

## Lecture 06

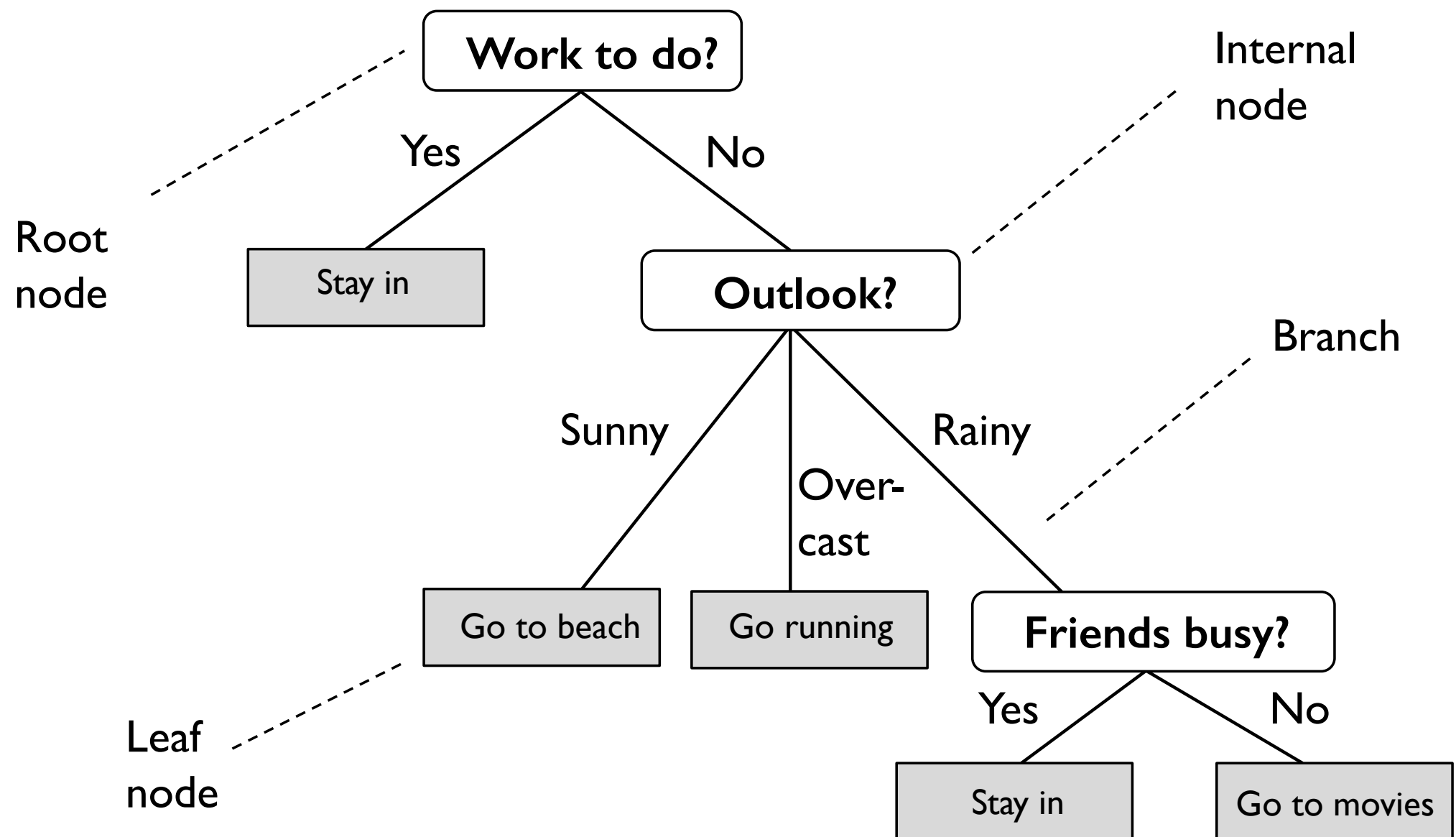
# Decision Trees

STAT 479: Machine Learning, Fall 2019

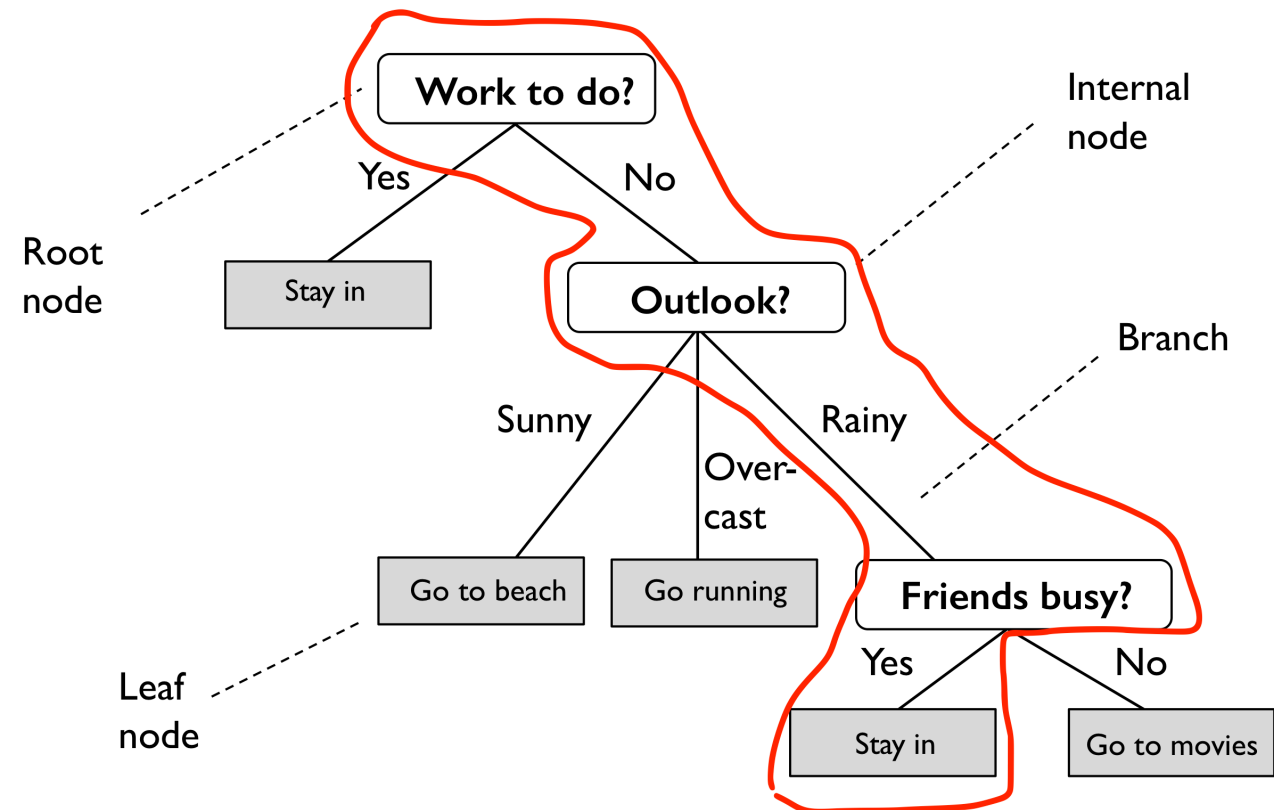
Sebastian Raschka

<http://stat.wisc.edu/~sraschka/teaching/stat479-fs2019/>

# Decision Tree Terminology



# Decision Trees as Rulesets



**IF**

\_\_\_\_\_

\_\_\_\_\_

**THEN**

\_\_\_\_\_

# Decision Trees and ML Categories

- Supervised vs. unsupervised learning algorithm
- Classification vs. regression
- Optimization method: \_\_\_\_\_
- Eager vs. lazy learning algorithm
- Batch vs. online learning algorithm
- Parametric vs. nonparametric model
- Deterministic vs. stochastic

# Recursion / Recursive Algorithms

```
1 def some_fun(x):  
2     if x == []:  
3         return 0  
4     else:  
5         return 1 + some_fun(x[1:])
```

What does this function do?

# Divide & Conquer Algorithms: Quicksort

```
1 def quicksort(array):
2     if len(array) < 2:
3         return array
4     else:
5         pivot = array[0]
6         smaller, bigger = [], []
7         for ele in array[1:]:
8             if ele <= pivot:
9                 smaller.append(ele)
10            else:
11                bigger.append(ele)
12        return quicksort(smaller) + [pivot] + quicksort(bigger)
```

# Divide & Conquer Algorithms: Quicksort

```
1 def quicksort(array):
2     if len(array) < 2:
3         return array
4     else:
5         pivot = array[0]
6         smaller, bigger = [], []
7         for ele in array[1:]:
8             if ele <= pivot:
9                 smaller.append(ele)
10            else:
11                bigger.append(ele)
12     return quicksort(smaller) + [pivot] + quicksort(bigger)
```



# Time complexity of quicksort:

$\mathcal{O}(\underline{\hspace{2cm}})$  ("on average")

```
1 def quicksort(array):
2     if len(array) < 2:
3         return array
4     else:
5         pivot = array[0]
6         smaller, bigger = [], []
7         for ele in array[1:]:
8             if ele <= pivot:
9                 smaller.append(ele)
10            else:
11                bigger.append(ele)
12        return quicksort(smaller) + [pivot] + quicksort(bigger)
```



# Array Sorting Algorithms

Algorithm	Time Complexity			Space Complexity
	Best	Average	Worst*	
<u>Quicksort</u>	$\Omega(n \log(n))$	$\Theta(n \log(n))$	$O(n^2)$	$O(\log(n))$
<u>Mergesort</u>	$\Omega(n \log(n))$	$\Theta(n \log(n))$	$O(n \log(n))$	$O(n)$
<u>Timsort</u>	$\Omega(n)$	$\Theta(n \log(n))$	$O(n \log(n))$	$O(n)$
<u>Heapsort</u>	$\Omega(n \log(n))$	$\Theta(n \log(n))$	$O(n \log(n))$	$O(1)$
<u>Bubble Sort</u>	$\Omega(n)$	$\Theta(n^2)$	$O(n^2)$	$O(1)$
<u>Insertion Sort</u>	$\Omega(n)$	$\Theta(n^2)$	$O(n^2)$	$O(1)$
<u>Selection Sort</u>	$\Omega(n^2)$	$\Theta(n^2)$	$O(n^2)$	$O(1)$
<u>Tree Sort</u>	$\Omega(n \log(n))$	$\Theta(n \log(n))$	$O(n^2)$	$O(n)$
<u>Shell Sort</u>	$\Omega(n \log(n))$	$\Theta(n(\log(n))^2)$	$O(n(\log(n))^2)$	$O(1)$
<u>Bucket Sort</u>	$\Omega(n+k)$	$\Theta(n+k)$	$O(n^2)$	$O(n)$
<u>Radix Sort</u>	$\Omega(nk)$	$\Theta(nk)$	$O(nk)$	$O(n+k)$
<u>Counting Sort</u>	$\Omega(n+k)$	$\Theta(n+k)$	$O(n+k)$	$O(k)$
<u>Cubesort</u>	$\Omega(n)$	$\Theta(n \log(n))$	$O(n \log(n))$	$O(n)$

<http://www.bigocheatsheet.com>

\* "worst" ~ inversely-sorted array

# Time Complexity ("Big-O")

Growing the tree:  $O(\dots)$

Tip: It can be shown that optimal split is on boundary between adjacent examples (similar feature value) with different class labels.

Fayyad, Usama Mohammad. "On the induction of decision trees for multiple concept learning." (1992).

# Time Complexity ("Big-O")

Querying the tree:  $\mathcal{O}(\dots)$

# Decision Tree in Pseudocode

GenerateTree( $\mathcal{D}$ ):

- if  $y = 1 \ \forall \ \langle \mathbf{x}, \mathbf{y} \rangle \in \mathcal{D}$  or  $y = 0 \ \forall \ \langle \mathbf{x}, \mathbf{y} \rangle \in \mathcal{D}$ :
  - return Tree
- else:
  - Pick best feature  $x_j$ :
    - $\mathcal{D}_0$  at Child<sub>0</sub> :  $x_j = 0 \ \forall \ \langle \mathbf{x}, \mathbf{y} \rangle \in \mathcal{D}$
    - $\mathcal{D}_1$  at Child<sub>1</sub> :  $x_j = 1 \ \forall \ \langle \mathbf{x}, \mathbf{y} \rangle \in \mathcal{D}$

return Node( $x_j$ , GenerateTree( $\mathcal{D}_0$ ), GenerateTree( $\mathcal{D}_1$ ))

# Generic Tree Growing Algorithm

- 1)** Pick the feature that, when parent node is split, results in the largest information gain
- 2)** Stop if child nodes are pure or information gain  $\leq 0$
- 3)** Go back to step 1 for each of the two child nodes

# Generic Tree Growing Algorithm

- 1) Pick the feature that, when parent node is split, results in the largest information gain
  - 2) Stop if child nodes are pure or information gain  $\leq 0$
  - 3) Go back to step 1 for each of the two child nodes
- How make predictions of features in dataset not sufficient to make child nodes pure?

# Design choices

- How to split
  - what measurement/criterion as measure of goodness
  - binary vs multi-category split
- When to stop
  - if leaf nodes contain only examples of the same class
  - feature values are all the same for all examples
  - statistical significance test

# ID3 -- Iterative Dichotomizer 3

- one of the earlier/earliest decision tree algorithms
- Quinlan, J. R. 1986. Induction of Decision Trees. Mach. Learn. 1, 1 (Mar. 1986), 81-106.
- cannot handle numeric features
- no pruning, prone to overfitting
- short and wide trees (compared to CART)
- maximizing information gain/minimizing entropy
- discrete features, binary and multi-category features



# C4.5

- continuous and discrete features
- Ross Quinlan 1993, Quinlan, J. R. (1993). C4.5: Programming for machine learning. *Morgan Kauffmann*, 38, 48.
- continuous is very expensive, because must consider all possible ranges
- handles missing attributes (ignores them in gain compute)
- post-pruning (bottom-up pruning)
- Gain Ratio

# CART

- Breiman, L. (1984). *Classification and regression trees*. Belmont, Calif: Wadsworth International Group.
- continuous and discrete features
- strictly binary splits (taller trees than ID3, C4.5)
- binary splits can generate better trees than C4.5, but tend to be larger and harder to interpret; k-attributes has a ways to create a binary partitioning
- variance reduction in regression trees
- Gini impurity, twoing criteria in classification trees
- cost complexity pruning

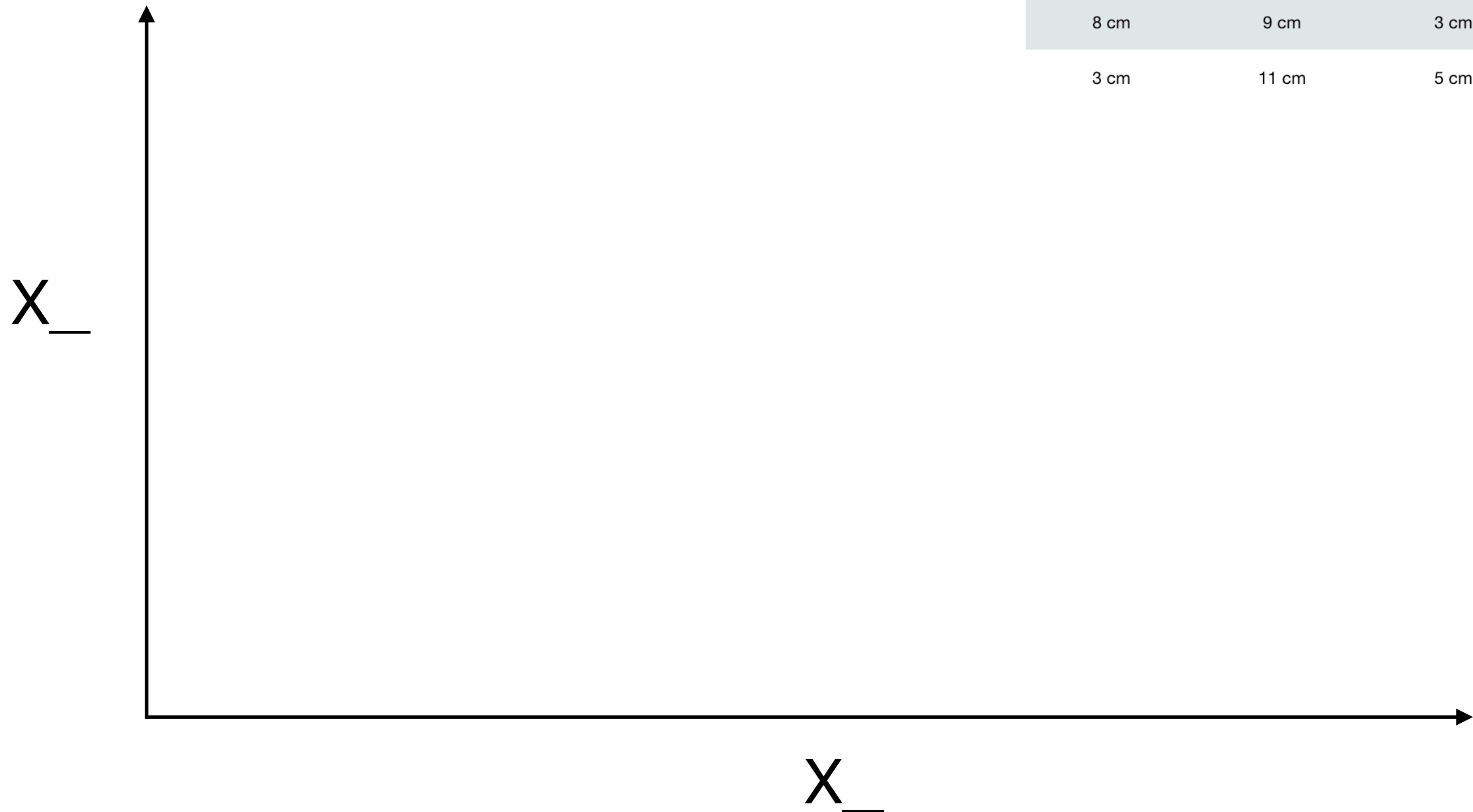
# Others

- CHAID (CHi-squared Automatic Interaction Detector); Kass, G. V. (1980). "An exploratory technique for investigating large quantities of categorical data". *Applied Statistics*. 29 (2): 119–127.
- MARS (Multivariate adaptive regression splines); Friedman, J. H. (1991). "Multivariate Adaptive Regression Splines". *The Annals of Statistics*. 19: 1
- C5.0 (patented)
- ...

# Finding a Decision Rule

$x_1$	$x_2$	$x_3$	$y$
6 cm	8 cm	9 cm	1
4 cm	11 cm	2 cm	0
6 cm	12 cm	4 cm	0
10 cm	9 cm	3 cm	1
5 cm	7 cm	8 cm	0
8 cm	9 cm	3 cm	1
3 cm	11 cm	5 cm	0

# Drawing a Decision Boundary



$x_1$	$x_2$	$x_3$	$y$
6 cm	8 cm	9 cm	1
4 cm	11 cm	2 cm	0
6 cm	12 cm	4 cm	0
10 cm	9 cm	3 cm	1
5 cm	7 cm	8 cm	0
8 cm	9 cm	3 cm	1
3 cm	11 cm	5 cm	0

# The Splitting Criterion

# Information Gain

$$GAIN(\mathcal{D}, x_j) = H(\mathcal{D}) - \sum_{v \in Values(x_j)} \frac{|\mathcal{D}_v|}{|\mathcal{D}|} H(\mathcal{D}_v)$$

# Shannon Entropy

Refer to lecture notes



# Entropy

$$H = - \sum_i p(i | x_j) \log_2(p(i | x_j))$$

# Gini Impurity

$$Gini = 1 - \sum_i (p(i | x_j)^2)$$

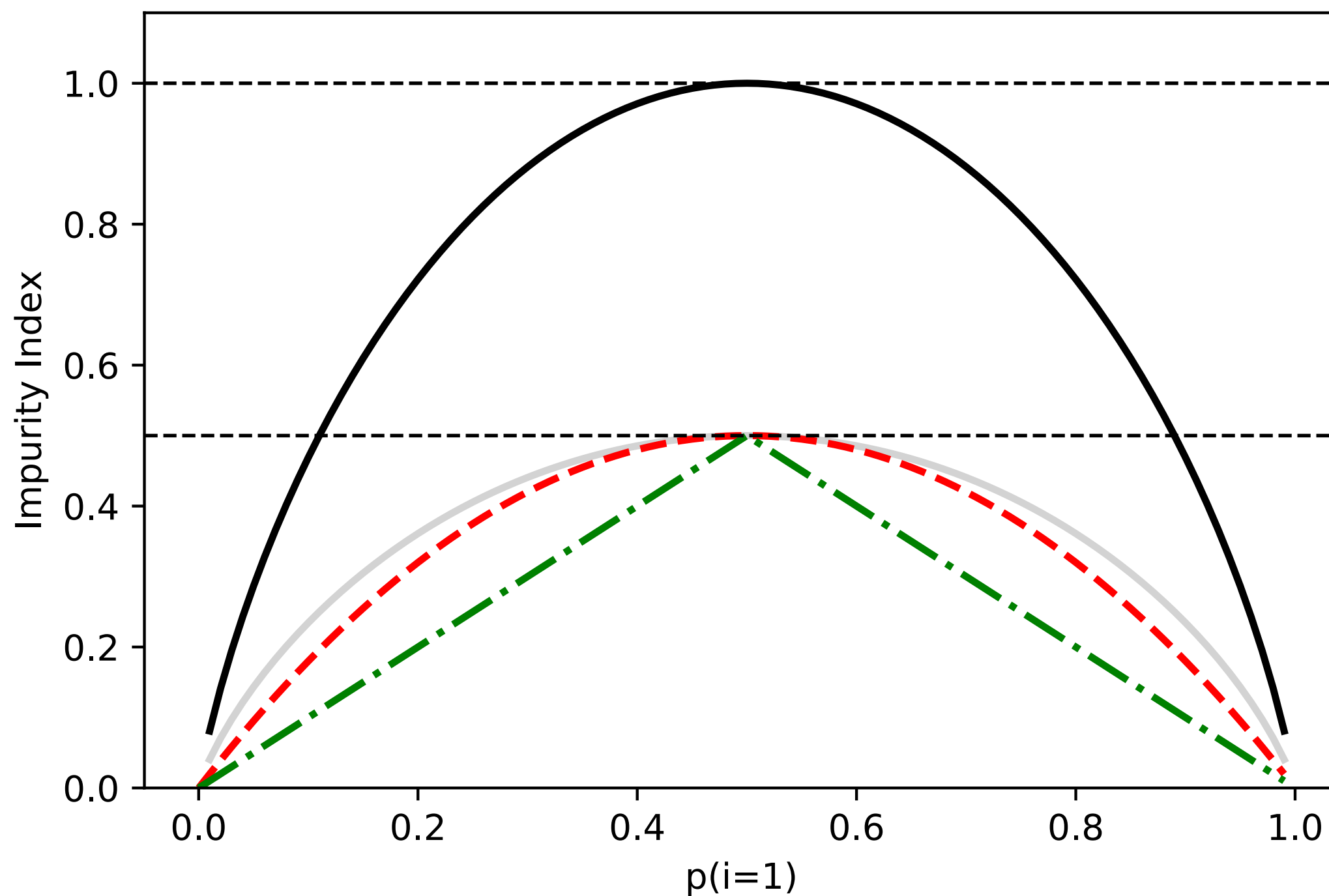
# Misclassification Error

$$ERR = \frac{1}{n} \sum_{i=1}^n L(\hat{y}^{[i]}, y^{[i]}),$$

$$L(\hat{y}, y) = \begin{cases} 0 & \text{if } \hat{y} = y, \\ 1 & \text{otherwise.} \end{cases}$$

# Misclassification Error

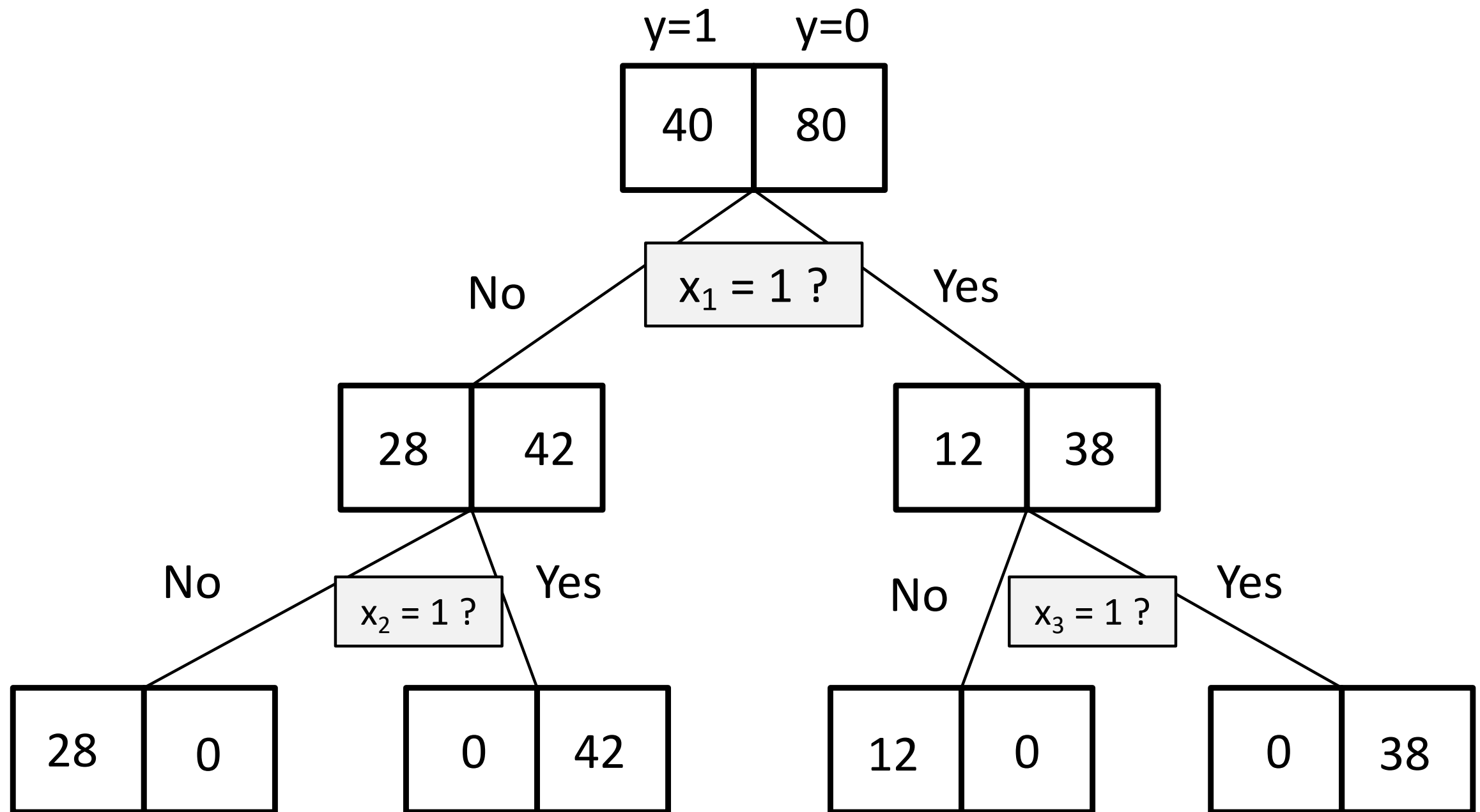
$$ERR = 1 - \max_i(p(i | x_j))$$



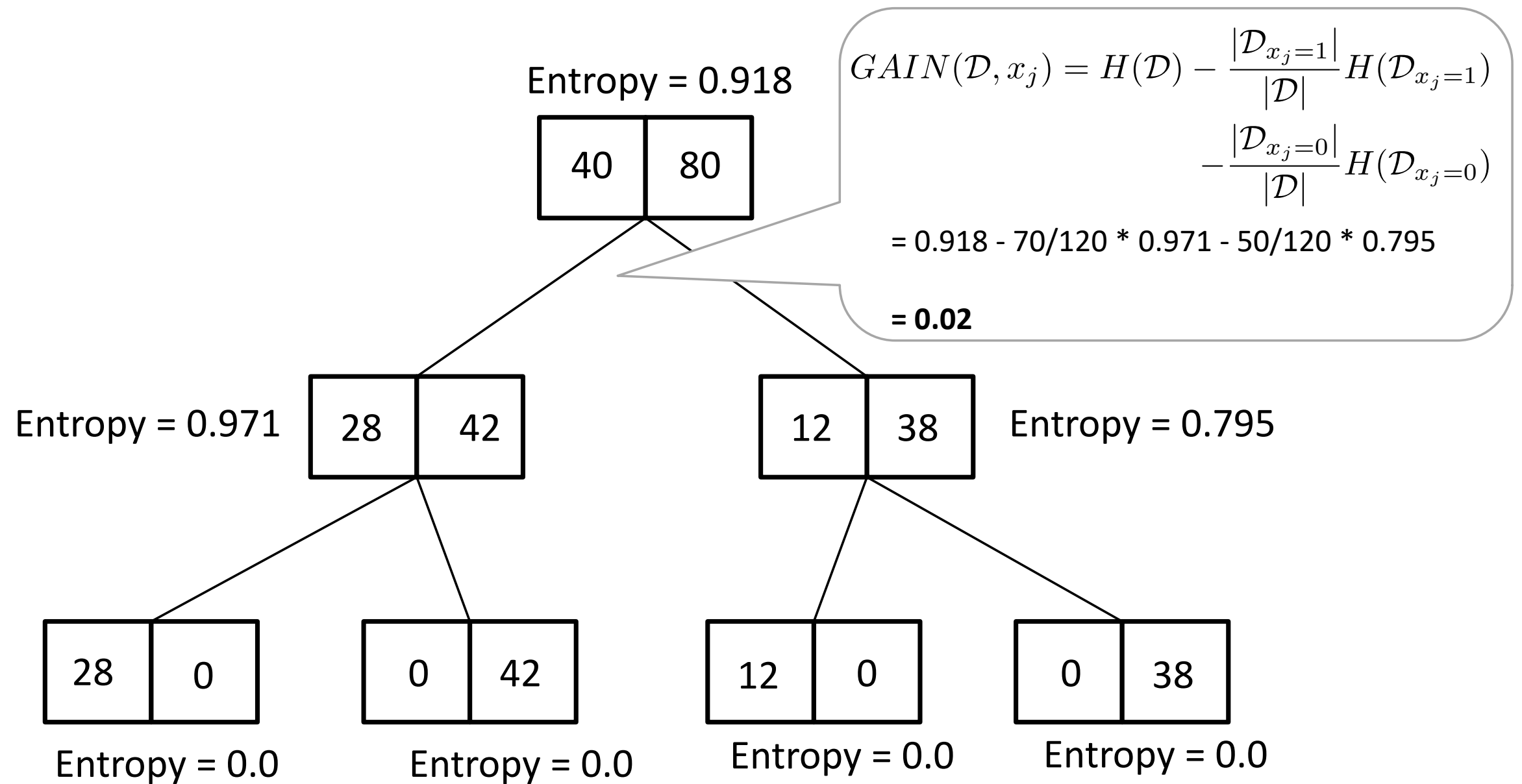
# **Why Growing Decision Trees via Entropy instead of Misclassification Error?**

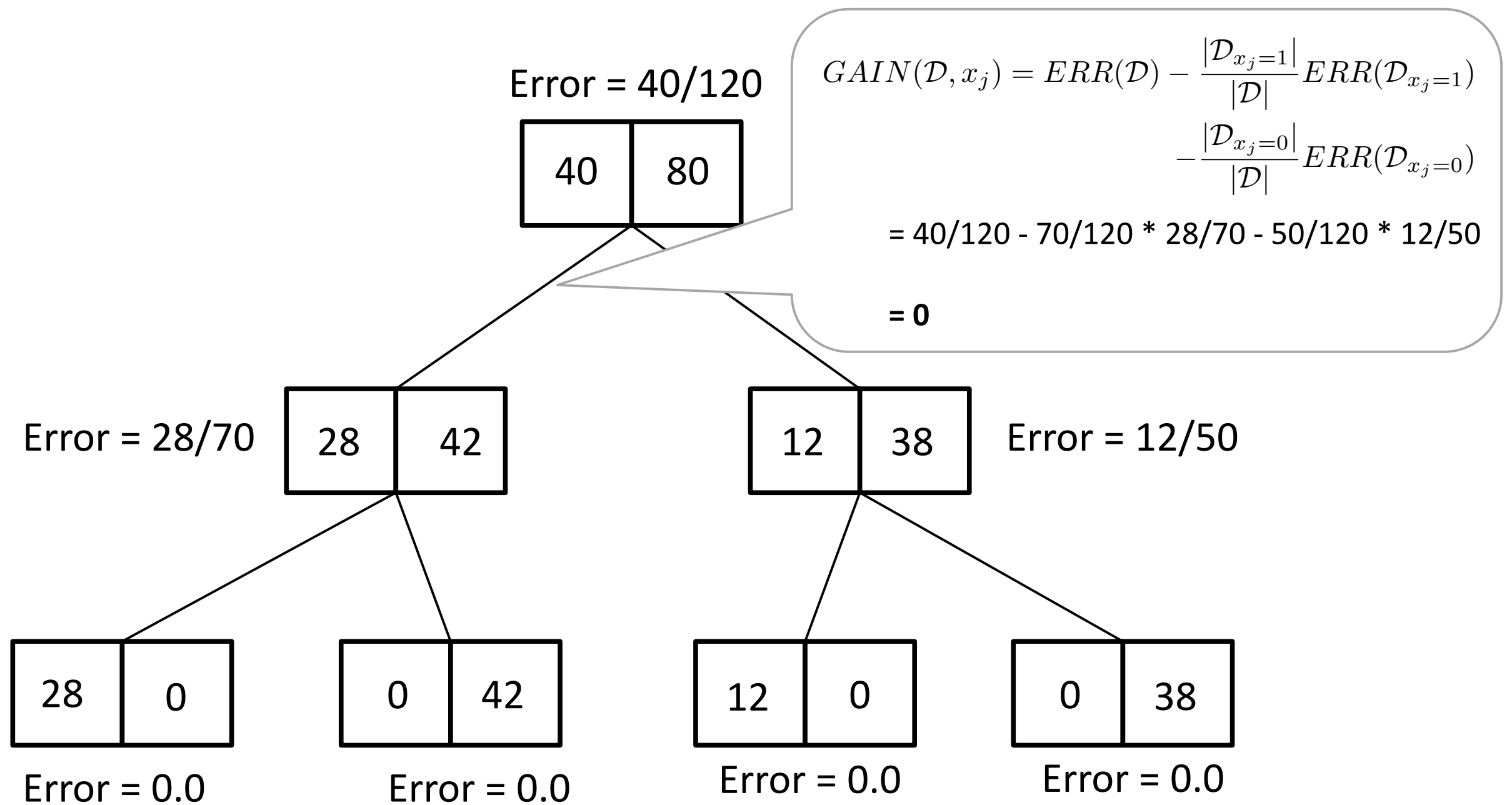
# Why Growing Decision Trees via Entropy instead of Misclassification Error?

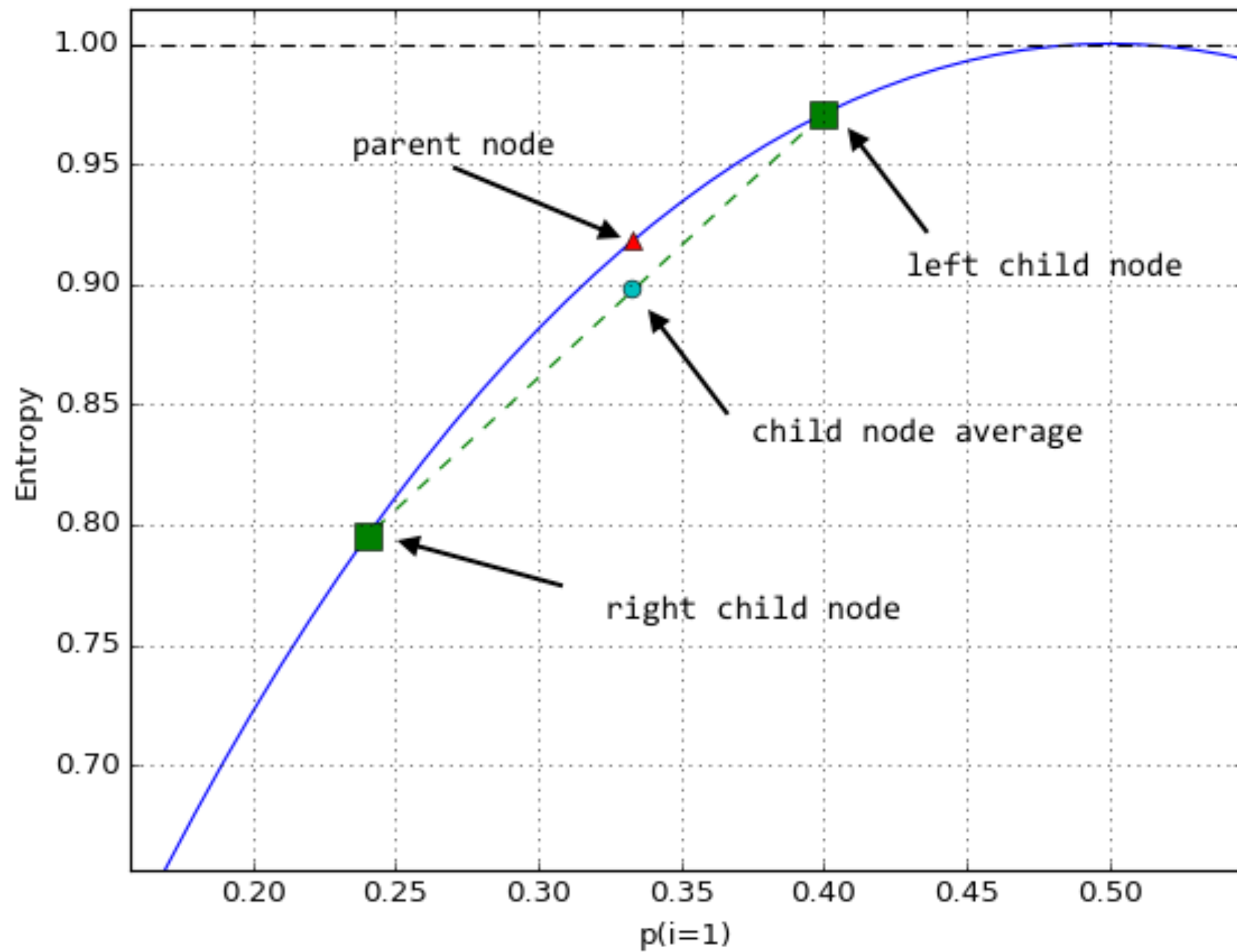
$$GAIN(\mathcal{D}, x_j) = I(\mathcal{D}) - \sum_{v \in Values(x_j)} \frac{|\mathcal{D}_v|}{|\mathcal{D}|} I(\mathcal{D}_v)$$











# Gain Ratio

$$\textit{GainRatio}(\mathcal{D}, x_j) = \frac{\textit{Gain}(\mathcal{D}, x_j)}{\textit{SplitInfo}(\mathcal{D}, x_j)}$$

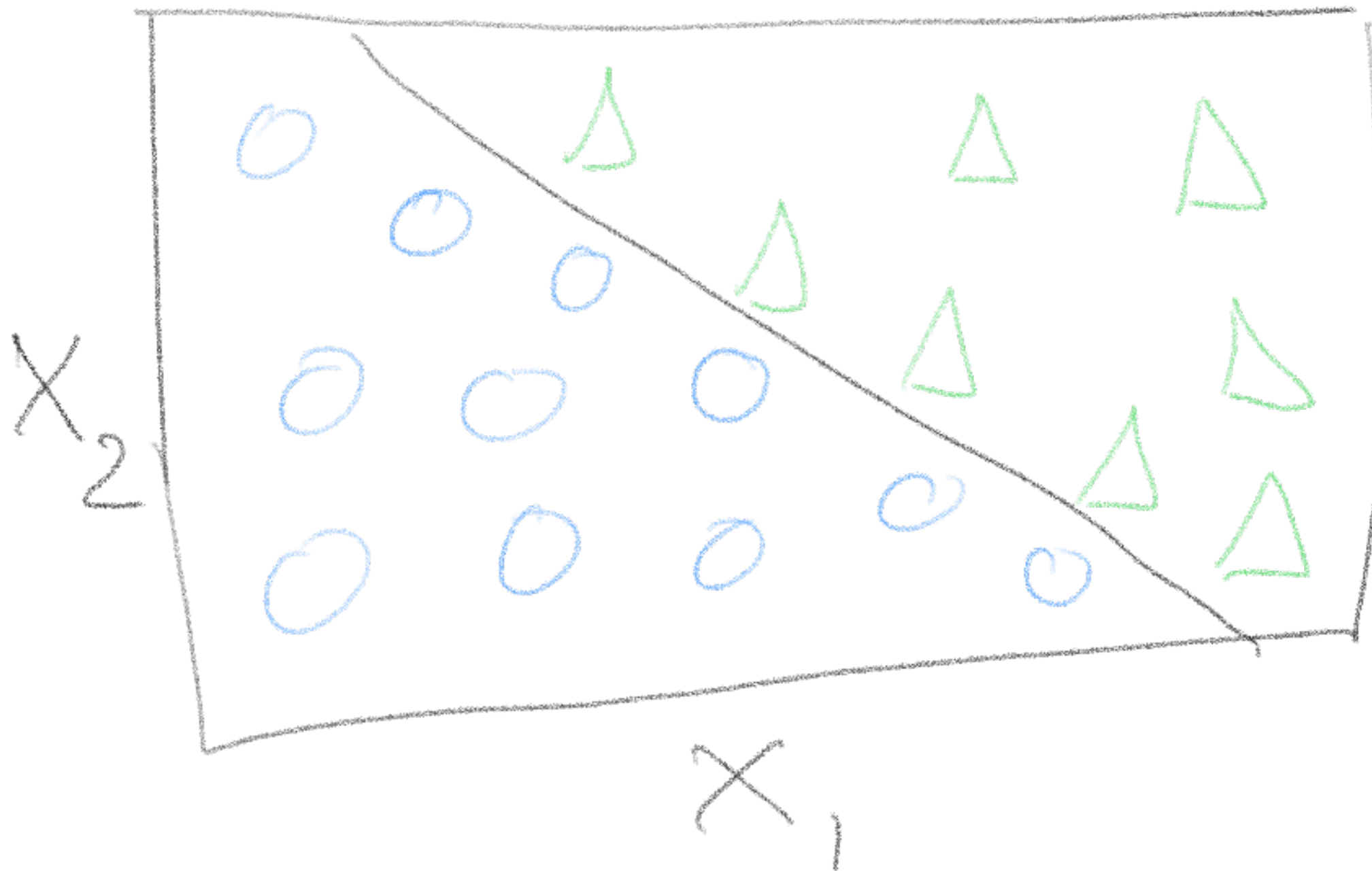
Quinlan 1986

where the split information measures the entropy of the feature:

$$\textit{SplitInfo}(\mathcal{D}, x_j) = - \sum_{v \in x_j} \frac{|\mathcal{D}_v|}{|\mathcal{D}|} \log_2 \frac{|\mathcal{D}_v|}{|\mathcal{D}|}$$

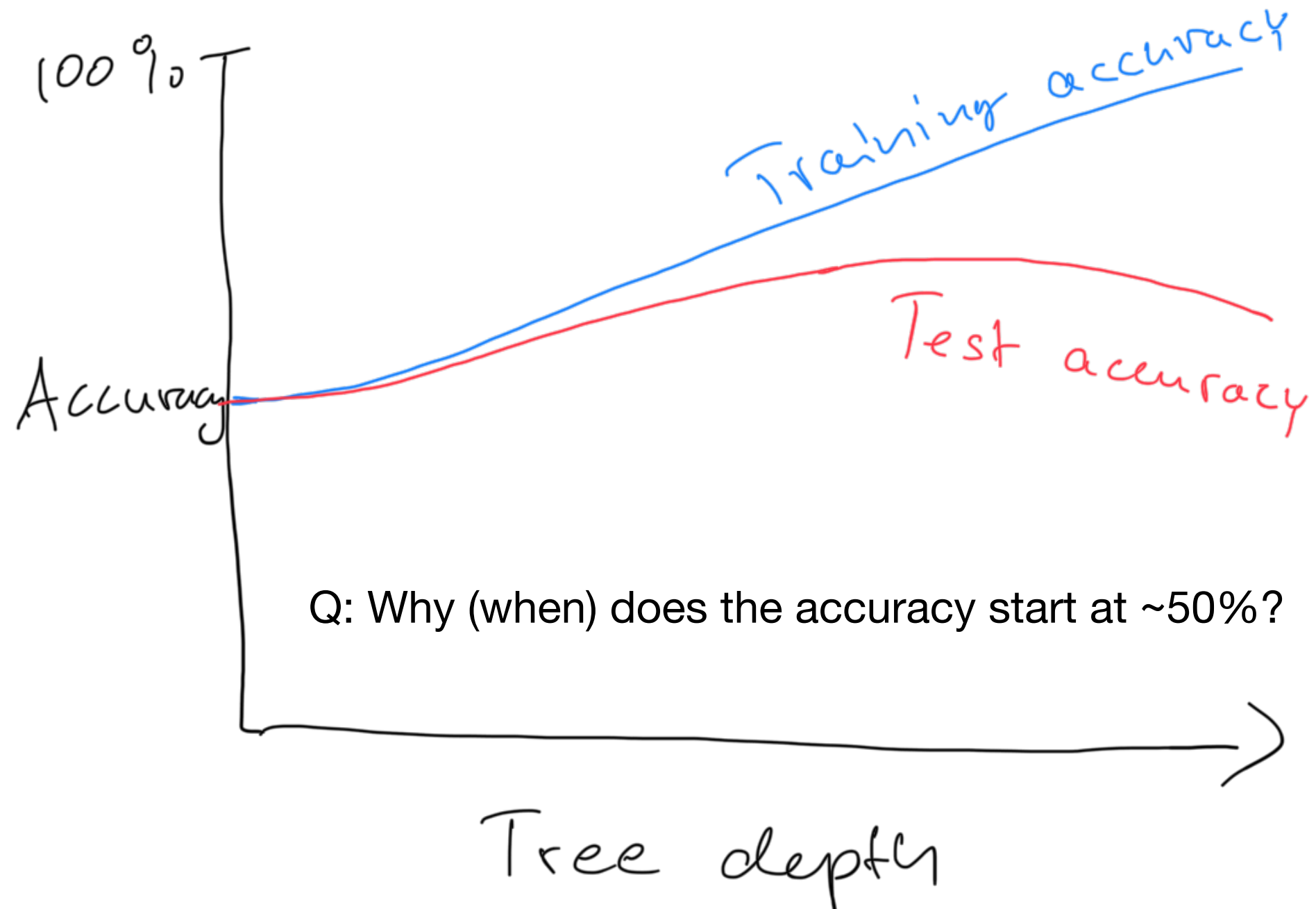
Penalizes splitting categorical attributes with many values (e.g., think date column, or really bad: row ID) via the split information

# Shortcomings



## How would the decision tree split look like?

# Overfitting



# Pre-Pruning

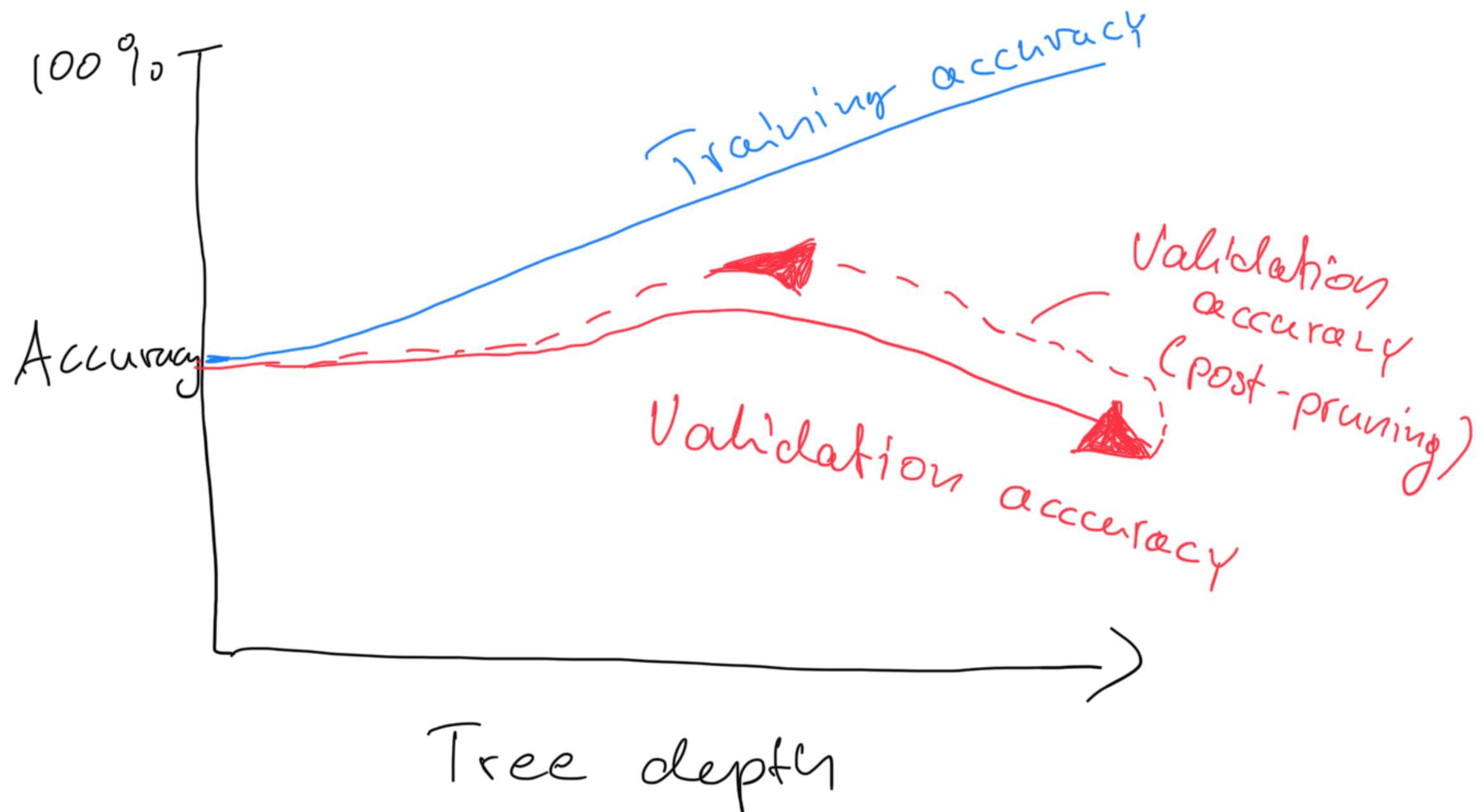
- Set a depth cut-off (maximum tree depth) *a priori*
- Cost-complexity pruning: , where is an impurity measure, is a tuning parameter, and is the total number of nodes.
- Stop growing if split is not statistically significant (e.g.,  $\chi^2$  test)
- Set a minimum number of data points for each node

# Post-Pruning

- Grow full tree first, then remove nodes, in C4.5
- Reduced-error pruning, remove nodes via validation set eval. (problematic for limited data)
- Can also convert trees to rules first and then prune the rules



# Post-Pruning



# Regression Trees

# Decision Tree Summary: Pros and Cons

- (+) Easy to interpret and communicate
- (+) Can represent "complete" hypothesis space
- (-) Easy to overfit
- (-) Elaborate pruning required
- (-) Expensive to just fit a "diagonal line"
- (-) Output range is bounded (dep. on training examples) in regression trees

# Reading Assignments

Python Machine Learning, 2nd Ed., Ch03 pg. 88-104

# Demo

06-trees\_demo.ipynb