

The heavy-tail phenomenon in SGD.

Fonction à minimiser, sur \mathbb{R}^d :

$$f(x) = \frac{1}{n} \sum_{i=1}^m f_i(x)$$

f_i différentiables
L-smooth

μ -convexes / μ -PL.

Minibatch-SGD:

$$x_{n+1} = x_n - \eta \times \widehat{\nabla f}(x_n) \quad \text{step-size } \eta.$$

avec $\widehat{\nabla f}(x_n) = \frac{1}{b} \sum_{i \in B} \nabla f_i(x_n).$

et $B \sim \text{Unif}(\{B \subseteq [m], |B|=b\}).$

Théorème Si les f_i sont μ -fortement convexes et L-smooth,
alors $\forall \eta < \frac{1}{2L}$

$$\mathbb{E}[|x_n - x_*|^2] \leq (1 - \eta\mu)^n |x_0 - x_*|^2 + \frac{2\eta}{\mu} \sigma^2$$

où σ^2 dépend de f et b : $\sigma^2 = \text{Var}(\widehat{\nabla f}(x_n)).$

En particulier, $(x_n - x_*)_n$ est une suite L^2 -bornée
donc tendue: quitte à extraire on peut supposer
que $x_n \xrightarrow[n \rightarrow \infty]{\text{loi}} x_*$

Que dire de la loi limite et de sa dépendance
en η, b, L, μ ?

(2)

Écrivons $\nabla f_i(x) \approx \nabla f_i(x_*) + \nabla^2 f_i(x_*) (x_n - x_*)$
 autour de x_* . Si x_* est un minimum de
 chaque f_i (interpolation) alors $\nabla f_i(x_*) = 0$.

On a donc :

$$x_{n+1} = x_n - \frac{\eta}{b} \left(\sum_{i \in B_n} \nabla f_i(x_n) \right)$$

$$\approx x_n - \frac{\eta}{b} \sum_{i \in B_n} \nabla^2 f_i(x_*) x_n + \frac{\eta}{b} \sum_{i \in B_n} \nabla^2 f_i(x_*) x_*$$

$$\approx \left(\text{Id} - \frac{\eta}{b} \sum_{i \in B_n} \nabla^2 f_i(x_*) \right) x_n + \left(\frac{\eta}{b} \sum_{i \in B_n} \nabla^2 f_i(x_*) x_* \right)$$

$$= A_n$$

matrice aléatoire
de taille $d \times d$

vecteur
aléatoire de
taille d .

$$= B_n.$$

On a donc $x_{n+1} \approx A_n x_n + B_n$. (1)

(1) est exacte dans le cas où les f_i sont toutes
 quadratiques. C'est le cas de la régression
 linéaire classique : training data (z_i, y_i)

avec $y_i = \langle x_*, z_i \rangle + \varepsilon$ $z_i \in \mathbb{R}^d$ $y_i \in \mathbb{R}$ $x \in \mathbb{R}^d$

loss : $f(x) = \frac{1}{2} \sum_{i=1}^n \underbrace{(\langle x, z_i \rangle - y_i)^2}_{f_i(x)}$

$$\nabla f_i(x) = (\langle x, z_i \rangle - y_i) z_i$$

$$\nabla^2 f_i(x) = z_i z_i^*.$$

La récursion (1) s'écrit alors

$$x_{n+1} = \left(I - \frac{\eta}{b} \sum_{i \in B_n} z_i z_i^T \right) x_n + \frac{\eta}{b} \sum_{i \in B_n} z_i z_i^T x^*$$

On suppose qu'on est dans un cadre de SGD online : les butchs (y_i, z_i) sont générés à chaque étape et jamais réutilisés.

hypothèses

$$z_i \sim N(0, I)$$

dimension 1

$$x_0 = 1$$

~~$x_0 = 1$~~

$$x_{n+1} = \left(1 - \frac{\eta}{b} K_n^2 \right) x_n + \frac{\eta}{b} K_n^2$$

$$\text{où } K_n^2 \sim \chi^2(b), \text{ iid de } (x_n).$$

Récursions affines et queue de proba.

On considère l'équation en loi $\boxed{X \stackrel{\text{loi}}{=} AX + B.} \quad (2)$

Si la récursion $x_{n+1} = A_n x_n + B_n$ converge en loi c'est forcément vers une solution de (2)

~~Plus~~ On supposera toujours que (AB) a une densité contre lebesgue. Par simplicité on suppose aussi que $A, B > 0$ mais les résultats s'étendent à tous les cas.

lemme 1 $\times \lim_{n \rightarrow \infty} \sum_{k=0}^n A_n A_{n-1} \dots A_{n-k} B_k$

(4)

la série converge si $E \ln|A| = \gamma < 0$.

Théorème (Kesten 1973-75)

(1) On suppose qu'il existe $s > 0$ tel que

(i) $E A^s = 1$ (ii) $E A^s \ln A < \infty$ (iii) $E B^s < \infty$.

Alors $\exists c > 0$ tq $P(X > x) \sim \frac{c}{x^s}$.

(2) Si A peut prendre des valeurs négatives, le même résultat est vrai.

(3) En plusieurs dimensions on pose $h(s) = \lim_{n \rightarrow \infty} (E \|A_n \dots A_1\|^s)^{\frac{1}{n}}$.

Si $\exists s$ tq $h(s) = 1$, $E \|A\|^s \ln \|A\| < \infty$, $E \|B\|^s < \infty$
alors $\forall u \in \mathbb{R}^n$, $u \neq 0$ $\exists c(u) > 0$ tq

$P(\langle u, X \rangle > x) \sim \frac{c(u)}{x^s}$.

Application K_n a tous ses moments exponentiels
donc $E |K_n|^a < \infty \forall a$.

$s \mapsto E |M_n|^s$ est str. \uparrow donc $\exists ! s$ tq
 $P(X_\infty > x) \sim c x^{-s}$.

En plus $s = f(\underline{a}, \underline{b}, \underline{d})$

s grand \Rightarrow queue fine
 s petit \Rightarrow queue lourde

What happens for FC networks / Deep Nets? (5)

$$f(\theta, x) = \sum_{i=1}^K \theta_i^1 \varphi(\langle \theta_i^2, x \rangle + \theta_i^3)$$

\hookrightarrow activation.

$$\theta_i^1 \in \mathbb{R}$$

$$\theta_i^2 \in \mathbb{R}^D$$

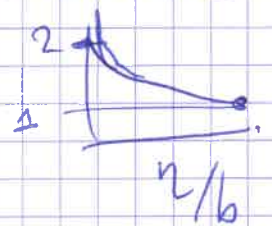
$$\theta_i^3 \in \mathbb{R} \text{ bias}$$

Idea when $K \rightarrow \infty$ the θ_i^2 become independent.

Empirically their distribution seems to become heavy-tailed.

SURTOUT pour les couches extrêmes les couches du milieu sont HT.

$$\alpha = f(\eta/b)$$



Conséquences sur la Modélisation

$$(SGD) \quad x_{n+1} - x_n = \eta \widehat{\nabla f}(x_n) \iff dx_t = -\nabla f(x_t) dt + \underset{\text{noise}}{\sum_t \epsilon_t}$$

$$\text{Langevin: } x_t \xrightarrow[n \rightarrow \infty]{\text{loi}} X \sim \frac{e^{-f}}{Z}$$

Si f a une croissance linéaire, e^{-f}/Z a des queues de distribution fines.

$$\text{Conclusion (SGD)} \iff dx_t = -\nabla f(x_t) dt + dL_t^\alpha$$

où L^α est un processus de Lévy.

Ultime remarque

Dans le cas où les x_i sont isotropiques :

$\exists \eta_1 < \eta_2$ tq :

$$\eta < \eta_1 \Rightarrow S > 2$$

$$\eta \in (\eta_1, \eta_2) \Rightarrow S \equiv 2$$

$$\eta > \eta_2 \Rightarrow S \in (0, 2)$$