

# Mathematical Foundations of Plant Semantics

Simon Castellan<sup>1</sup>, Jos Käfer<sup>2</sup>, Eric Tannier<sup>3</sup>

<sup>1</sup> Inria Rennes

<sup>2</sup> CNRS

<sup>3</sup> Inria Lyon

February 16th, 2023

## Another research

Can we do research that:

- ▶ solve a real need arising from society,
- ▶ empowers users without enslaving them (Illich's *conviviality*),
- ▶ is somewhat aligned with the ecological transition.

## Another research

Can we do research that:

- ▶ solve a real need arising from society,
- ▶ empowers users without enslaving them (Illich's *conviviality*),
- ▶ is somewhat aligned with the ecological transition.

A case study, PI@ntNet: can we improve on the points above?

- ▶ Hard to trust
- ▶ Very knowledgeable but a bad teacher
- ▶ What if it goes away?



# Modern botanist technology [Bonnier'1904]

**1<sup>er</sup> GROUPE :**

|   |   |   |
|---|---|---|
| ✕ Étamines et pétales réunis aux sépales par leur base. [En relevant le calice jusqu'à la base, on enlève en même temps les étamines et les pétales.] | <input type="checkbox"/> Feuilles épaisses et charnues, 6 à 20 pétales, ou même plus.   | — Plus de 5 sépales; carpelles libres..... <i>Crassulacées</i> , p. 110.  |
|   |   | — 5 sépales; carpelles soudés entre eux, sauf au sommet..... <i>FICOIDÉES</i> , p. 112.   |
| ✕ Étamines et pétales réunis aux sépales par leur base. [En relevant le calice jusqu'à la base, on enlève en même temps les étamines et les pétales.] | <input type="checkbox"/> Feuilles non charnues; 4 à 8 pétales rarement plus.  | <input type="checkbox"/> Calice d'un rouge vif, un peu en forme de toupie PU; pétales d'un rouge vif; arbuste à feuilles luisantes, coriaces, simples..... <i>GRANATÉES</i> , p. 103.   |
|   |   | <input type="checkbox"/> Arbuste odorant, à feuilles persistantes, coriaces, tout à fait entières et sans stipules; fleurs blanches..... <i>MYRTACÉES</i> , p. 107.   |
|   |   | <input type="checkbox"/> Arbuste à feuilles dentées opposées et sans stipules, à fleurs blanches très odorantes..... <i>PHILADELPHÉES</i> , p. 110.   |
| ✕ Étamines non réunies aux sépales et réunies entre elles, au moins à la base, ou disposées par groupes.  | * Arbre ou arbuste.   | <input type="checkbox"/> Pédicule soudé avec la bractée TI; 5 sépales libres, 5 pétales; feuilles molles, non persistantes..... <i>TILIACÉES</i> , p. 54.   |
|   |   | <input type="checkbox"/> Pédicule non soudé avec la bractée; sépales soudés; 3 à 8 pétales; feuilles assez coriaces, persistant pendant l'hiver..... <i>HESPÉRIDÉES</i> , p. 62.  |
|   | * Plante n'étant ni un arbre ni un arbuste.   | <input type="checkbox"/> Feuilles opposées, entières; 3 à 5 styles; étamines disposées par groupes [ex.: H. A.]..... <i>HYPERICINÉES</i> , p. 59.   |
|   |   | <input type="checkbox"/> Feuilles alternes. <ul style="list-style-type: none"> <li>= Fleurs en grappe allongée, [ex.: LL.]; pétales très divisés; calice simple..... <i>RÉSÉDACÉES</i>, p. 38.</li> <li>= Fleurs non en grappe allongée, souvent disposées à l'aisselle des feuilles; calice souvent double MS, AO..... <i>MALVACÉES</i>, p. 55.</li> </ul> |
| ✕ Étamines non réunies aux sépales et libres entre elles jusqu'à la base.   | <input type="checkbox"/> Plus de 16 pétales.  | <input type="checkbox"/> Plante flottant sur l'eau ou submergée, à feuilles dont le bord est entier [NL, NA, coupe de fleurs en long]..... <i>NYMPHÉACÉES</i> , p. 11.  |
|   |   | <input type="checkbox"/> Plante ni flottante ni submergée; feuilles profondément découpées. <ul style="list-style-type: none"> <li>✕ 4 sépales; feuilles plus ou moins arrondies, entières CP..... <i>CAPPARIDÉES</i>, p. 31.</li> <li>✕ 2 sépales tombant quand la fleur s'ouvre [ex.: PA, C]..... <i>PAPAVÉRACÉES</i>, p. 12.</li> </ul>                  |
|   |   | <input type="checkbox"/> Plus de 4 pétales. <ul style="list-style-type: none"> <li>• Feuilles entières ou faiblement dentées, souvent opposées..... <i>CISTINÉES</i>, p. 34.</li> <li>• Feuilles très découpées..... <i>BEVONGLACÉES</i>, p. 2.</li> </ul>  |
| <input type="checkbox"/> Moins de 16 pétales.   | <input type="checkbox"/> Pistil à un seul ovaire et pétales chiffonnées ou tordus sur eux-mêmes dans le bouton. <ul style="list-style-type: none"> <li>△ 4 pétales [ex.: P, C]..... <i>Alismacées</i>, p. 290.</li> </ul>                                   |   |
|   | <input type="checkbox"/> Pistil, en général, à plusieurs parties et pétales non chiffonnés ni tordus dans le bouton. <ul style="list-style-type: none"> <li>△ Fleur n'ayant pas à la fois 3 sépales et 3 pétales..... <i>Alismacées</i>, p. 290.</li> </ul> |   |

*Flore complète de la France et de la Suisse, pour trouver facilement les noms de plantes, SANS MOTS TECHNIQUES*

# The rabbit hole

*Can we automate the work of building determination keys and alleviate some of their limits?*

**Output:** offline app or even paper version of the key.

## Challenges:

- ▶ No open source morphological database.
- ▶ No formal description of plants.

# The rabbit hole

*Can we automate the work of building determination keys and alleviate some of their limits?*

**Output:** offline app or even paper version of the key.

## Challenges:

- ▶ No open source morphological database.
- ▶ No formal description of plants.
- ▶ Avoid over-engineering !
- ▶ Participative research to build data.

# The ID3 algorithm [Quinlan'86]

In what order should we ask the questions?



# The ID3 algorithm [Quinlan'86]

In what order should we ask the questions?

**Bayesian algorithm** using **information theory**:

1. Start with an **initial probability distribution**  $d$ .

# The ID3 algorithm [Quinlan'86]

In what order should we ask the questions?

**Bayesian algorithm** using **information theory**:

1. Start with an **initial probability distribution**  $d$ .
2. For every question  $q$ :
  - ▶ compute the average **information** *after*  $q$ .
3. Ask the question with the **largest information**.

# The ID3 algorithm [Quinlan'86]

In what order should we ask the questions?

**Bayesian algorithm** using **information theory**:

1. Start with an **initial probability distribution**  $d$ .
2. For every question  $q$ :
  - ▶ compute the average **information** *after*  $q$ .
3. Ask the question with the **largest information**.
4. **Update**  $d$  with the user answer and go back to (1).

# The ID3 algorithm [Quinlan'86]

In what order should we ask the questions?

**Bayesian algorithm** using **information theory**:

1. Start with an **initial probability distribution**  $d$ .
2. For every question  $q$ :
  - ▶ compute the average **information** *after*  $q$ .
3. Ask the question with the **largest information**.
4. **Update**  $d$  with the user answer and go back to (1).

Greedy and non-optimal, but good enough for now.

## An aside about information theory [Shannon'48]

**Information** is usually defined as the opposite of **entropy**:

$$\begin{aligned} \text{entropy} : \mathcal{D}(X) &\rightarrow \mathbb{R} \\ d &\mapsto - \sum_{x \in X} d(x) \times \log(d(x)) \end{aligned}$$

## An aside about information theory [Shannon'48]

**Information** is usually defined as the opposite of **entropy**:

$$\begin{aligned} \text{entropy} : \mathcal{D}(X) &\rightarrow \mathbb{R} \\ d &\mapsto - \sum_{x \in X} d(x) \times \log(d(x)) \end{aligned}$$

Extremal values:

$$\begin{aligned} \text{entropy}(\text{uniform}(\{1, \dots, n\})) &= \log(n) \\ \text{entropy}(\text{dirac}(x)) &= 0 \end{aligned}$$

# An aside about information theory [Shannon'48]

**Information** is usually defined as the opposite of **entropy**:

$$\begin{aligned} \text{entropy} : \mathcal{D}(X) &\rightarrow \mathbb{R} \\ d &\mapsto - \sum_{x \in X} d(x) \times \log(d(x)) \end{aligned}$$

Extremal values:

$$\begin{aligned} \text{entropy}(\text{uniform}(\{1, \dots, n\})) &= \log(n) \\ \text{entropy}(\text{dirac}(x)) &= 0 \end{aligned}$$

**Maximizing information**  $\leftrightarrow$  **minimizing entropy**

# Bayesian Update

## Input:

- ▶  $d \in \mathcal{D}(\text{Species})$  : prior knowledge of what species it may be.
- ▶  $o \in \text{Obs}$ : user answer (“red flower”)

## Output:

- ▶  $d' \in \mathcal{D}(\text{Species})$  : posterior knowledge



# Bayesian Update

## Input:

- ▶  $d \in \mathcal{D}(\text{Species})$  : prior knowledge of what species it may be.
- ▶  $o \in \text{Obs}$ : user answer (“red flower”)

## Output:

- ▶  $d' \in \mathcal{D}(\text{Species})$  : posterior knowledge

For  $s \in \text{Species}$ : (Bayes' Law)

$$d'(s) \propto d(s) \times \text{score}(s, o)$$

$\uparrow$   $\uparrow$   $\uparrow$

$P(\text{we see } s \mid \text{observing } o)$   $P(\text{we see } s)$   $P(\text{observing } o \mid \text{we see } s)$

## Input to the ID3 algorithm

```
(* Representation of the description of a species *)  
(* e.g. ``white flowers, simple lanceolate leaves'' *)  
type species
```

```
(* Representation of observation (user answers). *)  
type observation
```

```
(* Given an observation, how likely is  
   it to be a particular species ? *)  
val score : species -> observation -> float
```

How to describe species and observation ?

↪ score tangles species and observation.

# General shape of model

We distinguish two things:

- ▶ **Plant description**, which are certain and exhaustive  
“Tree with white flowers, large leaves ...“
- ▶ **Observation**, which are uncertain  
“Tree with whiteish flowers (maybe rose), not sure about leaves (its winter)”

# General shape of model

We distinguish two things:

- ▶ Plant: **Plant description**, which are certain and exhaustive  
“Tree with white flowers, large leaves ...”
- ▶ Obs: **Observation**, which are uncertain  
“Tree with whiteish flowers (maybe rose), not sure about leaves (its winter)”

We let  $\text{Obs} = \mathcal{D}(\text{Plant})$ .

# General shape of model

We distinguish two things:

- ▶ Plant: **Plant description**, which are certain and exhaustive  
“Tree with white flowers, large leaves ...”
- ▶ Obs: **Observation**, which are uncertain  
“Tree with whiteish flowers (maybe rose), not sure about leaves (its winter)”

We let  $\text{Obs} = \mathcal{D}(\text{Plant})$ .

$$P(\text{observing } o \mid \text{we see } s)$$

## General shape of model

We distinguish two things:

- ▶ Plant: **Plant description**, which are certain and exhaustive  
“Tree with white flowers, large leaves ...”
- ▶ Obs: **Observation**, which are uncertain  
“Tree with whiteish flowers (maybe rose), not sure about leaves (its winter)”

We let  $\text{Obs} = \mathcal{D}(\text{Plant})$ .

$$\begin{aligned} & P(\text{observing } o \mid \text{we see } s) \\ = & \sum_{p \in \text{Plant}} o(p) \times P(\text{observed plant is } p \mid \text{we see } s) \end{aligned}$$

# General shape of model

We distinguish two things:

- ▶ Plant: **Plant description**, which are certain and exhaustive  
“Tree with white flowers, large leaves ...”
- ▶ Obs: **Observation**, which are uncertain  
“Tree with whiteish flowers (maybe rose), not sure about leaves (its winter)”

We let  $\text{Obs} = \mathcal{D}(\text{Plant})$ .

$$\begin{aligned} & P(\text{observing } o \mid \text{we see } s) \\ = & \sum_{p \in \text{Plant}} o(p) \times \underbrace{P(\text{observed plant is } p \mid \text{we see } s)}_{s(p)} \end{aligned}$$

Thus we can also let  $\text{Species} = \mathcal{D}(\text{Plant})$  !

# General shape of model

We distinguish two things:

- ▶ Plant: **Plant description**, which are certain and exhaustive  
“Tree with white flowers, large leaves ...”
- ▶ Obs: **Observation**, which are uncertain  
“Tree with whiteish flowers (maybe rose), not sure about leaves (its winter)”

We let  $\text{Obs} = \mathcal{D}(\text{Plant})$ .

$$\begin{aligned} & P(\text{observing } o \mid \text{we see } s) \\ = & \sum_{p \in \text{Plant}} o(p) \times \underbrace{P(\text{observed plant is } p \mid \text{we see } s)}_{s(p)} \end{aligned}$$

Thus we can also let  $\text{Species} = \mathcal{D}(\text{Plant})$  !



## An example

```
type color = Red | Blue | White | Rose
type plant = { flower_color: color }
```

## An example

```
type color = Red | Blue | White | Rose
type plant = { flower_color: color }
```

```
type observation = plant distribution
let whiteish : observation = fun p ->
  match p.flower_color with
  | White -> 0.8
  | Rose -> 0.2
  | _ -> 0.
```

## An example

```
type color = Red | Blue | White | Rose
type plant = { flower_color: color }
```

```
type observation = plant distribution
let whiteish : observation = fun p ->
  match p.flower_color with
  | White -> 0.8
  | Rose -> 0.2
  | _ -> 0.
```

```
type species = plant distribution
let laurier_rose : species = fun p ->
  match p.flower_color with
  | Rose -> 0.8
  | White -> 0.2 (* cultivars *)
  | _ -> 0.0
```

# Where are we ?

We have all the conceptual ingredients to make it work:

- ▶ The decision tree **algorithm**
- ▶ Species and observation with **uncertainty**
- ▶ How to make species and observation **interact** (score)

# Where are we ?

We have all the conceptual ingredients to make it work:

- ▶ The decision tree **algorithm**
- ▶ Species and observation with **uncertainty**
- ▶ How to make species and observation **interact** (score)

However:

- ▶ What is the real type for `plant` ?
- ▶ Where do we get the data?
- ▶ If `plant` is big, isn't `score` intractable ?

# The project in a nutshell

1. How to have expert botanists write the type plant ?
2. How to describe the distribution probabilities for species ?
  - ▶ User-friendliness: we need a lot of data, hence of a lot of people!
  - ▶ Link formal/informal: bibliographical info linked to data

# The Flat Model – Plant

The naive approach:

$$\text{Plant} = \prod_{i \in I} C_i$$

where:

- ▶  $I$  is the set of **characters**
- ▶ Each  $C_i$  is a simple sum  $1 + \dots + 1$

# The Flat Model – Plant

The naive approach:

$$\text{Plant} = \prod_{i \in I} C_i$$

where:

- ▶  $I$  is the set of **characters**
- ▶ Each  $C_i$  is a simple sum  $1 + \dots + 1$

Can be described by a textual format:

```
flower-color = [ red blue white rose ... ];  
flower-petal-number = [ 1 2 3 ... ];  
leaf-structure = [ simple divided ];
```



# The Flat Model – Species

Distributions over Plant are inconvenient.

However, we have a correspondance:

$$\mathcal{D}(S \times T) \cong \mathcal{D}(S) \times \mathcal{D}(T)$$

# The Flat Model – Species

Distributions over Plant are inconvenient.

However, we have a correspondance:

$$\mathcal{D}(S \times T) \cong \mathcal{D}(S) \times \mathcal{D}(T)$$

Thus we use

$$\mathcal{D}(\text{Plant}) \cong \prod_{i \in I} \mathcal{D}(C_i)$$

## The Flat Model – Species

Distributions over Plant are inconvenient.  
However, we have a correspondance:

$$\mathcal{D}(S \times T) \cong \mathcal{D}(S) \times \mathcal{D}(T)$$

Thus we use

$$\mathcal{D}(\text{Plant}) \cong \prod_{i \in I} \mathcal{D}(C_i)$$

```
# laurier-rose.species  
flower-color = [ white: 0.8 rose: 0.2 ];  
# we can omit the rest and use uniform distribution
```

## The Flat Model – Species

Distributions over Plant are inconvenient.  
However, we have a correspondance:

$$\mathcal{D}(S \times T) \rightleftharpoons \mathcal{D}(S) \times \mathcal{D}(T)$$

Thus we use

$$\mathcal{D}(\text{Plant}) \rightleftharpoons \prod_{i \in I} \mathcal{D}(C_i)$$

```
# laurier-rose.species  
flower-color = [ white: 0.8 rose: 0.2 ];  
# we can omit the rest and use uniform distribution
```

With this representation, score can be computed in time  
 $O(C_1 + \dots + C_n)$  instead of  $O(C_1 \times \dots \times C_n)$  !

# Limits of the flat model

1. Plant is too **rigid**: what about plants without flowers?
2. Species cannot represent all probability distributions.  
No **correlation** between trait distributions.
3. Compositional structure on species?  
Merge different descriptions ...

## Going full algebraic type

What if we allow Plant to be an algebraic type?

```
plant := {  
  leaf;  
  flower;  
}.  
leaf := {  
  position: [ base | stem {disposition} ];  
  venation;  
  attachment;  
}.  
flower := {  
  inflorescence;  
  sex: [ unisexual | hermaphrodism ];  
  color: [ red | blue | white | rose ];  
}.
```

↪ Similar to ontologies (e.g. RDFS).

## What about Species?

We need to extend our **abstraction**:

$$\mathcal{D}(S \times T) \rightleftharpoons \mathcal{D}(S) \times \mathcal{D}(T)$$

to:

$$\mathcal{D}(S \oplus T) \rightleftharpoons$$

## What about Species?

We need to extend our **abstraction**:

$$\mathcal{D}(S \times T) \rightleftharpoons \mathcal{D}(S) \times \mathcal{D}(T)$$

to:

$$\mathcal{D}(S \oplus T) \rightleftharpoons [0, 1] \times \mathcal{D}(S) \times \mathcal{D}(T)$$

Example of a distribution:

```
{ leaf = { position = stem };  
  flower = { color = [ 0.8: rose | 0.2: white ] }  
}
```



## What about Species?

We need to extend our **abstraction**:

$$\mathcal{D}(S \times T) \rightleftharpoons \mathcal{D}(S) \times \mathcal{D}(T)$$

to:

$$\mathcal{D}(S \oplus T) \rightleftharpoons [0, 1] \times \mathcal{D}(S) \times \mathcal{D}(T)$$

Example of a distribution:

```
{ leaf = { position = stem };  
  flower = { color = [ 0.8: rose | 0.2: white ] }  
}
```

Abstract type interpretation:

$$\llbracket S \times T \rrbracket = \llbracket S \rrbracket \times \llbracket T \rrbracket \quad \llbracket S + T \rrbracket = [0, 1] \times \llbracket S \rrbracket \times \llbracket T \rrbracket.$$

# Biological species model of Linear Logic

We cannot still express **polymorphism**:

*simple basal leaves, compound stem leaves*

# Biological species model of Linear Logic

We cannot still express **polymorphism**:

*simple basal leaves, compound stem leaves*

↔ Can only say: *simple or compound, basal or compound.*

## Biological species model of Linear Logic

We cannot still express **polymorphism**:

*simple basal leaves, compound stem leaves*

$\rightsquigarrow$  Can only say: *simple or compound, basal or compound*. We add an exponential  $!S$  to the type language:

$$\llbracket !S \rrbracket = \mathcal{D}(S)$$

We can thus write:

```
[ 0.5 { leaf = { simple ; basal } }  
| 0.5 { leaf = { compound; stem } } ]
```

In the plant description, we can add modalities:

```
plant = { leaf!; flower!; stem }
```

# Conclusion

We have a working prototype for the flat model:

- ▶ An editor for entering species distribution
- ▶ Greedy algorithm implementation

# Conclusion

We have a working prototype for the flat model:

- ▶ An editor for entering species distribution
- ▶ Greedy algorithm implementation

Inria exploratory project *Back to the trees*:

- ▶ Extend the model
- ▶ Work with local associations to build a database
- ▶ Implement non-greedy algorithms ?



**WE NEED  
YOU**