# Lecture 06

# Decision Trees

STAT 451: Intro to Machine Learning, Fall 2021

Sebastian Raschka

# Lecture 6: Decision Trees
# Topics

**1. Intro to decision trees**

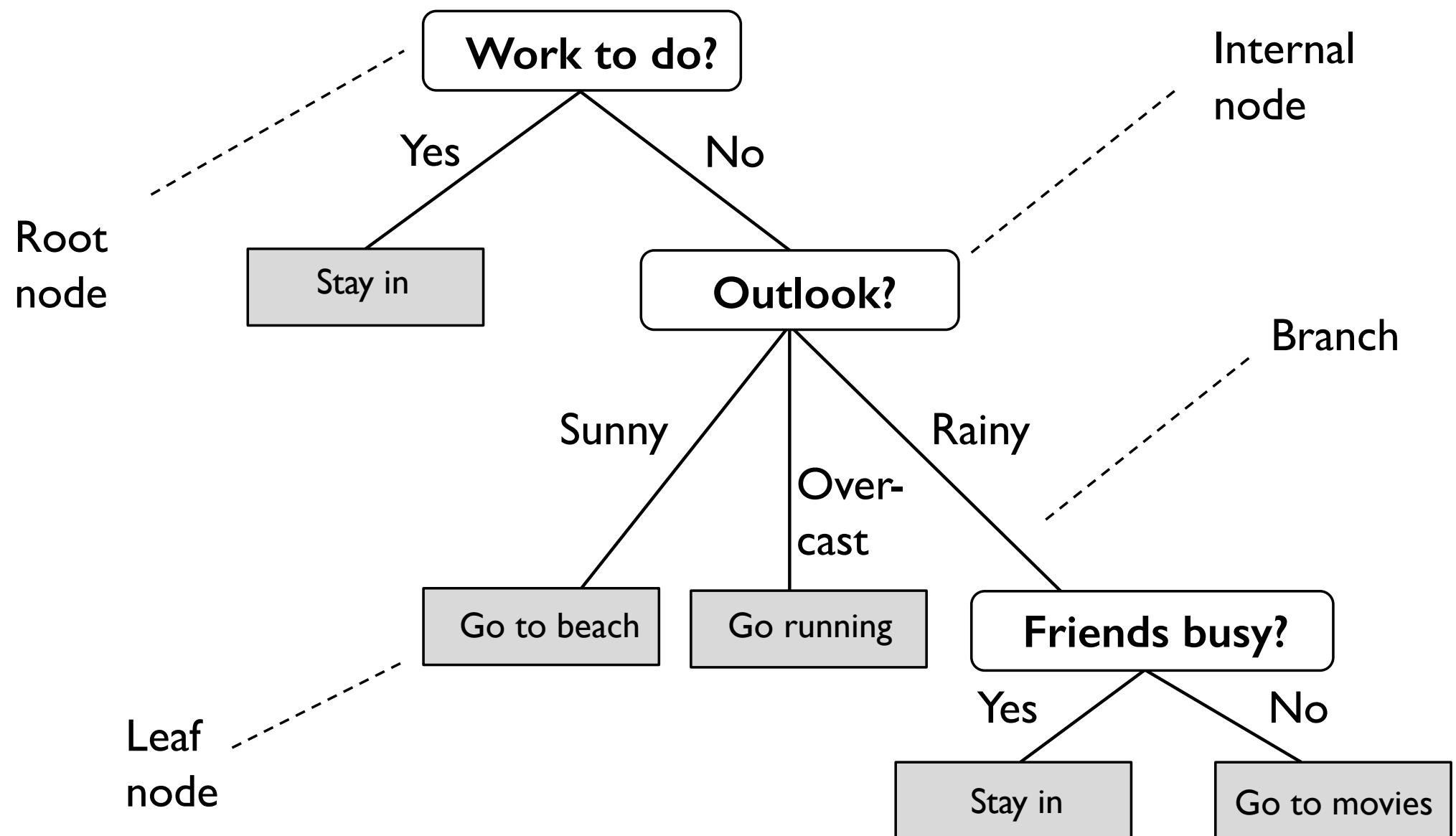2. Recursive algorithms & Big-O

3. Types of decision trees

4. Splitting criteria

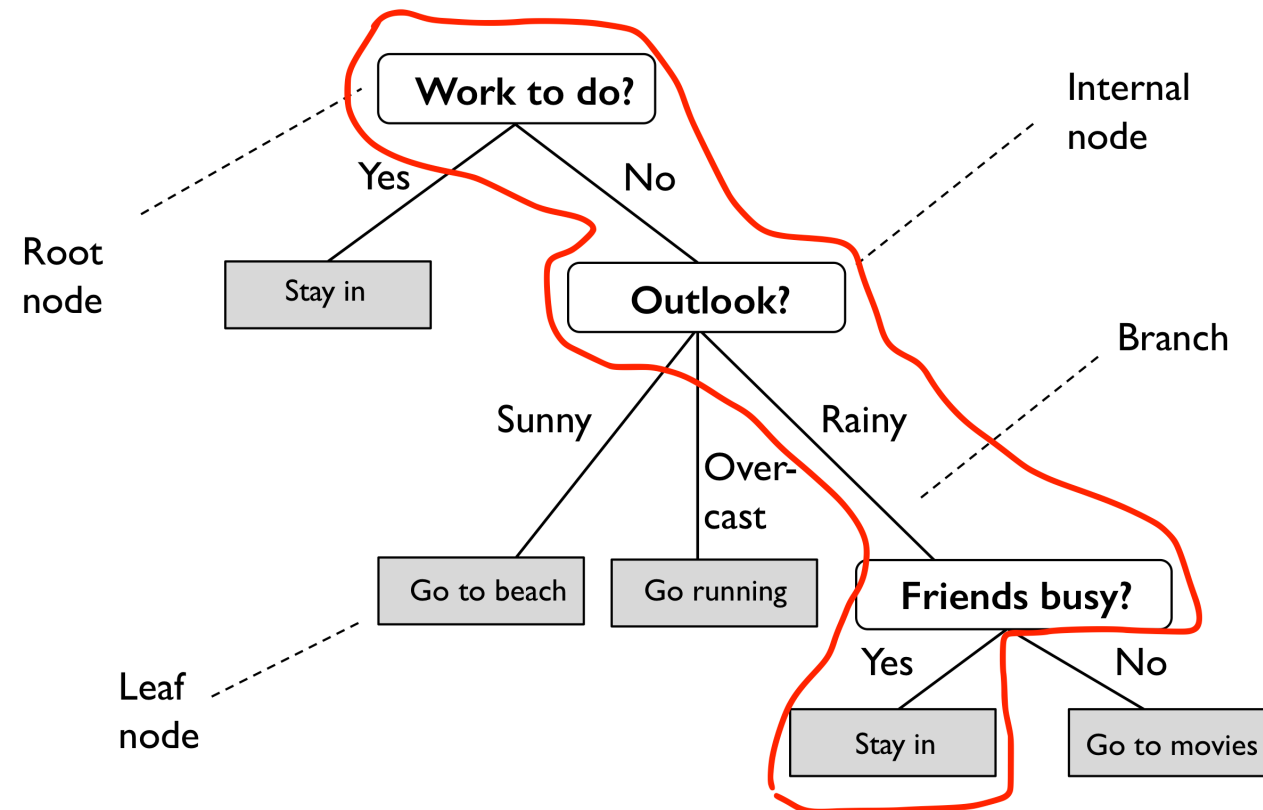5. Gini & Entropy vs misclassification error

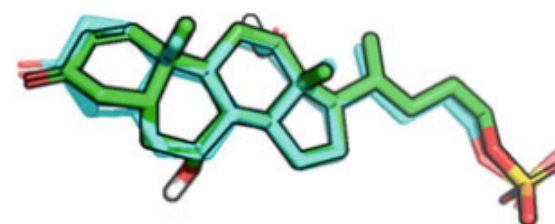6. Improvements & dealing with overfitting
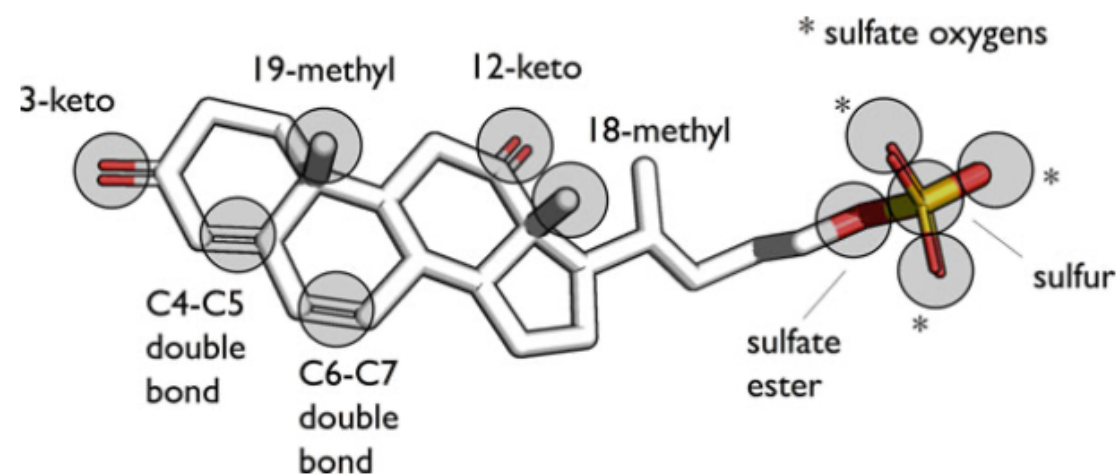
7. Code example

# Decision Tree Terminology

# Decision Trees as Rulesets



**IF**

**THEN**

* sulfate oxygens

3-keto
19-methyl
12-keto
18-methyl
C4-C5 double bond
C6-C7 double bond
sulfate ester
sulfur
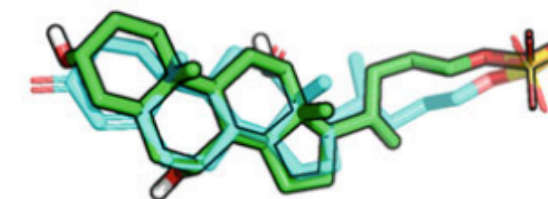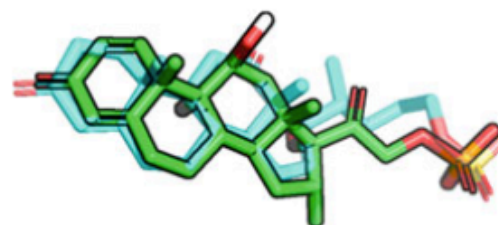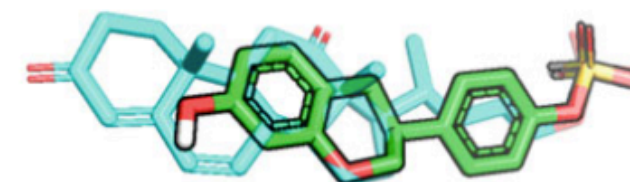
"ENE4:" 90.5 %

ZINC72400307: 90.4 %

ZINC03876071: 3.2 %

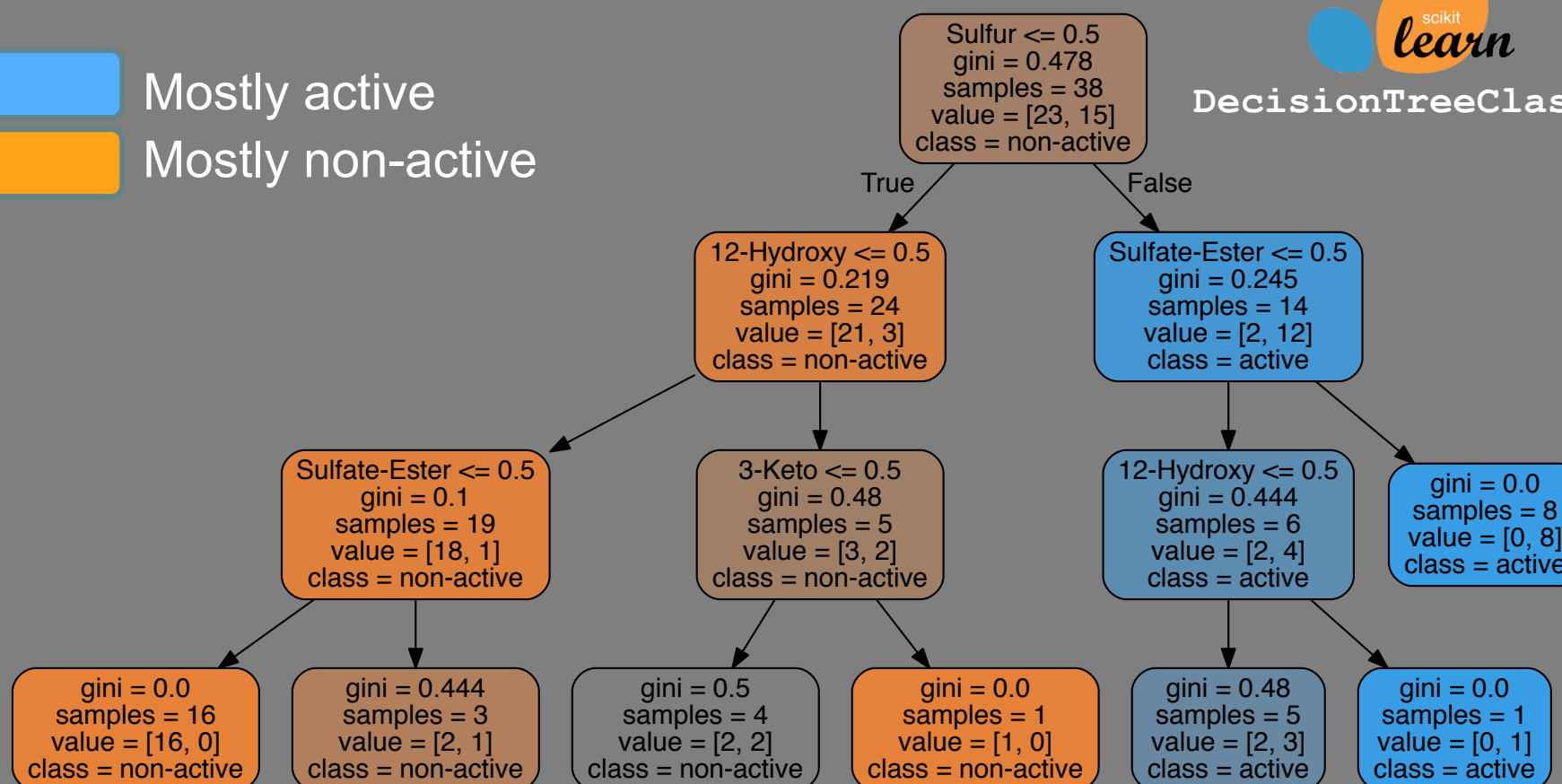ZINC40576706: 0.0 %

Mostly active
Mostly non-active

scikit learn
DecisionTreeClassifier

Sulfur <= 0.5
gini = 0.478
samples = 38
value = [23, 15]
class = non-active

True

False

12-Hydroxy <= 0.5
gini = 0.219
samples = 24
value = [21, 3]
class = non-active

Sulfate-Ester <= 0.5
gini = 0.245
samples = 14
value = [2, 12]
class = active

Sulfate-Ester <= 0.5
gini = 0.1
samples = 19
value = [18, 1]
class = non-active

3-Keto <= 0.5
gini = 0.48
samples = 5
value = [3, 2]
class = non-active

12-Hydroxy <= 0.5
gini = 0.444
samples = 6
value = [2, 4]
class = active

gini = 0.0
samples = 8
value = [0, 8]
class = active

gini = 0.0
samples = 16
value = [16, 0]
class = non-active

gini = 0.444
samples = 3
value = [2, 1]
class = non-active

gini = 0.5
samples = 4
value = [2, 2]
class = non-active

gini = 0.0
samples = 1
value = [1, 0]
class = non-active

gini = 0.48
samples = 5
value = [2, 3]
class = active

gini = 0.0
samples = 1
value = [0, 1]
class = active

Sebastian Raschka, Leslie A. Kuhn, Anne M. Scott, and Weiming Li (2018) *Computational Drug Discovery and Design: Automated Inference of Chemical Group Discriminants of Biological Activity from Virtual Screening Data.* Springer. ISBN: 978-1-4939-7755-0

# Decision tree-based diagnosis of coronary artery disease: CART model

Mohammad M. Ghiasi [a] ✉, Sohrab Zendehboudi [a], Ali Asghar Mohsenipour [b]

Show more ∨



Fig. 2. Graphical representation of the CART model (using all features) introduced for CAD diagnosis.
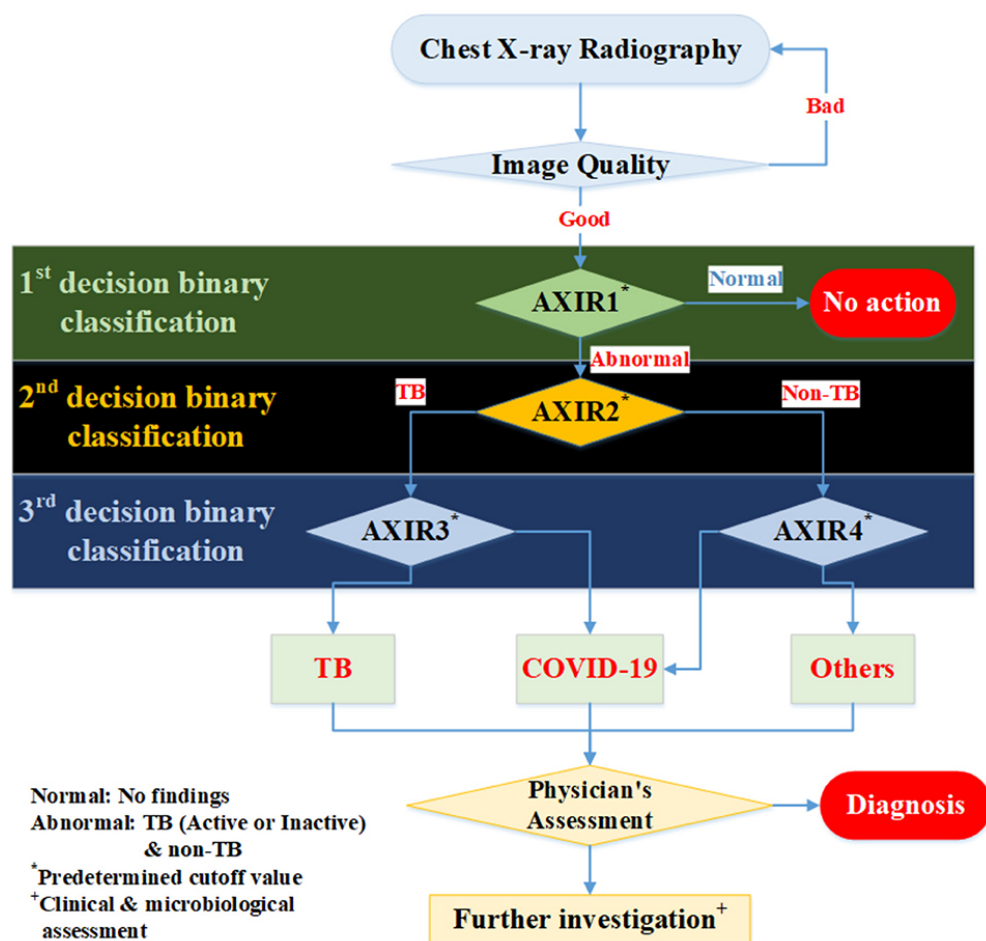
**https://www.sciencedirect.com/science/article/abs/pii/S0169260719308971**

# Deep Learning-Based Decision-Tree Classifier for COVID-19 Diagnosis From Chest X-ray Imaging

Seung Hoon Yoo[1], Hui Geng[1], Tin Lok Chiu[1], Siu Ki Yu[1], Dae Chul Cho[2], Jin Heo[2], Min Sung Choi[2], Il Hyun Choi[2], Cong Cung Van[3], Nguen Viet Nhung[3], Byung Jun Min[4*] and Ho Lee[5*]

**https://www.frontiersin.org/articles/10.3389/fmed.2020.00427/full**

**Random forests, adaptive boosting, gradient boosting**

# Lecture 6: Decision Trees
# Topics

1. Intro to decision trees

**2. Recursive algorithms & Big-O**

3. Types of decision trees

4. Splitting criteria

5. Gini & Entropy vs misclassification error

6. Improvements & dealing with overfitting

7. Code example

# 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



```python
def quicksort(array):
    if len(array) < 2:
        return array
    else:
        pivot = array[0]
        smaller, bigger = [], []
        for ele in array[1:]:
            if ele <= pivot:
                smaller.append(ele)
            else:
                bigger.append(ele)
        return quicksort(smaller) + [pivot] + quicksort(bigger)
```

# Time complexity of quicksort:

$$\mathcal{O}\ (\underline{\hspace{5cm}})$$

```
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 * | Worst |
| Quicksort | Ω(n log(n)) | Θ(n log(n)) | O(n^2) | O(log(n)) |
| Mergesort | Ω(n log(n)) | Θ(n log(n)) | O(n log(n)) | O(n) |
| Timsort | Ω(n) | Θ(n log(n)) | O(n log(n)) | O(n) |
| Heapsort | Ω(n log(n)) | Θ(n log(n)) | O(n log(n)) | O(1) |
| Bubble Sort | Ω(n) | Θ(n^2) | O(n^2) | O(1) |
| Insertion Sort | Ω(n) | Θ(n^2) | O(n^2) | O(1) |
| Selection Sort | Ω(n^2) | Θ(n^2) | O(n^2) | O(1) |
| Tree Sort | Ω(n log(n)) | Θ(n log(n)) | O(n^2) | O(n) |
| Shell Sort | Ω(n log(n)) | Θ(n(log(n))^2) | O(n(log(n))^2) | O(1) |
| Bucket Sort | Ω(n+k) | Θ(n+k) | O(n^2) | O(n) |
| Radix Sort | Ω(nk) | Θ(nk) | O(nk) | O(n+k) |
| Counting Sort | Ω(n+k) | Θ(n+k) | O(n+k) | O(k) |
| Cubesort | Ω(n) | Θ(n log(n)) | O(n log(n)) | O(n) |

http://www.bigocheatsheet.com       * "worst" ~ inversely-sorted array

# **Decision Tree in Pseudocode**

GenerateTree($\mathcal{D}$):

- if $y = 1 \ \forall \ \langle \mathbf{x}, \mathbf{y} \rangle \in \mathcal{D}$ or $y = 0 \ \forall \ \langle \mathbf{x}, y \rangle \in \mathcal{D}$ :

  - return Tree

- else:

  - Pick best feature $x_j$:

    - $\mathcal{D}_0$ at Child$_0$ : $x_j = 0 \ \forall \ \langle \mathbf{x}, y \rangle \in \mathcal{D}$
    - $\mathcal{D}_1$ at Child$_1$ : $x_j = 1 \ \forall \ \langle \mathbf{x}, y \rangle \in \mathcal{D}$

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

# Time Complexity ("Big-O")

Growing the tree: $\mathcal{O}(\;.\;.\;.$

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

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

Sorting a feature column : $\mathcal{O}(n \log n)$

Consider/sort $m$ feature columns: $\mathcal{O}(m)$

Decision tree has up to $n$ terminal leaf nodes

There are $2n - 1$ internal nodes in the tree

The number we split nodes is $2n - 1 - n = n - 1$

Complexity with re-sorting: $\mathcal{O}(mn^2 \log n)$

Complexity with sorting once & caching: $\mathcal{O}(mn \log n)$

# Time Complexity ("Big-O")

Querying the tree:    $\mathcal{O}(\,\ldots$

# Lecture 6: Decision Trees
# Topics

1. Intro to decision trees

2. Recursive algorithms & Big-O

**3. Types of decision trees**

4. Splitting criteria

5. Gini & Entropy vs misclassification error

6. Improvements & dealing with overfitting

7. Code example

# Decision Tree in Pseudocode

GenerateTree($\mathcal{D}$):

- if $y = 1 \; \forall \; \langle \mathbf{x}, \mathbf{y} \rangle \in \mathcal{D}$ or $y = 0 \; \forall \; \langle \mathbf{x}, y \rangle \in \mathcal{D}$ :

  - return Tree

- else:

  - Pick best feature $x_j$:

    - $\mathcal{D}_0$ at Child$_0$ : $x_j = 0 \; \forall \; \langle \mathbf{x}, y \rangle \in \mathcal{D}$
    - $\mathcal{D}_1$ at Child$_1$ : $x_j = 1 \; \forall \; \langle \mathbf{x}, y \rangle \in \mathcal{D}$

  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 <= 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 <= 0

3) Go back to step 1 for each of the two child nodes

- How make predictions if features in dataset are 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
  - maximum number of splits
  - score threshold or statistical significance test

# ID3 -- Iterative Dichotomizer 3

- one of the earlier 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)

- ...

# Lecture 6: Decision Trees
# Topics

1. Intro to decision trees

2. Recursive algorithms & Big-O

3. Types of decision trees

**4. Splitting criteria**

5. Gini & Entropy vs misclassification error

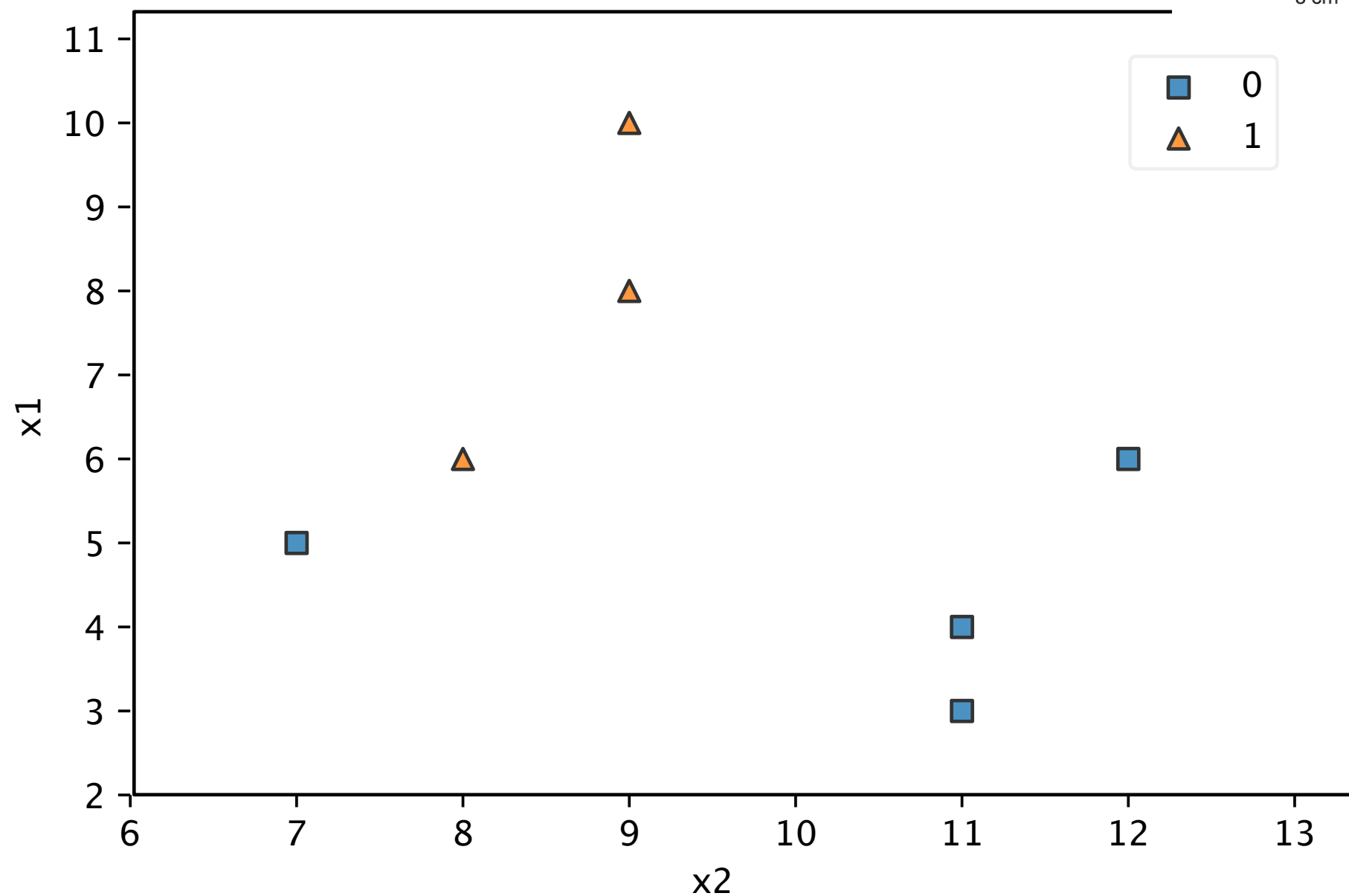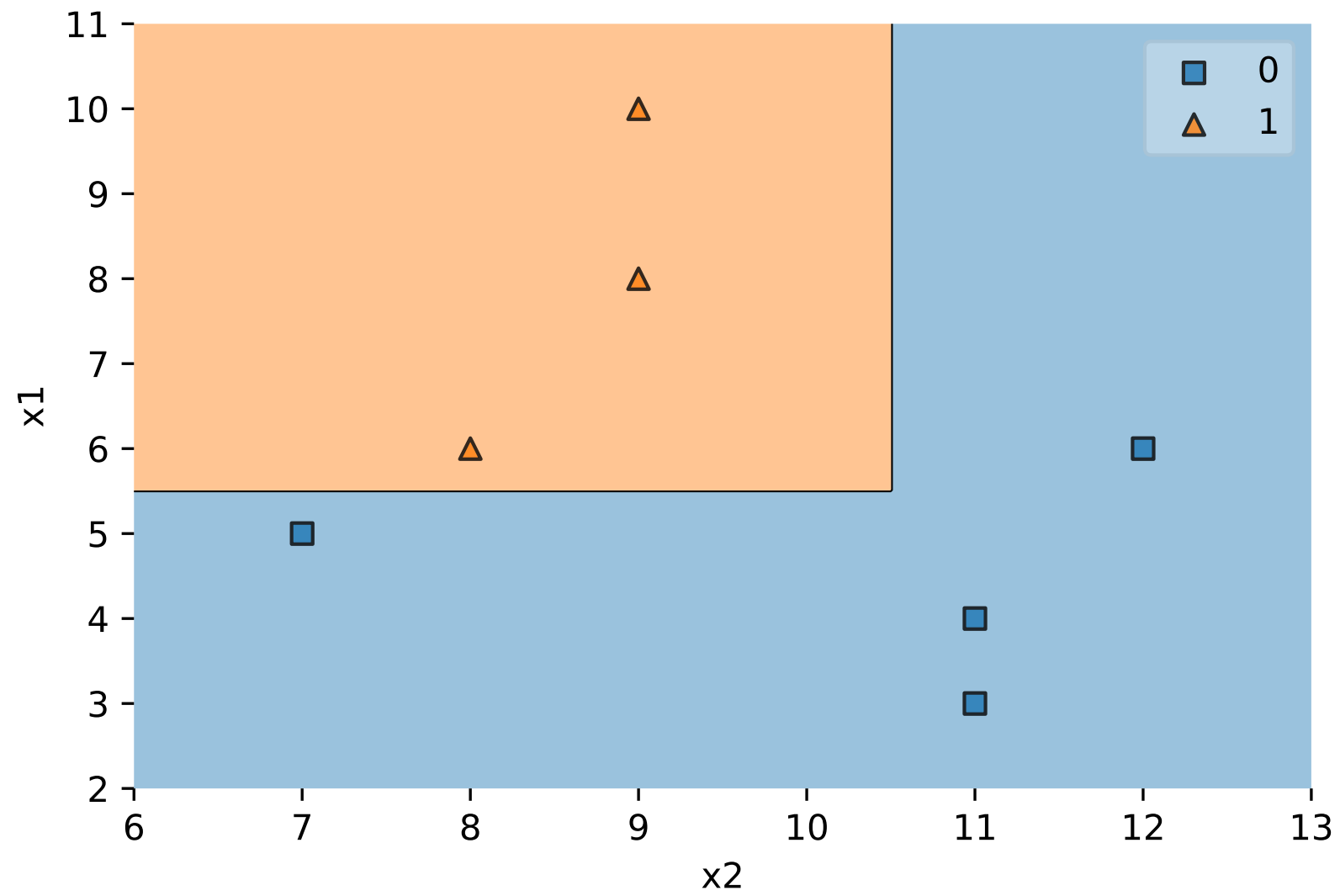6. Improvements & dealing with overfitting

7. Code example

# 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

| x₁ | x₂ | x₃ | 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 |

Decision tree diagram:

x1<=5.5
entropy=0.985
samples=7
value=[4,3]
class=Class0

True → entropy=0.0 / samples=3 / value=[3,0] / class=Class0

False → x2<=10.5 / entropy=0.811 / samples=4 / value=[1,3] / class=Class1

x2<=10.5 branches:
- entropy=0.0 / samples=3 / value=[0,3] / class=Class1
- entropy=0.0 / samples=1 / value=[1,0] / class=Class0

STAT 451: Intro to ML Lecture 6: Decision Trees

# The Splitting Criterion

# Information Gain

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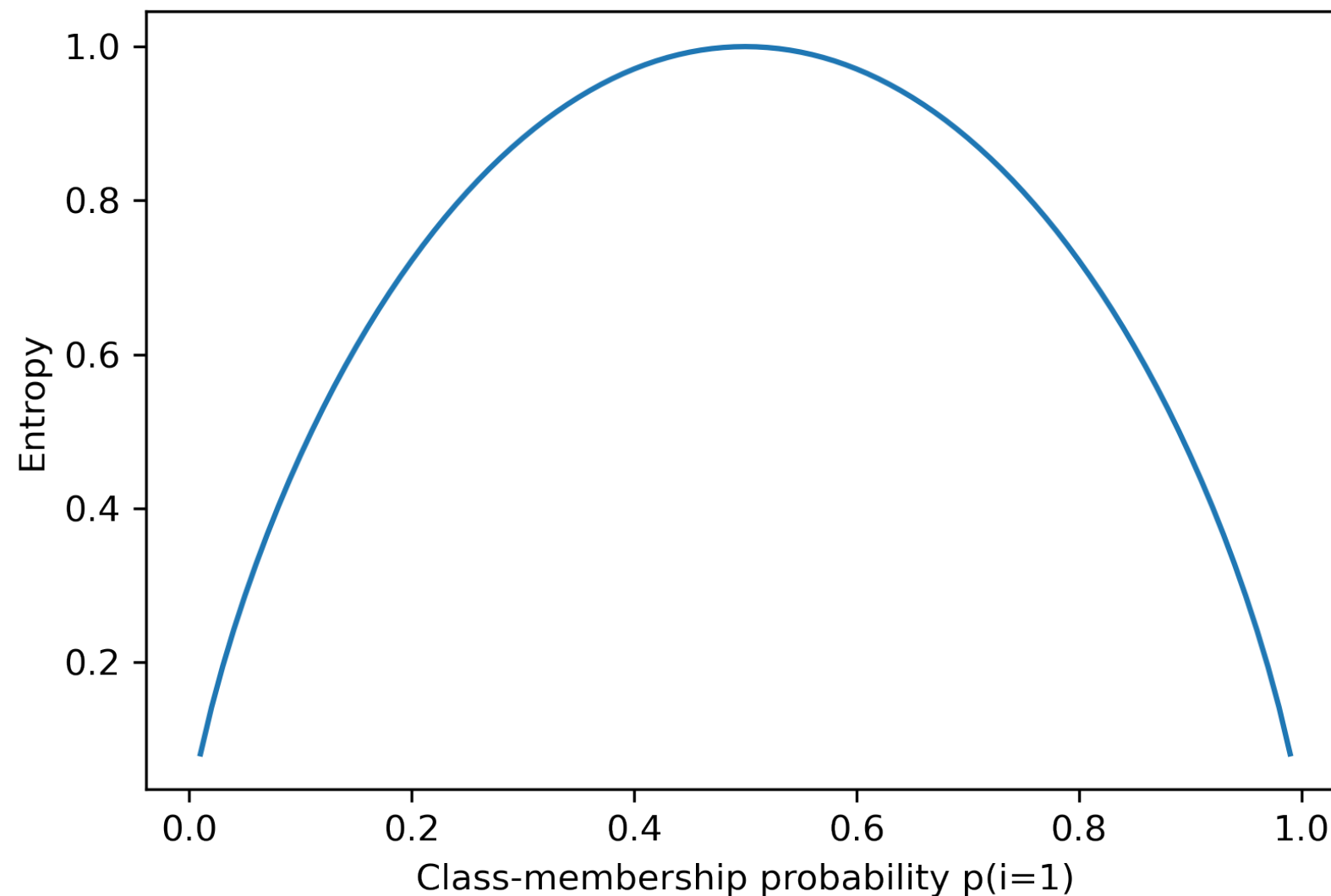
# (Shannon) Entropy

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

# (Shannon) Entropy

$$H = -\sum_{i} p(i \mid x_j)\log_2(p(i \mid x_j))$$

# Information Gain

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

# Gini Impurity

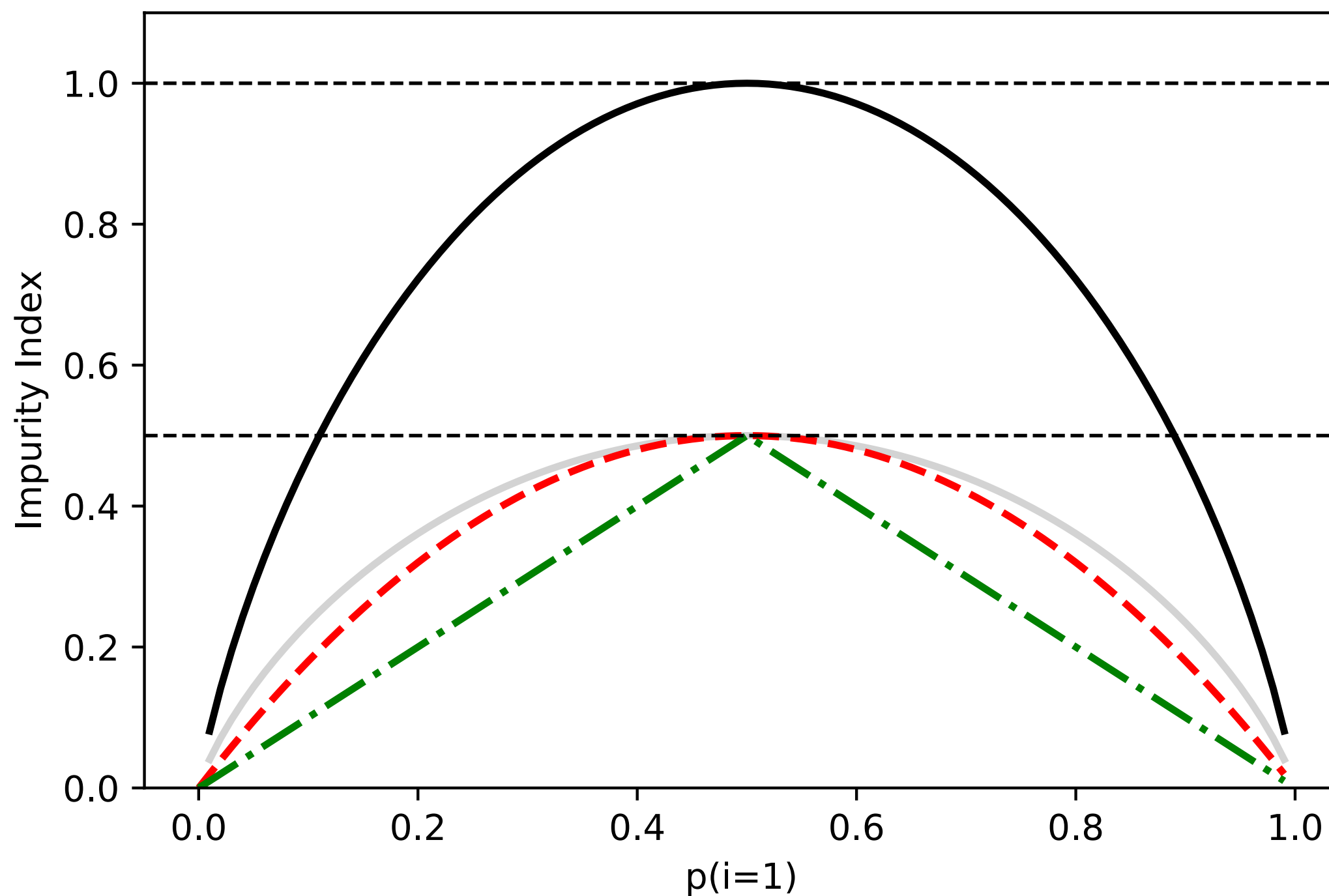$$Gini = 1 - \sum_i \left( p(i \mid x_j)^2 \right)$$

# Misclassification Error

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

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

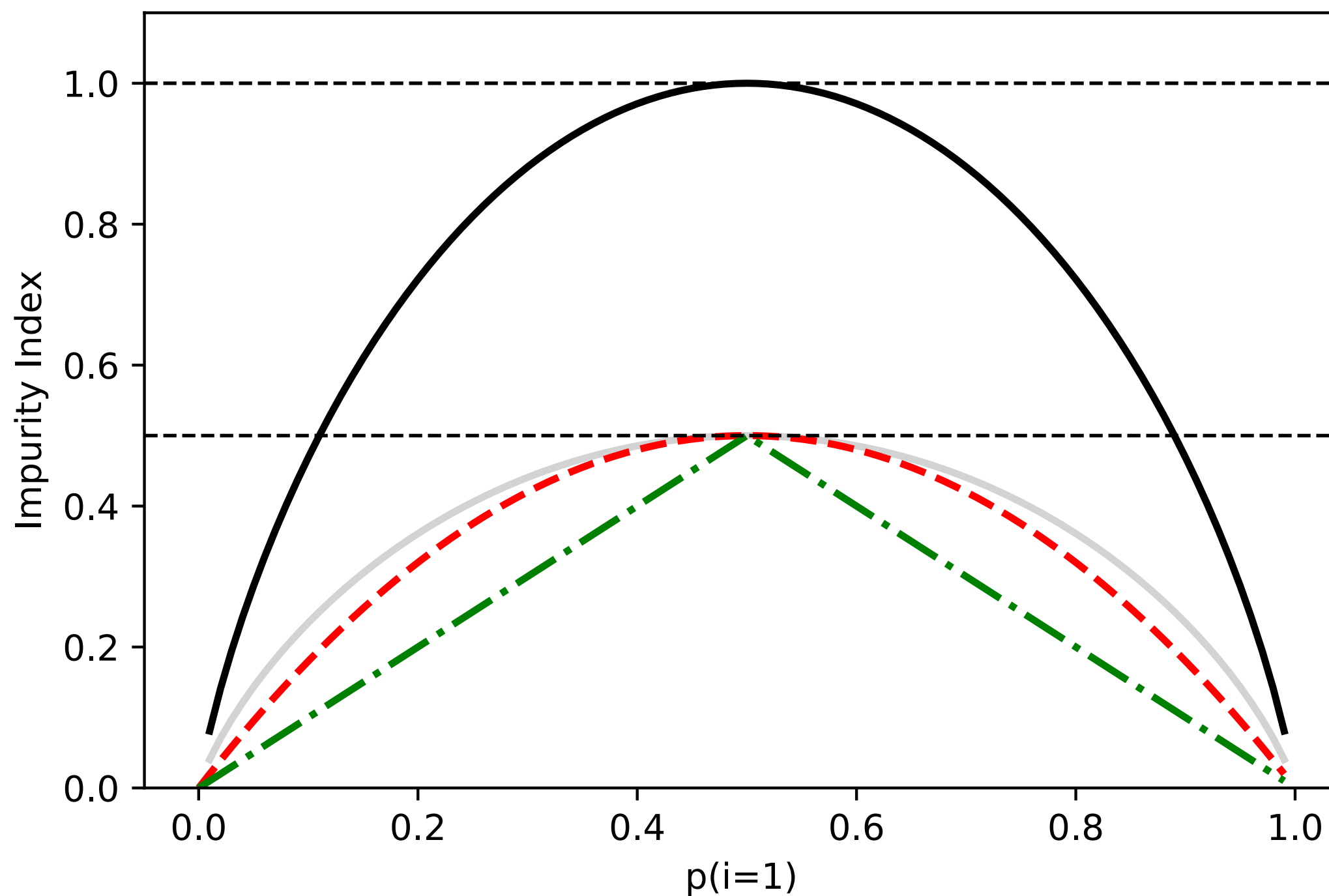# Misclassification Error

$$ERR = 1 - \max_{i}(p(i \mid x_j))$$

# Lecture 6: Decision Trees
# Topics

1. Intro to decision trees

2. Recursive algorithms & Big-O

3. Types of decision trees

4. Splitting criteria

**5. Gini & Entropy vs misclassification error**

6. Improvements & dealing with overfitting

7. Code example

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

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

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

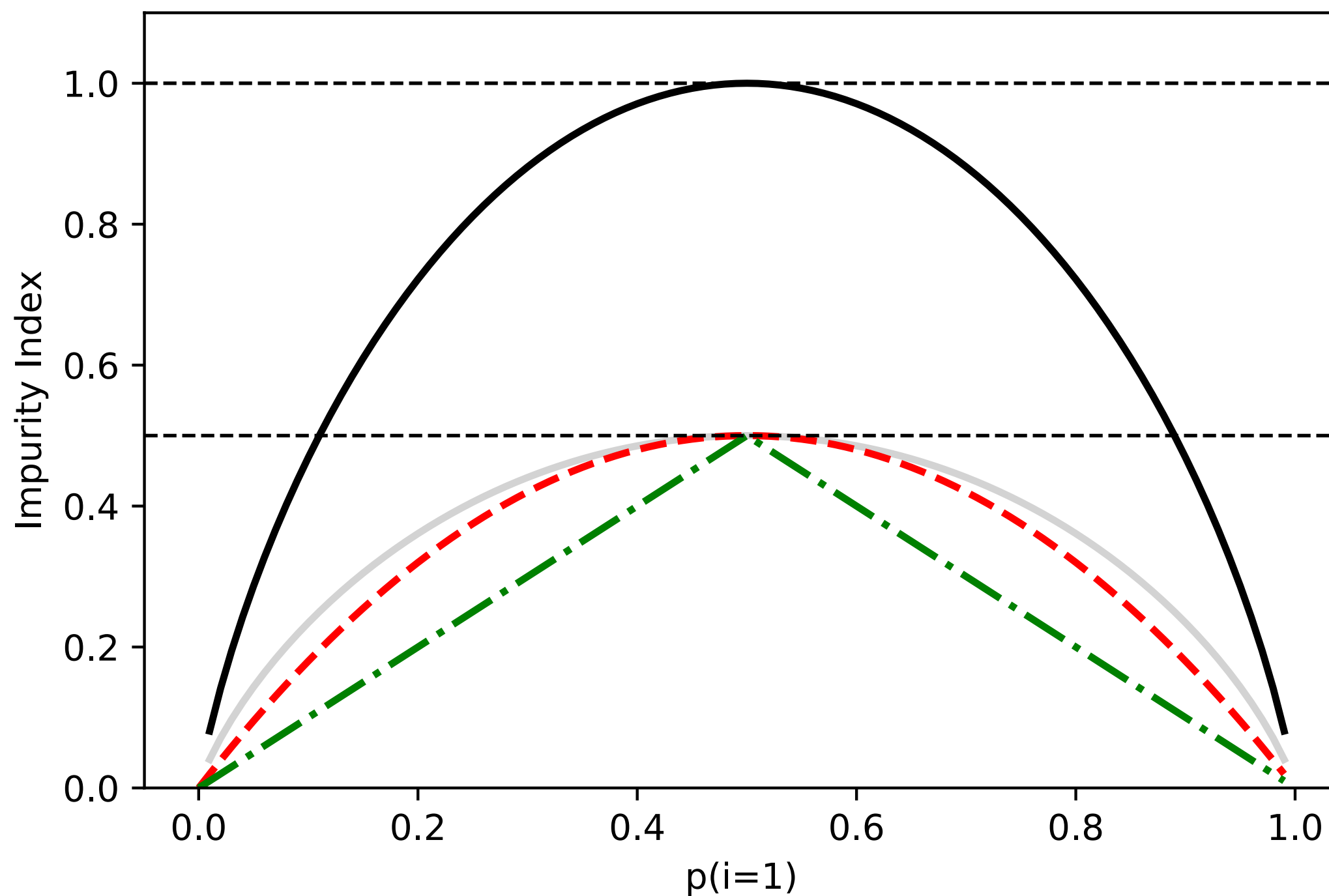# **Entropy**

$$H = -\sum_i p(i \,|\, x_j)\log_2(p(i \,|\, x_j))$$
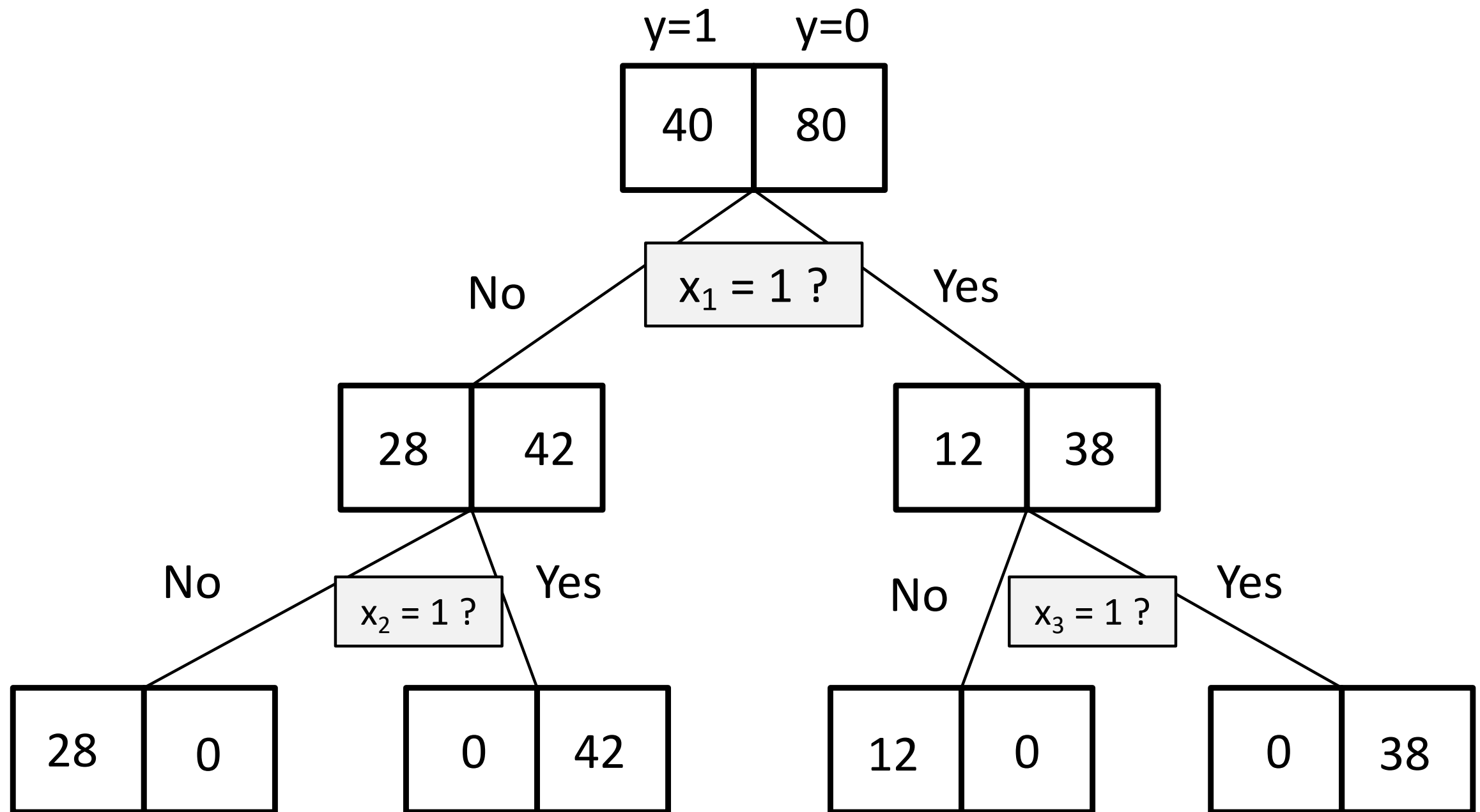
# Gini Impurity

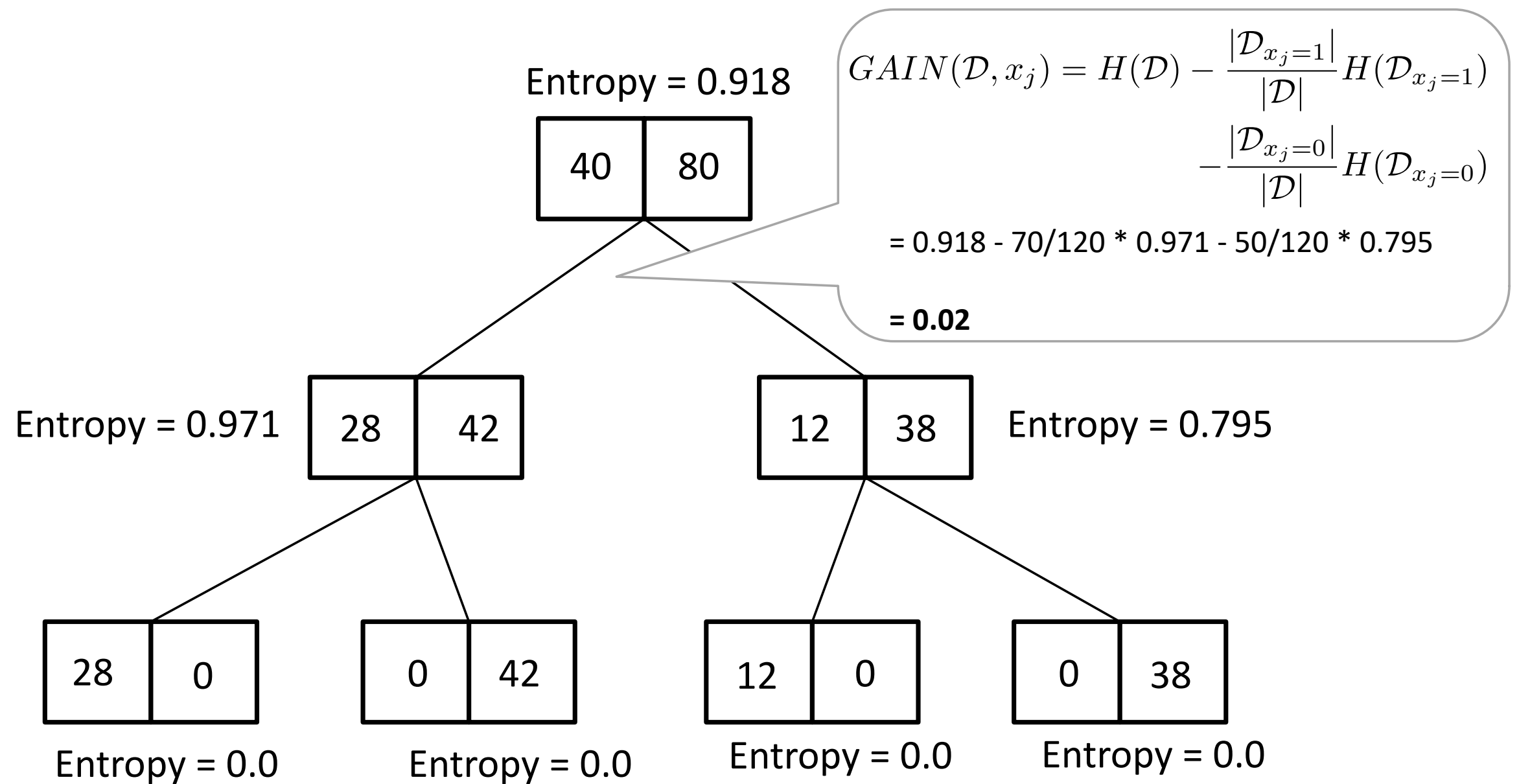$$Gini = 1 - \sum_i \left( p(i \mid x_j)^2 \right)$$

# Misclassification Error

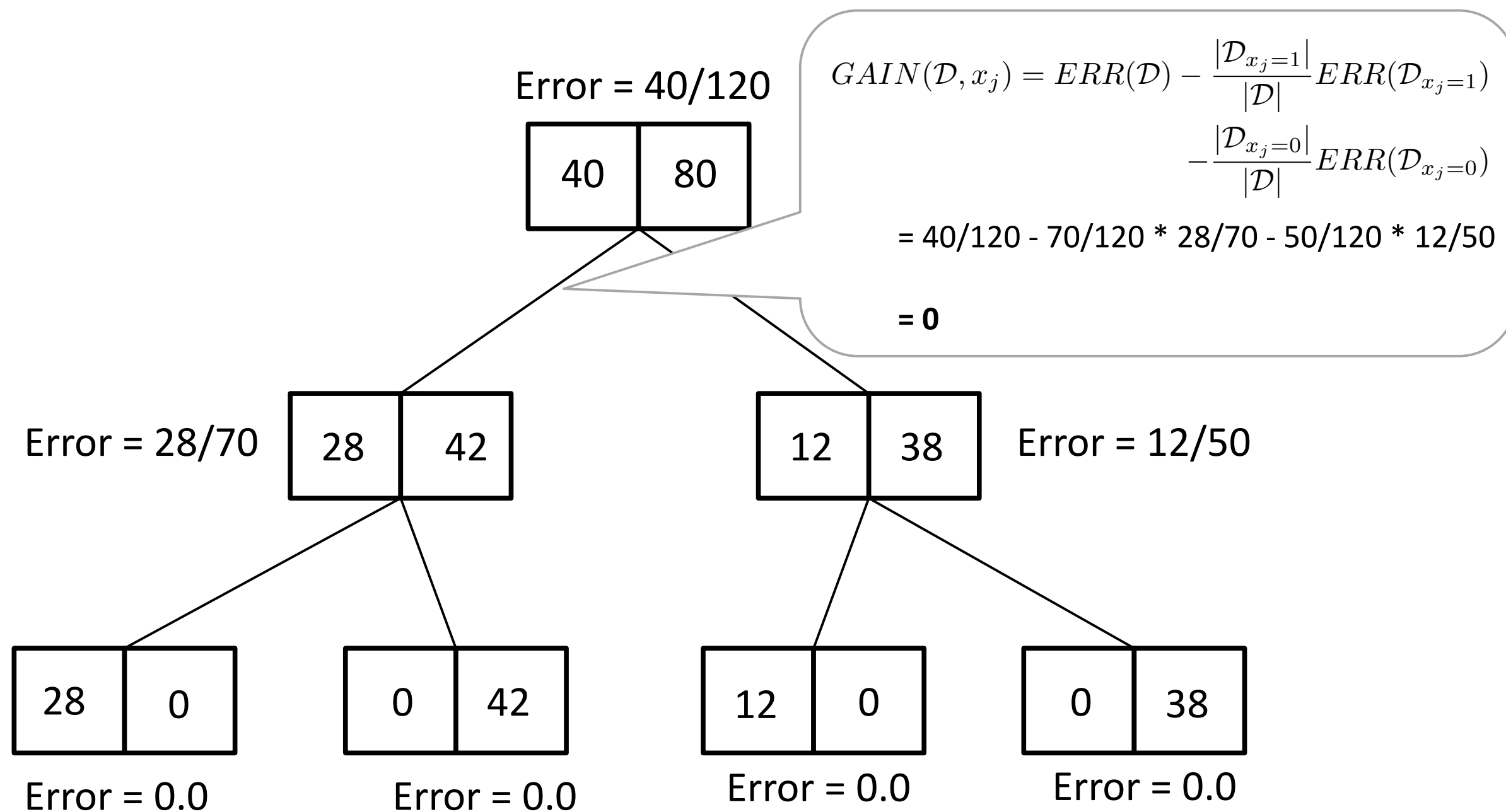$$ERR = \frac{1}{n} \sum_{i=1}^{n} L(\hat{y}^{[i]}, y^{[i]}),$$

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

Entropy = 0.918

| 40 | 80 |

$$GAIN(\mathcal{D}, x_j) = H(\mathcal{D}) - \frac{|\mathcal{D}_{x_j=1}|}{|\mathcal{D}|} H(\mathcal{D}_{x_j=1})$$

$$- \frac{|\mathcal{D}_{x_j=0}|}{|\mathcal{D}|} H(\mathcal{D}_{x_j=0})$$

= 0.918 - 70/120 * 0.971 - 50/120 * 0.795

**= 0.02**

Entropy = 0.971

| 28 | 42 |

| 12 | 38 |

Entropy = 0.795

| 28 | 0 |

Entropy = 0.0

| 0 | 42 |

Entropy = 0.0

| 12 | 0 |

Entropy = 0.0

| 0 | 38 |

Entropy = 0.0

Error = 40/120

$$GAIN(\mathcal{D}, x_j) = ERR(\mathcal{D}) - \frac{|\mathcal{D}_{x_j=1}|}{|\mathcal{D}|} ERR(\mathcal{D}_{x_j=1})$$

$$- \frac{|\mathcal{D}_{x_j=0}|}{|\mathcal{D}|} ERR(\mathcal{D}_{x_j=0})$$

= 40/120 - 70/120 * 28/70 - 50/120 * 12/50

**= 0**

| 40 | 80 |

Error = 28/70

| 28 | 42 |

| 12 | 38 |

Error = 12/50

| 28 | 0 |

Error = 0.0

| 0 | 42 |

Error = 0.0
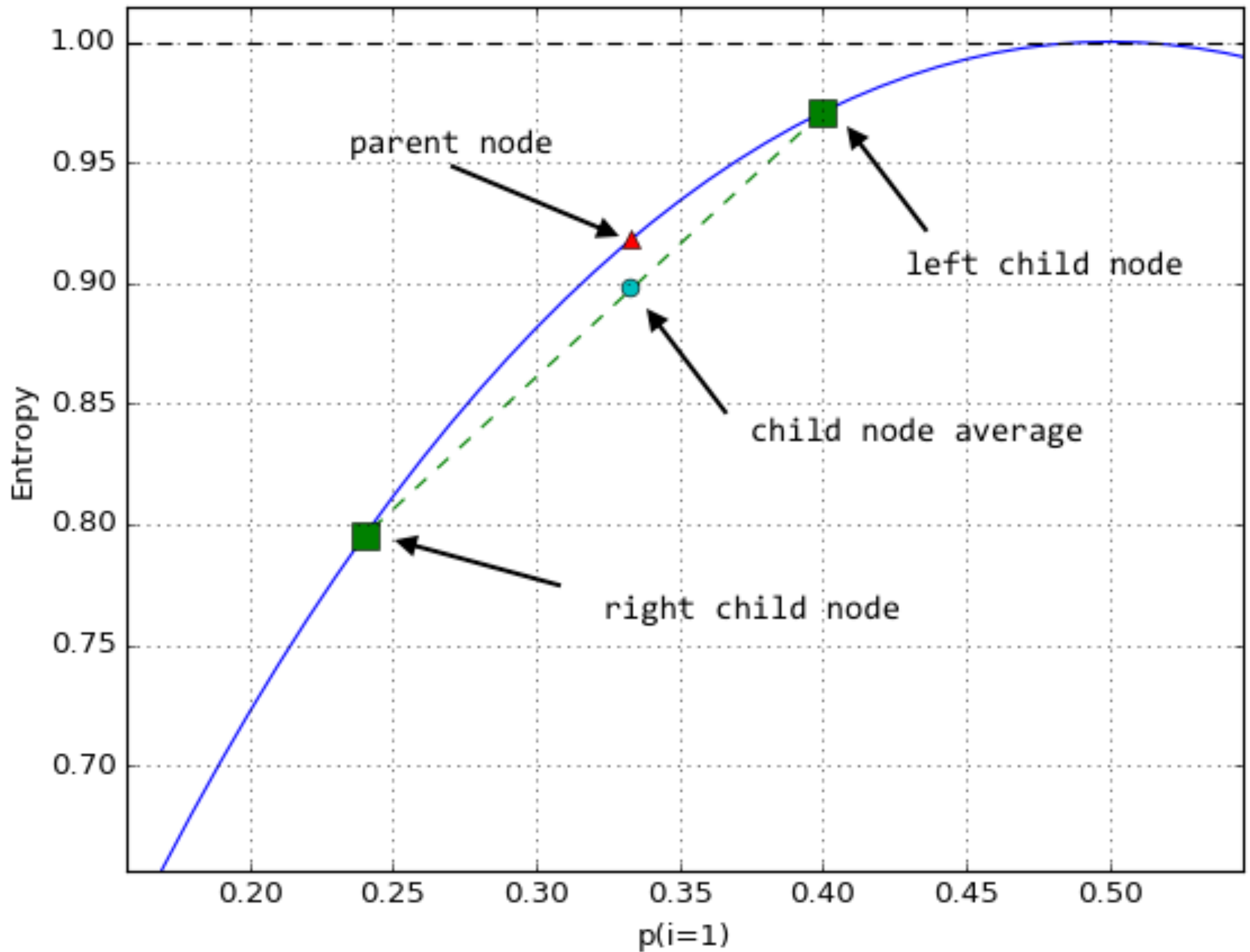
| 12 | 0 |

Error = 0.0

| 0 | 38 |

Error = 0.0

# Lecture 6: Decision Trees
# Topics

1.  Intro to decision trees

2. Recursive algorithms & Big-O

3. Types of decision trees

4. Splitting criteria

5. Gini & Entropy vs misclassification error

**6. Improvements & dealing with overfitting**

7. Code example

# Gain Ratio

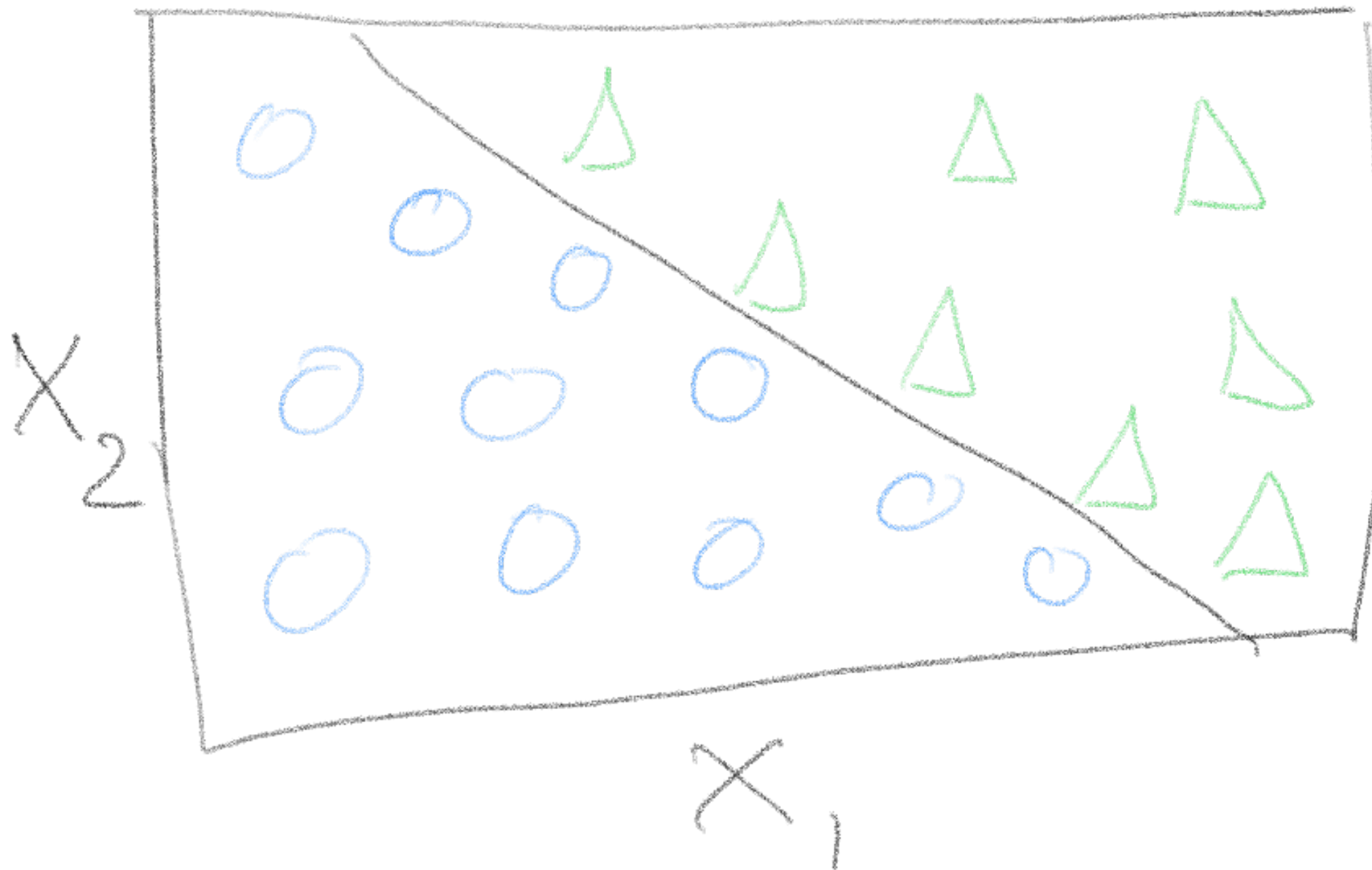$$GainRatio(\mathcal{D}, x_j) = \frac{Gain(\mathcal{D}, x_j)}{SplitInfo(\mathcal{D}, x_j)}$$

where the split information measures the entropy of the feature:

$$SplitInfo(\mathcal{D}, x_j) = -\sum_{v \in x_j} \frac{|\mathcal{D}_v|}{|\mathcal{D}|} \log_2 \frac{|\mathcal{D}_v|}{|\mathcal{D}|}$$
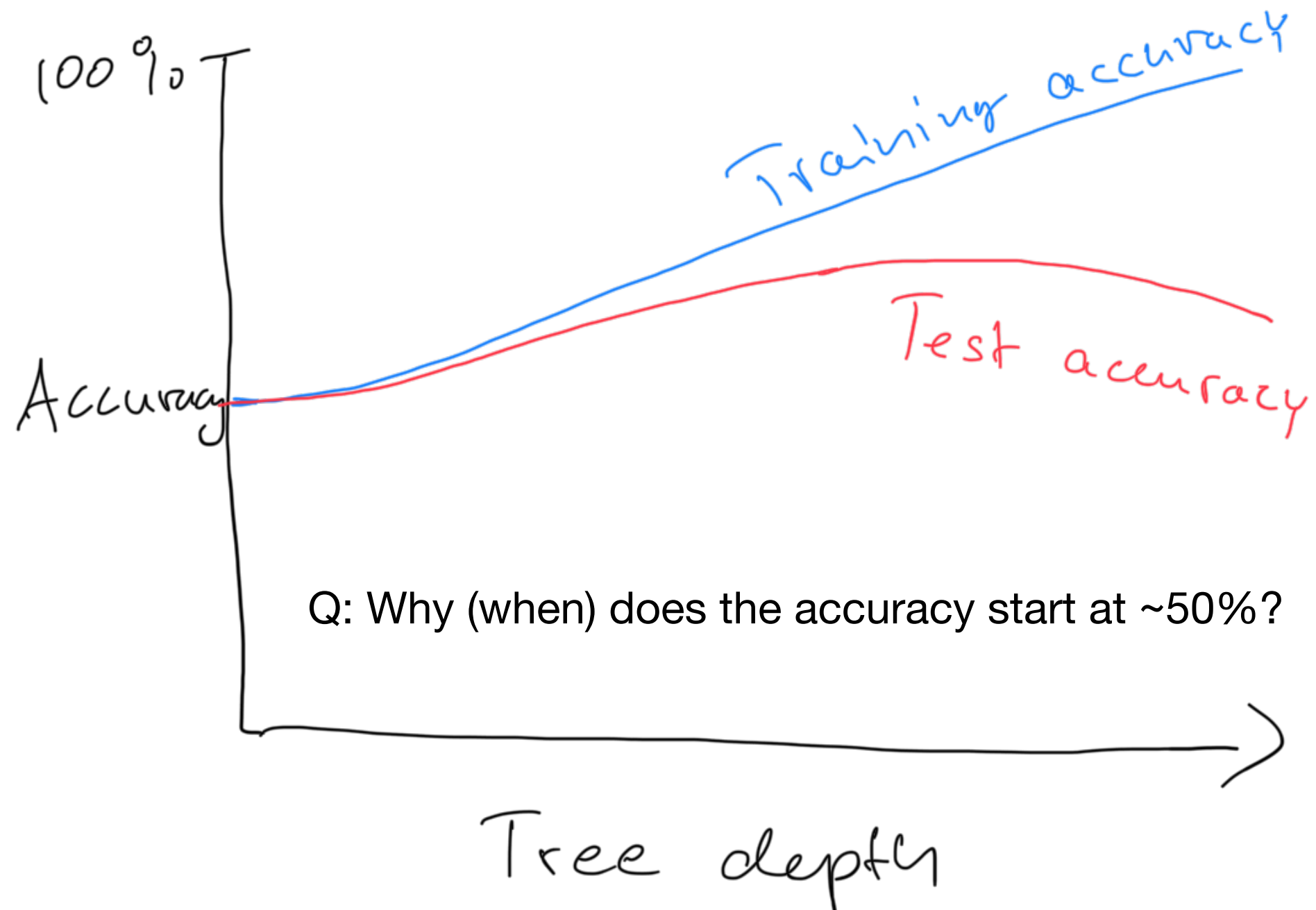
Penalizes splitting categorical attributes with many values (e.g., think date column, or really bad: row ID) via the split information

# Shortcomings



How would the decision tree split look like?

# Overfitting



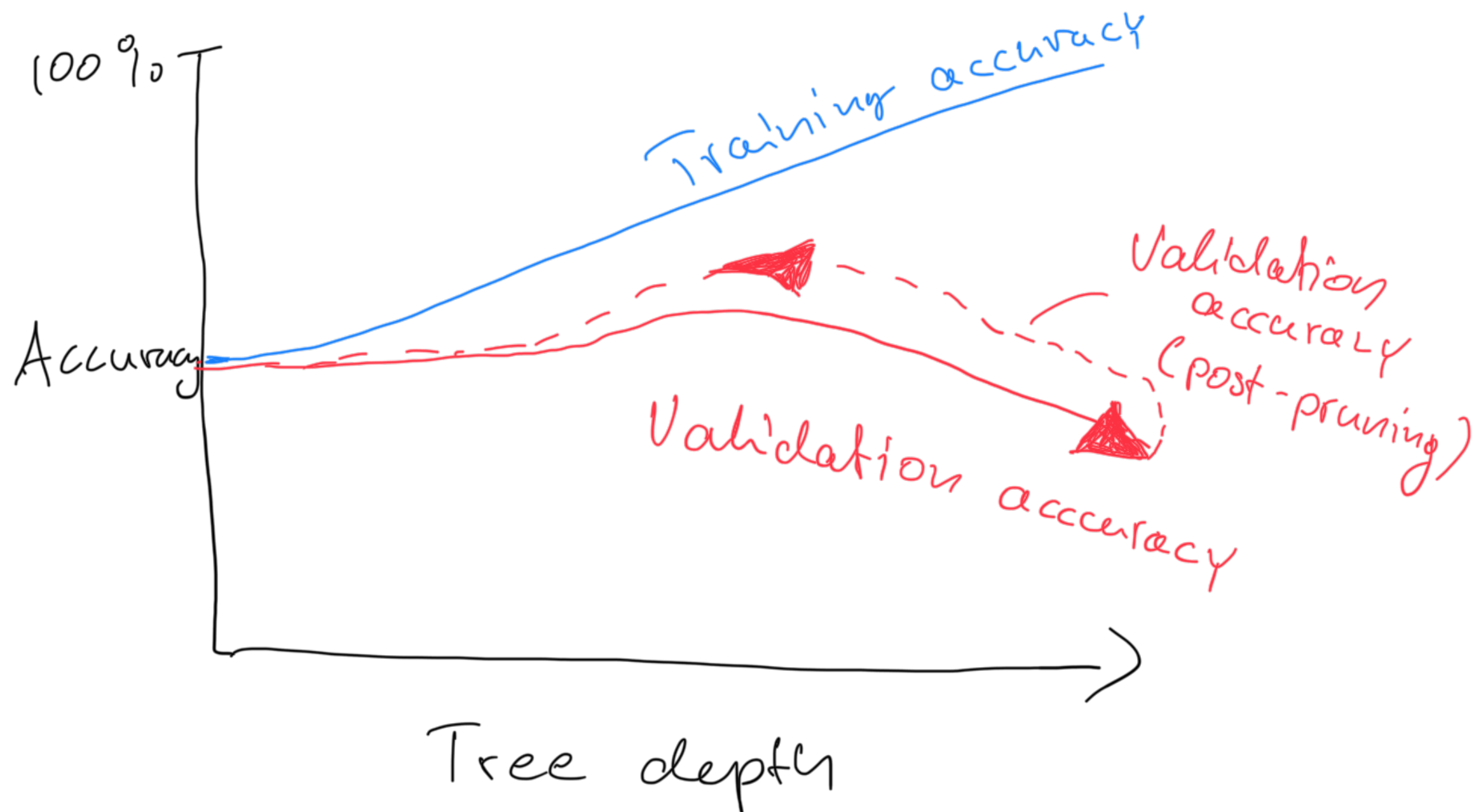Q: Why (when) does the accuracy start at ~50%?

# 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

# Lecture 6: Decision Trees
# Topics

1. Intro to decision trees

2. Recursive algorithms & Big-O

3. Types of decision trees

4. Splitting criteria

5. Gini & Entropy vs misclassification error

6. Improvements & dealing with overfitting

**7. Code example**

# Demo