# Introduction to machine learning in Hydrology

**Lazaro J. Perez & Marc Berghouse**

contact: lazaro.perez@dri.edu

SB
LAB

DRI

# Outline

- Decision Trees

- Decision Trees in Matlab

- Classification and Regression Learner App

# Decision trees

## Definition

A decision tree is a non-parametric supervised learning algorithm that is utilized for both classification and regression tasks

It has a hierarchical tree structure, which consists of a root node, branches, internal nodes, and leaf nodes
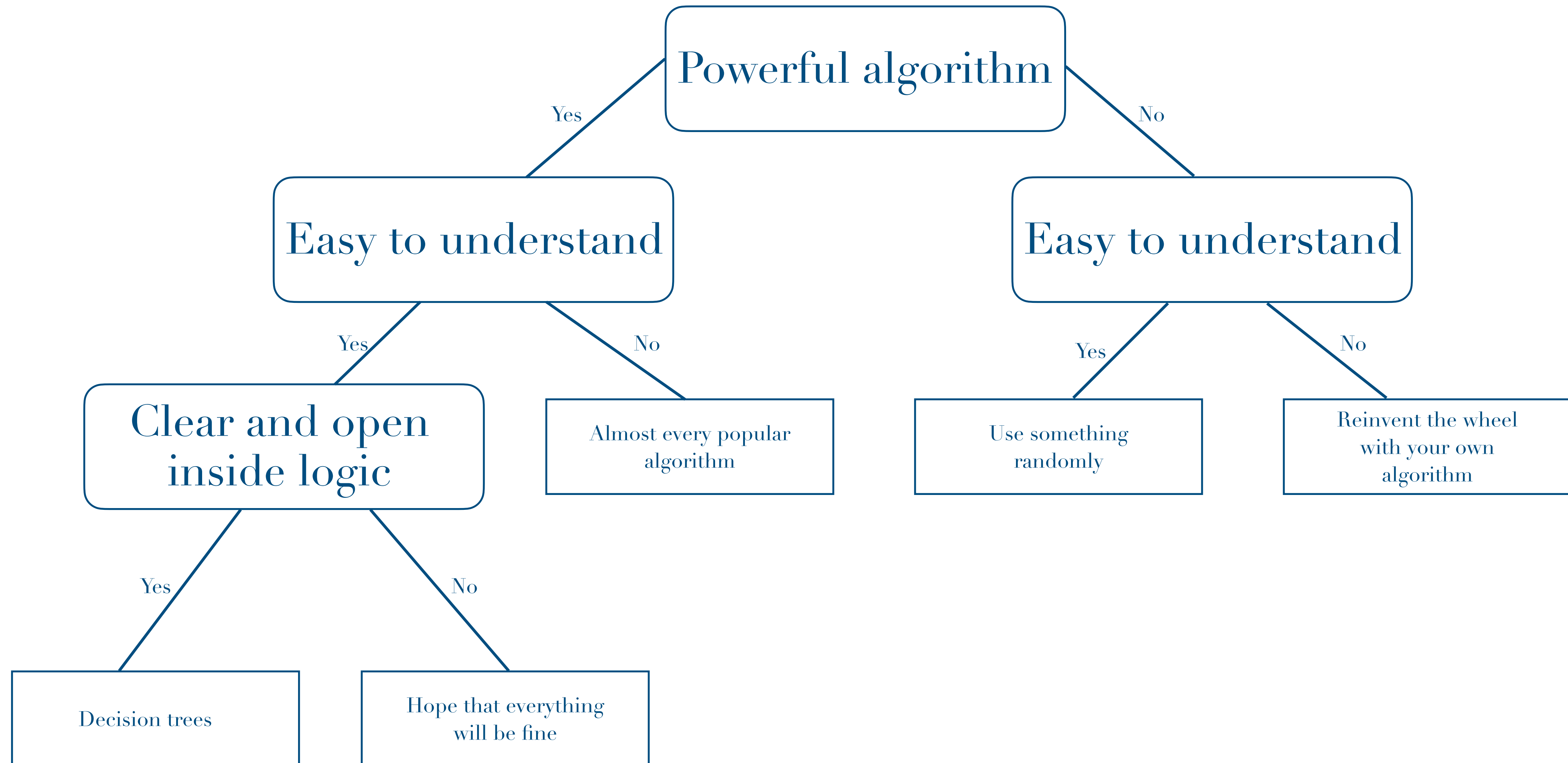
## General consensus

It's a binary tree that recursively splits the dataset until we're left with pure leaf nodes

A decision tree analysis is a divide-and-conquer approach to classification and regression
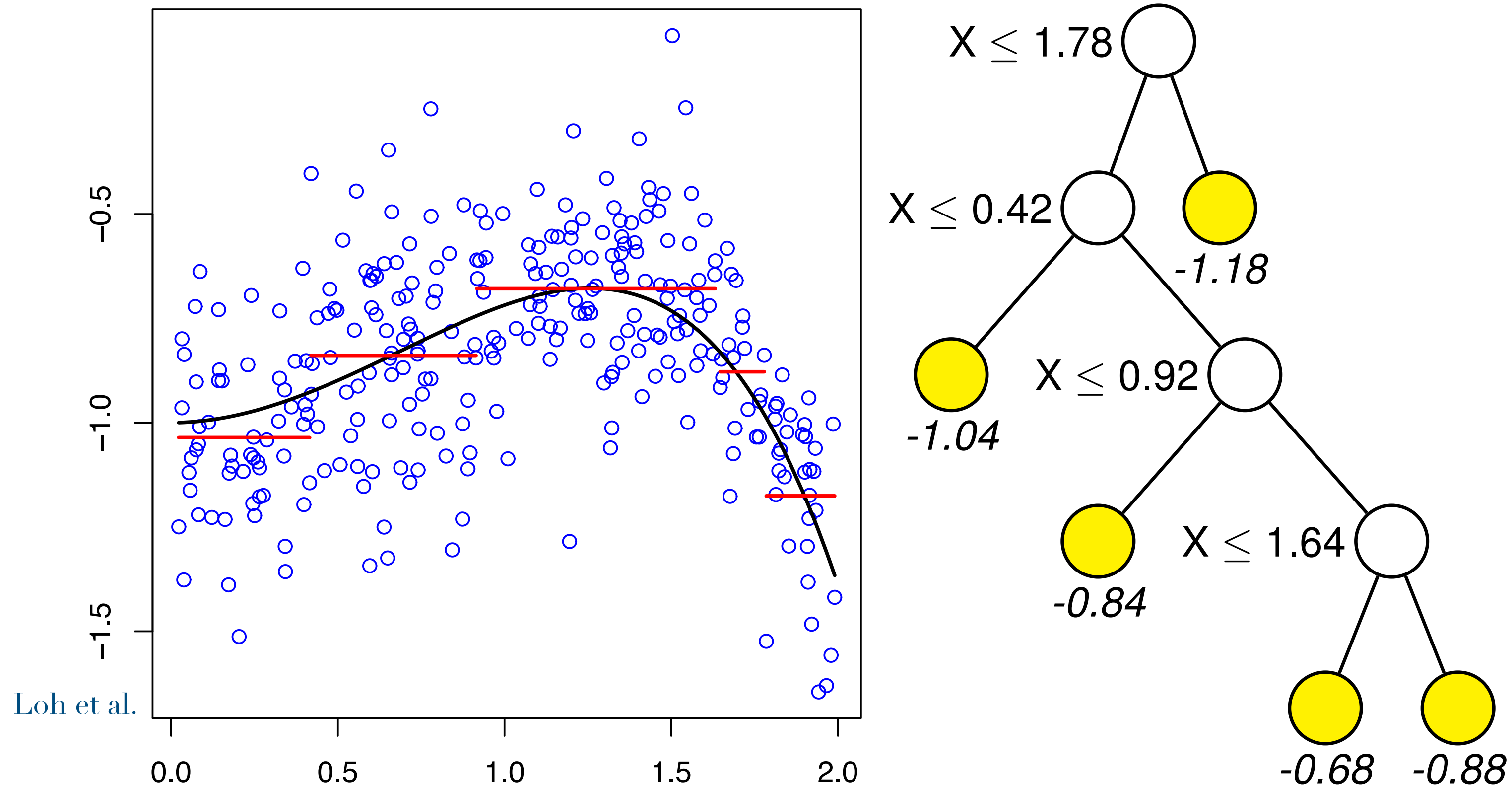
# Decision trees: Graphical representation

# Decision trees

## Development

First generation

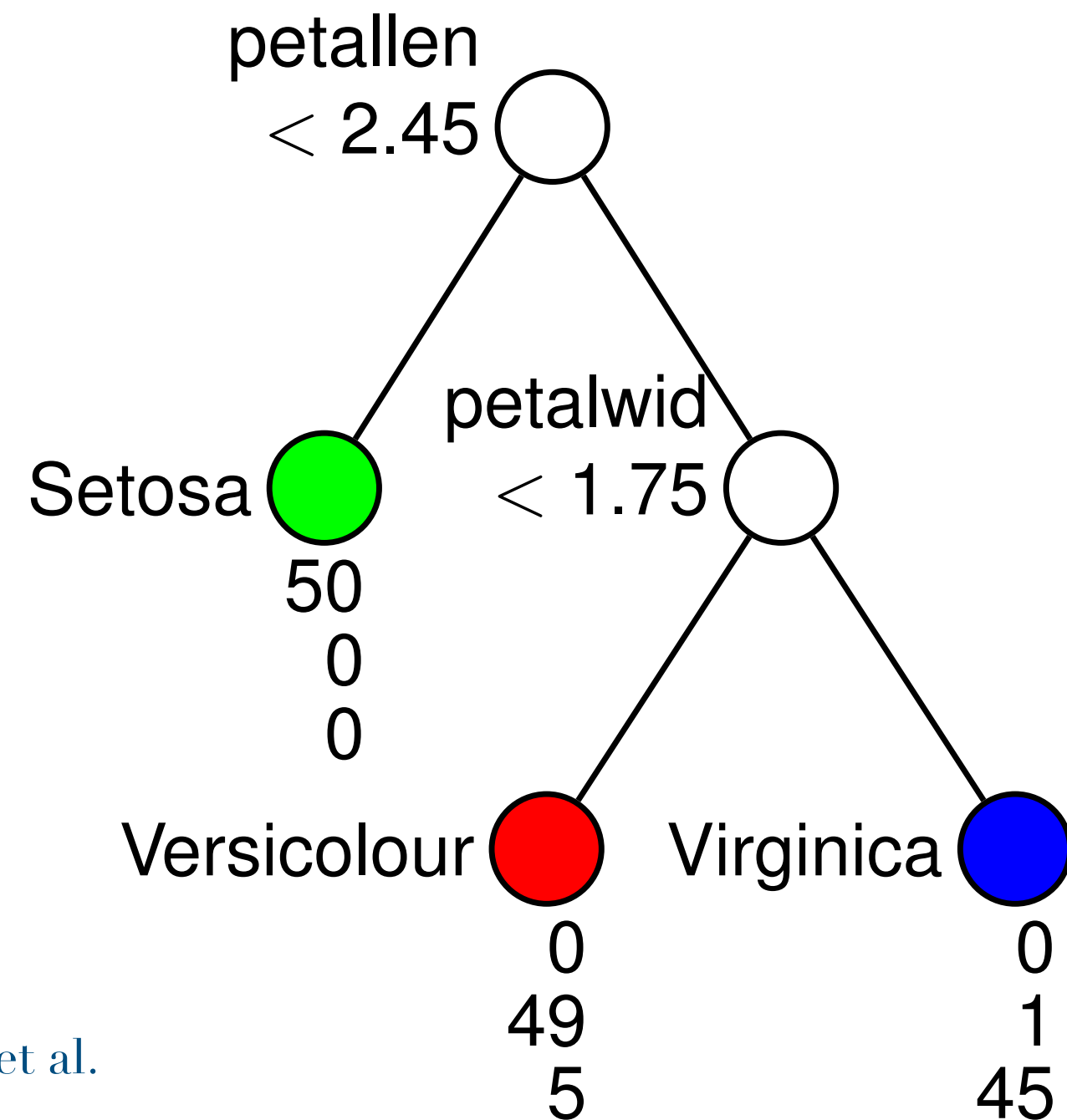Automatic interaction detection (AID): Morgan and Sonquist, 1963
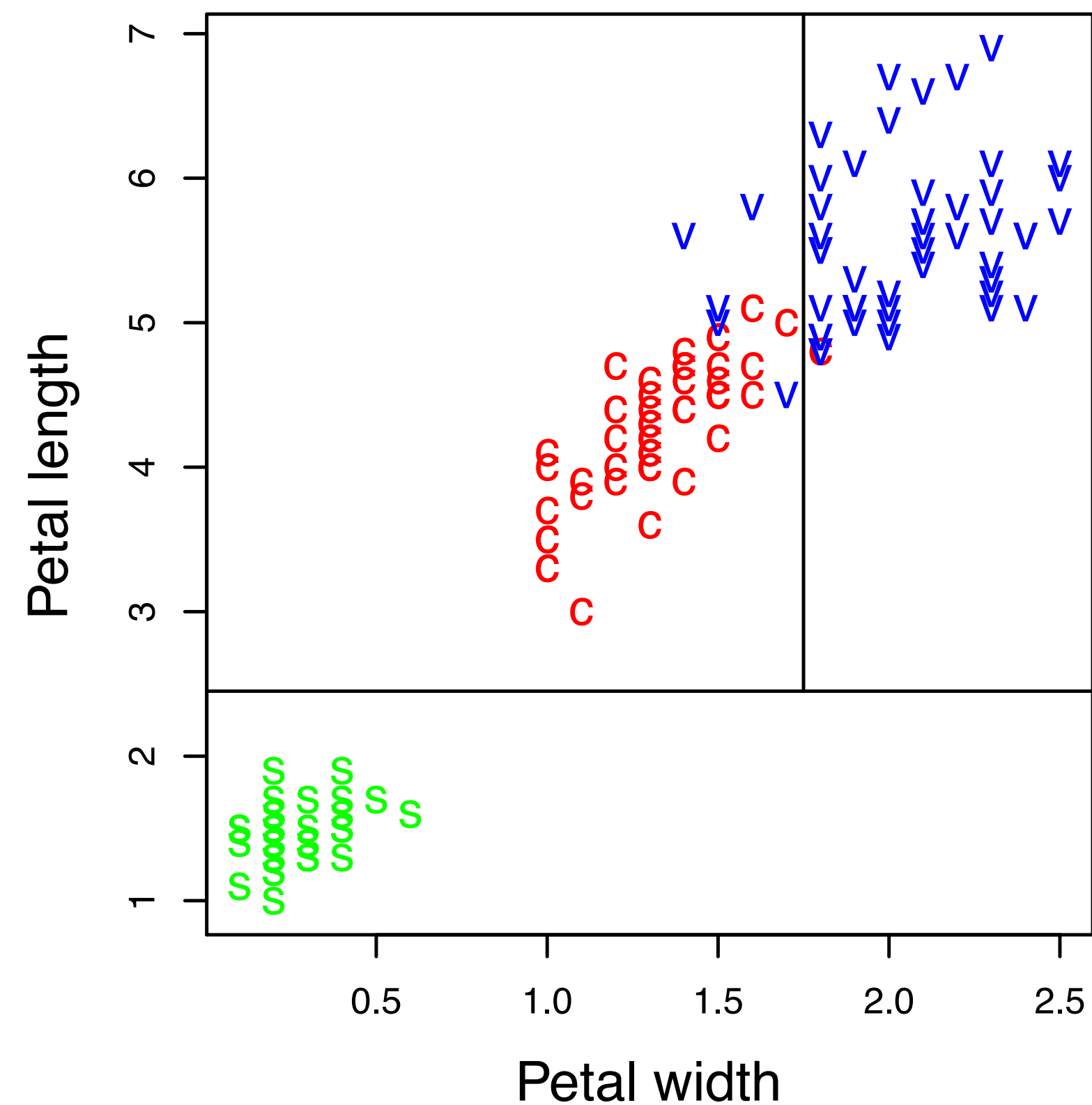


Loh et al.

# Decision trees

## Development

### First generation

Theta automatic interaction detection (THAID): Messenger and Mandell, 1972



Loh et al.

# Decision trees

**Development**

Second generation

Classification and regression tree (CART): Breiman et al., 1984

- Adds cross-validation
- Adds pruning

Third generation

Quick unbiased efficient statistical tree (QUEST): Loh and Shih, 1997

- Merge classes to get binary splits

Fourth generation

Random forest: Breiman et al., 2001

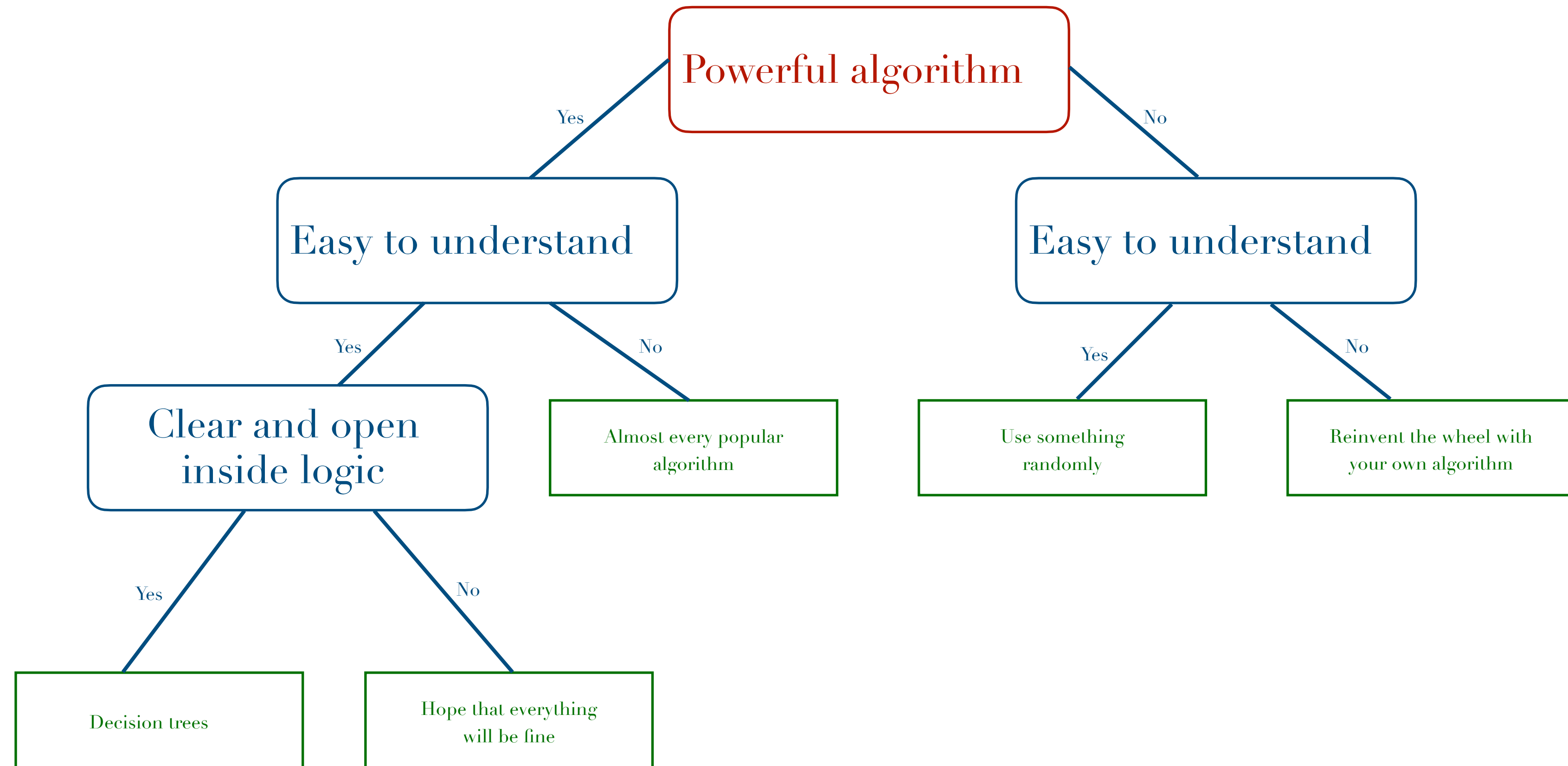- Final predicted value is average of values from the trees

# Decision trees: Key concepts

**Root: Top node**

**Node: Each object in a tree**

**Leaf node: Final node**

Powerful algorithm

Yes — Easy to understand

No — Easy to understand

Yes — Clear and open inside logic

No — Almost every popular algorithm

Yes — Use something randomly

No — Reinvent the wheel with your own algorithm

Yes — Decision trees

No — Hope that everything will be fine

# Decision trees: Splitting criterion

**Entropy**

$$S = -\sum_{k \in K} p(k) \log_2 p(k)$$

Entropy values fall between 0 and 1          Pure leaf nodes (only one class) has $S=0$

**Information gain**

$$\Theta = S(p) - \sum w_i S(c_i)$$

Information contained in a state

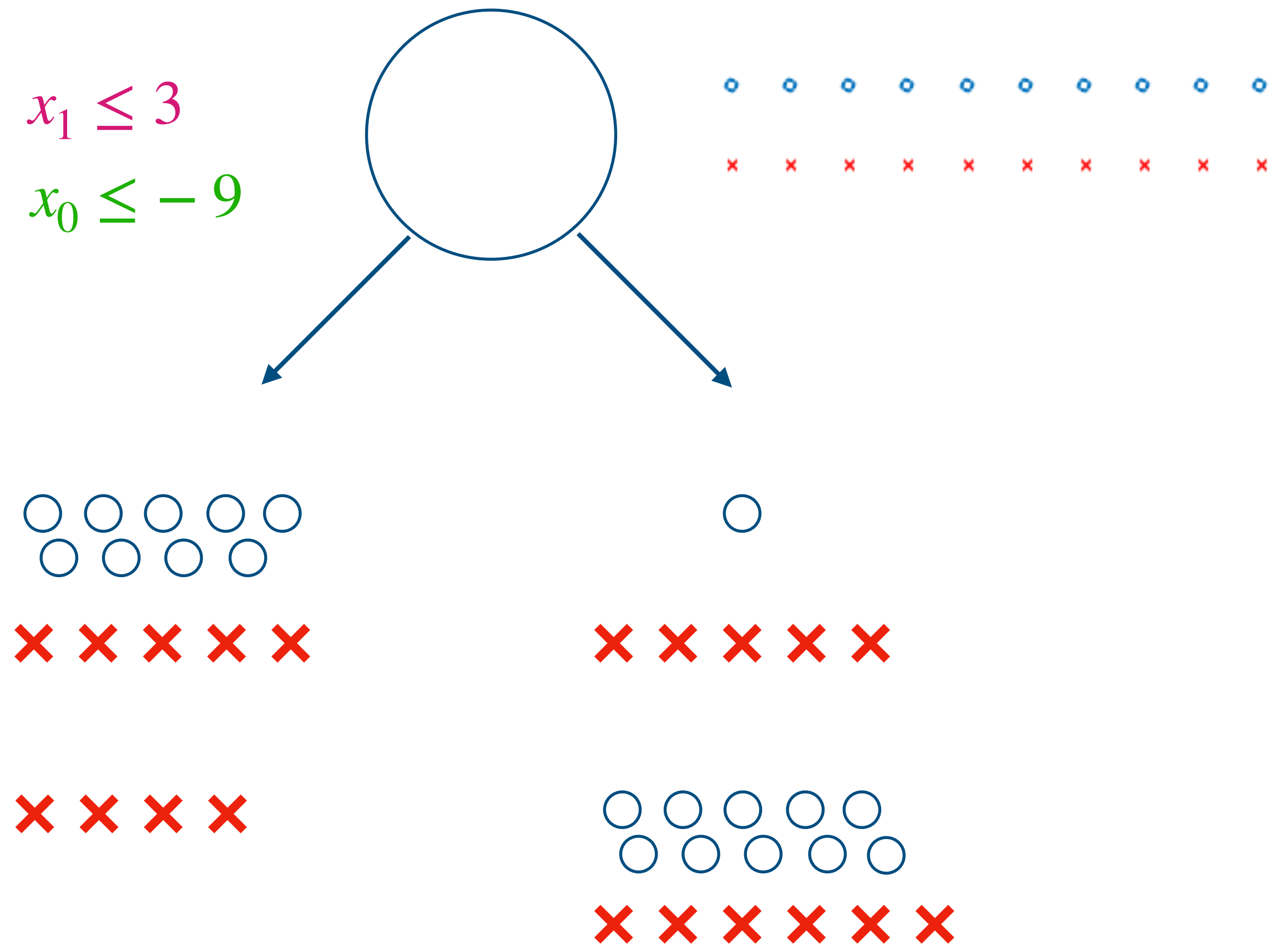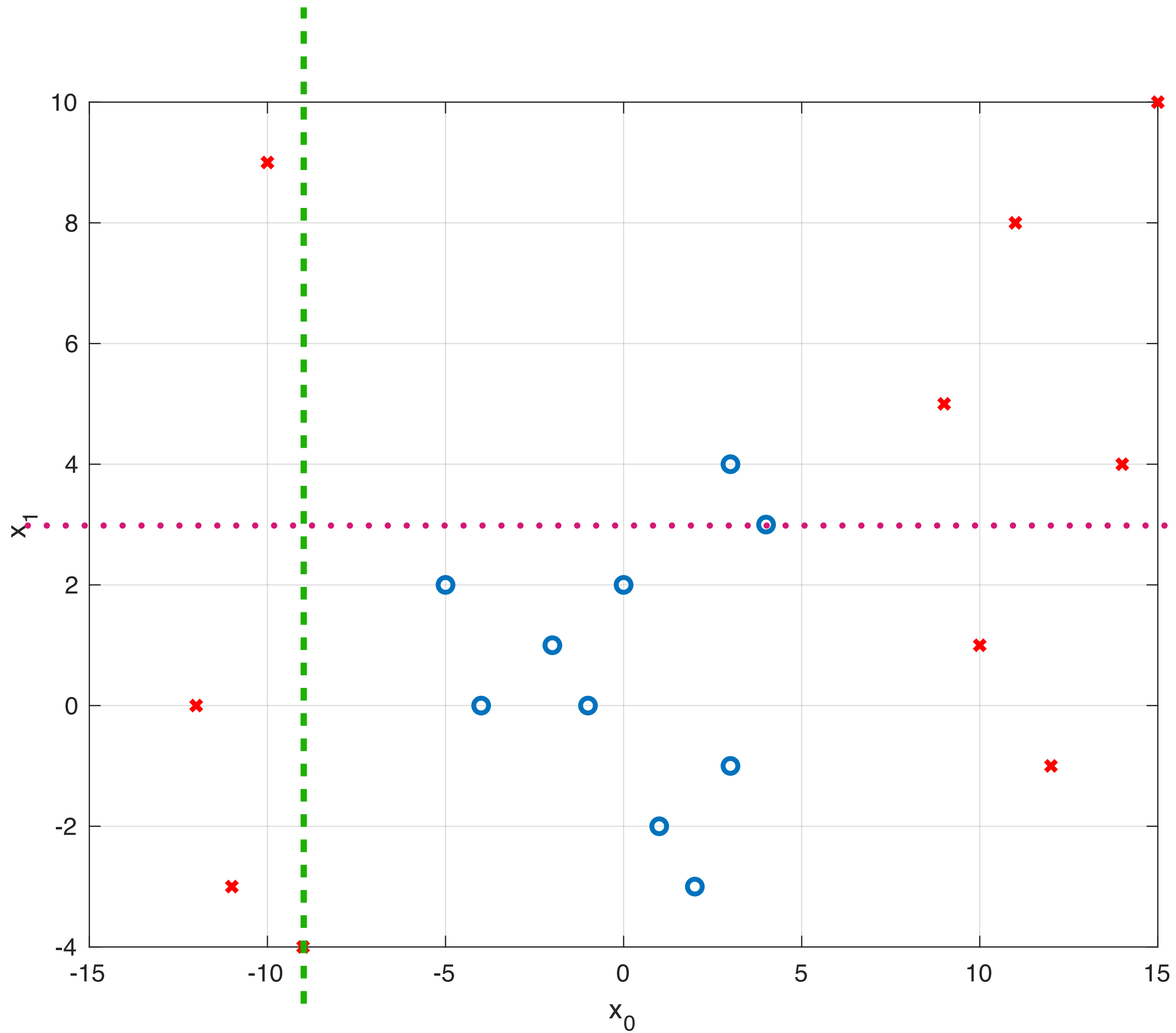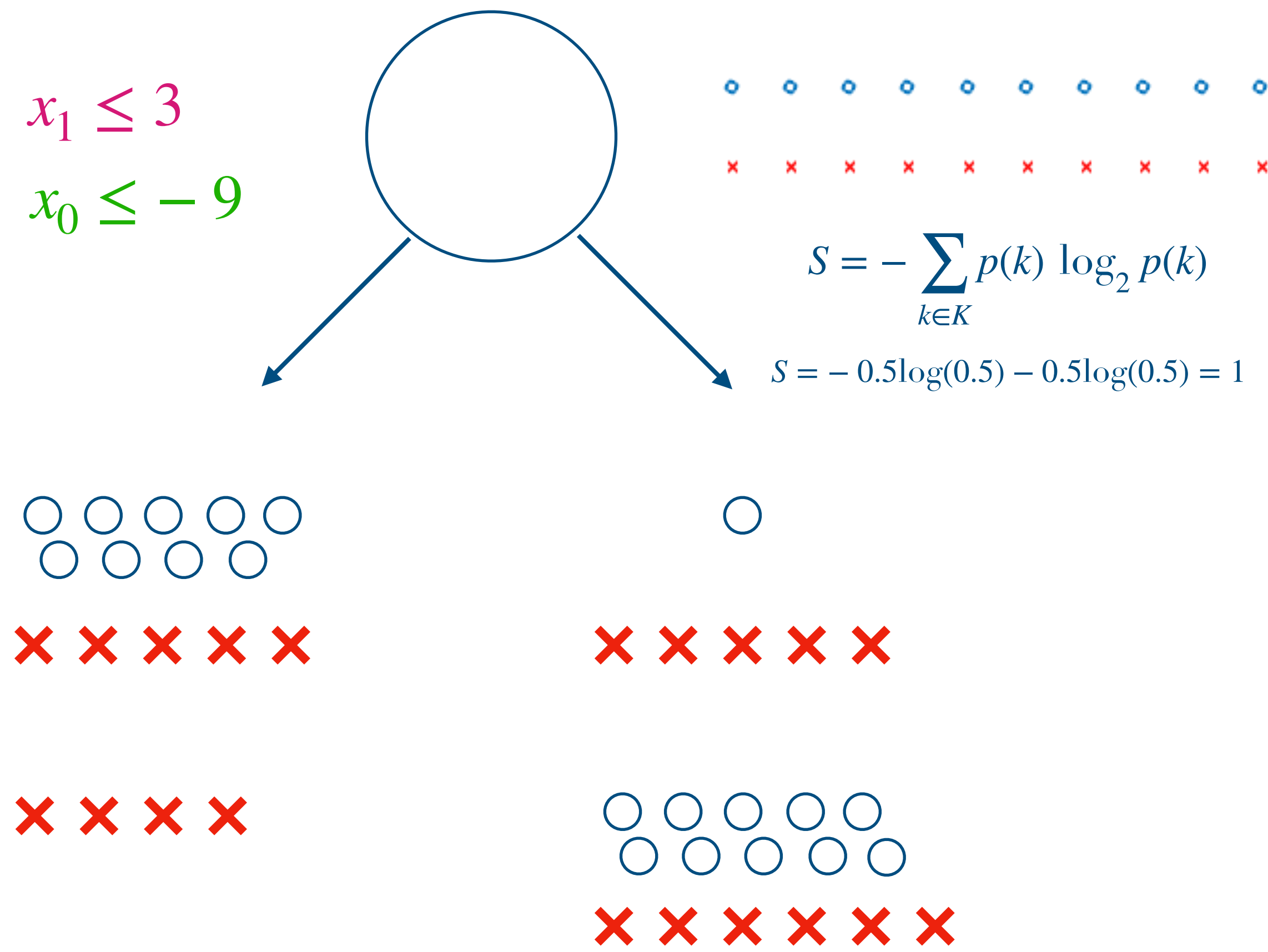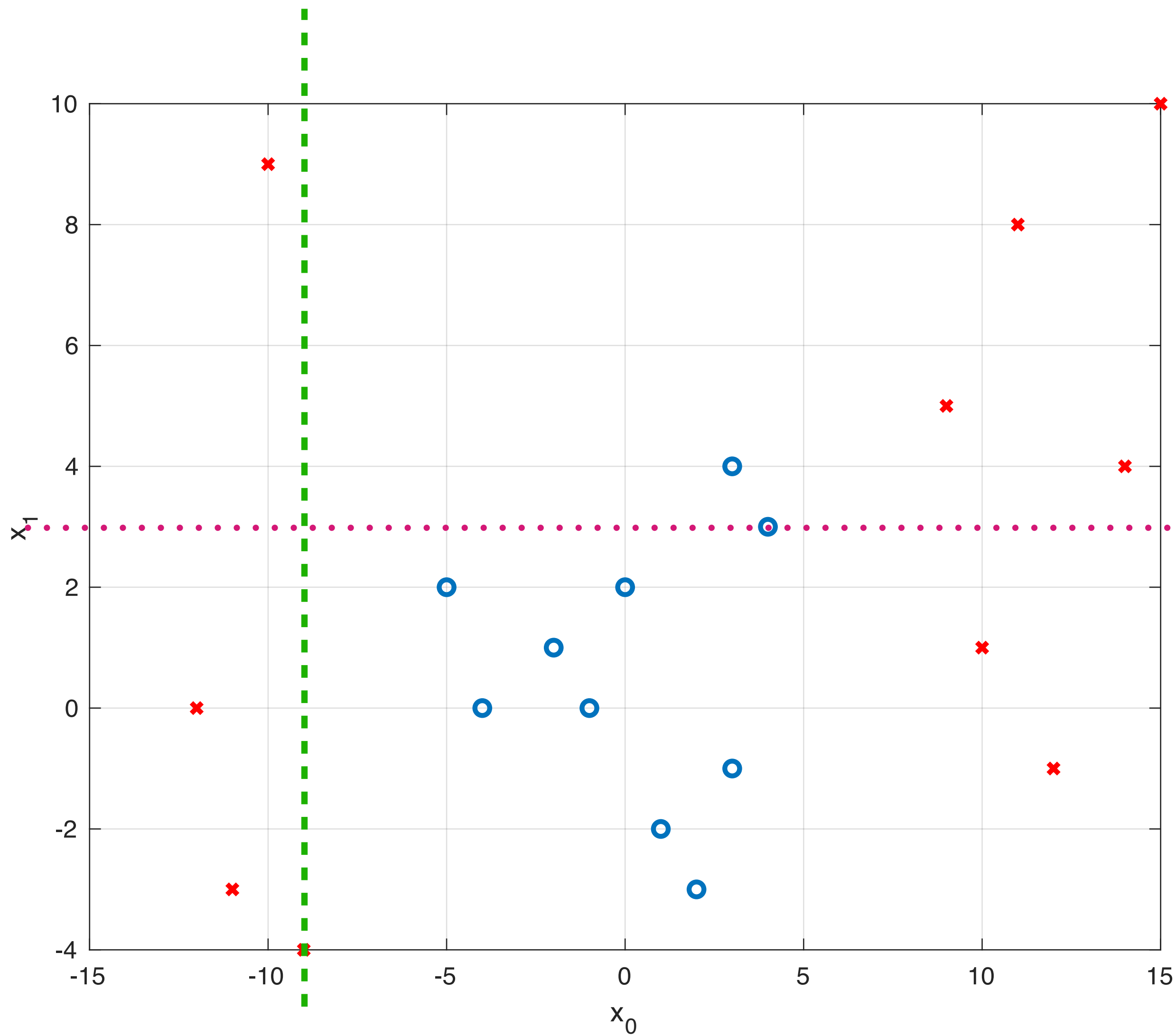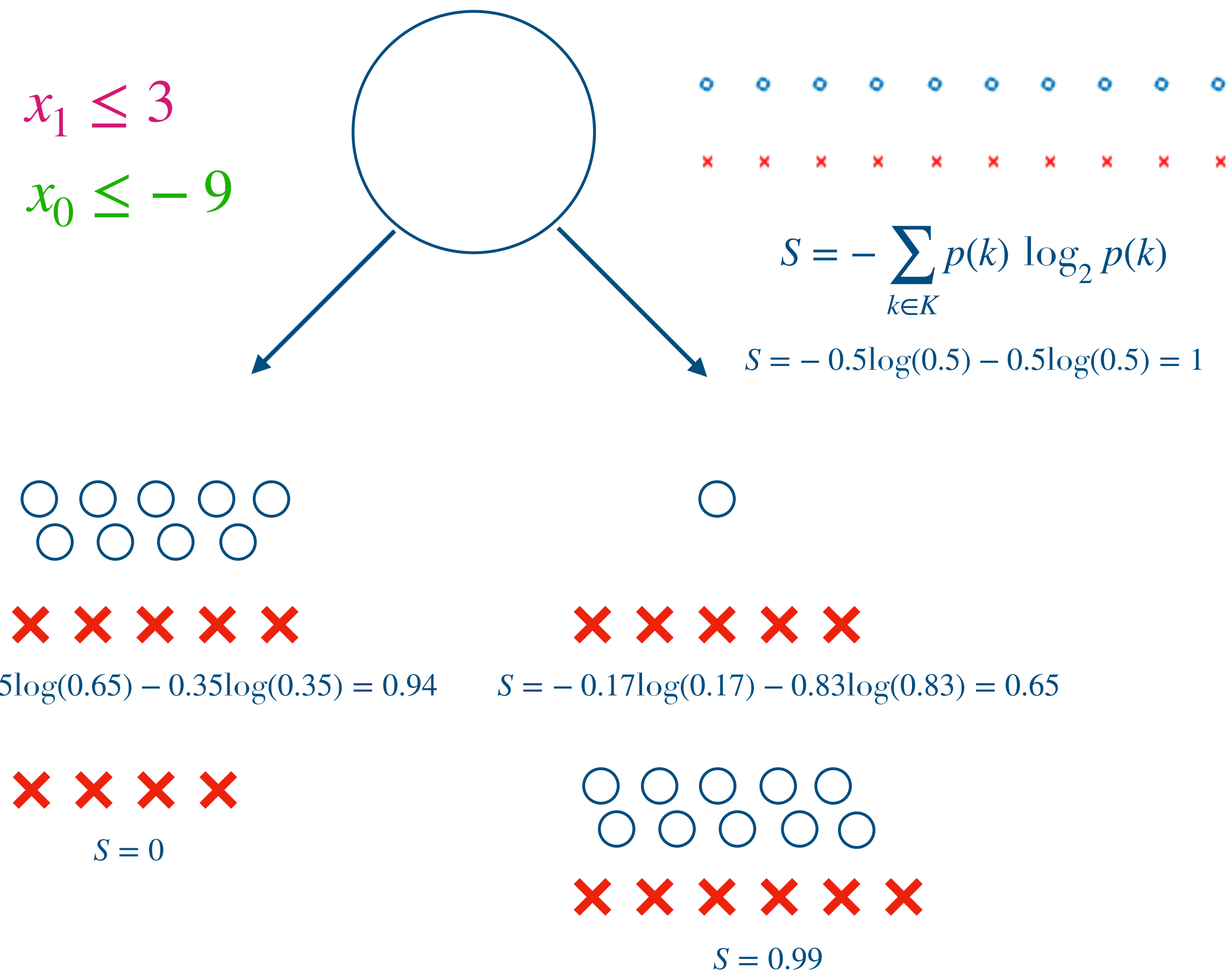The attribute with the highest information gain will produce the best split

# Decision trees: Example

# Decision trees: Example



$x_1 \leq 3$
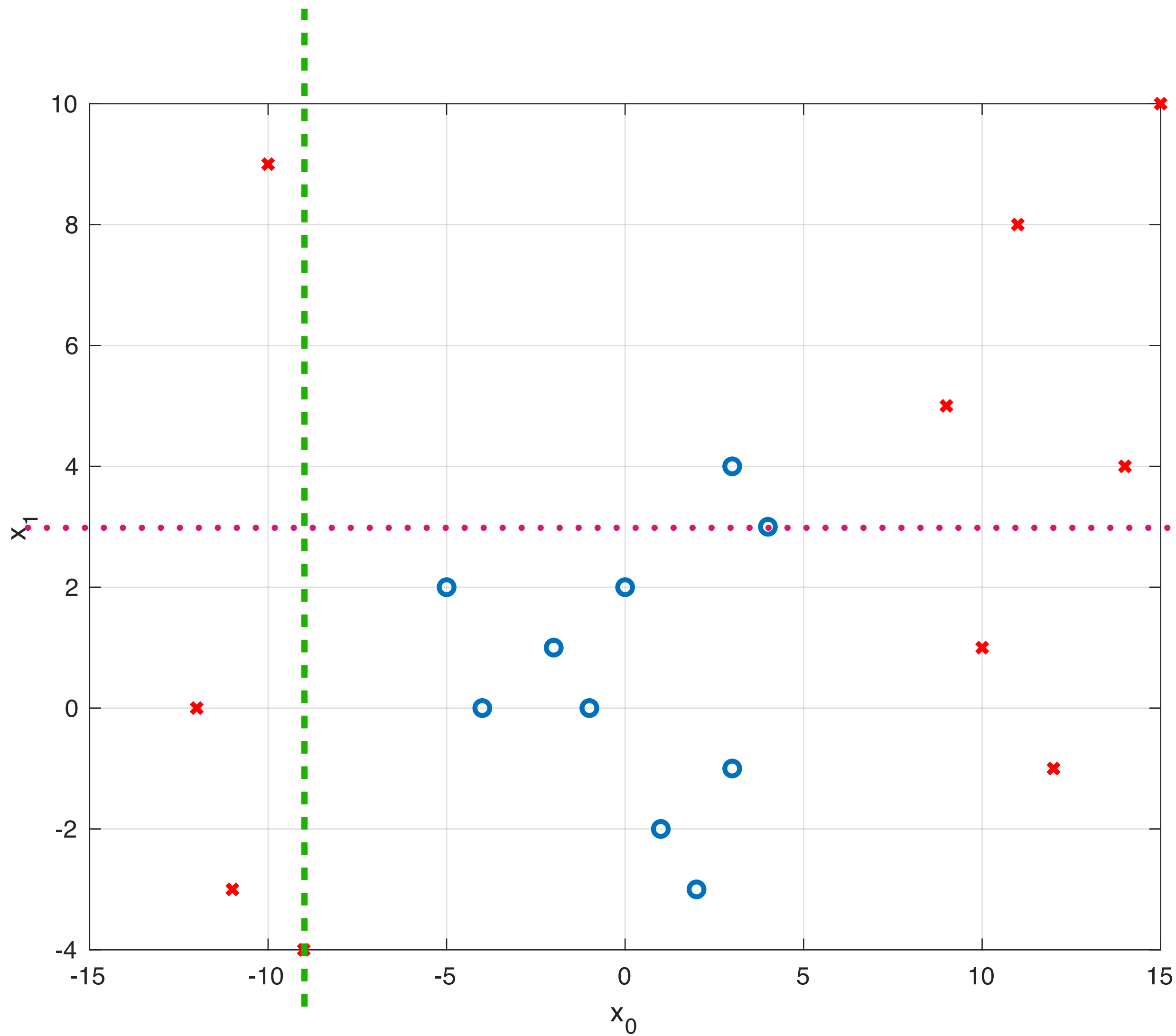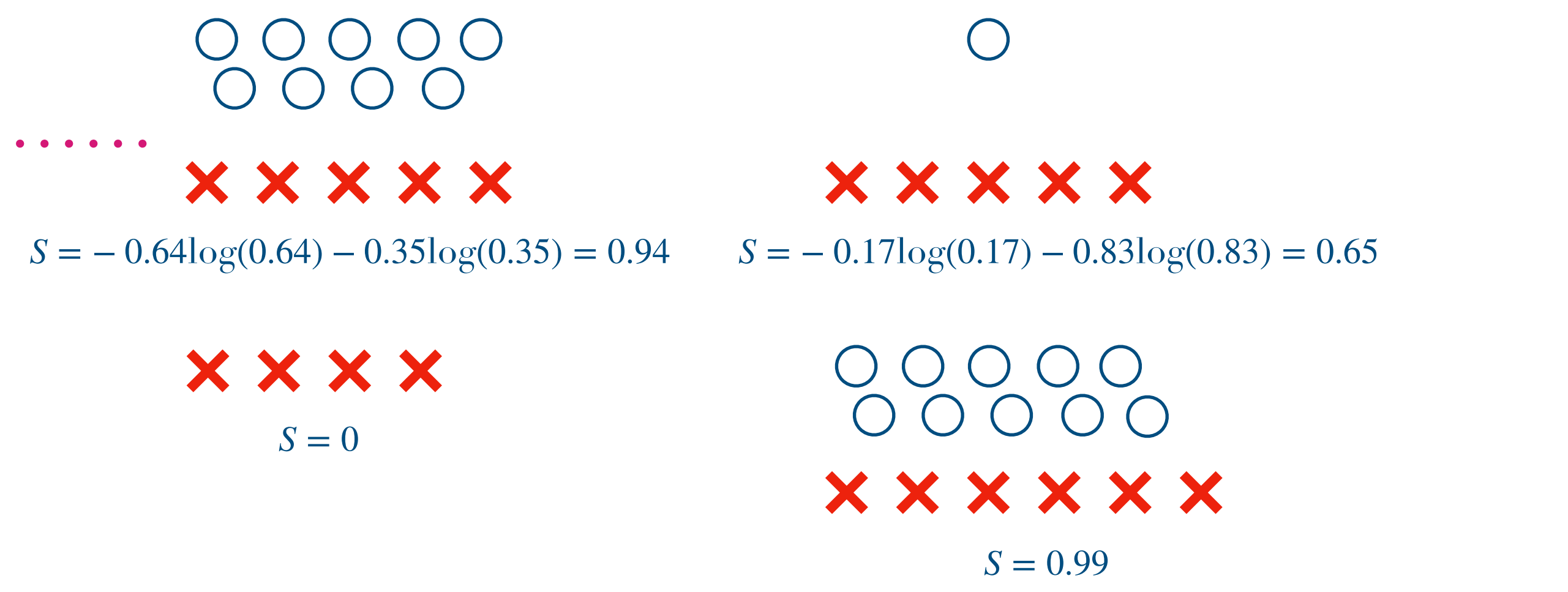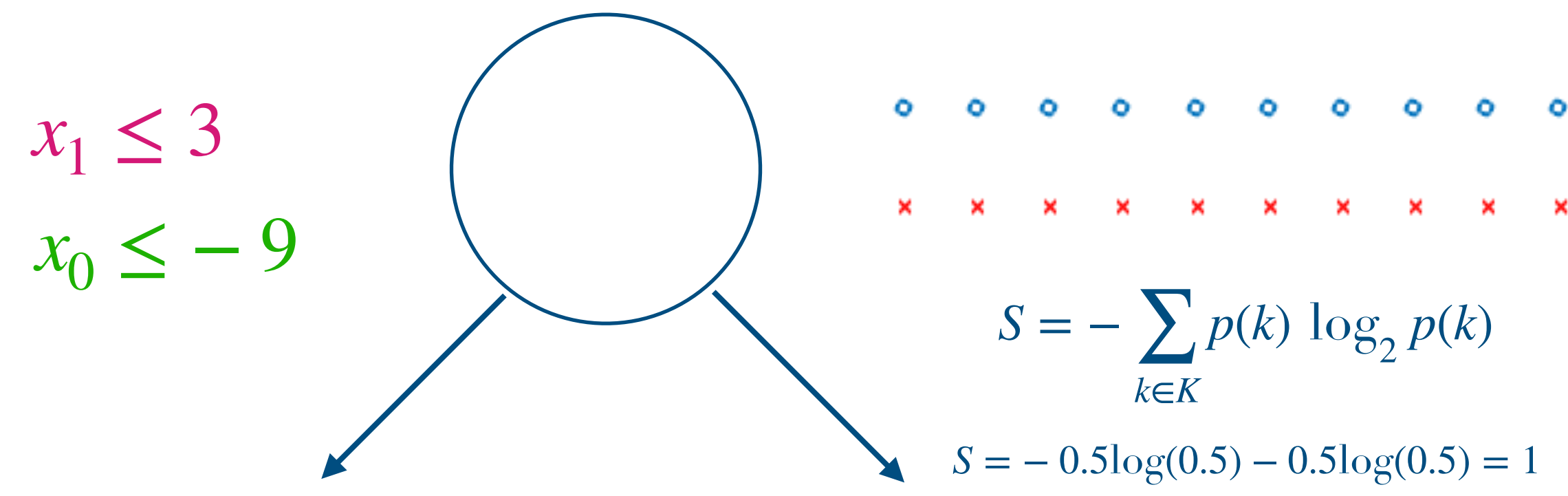
$x_0 \leq -9$

# Decision trees: Example

# Decision trees: Example



$x_1 \leq 3$
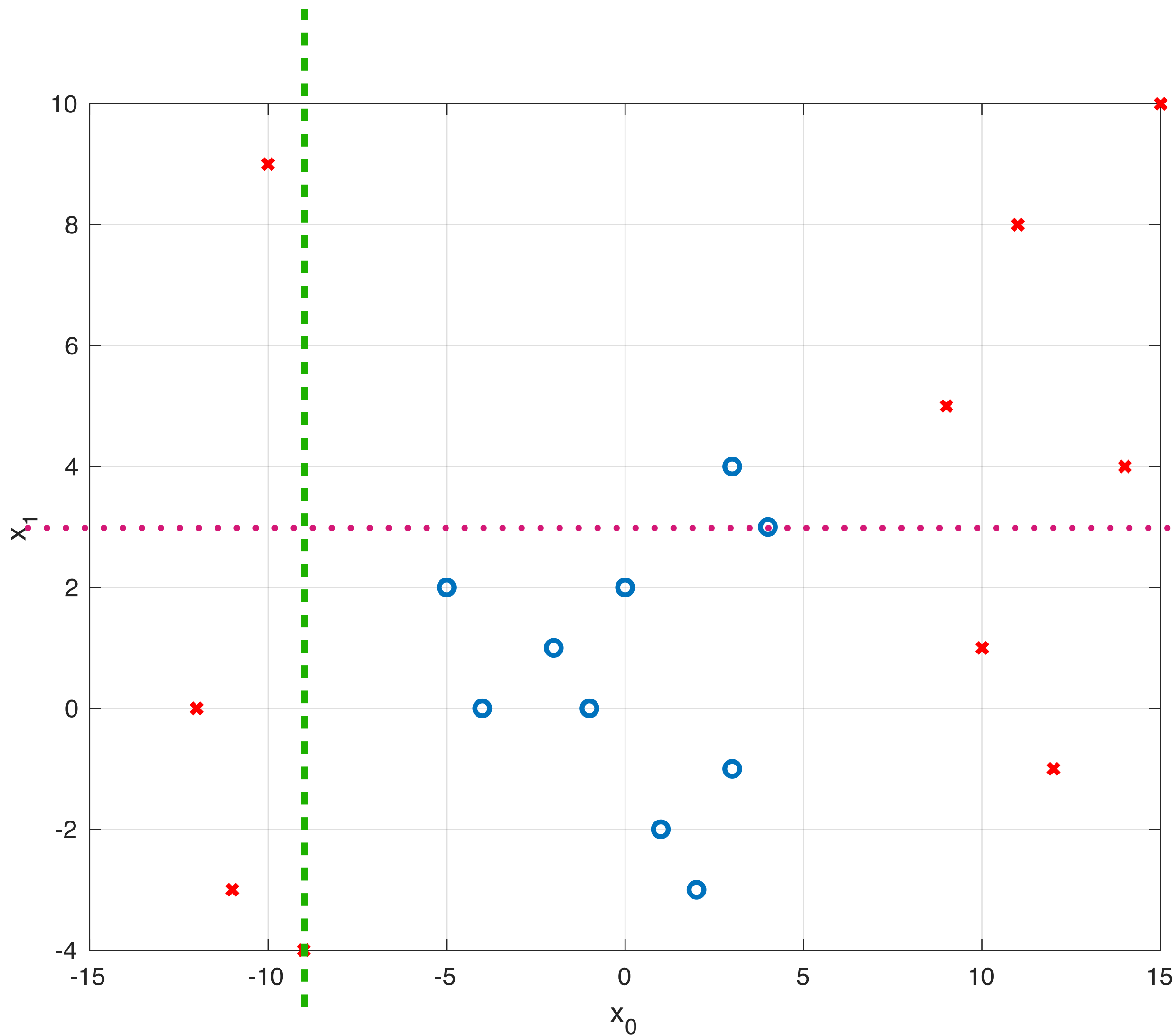
$x_0 \leq -9$

$$S = -\sum_{k \in K} p(k) \log_2 p(k)$$

$S = -0.5\log(0.5) - 0.5\log(0.5) = 1$

# Decision trees: Example



$x_1 \leq 3$

$x_0 \leq -9$

$$S = -\sum_{k \in K} p(k) \log_2 p(k)$$

$S = -0.5\log(0.5) - 0.5\log(0.5) = 1$

$S = -0.65\log(0.65) - 0.35\log(0.35) = 0.94$

$S = -0.17\log(0.17) - 0.83\log(0.83) = 0.65$

$S = 0$

$S = 0.99$

# Decision trees: Example



$x_1 \leq 3$

$x_0 \leq -9$

$$S = -\sum_{k \in K} p(k) \log_2 p(k)$$

$S = -0.5\log(0.5) - 0.5\log(0.5) = 1$

$S = -0.64\log(0.64) - 0.35\log(0.35) = 0.94$

$S = -0.17\log(0.17) - 0.83\log(0.83) = 0.65$
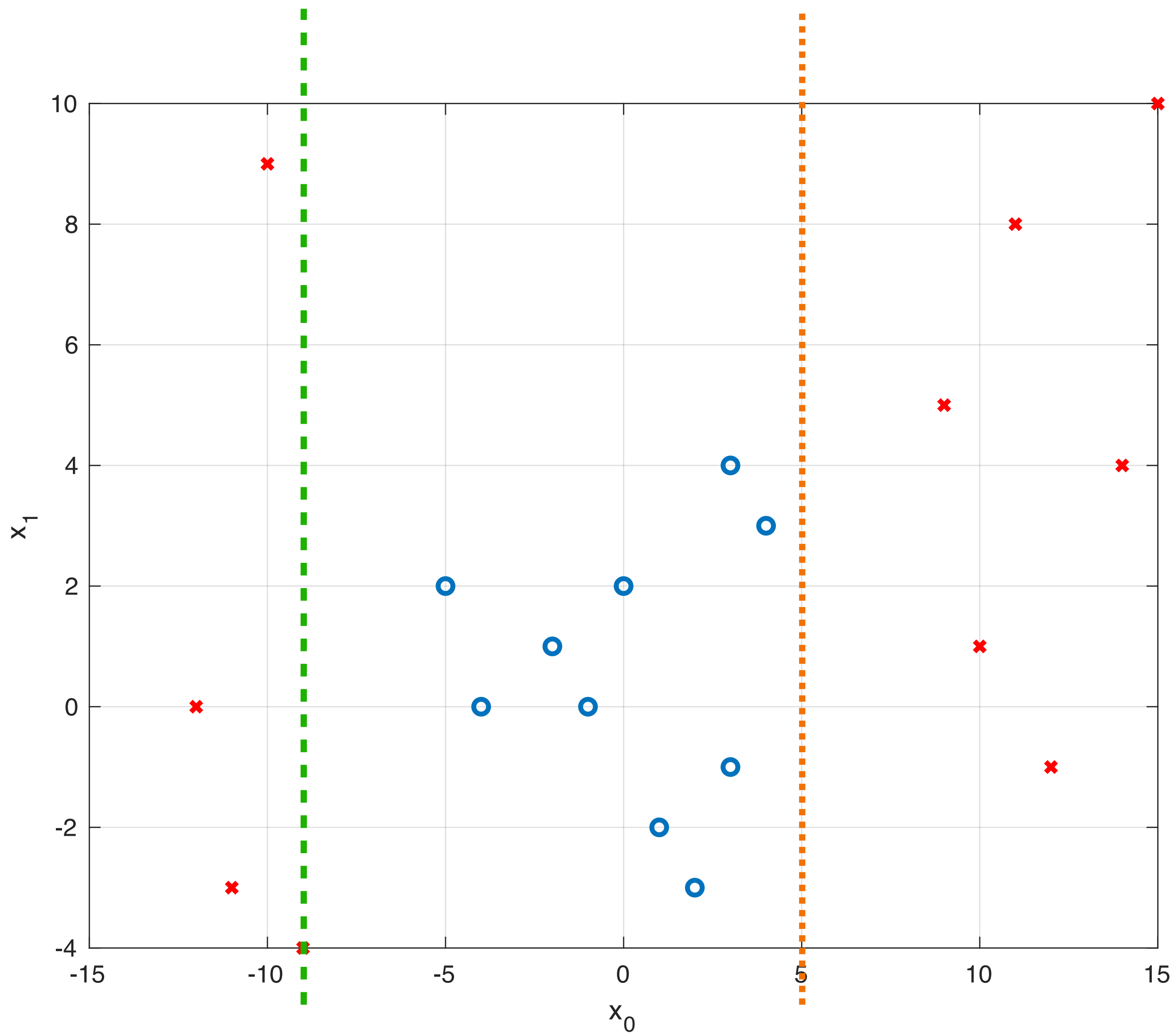
$S = 0$

$S = 0.99$

$$\Theta = S(p) - \sum w_i S(c_i)$$

$$\Theta_1 = 1 - \frac{14}{20} * 0.94 - \frac{6}{20} * 0.65 = 0.14$$

$$\Theta_2 = 1 - \frac{4}{20} * 0 - \frac{16}{20} * 0.99 = 0.2$$

# Decision trees: Example



$x_1 \leq 3$

$x_0 \leq -9$

$$S = -\sum_{k \in K} p(k) \log_2 p(k)$$

$S = -0.5\log(0.5) - 0.5\log(0.5) = 1$

$S = -0.64\log(0.64) - 0.35\log(0.35) = 0.94$

$S = -0.17\log(0.17) - 0.83\log(0.83) = 0.65$

$S = 0$

$S = 0.99$

$\Theta_2 > \Theta_1$

# Decision trees: Example

# Decision trees

## Advantages

Simple to understand: visual representations of decision trees make them easier to understand

Little to no data preparation

Flexible: can be leveraged for both classification and regression tasks

## Disadvantages

Prone to overfitting: Complex decision trees tend to overfit

High variance estimator: Small variations within data can produce a very different decision tree

# Decision trees in Matlab

Classification and regression tree (CART): Breiman et al., 1984

- Decision trees, or classification trees and regression trees, predict responses to data in Matlab

- Allows you to predict a response following the decisions in the tree from the root (beginning) node down to a leaf node

- The leaf node contains the response

- Classification trees give nominal responses, such as `'true'` or `'false'`

- Regression trees give numeric responses.

# Decision trees in Matlab

# Decision trees in Matlab