

Approximations stochastiques

Convergence et théorie asymptotique :

Le cadre des approximations stochastiques recouvre un spectre beaucoup plus large que ceux présentés les séances précédentes (gradient stochastique et on stochastique)

la situation générale est la suivante. On suppose que l'on cherche θ^*

$$\text{tg } h(\theta^*) = \mathbb{E}_{\theta^*} [H(\theta^*, Y)] = 0 \quad (*)$$

où H est connue mais la distribution de Y (qui peut dépendre de θ) ne l'est pas.

le fait est que l'on n'a pas accès directement à la fonction de coût h mais à des observations aléatoires $H(\theta, Y_n)$

ex: • minimiser / maximiser $Q(\theta)$ avec comme observations une version bruitée du gradient de Q : $H(\theta, Y_n) = \nabla Q(\theta) + Y_n$
où $(Y_n)_{n \geq 0}$ sont des réalisations d'un bruit additif.

• le gradient stochastique : $H(\theta, Y_n) = f(\theta, W_{n+1})$ où $W_{n+1} \sim P_\theta$

Vocabulaire: la fonction h est appelée le champ moyen.

L'idée des approximations stochastiques est de trouver θ^* itérativement en utilisant une mise à jour sous la forme:

$$\theta_n = \theta_{n-1} + \gamma_n h(\theta_{n-1}) + \gamma_n \eta_n \quad (**)$$

où $h(\theta) = \mathbb{E}_\theta [H(\theta, Y)]$ est notre fonction à annuler

et $\eta_n = H(\theta_{n-1}, Y_n) - h(\theta_{n-1})$ joue le rôle d'une perturbation aléatoire.

la raison de ce schéma apparaîtra au fur et à mesure puisque

(**) s'écrit en fait $\theta_n = \theta_{n-1} + \gamma_n H(\theta_{n-1}, Y_n)$ dans bcp de cas ce qui ressemble bien à notre cas du ∇ st.

Pour montrer la convergence de la suite $(\theta_n)_{n \geq 0}$, il est nécessaire de procéder en deux étapes :

- (1) Trouver des conditions générales sur une suite déterministe de $(\eta_n)_{n \geq 0}$ et sur h qui assure la convergence déterministe de $(O_n)_{n \geq 0}$
- (2) Montrer que ces conditions sont satisfaites avec probabilité 1 c'est à dire pour presque tout chemin $(\eta_n(\omega), \dots, \eta_n(\omega))$ du processus de bruit.

Le 1^{er} point est développé ci-dessous : sous des conditions de stabilité sur $(O_n)_n$ et des propriétés sur $(\eta_n)_{n \geq 0}$ (déterministe), nous allons voir que la convergence de la suite est fortement liée à la convergence de l'équation dite du champ moyen : $\frac{d\Theta_t}{dt} = h(\Theta_t)$

I Convergence ponctuelle dans le cas borné :

On suppose que notre terme de bruit (perturbation) peut s'écrire sous la forme : $\eta_n = \varepsilon_n + \gamma_n$ où :

- la série $\sum_{i=1}^n \gamma_i \varepsilon_i$ converge
- $\gamma_n \xrightarrow[n \rightarrow \infty]{} 0$

On va voir que ces conditions sont imposées par la convergence.

Hypothèses : $\mathcal{O} \subset \mathbb{R}^d$ est un ouvert dont on note $\partial \mathcal{O}$ la frontière.

(A) h est un champ de vecteur continu sur \mathcal{O}

\exists V une fonction positive (ou nulle) C^1 tq

(i) $\forall x \in \mathcal{O} \quad \langle \nabla V(x), h(x) \rangle \leq 0$

(ii) l'ensemble $S = \{x : \langle \nabla V(x), h(x) \rangle = 0\}$ satisfait $V(S)$ est d'intérieur vide.

Rq : la fonction V est connue sous le nom de fonction de Lyapunov pour le champ h . C'est une condition très courante et peu restrictive. Par exemple, si $h = -\nabla Q$ pour une fonction d'énergie ou de vraisemblance Q , le choix $V = Q$ fonctionne dès que Q est C^1 .

On va voir que l'hypothèse (A) va permettre d'assurer la convergence vers l'ensemble S . En général, S coïncidera avec $\{x : h(x) = 0\}$. De plus si $|S| < \infty$, on convergera vers un point de S .

le point (ii) semble compliqué à prouver, cependant le théorème de Sard (géométrie différentielle) peut souvent être utilisé. SP dit que pour une fonction $V \in C^d$ alors $V(\{\nabla V = 0\})$ est d'intérieur vide.

Nous allons maintenant donner les conditions de convergence de notre suite $(\theta_n)_n$ que nous allons écrire sous la forme: On donne

$$\theta_n = \theta_{n-1} + \gamma_n h(\theta_{n-1}) + \gamma_n e_n + \gamma_n x_n \quad (1)$$

avec $\gamma_n \geq 0$; $\sum_{p=1}^{\infty} \gamma_p = +\infty$ et $\lim_{n \rightarrow \infty} \gamma_n = 0$

h est défini sur $\mathcal{O} \subset \mathbb{R}^d$, en et x_n sont des perturbations et γ_n est appelé suite de gain (scalaires positive ou nulle)

Def: On dit que l'algorithme est A-stable si

θ_n reste dans un sous ensemble compact de \mathcal{O}

$\lim_{p \rightarrow \infty} \sum_{n=1}^p \gamma_n e_n$ existe

et $\lim_{n \rightarrow \infty} |\gamma_n| = 0$

théorème: On suppose que l'algorithme donné sous la forme (1) est A-stable. Sous l'hypothèse (A), la distance de θ_n à l'ensemble S notée $d(\theta_n, S)$ converge vers 0. En particulier si $|S| < \infty$, θ_n tend vers un point de S .

Preuve: Soit $F(\theta) = \langle \nabla V(\theta), h(\theta) \rangle$.

$$\text{Soit } \theta'_n = \theta_n + \sum_{i=n+1}^{\infty} \gamma_i e_i$$

$$\text{et } \delta'_n = - \sum_{i=n}^{\infty} \gamma_i e_i$$

$$\begin{aligned} \text{Donc } \theta'_{n+1} &= \theta_n + \delta'_{n+1} = \theta_{n-1} + \gamma_n h(\theta_{n-1}) + \gamma_n e_n + \gamma_n x_n - \delta'_{n+1} \\ &= \underbrace{\theta_{n+1} - \delta_n}_{\theta'_n} + \underbrace{\delta_n - \delta'_{n+1}}_{\gamma_n e_n} + h(\theta_{n-1}) + \gamma_n e_n + \gamma_n x_n \\ &= \theta'_{n-1} - \gamma_n e_n + h(\theta'_{n-1} + \delta_n) + \gamma_n e_n + \gamma_n x_n \\ &= \theta'_{n-1} + h(\theta'_{n-1} + \delta_n) + \gamma_n x_n \quad (\square) \end{aligned}$$

En posant $x'_n = x_n - h(\theta'_{n-1}) + h(\theta'_{n-1} + \delta_n)$ on obtient

$$\theta'_n = \theta'_{n-1} + r_n h(\theta'_{n-1}) + r_n r'_n.$$

(4)

De plus comme la série des r_n converge, la suite $r_n \xrightarrow{n \rightarrow \infty} 0$.

Par hypothèse $\theta_n \in K$ compact de S (A stabilité) donc

$\theta'_n \in K_0 \subset S$ compact aussi

Comme h est uniformément continue sur K_0 donc

$$h(\theta'_{n-1}) - h(\theta'_{n-1} + r_n) \xrightarrow{n \rightarrow \infty} 0$$

$$\text{Enfin } r_n \xrightarrow{n \rightarrow \infty} 0 \quad \text{Donc } r'_n \xrightarrow{n \rightarrow \infty} 0$$

la fonction V est C^1 donc en utilisant la formule de Taylor :

~~soit~~ $\exists \theta'' \in [\theta'_{n-1}; \theta'_n]$ tq

$$V(\theta'_n) = V(\theta'_{n-1}) + \underbrace{r_n}_{\substack{\text{car } \theta'_n - \theta'_{n-1}}} \langle \nabla V(\theta''), h(\theta'_{n-1}) + r'_n \rangle$$

$$\begin{aligned} &= V(\theta'_{n-1}) + r_n \langle h(\theta'_{n-1}); \nabla V(\theta'_{n-1}) \rangle \\ &\quad + r_n \langle h(\theta'_{n-1}); \nabla V(\theta'') - \nabla V(\theta'_{n-1}) \rangle \\ &\quad + r_n \langle \nabla V(\theta''), r'_n \rangle \end{aligned}$$

$$\text{Soit } r''_n = \langle h(\theta'_{n-1}); \nabla V(\theta'') - \nabla V(\theta'_{n-1}) \rangle + \langle \nabla V(\theta''), r'_n \rangle$$

$$\text{alors } |r''_n| \leq \|h\|_{L^\infty(K_0)} \|\nabla V(\theta'') - \nabla V(\theta'_{n-1})\| + \|\nabla V\|_{L^\infty(K_0)} |r'_n|$$

$$\text{or } r'_n \xrightarrow{n \rightarrow \infty} 0 \quad \text{et } \|\nabla V(\theta'') - \nabla V(\theta'_{n-1})\| \xrightarrow{n \rightarrow \infty} 0 \quad \text{car}$$

$$\theta'' \in [\theta'_n; \theta'_{n-1}] \quad \text{et} \quad |\theta'_{n-1} - \theta'_n| \leq r_n |h(\theta'_{n-1})| + r_n |r'_n| \xrightarrow{n \rightarrow \infty} 0$$

$\leq C \text{ sur } K_0 \text{ et } r_n \rightarrow 0$

donc $|\theta'_{n-1} - \theta''| \rightarrow 0$. Par continuité de ∇V car $V \in C^1$ on a le résultat.

$$\text{Donc } V(\theta'_n) = V(\theta'_{n-1}) + r_n F(\theta'_{n-1}) + r_n r''_n$$

avec $r''_n \xrightarrow{n \rightarrow \infty} 0$

Soit maintenant \mathcal{P} un voisinage de $\mathcal{S} \cap K_0$. On rappelle que

$$\mathcal{S} = \{F=0\} \quad \text{et} \quad F(\theta) \leq 0 \quad \forall \theta \in K_0 \quad \text{par (A)}$$

Par continuité de F sur K_0 :

$$\text{si } \theta_{n-1} \in K_0 \setminus \mathcal{P} \text{ alors } F(\theta_{n-1}) \leq -2\varepsilon$$

$$\text{si } \theta_{n-1} \notin K_0 \setminus \mathcal{P} \text{ i.e. } \theta_{n-1} \in \mathcal{P} \text{ alors } F(\theta_{n-1}) \leq C_1 \text{ étant une fonction } C^0 \text{ sur un compact.}$$

Donc pour n assez grand : $|r_n| \leq \varepsilon$ donc $\gamma_n |r_n| \leq \varepsilon \gamma_n$

$$(\Delta) \quad V(\theta'_n) \leq V(\theta_{n-1}) - \gamma_n \varepsilon \mathbb{1}_{\{\theta_{n-1} \notin \mathcal{P}\}} + \gamma_n C \mathbb{1}_{\{\theta_{n-1} \in \mathcal{P}\}}$$

$$\text{où } C = C_1 + \varepsilon.$$

\hookrightarrow qui prend en compte $\gamma_n r_n$ de la case cd

On va terminer la preuve en démontrant que $V(\theta'_n) \rightarrow V(\mathcal{S})$
et ensuite $\theta'_n \rightarrow \mathcal{S}$

$$\bullet \text{ Soit } A_\alpha = \{x \in \mathbb{R} : d(x, V(\mathcal{S} \cap K_0)) < \alpha\}$$

comme $\mathcal{S} \neq \emptyset$ A_α est une union d'intervalles de longueur 2α

Comme K_0 compact (i.e. borné), A_α est donc une union finie d'intervalles de longueur au moins 2α . (on peut avoir $\frac{1}{n} \rightarrow 0$ donc plus gd)

Soit α fixe, petit. On définit $\mathcal{P} = V^{-1}(A_\alpha)$

Dans ce cas, (Δ) définit une suite $u_n = V(\theta'_n)$ tq

$$u_n \leq u_{n-1} - \gamma_n \varepsilon \mathbb{1}_{\{u_{n-1} \notin A_\alpha\}} + \gamma_n C \mathbb{1}_{\{u_{n-1} \in A_\alpha\}}$$

A chaque fois que $u_{n-1} \notin A_\alpha$ $u_n \leq u_{n-1} - \gamma_n \varepsilon$

Puisque $\sum \gamma_i = +\infty$ u_n ne peut pas être une infinité de fois en dehors de A_α . ~~Donc~~ (car $u_n \in V(\theta'_n) \subset V(K_0)$ compact)

Donc à partir d'un certain rang $u_n \in A_\alpha$ i.e. $d(u_n, A_\alpha) = 0$

- $S \cap K_0$ est défini comme $S = F^{-1}(0)$ donc 'S fermé' ⑥
 donc $S \cap K_0$ fermé dans un compact donc compact
 Comme $V \in C^0$: $V(S \cap K_0)$ compact donc fermé et $V(S)$ est
 d'intérieur vide donc par des arguments topologiques

$$\begin{aligned} & \mu_n \xrightarrow{n \rightarrow \infty} \text{un élément de } V(S) \\ \text{càd } & V(\theta_{n+1}') \xrightarrow{n \rightarrow \infty} V(S) \end{aligned}$$

- Cette convergence de $V(\theta_n')$ implique que $\forall \varepsilon > 0 \exists n(\varepsilon)$ tq

$$\forall p > n \geq n(\varepsilon) \quad |V(\theta_n') - V(\theta_p')| \leq \varepsilon$$

Donc par (A) en iterant en θ_n' et θ_p'

$$\begin{aligned} & V(\theta_{p+1}') - V(\theta_p') \geq \gamma_p \varepsilon \mathbb{1}_{\{\theta_{p+1}' \notin \mathcal{D}\}} + \gamma_p C \mathbb{1}_{\{\theta_{p+1}' \in \mathcal{D}\}} \\ \Rightarrow & V(\theta_n') - V(\theta_p') \geq \varepsilon \sum_{k=n+1}^p \gamma_k \mathbb{1}_{\{\theta_{k+1}' \notin \mathcal{D}\}} + C \sum_{k=n+1}^p \gamma_k \mathbb{1}_{\{\theta_{k+1}' \in \mathcal{D}\}} \\ \text{D'où } & \varepsilon \sum_{k=n+1}^p \gamma_k \mathbb{1}_{\{\theta_{k+1}' \notin \mathcal{D}\}} - C \sum_{k=n+1}^p \gamma_k \mathbb{1}_{\{\theta_{k+1}' \in \mathcal{D}\}} \leq \varepsilon \quad \textcircled{a} \end{aligned}$$

Soit p tq $p > n \geq n(\varepsilon)$ et $\varepsilon(\gamma_{n+1} + \gamma_{n+2} + \dots + \gamma_p) > \varepsilon$
 (possible car la série diverge)

Donc $\exists k \in \{n, \dots, p\}$ tq $\theta_{k+1}' \in \mathcal{D}$

Puisque $\|\theta_{k+1}' - \theta_n'\|$ de l'ordre de $\frac{\varepsilon}{\varepsilon}$ par \textcircled{a} (on n'a plus que
 les termes $\notin \mathcal{D}$ donc que le 1er membre $\leq \varepsilon$)

$d(\theta_n', \mathcal{D}) \leq \theta(\frac{\varepsilon}{\varepsilon})$. Puisque ε arbitraire, on peut
 choisir $\varepsilon = \varepsilon^2$ et donc $d(\theta_n', \mathcal{D}) \xrightarrow{n \rightarrow \infty} 0$

Comme le voisinage \mathcal{D} de S est aussi arbitraire, on obtient

$$d(\theta_n', S) \xrightarrow{n \rightarrow \infty} 0$$

$$\Rightarrow \lim_{n \rightarrow \infty} d(\theta_n, S) = 0.$$

Remarque: La condition d'un algorithme A stable est l'hypothèse de base de la convergence de $(\Theta_n)_{n \geq 0}$. Le problème vient du fait de devoir supposer que Θ_n reste dans un compact.

On peut montrer qu'en utilisant une suite de compacts croissants et en supposant que $(\Theta_n)_n$ ne tend pas vers l'infini, on revient à la notion de A-stabilité. Je vous renvoie vers le poly de Bernard Delyon (Stochastic approximation with decreasing gain: convergence and asymptotic theory, (2000)) Cela utilise de manière plus fine la fonction de Lyapunov.

La condition devient:

(B) : $h \in C^0(D)$; $\exists V$ positive (ou nulle) C^1 , $\exists K \subset D$ compact tq :

- (i) $V(x) \rightarrow +\infty$ si $x \rightarrow \partial D$ ou $|x| \rightarrow +\infty$
- (ii) $\langle \nabla V(x), h(x) \rangle < 0$ si $x \notin K$.

Et le théorème :

théorème : On suppose la condition (B) vérifiée. On suppose qu'il existe un compact $K_0 \subset D$ tq

$\Theta_n \in K_0$ infiniment souvent

On suppose que $\forall M \in \mathbb{N}$

$\lim_{n \rightarrow \infty} \sum_{i=1}^n \gamma_i \mathbb{1}_{\{V(\Theta_{i-1}) \leq M\}}$ existe et $\lim_{n \rightarrow \infty} |n| \mathbb{1}_{\{V(\Theta_{i-1}) \leq M\}} = 0$

Alors l'algorithme est A-stable.

II Dans le cas de bruit aléatoire : cas particulier des algo de Robbins-Monro :

On s'est pour l'instant concentré sur le cas déterministe ; maintenant on regarde un cas plus général où le bruit est aléatoire et donc la convergence sera soit ps soit en proba (on aura ps !)

Pour comprendre et parce qu'ils reflètent bien le cas général, on se concentre sur les algorithmes dits de Robbins-Monro qui s'écrivent

$$\Theta_n = \Theta_{n-1} + \gamma_n H(\Theta_{n-1}, Y_n) \quad \text{où } Y_n \sim P_{\Theta_{n-1}} \text{ et tq} \quad (8)$$

$$P(Y_n \in A \mid Y_{n-1}, Y_{n-2}, \dots, \Theta_0) = P_{\Theta_{n-1}}(Y_n \in A)$$

L'algorithme cherche la solution de $E_\Theta [H(\Theta, Y)] = 0 = h(\Theta)$

De nouveau on réécrit cet équation en faisant apparaître h :

$$\Theta_n = \Theta_{n-1} + \gamma_n h(\Theta_{n-1}) + \gamma_n e_n \quad (R1)$$

$$\text{avec } e_n = H(\Theta_{n-1}, Y_n) - h(\Theta_{n-1})$$

On pose les hypothèses suivantes:

$$(R_0): \sum \gamma_n = +\infty \quad \sum \gamma_n^2 < \infty$$

(R1): $h \in C^0$; $S = \{\Theta : h(\Theta) = 0\}$ est fini et $\exists V$ continuellement différentiable tq

$$\begin{cases} \lim_{z \rightarrow \partial \mathcal{O}} V(z) = \infty \quad \text{non nécessaire si } (\Theta_n) \text{ reste ds un compact.} \\ \langle \nabla V(\Theta), h(\Theta) \rangle \leq 0 \\ \{ \langle \nabla V(\Theta), h(\Theta) \rangle = 0 \} = S \end{cases}$$

$$(R_2): K \text{ compact de } \mathcal{O} \quad \sup_{\Theta \in K} E_\Theta [|H(\Theta, Y)|^2] < \infty$$

Alors

Théorème: sous les hypothèses $(R_0), (R_1)$ et (R_2) , l'algorithme $(R1)$ converge vers Θ^* tq $h(\Theta^*) = 0$ avec probabilité 1.

Preuve: On regarde le cas compact. $\Theta_n \in K_0 \quad \forall n \geq 0$

$$\text{Soit } X_n = \gamma_n e_n \quad \text{et } S_n = \sum_{k=1}^n X_k$$

$$\begin{aligned} E_\Theta [|X_n|^2 \mid \mathcal{F}_{n-1}] &= \gamma_n^2 E_\Theta [|e_n|^2 \mid \mathcal{F}_{n-1}] \\ &\leq \gamma_n^2 \sup_{\Theta \in K_0} E_\Theta [|H(\Theta, Y) - h(\Theta)|^2 \mid \mathcal{F}_{n-1}] \\ &\leq \gamma_n^2 C_{K_0} \quad \text{car } h \in C^0 \text{ et } H \text{ aussi } + (R_2) \end{aligned}$$

$$\text{Donc } E_\Theta [|X_n|^2] \leq \gamma_n C_{K_0}$$

Comme $\sum \gamma_n^2 < \infty$ la suite S_n converge ps (thm de Proba)

Et on retombe sur les conditions du précédent théorème ($n \rightarrow \infty$)

Maintenant: étape suivante: (Y_n) chaîne de Markov!!! La suite à la séance 10 où on aura vu et manipulé les CN.