Deviation inequalities for Markov chains, with applications to SGD and empirical risk minimization

Pierre Alquier





High-Dimensional Statistical Modeling Team Seminar March 1st, 2022

Co-authors



Fan, X. and Alquier, P. and Doukhan, P. (2021). Deviation inequalities for stochastic approximation by averaging. Preprint arXiv:2102.08685.



Xiequan Fan

Tianjin University



Paul Doukhan

CY Cergy Paris Université

Objective

General problem in probability and statistics

$$\mathbb{P}\left\{\left|\frac{1}{n}\sum_{i=1}^{n}X_{i}-\frac{1}{n}\mathbb{E}\left(\sum_{i=1}^{n}X_{i}\right)\right|\geq x\right\}\leq ?$$

What can we expect? (1/2)

Chebyshev's inequality

$$\mathbb{P}\Big\{|U-\mathbb{E}(U)|\geq x\Big\}\leq \frac{\mathrm{Var}(U)}{x^2}.$$

In a first time, assume the X_i 's are independent, $\mathbb{E}(X_i) = \mu$ and $\mathrm{Var}(X_i) = \sigma^2$,

$$\mathbb{P}\left\{\left|\frac{1}{n}\sum_{i=1}^{n}X_{i}-\mu\right| \geq x\right\} = \frac{\operatorname{Var}\left(\sum_{i=1}^{n}X_{i}\right)}{n^{2}x^{2}}$$
$$= \frac{\sigma^{2}}{n^{2}x^{2}}.$$

But...



(Photo : Wikipedia).

What can we expect? (2/2)

$$\mathbb{P}\left\{\left|\frac{1}{n}\sum_{i=1}^{n}X_{i}-\mu\right|\geq x\right\}\leq \frac{\sigma^{2}}{nx^{2}}.$$

However, CLT:

$$\sqrt{\frac{n}{\sigma^2}}\left(\frac{1}{n}\sum_{i=1}^n X_i - \mu\right) \rightsquigarrow \mathcal{N}(0,1).$$

So, we expect:

$$\mathbb{P}\left\{\left|\frac{1}{n}\sum_{i=1}^{n}X_{i}-\mu\right|\geq x\right\}\simeq2\Phi\left(\frac{x\sqrt{n}}{\sigma}\right)\sim\frac{2\mathrm{e}^{-\frac{x^{2}n}{2\sigma^{2}}}}{\frac{x\sqrt{n}}{\sigma}\sqrt{2\pi}}.$$

Chernoff bound

Chernoff bound

$$\mathbb{P}\Big\{U - \mathbb{E}(U) \ge x\Big\} = \mathbb{P}\Big\{e^{s(U - \mathbb{E}(U))} \ge e^{sx}\Big\} \le \frac{\mathbb{E}\left(e^{s(U - \mathbb{E}(U))}\right)}{e^{sx}}.$$

$$\mathbb{P}\left\{\frac{1}{n}\sum_{i=1}^{n}X_{i}-\mu\geq x\right\}\leq \frac{\mathbb{E}\left(e^{\frac{s}{n}\sum_{i=1}^{n}(X_{i}-\mu)}\right)}{e^{sx}}$$
$$=e^{-sx}\prod_{i=1}^{n}\mathbb{E}\left(e^{\frac{s}{n}(X_{i}-\mu)}\right).$$

Hoeffding's inequality

$$\mathbb{P}\left\{\frac{1}{n}\sum_{i=1}^{n}X_{i}-\mu\geq x\right\}\leq \mathrm{e}^{-sx}\prod_{i=1}^{n}\mathbb{E}\left(\mathrm{e}^{\frac{s}{n}(X_{i}-\mu)}\right).$$

Hoeffding's lemma - U bounded : $\underline{a \leq U \leq b}$

$$\mathbb{E}\left(e^{s[U-\mathbb{E}(U)]}\right) \le e^{\frac{s^2(b-a)^2}{8}}.$$

Hoeffding's inequality

Assume the X_i 's are independent and $a \leq X_i \leq b$,

$$\mathbb{P}\left\{\left|\frac{1}{n}\sum_{i=1}^{n}X_{i}-\mu\right|\geq x\right\}\leq 2\mathrm{e}^{-\frac{2nx^{2}}{(b-a)^{2}}}.$$

McDiarmid's inequality

McDiarmid's inequality

Assume the X_i 's are independent and $f: \mathcal{X}^n \to \mathbb{R}$ such that

$$|f(x_1,\ldots,x_{i-1},x_i,x_{i+1},\ldots,x_n)-f(x_1,\ldots,x_{i-1},x_i',x_{i+1},\ldots,x_n)| \leq c.$$

then

$$\mathbb{P}\left\{\left|\frac{f(X_1,\ldots,X_n)-\mathbb{E}[f(X_1,\ldots,X_n)]}{n}\right|\geq x\right\}\leq 2\mathrm{e}^{-\frac{2x^2n}{c^2}}.$$

We recover Hoeffding for $f(x_1, ..., x_n) = \sum_{i=1}^n x_i$, c = (b-a).

Assumptions on moments

Hoeffding's lemma - U bounded : a < U < b

$$\mathbb{E}\left(e^{s[U-\mathbb{E}(U)]}\right) \le e^{\frac{s^2(b-a)^2}{8}}.$$

In general, why not assuming U satisfies such an inequality?

Definition - sub-Gaussian random variable *U*

$$\mathbb{E}\left(\mathrm{e}^{s[U-\mathbb{E}(U)]}\right) \le \mathrm{e}^{s^2C_0^2}$$

$$U$$
 sub-Gaussian $\Leftrightarrow \forall k \in \mathbb{N}, \mathbb{E}(|U|^{2k}) \leq k! C_1^k$.

Contents

- Deviation inequalities for time series : introduction
 - Why deviation inequalities?
 - Deviation inequalities for time series
- 2 Non-homogeneous Markov chains
 - Inequalities for non-homogeneous Markov chains
 - Applications in machine learning

Objective of this talk

Objective : for some time series $\{X_t, t=0,\ldots,\infty\}$

$$\mathbb{P}\left\{\left|\frac{1}{n}\sum_{t=1}^{n}X_{t}-\frac{1}{n}\mathbb{E}\left(\sum_{t=1}^{n}X_{t}\right)\right|\geq x\right\}\leq ?$$

$$\mathbb{P}\left\{\left|\frac{f(X_{1},\ldots,X_{n})-\mathbb{E}[f(X_{1},\ldots,X_{n})]}{n}\right|\geq x\right\}\leq ?$$

$$\mathbb{P}\left\{\frac{1}{n}\sum_{t=1}^{n}X_{t}-\mu\geq x\right\} \leq \frac{\mathbb{E}\left(e^{\frac{s}{n}\sum_{t=1}^{n}(X_{t}-\mu)}\right)}{e^{sx}}$$
$$=e^{-sx}\prod_{t=1}^{n}\mathbb{E}\left(e^{\frac{s}{n}(X_{t}-\mu)}\right)$$

Objective of this talk

Objective : for some time series $\{X_t, t = 0, \dots, \infty\}$

$$\mathbb{P}\left\{\left|\frac{1}{n}\sum_{t=1}^{n}X_{t}-\frac{1}{n}\mathbb{E}\left(\sum_{t=1}^{n}X_{t}\right)\right|\geq x\right\}\leq ?$$

$$\mathbb{P}\left\{\left|\frac{f(X_{1},\ldots,X_{n})-\mathbb{E}[f(X_{1},\ldots,X_{n})]}{n}\right|\geq x\right\}\leq ?$$



$$\mathbb{P}\left\{\frac{1}{n}\sum_{t=1}^{n}X_{t}\right\} \leq \frac{\mathbb{E}\left(e^{\frac{s}{n}\sum_{t=1}^{n}(X_{t}-\mu)}\right)}{e^{sx}}$$

$$= e^{-sx}\prod_{t=1}^{n}\mathbb{E}\left(e^{\frac{s}{n}(X_{t}-\mu)}\right)$$

Deviation for time series: an active research field













A remarkable result for Markov chains



Available online at www.sciencedirect.com ScienceDirect Sechaetic Processes and their Ambications 125 (2015) 60–60

stochastic processes and their applications

Deviation inequalities for separately Lipschitz functionals of iterated random functions

Jérôme Dedeckera.*. Xieguan Fanb

³ Université Paris Descartes, Sorbonne Paris Cité, Luboratoire MAPS and CNRS UMR 8145, 75016 Paris, France ³ Regularity Tears, Isriu and MAS Luboratory, Ecole Centrale Paris - Grande Voie des Vignes, 92295 Chicago, Mohler, Fasser

Received 11 February 2014; received in revised form 18 July 2014; accepted 2 August 2014

Abstract

We consider in X-valued Markov danix X_1, X_2, \dots, X_n belonging to a class of iterated random functions, which is "one easy contracting" with regoed to some distance d on $X_n^{-1}B^2$ is any separately distance d on $X_n^{-1}B^2$ is any separately distance d on d of d

MSC 60G42-60005-60ELS

Keywords: Berated random functions; Martingales; Exponential inequalities; Moment inequalities; Wasserstein distances

1 A class of iterated random functions

Let (Ω, A, \mathbb{P}) be a probability space. Let (\mathcal{X}, d) and (\mathcal{Y}, δ) be two complete separable metric spaces. Let $(\varepsilon_i)_{i\geq 1}$ be a sequence of independent and identically distributed (iid) \mathcal{Y} -valued

http://dx.doi.org/10.1016/j.spa.2014.08.001 0304-4149/S; 2014 Elsevier B.V. All rights reserved. study Markov chains of the form

$$X_n = F(X_{n-1}, \varepsilon_n)$$

 provide deviation inequalities when

$$\mathbb{E}\bigg\{d\Big(F(x,\varepsilon_n),F(x',\varepsilon_n)\Big)\bigg\}\leq \rho d(x,x')$$

for some
$$\rho < 1$$
.

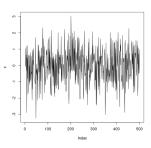
^{*} Corresponding author, Tel.: +33 1 83 94 88 72.

E-molf addresses: interne didector@maisdescures & (I. Dedector), favoiceam@hormail.com (X. Fan).

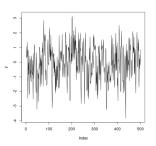
Example (1/2)

AR(1) process

$$X_n = F(X_{n-1}, \varepsilon_n) := \rho X_{n-1} + \varepsilon_n$$
$$|F(x, \varepsilon_n) - F(x', \varepsilon_n)| \le \rho |x - x'|$$



$$\rho = 0$$

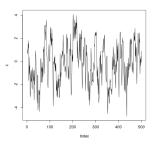


$$\rho = 0.5$$

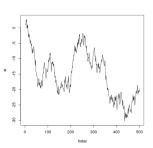
Example (2/2)

AR(1) process

$$X_n = F(X_{n-1}, \varepsilon_n) := \rho X_{n-1} + \varepsilon_n$$
$$|F(x, \varepsilon_n) - F(x', \varepsilon_n)| \le \rho |x - x'|$$



$$\rho = 0.8$$

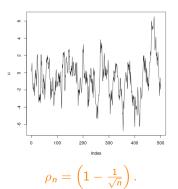


$$\rho = 1$$

What happens for non-homogeneous chains?

AR(1) process with varying coefficients

$$X_n = F_n(X_{n-1}, \varepsilon_n) := \rho_n X_{n-1} + \varepsilon_n$$



Inequalities for non-homogeneous Markov chains

- Deviation inequalities for time series : introduction
 - Why deviation inequalities?
 - Deviation inequalities for time series
- Non-homogeneous Markov chains
 - Inequalities for non-homogeneous Markov chains
 - Applications in machine learning

A class of non-homogeneous Markov chains

- X_n takes values in (\mathcal{X}, d) . Example : $\mathcal{X} = \mathbb{R}^d$, d large.
- (ε_n) are i.i.d random variables in (\mathcal{Y}, δ) .

Definition

- $\mathbb{E}\bigg\{d\Big(F_n(x,\varepsilon_n),F_n(x',\varepsilon_n)\Big)\bigg\}\leq \rho_n d(x,x').$

VAR with variying coefficients



Phillips, P.C.B. (1988). Regression theory for near integrated time series. Econometrica.

- $X_n \in \mathbb{R}^d$.
- (ε_n) are i.i.d $\mathcal{N}(0, \sigma^2 I_d)$.

- **3** $\tau_n = 1$.

Example: stochastic optimization

Minimize
$$L(x) = \sum_{i=1}^{N} \ell_i(x)$$

For I drawn uniformly in $\{1, \ldots, N\}$ with M elements,

$$\hat{\nabla}_n L(x) := \frac{1}{M} \sum_{i \in I} \nabla \ell_i(x).$$

Projected tochastic gradient descent (SGD) :

$$X_n = \Pi_{\mathcal{C}} \left[X_{n-1} - \frac{\gamma}{n^{\alpha}} \hat{\nabla}_n L(x) \right]$$

Projected stochastic gradient Langevin descent (SGLD) :

$$X_n = \Pi_{\mathcal{C}} \left[X_{n-1} - \frac{\gamma}{n^{\alpha}} \hat{\nabla}_n L(x) + \frac{\eta}{n^{\beta}} \varepsilon_n \right]$$

Example: SGD

Assume *L* is *m*-strongly convex and ∇L is ℓ -Lipschitz.

SGD -
$$\alpha \in [0,1]$$
, $\gamma > 0$

•
$$\rho_n \sim 1 - \frac{m\gamma}{n^{\alpha}}$$
 for $\alpha > 0$,
• $\rho_n = 1 - 2m\gamma + \ell^2 \gamma^2$ if $\alpha = 0$.

Example : SGLD

Assume L is m-strongly convex and ∇L is ℓ -Lipschitz.

SGLD -
$$\alpha, \beta \in [0,1]$$
, $\gamma, \eta > 0$, $\varepsilon_{\it n} \sim \mathcal{N}(0,1)$

•
$$\rho_n \sim 1 - \frac{m\gamma}{n^{\alpha}}$$
 for $\alpha > 0$,
• $\rho_n = 1 - 2m\gamma + \ell^2 \gamma^2$ if $\alpha = 0$.

Deviation inequality

Theorem (Proposition 3.1 in the paper) - $p \in [1, +\infty], d \in \mathbb{N}$

Assume $f: \mathcal{X}^n \to \mathbb{R}^d$ such that

$$|f(x_1,\ldots,x_i,\ldots,x_n)-f(x_1,\ldots,x_i',\ldots,x_n)|\leq d(x_i,x_i'),$$

 $\mathbb{E}_{\varepsilon_n}([\mathbb{E}_{\varepsilon_n'}\delta(\varepsilon_n,\varepsilon_n')]^k) \leq C_1^k k!$ and a similar condition for X_1 ,

$$\begin{split} & \mathbb{P}\left\{\left\|\frac{f(X_1,\ldots,X_n)-\mathbb{E}[f(X_1,\ldots,X_n)]}{n}\right\|_{\rho} \geq x\right\} \\ & \leq \left\{ \begin{array}{ll} \mathrm{e}^{-c_{p,d}nx} & \rho_n \leq \rho < 1,\, \tau_n \leq \frac{\tau}{n^{\alpha}},\, \alpha \in (0,1] \\ \mathrm{e}^{-c_{p,d}n(x1_{x>1}+x^21_{x\leq 1})} & \rho_n \leq 1-\frac{\rho}{n^{\alpha}},\, \tau_n \leq \frac{\tau}{n^{\alpha}},\, \alpha \in [0,1) \text{ or } \\ & \rho_n \leq 1-\frac{\rho}{n},\, \xi>0, \\ \mathrm{e}^{-c_{p,d}n^{1-2\alpha}x^2} & \rho_n \leq 1-\frac{\rho}{n^{\alpha}},\, (\tau_n \leq \tau \text{ or } \xi>0), \alpha \in (0,1/2). \end{array} \right. \end{split}$$

Proof technique

The proof technique relies on martingale decomposition :

$$f(X_1,\ldots,X_n)-\mathbb{E}[f(X_1,\ldots,X_n)]=\sum_{t=1}^n M_t$$

where

$$M_t = \mathbb{E}[f(X_1, \ldots, X_n)|X_1, \ldots, X_t] - \mathbb{E}[f(X_1, \ldots, X_n)|X_1, \ldots, X_{t-1}].$$

Conditional Chernoff:

$$\frac{\mathbb{E}\left(e^{\frac{s}{n}\sum_{t=1}^{n}M_{t}}\right)}{e^{sx}} = \frac{\mathbb{E}\left[e^{\frac{s}{n}\sum_{t=1}^{n-1}M_{t}}\mathbb{E}\left(e^{\frac{s}{n}M_{n}}|X_{1},\ldots,X_{n-1}\right)\right]}{e^{sx}}$$

Here the study of $\mathbb{E}\left(\mathrm{e}^{\frac{s}{n}M_n}|X_1,\ldots,X_{n-1}\right)$ requires some care...

Shameless name-dropping

In the paper, we provide an exhaustive list of inequalities, under various moment assumptions :

- exponential inequalities :
 - McDiarmid,
 - Hoeffding,
 - Bernstein.
- semi-exponential inequalities :
 - Fuk-Nagaev,
 - von Bahr-Esseen.
- moment inequalities :
 - Marcinkiewicz-Zygmund,
 - von Bahr-Esseen.

Applications

- Deviation inequalities for time series: introduction
 - Why deviation inequalities?
 - Deviation inequalities for time series
- Non-homogeneous Markov chains
 - Inequalities for non-homogeneous Markov chains
 - Applications in machine learning

Empirical risk minimization (1/2)

In the stationary case,

$$f(X_1,\ldots,X_n)=\frac{1}{n}\sum_{t=1}^n\ell(\theta,X_i)=R_n(\theta)$$

then

$$\mathbb{E}\left[f(X_1,\ldots,X_n)\right] = \mathbb{E}\left[\ell(\theta,X)\right] = R(\theta).$$

$$\mathbb{P}\left\{\left|R(\theta) - R_n(\theta)\right| \ge x\right\} \\
\le \begin{cases}
e^{-cnx} & \tau_n \le \frac{\tau}{n^{\alpha}}, \ \alpha \in (0, 1] \\
e^{-cn(x1_{x>1} + x^21_{x\le 1})} & \rho_n \le 1 - \frac{\rho}{n^{\alpha}}, \ \tau_n \le \frac{\tau}{n^{\alpha}}, \ \alpha \in [0, 1) \\
e^{-cn^{1-2\alpha}x^2} & \rho_n \le 1 - \frac{\rho}{n^{\alpha}}, \ \tau_n \le \tau, \ \alpha \in (0, 1/2).
\end{cases}$$

Empirical risk minimization (2/2)

ERM

$$\hat{\theta} = \arg\min_{\theta \in \Theta} R_n(\theta).$$

Say $Card(\Theta) = N$ is finite,

$$\mathbb{P}\bigg\{R(\hat{\theta}) \geq R_n(\hat{\theta}) + x\bigg\} \leq \begin{cases} Ne^{-cnx}, \\ Ne^{-cn(x1_{x>1} + x^21_{x\leq 1})}, \\ Ne^{-cn^{1-2\alpha}x^2}. \end{cases}$$

Application to SGLD (1/2)

L is *m*-strongly convex and ∇L is ℓ -Lipschitz.

SGLD -
$$\alpha, \beta \in [0, 1], \gamma, \eta \ge 0$$

$$X_n = \Pi_{\mathcal{C}} \left[X_{n-1} - \frac{\gamma}{n^{\alpha}} \hat{\nabla}_n L(x) + \frac{\eta}{n^{\beta}} \varepsilon_n \right], \quad \bar{X}_n = \frac{1}{n} \sum_{t=1}^n X_t.$$

For some
$$c_{p,d} = c_{p,d}(\ell, m)$$
,
$$\mathbb{P}\left\{\left\|\bar{X}_n - \mathbb{E}(\bar{X}_n)\right\|_p \ge x\right\}$$
$$\leq \begin{cases} e^{-c_{p,d}n(x\mathbf{1}_{x>1} + x^2\mathbf{1}_{x\leq 1})} & [0 \le \alpha \le \beta] \text{ or } [0 \le \alpha, \eta = 0] \\ e^{-c_{p,d}n^{1-2\alpha}x^2} & 0 < \alpha < \frac{1}{2}, \beta = 0. \end{cases}$$

Application to SGLD (2/2)

Assume in addition that L is Lipschitz.

Theorem - $\alpha \in [1/2, 1], \eta = 0$

$$\mathbb{E}\Big(\|\bar{X}_n-x^*\|_2^2\Big)\leq \frac{c}{n},\quad c=c(d,L,m).$$



Bach, F. and Moulines, E. (2011). Non-Asymptotic Analysis of Stochastic Approximation Algorithms for Machine Learning. *NIPS*.

Combine Bach & Moulines (2011) with our inequality

$$\mathbb{P}\left\{\left\|\bar{X}_n - x^*\right\|_2 \le \sqrt{\frac{c + \frac{1}{c_{2,d}}\log\left(\frac{1}{\delta}\right)}{n}}\right\} \ge 1 - \delta.$$