

Algorithmes d'apprentissage pour les chaînes de Markov cachées

Alexis Jacq

July 5, 2015

1 Un peu de philo...

Pour proposer une explication d'un phénomène, dans le but de le prédire ou bien de l'exploiter, l'Homme a recours à des modèles simplifiés de ce phénomène. Ces modèles consistent à projeter le phénomène dans un formalisme accessible afin de le rendre compréhensible et descriptible. Ainsi, les représentations cognitives du monde qui se forment dans nos cerveaux, nos langages, nos religions et enfin nos modèles scientifiques sont autant de modèles abordés pour s'entendre sur une compréhension – d'abord personnelle puis commune – de notre environnement.

L'apprentissage consiste à adapter ou bien à affiner les modèles choisis afin d'optimiser leur capacité à décrire les phénomènes observés. Toujours dans les mêmes exemples : au niveau du cerveau, l'apprentissage neuronal nous permet de cerner des concepts toujours plus pertinents et de mémoriser leur dynamique. Au niveau des langages, nous nous efforçons d'acquiescer, par les échanges verbaux, un vocabulaire toujours plus adapté à nos sociétés et nos mœurs, puis l'on suppose ce vocabulaire assez efficace pour l'utiliser et ainsi le tester. Dans les religions, nous cherchons à "grandir dans la foi", ce qui consiste en un perpétuel effort introspectif d'adaptation aux convictions spirituelles et d'acceptation des idéologies. Enfin, en science, la démarche est plus simple à décrire : tablier et supposer des modèles mécanistiques puis affiner leurs paramètres pour les adapter à la description du phénomène étudié. L'apprentissage se fait donc en deux étapes : l'**acceptation** du modèle, puis son **adaptation**.

Ces deux étapes ne sont absolument pas disjointes : en adaptant le modèle d'une certaine façon, on le modifie, et on en accepte un nouveau, etc... L'apprentissage peut donc être vu comme un permanent saut-mouton entre acceptation et adaptation : on accepte le modèle pour établir une prédiction, on adapte le modèle selon l'écart entre ce qui a été prédit et ce qui est advenu, puis on accepte le modèle adapté et on recommence etc... Il est important de voir que tous les algorithmes ou autres démarches visant un apprentissage reprennent ce schéma d'acceptation/adaptation.

2 Un peu de formalisme...

Nous nous plaçons ici dans le cadre d'un formalisme mathématique. Comme dans la science prétendante [?], nous nous intéressons à une *classe* de modèles très simple et très maniables : tant donné le phénomène

observ $X \rightarrow Y^{\text{obs}}$ o l'on cherche expliquer la variable d'intrt Y^{obs} avec la variable X , on tablit le modle suivant :

$$\mathcal{M}_P : X \mapsto Y$$

o \mathcal{M}_P dsigne le modle (dpendent du jeu de paramtre P), vu comme une application qui X (variable connue) associe Y (la prdiction). Ici, adapter le modle, c'est optimiser les paramtres P dans le but de minimiser l'erreur de prdiction, savoir l'"cart" entre Y prdit et Y^{obs} observ. Comme nous l'avons vu prcdemment, la dfinition de cet cart va beaucoup dpendre du contexte. Dans cette sance, nous allons nous intrresser un contexte stochastique. La mesure d'erreur de prdiction sera donc dfinie par la vraisemblance. Cette quantit a dj t introduite dans la sance prcdente ([?] section 2.3.1). Rappelons juste qu'il s'agit de la capacit du modle *expliquer* l'observation Y^{obs} . On la note \mathcal{L} (pour *likelihood* = vraisemblance) :

$$\mathcal{L}(\mathcal{M}_P, X, Y^{\text{obs}}) = \mathbb{P}[Y^{\text{obs}} | \mathcal{M}_P, X]$$

3 ...Maintenant on peut y aller

Dans cette sance, nous tudions un modle trs particulier : les chaînes de Markov caches (HMM pour Hidden Markov Model). Mais avant tout : qu'est-ce qu'une chane de Markov ?

3.1 Approche non rigoureuse

Prenons un processus purement alatoire, par exemple un "pile ou face". Pourquoi un *processus* ? Parce qu'on va faire plusieurs jets de suite, ce qui apporte une dimension temporelle. De ce fait, on va pouvoir les nommer par indice, par ex j_1 le premier jet, j_2 le deuxime etc... A chaque fois j_i vaut pile ou face selon le rsultat obtenu. Pourquoi *purement* alatoire ? En fait, un jet de pice n'est absolument pas alatoire : il dpend de la force et du point d'application du coup de pouce, de la position initiale dans l'espace, de la topographie du sol, peut-être même de la temprature... Mais notre chelle, tant incapable d'tablir ces liens de causalit, on prfre modliser le jet de pice comme un vnement ne dpendant d'absolument rien. Rien ne permet de prdire si le rsultat sera pile ou face.

Maintenant, regardons un autre processus. On continue jeter successivement des pices. Sauf que cette fois-ci, on compte un score. Supposons que l'on joue pile : si le rsultat d'un lancer donne pile, on gagne un point, sinon on perd un point. Indions l'tat de ce score pour tudier son avancement : s_1 pour le score aprs le premier jet, s_2 aprs le second etc... Ainsi, s_n reprsente le score au bout de n jets :

$$\begin{aligned} s_1 &= \mathbf{1}_{[j_1=\text{pile}]} - \mathbf{1}_{[j_1=\text{face}]} \\ s_2 &= s_1 + \mathbf{1}_{[j_2=\text{pile}]} - \mathbf{1}_{[j_2=\text{face}]} \\ &\dots \\ s_n &= s_{n-1} + \mathbf{1}_{[j_n=\text{pile}]} - \mathbf{1}_{[j_n=\text{face}]} \end{aligned}$$

Cette fois-ci, on voit que s_n dpend de s_{n-1} qui est connu l'instant n . En effet si l'instant $n-1$, le score est $s_{n-1} = 10$ on sait que s_n vaudra 11 ou 9 selon j_n . En outre, sans la connaissance de s_{n-1} on ne pourrai rien prdire de plus que " s_n sera compris entre $-n$ et $+n$ (inclus) et aura la mme

parit que n ".

Remarquons autre chose : si on connaît s_{n-1} , a ne sert rien de connaître s_{n-2} pour obtenir des informations propos de s_n . En effet, les informations apportées par s_{n-2} sont incluses dans les informations apportées par s_{n-1} . Par ex, si on connaît $s_{n-2} = 9$ on sait que s_n vaudra 7, 9, ou 11, tandis que si on connaît $s_{n-1} = 10$, on sait bien que $s_n \in \{7, 9, 11\}$ mais plus encore : $s_n \in \{9, 11\}$. On appelle cela "proprit d'oubli" ou bien "proprit de Markov". C'est dire que si on connaît s_{n-1} et s_{n-2} , on peut *oublier* s_{n-2} qui ne sert plus rien.

Plus formellement, on dit que la **loi** (c'est dire la nature du modle stochastique) de s_n **sachant** $s_{n-1}, s_{n-2}, \dots, s_2, s_1$ est la même que celle de s_n **sachant** seulement s_{n-1} .

Il est trs frquent de rencontrer des modles stochastiques vrifiant cette propriit. Le premier exemple, et sans doute le plus simple, est celui des **marches alatoires**. Il faut imaginer un homme complètement soûl errant dans une ville. Tantt il fait un pas vers la gauche, puis vers l'avant, puis un autre vers la gauche suivit d'un pas en arrire. chaque instant (les instants ici representent les pas successifs) la position courante seule donne autant d'information sur la position future que l'integralit de l'errance depuis la sortie du pub. En fait, les marches alatoires se gnralisent comme des accumulations additives de ralisation d'une variable alatoire. On peut donc les dfinir par rcurrence comme suit :

$$z_n = z_{n-1} + x_n$$

o z_n est la n -me position de la marche, et x_n la n -me ralisation d'une variable alatoire. Vous l'avez sans doute devin, le processus dfini par les scores de jets de pice est donc une marche alatoire. Mais attention : toutes les chaînes de Markov ne sont pas des marches alatoires ! En effet, on pourrait trs bien dfinir la chaîne suivante (cette rcurrence s'appelle le *critre fondamental des chaînes de Markov*, toute chaîne de Markov peut s'crire sous cette forme et n'importe quelle suite vrifiant cette propriit dfinie une chaîne de Markov):

$$z_n = f(z_{n-1}, x_n, n)$$

o f est mesurable de $\mathbb{R}^2 \times \mathbb{N}$ dans \mathbb{R} , mais a (les mesures et tout) on le verra plus loin, si a vous intresse. Pour l'instant, nous allons nous contenter de la dfinition suivante (je vous rassure, elle est on ne peut plus rigoureuse et suffisante) :

3.2 Approche rigoureuse des chaînes de Markov

Je vais copier-coller (Ha, si seulement on pouvait réellement copier-coller du LaTeX !) l'introduction aux chaînes de Markov du Pr. Jean Jacod que l'on retrouve en début de ses poly de cours Jussieu ([?],[?]) :

L'idée des chaînes de Markov est très simple : il s'agit d'une suite $(X_n)_{n \in \mathbb{N}}$ de variables aléatoires valeurs dans un espace mesurable (E, \mathcal{E}) , définies sur un espace de probabilité $(\Omega, \mathcal{F}, \mathbb{P})$, et telle que pour tout n on ait la propriété suivante :

$$\left. \begin{array}{l} \text{Conditionnellement à la valeur de } X_n, \text{ les variables } (X_0, \dots, X_{n-1}) \\ \text{d'une part, } (X_{n+1}, X_{n+2}, \dots) \text{ d'autre part, sont indépendantes.} \end{array} \right\} \quad (1)$$

Ce modèle permet de rendre compte d'un très grand nombre de situations concrètes.

Encore une fois, nous ne nous intéresserons pas pour l'instant à la théorie de la mesure (qui pourrait très bien – et ce serait une superbe idée – être la thématique d'un GT part entière). Donc pour vous expliquer la phrase ... valeurs dans un espace mesurable (E, \mathcal{E}) , définies sur un espace de probabilité $(\Omega, \mathcal{F}, \mathbb{P})$..., je vous dirais tout simplement : “cela veut dire que tout est bien définie et utilisable dans le contexte voulu”. C'est une phrase qui permet d'écarter tout phénomène pathologique exotique qui pourrait se glisser dans notre définition. Un peu comme une formule de marabout [?] pour repousser les mauvais esprits.

Cette définition se réécrit en termes plus symboliques :

Définition/ Propriété [*chaîne de Markov*] : Soit $(X_n)_{n \in \mathbb{N}}$ une suite de variables aléatoires définies *blabla* et valeurs dans l'espace mesurable *blabla*. $(X_n)_{n \in \mathbb{N}}$ définit une chaîne de Markov si et seulement si $\forall n$:

$$\mathbb{P}[X_{n+1}, X_{n+2}, \dots | X_0, \dots, X_{n-1}, X_n] = \mathbb{P}[X_{n+1}, X_{n+2}, \dots | X_n] \quad (2)$$

Cette dernière propriété (2) est tout bonnement la **propriété de Markov faible** “gnralise” (*faible* car il y a une version forte qui se généralise à un domaine d'application plus tendue (elle est super compliquée donc je ne prendrais pas le risque de l'introduire) et *gnralise* car elle concerne “tout l'avenir” sachant “tout le passé”). Dans ce chapitre, nous ne nous intéressons qu'aux chaînes de Markov **homogènes**. C'est à dire que la loi d'un événement futur X_{n+1} ou de n'importe quelle réunion d'événements futurs $(X_{n+1}, X_{n+2}, \dots)$ conditionnellement à un événement de départ X_n reste la même quelque soit l'instant n . En d'autres termes :

Propriété [*homogénéité*] :

$$\forall n, \quad \mathbb{P}[X_{n+1}, X_{n+2}, \dots | X_n] = \mathbb{P}[X_1, X_2, \dots | X_0] \quad (3)$$

On tablie donc la **rgle d'oubli** des chaînes de Markov homogènes :

propriété [règle d'oubli] :

$$\mathbb{P}[X_{n+1}, X_{n+2}, \dots | X_0, \dots, X_{n-1}, X_n] = \mathbb{P}[X_1, X_2, \dots | X_0] \quad (4)$$

Cette dernière propriété sera très importante par la suite. Notons que l'homogénéité peut aussi se voir comme la simplification du critère fondamentale : $(Z_n)_{n \in \mathbb{N}}$ est une chaîne de Markov homogène si et seulement si ses réalisations $(z_n)_{n \in \mathbb{N}}$ vérifient la récurrence suivante :

$$z_n = f(z_{n-1}, x_n)$$

où f est mesurable de \mathbb{R}^2 dans \mathbb{R} et les $(x_n)_{n \in \mathbb{N}}$ sont les réalisations d'une suite de variables aléatoires $(X_n)_{n \in \mathbb{N}}$ toutes de même loi et telle que toutes les unions possibles $(X_n)_{n \in \mathbb{I} \subset \mathbb{N}}$ soient indépendantes de Z_0 .

3.3 Chaînes de Markov cachées

Voilà, normalement on ne devrait pas avoir besoin de plus de notions que ça pour aborder les HMM. Pour ceux qui veulent approfondir les outils et ustensiles des chaînes de Markov, je conseillerai le poly du GT de François Bienvenu sur les modèles matriciels de populations [?] (ou bien, évidemment, un vrai bouquin de votre BU, mais ça sera probablement super drôle et ennuyeux alors que nos notes nous sont vraiment cool). Nous pouvons maintenant nous poser la question qui nous démangeait tant : *Mais où se cachent donc les chaînes de Markov cachées ?*

Bon, posons-nous une question un peu plus fine : Qu'est-ce qu'une chaîne de Markov *cachée* ? C'est avant tout un modèle. Et la meilleure façon d'introduire un modèle, c'est de présenter un phénomène très simple qui se modélise très bien avec. Je ne ferais pas la même erreur que celle de mon professeur de première année de master qui était allé chercher la probabilité d'apparition des lots CpG proximités des promoteurs dans l'ADN. Ce phénomène d'apparition de ces lots non-méthylés est très intéressant modéliser avec des chaînes de Markov cachées mais bon... Il y a vraiment plus simple comme introduction...

Prenons deux mecs bourrés à la sortie d'un bar. Nous allons faire la supposition suivante : l'un (Pierrot) est vraiment foutu et entame donc une marche aléatoire. L'autre (Robert) a un peu moins bu : il essaye comme il le peut de suivre Pierrot pour ne pas le perdre de vue. Mais bon, il en a quand même un bon coup dans le nez et les choses sont moins simples qu'elles en ont l'air : il n'y voit plus rien. Pierrot-aléatoire se sent rassuré de la présence bienveillante de son ami Robert qui essaye de le suivre. Mais il ne peut pas s'arrêter de marcher et a beaucoup de mal à contrôler ses jambes. Lorsqu'il s'écarte de son ami, il pousse un cri de détresse "HAAAAOOOO". Quand il s'en rapproche, il pousse un cri rassuré "HOOOOOOO". Robert utilise ces cris pour le repérer et le suivre.

Le problème est le suivant : Pierrot pousse des cris approximatifs et se trompe parfois.

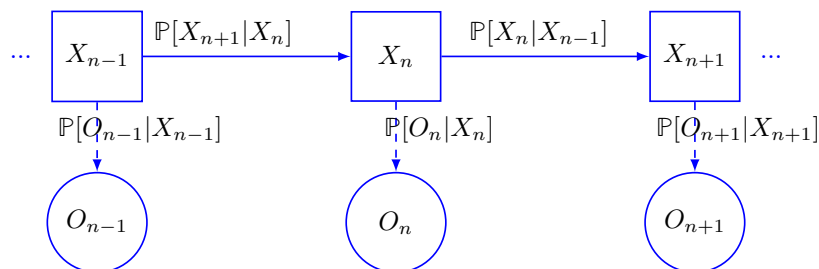
Formalisons cette situation : Il y a deux "tats" possibles pour Pierrot : soit il se rapproche (tat x_1), soit il s'écarte (tat x_2). Notons X_n la variable aléatoire qui représente l'état dans lequel se trouve Pierrot à l'instant n . Lorsque Pierrot prend une direction, il a du mal à s'en faire tout de suite.

Je veux dire par l que lorsqu'il commence s'écarter, il essaye de se sortir de cette situation mais c'est difficile. Disons qu'il a, chaque instant n o il s'écarter, une probabilit de revenir vers Robert $\mathbb{P}[X_{n+1} = x_1 | X_n = x_2] = 0.4$; Bien entendu cela implique la probabilit de continuer s'écarter $\mathbb{P}[X_{n+1} = x_2 | X_n = x_2] = 0.6$. D'autre part, si Pierrot sent qu'il se rapproche, il fait attention ne pas perdre ce bon plie. Sa probabilit de rester dans cet tat de rapprochement est donc leve $\mathbb{P}[X_{n+1} = x_1 | X_n = x_1] = 0.7$; ce qui implique $\mathbb{P}[X_{n+1} = x_2 | X_n = x_1] = 0.3$. Comme Pierrot est sacrement bourr, chaque instant il oublie le pass. Bref, $\mathbb{P}[X_{n+1} | X_n, X_{n-1} \dots X_0] = \mathbb{P}[X_1 | X_0]$. On est clairement face une chane de Markov deux tats. Plus haut, j'ai parl de marche alatoire : il s'agissait de chanes de Markov nombre infini d'tats (la variable alatoire representant l'tat l'instant n de la chane que je notais Z_n pouvait prendre n'importe quelle valeur entiere entre $-\infty$ et $+\infty$ selon la grandeur de n alors qu'ici notre variable X_n de la chane-Pierrot prend comme valeur x_1 ou bien x_2 quelque soit n).

Dans une HMM, les probabilities lies la chane dcrivent ce qu'on appelle les **lois de transition** : elles rgissent le passage d'un tat l'autre.

Pour ce qui est des cris que pousse Pierrot, crivons qu'il y a deux cris possibles : o_1 (HOOOOO) et o_2 (HAAAAO). Notons aussi O_n la variable alatoire qui represente le cri pouss l'instant n . Nous utilisons des \mathcal{O} parce qu'il s'agit d' \mathcal{O} bservations. Disons donc que la probabilit que Pierrot pousse le bon cri quand il est dans l'tat de rapprochement x_1 est $\mathbb{P}[O_{n+1} = o_1 | X_n = x_1] = 0.8$ et que la probabilit qu'il pousse le bon cri quand il s'écarter (x_2) est $\mathbb{P}[O_{n+1} = o_2 | X_n = x_2] = 0.7$. Ces probabilit sont appeles **lois d'mission**. Comme on suppose ici (et c'est une trs grosse supposition !) qu'il n'y a que deux cris possible, $\mathbb{P}[O_{n+1} = o_2 | X_n = x_1] = 0.2$ et $\mathbb{P}[O_{n+1} = o_1 | X_n = x_2] = 0.3$.

La manire la plus simple de reprsenter graphiquement ce processus est de dessiner un "peigne" : le temps s'écoule de gauche droite et de l'instant n l'instant $n + 1$ on passe de l'tat X_n l'tat X_{n+1} . chaque instant, un signal observable (cri) est mis (verticalement). Sur le graphe suivant, les tats sont represents par des carrs, et les signaux observables par des ronds :



Intressons-nous maintenant quelques petites propriets qui dcoulent d'une telle structure...

3.3.1 Quelques propriets qui dcoulent de cette structure

3.4 L'algorithme EM

3.5 L'algorithme de Viterbi

References

- [1] Jacq, A., Bienvenu, F. (2013) *Introduction à l'apprentissage*. sance GT n.8
- [2] Jacod, J. (2003) *Chaînes de Markov, Processus de Poisson et Applications*. Poly. de cours
- [3] Jacod, J. (2004) *Processus de Markov, application à la dynamique des populations*. Poly. de cours
- [4] Metalshine (2047) *Les distances n'existent pas*.
- [5] Bienvenu, F. (2013) *Introduction aux modèles matriciels de populations*. sance GT n.6,7