

Mésoschallenge 2014: Biodiversiton

Y. Laizet, A. Moreau, J.-M. Frigerio, Ph. Chaumeil, P. Gay, P. Ramet,
D. Sherman, A. Franc

Equipe Pleiade, INRIA/INRA, Bordeaux

JSCIA, 24 mars 2015

- 1 Contexte de biodiversité : en quoi les objets biologiques sont-ils différents ?
- 2 Quels sont les patterns et motifs, locaux et globaux ?
- 3 Domaine de l'évolution et l'écologie des communautés
- 4 Travaux sur les arbres de la forêts guyanaise
- 5 Génomes comme empreinte moléculaire de l'évolution
- 6 Pas de processus fonctionnels

- 1 Taxonomie moléculaire
- 2 Chaque individu a un *attribut* : une séquence (un mot de 500 lettres, alphabet de 4 lettres $w \in \{A, T, C, G\}^n$)
- 3 Histoire inférée par des phylogénies moléculaire (Tree of life)
- 4 Modèles statistiques atteignent leurs limites si $10^4 \sim 10^3$ individus (ML, bayésien)
- 5 Travail sur des distances
- 6 Distance d'édition (Levenstein, 1965)

- 1 Taxonomie moléculaire
- 2 Chaque individu a un *attribut* : une séquence (un mot de 500 lettres, alphabet de 4 lettres $w \in \{A, T, C, G\}^n$)
- 3 Histoire inférée par des phylogénies moléculaire (Tree of life)
- 4 Modèles statistiques atteignent leurs limites si $10^4 \sim 10^3$ individus (ML, bayésien)
- 5 Travail sur des distances
- 6 Distance d'édition (Levenstein, 1965)

Question

Quelle est la forme d'un nuage de points ?

Distance geometry, machine learning, manifold learning, ...

Multidimensional Scaling : le problème

- ▶ Soit un ensemble $V = \{1, n\} \subset \mathbb{N}$ de n objets
- ▶ Les distances deux à deux sont donnés : $d(i, j) := d_{ij}$
- ▶ On se donne un entier $r < n$
- ▶ On cherche n points $x_i \in \mathbb{R}^r$ tels que

$$d_x(i, j) := \|x_i - x_j\| \simeq d_{ij} \quad (1)$$

Résultat

Il existe une solution exacte : méthode spectrale

voir Izenman, 2007, Modern Multivariate Statistical Techniques

Remarque

Il est possible de reconstruire $\langle x_i, x_j \rangle$ à partir des $d_x(i, j) = \|x_i - x_j\|$.

Algorithm 1 pseudocode for Multidimensional Scaling

- 1: compute $C_{ij} = -\frac{1}{2}(d_{ij}^2 - d_i^2 - d_j^2 + d_{..}^2)$, $1 \leq i, j \leq n$
 - 2: compute x_k : $Cx_k = \lambda_k x_k$
 - 3: $X = (x_k)_k$
 - 4: $L = \text{diag}(\sqrt{\lambda_k})_k$
 - 5: $Y = XL$
-

- ▶ MDS : minimisation de la valeur absolue des écarts
- ▶ NLM (Sammon, 1969) : poids faible sur les distances grandes : minimisation sur les valeurs relatives
- ▶ Fonction à minimiser :

$$\phi = \frac{1}{c} \sum_{i < j} \frac{[d_{ij} - d_x(i, j)]^2}{d_{ij}} \quad (2)$$

où c est une constante de normalisation

Nonlinear mapping : formalisation

- On se donne un tableau de distances

$$D = [d_{ij}]_{i,j=1\dots n} \quad (3)$$

et une fonction objectif

$$\phi = \sum_{i < j} \omega(d_{ij}) [d_{ij} - d_x(i, j)]^2 \quad (4)$$

Problème

<i>D</i>	<i>étant donné</i>	
<i>p</i>	<i>étant donné</i>	
<i>ω</i>	<i>étant donné</i>	(5)
<i>trouver</i>	$X = \{x_1, \dots, x_n\} \subset \mathbb{R}^p$	
<i>tel que</i>	ϕ minimal	

- ▶ On peut définir

$$\begin{aligned}\phi(X) &= 2 \sum_{i < j} \omega(d_{ij}) [d_{ij} - \|x_i - x_j\|]^2 \\ &= \sum_i \phi_i(x_i) \quad \phi_i(x) = \sum_{j \neq i} \omega(d_{ij}) [d_{ij} - \|x - x_j\|]^2\end{aligned}\tag{6}$$

- ▶ Mais ...ce diagramme n'est pas commutatif

$$\begin{array}{ccc}(\dots, x_i, \dots, x_j, \dots) & \xrightarrow{t+1} & (\dots, x_i, \dots, x'_j, \dots) \\ \downarrow t+1 & \searrow & \downarrow t+2 \\ (\dots, x'_i, \dots, x_j, \dots) & \xrightarrow{t+2} & (\dots, x''_i, \dots, x''_j, \dots)\end{array}$$

Algorithm 2 pseudocode for sequential optimisation

- 1: Sequential optimization
 - 2: **for** $t = 1$ to s **do**
 - 3: **for** $i = 1$ to n **do**
 - 4: select x_i
 - 5: compute $\nabla_i = \nabla\phi|_{x_i}$
 - 6: find α such that $\phi_i(x_i + \alpha\nabla_i)$ minimum
 - 7: update $x_i \leftarrow x_i + \alpha\nabla_i$
 - 8: **end for**
 - 9: **end for**
-

- remplacer

$$\phi_i(x) = \sum_{j \neq i} \omega(d_{ij}) [d_{ij} - \|x - x_j\|]^2 \quad (7)$$

par

$$\phi_i(x) = \sum_{0 < d_{ij} < \theta} \omega(d_{ij}) [d_{ij} - \|x - x_j\|]^2 \quad (8)$$

- Alors, si $d_{ij} > \theta$, le diagramme est commutatif, et (i, j) peuvent être optimisés en parallèle (si $d_x(i, j) > \theta$)

Un exemple de parallélisation MPI avec les mathématiques discrètes

- 1 partition : choisir quels points envoyer en parallèle ou non
- 2 distribution : allouer les points aux cœurs
- 3 équilibre des charges

- ▶ On définit un seuil ϵ et on en déduit θ tel que

$$d_{ij} > \theta \implies \omega(d_{ij})[d_{ij} - d_x(i,j)]^2 < \epsilon \quad (9)$$

- ▶ On construit deux graphes $G = (V, E)$ et $G' = (V, E')$ où

- 1 V est l'ensemble des points $V = \{1, n\}$
- 2 $i \sim j$ dès que $d_{ij} > \theta$ pour G , $d_{ij} < \theta$ pour G' ,

- ▶ Alors, si $i \sim j$ dans E , on peut envoyer i et j en parallèle sur deux cœurs différents. De même un ensemble formé de un point par composante connexe de G' .

Résultat

S'il existe une clique de m points dans G , alors, on peut envoyer m points en parallèle sur m cœurs. De même s'il existe m composantes connexes dans G'

Partition par coloriage

- ▶ Cette partition est non optimale (cliques trop restrictives, et G' peut être connexe par linkage avec des points très éloignés).

Remarque

Il existe un lien avec le coloriage de graphes : trouver le nombre minimal de couleurs telles que deux sommets reliés par un lien sont de couleurs différentes

- ▶ Si deux points sont de même couleur, alors ils peuvent être envoyés en parallèle

Procédure

- 1 *Identifier les composantes connexes de G'*
- 2 *Colorier chaque composante connexe*
- 3 *Envoyer en parallèle une couleur complète de chaque composante*

Definition

Une *coloration* d'un graphe est une partition des sommets en ensembles deux à deux non adjacents. Le *nombre chromatique d'un graphe* est le nombre minimal d'éléments d'une coloration. Il est noté en général $\chi(G)$.

Remarque

Trouver une coloration d'un graphe est un problème NP-complet.

► Il existe en revanche des algorithmes approchés quadratiques. Un algorithme glouton permet de trouver une coloration avec $d + 1$ composantes si d est le degré maximal dans G .

Algorithm 3 pseudocode for parallel optimisation

- 1: Have $D = (d_{ij})_{i,j}$
- 2: Build $G = (V, E) : i \sim j \implies d(i, j) < \theta$
- 3: Compute the connected components $= (C_1, \dots, C_m)$
- 4: Compute a minimal coloration $C_k = \bigsqcup_{j \leq K} K_{kj}$
- 5: set $X = (x_1^0, \dots, x_n^0)$ (Initialize by *MDS*)
- 6: **for** $t = 1$ to s **do**
- 7: distrib = \emptyset
- 8: **for** $k = 1$ to K **do**
- 9: **for** $a = 1$ to m **do**
- 10: $X_{ak} = \{x_i \in C_a : \kappa(x_i) = k\}$
- 11: distrib = distrib $\cup X_{ak}$
- 12: **end for**
- 13: **end for**
- 14: mpiexec -np N optimize [param]
- 15: **end for**

► Ce qui a été fait

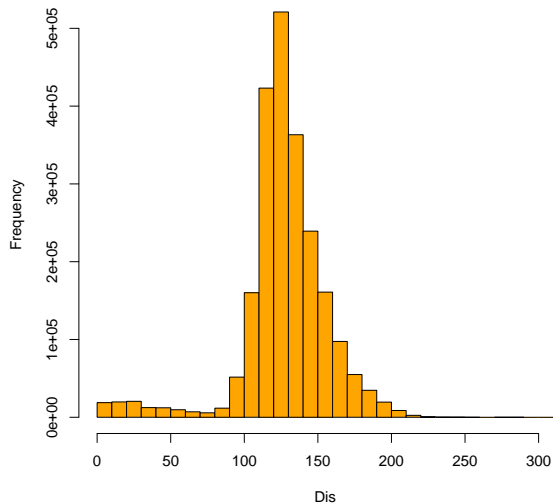
- 1 Poser le problème
- 2 Résoudre l'optimisation pour un point
- 3 Algorithme d'optimisation en séquentiel (moindres carrés alternés)
- 4 Paralléliser le code sur les points
 - 1 partition des points selon une coloration
 - 2 construire une distribution
 - 3 écrire un code en MPI-python

► Ce qui reste à faire

- 1 Equilibrer la distribution
- 2 Accélérer l'optimisation pour un point
- 3 Calculer la partition à chaque étape sur $d_x(i, j)$ et non une fois pour toute sur $d(i, j)$
- 4 Utiliser la hiérarchie cœurs / nœuds pour écrire une partie en open-MP (mémoire partagée)

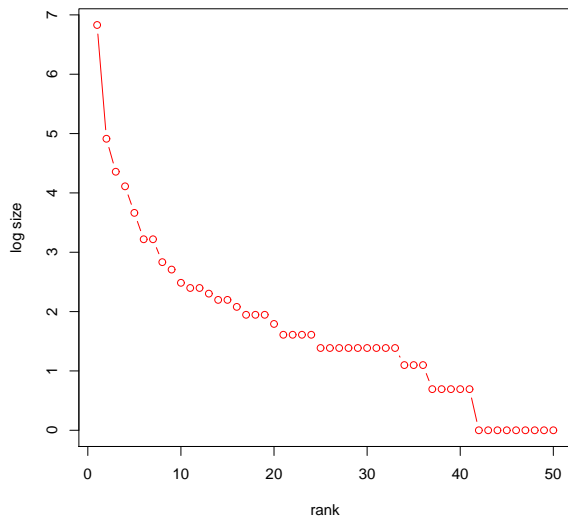
Histogramme des distances

Pairwise distances



Taille des composantes connexes

50 connected components

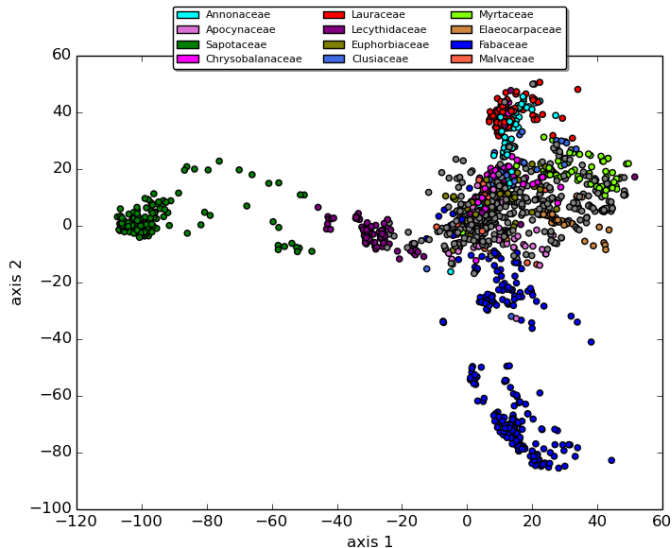


Apport : développement

Mise au point du code plus que calcul intensif

- 1 Pas de calcul très intensif actuellement (reste du niveau de la queue habituelle Tier 2 en nombre de CPU et *wall clock*)
- 2 Mais accès prioritaire d'une partie de la machine pour la mise au point des codes
- 3 Perspective de passage à l'échelle sur Tier 1 (plus grands jeux de données)

Taxonomie : MDS sur 1500 arbres de Guyane



Taxonomie : NLM sur 1500 arbres de Guyane

