

Contents

Chapter 1

Decision Trees

Let's start with training set:

Definition 1.0.1: Training set

is defined training set a *set of examples*, where: $\langle x^{(i)}, y^{(i)} \rangle$ where:

- i is the instance of the example
- $x^{(i)} \in X$ is the set of *input*
- $y^{(i)} \in Y$ is the set of *output*

the problem of machine learning is to find a function $h : X \rightarrow Y$ that approximates the real function $f : X \rightarrow Y$. We have two types of problems:

- **Classification:** Y is a discrete set of values (e.g. $\{0, 1\}$)
- **Regression:** Y is a continuous set of values (e.g. \mathbb{R})

1.1 Hypothesis space

In machine learning the *hypothesis space* H is defined as the set of all possible functions that can be used to approximate the real function $f : X \rightarrow Y$. Formally:

Definition 1.1.1: Hypothesis space

A hypothesis space H is defined as the set: $H = \{h | h : X \rightarrow Y\}$ where:

- h is a function (hypothesis) that maps input X to output Y
- X the input space (features, domain of data).
- Y the output space (labels, range of data).
- $|H|$ is the size of the hypothesis space (number of possible hypotheses)

this let us to define the model:

Definition 1.1.2: Model

A model is a way to compute a function $h \in H$ from the training set.

Example 1.1.1 (Decision tree)

A good day to play tennis? Our function F is:

$$F : Outlook \times Humidity \times Wind \times Temp \rightarrow PlayTennis?$$



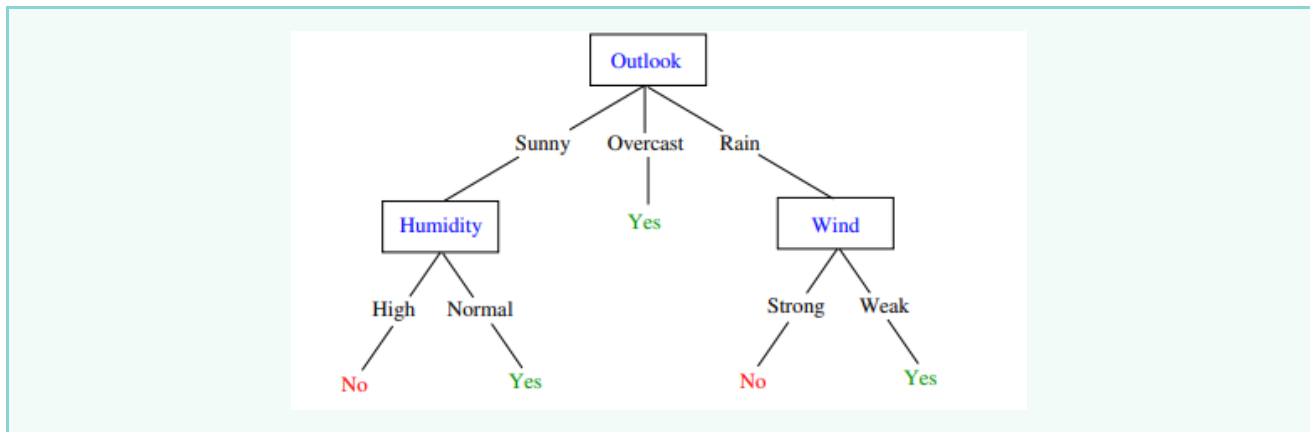
where:

- $Outlook \in \{Sunny, Overcast, Rain\}$
- $Humidity \in \{High, Normal\}$
- $Wind \in \{Weak, Strong\}$
- $PlayTennis? \in \{Yes, No\}$

Every node tests an attribute. Each branch corresponds to one of the possible values for that attribute. Each leaf node assigns a classification (Yes or No), in other words predicts the answer Y .

The problem configuration is the following:

- X is the set of all possible $x \in X$ that corresponds to a vector of attributes $(Outlook, Humidity, Wind, Temp)$
- Target function $f : X \rightarrow Y$ is the function that maps the attributes to the target variable $PlayTennis?$ (booleans)
- Hypothesis space $H = \{h|h : X \rightarrow Y\}$ is the set of all possible decision trees that can be constructed using the attributes in X to predict the target variable Y



1.1.1 Top-down inductive construction

Let $X = X_1 \times X_2 \cdots \times X_n$ where $X_i = \{\text{True}, \text{False}\}$

Can we represent, for instance, $Y = X_2 \wedge X_5$? or $Y = X_2 \wedge X_5 \vee (\neg X_3) \wedge X_4 \wedge X_1$?
and:

- do we have a decision tree for each h in the space hypothesis?
- if the tree exists, is it unique?
- if it is not unique, do we have a preference?

Theorem 1.1.1 Basta - Bonzo

Main loop:

- **Pick the "best" attribute X_i :** At the current node, choose which feature/attribute will best split the training data.
Best means: the attribute that gives the most information gain
- **Create a child node for each possible value of X_i :** for instance if attribute is "weather" with values "sunny", "rainy", "overcast", create three child nodes.
- **check if all examples in the child node are pure:** if all examples belong to the same class (e.g., all "yes" or all "no"), make that node a leaf node with that class label. If not repeat the process recursively for each child node.

1.1.2 Entropy

Definition 1.1.3: Entropy

The entropy $H(S)$ of a set of examples S is defined as:

$$H(S) = - \sum_{i=1}^n P(X = i) \log_2 P(X = i)$$

where:

- $P(X = i)$ is the proportion of examples in S that belong to class i
- n is the number of classes (the number of possible values of X)

In other words, Entropy measures the *degree of uncertainty* of the information. It is maximal when X is uniformly distributed (all classes have the same probability) and minimal (zero) when all examples belong to the same class (pure set)

Information Theory (Shannon 1948)

The entropy is the average amount of information produced by a stochastic source of data. The *information* is associated to the *probability* of each datum (the surprise element):

- An event with probability 1 (certain event) provides no information (no surprise): $I(1) = 0$.
- An event with probability 0 (impossible event) provides infinite information (really surprising): $I(0) = \infty$.
- Given two independent events A and B , the information provided by both events is the sum of the information provided by each event:

$$I(A \cap B) = I(A) + I(B)$$

So is natural defining

$$I(p) = -\log_2(p)$$

Code Theory (Shannon-Fano 1949, Huffman 1952)

The entropy is also related to the average number of bits required to transmit outcomes produced by a stochastic source process x .

Let suppose to have n events with same probability $p_i = \frac{1}{n}$. How many bits do we need to encode these events? The answer is $\log_2(n)$ bits. For instance, if we have 4 events, we need 2 bits to encode them:.

In this case:

$$H(X) = -\sum_{i=1}^n P(X=i) \log_2 P(X=i) = -\sum_{i=1}^n \frac{1}{n} \log_2 \frac{1}{n} = \log_2(n)$$

1.1.3 Information Gain

In a decision tree, the goal is to maximize the information gain during the execution of the algorithm. In other words, the final split should result in the minimum possible impurity. Here are the main formulas:

Theorem 1.1.2 Entropy of X

$$H(X) = -\sum_{i=1}^n P(X=i) \log_2 P(X=i)$$

Theorem 1.1.3 Conditional Entropy of X given a specific $Y = v$

$$H(X | Y = v) = -\sum_{i=1}^n P(X=i | Y=v) \log_2 P(X=i | Y=v)$$

This measures the entropy of X restricted to the subgroup where $Y = v$.

Theorem 1.1.4 Conditional Entropy of X given Y

$$H(X | Y) = \sum_{v=1}^m P(Y=v) H(X | Y=v)$$

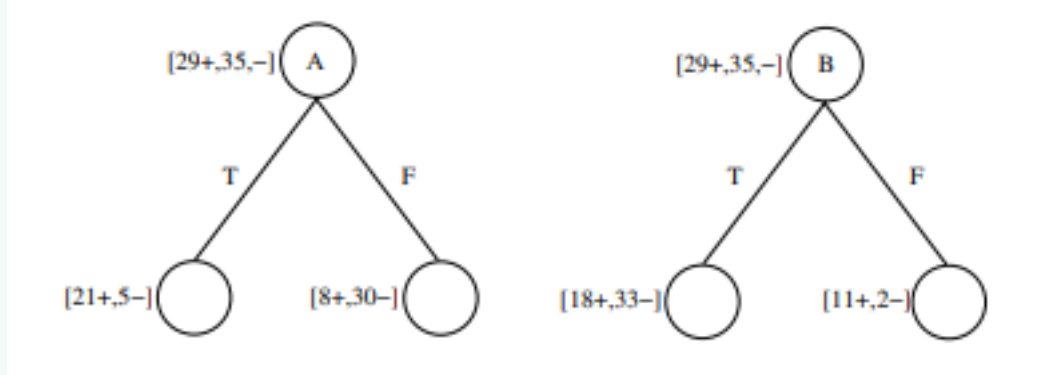
This is the generalization of ??, used to evaluate the utility of an attribute. It measures the average impurity that remains in X after splitting the data using all possible values of Y .

Theorem 1.1.5 Information Gain between X and Y

Here we are! $I(X, Y) = H(X) - H(X|Y) = H(Y) - H(Y|X)$

Example 1.1.2 (Information gain)

Let us measure the entropy reduction of the target variable Y due to some attribute X , that is the information gain $I(Y, X)$ between Y and X



$$H(Y) = -\frac{29}{64} \log_2\left(\frac{29}{64}\right) - \frac{35}{64} \log_2\left(\frac{35}{64}\right) = 0.994$$

$$H(Y | A = T) = -\frac{21}{26} \log_2\left(\frac{21}{26}\right) - \frac{5}{26} \log_2\left(\frac{5}{26}\right) = 0.706$$

$$H(Y | A = F) = -\frac{8}{38} \log_2\left(\frac{8}{38}\right) - \frac{30}{38} \log_2\left(\frac{30}{38}\right) = 0.742$$

$$H(Y | A) = 0.706 \cdot \frac{26}{64} + 0.742 \cdot \frac{38}{64} = 0.726$$

$$I(Y, A) = H(Y) - H(Y | A) = 0.994 - 0.726 = 0.288$$

$$H(Y | B) = 0.872$$

$$I(Y, B) = 0.122$$

Chapter 2

Overfitting

Let us consider the error of the hypothesis h

- on the training set, $error_{train}(h)$
- on the full data set \mathcal{D} , $error_{\mathcal{D}}(h)$

Definition 2.0.1: Overfitting

It's said that h *overfits* the training set if there exists another hypothesis h' such that:

$$error_{train}(h) < error_{train}(h')$$

but

$$error_{\mathcal{D}}(h) > error_{\mathcal{D}}(h')$$

These models (h and h') represent two different situations. The first corresponds to a model that fits the training dataset very closely, including its uncertainty and noise. The second is simpler: it captures only the general trend of the training data and avoids fitting the noise. As a consequence, the error with respect to the true data distribution \mathcal{D} is larger for the first model than for the second. The second one is better! Let's generalise.

But *We do not know* \mathcal{D}

2.1 Avoiding the over fitting

2.1.1 Detecting the Overfitting: validation set

For Detecting the Overfitting it's useful dividing the data available in two disjoint sets:

- **Training set:** set of data that the model *use for learning*. The tree is built by this data
- **validation set:** This set is not shown during the training- It's used as "test" for evaluating the accuracy of the model

2.1.2 Early stopping

This is a proactive strategy. Instead of let the tree grows until his major complexity, it's stopped first the possibility of Overfitting. The growing of a branch is stopped if these two conditions is verified:

- **The improvement is too small:** if a possible division of data produces a gain of information below a certain threshold, it means that it's not useful to continue
- **There are not enough data:** if a node contains a number of examples too much low, any decision taken would be statistically unreliable and probably based on noise. The tree stops to avoid creating rules based on coincidences.

2.1.3 Post - Pruning

This strategy is **reactive**. The decision tree is let grow completely on the training set, which may lead to overfitting, and then the useless or harmful branches are pruned.

Definition 2.1.1: Reduce-Error Post-Pruning

The *reduce-error post-pruning* technique works as follows:

- build the tree completely
- evaluate each branch using a validation set
- prune the branch whose removal improves accuracy the most
- repeat until no further pruning improves the accuracy

Chapter 3

Probabilistic approach

3.1 core idea

we have two main points of views:

- **traditional view:** we wanna to approximate a function $f : X \rightarrow X$
- **Probabilist view:** we wanna compute probabilities: $p : P(Y | X)$

3.1.1 Probs basics

Random variables

A random variables X represents an oyt come about which we're ncertain

Example 3.1.1 (Random variables)

- $X = \text{true}$ if a randomly drawn stdent is male
- $X =$ first name of the student
- $X = \text{true}$ if a randomly drawn stdent have the same birthday

Formal def:

Definition 3.1.1: Probs variables

the set Ω of the possible outcomes is called the sample space. It is said random variable a measurable function over Ω :

- Discrete: $\Omega \rightarrow \{m, f\}$
- Continuos: $\Omega \rightarrow \mathbb{R}$

Definition 3.1.2: Probs def

it is defined $P(X)$ is the fraction of times X is true in repeated runs of the same experiment.

Note:

The definition requires that all samples

Pay attention:

Wrong Concept 3.1: bad examples

Sample space, let Ω be a space made the possible sum:

$$\Omega = \{2, 3, 4, \dots, 12\}$$

Problem: not all sums are equally likely! It should be:

$$\begin{aligned} P(\text{sum} = 2) &= 1/11 \\ P(\text{sum} = 7) &= 1/11 \end{aligned}$$

but in reality:

- Sum = 2: can only happen one way: (1, 1)
- Sum = 7: can happen six ways: (1, 6), (2, 5), (3, 4), (4, 3), (5, 2), (6, 1)

so

$$P(\text{sum} = 2) \neq P(\text{sum} = 7)$$

A correct approach is

Claim 3.1.1 correct approach

Be $\Omega = (1, 1), (1, 2), (1, 3), \dots, (6, 5), (6, 6)$, where $|\Omega| = 36$ outcomes
each pair has equally probability = $\frac{1}{36}$
Now here is a correctly computing:

$$\begin{aligned} P(\text{sum} = 2) &= \frac{|(1,1)|}{36} = \frac{1}{36} \\ P(\text{sum} = 7) &= \frac{|(1,6),(2,5),(3,4),(4,3),(5,2),(6,1)|}{36} = \frac{6}{36} \end{aligned}$$

The Axioms of Probability Theory

These are the fundamental rules that make probability a "reasonable theory of uncertainty":

Axioms of probability theory

$$(1) \text{ Non-negativity: } 0 \leq P(A) \leq 1 \quad \text{for all events } A. \quad (3.1)$$

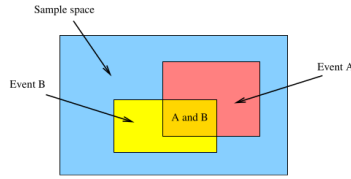
$$(2) \text{ Normalization: } P(\Omega) = 1. \quad (3.2)$$

$$(3) \text{ Countable additivity: } \text{If } A_1, A_2, \dots \text{ are disjoint, then } P\left(\bigcup_{i=1}^{\infty} A_i\right) = \sum_{i=1}^{\infty} P(A_i). \quad (3.3)$$

Then:

Corollary 3.1.1 consequences of the axioms

- Monotonicity: If $A \subseteq B$, then $P(A) \leq P(B)$
- Union rule (for two events): $P(A \cup B) = P(A) + P(B) - P(A \cap B)$
- $P(\text{True}) = 1$
- $P(\text{False}) = 0$



Derived theorems

Corollary 3.1.2 Complement Rule

$$P(\neg A) = 1 - P(A)$$

Dm:

$$P(A \cup \neg A) = P(A) + P(\neg A) - P(A \cap \neg A)$$

But:

$$P(A \cup \neg A) = P(\text{True}) = 1 \quad \text{and} \quad P(A \cap \neg A) = P(\text{False}) = 0$$

Therefore:

$$1 = P(A) + P(\neg A) - 0 \implies P(\neg A) = 1 - P(A) \quad \text{QED}$$



Corollary 3.1.3 Partition Rule

$$P(A) = P(A \cap B) + P(A \cap \neg B)$$

Proof:

$$\begin{aligned} A &= A \cap (B \cup \neg B) && [\text{since } B \cup \neg B \text{ is always True}] \\ &= (A \cap B) \cup (A \cap \neg B) && [\text{distributive law}] \end{aligned}$$

Hence,

$$\begin{aligned} P(A) &= P((A \cap B) \cup (A \cap \neg B)) \\ &= P(A \cap B) + P(A \cap \neg B) - P((A \cap B) \cap (A \cap \neg B)) \\ &= P(A \cap B) + P(A \cap \neg B) - P(\text{False}) \\ &= P(A \cap B) + P(A \cap \neg B) \end{aligned}$$



Multivalued Discrete Random Variables

Definition 3.1.3: k-value Discrete Random Variables

A random variable A is *k-valued discrete* if it takes exactly one value from

$$\{v_1, v_2, \dots, v_k\}.$$

Proposition 3.1.1 Key proprieties

1. **Mutual exclusivity:** For $i \neq j$,

$$P(A = v_i \cap A = v_j) = 0$$

2. **Exhaustiveness:**

$$P(A = v_1 \cup A = v_2 \cup \dots \cup A = v_k) = 1$$

Conditional Probability

Definition 3.1.4: Conditional probs

The Conditional probs of the event A *given* the event B is defined as the quantity

$$P(A | B) = \frac{P(A \cap B)}{P(B)}$$

Corollary 3.1.4 Chain rule

$$P(A \cap B) = P(B)P(A | B) = P(A)P(B | A)$$

Independent Events

Definition 3.1.5: Independent Events

Events A and B are independent when:

$$P(A | B) = P(A)$$

(Meaning: B provides no information about A .)

Corollary 3.1.5 consequences

- $P(A \cap B) = P(A)P(B)$ (from chain rule)
- $P(B|A) = P(B)$ (symmetry)

Bayes' Rule: The Heart of Probabilistic ML (ok chat... really?)

Theorem 3.1.1 Bayes' rule

Now we have Bayes rule

$$P(A | B) = \frac{P(A)P(B|A)}{P(B)}$$

Proof: It's true by the chain rule that: $P(A \cap B) = P(B)P(A | B)$. It's true also the reverse case $P(A \cap B) = P(A)P(B | A)$.

Since both expressions equal $P(A \cap B)$, they must equal each other:

$$P(A)P(B | A) = P(B)P(A | B)$$

that it's equal to

$$P(A | B) = \frac{P(A)P(B | A)}{P(B)}$$



Example 3.1.2 (The trousers problem)

Setup:

- 60% of students are boys, 40% are girls
- girls wear in the same number skirt and trousers
- boys only wear trousers

If we see a student wearing trousers, what is the probability that is a girl?

Solution: The probab a priori that a student is a girl is

$$P(G) = \frac{2}{5}$$

the probability that a student wears trousers is

$$P(T) = \frac{1}{5} + \frac{3}{5} = \frac{4}{5}$$

the probability that a student wear trousers, given that the student is a girl, is

$$P(T | G) = 1/2$$

So

$$P(G | T) = \frac{p(G)p(T | G)}{P(T)} = \frac{2/5 \cdot 1/2}{4/5} = 1/4$$

Q

Machine Learning Form

Machine Learning Form For discrete Y with values $\{y_1, y_2, \dots, y_m\}$ and X with values $\{x_1, x_2, \dots, x_n\}$:

$$P(Y = y_i | X = x_j) = \frac{P(Y = y_i) \cdot P(X = x_j | Y = y_i)}{P(X = x_j)}$$

Expanding the denominator:

$$\begin{aligned} P(X = x_j) &= \sum_i P(X = x_j, Y = y_i) \quad [\text{sum over all } Y \text{ values}] \\ &= \sum_i P(Y = y_i) \cdot P(X = x_j | Y = y_i) \quad [\text{chain rule}] \end{aligned}$$

Complete Bayes' Rule:

$$P(Y = y_i | X = x_j) = \frac{P(Y = y_i) \cdot P(X = x_j | Y = y_i)}{\sum_i P(Y = y_i) \cdot P(X = x_j | Y = y_i)}$$

Terminology:

$$\underbrace{P(Y | X)}_{\text{posterior}} = \frac{\overbrace{P(X | Y)}^{\text{likelihood}} \cdot \overbrace{P(Y)}^{\text{prior}}}{\underbrace{P(X)}_{\text{marginal}}}$$

- **Posterior** $P(Y | X)$: What we want – probability of Y given observed X
- **Likelihood** $P(X | Y)$: How likely is X if Y is true?
- **Prior** $P(Y)$: What we believed before seeing X
- **Marginal** $P(X)$: Overall probability of observing X (normalization constant)

Alternative form:

$$\text{Posterior} = \frac{\text{Likelihood} \cdot \text{Prior}}{\text{Marginal Likelihood}}$$

where:

$$\text{Marginal} = \sum_Y P(X | Y) \cdot P(Y)$$

The term “marginal” means we’ve **marginalized** (integrated/summed) over Y .