

Fundamental Concentration Inequalities

Guillaume Obozinski

Swiss Data Science Center
EPFL & ETH Zürich



RLSS, juillet 2019, Lille

What and why?

- **What is the concentration of measure phenomenon?**

This refers to the phenomenon that there are certain ways to combine random variables that produce r.v. that are *concentrated* around their expectation. One of the main case of interest are averages of independent variables.

- **Why do we need it for reinforcement learning?**

RL require to make decisions in the presence of “uncertain uncertainty”, r.v.s whose distributions are not known initially. This requires to be able to produce confidence intervals (or confidence regions) for these r.v. in the environment that are not yet know, but that are typically being learned in the RL algorithm.

- **Why is the central limit theorem not sufficient?**

The CLT only produces asymptotic CIs with an error which is a priori not quantified.

Union bound

Let A_1, A_2, \dots, A_k be events. We have

$$\mathbb{P}(A_1 \cup A_2 \cup \dots \cup A_k) \leq \mathbb{P}(A_1) + \mathbb{P}(A_2) + \dots + \mathbb{P}(A_k)$$

Proof.

$$\mathbb{E}[\mathbf{1}_{A_1 \cup A_2 \cup \dots \cup A_k}] \leq \mathbb{E}[\mathbf{1}_{A_1} + \mathbf{1}_{A_2} + \dots + \mathbf{1}_{A_k}].$$

□

Example

Let $X_t \sim \mathcal{N}(0, \sigma^2)$ (not necessarily independent)

$$\mathbb{P}(\max_t X_t > x) = \mathbb{P}(\bigcup_t \{X_t > x\}) \leq \sum_{t=1}^T \mathbb{P}(X_t > x) \leq T \exp\left(-\frac{x^2}{2\sigma^2}\right)$$

So with probability $1 - \delta$, we have

$$X \leq \sigma \sqrt{2 \log \frac{T}{\delta}}$$

Markov, Chebychev and Chernoff

Markov inequality

If $X \geq 0$ a.s. and $t > 0$, then
$$\mathbb{P}(X > t) \leq \frac{\mathbb{E}[X]}{t}$$

Chebychev inequality

$$\forall t > 0, \quad \mathbb{P}(X - \mathbb{E}[X] > t) \leq \frac{\text{Var}(X)}{t^2}$$

Chernoff inequality

$$\mathbb{P}(X > t) \leq \inf_{r \geq 0} e^{-rt} \mathbb{E}[e^{rX}]$$

Note that $r \mapsto \mathbb{E}[e^{rX}]$ is the *moment generating function* (MGF) of X .

Cramér-Chernoff Method

$\forall r > 0,$

$$\mathbb{P}(X > t) \leq e^{-rt} \mathbb{E}[e^{rX}] = \exp(\psi_X(r) - rt) \quad \text{for } \psi_X(r) = \log \mathbb{E}[e^{rX}]$$

ψ is the log MGF of X , aka *cumulant generating function* if $\mathbb{E}[X] = 0$.

Since this true for all $r \geq 0$ if

$$\psi_X^*(t) = \sup_{r \geq 0} rt - \psi_X(r),$$

then we have

$$\mathbb{P}(X > t) \leq \exp(-\psi_X^*(t))$$

- ψ_X^* is called the Cramér transform of X
- If $t \geq \mathbb{E}[X]$, then $\psi_X^*(t) = \sup_{r \in \mathbb{R}} rt - \psi_X(r)$,
i.e., ψ_X^* is the Fenchel-Legendre conjugate of ψ_X .

Applying the Cramér-Chernoff to the Gaussian

Let $X \sim \mathcal{N}(0, \sigma^2)$, then

$$\mathbb{E}[e^{rX}] = e^{\frac{r^2\sigma^2}{2}}, \quad \psi(r) = \frac{r^2\sigma^2}{2}, \quad \psi^*(t) = \frac{t^2}{2\sigma^2},$$

So that $1 - \Phi(t) := \mathbb{P}(X > t) \leq e^{-\psi^*(t)} = e^{-\frac{t^2}{2\sigma^2}}$.

But it is well-known that for all $t > 0$,

$$\left(\frac{1}{t} - \frac{1}{t^3}\right) \cdot \frac{1}{\sqrt{2\pi\sigma^2}} e^{-\frac{t^2}{2\sigma^2}} \leq 1 - \Phi(t) \leq \frac{1}{t} \cdot \frac{1}{\sqrt{2\pi\sigma^2}} e^{-\frac{t^2}{2\sigma^2}}.$$

In fact

$$\sup_{t \geq 0} (1 - \Phi(t)) e^{\frac{t^2}{2\sigma^2}} = \frac{1}{2}.$$

So the Cramér-Chernoff produces a relatively good bound.

MGF inequality for bounded r.v.

Bernoulli r.v. X

For $X_B \sim \text{Ber}(\theta)$, we have $\mathbb{E}[e^{sX_B}] = 1 - \theta + \theta e^s$

Any bounded r.v. X on $[0, 1]$

If $\mathbb{E}[X] = \theta$, $\forall s \in \mathbb{R}$, we have $\mathbb{E}[e^{sX}] \leq \mathbb{E}[(1 - X) + Xe^s] = 1 - \theta + \theta e^s$

So

$$\mathbb{E}[e^{sX}] \leq \mathbb{E}[e^{sX_B}] = 1 - \theta + \theta e^s.$$

And

$$\mathbb{E}[e^{s(X-\theta)}] \leq \mathbb{E}[e^{s(X_B-\theta)}] = (1 - \theta + \theta e^s) e^{-s\theta} = e^{\phi(s)},$$

with

$$\phi(s) := \log(1 - \theta + \theta e^s) - s\theta.$$

Key inequality (Hoeffding's Lemma)

Let $\phi(s) := \log(1 - \theta + \theta e^s) - s\theta$. We have $\phi(s) \leq \frac{s^2}{8}$.

Proof.

By Taylor-Lagrange $\phi(s) = \phi(0) + s\phi'(0) + \frac{s^2}{2}\phi''(t)$ with $t \in (0, s)$.

$$\phi'(t) + \theta = \frac{\theta e^t}{1 - \theta + \theta e^t} = \frac{1}{1 + \alpha e^{-t}} \quad \text{with} \quad \alpha = \frac{1-\theta}{\theta}.$$

$$\phi''(t) = \frac{\alpha e^{-t}}{(1 + \alpha e^{-t})^2} = \phi'(t)(1 - \phi'(t)) \leq \frac{1}{4} \quad \text{since} \quad \phi'(t) \leq 1.$$

So $\phi(0) = 0$, $\phi'(0) = 0$ and, by T.-L.,

$$\phi(s) \leq \frac{s^2}{2}\phi''(t) \leq \frac{s^2}{8}.$$



Bounded r.v. are sub-Bernoulli and thus sub-Gaussian

Let

- X be a r.v. on $[0, 1]$ with $\mathbb{E}[X] = \theta$
- $X_B \sim \text{Ber}(\theta)$
- $X_G \sim \mathcal{N}(0, \frac{1}{4})$
- $\phi(s) := \log(1 - \theta + \theta e^s) - s\theta$.

Then $\mathbb{E}[e^{s(X-\theta)}] \leq \mathbb{E}[e^{s(X_B-\theta)}] = e^{\phi(s)} \leq e^{\frac{s^2}{8}} = \mathbb{E}[e^{sX_G}]$.

Bounded r.v. are sub-Bernoulli and thus sub-Gaussian

Let

- X be a r.v. on $[0, 1]$ with $\mathbb{E}[X] = \theta$
- $X_B \sim \text{Ber}(\theta)$
- $X_G \sim \mathcal{N}(0, \frac{1}{4})$
- $\phi(s) := \log(1 - \theta + \theta e^s) - s\theta$.

Then $\forall s \geq 0$, $\mathbb{E}[e^{s(X-\theta)}] \leq \mathbb{E}[e^{s(X_B-\theta)}] = e^{\phi(s)} \leq e^{\frac{s^2}{8}} = \mathbb{E}[e^{sX_G}]$.

Now, let

- Y be a random variable on the interval $[a, b]$
- $X := \frac{Y - a}{b - a} \in [0, 1]$ so that $Y = (b - a)X + a$.
- $\tilde{Y} = Y - \mathbb{E}[Y]$, $\tilde{X} = X - \mathbb{E}[X]$, $\tilde{X}_B = X_B - \mathbb{E}[X_B]$,

We have $\tilde{Y} = (b - a)\tilde{X}$ and

$$\mathbb{E}[e^{s\tilde{Y}}] = \mathbb{E}[e^{s(b-a)\tilde{X}}] \leq \mathbb{E}[e^{s(b-a)\tilde{X}_B}] = e^{\phi(s(b-a))} \leq e^{\frac{s^2(b-a)^2}{8}}.$$

Hoeffding inequality

Let X_i be *independent* bounded r.v. such that

- $\mathbb{E}[X_i] = 0$ and X_i has support in $[a_i, b_i]$.

Let $\tau^2 := \frac{1}{n} \sum_i \tau_i^2$ with $\tau_i^2 := \frac{1}{4}(b_i - a_i)^2$. Note that $\text{Var}(X_i) \leq \tau_i^2$.

Then $\forall x \geq 0$, $\mathbb{P}(\bar{X} \geq x) \leq \exp\left(-\frac{nx^2}{2\tau^2}\right)$ with $\bar{X} := \frac{1}{n} \sum_{i=1}^n X_i$.

Proof. $\mathbb{P}(\sum_i X_i \geq nx) = \mathbb{P}\left(\exp\left(s \sum_i X_i\right) \geq \exp(snx)\right)$

$$\leq e^{-snx} \mathbb{E}\left[\prod_i e^{sX_i}\right] = e^{-snx} \prod_i \mathbb{E}\left[e^{sX_i}\right]$$
$$\stackrel{\text{sub-G}}{\leq} \exp\left(-snx + \frac{s^2}{8} \sum_i (b_i - a_i)^2\right)$$
$$= \exp\left(-snx + \frac{s^2}{2} n\tau^2\right)$$

Thus $\mathbb{P}(\sum_i X_i \geq nx) \leq \exp\left(-\frac{nx^2}{2\tau^2}\right)$ by setting $s = \frac{x}{n\tau^2} \geq 0$

which minimizes the RHS w.r.t. s . □

Comparing Hoeffding with the CLT

Let X_i be *independent* bounded r.v. such that

- $\mathbb{E}[X_i] = 0$ and X_i has support in $[a_i, b_i]$.
- Let $\tau^2 := \frac{1}{n} \sum_i \tau_i^2$ with $\tau_i^2 := \frac{1}{4}(b_i - a_i)^2$.
- Let $\sigma^2 := \frac{1}{n} \sum_i \sigma_i^2$ with $\sigma_i^2 = \text{Var}(X_i) \leq \tau_i^2$.

By the CLT:

$$\sqrt{n}\bar{X} \xrightarrow{(d)} X^* \quad \text{with} \quad X^* \sim \mathcal{N}(0, \sigma^2)$$

We can compare:

Hoeffding:

$$\mathbb{P}\left(\sqrt{n}\bar{X} > x\right) \leq \exp\left(-\frac{x^2}{2\tau^2}\right)$$

CLT:

$$\mathbb{P}\left(\sqrt{n}\bar{X} \geq x\right) \xrightarrow{n \rightarrow \infty} \mathbb{P}\left(X^* \geq x\right) \leq \frac{1}{x} \frac{1}{\sqrt{2\pi}} \exp\left(-\frac{x^2}{2\sigma^2}\right)$$

High probability statement of Hoeffding's inequality

As before let $\tau^2 = \frac{1}{n} \sum_{i=1}^n (b_i - a_i)^2$.

Hoeffding inequality

$$\mathbb{P}(\bar{X} > x) \leq \exp\left(-\frac{nx^2}{2\tau^2}\right)$$

By setting the RHS to δ , we obtain the following reformulation.

High probability statement:

With probability $1 - \delta$, $\bar{X} \leq \sqrt{\frac{\tau^2}{n} \cdot 2 \log\left(\frac{1}{\delta}\right)}$.

Or equivalently $\sum_{i=1}^n X_i \leq \sqrt{\sum_{i=1}^n (b_i - a_i)^2} \sqrt{2 \log\left(\frac{1}{\delta}\right)}$.

Sharper than Hoeffding: the Chernoff-Hoeffding inequality

If X_i are independent r.v. on $[0,1]$ with $\mathbb{E}[X_i] = \theta_i$, then

$$\mathbb{P}\left(\frac{1}{n} \sum_i X_i \geq q\right) \leq \exp\left(-n \text{KL}(q \parallel \theta)\right)$$

$$\text{with } \text{KL}(q \parallel \theta) = q \log \frac{q}{\theta} + (1 - q) \log \frac{1 - q}{1 - \theta}.$$

Proof

$$\begin{aligned} \mathbb{P}\left(\sum_i X_i \geq nq\right) &= \mathbb{P}\left(\exp\left(s \sum_i X_i\right) \geq \exp(snq)\right) \\ &\leq e^{-snq} \mathbb{E}\left[\prod_i e^{sX_i}\right] = e^{-snq} \prod_i \mathbb{E}\left[e^{sX_i}\right] \\ &= e^{-snq} \prod_i (1 - \theta_i + \theta_i e^s) \\ &\leq e^{-snq} (1 - \theta + \theta e^s)^n \quad \text{with } \theta = \frac{1}{n} \sum \theta_i, \end{aligned}$$

by the arithmetico-geometric inequality.

Let $\psi(s) = n \log(1 - \theta + \theta e^s)$. Then $\psi'(s^*) - nq = 0$ iff

$$\frac{\theta e^{s^*}}{1 - \theta + \theta e^{s^*}} = q \quad \Leftrightarrow \quad e^{s^*} = \frac{q}{1 - q} \frac{1 - \theta}{\theta}.$$

Sharper than Hoeffding: the Chernoff-Hoeffding inequality

We found $\psi'(s^*) - nq = 0$ iff

$$\frac{\theta e^{s^*}}{1 - \theta + \theta e^{s^*}} = q \quad \Leftrightarrow \quad e^{s^*} = \frac{q}{1 - q} \frac{1 - \theta}{\theta}.$$

$$\begin{aligned} \log \mathbb{P}(\sum_i X_i \geq nq) &\leq n \log \left(\frac{\theta e^{s^*}}{q} \right) - s^* nq \\ &\leq n \log \frac{\theta}{q} + s^* n(1 - q) = n \log \frac{\theta}{q} + n(1 - q) \left[\log \frac{1 - \theta}{1 - q} - \log \frac{\theta}{q} \right] \\ &= -nq \log \frac{q}{\theta} - n(1 - q) \log \frac{1 - q}{1 - \theta} = -n \text{KL}(q \parallel \theta) \end{aligned}$$

Bennett's inequality

Let X_i be *independent* bounded r.v. such that

- $\mathbb{E}[X_i] = 0$ and $\mathbb{P}(X_i \leq 1) = 1$.
- Let $\sigma^2 := \frac{1}{n} \sum_i \sigma_i^2$ with $\sigma_i^2 = \text{Var}(X_i) \leq \tau_i^2$.

Then

$$\mathbb{P}\left(\frac{1}{n} \sum_i X_i > x\right) \leq \exp\left(-n\sigma^2 h\left(\frac{x}{\sigma^2}\right)\right)$$

for $h(u) = (1+u)\log(1+u)$

Or equivalently

$$\mathbb{P}\left(\frac{1}{n} \sum_i X_i > x\right) \leq \exp\left(-n\left(\sigma^2 + x\right) \log\left(1 + \frac{x}{\sigma^2}\right)\right)$$

see, e.g. Boucheron et al. (2003) for a proof.

Bernstein's Inequality

Bennett's inequality: $\mathbb{P}\left(\frac{1}{n} \sum_i X_i > x\right) \leq \exp\left(-n\sigma^2 h\left(\frac{x}{\sigma^2}\right)\right)$

for $h(u) = (1+u)\log(1+u)$ but $h(u) \geq \frac{1}{2} \frac{u^2}{1+u/3}$ which implies

Bernstein's inequality

$$\mathbb{P}\left(\frac{1}{n} \sum_i X_i > x\right) \leq \exp\left(-\frac{nx^2}{2(\sigma^2 + x/3)}\right)$$

compare with Hoeffding's inequality

$$\mathbb{P}\left(\frac{1}{n} \sum_i X_i > x\right) \leq \exp\left(-\frac{nx^2}{2\tau^2}\right)$$

- If $x \ll \sigma^2$ this captures the right asymptotic variance
- If $\sigma^2 + x/3 \geq \tau^2$ then this is worse than Hoeffding
- But when $\sigma^2 + x/3 < \tau^2$ it captures relevant behavior for small σ^2
 - e.g. $\text{Bin}(n, \lambda/n) \rightarrow \text{Poisson}(\lambda)$ with tail in $e^{-\lambda}$.

High probability statement of Bernstein's inequality

Bernstein's inequality

$$\mathbb{P}\left(\frac{1}{n} \sum_i X_i > x\right) \leq \exp\left(-\frac{nx^2}{2(\sigma^2 + x/3)}\right)$$

By solving for x in $t = nx^2/(2(\sigma^2 + x/3))$ we get

$$x = \frac{t}{3n} + \sqrt{\frac{t^2}{9n^2} + \frac{2\sigma^2 t}{n}} \geq \frac{t}{3n} + \sqrt{\frac{2\sigma^2 t}{n}},$$

we get

$$\mathbb{P}\left(\frac{1}{n} \sum_i X_i > \sqrt{\frac{2\sigma^2 t}{n}} + \frac{t}{3n}\right) \leq e^{-t}$$

So that with probability $1 - \delta$, we have

$$\frac{1}{n} \sum_i X_i > \sqrt{\frac{2\sigma^2 \log(\frac{1}{\delta})}{n}} + \frac{\log(\frac{1}{\delta})}{3n}$$

References I

- Boucheron, S., Lugosi, G., and Bousquet, O. (2003). Concentration inequalities. In *Summer School on Machine Learning*, pages 208–240. Springer.
- Massart, P. (2003). Concentration inequalities and model selection. Lectures from the 33rd Saint-Flour Summer School on Probability Theory. *Lecture Notes in Mathematics*, 1896.
- Vershynin, R. (2018). *High-dimensional probability: An introduction with applications in data science*, volume 47. Cambridge University Press.