# DATA SCIENCE COURSE
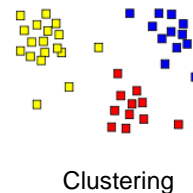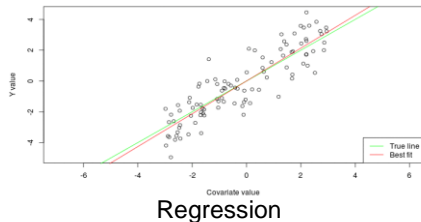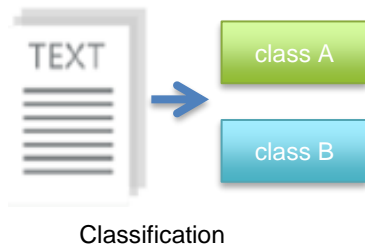
MACHINE LEARNING  DAY 03

# Machine Learning Types

*Machine Learning Basics*

- *Supervised*: learning with labeled data
  - Example: email classification, image classification
  - Example: regression for predicting real-valued outputs
- *Unsupervised*: discover patterns in unlabeled data
  - Example: cluster similar data points
- *Reinforcement learning*: learn to act based on feedback/reward
  - Example: learn to play Go

Classification

Regression

Clustering

# Decision Tree and Random Forest
## Let's play a Game!

# Decision Trees

| Gender | Age | App |
|--------|-----|-----|
| F | 15 | (Pokémon) |
| F | 25 | (WhatsApp) |
| M | 32 | (Snapchat) |
| F | 40 | (WhatsApp) |
| M | 12 | (Pokémon) |
| M | 14 | (Pokémon) |

Quiz 1

Answer

| Gender | Age | App |
|--------|-----|-----|
| F | 15 | Pokémon |
| F | 25 | WhatsApp |
| M | 32 | Snapchat |
| F | 40 | WhatsApp |
| M | 12 | Pokémon |
| M | 14 | Pokémon |

| Gender | Age | App |
|--------|-----|-----|
| F | 15 | Pokémon |
| F | 25 | WhatsApp |
| M | 32 | Snapchat |
| F | 40 | WhatsApp |
| M | 12 | Pokémon |
| M | 14 | Pokémon |

| Gender | Age | App |
|--------|-----|-----|
| F | 15 |  |
| F | 25 |  |
| M | 32 |  |
| F | 40 |  |
| M | 12 |  |
| M | 14 |  |

Age

<20    ≥20

| Gender | Age | App |
| --- | --- | --- |
| F | 25 | WhatsApp |
| M | 32 | Snapchat |
| F | 40 | WhatsApp |

Age
<20 / ≥20
Pokémon (<20)
Gender (≥20)
F / M
WhatsApp (F) / Snapchat (M)

# Student Admissions

# Definition

- A tree-like model that illustrates series of events leading to certain decisions
- Each node represents a test on an attribute and each branch is an outcome of that test

Who to loan?

- Not a student
- 45 years old
- Medium income
- Fair credit record

- Student
- 27 years old
- Low income
- Excellent credit record

# Definition

- A tree-like model that illustrates series of events leading to certain decisions
- Each node represents a test on an attribute and each branch is an outcome of that test

## Who to loan?

- Not a student
- 45 years old
- Medium income
- Fair credit record
- ➢ Yes

- Student
- 27 years old
- Low income
- Excellent credit record

# Definition

- A tree-like model that illustrates series of events leading to certain decisions
- Each node represents a test on an attribute and each branch is an outcome of that test

## Who to loan?

- Not a student
- 45 years old
- Medium income
- Fair credit record
- ➢ Yes

- Student
- 27 years old
- Low income
- Excellent credit record
- ➢ No

# Decision Tree Learning

- We use labeled data to obtain a suitable decision tree for future predictions
  - We want a decision tree that works well on unseen data, while asking as few questions as possible

| Outlook | Temperature | Humidity | Wind | Play Tennis? |
|---------|-------------|----------|------|--------------|
| Sunny | Hot | High | Weak | No |
| Sunny | Hot | High | Strong | No |
| Overcast | Hot | High | Weak | Yes |
| Rainy | Mild | High | Weak | Yes |
| Rainy | Cool | Normal | Weak | Yes |
| Rainy | Cool | Normal | Strong | No |
| Overcast | Cool | Normal | Strong | Yes |
| Sunny | Mild | High | Weak | No |
| Sunny | Cool | Normal | Weak | Yes |
| Rainy | Mild | Normal | Weak | Yes |
| Sunny | Mild | Normal | Strong | Yes |
| Overcast | Mild | High | Strong | Yes |
| Overcast | Hot | Normal | Weak | Yes |
| Rainy | Mild | High | Strong | No |

# Decision Tree Learning

- Basic step: choose an attribute and, based on its values, split the data into smaller sets
  - ➢ Recursively repeat this step until we can surely decide the label

| Outlook | Temperature | Humidity | Wind | Play Tennis? |
|---|---|---|---|---|
| Sunny | Hot | High | Weak | No |
| Sunny | Hot | High | Strong | No |
| Overcast | Hot | High | Weak | Yes |
| Rainy | Mild | High | Weak | Yes |
| Rainy | Cool | Normal | Weak | Yes |
| Rainy | Cool | Normal | Strong | No |
| Overcast | Cool | Normal | Strong | Yes |
| Sunny | Mild | High | Weak | No |
| Sunny | Cool | Normal | Weak | Yes |
| Rainy | Mild | Normal | Weak | Yes |
| Sunny | Mild | Normal | Strong | Yes |
| Overcast | Mild | High | Strong | Yes |
| Overcast | Hot | Normal | Weak | Yes |
| Rainy | Mild | High | Strong | No |

**Outlook**

# Decision Tree Learning

- Basic step: choose an attribute and, based on its values, split the data into smaller sets
  - ➢ Recursively repeat this step until we can surely decide the label

**Outlook = Sunny**

| Temperature | Humidity | Wind | Play Tennis? |
|---|---|---|---|
| Hot | High | Weak | No |
| Hot | High | Strong | No |
| Mild | High | Weak | No |
| Cool | Normal | Weak | Yes |
| Mild | Normal | Strong | Yes |

**Outlook = Overcast**

| Temperature | Humidity | Wind | Play Tennis? |
|---|---|---|---|
| Hot | High | Weak | Yes |
| Cool | Normal | Strong | Yes |
| Mild | High | Strong | Yes |
| Hot | Normal | Weak | Yes |

**Outlook = Rainy**

| Temperature | Humidity | Wind | Play Tennis? |
|---|---|---|---|
| Mild | High | Weak | Yes |
| Cool | Normal | Weak | Yes |
| Cool | Normal | Strong | No |
| Mild | Normal | Weak | Yes |
| Mild | High | Strong | No |

**Outlook**

Sunny    Overcast    Rainy

?    ?    ?

# Decision Tree Learning

- Basic step: choose an attribute and, based on its values, split the data into smaller sets
  - ➤ Recursively repeat this step until we can surely decide the label

**Outlook = Sunny**

| Temperature | Humidity | Wind | Play Tennis? |
|---|---|---|---|
| Hot | High | Weak | No |
| Hot | High | Strong | No |
| Mild | High | Weak | No |
| Cool | Normal | Weak | Yes |
| Mild | Normal | Strong | Yes |

**Outlook = Overcast**

| Temperature | Humidity | Wind | Play Tennis? |
|---|---|---|---|
| Hot | High | Weak | Yes |
| Cool | Normal | Strong | Yes |
| Mild | High | Strong | Yes |
| Hot | Normal | Weak | Yes |

**Outlook = Rainy**

| Temperature | Humidity | Wind | Play Tennis? |
|---|---|---|---|
| Mild | High | Weak | Yes |
| Cool | Normal | Weak | Yes |
| Cool | Normal | Strong | No |
| Mild | Normal | Weak | Yes |
| Mild | High | Strong | No |

# Decision Tree Learning

- Basic step: choose an attribute and, based on its values, split the data into smaller sets
  - ➢ Recursively repeat this step until we can surely decide the label
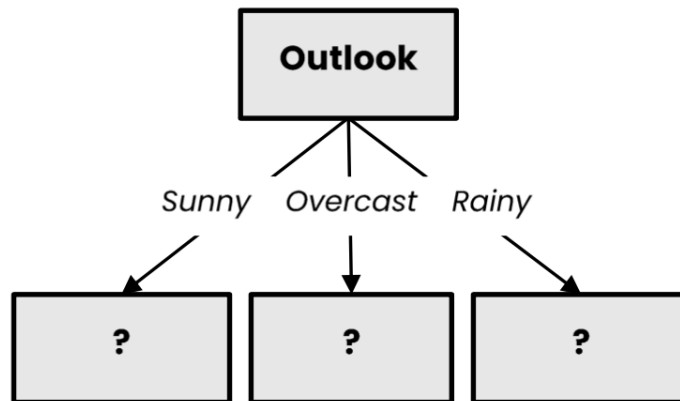
**Outlook = Sunny**

**Humidity = High**

| Temperature | Wind | Play Tennis? |
|---|---|---|
| Hot | Weak | No |
| Hot | Strong | No |
| Mild | Weak | No |

**Humidity = Normal**

| Temperature | Wind | Play Tennis? |
|---|---|---|
| Cool | Weak | Yes |
| Mild | Strong | Yes |

**Outlook = Overcast**

| Temperature | Humidity | Wind | Play Tennis? |
|---|---|---|---|
| Hot | High | Weak | Yes |
| Cool | Normal | Strong | Yes |
| Mild | High | Strong | Yes |
| Hot | Normal | Weak | Yes |

**Outlook = Rainy**

**Wind = Strong**

| Temperature | Humidity | Play Tennis? |
|---|---|---|
| Cool | Normal | No |
| Mild | High | No |

**Wind = Weak**

| Temperature | Humidity | Play Tennis? |
|---|---|---|
| Mild | High | Yes |
| Cool | Normal | Yes |
| Mild | Normal | Yes |

Tree:

Outlook
- Sunny → Humidity
  - High → ?
  - Normal → ?
- Overcast → Yes
- Rainy → Wind
  - Strong → ?
  - Weak → ?

# Decision Tree Learning

- Basic step: choose an attribute and, based on its values, split the data into smaller sets
  - ➢ Recursively repeat this step until we can surely decide the label
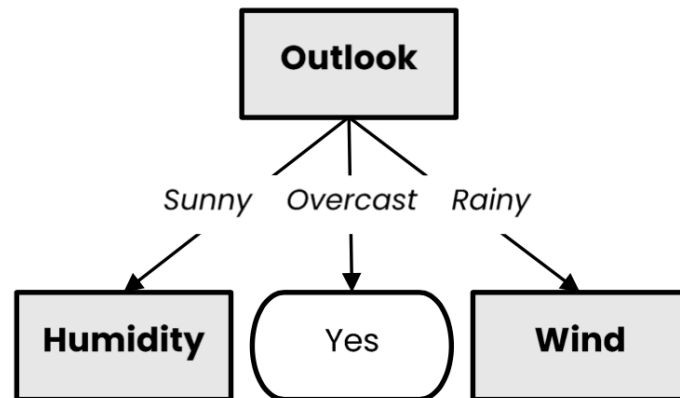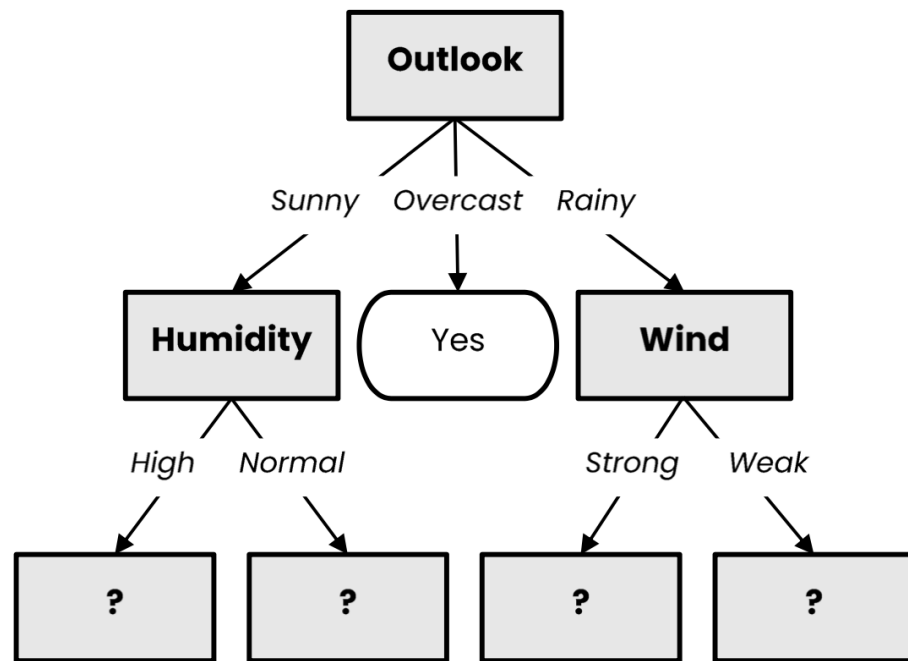
**Outlook = Sunny**

**Humidity = High**

| Temperature | Wind | Play Tennis? |
|---|---|---|
| Hot | Weak | No |
| Hot | Strong | No |
| Mild | Weak | No |

**Humidity = Normal**

| Temperature | Wind | Play Tennis? |
|---|---|---|
| Cool | Weak | Yes |
| Mild | Strong | Yes |

**Outlook = Overcast**

| Temperature | Humidity | Wind | Play Tennis? |
|---|---|---|---|
| Hot | High | Weak | Yes |
| Cool | Normal | Strong | Yes |
| Mild | High | Strong | Yes |
| Hot | Normal | Weak | Yes |

**Outlook = Rainy**

**Wind = Strong**

| Temperature | Humidity | Play Tennis? |
|---|---|---|
| Cool | Normal | No |
| Mild | High | No |

**Wind = Weak**

| Temperature | Humidity | Play Tennis? |
|---|---|---|
| Mild | High | Yes |
| Cool | Normal | Yes |
| Mild | Normal | Yes |

# What is a good attribute?

| Does it fly? | Color | Class |
|---|---|---|
| No | Brown | **Mammal** |
| No | White | **Mammal** |
| Yes | Brown | **Bird** |
| Yes | White | **Bird** |
| No | White | **Mammal** |
| No | Brown | **Bird** |
| Yes | White | **Bird** |

**Color**

Brown    White

33% Mammal
66% Bird

50% Mammal
50% Bird

**Fly?**

Yes    No

100% Bird

75% Mammal
25% Bird

- Which attribute provides better splitting?
- Why?
  - ➤ Because the resulting subsets are more pure
  - ➤ Knowing the value of this attribute gives us more information about the label
    (the entropy of the subsets is lower)

# Information Gain

# Entropy

- Entropy measures the degree of randomness in data



**Low entropy**    **High entropy**

- For a set of samples $X$ with $k$ classes:

$$entropy(X) = -\sum_{i=1}^{k} p_i \log_2(p_i)$$

where $p_i$ is the proportion of elements of class $i$

- Lower entropy implies greater predictability!

# Information Gain

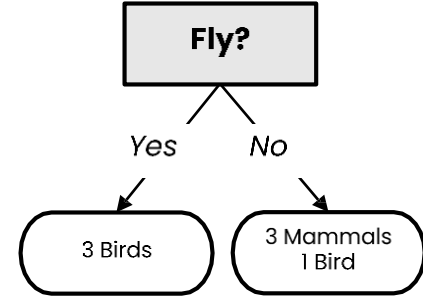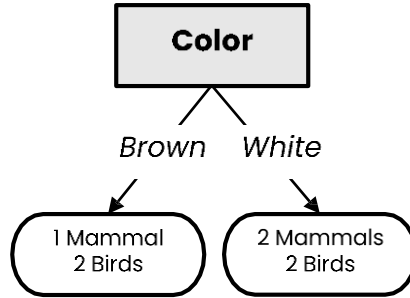- The information gain of an attribute a is the expected reduction in entropy due to splitting on values of a:

$$gain(X, a) = entropy(X) - \sum_{v \in Values(a)} \frac{|X_v|}{|X|} entropy(X_v)$$

where $X_v$ is the subset of $X$ for which $a = v$
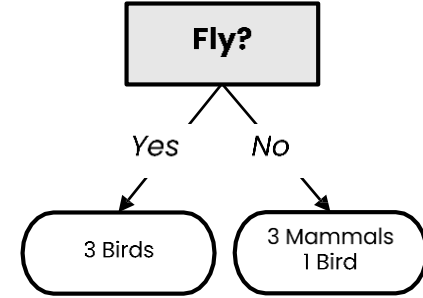
# Best attribute = highest information gain

| Does it fly? | Color | Class |
|---|---|---|
| No | Brown | Mammal |
| No | White | Mammal |
| Yes | Brown | Bird |
| Yes | White | Bird |
| No | White | Mammal |
| No | Brown | Bird |
| Yes | White | Bird |

**Color**

Brown    White

1 Mammal
2 Birds

2 Mammals
2 Birds

**Fly?**

Yes    No

3 Birds

3 Mammals
1 Bird

$$entropy\,(X) = -p_{\mathrm{mammal}} \log_2 p_{\mathrm{mammal}} - p_{\mathrm{bird}} \log_2 p_{\mathrm{bird}} = -\frac{3}{7}\log_2\frac{3}{7} - \frac{4}{7}\log_2\frac{4}{7} \approx 0.985$$

# Best attribute = highest information gain

| Does it fly? | Color | Class |
|---|---|---|
| No | Brown | Mammal |
| No | White | Mammal |
| Yes | Brown | Bird |
| Yes | White | Bird |
| No | White | Mammal |
| No | Brown | Bird |
| Yes | White | Bird |

**Color**

Brown / White

1 Mammal 2 Birds | 2 Mammals 2 Birds
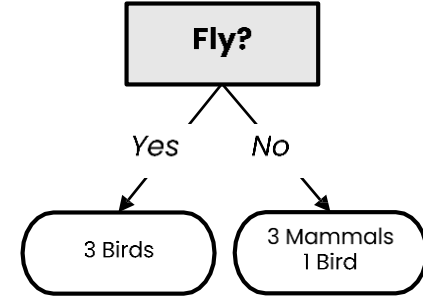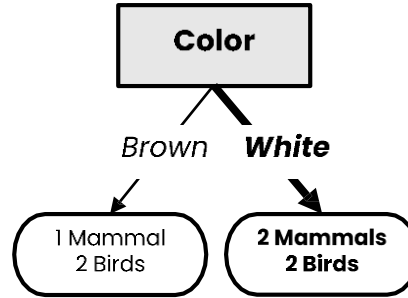
**Fly?**

Yes / No

3 Birds | 3 Mammals 1 Bird

$$entropy\,(X) = -p_{\text{mammal}} \log_2 p_{\text{mammal}} - p_{\text{bird}} \log_2 p_{\text{bird}} = -\frac{3}{7} \log_2 \frac{3}{7} - \frac{4}{7} \log_2 \frac{4}{7} \approx 0.985$$

$$entropy\,(X_{color=brown}) = -\frac{1}{3} \log_2 \frac{1}{3} - \frac{2}{3} \log_2 \frac{2}{3} \approx 0.918$$

# Best attribute = highest information gain

| Does it fly? | Color | Class |
|---|---|---|
| No | Brown | Mammal |
| No | White | Mammal |
| Yes | Brown | Bird |
| Yes | White | Bird |
| No | White | Mammal |
| No | Brown | Bird |
| Yes | White | Bird |

**Color**

Brown    **White**

```
1 Mammal      2 Mammals
2 Birds       2 Birds
```

**Fly?**

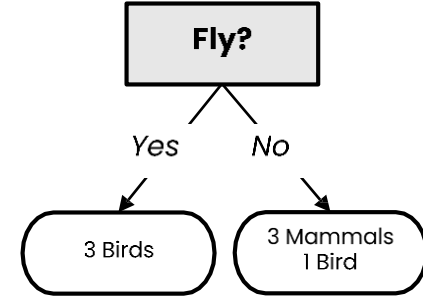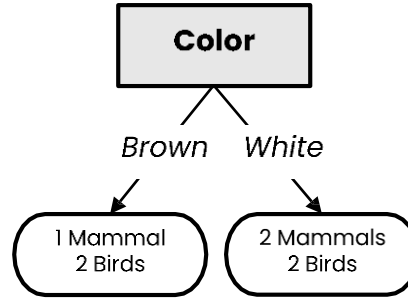Yes    No

```
3 Birds       3 Mammals
              1 Bird
```

$$entropy\ (X) = -p_{\text{mammal}} \log_2 p_{\text{mammal}} - p_{\text{bird}} \log_2 p_{\text{bird}} = -\frac{3}{7} \log_2 \frac{3}{7} - \frac{4}{7} \log_2 \frac{4}{7} \approx 0.985$$

$$entropy\ (X_{color=brown}) = -\frac{1}{3} \log_2 \frac{1}{3} - \frac{2}{3} \log_2 \frac{2}{3} \approx 0.918 \qquad entropy\ (X_{color=white}) = 1$$

# Best attribute = highest information gain

| Does it fly? | Color | Class |
|---|---|---|
| No | Brown | Mammal |
| No | White | Mammal |
| Yes | Brown | Bird |
| Yes | White | Bird |
| No | White | Mammal |
| No | Brown | Bird |
| Yes | White | Bird |

**Color**

Brown    White

1 Mammal
2 Birds

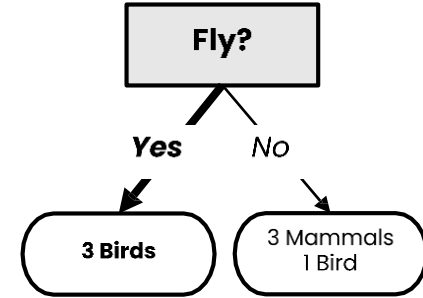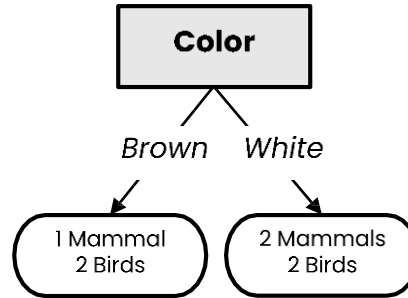2 Mammals
2 Birds

**Fly?**

Yes    No

3 Birds

3 Mammals
1 Bird

$$entropy\,(X) = -\,p_{\text{mammal}} \log_2 p_{\text{mammal}} - p_{\text{bird}} \log_2 p_{\text{bird}} = -\frac{3}{7}\log_2\frac{3}{7} - \frac{4}{7}\log_2\frac{4}{7} \approx 0.985$$

$$entropy\,(X_{color=brown}) = -\frac{1}{3}\log_2\frac{1}{3} - \frac{2}{3}\log_2\frac{2}{3} \approx 0.918 \qquad entropy\,(X_{color=white}) = 1$$

$$\boldsymbol{gain\,(X, color) = 0.985 - \frac{3}{7} \cdot 0.918 - \frac{4}{7} \cdot 1 \approx 0.020}$$

# Best attribute = highest information gain

| Does it fly? | Color | Class |
|---|---|---|
| No | Brown | Mammal |
| No | White | Mammal |
| Yes | Brown | Bird |
| Yes | White | Bird |
| No | White | Mammal |
| No | Brown | Bird |
| Yes | White | Bird |

**Color**

Brown / White

1 Mammal 2 Birds

2 Mammals 2 Birds

**Fly?**

**Yes** / No

**3 Birds**

3 Mammals 1 Bird

$$entropy\,(X) = -p_{\mathrm{mammal}} \log_2 p_{\mathrm{mammal}} - p_{\mathrm{bird}} \log_2 p_{\mathrm{bird}} = -\frac{3}{7} \log_2 \frac{3}{7} - \frac{4}{7} \log_2 \frac{4}{7} \approx 0.985$$

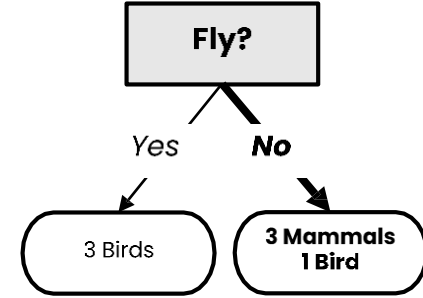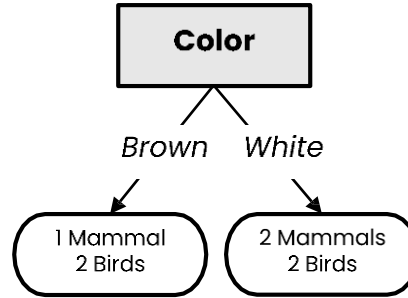$$entropy\,(X_{color=brown}) = -\frac{1}{3} \log_2 \frac{1}{3} - \frac{2}{3} \log_2 \frac{2}{3} \approx 0.918 \qquad entropy\,(X_{color=white}) = 1$$

$$\boldsymbol{gain\,(X, color) = 0.985 - \frac{3}{7} \cdot 0.918 - \frac{4}{7} \cdot 1 \approx 0.020}$$

$$entropy\,(X_{fly=yes}) = 0$$

# Best attribute = highest information gain

| Does it fly? | Color | Class |
|---|---|---|
| No | Brown | Mammal |
| No | White | Mammal |
| Yes | Brown | Bird |
| Yes | White | Bird |
| No | White | Mammal |
| No | Brown | Bird |
| Yes | White | Bird |

**Color**

Brown / White

1 Mammal 2 Birds

2 Mammals 2 Birds

**Fly?**

Yes / *No*

3 Birds

**3 Mammals 1 Bird**

$$entropy\,(X) = -p_{\text{mammal}} \log_2 p_{\text{mammal}} - p_{\text{bird}} \log_2 p_{\text{bird}} = -\frac{3}{7} \log_2 \frac{3}{7} - \frac{4}{7} \log_2 \frac{4}{7} \approx 0.985$$

$$entropy\,(X_{color=brown}) = -\frac{1}{3} \log_2 \frac{1}{3} - \frac{2}{3} \log_2 \frac{2}{3} \approx 0.918 \qquad entropy\,(X_{color=white}) = 1$$
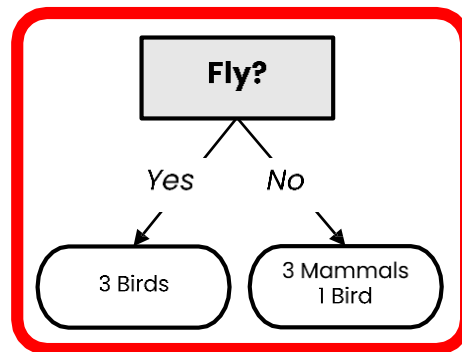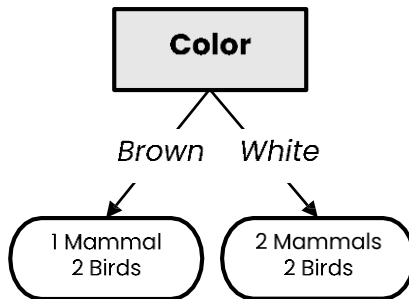
$$\boldsymbol{gain\,(X, color) = 0.985 - \frac{3}{7} \cdot 0.918 - \frac{4}{7} \cdot 1 \approx 0.020}$$

$$entropy\,(X_{fly=yes}) = 0 \qquad\qquad entropy\,(X_{fly=no}) = -\frac{3}{4} \log_2 \frac{3}{4} - \frac{1}{4} \log_2 \frac{1}{4} \approx 0.811$$

# Best attribute = highest information gain

In practice, we compute $entropy(X)$ only once!

| Does it fly? | Color | Class |
|---|---|---|
| No | Brown | Mammal |
| No | White | Mammal |
| Yes | Brown | Bird |
| Yes | White | Bird |
| No | White | Mammal |
| No | Brown | Bird |
| Yes | White | Bird |



$$entropy\,(X) = -\,p_{\mathrm{mammal}} \log_2 p_{\mathrm{mammal}} - p_{\mathrm{bird}} \log_2 p_{\mathrm{bird}} = -\frac{3}{7}\,\log_2\frac{3}{7} - \frac{4}{7}\log_2\frac{4}{7} \approx 0.985$$

$$entropy\,(X_{color=brown}) = -\frac{1}{3}\,\log_2\frac{1}{3} - \frac{2}{3}\log_2\frac{2}{3} \approx 0.918 \qquad entropy\,(X_{color=white}) = 1$$

$$\boldsymbol{gain\,(X, color) = 0.985 - \frac{3}{7} \cdot 0.918 - \frac{4}{7} \cdot 1 \approx 0.020}$$

$$entropy\,(X_{fly=yes}) = 0 \qquad\qquad entropy\,(X_{fly=no}) = -\frac{3}{4}\,\log_2\frac{3}{4} - \frac{1}{4}\log_2\frac{1}{4} \approx 0.811$$

$$\boldsymbol{gain\,(X, fly) = 0.985 - \frac{3}{7} \cdot 0 - \frac{4}{7} \cdot 0.811 \approx \boxed{0.521}}$$

# Gini Impurity

# Gini Impurity

- Gini impurity measures how often a randomly chosen example would be incorrectly labeled if it was randomly labeled according to the label distribution



**Error of classifying randomly picked fruit with randomly picked label**
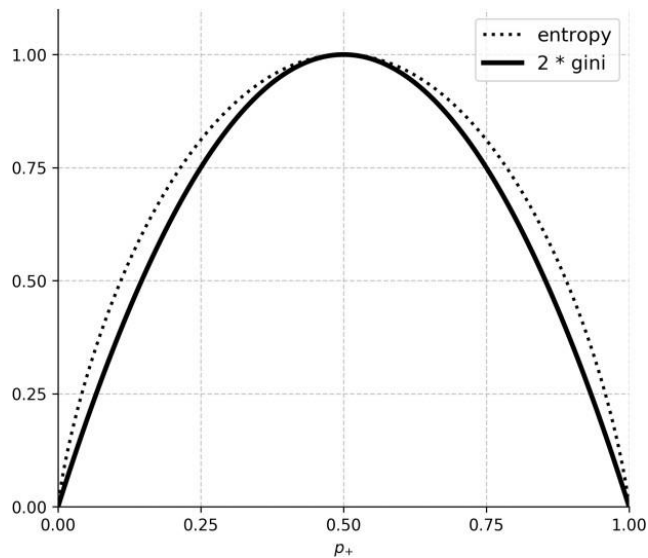
- For a set of samples $X$ with $k$ classes:

$$gini(X) = 1 - \sum_{i=1}^{k} p_i^2$$

where $p_i$ is the proportion of elements of class $i$

- Can be used as an alternative to entropy for selecting attributes!

# Entropy versus Gini Impurity

- Entropy and Gini Impurity give similar results in practice
  - ➢ They only disagree in about 2% of cases
    "Theoretical Comparison between the Gini Index and Information Gain Criteria" [Răileanu & Stoffel, AMAI 2004]
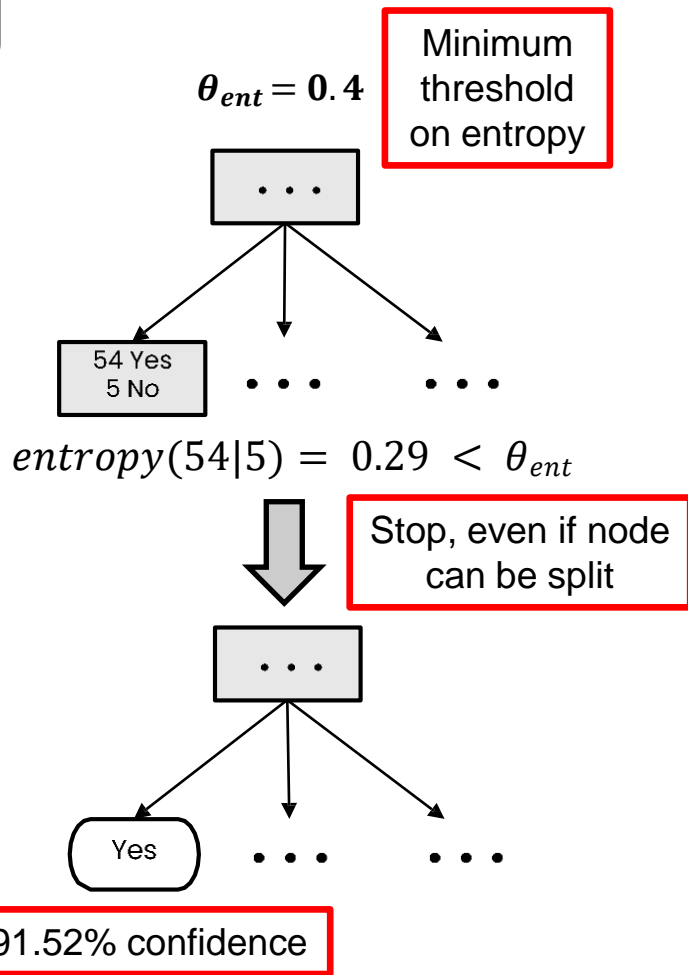  - ➢ Entropy might be slower to compute, because of the log

# Pruning

# Pruning

- Pruning is a technique that reduces the size of a decision tree by removing branches of the tree which provide little predictive power
- It is a **regularization** method that reduces the complexity of the final model, thus reducing overfitting
  - Decision trees are prone to overfitting!

- Pruning methods:
  - Pre-pruning: Stop the tree building algorithm before it fully classifies the data
  - Post-pruning: Build the complete tree, then replace some non-leaf nodes with leaf nodes if this improves validation error

# Pre-pruning

- Pre-pruning implies early stopping:
  - ➤ If some condition is met, the current node will not be split, even if it is not 100% pure
  - ➤ It will become a leaf node with the label of the majority class in the current set

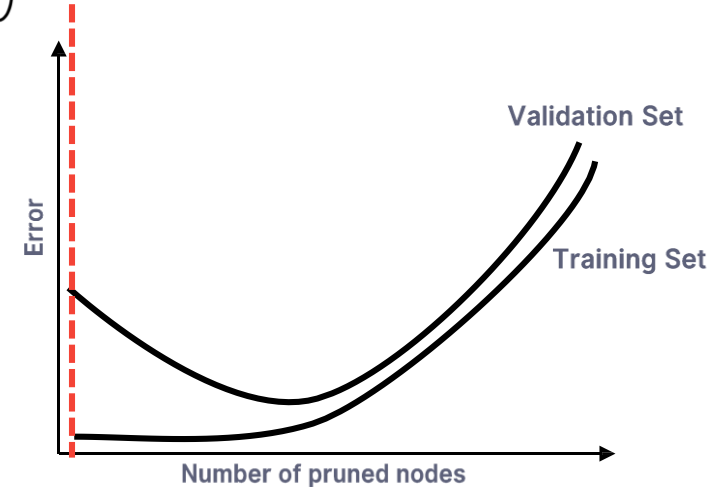(the class distribution could be used as prediction confidence)

- Common stopping criteria include setting a threshold on:
  - ➤ Entropy (or Gini Impurity) of the current set
  - ➤ Number of samples in the current set
  - ➤ Gain of the best-splitting attribute
  - ➤ Depth of the tree

$\theta_{ent} = 0.4$

Minimum threshold on entropy

. . .

54 Yes
5 No

. . .     . . .

$entropy(54|5) = 0.29 < \theta_{ent}$

Stop, even if node can be split

. . .

Yes     . . .     . . .
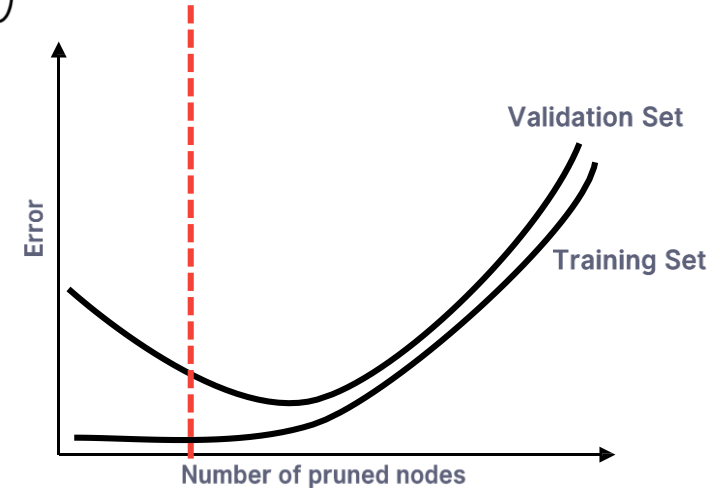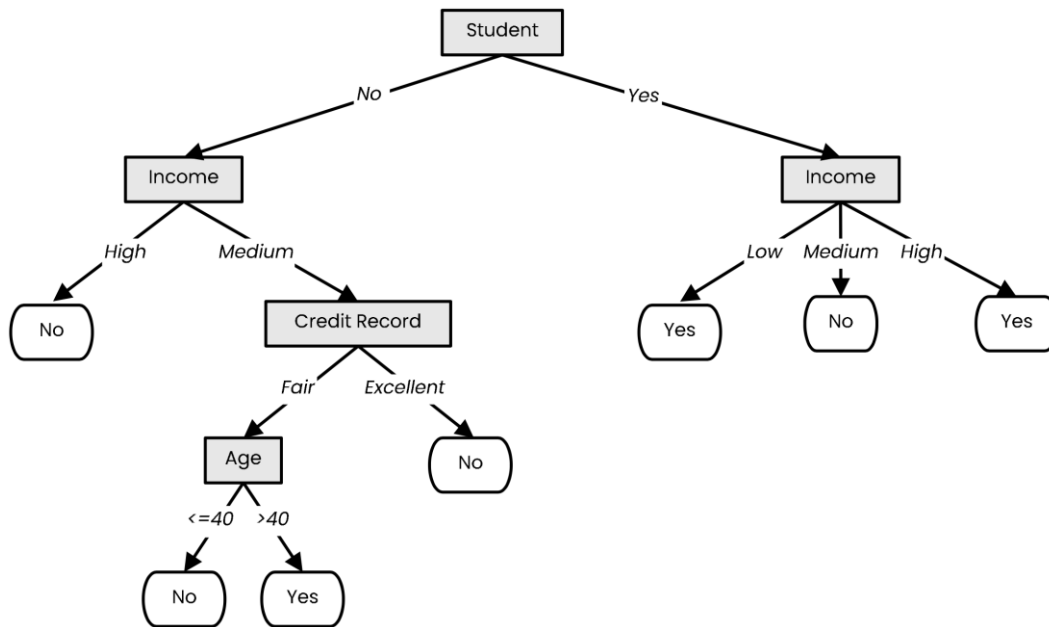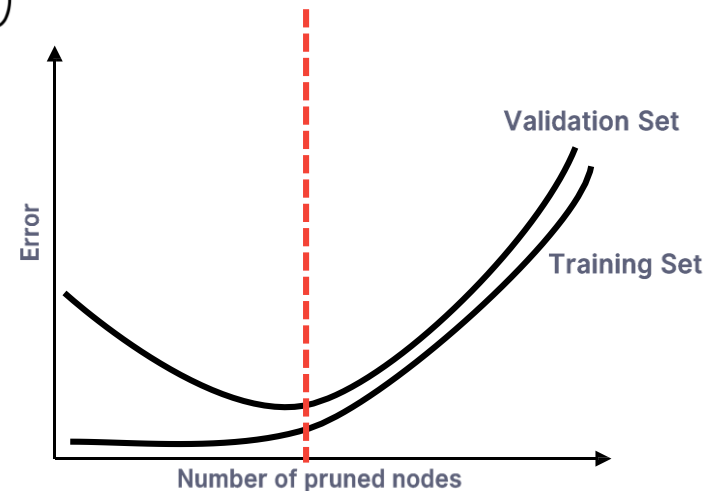
91.52% confidence

# Post-pruning



- Prune nodes in a bottom-up manner, if it decreases validation error

# Post-pruning



- Prune nodes in a bottom-up manner, if it decreases validation error
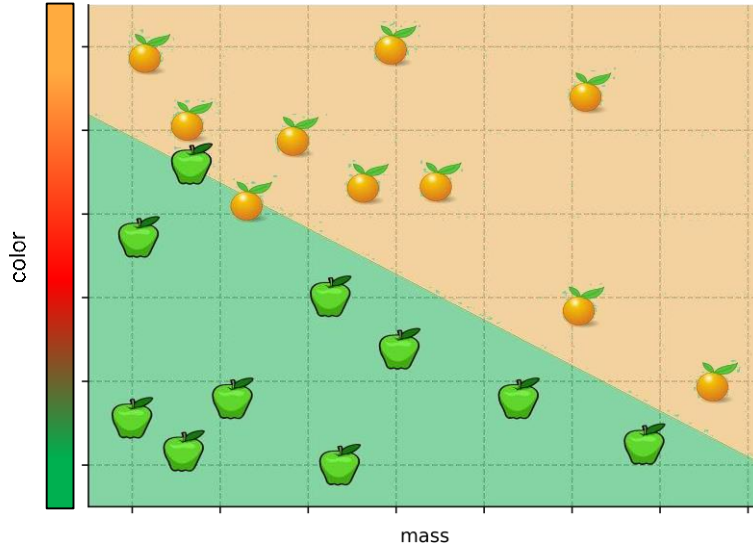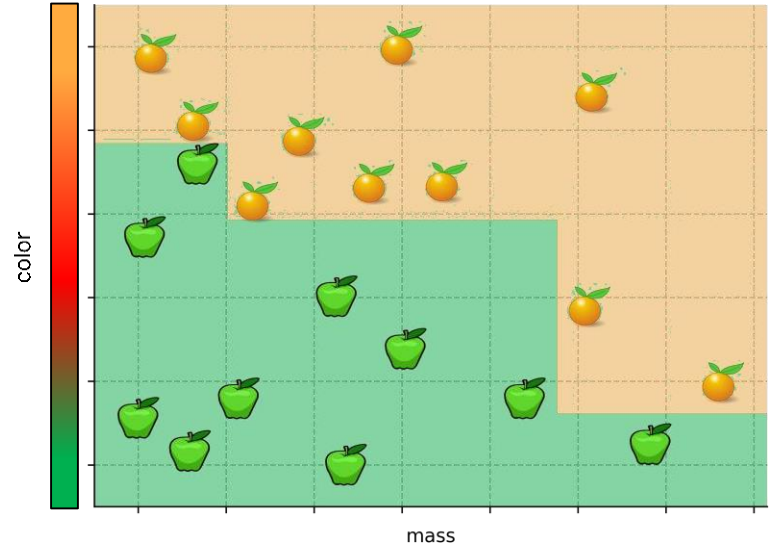
# Post-pruning



- Prune nodes in a bottom-up manner, if it decreases validation error

# Decision Boundaries

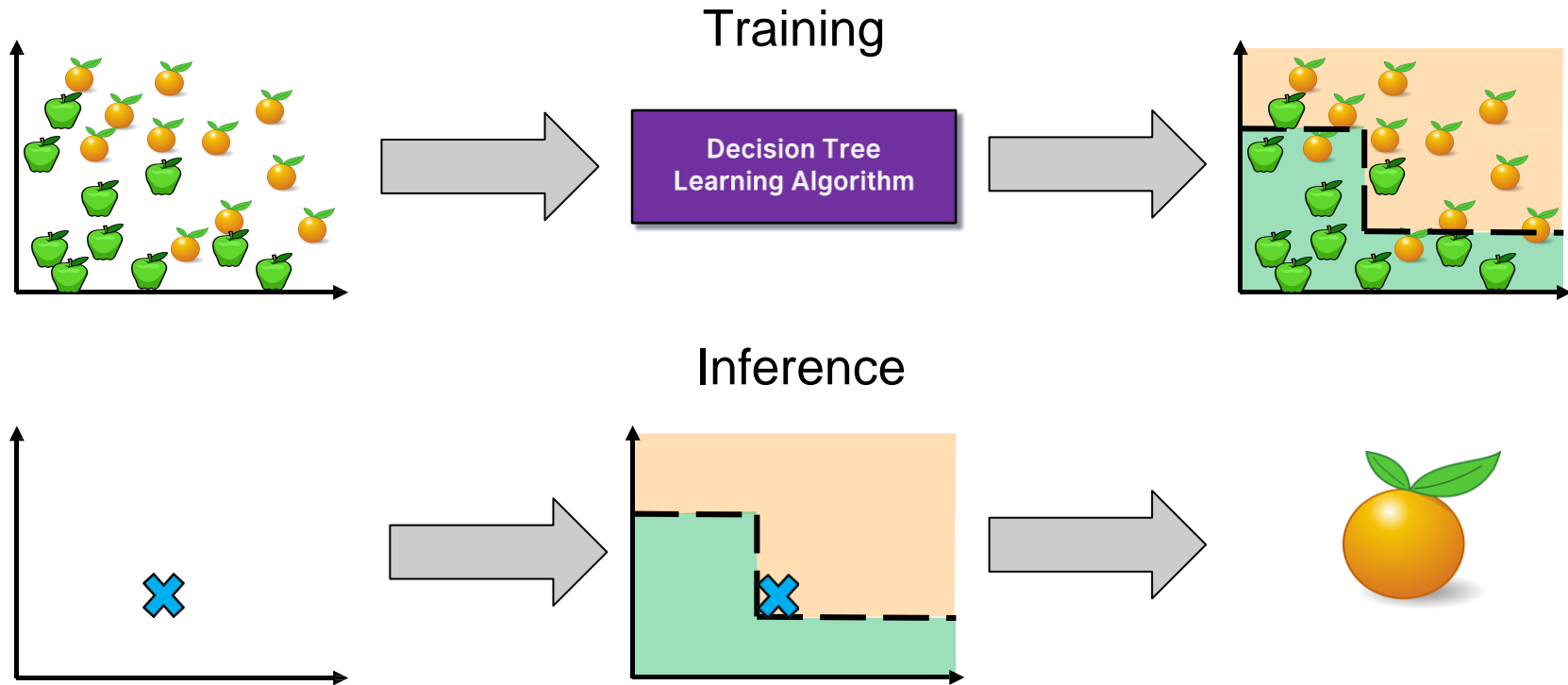- Decision trees produce non-linear decision boundaries



Logistec Regression

Decision Tree
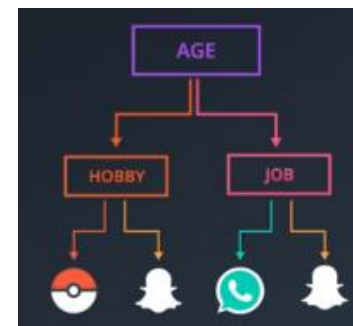
# Decision Trees: Training and Inference

# Random Forests
## (Ensemble learning with decision trees)

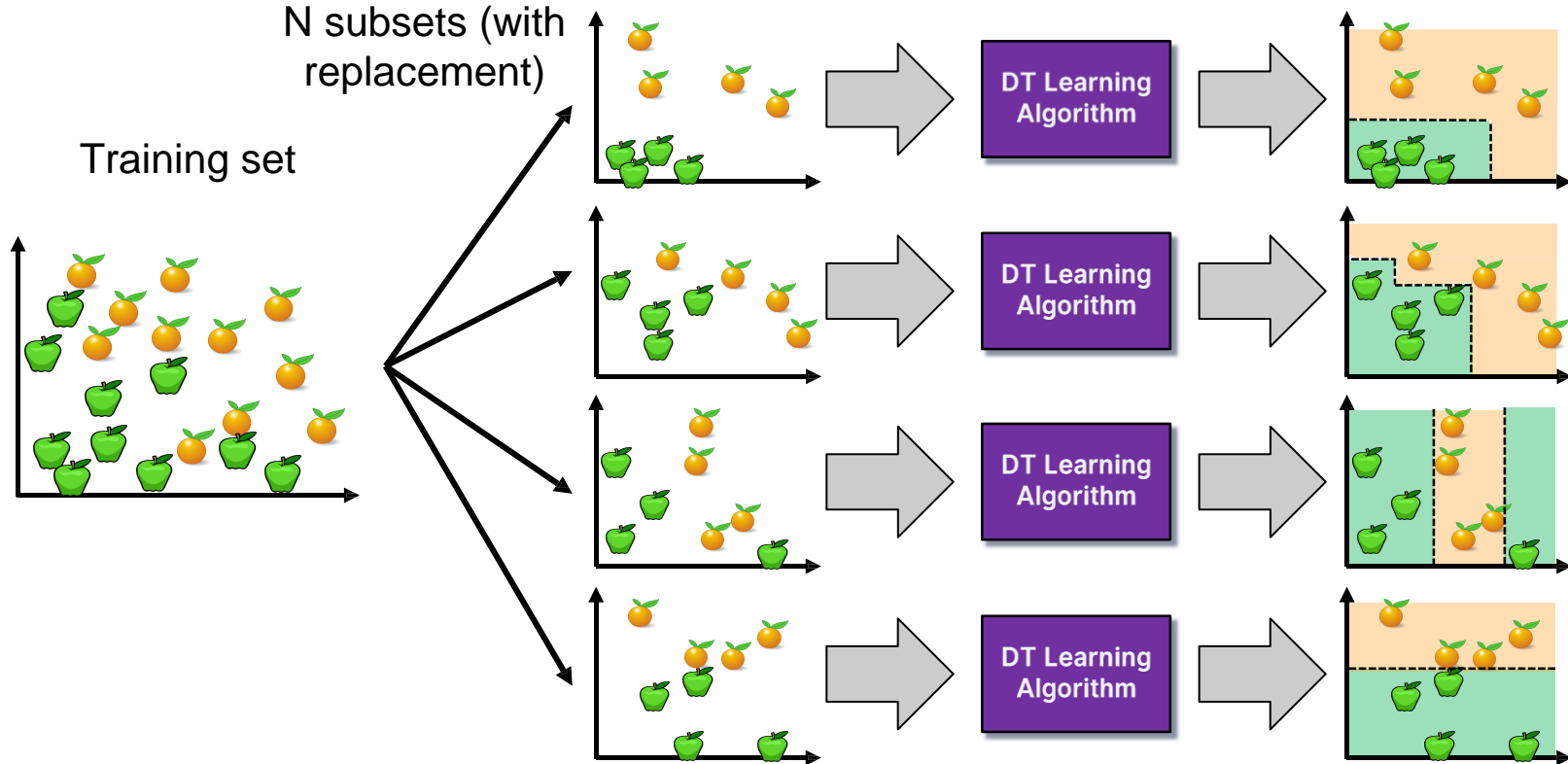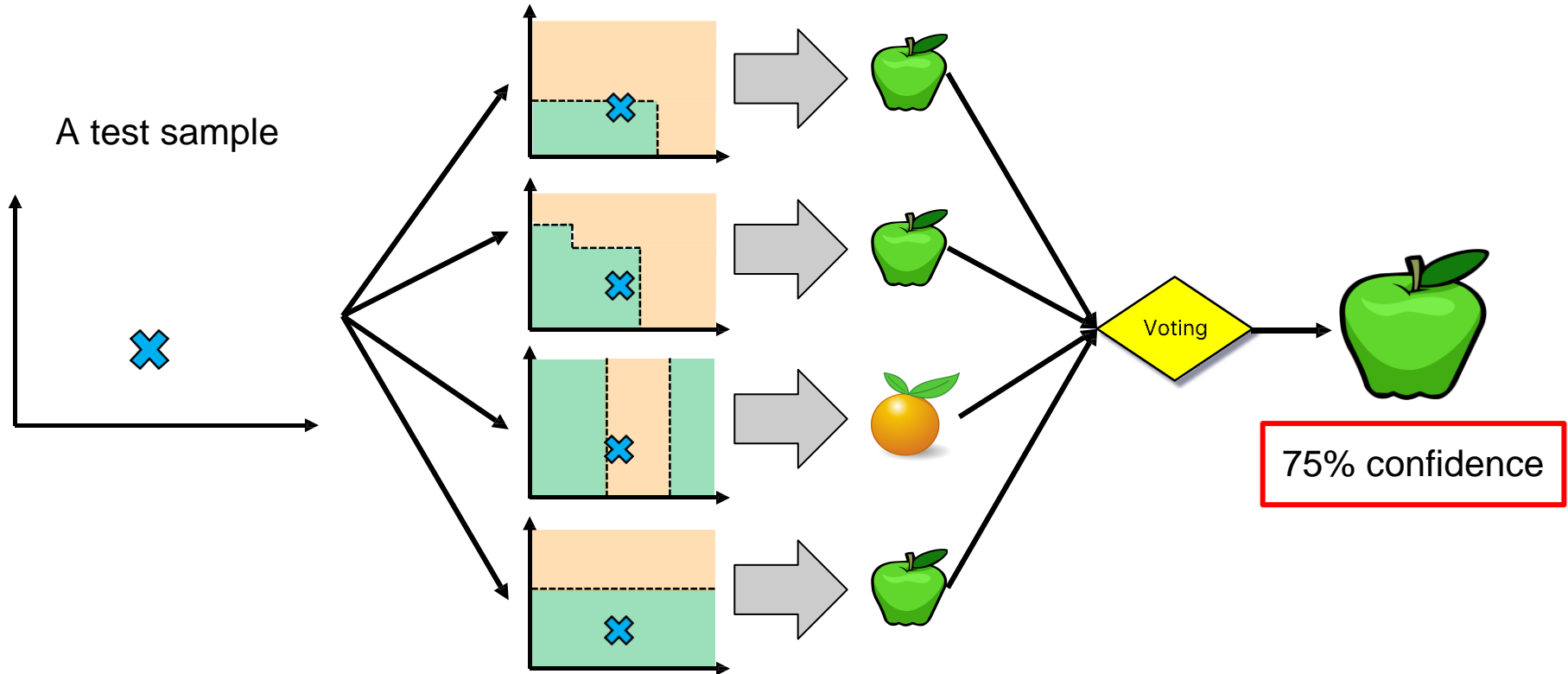| Gender | Age | Location | Platform | Job | Hobby | App |
|--------|-----|----------|----------|-----|-------|-----|
| F | 30 | US | IOS | School | Games | Whatsapp |
| F | 11 | France | Android | Work | Tennis | Pokemon Go |
| M | 16 | Chile | IOS | Temp | Tennis | Snapchat |
| F | 15 | China | IOS | Retired | Chess | Whatsapp |
| M | 25 | Us | Android | School | Games | Snapchat |
| M | 32 | Us | IOS | School | Tennis | Whatsapp |
| F | 40 | Egypt | Android | Work | Chess | Snapchat |
| M | 12 | France | Android | Temp | Tennis | Whatsapp |
| M | 14 | Australia | Android | School | chess | Pokemon Go |

Random Forests

# Random Forests

- Random Forests:
  - ➤ Instead of building a single decision tree and use it to make predictions, build many slightly different trees and combine their predictions

- We have a single data set, so how do we obtain slightly different trees?
  1. Bagging (**B**ootstrap **Agg**regat**ing**):
  - ➤ Take random subsets of data points from the training set to create N smaller data sets
  - ➤ Fit a decision tree on each subset

  2. Random Subspace Method (also known as Feature Bagging):
  - ➤ Fit N different decision trees by constraining each one to operate on a random subset of features
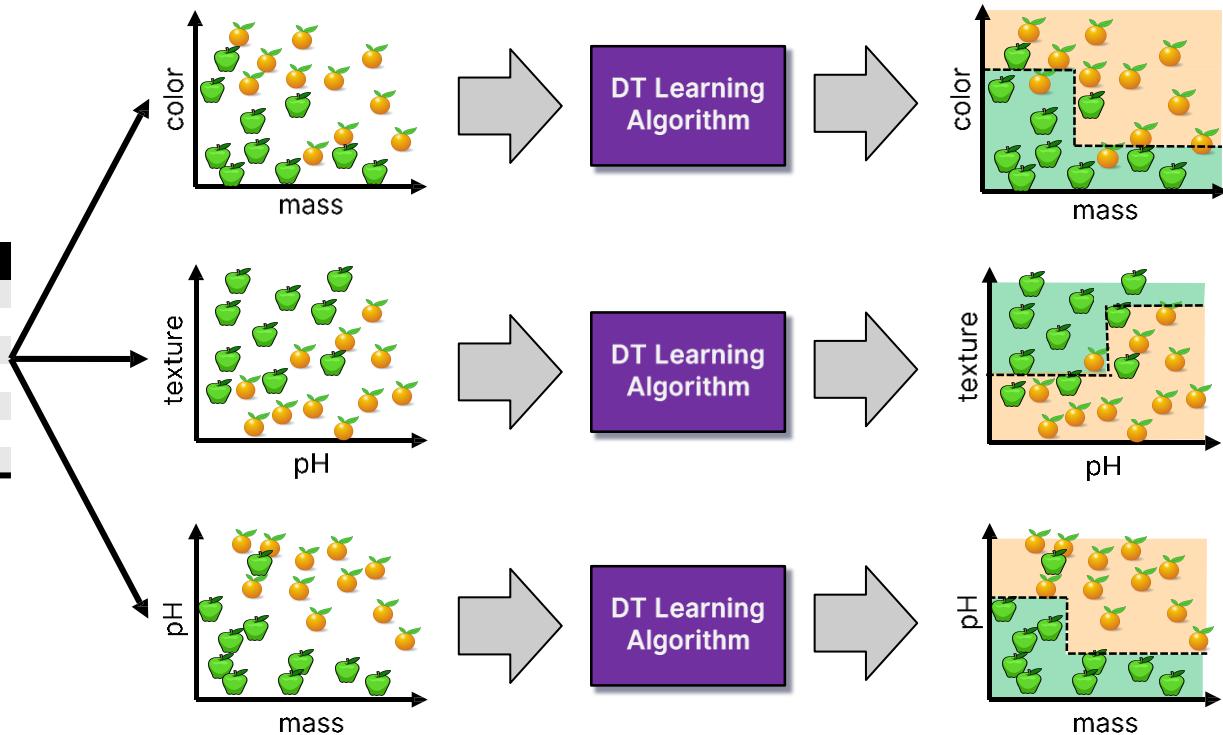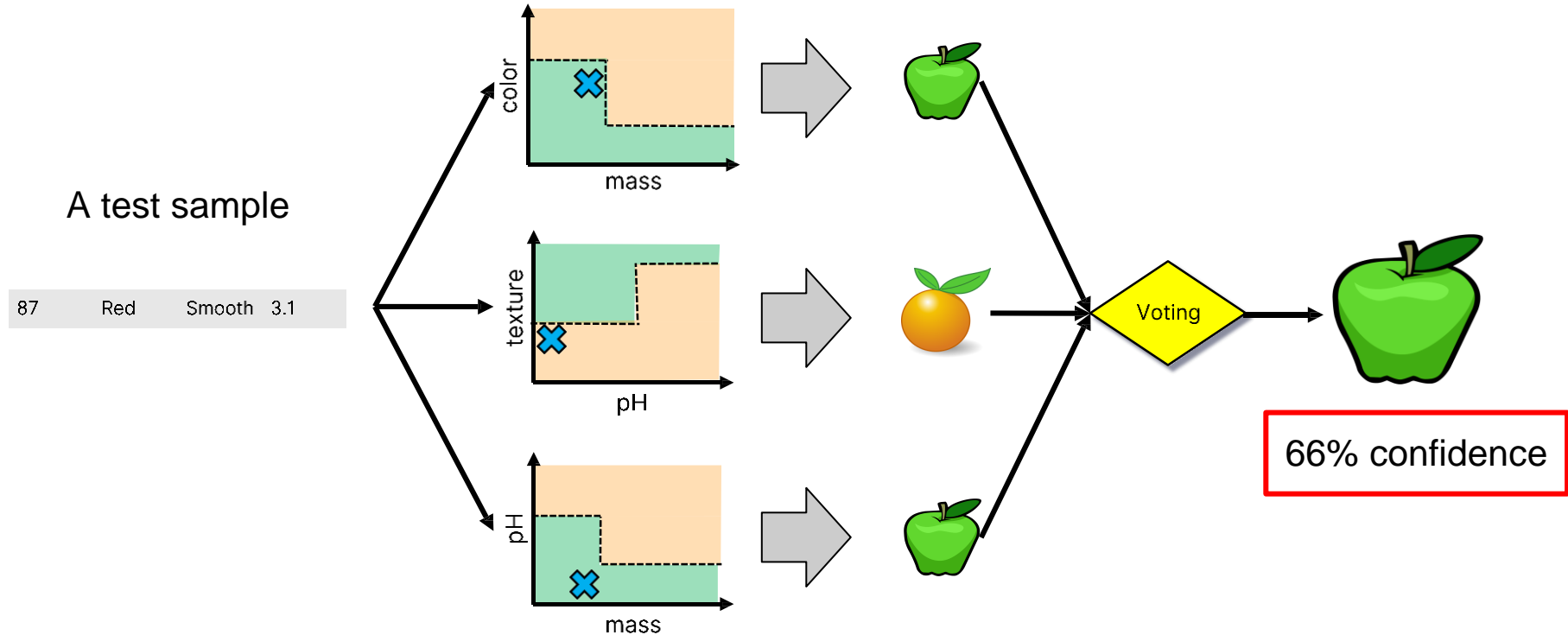
# Bagging at training time

N subsets (with replacement)

Training set



DT Learning Algorithm

DT Learning Algorithm

DT Learning Algorithm

DT Learning Algorithm

# Bagging at inference time

A test sample

75% confidence

# Random Subspace Method at training time

# Random Subspace Method at inference time
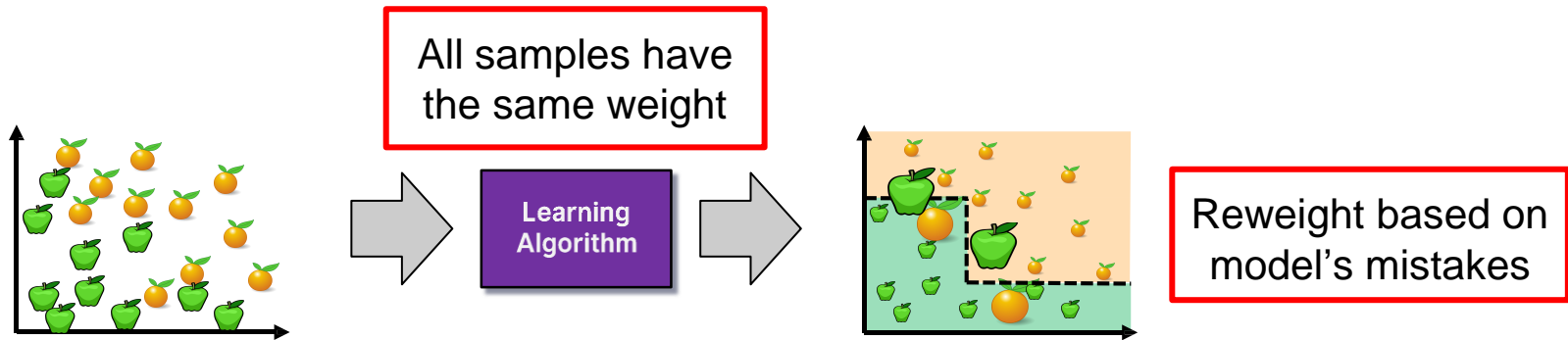
# Random Forests

# Ensemble Learning

- Ensemble Learning:
  - ➢ Method that combines multiple learning algorithms to obtain performance improvements over its components

- **Random Forests** are one of the most common examples of ensemble learning

- Other commonly-used ensemble methods:
  - ➢ Bagging: multiple models on random subsets of data samples
  - ➢ Random Subspace Method: multiple models on random subsets of features
  - ➢ Boosting: train models iteratively, while making the current model focus on the mistakes of the previous ones by increasing the weight of misclassified samples
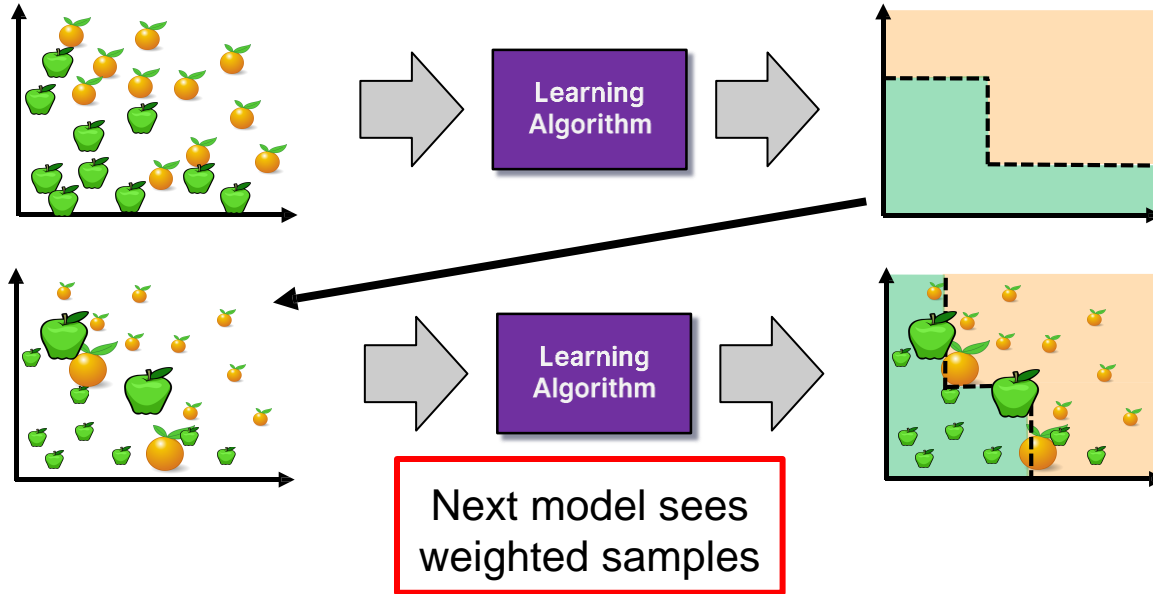
# Boosting



All samples have the same weight

Learning Algorithm

# Boosting

# Boosting



Next model sees weighted samples

# Boosting



Reweight based on current model's mistakes

# Boosting

# Boosting

# Summary

- Ensemble Learning methods combine multiple learning algorithms to obtain performance improvements over its components

- Commonly-used ensemble methods:
  - ➢ Bagging (multiple models on random subsets of data samples)
  - ➢ Random Subspace Method (multiple models on random subsets of features)
  - ➢ Boosting (train models iteratively, while making the current model focus on the mistakes of the previous ones by increasing the weight of misclassified samples)
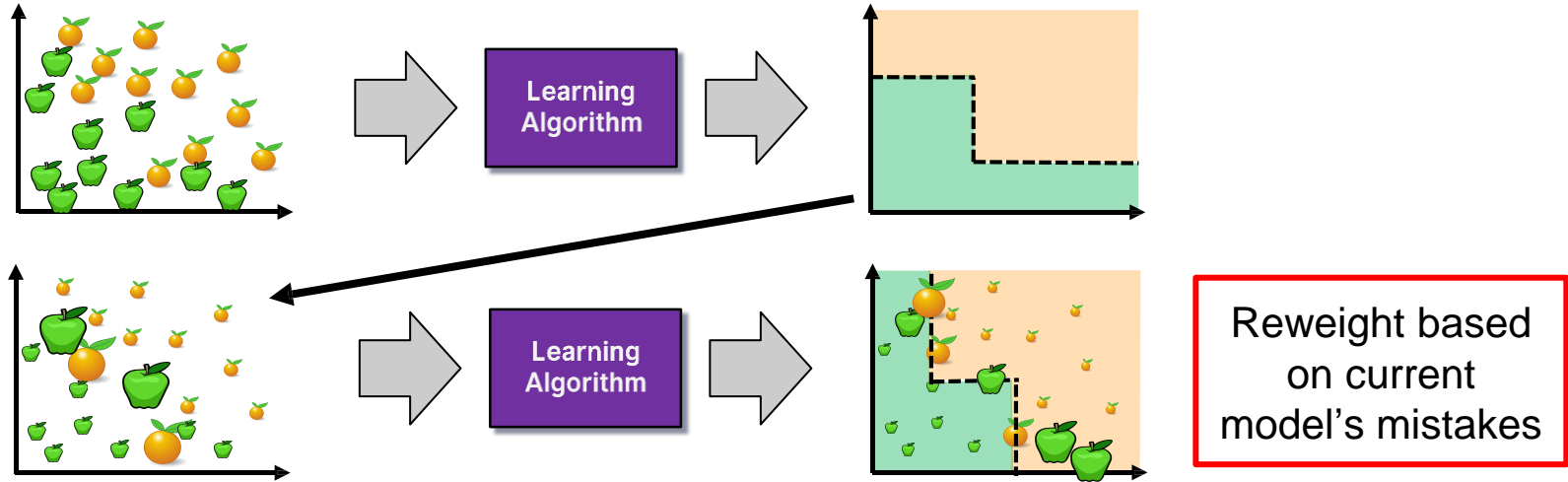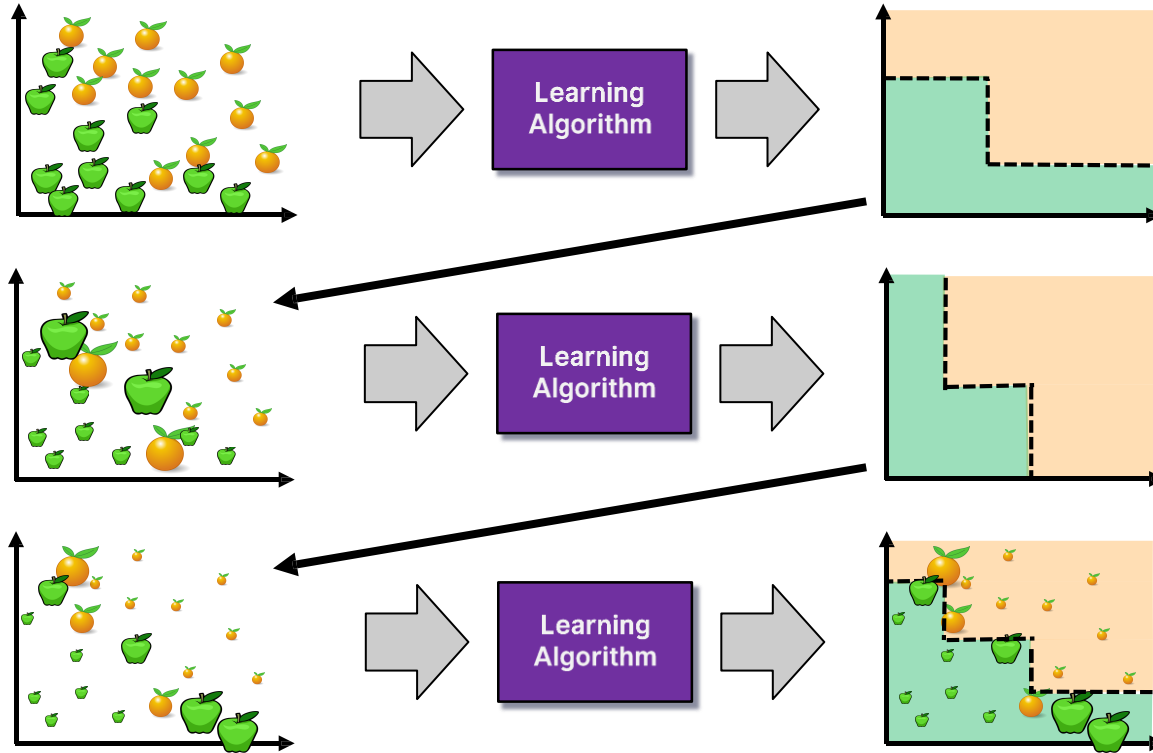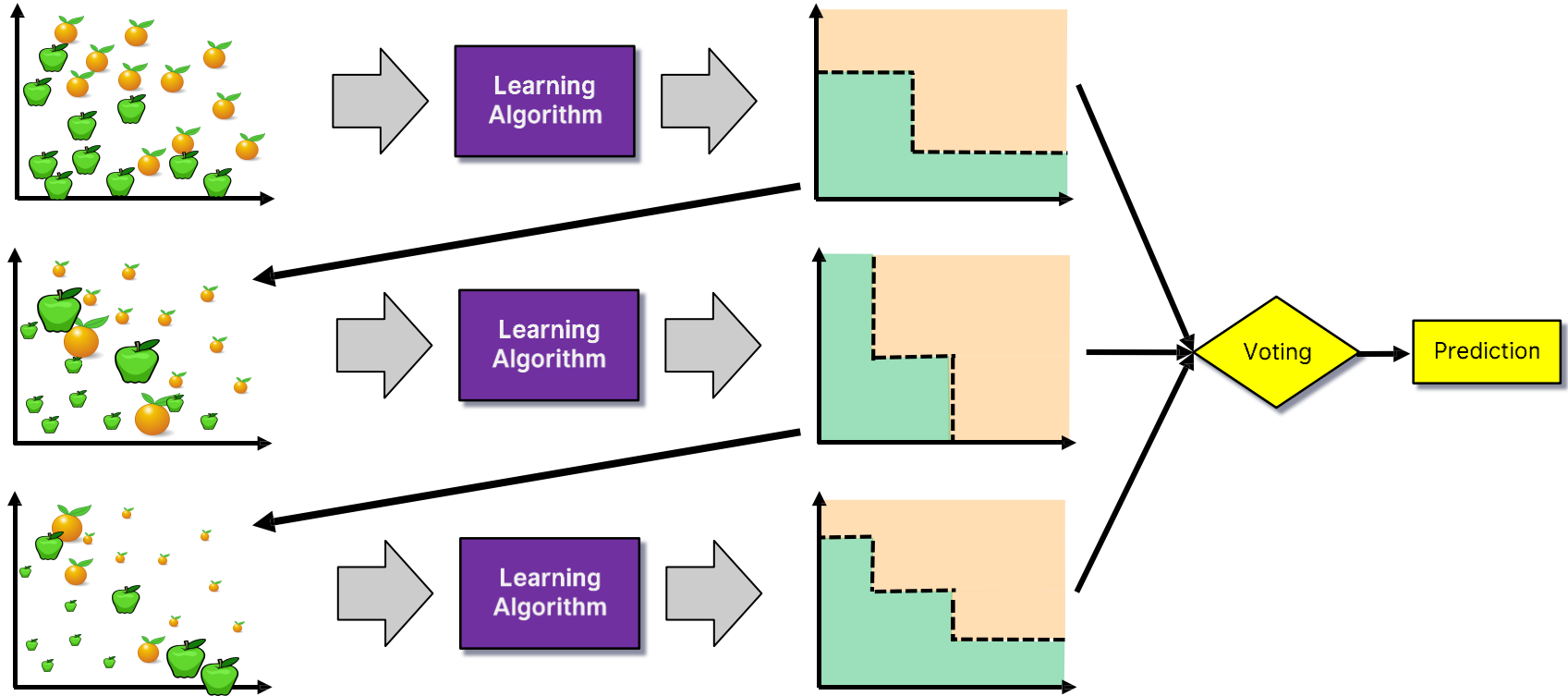
- **Random Forests** are an ensemble learning method that employ decision tree learning to build multiple trees through **bagging** and **random subspace method**.
  - ➢ They rectify the overfitting problem of decision trees!

# Decision Trees and Random Forest (Python)

```python
from sklearn.tree import DecisionTreeClassifier
from sklearn.ensemble import RandomForestClassifier

clf = DecisionTreeClassifier(criterion = "entropy", min_samples_leaf = 3)
# Lots of parameters: criterion = "gini" / "entropy";
#                      max_depth;
#                      min_impurity_split;

clf.fit(X, y) # It can only handle numerical attributes!
# Categorical attributes need to be encoded, see LabelEncoder and OneHotEncoder

clf.predict([x]) # Predict class for x

clf.feature_importances_ # Importance of each feature
clf.tree_ # The underlying tree object

clf = RandomForestClassifier(n_estimators = 20) # Random Forest with 20 trees
```