

Cours MDP et applications.

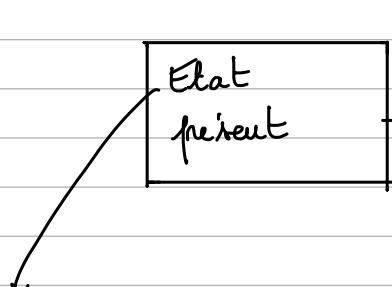
Chapitre 1.

Chaînes de Markov contrôlées

Philosophie générale :

1. Modélisation

Chaînes de Markov contrôlées.



→ évolution aléatoire de type Markov

action : décision prise par l'individu en fonction de ce qu'il observe et en vue de maximiser un critère de satisfaction

nouvel état de l'individu

état peut être :

- grandeur physique
- grandeur économique
- état dans un jeu

⋮

généralement, il s'agit de l'état d'un individu

Calcul de la nouvelle satisfaction instantanée

Très nombreuses applications : finance, économie mais aussi biologie (ex: épidémiologie) ...

Paradigme : Dans cette phase de modélisation, on suppose **CONNUES** toutes les quantités régissant l'évolution du système :

(les réalisations * forme de transitions aléatoires ne sont pas connues, mais la

distribution, statistique et connue en fonction des états occupés et des actions suivies) * forme des récompenses (ou satisfactions) en fonction des états occupés et des actions suivies.

Naturellement, les décisions ne sont pas connues à l'avance.

OBJECTIF : caractériser mathématiquement les meilleures décisions.

d. Apprentissage

En pratique : on ne connaît pas les quantités mathématiques nécessitant le modèle. Les seules choses observées sont les états occupés, les actions décidées et les récompenses reçues.

QUESTION : COMMENT APPRENDRE LA MEILLEURE DÉCISION A "JOUER" EN FONCTION DES OBSERVATIONS ?

→ problématique typique de l'apprentissage statistique. On parle d'APPRENTISSAGE par RÉINFORCEMENT (une des branches de l'apprentissage avec l'apprentissage supervisé et l'apprentissage non supervisé)

intuitivement, la "pertinence" de la décision affise doit évoluer de façon croissante avec le nombre d'observations

Idée générale : pour répondre à 2, il faut bien comprendre 1.

I quelques rafles sur les chaînes de Markov



! Pas de contrôle. Le cadre correspond à

celui étudié en cours de processus stochastiques.

1) principe général

Une chaîne de Markov décrit l'évolution (au cours du temps) d'un système **aléatoire** sans mémoire. Il s'agit de la version probabiliste d'une suite récursive de la forme

$$(*) \quad x_{n+1} = f(n, x_n)$$

Diagramme illustrant la relation entre les états et la fonction f :

- Etat de l'individu à l'instant $n+1$: x_{n+1}
- Fonction "DETERMINISTE" expliquant le passage de l'état courant à l'état suivant : $f(n, \cdot)$
- Etat de l'individu à l'instant n : x_n

Remarques:

- a) la notion d'état est ici la même que la notion d'état dans le paragraphe introductif

b) Ici, il n'y a pas possibilité pour l'individu de prendre une décision

c) Aucun AET à ce stade. Le passage de l'état x_n à l'état x_{n+1} est déterministe

Définition: On appelle espace d'états l'ensemble S des valeurs pour les états de l'individu, à n'importe quel instant. Autrement dit, on suppose, dans $(*)$,

a) $x_0 \in S$

b) $\forall n \in \mathbb{N}, \quad f(n, \cdot) : S \rightarrow S, \text{ i.e}$

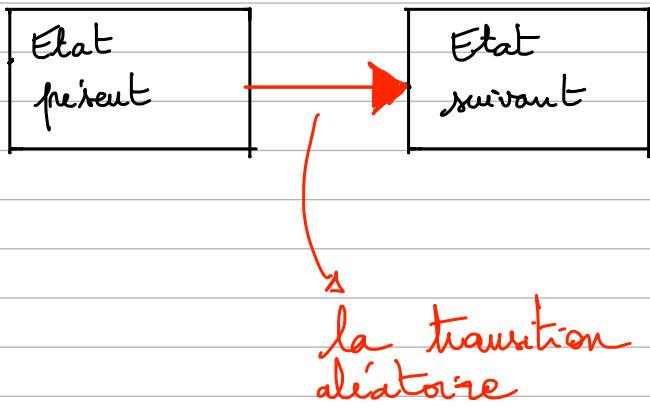
$$f : \mathbb{N} \times S \longrightarrow S$$

Dans la suite, on supposera S un σ -dénombrable et souvent fini.



A ce stade, pas d'aléa'

En présence d'un aléa, il y a une incertitude sur ce que vaut x_{n+1} , même en connaissance de x_n . On peut donc reprendre le graphique introduit.



2) matrice de transition

la fonction $f(n, \cdot)$ est remplacée par une matrice décrivant les probabilités de transiter d'un état à un autre

(**) $\left(f(n, i) \right)_{i \in S}$ avec $\left(P(n, i; j) \right)_{i, j \in S}$

Définition: Une collection de réels $\left(q(i, j) \right)_{i, j \in S}$ forme une matrice de transition ($n \in S$) si

$\forall i \in S, q(i, \cdot)$ forme une collection de poids de probabilité sur S , i.e.

a) $\forall j \in S, q(i, j) \geq 0$

b) $\sum_{j \in S} q(i, j) = 1$.

Application: On suppose donc que, dans (**),

$\forall n \in \mathbb{N}$, $P(n, \cdot, \cdot)$ est une matrice de transition.

$(\Omega, \mathcal{F}, \mathbb{P})$ définition d'une chaîne de Markov
 $(\mathbb{S}, \mathcal{A}, \mathbb{P})$ espace de probabilité

Définition: Étant une suite $(P(n, \cdot, \cdot))_{n \in \mathbb{N}}$ de matrices de transition sur S , une collection (X_n) de variables aléatoires à valeurs dans S (formant ainsi un processus indexé par le temps discret) constitue une chaîne de Markov de matrices de transition $(P(n, \cdot, \cdot))_{n \in \mathbb{N}}$ si l'une des définitions suivantes est vérifiée:

a) $\forall n \in \mathbb{N}, \forall i_0, \dots, i_m, j \in S,$

$$\begin{aligned} \mathbb{P}\{X_0 = i_0, \dots, X_m = i_m, X_{m+1} = j\} \\ = P(m, i_m, j) \quad \mathbb{P}\{X_0 = i_0, \dots, X_m = i_m\} \end{aligned}$$

b) $\forall n \in \mathbb{N}, \forall i_0, \dots, i_m, j \in S$ tels que

$$\mathbb{P}\{X_0 = i_0, \dots, X_m = i_m\} > 0$$

$$\mathbb{P}\{X_{m+1} = j \mid X_0 = i_0, \dots, X_m = i_m\}$$

$$= \mathbb{P}\{X_{m+1} = j \mid X_m = i_m\}$$

$$= \sum (n, i_m, j).$$

On peut généraliser cette définition en utilisant le concept de filtration

Définition: Étant donnée \mathcal{F} une sous- σ -algèbre de \mathcal{E} , et $X: \Omega \rightarrow E$ une variable aléatoire à valeurs dans E dénombrable, on appelle la conditionnelle de X sachant \mathcal{F} la collection

de variables aléatoires :

$$\left(\mathbb{P}(X=i_j | \mathcal{G}_f) \right)_{i \in E}$$

$\underbrace{\quad}_{\mathbb{E} \left(\mathbf{1}_{\{X=i_j\}} | \mathcal{G}_f \right)}$

Si $\mathcal{G}_f = \sigma(Y)$ où $Y: \Omega \rightarrow F$ dénombrable abs sur $\{Y=j\}$, $j \in F$ tel que $\mathbb{P}\{Y=j\} > 0$, on a (TD)

$$\mathbb{E} \left(\mathbf{1}_{\{X=i_j\}} | \mathcal{G}_f \right) = \mathbb{P}\{X=i_j | Y=j\}$$

Autrement dit

$$\mathbb{E} \left(\mathbf{1}_{\{X=i_j\}} | \mathcal{G}_f \right) = \sum_{f \in F} \mathbb{P}\{X=i_j | Y=f\} \mathbf{1}_{\{Y=f\}}$$

$\underbrace{\quad}_{\mathbb{P}\{Y=f\} = 0}$

\mathcal{G}_f la définition est arbitraire

Définition : On appelle filtration sur $(\Omega, \mathcal{S}, \mathbb{P})$ toute collection croissante de sous-tribus de \mathcal{S} .
 - autrement dit $(\mathcal{G}_n)_{n \in \mathbb{N}}$

$$\mathcal{G}_n \subset \mathcal{G}_{n+1}, \quad n \in \mathbb{N}.$$

Application :

Définition : Étant une suite $(P(n, \cdot, \cdot))_{n \in \mathbb{N}}$ de matrices de transition sur S et une filtration $(\mathcal{G}_n)_{n \in \mathbb{N}}$ sur $(\Omega, \mathcal{S}, \mathbb{P})$ une collection $(X_n)_{n \in \mathbb{N}}$ de variables aléatoires à valeurs dans S constiue une $(\mathcal{G}_n)_{n \in \mathbb{N}}$ -chaîne de Markov de matrices de transition $(P(n, \cdot, \cdot))_{n \in \mathbb{N}}$.

- $\forall n \in \mathbb{N}$, X_n est \mathcal{S}_n^1 -mesurable

- $\forall n \in \mathbb{N}$, $\forall i \in S$,

$$\mathbb{P}\left(\{X_{n+1} = i\} \mid \mathcal{S}_n^1\right) = P(n, X_n, i)$$

(P p.l.)

Remarque: toute chaîne de Markov au sens de la première définition forme une $(\mathcal{S}_n^1)_{n \in \mathbb{N}}$ -chaîne de Markov pour

$$\mathcal{S}_n^1 := \langle (X_0, \dots, X_n) \rangle.$$

en effet

$$\text{sur } \{X_n = j\} \quad (\text{avec } \mathbb{P}\{X_n = j\} > 0)$$

$$\mathbb{P}\left(\{X_{n+1} = i\} \mid \mathcal{S}_n^1\right) = \mathbb{P}\left(\{X_{n+1} = i\} \mid \{X_n = j\}\right)$$

def loi conditionnelle $= P(n, j, i) = P(n, X_n, i).$

def chaîne de Markov

4) Construction canonique.

On donne ici la construction canonique d'une chaîne de Markov généralisant la définition d'une chaîne déterministe sans mémoire donnée dans §1.

Proposition: Soit $f: N \times S \times [0, 1] \rightarrow S$ une fonction mesurable.

Soyons $x_0: \Omega \rightarrow S$ et $(\mathbb{I}_n)_{n \in \mathbb{N}, \omega \in \Omega}$ une w.a. et une suite de w.a. telles que

i) $\mathbb{I}_n \sim \text{Unif}[0, 1]$ pour chaque $n \in \mathbb{N}$

- ii) (U_n) ensemble suite de r.v.a. \perp
 iii) $X_0 \perp (U_n)$ alors la suite $(X_n)_{n \in \mathbb{N}}$ définie par

$$X_{n+1} = f(n, X_n, U_{n+1})$$

est une chaîne de Markov de matrices de transition

$$\left(P(n, i; j) := \mathbb{P} \{ f(n, i, U_n) = j \}_{\substack{i, j \in S \\ n \in \mathbb{N}}} \right)$$

Remarque : (a) ii) et iii) sont \Leftrightarrow à, pour tout $n \in \mathbb{N}$, X_0, U_1, \dots, U_n sont \perp .

(b) démonstration faite en TD

En réalité, il existe une réciproque (preuve également faite en TD)

Proposition : Toute chaîne de Markov (de matrices de transition $(P(n, i; j))_{n, i, j}$ données) admet une forme canonique

Indice ; $f(n, i, \cdot)$ fonction quantile de la loi $(P(n, i, j))_{j \in S}$.

II Chaînes contrôlées

Nous avons maintenant tous les éléments pour définir un système contrôlé, analogue des chaînes de Markov

1) Cas déterministe

Commençons par examiner ce que serait la version déterministe

Dans l'écriture

$$x_{n+1} = f(n, x_n)$$

on dit **AJOUTER** la décision prise par l'individu à l'instant n .

Définition: On appelle espace d'actions un ensemble A contenant toutes les actions (ou tous les choix) possibles pour l'individu à n'importe quel instant. Dans un but de simplification, on suppose de fait que A est indépendant du temps (mais cela pourrait être généralisé au cas où les actions possibles varient avec le temps).

Dans la suite, A est au moins au + dénombrable, mais on peut aussi avoir des cas où A est un Borelien de \mathbb{R}^d , pour un certain $d \geq 1$. On fait donc cette hypothèse ici.

Nous sommes donc amenés à considérer

$$f: \mathbb{N} \times S \times A \longrightarrow S$$

mesurable (i.e., $\forall (n, i) \in \mathbb{N} \times S$, $f(n, i, \cdot): A \rightarrow S$ est mesurable au sens clair que

$$\forall j \in S, \{a \in A : f(n, i, a) = j\} \in \mathcal{B}(\mathbb{R}^d)$$

(démonstration à faire en TD)

Définition: Une suite $(x_n, a_n)_{n \in \mathbb{N}}$ à valeurs dans $S \times A$ forme une chaîne sans mémoire contrôlée par $(a_n)_{n \in \mathbb{N}}$ et dirigée par f si

$$\forall n \in \mathbb{N} \quad x_{n+1} = f(n, x_n, a_n).$$

Attention: Il n'y a pas de dynamique pour $(a_n)_{n \in \mathbb{N}}$:
 — chaque a_n est choisi librement par l'individu
 dont l'évolution est représentée par la chaîne.

2) matrices de transition contrôlées.

On veut maintenant généraliser la définition donnée dans le §1 précédent au cas stochastique en nous appuyant sur le concept de matrices de transition vu dans la section I.

Définition. On appelle matrices de transition contrôlées toute application mesurable

$$P: \mathbb{N} \times S \times A \times S \rightarrow [0, 1]$$

telle que pour tout $(n, a) \in \mathbb{N} \times A$ la

$\left(P(n, i, a, j) \right)_{i, j}$ forme une matrice de transition

c'est à dire pour tout $(n, y, a) \in \mathbb{N} \times S \times A$

$P(n, i, a, \cdot)$ est une famille de poids de probabilité sur S .

Remarque: On dit que la collection est homogène si

$$P(n, i, a, j) = P(i, a, j)$$

de sorte que $P: S \times A \times S \rightarrow [0, 1]$ mesurable
 de sorte que $P(\cdot, a, \cdot)$ forme une matrice de transition pour tout $a \in A$.

Interprétation: Si l'individu est dans l'état i à l'instant n, alors en choisissant l'action a et, il a probabilité $P(n, i, a, j)$ de se retrouver dans l'état j à l'instant suivant n+1.

Exemple: On considère un service informatique soumis régulièrement à des attaques de hackers. Le service a deux états



À chaque instant $n \in \mathbb{N}$

En état 0 : il y a deux choix

- NE RIEN FAIRE . La protection du système est dégradée, mais cela ne coûte rien au manager du système !

- METTRE à JOUR la protection . La protection du système est renforcée, mais il y a des frais de mise à jour

En état 1 il y a deux choix

- NE RIEN FAIRE . Le système fonctionne en mode dégradé. Cela ne coûte rien en termes de réparation, mais cela a un coût sur les activités autour .

- REPARER le système . Cela a un coût pour le service informatique mais préserve les activités périphériques .

On a donc $A = \{0, 1\}$

On ne fait rien \xrightarrow{s} \xrightarrow{r} On agit

et on voit

$$P(0, 0, 1) = p \quad \text{ou } p \in \{0, 1\}.$$

$$P(0, 0, 0) = 1-p$$

$$P(0, 1, 1) = q$$

$$P(0, 1, 0) = 1-q$$

Clairement $p > q$.

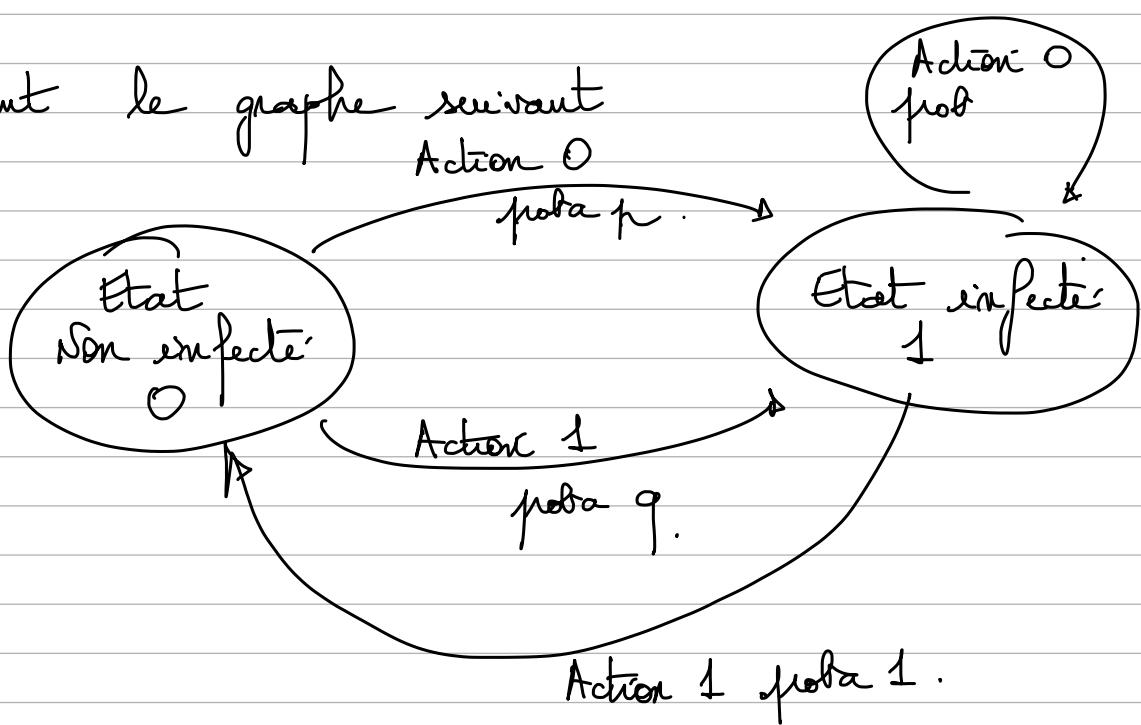
$$P(1, 0, 1) = 1$$

$$P(1, 0, 0) = 0$$

$$P(1, 1, 0) = 1$$

$$P(1, 1, 1) = 0.$$

On obtient le graphe suivant



Remarque: le modèle est simple. On oublie la protection passée dans la probabilité d'être affecté par une attaque.

On pourrait ajouter un état: "protégé" mais pas mis à jour.

3) Chaînes de Markov contrôlées

On veut maintenant associer une suite d'états et une suite d'actions aux matrices de transition contrôlées. Il y a néanmoins une subtilité: quelles sont les informations disponibles à l'instant n (ou à l'observateur) pour prendre la décision à un instant donné? Il s'agit de fait d'une question de mesurabilité.

On commence par la définition suivante, qui généralise celle introduite pour les chaînes de Markov.

Définition: Étant donnée une filtration $(\mathcal{F}_n)_{n \in \mathbb{N}}$ et une famille de matrices de transition contrôlées

$$P: N \times S \times A \ni (n, i, a) \mapsto P(n, i; a) \in [0, 1]$$

une suite $(X_n, A_n)_{n \in \mathbb{N}}$ de variables aléatoires à valeurs dans $S \times A$ est une $(\mathcal{F}_n)_{n \in \mathbb{N}}$ - chaîne de Markov contrôlée associée à P si

- $\forall n \in \mathbb{N}$, (X_n, A_n) est \mathcal{F}_n -mesurable

- $\forall n \in \mathbb{N}$, $\forall i \in S$

$$P(X_{n+1} = j | \mathcal{F}_n) = P(n, X_n, A_n, i)$$

Illustrons la dernière propriété dans le cas où A est un + dénombrable.

Proposition: Si A est un + dénombrable.

$\forall n \in \mathbb{N} \quad \forall (i_0, a_0, i_1, a_1, \dots, i_m, a_m) \in \underbrace{S \times A \times S \times A \times \dots}_{m \text{ fois}}$

$\forall j \in S$

$$\mathbb{P} \left\{ X_0 = i_0, A_0 = a_0, \dots, X_m = i_m, A_m = a_m, X_{m+1} = j \right\}$$

$$= \mathbb{P}(n, i_m, a_m, j) \quad \mathbb{P} \left\{ X_0 = i_0, A_0 = a_0, \dots, X_m = i_m, A_m = a_m \right\}.$$

Preuve: On a

$$\{X_0 = i_0, A_0 = a_0, \dots, X_m = i_m, A_m = a_m\} \in \mathcal{F}_m^1$$

$$\mathbb{E} \left\{ \mathbb{1}_{\{X_{m+1} = j\}} \mid \{X_0 = i_0, A_0 = a_0, \dots, X_m = i_m, A_m = a_m\} \right\}$$

$$= \mathbb{E} \left\{ \mathbb{E} \left(\mathbb{1}_{\{X_{m+1} = j\}} \mid \mathcal{F}_m^1 \right) \mid \{X_0 = i_0, A_0 = a_0, \dots, j\} \right\},$$

$$\mathbb{P}(n, i_m, a_m, j).$$

$$\mathbb{P}(n, i_m, a_m, j)$$

$$= \mathbb{P}(n, i_m, a_m, j) \quad \mathbb{P} \left\{ X_0 = i_0, A_0 = a_0, \dots, j \right\}.$$

Dans le cas où A est non-dénombrable, on ne peut plus procéder comme ci-dessus: les événements peuvent être de probabilité nulle (en particulier si A_0, \dots, A_m sont à droite).

Proposition: Dans le cas général, si B_0, \dots, B_m sont des ensembles de A , et (i_0, \dots, i_m) sont dans S , alors

$$\begin{aligned} & \mathbb{P} \left(\{X_0 = i_0, A_0 \in B_0, X_1 = i_1, A_1 \in B_1, \dots, X_m = i_m, A_m \in B_m, X_{m+1} = j\} \right) \\ &= \int_{S^{m+1} \times (\mathbb{R}^d)^{m+1}} \mathbb{P}(n, i_m, a_m) \quad \mathbb{1}_{B_0 \times \dots \times B_m}(a_0, \dots, a_m) \quad \mathbb{1}_{\{i_0, \dots, i_m\}}(j_0, \dots, j_m) \\ & \qquad \qquad \qquad d\mathbb{P}_{(X_0, A_0, \dots)}(j_0, a_0, \dots, j_m, a_m) \end{aligned}$$

Preuve. Comme précédemment

$$\{x_0=i_0, A_0 \in \mathcal{B}_0, \dots, x_m=i_m, A_m \in \mathcal{B}_m\} \in \mathfrak{F}_m$$

$$\begin{aligned} & \mathbb{P}\left(\{x_0=i_0, A_0 \in \mathcal{B}_0, x_1=i_1, A_1 \in \mathcal{B}_1, \dots, x_m=i_m, A_m \in \mathcal{B}_m, x_{m+1}=j\}\right) \\ &= \mathbb{E} \underbrace{\{P(a, x_m, A_m, j)\}}_{\{x_0=i_0, A_0 \in \mathcal{B}_0, \dots, x_m=i_m, A_m \in \mathcal{B}_m\}} \end{aligned}$$

$$P(a, x_m, A_m, j)$$

Application. Si $\mathbb{P}\{x_0=i_0, \dots, x_m=i_m\} > 0$, on appelle la conditionnelle de (A_0, \dots, A_m) sachant $x_0=i_0, \dots, x_m=i_m$ la mesure de probabilité

$$\begin{aligned} & \mathbb{P}_{(A_0, \dots, A_m) | x_0=i_0, \dots, x_m=i_m}(B_0 \times \dots \times B_m) \\ &= \mathbb{P}\{(A_0, \dots, A_m) \in B_0 \times \dots \times B_m | x_0=i_0, \dots, x_m=i_m\}. \end{aligned}$$

Alors

$$\mathbb{P}\{x_0=i_0, A_0=a_0, \dots, x_m=i_m, A_m=a_m, x_{m+1}=j\}$$

$$= \int_{\mathcal{G}^{m+1}} P(a, i_m, a_m) P_{(A_0, \dots, A_m) | x_0=i_0, \dots, x_m=i_m}(da_0, \dots, da_m).$$

4) Mesurabilité des actions.

On distingue 3 types de propriété de mesurabilité pour $(A_n)_{n \in \mathbb{N}}$.

Définition 1. On dit que les actions $(A_n)_{n \in \mathbb{N}}$ s'écrivent de forme fermée dépendant du 'pas' si

$\forall n \in \mathbb{N}$, A_n est $\sigma(x_0, \dots, x_n)$ measurable

c'est à dire $A_n = f_n(x_0, \dots, x_n)$

On choisit l'action à l'instant n en fonction des observations des états passés de 0 à n .

Définition 2: On dit que les actions $(A_n)_{n \in \mathbb{N}}$ s'écrivent sous forme fermée markoviennne si:

$\forall n, A_n$ est $\sigma(X_n)$ -mesurable

c'est à dire $A_n = f_n(X_n)$

On choisit l'action à l'instant n en fonction de la seule observation de l'état à l'instant n .

Définition 3: On dit que les actions $(A_n)_{n \in \mathbb{N}}$ s'écrivent sous forme mixte markoviennne si il existe une suite $(V_n)_{n \in \mathbb{N}}$ de va l de loi uniforme sur $[0,1]$ telle que

- $\forall n \in \mathbb{N}, V_n$ est \mathcal{E}_n -mesurable
- $\forall n \in \mathbb{N}, V_n \perp \sigma(X_n) \vee \mathcal{E}_{n-1}$.
- $\forall n \in \mathbb{N}, A_n$ est $\sigma(X_n, V_n)$ mesurable.

Interprétation: à l'instant n , on a obtenu, en plus de l'observation de l'état X_n à l'instant n un tirage V_n pour décider de l'action à suivre

Exercice (TD): On suppose A fini.

À chaque instant n , on te donne

$$f_n: S \rightarrow \underbrace{\mathcal{L}(A)}_{\text{espace de probabilité sur } A}$$

(v l'espace de probabilité sur A)

Montrer que f_n induit une suite d'actions

markoviennes mixtes.

Indice : On appelle $(V_n)_{n \in \mathbb{N}}$ suite de r.a. 1 de loi uniforme sur $[0,1]$ et $\mathcal{L}(V_n)$ la loi de la $(V_n)_{n \in \mathbb{N}}$.

Pour $\mu \in \mathcal{L}(A)$, on appelle $\mathcal{Q}(\mu, \cdot)$ la fonction de $[0,1]$ dans A la fonction suivante.

$$\mathcal{Q}(\mu, u) = \sum_{i=1}^{|A|} a_i \mathbf{1}_{[\mu(a_1) + \dots + \mu(a_{i-1}), \mu(a_1) + \dots + \mu(a_i))}(u)$$

\star

$u \in [0,1] \quad \text{où } A = \{a_1, \dots, a_M\}$

Alors on construit, pour x_0 donné,

$$t_0 = \mathcal{Q}(f_0(x_0), V_0)$$

$$x_1 \sim P(0, x_0, t_0, \cdot)$$

$$t_1 = \mathcal{Q}(f_1(x_1), V_1)$$

$$x_2 \sim P(1, x_1, t_1, \cdot)$$

:

$$\text{et } \xi'_n = \sigma(x_0, t_0, \dots, x_n, t_n).$$

Exercice (TD)

On se donne

$$(f_n: S \rightarrow \mathcal{L}(A))_{n \in \mathbb{N}}$$

$$F: \mathbb{N} \times S \times A \times [0,1] \rightarrow S$$

$(x_0, (V_n)_{n \in \mathbb{N}}, (U_n)_{n \in \mathbb{N}, \text{def}})$ collection de r.a. 1

avec $V_n \sim \text{Unif}[0,1]$ et $U_n \sim \text{Unif}[0,1]$. On construit

$$t_0 = \mathcal{Q}(f_0(x_0), V_0)$$

$$x_1 = F(0, s_0, t_0, U_1)$$

$$A_1 = g(f_1(x_1), v_1)$$

$$x_2 = F(1, S_1, A_1, u_1)$$

⋮

montrer que l'on obtient une chaîne de markov contrôlée avec

$$\tilde{x}_n = \sigma(x_0, v_0, u_1, \dots, u_n, v_n).$$

5) récompenses.

Etant donnée une chaîne de markov contrôlée $(x_n, a_n)_{n \in \mathbb{N}}$ comme précédemment, on associe

R_n : \mathbb{R}_m^m -measurable représentant la récompense à l'instant n

Typiquement

$$R_n = r_n(x_n, a_n).$$

Objectif : Maximiser les récompenses moyennes accumulées.
Ce sera le chapitre 2.

III Exemples (TD)

1) Approvisionnement de stock.

Un centre de livraison s'approvisionne de la façon suivante

- s_0 : stock initial

- A l'instant $n \in \mathbb{N}$

s_n : état du stock (entier dans $\{0, \dots, N\}$)

a_n : approvisionnement (entier dans $\{0, \dots, N\}$)

• Entre n et $n+1$

D_{n+1} : demande aléatoire (independante des observations passées) de loi homogène (en $\lambda_0, \dots, \lambda_D$)
la demande est satisfaite si $S_n + A_n \geq D_{n+1}$
Sinon elle est satisfaite dans la mesure du possible.

Ecrire la chaîne de Markov contrôlée associée

Indice: $S_{n+1} = \max(S_n + A_n - D_{n+1}, 0)$.

2) Bandit multi-bras.

Un joueur observe les résultats de K machines à sous différentes. À chaque étape il en choisit une et une seule et la joue : les $K-1$ autres machines ne sont pas jouées.

À chaque partie n , on connaît les états des K machines. Un seul état parmi les K est alors actualisé.

On suppose que la machine $i \in \{1, \dots, K\}$ peut être dans les états d'un ensemble S^i . La machine i passe d'un état s_i à un état s'_i avec une probabilité:

$$P(i; s_i, s'_i).$$

Ecrire une chaîne de Markov contrôlée avec

$S^1 \times \dots \times S^K$ comme espace d'états

$\{1, \dots, K\}$ comme espace d'actions.

3) Sélection d'un film sur Netflix

Un abonné Netflix sélectionne un film avec la règle suivante : il visualise les propositions de façon séquentielle et ne peut choisir que le film dont il est en train de consulter la notice (pas de retour en arrière).

$$S = \{0,1\}$$



1 si le film consulté est le meilleur choix (au regard des goûts de l'abonné) parmi tous ceux consultés jusqu'ici
0 sinon

$$A = \{0,1\}$$

1 si on choisit le film dont on consulte la notice, 0 sinon

On obtient donc une chaîne de Markov contrôlée.
Montrer que les transitions sont

$$P(n, i, \alpha, 1) = \frac{1}{n+1}.$$

$$\begin{aligned} & \cdot \quad \alpha \in \{0,1\} \\ & \alpha \in \{0,1\} \end{aligned}$$