

Régression logistique multiple et prédiction des facteurs de risque concernant la survie en soins intensifs.

changer le titre

Julia Guerra ¹, Maxime Jaunatre ², Ellie Tideswell ³ | Master 2 BEE Grenoble
Mail ¹, Mail ² Mail ³ | 29 novembre 2019

Todo list

changer le titre	1
jeu de donnée	2
algo d'apprentissage	2
mettre exemple de matrice	2
a voir si on choisit bien ce modèle ou non	3
fig sensi	3
figure sensispeci	3
Parametrages de viterbi - transition entre modèles + et -	3
fig smoothing	4

L'avancée des techniques de séquençage a permis d'obtenir de grandes quantités d'informations génomiques. Durant ces dernières années, la génétique a embrassé les outils mathématiques de modélisation, parvenant à une caractérisation statistique des données de séquençage. Cette approche a permis de décrire avec précision les motifs observés dans des régions étudiées, et de prédire leurs présences dans des séquences encore inconnues (Wu *et al.*, 2010). Une des méthodes les plus connues dans ce domaine est l'utilisation des chaînes de Markov, introduites premièrement par Churchill (1992) pour l'analyse de séquences génomiques puis par Durbin *et al.* (1998) pour la détection de régions CGI.

Dans le génome des deutérostomiens, la fréquence du dinucléotide C-G est moins importante qu'attendu sous une distribution aléatoire indépendante des quatre bases azotées. Ceci est une conséquence des mécanismes de protection contre la mutation spontanée du génome. Cependant, dans certaines régions de l'ADN nommées îlots CpG (ou CGI) ce processus de mutation est évolutivement réprimé et donc la fréquence des dinucléotides C-G est donc plus élevée; par exemple aux alentours de

certain promoteurs (Haque *et al.*, 2011, Saxonov *et al.*, 2006, Wu *et al.*, 2010).

La grande variabilité dans la taille, la composition et l'emplacement de ces CGI rend difficile leurs définitions et donc l'établissement d'un algorithme unique permettant leurs détections indubitables (Wu *et al.*, 2010). Ainsi, les modèles de Markov permettent de modéliser les fréquences des nucléotides en fonction de séquences déjà connues; ces séquences contenant ou non des CGI. En supplément des chaînes de Markov simples, il existe aussi les chaînes de Markov cachées (HMM, Churchill (1992)) : ces dernières décrivent de nombreux processus réels qui suivent un modèle de Markov, mais qui ne sont pas observables. Une chaîne de Markov cachée permettrait donc l'utilisation d'un seul modèle pour identifier un nucléotide (l'observation) et si ce dernier est à l'intérieur d'un îlot CpG ou non (aussi appelé l'état de la région). Les HMM permettent ainsi d'augmenter la résolution de l'analyse, c'est-à-dire de détecter l'emplacement des régions CGI à l'intérieur des séquences.

Matériel et méthodes

0.1 Modèles de Markov simples

jeu de donnée

Les modèles de Markov réalisés dans cette étude ont été construits à partir de deux jeux de séquences de souris (*Mus musculus*). Ces jeux de séquences avaient été caractérisés en avance comme contenant des îlots CpG (on notera “CpG+”), ou ne contenant pas d’îlots CpG (“CpG-”). Les deux jeux de séquences “app”, pour la construction des modèles CpG+ et CpG-, contenaient 1160 et 5755 séquences respectivement. Des jeux supplémentaires “test” également caractérisés comme CpG+ ou CpG- (1163 et 5137 séquences) ont servi à évaluer la performance des modèles.

algo d'apprentissage

Dans un premier temps, les fréquences relatives d’observation des bases A, C, G, T ont été calculées pour la totalité des séquences de chaque jeu de données “app” (R, fonction `count` du package `seqinr`; Charif & Lobry (2007)). Ces données ont permis de construire la matrices de probabilité A du modèle d’ordre 0 (M0). Le terme “ordre” fait référence au nombre de bases précédentes conditionnant la probabilité de présence de la base étudiée. De cette manière, le M0 considère la probabilité d’occurrence de chaque base comme une variable aléatoire (équation 1) dont les probabilités d’occurrence (équation 2) sont différentes. En plus, des résultats différents sont attendus en fonction de la nature CpG+ ou CpG- des séquences; c’est pourquoi deux matrices de probabilité A+ et A- ont été construites, provenant respectivement des comptages du jeu CpG+ et CpG-.

Le modèle de Markov d’ordre 1 (M1) rassemble les occurrences de chaque base en fonction de la base pré-

cédente. Les matrices de transition q1+ et q1- sont matrices 4x4 qui ont été donc construites à partir des comptages de chaque couple de bases. Vu qu’elle ne peut pas dépendre d’une base précédente, la probabilité d’occurrence de chaque base initiale a été considérée comme une variable aléatoire (équation 3) à probabilités équivalentes (équation 4). Ce protocole de construction de modèle a été refait pour l’ordre 2, obtenant une matrice 16x4. Pareil pour l’ordre 3 (matrice 64 x 4), l’ordre 4 ... jusqu’à l’ordre 5. Pour les modèles d’ordre supérieur 0, les lignes de la matrice ont été rangées de sorte que la somme de chaque ligne soit égal à 1, à cause de la nature conditionnelle des probabilités, comme dans l’exemple suivant (table 1).

$$Y \in B; B = a, c, g, t \quad (1)$$

$$P(Y_i = k) \forall k \in B \quad (2)$$

$$X \in B; B = a, c, g, t \quad (3)$$

$$P(A) = P(C) = P(G) = P(T) = P(X \in B) = \frac{1}{4} \quad (4)$$

mettre exemple de matrice

A partir des matrices de transition, on peut calculer la log-Vraisemblance d’une séquence sous un modèle MX correspondant comme la somme du log de la probabilité de premières bases (région de taille égale à l’ordre) avec la somme du produit de la matrice de transition par la matrice d’occurrence des mots dans la séquence (voir équation 5).

	a	c	g	t
a	0.29	0.21	0.30	0.20
c	0.26	0.30	0.17	0.27
g	0.24	0.27	0.30	0.20
t	0.18	0.26	0.28	0.29

TABLE 1 – Matrice de transition du modèle CpG+ d’ordre 1

$$P_+(Sequence) = \log[P_{initial}(mot_{initial})] + \sum_{i=1}^{n=4^{ordre}+1} \log[P_i(mot_i) \cdot N_i(mot_i)] \quad (5)$$

Choix du meilleur modele

La performance du M1 a été testée sur les deux jeux de séquences de test. La log-vraisemblance de chaque séquence a été calculé pour chaque modèle (CpG+ et

CpG-) et la séquence est donc associée à l’état pour lequel la log-vraisemblance est la plus grande. Pour le jeu de données CpG+, les séquences caractérisées comme CpG+ sont considérées comme vrais positifs (VP) et les séquences caractérisées CpG- comme faux négatifs

(FN). Pour le jeu de données CpG-, les séquences caractérisées comme CpG+ sont considérées comme faux positifs (FP) et celles caractérisées comme CpG- comme vrais négatifs (VN). Le même protocole a été suivi pour tester la performance des modèles 1 à 6.

La spécificité et la sensibilité de chaque modèle ont été calculées à partir de ces résultats, selon les équations illustrées en 6 et 7.

$$Sensitivity = \frac{VP}{VP + FN} \quad (6)$$

$$Specificity = \frac{VN}{VN + FP} \quad (7)$$

Ce processus, répété pour toutes les combinaisons de Mi+/Mj- (avec i et j allant de 0 à 5), a permis de connaître la meilleure combinaison de modèle. Pour l'obtenir, les données de sensibilité et spécificité pour les modèles ont été sommées entre elles. La combinaison d'ordres portant la valeur maximale étant la valeur (5,4) de la matrice

a voir si on choisit bien ce modele ou non

; les calculs de la chaîne de Markov cachée ont été réalisés sur un modèle d'ordre 5 pour les séquences CpG+ et un modèle d'ordre 4 pour les séquences CpG-.

fig sensi

figure sensispeci

0.2 Modèles de Markov cachés

Afin d'augmenter la résolution de l'analyse des séquences et donc de détecter plutôt l'emplacement des régions CGI (îlots CpG+) que la catégorie de toute la séquence les modèles de Markov cachés ont été utilisés. Les HMC permettent d'inclure les effets d'un processus sous-jacent au processus principal observé. Dans notre cas, la chaîne cachée sert à modéliser les transitions entre régions CGI et régions pas CGI à l'intérieur des séquences, ce qui n'est pas directement observable; le processus principale MM centré toujours sur les observations de bases à chaque position. HERE : mettre en ordre / propre la description des calculs HMC Par rapport aux modèles de Markov simples, un troisième paramètre doit donc être calculé (ainsi que les proba-

bilités de transition entre états et les distribution initiales), ceci étant les probabilités d'émission des observations. Ces probabilités sont calculées à partir des valeurs moyennes des longueurs des îlots CpG, et des régions entre ces îlots. Les trois matrices (contenant des probabilités initiales, de transition et d'émission) ont la particularité d'avoir les lignes stochastiques. Ceci signifie que la somme des éléments de la ligne est égale à 1.

0.3 L'algorithme de Viterbi

Afin de trouver la séquence optimale d'états qui correspond à une séquence donnée d'observations, il est possible d'utiliser une fenêtre glissante (un algorithme naïf), dans laquelle les log vraisemblances sont calculés pour des segments de bases d'une longueur donnée. Bien que facile à implémenter, les résultats (en terme des prédictions des CGI) dépendent de la taille de la fenêtre choisi, ceci peut représenter un biais de cette méthode. Une façon alternative peut être l'algorithme de Viterbi, un exemple de la programmation dynamique, qui permet d'identifier la séquence qui maximise la probabilité de générer les observations (Pardoux). Le chemin le plus probable étant donné un modèle est déterminé via une procédure récursive. L'algorithme de Viterbi est décrit comme suit :

$P(\text{next observation}|\text{next state}) \max \text{ current state}$
 $P(\text{next state}|\text{current state}) \max \text{ previous states}$
 $P(\text{previous states, current state}|\text{observations so far})$

ou : $P(O_{t+1}|S_{t+1}) \max_{st} P(S_{t+1}|st) \max_{s1 \dots st-1} P(s1 \dots st-1, st|O1 \dots Ot)$

Parametrages de viterbi - transition entre modèles + et -

0.4 Smoothing

La technique de "Smoothing" représente une technique mathématique qui enlève la variabilité parmi les données, impliquant souvent la redistribution du poids entre des régions de haute probabilité, et des régions de "zéro probabilité" (Boodidi 2007). Dans le cadre de cette étude, le "Smoothing" revient donc à lisser la caractérisation des différentes régions en les ré-assignant selon 2

procédés successifs. Dans un premier temps, les régions de longueur inférieur à un certain seuil (S) sont assignée à une nouvelle catégorie “Ambiguous”, en vert dans la figure Cette première étape comporte également un algorithme qui compile ces nouvelles régions en une seule quand elles se suivent dans la séquence (voir bases 9 à 13), afin de mesurer la longueur de cette nouvelle région dont la catégorie est devenue unique.

Le second procédé vérifie la longueur de ces nouvelles régions ambiguës et leurs situations sur la séquence. En effet, il arrive qu’une région de petite taille soit considérée comme ambiguë entre deux régions d’une même catégorie (voir base 5). On peut donc supposer qu’il s’agit de bruit et que cette région est probablement de la même catégorie que celles qui l’entoure. Ainsi, le second procédé de smoothing va ré-assigner des régions ambiguës si leurs tailles sont inférieures à un seuil et que les régions bordantes sont de même nature. On note que l’algorithme de ‘Smoothing’ ne peut être utilisé qu’avec le second procédé, car en l’absence de régions ambiguës aucune ré-assignation vers CpG+ ou CpG- n’est possible.

Résultats

mus1

table tronquée figure tronquée description du chromosome (nombre d’îlots cpg, taille des cpg, fenetre de smooth)

mus2

description

mus3

description

Discussion

nos résultats sont badass moduler le choix de smoothing améliorations des algos (coder en autre langage), parallélisation du choix de meilleur modele

Ressources

Ce document est disponible en ligne sous format “.Rnw”, contenant tout le code nécessaire à la reproduction de l’analyse, réalisée avec un script en langage R (R.Team, 2017), ainsi que le jeu de données de départ. L’ensemble est situé sur Github : <https://github.com/gowachin/BeeMarkov> et peut être installé sur R via les commandes suivantes.

```
> # NOT RUN
> library(devtools)
> install_github("gowachin/BeeMarkov")
> library(BeeMarkov)
```

Bibliographie

- Charif, Delphine, & Lobry, Jean R. 2007. SeqinR 1.0-2 : A Contributed Package to the R Project for Statistical Computing Devoted to Biological Sequences Retrieval and Analysis.
- Churchill, Gary A. 1992. Hidden Markov chains and the analysis of genome structure. *Computers and Chemistry*, **16**(2), 107–115.
- Durbin, Richard, Eddy, Sean R., Krogh, Anders, & Mitchison, Graeme. 1998. Biological sequence analysis. *Biological sequence analysis*.
- Haque, A. N.A., Hossain, M. E., Haque, M. E., Hasan, M. M., Malek, M. A., Rafii, M. Y., & Shamsuzzaman, S. M. 2011. CpG islands and the regulation of transcription. *GENES & DEVELOPMENT*, **25**(1), 1010–1022.
- R.Team. 2017. R : A language and environment for statistical computing (Version 3.4. 2)[Computer software]. *Vienna, Austria : R Foundation for Statistical Computing*.
- Saxonov, Serge, Berg, Paul, & Brutlag, Douglas L. 2006. A genome-wide analysis of CpG dinucleotides in the human genome distinguishes two distinct classes of promoters. *Proceedings of the National Academy of Sciences of the United States of America*, **103**(5), 1412–1417.

Wu, Hao, Caffo, Brian, Jaffee, Harris A., Irizarry, Rafael A., & Feinberg, Andrew P. 2010. Redefining CpG

islands using hidden Markov models. *Biostatistics*, **11**(3), 499–514.

Annexes

mus1

table figure

mus2

table figure

mus3

table figure

scripts