

# A limit conditional theorem for random walks, Applications to IS

Conditioned random walks, Trondheim, 2012

Michel Broniatowski, Virgile Caron

LSTA, Université Paris 6, France

June 27th, 2012

$X_1^n := (X_1, \dots, X_n)$ ,  $X_i$ 's i.i.d. on  $\mathbb{R}^d$  with density  $p$ .

$f : \mathbb{R}^d \rightarrow \mathbb{R}^s$

local condition

$$\frac{1}{n} \sum f(X_i) = a_n \quad (\text{Local})$$

$(a_n)_{n \geq 1}$  is a convergent sequence.

Either  $a_n \rightarrow EX$  (LLN range) or  $a_n \rightarrow a \neq EX$  (LD).

Question What about the distribution of  $(X_1, \dots, X_{k_n})$  given Local when

$$k_n/n \rightarrow 1 \quad \text{long runs}$$

$$n - k_n \rightarrow \infty \quad \text{not all}$$

Notation  $p(X_1^k = x_1^k \mid \frac{1}{n} \sum f(X_i) = a_n) =: p_{a_n}(X_1^k = x_1^k)$ .

Take  $f(x) = x$ ,  $d = s = 1$ ,  $k = 1$ .

$$\frac{1}{n} \sum X_i = a$$

Gibbs conditioning principle (local form) with  $a_n$  fixed  $= a$ .

$$\phi(t) \quad : \quad = E \exp tX$$

$$m(t) = \frac{d}{dt} \log \phi(t), \quad s^2(t) = \frac{d^2}{dt^2} \log \phi(t), \quad \mu_3(t) := \frac{d^3}{dt^3} \log \phi(t)$$

$t$  such that  $m(t) = a$

$$\pi^a(x) := \frac{\exp tx}{\phi(t)} p(x)$$

Recall  $E_{\pi^a}(X) = a$ ,  $\text{Var}_{\pi^a}(X) = s^2(t)$ . Then (extensions (Diaconis and Friedman (1988)  $k = o(n)$ , Dembo and Zeitouni (1995)  $\limsup_n k/n < 1$ )

$$\int |p_a - \pi^a| dx \rightarrow 0 \text{ as } n \rightarrow \infty$$
$$\sup_{A \in \mathcal{B}(\mathbb{R})} P_a(A) - \Pi^a(A) \rightarrow 0.$$

Relevance of this question for IS?

$$\begin{aligned} p_A \left( X_1^k = x_1^k \right) & : = p \left( X_1^k = x_1^k \mid \frac{1}{n} \sum f(X_i) \in A \right) \\ & \sim \int_A g_u \left( X_1^k = x_1^k \right) p \left( \frac{1}{n} \sum f(X_i) = u \mid \frac{1}{n} \sum f(X_i) \in A \right) \\ & = : g_A \left( X_1^k = x_1^k \right). \end{aligned}$$

a mixture of the approximating densities with weights charging  $A$ . (no dominating point in this approach).

If  $g_u$  is "good"

$$p_u \left( X_1^k = x_1^k \right) = g_u \left( X_1^k = x_1^k \right) (1 + o(1))$$

then we expect

$$p_A \left( X_1^k = x_1^k \right) = g_A \left( X_1^k = x_1^k \right) (1 + o(1))$$

which with  $k_n$  "close to  $n$ " is OK when sampling under  $g_A$  is possible.

A recursive approximation (turn to  $f(x) = x$ ,  $d = s = 1$  for convenience)

$$p_{a_n} \left( X_1^k = x_1^k \right) = \prod_{i=0}^{k-1} p \left( X_{i+1} = x_{i+1} \mid S_1^n = na_n, X_1^i = x_1^i \right)$$

Now for any  $\alpha$  (invariance of conditional densities under any tilting)

$$p \left( X_{i+1} = x_{i+1} \mid S_1^n = na_n, X_1^i = x_1^i \right) = \pi^\alpha \left( X_{i+1} = x_{i+1} \mid S_1^n = na_n, X_1^i = x_1^i \right)$$

$$\begin{aligned} & \pi^{m_i} (X_{i+1} = x_{i+1} | S_1^n = na_n, X_1^i = x_1^i) \\ = & \pi^{m_i} (X_{i+1} = x_{i+1}) \frac{\pi^{m_i} (S_{i+1}^n = na_n - s_1^i)}{\pi^{m_i} (S_i^n = na_n - s_1^{i-1})} \end{aligned}$$

with  $s_1^i := x_1 + \dots + x_i$ , with  $m_i$  making the ratio simple to evaluate. A precise evaluation of the dominating terms in this latest expression is needed in order to handle the product in the joint density. Center, reduce, Edgeworth expansions with

$$m_i := m(t_i) := \frac{1}{n - i + 1} (na_n - s_1^{i-1})$$

. **Pb: the orders of magnitude of the  $x_i$ 's in order to control the  $k_n$  products  $\implies$  under the conditional sampling** Develop some maximal inequalities. Use Edgeworth expansions, etc.

Instead of arbitrary  $x_i$ 's consider  $Y_i$ 's random under the conditional distribution **(typical paths)**.

Order of magnitude:

$$\max_{1 \leq i \leq n} |Y_i| = O_{P_{a_n}}(\log n)$$

(not large, but  $\rightarrow \infty$ )

**Theorem:**

$$\begin{aligned} p_{a_n} \left( X_1^{k_n} = Y_1^{k_n} \right) &= g_{a_n} \left( Y_1^{k_n} \right) \left( 1 + o_{P_{a_n}}(\delta_n) \right) \\ g_{a_n} \left( Y_1^{k_n} \right) &= p_{a_n} \left( X_1^{k_n} = Y_1^{k_n} \right) \left( 1 + o_{G_{a_n}}(\delta_n) \right) \end{aligned}$$

For IS: this is OK (sample under  $g_u$ )

**Remark:** Implies

$$\sup_{B \in \mathcal{B}(\mathbb{R}^{k_n})} P_{a_n}(B) - G_{a_n}(B) \rightarrow 0.$$

$$g_t(y_1^k) := \prod_{i=0}^{k-1} g_i(y_{i+1} | y_1^i).$$

$$g_{i+1}(y_{i+1} | y_1^i) = C_i p(y_{i+1} | \alpha\beta + a_n, \alpha, y_{i+1})$$

where  $n(\mu, \tau, x)$  is the normal density with mean  $\mu$  and variance  $\tau$  at  $x$ .  
Here

$$\alpha = s_{i,n}^2 (n - i - 1), \quad \beta = t_{i,n} + \frac{\mu_3^{(i,n)}}{2s_{i,n}^4 (n - i - 1)}$$

The terms in  $\alpha$  and  $\beta$  depend on the past values and on the m.g.f. of  $X$ .

**Remark:** coincides with the exact gaussian conditional density in the gaussian case, with  $k$  up to  $n$ . For fixed  $k$ : coincides with the usual tilted+rate. When  $k_n$  large, the quadratic term in the  $g_{i+1}(y_{i+1} | y_1^i)g_{i+1}(y_{i+1} | y_1^i)$  is dominant (reduces the variance). When conditioning on an average of the  $f(X_i)$ 's a change in the formula.



## Consequence

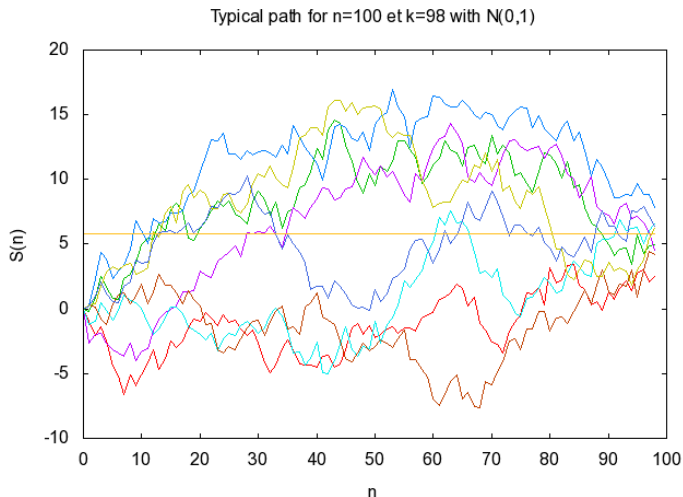
$$\begin{aligned} p_A \left( X_1^{k_n} = Y_1^{k_n} \right) &= g_A \left( Y_1^{k_n} \right) \left( 1 + o_{P_A}(\delta_n) \right) \\ g_A \left( Y_1^{k_n} \right) &= p_A \left( X_1^{k_n} = Y_1^{k_n} \right) \left( 1 + o_{G_A}(\delta_n) \right) \end{aligned}$$

for any "thick"  $A \in \mathcal{B}(\mathbb{R})$  and

$$\sup_{B \in \mathcal{B}(\mathbb{R}^{k_n})} P_A(B) - G_A(B) \rightarrow 0.$$

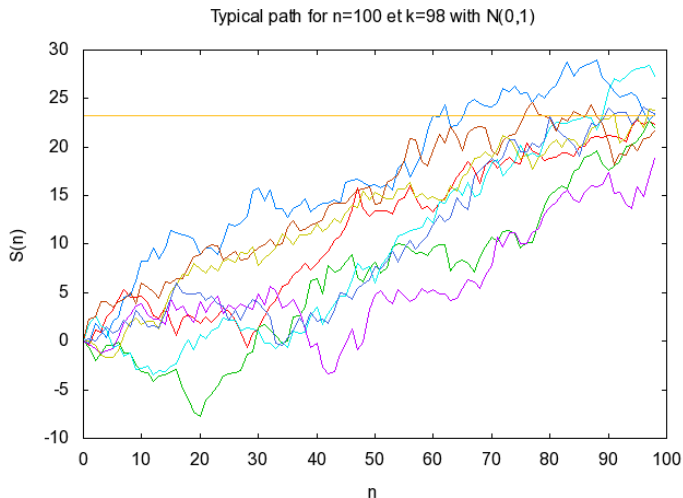
# Trajectories

Normal Case, CLT



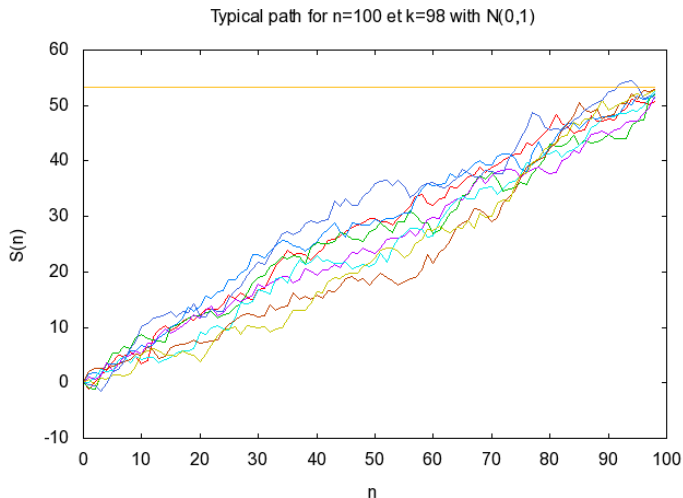
# Trajectories

Normal Case, MD



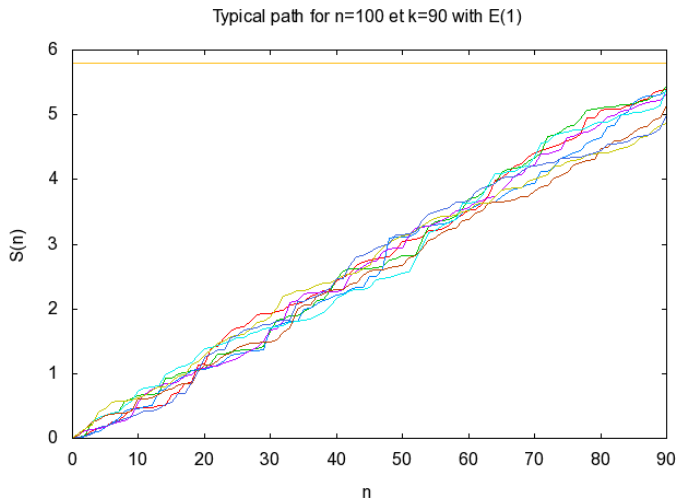
# Trajectories

Normal Case, LDP



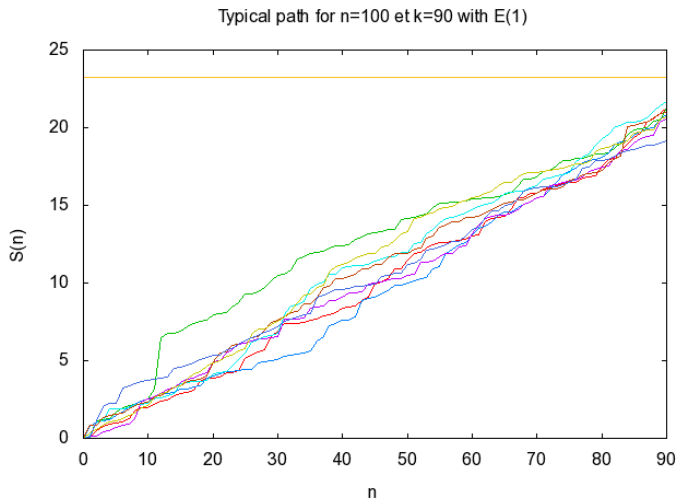
# Trajectories

Exponential Case, CLT



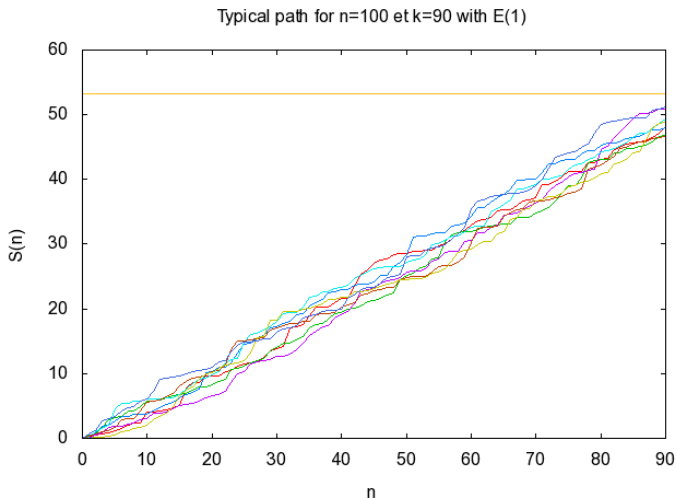
# Trajectories

Exponential Case, MD



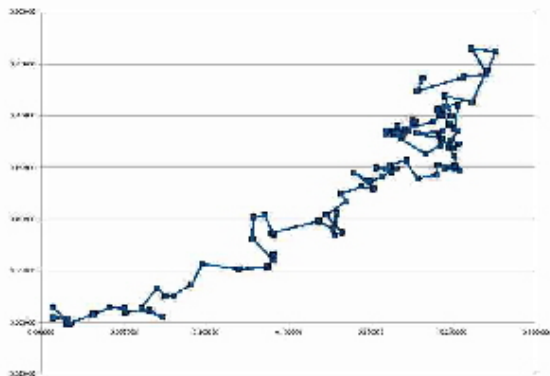
# Trajectories

Exponential Case, LD



# Trajectories

Normal Case, R2, MD





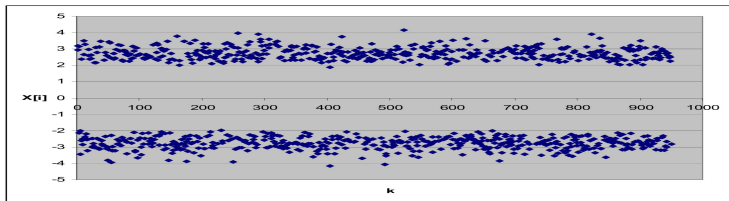
$f : \mathbb{R}^d \rightarrow \mathbb{R}$ .  $\mathbf{X}_1, \dots, \mathbf{X}_n$  i.i.d.  $f(\mathbf{X}_1)$  light tails.  $a_n$  large (high level,  $a_n \gg Ef(\mathbf{X}_1)$ ). **Draw**  $x : f(x)$  **of order**  $a_n$ .

For

$$P \left( \text{all the } \mathbf{X}'_i \text{'s } \approx a_n \mid \frac{1}{n} \sum \mathbf{X}_i > a_n \right) \rightarrow 1$$

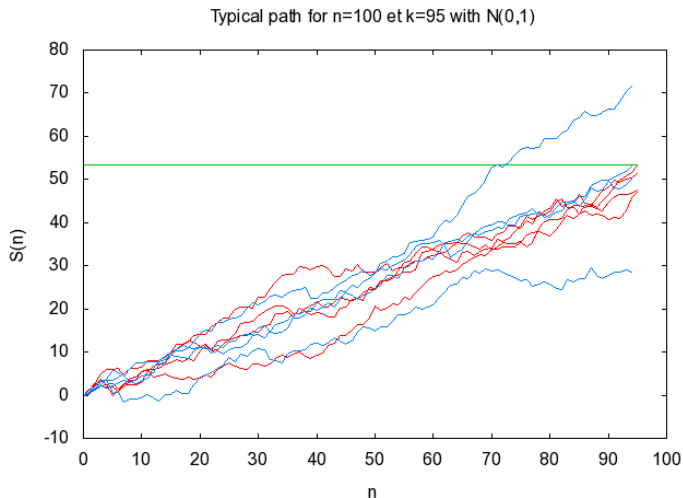
$f(x) = x^2$ ,  $\mathcal{L}(\mathbf{X}) = \text{Symmetric Weibull shape parameter } 2$ ,  $a_n = 10$ ,  
 $n = 1000$

No more Gibbs equivalent in this asymptotics (B-Cao, 2012 Arxiv)



# Trajectories

Blue:Tilted, Red:Adaptative Method



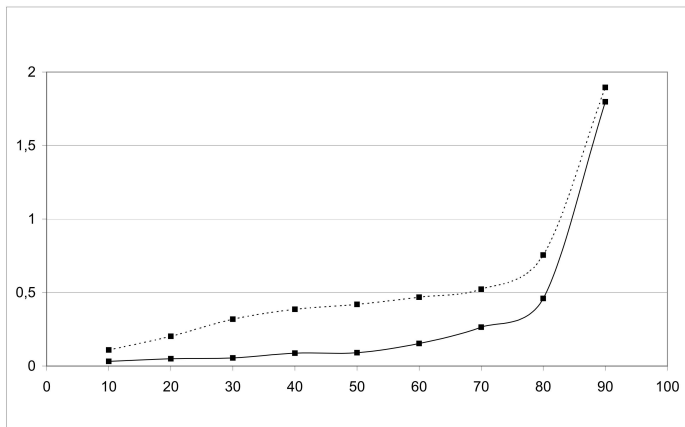
How long are good conditioned sampled runs?

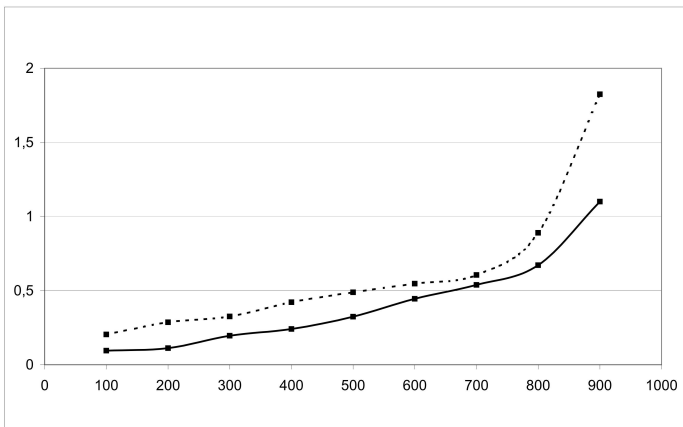
An empirical benchmark: the relative error as a function of  $k_n$

$$RE(k) := E_{G_u} \frac{|p_u(Y_1^k) - g_u(Y_1^k)|}{p_u(Y_1^k)}$$

to be estimated.

Remark: When  $A = (a, \infty)$ , the same indicators for  $k_n$





Evaluate  $k = k_n$

- Simulate  $u$  in  $A$  following

$$L\left(\frac{1}{n}\sum f(X_i) \middle| \frac{1}{n}\sum f(X_i) \in A\right)$$

(Approximate the distribution (ex: LDP  $\implies$  use Petrov, etc), or Metropolis-Hastings)

Evaluate  $k = k_n$

- Simulate  $u$  in  $A$  following

$$L \left( \frac{1}{n} \sum f(X_i) \middle| \frac{1}{n} \sum f(X_i) \in A \right)$$

(Approximate the distribution (ex: LDP  $\implies$  use Petrov, etc), or Metropolis-Hastings)

- Simulate  $Y_1^k$  with density  $g_u \sim p \left( \cdot \middle| \frac{1}{n} \sum f(X_i) = u \right)$ .

Evaluate  $k = k_n$

- Simulate  $u$  in  $A$  following

$$L \left( \frac{1}{n} \sum f(X_i) \middle| \frac{1}{n} \sum f(X_i) \in A \right)$$

(Approximate the distribution (ex: LDP  $\implies$  use Petrov, etc), or Metropolis-Hastings)

- Simulate  $Y_1^k$  with density  $g_u \sim p \left( \cdot \middle| \frac{1}{n} \sum f(X_i) = u \right)$ .
- Simulate  $Y_{k+1}^n$  with tilted density at point  $m_k$



Evaluate  $k = k_n$

- Simulate  $u$  in  $A$  following

$$L \left( \frac{1}{n} \sum f(X_i) \middle| \frac{1}{n} \sum f(X_i) \in A \right)$$

(Approximate the distribution (ex: LDP  $\implies$  use Petrov, etc), or Metropolis-Hastings)

- Simulate  $Y_1^k$  with density  $g_u \sim p(\cdot | \frac{1}{n} \sum f(X_i) = u)$ .
- Simulate  $Y_{k+1}^n$  with tilted density at point  $m_k$
- Evaluate the IS ratio

$$\frac{\prod_{i=1}^n p(Y_i)}{g_u(Y_1^k) \prod_{i=k+1}^n \pi^{m_k}(Y_i)} \mathbf{1}_A \left( \frac{1}{n} \sum f(Y_i) \right)$$

Evaluate  $k = k_n$

- Simulate  $u$  in  $A$  following

$$L \left( \frac{1}{n} \sum f(X_i) \middle| \frac{1}{n} \sum f(X_i) \in A \right)$$

(Approximate the distribution (ex: LDP  $\implies$  use Petrov, etc), or Metropolis-Hastings)

- Simulate  $Y_1^k$  with density  $g_u \sim p(\cdot | \frac{1}{n} \sum f(X_i) = u)$ .
- Simulate  $Y_{k+1}^n$  with tilted density at point  $m_k$
- Evaluate the IS ratio

$$\frac{\prod_{i=1}^n p(Y_i)}{g_u(Y_1^k) \prod_{i=k+1}^n \pi^{m_k}(Y_i)} 1_A \left( \frac{1}{n} \sum f(Y_i) \right)$$

- Repeat from top  $L$  times

Evaluate  $k = k_n$

- Simulate  $u$  in  $A$  following

$$L \left( \frac{1}{n} \sum f(X_i) \middle| \frac{1}{n} \sum f(X_i) \in A \right)$$

(Approximate the distribution (ex: LDP  $\implies$  use Petrov, etc), or Metropolis-Hastings)

- Simulate  $Y_1^k$  with density  $g_u \sim p(\cdot | \frac{1}{n} \sum f(X_i) = u)$ .
- Simulate  $Y_{k+1}^n$  with tilted density at point  $m_k$
- Evaluate the IS ratio

$$\frac{\prod_{i=1}^n p(Y_i)}{g_u(Y_1^k) \prod_{i=k+1}^n \pi^{m_k}(Y_i)} 1_A \left( \frac{1}{n} \sum f(Y_i) \right)$$

- Repeat from top  $L$  times
- Average  $\implies \widetilde{P}_n$

**Properties** Standard IS (i.i.d. replications under the tilted at dominating point). In LDP, for

$$\frac{1}{n} \sum f(X_i) > a$$

$\widehat{Var}_{\pi^a} \widehat{P}_n$  proportional to  $\sqrt{n}$  (e.g. Sadowsky and Bucklew)

$\widehat{Var}_{g_A} \widehat{P}_n$  proportional to  $\sqrt{n-k}$  (optimal on the  $k_n$  first summands)

$\widetilde{P}_n$  has a small asymptotic variability when evaluated on classes of subsets of  $\mathbb{R}^n$  whose probability goes to 1 under the sampling  $g_A$ .

**Defaults:** many steps, time run (obviously not worth for standard cases)

Accuracy (no dominating point)?

Compare with other methods in quasi-standard cases

$X_1, \dots, X_{100}$  where  $X_1$  has a normal distribution  $N(0.05, 1)$  and let

$$\mathcal{E}_{100} := \left\{ x_1^{100} : \frac{|x_1 + \dots + x_{100}|}{100} > 0.28 \right\}$$

$$P_{100} = P((X_1, \dots, X_{100}) \in \mathcal{E}_{100}) = 0.01120.$$

simple **disymmetric case**. The standard i.i.d. IS scheme introduces the dominating point  $a = 0.28$  and the family of i.i.d. tilted r.v.'s with common  $N(a, 1)$  distribution. The resulting estimator of  $P_{100}$  is 0,01074 (with  $L = 1000$ ), indicating that the event  $S_{1,100}/100 < -0.28$  is ignored in the evaluation of  $P_{100}$ . Also the hit rate is of order 50%. It can also be seen that  $S_1^{100}/100 < -0.28$  is never visited through the procedure.

**Comparison with the cross entropy method.** The sampling distribution is chosen as a normal one with variance 1, as adapted to this situation; the mean is estimated recursively through Kullback minimisation. When the initialisation mean is close to 0.28 then the performance is similar to the classical IS scheme, since the successive means keep close to 0.28; at the contrary when it is defined close to -0.28 the sequence of sampling distributions tend to concentrate around  $N(-.28, 1)$  and the resulting estimate produces a relative error of order 100%. Indeed it is roughly  $|(10^{-4} - 10^{-2}) / 10^{-2}|$  since  $P_{\mathcal{E}_{100}^-} \sim 10^{-4}$  where

$$\mathcal{E}_{100}^- := \left\{ x_1^{100} : \frac{x_1 + \dots + x_{100}}{100} < -0.28 \right\}.$$

Rôle of point conditioning in stats:  
sufficiency

$$p_{\theta} (x_1^n | t(x_1^n)) \text{ independent upon } \theta$$

Rao-Blackwell:  $S(X_1^n)$  estimator of  $\theta$ .  $t(x_1^n)$  any statistics

$$MSE (E (S(X_1^n) | t(X_1^n))) \leq MSE (S(X_1^n))$$

Optimality when  $t$  is sufficient (Lehman-Scheffé) **PB: estimate**  
 $E (S(X_1^n) | T(X_1^n))$ .

Observe  $x_1^n$ , and  $t(x_1^n)$ ; in exponential families  $t(x_1^n) = \frac{1}{n} \sum t(x_i) = t_{obs}$  (converges under  $\theta_0$ ).

Choose  $k_n$

Simulate according to  $g_{t_{obs}}(\cdot) \sim p_\theta(\cdot | \frac{1}{n} \sum t(X_i) = \frac{1}{n} \sum t(x_i))$  (ind upon  $\theta$ )

Estimate  $E(S(X_1^n) | t(X_1^n) = t(x_1^n))$  averaging the values of  $S$  on the  $k_n$  realizations under  $g_{t_{obs}}$

**Remark**  $t$  sufficient for  $g_{t_{obs}}(\cdot)$ , so put any  $\theta$  in the definition of  $g_{t_{obs}}$