Updated 10/22/2025

Part 3: Applications L12: Decision trees (1/2)

Wolfgang Gatterbauer

cs7840 Foundations and Applications of Information Theory (fa25)

https://northeastern-datalab.github.io/cs7840/fa25/

10/20/2025

Pre-class conversations

- Last class recapitulation
- Piazza: is it cumbersome to find comments on your scribes? Any suggestion for a better system?
- Projects: I have time to talk today after class and also THU

- Today:
 - Decision trees
 - backed in: Occam, MDL, fun connections to topics you have seen before here or elsewhere

Decision trees

Formal setup

EXAMPLE: Classifying days based on weather conditions.

| /Index | (not part |
|----------|-----------|
| , | "data") |
| | |

| day | (O)utlook | (T)emp. | (H)umidity | (W)ind | (C)lass |
|-----|-----------|---------|------------|--------|---------|
| 1 | sunny | hot | high | weak | no |
| 2 | sunny | hot | high | strong | no |
| 3 | overcast | hot | high | weak | yes |
| 4 | rain | mild | high | weak | yes |
| 5 | rain | cool | normal | weak | yes |
| 6 | rain | cool | normal | strong | no |
| 7 | overcast | cool | normal | strong | yes |
| 8 | sunny | mild | high | weak | no |
| 9 | sunny | cool | normal | weak | yes |
| 10 | rain | mild | normal | weak | yes |
| 11 | sunny | mild | normal | strong | yes |
| 12 | overcast | mild | high | strong | yes |
| 13 | overcast | hot | normal | weak | yes |
| 14 | rain | mild | high | strong | no |

- What are the features (predictors)?
- What is the target concept (the response variable)?



Formal setup

EXAMPLE: Classifying days based on weather conditions.

Class label y_i denotes weather a particular event happened (here whether tennis game happened)

Columns denote m = 4 features $\{X_j\}_{j=1}^m$. Domain \mathcal{X}_H of feature X_H is {high, normal}

| • | Concrete | Problem | Setting | in | this | Examp | le: |
|---|----------|---------|---------|----|------|-------|-----|
|---|----------|---------|---------|----|------|-------|-----|

| | Predictors | | | | Response | | |
|-----|------------|---------|----------|----|----------|---------|-------------------------------------|
| day | (O)utlook | (T)emp. | (H)umidi | ty | (W)ind | (C)lass | |
| 1 | sunny | hot | high | | weak | no | |
| 2 | sunny | hot | high | | strong | no | |
| 3 | overcast | hot | high | | weak | yes | |
| 4 | rain | mild | high | | weak | yes | $\langle \mathbf{x}_4, y_4 \rangle$ |
| 5 | rain | cool | normal | | weak | yes | A |
| 6 | rain | cool | normal | | strong | no | |
| 7 | overcast | cool | normal | | strong | yes | |
| 8 | sunny | mild | high | | weak | no | |
| 9 | sunny | cool | normal | | weak | yes | |
| 10 | rain | mild | normal | | weak | yes | |
| 11 | sunny | mild | normal | | strong | yes | |
| 12 | overcast | mild | high | | strong | yes | |
| 13 | overcast | hot | normal | | weak | yes | |
| 14 | rain | mild | high | | strong | no | |

Input:

• Output: ?

Rows denote labeled instances $\langle \mathbf{x}_i, y_i \rangle$.

Formal setup

EXAMPLE: Classifying days based on weather conditions.

Class label y_i denotes weather a particular event happened (here whether tennis game happened)

Columns denote m=4 features $\{X_j\}_{j=1}^m$. Domain \mathcal{X}_H of feature X_H is {high, normal}/

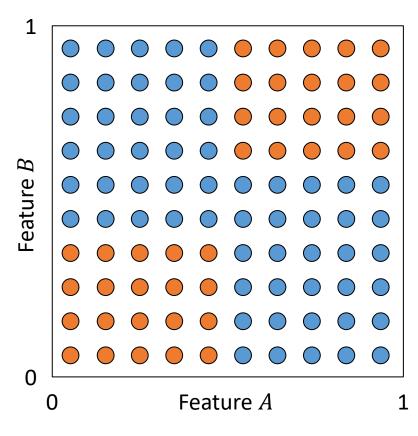
| | | Predictors | | | | Response | |
|-----|-----------|------------|----------|----|--------|----------|-------------------------------------|
| day | (O)utlook | (T)emp. | (H)umidi | ty | (W)ind | (C)lass | |
| 1 | sunny | hot | high | | weak | no | |
| 2 | sunny | hot | high | | strong | no | |
| 3 | overcast | hot | high | | weak | yes | |
| 4 | rain | mild | high | | weak | yes | $\langle \mathbf{x}_4, y_4 \rangle$ |
| 5 | rain | cool | normal | | weak | yes | A |
| 6 | rain | cool | normal | | strong | no | |
| 7 | overcast | cool | normal | | strong | yes | |
| 8 | sunny | mild | high | | weak | no | |
| 9 | sunny | cool | normal | | weak | yes | |
| 10 | rain | mild | normal | | weak | yes | |
| 11 | sunny | mild | normal | | strong | yes | |
| 12 | overcast | mild | high | | strong | yes | |
| 13 | overcast | hot | normal | | weak | yes | |
| 14 | rain | mild | high | | strong | no | |

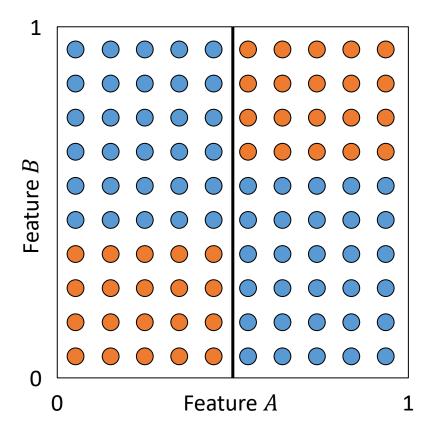
Rows denote labeled instances $\langle \mathbf{x}_i, y_i \rangle$.

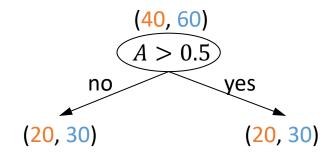
- Concrete Problem Setting in this Example:
 - Set of possible instances $\mathcal{X} = \mathcal{X}_0 \times \mathcal{X}_T \times \mathcal{X}_H \times \mathcal{X}_W$
 - Set of possible labels $\mathcal{Y} = \{\text{yes, no}\}\$ with size $k = |\mathcal{Y}| = 2$ (binary)
 - Unknown target function $f: \mathcal{X} \to \mathcal{Y}$
 - Set of function hypotheses $H = \{h \mid h: X \rightarrow Y\}$
- Input: n training examples of unknown target function f

$$\{\langle \mathbf{x}_i, y_i \rangle\}_{i=1}^n = \{\langle \mathbf{x}_1, y_1 \rangle, \dots, \langle \mathbf{x}_n, y_n \rangle\}$$

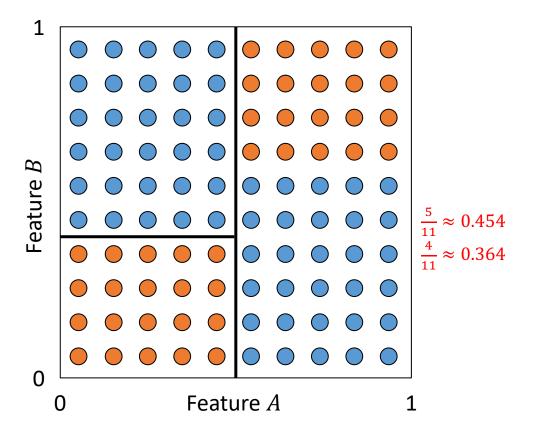
• Output: Hypothesis $h \in H$ that best approximates f

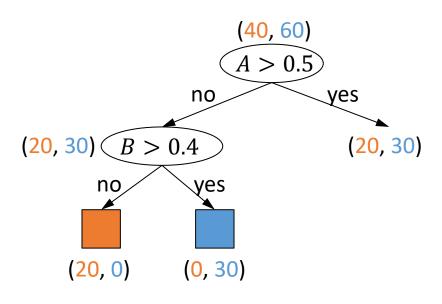


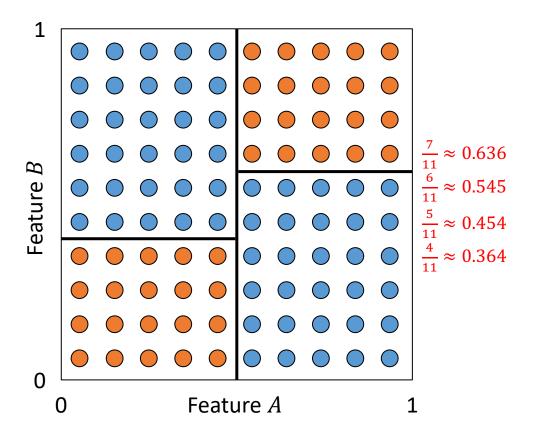


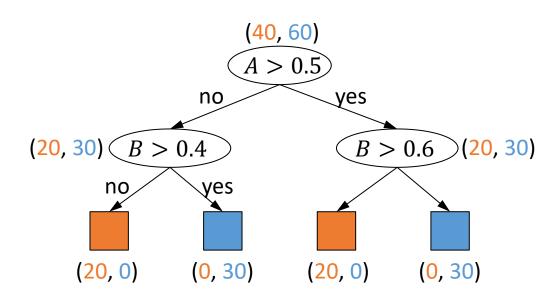


(This is not what a practical algorithm with do. We come back to that later.)

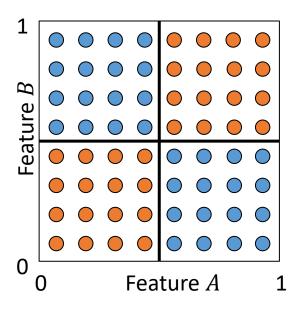


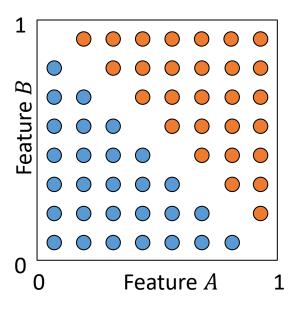


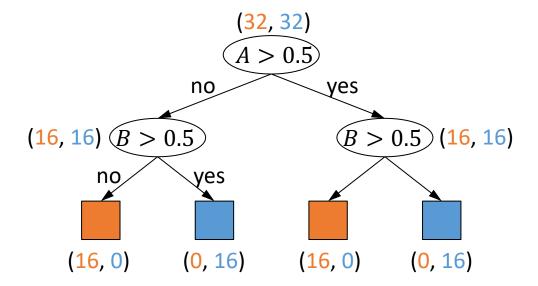




Rectilinear decision boundaries

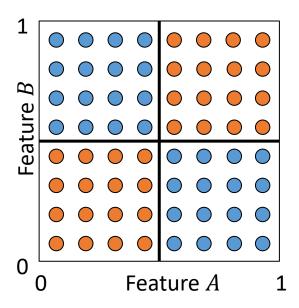


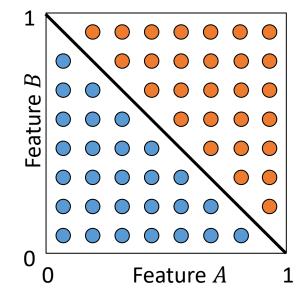


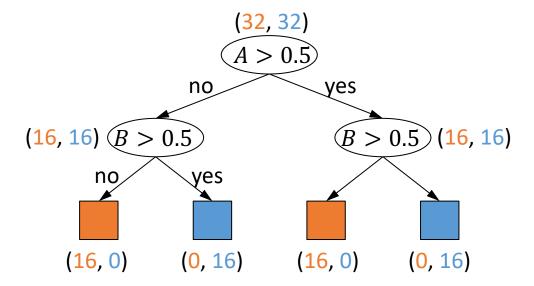


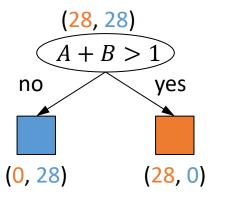


Rectilinear vs. oblique decision boundaries









Rectilinear vs. oblique decision boundaries

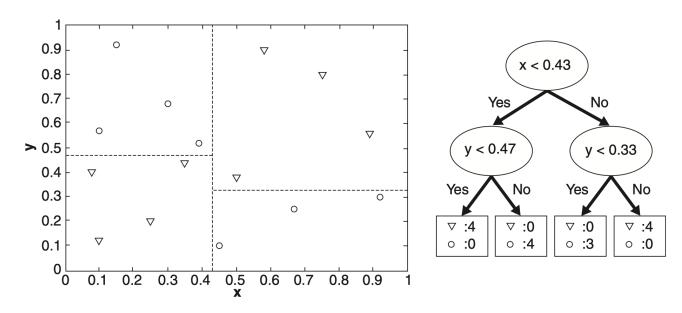


Figure 3.20. Example of a decision tree and its decision boundaries for a two-dimensional data set.

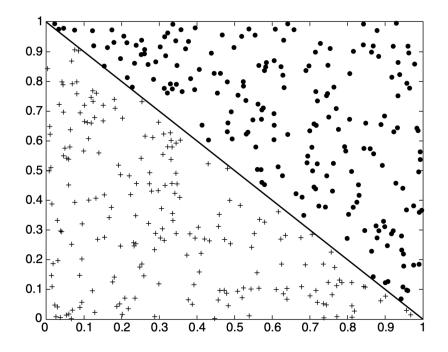


Figure 4.21. Example of data set that cannot be partitioned optimally using test conditions involving single attributes.

Rectilinear vs. oblique decision boundaries

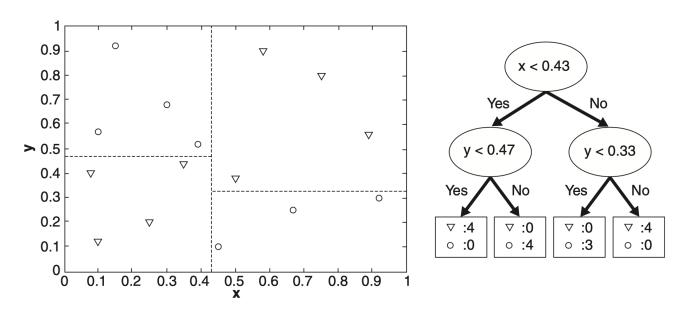


Figure 3.20. Example of a decision tree and its decision boundaries for a two-dimensional data set.

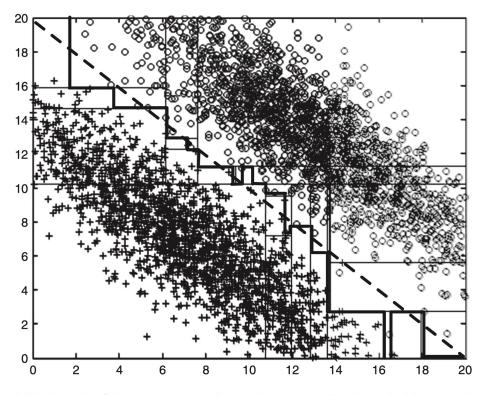


Figure 3.21. Example of data set that cannot be partitioned optimally using a decision tree with single attribute test conditions. The true decision boundary is shown by the dashed line.

Decision Trees

- Decision trees are a popular technique for both <u>classification</u> (categorical values) and <u>regression</u> (numerical values) tasks.
- During training (when developing the tree):
 - The training data is recursively partitioned into smaller yet <u>increasingly homogenous</u> regions
 - Each internal node tests an attribute to direct the example toward <u>a leaf</u>, which predicts the <u>majority class</u> (or mean value) of the examples contained in its region.
- During prediction (when classifying an instance):
 - The instance is routed from the root to a leaf according to the outcomes of the tests along the path. The leaf's label becomes the tree's prediction.

How to split?

Constructing Decision Trees (DTs)

- Given some training data, what is the "optimal" DT?
 - With optimal, we mean here the "smallest" that fits the data perfectly
- In general, finding an optimal DT is called intractable (NP-hard)
 - There are exponentially many DT's that could be constructed from a given set of attributes (exponential in number of attributes)
- In practice, we use greedy heuristics to construct a good DT (making a series of locally optimal decisions)
- Hunt's algorithm: a decision tree is grown in a recursive fashion by partitioning the training records into successively "purer" subsets
 - We can use different notions of "impurity"

Practical hardness of optimal decision trees?

CONSTRUCTING OPTIMAL BINARY DECISION TREES IS NP-COMPLETE*

Laurent HYAFIL

IRIA - Laboria, 78150 Rocquencourt, France

and

Ronald L. RIVEST

Dept. of Electrical Engineering and Computer Science, M.I.T., Cambridge, Massachusetts 02139, USA

Information Processing Letters'76

Effectiveness with modern ILP solvers. Problem: we may need exponentially many statistics over the data (but exponential only in number of attribute not data size)



Cp. to Shannon-Fano top-down vs. Huffman optimal bottom-up!

We demonstrate that constructing optimal binary decision trees is an NP-complete problem, where an optimal tree is one which minimizes the expected number of tests required to identify the unknown object.

Let $p(x_i)$ be the length of the path from the root of the tree to the terminal node naming x_i , that is, the number of tests required to identify x_i . Then the cost of this tree is merely the external path length, that is, $\sum_{x_i \in X} p(x_i)$. This model is identical to that studied by Garey [3].

The decision tree problem $DT(\mathcal{I}, X, w)$ is to determine whether there exists a decision tree with cost less than or equal to w, given \mathcal{I} and X.

To show that DT is NP-complete, we show that EC3 α DT, where EC3 is the problem of finding an exact cover for a set X, and where each of the subsets available for use contains exactly 3 elements. More

Practical hardness of optimal decision trees?

CONSTRUCTING OPTIMAL BINARY DECISION TREES IS NP-COMPLETE*

Laurent HYAFIL

IRIA - Laboria, 78150 Rocquencourt, France

and

Ronald L. RIVEST

Dept. of Electrical Engineering and Computer Science, M.I.T., Cambridge, Massachusetts 02139, USA

Information Processing Letters'76

Mach Learn (2017) 106:1039–1082 DOI 10.1007/s10994-017-5633-9

Optimal classification trees

Dimitris Bertsimas¹ · Jack Dunn²

7 Conclusions

In this paper, we have revisited the classical problem of decision tree creation under a modern MIO lens. We presented a novel MIO formulation for creating optimal decision trees that captures the discrete nature of this problem and motivates two new methods, OCT and OCT-H.

Experiments with synthetic data provide strong evidence that optimal decision trees can better recover the true generating decision tree in the data, contrary to the popular belief that such optimal methods will just overfit the training data at the expense of generalization ability.

Exploiting the astonishing progress of MIO solvers in the last 25 years, our comprehensive computational experiments showed that both the OCT and OCT-H problems are both practically solvable and deliver solutions that outperform the classical state-of-the-art methods, with average absolute improvements in out-of-sample accuracy over CART of 1–2% for OCT and 3–5% for OCT-H across all datasets, depending on the depth of the tree. We also provided guidance for each method that predict when each is most likely to outperform CART.

Hunt's algorithm: Top-down induction of Decision trees

- Create a root node x; assign it all training examples: D_x
- Repeat {
 - If all records in D_x belong to the same class, then make x a leaf node and assign it the class label
 - Else:
 - Choose an attribute A that partitions the training records at node x into the "purest" subsets
 - For each value v of A, create a new child node $X_{A=v}$ and assign it the training examples $D_{A=v}$
 - Choose a non-leaf node x
- Until all nodes are leaves

There are different ways to measure "impurity" and we will discuss in a moment variants of how we could measure that

2.1 Divide and conquer

The skeleton of Hunt's method for constructing a decision tree from a set T of training cases is elegantly simple. Let the classes be denoted $\{C_1, C_2, \ldots, C_k\}$. There are three possibilities:

*****..

[Quinlan'93]

"Tan, Steinbach, Kumar. Introduction to Data Mining-Pearson" mentions "Hunt's algorithm" but does not leave a citation. Quinlan himself in his 1993 book on C4.5 refers to "Hunt, Marin, Stone. Experiments in induction. Academic press, 1966".

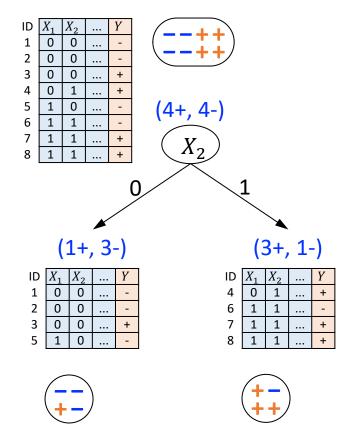
Gatterbauer. Foundations and Applications of Information Theory: https://northeastern-datalab.github.io/cs7840/

To determine how well a test condition performs, we compare the "impurity" of the parent node (before splitting) with the "impurity" of the child nodes (after splitting).

Call the impurity at node N:I(N) currently a black box fct

Impurity before: I(N)

Impurity after: ?



To determine how well a test condition performs, we compare the "impurity" of the parent node (before splitting) with the "impurity" of the child nodes (after splitting).

Call the impurity at node N: I(N)

Impurity before: I(N)

average) impurity

Impurity after:
$$\sum_{C \in \text{children}(N)} p_C \cdot I(C) = \mathbb{E}_{p(C)}[I(C)]$$

infor. gain (IG): $I(N) - \mathbb{E}_{p(C)}[I(C)]$

(3+, 1-)

(4+, 4-)





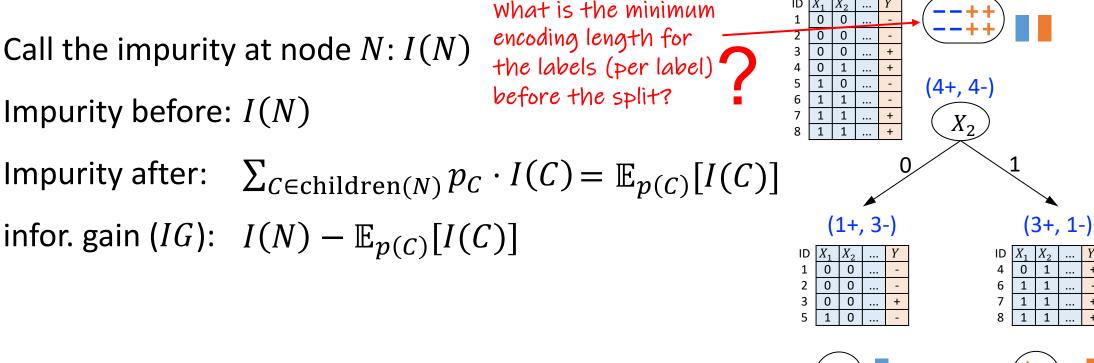
Have we seen that before

To determine how well a test condition performs, we compare the "impurity" of the parent node (before splitting) with the "impurity" of the child nodes (after splitting).

Call the impurity at node N: I(N)

Impurity before: I(N)

encoding length for



To determine how well a test condition performs, we compare the "impurity" of the parent node (before splitting) with the "impurity" of the child nodes (after splitting).

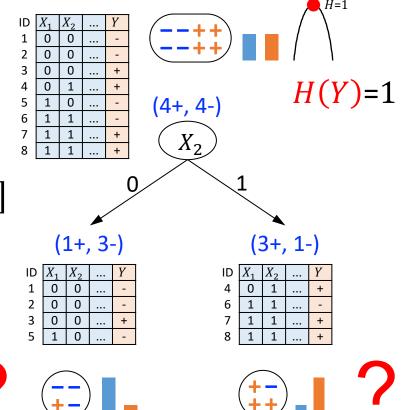
Call the impurity at node N: I(N) = H(Y)

(per label)

Impurity before: I(N)

Impurity after: $\sum_{C \in \text{children}(N)} p_C \cdot I(C) = \mathbb{E}_{p(C)}[I(C)]$

infor. gain (IG): $I(N) - \mathbb{E}_{p(C)}[I(C)]$



To determine how well a test condition performs, we compare the "impurity" of the parent node (before splitting) with the "impurity" of the child nodes (after splitting).

(per label)

Call the impurity at node N: I(N) = H(Y)

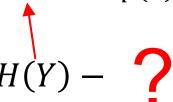
Impurity before: I(N)

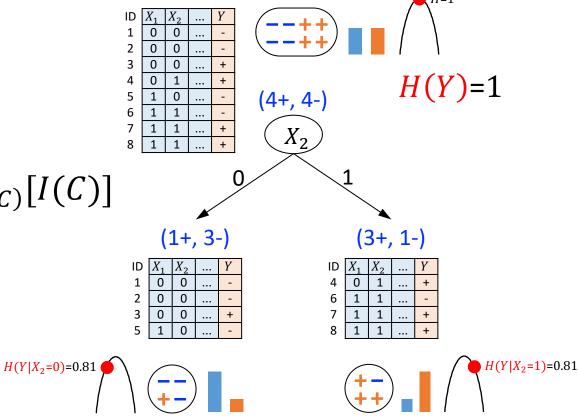
Impurity after: $\sum_{C \in \text{children}(N)} p_C \cdot I(C) = \mathbb{E}_{p(C)}[I(C)]$

infor. gain (IG): $I(N) - \mathbb{E}_{p(C)}[I(C)]$

"information gain"

IG or Δ_{info} :





To determine how well a test condition performs, we compare the "impurity" of the parent node (before splitting) with the "impurity" of the child nodes (after splitting).

(per label)

 $H(Y|X_2=0)=0.81$

Call the impurity at node N: I(N) = H(Y)

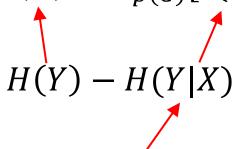
Impurity before: I(N)

Impurity after: $\sum_{C \in \text{children}(N)} p_C \cdot I(C) = \mathbb{E}_{p(C)}[I(C)]$

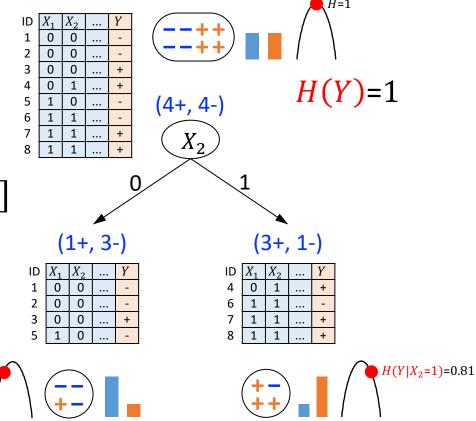
infor. gain (IG): $I(N) - \mathbb{E}_{p(C)}[I(C)]$

"information gain"

IG or Δ_{info} :

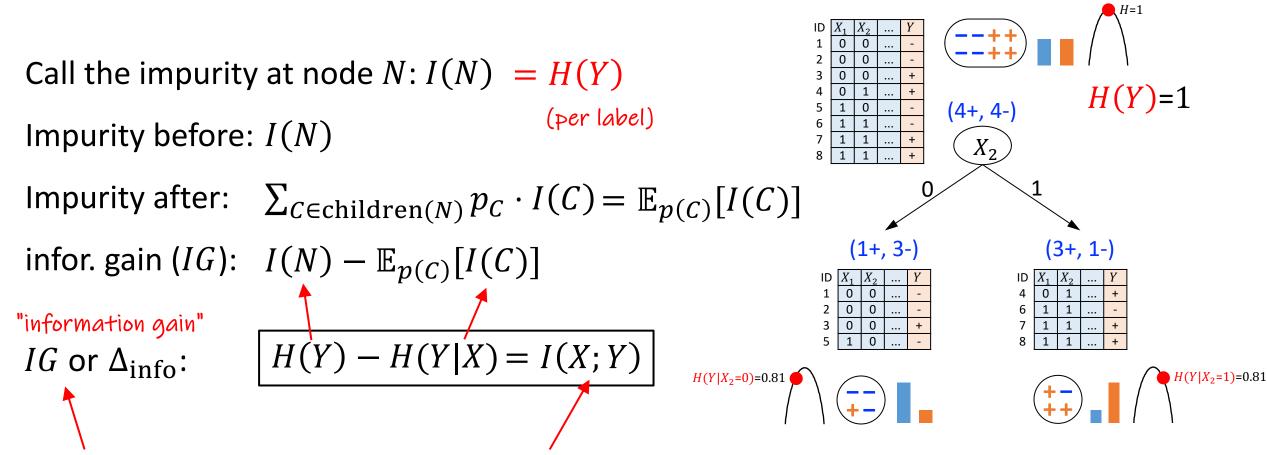


Conditional entropy: the amount of information needed to describe the outcome of RV Y given that we know the value of another RV X.



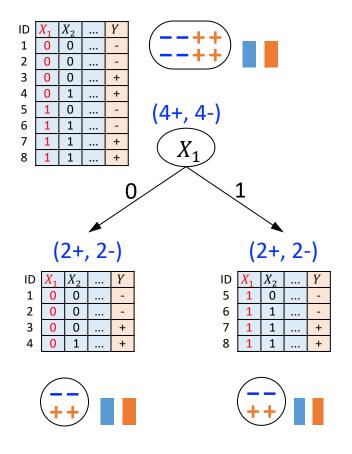
$$H(Y|X) = \sum_{x} p(x) \cdot H(Y|X = x) = \mathbb{E}_{p(x)}[H(Y|X = x)]$$

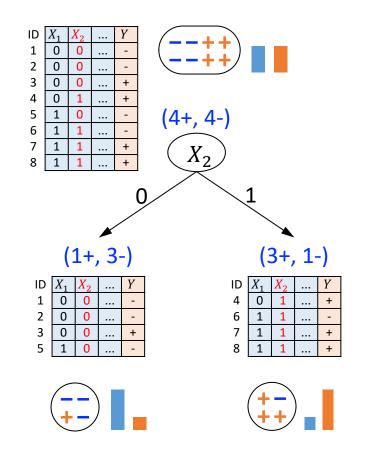
To determine how well a test condition performs, we compare the "impurity" of the parent node (before splitting) with the "impurity" of the child nodes (after splitting).



 Δ_{info} (information gain gain IG) is the reduction of class label entropy H(Y) from the parent (i.e. the training data in a branch) to the average entropies of the children (i.e. the new partitions constructed from the values of variable X.

"Impurity" reduction

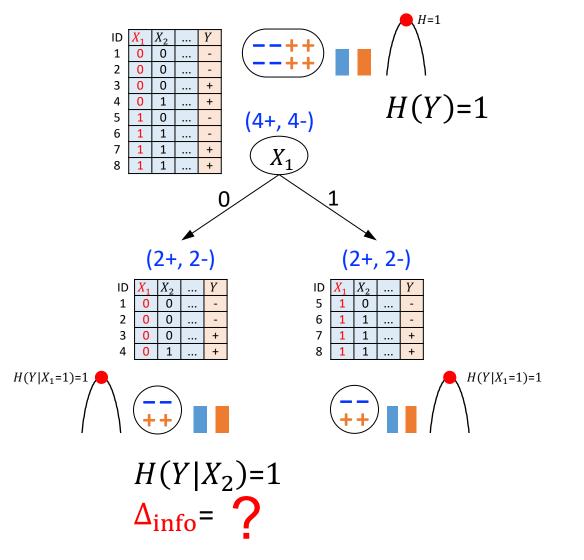


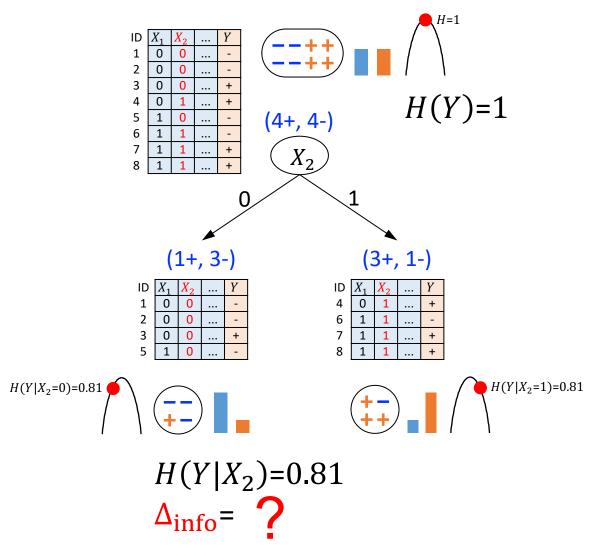




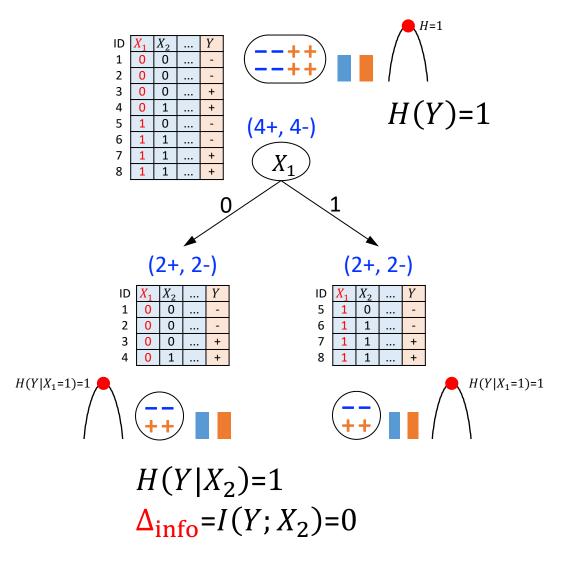
Which variable would you choose ?

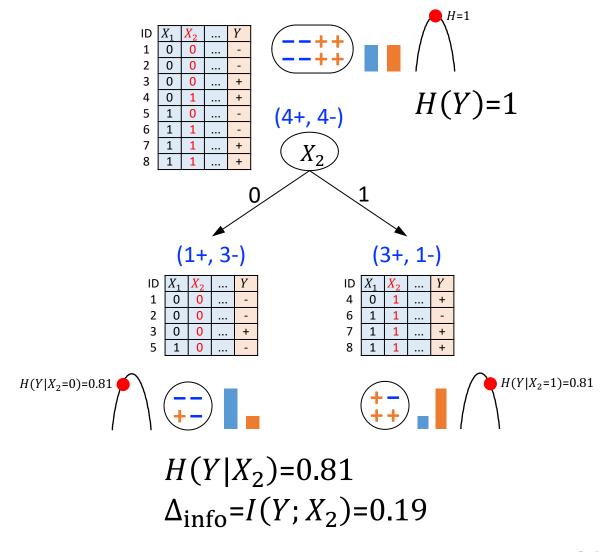
"Impurity" reduction, measured by entropy



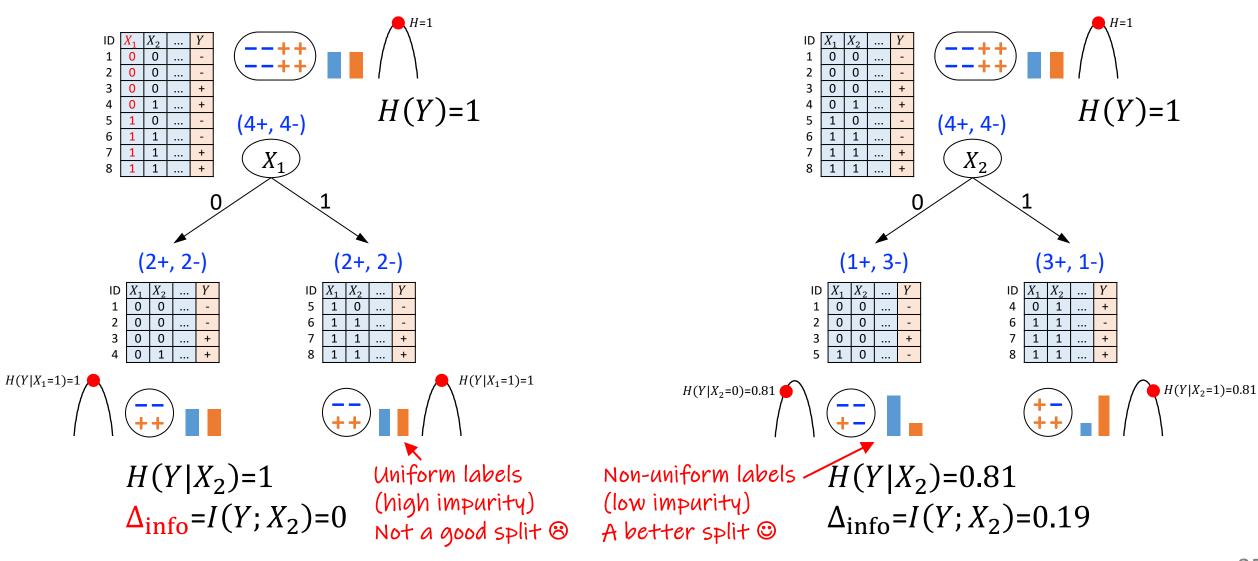


"Impurity" reduction, measured by entropy



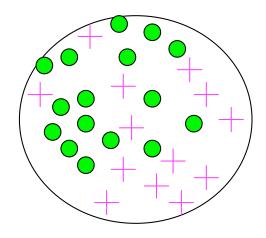


"Impurity" reduction, measured by entropy

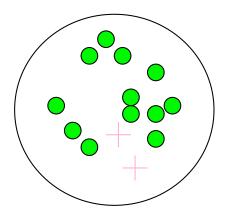


Impurity

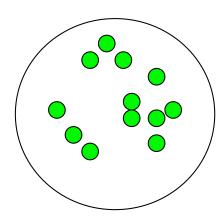
Very impure group



Less impure

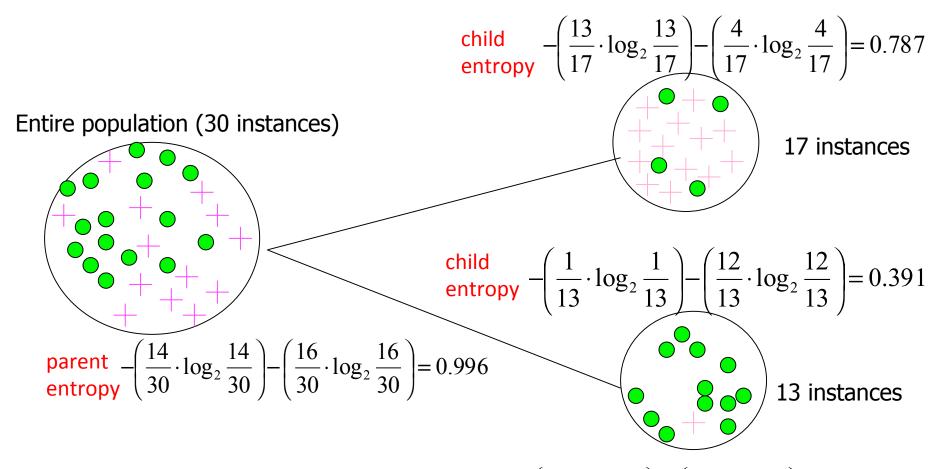


Minimum impurity



Calculating Information Gain

Information Gain = entropy(parent) - [average entropy(children)]



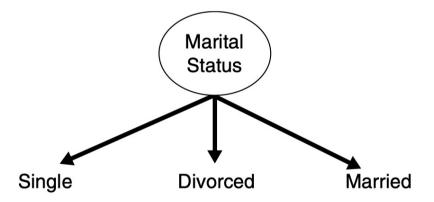
(Weighted) Average Entropy of Children =
$$\left(\frac{17}{30} \cdot 0.787\right) + \left(\frac{13}{30} \cdot 0.391\right) = 0.615$$

Information Gain = 0.996 - 0.615 = 0.38

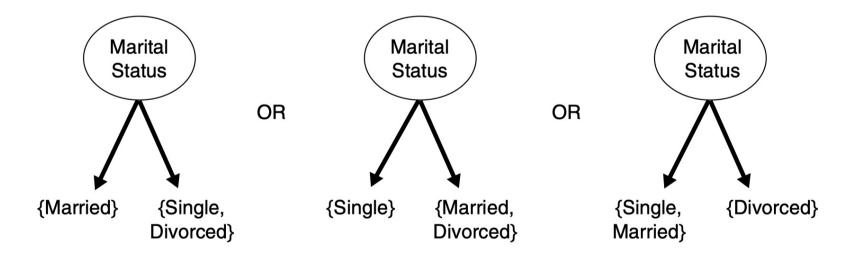
38

Test conditions for nominal attributes

Multiway split



Binary split by grouping attributes



Impurity measures (in addition to entropy)

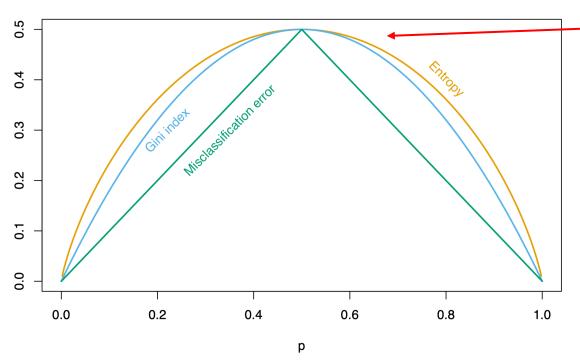
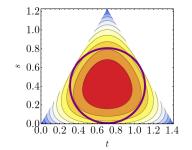


FIGURE 9.3. Node impurity measures for two-class classification, as a function of the proportion p in class 2. Cross-entropy has been scaled to pass through (0.5, 0.5).

Entropy =
$$-\sum_{i=0}^{c-1} p_i(t) \log_2 p_i(t)$$
,
Gini index = $1 - \sum_{i=0}^{c-1} p_i(t)^2$,
Classification error = $1 - \max[p_i(t)]$,

| Another w | au to think | about the | u-axic ic | I(S) |
|--------------------|---------------|---------------|------------|-----------------|
| 7 (10 () 10 () 0 | ory to trivia | . DIDOUT TYPO | 9 01013 13 | $\max_{S} I(S)$ |

| Entropy | Gini index |
|--|--|
| measures the amount of uncertainty (or | measures the probability of misclassifying a |
| randomness) in a set / can be interpreted as | randomly chosen element in a set / can be |
| the average amount of information needed | interpreted as the expected error rate in a |
| to specify the class of an instance. | classifier. |
| The range of entropy is [0, lg(c)], where c is | The range of the Gini index is [0, 1-1/c] |
| the number of classes. | (often incorrectly stated as [0,1]) |
| It has a bias toward selecting splits that | It has a bias toward selecting splits that |
| result in a higher reduction of uncertainty | result in a more balanced (equally sized) |
| (distinguishes more between highly impure | distribution of classes. |
| and moderately impure splits, better for | |
| imbalanced datasets) | |
| Entropy was proposed by Quinlan (ID3 and | Gini index is typically used in CART |
| C4.5 algorithms). Earlier paper from 1979 is | ("Classification and Regression Trees") |
| hard to find, usual citation: "Quinlan. | "Breiman, Friedman, Olshen, Stone: |
| Induction of Decision Trees, 1986." | Classification and Regression Trees, 1984." |
| https://doi.org/10.1007/BF00116251 | https://doi.org/10.1201/9781315139470 |



Simple interesting scribe: create a notebook and figures that compare ternary Gini vs entropy function over the probability simplex. How are the derivatives close to "pure sets"? Two starter posts:

 $\frac{https://physics.stackexchange.com/questions/363545/what-is-the-relation-between-linear-purity-and-von-neumann-entropy-of-a-state}{}$

 $\frac{\text{https://math.stackexchange.com/questions/3791203/what-do-the-level-sets-of-the-shannon-entropy-look-like}{\text{https://math.stackexchange.com/questions/3791203/what-do-the-level-sets-of-the-shannon-entropy-look-like}{\text{https://math.stackexchange.com/questions/3791203/what-do-the-level-sets-of-the-shannon-entropy-look-like}{\text{https://math.stackexchange.com/questions/3791203/what-do-the-level-sets-of-the-shannon-entropy-look-like}{\text{https://math.stackexchange.com/questions/3791203/what-do-the-level-sets-of-the-shannon-entropy-look-like}{\text{https://math.stackexchange.com/questions/3791203/what-do-the-level-sets-of-the-shannon-entropy-look-like}{\text{https://math.stackexchange.com/questions/3791203/what-do-the-level-sets-of-the-shannon-entropy-look-like}{\text{https://math.stackexchange.com/questions/3791203/what-do-the-level-sets-of-the-shannon-entropy-look-like}{\text{https://math.stackexchange.com/questions/3791203/what-do-the-level-sets-of-the-shannon-entropy-look-like}{\text{https://math.stackexchange.com/questions/stackexchange.com/q$

Gini index and "logical entropy"

When there are point probabilities $p = (p_1, ..., p_n)$ for p_j as the probability of the outcome $u_j \in U$ with $\sum_{j=1}^n p_j = 1$, then $\Pr(B_i) = \sum \{p_j : u_j \in B_i\}$ in the formula for logical entropy. This also gives the definition of logical entropy for any probability distribution $p = (p_1, ..., p_n)$,

$$h(p) = 1 - \sum_{j=1}^{n} p_j^2.$$
 (2.3)

$$1 = 1^{2} = (p_{1} + \dots + p_{n})(p_{1} + \dots + p_{n}) = \sum_{j=1}^{n} p_{i}^{2} + \sum_{j \neq k} p_{j} p_{k}$$
(2.4)

so that:

$$h(p) = 1 - \sum_{j=1}^{n} p_i^2 = \sum_{j=1}^{n} p_j (1 - p_j) = \sum_{j \neq k} p_j p_k = 2 \sum_{j < k} p_j p_k$$
 (2.5)

Gini index and "logical entropy"

1.5 Brief History of the Logical Entropy Formula

The logical entropy formula $h(p) = \sum_i p_i (1-p_i) = 1 - \sum_i p_i^2$ is the probability of getting distinct values $u_i \neq u_j$ in two independent samplings of the random variable u. The complementary measure $1 - h(p) = \sum_i p_i^2$ is the probability that the two drawings yield the same value from U. Thus $1 - \sum_i p_i^2$ is a measure of heterogeneity or diversity in keeping with our theme of information as distinctions, while the complementary measure $\sum_i p_i^2$ is a measure of homogeneity or concentration. Historically, the formula can be found in either form depending on the particular context. The p_i 's might be relative shares such as the relative share of organisms of the ith species in some population of organisms, and then the interpretation of p_i as a probability arises by considering the random choice of an organism from the population.

According to I. J. Good, the formula has a certain naturalness: "If p_1, \ldots, p_t are the probabilities of t mutually exclusive and exhaustive events, any statistician of this century who wanted a measure of homogeneity would have take about two seconds to suggest $\sum p_i^2$ which I shall call ρ ." [13, p. 561] As noted by Bhargava and Uppuluri [4], the formula $1 - \sum p_i^2$ was used by Gini in 1912 [10] as a measure of "mutability" or diversity. But another development of the formula (in the complementary form) in the early twentieth century was in cryptography. The American cryptologist, William F. Friedman, devoted a 1922 book [9] to the "index of coincidence" (i.e., $\sum p_i^2$). Solomon Kullback (see the Kullback-Leibler divergence treated later) worked as an assistant to Friedman and wrote a book on cryptology which used the index [16].

"Logical entropy"

2.2 Logical Entropy, Not Shannon Entropy, Is a (Non-negative) Measure

• •

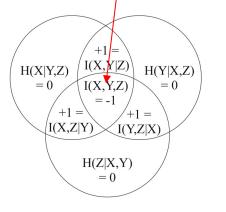
As we will see, for three or more random variables, the Shannon mutual information can have negative values—which has no known interpretation.

Recall that this concern can be easily avoided by more careful notation and *not* using the terminology of "mutual information" for what we called the "interaction information"

4.2 An Example of Negative Mutual Information for Shannon Entropy

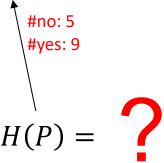
Norman Abramson gives an example [1, pp. 130–131] where the Shannon mutual information of three variables is negative.³ William Feller gives a similar concrete example that we will use [11, Exercise 26, p. 143]. Any probability theory textbook example to show that pair-wise independence does not imply mutual independence for three or more random variables would do as well.

Fig. 4.6 Negative 'area' I(X, Y, Z) in Venn diagram



End-to-end "Tennis" Example by Tom Mitchell

| | | Pred | ictors | | Response | |
|-----|-----------|---------|------------|--------|----------|-------------------------------------|
| day | (O)utlook | (T)emp. | (H)umidity | (W)ind | (P)lay | |
| 1 | sunny | hot | high | weak | no | |
| 2 | sunny | hot | high | strong | no | |
| 3 | overcast | hot | high | weak | yes | |
| 4 | rain | mild | high | weak | yes | $\langle \mathbf{x}_4, y_4 \rangle$ |
| 5 | rain | cool | normal | weak | yes | |
| 6 | rain | cool | normal | strong | no | |
| 7 | overcast | cool | normal | strong | yes | |
| 8 | sunny | mild | high | weak | no | |
| 9 | sunny | cool | normal | weak | yes | |
| 10 | rain | mild | normal | weak | yes | |
| 11 | sunny | mild | normal | strong | yes | |
| 12 | overcast | mild | high | strong | yes | |
| 13 | overcast | hot | normal | weak | yes | |
| 14 | rain | mild | high | strong | no | |
| | | • | | • | | |



Example: Deciding whether to play or not to play tennis on a Saturday (binary classification) Columns denote 4 features X_i . Rows denote labeled instances $\langle \mathbf{x}_i, y_i \rangle$. Play denotes the classification.

| | | Response | | | |
|-----|-----------|----------|------------|--------|--------|
| day | (O)utlook | (T)emp. | (H)umidity | (W)ind | (P)lay |
| 1 | sunny | hot | high | weak | no |
| 2 | sunny | hot | high | strong | no |
| 3 | overcast | hot | high | weak | yes |
| 4 | rain | mild | high | weak | yes |
| 5 | rain | cool | normal | weak | yes |
| 6 | rain | cool | normal | strong | no |
| 7 | overcast | cool | normal | strong | yes |
| 8 | sunny | mild | high | weak | no |
| 9 | sunny | cool | normal | weak | yes |
| 10 | rain | mild | normal | weak | yes |
| 11 | sunny | mild | normal | strong | yes |
| 12 | overcast | mild | high | strong | yes |
| 13 | overcast | hot | normal | weak | yes |
| 14 | rain | mild | high | strong | no |

What happens if we split by attribute W



#no: 5
#yes: 9
$$H(P) = H\left(\frac{9}{14}, \frac{5}{14}\right) = 0.940$$

Example: Deciding whether to play or not to play tennis on a Saturday (binary classification) Columns denote 4 features X_i . Rows denote labeled instances $\langle \mathbf{x}_i, y_i \rangle$. Play denotes the classification.

| | | Pred | ictors | | Response |
|-----|-----------|---------|------------|--------|----------|
| day | (O)utlook | (T)emp. | (H)umidity | (W)ind | (P)lay |
| 1 | sunny | hot | high | weak | no |
| 3 | overcast | hot | high | weak | yes |
| 4 | rain | mild | high | weak | yes |
| 5 | rain | cool | normal | weak | yes |
| 8 | sunny | mild | high | weak | no |
| 9 | sunny | cool | normal | weak | yes |
| 10 | rain | mild | normal | weak | yes |
| 13 | overcast | hot | normal | weak | yes |
| 2 | sunny | hot | high | strong | no |
| 6 | rain | cool | normal | strong | no |
| 7 | overcast | cool | normal | strong | yes |
| 11 | sunny | mild | normal | strong | yes |
| 12 | overcast | mild | high | strong | yes |
| 14 | rain | mild | high | strong | no |

What happens if we split by attribute W?

$$H(P|W) =$$

$$I(P;W) =$$

now partitioned by W
$$H(P) = H\left(\frac{9}{14}, \frac{5}{14}\right) = 0.940$$

Example: Deciding whether to play or not to play tennis on a Saturday (binary classification) Columns denote 4 features X_i . Rows denote labeled instances $\langle \mathbf{x}_i, y_i \rangle$. Play denotes the classification.

| | | Pred | ictors | | Response |
|-----|-----------|---------|------------|--------|----------|
| day | (O)utlook | (T)emp. | (H)umidity | (W)ind | (P)lay |
| 1 | sunny | hot | high | weak | no |
| 3 | overcast | hot | high | weak | yes |
| 4 | rain | mild | high | weak | yes |
| 5 | rain | cool | normal | weak | yes |
| 8 | sunny | mild | high | weak | no |
| 9 | sunny | cool | normal | weak | yes |
| 10 | rain | mild | normal | weak | yes |
| 13 | overcast | hot | normal | weak | yes |
| 2 | sunny | hot | high | strong | no |
| 6 | rain | cool | normal | strong | no |
| 7 | overcast | cool | normal | strong | yes |
| 11 | sunny | mild | normal | strong | yes |
| 12 | overcast | mild | high | strong | yes |
| 14 | rain | mild | high | strong | no |

What happens if we split by attribute W?

$$\mathbb{E}_{p(v)}[H(P|W=v)]$$

$$H(P|W) = \sum_{v} p(v) \cdot H(P|W=v)$$

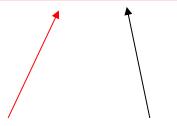
$$H(P|W=\text{weak}) =$$

$$H(P|W=\text{strong}) =$$

$$I(P;W) =$$
 now partitioned by W
$$H(P) = H\left(\frac{9}{14}, \frac{5}{14}\right) = 0.940$$

Example: Deciding whether to play or not to play tennis on a Saturday (binary classification) Columns denote 4 features X_i . Rows denote labeled instances $\langle \mathbf{x}_i, y_i \rangle$. Play denotes the classification.

| | | Predictors | | | | | | |
|-----|-----------|------------|------------|--------|--------|--|--|--|
| day | (O)utlook | (T)emp. | (H)umidity | (W)ind | (P)lay | | | |
| 1 | sunny | hot | high | weak | no | | | |
| 3 | overcast | hot | high | weak | yes | | | |
| 4 | rain | mild | high | weak | yes | | | |
| 5 | rain | cool | normal | weak | yes | | | |
| 8 | sunny | mild | high | weak | no | | | |
| 9 | sunny | cool | normal | weak | yes | | | |
| 10 | rain | mild | normal | weak | yes | | | |
| 13 | overcast | hot | normal | weak | yes | | | |
| 2 | sunny | hot | high | strong | no | | | |
| 6 | rain | cool | normal | strong | no | | | |
| 7 | overcast | cool | normal | strong | yes | | | |
| 11 | sunny | mild | normal | strong | yes | | | |
| 12 | overcast | mild | high | strong | yes | | | |
| 14 | rain | mild | high | strong | no | | | |



What happens if we split by attribute W?

$$\mathbb{E}_{p(v)}[H(P|W = v)]$$

$$H(P|W) = \sum_{v} p(v) \cdot H(P|W = v)$$

$$H(P|W = \text{weak}) = H\left(\frac{2}{8}, \frac{6}{8}\right) = 0.811$$

$$H(P|W = \text{strong}) = H\left(\frac{3}{6}, \frac{3}{6}\right) = 1$$

$$H(P|W) = \frac{8}{14} \cdot 0.811 + \frac{6}{14} \cdot 1 = 0.892$$

$$I(P; W) = \sum_{v} p(v) \cdot H(P|W = v)$$

$$H(P|W = v)$$

$$H($$

now partitioned by W
$$H(P) = H\left(\frac{9}{14}, \frac{5}{14}\right) = 0.940$$

Example: Deciding whether to play or not to play tennis on a Saturday (binary classification) Columns denote 4 features X_i . Rows denote labeled instances $\langle \mathbf{x}_i, y_i \rangle$. Play denotes the classification.

| | | Predictors | | | | | | |
|-----|-----------|------------|------------|--------|--------|--|--|--|
| day | (O)utlook | (T)emp. | (H)umidity | (W)ind | (P)lay | | | |
| 1 | sunny | hot | high | weak | no | | | |
| 3 | overcast | hot | high | weak | yes | | | |
| 4 | rain | mild | high | weak | yes | | | |
| 5 | rain | cool | normal | weak | yes | | | |
| 8 | sunny | mild | high | weak | no | | | |
| 9 | sunny | cool | normal | weak | yes | | | |
| 10 | rain | mild | normal | weak | yes | | | |
| 13 | overcast | hot | normal | weak | yes | | | |
| 2 | sunny | hot | high | strong | no | | | |
| 6 | rain | cool | normal | strong | no | | | |
| 7 | overcast | cool | normal | strong | yes | | | |
| 11 | sunny | mild | normal | strong | yes | | | |
| 12 | overcast | mild | high | strong | yes | | | |
| 14 | rain | mild | high | strong | no | | | |

What happens if we split by attribute W?

$$\mathbb{E}_{p(v)}[H(P|W=v)]$$

$$H(P|W) = \sum_{v} p(v) \cdot H(P|W=v)$$

$$H(P|W=\text{weak}) = H\left(\frac{2}{8}, \frac{6}{8}\right) = 0.811$$

$$H(P|W=\text{strong}) = H\left(\frac{3}{6}, \frac{3}{6}\right) = 1$$

$$H(P|W) = \frac{8}{14} \cdot 0.811 + \frac{6}{14} \cdot 1 = 0.892$$

$$I(P;W) = H(P) - H(P|W) = 0.048$$

now partitioned by W $H(P) = H\left(\frac{9}{14}, \frac{5}{14}\right) = 0.940$

Example: Deciding whether to play or not to play tennis on a Saturday (binary classification) Columns denote 4 features X_i . Rows denote labeled instances $\langle \mathbf{x}_i, y_i \rangle$. Play denotes the classification.

| | | Response | | | |
|-----|-----------|----------|------------|--------|--------|
| day | (O)utlook | (T)emp. | (H)umidity | (W)ind | (P)lay |
| 1 | sunny | hot | high | weak | no |
| 2 | sunny | hot | high | strong | no |
| 3 | overcast | hot | high | weak | yes |
| 4 | rain | mild | high | weak | yes |
| 5 | rain | cool | normal | weak | yes |
| 6 | rain | cool | normal | strong | no |
| 7 | overcast | cool | normal | strong | yes |
| 8 | sunny | mild | high | weak | no |
| 9 | sunny | cool | normal | weak | yes |
| 10 | rain | mild | normal | weak | yes |
| 11 | sunny | mild | normal | strong | yes |
| 12 | overcast | mild | high | strong | yes |
| 13 | overcast | hot | normal | weak | yes |
| 14 | rain | mild | high | strong | no |

Now we calculate mutual information (aka information gain) between P and the other attributes

$$I(P; H) = 0.152$$

 $I(P; T) = 0.029$
 $I(P; O) = 0.246$

$$I(P; W) = H(P) - H(P|W) = 0.048$$

Which attribute do we pick?



Example: Deciding whether to play or not to play tennis on a Saturday (binary classification) Columns denote 4 features X_i . Rows denote labeled instances $\langle \mathbf{x}_i, y_i \rangle$. Play denotes the classification.

| | | Response | | | |
|-----|-----------|----------|------------|--------|--------|
| day | (O)utlook | (T)emp. | (H)umidity | (W)ind | (P)lay |
| 1 | sunny | hot | high | weak | no |
| 2 | sunny | hot | high | strong | no |
| 3 | overcast | hot | high | weak | yes |
| 4 | rain | mild | high | weak | yes |
| 5 | rain | cool | normal | weak | yes |
| 6 | rain | cool | normal | strong | no |
| 7 | overcast | cool | normal | strong | yes |
| 8 | sunny | mild | high | weak | no |
| 9 | sunny | cool | normal | weak | yes |
| 10 | rain | mild | normal | weak | yes |
| 11 | sunny | mild | normal | strong | yes |
| 12 | overcast | mild | high | strong | yes |
| 13 | overcast | hot | normal | weak | yes |
| 14 | rain | mild | high | strong | no |

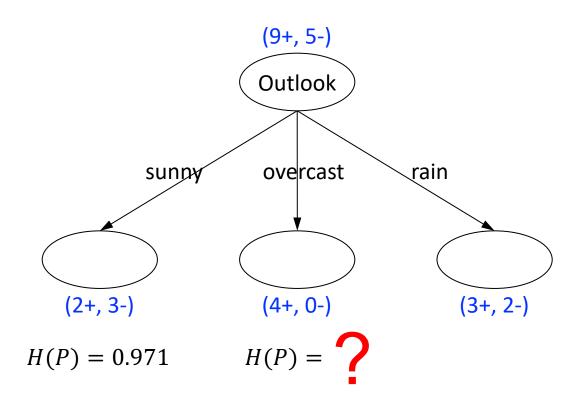
Now we calculate mutual information (aka information gain) between P and the other attributes

$$I(P; W) = I(P; W) = I(P; W) = I(P; T) = 0.029$$
 $I(P; O) = 0.246$
We pick the h

$$I(P; W) = H(P) - H(P|W) = 0.048$$

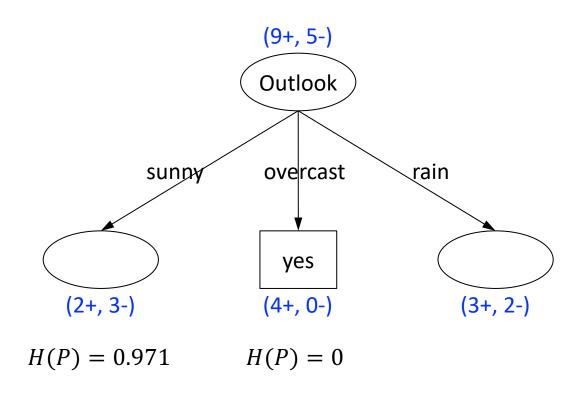
We pick the attribute with the highest information gain

| | | Response | | | |
|-----|-----------|----------|------------|--------|--------|
| day | (O)utlook | (T)emp. | (H)umidity | (W)ind | (P)lay |
| 1 | sunny | hot | high | weak | no |
| 2 | sunny | hot | high | strong | no |
| 8 | sunny | mild | high | weak | no |
| 9 | sunny | cool | normal | weak | yes |
| 11 | sunny | mild | normal | strong | yes |
| 3 | overcast | hot | high | weak | yes |
| 7 | overcast | cool | normal | strong | yes |
| 12 | overcast | mild | high | strong | yes |
| 13 | overcast | hot | normal | weak | yes |
| 4 | rain | mild | high | weak | yes |
| 5 | rain | cool | normal | weak | yes |
| 6 | rain | cool | normal | strong | no |
| 10 | rain | mild | normal | weak | yes |
| 14 | rain | mild | high | strong | no |





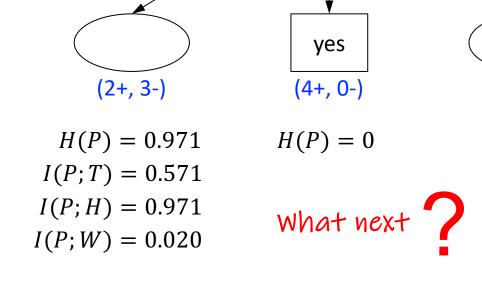
| | | Response | | | |
|-----|-----------|----------|------------|--------|--------|
| day | (O)utlook | (T)emp. | (H)umidity | (W)ind | (P)lay |
| 1 | sunny | hot | high | weak | no |
| 2 | sunny | hot | high | strong | no |
| 8 | sunny | mild | high | weak | no |
| 9 | sunny | cool | normal | weak | yes |
| 11 | sunny | mild | normal | strong | yes |
| 3 | overcast | hot | high | weak | yes |
| 7 | overcast | cool | normal | strong | yes |
| 12 | overcast | mild | high | strong | yes |
| 13 | overcast | hot | normal | weak | yes |
| 4 | rain | mild | high | weak | yes |
| 5 | rain | cool | normal | weak | yes |
| 6 | rain | cool | normal | strong | no |
| 10 | rain | mild | normal | weak | yes |
| 14 | rain | mild | high | strong | no |





Example: Deciding whether to play or not to play tennis on a Saturday (binary classification) Columns denote 4 features X_i . Rows denote labeled instances $\langle \mathbf{x}_i, y_i \rangle$. Play denotes the classification.

| | | Predictors | | | | | | |
|-----|-----------|------------|------------|--------|--------|--|--|--|
| day | (O)utlook | (T)emp. | (H)umidity | (W)ind | (P)lay | | | |
| 1 | sunny | hot | high | weak | no | | | |
| 2 | sunny | hot | high | strong | no | | | |
| 8 | sunny | mild | high | weak | no | | | |
| 9 | sunny | cool | normal | weak | yes | | | |
| 11 | sunny | mild | normal | strong | yes | | | |
| 3 | overcast | hot | high | weak | yes | | | |
| 7 | overcast | cool | normal | strong | yes | | | |
| 12 | overcast | mild | high | strong | yes | | | |
| 13 | overcast | hot | normal | weak | yes | | | |
| 4 | rain | mild | high | weak | yes | | | |
| 5 | rain | cool | normal | weak | yes | | | |
| 6 | rain | cool | normal | strong | no | | | |
| 10 | rain | mild | normal | weak | yes | | | |
| 14 | rain | mild | high | strong | no | | | |



sunny

(9+, 5-)

Outlook

overcast

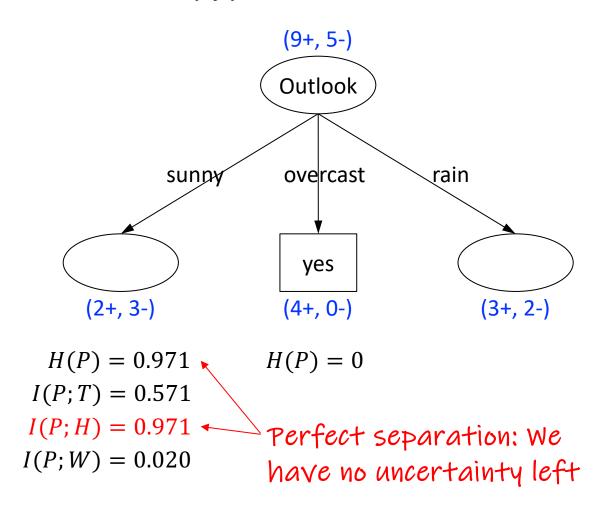
rain

(3+, 2-)

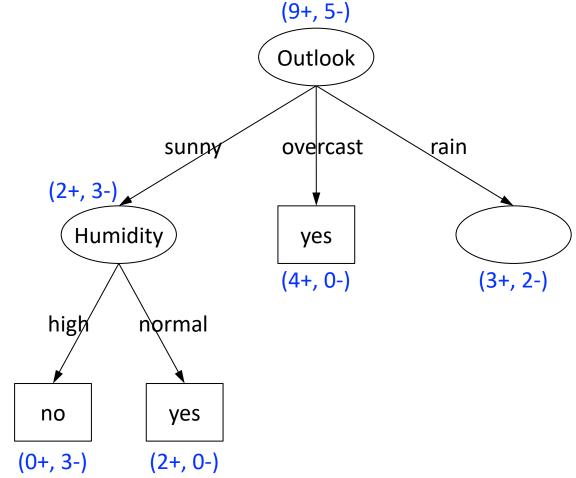
now partitioned by O

| | | Predictors | | | | | |
|-----|-----------|------------|------------|--------|--------|--|--|
| day | (O)utlook | (T)emp. | (H)umidity | (W)ind | (P)lay | | |
| 1 | sunny | hot | high | weak | no | | |
| 2 | sunny | hot | high | strong | no | | |
| 8 | sunny | mild | high | weak | no | | |
| 9 | sunny | cool | normal | weak | yes | | |
| 11 | sunny | mild | normal | strong | yes | | |
| 3 | overcast | hot | high | weak | yes | | |
| 7 | overcast | cool | normal | strong | yes | | |
| 12 | overcast | mild | high | strong | yes | | |
| 13 | overcast | hot | normal | weak | yes | | |
| 4 | rain | mild | high | weak | yes | | |
| 5 | rain | cool | normal | weak | yes | | |
| 6 | rain | cool | normal | strong | no | | |
| 10 | rain | mild | normal | weak | yes | | |
| 14 | rain | mild | high | strong | no | | |

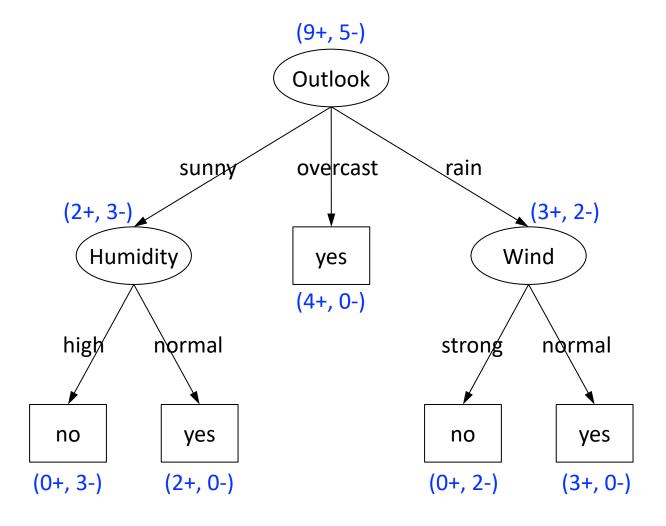




| | | Pred | ictors | | Response |
|-----|-----------|-------------|------------|-------------|----------|
| day | (O)utlook | (T)emp. | (H)umidity | (W)ind | (P)lay |
| 1 | sunny | hot | high | weak | no |
| 2 | sunny | hot | high | strong | no |
| 8 | sunny | mild | high | weak | no |
| 9 | sunny | cool | normal▲ | weak | yes |
| 11 | sunny | mild | normal | strong | yes |
| 3 | overcast | hot | high | weak | yes |
| 7 | overcast | cool | normal | strong | yes |
| 12 | overcast | mild | high | strong | yes |
| 13 | overcast | hot | normal | weak | yes |
| 4 | rain | mild | high | weak | yes |
| 5 | rain | cool | normal | weak | yes |
| 6 | rain | cool | normal | strong | no |
| 10 | rain | mild | normal | weak | yes |
| 14 | rain | mild | high | strong | no |
| | | fu | rther | ı Partit | -ioned b |
| И | ow par | - tition | ied by (|) | |



| | | Response | | | |
|-----|-----------|----------|------------|--------|--------|
| day | (O)utlook | (T)emp. | (H)umidity | (W)ind | (P)lay |
| 1 | sunny | hot | high | weak | no |
| 2 | sunny | hot | high | strong | no |
| 8 | sunny | mild | high | weak | no |
| 9 | sunny | cool | normal | weak | yes |
| 11 | sunny | mild | normal | strong | yes |
| 3 | overcast | hot | high | weak | yes |
| 7 | overcast | cool | normal | strong | yes |
| 12 | overcast | mild | high | strong | yes |
| 13 | overcast | hot | normal | weak | yes |
| 6 | rain | cool | normal | strong | no |
| 14 | rain | mild | high | strong | no |
| 4 | rain | mild | high | weak | yes |
| 5 | rain | cool | normal | weak | yes |
| 10 | rain | mild | normal | weak | yes |



Gain ratio

Example: Deciding whether to play or not to play tennis on a Saturday (binary classification) Columns denote 4 features X_i . Rows denote labeled instances $\langle \mathbf{x}_i, y_i \rangle$. Play denotes the classification.

| | | | Response | | |
|-----|-----------|---------|------------|--------|--------|
| day | (O)utlook | (T)emp. | (H)umidity | (W)ind | (P)lay |
| 1 | sunny | hot | high | weak | no |
| 2 | sunny | hot | high | strong | no |
| 3 | overcast | hot | high | weak | yes |
| 4 | rain | mild | high | weak | yes |
| 5 | rain | cool | normal | weak | yes |
| 6 | rain | cool | normal | strong | no |
| 7 | overcast | cool | normal | strong | yes |
| 8 | sunny | mild | high | weak | no |
| 9 | sunny | cool | normal | weak | yes |
| 10 | rain | mild | normal | weak | yes |
| 11 | sunny | mild | normal | strong | yes |
| 12 | overcast | mild | high | strong | yes |
| 13 | overcast | hot | normal | weak | yes |
| 14 | rain | mild | high | strong | no |

I was missing a better predictor. Which one



$$\Delta_{\rm info}$$

$$I(P; O) = 0.246$$

 $I(P; H) = 0.152$
 $I(P; W) = 0.048$
 $I(P; T) = 0.029$

$$I(P; H) = 0.152$$

$$I(P;W) = 0.048$$

$$I(P;T) = 0.029$$

$$H(P) = 0.940$$

Example: Deciding whether to play or not to play tennis on a Saturday (binary classification) Columns denote 4 features X_i . Rows denote labeled instances $\langle \mathbf{x}_i, y_i \rangle$. Play denotes the classification.

| | Predictors | | | | | | |
|-------|------------|---------|------------|--------|--------|--|--|
| (D)ay | (O)utlook | (T)emp. | (H)umidity | (W)ind | (P)lay | | |
| 1 | sunny | hot | high | weak | no | | |
| 2 | sunny | hot | high | strong | no | | |
| 3 | overcast | hot | high | weak | yes | | |
| 4 | rain | mild | high | weak | yes | | |
| 5 | rain | cool | normal | weak | yes | | |
| 6 | rain | cool | normal | strong | no | | |
| 7 | overcast | cool | normal | strong | yes | | |
| 8 | sunny | mild | high | weak | no | | |
| 9 | sunny | cool | normal | weak | yes | | |
| 10 | rain | mild | normal | weak | yes | | |
| 11 | sunny | mild | normal | strong | yes | | |
| 12 | overcast | mild | high | strong | yes | | |
| 13 | overcast | hot | normal | weak | yes | | |
| 14 | rain | mild | high | strong | no | | |

The day has the highest mutual information

$$\Delta_{\rm info}$$

$$I(P; D) = 0.940$$

$$I(P; O) = 0.246$$

$$I(P;H) = 0.152$$

$$I(P; W) = 0.048$$

 $I(P; T) = 0.029$

$$I(P;T) = 0.029$$

$$H(P) = 0.940$$

Gain ratio

Disadvantage of information gain: It prefers attributes with <u>large number of values</u> that split the data into <u>small, pure subsets</u>

Quinlan's gain ratio (introduced with C4.5) uses normalization on the splitting criterion, i.e. it takes into account the number of outcomes produced by the attribute test condition and their distribution.

The **gain ratio** penalizes attributes such as Date by incorporating a term, called **split information**, that is sensitive Gain to how broadly and uniformly the attribute splits the data

Gain ratio =
$$\frac{\Delta_{info}}{split info}$$

The "split information" is just the entropy of the "split distribution", i.e. the distribution of the values of the attribute on which we split

If all are values balanced and we have k values, then $= \ln(k)$

Example: Deciding whether to play or not to play tennis on a Saturday (binary classification) Columns denote 4 features X_i . Rows denote labeled instances $\langle \mathbf{x}_i, y_i \rangle$. Play denotes the classification.

| | Predictors | | | | | | |
|-------|------------|---------|------------|--------|--------|--|--|
| (D)ay | (O)utlook | (T)emp. | (H)umidity | (W)ind | (P)lay | | |
| 1 | sunny | hot | high | weak | no | | |
| 2 | sunny | hot | high | strong | no | | |
| 3 | overcast | hot | high | weak | yes | | |
| 4 | rain | mild | high | weak | yes | | |
| 5 | rain | cool | normal | weak | yes | | |
| 6 | rain | cool | normal | strong | no | | |
| 7 | overcast | cool | normal | strong | yes | | |
| 8 | sunny | mild | high | weak | no | | |
| 9 | sunny | cool | normal | weak | yes | | |
| 10 | rain | mild | normal | weak | yes | | |
| 11 | sunny | mild | normal | strong | yes | | |
| 12 | overcast | mild | high | strong | yes | | |
| 13 | overcast | hot | normal | weak | yes | | |
| 14 | rain | mild | high | strong | no | | |

The day has the highest mutual information

| $\Delta_{	ext{info}}$ | split info | gain ratio |
|----------------------------|--------------|--|
| $I(P; \mathbf{D}) = 0.940$ | H(D) = 3.81 | 0.247 |
| I(P; O) = 0.246 | H(0) = 1.58 | 0.156 |
| I(P;H) = 0.152 | H(H) = 1 | 0.152 |
| I(P;W) = 0.048 | H(W) = 0.99 | 0.048 |
| I(P;T) = 0.029 | H(T) = 1.56 | 0.019 |
| H(P) = 0.940 | Gain ratio = | $=\frac{\Delta_{\mathrm{info}}}{\mathrm{splitinfo}}$ |

The normalization still does *not* help here. It would help if the data set was bigger as H(D) grows with the size of the dataset, while H(O) would stay the same.

Updated 10/23/2025

Part 3: Applications L13: Decision trees (2/2) and MDL

Wolfgang Gatterbauer

cs7840 Foundations and Applications of Information Theory (fa25)

https://northeastern-datalab.github.io/cs7840/fa25/

10/23/2025

Pre-class conversations

- Last class recapitulation
- Projects: I have time to talk today after class

- Today:
 - Decision trees
 - Occam, MDL
 - Logistic regression

The Parity Function

Expressiveness

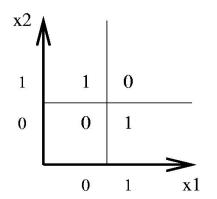
Decision trees have a variable-sized hypothesis space

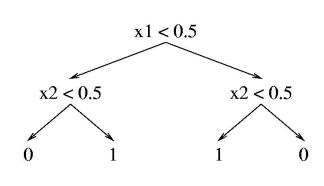
 As the #nodes (or depth) increases, the hypothesis space grows

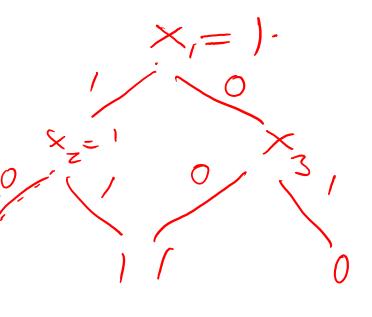
 Depth 1 ("decision stump"): can represent any boolean function of one feature

– Depth 2: any boolean fn of two features; some involving three features (e.g., $(x_1 \land x_2) \lor (\neg x_1 \land \neg x_3)$)

etc.



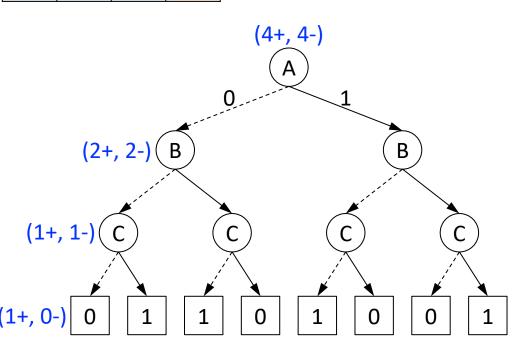






| Predict | Resp. | | |
|---------|-------|---|---|
| Α | В | С | Υ |
| 0 | 0 | 0 | 0 |
| 0 | 0 | 1 | 1 |
| 0 | 1 | 0 | 1 |
| 0 | 1 | 1 | 0 |
| 1 | 0 | 0 | 1 |
| 1 | 0 | 1 | 0 |
| 1 | 1 | 0 | 0 |
| 1 | 1 | 1 | 1 |

Decision Tree for parity function of 3 Boolean attributes: true iff an odd number of variables are true



$$H(Y) = ?$$

$$H(Y|A) = ?$$

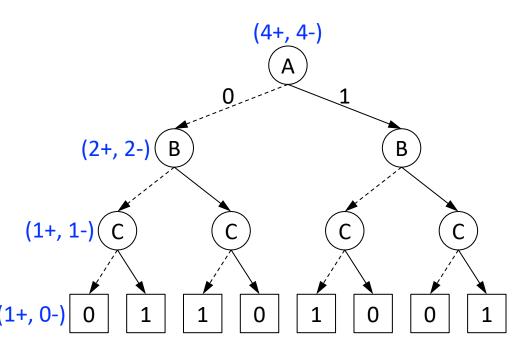
$$H(Y|B) = ?$$

$$H(Y|A,B) = ?$$

$$H(Y|A,B,C) = ?$$

| Predict | Resp. | | |
|---------|-------|---|---|
| Α | В | С | Υ |
| 0 | 0 | 0 | 0 |
| 0 | 0 | 1 | 1 |
| 0 | 1 | 0 | 1 |
| 0 | 1 | 1 | 0 |
| 1 | 0 | 0 | 1 |
| 1 | 0 | 1 | 0 |
| 1 | 1 | 0 | 0 |
| 1 | 1 | 1 | 1 |

Decision Tree for parity function of 3 Boolean attributes: true iff an odd number of variables are true



$$H(Y) = 1$$

$$H(Y|A) = 0.5 \cdot 1 + 0.5 \cdot 1 = 1$$

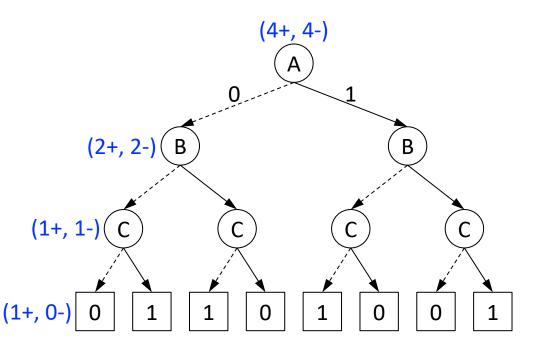
$$H(Y|B) = 0.5 \cdot 1 + 0.5 \cdot 1 = 1$$

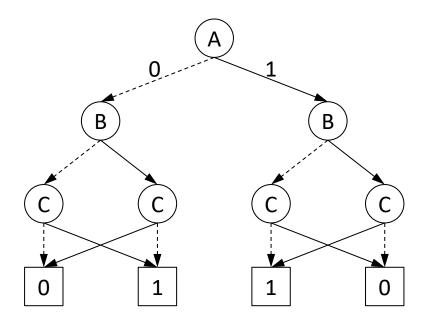
$$H(Y|A,B) = 0.5 \cdot 1 + 0.5 \cdot 1 = 1$$

$$H(Y|A,B,C) = 0$$

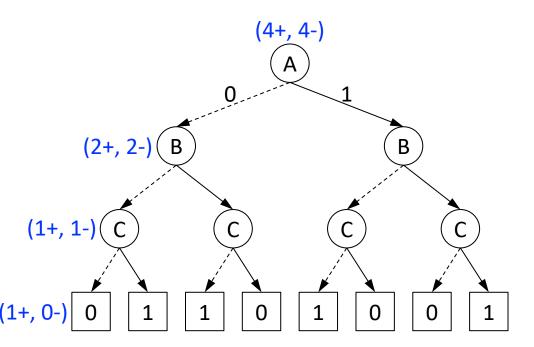
Decision trees vs. circuits

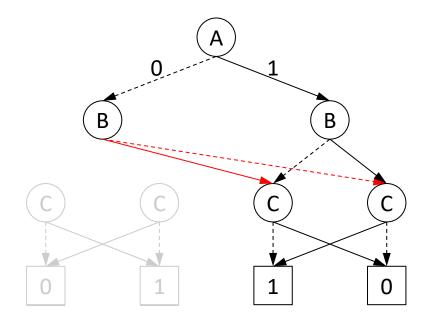
| Predict | Resp. | | |
|---------|-------|---|---|
| Α | В | С | Υ |
| 0 | 0 | 0 | 0 |
| 0 | 0 | 1 | 1 |
| 0 | 1 | 0 | 0 |
| 0 | 1 | 1 | 1 |
| 1 | 0 | 0 | 0 |
| 1 | 0 | 1 | 1 |
| 1 | 1 | 0 | 0 |
| 1 | 1 | 1 | 1 |





| Predict | Resp. | | |
|---------|-------|---|---|
| Α | В | С | Υ |
| 0 | 0 | 0 | 0 |
| 0 | 0 | 1 | 1 |
| 0 | 1 | 0 | 0 |
| 0 | 1 | 1 | 1 |
| 1 | 0 | 0 | 0 |
| 1 | 0 | 1 | 1 |
| 1 | 1 | 0 | 0 |
| 1 | 1 | 1 | 1 |

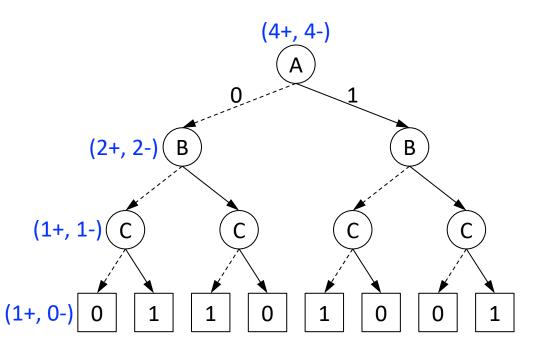


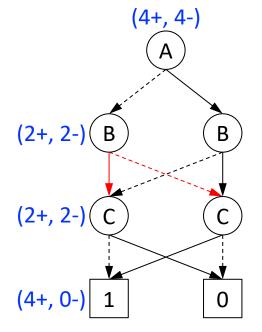


| Predict | Predictors | | |
|---------|------------|---|---|
| Α | В | C | Υ |
| 0 | 0 | 0 | 0 |
| 0 | 0 | 1 | 1 |
| 0 | 1 | 0 | 0 |
| 0 | 1 | 1 | 1 |
| 1 | 0 | 0 | 0 |
| 1 | 0 | 1 | 1 |
| 1 | 1 | 0 | 0 |
| 1 | 1 | 1 | 1 |

The DT grows exponentially with the number of attributes (linearly in the size of the truth table).

The OBDD (Ordered Binary Decision Diagrams) grows linearly in number of attributes (exponentially more succinct than truth table)





Overfitting

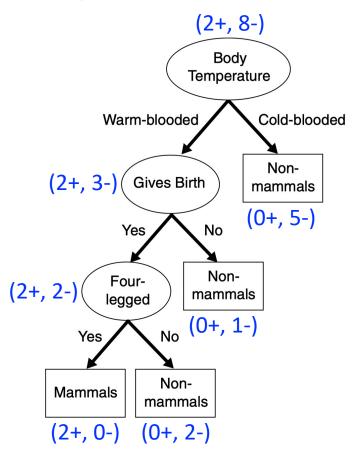
Overfitting due to presence of noise

Training set

name
porcupine
cat
bat
whale
salamander
komodo dragon
python
salmon
eagle
guppy

| P | Predictors | | | | |
|------------|-------------|--------|--------|--|--|
| Body Temp. | Gives Birth | 4 legs | Mammal | | |
| warm | yes | yes | yes | | |
| warm | yes | yes | yes | | |
| warm | yes | no | no | | |
| warm | yes | no | no | | |
| cold | no | yes | no | | |
| cold | no | yes | no | | |
| cold | no | no | no | | |
| cold | no | no | no | | |
| warm | no | no | no | | |
| cold | yes | no | no | | |

DT 1



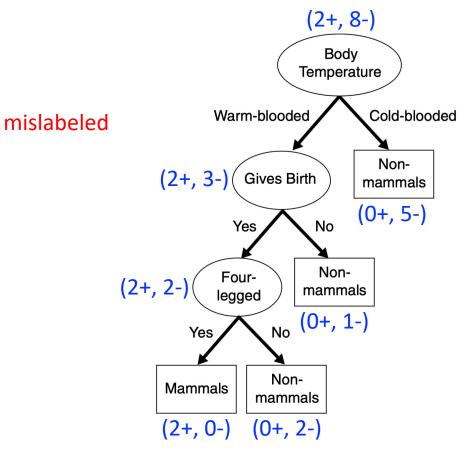
0% training error

Overfitting due to presence of noise

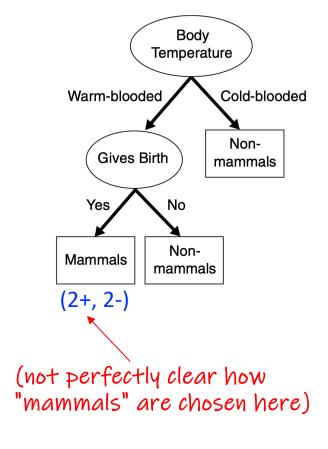
Training set

| | Predictors | | | Label |
|---------------|------------|-------------|--------|--------|
| name | Body Temp. | Gives Birth | 4 legs | Mammal |
| porcupine | warm | yes | yes | yes |
| cat | warm | yes | yes | yes |
| bat | warm | yes | no | no |
| whale | warm | yes | no | no |
| salamander | cold | no | yes | no |
| komodo dragon | cold | no | yes | no |
| python | cold | no | no | no |
| salmon | cold | no | no | no |
| eagle | warm | no | no | no |
| guppy | cold | yes | no | no |

DT 1



DT 2



0% training error

20% training error

Overfitting due to presence of noise

Training set

| | Predictors | | | Label |
|---------------|------------|-------------|--------|--------|
| name | Body Temp. | Gives Birth | 4 legs | Mammal |
| porcupine | warm | yes | yes | yes |
| cat | warm | yes | yes | yes |
| bat | warm | yes | no | no |
| whale | warm | yes | no | no |
| salamander | cold | no | yes | no |
| komodo dragon | cold | no | yes | no |
| python | cold | no | no | no |
| salmon | cold | no | no | no |
| eagle | warm | no | no | no |
| guppy | cold | yes | no | no |

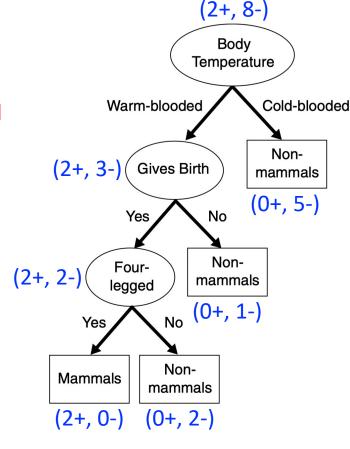
mislabeled

Test set

human
pigeon
elephant
leopard shark
turtle
penguin
eel
dolphin
spiny anteater
gila monster

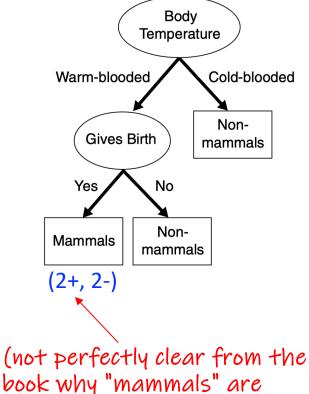
| yes | no | yes |
|-----|----------------------------|--|
| no | no | no |
| yes | yes | yes |
| yes | no | no |
| no | yes | no |
| no | no | no |
| no | no | no |
| yes | no | yes |
| no | yes | yes |
| no | yes | no |
| | no yes yes no no no yes no | no no yes yes yes no no yes no no no no no yes no no yes no no yes no no yes |

DT 1



0% training error 30% test error

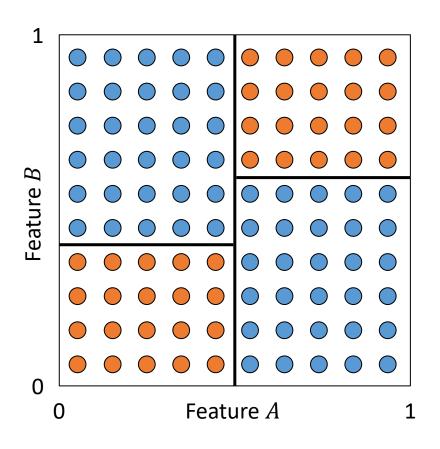
DT 2

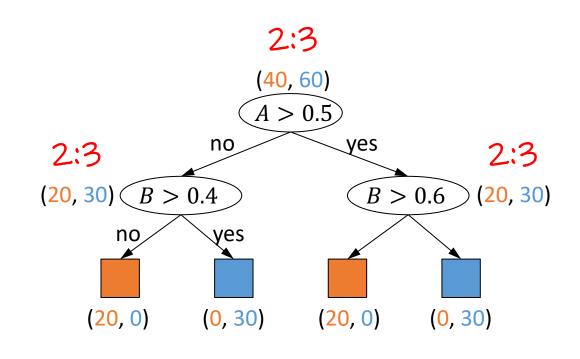


book why "mammals" are chosen here)

20% training error10% test error

Practical considerations

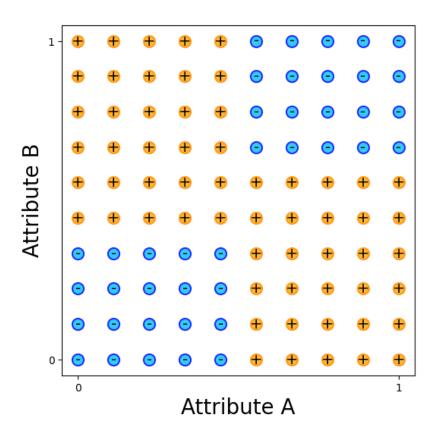




What will happen if we apply information gain



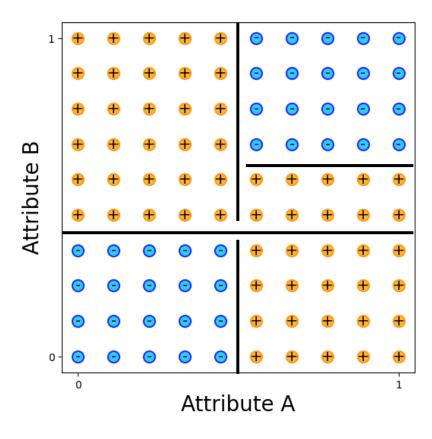
Decision Tree Classification

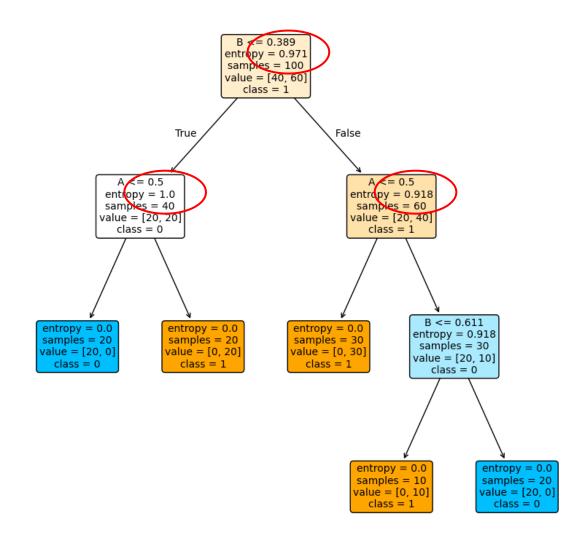


What will happen if we apply information gain



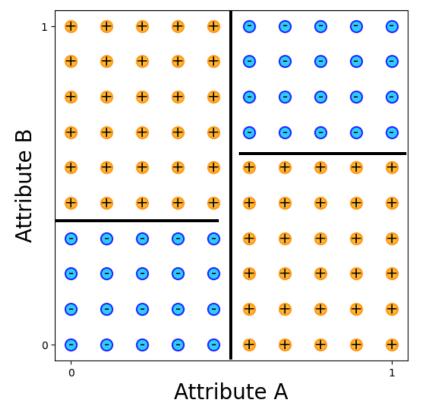
Decision Tree Classification





Decision Tree Classification

as of October 2024



I could not figure out how to force scikit-learn to create the original smaller decision tree shown on the left (max_depth did not help either).

Can you?





While you can't directly force a decision tree in scikit-learn to use a particular attribute first, you can influence its behavior by:

1. Feature Engineering:

Create a new feature:

Combine the attribute you want to prioritize with other features or create a new feature based on its transformations. This can increase its importance in the decision-making process.

Scale the feature:

If the attribute has a different scale compared to other features, scaling it can make it more prominent in the tree's decision-making.

2. Hyperparameter Tuning:

max depth

Limiting the maximum depth of the tree can prevent it from exploring deeper levels where your desired attribute might be used.

min_samples_split:

This parameter sets the minimum number of samples required to split an internal node. Increasing this value can force the tree to consider attributes with higher information gain earlier.

min_samples_leaf:

This parameter sets the minimum number of samples required to be at a leaf node. Increasing this value can have a similar effect to increasing min samples split.

3. Custom Splitting Criteria:

Implement your own splitting criterion: You can write a custom function to
calculate the splitting criterion, giving more weight to the attribute you want to
prioritize.

However, keep in mind that:

- Decision trees are designed to find the best splits based on the data. Forcing a specific attribute might lead to a suboptimal model.
- The importance of an attribute depends on its relationship with the target variable. If the attribute is not strongly correlated with the target, it might not be used even if you try to force it.

Here's an example of how to use feature engineering to influence the decision tree: