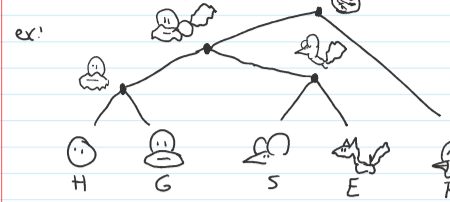


Phylogénétique

• Étude de relations évolutives entre objets
 ↳ objets appelés taxa (pluriel de taxon)



ex: langues



Suppositions:

- évolution en arbre
- chaque taxon ancestral donne lieu à 2 taxa enfants (arbres binaire)
- feuilles = taxa contemporains (d'aujourd'hui)
- nœuds internes = taxa ancestraux inconnus

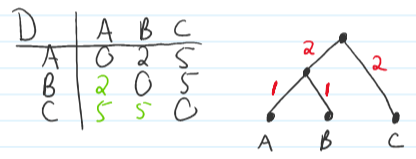
Entrée: feuilles (séquences, distances, ...)
 Sortie: arbre
 ↳ étiquette sur les arêtes qui représente un temps d'évol.
 ↳ état des ancêtres

Approches:

- ① par distances
- ② par caractères
- ③ statistiques

① Approches par distances

ex: 3 taxa A, B, C $D_{x,y}$ = distance entre X et Y



• Algorithme UPGMA
 ↳ unweighted pair group with arithmetic mean

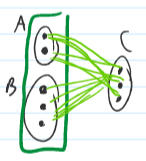
Entrée: matrice D (rangées = taxa, colonnes = taxa)
 $D_{x,y}$ = dist. X et Y symétrique

- 1) choisir A, B qui minimisent $D_{A,B}$
- 2) ajouter un nouveau taxon qui va appeler (A,B)
 ↳ ajouter une rangée + colonne à D pour (A,B)
- 3) retirer les rangées A et B
- 4) pour chaque taxon C restant calculer $D_{(A,B),C}$

$$D_{(A,B),C} = \text{moyenne des dists entre les taxa (A,B) et C}$$

$$= \frac{1}{\# \text{dists}} \cdot \sum_{x \in (A,B)} \sum_{y \in C} D_{x,y}$$

$$= \frac{1}{|A \cup B| \cdot |C|} \cdot \sum_{x \in (A,B)} \sum_{y \in C} D_{x,y}$$



ex: D

	A	B	C	D	E	AUB	DUE	AUB UC
A	0	2	4	6	6			
B	2	0	4	6	6			
C			0	6	6	4	6	
D				0	6	6		
E					0	6		
AUB						0	6	
DUE							0	

T

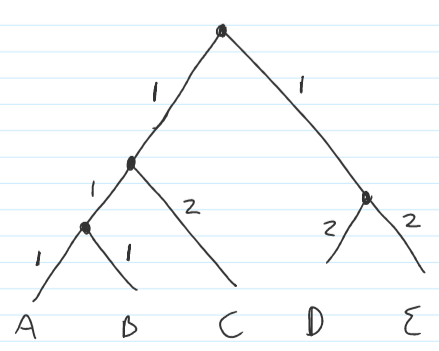
6+6+6+6
4

5) goto 1) jusqu'à ce qu'il reste 1 taxon

Q: peut-on étiquetter les branches de façon à ce que dans l'arbre T,

$$\text{dist}_T(X,Y) = D_{X,Y} \quad \text{dist}_T(X,Y) = \text{somme des poids des branches sur le chemin } X,Y \text{ dans } T$$

A	0	2	4	6	6
B	2	0	4	6	6
C	4	4	0	6	6
D	6	6	6	0	4
E	6	6	6	4	0



Un tel étiquetage fonctionne si la distance D (matrice) est ultramétrique.

• Une distance D est ultramétrique s'il existe un arbre

Un tel éligibilité fonctionne si la distance D (matrice) est ultramétrique.

- Une distance D est ultramétrique s'il existe un arbre T qui satisfait les conditions:
 - ① Pour tous taxa u, v , $D_{u,v} = \text{dist}_T(u, v)$
 - ② Soit r la racine de T . Alors \forall taxa u, v , $\text{dist}_T(r, u) = \text{dist}_T(r, v)$

Théorème: si D est ultramétrique, alors UPGMA reconstruit un arbre capable de satisfaire ① et ②

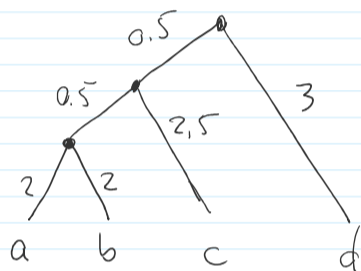
Dans les données réelles, les distances ne sont jamais ultramétrique.

Une autre approche, Neighbor-Joining, permet d'alléger ces conditions. Cette approche se concentre sur les distances additives.

- Une distance D est additive s'il existe un arbre T tel que
 - ① $\forall u, v$, $D_{u,v} = \text{dist}_T(u, v)$.

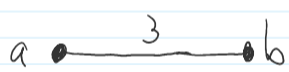
Q: comment vérifier qu'une distance est additive?

a	0	4	5	6
b		0	5	6
c			0	6
d				0

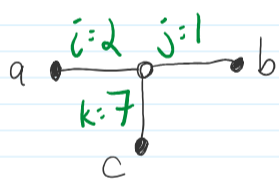


a	0	3	9	6
b		0	8	8
c			0	5
d				0

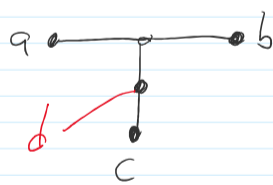
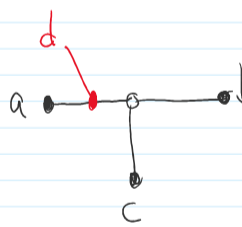
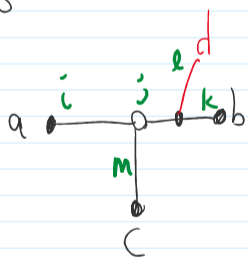
]? arbre additif? Cet arbre n'a pas besoin d'être enraciné.



$$\begin{aligned} i+j &= 3 \\ i+k &= 9 \\ j+k &= 8 \end{aligned}$$



$$\begin{aligned} i+j+l &= 6 \\ i+j+k &= 3 \\ &\vdots \end{aligned}$$

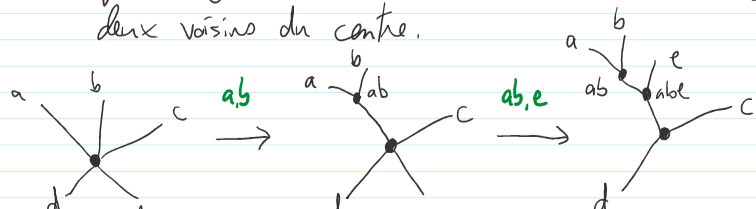
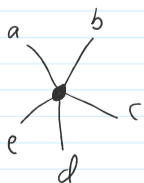


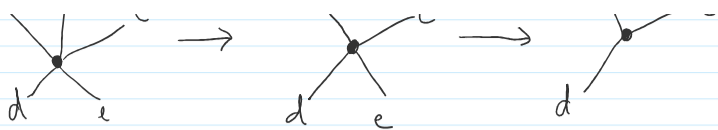
- Remarque: si une distance D est additive, alors il existe un seul arbre unique qui peut représenter les distances.

- L'algo Neighbor-Joining trouve un arbre additif s'il existe. S'il n'existe pas, l'algo approxime "bien" les distances.

Idee: démarrer d'un arbre étoile

- raffiner cet arbre une étape à la fois, en joignant successivement deux voisins du centre.





- On arrête quand le centre a 3 voisins.
- ↳ arbre non-enraciné.

Q: comment choisir les voisins du centre à joindre? (Complicé)

algoNJ(D)

soit T l'arbre étoile avec 1 feuille par tron

$c \leftarrow$ centre de T

tant que c a plus que 3 voisins

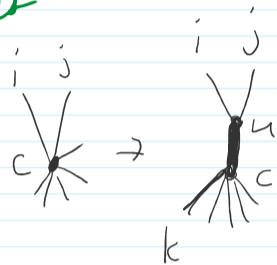
parmi les m voisins de c, choisir les deux voisins i, j qui minimisent

i, j proches $(m-2)D_{i,j} - (\sum_{k=1}^m D_{i,k} + \sum_{k=1}^m D_{j,k})$ *i et j bins des autres*

Joindre i, j en créant un nouveau nœud u

Mettre à jour $D_{u,k} \forall k \neq i, j$

$$D_{u,k} = \frac{1}{2} [D_{i,k} + D_{j,k} - D_{i,j}]$$



x

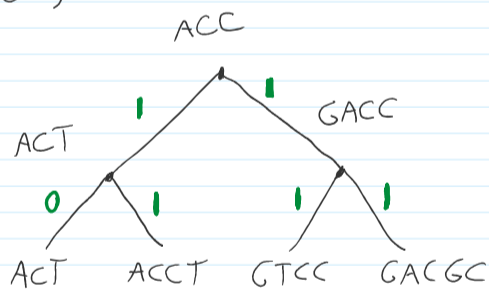
Approches par caractères

Entrée: séquences S_1, S_2, \dots, S_m

Sortie: arbre T avec

- feuilles étiquetées par les S_i
- tout nœud interne v soit étiqueté par une séquence $s(v)$
- minimiser la somme des coûts sur les branches (maximiser un score)

ex: S_1 ACT
 S_2 ACCT
 S_3 GTCC
 S_4 GACGC



recherche Arbre (S_1, S_2, \dots, S_m)

$\min T = \text{null}, \min \text{Cost} = \infty$

pour chaque arbre T avec les S_i aux feuilles

- trouver le meilleur étiquetage des nœuds internes de T, qui minimise $\text{cost}(T)$

- si $\text{cost}(T) < \min \text{Cost}$

$\min T = T$
 $\min \text{Cost} = \text{cost}(T)$

return T

Complexité: $\Omega((2n-3)!!) = \Omega((2n-3) \cdot (2n-5) \cdot (2n-7) \cdot \dots \cdot 5 \cdot 3 \cdot 1)$
au moins $= \Omega(\text{absurde})$

Problème de la petite parcimonie

- Parcimonie: l'hypothèse + simple = préférée
 = demande # min de changements

• Parcimonie: l'hypothèse + simple = préférée
 = demande # min de changements

• "Petite": sous-problème de la parcimonie générale
 général: T est inconnu
 petite: T est connu

Entrée: séquences S_1, S_2, \dots, S_m , arbre T avec S_i aux feuilles
 Sortie: étiquetage $s(v)$ des nœuds de T qui minimise

$$\sum_{uv \in E(T)} d(s(u), s(v))$$

arêtes

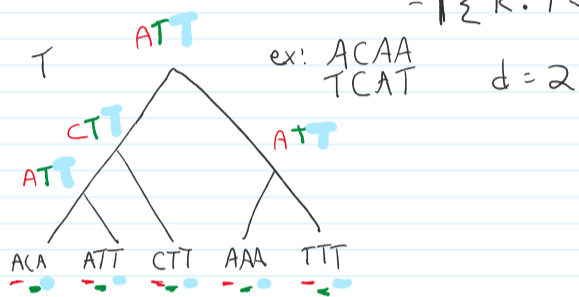
où d est une distance.

Version simple: d = distance de Hamming

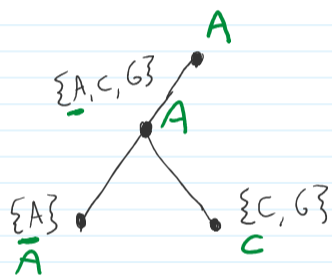
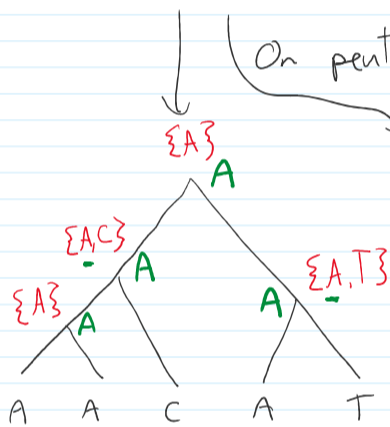
• S_1, S_2, \dots, S_m sont de la même longueur

• $d(S_i, S_j)$ = nombre de positions avec caractère différents

$$= |\{k : 1 \leq k \leq |S_i| \text{ et } S_i[k] \neq S_j[k]\}|$$



On peut traiter chaque pos. une à la fois



// si parent rien en commun, choix arbitraire

$C = \text{dict}()$ // $C[v]$ = caractères candidats au nœud v

getCandidats(v)

si v = feuille

| $C[v]$ = caractère assigné à v

sinon

| Soit v_1, v_2 les enfants de v

| getCandidats(v_1)

| getCandidats(v_2)

| si $C[v_1] \cap C[v_2] \neq \emptyset$

| | $C[v] = C[v_1] \cap C[v_2]$

| sinon

| | $C[v] = C[v_1] \cup C[v_2]$

x

getCandidat(racine)

// phase 2

getEtiquette(v)

si v = racine

| $s(v)$ = élément arbitraire de $C[v]$

sinon si v ≠ feuille

| soit p le parent de v

| si $\langle p \rangle \in C[v]$

Complexité:

- chaque nœud v est visité une fois dans l'exécution de cet algo

- temps pour traiter un nœud $O(|\Sigma|)$

Σ = alphabet pour calculer n, U

$\Rightarrow O(n|\Sigma|)$

n = nb de nœuds

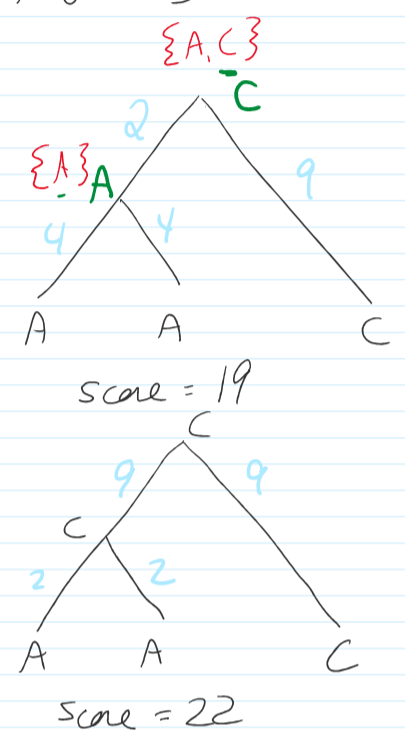
soit p le parent de v
 si $s(p) \in C[v]$
 $| s(v) = s(p)$
 sinon
 $| s(v) = \text{élem. arbitraire de } C[v]$

\times
 si ($v \neq$ feuille)
 $| \text{get étiquette}(v_1) \quad // v_1, v_2 \text{ enfants de } v$
 $| \text{get étiquette}(v_2)$

Version générale: matrice de coûts entre caractères
 $d(s_i, s_j) = \sum_{k=1}^{|s_i|} M[s_i[k], s_j[k]]$

ex: $M:$

	A	C	G	T
A	4	2	1	2
C	2	9	3	1
G	1	3	5	2
T	2	1	2	4

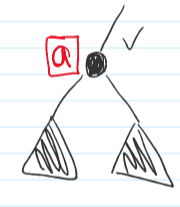


Si on veut maximiser la somme des scores sur les branches, ceci n'est pas un bon exemple

Avec un seul caractère, pour trouver un étiquetage qui maximise

$$\sum_{u \in E(T)} M[s(u), s(v)]$$

on utilise la prog. dynamique.



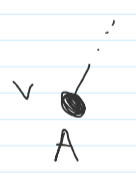
On définit, pour chaque nœud v et chaque $a \in \Sigma$,

$$V[v, a] = \text{score maximum du sous-arbre enraciné en } v \text{ avec la contrainte que } s(v) = a$$

On cherche, comme score optimal, l'entrée $\max_{a \in \Sigma} V[\text{racine}, a]$

Cas de base, $v =$ feuille

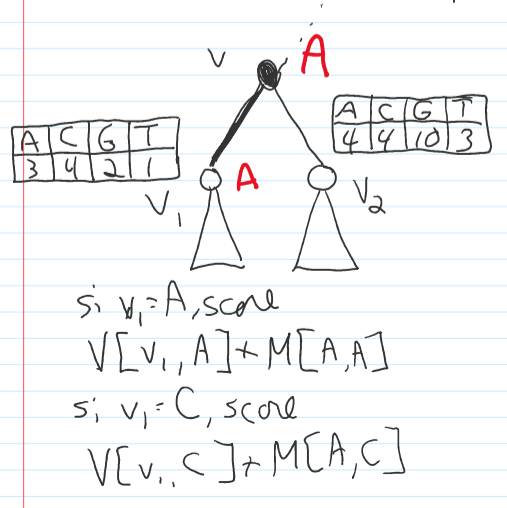
$$V[v, a] = \begin{cases} 0 & \text{si } a \text{ est le caractère assigné à la feuille} \\ -\infty & \text{sinon} \end{cases}$$



Récurrance

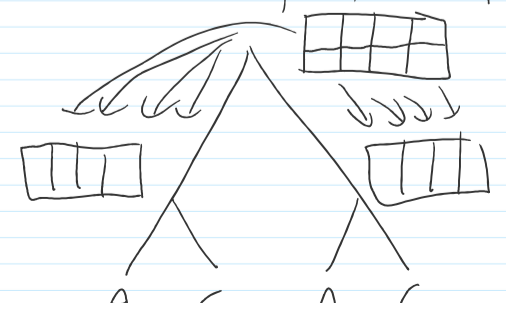
Soit v_1, v_2 les enfants de v

$$V[v, a] = \max_{b \in \Sigma} [V[v_1, b] + M[a, b]] + \max_{b \in \Sigma} [V[v_2, b] + M[a, b]]$$

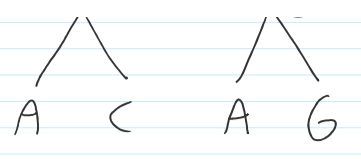


$M:$

	A	C	G	T
A	4	2	1	2
C	2	9	3	1
G	1	3	5	2
T	2	1	2	4



$V[v, c] + M[A, C]$



- Calculer toutes les entrées $V[v, a]$, il faut un temps "nb d'entrées dans V" x "temps pour calculer une entrée"
 "nb de nœuds" x "nb de symboles" x "nb de symboles"
 $= O(n \cdot |\Sigma|^2)$

- Algo de N, U : algo de Fitch
- Algo de prog. dyn. : algo de Sankoff - Rousseau

En pratique pour trouver le meilleur arbre T

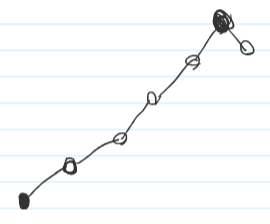
$T \leftarrow$ arbre Neighbor Joining (S_1, \dots, S_m)

fini = false
 tant que pas fini

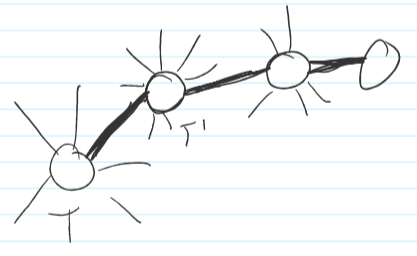
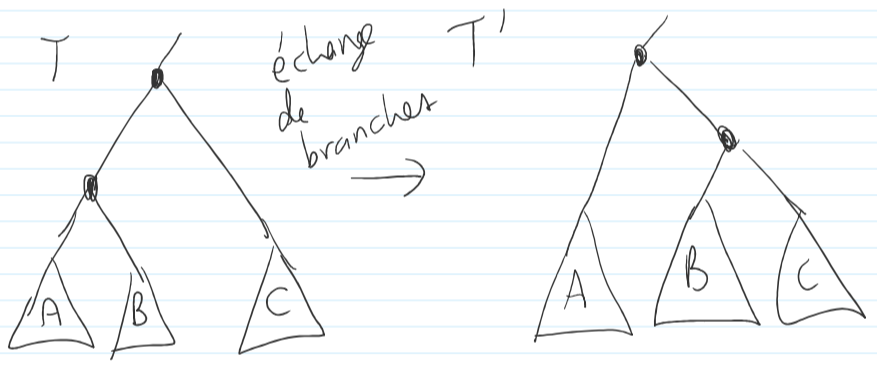
pour chaque voisin T' de T

si $cost(T') < cost(T)$
 $T \leftarrow T'$

si T n'a pas été mis à jour, fini = true



Voisin d'un arbre :



À retenir : • approches par dist. : UPGMA (en détail)
 Neighbor Joining (principe, ultramétrique vs additive)

• approches par caractères : algo N, U min. Hamming
 algo prog. dyn. max score M
 contexte général (trouver T)