

le gradient stochastique

1 Du déterministe à l'aléatoire:

On considère le problème suivant: $\min_{u \in U} J(u)$ où $\exists j$ tq j soit une fonction dépendant d'une variable aléatoire W définie sur $(\Omega, \mathcal{A}, \mathbb{P})$ et à valeurs dans \mathcal{W} muni d'une tribu et $J(u) = \mathbb{E}[j(u, W)]$.

le problème s'écrit donc trouver $\min_{u \in U} \mathbb{E}[j(u, W)]$ (1)

D'un point de vue purement théorique, on peut calculer, en tout point de u sous certaines conditions de régularité, les valeurs de $J(u)$ et de son gradient $\nabla_u J$. C'est un problème d'intégrale à paramètre. On se ramène donc à un problème de l'énumération pour lequel plusieurs méthodes de minimisations (gradient, Newton, etc) peuvent être appliquées.

Cependant d'un point de vue pratique, chaque évaluation de $J(u)$ et $\nabla_u J$ nécessite le calcul d'une espérance, ce qui peut être coûteux en tps de calcul surtout si W est de grande dimension.

→ Idée du gradient stochastique

L'idée est de faire évoluer simultanément le calcul de \mathbb{E} et la méthode de descente de gradient en s'appuyant sur la méthode de Monte Carlo qui donne une approximation de \mathbb{E} .

On tire aléatoirement une suite i.i.d de réalisations de W : (w_1, \dots, w_k, \dots) selon sa loi (celle sous laquelle est calculée \mathbb{E}). On génère une suite $(u_k)_k$ telle que $u_k \rightarrow u^*$ solution de (1). la mise à jour de u_k se fait avec une valeur simulée de w_k par un pas de gradient de la fonction j .

Il faut bien sûr s'assurer de la convergence de cet algorithme en particulier il faut vérifier que la variation en u est lente pour que la moyennisation en W ait un sens.

Voici l'algorithme en pseudo-code :

On se donne u_0

$k \rightarrow k+1$: on simule $w^{k+1} \sim P_w$ et $u_{k+1} = u_k - \varepsilon_k \nabla_{u_k} j(u_k, w_{k+1})$
où $(\varepsilon_k)_k$ suite $\searrow 0$

→ Interprétation intuitive :

L'idée est de profiter des itérations d'optimisation pour effectuer le calcul de l'espérance :

$$u_{k+1} = u_k - \varepsilon_k \nabla_{u_k} j(u_k, w_{k+1})$$

On somme k fois cette formule à partir d'un indice k_0 fixé

$$u_{k+k_0} = u_{k_0} - \sum_{\ell=0}^{k-1} \varepsilon_{k_0+\ell} \nabla_{u_{k_0+\ell}} j(u_{k_0+\ell}, w_{k_0+\ell+1})$$

Si on sait que $u \mapsto \nabla_u J(u, w)$ est "suffisamment régulière" et que $|u_{k_0+\ell} - u_{k_0+\ell'}|$ "petit" (à définir évidemment de manière plus rigoureuse plus tard) alors :

$$\begin{aligned} u_{k+k_0} &\approx u_{k_0} - \sum_{\ell=0}^{k-1} \varepsilon_{k_0+\ell} \nabla_u j(u_{k_0}, w_{k_0+\ell+1}) \\ &\approx u_{k_0} - \left(\sum_{\ell=0}^{k-1} \varepsilon_{k_0+\ell} \right) \nabla J(u_{k_0}) \end{aligned}$$

En utilisant le résultat ci-dessus (*)

Cela donne une étape de descente de gradient sur J . Le gradient stochastique utilise les itérations de l'optimisation pour reconstituer l'espérance du gradient du critère (et non l'espérance du critère lui-même).

→ Rq : s'adapte au cas sous contrainte $u \in U$ admissible convexe fermé $\neq \emptyset$ par projection

(*) Soit $(x_k)_k$ une suite dans un Hilbert X qui converge en moyenne vers μ :

$$\lim_{N \rightarrow \infty} \left(\frac{1}{N} \sum_{k=1}^N x_k \right) = \mu$$

convergence en moyenne

Soit $(p_k)_k$ une suite de réels positifs $\searrow 0$. On suppose que la suite $(\varepsilon_k)_k$ définie par $\varepsilon_k = k(p_k - p_{k+1}) \geq 0$ est tq ε_k soit le terme général d'une série divergente. Alors :

$$\lim_{N \rightarrow \infty} \frac{\sum_{k=1}^N p_k x_k}{\sum_{k=1}^N p_k} = \mu$$

Preuve du théorème de Césaro en moyenne:

On se ramène au cas $\mu = 0$ en translatant x_k par μ .

Soit $y_k = \frac{1}{k} \sum_{\ell=1}^k x_\ell$. Par hypothèse $y_k \xrightarrow[k \rightarrow \infty]{} \mu = 0$

Comme $p_k \rightarrow 0$ et $E_k = k(p_k - p_{k+1})$ on obtient par somme télescopique

$$p_k = \sum_{\ell=k}^{\infty} \frac{E_\ell}{\ell}$$

Soit $N \in \mathbb{N}^*$, $\forall k \leq N$

$$p_k^{(0)} = p_N^k + p_{N+1} \quad \text{avec} \quad p_N^k = \sum_{\ell=k}^N \frac{E_\ell}{\ell}$$

Avec ces notations:

$$\sum_{k=1}^N E_k y_k = \sum_{k=1}^N \frac{E_k}{k} \sum_{\ell=1}^k x_\ell = \sum_{k=1}^N x_k \sum_{k=\ell}^N \frac{E_k}{k} = \sum_{k=1}^N p_N^k x_k \quad (..)$$

Ceci est vrai $\forall (x_k)_k$ donc on prend le cas particulier $x_k = 1 \forall k$ et on obtient $\sum_{k=1}^N E_k = \sum_{k=1}^N p_N^k$. (....)

On obtient donc:

$$\left\| \frac{\sum_{k=1}^N p_k x_k}{\sum_{k=1}^N p_k} \right\| = \frac{\left\| \sum_{k=1}^N p_N^k x_k + p_{N+1} \sum_{k=1}^N x_k \right\|}{\sum_{k=1}^N p_N^k + N p_{N+1}} \quad \text{par (.)}$$

$$\leq \frac{\left\| \sum_{k=1}^N p_N^k x_k \right\|}{\sum_{k=1}^N p_N^k + N p_{N+1}} + \frac{\left\| p_{N+1} \sum_{k=1}^N x_k \right\|}{\sum_{k=1}^N p_N^k + N p_{N+1}} \quad (\text{inégalité triangulaire})$$
$$\leq \frac{\left\| \sum_{k=1}^N p_N^k x_k \right\|}{\sum_{k=1}^N p_N^k} + \frac{\left\| p_{N+1} \sum_{k=1}^N x_k \right\|}{N p_{N+1}} \quad (\text{car } p_k \geq 0 \forall k)$$

$$\leq \left\| \frac{\sum_{k=1}^N E_k y_k}{\sum_{k=1}^N E_k} \right\| + \underbrace{\frac{1}{N} \left\| \sum_{k=1}^N x_k \right\|}_{\rightarrow 0 \text{ par hypothèse}} \quad \text{par (..) et (....)}$$

Or E_k est le terme général d'une série divergente donc le théorème de Césaro classique implique que le 1^{er} terme de droite $\xrightarrow[N \rightarrow \infty]{} 0$

Soit $(u_n) \in E$ un R.V.E.C. normé, $\lambda_n > 0$ tq $\sum \lambda_n$ diverge. Si $u_n \rightarrow \ell$

$$\text{alors } \sigma_n = \frac{\sum \lambda_k u_k}{\sum \lambda_k} \rightarrow \ell$$

Le gradient stochastique:

On va ici voir le cadre probabiliste de cet algorithme et l'étude de ses propriétés, en particulier sa convergence.

1) algorithme:

Soit $(\Omega, \mathcal{F}, \mathbb{P})$ un espace de probabilité et w définie sur Ω à valeurs dans \mathcal{W} . On note μ la loi de w . Soit \mathcal{U} un Hilbert et \mathcal{U}_c une partie convexe fermée non vide de \mathcal{U} . Soit $j: \mathcal{U} \times \mathcal{W} \rightarrow \bar{\mathbb{R}}$ et $J = \mathbb{E}_\mu[j(\cdot, w)]$.

On cherche u^* tq $u^* = \underset{u \in \mathcal{U}_c}{\operatorname{argmin}} J(u)$

Sous les hypothèses classiques de convexité de J et de différentiabilité, si on se résout à calculer J et ∇J sur une trajectoire alors on peut utiliser l'algorithme de descente de gradient qui converge à vitesse $\frac{1}{k}$ (l'existence d'un ensemble $\mathcal{U}^* = \{\underset{u \in \mathcal{U}_c}{\operatorname{argmin}} J(u)\}$ étant supposé $\neq \emptyset$). Ici on essaie d'inter les calculs des espérances.

Algorithme en pseudo-code:

```
• On se donne  $u_0 \in \mathcal{U}_c$  ainsi qu'une suite de pas  $(\epsilon_k)_{k \geq 0} \geq 0$   
• Pour  $k=1$ : end  
2)  $w_{k+1} \sim \mu$  (Réalisation  $1^{te}$  de  $w_1, \dots, w_k$ )  
   Calculer  $\nabla_u j(u_k, w_{k+1})$   
    $u_{k+1} = \operatorname{proj}_{\mathcal{U}_c} [u_k - \epsilon_k \nabla_u j(u_k, w_{k+1})]$   
fin.
```

Rq1: On verra plus loin le critère d'arrêt (end); on se doute que ce ne pourra pas être comme pour la descente déterministe dû aux aléas.

Rq2: Il faut être capable de simuler $(w_k)_k$ la suite iid $\sim \mu$. On verra plus tard (cours 2-3) quelques méthodes de simulations et pas forcément iid mais avec des garanties quand même.

On va regarder un exemple qui illustre bien ce que l'on fait et présente en contexte un peu plus général aussi.

Exemple: Soit $W: \Omega \rightarrow \mathbb{R}$ l'... la question est calculer $E_\mu[W]$;

la manière connue pour approcher $E[W]$ est d'utiliser la formule de Monte Carlo:

$$E[W] \approx \frac{1}{k} \sum_{j=1}^k W_j \quad \text{où } W_j \sim \mu. \text{ i.i.d.}$$

On note $U_{k+1} = \frac{1}{k+1} \sum_{j=1}^{k+1} W_j$. On sait par la loi forte des grands nombres

que $U_k = \frac{1}{k} \sum_{j=1}^k W_j \xrightarrow[k \rightarrow \infty]{p.s.} E_\mu[W]$ quand $W_j \sim \mu$.

On va réécrire $(U_k)_k$ sous une autre forme:

$$\begin{aligned} U_{k+1} &= \frac{1}{k+1} \sum_{j=1}^k W_j + \frac{W_{k+1}}{k+1} = \frac{1}{k} \sum_{j=1}^k W_j - \frac{1}{k+1} \left(\frac{1}{k} \sum_{j=1}^k W_j - W_{k+1} \right) \\ &= U_k - \frac{1}{k+1} (U_k - W_{k+1}) \end{aligned}$$

Soit $E_k = \frac{1}{k+1}$ et $j(u, w) = \frac{1}{2} (u - w)^2$; on obtient alors que

$$U_{k+1} = U_k - E_k \nabla_u j(U_k, W_{k+1})$$

Or l'espérance d'une variable aléatoire est la valeur minimale du critère de dispersion d'un nuage de points: $E[W] = \operatorname{argmin}_{u \in \mathbb{R}} \frac{1}{2} (u - w)^2$

Donc la méthode d'approximation de l'espérance par somme de Monte Carlo est en fait un gradient stochastique sur la fonctionnelle de dispersion.

On peut avec cet exemple faire quelques remarques:

- On a une idée de la suite (E_k) dans ce cas précis $E_k = \frac{1}{k} \rightarrow 0$ mais pas trop vite... ($\sum \frac{1}{k} = +\infty$)
- la méthode de Monte Carlo converge p.s., ce qui nous laisse espérer de même au moins sous certaines conditions pour le gradient stochastique.
- Il existe aussi des TCL pour les sommes de MC donc peut-être aussi pour le GS.

Rappels de modes de convergence d'une s.a.:

• p.s.: $\forall \varepsilon > 0 \quad \lim_{k \rightarrow \infty} P\left(\sup_{n \geq k} \|U_n - U^*\| > \varepsilon\right) = 0$

• en P: $\forall \varepsilon > 0 \quad \lim_{k \rightarrow \infty} P(\|U_k - U^*\| > \varepsilon) = 0$

• en E: $\lim_{k \rightarrow \infty} E[\|U_k - U^*\|] = 0$

))
cette dernière
))

2) Résultats de Convergence:

On va énoncer et démontrer un thm de convergence en moyenne quadratique car simple à prouver. Mais les conditions sont très restrictives donc on énoncera les thm généraux de convergence pour les Approximations stochastiques qui laisseront voir ce que l'on espère convergence ps et TCL.

H1 la variable aléatoire $j(u, w) : \mathcal{X} \rightarrow \mathbb{R}$ est mesurable et son espérance existe $\forall u \in \mathcal{U}_c$.

H2 la fonction $u \mapsto j(u, w)$ est convexe, semi C^0 inférieurement, à valeur dans $\bar{\mathbb{R}}$, différentiable $\forall w \in W$

H3 $\exists m > 0$ tq $\forall u \in \mathcal{U}_c, \forall w \in W \quad \|\nabla_u j(u, w)\| \leq m$

H4 Le probl $\min_{u \in \mathcal{U}_c} J(u)$ admet un ensemble non vide de solutions \mathcal{U}^* qui vérifie la relation: $\forall u \in \mathcal{U}_c \quad J(u) - J^* \geq c \cdot \text{dist}_{\mathcal{U}^*}(u)^2$
où $J^* = \min_{u \in \mathcal{U}_c} J(u)$ et $\text{dist}_{\mathcal{U}^*}$ est la fonction distance à \mathcal{U}^*

H5 (ε_k) vérifie $\varepsilon_k > 0 \forall k$ et $\sum_{k \in \mathbb{N}} \varepsilon_k = +\infty$, $\sum_{k \in \mathbb{N}} \varepsilon_k^2 < \infty$

Théorème de convergence en moyenne quadratique:

Sous les hypothèses **(H1-5)**, la suite $(U_k)_k$ de v.a générée par le GS converge en moyenne quadratique vers \mathcal{U}^* :

$$\lim_{k \rightarrow \infty} \mathbb{E} [\text{dist}_{\mathcal{U}^*}(U_k)^2] = 0$$

Quelques remarques avant la preuve:

- convergence L^2 assez faible...
- H3 très restrictive: peu ex ne prend pas en cpte le cas des sommes de N.C
- H5: très très courante et se retrouve dans presque tous les thm de convergence des approximations stochastiques qui sont des algorithmes s'écrivant sous la forme $\Theta_n = \Theta_{n-1} + \gamma_n \cdot h(\Theta_{n-1}) + \gamma_n \cdot \eta_n$
où $h(\theta) = \mathbb{E}[H(\theta, W)]$ et la fonction que l'on veut annuler et $\eta_n = H(\Theta_{n-1}, X_n) - h(\Theta_{n-1})$ joue le rôle d'une perturbation aléatoire

Preuve du thm de convergence :

U^* étant un convexe fermé (car j convexe c'est un convexe fermé), la projection sur U^* est bien définie. Soit $(u_k)_k$ la suite produite par l'algorithme (GS) avec une réalisation de la suite $(w_k)_k$. Soit $\bar{u}_k = \text{proj}_{U^*}(u_k)$. On a alors $\text{dist}_{U^*}(u_k)^2 = \|u_k - \bar{u}_k\|^2$. Soit $d_k = \text{dist}_{U^*}(u_k)^2$.

$$d_{k+1} = \|u_{k+1} - \bar{u}_{k+1}\|^2 \leq \|u_{k+1} - \bar{u}_k\|^2 \quad \text{car } d_{k+1} \text{ est la distance minimale de } u_{k+1} \text{ à } U^*.$$

$$\leq \|\text{proj}_{U^*}(u_k - \varepsilon_k \nabla_u j(u_k, w_{k+1})) - \bar{u}_k\|^2 \quad \text{par def de } u_{k+1}$$

$$\leq \|u_k - \varepsilon_k \nabla_u j(u_k, w_{k+1}) - \bar{u}_k\|^2 \quad \text{car la projection est contractante}$$

$$\leq d_k + \varepsilon_k^2 m^2 - 2\varepsilon_k \langle u_k - \bar{u}_k ; \nabla_u j(u_k, w_{k+1}) \rangle \quad \text{developpement du et } \nabla j \text{ borné}$$

On prend l'espérance conditionnelle sachant la filtration \mathcal{F}_k

$$\mathbb{E}[d_{k+1} | \mathcal{F}_k] \leq \mathbb{E}[d_k | \mathcal{F}_k] + \varepsilon_k^2 m^2 - 2\varepsilon_k \mathbb{E}[\langle u_k - \bar{u}_k ; \nabla_u j(u_k, w_{k+1}) \rangle | \mathcal{F}_k]$$

$$\leq d_k + \varepsilon_k^2 m^2 - 2\varepsilon_k \langle u_k - \bar{u}_k ; \mathbb{E}[\nabla_u j(u_k, w_{k+1}) | \mathcal{F}_k] \rangle$$

Comme j est différentiable de gradient en u borné :

$$\mathbb{E}[\nabla_u j(u_k, w_{k+1}) | \mathcal{F}_k] = \nabla_u \mathbb{E}[j(u_k, w_{k+1}) | \mathcal{F}_k]$$

$$= \nabla_u J(u_k)$$

$$\text{D'où } \mathbb{E}[d_{k+1} | \mathcal{F}_k] \leq d_k + \varepsilon_k^2 m^2 - 2\varepsilon_k \langle u_k - \bar{u}_k ; \nabla_u J(u_k) \rangle$$

$$\leq d_k + \varepsilon_k^2 m^2 - 2\varepsilon_k (J(u_k) - J^*) \quad \text{par les propriétés de convexité}$$

$$\leq d_k (1 - 2\varepsilon_k c) + \varepsilon_k^2 m^2 \quad \text{par H}_4$$

En reprenant \mathbb{E} on obtient,

$$\mathbb{E}[d_{k+1}] \leq (1 - 2\varepsilon_k c) \mathbb{E}[d_k] + \varepsilon_k^2 m^2$$

On montre alors par récurrence que $\forall k$ suffisamment grand: $\forall n \in \mathbb{N}^*$

$$\mathbb{E}[d_{k+n+1}] \leq \left[\prod_{\ell=0}^n (1 - 2\varepsilon_{k+\ell} c) \right] \mathbb{E}[d_k] + \left(\sum_{\ell=0}^n \varepsilon_{k+\ell}^2 \right) m^2 \quad \text{CRPC}$$

$\rightarrow 0 \text{ (1)}$ $\rightarrow 0 \text{ car } \sum_{\ell \in \mathbb{N}} \varepsilon_{k+\ell}^2 < \infty$

(1) car soit k_0 tq $\forall \ell \geq k_0$ $0 < (1 - 2\varepsilon_{k+\ell} c) \leq 1$

existe car $\varepsilon_k \rightarrow 0$

$$\text{donc } p_k = \prod_{\ell=1}^k (1 - 2\varepsilon_{k-\ell} c) = C \prod_{\ell=k_0}^k (1 - 2\varepsilon_{k-\ell} c) = C \gamma_k. \quad \forall k > 0 \quad \gamma_k > 0$$

donc convergence. De plus $\log \gamma_k = \sum_{\ell=k_0}^k \log(1 - 2\varepsilon_{k-\ell} c) \leq -2c \sum_{\ell=k_0}^k \varepsilon_{k-\ell}$

$\rightarrow -\infty$

$$\Rightarrow \gamma_k \rightarrow 0$$

Théorème (Vitesse de convergence)

Sous les mêmes hypothèses que précédemment, on choisit $\varepsilon_k = \frac{1}{ck + \frac{m^2}{cd^0}}$

avec $d_0 = \text{dist}_{\mathcal{H}}(u_0)^2$ on obtient la borne de la vitesse de convergence

$$\mathbb{E}[\text{dist}_{\mathcal{H}}(u_k)^2] \leq \frac{1}{\frac{c^2}{m^2}k + \frac{1}{d_0}} \quad \forall k \in \mathbb{N}$$

Preuve: On repart de l'inégalité suivante:

$$\mathbb{E}[d_{k+1}] \leq (1 - 2\varepsilon_k c) \mathbb{E}[d_k] + \varepsilon_k^2 m^2$$

Si $\varepsilon_k = \frac{\gamma}{\alpha k + \beta}$ alors on montre (par récurrence) que si $\alpha = \frac{c^2}{m^2}$, $\beta = \frac{1}{d_0}$

$$\text{et } \gamma = \frac{c}{m^2} : (\alpha k + \beta) \mathbb{E}[d_k] \leq 1 \quad \text{c.q.p.d.}$$

3) Considérations pratiques de l'algorithme (GS):

- Critère d'arrêt: * Ne peut être, comme pour la GD , basé sur $\|u_{k+1} - u_k\|$ car du fait que $\varepsilon_k \rightarrow 0$ et $|D_j| \leq m$ cette différence $\rightarrow 0$ par construction.

* pas non plus sur D_j qui n'est pas informatif sur ∇J .

* une approximation de $\mathbb{E}[D_j]$ peut être utilisée

* le plus souvent un nb fixé par l'utilisateur.

- suite ε_k : en théorie pour avoir la plus grande vitesse asymptotique il faut $\varepsilon_k = \frac{1}{k}$. Mais en pratique (et pour le TCC aussi) $\frac{1}{\varepsilon^\alpha}$, $\alpha \in]\frac{1}{2}, 1]$, est parfois meilleur pour avoir un poids assez important de D_j % à la précédente valeur. De plus, on peut aussi avoir besoin de ce qu'on appelle un laps de chauffe pour retirer les 1^{ères} valeurs aberrantes.

On va maintenant regarder un cadre un peu plus général qui nous donnera des théorèmes de convergence plus forts (on attend p.s.) et en TCC.