

Repliement protéique simplifié via un algorithme de Monte Carlo.

Basé sur l'article : Thachuk C, Shmygelska A, Hoos HH. *A replica exchange Monte Carlo algorithm for protein folding in the HP model. BMC Bioinformatics. 2007 Sep 17;8:342. PubMed PMID: 17875212; PubMed Central PMCID: PMC2071922.*

Introduction:

Les progrès technologiques ont permis ces dernières années d'obtenir rapidement et en masses des données de séquençages protéique. Cependant, le séquençage ne produit que la séquences d'acides aminés qui constitue la protéine d'intérêt, mais ne donne à priori aucune information sur sa structures secondaire, tertiaire et quaternaire. Un des objectifs lié à l'arrivage massif de données est de pouvoir prédire les repliements protéiques rapidement sans avoir à effectué des analyses en laboratoire comme des cristallographies au rayon X ou des imageries à résonance magnétiques. Pour ce faire, il est possible de créer des modèles de repliement aléatoire afin de tenter d'extrapoler au mieux un repliement protéique. Depuis plusieurs années déjà, différentes équipes ont tenté de créer des modèles capable de prédire des repliement de protéines. Parmi les différents travaux effectués, ceux de Thachuk et son équipe sont basé sur un algorithme de Monte Carlo.

Mon travail ici, a constitue en une inspiration de l'étude de Thachuk *et al.* (2007), afin d'implémenter un modèle simple de repliement protéique avec des algorithme de Monte Carlo.

Matériel et méthodes:

Tout les fichiers et documents sont disponible sur GitHub à cette adresse (https://github.com/nicolassilva/MonteCarloAlgorithm_Search). Sur le git, est également disponible un fichier .yml avec lequel la création d'un environnement conda identique à celui dans lequel le projet a été réalisé est possible.

L'algorithme créer ici s'effectuera avec un modèle de séquence "Polaire/Hydrophobe" (Dill 1985). Ce modèle classifie les acides aminés qui constituent une protéine en deux catégorie : les acides aminés polaires et les acides aminés hydrophobe. La séquence protéiques est donc modélisé par un enchainement d'acides aminés H et P. Pour ce modèle, une séquence aléatoire a été utilisée et est disponible dans le répertoire *data* du git. La séquence protéique est ensuite projetée en ligne droite sur une grille en 2D.

L'algorithme de Monte Carlo implémenté ici consiste en la prise aléatoire d'un avide aminé de la séquence et à lui faire effectuer un mouvement tout aussi aléatoire

Annexes

Structures du programme:

Le code python implémenté pour l'algorithme de Recherche de Monte Carlo à été construit de la manière suivante:

- Une classe **File** avec une méthode *readFile* permettant de lire le fichier contenant la chaîne d'acides aminés.
- Une classe **Array** avec trois méthodes. La première *createArray* qui permet de créer une matrice de deux fois la taille de la séquence étudié ; la seconde *energy* qui permet de calculer l'énergie d'une conformation ; et la troisième *writeFile* qui permet d'écrire dans un fichier la matrice avec la conformation finale obtenue.
- La séquence est importé et une matrice est créée avec la séquence au centre de celle-ci. Les coordonnées des acides aminés dans la matrice sont alors stockés dans un dictionnaire, et l'énergie de la conformation est calculée.
- Commence ensuite l'algorithme de Monte Carlo. Un acide aminé est tiré au hasard dans la séquence, ainsi qu'un des 8 mouvements modélisés. Si la case d'arrivée du dit acide aminé est disponible (c'est-à-dire pas déjà occupée par un autre acide aminé) alors le déplacement est autorisé et les autres acides aminés suivant déplacés aussi.
- Une nouvelle matrice avec la nouvelle séquence est créée et l'énergie associée à cette nouvelle conformation est calculée. Si la nouvelle énergie est inférieure ou égale à l'énergie de l'ancienne conformation, alors la nouvelle conformation est conservée. Sinon le déplacement est annulé.
- L'algorithme de recherche de Monte Carlo recommence ensuite autant de fois que le *nstep*.
- A la fin du script, la matrice ainsi que l'énergie de la dernière conformation conservé est affichée. Et la matrice est sauvegarder dans un fichier *.txt*.

Exemple d'utilisation du programme:

Dans un premier temps il faut obtenir la chaîne protéique d'intérêt sous forme d'enchaînement de 'h' et de 'p'. Il doit s'agir d'un fichier *.txt*.
Il faut ensuite télécharger le fichier *monteCarloSearch.yml* de mon git et créer l'environnement conda associé. De cet environnement il faut alors exécuter le script *python monteCarlo_Search.py* en passant comme argument le nom du fichier et le nombre d'itération (*nstep*). Un fichier *results.txt* contenant la matrice avec la conformation finale est ensuite créé dans un dossier *results*.

Exemple : `python3 monteCarlo_Search.py sequence.txt 500`