Constraint Programming and Biology: Phylogenetic trees

Agostino Dovier

Dept. Math and Computer Science, Univ. of Udine, Italy

ACP Summer School in Constraint Programming Wrocław, September 2012

Agostino Dovier (DIMI, UDINE Univ.)

CP and Biology

A
A
B
A
B
A
B
A
B
A
B
A
B
A
B
A
B
A
B
A
B
A
B
A
B
A
B
A
B
A
B
A
B
A
B
A
B
A
B
A
B
A
B
A
B
A
B
A
B
A
B
A
B
A
B
A
B
A
B
A
B
A
B
A
B
A
B
A
B
A
B
A
B
A
B
A
B
A
B
A
B
A
B
A
B
A
B
A
B
A
B
A
B
A
B
A
B
A
B
A
B
A
B
A
B
A
B
A
B
A
B
A
B
A
B
A
B
A
B
A
B
A
B
A
B
A
B
A
B
A
B
A
B
A
B
A
B
A
B
A
B
A
B
A
B
A
B
A
B
A
B
A
B
A
B
A
B
A
B
A
B
A
B
A
B
A
B
A
B
A
B
A
B
A
B
A
B
A
B
A
B
A
B
A
B
A
B
A
B
A
B
A
B
A
B
A
B
A
B
A
B
A
B
A
B
A
B
A
B
A
B
A
B
A
B
A
B
A
B
A
B
A
B
A
B
A
B

Phylogenetic trees Basics

- The starting point is a set *L* of elementary (taxonomic) units, known as taxa (e.g.,
 - $L = \{English, German, French, Spanish, Italian\}$ or
 - $L = \{ dog, cat, horse, chicken \} \}$
- A set *I* of characters is assigned to each element of *L* (e.g., characters "hand" and "father", or characters "number of legs", "length of the tail", etc.)
- Characters can have values (e.g. {1 (hand), 2 (mano/main)} for "hand" and {1 (father/padre), 2 (vater/pere)} for "father") Each element in L is assigned a value.
- Let us focus on Boolean characters
- Reconstruction of phylogenies is the first step of reconstructing the evolutionary history of the set *L* of taxa.

・ロト ・同ト ・ヨト ・ヨト

Example (Grapewine Genome Characterization — Nature 449, Sept. 2007)



Figure 3 | Positions of the polyploidization events in the evolution of plants with a sequenced genome. Each star indicates a WGD (tetraploidization) event on that branch. The question mark indicates that ancient events are visible in the rice genome that would require other monocotyledon genome sequences to be resolved. The formation of the palaeo-hexaploid ancestral genome occurred after divergence from monocotyledons and before the radiation of the Eurosids.

Agostino Dovier (DIMI, UDINE Univ.)

CP and Biology

Wrocław, September 2012 3 / 21

Example (Indo-European languages phylogeny-Erdem et al. 2003)



Agostino Dovier (DIMI, UDINE Univ.)

Example (Indo-European languages phylogeny–Erdem et al. 2003)

Hittite (HI), Luvian (LU), Lycian (LY), Tocharian A (TA), Tocharian B (TB), Vedic (VE), Avestan (AV), Old Persian (PE), Classical Armenian (AR), Ancient Greek (GK), Latin (LA), Oscan (OS), Umbrian (UM), Gothic (GO), Old Norse (ON), Old English (OE), Old High German (OG), Old Irish (OI), Welsh (WE), Old Church Slavonic (OC), Old Prussian (PR), Lithuanian (LI), Latvian (LT), and Albanian (AL).

< ロ > < 同 > < 回 > < 回 > < 回 > <

A phylogeny

for a set *L* of taxa is a

- finite binary tree (V, E) (rooted by some v ∈ V) with leaves L ⊆ V (taxa=leaves, with a slight abuse of notation)
- along with two finite sets *I* and *S* and a function $f : L \times I \longrightarrow S$.
- Intuitively, L are the taxonomic units, V \ L describes their ancestral units and E genetic relationships between them.
- *I* is the set of characters, and *S* is the set of their values (also knows are states of these characters)
- *f* labels every leaf $v \in L$ by assigning a state to each character $i \in I$

< ロ > < 同 > < 回 > < 回 > < 回 > <

Example (from Erdem 2011)



A phylogeny (V, E, L, I, S, f) where $L = \{\text{English}, \text{German}, \text{French}, \text{Spanish}, \text{Italian}\}$ (taxa) $I = \{\text{Hand}, \text{Father}\}$ (characters), $S = \{1, 2\}$ (states).

Agostino Dovier (DIMI, UDINE Univ.)

CP and Biology

Wrocław, September 2012 6 / 21

- A character *i* ∈ *I* is compatible with a phylogeny (*V*, *E*, *L*, *I*, *S*, *f*) if there is a function *g* : *V* × {*i*} → *S* s.t.
 - ✓ $(\forall v \in L)(g(v, i) = f(v, i))$ (*g* extends *f*) and
 - ✓ Let us denote $V_{i,s} = \{x \in V : g(x,i) = s\}$.

A character *i* ∈ *I* is compatible with a phylogeny (*V*, *E*, *L*, *I*, *S*, *f*) if there is a function *g* : *V* × {*i*} → *S* s.t.

✓
$$(\forall v \in L)(g(v, i) = f(v, i))$$
 (g extends f) and

✓ Let us denote
$$V_{i,s} = \{x \in V : g(x,i) = s\}$$
.
Then $(\forall s \in S)(V_{i,s} \neq \emptyset \Rightarrow$

the subgraph $(V_{i,s}, E \cap V_{i,s}^2)$ of (V, E) is a rooted tree.)

Agostino Dovier (DIMI, UDINE Univ.)

< ロ > < 同 > < 回 > < 回 > < 回 > <

Example (from Erdem 2011)



Character Hand is compatible with the above tree

Agostino Dovier (DIMI, UDINE Univ.)

CP and Biology

- ₹ € → Wrocław, September 2012 8/21

∃ >

Example (from Erdem 2011)



Character Hand is compatible with the above tree

Agostino Dovier (DIMI, UDINE Univ.)

CP and Biology

A 30 b Wrocław, September 2012 8/21

-

э

Example (from Erdem 2011)



Character Hand is compatible with the above tree

Agostino Dovier (DIMI, UDINE Univ.)

CP and Biology

- ₹ € → Wrocław, September 2012 8/21

-

Example (from Erdem 2011)



Character Father is incompatible with the above tree

Agostino Dovier (DIMI, UDINE Univ.)

CP and Biology

Wrocław, September 2012 9 / 21

Example (from Erdem 2011)



Character Father is incompatible with the above tree

Agostino Dovier (DIMI, UDINE Univ.)

CP and Biology

Wrocław, September 2012 9 / 21

Example (from Erdem 2011)



Character Father is incompatible with the above tree

Agostino Dovier (DIMI, UDINE Univ.)

CP and Biology

Wrocław, September 2012 9 / 21

- A character *i* ∈ *I* is compatible with a phylogeny (*V*, *E*, *L*, *I*, *S*, *f*) if there is a function *g* : *V* × {*i*} → *S* s.t.
 - ✓ $(\forall v \in L)(g(v, i) = f(v, i))$ (*g* extends *f*) and
 - ✓ Let us denote $V_{i,s} = \{x \in V : g(x,i) = s\}$. Then $(\forall s \in S)(V_{i,s} \neq \emptyset \Rightarrow$

the subgraph $(V_{i,s}, E \cap V_{i,s}^2)$ of (V, E) is a rooted tree.)

- Otherwise it is incompatible
- The above (sub-tree) requirement implicitly states that when a character changes (in the evolution) it never go back to the previous value (Camin-Sokal). Moreover, that the "change" occurs in a unique place (Dollo).
- cladistic characters are those that have a temporal order (ancestral/derived), otherwise they are qualitative

-

・ロッ ・雪 ・ ・ ヨ ・ ・

k-INCOMPATIBILITY PROBLEM

Given sets *L* (taxa/leaves), *I* (characters), and *S* (states), a function $f : L \times I \longrightarrow S$, and $k \in \mathbb{N}$, decide the existence of a phylogeny (V, E, L, I, S, f) with at most *k* incompatible characters.

< ロ > < 同 > < 回 > < 回 >

k-INCOMPATIBILITY PROBLEM

Given sets *L* (taxa/leaves), *I* (characters), and *S* (states), a function $f : L \times I \longrightarrow S$, and $k \in \mathbb{N}$, decide the existence of a phylogeny (V, E, L, I, S, f) with at most *k* incompatible characters.

This problem is NP-complete (Day, Sankoff 1986). We'll see a sketch of the proof later.

< ロ > < 同 > < 回 > < 回 >

k-INCOMPATIBILITY PROBLEM

Given sets *L* (taxa/leaves), *I* (characters), and *S* (states), a function $f : L \times I \longrightarrow S$, and $k \in \mathbb{N}$, decide the existence of a phylogeny (V, E, L, I, S, f) with at most *k* incompatible characters.

This problem is NP-complete (Day, Sankoff 1986). We'll see a sketch of the proof later.

Search space (exercise): 1) How many binary trees are there with n leaves? 2) Now, leaves can be labeled by the n taxa: multiply by n! 3) For phylogenies symmetries must be removed (left/right son is the same). How many trees remains? Is there a compact formula?

Agostino Dovier (DIMI, UDINE Univ.)

Phylogenetic trees

Alternative definitions

Given the set L = {a₁,..., a_n} of taxa, and characters
 I = {i₁,..., i_m}, the values of the function *f* can be compactly expressed by a matrix X:



• In the previous example:

	English	German	French	Spanish	Italian		
Hand	1	1	2	2	2		
Father	1	2	2	1	1		

Agostino Dovier (DIMI, UDINE Univ.)

12/21

Phylogenetic trees

Alternative definitions

- Focusing on Boolean values, X is a $m \times n \{0, 1\}$ matrix.
- Each column is a vector of $\{0, 1\}^m$
- Given two vectors in {0,1}^m they are connected if they differ in exactly one element (their Hamming distance d_H is 1—e.g., (0,0,1,0,1,1,0) and (0,0,1,1,1,1,0))
- {0,1}^{*m*} with the above notion generates an hypercube graph. There is always a path between two vectors.

Phylogenetic trees

Alternative definitions

- Focusing on Boolean values, X is a $m \times n \{0, 1\}$ matrix.
- Each column is a vector of $\{0, 1\}^m$
- Given two vectors in {0,1}^m they are connected if they differ in exactly one element (their Hamming distance d_H is 1—e.g., (0,0,1,0,1,1,0) and (0,0,1,1,1,1,0))
- {0,1}^{*m*} with the above notion generates an hypercube graph. There is always a path between two vectors.
- $V \subseteq \{0,1\}^m$ identifies a graph (V, E_V) , where $E_V = \{\{x,y\} \subseteq V \mid d_H(x,y) = 1\}.$
- If (V, E_V) is acyclic, we can easily build a rooted tree.

(日)

Phylogenetic trees

Alternative definitions

- Focusing on Boolean values, X is a $m \times n$ {0, 1} matrix.
- Each column is a vector of $\{0, 1\}^m$
- Given two vectors in {0,1}^m they are connected if they differ in exactly one element (their Hamming distance d_H is 1—e.g., (0,0,1,0,1,1,0) and (0,0,1,1,1,1,0))
- {0,1}^{*m*} with the above notion generates an hypercube graph. There is always a path between two vectors.
- $V \subseteq \{0, 1\}^m$ identifies a graph (V, E_V) , where $E_V = \{\{x, y\} \subseteq V \mid d_H(x, y) = 1\}.$
- If (V, E_V) is acyclic, we can easily build a rooted tree.
- The tree is a phylogeny if, starting from the root, (cladistic) character changes 1 → 0 are forbidden (Camin and Sokal 1965)
- Furthermore, one can require that the change of a character occurs in a unique edge (Dollo)

Agostino Dovier (DIMI, UDINE Univ.)

CP and Biology

Phylogenetic trees

Example (Camin and Sokal)



Agostino Dovier (DIMI, UDINE Univ.)

CP and Biology

Wrocław, September 2012 14 / 21

э

< ロ > < 同 > < 回 > < 回 >

Example (Camin and Sokal)



< ∃⇒

Example (Camin and Sokal)



< ∃⇒

< A

Example (Camin and Sokal)



< ∃→

Example (Camin and Sokal)



э

Example (Camin and Sokal)



< ∃→

Phylogenetic trees Another Problem

QCS/QDO (Qualitative Camin-Sokal/Dollo)

Given a set *L* of taxa, a set *I* of characters, and a number *k*, is there *V* s.t. |V| = k and $L \subseteq V \subseteq \{0, 1\}^n$ such that (V, E_V) is acyclic and connected (there is a path between each pair of nodes) and the induced tree is a phylogeny?

< ロ > < 同 > < 回 > < 回 >

Phylogenetic trees Another Problem

QCS/QDO (Qualitative Camin-Sokal/Dollo)

Given a set *L* of taxa, a set *I* of characters, and a number *k*, is there *V* s.t. |V| = k and $L \subseteq V \subseteq \{0, 1\}^n$ such that (V, E_V) is acyclic and connected (there is a path between each pair of nodes) and the induced tree is a phylogeny?

This problem (s) is NP-complete (Day, Johnson, Sankoff 1986). Reduction from vertex cover (see also the Haplotype Inference proof).

< ロ > < 同 > < 回 > < 回 > < 回 > <

NP completeness of the *k*-Incompatibility Problem

An incompatible gadget

Consider two taxa and three/four characters.



Agostino Dovier (DIMI, UDINE Univ.)

CP and Biology

Wrocław, September 2012 16 / 21

4 3 5 4 3

NP completeness of the *k***-Incompatibility Problem**

An incompatible gadget

Consider two taxa and three/four characters.



Agostino Dovier (DIMI, UDINE Univ.)

CP and Biology

Wrocław, September 2012 16 / 21

4 3 5 4 3

NP completeness of the *k***-Incompatibility Problem**

An incompatible gadget

Consider two taxa and three/four characters.



It can be proved that if the three characters contain pairs (0, 1), (1, 0), (1, 1) then the characters are not compatible.

Agostino Dovier (DIMI, UDINE Univ.)

CP and Biology

Wrocław, September 2012 16 / 21

Complexity

NP completeness of the k-Incompatibility Problem **Reduction (sketch)**

The property is exploited in the reduction of k-clique to k-Incompatibility.



$$|L| = 3 \frac{|V|(|V|-1)}{2}$$

Vertexes of a non-edge cannot stay both in a clique. Add the incompatible gadget to all non-edges.



Agostino Dovier (DIMI, UDINE Univ.)

Complexity

NP completeness of the k-Incompatibility Problem **Reduction (sketch)**

The property is exploited in the reduction of k-clique to k-Incompatibility.



$$L| = 3\frac{|V|(|V|-1)}{2}$$

Vertexes of a non-edge cannot stay both in a clique. Add the incompatible gadget to all non-edges. Set all other cells to 0

(1,2)		((1,3)		(1,4)		(2,3)		(2,4)			(3,4)						
1	0	0	0	0	0	0	1	1	0	0	0	0	0	0	0	0	0	0
2	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0
3	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	1	1	0
4	0	0	0	0	0	0	0	1	1	0	0	0	0	0	0	0	1	1

Encoding

Binary tree representation

 Input vector L of n elements (taxa) each of them characterized by a m-tuple of Boolean (character) values. E.g. m = 3, n = 4:

L = [[0, 1, 1], [1, 0, 0], [1, 1, 0], [1, 0, 1]]

- The Tree can be represented by a FD vector of t elements valued in 1,..., t + 1. Tree[i] = j means that node i is a son of node j. For the root r, Tree[r] = t + 1.
- To force the tree to be binary, for i = 1, ..., n we state

 $count(i, Tree, \leq, 2)$

- Taxa are the leaves of the tree. To avoid symmetries, wlog,
 - \checkmark leaves are nodes 1–*n*

✓ Tree[1] =
$$n + 1$$

✓ Tree[t] = t + 1 (t is the root)

For
$$i, j \in \{1, \dots, t\}$$
: $i < j \rightarrow \text{Tree}[i] \leq \text{Tree}[j]$

Agostino Dovier (DIMI, UDINE Univ.)

Encoding Hypercube tree

- Each node of the tree is assigned a *m*-tuple of Boolean Values. This is stored in a vector Chars.
- Chars[1]–Chars[*n*] are assigned using the input *L*. Values for internal nodes must be computed.
- For *i* < *j*, if Tree[*i*] = *j*, the Hamming difference of the corresponding tuples is 1. Precisely:

$$\text{Tree}[i] = j \rightarrow \left(\sum_{\ell=1}^{m} |\text{Chars}[i][\ell] - \text{Chars}[j][\ell]|\right) = 1$$

Agostino Dovier (DIMI, UDINE Univ.)

< ロ > < 同 > < 回 > < 回 > < 回 > <

Encoding Hypercube tree

- Each node of the tree is assigned a *m*-tuple of Boolean Values. This is stored in a vector Chars.
- Chars[1]–Chars[*n*] are assigned using the input *L*. Values for internal nodes must be computed.
- For *i* < *j*, if Tree[*i*] = *j*, the Hamming difference of the corresponding tuples is 1. Precisely:

$$\text{Tree}[i] = j \rightarrow \left(\sum_{\ell=1}^{m} |\text{Chars}[i][\ell] - \text{Chars}[j][\ell]|\right) = 1$$

- Actually, we can either relax the above constraint to ≤ 1 (see e.g. hand/father example, italian and spanish) or (alternatively)
- Add the redundant constraint

AllDifferentTuples(Chars)

Agostino Dovier (DIMI, UDINE Univ.)

CP and Biology

< ロ > < 同 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ >

Encoding *k*-incompatibility

- We need to state that a character changes (actually, increases) in at most one node. This makes the tree compatible with that character.
- Let Comp be a vector of *m* elements (one per character).
- For i < j, let $F_{i,j} = 1$ if Tree[i] = j, $F_{i,j} = 0$ otherwise.
- Then, for $\ell = 1, \ldots, m$ (and $i, j = 1, \ldots, n$:

$$\mathsf{Comp}[\ell] = \sum_{i < j} F_{i,j}(\mathsf{Chars}[i][\ell] - \mathsf{Chars}[j][\ell])$$

- Basically, after variable instantiation, Comp[l] will contain the number of changes of character l in the tree.
- The number of values different from 1 and 0 in Comp is forced to be less than or equal to *k*.
- See the program phylo.pl. Call the goal as: :phylo(List,0),phylo_decision(List,1,9,Tree).

Agostino Dovier (DIMI, UDINE Univ.)

CP and Biology

References

(Some) References

- Day, W.H.E., Johnson D.S., Sankoff, D. The Computational complexity of Inferring Rooted Phylogenies by Parsimony. Math. Biosciences 81:33-42, 1986.
- Day, W.H.E., Sankoff, D. Computational complexity of Inferring Phylogenies by Compatibility. Systematic Zoology 35(2):224–229, 1986.
- Erdem E., Lifschitz V., Nakhleh L., Ringe D. Reconstructing the Evolutionary History of Indo-European Languages Using Answer Set Programming. PADL 2003 160-176
- Thomas Schiex et al. Papers on complex pedigree reconstructions using weighted constraint satisfaction. In WCB 05, WCB 06, WCB 07,
- Moore N.C.A., and Prosser P. The Ultrametric Constraint and its Application to Phylogenetics. J. Artif. Intell. Res. 32:901–938, 2008 (also in WCB 06)
- Erdem E. Applications of Answer Set Programming in Phylogenetic Systematics MG65, LNCS 6565, 2011.
- Le Tiep, Nguyen Hieu, Pontelli Enrico, and Cao Son Tran. ASP at Work: An ASP Implementation of PhyloWS. ICLP 2012, LIPICS vol 17. (also in WCB 12)

・ロット (四) (日) (日) (日)