

# Inferring the presence of metabolites

March 10, 2023

## 1 Model

### 1.1 Definition

We seek to infer the presence or absence of  $M$  metabolites in  $S$  species. Let us define a vector of vectors  $\vec{v} = \{\vec{v}_1, \dots, \vec{v}_I\}$  where  $\vec{v}_i = \{v_1, \dots, v_{J_i}\}$  and  $J_i$  is the length of vector  $\vec{v}_i$ . Since we are interested in at least species and metabolites, we expect  $I \geq 2$ . Hence for this model we expect  $\vec{v} = \{\vec{s}, \vec{m}, \dots\}$  with  $\vec{s} = \{s_1, \dots, s_S\}$  all species and  $\vec{m} = \{m_1, \dots, m_M\}$  all metabolites. Define a vector  $\vec{w}$  where  $\vec{w} = \vec{v} \setminus \{\vec{s}, \vec{m}\}$ .

Let us further define  $\vec{c}$  a vector containing a single element of  $\vec{v}_i$  for each  $\vec{v}_i$  in  $\vec{v}$ . Similarly we have  $\vec{b} = \vec{c} \setminus \{s, m\}$ . We denote by  $x_{\vec{c}} = x_{sm\vec{b}}$  whether metabolite  $m$  is present ( $x_{sm\vec{b}} = 1$ ) or absent ( $x_{sm\vec{b}} = 0$ ) in species  $s$  in properties  $\vec{b} \in \vec{w}$ .

We denote  $\vec{x}_{s\vec{b}} = (x_{s1\vec{b}}, \dots, x_{sM\vec{b}})$  the vector of molecules present in properties  $\vec{b}$  of a specific species  $s$ . Finally we define  $\vec{x}_s = (\vec{x}_{s1\vec{w}}, \dots, \vec{x}_{sM\vec{w}})$  the vector of presence/absence of all molecules for each element of each property in  $\vec{w}$  for species  $s$ .

We assume that related species share a similar set of metabolites and that metabolites related in their synthesis share a similar distribution across species. Let  $\mathbb{P}(x_{sm} = 1 | \mu_m, \epsilon_{sm\vec{b}})$  be the probability with which metabolite  $m$  is present in species  $s$ . We then assume that

$$\text{logit } \mathbb{P}(x_{sm} = 1 | \mu_m, \epsilon_{sm\vec{b}}) = \mu_m + \epsilon_{sm\vec{b}} \quad (1)$$

where  $\mu_m$  is a metabolite-specific intercept and  $\epsilon_{sm\vec{b}}$  is normally distributed with mean 0 and co-variance :

$$\text{cov}(\epsilon_{\vec{c}}, \epsilon_{\vec{c}'} ) = \sum_{i=1}^I \sum_{d \neq i} \sum_{f=1}^F \alpha_{df} \sigma_{c_i c'_i}^{(f)} \quad (2)$$

with  $I$  the length of vector  $\vec{v}$ ,  $\sigma$  a known measure of covariance and  $\alpha$  a positive scalar.  $\sigma_{c_i c'_i}^{(f)}$  is defined as a known measure of covariance between property  $c_i$  and  $c'_i$  when looking at feature  $f$ .  $f$  being a feature at which the variance is measured. For example, this would allow to discriminate the variance of two species when looked at the "phenotype", or "environment" level or any arbitrary feature one is interested in.

#### 1.1.1 Updating parameters

To update our parameters using Metropolis-Hastings algorithm, we need to propose a suitable ration  $h$  for each parameter to be updated. We know that :

$$P(\vec{\alpha} | \vec{x}_s, \vec{\mu}) \propto P(\vec{x}_s | \vec{\mu}, \vec{\alpha}) P(\vec{\alpha}) \quad (3)$$

and,

$$P(\vec{\mu} | \vec{x}_s, \vec{\alpha}) \propto P(\vec{x}_s | \vec{\mu}, \vec{\alpha}) P(\vec{\mu}) \quad (4)$$

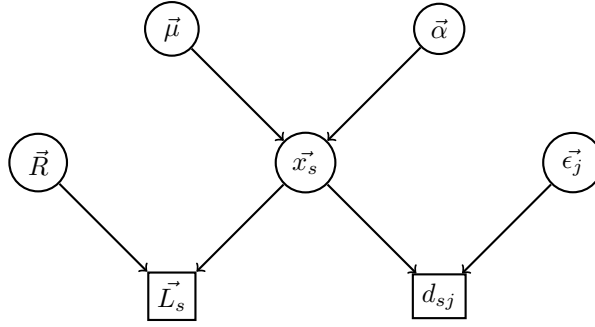


Figure 1: Potential DAG of the model.

It is then possible to update each parameter in turn giving :

$$q_\alpha = \min \left[ 1, \frac{P(\vec{x}_s | \vec{\mu}, \vec{\alpha})}{P(\vec{x}_s | \vec{\mu}, \vec{\alpha})} \cdot \frac{P(\vec{\alpha})}{P(\vec{\alpha})} \right] \quad \text{and,} \quad q_\mu = \min \left[ 1, \frac{P(\vec{x}_s | \vec{\mu}, \vec{\alpha})}{P(\vec{x}_s | \vec{\mu}, \vec{\alpha})} \cdot \frac{P(\vec{\mu})}{P(\vec{\mu})} \right] \quad (5)$$

Furthermore, we have two origins of data, mass spectrometry data and the Lotus database. We denote  $d_{sj}$  the  $j^{\text{th}}$  mass spectrometry run for species  $s$  and  $\vec{L}_s$  all molecules assigned to species  $s$  present in the LOTUS database. Finally we define  $R$ , a function representing the research effort produced for either a species  $s$  or a specific molecule  $m$ . We also define  $\vec{\epsilon}_j$  a vector of error that is specific for each mass-spectrometry run.

A DAG of the model can be seen in Figure 1.

With  $\vec{\mu} = (\mu_1, \dots, \mu_M)$

## 1.2 LOTUS database

Since LOTUS database [1] has no properties in  $\vec{\epsilon}$  other than species and molecules, we denote

$$P(\vec{L}_s | \vec{x}_s, \vec{R}) = P(\vec{L}_s | \vec{\xi}_s, \vec{R}),$$

with  $\vec{\xi}_s = (\xi_{s1}, \dots, \xi_{sM})$  the vector of presence/absence of all molecules  $M$  in species  $s$ . Furthermore,  $\xi_{sm} = \min(1, \sum_t x_{smt})$  the minimum between 1 and the sum of presence or absence of a molecule across all tissues.

More generally, we would have

$$\xi_{sm} = \min \left[ 1, \sum_{i=1}^C \sum_{i \neq \{s, m\}} \sum_{j=1}^J x_{ij} \right] \quad (6)$$

## 1.3 MS data

Let  $\vec{d}_{sj} = (d_{sj1}, \dots, d_{sjM})$  be the presence-absence vector of each metabolite  $m$  obtained with mass-spectrometry run  $j = 1, \dots, J_s$  performed on species  $s$ . Assuming a false-positive and false-negative error rates  $\epsilon_{01}$  and  $\epsilon_{10}$ , respectively, we have

$$\mathbb{P}(\vec{d}_{sj} | \vec{x}, \epsilon_{01}, \epsilon_{10}) = \prod_m \left[ x_{sm} \left( \epsilon_{10}^{1-d_{sjm}} (1 - \epsilon_{10})^{d_{sjm}} \right) + (1 - x_{sm}) \left( \epsilon_{01}^{d_{sjm}} (1 - \epsilon_{01})^{1-d_{sjm}} \right) \right].$$

## 2 Idea and concept

We consider two sets of data informative about  $\vec{x}$ : i) Presence-absence data obtained with mass-spectrometry and ii) presence-only reports of specific metabolites in specific specie. Let  $\vec{d}_{sj} = (d_{sj1}, \dots, d_{sjM})$  be the presence-absence vector of each metabolite  $m$  obtained with mass-spectrometry run  $j = 1, \dots, J_s$  performed on species  $s$ . Assuming a false-positive and false-negative error rates  $\epsilon_{01}$  and  $\epsilon_{10}$ , respectively, we have

$$\mathbb{P}(\mathbf{d}_{sj}|\mathbf{x}, \epsilon_{01}, \epsilon_{10}) = \prod_m \left[ x_{sm} \left( \epsilon_{10}^{1-d_{sjm}} (1 - \epsilon_{10})^{d_{sjm}} \right) + (1 - x_{sm}) \left( \epsilon_{01}^{d_{sjm}} (1 - \epsilon_{01})^{1-d_{sjm}} \right) \right].$$

To model the presence only data, it must be put in relation to the expected research effort. Let  $p_{sm}$  denote the known number of presence-only reports for metabolite  $m$  in species  $s$  and  $n_{sm}$  the unknown number of research projects that aimed at discovering metabolite  $m$  in species  $s$ . Assuming a false-positive and false-negative error rates  $\pi_{01}$  and  $\pi_{10}$ , respectively, we have

$$\mathbb{P}(p_{sm}|n_{sm}, \pi_{01}, \pi_{10}) =$$

### 3 Formulation

We seek to infer the presence or absence of  $M$  metabolites in tissue  $T$  in  $S$  species. We denote by  $x_{smt}$  whether metabolite  $m = 1, \dots, M$  is present ( $x_{smt} = 1$ ) or absent ( $x_{smt} = 0$ ) in tissue  $t = 1, \dots, T$  in species  $s = 1, \dots, S$ . We denote  $\vec{x}_{st} = (x_{st1}, \dots, x_{stM})$  the vector of molecules present in a tissue  $T$  of a specific species  $S$ .

Let us further denote  $\vec{x}_s = (x_{s1}, \dots, x_{sT}) = (x_{s11}, \dots, x_{s1M}, x_{s21}, \dots, x_{sTM})$  the vector of presence/absence of all molecules across all tissues for species  $s$ .

We have two origins of data, mass spectrometry data and the Lotus database. We denote  $d_{sj}$  the  $j^{\text{th}}$  mass spectrometry run for species  $s$  and  $\vec{L}_s$  all molecules assigned to species  $s$  present in the LOTUS database. Finally we define  $R$ , a function representing the research effort produced for either a species  $s$  or a specific molecule  $m$ . We also define  $\vec{\epsilon}_j$  a vector of error that is specific for each mass-spectrometry run.

A DAG of the model can be seen in Figure 1.  $\vec{\mu}$  being the vector of the average presence/absence of each molecule across all species and  $\vec{\alpha}$  represents the vector of all known measures of covariates between species and molecules.

Since the tissue specific origin of a molecule is not known in the LOTUS database [1], we denote

$$P(\vec{L}_S|\vec{x}_s, R) = P(\vec{L}_S|\vec{\xi}_s, \vec{R}),$$

with  $\vec{\xi}_s = (\xi_{s1}, \dots, \xi_{sM})$  the vector of presence/absence of all molecules  $M$  in species  $s$ . Furthermore,  $\xi_{sm} = \min(1, \sum_t x_{smt})$  the minimum between 1 and the sum of presence or absence of a molecule across all tissues.

We can also denote  $\vec{R} = (R_{11}, \dots, R_{1M}, R_{21}, \dots, R_{SM})$  the vector of *research effort* of all molecules across all species.

The probability of having a molecule present in the LOTUS database not only depends on the presence/absence of that molecule in a species but also on the research effort done for a specific molecule or species. We thus have  $R_{sm} = f(n_s, n_m)$  with  $R_{sm} \in [0, 1]$  and where  $n_s$  and  $n_m$  are the number of scientific papers that relate respectively the species or the molecules of interest. We thus have the following matrix :

$$\begin{matrix} L_{sm} = NA & L_{sm} = 1 \\ x_{sm} = 0 & \begin{pmatrix} 1 & 0 \\ 1 - R_{sm} & R_{sm} \end{pmatrix} \\ x_{sm} = 1 & \end{matrix}$$

Tissue of origin is usually known in mass spectrometry analysis. From Figure 1, we have  $P(d_{sj}|\vec{x}_s, \vec{\epsilon}_j) = P(d_{sj}|x_{st(d_{sj})}, \vec{\epsilon}_j)$  where  $t(d_{sj})$  reflects the tissue from which mass spectrometry run  $j$  in species  $s$  was sampled.

We thus have:

$$P(\mathbf{d}, \mathbf{x}, \boldsymbol{\mu}, \boldsymbol{\alpha}) = P(\boldsymbol{\mu})P(\boldsymbol{\alpha}) \prod_{s=1}^S P(\vec{x}_s|\boldsymbol{\mu}, \boldsymbol{\alpha}) \prod_{j=1}^J P(d_{sj}|x_{st(d_{sj})}, \vec{\epsilon}_j) \quad (7)$$

Similarly, for LOTUS database we have :

$$P(\mathbf{L}, \mathbf{x}, \boldsymbol{\mu}, \boldsymbol{\alpha}) = P(\boldsymbol{\mu})P(\boldsymbol{\alpha})P(\mathbf{R}) \prod_{s=1}^S P(\vec{x}_s|\boldsymbol{\mu}, \boldsymbol{\alpha})P(L_s|\vec{x}_s, \mathbf{R}) \quad (8)$$

### 3.1 Decreasing complexity

If one is interested in the presence/absence of a molecule not in the species but on the Genus/Order level, one can easily remodel the previous model such as  $\vec{x}_g = (x_{g1}, \dots, x_{gM})$  and where  $x_{gm} = \min(1, \sum_{s \in g} \sum_t^T x_{smt})$ .

## References

- [1] Adriano Rutz, Maria Sorokina, Jakub Galgonek, Daniel Mietchen, Egon Willighagen, Arnaud Gaudry, James G Graham, Ralf Stephan, Roderic Page, Jiří Vondrášek, Christoph Steinbeck, Guido F Pauli, Jean-Luc Wolfender, Jonathan Bisson, and Pierre-Marie Allard. The LOTUS initiative for open knowledge management in natural products research. *eLife*, 11:e70780, May 2022. doi:10.7554/eLife.70780.