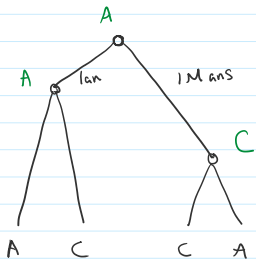


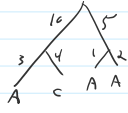
Inference statistique

18 novembre 2022 10:24



En principe, chaque branche doit avoir ses propres coûts, qui dépendent de leur longueur et probabilité.

- On suppose qu'on nous donne
- un arbre T , feuille / caractère
 - longueurs sur branches
 - chaque branche a une matrice de probab. de mutation qui lui est propre



$\forall uv \in E(T), Pr_{uv}[a,b] = \text{prob. que } a \text{ mute en } b \text{ sur } uv$

A → A	4/5
A → C	1/5
C → A	1/3
C → C	2/3

A → A	2/3
A → C	1/3
C → A	1/4
C → C	3/4

Supposons $\Sigma = \{A, C\}$

raisonnable $(T) \approx Pr[\text{séquences} | T]$

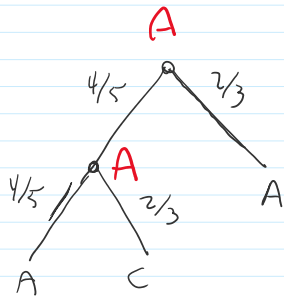
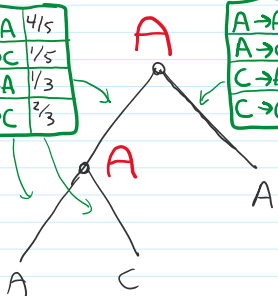
$= \sum_{\text{étiquetages possibles de } T} Pr[\text{l'étiquetage } \pi \text{ en lieu}]$

$= Pr \left[\begin{array}{c} A \\ / \quad \backslash \\ A \quad C \end{array} \right] + Pr \left[\begin{array}{c} A \\ / \quad \backslash \\ C \quad A \end{array} \right] + Pr \left[\begin{array}{c} C \\ / \quad \backslash \\ A \quad A \end{array} \right] + Pr \left[\begin{array}{c} C \\ / \quad \backslash \\ C \quad A \end{array} \right]$

ex:

A → A	4/5
A → C	1/5
C → A	1/3
C → C	2/3

A → A	2/3
A → C	1/3
C → A	1/4
C → C	3/4



$Pr \left[\begin{array}{c} A \\ / \quad \backslash \\ A \quad C \end{array} \right] = \frac{4}{5} \cdot \frac{2}{3} \cdot \frac{4}{5} \cdot \frac{2}{3} \cdot \frac{1}{2}$

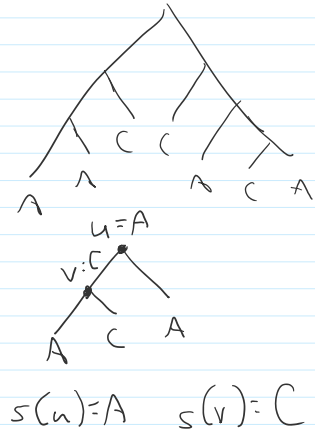
prob. d'avoir A à la racine

Prob. que tous les changements prédits sur les branches

survivent = $\frac{2}{3} \cdot \frac{4}{5} \cdot \frac{4}{5} \cdot \frac{1}{5} = \frac{32}{375}$

Il faut aussi considérer la prob que A apparaisse à la racine. Souvent, on attribue de façon uniforme, i.e. $Pr[A \text{ racine}] = \frac{1}{2} = \frac{1}{|\Sigma|}$

- En général, T a n feuilles (1 caractère)
accès à $\text{Pr}_{uv}[a, b]$
- Soit E l'ensemble des étiquetages possibles
des nœuds internes de T
↳ $s \in E \rightarrow \forall u \in V(T), s(v) = \text{caractère de } v$
- $\forall a \in \Sigma, \pi(a) = \text{Prob. d'observer } a \text{ à la racine}$

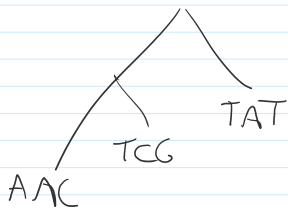


• On cherche

$$\sum_{s \in E} \left(\prod_{uv \in E(T)} \text{Pr}_{uv}[s(u), s(v)] \right) \cdot \pi(s(\text{racine}))$$

produit itéré

Ceci est pour un seul caractère: hypothèse d'indépendance entre les positions



→ la vraisemblance de T
= produit des vraisemblances par caractère

Calculer $\sum_{s \in E} (\dots)$ peut prendre un temps $O(|\Sigma|^I)$ où $I = \# \text{ nœuds interne}$
naïvement

Il existe un algo en temps polynomial pour calculer $\sum_{s \in E} (\dots)$
par prog. dynamique (algo de Felsenstein)