

An integrative model for the protein inference problem

Emile Benoist¹, Guillaume Fertin¹, Géraldine Jean¹,
Dominique Tessier²

1.Nantes Université, LS2N (UMR 6004), Nantes, France

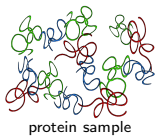
2.Inrae, BIA (UR 1268), Nantes, France

{emile.benoist,geraldine.jean,guillaume.fertin}@univ-nantes.fr,

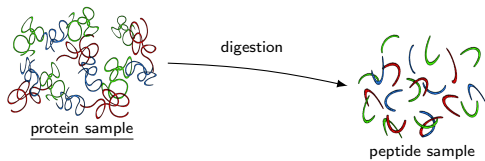
dominique.tessier@inrae.fr

Thursday, November 17

The protein inference problem

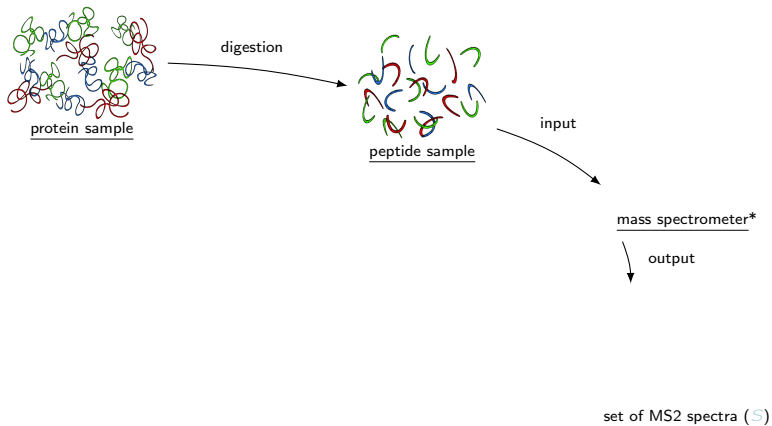


The protein inference problem

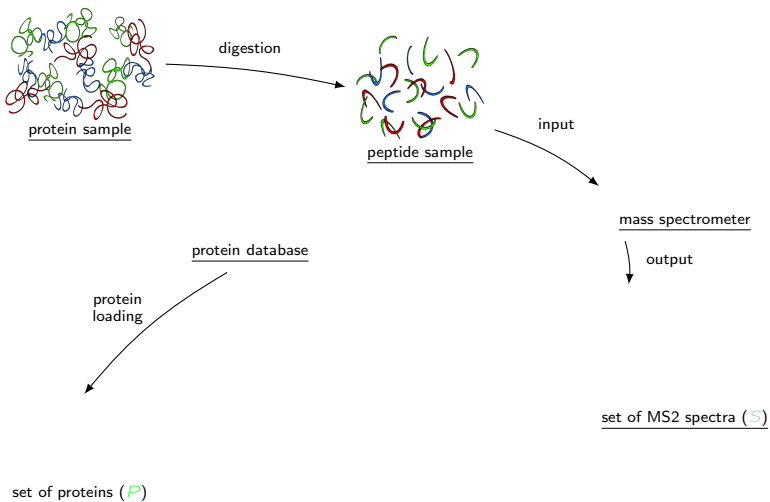


The protein inference problem

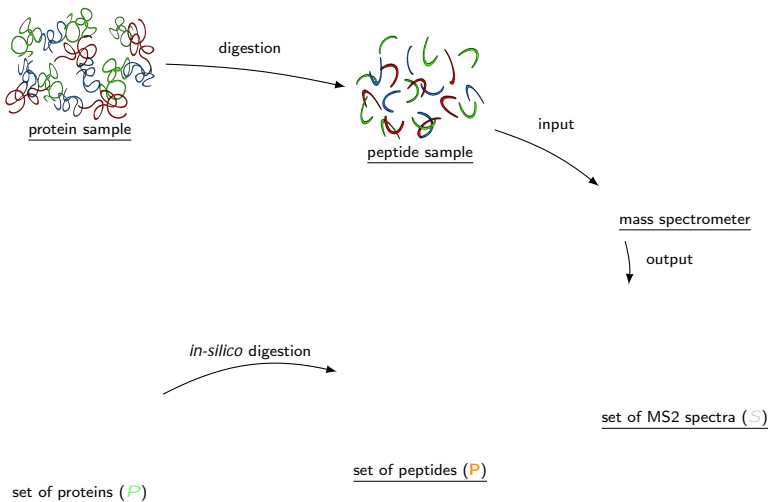
*https://fr.wikipedia.org/wiki/Spectrometrie_de_masse



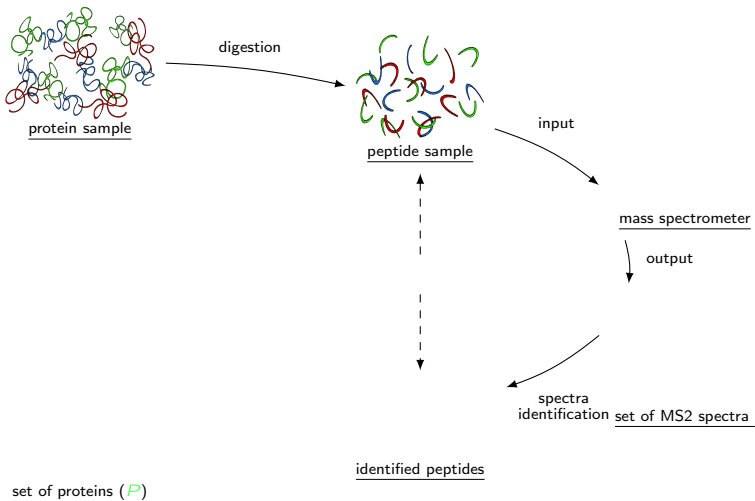
The protein inference problem



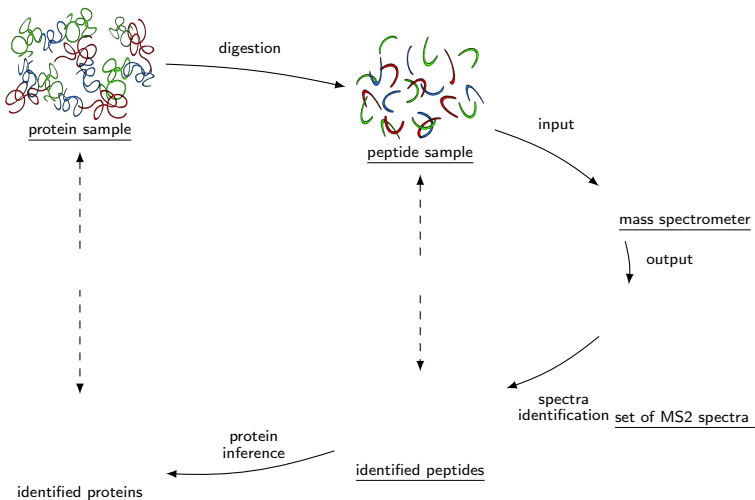
The protein inference problem



The protein inference problem



The protein inference problem



Existing models

Parsimonious models (Occam's razor principle):

Infers the minimum number of proteins that can "explain" all identified peptides [IDPicker, DBParser, MassSieve, LDFA]

Optimistic models:

Infers all proteins meeting a simple criterion [DTASelect]

Statistical models:

Compute for each candidate protein, the probability that the protein is present in the sample [ProteinProphet, MSBayesPro, ProteinLP]

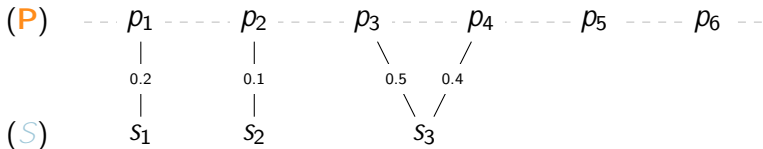
Improvement idea

Each protein inference model relies on a set of identified peptides, themselves identified from all the MS2 spectra independently of the proteins of the database

Improvement idea

Each protein inference model relies on a set of identified peptides, themselves identified from all the MS2 spectra independently of the proteins of the database

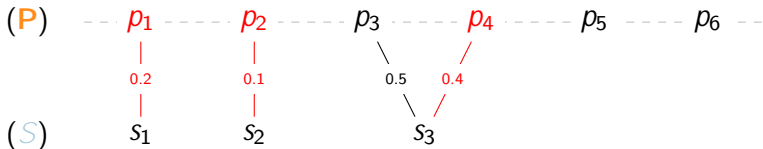
The origin of the peptides (proteins from the database) is an interesting source of information :



Improvement idea

Each protein inference model relies on a set of identified peptides, themselves identified from all the MS2 spectra independently of the proteins of the database

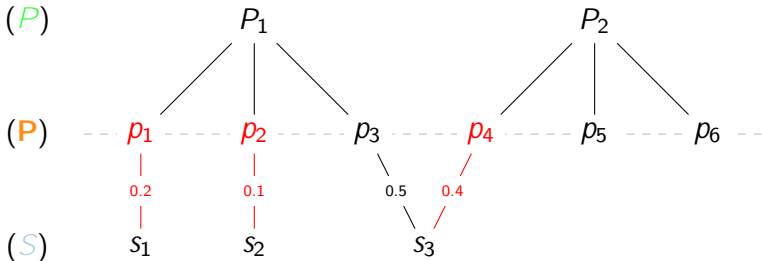
The origin of the peptides (proteins from the database) is an interesting source of information :



Improvement idea

Each protein inference model relies on a set of identified peptides, themselves identified from all the MS2 spectra independently of the proteins of the database

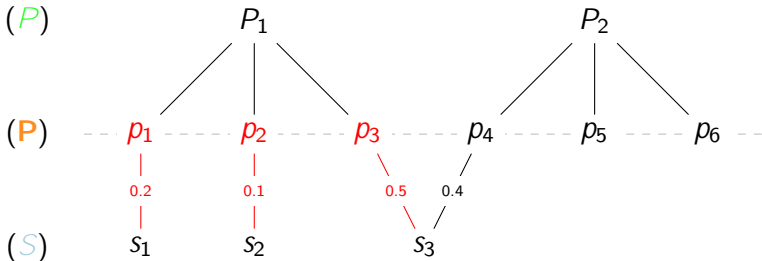
The origin of the peptides (proteins from the database) is an interesting source of information :



Improvement idea

Each protein inference model relies on a set of identified peptides, themselves identified from all the MS2 spectra independently of the proteins of the database

The origin of the peptides (proteins from the database) is an interesting source of information :



The Global Spectra Interpretation (GSI)

Proteins (\mathcal{P}):

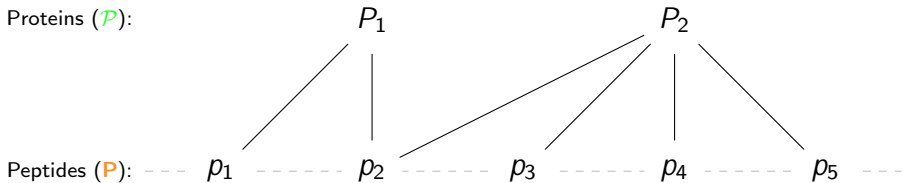
P_1

P_2



Proteins loaded from a database

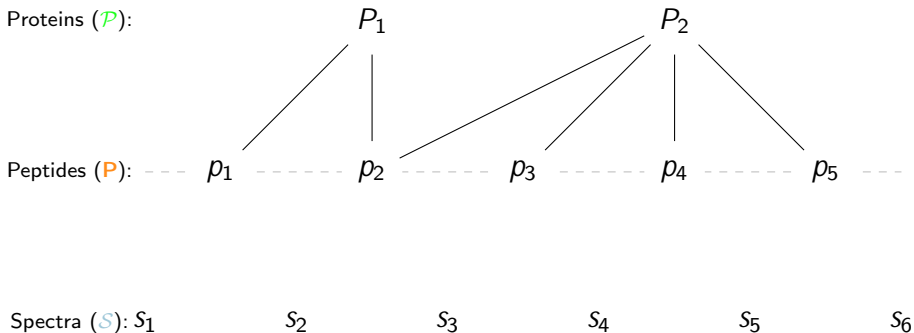
The Global Spectra Interpretation (GSI)



In-silico digestion (trypsin) of each protein

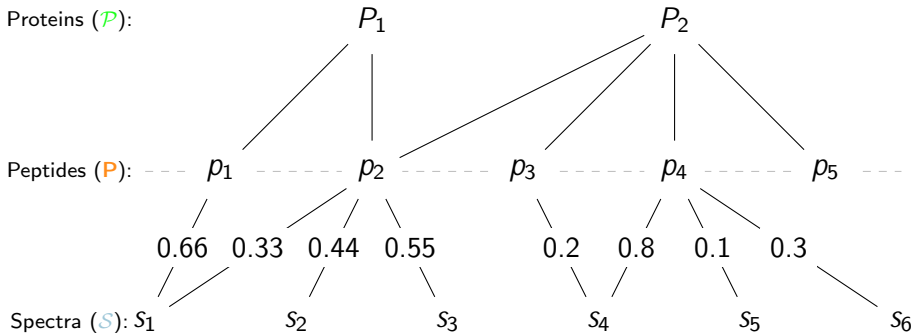
The Global Spectra Interpretation (GSI)

Proteins (\mathcal{P}):



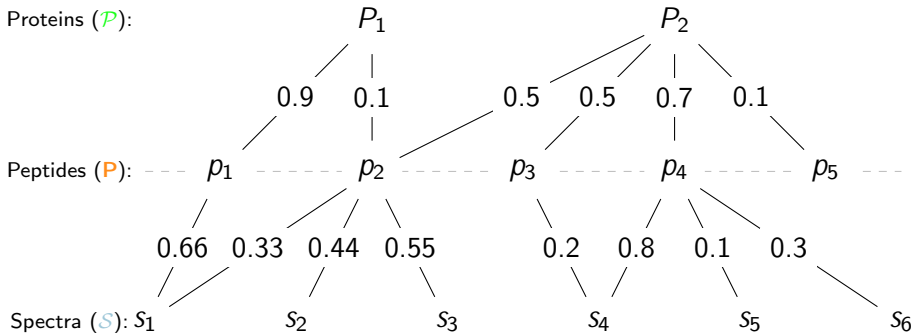
The set of MS2 spectra obtained by mass spectrometry

The Global Spectra Interpretation (GSI)



A set of scores (PSM) provide by software comparing MS2 spectra and peptides

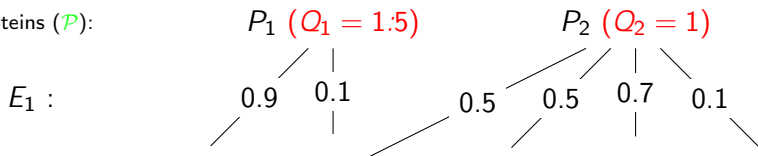
The Global Spectra Interpretation (GSI)



$P(P_i; p_j)$: the probability that a spectrum corresponding to p_j exists in \mathcal{S} given that P_i has an abundance of 1 in the sample

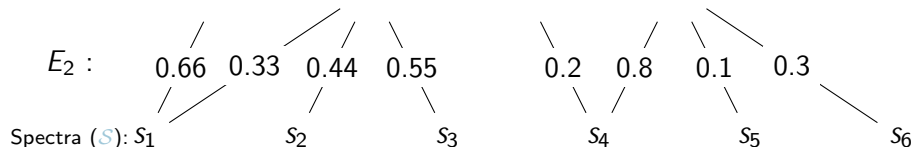
Solution example

Proteins (\mathcal{P}):



Peptides (\mathcal{P}):

----- p_1 ----- p_2 ----- p_3 ----- p_4 ----- p_5 -----



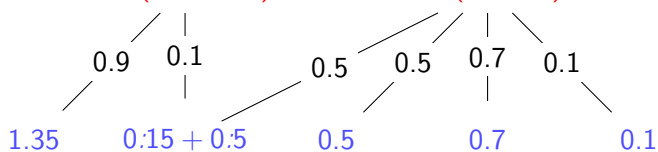
A quantity (or abundance) for each protein in \mathcal{P} . It represents the assumed abundance of the protein in the sample

Solution example

Proteins (\mathcal{P}):

P_1 ($Q_1 = 1.5$) P_2 ($Q_2 = 1$)

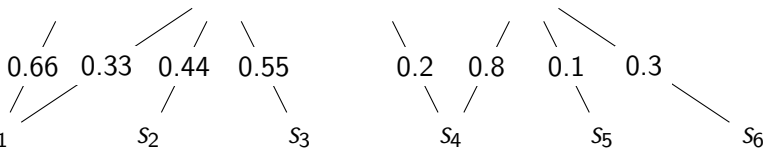
E_1 :



Peptides (\mathcal{P}):

----- p_1 ----- p_2 ----- p_3 ----- p_4 ----- p_5 -----

E_2 :



Spectra (\mathcal{S}):

$$\text{For each peptide } p_j \text{ in } \mathcal{P} : \bar{q}_j = \sum_{(P_i; p_j)} P(P_i; p_j) Q_i$$

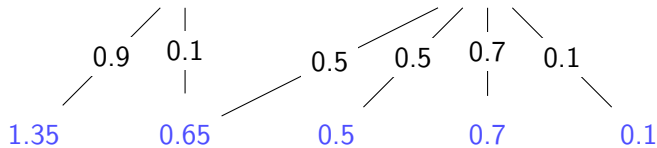
Solution example

Proteins (\mathcal{P}):

P_1 ($Q_1 = 1.5$)

P_2 ($Q_2 = 1$)

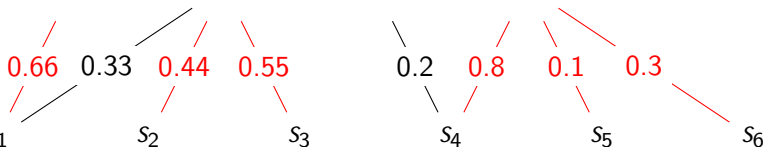
E_1 :



Peptides (\mathcal{P}):

ρ_1 ρ_2 ρ_3 ρ_4 ρ_5

E_2 :



Spectra (\mathcal{S}): S_1

S_2

S_3

S_4

S_5

S_6

One PSM for each spectrum in \mathcal{S}

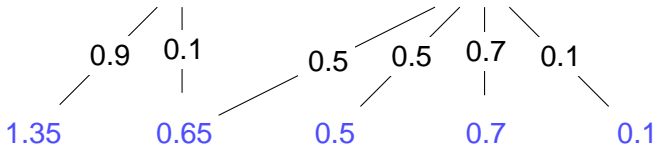
Solution example

Proteins (P):

P_1 ($Q_1 = 1:5$)

P_2 ($Q_2 = 1$)

E_1 :



Peptides (P):

p_1 p_2 p_3 p_4 p_5

1 2 0 3 0

E_2 :



Spectra (S): s_1

s_2

s_3

s_4

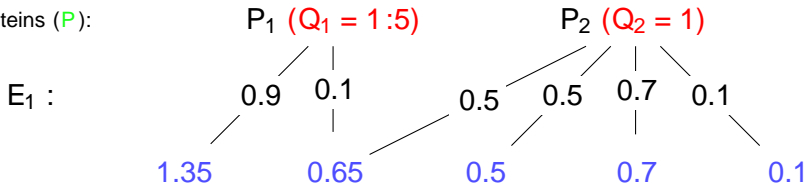
s_5

s_6

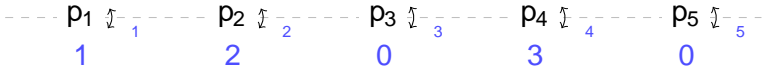
For each peptide p_j in P : q_j = number of selected (red) edges adjacent to p_j

Solution example

Proteins (P):



Peptides (P):



E_2 :



Spectra (S):

s_1 s_2 s_3 s_4 s_5 s_6

$$j = \overline{j} \quad \underline{q_j} \quad (\quad \underline{1} = j1:35 \quad 1j = 0:35)$$

Objective function

Minimize the difference () between the number of times a peptide is identified and the number of times it is expected to be identified given the quantities of the proteins

Objective function

Minimize the difference () between the number of times a peptide is identified and the number of times it is expected to be identified given the quantities of the proteins

Minimize the score on the chosen edges (red edges)

Objective function

Minimize the difference () between the number of times a peptide is identified and the number of times it is expected to be identified given the quantities of the proteins

Minimize the score on the chosen edges (red edges)

$$\min z = \sum_{i=1}^m x_i + \sum_{k=1}^n s_k$$

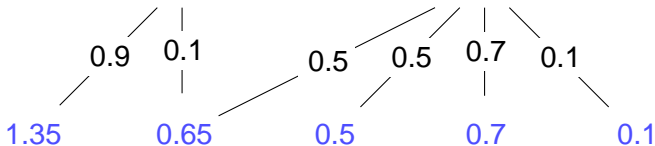
m : number of peptide nodes, n : number of spectrum nodes, s_k : the weighting factors

Objective function

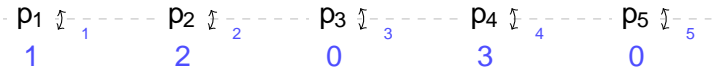
Proteins (P):

$P_1 (Q_1 = 1:5)$ $P_2 (Q_2 = 1)$

E_1 :



Peptides (P):



E_2 :



Spectra (S): S_1 S_2 S_3 S_4 S_5 S_6

$$\begin{aligned}
 z &= 1 \cdot 0.35 + 1.35 + 0.5 + 2 \cdot 3 + 0.1 + 2 \cdot 0.66 + 0.44 + 0.55 + 0.8 + 0.1 + 0.3 \\
 &= 4:6 \quad 1 + 2:85 \quad 2
 \end{aligned}$$

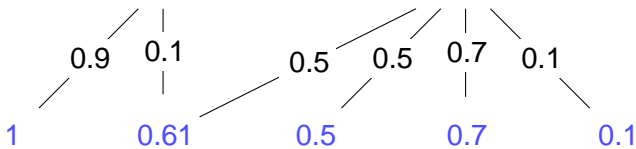
Objective function

Proteins (P):

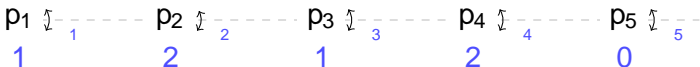
$P_1 (Q_1 = 1:11)$

$P_2 (Q_2 = 1)$

E_1 :



Peptides (P):



E_2 :



Spectra (S):



$$z = \begin{matrix} 1 & 0+1:39+0:5+1:3+0:1 & + & 2 & 0:66+0:44+0:55+0:2+0:1+0:3 \\ = & 3:29 & & & 1+2:25 & 2 \end{matrix}$$

In what way the GSI is an integrative model?

One of our goals is to evolve the model by taking into account other sources of information in order to improve its performance

In what way the GSI is an integrative model?

One of our goals is to evolve the model by taking into account other sources of information in order to improve its performance

The probabilities on the edges between proteins and peptides can carry a lot of information :

- the probability of the peptide to be ionized

- the probability of a missed cleavage around the peptide

- the probability that the peptide will be affected by

- modifications

- ...

In what way the GSI is an integrative model?

One of our goals is to evolve the model by taking into account other sources of information in order to improve its performance

The probabilities on the edges between proteins and peptides can carry a lot of information :

- the probability of the peptide to be ionized

- the probability of a missed cleavage around the peptide

- the probability that the peptide will be affected by

- modifications

- ...

For now, this probability is fixed to 1 (each peptide has the same probability to be identified).

GSI complexity

The Global Spectra Interpretation is an NP-hard problem (reduction from the RXC3 : a restricted version of the Exact-Set-Cover problem) :
under the assumption that $P \neq NP$, there is no polynomial algorithm that can solve the GSI problem

GSI complexity

The Global Spectra Interpretation is an NP-hard problem (reduction from the RXC3 : a restricted version of the Exact-Set-Cover problem) :
under the assumption that $P \neq NP$, there is no polynomial algorithm that can solve the GSI problem

The GSI is FPT with respect to the number of ambiguous edges (an edge adjacent to an ambiguous spectrum) :
there exist an algorithm that solve the GSI in a time that is exponential only on the number of ambiguous edges

GSI complexity

The Global Spectra Interpretation is an NP-hard problem (reduction from the RXC3 : a restricted version of the Exact-Set-Cover problem) :
under the assumption that $P \neq NP$, there is no polynomial algorithm that can solve the GSI problem

The GSI is FPT with respect to the number of ambiguous edges (an edge adjacent to an ambiguous spectrum) :
there exist an algorithm that solve the GSI in a time that is exponential only on the number of ambiguous edges

Moving from theory to practice!

Using the mixed-integer linear programming to solve the GSI problem

Mixed-integer linear programming problem

$$\min z = \sum_{j \in P} p_j x_j + \sum_{(p_j; s_k) \in E_2} s_k x_{j;k}$$

subject to :

$$\text{C.1} \quad \sum_{j \in P} \bar{q}_j x_j \leq 8p_j \quad P$$

$$\text{C.2} \quad \sum_{j \in P} q_j x_j \leq 8p_j \quad P$$

$$\text{C.3} \quad x_{j;k} = 1 \quad (p_j; s_k) \in E_2 \quad S$$

$$\bar{q}_j = \sum_{(P_i; p_j) \in E_1} Q_i \quad P \quad 8p_j \quad P$$

$$q_j = \sum_{(p_j; s_k) \in E_2} x_{j;k} \quad P \quad 8p_j \quad P$$

Q_i : the quantity assigned to the protein P_i

\bar{q}_j : the difference between \bar{q}_j and q_j

$x_{j;k}$: 1 if the edge $(p_j; s_k)$ is selected, 0 otherwise

Explanation of constraints

For each $p_j \in P$, x_j must be equal to q_j or \bar{q}_j

Absolute value is not a linear function but can be expressed as follows :

$$\begin{array}{l}
 \text{C.1} \quad x_j - \bar{q}_j \leq q_j - 8p_j \quad p_j \in P \\
 \text{C.2} \quad x_j - q_j \leq \bar{q}_j - 8p_j \quad p_j \in P
 \end{array}
 \quad \vee \quad
 \begin{array}{l}
 x_j - q_j \leq \bar{q}_j - 8p_j \quad p_j \in P \\
 x_j - \bar{q}_j \leq q_j - 8p_j \quad p_j \in P
 \end{array}$$

Explanation of constraints

For each $p_j \in P$, x_j must be equal to $q_j - \bar{q}_j$

Absolute value is not a linear function but can be expressed as follows :

$$\begin{aligned}
 \text{C.1} \quad & x_j - \bar{q}_j \leq q_j - 8p_j \quad p_j \in P \\
 \text{C.2} \quad & x_j - q_j \leq \bar{q}_j - 8p_j \quad p_j \in P
 \end{aligned}
 \quad \Leftrightarrow \quad
 \begin{aligned}
 & x_j - q_j \leq \bar{q}_j - 8p_j \quad p_j \in P \\
 & + \\
 \text{min } z = & \sum_{p_j \in P} x_j + \dots = \sum_{p_j \in P} (x_j - q_j - \bar{q}_j + 8p_j)
 \end{aligned}$$

Explanation of constraints

For each $p_j \in P$, j must be equal to q_j or \bar{q}_j

Absolute value is not a linear function but can be expressed as follows :

$$\begin{aligned}
 \text{C.1} \quad & j - \bar{q}_j - \frac{q_j}{8p_j} \in P \\
 \text{C.2} \quad & j - \frac{q_j}{8p_j} - \bar{q}_j \in P \\
 \min z = & \sum_{p_j \in P} \dots \Rightarrow j = \frac{q_j}{8p_j} - \bar{q}_j \in P
 \end{aligned}$$

Each spectrum must be identified to exactly one peptide :

$$\text{C.3} \quad \sum_{(p_j; s_k) \in E_2} x_{p_j; s_k} = 1 \quad \forall s_k \in S$$

Some running time

Fixed parameters :

the number of proteins : 10375

the number of peptides : 370759

Variable parameters :

the number of spectra

the number of edges per spectrum

the scores on the edges

the values of α_1 and α_2

Some running time

Fixed parameters :

the number of proteins : 10375

the number of peptides : 370759

Variable parameters :

the number of spectra

the number of edges per spectrum

the scores on the edges

the values of α_1 and α_2

A first series of tests was carried out on an instance containing only one edge per spectrum. Even with a high number of spectra, resolution times do not exceed a few seconds (10 seconds).

Some running time

2 edges per spectrum :

3 edges per spectrum :

The higher the ambiguity (scores close to each other), the higher the execution times

number of proteins : 10735 / number of peptides : 370759 / $\tau_1 = 0:5$ / $\tau_2 = 0:5$

Some running time

Scores = 0.5 / 0.5 :

Scores = 0.2 / 0.8 :

The higher t_1 compared to t_2 , the higher the execution times

number of proteins : 10735 / number of peptides : 370759 / number of edges per spectrum : 2

What about performance

Two questions remain open :

How well the model is able to correct errors of a spectrum identification tool ?

What about performance

Two questions remain open :

How well the model is able to correct errors of a spectrum identification tool ?

Is the model capable of correctly retrieving proteins from a sample?

What about performance

Two questions remain open :

How well the model is able to correct errors of a spectrum identification tool ?

Is the model capable of correctly retrieving proteins from a sample?

The next few months will consist of answering these 2 questions by performing tests on real data sets.

Conclusion

Contribution :

We propose a new model for the protein inference problem that can keep several candidate peptides per spectrum

We demonstrate the scaling of our model with an exact resolution

Work to do :

Perform the tests on the performance of the model

Integration of information through probabilities on edges between proteins and peptides

Evolve the model by the integration of new source of information

Thank you for your attention

