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



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

Pierre Alquier, RIKEN AIP Deviation inequalities for Markov chains

### 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

Pierre Alquier, RIKEN AIP Deviation inequalities for Markov chains

Why deviation inequalities ? Deviation inequalities for time series

### Objective

#### General problem in probability and statistics

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

Why deviation inequalities ? Deviation inequalities for time series

### What can we expect ? (1/2)

Chebyshev's inequality

$$\mathbb{P}\Big\{|U - \mathbb{E}(U)| \ge x\Big\} \le rac{\operatorname{Var}(U)}{x^2}.$$

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

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

Why deviation inequalities ? Deviation inequalities for time series

#### But...



(Photo : Wikipedia).

Why deviation inequalities ? Deviation inequalities for time series

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 :

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

Why deviation inequalities ? Deviation inequalities for time series

### Chernoff bound

#### Chernoff bound

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

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

Why deviation inequalities ? Deviation inequalities for time series

### Hoeffding's inequality

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

Hoeffding's lemma - U bounded :  $a \le U \le b$ 

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

#### Hoeffding's inequality

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

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

Why deviation inequalities ? Deviation inequalities for time series

### 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  $f(x_1, \ldots, x_n) = \sum_{i=1}^n x_i$ , c = (b-a).

Why deviation inequalities ? Deviation inequalities for time series

#### Assumptions on moments

Hoeffding's lemma - U bounded :  $a \le U \le b$ 

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

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

Definition - sub-Gaussian random variable U

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

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

Why deviation inequalities ? Deviation inequalities for time series

#### 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

Why deviation inequalities ? Deviation inequalities for time series

### Objective of this talk

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

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

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

Why deviation inequalities ? Deviation inequalities for time series

### Objective of this talk

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

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



Why deviation inequalities ? Deviation inequalities for time series

### Deviation for time series : an active research field



Pierre Alquier, RIKEN AIP Deviation inequalities for Markov chains

Why deviation inequalities ? Deviation inequalities for time series

#### A remarkable result for Markov chains

processes and their applications





#### Deviation inequalities for separately Lipschitz functionals of iterated random functions

#### Jérôme Dedecker<sup>a,\*</sup>, Xiequan Fan<sup>b</sup>

<sup>8</sup>Université Paris Descarier, Scohome Paris Cite, Laboratoire MAPS and CRRS UMR 8145, 75016 Paris, France <sup>8</sup>Regularity Team, Juria and MIS Laboratory, Ecolic Centrale Paris - Grande Viste des Vignes, 92295 Chietono, Molder, Foure

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

#### Abstract

We consider an X-valued Matter data X<sub>1</sub>, X<sub>2</sub>,...,X<sub>k</sub> belonging to a data of iterator dandom function, which "itera exportanting" with represent some distance data  $K \neq k$  is any segment Lipschie function with negoci to  $A_k$  are as a well known decomposition of  $S_k = f(X_1, \ldots, X_k) = B_k$  $(X_1, \ldots, X_k) = B_k$  ( $X_1, \ldots, X_k$ ) are and maturized  $X_k$  is a straight of the set of

MSC: 60G42; 60005; 60EL5

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

#### 1. A class of iterated random functions

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

\* Corresponding author, Tel: +33 1 83 94 58 72. E-moil addresses: jerome dedecker@parisdescartes.fr (J. Dedecker), fansieguan@hormail.com (X. Fan).

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

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

for some  $\rho < 1$ .

Why deviation inequalities ? Deviation inequalities for time series

## Example (1/2)

#### AR(1) process

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



 $\rho = 0$ 

$$\rho = 0.5$$

Pierre Alquier, RIKEN AIP Deviation inequalities for Markov chains

Why deviation inequalities ? Deviation inequalities for time series

## Example (2/2)

#### AR(1) process

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



Pierre Alquier, RIKEN AIP Deviation inequalities for Markov chains

Why deviation inequalities ? Deviation inequalities for time series

### 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

#### 2 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.
- $(\varepsilon_n)$  are i.i.d random variables in  $(\mathcal{Y}, \delta)$ .

#### Definition

• 
$$X_n = F_n(X_{n-1}, \varepsilon_n).$$
  
•  $\mathbb{E}\left\{d\left(F_n(x, \varepsilon_n), F_n(x', \varepsilon_n)\right)\right\} \le \rho_n d(x, x').$   
•  $d\left(F_n(x, y), F_n(x, y')\right) \le \tau_n \delta(y, y') + \xi_n.$ 

Inequalities for non-homogeneous Markov chains Applications in machine learning

### VAR with variying coefficients



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

• 
$$X_n \in \mathbb{R}^d$$
.

• 
$$(\varepsilon_n)$$
 are i.i.d  $\mathcal{N}(0, \sigma^2 I_d)$ .

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

$$\rho_n = ||A_n||_{\text{op}} = \sup_{x \neq 0} \frac{||A_n x||}{||x||} \xrightarrow[n \to \infty]{} 1$$

3 
$$\tau_n = 1, \ \xi_n = 0.$$

### 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) :

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

• Projected stochastic gradient Langevin descent (SGLD) :

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

### Example : SGD

# Assume L is *m*-strongly convex, *M*-Lipschitz and $\nabla L$ is $\ell$ -Lipschitz.

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

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

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

### Example : SGLD

# Assume L is *m*-strongly convex, *M*-Lipschitz and $\nabla L$ is $\ell$ -Lipschitz.

### SGLD - $\alpha, \beta \in [0, 1], \gamma, \eta > 0, \varepsilon_n \sim \mathcal{N}(0, 1)$ **1** $X_n = F_n(X_{n-1}, \varepsilon_n) = \prod_{\mathcal{C}} \left[ X_{n-1} - \frac{\gamma}{n^{\alpha}} \hat{\nabla}_n L(X_{n-1}) + \frac{\eta}{n^{\beta}} \varepsilon_n \right].$ **2** $\bullet \rho_n \sim 1 - \frac{m\gamma}{n^{\alpha}} \text{ for } \alpha > 0,$ $\bullet \rho_n = 1 - 2m\gamma + \ell^2 \gamma^2 \text{ if } \alpha = 0.$ **3** $\xi_n = \frac{2\gamma M}{n^{\alpha}}, \tau_n = \frac{\eta}{n^{\beta}}.$

### 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

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

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

$$\mathbb{P}\left\{ \left\| \frac{f(X_1,\ldots,X_n) - \mathbb{E}[f(X_1,\ldots,X_n)]}{n} \right\|_p \ge x \right\}$$
  
$$\leq \left\{ \begin{array}{l} e^{-c_{p,d}nx} & \rho_n \le 1 - \rho < 1, \ \tau_n + \xi_n \le \frac{\tau}{n^{\alpha}}, \ \alpha \in (0,1] \\ e^{-c_{p,d}n(x_{1_{x>1}+x^2 \mathbf{1}_{x\le 1}})} & \rho_n \le 1 - \frac{\rho}{n^{\alpha}}, \ \tau_n + \xi_n \le \frac{\tau}{n^{\alpha}}, \ \alpha \in [0,1), \\ e^{-c_{p,d}n^{1-2\alpha}x^2} & \rho_n \le 1 - \frac{\rho}{n^{\alpha}}, \ \tau_n + \xi_n \le \tau, \alpha \in (0,1/2). \end{array} \right\}$$

## 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

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

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

Here the study of  $\mathbb{E}\left(\mathrm{e}^{rac{s}{n}M_n}|X_1,\ldots,X_{n-1}
ight)$  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

#### 2 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

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

$$\mathbb{P}\left\{ \left| R(\theta) - R_n(\theta) \right| \ge x \right\} \le \begin{cases} e^{-cnx}, \\ e^{-cn(x\mathbf{1}_{x>1} + x^2\mathbf{1}_{x\leq 1})}, \\ e^{-cn^{1-2\alpha}x^2}. \end{cases}$$

Inequalities for non-homogeneous Markov chains Applications in machine learning

### Empirical risk minimization (2/2)

#### ERM

$$\hat{\theta} = \underset{\theta \in \Theta}{\operatorname{arg\,min}} R_n(\theta).$$

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

$$\mathbb{P}\left\{R(\hat{\theta}) \geq R_n(\hat{\theta}) + x\right\} \leq \begin{cases} N e^{-cnx}, \\ N e^{-cn(x_{1_{x>1}} + x^2 \mathbf{1}_{x\leq 1})}, \\ N e^{-cn^{1-2\alpha}x^2}. \end{cases}$$

Inequalities for non-homogeneous Markov chains Applications in machine learning

### Application to SGLD (1/2)

*L* is *m*-strongly convex, *M*-Lipschitz and  $\nabla L$  is  $\ell$ -Lipschitz.

SGLD - 
$$\alpha \in (0, 1), \beta < \alpha, \gamma > 0, \eta \ge 0$$
  
$$X_n = \prod_{\mathcal{C}} \left[ X_{n-1} - \frac{\gamma}{n^{\alpha}} \hat{\nabla}_n L(X_{n-1}) + \frac{\eta}{n^{\beta}} \varepsilon_n \right], \quad \bar{X}_n = \frac{1}{n} \sum_{t=1}^n X_t.$$

For some 
$$c_{p,d} = c_{p,d}(\ell, m, M)$$
,  
 $\mathbb{P}\left\{ \left\| \bar{X}_n - \mathbb{E}(\bar{X}_n) \right\|_p \ge x \right\} \le e^{-c_{p,d}n(x\mathbf{1}_{x>1}+x^2\mathbf{1}_{x\le 1})}.$ 

Inequalities for non-homogeneous Markov chains Applications in machine learning

### Application to SGLD (2/2)

#### Theorem - Moulines and Bach 2011

$$\mathbb{E}\|\bar{X}_n - x^*\|_2^2 \le \frac{C_0}{n}$$

Moulines, E. and Bach, F. (2011). Non-asymptotic analysis of stochastic approximation algorithms for machine learning. *NIPS*.

#### Combine with our inequality

$$\mathbb{P}\left\{\left\|\bar{X}_n-x^*\right\|_2\leq \sqrt{\frac{C_0+\frac{1}{c_{2,d}}\log\left(\frac{1}{\delta}\right)}{n}}\right\}\geq 1-\delta.$$