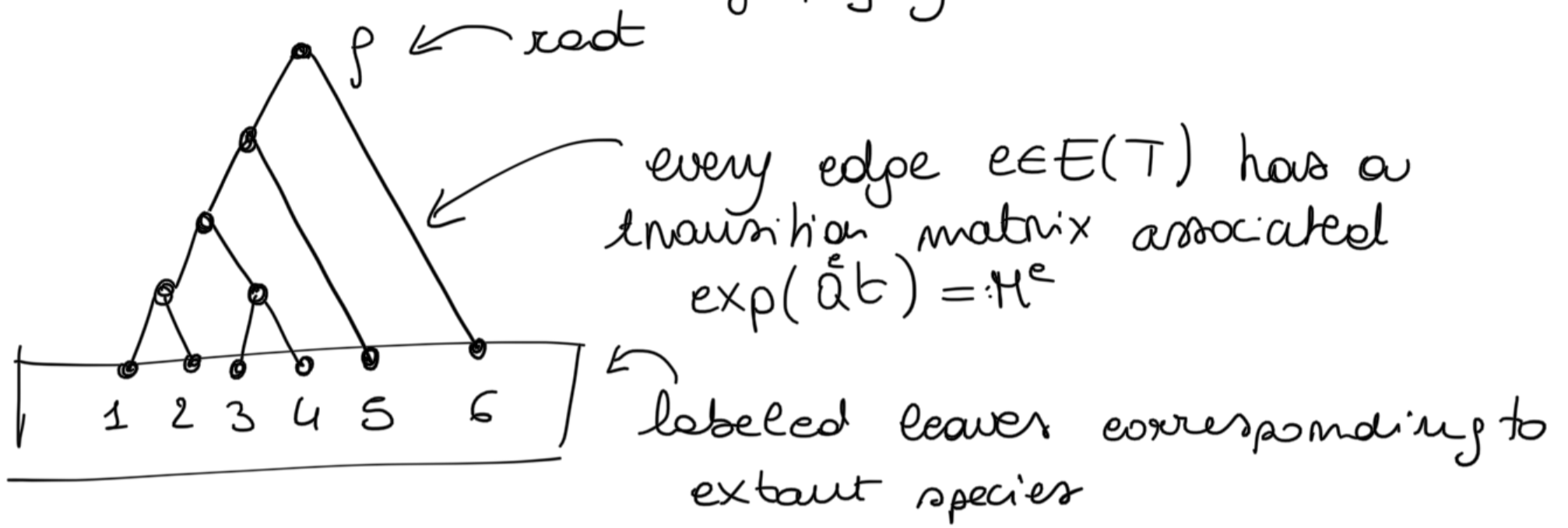


GROUP-BASED MODELS

07-05-21
Francesca

We build on the shoulders of phylogenetic models



ex 1: CFN model

$$Q^{\text{CFN}} = \begin{pmatrix} -\alpha & \alpha \\ \alpha & -\alpha \end{pmatrix}$$

↑ rate matrix

$$\exp(Q^{\text{CFN}} t) = \begin{bmatrix} \frac{1}{2}(1 + e^{-2\alpha t}) & \frac{1}{2}(1 - e^{-2\alpha t}) \\ \frac{1}{2}(1 - e^{-2\alpha t}) & \frac{1}{2}(1 + e^{-2\alpha t}) \end{bmatrix}$$

↑ transition matrix

ex 2: K3P model

$$Q^{\text{K3P}} = \begin{bmatrix} -\alpha - \beta - \gamma & \beta & \alpha & \gamma \\ \beta & -\alpha - \beta - \gamma & \gamma & \alpha \\ \alpha & \gamma & -\alpha - \beta - \gamma & \beta \\ \gamma & \alpha & \beta & -\alpha - \beta - \gamma \end{bmatrix}$$

↓

$$\exp(Q^{\text{K3P}} t) = \begin{bmatrix} a & b & c & d \\ b & a & d & c \\ c & d & a & b \\ d & c & b & a \end{bmatrix}$$

$$a = \frac{1}{4} \left(1 + e^{-2(\alpha+\beta)t} + e^{-2(\alpha+\gamma)t} + e^{-2(\beta+\gamma)t} \right)$$

$$b = \frac{1}{4} \left(1 - e^{-2(\alpha+\beta)t} - e^{-2(\alpha+\gamma)t} + e^{-2(\beta+\gamma)t} \right)$$

$$c = \frac{1}{4} \left(1 - e^{-2(\alpha+\beta)t} + e^{-2(\alpha+\gamma)t} - e^{-2(\beta+\gamma)t} \right)$$

$$d = \frac{1}{4} \left(1 + e^{-2(\alpha+\beta)t} - e^{-2(\alpha+\gamma)t} - e^{-2(\beta+\gamma)t} \right)$$

We end up with a model that is not so easy to describe!

What? The story of a simplification

How? Put some additional structure on the model in order to have a "nicer" description.

Def: A phylogenetic model is called a GROUP-BASED model if for every transition matrix M^e associated to $e \in E(T)$ there is a function $f^e: G \rightarrow \mathbb{R}$ from a finite abelian group $(G, +)$ s.t

$$M_{gh}^e = M^e(g, h) = f^e(g-h)$$

Note: In most of the cases $G = \mathbb{Z}_2$ or $G = \mathbb{Z}_2 \times \mathbb{Z}_2$

ex 1: CFN is a group based model

$$\exp(Q^{CFN} t) = \begin{bmatrix} \frac{1}{2}(1+e^{-2\alpha t}) & \frac{1}{2}(1-e^{-2\alpha t}) \\ \frac{1}{2}(1-e^{-2\alpha t}) & \frac{1}{2}(1+e^{-2\alpha t}) \end{bmatrix} \quad \begin{array}{l} \{A, C\} \rightarrow 0 \\ \{G, T\} \rightarrow 1 \end{array}$$

Define $f: \mathbb{Z}_2 \rightarrow \mathbb{R}$ as

$$f(0) = \frac{1}{2}(1+e^{-2\alpha t}) \quad f(1) = \frac{1}{2}(1-e^{-2\alpha t})$$

$$\Rightarrow M_{00} = f(0-0) = f(0) \quad M_{01} = f(0-1) = f(1) = M_{10}$$

$$M_{11} = f(1-1) = f(0)$$

ex 2: K3P is a group-based model

$$\exp(Q^{K3P} t) = \begin{bmatrix} a & b & c & d \\ b & a & d & c \\ c & d & a & b \\ d & c & b & a \end{bmatrix} \quad \begin{array}{l} A \rightarrow (0,0) \\ C \rightarrow (0,1) \\ G \rightarrow (1,0) \\ T \rightarrow (1,1) \end{array}$$

$$a = \frac{1}{4} \left(1 + e^{-2(\alpha+\beta)t} + e^{-2(\alpha+\gamma)t} + e^{-2(\beta+\gamma)t} \right)$$

$$b = \frac{1}{4} \left(1 - e^{-2(\alpha+\beta)t} - e^{-2(\alpha+\gamma)t} + e^{-2(\beta+\gamma)t} \right)$$

$$c = \frac{1}{4} \left(1 - e^{-2(\alpha+\beta)t} + e^{-2(\alpha+\gamma)t} - e^{-2(\beta+\gamma)t} \right)$$

$$d = \frac{1}{4} \left(1 + e^{-2(\alpha+\beta)t} - e^{-2(\alpha+\gamma)t} - e^{-2(\beta+\gamma)t} \right)$$

Define $f: \mathbb{Z}_2 \times \mathbb{Z}_2 \longrightarrow \mathbb{R}$ as

$$f(0,0) = a, \quad f(0,1) = b, \quad f(1,0) = c, \quad f(1,1) = d$$

$$M_{(0,1),(1,1)} = f(0,1) - f(1,1) = f(1,0) = c$$

Why not $G = \mathbb{Z}_4$? For example, $\bar{f}: \mathbb{Z}_4 \longrightarrow \mathbb{R}$ s.t

$$\bar{f}(0) = a, \quad \bar{f}(1) = b, \quad \bar{f}(2) = c, \quad \bar{f}(3) = d$$

$$M_{0,1} \neq \bar{f}(0-1) = \bar{f}(-1) = \bar{f}(3) = d$$

↑ the issue is that elements in \mathbb{Z}_4 are not inverses of themselves, so the symmetry of the matrix is not guaranteed.

Def: A 1-dimensional character of a group G is a group homomorphism $\chi: G \longrightarrow \mathbb{C}^*$
 ↑ additive ↑ multiplicative

Prop: $\hat{G} = \{ \chi: G \longrightarrow \mathbb{C}^* \mid \chi \text{ homo} \}$ is a group with the operation $(\chi_1 \cdot \chi_2)(g) = \chi_1(g)\chi_2(g)$.
 Moreover, $G \cong \hat{\hat{G}}$ if G is a finite abelian group.

Pf

It is not difficult to check that the operation is well defined and that \hat{G} is a group.

We now show that $G \cong \hat{\hat{G}}$. Given that G is a finite abelian group, $G \cong \mathbb{Z}_{p_1^{\alpha_1}} \times \dots \times \mathbb{Z}_{p_r^{\alpha_r}}$, for p_1, \dots, p_r primes and $\alpha_1, \dots, \alpha_r \in \mathbb{N}$.

It is also easy to see that $\hat{G \times H} = \hat{G} \times \hat{H}$, hence it

suffices to show the isomorphism for $G \cong \mathbb{Z}_m$.
 In this case, G is cyclic and $\chi: G \rightarrow \mathbb{C}^*$ is
 characterized by where it sends a generator, which will
 go to a m -th root of unity (since $\chi(g)^m = \chi(mg) =$
 $= \chi(0) = 1 \Rightarrow \chi(g)$ m -th root of 1). This correspondence
 describe the isomorphism desired.

Def: Discrete Fourier transform of $f: G \rightarrow \mathbb{C}$ is
 $\hat{f}: \hat{G} \rightarrow \mathbb{C}$ defined by

$$\hat{f}(\chi) = \sum_{g \in G} f(g) \chi(g)$$

Intuition: $f \longleftrightarrow A = (a_g)_{g \in G} \in \mathbb{C}^G \Rightarrow a_g \in \mathbb{C}$

$\hat{f} \longleftrightarrow B = (b_\chi)_{\chi \in \hat{G}} \in \mathbb{C}^{\hat{G}} \Rightarrow b_\chi \in \mathbb{C}$

$$B = \begin{bmatrix} b_{\chi_1} \\ \vdots \\ b_{\chi_{|G|}} \end{bmatrix}$$

A

$|G| \times |G|$ matrix

Fourier coordinates
 \updownarrow
 phase space

Th: $f(g) = \frac{1}{|G|} \sum_{\chi \in \hat{G}} \hat{f}(\chi) \overline{\chi(g)}$

Def: $f_1, f_2: G \rightarrow \mathbb{C}$ functions

$$(f_1 * f_2)(g) = \sum_{h \in G} f_1(h) f_2(g-h)$$

Prop: $\widehat{f_1 * f_2}(\chi) = \hat{f}_1(\chi) \hat{f}_2(\chi)$

$g-h = g'$

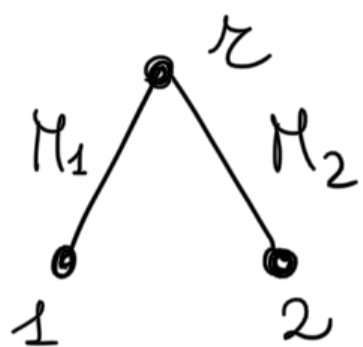
pt $\widehat{f_1 * f_2}(\chi) = \sum_{g \in G} \sum_{h \in G} f_1(h) f_2(g-h) \chi(g) =$

$$= \sum_{g' \in G} \sum_{h \in G} f_1(h) f_2(g') \chi(g'+h) =$$

$$= \left(\sum_{h \in G} f_1(h) \chi(h) \right) \left(\sum_{g' \in G} f_2(g') \chi(g') \right)$$

$$= f_1(x) f_2(x)$$

ex T₂:



Let π be the root distribution
 $\pi: G \rightarrow \mathbb{R}$

The probability of observing (g_1, g_2) on the leaves:

$$p(g_1, g_2) = \sum_{h \in G} \pi(h) M_1(h, g_1) M_2(h, g_2) =$$

$$= \sum_{h \in G} \pi(h) f_1(g_1 - h) f_2(g_2 - h)$$

From the phase space:

$$\begin{cases} g_1 - h = g'_1 \\ g_2 - h = g'_2 \end{cases}$$

$$\hat{p}(x_1, x_2) = \sum_{(g_1, g_2) \in G^2} \sum_{h \in G} \pi(h) f_1(g_1 - h) f_2(g_2 - h) x_1(g_1) x_2(g_2) =$$

$$= \sum_{g'_1 \in G} \sum_{g'_2 \in G} \sum_{h \in G} \pi(h) f_1(g'_1) f_2(g'_2) x_1(g'_1 + h) x_2(g'_2 + h) =$$

$$= \left(\sum_{h \in G} \pi(h) x_1(h) x_2(h) \right) \left(\sum_{g'_1 \in G} f_1(g'_1) x_1(g'_1) \right) \left(\sum_{g'_2 \in G} f_2(g'_2) x_2(g'_2) \right)$$

$$\uparrow x_i(g'_i + h) = x_i(g'_i) x_i(h), \quad \forall i = 1, 2$$

$$= \hat{\pi}(x_1 x_2) \hat{f}_1(x_1) \hat{f}_2(x_2)$$

↑ In the phase space the probability $p(g_1, g_2)$ has this "monomial" form, which leads to the following result.

Prop: Consider the map

$$\varphi_G: \mathbb{C}^G \times \mathbb{C}^G \times \mathbb{C}^G \longrightarrow \mathbb{C}^{G \times G}$$

$$\left((a_g)_{g \in G}, (b_g)_{g \in G}, (c_g)_{g \in G} \right) \longmapsto \left(a_{g_1} b_{g_2} c_{g_1 + g_2} \right)_{(g_1, g_2) \in G^2}$$

The Zariski closure of the image of φ_G is equal to the Zariski closure of the cone over the image of the group-based model of ex T₂.

Rem: φ_G is monomial map, so the Zariski closure of its image is a toric variety.

pf (idea)

Observe that in the phase space

$$\hat{p}(x_1, x_2) = \hat{f}_1(x_1) \hat{f}_2(x_2) \hat{\pi}(x_1 x_2)$$

which is exactly of the same form of the image of φ_G . However, f_1, f_2, π are subjects to the following constraints:

$$\sum_{g \in G} f_1(g) = \sum_{g \in G} f_2(g) = \sum_{g \in G} \pi(g)$$

$$\Rightarrow \hat{f}_1(\text{id}) = \hat{f}_2(\text{id}) = \hat{\pi}(\text{id}) = 1, \text{ since}$$

$$\text{id}(g) = 1 \quad \forall g \in G \text{ and thus}$$

$$\hat{f}_i(\text{id}) = \sum_{g \in G} f_i(g) \text{id}(g) = 1 \quad (\text{same for } \pi)$$

$$\hat{f}_1 = \begin{bmatrix} 1 \\ a_1 \\ \vdots \\ a_m \end{bmatrix}, \quad \hat{f}_2 = \begin{bmatrix} 1 \\ b_1 \\ \vdots \\ b_m \end{bmatrix}, \quad \hat{\pi} = \begin{bmatrix} 1 \\ c_1 \\ \vdots \\ c_m \end{bmatrix} \in U_0 \cong \mathbb{R}^m$$

" $\{(x_0: \dots: x_m) | x_0 \neq 0\}$

So we actually have a map of the form of φ_G restricted to a dense subset $U_0 \times U_0 \times U_0$.

In order for it to be well-defined, we need to take the cone of the image (or in other terms to "projectivize"). Then we can conclude by taking the Zariski closures of both and conclude that they must be equal.

Up to now we have considered only group-based model with no additional assumptions on the root distribution or on the transition matrices. They can be collected under the name of GENERIC group-based models.

ex: Consider the generic group-based model on $\mathbb{Z}_2 \times \mathbb{Z}_2$

with identification:

$$A = (0,0), C = (0,1), G = (1,0), T = (1,1)$$

$$\varphi_G(a_A, a_C, a_G, a_T, b_A, b_C, b_G, b_T, c_A, c_C, c_G, c_T) =$$

$$= \begin{bmatrix} a_A b_A c_A & a_A b_C c_C & a_A b_G c_G & a_A b_T c_T \\ a_C b_A c_C & a_C b_C c_C & a_C b_G c_T & a_C b_T c_G \\ a_G b_A c_G & a_G b_C c_T & a_G b_G c_A & a_G b_T c_C \\ a_T b_A c_T & a_T b_C c_G & a_T b_G c_C & a_T b_T c_A \end{bmatrix}$$


Macaulay2 \rightarrow vanishing ideal, to check it in tonic

Example of generator: $q_{AA} q_{CG} q_{TC} - q_{TA} q_{CC} q_{AG}$ in Fourier coordinates.

It is possible to go back to standard coordinates, for example:

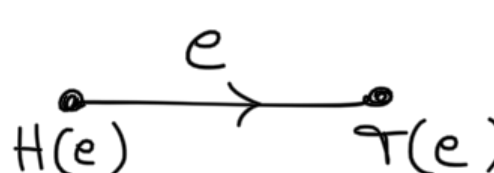
$$q_{CG} = \sum_{(i,j), (k,e) \in \mathbb{Z}_2 \times \mathbb{Z}_2} \underbrace{\chi_C(i,j) \chi_G(k,e)}_{(-1)^{j+k}} p((i,j), (k,e))$$

with $\chi_A(i,j) = 1$, $\chi_C(i,j) = (-1)^j$, $\chi_G(i,j) = (-1)^i$, $\chi_T(i,j) = (-1)^{i+j}$

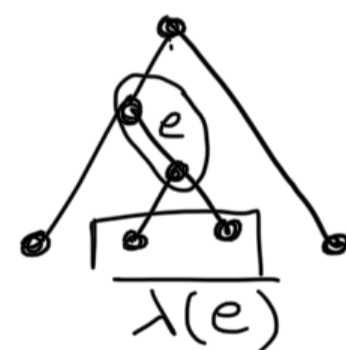
We generalise the case  to the model represented by a rooted tree with m -leaves T .

• $\pi: G \rightarrow \mathbb{R}$ root distribution

• $M^e(g,h)$ transition matrix

• $\forall e \in E(T)$ 

• $\lambda(e) =$ set of leaves lying below e



• $\mathcal{Y}_{int} = V(T) \setminus \{\text{leaf}\}$

The probability of seeing the tuple (g_1, \dots, g_m) on the leaves is:

$$p(g_1, \dots, g_m) = \sum_{g \in G^{\mathcal{Y}_{int}}} \pi(g^x) \prod_{e \in E(T)} M^e(g_{T(e)}, g_{H(e)}) =$$

$$= \sum_{g \in G^{\text{int}}} \pi(g_r) \prod_{e \in E(T)} f^e(g_{H(e)} - g_{T(e)})$$

Prop: Let $p: G^m \rightarrow \mathbb{R}$ denote the probability distribution of a group-based model on a m -leaf tree. Then the Fourier transform of p is

$$\hat{p}(\chi_1, \dots, \chi_m) = \hat{\pi} \left(\prod_{i=1}^m \chi_i \right) \prod_{e \in E(T)} \hat{f}^e \left(\prod_{i \in \lambda(e)} \chi_i \right)$$

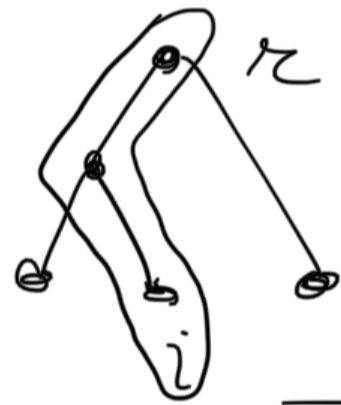
pf

$$\hat{p}(\chi_1, \dots, \chi_m) = \sum_{(g_1, \dots, g_m) \in G^m} \sum_{g \in G^{\text{int}}} \pi(g_r) \prod_{e \in E(T)} f^e(g_{H(e)} - g_{T(e)}) \prod_{i=1}^m \chi_i(g_i) =$$

change of variable: $g_e = g_{H(e)} - g_{T(e)} \quad \forall e \in E(T)$

$$\Rightarrow g_i = g_r + \sum_{e \in P(i)} g_e$$

where $P(i) = \{ \text{edges connecting the root to the } i\text{-th leaf} \}$

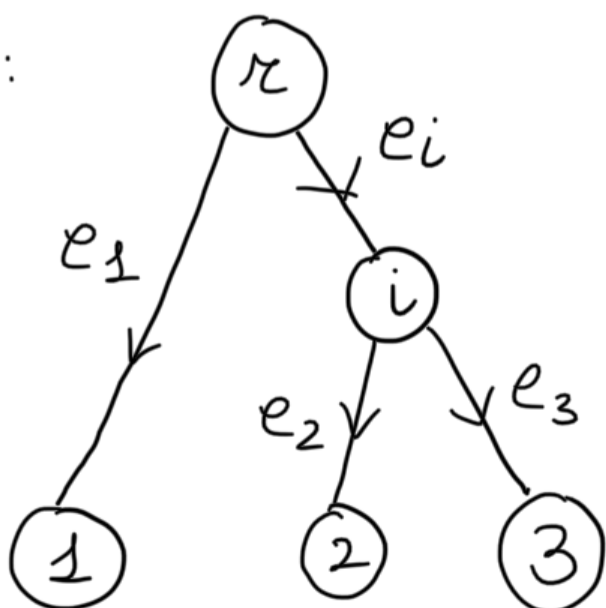


$$= \sum_{g_r \in G} \sum_{g \in G^{E(T)}} \pi(g_r) \prod_{e \in E(T)} f^e(g_e) \prod_{i=1}^m \chi_i \left(g_r + \sum_{e \in P(i)} g_e \right) =$$

$$\stackrel{*}{=} \left(\sum_{g_r \in G} \pi(g_r) \prod_{i=1}^m \chi_i(g_r) \right) \left(\prod_{e \in E(T)} \sum_{g_e \in G} f^e(g_e) \prod_{i \in \lambda(e)} \chi_i(g_e) \right) =$$

$$= \hat{\pi} \left(\prod_{i=1}^m \chi_i \right) \prod_{e \in E(T)} \hat{f}^e \left(\prod_{i \in \lambda(e)} \chi_i \right)$$

ex:



$$\hat{p}(x_1, x_2, x_3) = \sum_{g_r \in G} \sum_{g_{e_1}, g_{e_2}, g_{e_3}} \pi(g_r) f^{e_1}(g_{e_1}) f^{e_2}(g_{e_2}) f^{e_3}(g_{e_3}) x_1(g_r + g_{e_1}) x_2(g_r + g_{e_1} + g_{e_2}) x_3(g_r + g_{e_1} + g_{e_2} + g_{e_3})$$

$$\stackrel{*}{=} \left(\sum_{g_r \in G} \pi(g_r) x_1(g_r) x_2(g_r) x_3(g_r) \right) \cdot \left(\sum_{g_{e_1}} f^{e_1}(g_{e_1}) x_1(g_{e_1}) \right) \left(\sum_{g_{e_2}} f^{e_2}(g_{e_2}) x_2(g_{e_2}) \right) \left(\sum_{g_{e_3}} f^{e_3}(g_{e_3}) x_3(g_{e_3}) \right)$$

$$= \hat{\pi}(x_1 x_2 x_3) \hat{f}^{e_1}(x_1) \hat{f}^{e_2}(x_2) \hat{f}^{e_3}(x_3)$$

Def: In the Fourier coordinates, the Zariski closure of the cone over the perennic group-based model on m -leaves is equal to the Zariski closure of the image of

$$\varphi^{G, T}: \mathbb{C}^{|G| \times (|E(T)| + 1)} \longrightarrow \mathbb{C}^{|G|^m}$$

$$\left(\left(a_g^e \right)_{e \in E(T), g \in G}, \left(a_g^r \right)_{g \in G} \right) \longmapsto a_{\sum_{i \in [m]} g_i}^r \prod_{e \in E(T)} a_{\sum_{i \in [m]} g_i}^e$$

ex: $T = (\heartsuit)$

$$\varphi^{G, T}(a) = a_{g_1 + g_2 + g_3}^r a_{g_2 + g_3}^{e_1} a_{g_1}^{e_2} a_{g_2}^{e_3} a_{g_3}^{e_3}$$

Up to now we have studied the perennic group-based model. Now we see what happens when we add an hp on the root distribution. In particular, it is often the case that the root distribution is assumed to be the stationary distribution.

It can be seen that in all the cases that we have encountered it coincides with the uniform distribution.

Def: A stationary distribution π is s.t. $\pi \exp(Qt) = \pi$

$\forall \alpha$.

ex: Consider the JC^{69} model with the uniform distribution

$$\pi(A) = \pi(C) = \pi(G) = \pi(T) = \frac{1}{4}$$

$$\begin{bmatrix} \frac{1}{4} & \frac{1}{4} & \frac{1}{4} & \frac{1}{4} \end{bmatrix} \frac{1}{4} \begin{bmatrix} 1+3e^{-4t\alpha} & 1-e^{-4t\alpha} & 1-e^{-4t\alpha} & 1-e^{-4t\alpha} \\ 1-e^{-4t\alpha} & 1+3e^{-4t\alpha} & 1-e^{-4t\alpha} & 1-e^{-4t\alpha} \\ 1-e^{-4t\alpha} & 1-e^{-4t\alpha} & 1+3e^{-4t\alpha} & 1-e^{-4t\alpha} \\ 1-e^{-4t\alpha} & 1-e^{-4t\alpha} & 1-e^{-4t\alpha} & 1+3e^{-4t\alpha} \end{bmatrix} =$$

"exp($Q^{JC}t$)"

$$= \begin{bmatrix} \frac{1}{4} & \frac{1}{4} & \frac{1}{4} & \frac{1}{4} \end{bmatrix}$$

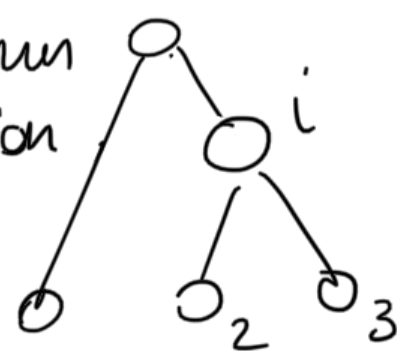
Prop: $f: G \rightarrow \mathbb{C}$ constant defined by $f(g) = \frac{1}{|G|}$

Then $\hat{f}(X) = \begin{cases} 1 & \text{if } X = id \\ 0 & \text{otherwise} \end{cases}$

Cor: In the Fourier coordinates, the Zariski closure of the image of the generic group-based model with uniform node distribution is equal to the Zariski closure of the image of

$$\varphi^{G,T}: \begin{matrix} \mathbb{C}^{|G| \times E(T)} \\ (a_g)_{\substack{e \in E(T) \\ g \in G}} \end{matrix} \longrightarrow \begin{matrix} \mathbb{C}^{|G|^m} \\ \prod_{e \in E(T)} \prod_{i \in \lambda(e)} a_{\sum g_i}^e \end{matrix} \begin{matrix} \text{if } \sum_{i \in [m]} g_i = 1 \\ 0 \text{ otherwise} \end{matrix}$$

ex: $T = (\heartsuit)$ with JC with uniform node distribution



$$f^e(A) = a$$

$$f^e(C) = f^e(G) = f^e(T) = b$$

a and b are functions of α and α

$$\Rightarrow \hat{f}^e(A) = a + 3b, \quad \hat{f}^e(C) = \hat{f}^e(G) = \hat{f}^e(T) = a - b$$

$$\Rightarrow a_c^e = a_G^e = a_T^e \quad \forall c \in E(T)$$

Using these constraints we get linear invariants.

$$q_{CGT} = \begin{matrix} e_1 & e_2 & e_3 & e_i \\ a_c & a_G & a_T & a_{\underbrace{G+T}_C} \end{matrix} = \begin{matrix} C = (0, 1) \\ G = (1, 0) \\ T = (1, 1) \end{matrix}$$

$$= \begin{matrix} e_1 & e_2 & e_3 & e_i \\ a_c & a_c & a_c & a_c \end{matrix}^C$$

$$q_{GTC} = \begin{matrix} e_1 & e_2 & e_3 & e_i \\ a_G & a_T & a_c & a_{\underbrace{T+G}_C} \end{matrix} = \begin{matrix} a_c & a_c & a_c & a_c \end{matrix}^e$$

$\Rightarrow q_{CGT} - q_{GTC}$ is a linear invariant for the model