# Random Forest and AdaBoost
# (Warm-up Class)
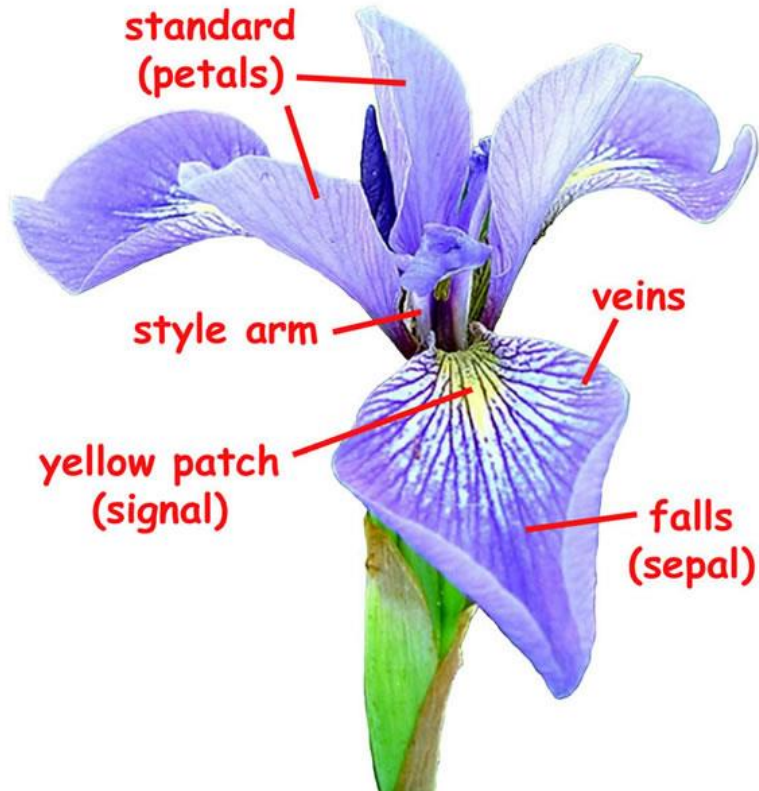
Quang-Vinh Dinh
Ph.D. in Computer Science

Year 2023

# Random Forest

Quang-Vinh Dinh
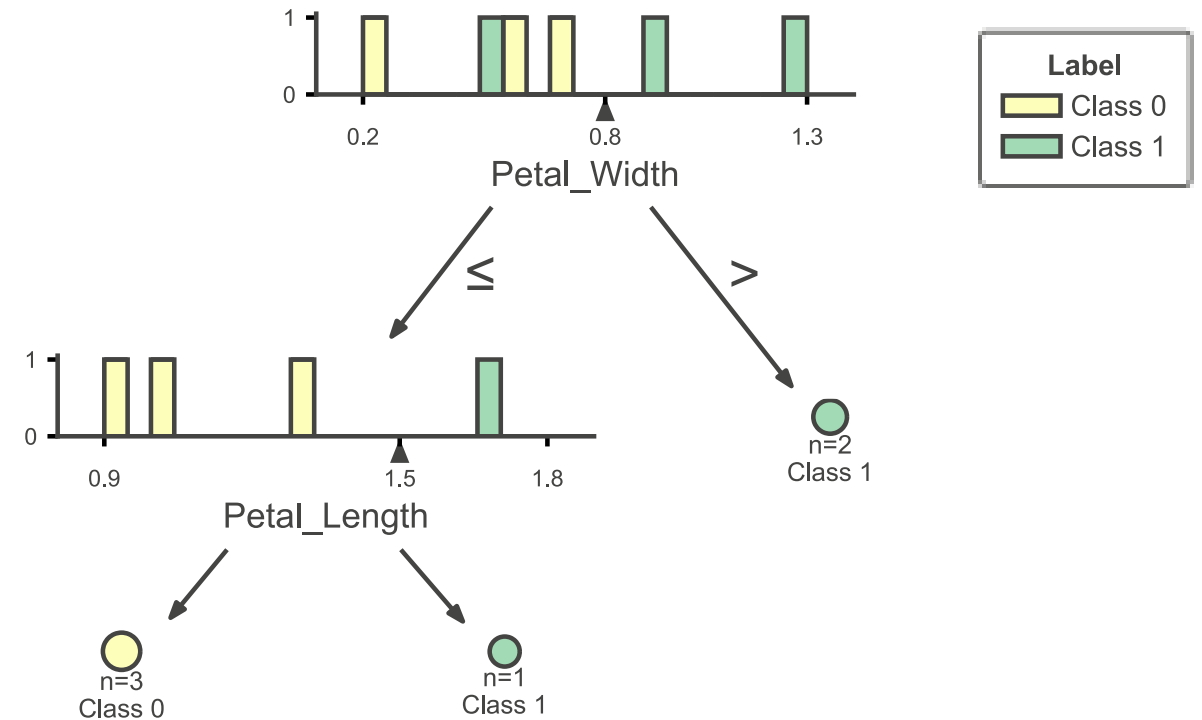Ph.D. in Computer Science
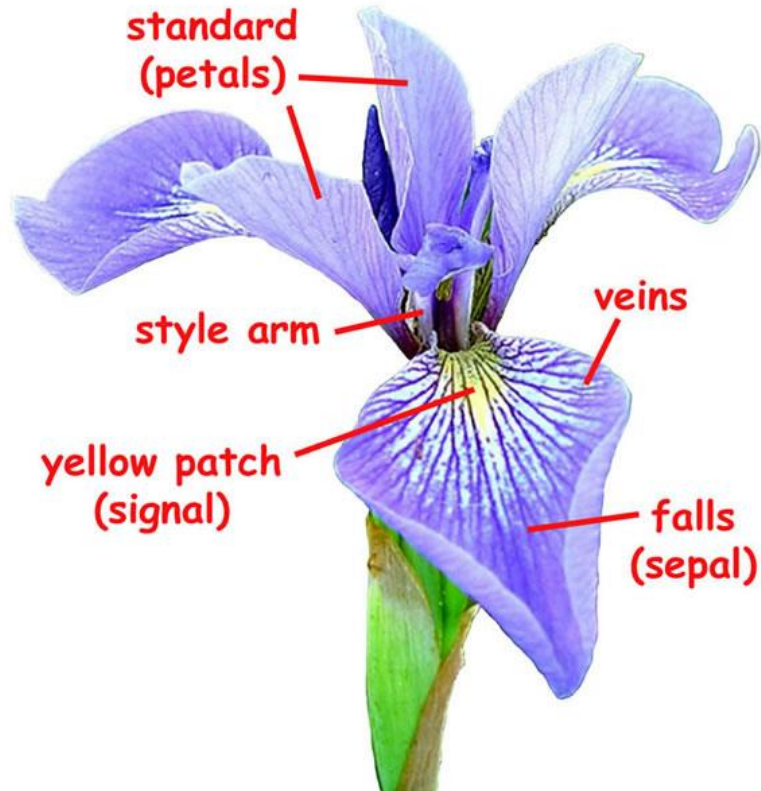
# Decision Tree

| Petal_Length | Petal_Width | Label |
|:---:|:---:|:---:|
| 1 | 0.2 | 0 |
| 1.3 | 0.6 | 0 |
| 0.9 | 0.7 | 0 |
| 1.7 | 0.5 | 1 |
| 1.8 | 0.9 | 1 |
| 1.2 | 1.3 | 1 |

❖ **Observation**



https://www.fs.usda.gov/wildflowers/
beauty/iris/flower.shtml



1

# Decision Tree

| Petal_Length | Petal_Width | Sepal_length | Label |
|:---:|:---:|:---:|:---:|
| 1 | 0.2 | 5.1 | 0 |
| 1.3 | 0.6 | 4.9 | 0 |
| 0.9 | 0.7 | 4.7 | 0 |
| 1.7 | 0.5 | 4.8 | 1 |
| 1.8 | 0.9 | 6.6 | 1 |
| 1.2 | 1.3 | 5.2 | 1 |

❖ **Observation**



https://www.fs.usda.gov/wildflowers/
beauty/iris/flower.shtml

# Decision Tree

| Petal_Length | Petal_Width | Sepal_length | Sepal_Width | Label |
|---|---|---|---|---|
| 1 | 0.2 | 5.1 | 3.5 | 0 |
| 1.3 | 0.6 | 4.9 | 3 | 0 |
| 0.9 | 0.7 | 4.7 | 3.2 | 0 |
| 1.7 | 0.5 | 4.8 | 2.8 | 1 |
| 1.8 | 0.9 | 6.6 | 3.3 | 1 |
| 1.2 | 1.3 | 5.2 | 2.4 | 1 |

❖ **Observation**



standard (petals)

style arm

veins

yellow patch (signal)

falls (sepal)

https://www.fs.usda.gov/wildflowers/
beauty/iris/flower.shtml

problem 1 $\Longrightarrow$ problem 2

# RF-Classification

| Petal_Length | Petal_Width | Label |
|:---:|:---:|:---:|
| 1 | 0.2 | 0 |
| 1.3 | 0.6 | 0 |
| 0.9 | 0.7 | 0 |
| 1.7 | 0.5 | 1 |
| 1.8 | 0.9 | 1 |
| 1.2 | 1.3 | 1 |

❖ **Simple IRIS**



Create a new dataset

Select features randomly

Build decision tree

Enough?

6

# RF-Classification

❖ **Simple IRIS**

1

Create a new dataset

Select features randomly

Build decision tree

Enough?

| Petal_Length | Petal_Width | Label |
|---|---|---|
| 1 | 0.2 | 0 |
| 1.3 | 0.6 | 0 |
| 0.9 | 0.7 | 0 |
| 1.7 | 0.5 | 1 |
| 1.8 | 0.9 | 1 |
| 1.2 | 1.3 | 1 |

| Petal_Length | Petal_Width | Label |
|---|---|---|
| 1 | 0.2 | 0 |
| 1.3 | 0.6 | 0 |
| 1 | 0.2 | 0 |
| 1.8 | 0.9 | 1 |
| 1.8 | 0.9 | 1 |
| 1.2 | 1.3 | 1 |

# RF-Classification

❖ **Simple IRIS**

1



| Petal_Length | Petal_Width | Label |
|:---:|:---:|:---:|
| 1 | 0.2 | 0 |
| 1.3 | 0.6 | 0 |
| 0.9 | 0.7 | 0 |
| 1.7 | 0.5 | 1 |
| 1.8 | 0.9 | 1 |
| 1.2 | 1.3 | 1 |

| Petal_Length | Label |
|:---:|:---:|
| 1 | 0 |
| 1.3 | 0 |
| 1 | 0 |
| 1.8 | 1 |
| 1.8 | 1 |
| 1.2 | 1 |

# RF-Classification

❖ **Simple IRIS**

1

| Petal_Length | Label |
|--------------|-------|
| 1 | 0 |
| 1.3 | 0 |
| 1 | 0 |
| 1.8 | 1 |
| 1.8 | 1 |
| 1.2 | 1 |

# RF-Classification

❖ **Simple IRIS**

2



Create a new dataset

Select features randomly

Build decision tree

Enough?

| Petal_Length | Petal_Width | Label |
|---|---|---|
| 1 | 0.2 | 0 |
| 1.3 | 0.6 | 0 |
| 0.9 | 0.7 | 0 |
| 1.7 | 0.5 | 1 |
| 1.8 | 0.9 | 1 |
| 1.2 | 1.3 | 1 |

| Petal_Length | Petal_Width | Label |
|---|---|---|
| 1.3 | 0.6 | 0 |
| 1.3 | 0.6 | 0 |
| 0.9 | 0.7 | 0 |
| 0.9 | 0.7 | 0 |
| 1.8 | 0.9 | 1 |
| 1.2 | 1.3 | 1 |

# RF-Classification

❖ **Simple IRIS**

2

Create a new dataset

Select features randomly

Build decision tree

Enough?

| Petal_Length | Petal_Width | Label |
|---|---|---|
| 1 | 0.2 | 0 |
| 1.3 | 0.6 | 0 |
| 0.9 | 0.7 | 0 |
| 1.7 | 0.5 | 1 |
| 1.8 | 0.9 | 1 |
| 1.2 | 1.3 | 1 |

| Petal_Width | Label |
|---|---|
| 0.6 | 0 |
| 0.6 | 0 |
| 0.7 | 0 |
| 0.7 | 0 |
| 0.9 | 1 |
| 1.3 | 1 |

# RF-Classification

❖ **Simple IRIS**

2



| Petal_Width | Label |
|:---:|:---:|
| 0.6 | 0 |
| 0.6 | 0 |
| 0.7 | 0 |
| 0.7 | 0 |
| 0.9 | 1 |
| 1.3 | 1 |

Create a new dataset

Select features randomly

Build decision tree

Enough?

x[0] <= 0.8
entropy = 0.918
samples = 6
value = [4, 2]

entropy = 0.0
samples = 4
value = [4, 0]

entropy = 0.0
samples = 2
value = [0, 2]

Petal_Width

≤

>

n=4
Class 0

n=2
Class 1

**Label**
Class 0
Class 1

# RF-Classification

❖ **Simple IRIS**

3

Create a new dataset

Select features randomly

Build decision tree

Enough?

| Petal_Length | Petal_Width | Label |
|---|---|---|
| 1 | 0.2 | 0 |
| 1.3 | 0.6 | 0 |
| 0.9 | 0.7 | 0 |
| 1.7 | 0.5 | 1 |
| 1.8 | 0.9 | 1 |
| 1.2 | 1.3 | 1 |

| Petal_Length | Petal_Width | Label |
|---|---|---|
| 1 | 0.2 | 0 |
| 1.3 | 0.6 | 0 |
| 1.2 | 1.3 | 1 |
| 1.8 | 0.9 | 1 |
| 1.8 | 0.9 | 1 |
| 1.2 | 1.3 | 1 |

# RF-Classification

❖ **Simple IRIS**

3



| Petal_Length | Petal_Width | Label |
|---|---|---|
| 1 | 0.2 | 0 |
| 1.3 | 0.6 | 0 |
| 0.9 | 0.7 | 0 |
| 1.7 | 0.5 | 1 |
| 1.8 | 0.9 | 1 |
| 1.2 | 1.3 | 1 |

| Petal_Length | Label |
|---|---|
| 1 | 0 |
| 1.3 | 0 |
| 1.2 | 1 |
| 1.8 | 1 |
| 1.8 | 1 |
| 1.2 | 1 |

Create a new dataset

Select features randomly

Build decision tree

Enough?

# RF-Classification

❖ **Simple IRIS**

3



```
x[0] <= 1.1
entropy = 0.918
samples = 6
value = [2, 4]
```

```
entropy = 0.0
samples = 1
value = [1, 0]
```

```
x[0] <= 1.25
entropy = 0.722
samples = 5
value = [1, 4]
```

```
entropy = 0.0
samples = 2
value = [0, 2]
```

```
entropy = 0.918
samples = 3
value = [1, 2]
```

| Petal_Length | Label |
|:---:|:---:|
| 1 | 0 |
| 1.3 | 0 |
| 1.2 | 1 |
| 1.8 | 1 |
| 1.8 | 1 |
| 1.2 | 1 |

# RF-Classification

❖ **Simple IRIS**

inference

Petal_Length = 1.7

Petal_Width = 0.8

| Petal_Length | Petal_Width | Label |
|---|---|---|
| 1 | 0.2 | 0 |
| 1.3 | 0.6 | 0 |
| 0.9 | 0.7 | 0 |
| 1.7 | 0.5 | 1 |
| 1.8 | 0.9 | 1 |
| 1.2 | 1.3 | 1 |

# RF-Classification

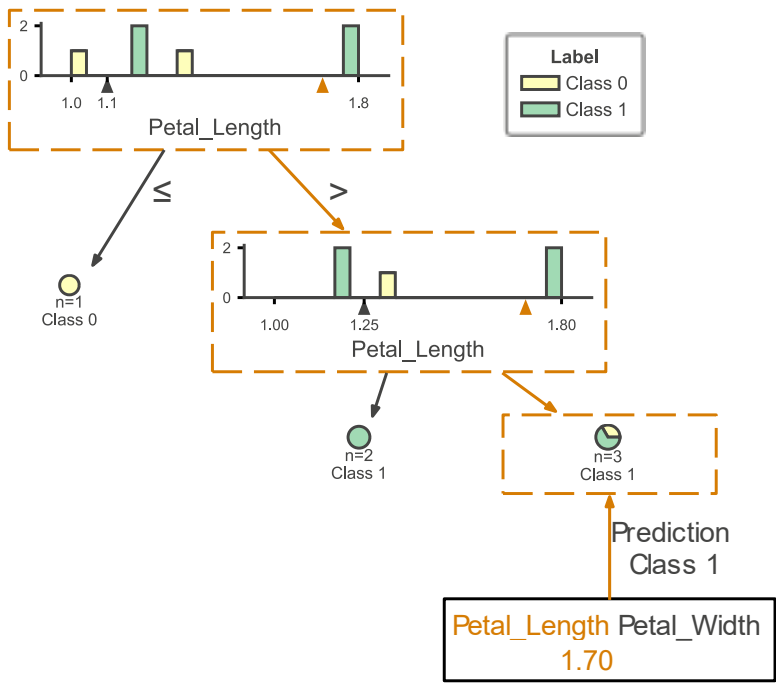❖ **Using sklearn**
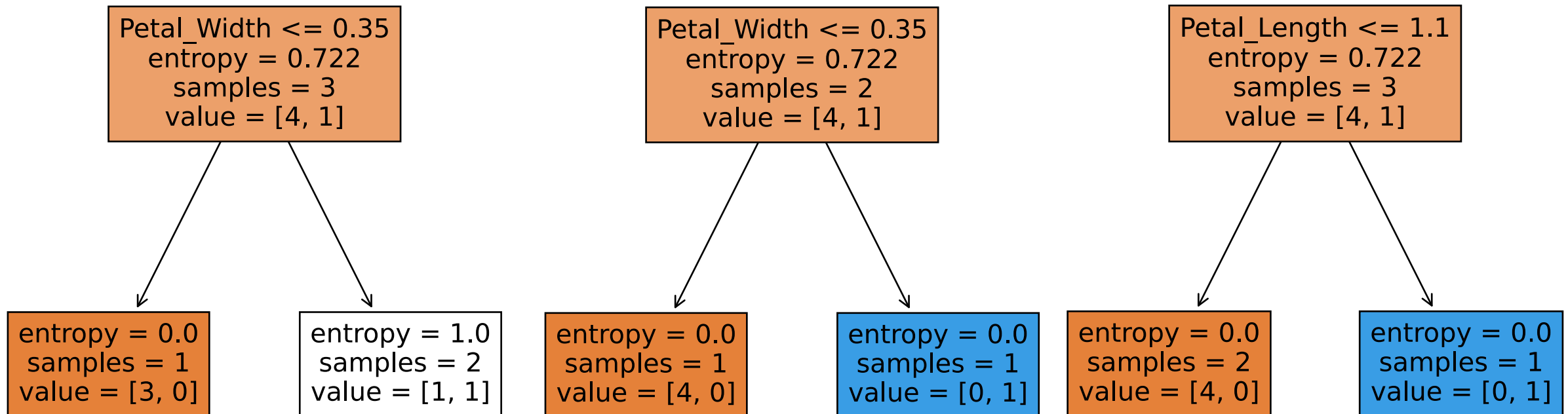
Another experiment

```
1  rf_classifier = RandomForestClassifier(n_estimators=3,
2                                          max_features=1,
3                                          max_depth=1,
4                                          criterion='entropy',
5                                          max_samples=5)
6  rf_classifier.fit(x_data, y_train)
7  rf_classifier.predict(np.array([[2.7, 0.8]]))
```
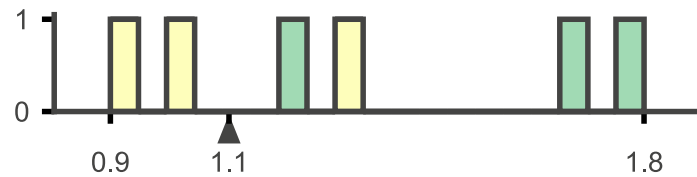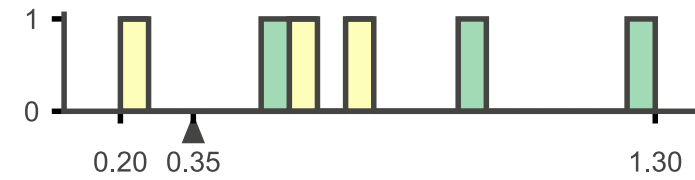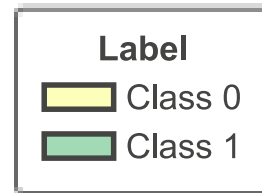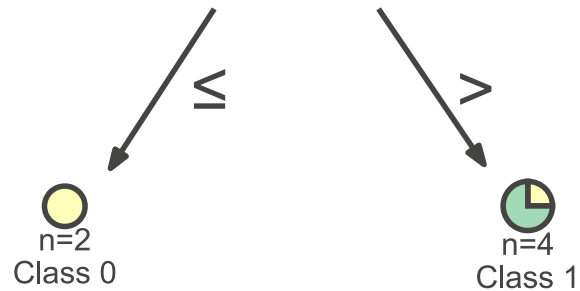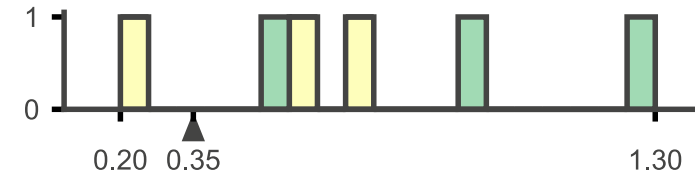
```
array([0])
```

# RF-Classification

❖ **Using sklearn**

Using all the training samples

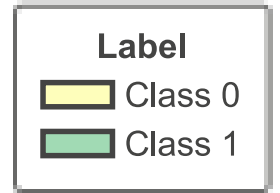| Outlook | Temp | Humidity | Wind | Play Tennis |
|---|---|---|---|---|
| Sunny | Hot | High | Weak | No |
| Sunny | Hot | High | Strong | No |
| Overcast | Hot | High | Weak | Yes |
| Rain | Mild | High | Weak | Yes |
| Rain | Cool | Normal | Weak | Yes |
| Rain | Cool | Normal | Strong | No |
| Overcast | Cool | Normal | Strong | Yes |
| Sunny | Mild | High | Weak | No |
| Sunny | Cool | Normal | Weak | Yes |
| Rain | Mild | Normal | Weak | Yes |
| Sunny | Mild | Normal | Strong | Yes |
| Overcast | Mild | High | Strong | Yes |
| Overcast | Hot | Normal | Weak | Yes |
| Rain | Mild | High | Strong | No |

**Entropy:**

$$E(S) = -\sum_{c \in C} p_c log_2 p_c$$

**Information Gain**

$$IG(S,F) = E(S) - \sum_{f \in F} \frac{|S_f|}{|S|} E(S_f)$$

Decision Tree

$$S = \{9: Yes, 5: No\} \longrightarrow E(S) = -\frac{9}{14} log_2\left(\frac{9}{14}\right) - \frac{5}{14} log_2\left(\frac{5}{14}\right) = 0.94$$

$$S_{weak} = \{6: Yes, 2: No\} \longrightarrow E(S_{weak}) = -\frac{6}{8} log_2\left(\frac{6}{8}\right) - \frac{2}{8} log_2\left(\frac{6}{8}\right) = 0.811$$

$$S_{Strong} = \{3: Yes, 3: No\} \longrightarrow E(S_{Strong}) = -\frac{3}{6} log_2\left(\frac{3}{6}\right) - \frac{3}{6} log_2\left(\frac{3}{6}\right) = 1$$

$$\Rightarrow Gain(S, Wind) = E(S) - \frac{8}{14} E(S_{weak}) - \frac{6}{14} E(S_{Strong})$$

**Category = 2**

$$= 0.94 - \frac{8}{14} * 0.811 - \frac{6}{14} * 1 = 0.048$$

**Category = 3 > 2 → Combine →**
Option_1: Sunny - (Overcast, Rain)
Option_2: Overcast - (Sunny, Rain)
Option_3: Rain – (Sunny, Overcast)

$$Gain(S, Outlook) = max \begin{cases} IG(S, Option\_1) = 0.102 \\ IG(S, Option\_2) = 0.226 \\ IG(S, Option\_3) = 0.003 \end{cases}$$

$$S_{Sunny} = \{2: Yes, 3: No\} \longrightarrow E(S_{Sunny}) = 0.97$$

$$S_{Overcast,Rain} = \{7: Yes, 2: No\} \longrightarrow E(S_{Overcast,Rain}) = 0.764$$

$$IG(S, Option\_1)$$

$$= E(S) - \frac{5}{14} E(S_{Sunny}) - \frac{9}{14} E(S_{Overcast,Rain})$$

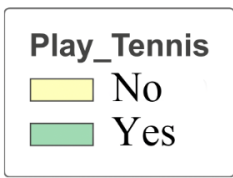$$= 0.94 - \frac{5}{14} * 0.97 - \frac{9}{14} * 0.764 = 0.102$$

Gain(S, Outlook) = 0.226

Gain(S, Temp) = 0.015

Gain(S, Humidity) = 0.151

Gain(S, Wind) = 0.048

**Choose Outlook with highest Gain score for root node**

**Option_2 is used to split**

19

Decision Tree

20

Training phase

Test = <outlook=Sunny, temperature=Hot, humidity=High, Wind=Weak>

Tree 1

Tree 2

Tree 3

predict

predict

predict

No

Yes

Yes

Voting

Final predict: Yes

Test phase

```python
from sklearn.ensemble import RandomForestClassifier

classifier = RandomForestClassifier(n_estimators=3,
                                     max_features=2,
                                     criterion='entropy',
                                     max_samples=10)

classifier.fit(X, y)
```

| Outlook | Temp | Humidity | Wind | Play Tennis |
|---------|------|----------|------|-------------|
| Sunny | Hot | High | Weak | No |
| Sunny | Hot | High | Strong | No |
| Overcast | Hot | High | Weak | Yes |
| Rain | Mild | High | Weak | Yes |
| Rain | Cool | Normal | Weak | Yes |
| Rain | Cool | Normal | Strong | No |
| Overcast | Cool | Normal | Strong | Yes |
| Sunny | Mild | High | Weak | No |
| Sunny | Cool | Normal | Weak | Yes |
| Rain | Mild | Normal | Weak | Yes |
| Sunny | Mild | Normal | Strong | Yes |
| Overcast | Mild | High | Strong | Yes |
| Overcast | Hot | Normal | Weak | Yes |
| Rain | Mild | High | Strong | No |

Tree 1



23

```
from sklearn.ensemble import RandomForestClassifier

classifier = RandomForestClassifier(n_estimators=3,
                                     max_features=2,
                                     criterion='entropy',
                                     max_samples=10)

classifier.fit(X, y)
```

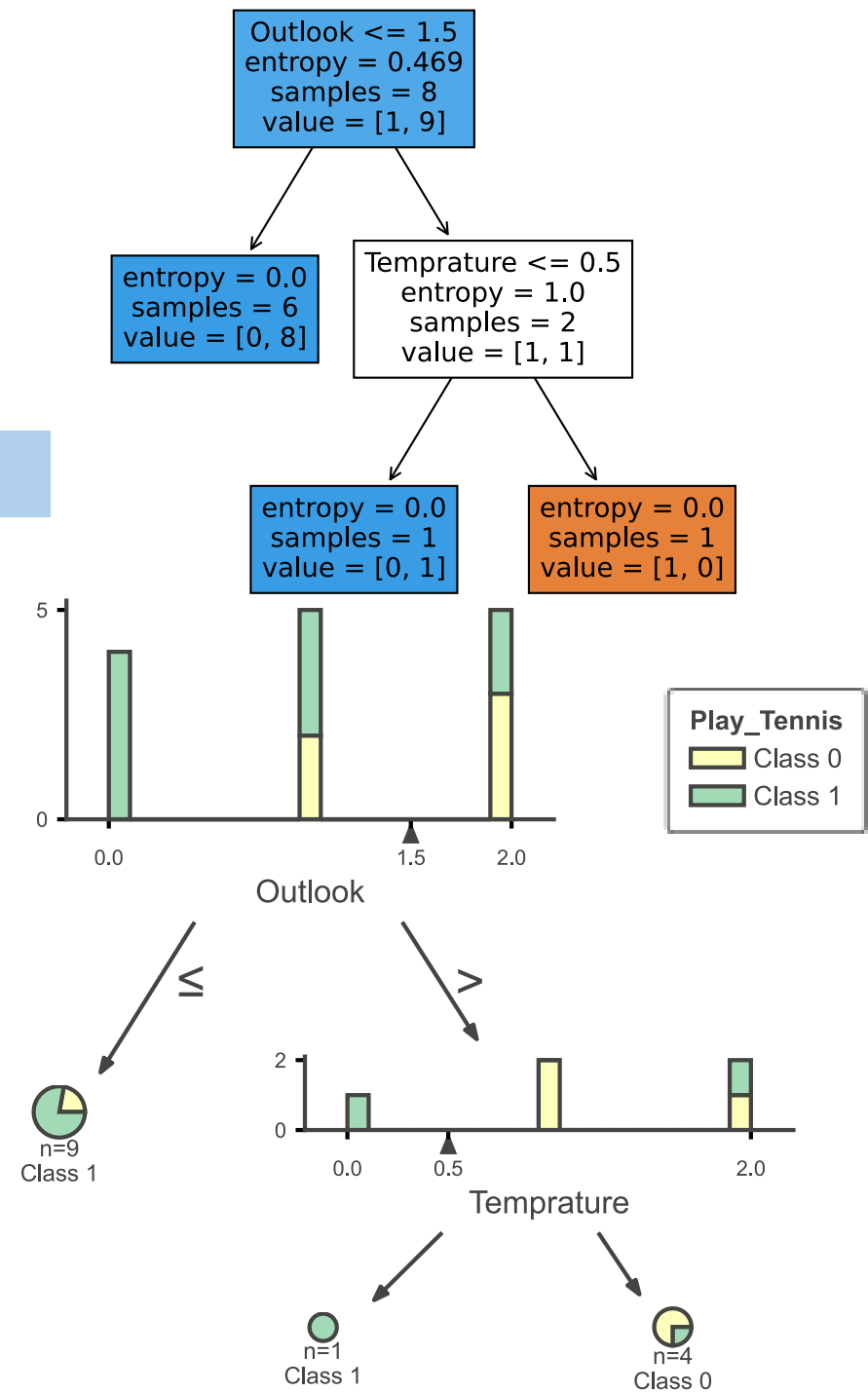| Outlook | Temp | Humidity | Wind | Play Tennis |
|---------|------|----------|------|-------------|
| Sunny | Hot | High | Weak | No |
| Sunny | Hot | High | Strong | No |
| Overcast | Hot | High | Weak | Yes |
| Rain | Mild | High | Weak | Yes |
| Rain | Cool | Normal | Weak | Yes |
| Rain | Cool | Normal | Strong | No |
| Overcast | Cool | Normal | Strong | Yes |
| Sunny | Mild | High | Weak | No |
| Sunny | Cool | Normal | Weak | Yes |
| Rain | Mild | Normal | Weak | Yes |
| Sunny | Mild | Normal | Strong | Yes |
| Overcast | Mild | High | Strong | Yes |
| Overcast | Hot | Normal | Weak | Yes |
| Rain | Mild | High | Strong | No |

Tree 2



24

Tree 3

| Outlook | Temp | Humidity | Wind | Play Tennis |
|---|---|---|---|---|
| Sunny | Hot | High | Weak | No |
| Sunny | Hot | High | Strong | No |
| Overcast | Hot | High | Weak | Yes |
| Rain | Mild | High | Weak | Yes |
| Rain | Cool | Normal | Weak | Yes |
| Rain | Cool | Normal | Strong | No |
| Overcast | Cool | Normal | Strong | Yes |
| Sunny | Mild | High | Weak | No |
| Sunny | Cool | Normal | Weak | Yes |
| Rain | Mild | Normal | Weak | Yes |
| Sunny | Mild | Normal | Strong | Yes |
| Overcast | Mild | High | Strong | Yes |
| Overcast | Hot | Normal | Weak | Yes |
| Rain | Mild | High | Strong | No |

```
from sklearn.ensemble import RandomForestClassifier

classifier = RandomForestClassifier(n_estimators=3,
                                     max_features=2,
                                     criterion='entropy',
                                     max_samples=10)

classifier.fit(X, y)
```
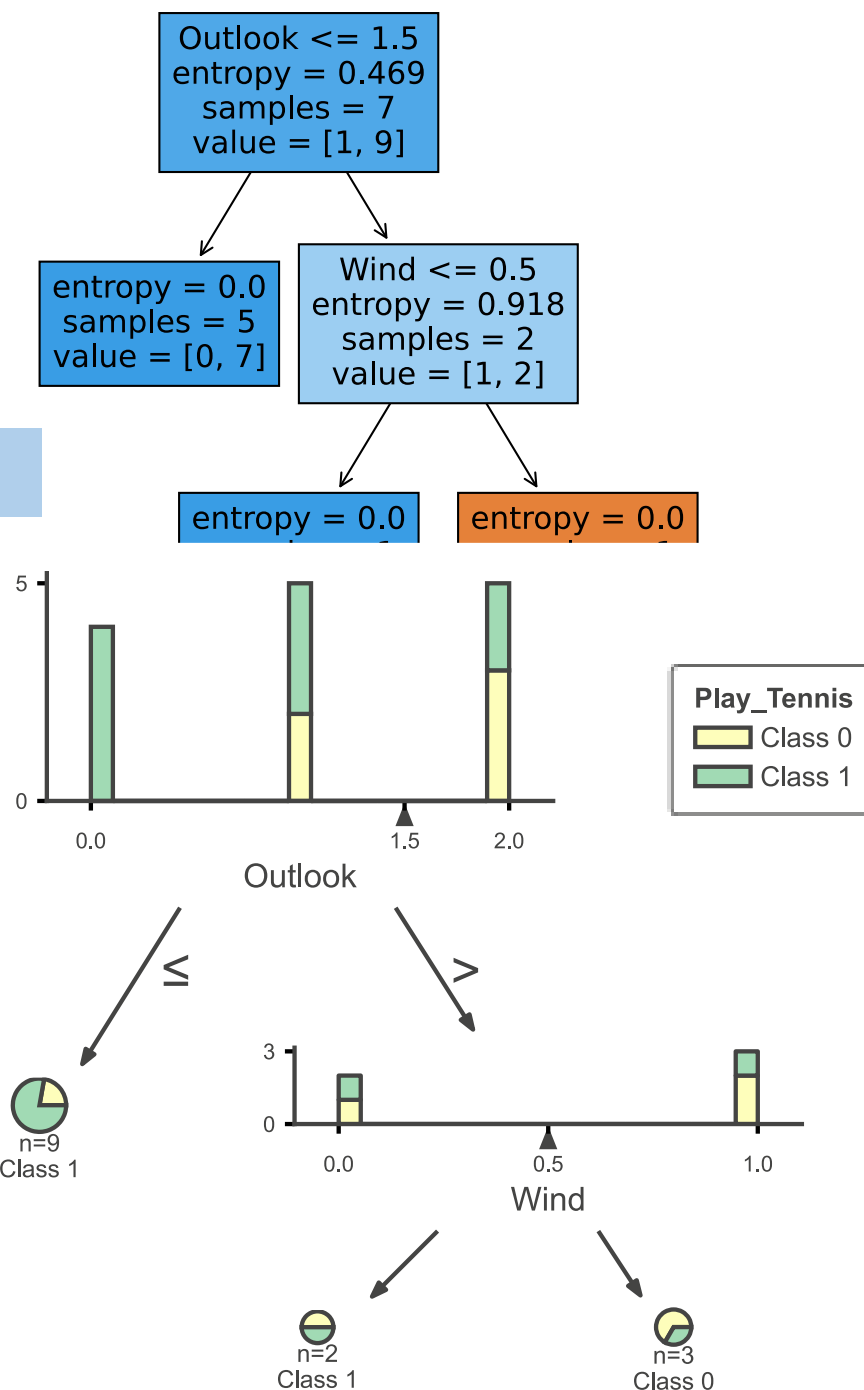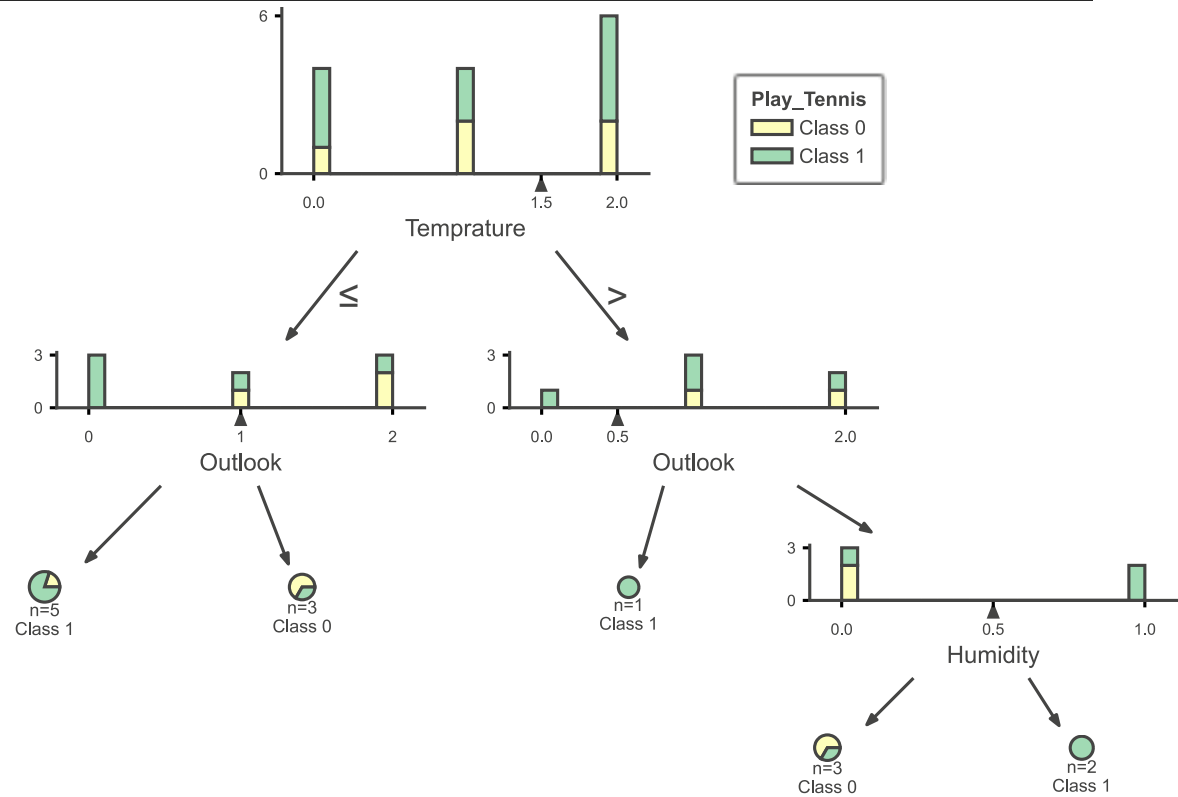


25

# Decision Tree - Regression

❖ **Salary prediction**

| Experience | Salary |
|------------|--------|
| 1 | 0 |
| 1.5 | 0 |
| 2 | 0 |
| 2.5 | 0 |
| 3 | 60 |
| 3.5 | 64 |
| 4 | 55 |
| 4.5 | 61 |
| 5 | 66 |
| 5.5 | 83 |
| 6 | 93 |
| 6.5 | 91 |
| 7 | 98 |
| 7.5 | 101 |



**When Experience = 5.3,**

**Salary = ?**

| Experience | Salary |
|------------|--------|
| 1 | 0 |
| 1.5 | 0 |
| 2 | 0 |
| 2.5 | 0 |
| 3 | 60 |
| 3.5 | 64 |
| 4 | 55 |
| 4.5 | 61 |
| 5 | 66 |
| 5.5 | 83 |
| 6 | 93 |
| 6.5 | 91 |
| 7 | 98 |
| 7.5 | 101 |

| Experience | Salary |
|------------|--------|
| 1 | 0 |



$$\mu_L = \frac{1}{|L|} \sum_i L_i = 0$$

$$mse_L = \frac{1}{|L|} \sum_i (L_i - \mu)^2 = 0$$

| Experience | Salary |
|------------|--------|
| 1.5 | 0 |
| 2 | 0 |
| 2.5 | 0 |
| 3 | 60 |
| 3.5 | 64 |
| 4 | 55 |
| 4.5 | 61 |
| 5 | 66 |
| 5.5 | 83 |
| 6 | 93 |
| 6.5 | 91 |
| 7 | 98 |
| 7.5 | 101 |

$$a_{mse} = \frac{|L|}{|S|} mse_L + \frac{|R|}{|S|} mse_R$$

$$= \frac{1}{14} * 0 + \frac{13}{14} * 1275.15$$

$$= 1184.07$$

$$\mu = \frac{1}{|S|} \sum_i S_i = 55.14$$

$$mse = \frac{1}{|S|} \sum_i (S_i - \mu)^2 = 1417.97$$

$$\mu_R = \frac{1}{|R|} \sum_i R_i = 59.38$$

$$mse_R = \frac{1}{|R|} \sum_i (R_i - \mu)^2 = 1275.15$$

| Experience | Salary |
|------------|--------|
| 1 | 0 |
| 1.5 | 0 |
| 2 | 0 |
| 2.5 | 0 |
| 3 | 60 |
| 3.5 | 64 |
| 4 | 55 |
| 4.5 | 61 |
| 5 | 66 |
| 5.5 | 83 |
| 6 | 93 |
| 6.5 | 91 |
| 7 | 98 |
| 7.5 | 101 |

| Experience | Salary |
|------------|--------|
| 1 | 0 |
| 1.5 | 0 |
| 2 | 0 |
| 2.5 | 0 |

$$\mu_L = \frac{1}{|L|} \sum_i L_i = 0$$

$$mse_L = \frac{1}{|L|} \sum_i (L_i - \mu)^2 = 0$$

| Experience | Salary |
|------------|--------|
| 3 | 60 |
| 3.5 | 64 |
| 4 | 55 |
| 4.5 | 61 |
| 5 | 66 |
| 5.5 | 83 |
| 6 | 93 |
| 6.5 | 91 |
| 7 | 98 |
| 7.5 | 101 |

$$a_{mse} = \frac{|L|}{|S|} mse_L + \frac{|R|}{|S|} mse_R$$

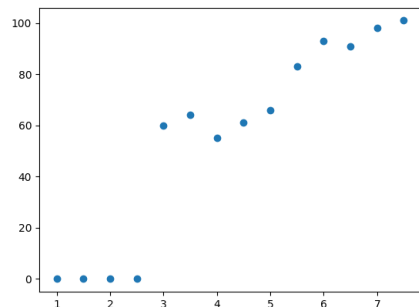$$= \frac{4}{14} * 0 + \frac{10}{14} * 282.35$$

$$= 201.68$$

$$\mu = \frac{1}{|S|} \sum_i S_i = 55.14$$

$$mse = \frac{1}{|S|} \sum_i (S_i - \mu)^2 = 1417.97$$

$$\mu_R = \frac{1}{|R|} \sum_i R_i = 77.2$$

$$mse_R = \frac{1}{|R|} \sum_i (R_i - \mu)^2 = 282.35$$

| Experience | Salary |
|:---:|:---:|
| 1 | 0 |
| 1.5 | 0 |
| 2 | 0 |
| 2.5 | 0 |
| 3 | 60 |
| 3.5 | 64 |
| 4 | 55 |
| 4.5 | 61 |
| 5 | 66 |
| 5.5 | 83 |
| 6 | 93 |
| 6.5 | 91 |
| 7 | 98 |
| 7.5 | 101 |

$a_{mse} = 1184.07$

$a_{mse} = 911.19$

$a_{mse} = 588.68$

$a_{mse} = 201.68$

$a_{mse} = 383.92$

$a_{mse} = 526.52$

$a_{mse} = 543.51$

$a_{mse} = 575.09$

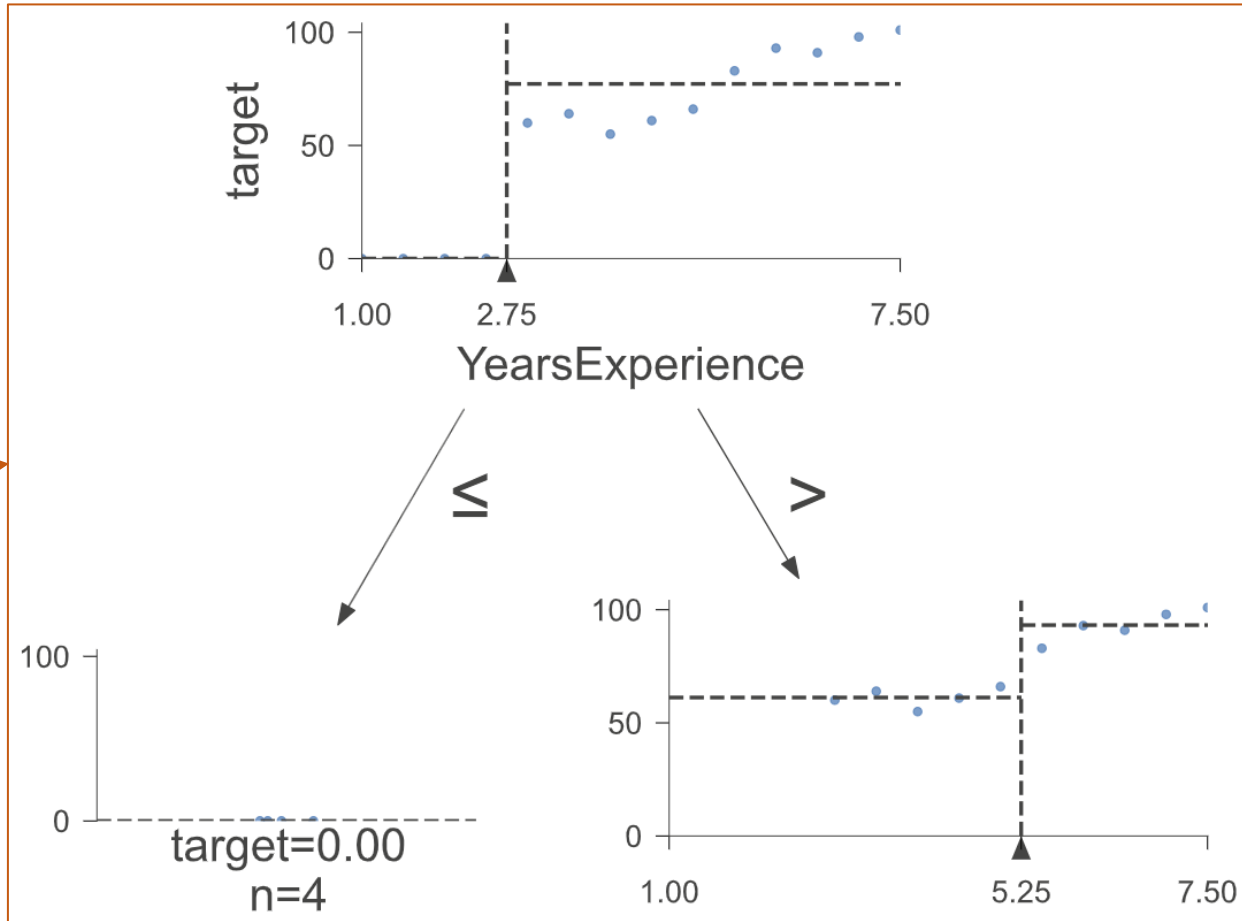$a_{mse} = 613.34$

$a_{mse} = 758.4$

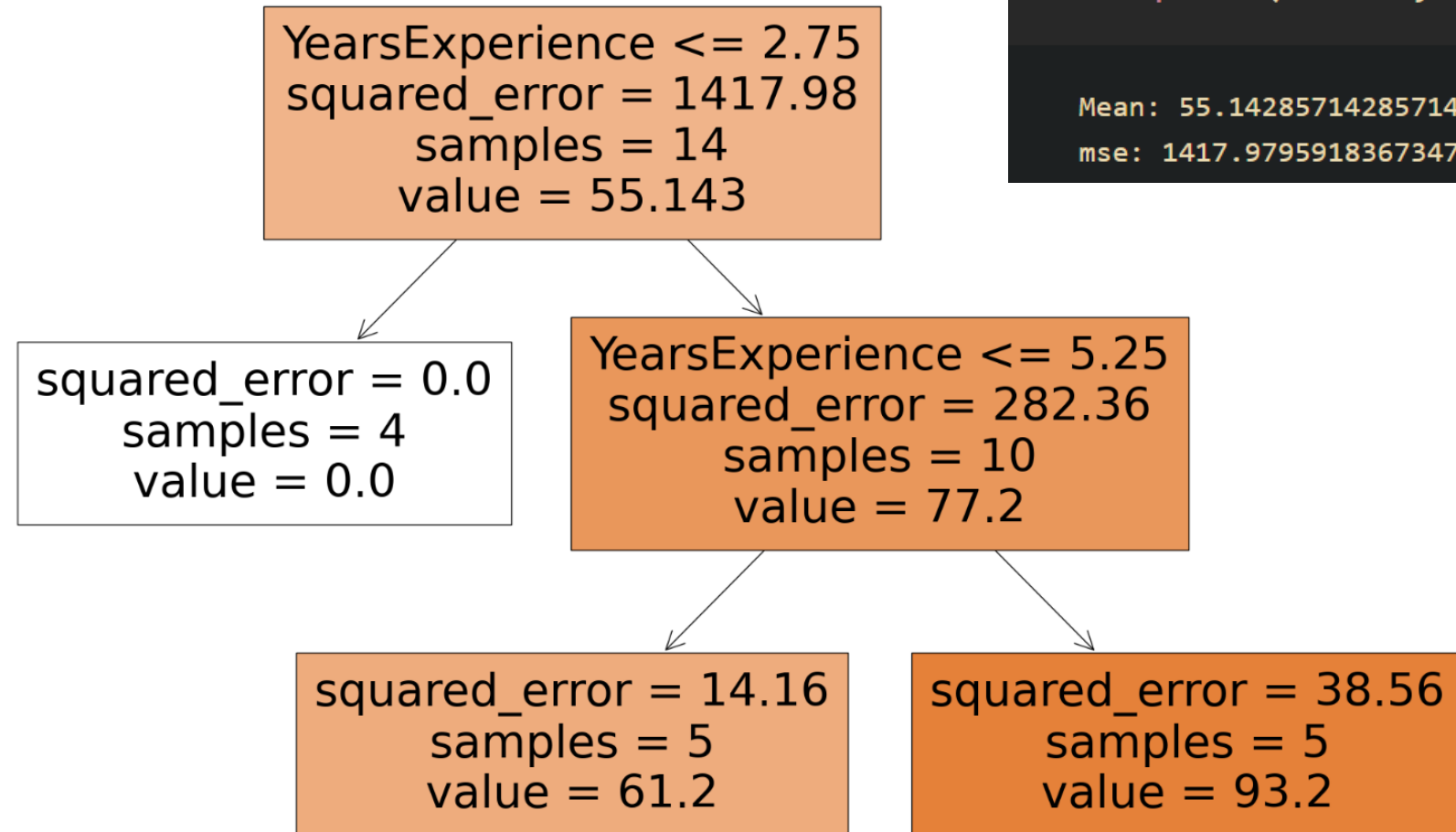$a_{mse} = 947.73$

$a_{mse} = 1090.05$

$a_{mse} = 1256.21$



| Experience | Salary |
|:---:|:---:|
| 1 | 0 |
| 1.5 | 0 |
| 2 | 0 |
| 2.5 | 0 |

| Experience | Salary |
|:---:|:---:|
| 3 | 60 |
| 3.5 | 64 |
| 4 | 55 |
| 4.5 | 61 |
| 5 | 66 |
| 5.5 | 83 |
| 6 | 93 |
| 6.5 | 91 |
| 7 | 98 |
| 7.5 | 101 |

29

# Decision Tree - Regression

```
1  y_mean = y.mean()
2  print('Mean:', y_mean)
3
4  diff = (y - y_mean)**2
5  mse = diff.sum()/14
6  print('mse:', mse)
```

```
Mean: 55.142857142857146
mse: 1417.9795918367347
```
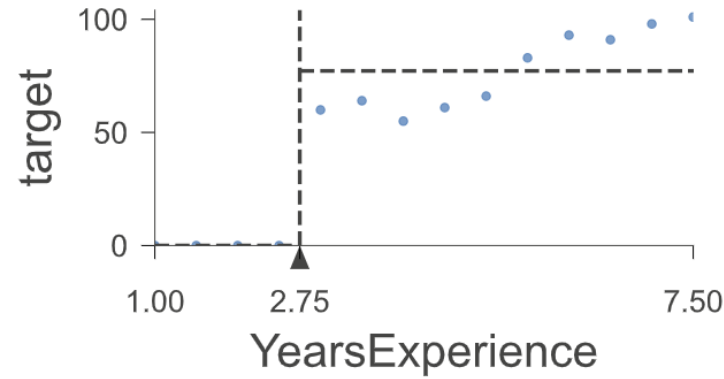
❖ **Salary prediction**
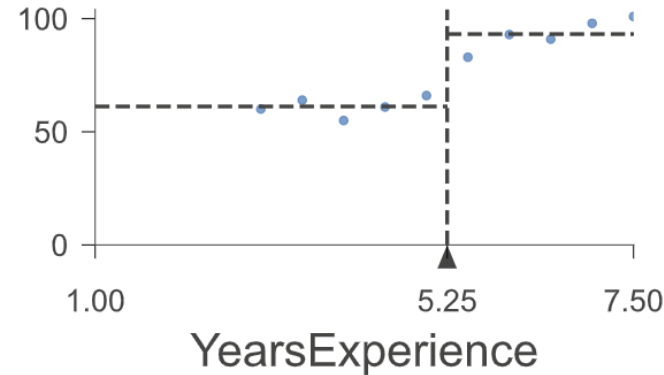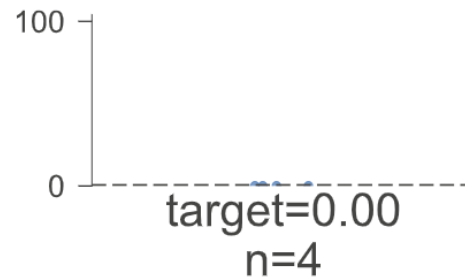
| Experience | Salary |
|------------|--------|
| 1 | 0 |
| 1.5 | 0 |
| 2 | 0 |
| 2.5 | 0 |
| 3 | 60 |
| 3.5 | 64 |
| 4 | 55 |
| 4.5 | 61 |
| 5 | 66 |
| 5.5 | 83 |
| 6 | 93 |
| 6.5 | 91 |
| 7 | 98 |
| 7.5 | 101 |

```
YearsExperience <= 2.75
squared_error = 1417.98
samples = 14
value = 55.143
```

```
squared_error = 0.0
samples = 4
value = 0.0
```

```
YearsExperience <= 5.25
squared_error = 282.36
samples = 10
value = 77.2
```

```
squared_error = 14.16
samples = 5
value = 61.2
```

```
squared_error = 38.56
samples = 5
value = 93.2
```

30

Decision Tree Regression

| Experience | Salary |
|---|---|
| 1 | 0 |
| 1.5 | 0 |
| 2 | 0 |
| 2.5 | 0 |
| 3 | 60 |
| 3.5 | 64 |
| 4 | 55 |
| 4.5 | 61 |
| 5 | 66 |
| 5.5 | 83 |
| 6 | 93 |
| 6.5 | 91 |
| 7 | 98 |
| 7.5 | 101 |

| Experience | Salary |
|---|---|
| 3 | 60 |
| 3.5 | 64 |
| 4 | 55 |
| 4.5 | 61 |
| 5 | 66 |
| 5.5 | 83 |
| 6 | 93 |
| 6.5 | 91 |
| 7 | 98 |
| 7.5 | 101 |

| Experience | Salary |
|---|---|
| 1 | 0 |
| 1.5 | 0 |
| 2 | 0 |
| 2.5 | 0 |

target=0.00
n=4

| Experience | Salary |
|---|---|
| 3 | 60 |
| 3.5 | 64 |
| 4 | 55 |
| 4.5 | 61 |
| 5 | 66 |

target=61.20
n=5

target=93.20
n=5

| Experience | Salary |
|---|---|
| 5.5 | 83 |
| 6 | 93 |
| 6.5 | 91 |
| 7 | 98 |
| 7.5 | 101 |

31

# Random Forest Regression

```
1  dt_regressor = RandomForestRegressor(n_estimators=3,
2                                         max_depth=2)
3  dt_regressor.fit(X, y)
```

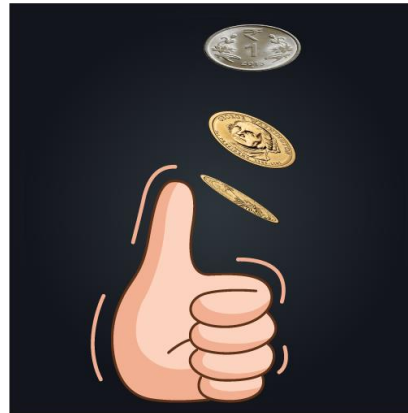❖ **Salary prediction**

# Random Forest

❖ **Bernoulli Random variables**

A numerical description of the outcome
of a statistical experiment

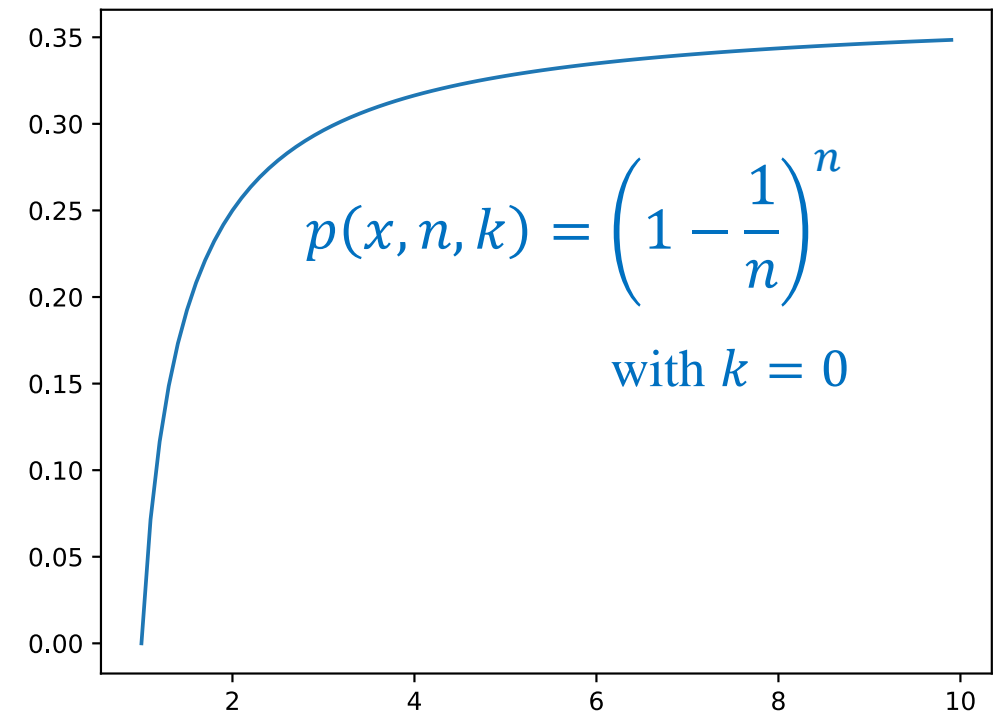$$p(x) = p\{X = x\} = \begin{cases} p & when\ x = 1 \\ 1 - p & when\ x = 0 \end{cases}$$

$$p(x, n, k) = C_n^k \left(\frac{1}{n}\right)^k \left(1 - \frac{1}{n}\right)^{n-k}$$

**Toss a coin**

Sample space: S = {tail, head}

X= {0, 1}

$$p(x, n, k) = \left(1 - \frac{1}{n}\right)^n$$

with $k = 0$

# **Adaptive Boosting**
# **(Warm-up Class)**

Quang-Vinh Dinh
Ph.D. in Computer Science

*Year 2023*

# AdaBoost

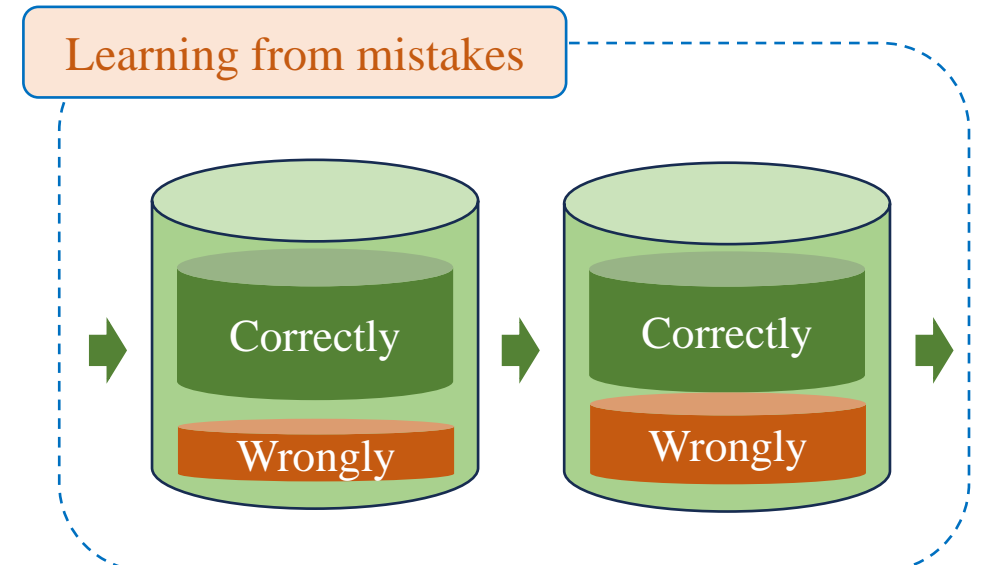❖ **Idea**
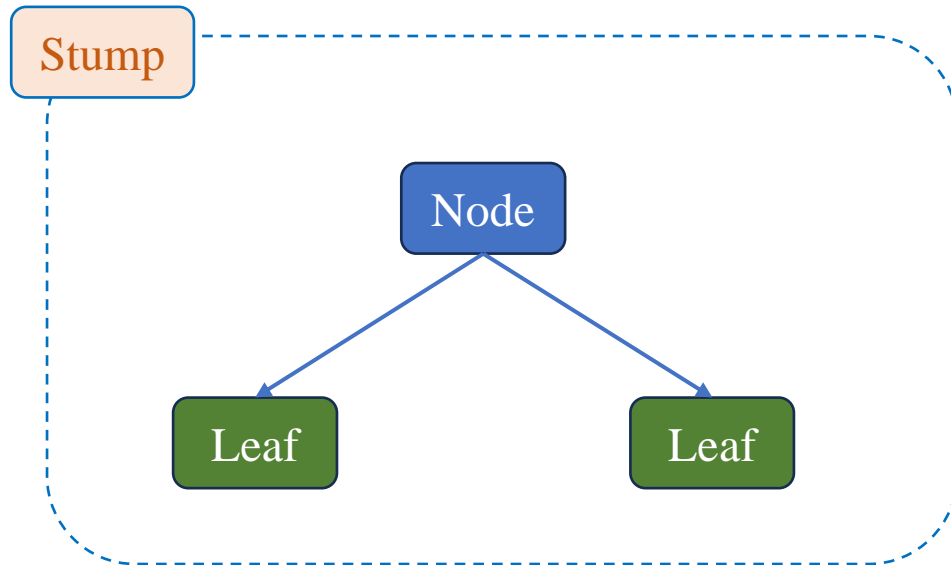


https://www.fs.usda.gov/wildflowers/beauty/iris/flower.shtml
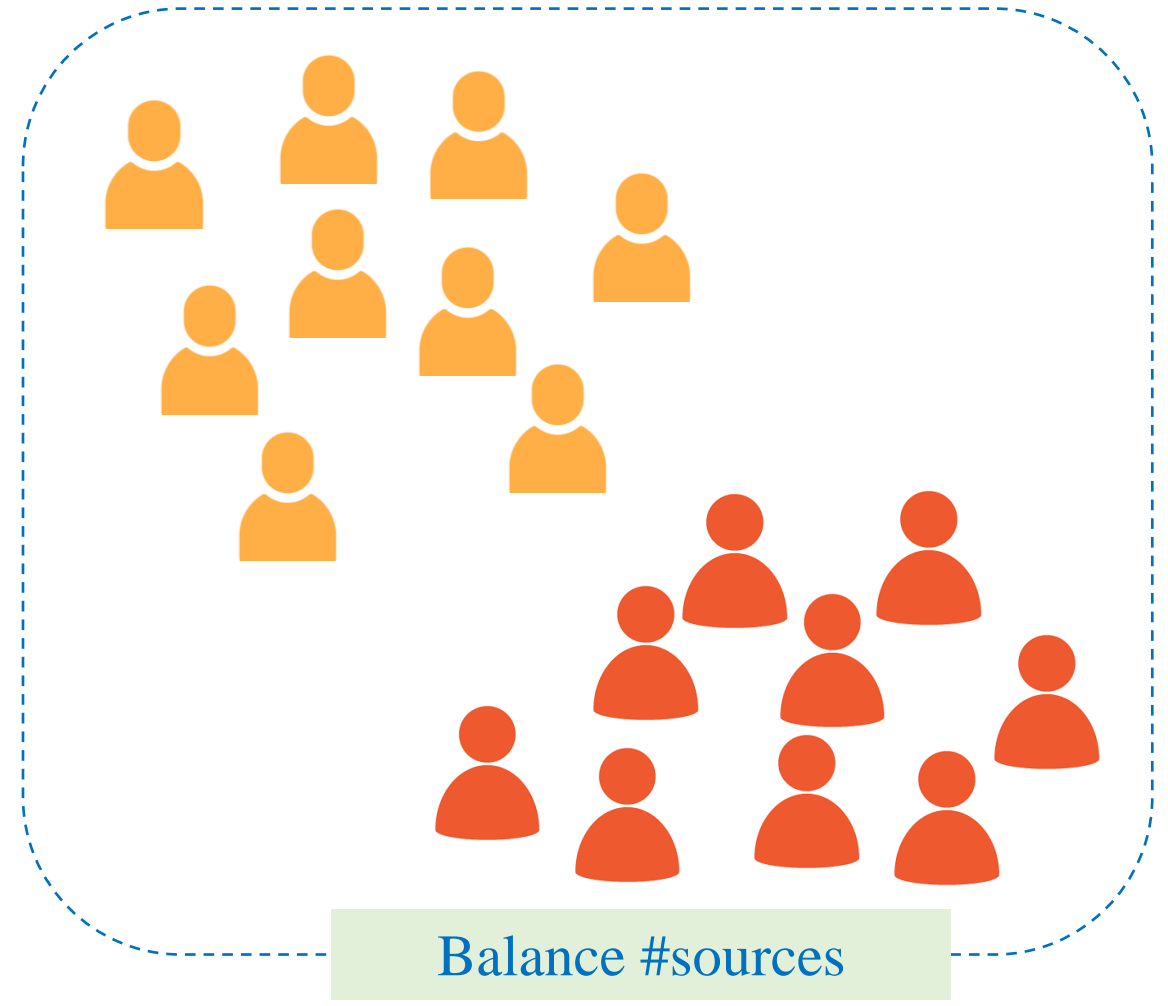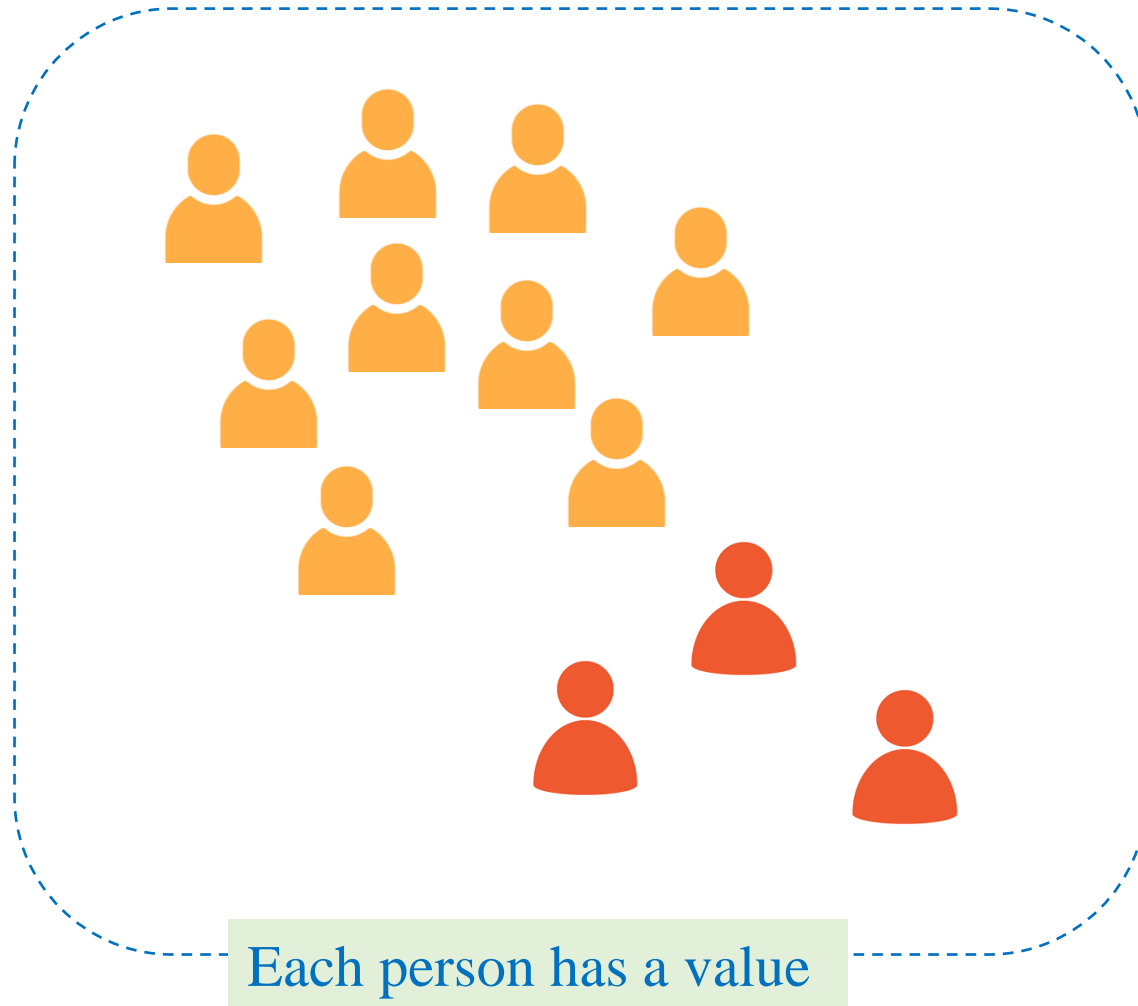
# AdaBoost

❖ **Discussion**

1) Are a wide range of features used to build a forest?
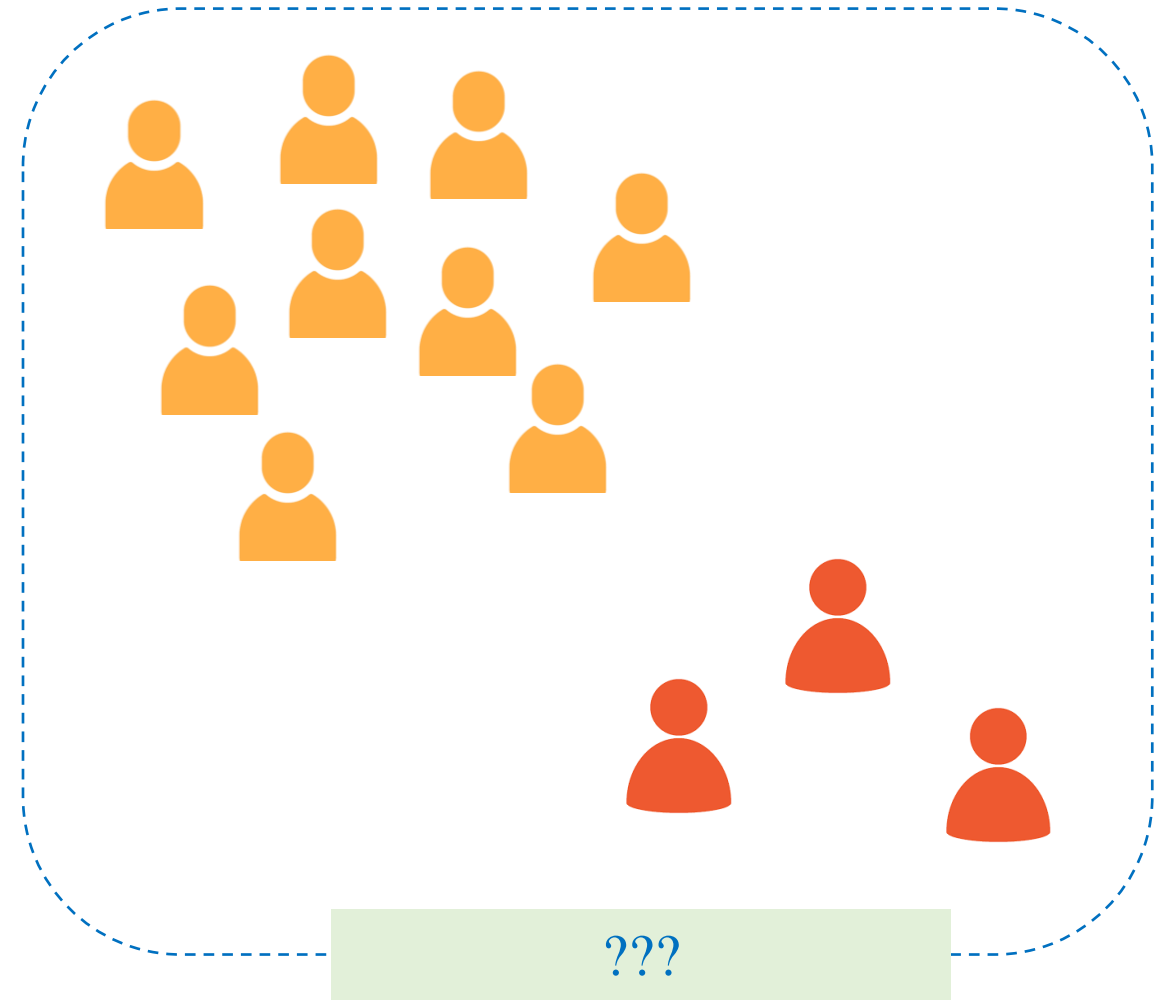
2) How to create a new dataset?

# AdaBoost

❖ **Idea**

How to balance the two groups' values?



Each person has a value

Balance #sources

# AdaBoost

❖ **Idea**

Any other ideas?



Each person has a value

???

# AdaBoost

❖ Create a new dataset

| Petal_Length | Petal_Width | Label |
|---|---|---|
| 1 | 0.2 | 0 |
| 1.3 | 0.6 | 0 |
| 0.9 | 0.7 | 0 |
| 1.7 | 0.5 | 1 |
| 1.8 | 0.9 | 1 |
| 1.2 | 1.3 | 1 |

| Petal_Length | Petal_Width | Label |
|---|---|---|
| 1 | 0.2 | 0 |
| 1.3 | 0.6 | 0 |
| 0.9 | 0.7 | 0 |
| 1.7 | 0.5 | 1 |
| 1.8 | 0.9 | 1 |
| 1.2 | 1.3 | 1 |

| Petal_Length | Petal_Width | Label | Evaluation |
|---|---|---|---|
| 1 | 0.2 | 0 | T |
| 1.3 | 0.6 | 0 | F |
| 0.9 | 0.7 | 0 | T |
| 1.7 | 0.5 | 1 | T |
| 1.8 | 0.9 | 1 | F |
| 1.2 | 1.3 | 1 | T |

| Petal_Length | Petal_Width | Label | Evaluation |
|---|---|---|---|
| 1 | 0.2 | 0 | T |
| 0.9 | 0.7 | 0 | T |
| 1.7 | 0.5 | 1 | T |
| 1.2 | 1.3 | 1 | T |

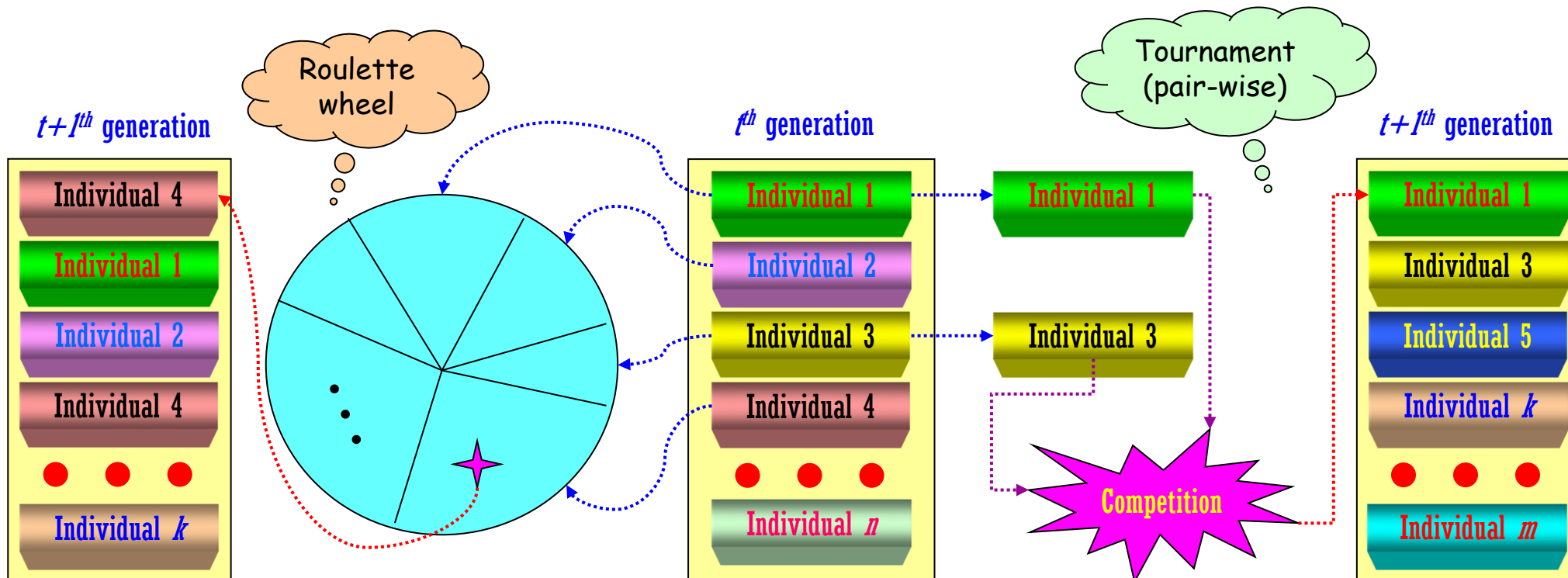| Petal_Length | Petal_Width | Label | Evaluation |
|---|---|---|---|
|  |  |  |  |
| 1.3 | 0.6 | 0 | F |
| 1.3 | 0.6 | 0 | F |
|  |  |  |  |
| 1.8 | 0.9 | 1 | F |
| 1.8 | 0.9 | 1 | F |

# Ideas from Genetic Algorithms

- **Roulette Wheel Selection**
  - ❖ The probability of selecting a given chromosome is proportional to its fitness
- **Tournament Selection**
  - ❖ Combine the fitness proportional concept with the random selection

# AdaBoost

❖ Create a new dataset

Add more randomness

| Petal_Length | Petal_Width | Label |
|---|---|---|
| 1 | 0.2 | 0 |
| 1.3 | 0.6 | 0 |
| 0.9 | 0.7 | 0 |
| 1.7 | 0.5 | 1 |
| 1.8 | 0.9 | 1 |
| 1.2 | 1.3 | 1 |

| Petal_Length | Petal_Width | Label |
|---|---|---|
| 1 | 0.2 | 0 |
| 1.3 | 0.6 | 0 |
| 0.9 | 0.7 | 0 |
| 1.7 | 0.5 | 1 |
| 1.8 | 0.9 | 1 |
| 1.2 | 1.3 | 1 |

| Petal_Length | Petal_Width | Label | Evaluation |
|---|---|---|---|
| 1 | 0.2 | 0 | T |
| 1.3 | 0.6 | 0 | F |
| 0.9 | 0.7 | 0 | T |
| 1.7 | 0.5 | 1 | T |
| 1.8 | 0.9 | 1 | F |
| 1.2 | 1.3 | 1 | T |

normalize

| Petal_Length | Petal_Width | Label | Evaluation | Score | Probability |
|---|---|---|---|---|---|
| 1 | 0.2 | 0 | T | 1 | 0.125 |
| 1.3 | 0.6 | 0 | F | 2 | 0.25 |
| 0.9 | 0.7 | 0 | T | 1 | 0.125 |
| 1.7 | 0.5 | 1 | T | 1 | 0.125 |
| 1.8 | 0.9 | 1 | F | 2 | 0.25 |
| 1.2 | 1.3 | 1 | T | 1 | 0.125 |

| Petal_Length | Petal_Width | Label |
|---|---|---|
| 1 | 0.2 | 0 |
| 1.3 | 0.6 | 0 |
| 1.2 | 1.3 | 1 |
| 1.7 | 0.5 | 1 |
| 1.8 | 0.9 | 1 |
| 1.8 | 0.9 | 1 |

40

# AdaBoost

❖ Create a new dataset

| Petal_Length | Petal_Width | Label |
|---|---|---|
| 1 | 0.2 | 0 |
| 1.3 | 0.6 | 0 |
| 0.9 | 0.7 | 0 |
| 1.7 | 0.5 | 1 |
| 1.8 | 0.9 | 1 |
| 1.2 | 1.3 | 1 |

| Petal_Length | Petal_Width | Label | Evaluation | Score | Probability |
|---|---|---|---|---|---|
| 1 | 0.2 | 0 | T | 1 | 0.125 |
| 1.3 | 0.6 | 0 | F | 2 | 0.25 |
| 0.9 | 0.7 | 0 | T | 1 | 0.125 |
| 1.7 | 0.5 | 1 | T | 1 | 0.125 |
| 1.8 | 0.9 | 1 | F | 2 | 0.25 |
| 1.2 | 1.3 | 1 | T | 1 | 0.125 |

Case 1

| Petal_Length | Petal_Width | Label |
|---|---|---|
| 1 | 0.2 | 0 |
| 1.3 | 0.6 | 0 |
| 1.2 | 1.3 | 1 |
| 1.7 | 0.5 | 1 |
| 1.8 | 0.9 | 1 |
| 1.8 | 0.9 | 1 |

| Petal_Length | Petal_Width | Label | Evaluation | Score | Probability |
|---|---|---|---|---|---|
| 1 | 0.2 | 0 | T | 1 | 0.142 |
| 1.3 | 0.6 | 0 | T | 1 | 0.142 |
| 1.2 | 1.3 | 1 | F | 2 | 0.29 |
| 1.7 | 0.5 | 1 | T | 1 | 0.142 |
| 1.8 | 0.9 | 1 | T | 1 | 0.142 |
| 1.8 | 0.9 | 1 | T | 1 | 0.142 |

# AdaBoost

❖ Create a new dataset

| Petal_Length | Petal_Width | Label |
|---|---|---|
| 1 | 0.2 | 0 |
| 1.3 | 0.6 | 0 |
| 0.9 | 0.7 | 0 |
| 1.7 | 0.5 | 1 |
| 1.8 | 0.9 | 1 |
| 1.2 | 1.3 | 1 |

| Petal_Length | Petal_Width | Label | Evaluation | Score | Probability |
|---|---|---|---|---|---|
| 1 | 0.2 | 0 | T | 1 | 0.125 |
| 1.3 | 0.6 | 0 | F | 2 | 0.25 |
| 0.9 | 0.7 | 0 | T | 1 | 0.125 |
| 1.7 | 0.5 | 1 | T | 1 | 0.125 |
| 1.8 | 0.9 | 1 | F | 2 | 0.25 |
| 1.2 | 1.3 | 1 | T | 1 | 0.125 |

Case 2

| Petal_Length | Petal_Width | Label |
|---|---|---|
| 1 | 0.2 | 0 |
| 1.3 | 0.6 | 0 |
| 1.2 | 1.3 | 1 |
| 1.7 | 0.5 | 1 |
| 1.8 | 0.9 | 1 |
| 1.8 | 0.9 | 1 |

| Petal_Length | Petal_Width | Label | Evaluation | Score | Probability |
|---|---|---|---|---|---|
| 1 | 0.2 | 0 | T | 1 | 0.142 |
| 1.3 | 0.6 | 0 | F | 2 | 0.29 |
| 1.2 | 1.3 | 1 | T | 1 | 0.142 |
| 1.7 | 0.5 | 1 | T | 1 | 0.142 |
| 1.8 | 0.9 | 1 | T | 1 | 0.142 |
| 1.8 | 0.9 | 1 | T | 1 | 0.142 |

Problem and solution?

Error < 0.5

For incorrect samples

When a model is good (small error),
scaled weights should increase/decrease slightly/significantly?

$f(x) = x$

$f(x) = e^{-x}$

$f(x) = -x$

$f(x) = e^{x}$

*error x*

# AdaBoost

For incorrect samples

When a model has a small error, increase significantly



$$f(x) = \frac{1-x}{x}$$

$$g(x) = \sqrt{\frac{1-x}{x}}$$

*error x*

*error x*

# AdaBoost

For correct samples

Decrease significantly

$$k(x) = \frac{1}{f(x)} = \frac{x}{1-x}$$

$$h(x) = \frac{1}{g(x)} = \sqrt{\frac{x}{1-x}}$$

# ❖ Create a new dataset

| Petal_Length | Petal_Width | Label | Probability |
|---|---|---|---|
| 1 | 0.2 | 0 | 0.166 |
| 1.3 | 0.6 | 0 | 0.166 |
| 0.9 | 0.7 | 0 | 0.166 |
| 1.7 | 0.5 | 1 | 0.166 |
| 1.8 | 0.9 | 1 | 0.166 |
| 1.2 | 1.3 | 1 | 0.166 |

| Petal_Length | Petal_Width | Label | Probability |
|---|---|---|---|
| 1 | 0.2 | 0 | 0.166 |
| 1.3 | 0.6 | 0 | 0.166 |
| 0.9 | 0.7 | 0 | 0.166 |
| 1.7 | 0.5 | 1 | 0.166 |
| 1.8 | 0.9 | 1 | 0.166 |
| 1.2 | 1.3 | 1 | 0.166 |

■ True

■ False

$$g(E) = \sqrt{\frac{1-E}{E}}$$

$$p_i = p_i g(E)$$

$$= 0.166 * 1.41 = 2.347$$

Update

| Petal_Length | Petal_Width | Label | Probability |
|---|---|---|---|
| 1 | 0.2 | 0 | 1.17 |
| 1.3 | 0.6 | 0 | 2.347 |
| 0.9 | 0.7 | 0 | 1.17 |
| 1.7 | 0.5 | 1 | 1.17 |
| 1.8 | 0.9 | 1 | 2.347 |
| 1.2 | 1.3 | 1 | 1.17 |

$$h(E) = \sqrt{\frac{E}{1-E}}$$

$$p_i = p_i h(E)$$

$$= 0.166 * 0.707 = 1.17$$

46

# ❖ Create a new dataset

| Petal_Length | Petal_Width | Label | Probability |
|:---:|:---:|:---:|:---:|
| 1 | 0.2 | 0 | 0.166 |
| 1.3 | 0.6 | 0 | 0.166 |
| 0.9 | 0.7 | 0 | 0.166 |
| 1.7 | 0.5 | 1 | 0.166 |
| 1.8 | 0.9 | 1 | 0.166 |
| 1.2 | 1.3 | 1 | 0.166 |

| Petal_Length | Petal_Width | Label | Probability |
|:---:|:---:|:---:|:---:|
| 1 | 0.2 | 0 | 0.166 |
| 1.3 | 0.6 | 0 | 0.166 |
| 0.9 | 0.7 | 0 | 0.166 |
| 1.7 | 0.5 | 1 | 0.166 |
| 1.8 | 0.9 | 1 | 0.166 |
| 1.2 | 1.3 | 1 | 0.166 |

True

False

Update

$$g(E) = \sqrt{\frac{1-E}{E}}$$

$$p_i = p_i g(E)$$

$$= 0.166 * 1.41 = 2.347$$

$$h(E) = \sqrt{\frac{E}{1-E}}$$

$$p_i = p_i h(E)$$

$$= 0.166 * 0.707 = 1.17$$

Normalized

| Petal_Length | Petal_Width | Label | Probability |
|:---:|:---:|:---:|:---:|
| 1 | 0.2 | 0 | 0.124 |
| 1.3 | 0.6 | 0 | 0.25 |
| 0.9 | 0.7 | 0 | 0.124 |
| 1.7 | 0.5 | 1 | 0.124 |
| 1.8 | 0.9 | 1 | 0.25 |
| 1.2 | 1.3 | 1 | 0.124 |

47

❖ Create a new dataset

| Petal_Length | Petal_Width | Label | Probability |
|---|---|---|---|
| 1 | 0.2 | 0 | 0.166 |
| 1.3 | 0.6 | 0 | 0.166 |
| 0.9 | 0.7 | 0 | 0.166 |
| 1.7 | 0.5 | 1 | 0.166 |
| 1.8 | 0.9 | 1 | 0.166 |
| 1.2 | 1.3 | 1 | 0.166 |

| Petal_Length | Petal_Width | Label | Probability |
|---|---|---|---|
| 1 | 0.2 | 0 | 0.166 |
| 1.3 | 0.6 | 0 | 0.166 |
| 0.9 | 0.7 | 0 | 0.166 |
| 1.7 | 0.5 | 1 | 0.166 |
| 1.8 | 0.9 | 1 | 0.166 |
| 1.2 | 1.3 | 1 | 0.166 |

■ True

■ False

$$g(E) = \sqrt{\frac{1-E}{E}}$$

$$p_i = p_i g(E) = p_i e^{\ln(g(E))}$$

$$= p_i e^{\ln\left(\sqrt{\frac{1-E}{E}}\right)} = p_i e^{\frac{1}{2}\ln\left(\frac{1-E}{E}\right)}$$

Update

$$h(E) = \sqrt{\frac{E}{1-E}}$$

$$p_i = p_i h(E) = p_i e^{\ln(h(E))}$$

$$= p_i e^{\ln\left(\sqrt{\frac{E}{1-E}}\right)} = p_i e^{-\frac{1}{2}\ln\left(\frac{1-E}{E}\right)}$$

Normalized

| Petal_Length | Petal_Width | Label | Probability |
|---|---|---|---|
| 1 | 0.2 | 0 | 0.124 |
| 1.3 | 0.6 | 0 | 0.25 |
| 0.9 | 0.7 | 0 | 0.124 |
| 1.7 | 0.5 | 1 | 0.124 |
| 1.8 | 0.9 | 1 | 0.25 |
| 1.2 | 1.3 | 1 | 0.124 |

# AdaBoost

For incorrect samples

Increase significantly



$$g(x) = \sqrt{\frac{1-x}{x}}$$

$$g(x) = e^{\frac{1}{2}\ln\left(\frac{1-x}{x}\right)}$$

error $x$

error $x$

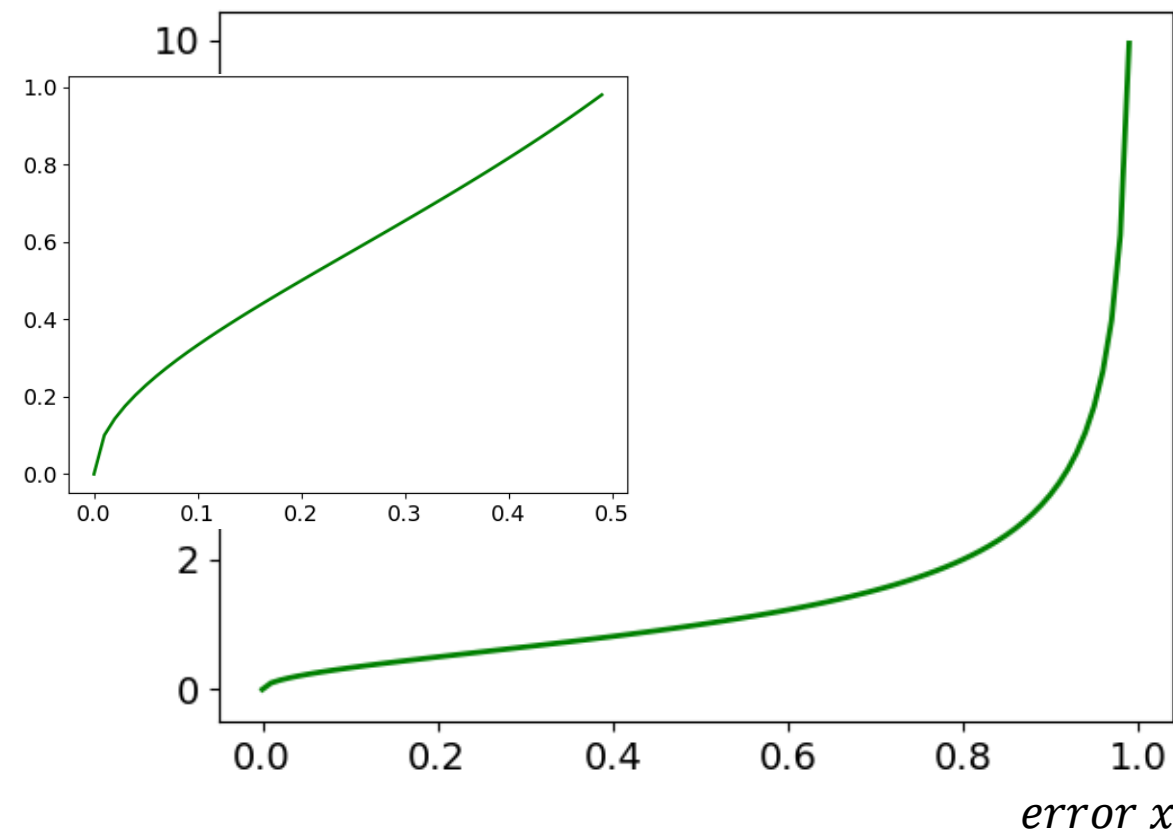https://mbernste.github.io/files/notes/AdaBoost.pdf

49

# AdaBoost

For correct samples

Decrease significantly

$$h(x) = \frac{1}{g(x)} = \sqrt{\frac{x}{1-x}}$$

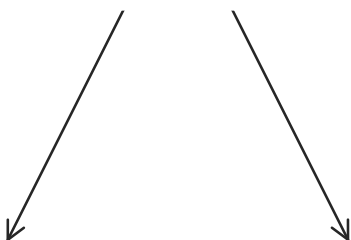$$h(x) = e^{-\frac{1}{2}\ln\left(\frac{1-x}{x}\right)}$$

# AdaBoost

**Implementation using sklearn**

| Petal_Length | Petal_Width | Label |
|--------------|-------------|-------|
| 1 | 0.2 | 0 |
| 1.3 | 0.6 | 0 |
| 0.9 | 0.7 | 0 |
| 1.7 | 0.5 | 1 |
| 1.8 | 0.9 | 1 |
| 1.2 | 1.3 | 1 |

```
dt_classifier = AdaBoostClassifier(n_estimators=3)
dt_classifier.fit(x_data, y_train)
```

Petal_Length <= 1.5
gini = 0.5
samples = 6
value = [0.5, 0.5]

Petal_Width <= 0.8
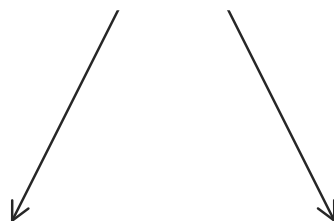gini = 0.5
samples = 6
value = [0.5, 0.5]

Petal_Length <= 1.5
gini = 0.5
samples = 6
value = [0.5, 0.5]

gini = 0.375
samples = 4
value = [0.5, 0.167]

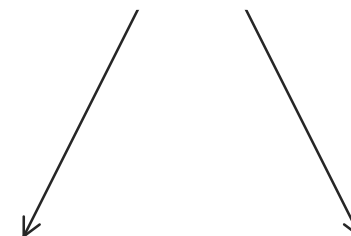gini = 0.0
samples = 2
value = [0.0, 0.333]

gini = 0.0
samples = 4
value = [0.5, 0.0]

gini = 0.0
samples = 2
value = [0.0, 0.5]

gini = 0.0
samples = 4
value = [0.5, 0.0]

gini = 0.0
samples = 2
value = [0.0, 0.5]