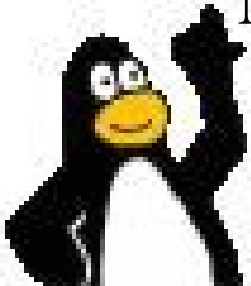


On phylogenetic trees – algebraic geometer's view

Joint project with Weronika Buczyńska
and other students

Microsoft Research Center – Trento, Feb. 16th, 2006

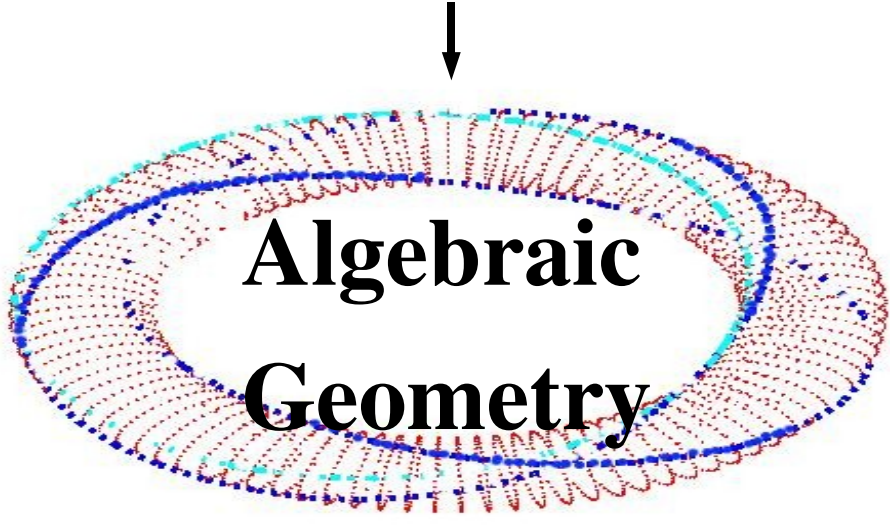
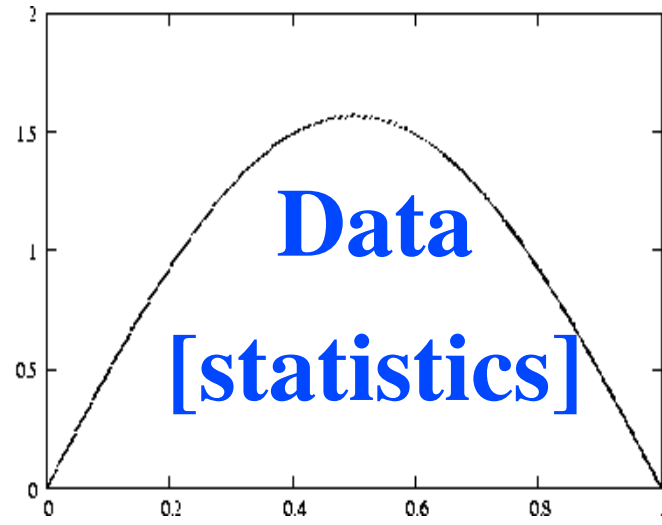


What about algebraic geometry?

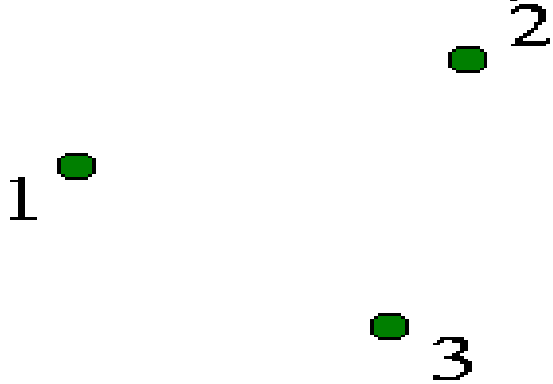


$I = \text{ideal}(x_3^4 - x_1^7, x_2^4 - x_1^6, x_2^3 - x_1^5, x_2^{19} - x_1^{31},$
 $x_2^{18} - x_1^{30}, x_2^{17} - x_1^{29}, x_2^{16} - x_1^{28}, x_3^{15} - x_1^{27},$
 $x_3^{14} - x_1^{26}, x_3^{13} - x_1^{25}, x_3^{12} - x_1^{24}, x_4^{11} - x_1^{23},$
 $x_4^{10} - x_1^{22}, x_4^9 - x_1^{21}, x_4^8 - x_1^{20}, x_2^7 - x_1^{32}, x_3^6 -$
 $x_1^{32}, x_4^5 - x_1^{32}, x_6^{19} - x_4^{31}, x_5^{19} - x_3^{31}, x_6^{18} -$
 $x_4^{30}, x_5^{18} - x_3^{30}, x_6^{17} - x_4^{29}, x_5^{17} - x_3^{29}, x_6^{16} -$
 $x_4^{28}, x_5^{16} - x_3^{28}, x_7^{15} - x_4^{27}, x_5^{15} - x_2^{27}, x_7^{14} -$
 $x_4^{26}, x_5^{14} - x_2^{26}, x_7^{13} - x_4^{25}, x_5^{13} - x_2^{25}, x_7^{12} -$
 $x_4^{24}, x_5^{12} - x_2^{24}, x_7^{11} - x_4^{23}, x_5^{11} - x_2^{23}, x_7^{10} -$
 $x_3^{22}, x_6^{10} - x_2^{22}, x_7^9 - x_3^{21}, x_6^9 - x_2^{21}, x_7^8 - x_3^{20},$
 $x_6^8 - x_2^{20}, x_6^7 - x_4^{32}, x_5^7 - x_3^{32}, x_5^6 - x_2^{32}, x_7^{18} -$
 $x_6^{19}, x_5^{18} - x_4^{19}, x_3^{18} - x_1^{19}, x_1^{17} - x_1^{10^*} 19, x_9^{17} -$
 $x_8^{19}, x_5^{16} - x_4^{17}, x_3^{16} - x_1^{17}, x_1^{16} - x_1^{10^*} 18, x_9^{16} -$
 $x_8^{18}, x_3^{14} - x_1^{12^*} 15, x_1^{14} - x_1^{10^*} 15, x_1^{13} - x_1^{8^*} 15, x_1^{12} -$
 $x_9^{14}, x_1^{10^*} 12 - x_8^{14}, x_9^{11} - x_1^{10^*} 11, x_9^{10} - x_1^{18^*} 31, x_1^{19^*} 29 -$
 $x_1^{17^*} 31, x_1^{19^*} 28 - x_1^{16^*} 31, x_7^{31} - x_1^{19^*} 32, x_1^{18^*} 29 - x_1^{16^*} 31, x_1^{17^*} 30 -$
 $x_1^{16^*} 31, x_1^{18^*} 28 - x_1^{16^*} 30, x_1^{19^*} 26 - x_1^{18^*} 27, x_1^{15^*} 30 - x_1^{14^*} 31,$
 $x_1^{19^*} 24 - x_1^{18^*} 25, x_1^{13^*} 30 - x_1^{12^*} 31, x_7^{30} - x_1^{18^*} 32, x_1^{17^*} 28 - x_1^{16^*} 29,$
 $x_1^{19^*} 22 - x_1^{17^*} 23, x_1^{11^*} 29 - x_1^{10^*} 31, x_1^{19^*} 20 - x_1^{17^*} 21, x_9^{29} - x_1^{8^*} 31,$
 $x_7^{29} - x_1^{17^*} 32, x_1^{17^*} 26 - x_1^{16^*} 27, x_8^{14^*} 19^* 32, x_1^{21^*} 25^* 31 -$
 $x_9^{13^*} 19^* 32,$

Relations
[algebra]



Example 1: independent variables



Relations in the spaces of triples

$$P_{\alpha\alpha\alpha}P_{\beta\beta\beta} = P_{\alpha\alpha\beta}P_{\beta\beta\alpha} = P_{\alpha\beta\alpha}P_{\beta\alpha\beta} = P_{\beta\alpha\alpha}P_{\alpha\beta\beta}$$

$$P_{\alpha\alpha\alpha}P_{\beta\beta\alpha} = P_{\beta\alpha\alpha}P_{\alpha\beta\alpha}$$

$$P_{\alpha\alpha\alpha}P_{\beta\alpha\beta} = P_{\beta\alpha\alpha}P_{\alpha\alpha\beta}$$

$$P_{\alpha\alpha\alpha}P_{\alpha\beta\beta} = P_{\alpha\beta\alpha}P_{\alpha\alpha\beta}$$

$$P_{\beta\beta\beta}P_{\alpha\alpha\beta} = P_{\alpha\beta\beta}P_{\beta\alpha\beta}$$

$$P_{\beta\beta\beta}P_{\alpha\beta\alpha} = P_{\alpha\beta\beta}P_{\beta\beta\alpha}$$

$$P_{\beta\beta\beta}P_{\beta\alpha\alpha} = P_{\beta\alpha\beta}P_{\beta\beta\alpha}$$

Consider 3 species with complementary features α and β which occur with respective probability p^i_α and p^i_β where $i=1,2,3$. Want to understand the distribution of features α and β in this 3-element system.

If the species are independent then the respective probability is as follows:

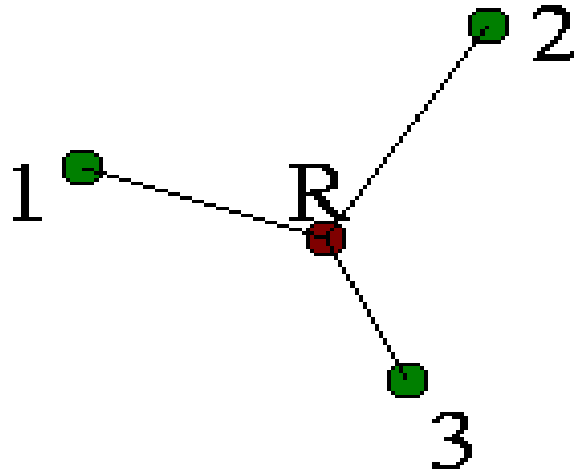
$$P_{\alpha\alpha\alpha} = p^1_\alpha p^2_\alpha p^3_\alpha \quad P_{\beta\beta\beta} = p^1_\beta p^2_\beta p^3_\beta$$

$$P_{\alpha\alpha\beta} = p^1_\alpha p^2_\alpha p^3_\beta \quad P_{\beta\beta\alpha} = p^1_\beta p^2_\beta p^3_\alpha$$

$$P_{\alpha\beta\alpha} = p^1_\alpha p^2_\beta p^3_\alpha \quad P_{\beta\alpha\beta} = p^1_\beta p^2_\alpha p^3_\beta$$

$$P_{\beta\alpha\alpha} = p^1_\beta p^2_\alpha p^3_\alpha \quad P_{\alpha\beta\beta} = p^1_\alpha p^2_\beta p^3_\beta$$

Example 2: variables in a tree



Now assume that the 3 species had a common ancestor R who had features α and β with probability r_α and r_β (which is not known) and each feature is duplicated in i -th species with probability a_i and changed to the other one with probability b_i

We find out that

If $r_\alpha \neq r_\beta$ then via a linear change of coordinates depending on r_α, r_β we arrive to previous equations, otherwise

$$P_{\alpha\alpha\alpha} = r_\alpha a_1 a_2 a_3 + r_\beta b_1 b_2 b_3 \quad P_{\beta\beta\beta} = r_\alpha b_1 b_2 b_3 + r_\beta a_1 a_2 a_3$$

$$P_{\alpha\alpha\beta} = r_\alpha a_1 a_2 b_3 + r_\beta b_1 b_2 a_3 \quad P_{\beta\beta\alpha} = r_\alpha b_1 b_2 a_3 + r_\beta a_1 a_2 b_3$$

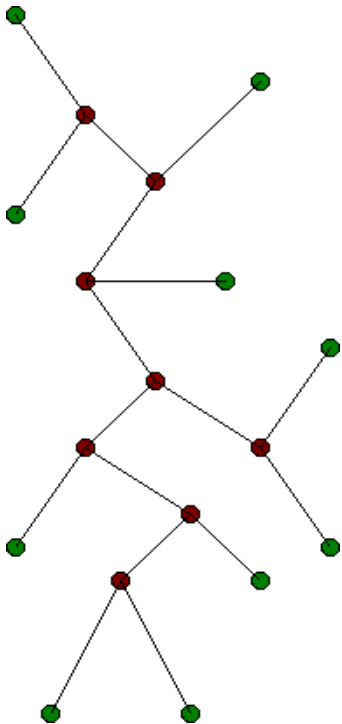
$$P_{\alpha\beta\alpha} = r_\alpha a_1 b_2 a_3 + r_\beta b_1 a_2 b_3 \quad P_{\beta\alpha\beta} = r_\alpha b_1 a_2 b_3 + r_\beta a_1 b_2 a_3$$

$$P_{\beta\alpha\alpha} = r_\alpha b_1 a_2 a_3 + r_\beta a_1 b_2 b_3 \quad P_{\alpha\beta\beta} = r_\alpha a_1 b_2 b_3 + r_\beta b_1 a_2 a_3$$

$$P_{\alpha\alpha\alpha} = P_{\beta\beta\beta} \quad P_{\alpha\alpha\beta} = P_{\beta\beta\alpha}$$

$$P_{\alpha\beta\alpha} = P_{\beta\alpha\beta} \quad P_{\beta\alpha\alpha} = P_{\alpha\beta\beta}$$

Definitions and assumptions



$$\begin{bmatrix} a & b \\ b & a \end{bmatrix}$$

- A tree T is 1-connected graph with the set of edges E and vertices $V = \text{Nodes} + \text{Leaves}$
- At each node we have 3 edges (*3-valent trees*).
- We consider propagation of two features α and β attached to each vertex of T (*binary trees*).
- The initial value of probability at a root is uniform and the change of the probability along every edge is symmetric (*symmetric trees*).
- We observe the distribution of features α and β only on leaves.

Varieties in projective space

Note that all relations in example 1 and 2 are homogeneous and therefore they define a set (variety) in a projective space with coordinates

$$[P_{\alpha\alpha\alpha}, P_{\alpha\alpha\beta}, P_{\alpha\beta\alpha}, P_{\beta\alpha\alpha}, P_{\beta\beta\alpha}, P_{\beta\alpha\beta}, P_{\alpha\beta\beta}, P_{\beta\beta\beta}]$$

Projective space = lines in a vector space. We will consider projective spaces whose coordinates are parametrized by possible distributions of features α and β on all leaves, let us call them the space of parameters.

Projective variety = zero set of homogeneous polynomials in coordinates of the vector space.

Geometric model of a phylogenetic tree

Geometric model of a phylogenetic tree T :

projective variety $X(T)$ in a projective space of parameters which describes geometric locus of probability distributions for a propagation of these features according to the tree T .

In general $X(T)$ is defined by formulas similar to these of Example 2:

$$P\left(X_l = a_{i(l)} : l \in L\right) = \sum_H \left(P\left(X_r = a_{i(r)}\right) \prod_e A_{i(s(e)), i(f(e))}^e \right)$$

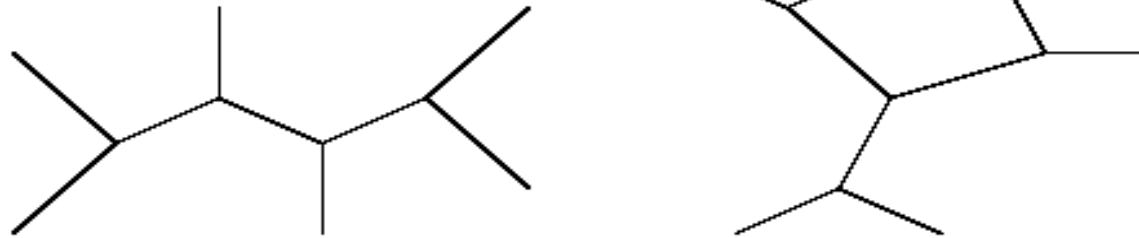
Using algebraic geometry

Understand relation of trees and their geometric models:
e.g. find intrinsic invariants of geometric models which distinguish their respective trees.

Example: Hilbert polynomial of a projective variety $X=X(T)$ in \mathbf{P}^N allows (under suitable conditions on X) to predict the number and degree of relations between the distributions of features on leaves.

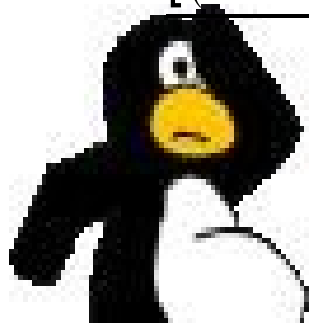
Example: two 3-valent trees with 6 leaves

Can we distinguish these trees by looking at their models?

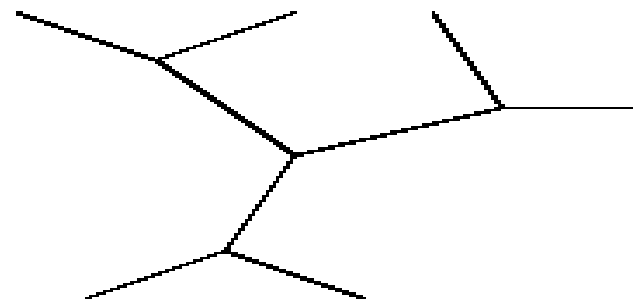
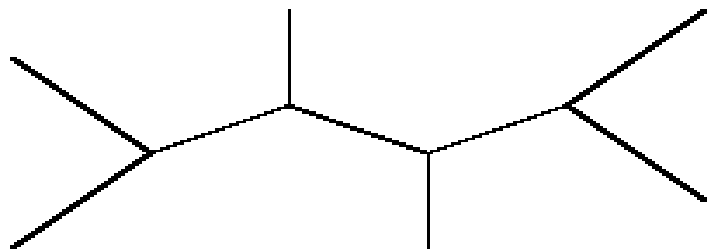


Hilbert polynomial does not distinguish these trees

$$\frac{\chi(t) \bar{\Gamma}(t+2)(t+3)(31t^6 + 374t^5 + 1942t^4 + 5616t^3 + 9511t^2 + 8988t + 3780)}{22680}$$



Example, continued: are the models distinguishable?



```
[jarekw@Hermiona jarekw]$ polymake
research/drzewa/maximapolymake/3caterpillar.poly
F2_VECTOR
F2_VECTOR
32 480 2400 6144 9312 8832 5280 1920 384
480 240 2400 9456 19904 24896 19104 8816 2240
2400 2400 760 5944 18976 32408 32168 18616 5824
6144 9456 5944 1316 8384 21648 29112 21552 8336
9312 19904 18976 8384 1392 7184 14584 14576 7176
8832 24896 32408 21648 7184 940 3816 5752 3816
5280 19104 32168 29112 14584 3816 406 1224 1224
1920 8816 18616 21552 14576 5752 1224 108 216
384 2240 5824 8336 7176 3816 1224 216 16
```

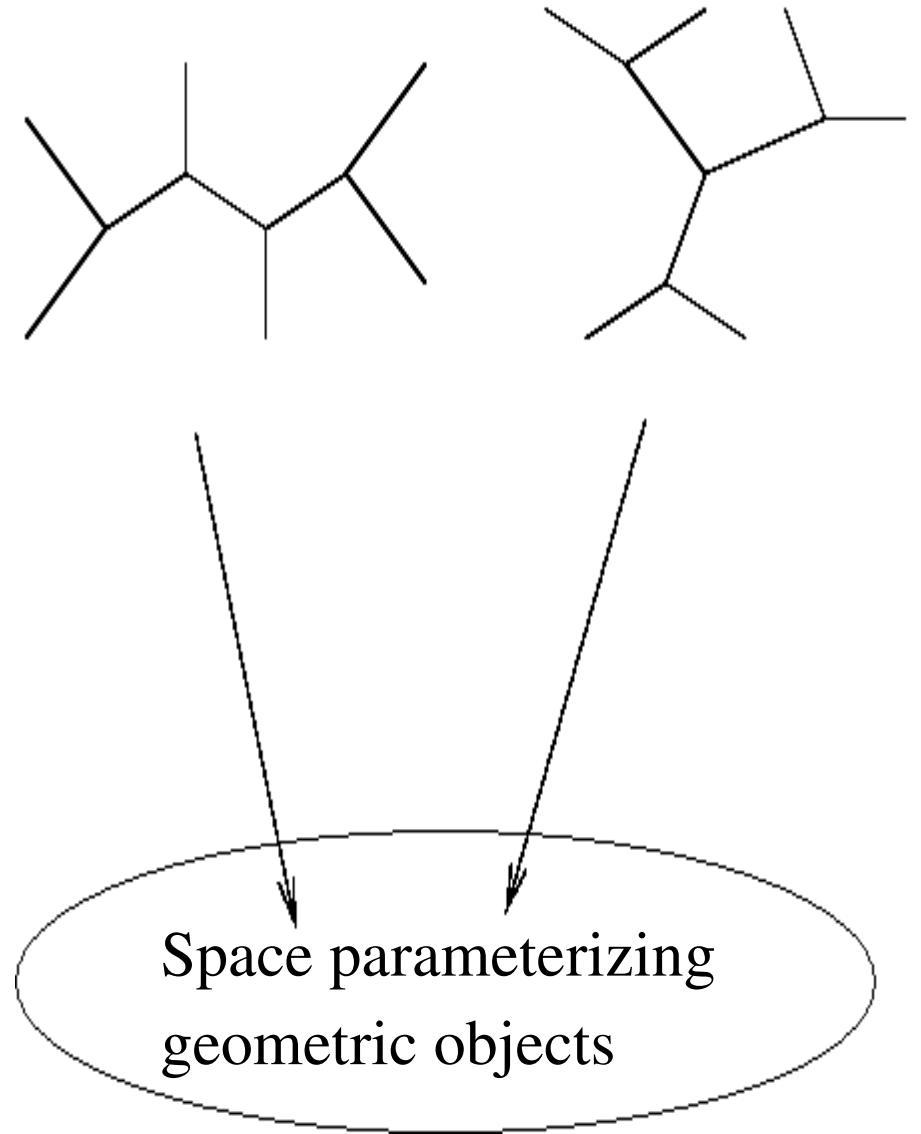
```
[jarekw@Hermiona jarekw]$ polymake
research/drzewa/maximapolymake/2star.poly
F2_VECTOR
F2_VECTOR
32 480 2400 6144 9312 8832 5280 1920 384
480 240 2400 9456 19920 24960 19200 8880 2256
2400 2400 760 5944 19008 32552 32408 18792 5872
6144 9456 5944 1316 8400 21744 29308 21720 8388
9312 19920 19008 8400 1392 7200 14640 14640 7200
8832 24960 32552 21744 7200 940 3820 5760 3820
5280 19200 32408 29308 14640 3820 406 1224 1224
1920 8880 18792 21720 14640 5760 1224 108 216
384 2256 5872 8388 7200 3820 1224 216 16
```

Conclusion: the geometric models are not isomorphic!!



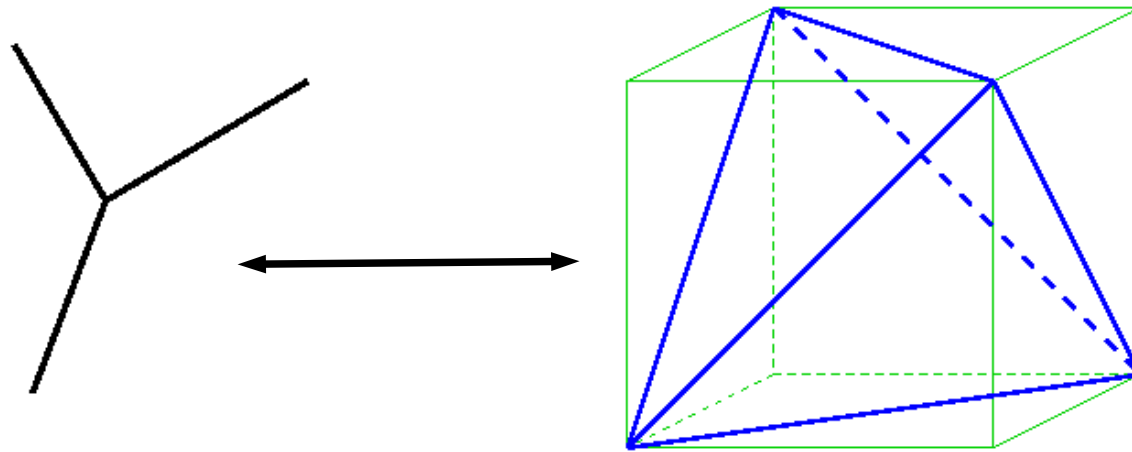
Explanation: deformation.

Theorem. Geometric models of any two binary symmetric 3-valent trees with the same number of leaves are deformation equivalent in the projective space of parameters. In particular equations defining one of them can be deformed to these defining any other one.



Tools: toric geometry

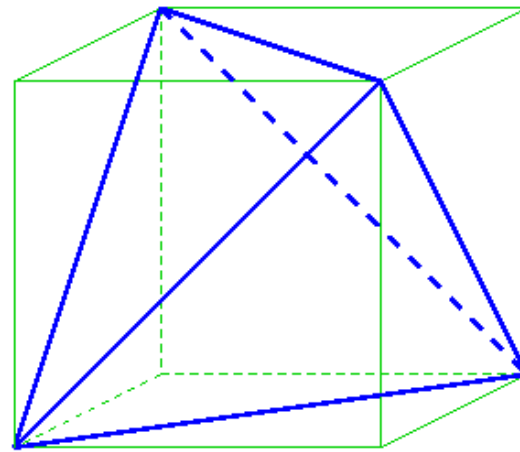
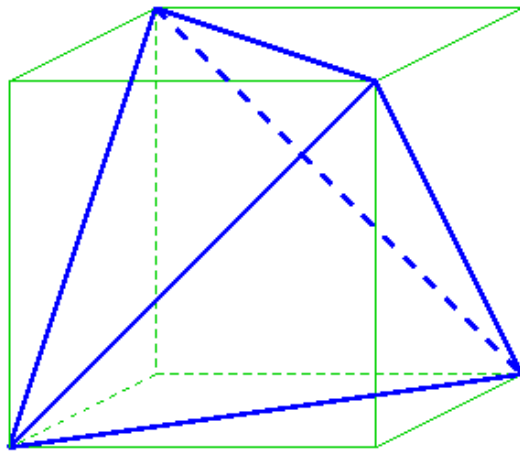
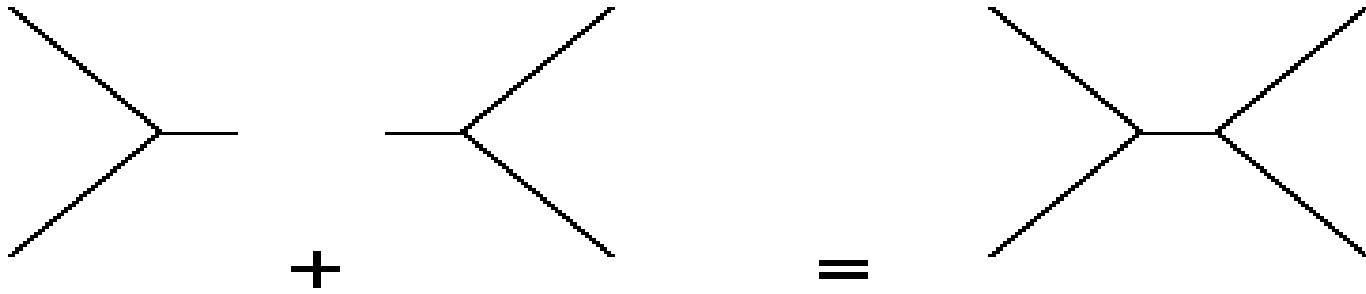
The geometric models of binary symmetric trees are *toric* hence they can be *described* by polytopes with vertices in a lattice.



Two lattices and a polytope

- Given a tree T we consider a free abelian group spanned on the edges of T , that is $M = \bigoplus_{e \in E} \mathbf{Z}e$
- The vertices of T live in the dual lattice $N = \text{Hom}(M, \mathbf{Z})$
 $V = \{v \in M : v(e) = 1 \text{ if } v \text{ is in edge } e \text{ and } v(e) = 0 \text{ otherwise}\}$
- In the lattice M consider a polytope $\Delta(T)$ with vertices in $\{u = \sum m(e)e : m(e) = 0, 1, v(u) = 0 \pmod{2}, \text{ for all nodes } v\}$
- The polytope $\Delta(T)$ provides a complete description of the variety $X(T)$.

Building up trees: fiber products of polytopes



Conclusion: what about algebraic geometry?

Let $X=X(T)$ in \mathbf{P}^N be a geometric model of a 3-valent binary symmetric tree. Then:

- $(X, O(1))$ is projectively normal
- X has terminal Gorenstein singularities
- X is a Fano variety of index 4, i.e. $K_X=O(-4)$

In dimension 3 any Fano variety with terminal Gorenstein singularities can be smoothed (to a *classical* Fano manifold) [Namikawa]. Does a similar phenomenon hold in higher dimensions?