Bregman Deviations of Generic Exponential Families

Patrick Saux¹

(Joint work with Sayak Ray Chowdhury², Odalric-Ambrym Maillard¹ and Aditya Gopalan³)

¹ Univ. Lille, Inria, CNRS, Centrale Lille, UMR 9189 - CRIStAL, F-59000, Lille, France ² Boston University, Boston, Massachussetts, United States

³ Indian Institute of Science, Bangalore, India





Table of Contents



Time-uniform concentration: method of mixture

Bregman uniform concentration for generic exponential families







Let X_1, \ldots, X_t i.i.d random variables distributed as $X \sim \nu \in \mathcal{P}(\mathbb{R}^d)$. Where is $= \mathbb{E}[X]$?



t=1000



We want confidence sets, not just estimation!

We want **confidence sets**, not just estimation! For $\delta \in (0, 1)$ and $t \in \mathbb{N}$, we need a set C_t^{δ} such that $\mathbb{P}\left(\mu \in C_t^{\delta}\right) \ge 1 - \delta$. Asymptotically:

$$C_t^{\delta} = \left[\widehat{\mu}_t - \frac{\sigma^2}{\sqrt{t}} \Phi^{-1}(1 - \delta/2), \ \widehat{\mu}_t + \frac{\sigma^2}{\sqrt{t}} \Phi^{-1}(1 - \delta/2)\right]$$

is such that

$$\lim_{t \to +\infty} \mathbb{P}\left(\mu \in \mathcal{C}_t^{\delta}\right) = 1 - \delta \,,$$

where $\widehat{\mu}_t = \frac{1}{t} \sum_{s=1}^t X_s$.

X Does not tell us anything about the small sample size regime...

Nonasymptotic confidence sets

We need some assumptions...

Sub- ψ distributions:

$$orall \lambda \in \mathcal{I} \subseteq \mathbb{R}_+, \ \log \mathbb{E}_{\boldsymbol{X} \sim
u} \left[e^{\lambda(\boldsymbol{X} - \mu)}
ight] \leq \psi(\lambda) \,.$$

Proposition (Chernoff bound)

$$\mathbb{P}\left(\widehat{\mu}_t - \mu \geq \psi_*^{-1}\left(\frac{1}{t}\log\frac{1}{\delta}\right)\right) \leq \delta\,,$$

where $\psi_*(u) = \sup_{\lambda \in \mathcal{I}} \lambda u - \psi(\lambda)$ is the Fenchel-Legendre conjugate of ψ .

Examples:

•
$$\nu = \mathcal{N}(\mu, \sigma^2), \ \psi(\lambda) = \sigma^2 \lambda^2 / 2$$

- ν has bounded support, $\psi(\lambda) = diam(Supp(\nu))^2 \lambda^2/8$.
- $\nu = \chi^2(k), \ \psi(\lambda) = (1 2\lambda)^{-k/2}, \ \mathcal{I} = (0, \frac{1}{2}).$

Nonasymptotic confidence sets

Proof:

$$\mathbb{P}\left(\sum_{s=1}^{t} X_{s} - \mu \geq tu\right) = \mathbb{P}\left(\prod_{s=1}^{t} e^{\lambda(X_{s} - \mu)} \geq e^{t\lambda u}\right)$$
$$\leq e^{-t\lambda u} \mathbb{E}\left[\prod_{s=1}^{t} e^{\lambda(X_{s} - \mu)}\right] \quad (Markov)$$
$$= e^{-t\lambda u} \prod_{s=1}^{t} \mathbb{E}\left[e^{\lambda(X_{s} - \mu)}\right] \quad (independence)$$
$$\leq e^{-t\lambda u + t\psi(\lambda)} \quad (sub-\psi),$$

then optimise in $\lambda \in \mathcal{I}$.

Nonasymptotic confidence sets

Other possible assumptions:

- Fully parametric
 - If you know the quantiles, use them!
- Bounded
 - With control of moments: Bennett, Berstein Boucheron et al. [2013], Bentkus [2004],
 - With empirical estimators of moments: Bernstein Maurer and Pontil [2009], Bentkus Kuchibhotla and Zheng [2021].
 - ▶ With only boundedness: Phan et al. [2021].
- Self-normalised sums
 - Bercu et al. [2015], Bercu and Touati [2019]

Is this all for mean estimation?

Detour: stochastic bandit





(a) 15% chance to win 10 \in

(b) 5% chance to win 10€



The casino does not tell you which one is the winner!

- Optimism: play the machine with highest plausible reward.
- More formally:
 - ▶ play argmax_{k=1,2} $U_{\tau_k}^k(\delta)$,
 - τ_k : number of pulls to machine k (random stopping time),
 - $U_{\tau_k}^k(\delta)$ is a measurable function of τ_k samples from machine k such that

$$\mathbb{P}\left(\mu \geq U^k_{\boldsymbol{\tau}_k}(\delta)\right) \leq \delta.$$

Concentration bound for sample of random size!









Machine 1: 5 pulls. Machine 2: 30 pulls.





Machine 1: 100 pulls. Machine 2: 30 pulls.

Table of Contents



2 Time-uniform concentration: method of mixture

3 Bregman uniform concentration for generic exponential families

Martingale and stopping time

- Stopping times are hard to deal with...
- ✤ ... but go well with martingales!

For $t \in \mathbb{N}$, assume we know an invertible, nondecreasing $F_t \colon \mathbb{R} \to \mathbb{R}_+$ s.t

(i)
$$M_t = F_t(S_t)$$
 defines a supermartingale adapted to the filtration $\mathcal{F}_t = \sigma(X_s, s \leq t)$, with $S_t = \sum_{s=1}^t X_s - \mu$;

(*ii*) $\mathbb{E}[M_0] \leq 1$.

Then for any \mathcal{F} -stopping time τ ,

$$\mathbb{P}\left(S_{\tau} \geq F_{\tau}^{-1}\left(\frac{1}{\delta}\right)\right) = \mathbb{P}\left(M_{\tau} \geq \frac{1}{\delta}\right)$$

$$\leq \delta \mathbb{E}\left[M_{\tau}\right] \quad (Markov)$$

$$\leq \delta \mathbb{E}\left[M_{0}\right] \quad (Doob)$$

$$\leq \delta .$$

To find a martingale

The sub- ψ assumption is really a martingale condition:

$$\forall \lambda \in \mathcal{I}, \ M^{\lambda} = \left(e^{\lambda S_t - t\psi(\lambda)}\right)_{t \in \mathbb{N}} \text{ is a nonnegative supermartingale}.$$

Proposition

For any \mathcal{F} -stopping time τ ,

$$\forall \lambda \in \mathcal{I}, \ \mathbb{P}\left(\widehat{\mu}_{\tau} - \mu \geq \frac{1}{\lambda}\left(\frac{1}{\tau}\log\frac{1}{\delta} + \psi(\lambda)\right)\right) \leq \delta$$

Cannot optimise in λ for all values of τ ! (optimising for a fixed time t_0 recovers the Fenchel-Legendre formula).

To find a good martingale

For any mixture density $q(\lambda)$ over \mathcal{I} ,

$$M = \left(\int_{\mathcal{I}} e^{\lambda S_t - t\psi(\lambda)} q(\lambda) d\lambda
ight)_{t \in \mathbb{N}}$$
 is also a nonnegative supermartingale.

Proposition (Chernoff-Laplace mixture bound (sub-Gaussian))
If
$$\psi(\lambda) = \frac{\sigma^2 \lambda^2}{2}$$
, then for any \mathcal{F} -stopping time τ and any $c > 0$,

$$\mathbb{P}\left(\widehat{\mu}_{\tau} - \mu \ge \sigma \sqrt{\frac{2\left(1 + \frac{c}{\tau}\right)\log\left(\frac{\sqrt{\frac{\tau}{c} + 1}}{\delta}\right)}{\tau}}\right) \le \delta.$$

Remark: this corresponds to the mixture distribution $\mathcal{N}\left(0, \frac{1}{c}\right)$.

Table of Contents



Time-uniform concentration: method of mixture



Parametric family indexed by $\theta \in \Theta$ (open set) of distributions ν_{θ} over \mathbb{R}^d given by

$$rac{d
u_{ heta}}{d
u_{ heta_{\circ}}}(x) = h(x)e^{\langle heta, F(x)
angle - \mathcal{L}(heta)} \,.$$

• *F*: feature function (of $x \in \mathbb{R}^d$),

- \mathcal{L} : log-partition function (of $\theta \in \Theta$), convex, twice differentiable.
 - Assume det $\nabla^2 \mathcal{L}(\theta) > 0$ for all $\theta \in \Theta$.

Bregman divergence:

$$\begin{split} \mathcal{B}_{\mathcal{L}}(\theta',\theta) &= \mathcal{L}(\theta') - \mathcal{L}(\theta) - \langle \theta' - \theta, \nabla \mathcal{L}(\theta) \rangle \\ &= \mathsf{KL}\left(\nu_{\theta} \| \nu_{\theta'}\right) \,. \end{split}$$

Examples

Gaussian $\mathcal{N}\left(\mu,\sigma^2\right)$ with known variance σ^2

$$egin{aligned} & heta &= \mu, \Theta = \mathbb{R} \,, \ & \mathcal{B}_{\mathcal{L}}(heta', heta) &= rac{(heta' - heta)^2}{2\sigma^2} \end{aligned}$$

Gaussian $\mathcal{N}\left(\mu,\sigma^{2}\right)$

$$egin{aligned} & heta = \left(rac{\mu}{\sigma^2}, -rac{1}{2\sigma^2}
ight)^ op, \Theta = \mathbb{R} imes \mathbb{R}^*_- \,, \ &\mathcal{B}_{\mathcal{L}}(heta', heta) = rac{1}{2}\lograc{ heta_2}{ heta'_2} + rac{ heta'_2}{2 heta_2} - heta'_2 \left(rac{ heta'_1}{2 heta'_2} - rac{ heta_1}{2 heta_2}
ight)^2 - rac{1}{2}. \end{aligned}$$

Bernoulli $\mathcal{B}(p)$

$$egin{aligned} & heta = eta, \Theta = (0,1)\,, \ &\mathcal{B}_\mathcal{L}(heta', heta) = heta \log rac{ heta}{ heta'} + (1- heta) \log rac{1- heta}{1- heta'} \end{aligned}$$

Lemma

For $\theta \in \Theta$ and λ s.t $\theta + \lambda \in \Theta$,

$$\log \mathbb{E}_{ heta} \left[e^{\langle \lambda, F(X) - \mathbb{E}_{ heta}[F(X)]
angle}
ight] = \mathcal{B}_{\mathcal{L}}(heta + \lambda, heta) \,.$$

Consequence: let $\hat{\mu}_t = \frac{1}{t} \sum_{s=1}^t F(X_s)$ and $\mu = \mathbb{E}_{\theta} [F(X)]$, then

$$\forall \lambda, \ M^{\lambda} = \left(e^{\langle \lambda, t(\hat{\mu}_t - \mu) \rangle - t \mathcal{B}_{\mathcal{L}}(\theta + \lambda, \theta)} \right)_{t \in \mathbb{N}} \text{ is a nonnegative (super)martingale.}$$

Mixture: for c > 0,

$$egin{aligned} q_{ heta}(\lambda|c) \propto e^{\langle heta+\lambda, c
abla \mathcal{L}(heta)
angle - c \mathcal{L}(heta)}\,, \ M_t &= \int M_t^\lambda q_{ heta}(\lambda|c) d\lambda\,. \end{aligned}$$

Bregman-Laplace confidence set

Regularised parameter estimate:

$$\widehat{\theta}_{t,c}(\theta) = \nabla \mathcal{L}^{-1}\left(rac{t}{t+c}\widehat{\mu}_t + rac{c}{t+c}\mathcal{L}(\theta)
ight).$$

Bregman information gain:

$$\gamma_{t,c}(heta) = \log rac{\int_{\Theta} e^{-c\mathcal{B}_{\mathcal{L}}(heta', heta)} d heta'}{\int_{\Theta} e^{-(t+c)\mathcal{B}_{\mathcal{L}}(heta',\widehat{ heta}_{t,c}(heta))} d heta'}\,.$$

Theorem (Bregman-Laplace mixture bound for exponential families)

For any \mathcal{F} -stopping time τ and any c > 0,

$$\mathbb{P}\left((\tau+c)\mathcal{B}_{\mathcal{L}}\left(\theta,\widehat{\theta}_{\tau,c}(\theta)\right)\geq\log\frac{1}{\delta}+\gamma_{\tau,c}(\theta)\right)\leq\delta$$

Remarks

$$\mathbb{P}\left((\tau+c)\mathcal{B}_{\mathcal{L}}\left(\theta,\widehat{\theta}_{\tau,c}(\theta)\right)\geq\log\frac{1}{\delta}+\gamma_{\tau,c}(\theta)\right)\leq\delta$$

• Laplace's method for approximating integrals: when $t \to +\infty$,

$$\gamma_{t,c}(\theta) = rac{\dim \Theta}{2} \log \left(1 + rac{t}{c}\right) + \mathcal{O}(1) \,.$$

Implicit confidence set...

...but convex: easy numerical solution.

Numerical experiments



Numerical experiments



Figure: Gaussian (mean and variance) for $t \in \{10, 25, 50, 100\}$ observations

Numerical experiments



Comparison of median confidence envelopes around the mean for $\mathcal{B}(0.8)$. Grey lines are trajectories of empirical means $\hat{\mu}_n$.

Examples

Table 2: Summary of Bregman confidence sets given by Theorem 3.2 for representative families Throughout, the following notations are used:

$S_n = \sum_{t=1}^n X_t$, $\hat{\mu}_n = \frac{S_n}{n}$, $Z_n(\mu, \sigma) = \frac{1}{\sigma^2} \sum_{t=1}^n (X_t - \mu)^2$,
$S_n^{(k)} = \sum_{t=1}^n X_t^k, L_n = \sum_{t=1}^n \log X_t, K_n = \sum_{t=1}^n \log X_t/2,$
$I(a, b) = \int_{-\infty}^{+\infty} e^{-ae^{\theta}+b\theta} d\theta,$
$J(a, b) = \int_{0}^{\infty} \exp \left(-a \log \Gamma\left(\frac{k}{2}\right) + b\frac{k}{2}\right) dk$ (if $k \in \mathbb{R}_{+}$),
$J(a, b) = \sum_{k'=1}^{\infty} \exp \left(-a \log \Gamma \left(\frac{k'}{2}\right) + b \frac{k'}{2}\right)$ (if $k \in \mathbb{N}$).

Name	Parameters	Formula
Gaussian	$\mu \in \mathbb{R}$	$\frac{1}{n+c}\frac{(S_n-n\mu)^2}{2\sigma^2}\leqslant \log\frac{1}{\delta}+\frac{1}{2}\log\frac{n+c}{c}$
Gaussian	$\sigma \in \mathbb{R}_+$	$\begin{array}{l} \frac{Q_n(\mu)}{2\sigma^2} - \left(\frac{n+\varepsilon}{2} + 1\right) \log \left(\frac{Q_n(\mu)}{2\sigma^2} + \frac{\varepsilon}{2}\right) \\ \leqslant \log \frac{1}{\delta} + \log \frac{\Gamma\left(\frac{\delta}{2}\right)}{\Gamma\left(\frac{n+\varepsilon}{2}\right)} - \log \left(\frac{n+\varepsilon}{2}\right) - \frac{\varepsilon}{2} \log \frac{\varepsilon}{2} \end{array}$
Gaussian	$\begin{array}{l} \mu \in \mathbb{R} \\ \sigma \in \mathbb{R}_+ \end{array}$	$\begin{split} & \frac{1}{2}Z_n(\mu,\sigma) - \frac{n+c+3}{2}\log\left(\frac{n}{n+c}Z_n(\hat{\mu},\sigma) + \frac{c}{n+c}Z_n(\mu,\sigma) + c\right) \\ & \leq \log\frac{1}{\delta} - \frac{n}{2}\log 2 - \left(\frac{c}{2}+2\right)\log c + \frac{1}{2}\log\left(n+c\right) \\ & + \log\Gamma\left(\frac{c+3}{2}\right) - \log\Gamma\left(\frac{n+c+3}{2}\right) \end{split}$
Bernoulli	$\mu \in [0,1]$	$\begin{array}{l} S_n \log \frac{1}{\mu} + (n - S_n) \log \frac{1}{1 - \mu} + \log \frac{\Gamma(S_n + c\mu)\Gamma(n - S_n + c(1 - \mu))}{\Gamma(c\mu)\Gamma(c(1 - \mu))} \\ \leqslant \log \frac{1}{\delta} + \log \frac{\Gamma(n + c)}{\Gamma(c)} \end{array} - \end{array}$
Exponential	$\mu \in \mathbb{R}_+$	$ \begin{split} & \frac{S_n}{\mu} - (n+c+1) \log \left(\frac{S_n}{\mu} + c \right) \\ & \leqslant \log \frac{1}{\theta} + \log \frac{\Gamma(c)}{\Gamma(n+c)} - \log(n+c) - c \log c \end{split} $

Chi-square	$k \in \mathbb{N}$ or $k \in \mathbb{R}_+$	$ \begin{array}{l} n \log \Gamma\left(\frac{k}{2}\right) - \frac{k}{2}K_n - \log J\left(c, c\psi_0\left(\frac{k}{2}\right)\right) \\ + \log J\left(n + c, K_n + c\psi_0\left(\frac{k}{2}\right)\right) \leqslant \log \frac{1}{\delta} \end{array} $
Poisson	$\lambda \in \mathbb{R}_+$	$\begin{array}{l} n\lambda\!-\!S_n\log\lambda \\ \leqslant \log\frac{1}{\delta}\!+\!\log I\left(c,c\lambda\right) - \log I\left(n\!+\!c,S_n\!+\!c\lambda\right) \end{array}$
Pareto	$\alpha \in \mathbb{R}$	$\begin{aligned} \alpha L_n &- (n+c+1)\log\left(\alpha L_n + c\right) \\ &\leqslant \log \frac{1}{\delta} + \log \frac{\Gamma(c)}{\Gamma(n+c)} - \log(n+c) - c\log c \end{aligned}$
Weibull	$\lambda \in \mathbb{R}_+$	$ \begin{split} & \frac{S_{\lambda}^{(k)}}{\lambda^{k}} - (n+c+1)\log\left(\frac{S_{\lambda}^{(k)}}{\lambda^{k}} + c\right) \\ & \leqslant \log \frac{1}{\delta} + \log \frac{\Gamma(c)}{\Gamma(n+c)} - \log(n+c) - c\log c \end{split} $
Gamma	$\lambda \in \mathbb{R}_+$	$\begin{array}{l} \frac{S_n}{\lambda} - ((n+c)k+1)\log(\frac{S_n}{\lambda}+ck) \\ \leqslant \log \frac{1}{\delta} + \log \frac{\Gamma(ck)}{\Gamma((n+c)k)} - \log((n+c)k) - ck\log ck \end{array}$

Questions?



References I

- V. Bentkus. On hoeffding's inequalities. The Annals of Probability, 32(2):1650-1673, 2004.
- B. Bercu and T. Touati. New insights on concentration inequalities for self-normalized martingales. Electronic Communications in Probability, 24:1–12, 2019.
- B. Bercu, B. Delyon, and E. Rio. Concentration inequalities for sums and martingales. Springer, 2015.
- S. Boucheron, G. Lugosi, and P. Massart. <u>Concentration inequalities: A nonasymptotic theory of</u> independence. Oxford university press, 2013.
- A. K. Kuchibhotla and Q. Zheng. Near-optimal confidence sequences for bounded random variables. In International Conference on Machine Learning, pages 5827–5837. PMLR, 2021.
- A. Maurer and M. Pontil. Empirical bernstein bounds and sample variance penalization. 2009.
- M. Phan, P. Thomas, and E. Learned-Miller. Towards practical mean bounds for small samples. In International Conference on Machine Learning, pages 8567–8576. PMLR, 2021.

Table of Contents



A bad idea: union bound

Say we have:

$$orall t \in \{1,\ldots,T\}, \ \mathbb{P}\left(\mu \in \mathcal{C}_t^\delta
ight) \geq 1-\delta$$
 .

Then:

$$\mathbb{P}\left(\forall t \in \{1, \dots, T\}, \ \mu \in \mathcal{C}_t^{\delta}\right) = 1 - \mathbb{P}\left(\bigcup_{t=1}^T \{\mu \notin \mathcal{C}_t^{\delta}\}\right)$$
$$\geq 1 - \sum_{t=1}^T \mathbb{P}\left(\mu \notin \mathcal{C}_t^{\delta}\right)$$
$$\geq 1 - T\delta.$$

A bad idea: union bound

Say we have:

$$orall t \in \{1,\ldots,T\}, \ \mathbb{P}\left(\mu \in \mathcal{C}_t^\delta
ight) \geq 1-\delta$$
 .

Then:

$$\mathbb{P}\left(\forall t \in \{1, \dots, T\}, \ \mu \in \mathcal{C}_t^{\delta}\right) = 1 - \mathbb{P}\left(\bigcup_{t=1}^T \{\mu \notin \mathcal{C}_t^{\delta}\}\right)$$
$$\geq 1 - \sum_{t=1}^T \mathbb{P}\left(\mu \notin \mathcal{C}_t^{\delta}\right)$$
$$\geq 1 - T\delta.$$

Proposition (Union bound)

 $(\mathcal{C}_t^{\delta/T})_{t=1}^T$ is a uniform confidence sequence at level δ .

Proposition (Doob's maximal inequality)

Let $(S)_{t\in\mathbb{N}}$ be a **nonnegative supermartingale** w.r.t a filtration $(\mathcal{F}_t)_{t\in\mathbb{N}}$. Then for all $p \ge 1$ and $\epsilon > 0$, it holds for all $T_0 < T$ that

$$\mathbb{P}\left(\max_{T_0 \leq t \leq T} S_t \geq \epsilon\right) \leq \frac{\mathbb{E}\left[S_{T_0}^{p}\right]}{\epsilon^{p}}$$

Idea: apply union bound over a geometric grid $(t_k)_{k \in \mathbb{N}}$ with $t_k = (1 + \eta)^k$.

Proposition (Doob's maximal inequality)

Let $(S)_{t\in\mathbb{N}}$ be a **nonnegative supermartingale** w.r.t a filtration $(\mathcal{F}_t)_{t\in\mathbb{N}}$. Then for all $p \ge 1$ and $\epsilon > 0$, it holds for all $T_0 < T$ that

$$\mathbb{P}\left(\max_{T_0 \leq t \leq T} S_t \geq \epsilon\right) \leq \frac{\mathbb{E}\left[S_{T_0}^{p}\right]}{\epsilon^{p}}$$

Idea: apply union bound over a geometric grid $(t_k)_{k \in \mathbb{N}}$ with $t_k = (1 + \eta)^k$.

Proposition (Geometric time-peeling)

Let au a $(\mathcal{F}_t)_{t\in\mathbb{N}}$ -stopping time. Then for $\eta>0$ and $\delta\in(0,1)$,

$$\mathbb{P}\left(\widehat{\mu}_{\tau} - \mu \geq \psi_*^{-1}\left(\frac{1 + \eta}{\tau} \log\left(\frac{\log(\tau) \log\left((1 + \eta)\tau\right)}{\delta \log^2(1 + \eta)}\right)\right) \leq \delta$$