Hence we calculate first the gini of leaves for each feature then the total gini.

In 1$^{st}$ iteration →

G_total_F1 = **0.416**

G_total_F2 = 0.444

G_total_F3 = 0.475

G_total_F4 = **0.416**

As F1 and F4 have same gini coefficients, we choose F4 for simple split.

**Iteration 1:**

- F1: Weather, F2: Temperature, F3: Humidity, F4: Wind
- $Gini\_leaves = 1-(P(Yes))^2 - (P(No))^2$
- Total Gini of node $= \sum weight \ of \ leaf \times Gini\_leaf$
  - $\times$ instance in leaf
  - $\times$ total instance in node

- iterate ① → Calculate Gini for each feature & take the lowest

① F1

| F1 | | |
|---|---|---|
| $\neq$③ clouds | ③ Sunny $\neq$④ | Rainy $\neq$④ |

| Hiking | | Hiking | | Hikin | |
|---|---|---|---|---|---|
| Yes | No | Yes | No | Yes | No |
| 2 | 1 | 1 | 2 | 1 | 3 |

$G = 1-(\frac{2}{3})^2 -(\frac{1}{3})^2 = 0.444$

$G_2 = 1- (\frac{2}{3})^2 -(\frac{1}{3})^2 = 0.44$

$G_3 = 1-(\frac{1}{4})^2-(\frac{3}{4})^2$

$= 0.375$

$Gini = \frac{3}{10} \times 0.444 + \frac{3}{10} \times 0.444 + \frac{4}{10} \times 0.375$

$= \boxed{0.416}$

② F2

| F2 | | |
|---|---|---|
| $\neq$④ Hot | ⑤ Mild | Cool ① |

| Hiking | | Hiking | | Hiking | |
|---|---|---|---|---|---|
| Yes | No | Yes | No | Yes | No |
| 2 | 2 | 3 | 2 | 0 | 1 |

$G_1 = 1-(\frac{1}{2})^2-(\frac{1}{2})^2$

$= 0.5$

$G_2 = 1-(\frac{3}{5})^2 (\frac{2}{5})^2$

$= 0.48$

$G_3 = 0$ (no impurity)

$= 0$

$Gini = \frac{4}{10} \times 0.5 + \frac{5}{10} \times 0.48 + \frac{1}{10} \times 0 = \boxed{0.444}$

③ F3

| F3 | | |
|---|---|---|
| $\neq$⑦ High | | Normal $\neq$③ |

| Hiking | | Hiking | |
|---|---|---|---|
| Yes | No | Yes | No |
| 3 | 4 | 2 | 1 |

$G_1 = 1-(\frac{3}{7})^2+(\frac{4}{7})^2$

$= 0.489$

$G = 1-(\frac{2}{3})^2 -(\frac{1}{3})^2$

$= 0.444$

$Gini = \frac{7}{10} \times 0.489 + \frac{3}{10} \times 0.444 = \boxed{0.475}$

④ F4

| F4 | | |
|---|---|---|
| $\neq$④ weak | | Strong $\neq$⑥ |

| Hiking | | Hiking | |
|---|---|---|---|
| Yes | No | Yes | No |
| 3 | 1 | 2 | 4 |

$G = 1-(\frac{3}{4})^2-(\frac{1}{4})^2 = 0.375$

$G = 1-(\frac{2}{6})^2+(\frac{4}{6})^2$

$= 0.444$

$Gini = \frac{4}{10} \times 0.375 + \frac{6}{10} \times 0.444$

$= \boxed{0.4164}$

We can find there're more impurities in both leaves so both leaves need further split and hence we try gini with the remaining 3 features in both leaves.

# Iteration 2:

The 4 gini values are: $f_1 = f_4 = 0.4164$ ✓

$f_2 = 0.444$    $f_3 = 0.475$

so we choose for $f_1$ or $f_4$ as they're same gini; for simplicity of sketch chose [f4]

. iterate ② → for feature ④; each leaf in it we shall check if it has impurities. Since both leaves have impurities so we shall calculate the gini for the other 3 features in accordance with feature ④ & choose lowest

f4

**f4** — weak / Strong

Hiking: Yes 3 / No (4) — impurity    impurity — Yes 2 / Hiking No 4

**f1**

f1 — Cloudy ② / Sunny ① / Rain ⓪

| | Yes | No |
|---|---|---|
| | 1 | 1 |
| | 1 | 0 |
| | 1 | 0 |

↳ G = 0.5    ↳ G = 0    ↳ G = 0

total $G = 2/4 \times 0.5 + 0 = \boxed{0.25}$

**f1** (right) — Cloudy ① / Sunny ② / Rain ③    ≠6

| | Yes | No |
|---|---|---|
| | 1 | 0 |
| | 1 | 1 |
| | 0 | 3 |

↳ G = 0    ↳ G = 0.5    ↳ G = 0.

$G = 2/6 \times 0.5 = \boxed{0.167}$

**f2**

f2 — hot ③ / mild ① / Cool ⓪

| | Yes | No |
|---|---|---|
| | 2 | 1 |
| | 1 | 0 |
| | 0 | 0 | → G=0

$G = 1 - (1/2)^2 - (1/2)^2 = 0.444$    ↳ G = 0

total $G = 3/4 \times 0.444 + 0 = \boxed{0.333}$

**f2** (right) — hot ⓪ / mild ④ / Cool ①    ≠6

| | Yes | No |
|---|---|---|
| | 0 | 1 |
| | 2 | 2 |
| | 0 | 1 |

↳ G=0    ↳ G = 0.5    ↳ G = 0.

total $G = 4/6 \times 0.5 = \boxed{0.333}$

**f3**

f3 — high ③ / Normal ①    ≠④

| | Yes | No |
|---|---|---|
| | 2 | 1 |
| | 1 | 0 | ↳ G=0

$G = 1 - (2/3)^2 - (1/3)^2 = 0.444$

$G = 0.444 \times 3/4 + 0 = \boxed{0.333}$

**f3** (right) — high ④ / Normal ②    ≠6

| | Yes | No |
|---|---|---|
| | 1 | 3 |
| | 1 | 1 |

$G = 1 - (1/4)^2 - (3/4)^2 = 3/8$    ↳ G = 0.5

total $G = 4/6 \times 3/8 + 2/6 \times 0.5 = \boxed{0.4167}$

In 2$^{nd}$ iteration →

| Wind (F4) | Weak | Strong |
|---|---|---|
| G_total_F1 | **0.25** | **0.167** |
| G_total_F2 | 0.333 | 0.333 |
| G_total_F3 | 0.333 | 0.4167 |

Here we can observe that for weak wind the best subsequent split will be by the weather (F1) and same goes for the strong wind as they have the least total gini.

However, there still few impurities after this split so we need one more iteration for splitting.

**Iteration 3:**

③

→ for . weak wind (f4) — G_f1 = 0.25 ✓  ∴ we'll Split by
         G_f2 = 0.333  (f1)
         G_f3 = 0.333

• Strong wind (f4) — G_f1 = 0.167 ✓  ∴ we'll Split by
         G_f2 = 0.333  (f1)
         G_f3 = 0.4167

f4
   weak             Strong

**Left (weak):** f1 → Cloudy / Sunny / Rain
- Cloudy: Yes 1 No 1
- Sunny: Yes 1 No 0
- Rain: Yes 1 No 0

**Right (Strong):** f1 → cloudy / Sunny / Rain
- cloudy: Yes 1 No 0
- Sunny: Yes 1 No 1
- Rain: Yes 0 No 2

→ As we observe Here ; Cloudy weak-wind weather Carry impurity (Try gain to Split further)

→ Here the Sunny Strong-wind weather hase in purity

• f4 = weak && f1 = Cloudy

• f4 = Strong && f1 = Sunny

f2
② hot   mild ◎   Cool ◎
- hot: Yes 1 No 1
- mild: Yes 0 No 0
- Cool: Yes 0 No 0

$G = 1 - \frac{1}{2}^2 - \frac{1}{2}^2$
= 0.5
total G = ½ * 0.5 + 0 = 0.5 *②

f2
hot   mild   Cool
- hot: Yes 0 No 1 → G=0
- mild: Yes 1 No 0 → G=0
- Cool: Yes 0 No 0
total G = 0 + 0 + 0 = 0

f3
High    Normal
- High: Yes 0 No 1 → G=0
- Normal: Yes 1 No 0 → G=0
G = 0

f3
High    Normal
- High: Yes 0 No 1 → G=0
- Normal: Yes 1 No 0 → G=0
total G = 0 + 0 = 0

* G_f2 = 0.5
  G_f3 = 0 ✓ Split by f3

* G_f2 = G_f3 = 0
→ for Simple Sketch → f3

| Wind (F4) & Weather (F1) | Weak & cloudy | Strong & Sunny |
|---|---|---|
| G_total_F2 | 0.5 | 0 |
| G_total_F3 | 0 | 0 |

For the weak wind and cloudy weather, we can see that F3 (Humidity) yields no impurities.

Also, for the strong wind and sunny weather, we can observe that both features yield no impurities.

We choose F3 for both splits and this is the final tree that we get.

Q2)

$$H(T) = -5/10 \log_2(5/10) - 5/10 \log_2(5/10) = 1$$

$$H(T|Weather) = 3/10\left(-2/3 \log_2(2/3) - 1/3 \log_2(1/3)\right) +$$

$$3/10\left(-2/3 \log_2(2/3) - 1/3 \log_2(1/3)\right) +$$

$$4/10\left(-1/4 \log_2(1/4) - 3/4 \log_2(3/4)\right) \simeq 0.275 + 0.275 + 0.325$$
$$\simeq 0.875$$

$$H(T|Temp) = 4/10\left(-3/4 \log_2(3/4) - 3/4 \log_2(3/4)\right) +$$

$$5/10\left(-3/5 \log_2(3/5) - 2/5 \log_2(2/5)\right) +$$

$$1/10\left(-1/1 \log_2(1) - 0/1 \log_2(0)\right) = 0.4 + 0.485 + 0$$
$$\simeq 0.885$$

$$H(T|Hum) = \frac{7}{10}\left(-\frac{3}{7} \log_2(3/7) - 4/7 \log_2(4/7)\right) +$$

$$\frac{3}{10}\left(-2/3 \log_2(2/3) - 1/3 \log_2(1/3)\right) \simeq 0.690 + 0.275$$
$$\simeq 0.965$$

$$H(T|Wind) = 4/10\left(-3/4 \log_2(3/4) - 1/4 \log_2(1/4)\right) +$$

$$6/10\left(-2/6 \log_2(2/6) - 4/6 \log_2(4/6)\right)$$

$$\simeq 0.325 + 0.551 \simeq 0.876$$

$IG(T|Weather) = 1 - 0.875 = \boxed{0.125}$ » Highest IG
$IG(T|Temp) = 1 - 0.885 = 0.115$
$IG(T|Hum) = 1 - 0.965 = 0.035$
$IG(T|Wind) = 1 - 0.876 = 0.124$

→ cloudy

| Temperature | Humidity | wind | Hiking |
|---|---|---|---|
| Hot | High | weak | No |
| Hot | Normal | weak | Yes |
| Mild | High | strong | Yes |

$$H(T) = -\frac{2}{3} \log_2 \left(\frac{2}{3}\right) - \frac{1}{3} \log_2 \left(\frac{1}{3}\right) \simeq 0.918$$

$$H(T|Temp) = \frac{2}{3} \left(-\frac{1}{2} \log_2 \left(\frac{1}{2}\right) - \frac{1}{2} \log_2 \left(\frac{1}{2}\right)\right) \simeq 0.667$$

$$H(T|Hum) = \frac{2}{3} \left(-\frac{1}{2} \log_2 \left(\frac{1}{2}\right) - \frac{1}{2} \log_2 \left(\frac{1}{2}\right)\right) \simeq 0.667$$

$$H(T|wind) = \frac{2}{3} \left(-\frac{1}{2} \log_2 \left(\frac{1}{2}\right) - \frac{1}{2} \log_2 \left(\frac{1}{2}\right)\right) \simeq 0.667$$

$$IG(T|Temp) = IG(T|Hum) = IG(T|wind)$$
$$\simeq 0.918 - 0.667 \simeq 0.251$$

since all information gain are equal Then I will choose
The first predictor ( Temperature )

→ sunny

| Temperature | Humidity | wind | Hiking |
|-------------|----------|--------|--------|
| HoT | High | weak | Yes |
| HoT | High | STrong | No |
| Mild | Normal | sTrong | Yes |

$$H(T) = -\frac{2}{3} \log_2 \left(\frac{2}{3}\right) - \frac{1}{3} \log_2 \left(\frac{1}{3}\right) \approx 0.918$$

$$H(T|Temp) = \frac{2}{3}\left(-\frac{1}{2} \log_2 \left(\frac{1}{2}\right) - \frac{1}{2} \log_2 \left(\frac{1}{2}\right)\right) \approx 0.667$$

$$H(T|Hum) = \frac{2}{3}\left(-\frac{1}{2} \log_2 \left(\frac{1}{2}\right) - \frac{1}{2} \log_2 \left(\frac{1}{2}\right)\right) \approx 0.667$$
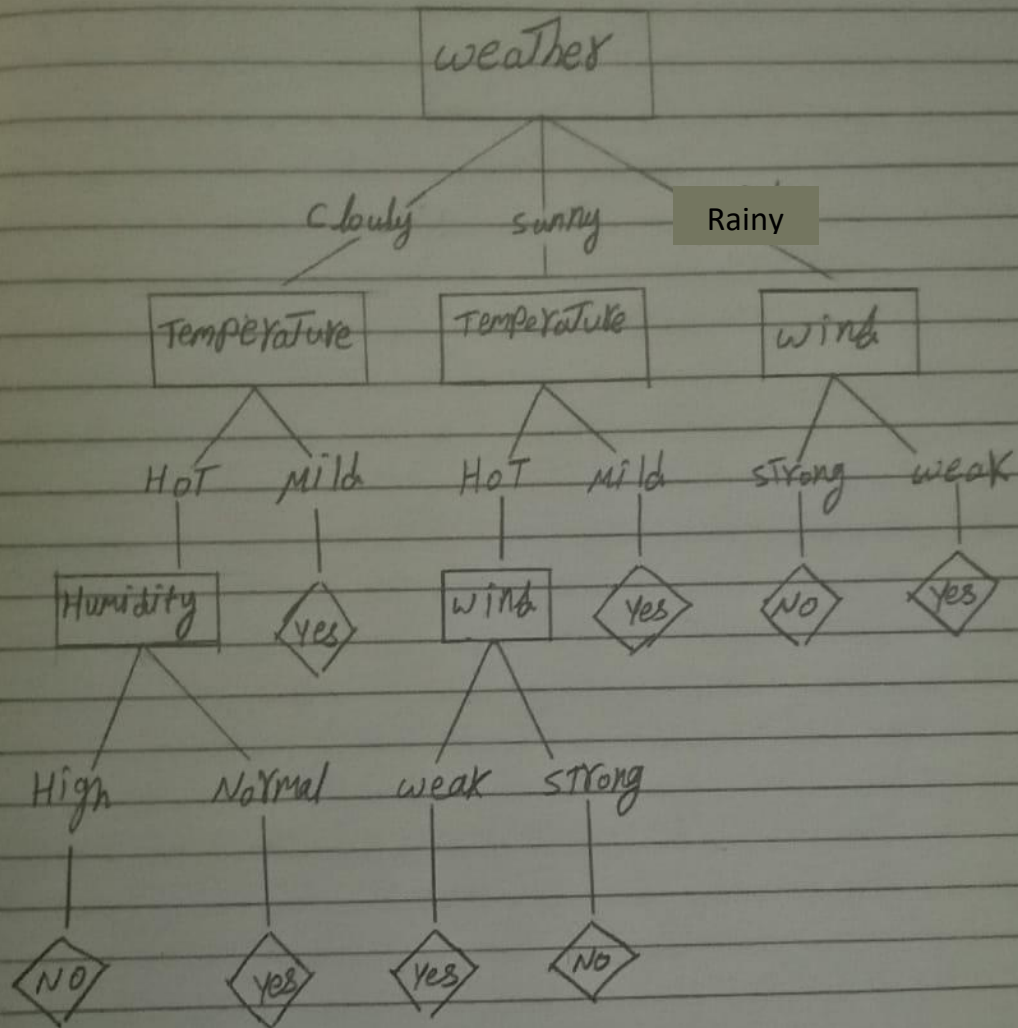
$$H(T|wind) = \frac{2}{3}\left(-\frac{1}{2} \log_2 \left(\frac{1}{2}\right) - \frac{1}{2} \log_2 \left(\frac{1}{2}\right)\right) \approx 0.667$$

$$IG(T|Temp) = H(T|Hum) = H(T|wind) \approx 0.918 - 0.667 \approx 0.251$$

since all IG are equal Then I'll choose The first
   predicTor (TemperaTure)

$$\boxed{4}$$

→ Rainy

| Temperature | Humidity | wind | Hiking | |
|---|---|---|---|---|
| Mild | High | sTrong | No | |
| Mild | High | weak | Yes | |
| cool | Normal | sTrong | No | |

$$H(T) = -\tfrac{1}{3}\log_2(\tfrac{1}{3}) - \tfrac{2}{3}\log_2(\tfrac{2}{3}) \simeq 0.918$$

$$H(T \mid Temp) = \tfrac{2}{3}\left(-\tfrac{1}{2}\log_2(\tfrac{1}{2}) - \tfrac{1}{2}\log_2(\tfrac{1}{2})\right) \simeq 0.667$$

$$H(T \mid Hum) = \tfrac{2}{3}\left(-\tfrac{1}{2}\log_2(\tfrac{1}{2}) - \tfrac{1}{2}\log_2(\tfrac{1}{2})\right) \simeq 0.667$$

$$H(T \mid wind) = \tfrac{2}{3}\left(0 - \tfrac{2}{2}\log_2(\tfrac{2}{2})\right) + \tfrac{1}{3}\left(-\tfrac{1}{1}\log_2(1) - 0\right) = 0$$

$$IG(T \mid Temp) = IG(T \mid Hum) \simeq 0.251$$

$$IG(T \mid wind) = 0.918 - 0 = \boxed{0.918} \Rightarrow highesT\ IG$$

5)

weather

Clouly    sunny    Rainy

Temperature    Temperature    wind

HoT    Mild    HoT    Mild    strong    weak

Humidity    yes    wind    yes    NO    yes

High    Normal    weak    strong

NO    yes    yes    NO
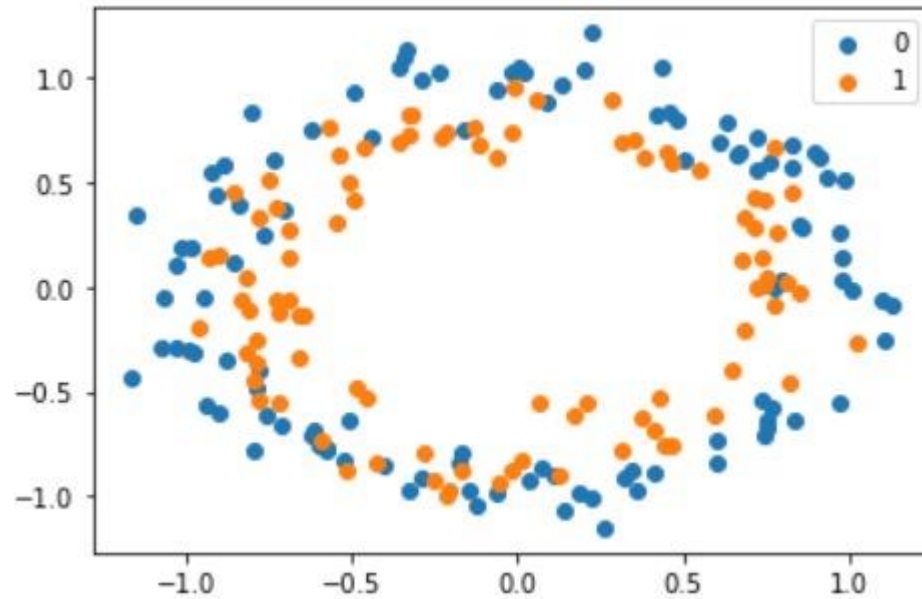
**Q3)** [1]

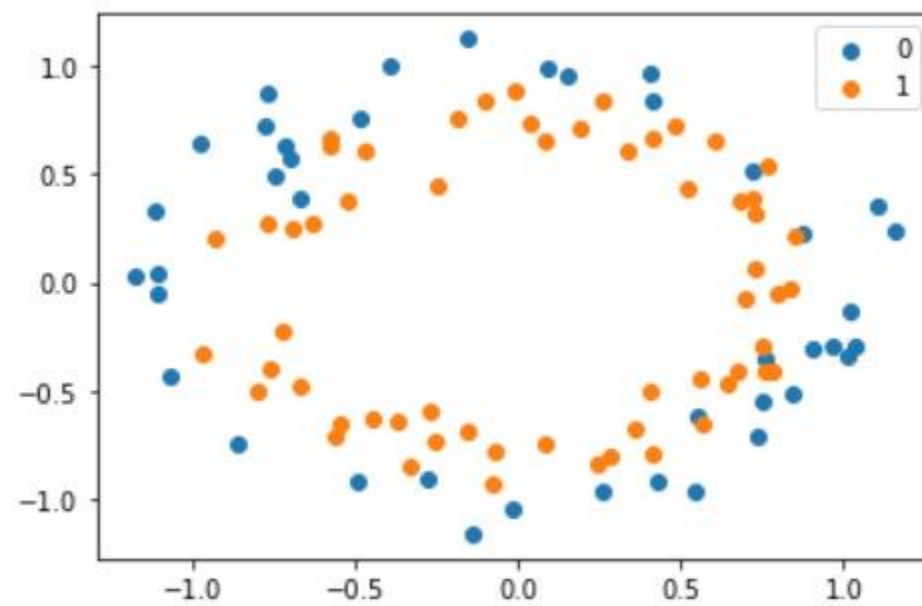| | Gini Index | Information Gain |
|---|---|---|
| **Advantages** | Easy to implement | Helpful in exploratory analysis |
| | Computational non-extensive | Sometimes Outperforms gini in data imbalance |
| | | Better in exponential data distributions |
| **Disasdvantage** | Doesn't perform well at some conditions | Computational extensive (Log) |
| | | prefer splits that result in large number of partitions, each being small but pure. |

# Part #2

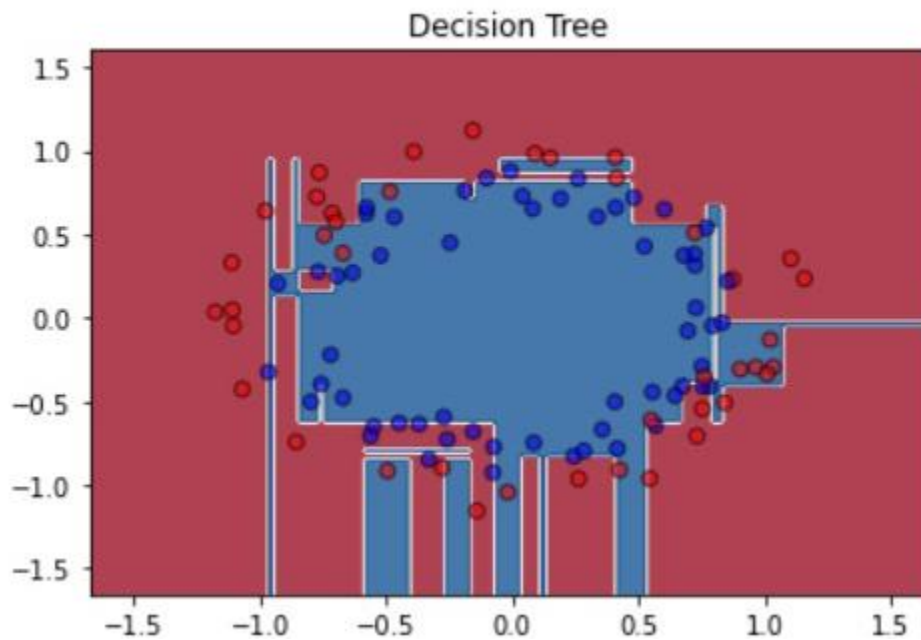**# Programming Questions**

**Q4.1)**

**Train Data:**



**Test Data:**

**Decision Boundary:**

The Accuracy = 0.6060606060606061
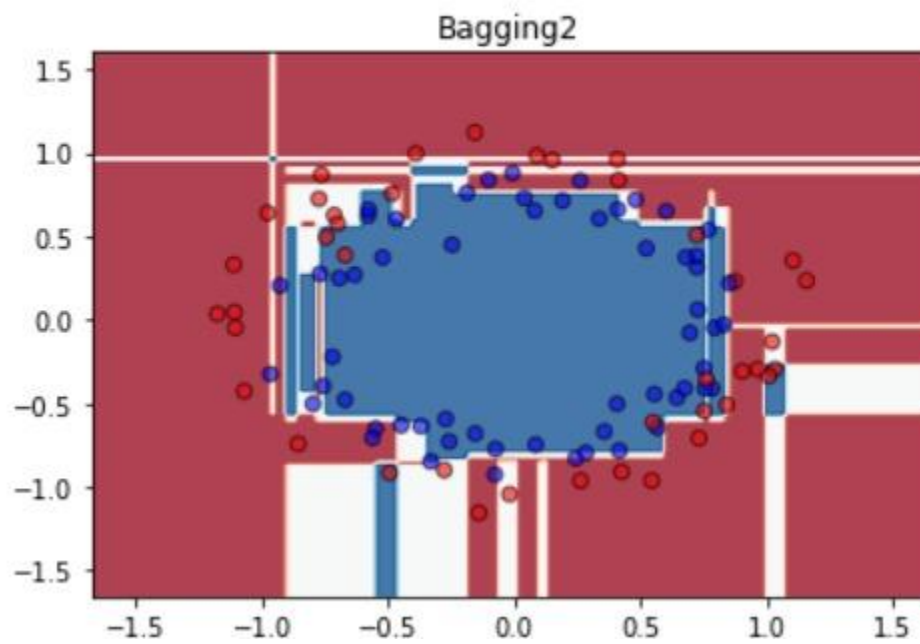
**Decision Tree**



## Q4.2)    Code:

```
7  num_estimator = [2,5,15,20]
8  # bsX = random.choices(trX, k = len(trX))
9  # bsY = random.choices(trY, k = len(trY))
10 # bsX, bsY = resample(trX, trY, random_state=rs)
11 # print(bsX)
12 h = .02
13 x_min, x_max = teX[:, 0].min() - .5, teX[:, 0].max() + .5
14 y_min, y_max = teX[:, 1].min() - .5, teX[:, 1].max() + .5
15 xx, yy = np.meshgrid(np.arange(x_min, x_max, h), np.arange(y_min, y_max, h))
16 for i in num_estimator:
17     df_prediction = pd.DataFrame()
18     Z=pd.DataFrame()
19     for n in range(i):
20         idx = random.choices(range(len(trX)), k = len(trX))
21         bsX=trX[idx]
22         bsY=trY[idx]
23         est =  DecisionTreeClassifier(random_state=rs)
24         clf = est.fit(list(bsX), list(bsY))
25         predY = clf.predict(teX)
26         Z[str(n)] = clf.predict_proba(np.c_[xx.ravel(), yy.ravel()])[:, 1]
27         df_prediction['M'+str(n)] = predY
28     df_prediction['FinalPredict'] = df_prediction.mode(axis=1)[0]
29     Z['vote']=Z.mean(axis=1)
30
31     dtAccuracy = accuracy_score(teY, df_prediction['FinalPredict'])
32     print("The Accuracy = ",dtAccuracy)
33     plotEstimator(trX, trY, teX, teY, est,'Bagging'+str(i),np.array(Z['vote']))
```
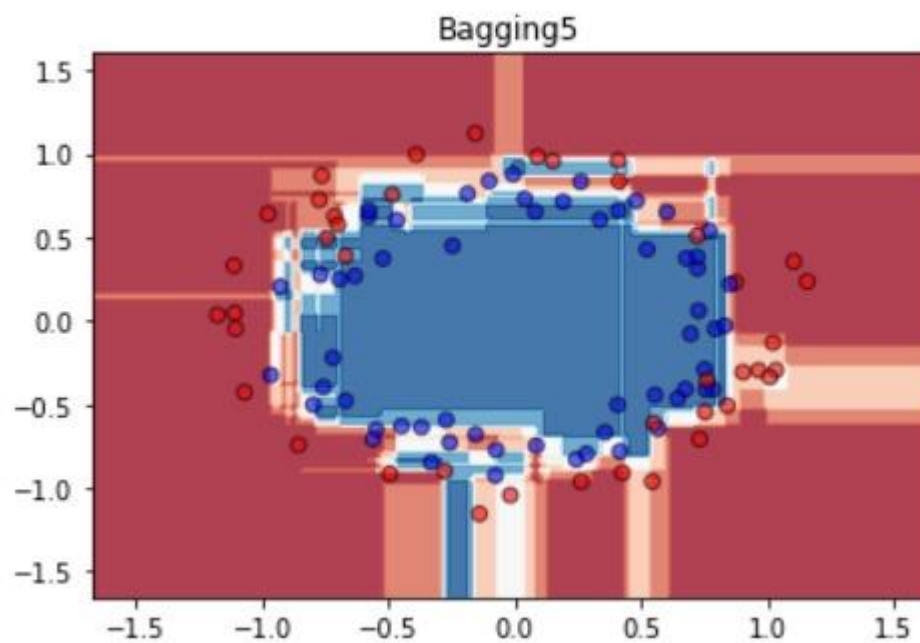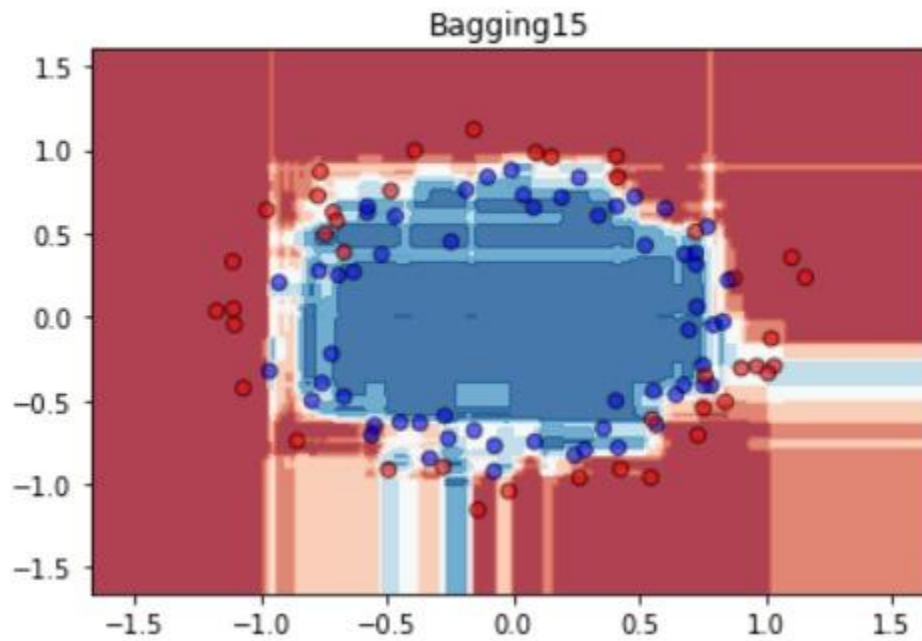
## 2 Bags

The Accuracy =  0.696969696969697



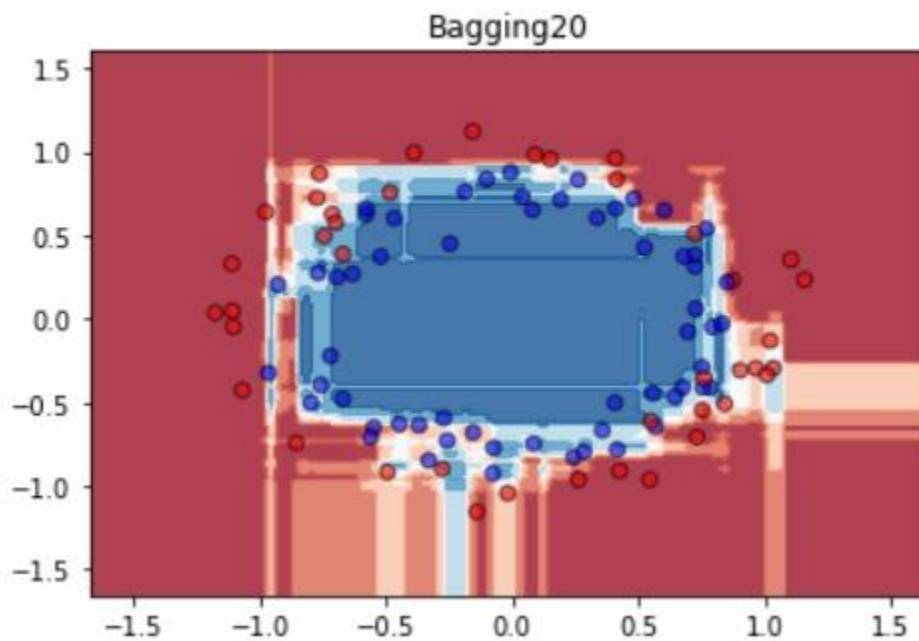Bagging2

## 5 Bags

The Accuracy =  0.7171717171717171



Bagging5

## 15 Bags

The Accuracy = 0.7070707070707071



Bagging15

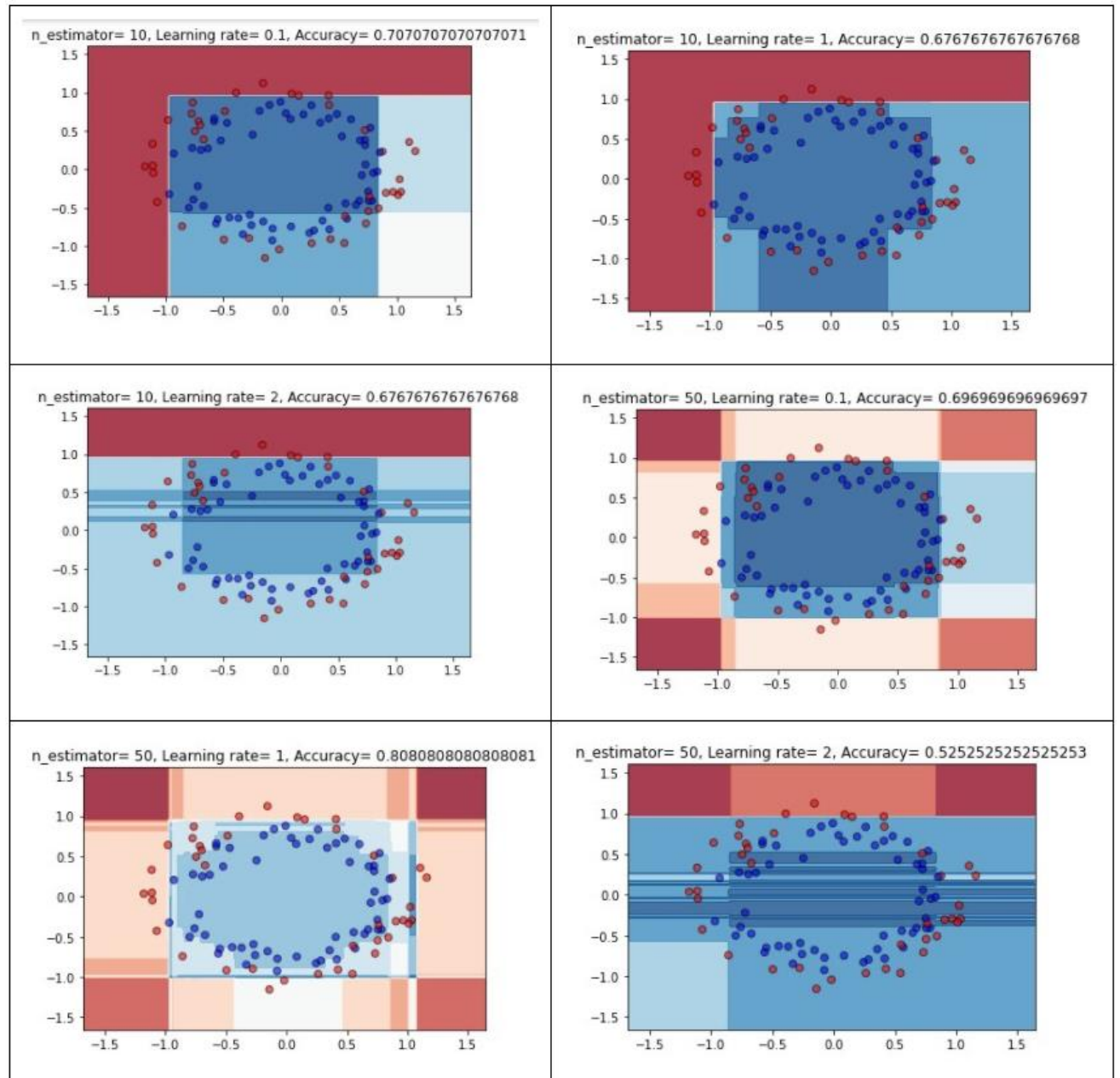## 20 Bags

The Accuracy = 0.7474747474747475



Bagging20

## Q5)

A single tree highly overfits the data. Using **'t' trees** ensures that each tree uses **different sets of data and variables**, where each tree yields to low bias and high variance. Hence, taking the average over all trees (ensembling: Bagging) or taking the **most frequent prediction** among the trees, each fitted to a subset of the original data set, we arrive to one bagged predictor, where the mean/most frequent prediction will be more stable and **less overfit**. [2]

The strength of bagging lies in the fact that it ensures that all trees are **different**. Since, joining several "weak learners" to provide a "strong learning" results in a smoother, and less wild variance in the model. [3]

Thus, the variance of the Random Forest (Bagging) is smaller compared to the variance of a single Decision Tree. [4]

So, in brief bagging overcomes overfitting by creating random subsets of the features and building smaller trees using those subsets.

**Q6)**



n_estimator= 10, Learning rate= 0.1, Accuracy= 0.7070707070707071

n_estimator= 10, Learning rate= 1, Accuracy= 0.6767676767676768

n_estimator= 10, Learning rate= 2, Accuracy= 0.6767676767676768

n_estimator= 50, Learning rate= 0.1, Accuracy= 0.696969696969697

n_estimator= 50, Learning rate= 1, Accuracy= 0.8080808080808081

n_estimator= 50, Learning rate= 2, Accuracy= 0.5252525252525253

n_estimator= 100, Learning rate= 0.1, Accuracy= 0.6767676767676768

n_estimator= 100, Learning rate= 1, Accuracy= 0.8080808080808081

n_estimator= 100, Learning rate= 2, Accuracy= 0.6464646464646465

n_estimator= 200, Learning rate= 0.1, Accuracy= 0.696969696969697

n_estimator= 200, Learning rate= 1, Accuracy= 0.797979797979798

n_estimator= 200, Learning rate= 2, Accuracy= 0.6868686868686869

# References:

[1] machine learning - When should I use Gini Impurity as opposed to Information Gain (Entropy)? - Data Science Stack Exchange

[2] (8) How does bagging avoid overfitting in Random Forest classification? - Quora

[3] Why does "bagging" in machine learning decrease variance? (techopedia.com)

[4] Understanding the Effect of Bagging on Variance and Bias visually | by Dr. Robert Kübler | Towards Data Science