

Algorithmes de gradient stochastiques moyennés

A. Godichon-Baggioni

Algorithme de gradient stochastique moyenné

DÉFINITION

Algorithme moyenné :

$$\bar{m}_n = \frac{1}{n+1} \sum_{k=0}^n m_k$$

où les m_k sont les estimateurs de gradient stochastique.

Ecriture récursive :

$$m_{n+1} = m_n - \gamma_{n+1} \nabla_{hg}(X_{n+1}, m_n)$$
$$\bar{m}_{n+1} = \bar{m}_n + \frac{1}{n+2} (m_{n+1} - \bar{m}_n).$$

avec $\gamma_n = c_\gamma n^{-\alpha}$ et $\alpha \in (1/2, 1)$.

LEMME DE TOEPLITZ

Lemme

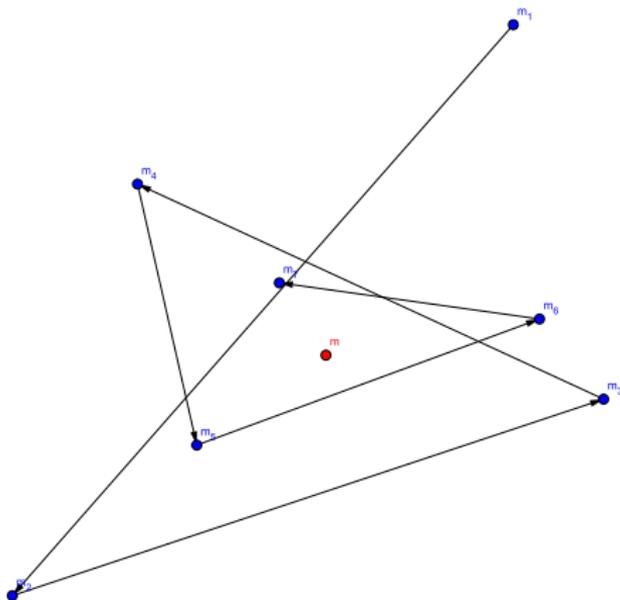
Soit (a_n) positive telle que $\sum_{n \geq 0} a_n = +\infty$ et X_n une suite de variables aléatoires convergeant presque sûrement vers X . Alors

$$\frac{1}{\sum_{k=0}^n a_k} \sum_{k=0}^n a_k X_k \xrightarrow[n \rightarrow +\infty]{p.s.} X.$$

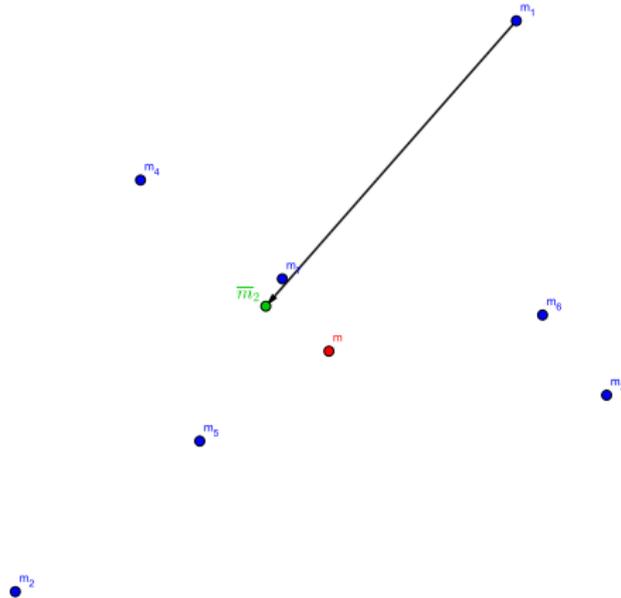
Application :

$$m_n \xrightarrow[n \rightarrow +\infty]{p.s.} m \quad \implies \quad \bar{m}_n \xrightarrow[n \rightarrow +\infty]{p.s.} m$$

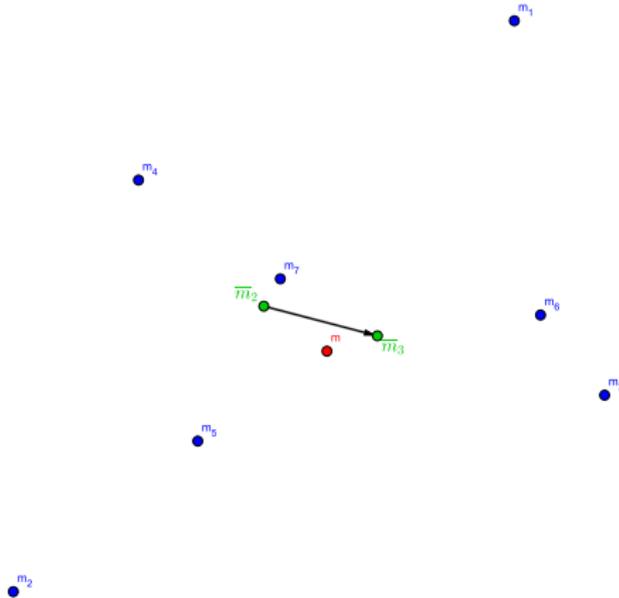
COMMENT ÇA MARCHE ?



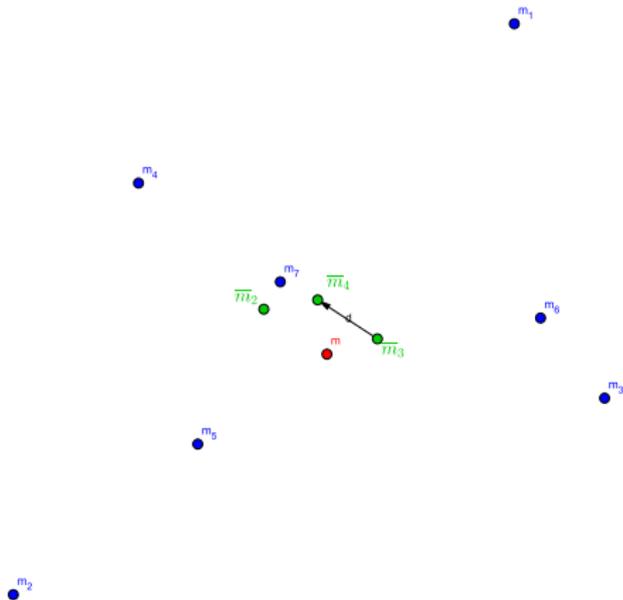
COMMENT ÇA MARCHE ?



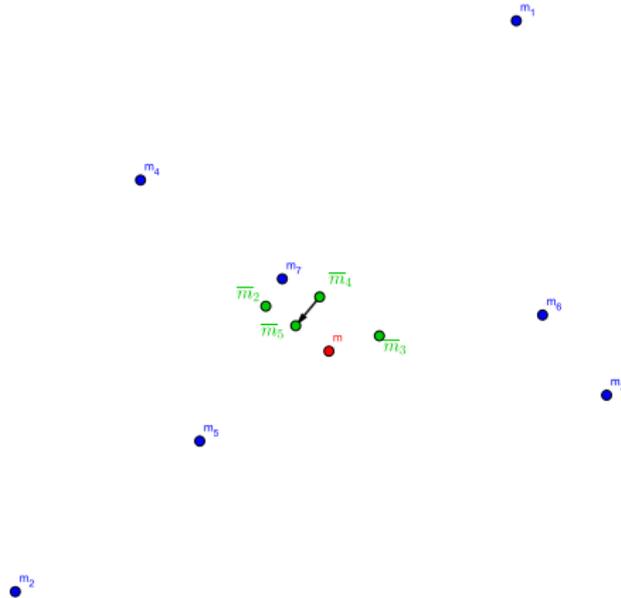
COMMENT ÇA MARCHE ?



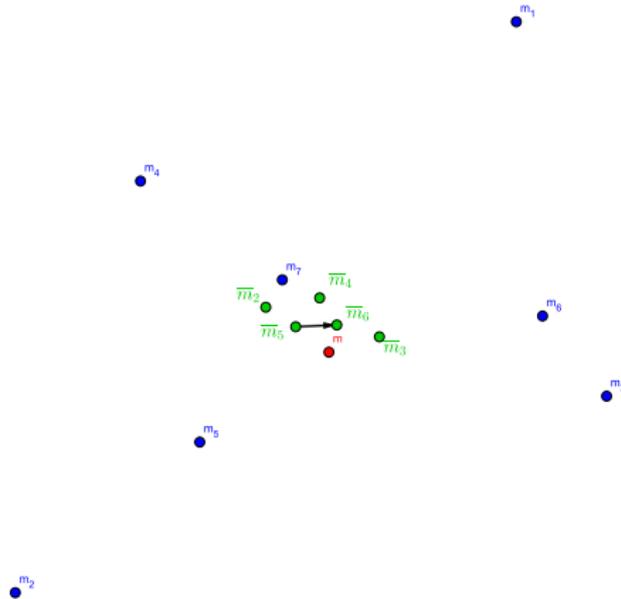
COMMENT ÇA MARCHE ?



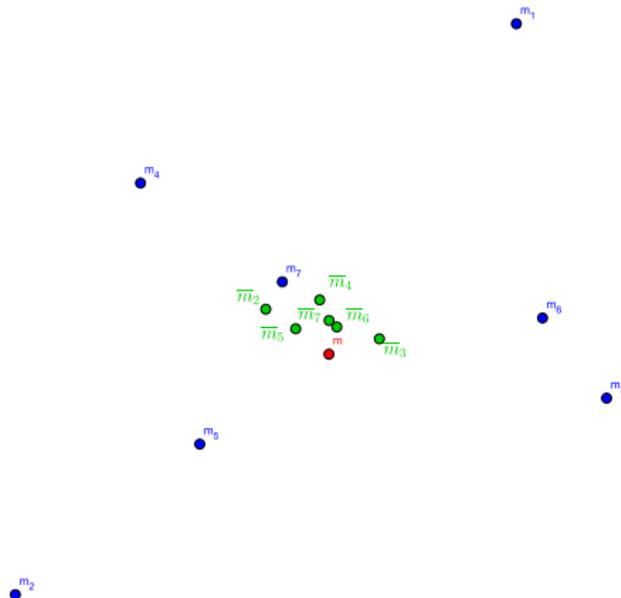
COMMENT ÇA MARCHE ?



COMMENT ÇA MARCHE ?



COMMENT ÇA MARCHE ?



Vitesse de convergence

UN COROLLAIRE DU LEMME DE TOEPLITZ

Corollaire

Soit (X_n) une suite de variables aléatoires positives et (a_n) une suite positive telles que

$$X_n = o(a_n) \quad p.s.$$

Alors

$$\sum_{k=1}^n X_k = O\left(\sum_{k=1}^n a_k\right) \quad a.s.$$

CADRE (GRADIENT STOCHASTIQUE)

(PS1) Il existe $\eta > \frac{1}{\alpha} - 1$ et $C_\eta \geq 0$ tels que

$$\mathbb{E} \left[\|\nabla_h g(X, h)\|^{2+2\eta} \right] \leq C_\eta \left(1 + \|h - m\|^{2+2\eta} \right)$$

(PS2) G est deux fois continûment différentiable et

$$\lambda_{\min} := \lambda_{\min} (\nabla^2 G(m)) > 0.$$

$$\mathbf{(PS1)} \text{ et } \mathbf{(PS2)} \implies \|m_n - m\|^2 = O\left(\frac{\ln n}{n^\alpha}\right) \quad p.s.$$

CADRE

(PS3) Il existe $\eta > 0$ et $C_\eta \geq 0$ t.q pour tout $h \in \mathcal{B}_\eta := \mathcal{B}(m, \eta)$,

$$\|\nabla G(h) - \nabla^2 G(m)(h - m)\| \leq C_\eta \|h - m\|^2$$

L'hypothèse **(PS3)** est vérifiée si $\nabla^2 G(\cdot)$ est Lipschitz sur \mathcal{B}_η .

VITESSE DE CONVERGENCE

Théorème

On suppose que les hypothèses (PS1) à (PS3) sont vérifiées. Alors, pour tout $\delta > 0$,

$$\|\bar{m}_n - m\|^2 = o\left(\frac{(\ln n)^{1+\delta}}{n}\right) \quad p.s.$$

PREUVE

La preuve repose sur le résultat suivant :

Théorème

Soit (ξ_k) une suite de différences de martingale telle que

$\mathbb{E} \left[\|\xi_k\|^2 \mid \mathcal{F}_{k-1} \right] \leq C$. Alors, pour tout $\delta > 0$,

$$\left\| \sum_{k=1}^n \xi_k \right\|^2 = o(n(\ln n)^{1+\delta}) \quad p.s.$$

EFFICACITÉ ASYMPTOTIQUE

(PS4) La fonction $\Sigma : \mathbb{R}^d \rightarrow \mathcal{M}_d(\mathbb{R})$ définie par

$$\Sigma(h) = \mathbb{E} [\nabla_h g(X, h) \nabla_h g(X, h)^T]$$

est continue en m .

Théorème

*On suppose que les hypothèses **(PS1)** à **(PS4)** sont vérifiées. Alors*

$$\sqrt{n} (\bar{m}_n - m) \xrightarrow[n \rightarrow +\infty]{\mathcal{L}} \mathcal{N} (0, H^{-1} \Sigma H^{-1})$$

avec $H = \nabla^2 G(m)$ et $\Sigma = \Sigma(m)$.

Régression linéaire

L'ALGORITHME

Algorithme moyenné :

$$\theta_{n+1} = \theta_n + \gamma_{n+1} (Y_{n+1} - X_{n+1}^T \theta_n) X_{n+1}$$

$$\bar{\theta}_{n+1} = \bar{\theta}_n + \frac{1}{n+2} (\theta_{n+1} - \bar{\theta}_n)$$

avec $\bar{\theta}_0 = \theta_0$.

VITESSE DE CONVERGENCE

Théorème

On suppose qu'il existe $\eta > \frac{1}{\alpha} - 1$ tel que X et ϵ admettent des moments d'ordre $4 + 4\eta$ et $2 + 2\eta$. Alors pour tout $\delta > 0$,

$$\|\bar{\theta}_n - \theta\|^2 = o\left(\frac{(\ln n)^{1+\delta}}{n}\right) p.s. \quad \text{et} \quad \sqrt{n}(\bar{\theta}_n - \theta) \xrightarrow[n \rightarrow +\infty]{\mathcal{L}} \mathcal{N}(0, \sigma^2 H^{-1})$$

SIMULATIONS

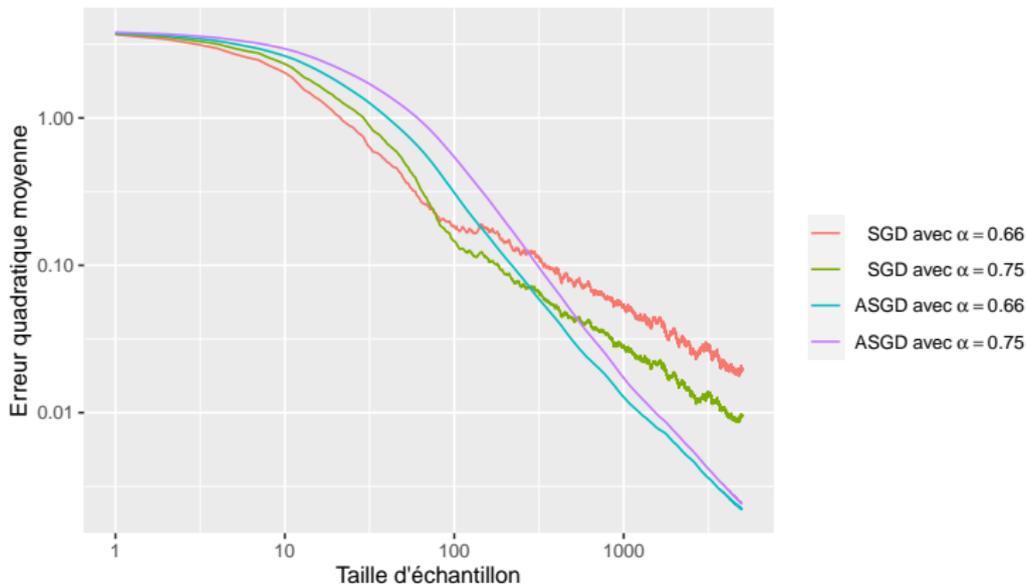


FIGURE – Evolution de l’erreur quadratique moyenne de l’estimateur de gradient θ_n (SGD) et de sa version moyennée $\bar{\theta}_n$ (ASGD) en fonction de la taille d’échantillon n dans le cadre de la régression linéaire.

TESTER $H_0 : \theta = \theta_0$ "EN LIGNE"

Réécriture du TLC : Sous H_0 ,

$$\sqrt{n} \frac{(\bar{\theta}_n - \theta_0)^T H (\bar{\theta}_n - \theta_0)}{\sigma^2} \xrightarrow[n \rightarrow +\infty]{\mathcal{L}} \chi_d^2$$

Application : Soit \bar{H}_n et $\hat{\sigma}_n^2$ des estimateurs consistants. Alors

$$C_n := \sqrt{n} \frac{(\bar{\theta}_n - \theta_0)^T \bar{H}_n (\bar{\theta}_n - \theta_0)}{\hat{\sigma}_n^2} \xrightarrow[n \rightarrow +\infty]{\mathcal{L}} \chi_d^2$$

SIMULATIONS

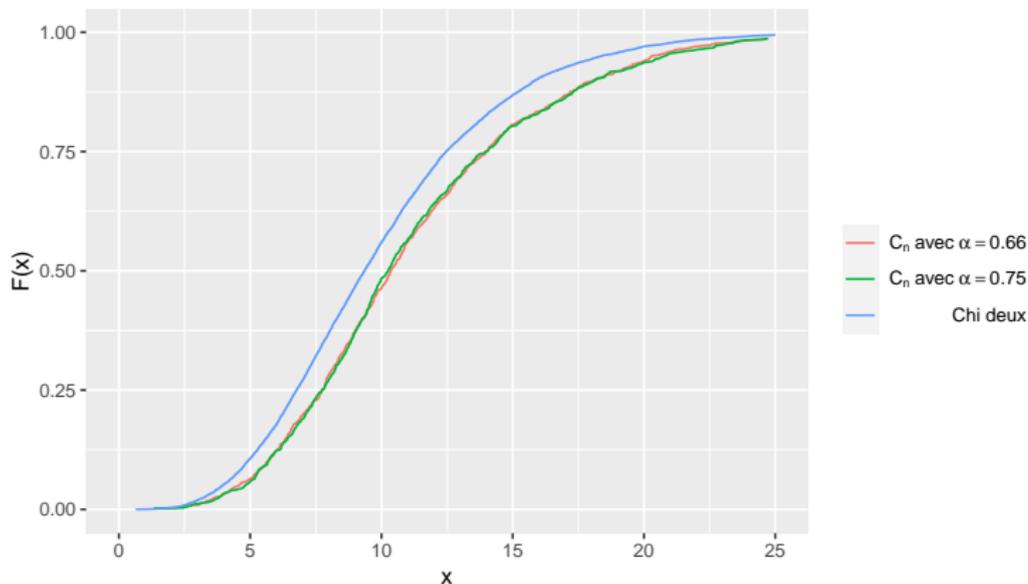


FIGURE – Comparaison de la fonction de répartition de C_n avec $n = 5000$, pour $\alpha = 0.66$ et $\alpha = 0.75$, et de celle d'une Chi 2 à 10 degrés de liberté dans le cadre du modèle linéaire.

Régression logistique

L'ALGORITHME

Algorithme moyenné :

$$\theta_{n+1} = \theta_n + \gamma_{n+1} (Y_{n+1} - \pi(X_{n+1}^T \theta_n)) X_{n+1}$$

$$\bar{\theta}_{n+1} = \bar{\theta}_n + \frac{1}{n+2} (\theta_{n+1} - \bar{\theta}_n)$$

avec $\bar{\theta}_0 = \theta_0$ et $\pi(x) = \frac{e^x}{1+e^x}$.

VITESSE DE CONVERGENCE

Théorème

On suppose que X admet un moment d'ordre 4. Alors pour tout $\delta > 0$,

$$\|\bar{\theta}_n - \theta\|^2 = o\left(\frac{(\ln n)^{1+\delta}}{n}\right) p.s. \quad \text{et} \quad \sqrt{n}(\bar{\theta}_n - \theta) \xrightarrow[n \rightarrow +\infty]{\mathcal{L}} \mathcal{N}(0, H^{-1})$$

SIMULATIONS

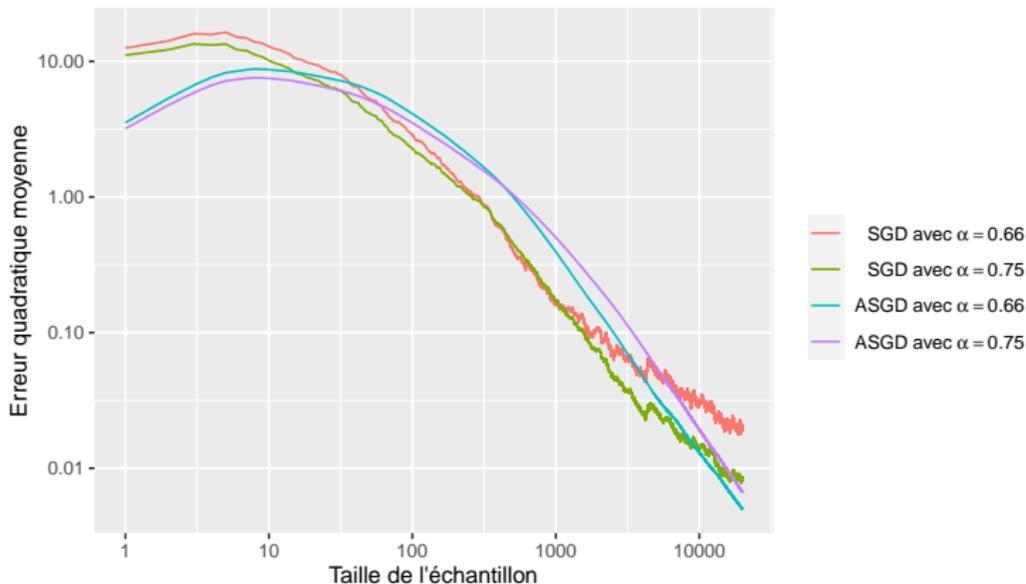


FIGURE – Evolution de l’erreur quadratique moyenne par rapport à la taille de l’échantillon des estimateurs de gradients θ_n (SGD) et de leurs versions moyennées $\bar{\theta}_n$ (ASGD) dans le cadre de la régression logistique.

TESTER $H_0 : \theta = \theta_0$ "EN LIGNE"

Réécriture du TLC : Sous H_0 ,

$$\sqrt{n} (\bar{\theta}_n - \theta_0)^T H (\bar{\theta}_n - \theta_0) \xrightarrow[n \rightarrow +\infty]{\mathcal{L}} \chi_d^2$$

Application : Soit \bar{H}_n un estimateur consistant de H . Alors

$$C_n := \sqrt{n} (\bar{\theta}_n - \theta_0)^T \bar{H}_n (\bar{\theta}_n - \theta_0) \xrightarrow[n \rightarrow +\infty]{\mathcal{L}} \chi_d^2$$

SIMULATIONS

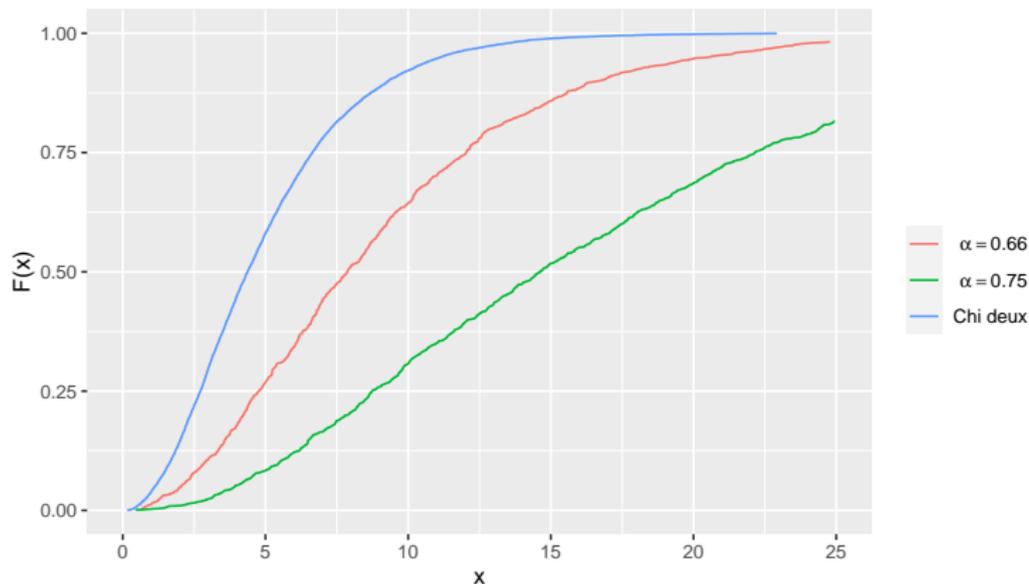


FIGURE – Comparaison de la fonction de répartition de C_n , avec $n = 20000$ et $\alpha = 0.66$ ou $\alpha = 0.75$, et de la fonction de répartition d'une Chi deux à 5 degrés de liberté dans le cadre de la régression logistique.

MOYENNÉ PONDÉRÉ

$$\bar{m}_n = \frac{1}{\sum_{k=1}^n \log(k+1)^w} \sum_{k=1}^n \log(k+1)^w m_k$$

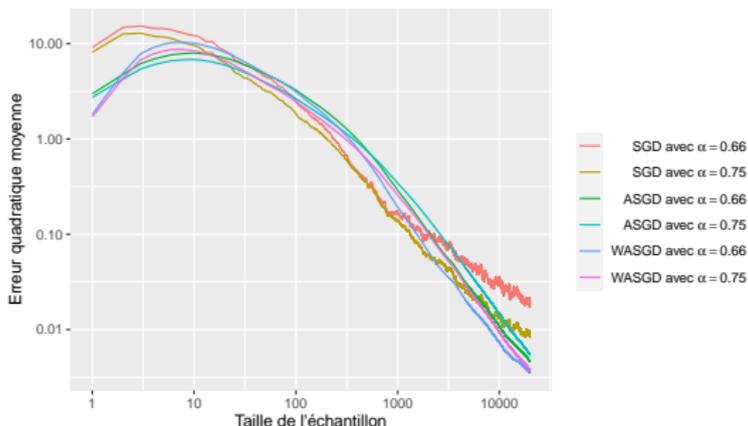


FIGURE – Evolution de l'erreur quadratique moyenne par rapport à la taille de l'échantillon des estimateurs de gradients θ_n (SGD) et de leurs versions moyennées $\bar{\theta}_n$ (ASGD) dans le cadre de la régression logistique.

EXERCICE

► Modèle linéaire :

On prend $\theta = (-2, -1, 0, 1, 2)^T$, $X_i \sim \mathcal{N}(0, I_5)$ et $\epsilon_i \sim \mathcal{N}(0, 1)$.

1. Générer un échantillon $(X_i, Y_i)_{i=1, \dots, n}$ avec $n = 10000$.
2. Ecrire une fonction qui ressorte l'ensemble de tous les estimateurs $\theta_i, i = 0, \dots, n$ et de leurs versions moyennées $\bar{\theta}_i$.
3. Tracer les courbes des évolutions des erreurs quadratiques.
4. Générer maintenant 50 échantillons et tracer les courbes des évolutions de des erreurs quadratiques moyennes.

► Faire de même pour la régression logistique.

- Revenir à l'exemple de la régression linéaire mais en prenant $X \sim \mathcal{N}(0, D)$ avec $D = \text{diag}(10^{-2}, 10^{-1}, 1, 10, 10^2)$. Regarder les évolutions des erreurs quadratiques moyennes pour les estimateurs de gradient stochastique et leurs versions moyennées.