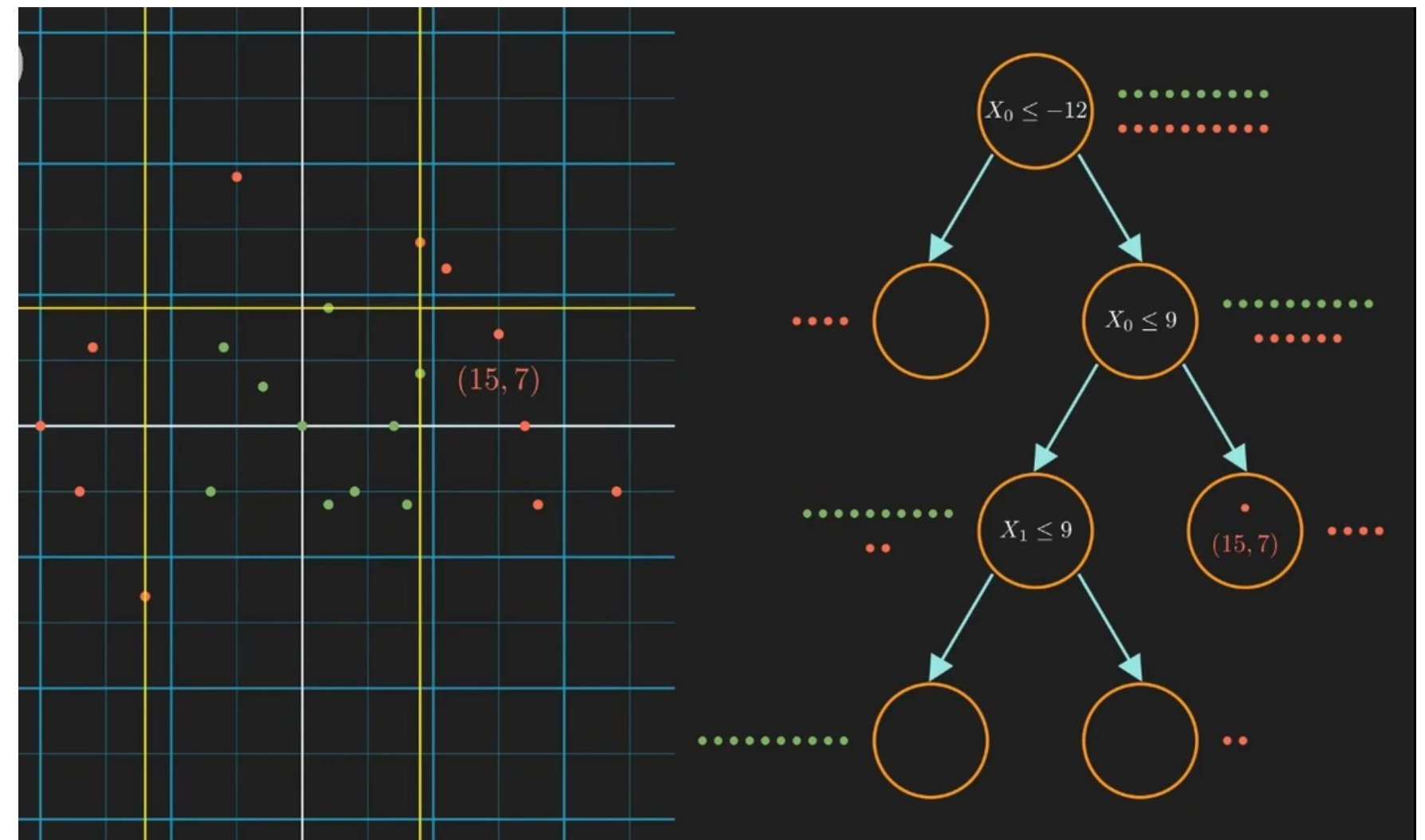


Decision Tree

A **decision tree** is a **supervised machine learning** algorithm that uses a **tree-like structure** to model decisions and predict outcomes. It is a flowchart-like structure where **internal nodes represent features or attributes**, **branches represent decision rules**, and **leaf nodes represent the final predictions or outcomes**.

The construction of a decision tree involves recursively partitioning the dataset based on different features and splitting criteria. The algorithm selects the best feature to split the data, based on metrics such as **information gain** or **Gini impurity**, which measure the purity or homogeneity of the subsets resulting from the split.



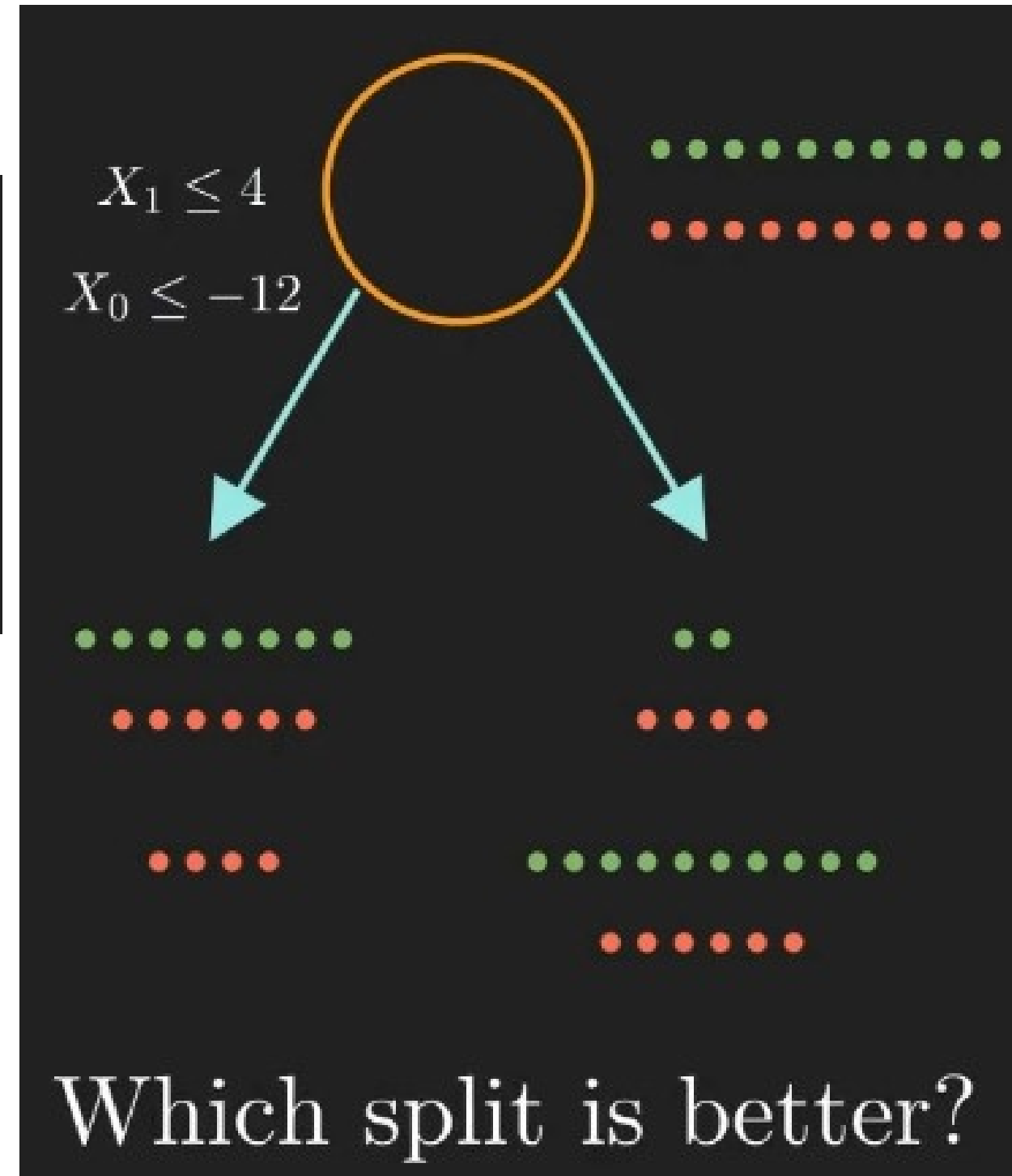
Algorithm 1 : Entropy and Information Gain

$$Entropy = \sum - p_i \log(p_i)$$

p_i = probability of class i

Information Gain =
Entropy(S) - (Weight *
Entropy(each feature)

Weight of a feature = Number of samples in the feature/
Total samples before split



Algorithm 2 : Gini Impurity Index and Gini Gain

$$\text{Gini Impurity Index} = 1 - \sum P(i)^2$$

Gini Gain =
Gini impurity before split
- sum (Weight * Gini
impurity of each subset)

How the decision tree would have processed?

- ① Initial gini index of column "Label".
4 spam, 6 non-spam.

$$\text{Gini} = 1 - \left(\frac{4}{10}\right)^2 - \left(\frac{6}{10}\right)^2 = \underline{\underline{0.48}}$$

- ② Based on Column "Sender".

Sender = 0. and $\neq 0$.

Labels: 1, 1, 0, 0 Labels: 0, 1, 1, 0, 0, 0

$$1^{\text{st}} \text{ Gini} = 1 - 0.25 - 0.25 = 0.5$$

$$2^{\text{nd}} \text{ Gini} = 1 - 0.111 - 0.444 = 0.444$$

$$\text{Gini Gain} = 0.48 - \frac{4}{10} \times 0.5 - \frac{6}{10} \times 0.444 \\ = \underline{\underline{0.0136}}$$

Sender = 1 and $\neq 0$

Labels: 0, 1, 0 Labels: 1, 1, 1, 0, 0, 0, 0

$$1^{\text{st}} \text{ Gini} = 1 - \left(\frac{1}{3}\right)^2 - \left(\frac{2}{3}\right)^2 = 0.444$$

$$2^{\text{nd}} \text{ Gini} = 1 - \left(\frac{3}{7}\right)^2 - \left(\frac{4}{7}\right)^2 = 0.489$$

$$\text{Gini Gain} = 0.48 - \frac{3}{10} \times 0.444 - \frac{7}{10} \times 0.489 \\ = -0.0045$$

Sender = 2 and $\neq 0$

- ③ Based on Column "Subject".

Subject = 0, 1, 2.

- ⑤ Based on "Attachment" = 0 or 1.

- ④ Based on Column Body Length

Sort it: 40, 50, 55, 60, 70, 75, 80, 90, 100, 120.

Splits at mid points.

Body length ≤ 45 & > 45

Labels: 0

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

Labels: (4)(1s), (5)(0s).

$$\text{Gini} = 1 - \left(\frac{4}{9}\right)^2 - \left(\frac{5}{9}\right)^2 \\ = 0.493$$

$$\text{Gini Gain} = 0.48 - 0 \times \frac{1}{10} - 0.493 \times \frac{9}{10} \\ = \underline{\underline{0.035}}$$

$$B.L \leq 52.5$$

$$B.L \leq 57.5$$

$$B.L \leq 65$$

Labels: 0, 0, 0, 0

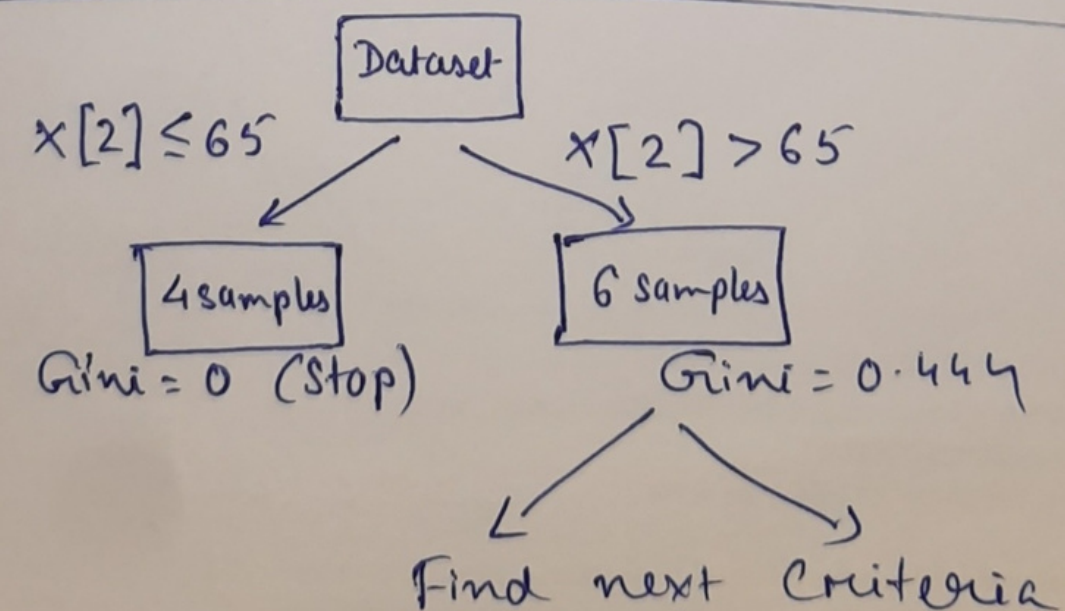
$$\text{Gini} = 0$$

$$B.L > 65$$

1 1 1 1 0 0

$$\text{Gini} = 0.444$$

$$\text{Gini Gain} = 0.48 - 0.444 \times 0.6 = \underline{\underline{0.214}}$$



The decision tree algorithm

- Select the best feature: The first step is to select the feature that best splits the data into classes or reduces the variance of the target variable for regression tasks. The most common measure for selecting the best feature is the information gain or the Gini impurity.
- Split the data: Once the best feature is selected, the data is split into subsets based on the values of the selected feature. Each subset corresponds to a child node of the current node in the tree.
- Recursively repeat steps 1 and 2: The above steps are repeated recursively for each child node until a stopping criterion is met. This stopping criterion can be a predefined maximum depth of the tree, a minimum number of samples required to split an internal node or a minimum number of samples required to be at a leaf node.
- Predict the target variable: Once the decision tree is built, the target variable can be predicted by traversing the tree from the root node to a leaf node based on the decision rules at each node.
- Prune the tree: To avoid overfitting, the tree can be pruned by removing unnecessary nodes or subtrees that do not improve the performance on a validation set.

Bayes

The **Naive Bayes theorem** is a statistical technique used for classification and prediction tasks. It is based on Bayes' theorem, which describes the probability of an event occurring given that another related event has already occurred.

The Naive Bayes algorithm assumes that the features (or attributes) used to predict the class of an instance are independent of each other. This assumption is often simplified as "naive" because features are often correlated to some extent in real-world scenarios.

The Naive Bayes algorithm works by calculating the probability of a given class, given a set of input features. It does this by multiplying the conditional probabilities of each feature given the class. The class with the highest probability is then assigned as the prediction for that input.

Here's the formula for Naive Bayes:

$$P(\text{class} \mid \text{features}) = (P(\text{class}) * P(\text{features} \mid \text{class})) / P(\text{features})$$

- $P(\text{class} \mid \text{features})$ is the probability of the class given the input features
- $P(\text{class})$ is the prior probability of the class (before observing the input features)
- $P(\text{features} \mid \text{class})$ is the conditional probability of the features given the class
- $P(\text{features})$ is the prior probability of the input features (before knowing the class)

To use Naive Bayes for classification, we need to train the model first. This involves estimating the prior probabilities and conditional probabilities from a labeled dataset. Then, when new instances are presented, we can use these probabilities to calculate the probability of each class and make predictions.

Naive Bayes has proven to be an effective algorithm for text classification, spam filtering, and other tasks where the input features are discrete and independent. However, it may not work well when the features are highly correlated, or when the distribution of the features is complex and not well understood.

Press Esc to exit full screen

Estimate the value of Y given that $X = (0, 2)$

X_1	X_2	Y
0	0	0
0	1	1
1	2	1
0	0	1
2	2	0
1	1	0
0	2	1
2	0	0
2	1	0
1	0	0

Let's compute $P(Y = 0|X = (0, 2))$ and $P(Y = 1|X = (0, 2))$...

$$P(Y = 0) = \frac{\#Y = 0}{\#Y = 0 + \#Y = 1} = \frac{6}{10}$$

$$P(Y = 1) = \frac{\#Y = 1}{\#Y = 0 + \#Y = 1} = \frac{4}{10}$$

$$P(X = (0, 2)|Y = 1) = \frac{1}{4}$$

$$P(X = (0, 2)|Y = 0) = 0$$

Estimated value of Y is 1 given $X = (0, 2)$

$$P(Y|X) = P(Y)*P(X|Y)/P(X)$$

X_1	X_2	Y
0	0	0
0	1	1
1	2	1
0	0	1
2	2	0
1	1	0
0	2	1
2	0	0
2	1	0
1	0	0

What if the features do not match our required test data exactly?

We Naively assume that our features are independent.

Let's compute $P(Y = 0|X = (0, 2))$ and $P(Y = 1|X = (0, 2)) \dots$

$$P(Y = 0) = \frac{\#Y = 0}{\#Y = 0 + \#Y = 1} = \frac{6}{10}$$

$$P(Y = 1) = \frac{\#Y = 1}{\#Y = 0 + \#Y = 1} = \frac{4}{10}$$

$$P(X = (0, 2)|Y = 1) = P(X_1 = 0|Y = 1) * P(X_2 = 2|Y = 1) = \frac{3}{4} * \frac{2}{4}$$

$$P(X = (0, 2)|Y = 0) = P(X_1 = 0|Y = 0) * P(X_2 = 2|Y = 0) = \frac{1}{6} * \frac{1}{6}$$

$$\frac{3}{4} * \frac{2}{4} > \frac{1}{6} * \frac{1}{6}$$

So, the estimated value of Y is 1

X_1	X_2	Y
0	0	0
0	1	1
1	2	1
0	0	1
2	2	0
1	1	0
0	2	1
2	0	0
2	1	0
1	0	0

Dealing with Continuous features

1. Discretization

$$age \in (11, 50) \longrightarrow age \in \{1, 2, 3, 4\}$$