# Concentration Inequalities with Machine Learning Applications

Stéphane Boucheron

LRI UMR CNRS 8623 Université Paris-Sud

# Organization

- From exponential inequalities to the concentration of measure phenomenon
- Concentration inequalities using the entropy method
- Learning-theoretical applications
- Moment inequalities using the generalized entropy method

# From exponential inequalities to concentration

## Lecture I: from exponential inequalities to concentration

- Introduction and Motivations
- Path to Bernstein inequality
- Martingales with bounded increments
- Efron-Stein inequality

# Lecture I: Roadmap

The introduction describes traditional exponential inequalities (Hoeffding/Bernstein) as non-asymptotic counterparts to limit theorems for sums of independent random variables. Concentration inequalities are presented as upper-bounds on tail probabilities for functions of many independent random variables. The scope of concentration inequalities is illustrated on a combinatorial optimization problem.

The path to Bernstein inequality is described in detail, stressing the fact that good bounds on the Log-Laplace transform of a random variable provide exponential bounds on the tail probabilities. The main topic of this course will be the derivation of Bernstein-like inequalities for general functions.

Martingales methods provide a general recipe for constructing Bernstein-like inequalities. The exponential super-martingale associated with martingales with bounded increments allows to refine the celebrated bounded-differences inequality. Despite and because of their generality, using martingale methods may be quite difficult. This has prompted the search for more user-friendly methods as (for example) the entropy method.

The first step in the entropy method is illustrated by the Efron-Stein inequality. The latter inequality provides a general and often tight upper-bound on the the variance of general functions of independent random variables. The Efron-Stein bound is first illustrated on a combinatorial optimization problem.

Basic concern of Probability Theory : sums of independent bounded random variables

$$(X_i)_{i \le n}$$
 i.i.d. in  $[-1, 1]$   $\mathbb{E}[X_i] = 0$   $\mathbb{E}[X_i^2] \le v$ 

$$Z \stackrel{\Delta}{=} \sum_{i=1}^{n} X_i$$

Basic concern of Probability Theory : sums of independent bounded random variables

$$(X_i)_{i < n}$$
 i.i.d. in  $[-1, 1]$   $\mathbb{E}[X_i] = 0$   $\mathbb{E}[X_i^2] \le v$ 

$$Z \stackrel{\Delta}{=} \sum_{i=1}^{n} X_i$$

Law of Large Numbers (LLN):

 $\forall \epsilon > 0, \quad \lim_n \mathbb{P}\left\{ |Z - \mathbb{E}[Z]| > n\epsilon \right\} = 0$ 

Central Limit Theorem (CLT, Invariance principle)

$$\forall x \in \mathbb{R} \quad \lim_{n} \mathbb{P}\left\{\frac{1}{\sqrt{n}}\left(Z - \mathbb{E}[Z]\right) < x\right\} = \int_{-\infty}^{x} \frac{1}{\sqrt{2\pi\nu}} e^{-\frac{x^2}{2\nu^2}} dx$$

LLN and CLT : asymptotic statements.

 $\mathbb{R}$  Need for statements dealing with fixed values of n.

Central Limit Theorem (CLT, Invariance principle)

$$\forall x \in \mathbb{R} \quad \lim_{n} \mathbb{P}\left\{\frac{1}{\sqrt{n}}\left(Z - \mathbb{E}[Z]\right) < x\right\} = \int_{-\infty}^{x} \frac{1}{\sqrt{2\pi\nu}} e^{-\frac{x^2}{2\nu^2}} dx$$

#### LLN and CLT : asymptotic statements. $\square$ Need for statements dealing with fixed values of n.





Probability Theory provides informations about the rate of convergence in the CLT. (Berry-Esseen, 1943)

$$\mathbb{E}[|X_i|^3] < \infty \quad \to \sup_x \left| \mathbb{P}\left\{ \frac{1}{\sqrt{n}} \left( Z - \mathbb{E}[Z] \right) < x \right\} - \int_{-\infty}^x \frac{1}{\sqrt{2\pi v}} e^{-\frac{x^2}{2v^2}} dx \right| \le \frac{C}{\sqrt{n}}$$

Berry-Esseen Theorems say nothing meaningful about

$$\left| \mathbb{P}\left\{ \left( Z - \mathbb{E}[Z] \right) > nx \right\} - \int_{x}^{\infty} \frac{1}{\sqrt{2\pi v}} e^{-\frac{nx^{2}}{2v^{2}}} dx \right|$$
arge deviations depend upon the P.D. of  $X_{i}$ .

Probability Theory provides informations about the rate of convergence in the CLT. (Berry-Esseen, 1943)

$$\mathbb{E}[|X_i|^3] < \infty \quad \to \sup_x \left| \mathbb{P}\left\{ \frac{1}{\sqrt{n}} \left( Z - \mathbb{E}[Z] \right) < x \right\} - \int_{-\infty}^x \frac{1}{\sqrt{2\pi v}} e^{-\frac{x^2}{2v^2}} dx \right| \le \frac{C}{\sqrt{n}}$$

Berry-Esseen Theorems say nothing meaningful about

$$\left| \mathbb{P}\left\{ \left( Z - \mathbb{E}[Z] \right) > nx \right\} - \int_{x}^{\infty} \frac{1}{\sqrt{2\pi v}} e^{-\frac{nx^{2}}{2v^{2}}} dx$$
  
Large deviations depend upon the P.D. of  $X_{i}$ .

If X is Gaussian: 
$$\mathbb{P}\left\{X \ge t\right\} = \int_{t}^{\infty} \frac{1}{x\sqrt{2\pi v}} x e^{-\frac{x^2}{2v}} dx$$
  
$$= \left[\frac{-v}{x} \frac{1}{\sqrt{2\pi v}} e^{-\frac{x^2}{2v}}\right]_{t}^{\infty} + \int_{t}^{\infty} \frac{\sqrt{v}}{\sqrt{2\pi}} \frac{1}{x^2} e^{-\frac{x^2}{2v}} dx$$
$$\leq \frac{\sqrt{v}}{\sqrt{2\pi t}} e^{-\frac{t^2}{2v}}$$

Gaussian tail bounds:

$$\mathbb{P}\left\{X \ge t\right\} \le \frac{\sqrt{v}}{\sqrt{2\pi t}}e^{-\frac{t^2}{2v}}$$

What about sums of i.i.d. random variables ? **Exponential inequalities** provide a widely applicable answer to those questions:

If Z is a sum of n independent centered [-1, 1]-variables:

 $\mathbb{P}\left\{Z \ge \mathbb{E}[Z] + t\right\} \le \exp\left(-\frac{t^2}{2n}\right) \qquad \text{Hoeffding}$  $\mathbb{P}\left\{Z \ge \mathbb{E}[Z] + t\right\} \le \exp\left(-\frac{t^2}{2(\operatorname{Var}[Z] + t/3)}\right) \qquad \text{Bernstein.}$ 

Gaussian-like tail bounds.

In Bernstein inequality, *n* does not appear in the exponent,

Bernstein inequality looks dimension-free.

## **Motivations**

"While exponential inequalities for sums of independent random variables are at the core of classical probabilities, the new abstract inequalities are far-reaching extensions that apply to considerably more general functions."

M. Talagrand. Ann. Probability, 1996

## **Motivations**

"While exponential inequalities for sums of independent random variables are at the core of classical probabilities, the new abstract inequalities are far-reaching extensions that apply to considerably more general functions."

M. Talagrand. Ann. Probability, 1996

"Any function of many independent random variables that depend on many of them but not too much on each of them is essentially constant."

M. Talagrand. Inventiones Mathematicae, 1996

## **Motivations**

"While exponential inequalities for sums of independent random variables are at the core of classical probabilities, the new abstract inequalities are far-reaching extensions that apply to considerably more general functions."

M. Talagrand. Ann. Probability, 1996

"Any function of many independent random variables that depend on many of them but not too much on each of them is essentially constant."

M. Talagrand. Inventiones Mathematicae, 1996

- What's in depend not too much on each of the variables ?
- How general should be the considerably more general functions ?
- Combinatorial optimization problems provide a cavalcade of illustrations of the concentration of measure phenomenon.

 $X_i$  are intervals from [0, 1]: extremities are picked by picking random numbers from [0, 1]

A packing is a set of pairwise disjoint intervals. *Z*: maximum cardinality of a packing extracted from  $X_1, \ldots, X_n$ 



*Z*: maximum cardinality of a packing extracted from  $X_1, \ldots, X_n$ 



 $\square$ : a maximum packing can be constructed (in  $\Theta(n \log n)$  operations) by a greedy algorithm.

1 sort intervals according their right-most extremities in ascending order.



- 2 the initial packing contains the first interval
- 3 scan the remaining intervals in ascending order
- 4 add the current interval to the packing if its left-most extremity is larger than the right-most extremity of the last interval in the packing.

- How does  $\mathbb{E}[Z]$  behave as  $\sqrt{n}$  goes to infinity?
- Law of large numbers? CLT? Asymptotic theorems ?
- **D** Tail behavior for fixed n ?
- $\checkmark$  No obvious way to represent Z as a sum of independent random variables.
- Z is a kind of "much more general function" of many independent random variables does not depend to much on each of them.



1 interval  $\leftrightarrow$  1 point (x, y) in  $[0, 1]^2$  with  $x \leq y$ .

Example: 10000 random intervals, Z = 114,  $\mathbb{E}[Z] \approx \frac{2}{\sqrt{\pi}}\sqrt{n} \approx 113$ 

Z satisfies a Central Limit Theorem:

 $(Z - \mathbb{E}[Z])/n^{1/4}$  is asymptotically Gaussian with variance  $2(4 - \pi)/\pi^{3/2}$ .

If  $Z \ge t$ , then there exists a sequence of t intervals  $X_{i_1}, X_{i_2}, \ldots, X_{i_t}$  that witnesses this fact.

This is just a sequence of *t* disjoint intervals !

A set property (P) is hereditary if every subset of a set satisfying (P) also satisfies (P)

Z is a configuration function: Z is the largest cardinality of a subset of the random variables that satisfy some hereditary property.

Concentration inequalities assert that:

$$\mathbb{P}\left\{Z \ge \mathbb{E}[Z] + t\right\} \le \exp\left(-\frac{t^2}{2(\mathbb{E}[Z] + t/3)}\right)$$
$$\mathbb{P}\left\{Z \le \mathbb{E}[Z] - t\right\} \le \exp\left(-\frac{t^2}{2\mathbb{E}[Z]}\right).$$

The concentration of measure perspective deals with deviations around expectation/median by focusing on the structure of increments of the functional under consideration.

 $X_1, \ldots, X_n, X_i \in [-1, 1]$ , independent,  $\mathbb{E}[X_i] = 0$ ,  $Var[X_i] = v$ .

#### Markov inequality.

Z an  $\mathbb{R}$ -valued random variable.

f a positive measurable function, such that f is non-decreasing on  $[t,\infty)$ .

$$\mathbb{P}\{Z > t\} = \mathbb{P}\{f(Z) \ge f(t)\} \quad f \text{ is non-decreasing.}$$
$$= \mathbb{E}\left[\mathbb{1}_{f(Z) \ge f(t)}\right]$$
$$\leq \mathbb{E}\left[\mathbb{1}_{f(Z) \ge f(t)} \frac{f(Z)}{f(t)}\right]$$
$$\leq \mathbb{E}\left[\frac{f(Z)}{f(t)}\right]$$

 $X_1, \ldots, X_n, X_i \in [-1, 1]$ , independent,  $\mathbb{E}[X_i] = 0$ ,  $Var[X_i] = v$ .

#### Markov inequality.

Z an  $\mathbb{R}$ -valued random variable.

f a positive measurable function, such that f is non-decreasing on  $[t,\infty)$ .

$$\mathbb{P}\{Z > t\} \le \mathbb{E}\Big[\frac{f(Z)}{f(t)}\Big]$$

Examples:

$$f(x) = x^{2}, \qquad \mathbb{P}\left\{Z - \mathbb{E}[Z] > t\right\} \leq \mathbb{E}\left[\frac{(Z - \mathbb{E}[Z])^{2}}{t^{2}}\right] = \frac{\operatorname{Var}(Z)}{t^{2}}$$

$$f(x) = \exp(\lambda x) \text{ where } \lambda > 0 \qquad \mathbb{P}\left\{Z > t\right\} \leq \exp\left(-\left(\lambda t - \log \mathbb{E}[e^{-\lambda Z}]\right)\right)$$

$$f(x) = |x|^{q} \qquad \mathbb{P}\left\{Z > t\right\} \leq \mathbb{E}\left[\frac{|Z|^{q}}{t^{q}}\right] \leq \left(\frac{||Z||_{q}}{t}\right)^{q}$$

IN Markov inequality relates tail-behavior and integrability.

 $X_1, \ldots, X_n, X_i \in [-1, 1]$ , independent,  $\mathbb{E}[X_i] = 0$ ,  $Var[X_i] = v$ .

Exponential Markov inequalities for sums of independent random variables.  $Z = \sum_{i=1}^{n} X_i$ 

$$\mathbb{P}\left\{Z > t\right\} \leq e^{-\lambda t} \mathbb{E}\left[e^{\lambda \sum_{i} X_{i}}\right]$$
$$\leq e^{-\lambda t} \mathbb{E}\left[\prod_{i} e^{\lambda X_{i}}\right]$$
$$\leq e^{-\lambda t} \prod_{i} \mathbb{E}\left[e^{\lambda X_{i}}\right] \quad \text{by independence}$$

 $X_1, \ldots, X_n, X_i \in [-1, 1]$ , independent,  $\mathbb{E}[X_i] = 0$ ,  $Var[X_i] = v$ .

Exponential Markov inequalities for sums of independent random variables.  $Z = \sum_{i=1}^{n} X_i$ 

$$\mathbb{P}\Big\{Z > t\Big\} \le e^{-\lambda t} \prod_{i} \mathbb{E}\Big[e^{\lambda X_{i}}\Big]$$

$$\begin{split} \prod_{i} \left( \mathbb{E} \Big[ e^{\lambda X_{i}} \Big] \right) &= \prod_{i} \left( \mathbb{E} \Big[ \sum_{k=0}^{\infty} \frac{\lambda^{k} X_{i}^{k}}{k!} \Big] \right) \\ &\leq \prod_{i} \left( \mathbb{E} \Big[ 1 + \sum_{k=2}^{\infty} \frac{\lambda^{k} |X_{i}|^{k}}{k!} \Big] \right) \text{ centering} \\ &\leq \prod_{i} \left( 1 + v \sum_{k=2}^{\infty} \frac{\lambda^{k}}{k!} \right) \text{ variance, boundedness} \\ &= \prod_{i} \left( 1 + v(e^{\lambda} - \lambda - 1) \right) \\ &\leq \exp \Big( \operatorname{Var}(Z)(e^{\lambda} - \lambda - 1) \Big) \end{split}$$

 $X_1, \ldots, X_n, X_i \in [-1, 1]$ , independent,  $\mathbb{E}[X_i] = 0$ ,  $Var[X_i] = v$ .

Exponential Markov inequalities for sums of independent random variables.  $Z = \sum_{i=1}^{n} X_i$ With  $\tau^*(\lambda) = e^{\lambda} - \lambda - 1$ ,

$$\mathbb{P}\left\{Z > t\right\} \le e^{-\lambda t} \exp\left(\operatorname{Var}(Z)\tau^*(\lambda)\right)$$

Convex duality:  $\tau(x) = \sup_{\lambda} [\lambda x - \tau^*(\lambda)] = (x+1)\log(x+1) - x$ 

$$\mathbb{P}\left\{Z > t\right\} \leq \exp\left(-\sup_{\lambda} \left[\lambda t - \operatorname{Var}(Z)\tau^{*}(\lambda)\right]\right)$$
$$\leq \exp\left(-\operatorname{Var}(Z)\tau\left(\frac{t}{\operatorname{Var}(Z)}\right)\right) \text{ Bennett}$$
$$\leq \exp\left(-\frac{t^{2}}{2\left(\operatorname{Var}(Z) + t/3\right)}\right) \text{ Bernstein}$$

Doob's embedding.

$$M_i = \mathbb{E}\Big[f(X_1, X_2, \dots, X_n) \mid X_1^i\Big] = \mathbb{E}[Z \mid X_1^i]$$

 $M_i$  is  $\sigma(X_1, \ldots, X_i)$ -measurable.

$$M_0 = \mathbb{E}[Z] \qquad M_n = Z$$

$$\mathbb{E}\Big[M_{i+1} \mid X_1^i\Big] = M_i$$

The sequence  $(M_i)_{i=0}^n$  is an  $(\sigma(X_1^n))_{i=0,...,n}$ -adapted martingale.

INF Doob's embedding allows to represent a general function of many random variables as the last value of a martingale. The increments  $M_{i+1} - M_i$  reflect the sensitivity of f with respect to its arguments.

Martingale with bounded increments.

Assumption: for all i,  $|M_{i+1} - M_i| \le 1$  and  $\mathbb{E}[Z] = 0$ 

Increasing process associated with martingale  $(M_i)$ .

$$\langle M \rangle_i = \sum_{j=1}^i \mathbb{E}[(M_j - M_{j-1})^2 \mid X_1^{j-1}]$$

 $\operatorname{Var}(M_n) = \operatorname{Var}(Z) = \mathbb{E}[\langle M \rangle_n]$ The process  $\left( \exp\left(\lambda M_i - \tau^*(\lambda) \langle M \rangle_i \right) \right)_i$  is an  $(\sigma(X_1^i))_i$ -adapted super-martingale:

$$\mathbb{E}\left[\left(\exp\left(\lambda M_{i}-\tau^{*}(\lambda)\langle M\rangle_{i}\right)\right)\mid X_{1}^{i-1}\right]\leq \exp\left(\lambda M_{i-1}-\tau^{*}(\lambda)\langle M\rangle_{i-1}\right)$$

#### Bounded-differences inequality (McDiarmid).

Assumption:  $\langle M \rangle_n \leq c$  with probability 1.

$$\mathbb{E}\left[e^{\lambda(Z-\mathbb{E}[Z])}\right] \le \exp(c\tau^*(\lambda))$$

As increments are bounded by 1, we have  $\langle M \rangle_n \leq n$ .

Azuma inequality (Hoeffding inequality for martingales with bounded increments)

$$\mathbb{P}\left\{Z - \mathbb{E}[Z] \ge t\right\} \le \exp\left(-\frac{t^2}{2n}\right)$$

○ A worst-case upper bound on  $\langle M \rangle_n$  may be excessively conservative.
 It is highly desirable to take advantage of the integrability of  $\langle M \rangle_n$  in order to derive tight tail bounds for martingale with bounded increments.

#### Bernstein inequality for martingales. Stopping

T is a stopping time with respect to  $(\sigma(X_1^i))_i$  if  $T \leq i$  is  $\sigma(X_1^i)$ -measurable.

The optional sampling theorem entails that

$$\left(\exp\left(\lambda M_0 - \tau^*(\lambda)\langle Z\rangle_0\right), \exp\left(\lambda M_T - \tau^*(\lambda)\langle M\rangle_T\right)\right)$$

is a super-martingale.

For a (finite) stopping time  $T \quad \mathbb{E}\left[\exp\left(\lambda Z_T - \tau^*(\lambda)\langle Z \rangle_T\right)\right] \leq 1$ Let *V* denote a fixed positive quantity

 $T = \min \left\{ n, \min \left\{ j : \langle M \rangle_{j+1} > V \right\} \right\}$  is a stopping time.

#### Bernstein inequality for martingales. Stopping

T is a stopping time with respect to  $(\sigma(X_1^i))_i$  if  $T \leq i$  is  $\sigma(X_1^i)$ -measurable.

$$\mathbb{P}\left\{Z \ge t \wedge T = n\right\} \le \mathbb{P}\left\{M_T \ge t\right\}$$
  
$$\le \mathbb{P}\left\{\lambda M_T - \tau^*(\lambda) \langle M \rangle_T \ge \lambda t - \tau^*(\lambda) \langle M \rangle_T\right\}$$
  
$$\le \exp\left(-\lambda t + \tau^*(\lambda)V\right)$$
  
$$\le \exp\left(-V\tau\left(\frac{t}{V}\right)\right) \text{ Optimizing with respect to } \lambda$$
  
$$\le \exp\left(-\frac{t^2}{2(V+t/3)}\right)$$

#### Bernstein inequality for martingales. Stopping

T is a stopping time with respect to  $(\sigma(X_1^i))_i$  if  $T \leq i$  is  $\sigma(X_1^i)$ -measurable.

$$\mathbb{P}\left\{Z \ge t \wedge T = n\right\} \le \exp\left(-\frac{t^2}{2(V+t/3)}\right)$$

$$\mathbb{P}\left\{Z \ge t\right\} \le \mathbb{P}\left\{\langle M \rangle_n \ge V\right\} + \exp\left(-\frac{t^2}{2(V+t/3)}\right)$$

The stopping-time trick allows to search for trade-offs and just requires the martingale to be smooth on average. It is tight for moderate deviations.

 $\odot$  The conditional variance process may be as difficult to analyze as Z...

🖙 Efron-Stein estimates will serve as surrogates for conditional variance process.

 $I_{1}, \ldots, X_{n}$ : independent random variables.

- $I_{1}, \ldots, X_{n}$ : independent random variables.
- $I = f(X_1, \ldots, X_n)$ : function of *many* independent random variables.

- $I_{1}, \ldots, X_{n}$ : independent random variables.
- $J = f(X_1, \ldots, X_n)$ : function of *many* independent random variables.
- $X'_1, \ldots, X'_n: independent copies of X_1, \ldots, X_n.$

- $I_{1}, \ldots, X_{n}$ : independent random variables.
- $J = f(X_1, \ldots, X_n)$ : function of *many* independent random variables.
- $X'_1, \ldots, X'_n: independent copies of X_1, \ldots, X_n.$ 
  - Two perturbations of Z :

 $X_i$  is replaced by an independent copy  $X'_i$ :

 $Z^{(i)} = f(X_1, \dots, X_{i-1}, X'_i, X_{i+1}, \dots, X_n)$ 

 $Z_i$  does not depend on  $X_i$ .

 $Z_i = f_i(X_1, \ldots, X_{i-1}, X_{i+1}, \ldots, X_n).$ 

- $I_{1}, \ldots, X_{n}$ : independent random variables.
- $= f(X_1, \ldots, X_n)$ : function of *many* independent random variables.
- $X'_1, \ldots, X'_n : independent copies of <math>X_1, \ldots, X_n .$ 
  - Two perturbations of Z :

 $X_i$  is replaced by an independent copy  $X'_i$ :

 $Z^{(i)} = f(X_1, \dots, X_{i-1}, X'_i, X_{i+1}, \dots, X_n)$ 

 $Z_i$  does not depend on  $X_i$ .

 $Z_i = f_i(X_1, \ldots, X_{i-1}, X_{i+1}, \ldots, X_n).$ 

 $V^+$  and V allow to monitor the sensitivity of Z

$$V^{+} = \sum_{i} \mathbb{E} \Big[ (Z - Z^{(i)})^{2} \mathbb{1}_{Z > Z^{(i)}} \mid X_{1}^{n} \Big]$$
$$V = \sum_{i} (Z - Z_{i})^{2}.$$
# Conventions

- $I_{1}, \ldots, X_{n}$ : independent random variables.
- $J = f(X_1, \ldots, X_n)$ : function of *many* independent random variables.
- $X'_1, \ldots, X'_n: independent copies of X_1, \ldots, X_n.$ 
  - Two perturbations of Z :

 $X_i$  is replaced by an independent copy  $X'_i$ :

 $Z^{(i)} = f(X_1, \dots, X_{i-1}, X'_i, X_{i+1}, \dots, X_n)$ 

 $Z_i$  does not depend on  $X_i$ .

 $Z_i = f_i(X_1, \ldots, X_{i-1}, X_{i+1}, \ldots, X_n).$ 

 $V^+$  and V allow to monitor the sensitivity of Z

$$V^{+} = \sum_{i} \mathbb{E} \Big[ (Z - Z^{(i)})^{2} \mathbb{1}_{Z > Z^{(i)}} \mid X_{1}^{n} \Big]$$
$$V = \sum_{i} (Z - Z_{i})^{2}.$$

$$\operatorname{Var}[Z] \le \mathbb{E}[V^+] = \begin{cases} \frac{1}{2} \sum_i \mathbb{E}\left[ (Z - Z^{(i)})^2 \right] = \sum_i \operatorname{Var}[Z \mid X_1^{i-1}, X_{i+1}^n] \\ \sum_i \mathbb{E}\left[ (Z - Z^{(i)})^2 \mathbb{1}_{Z > Z^{(i)}} \right] \end{cases}$$

As  $\operatorname{Var}[Z \mid X_1^{i-1}, X_{i+1}^n] \leq \mathbb{E}[(Z - Z_i)^2]$ , for any  $Z_i$  that is  $X_1^{i-1}, X_{i+1}^n$ -measurable,

#### $\operatorname{Var}[Z] \leq \mathbb{E}[V^+] \leq \mathbb{E}[V].$

E-S inequality is tight: for sums of independent random variables, it becomes an inequality.

E-S inequality is (sometimes) called the Jacknife estimate of variance.

E-S inequality is also called the tensorization property of variance.

 $\odot$  E-S inequality may be poor: let  $Z = \prod_i X_i$  where  $X_i$  are Bernouilli random variables with expectation p. Var $[Z] = p^n(1-p^n)$  while the Efron-Stein estimate equals  $np^n(1-p)$  !!

 $\operatorname{Var}[Z] \leq \mathbb{E}[V^+] \leq \mathbb{E}[V].$ 

Proof of Efron-Stein inequality.(I)

w.l.o.g. assume  $\mathbb{E}[Z] = 0$ .

$$\mathbb{E}[Z^2] = \mathbb{E}\Big[\big(\sum_{i=1}^n \mathbb{E}[Z \mid X_1^i] - \mathbb{E}[Z \mid X_1^{i-1}]\big)^2\Big] \text{ martingale decomposition of } Z$$
$$= \sum_{i=1}^n \mathbb{E}\Big[\big(\mathbb{E}[Z \mid X_1^i] - \mathbb{E}[Z \mid X_1^{i-1}]\big)^2\Big] \text{ orthogonality martingale increments}$$
$$= \sum_{i=1}^n \mathbb{E}_{X_1^i}\Big[\Big(\mathbb{E}_{X_{i+1}^n}[Z \mid X_1^i] - \mathbb{E}_{X_{i+1}^n}[\mathbb{E}_{X_i}[Z \mid X_1^{i-1}]]\Big)^2\Big]$$

 $\operatorname{Var}[Z] \leq \mathbb{E}[V^+] \leq \mathbb{E}[V].$ 

Proof of Efron-Stein inequality.(I)

w.l.o.g. assume  $\mathbb{E}[Z] = 0$ .

$$\mathbb{E}[Z^2] = \mathbb{E}\left[\left(\sum_{i=1}^n \mathbb{E}[Z \mid X_1^i] - \mathbb{E}[Z \mid X_1^{i-1}]\right)^2\right] \text{ martingale decomposition of } Z$$
$$= \sum_{i=1}^n \mathbb{E}\left[\left(\mathbb{E}[Z \mid X_1^i] - \mathbb{E}[Z \mid X_1^{i-1}]\right)^2\right] \text{ orthogonality martingale increments}$$

(conditional) Jensen inequality,  $x^2$  convex

$$\leq \sum_{i=1}^{n} \mathbb{E}_{X_{1}^{i-1}, X_{i+1}^{n}} \left[ \mathbb{E}_{X_{i}} \left[ (Z - \mathbb{E}_{X_{i}}[Z \mid X_{1}^{i-1}, X_{i+1}^{n}])^{2} \mid X_{1}^{i-1}, X_{i+1}^{n} \right] \right]$$
$$= \sum_{i=1}^{n} \operatorname{Var}[Z \mid X_{1}^{i-1}, X_{i+1}^{n}] = \mathbb{E}[V^{+}]$$

#### Where does convexity show up ?

#### Jensen inequality in large spaces.

 $\mathcal X$  : a locally convex Haussdorf topological vector space, with dual  $\mathcal X^*$ 

Rockafellar duality Lemma. f: a lower-semi-continuous convex function from  $\mathcal{X}$  on  $(-\infty,\infty]$ 

if g is defined on  $\mathcal{X}^*$  as  $g(y) = \sup_{x \in \mathcal{X}} y(x) - f(x)$ , then  $f(x) = \sup_{y \in \mathcal{X}^*} y(x) - g(y)$ 

 $\operatorname{Var}[Z] \leq \mathbb{E}[V^+] \leq \mathbb{E}[V].$ 

#### Proof of Efron-Stein inequality.(II)

Variance maps  $\mathbb{L}_2$  toward  $\mathbb{R}^+$ . Variance is convex !

duality 
$$\operatorname{Var}[Z] = \sup_{T \in \mathbb{L}_2^+} \left[ 2\mathbb{E}[TZ] - \operatorname{Var}[T] \right]$$

$$\operatorname{Var}[Z \mid X_{2}] = \mathbb{E}_{X_{2}} \left[ \sup_{T \in \mathbb{L}_{2}(X_{1})} \left\{ 2\mathbb{E}_{X_{1}}[T(X_{1})Z] - \operatorname{Var}[T] \right\} \right]$$

$$\geq \sup_{T \in \mathbb{L}_{2}(X_{1})} \mathbb{E}_{X_{2}} \left[ 2\mathbb{E}_{X_{1}}[T(X_{1})Z] - \operatorname{Var}[T] \right]$$

$$= \sup_{T \in \mathbb{L}_{2}(X_{1})} \mathbb{E}_{X_{1}} \left[ 2\mathbb{E}_{X_{2}}[Z]T_{X_{1}} - \operatorname{Var}[T] \right]$$

$$= \operatorname{Var}[\mathbb{E}[Z \mid X_{1}]]$$

 $\operatorname{Var}[Z] \leq \mathbb{E}[V^+] \leq \mathbb{E}[V].$ 

Proof of Efron-Stein inequality.(II)

duality 
$$\operatorname{Var}[Z] = \sup_{T \in \mathbb{L}_2^+} 2\mathbb{E}[TZ] - \operatorname{Var}[T]$$
  
 $\rightarrow \operatorname{Var}[Z \mid X_1, \dots, X_{i-1}] \leq \operatorname{Var}[Z \mid X_1, \dots, X_{i-1}, X_{i+1}, \dots, X_n]$ 

$$\begin{aligned} \mathsf{Var}[Z] &= \mathbb{E}\Big[\big(Z - \mathbb{E}[Z \mid X_1]\big)^2\Big] + \mathbb{E}\Big[\big(\mathbb{E}_{X_2}[Z \mid X_1] - \mathbb{E}[Z]\big)^2\Big] \\ &= \mathbb{E}_{X_1}\Big[\mathbb{E}_{X_2}\big[\big(Z - \mathbb{E}[Z \mid X_1]\big)^2\big]\Big] + \mathbb{E}_{X_1}\Big[\big(\mathbb{E}_{X_2}[Z \mid X_1] - \mathbb{E}[Z]\big)^2\Big] \\ &\leq \mathsf{Var}[Z \mid X_1] + \mathsf{Var}[\mathbb{E}[Z \mid X_1]] \\ &\leq \mathsf{Var}[Z \mid X_1] + \mathsf{Var}[Z \mid X_2] \end{aligned}$$

# **Application of Efron-Stein inequality**

#### Back to Random interval packings

 $Z_i$  is the size of the largest packing that can be constructed when the *i*th interval is removed.

 $Z_i = Z$  if the *i*th interval does not belong to a witness of the value of Z, otherwise  $Z_i \ge Z - 1$ .

$$0 \leq Z - Z_i \leq 1$$
 and  $\sum_i (Z - Z_i) \leq Z$ 

 $V \leq Z \rightarrow \operatorname{Var}[Z] \leq \mathbb{E}[Z]$ 

If  $\mathbb{E}[Z] \nearrow \infty$  with n, E-S and Chebyshev imply a law of large numbers for  $Z/\mathbb{E}[Z]$ .

The Efron-Stein estimate is within  $\pi/(4-\pi) \approx 3.66$  from truth (asymptotically)!!

 $\mathbb{R}$   $V \leq Z$  also holds for all configuration functions, this includes VC-dimension, fat-shattering dimension, VC-entropy, conditional Rademacher averages.

# Lecture II: The entropy method

- From exponential inequalities to the concentration of measure phenomenon
- Concentration inequalities using the entropy method
  - Entropy
  - Tensorization
  - Modified Logarithmic Sobolev inequalities
  - Exponential Efron-Stein inequalities
  - Convex distance inequality
- Learning-theoretical applications
- Moment inequalities using the generalized entropy method

### Lecture II: where are we?



# Lecture II: where are we ?

Concentration inequalities are aimed to extend the classical exponential inequalities for sums of independent random variables, to functions of independent random variables. Such functions show up in learning theory as the sumpremum of the deviations between the true risk and the empirical risk, the empirical VC-dimension, the empirical VC entropy, the eigenvalues of the Gram matrix.... They play an important role in the estimation of the generalization error and have proved useful in the design of model selection strategies. Even when a function cannot be coded as a sum of independent random variables, it may be usually represented as a sum of martingale increments. If the increments happen to be bounded, martingale extensions of the classical inequalities allow (in principle) to derive exponential inequalities for the function under concern and the moments of the increasing process associated to the martingale representation.

Unfortunately the increasing process is not so easy to analyze. Efron-Stein estimates of the variance of the function of interest have proved to be reasonable surrogates for the increasing process. Indeed, in the second lecture, we will relate the exponential moments of the function under concern and the exponential moments of the Efron-Stein estimates. This will be carried out using the entropy method.

Efron-Stein estimates have often turned out to be surprisingly easy to analyze. This is exemplified by a derivation of Talagrand's convex distance inequality.

# Lecture II: roadmap

The entropy of a random variable is defined and elementary properties are presented. Then the Gaussian logarithmic-Sobolev inequality is described. It is shown to imply a tight upper-bound on the log-Laplace transform of smooth functions of Gaussian random variables. A remarkable feature of the Gaussian logarithmic-Sobolev inequality lies in the fact that it is dimension-free.

The derivation of modified logarithmic-Sobolev inequalities in product spaces is carried in two steps: tensorization of entropy is derived by resorting to convexity arguments that proved useful when deriving the Efron-Stein inequality; then an *energy* that will prove to be adequate for general product spaces is defined. The Efron-Stein upper-bound on variance shows up in the right-hand-side of modified logarithmic-Sobolev inequalities. Under various integrability conditions of the Efron-Stein upper-bound of variance, Bernstein-like inequalities for *functions of many independent random variables that depend not too much on each of them* are derived.

Bernstein-like inequalities for self-bounded functionals are derived. A derivation of Talagrand's convex distance inequality is finally presented.

To get exponential tail bounds for Z  $\rightarrow$  control  $\mathbb{E}[\exp(\lambda Z)]$  or  $\log \mathbb{E}[\exp(\lambda(Z - \mathbb{E}[Z]))]$ .

Control over log-Laplace transform of smooth functions of Gaussian random variables can be obtained from functional inequalities known as logarithmic Sobolev inequalities.

Entropy of a (positive) function f > 0 with respect to a probability distribution

 $\operatorname{Ent}(f) \stackrel{\Delta}{=} \mathbb{E}[f \log f] - \mathbb{E}[f] \log \mathbb{E}[f] = \mathbb{E}[\Phi(f)] - \Phi(\mathbb{E}[f])$  with  $\Phi(x) = x \log x$ .

Jensen inequality  $\rightarrow$  Ent $(f) \ge 0$ .

To get exponential tail bounds for Z  $\rightarrow$  control  $\mathbb{E}[\exp(\lambda Z)]$  or  $\log \mathbb{E}[\exp(\lambda(Z - \mathbb{E}[Z]))]$ .

Control over log-Laplace transform of smooth functions of Gaussian random variables can be obtained from functional inequalities known as logarithmic Sobolev inequalities.

Entropy of a (positive) function f > 0 with respect to a probability distribution

$$\operatorname{Ent}(f) \stackrel{\Delta}{=} \mathbb{E}[f \log f] - \mathbb{E}[f] \log \mathbb{E}[f] = \mathbb{E}[\Phi(f)] - \Phi(\mathbb{E}[f])$$
 with  $\Phi(x) = x \log x$ .

Jensen inequality  $\rightarrow$  Ent $(f) \ge 0$ .

$$\lim_{a \nearrow \infty} \operatorname{Ent}[(f+a)^2] = 4\operatorname{Var}[f] .$$

To get exponential tail bounds for Z  $\rightarrow$  control  $\mathbb{E}[\exp(\lambda Z)]$  or  $\log \mathbb{E}[\exp(\lambda(Z - \mathbb{E}[Z]))]$ .

Control over log-Laplace transform of smooth functions of Gaussian random variables can be obtained from functional inequalities known as logarithmic Sobolev inequalities.

Entropy of a (positive) function f > 0 with respect to a probability distribution

$$\mathsf{Ent}(f) \stackrel{\Delta}{=} \mathbb{E}[f \log f] - \mathbb{E}[f] \log \mathbb{E}[f] = \mathbb{E}[\Phi(f)] - \Phi(\mathbb{E}[f])$$
 with  $\Phi(x) = x \log x$ .

Jensen inequality  $\rightarrow$  Ent $(f) \ge 0$ .

Gaussian Logarithmic-Sobolev Inequality (Gross, 1975)



To get exponential tail bounds for Z  $\rightarrow$  control  $\mathbb{E}[\exp(\lambda Z)]$  or  $\log \mathbb{E}[\exp(\lambda(Z - \mathbb{E}[Z]))]$ .

Control over log-Laplace transform of smooth functions of Gaussian random variables can be obtained from functional inequalities known as logarithmic Sobolev inequalities.

Entropy of a (positive) function f > 0 with respect to a probability distribution

$$\mathsf{Ent}(f) \stackrel{\Delta}{=} \mathbb{E}[f \log f] - \mathbb{E}[f] \log \mathbb{E}[f] = \mathbb{E}[\Phi(f)] - \Phi(\mathbb{E}[f])$$
 with  $\Phi(x) = x \log x$ .

Jensen inequality  $\rightarrow$  Ent $(f) \ge 0$ .

Gaussian Logarithmic-Sobolev Inequality (Gross, 1975)

$$\mathsf{Ent}(f^2) \le 2\mathbb{E}\big[ \| \bigtriangledown f \|^2 \big] \xrightarrow{\bullet} \mathsf{Ent}\big(e^{\lambda f}\big) \le \frac{\lambda^2}{2}\mathbb{E}\big[ \| \bigtriangledown f \|^2 e^{\lambda f} \big].$$

This step critically depends on the chain rule for derivation.

Herbst argument: derive a differential inequality satisfied by the log-Laplace transform.

Assumption:  $\| \bigtriangledown f \|_{\infty} \leq 1$ . Notation :  $H(\lambda) = \frac{1}{\lambda} \log \mathbb{E} \left[ e^{\lambda Z} \right]$ 

 $\frac{1}{\lambda^2} \operatorname{Ent}(e^{\lambda Z}) = H'(\lambda) \mathbb{E}[e^{\lambda Z}] \quad \text{and} \quad \lim_{\lambda \to 0^+} H(\lambda) = \mathbb{E}[Z]$ 

$$\begin{array}{ll} H'(\lambda) & \leq & \displaystyle\frac{1}{2} & \mbox{Gaussian logarithmic-Sobolev inequality} \\ H(\lambda) - H(0^+) & \leq & \displaystyle\frac{\lambda}{2} & \mbox{Integration} \\ \log \mathbb{E}\left[e^{\lambda(Z - \mathbb{E}[Z])}\right] & \leq & \displaystyle\frac{\lambda^2}{2} & \mbox{Rewriting.} \end{array}$$

$$\mathbb{P}\left\{f(X_1,\ldots,X_n) \ge \mathbb{E}[f] + t\right\} \le \exp\left(-\frac{t^2}{2\|\bigtriangledown f\|_{\infty}^2}\right)$$

Herbst argument: derive a differential inequality satisfied by the log-Laplace transform.

Assumption:  $\| \bigtriangledown f \|_{\infty} \leq 1$ .  $\frac{1}{\lambda^2} \operatorname{Ent}(e^{\lambda Z}) = H'(\lambda) \mathbb{E}[e^{\lambda Z}] \text{ and } \lim_{\lambda \to 0^+} H(\lambda) = \mathbb{E}[Z]$  $\log \mathbb{E}\left[e^{\lambda(Z - \mathbb{E}[Z])}\right] \leq \frac{\lambda^2}{2}$ 

$$\mathbb{P}\left\{f(X_1,\ldots,X_n) \ge \mathbb{E}[f] + t\right\} \le \exp\left(-\frac{t^2}{2\|\bigtriangledown f\|_{\infty}^2}\right)$$

Rev Any smooth function of a Gaussian random variable enjoys Gaussian tail-behavior. The dimension of the Gaussian vector does not appear in tail-behavior.

The Gaussian concentration inequality is dimension-free.

# Variation on Herbst's argument

The Gaussian logarithmic Sobolev inequality provides with concentration inequalities when  $|| \nabla F ||^2$  is (somewhat) exponentially integrable.

Decoupling inequality
$$\mathbb{E}\left[\lambda W e^{\lambda Z}\right] \leq \operatorname{Ent}\left[e^{\lambda Z}\right] + \mathbb{E}\left[e^{\lambda Z}\right] \log \mathbb{E}\left[e^{\lambda W}\right]$$
Variational characterization of Entropy:
$$\mathbb{E}\left[\frac{f}{\mathbb{E}[f]}g\right] \leq \operatorname{Ent}\left[\frac{f}{\mathbb{E}[f]}\right] + \log \mathbb{E}\left[e^{g}\right]$$

Applying the decoupling inequality to the Gaussian logarithmic Sobolev inequality:

$$\operatorname{Ent}\left[e^{\lambda f}\right] \leq \frac{\theta \lambda}{2} \left(\operatorname{Ent}\left[e^{\lambda f}\right] + \mathbb{E}\left[e^{\lambda f}\right] \log \mathbb{E}\left[e^{\frac{\lambda || \bigtriangledown f ||^{2}}{\theta}}\right]\right)$$

$$\downarrow$$

$$\frac{1}{\lambda^{2}} \frac{\operatorname{Ent}\left[e^{\lambda f}\right]}{\mathbb{E}\left[e^{\lambda f}\right]} \leq \frac{\theta}{2\lambda(1-\lambda\theta/2)} \log \mathbb{E}\left[e^{\frac{\lambda || \bigtriangledown f ||^{2}}{\theta}}\right]$$

#### Variation on Herbst's argument

Notations:

$$H(\lambda) = \frac{1}{\lambda} \log \mathbb{E}\left[e^{\lambda f}\right] \quad \text{and} \quad G(\lambda) = \log \mathbb{E}\left[e^{\lambda \|\nabla f\|^{2}}\right],$$
$$\frac{1}{\lambda^{2}} \frac{\mathsf{Ent}\left[e^{\lambda f}\right]}{\mathbb{E}\left[e^{\lambda f}\right]} \leq \frac{\theta}{2\lambda(1-\lambda\theta/2)} \log \mathbb{E}\left[e^{\frac{\lambda \|\nabla f\|^{2}}{\theta}}\right]$$

translates into:

$$H'(\lambda) \le \frac{\theta}{2\lambda(1-\lambda\theta/2)}G(\lambda/\theta)$$

Convexity of  $G \rightarrow \text{RHS}$  is non-decreasing with respect to  $\lambda$ 

$$H(\lambda) - H(0) \le \int_{s=0}^{\lambda} \frac{\theta}{2s(1-s\theta/2)} G(s/\theta) ds$$
$$\le \frac{\theta}{2(1-\lambda\theta/2)} G(\lambda/\theta)$$

# Variation on Herbst's argument

$$\log \mathbb{E}\left[e^{\lambda(f(X) - \mathbb{E}[f])}\right] \leq \frac{\lambda\theta}{2(1 - \lambda\theta/2)} \log \mathbb{E}\left[e^{\lambda \|\nabla f\|^2/\theta}\right]$$

If  $\| \nabla f \|_{\infty} \leq L$ , taking  $\theta$  toward 0, we recover sub-Gaussian behavior.

$$\log \mathbb{E}\left[e^{\lambda(f(X) - \mathbb{E}[f])}\right] \le \frac{\lambda^2 L^2}{2}$$

If  $\| \nabla f \| = C \|X\|$  (sub-quadratic function),  $\log \mathbb{E}\left[e^{\lambda \|\nabla f\|^2/\theta}\right] < \infty$  for a non-trivial range of values of  $\lambda$ .

In the sequel, we will try to reproduce this line of reasoning for product distributions.

# **Toward modified Logarithmic Sobolev inequalities**

The Gaussian logarithmic Sobolev inequality is (almost) a characterization of the Gaussian distribution.

In order to derive tail bounds for other kinds of distributions using analogues of Herbst argument, we need to develop:

• modified logarithmic Sobolev inequalities tailored to the distributions we have in mind (exponential, Poisson, product distributions,...)

• analogues of Herbst argument.

In Learning Theory we are interested in (modified) logarithmic Sobolev inequalities for product measures

# **Toward modified Logarithmic Sobolev inequalities**

The Gaussian logarithmic Sobolev inequality is (almost) a characterization of the Gaussian distribution.

In order to derive tail bounds for other kinds of distributions using analogues of Herbst argument, we need to develop:

• modified logarithmic Sobolev inequalities tailored to the distributions we have in mind (exponential, Poisson, product distributions,...)

• analogues of Herbst argument.

In Learning Theory we are interested in (modified) logarithmic Sobolev inequalities for product measures

General form:  $\operatorname{Ent}(f^2) \leq \mathcal{E}[f]$ ,

where  $\mathcal{E}[f]$  is an energy-like functional.

# **Toward modified Logarithmic Sobolev inequalities**

Solution Strategy The Gaussian logarithmic Sobolev inequality is (almost) a characterization of the Gaussian distribution.

In order to derive tail bounds for other kinds of distributions using analogues of Herbst argument, we need to develop:

• modified logarithmic Sobolev inequalities tailored to the distributions we have in mind (exponential, Poisson, product distributions,...)

• analogues of Herbst argument.

In Learning Theory we are interested in (modified) logarithmic Sobolev inequalities for product measures

General form:  $\operatorname{Ent}(f^2) \leq \mathcal{E}[f]$ ,

where  $\mathcal{E}[f]$  is an energy-like functional.

Balanced Bernouilli Log-Sobolev Inequality (Gross, 1975)

$$\operatorname{Ent}(f^2) \leq \frac{1}{2} \mathbb{E}[(Df)^2]$$
 with  $Df \stackrel{\Delta}{=} f(1) - f(-1)$ .

# **Tensorization of entropy**

Derivation of modified LS inequalities proceeds in two steps:

Tensorization. (analogue of Efron-Stein inequality)

Notation 
$$\operatorname{Ent}[Z \mid X_1, \dots, X_{i-1}, X_{i+1}, \dots, X_n]$$
  
=  $\mathbb{E}_{X_1, \dots, X_{i-1}, X_{i+1}, \dots, X_n} \left[ \operatorname{Ent}[Z(x_1, \dots, x_{i-1}, X_i, x_{i+1}, \dots, x_n)] \right]$ 

$$\operatorname{Ent}[Z] \le \sum_{i=1}^{n} \operatorname{Ent}[Z \mid X_{1}^{i-1}, X_{i+1}^{n}]$$

Upper-bounding of entropies with respect to single random variables.

 $\log x \leq x-1 \quad > \quad -\mathbb{E}[Y] \log \mathbb{E}[Y] \leq -\mathbb{E}[Y] \log u - \mathbb{E}[Y] + u \quad \forall u > 0.$ 

# **Tensorization of entropy**

The tensorization step is generic.

Its relies on a convexity property of the functional  $Z \mapsto \text{Ent}[Z]$  on  $\mathbb{L}^1_+$ . Not to be confused with the convexity of  $x \mapsto x \log x$ .

A representation formula for Entropy:

$$\mathsf{Ent}[Z] = \sup_{T \in \mathbb{L}_1^+} \mathbb{E}\Big[ Z \log \frac{T}{\mathbb{E}[T]} \Big]$$

 $\operatorname{Ent}[Z \mid X_2] \ge \operatorname{Ent}[\mathbb{E}[Z \mid X_1]] \quad 2 \text{ variables}$  $\operatorname{Ent}[Z \mid X_1, \dots, X_{i-1}, X_{i+1}, \dots, X_n] \ge \operatorname{Ent}[\mathbb{E}[Z \mid X_1, \dots, X_{i-1}] \mid X_{i+1}, \dots, X_n]$ 

The conditional distribution of Z with respect  $X_1, \ldots, X_{i-1}$  is a convex combination of conditional distributions with respect to  $X_1, \ldots, X_{i-1}, X_{i+1}, \ldots, X_n$ 

#### **Tensorization of entropy**

$$\begin{aligned} &\operatorname{Ent}[Z(X_{1}, X_{2})] \\ &= \sup_{T} \mathbb{E} \left[ Z \log \frac{T(X_{1}, X_{2})}{\mathbb{E}[T]} \right] \\ &= \sup_{T} \mathbb{E} \left[ Z \log \frac{T}{\mathbb{E}_{X_{2}}[T]} + Z \log \frac{\mathbb{E}_{X_{2}}[T]}{\mathbb{E}[T]} \right] \\ &\leq \sup_{T} \mathbb{E} \left[ Z \log \frac{T}{\mathbb{E}_{X_{2}}[T]} \right] + \sup_{T} \mathbb{E}_{X_{2}} \left[ \mathbb{E}_{X_{1}} \left[ Z \log \frac{\mathbb{E}_{X_{2}}[T]}{\mathbb{E}[T]} \right] \right] \\ &\leq \mathbb{E}_{X_{1}} \left[ \sup_{T \in \mathbb{L}_{1}^{+}(X_{2})} \mathbb{E}_{X_{2}} \left[ Z \log \frac{T}{\mathbb{E}_{X_{2}}[T]} \right] \right] + \mathbb{E}_{X_{2}} \sup_{T \in \mathbb{L}_{1}^{+}(X_{1})} \left[ \mathbb{E}_{X_{1}} \left[ Z \log \frac{T}{\mathbb{E}_{X_{1}}[T]} \right] \right] \\ &\leq \operatorname{Ent}[Z \mid X_{1}] + \operatorname{Ent}[Z \mid X_{2}]. \end{aligned}$$

The general formula follows by induction on the number of variables.

$$\underbrace{X_1, \dots, X_{n-1}}_{Y_1} \quad X_n$$

Solution When dealing with the variance, the proof of the Efron-Stein inequality reduced to the proof of the tensorization property of the variance. Here, the tensorization property of entropy leaves us with a sum of one-dimensional conditional entropies.

More work is needed in order to get an energy-like term.

$$\int \log x \le x - 1 \quad \to \quad -\mathbb{E}[Y] \log \mathbb{E}[Y] \le -\mathbb{E}[Y] \log u - \mathbb{E}[Y] + u \quad \forall u > 0.$$
$$\operatorname{Ent}[f(X)] \quad \le \quad \mathbb{E}\left[f(X) \log \frac{f(X)}{u} - (f(X) - u)\right] \quad \forall u > 0.$$

$$\log x \le x - 1 \quad \to \quad -\mathbb{E}[Y] \log \mathbb{E}[Y] \le -\mathbb{E}[Y] \log u - \mathbb{E}[Y] + u \quad \forall u > 0.$$
$$\mathsf{Ent}[f(X)] \quad \le \quad \mathbb{E}\left[f(X) \log \frac{f(X)}{u} - (f(X) - u)\right] \quad \forall u > 0.$$

Let  $g_i$  denote a function of the n-1 random variables:  $X_1^{i-1}, X_{i+1}^n$ .

$$\mathsf{Ent}[f \mid X_1^{i-1}, X_{i+1}^n] \leq \mathbb{E}_{X_1^{i-1}, X_{i+1}^n} \left[ \mathbb{E}_{X_i} \left[ f \log \frac{f}{g_i} - \left( f - g_i \right) \mid X_1^{i-1}, X_{i+1}^n \right] \right]$$

$$\log x \le x - 1 \quad \to \quad -\mathbb{E}[Y] \log \mathbb{E}[Y] \le -\mathbb{E}[Y] \log u - \mathbb{E}[Y] + u \quad \forall u > 0.$$
$$\mathsf{Ent}[f(X)] \quad \le \quad \mathbb{E}\Big[f(X) \log \frac{f(X)}{u} - (f(X) - u)\Big] \quad \forall u > 0.$$

Let g denote a function of n - 1 random variables. Combining with the tensorization property of Entropy:

$$\mathsf{Ent}[f] \leq \sum_{i} \mathbb{E}_{X_{1}^{i-1}, X_{i+1}^{n}} \left[ \mathbb{E}_{X_{i}} \left[ f \log \frac{f}{g} - (f - g) \right] = X_{1}^{i-1}, X_{i+1}^{n} \right]$$

re entails the Gaussian logarithmic Sobolev inequality, the Poissonian logarithmic Sobolev inequality...

$$\log x \le x - 1 \quad \to \quad -\mathbb{E}[Y] \log \mathbb{E}[Y] \le -\mathbb{E}[Y] \log u - \mathbb{E}[Y] + u \quad \forall u > 0.$$
$$\mathsf{Ent}[f(X)] \quad \le \quad \mathbb{E}\Big[f(X) \log \frac{f(X)}{u} - (f(X) - u)\Big] \quad \forall u > 0.$$

Specializing to functions  $f(X_1^n) = e^{\lambda Z}$ , letting  $g(X_1^{i-1}, X_{i+1}^n) = e^{\lambda Z_i}$ ,

$$\operatorname{Ent}\left[e^{\lambda Z}\right] \leq \sum_{i} \mathbb{E}_{X_{1}^{i-1}, X_{i+1}^{n}} \left[ \mathbb{E}_{X_{i}}\left[e^{\lambda Z}\left(\lambda Z - \lambda Z_{i}\right) - \left(e^{\lambda Z} - e^{\lambda Z_{i}}\right)\right] \right]$$
$$\leq \sum_{i} \mathbb{E}_{X_{1}^{i-1}, X_{i+1}^{n}} \left[ \mathbb{E}_{X_{i}}\left[e^{\lambda Z}\left(e^{\lambda(Z_{i}-Z)} - \lambda(Z_{i}-Z) - 1\right)\right] \right]$$
$$\leq \sum_{i} \mathbb{E}_{X_{1}^{i-1}, X_{i+1}^{n}} \left[ \mathbb{E}_{X_{i}}\left[e^{\lambda Z}\tau^{*}(\lambda(Z_{i}-Z))\right] \right]$$

where  $Z_i$  is a measurable function of  $X_1, \ldots, X_{i-1}, X_{i+1}, \ldots, X_n$ , and  $\tau^*(x) \stackrel{\Delta}{=} e^x - x - 1$ .

where  $Z_i$  is a measurable function of  $X_1, \ldots, X_{i-1}, X_{i+1}, \ldots, X_n$ , and  $\tau^*(x) \stackrel{\Delta}{=} e^x - x - 1$ .

 $If Z \ge Z_i and \lambda \ge 0 as e^x - x - 1 \le x^2/2 for x \le 0$ 

$$\operatorname{Ent}\left[e^{\lambda Z}\right] \leq \frac{\lambda^2}{2} \mathbb{E}\left[e^{\lambda Z} \sum_{i} (Z - Z_i)^2\right]$$

$$\begin{aligned} & \blacktriangleright \text{Let } V = \sum_{i} (Z - Z_{i})^{2}, \\ & \quad \text{Ent} \Big[ e^{\lambda Z} \Big] \leq \frac{\lambda^{2}}{2} \mathbb{E} \Big[ V e^{\lambda Z} \Big] \ \rightarrow \ \log \mathbb{E} \Big[ e^{\lambda (Z - \mathbb{E}[Z])} \Big] \leq \frac{\lambda \theta}{2(1 - \lambda \theta/2)} \log \mathbb{E} \Big[ e^{\lambda V/\theta} \Big] \end{aligned}$$

INFIGURE Not only the second moment of Z is related to some moment of V, (Efron-Stein) but also the log-Laplace transforms are connected.

$$\operatorname{Ent}\left[e^{\lambda Z}\right] \leq \operatorname{\mathbb{E}}\left[e^{\lambda Z}\sum_{i}\tau^{*}(\lambda(Z_{i}-Z))\right]$$

$$\operatorname{Self-bounded functions}$$

$$\operatorname{If} Z_{i} \leq Z \leq Z_{i}+1, \text{ as } \tau^{*} \text{ is convex:}$$

$$\tau^{*}(-\lambda(Z-Z_{i})) \leq (1+Z-Z_{i})\tau(0)+\tau^{*}(-\lambda)(Z-Z_{i})$$

$$= \tau^{*}(-\lambda)(Z-Z_{i})$$

$$\operatorname{Ent}\left[e^{\lambda Z}\right] \leq \frac{\lambda^{2}}{2}\operatorname{\mathbb{E}}\left[e^{\lambda Z}\tau^{*}(-\lambda)(\sum_{i}Z-Z_{i})\right]$$

• If furthermore  $\sum_{i} (Z - Z_i) \leq Z$   $\operatorname{Ent}\left[e^{\lambda Z}\right] \leq \tau^*(-\lambda) \mathbb{E}\left[Ze^{\lambda Z}\right]$ 

Stating and solving the differential inequality:

 $(\lambda$ 

$$\begin{split} \mathbb{E}\Big[\lambda Z e^{\lambda Z}\Big] &- \mathbb{E}\Big[e^{\lambda Z}\Big]\log\mathbb{E}\Big[e^{\lambda Z}\Big] \leq \tau^*(-\lambda)\mathbb{E}\Big[Z e^{\lambda Z}\Big]\\ & \text{multiplying by } \exp(-\lambda\mathbb{E}[Z])\\ & \text{regrouping}\\ & \text{dividing by } F(\lambda) = \mathbb{E}\Big[\exp\big(\lambda(Z - \mathbb{E}[Z])\big)\Big]\\ &- \tau^*(\lambda)\big)\frac{F'(\lambda)}{F(\lambda)} - \log F(\lambda) \leq \tau^*(\lambda)\mathbb{E}[Z] \end{split}$$

Note that  $\mathbb{E}[Z]\tau^*(\lambda)$  is a solution of the differential equation

$$(\lambda - \tau^*(\lambda))f'(\lambda) - f(\lambda) = \tau^*(\lambda)\mathbb{E}[Z]$$

One can check that  $\log F(\lambda) = \mathbb{E}[Z]\tau^*(\lambda)$  is the largest solution of the differential inequality, by showing that for all solutions  $h(\lambda) = (e^{\lambda} - 1)g(\lambda)$  of

$$(\lambda - \tau^*(\lambda))h'(\lambda) - h(\lambda) \le 0$$

are non-positive.

#### Self-bounded functions

• If  $Z_i \leq Z \leq Z_i + 1$ , as  $\tau^*$  is convex:

$$\tau^*(-\lambda(Z-Z_i)) \leq (1+Z-Z_i)\tau(0) + \tau^*(-\lambda)(Z-Z_i)$$
$$= \tau^*(-\lambda)(Z-Z_i)$$

$$\operatorname{Ent}\left[e^{\lambda Z}\right] \leq \frac{\lambda^2}{2} \mathbb{E}\left[e^{\lambda Z} \tau^*(-\lambda) \left(\sum_i Z - Z_i\right)\right]$$

• If furthermore  $\sum_{i} (Z - Z_i) \leq Z$   $\operatorname{Ent}\left[e^{\lambda Z}\right] \leq \tau^*(-\lambda) \mathbb{E}\left[Ze^{\lambda Z}\right]$ 

 $\log \mathbb{E}\left[e^{\lambda(Z-\mathbb{E}[Z])}\right] \le \mathbb{E}[Z]\tau^*(\lambda)$ 

$$\mathbb{P}\left\{Z \ge \mathbb{E}[Z] + t\right\} \le \exp\left(-\frac{t^2}{2(\mathbb{E}[Z] + t/3)}\right).$$
Rather than  $Z_i$ , we may consider  $Z^{(i)}$ , where  $Z^{(i)} = f(X_1, \ldots, X_{i-1}, X'_i, X_{i+1}, \ldots, X_n)$ Z = Z and  $Z^{(i)}$  are identically distributed.

$$\begin{aligned} \operatorname{Ent} \begin{bmatrix} e^{\lambda Z} \end{bmatrix} &\leq \sum_{i=1}^{n} \mathbb{E} \left[ e^{\lambda Z} \tau^* (-\lambda(Z - Z^{(i)})) \right] \\ &\leq \sum_{i=1}^{n} \mathbb{E} \left[ e^{\lambda Z} \tau^* (-\lambda(Z - Z^{(i)})) \mathbb{1}_{Z > Z^{(i)}} \right] \\ &+ \sum_{i=1}^{n} \mathbb{E} \left[ e^{\lambda Z} \tau^* (-\lambda(Z - Z^{(i)})) \mathbb{1}_{Z > Z^{(i)}} \right] \\ &\leq \sum_{i=1}^{n} \mathbb{E} \left[ e^{\lambda Z} \tau^* (-\lambda(Z - Z^{(i)})) \mathbb{1}_{Z > Z^{(i)}} \right] \\ &+ \sum_{i=1}^{n} \mathbb{E} \left[ e^{\lambda Z} (\tau^* (-\lambda(Z^{(i)} - Z)) \mathbb{1}_{Z > Z^{(i)}} \right] \\ &\leq \sum_{i=1}^{n} \mathbb{E} \left[ e^{\lambda Z} \left( e^{-\lambda(Z - Z^{(i)})} - 1 \right) \lambda(Z^{(i)} - Z) \mathbb{1}_{Z > Z^{(i)}} \right] \end{aligned}$$

Symmetrized version of modified logarithmic Sobolev inequalities.

$$\operatorname{Ent}\left[e^{\lambda Z}\right] \leq \sum_{i=1}^{n} \mathbb{E}\left[e^{\lambda Z}\psi(-\lambda(Z-Z^{(i)}))\mathbb{1}_{Z>Z^{(i)}}\right]$$

where  $Z^{(i)}$  is a measurable function of  $X_1, \ldots, X_{i-1}, X_{i+1}, \ldots, X_n$ , and  $X'_i$  while

$$\psi(x) = e^x \tau^*(-x) + \tau^*(x) = x(e^x - 1).$$

 $igsquir \psi()$  is convex and for x>0,  $\psi(-x)\leq x^2,$ 

for 
$$\lambda > 0$$
  $\operatorname{Ent}\left[e^{\lambda Z}\right] \leq \lambda^2 \mathbb{E}\left[e^{\lambda Z}V^+\right]$ 

Where  $\mathbb{E}[V^+]$  is the Efron-Stein upper-bound on Var[Z].

For  $x \leq 0$ ,  $\psi(x) \leq x^2$ . The modified logarithmic Sobolev inequality implies that for  $\lambda \geq 0$ :

$$\operatorname{Ent}\left[e^{\lambda Z}\right] \leq \lambda^{2} \mathbb{E}\left[e^{\lambda Z} V_{+}\right]$$
$$\leq \theta \lambda \left(\operatorname{Ent}\left[e^{\lambda Z}\right] + \mathbb{E}\left[e^{\lambda Z}\right] \log \mathbb{E}\left[e^{\lambda V_{+}/\theta}\right]\right).$$

which translates into:

$$\frac{1}{\lambda^2} \frac{\mathsf{Ent}\left[e^{\lambda Z}\right]}{\mathbb{E}\left[e^{\lambda Z}\right]} \leq \frac{\theta}{\lambda(1-\theta\lambda)} \log \mathbb{E}\left[e^{\lambda V_+/\theta}\right]$$

This has exactly the same form as the inequality encountered while carrying the extended Herbst argument !

An exponential Efron-Stein inequality.  $\theta > 0$  and  $\lambda \in (0, 1/\theta)$ ,

$$\log \mathbb{E}\left[\exp(\lambda(Z - \mathbb{E}[Z]))\right] \le \frac{\lambda\theta}{1 - \lambda\theta} \log \mathbb{E}\left[\exp\left(\frac{\lambda V_+}{\theta}\right)\right].$$

Solution The exponential integrability of Z is actually related to the exponential integrability of the Efron-Stein estimate(s) of variance. If we can bound  $\log \mathbb{E}\left[\exp\left(\frac{\lambda V_{+}}{\theta}\right)\right]$ , then we can use an exponential Markov inequality in order to get a Bernstein-like inequality for Z.

- Modus operandi: check whether
- V or V is constant.
- V or  $V^+$  is upper-bounded by aZ + b,
- V or  $V^+$  is simpler than Z.

Done or to be done for:

- $\bullet \sup_{f} R[f] R_{\mathsf{emp}}(f).$
- Empirical VC or fat-shattering dimension.
- Empirical VC entropy.
- Number of support vectors.



 $I \subseteq \mathcal{X}^n$ 



Control by convex distance provides sharp deviation inequalities around the median, for

- 1. configuration functions,
- 2. euclidean traveling salesman,
- 3. bin-packing,
- 4. minimum spanning trees.

Control by convex distance can be recovered from exponential Efron-Stein inequalities.

 $\int d_T(\cdot, \cdot)$  can be represented as a saddle point.  $\mathcal{M}(A)$  : set of probabilities on A.

$$d_T(X_1^n, A) = \inf_{\nu \in \mathcal{M}(A)} \sup_{\alpha: \|\alpha\|_2 \le 1} \sum_j \alpha_j \mathbb{E}_{\nu}[\mathbb{1}_{X_j \ne Y_j}]$$
$$= \sup_{\alpha: \|\alpha\|_2 \le 1} \inf_{\nu \in \mathcal{M}(A)} \sum_j \alpha_j \mathbb{E}_{\nu}[\mathbb{1}_{X_j \ne Y_j}]$$

 $d_T(\cdot, \cdot)$  can be represented as a saddle point.  $\mathcal{M}(A)$  : set of probabilities on A.

$$d_T(X_1^n, A) = \inf_{\nu \in \mathcal{M}(A)} \sup_{\alpha: \|\alpha\|_2 \le 1} \sum_j \alpha_j \mathbb{E}_{\nu}[\mathbbm{1}_{X_j \ne Y_j}]$$
$$= \sup_{\alpha: \|\alpha\|_2 \le 1} \inf_{\nu \in \mathcal{M}(A)} \sum_j \alpha_j \mathbb{E}_{\nu}[\mathbbm{1}_{X_j \ne Y_j}]$$

Sion Minimax Theorem  $f: \mathcal{X} \times \mathcal{Y} \to \mathbb{R}$ 

convex and lower-semi-continuous with respect to xconcave and upper-semi-continuous with respect to y $\mathcal{X}$  convex and compact

 $\inf_{x} \sup_{y} f(x,y) = \sup_{y} \inf_{x} f(x,y) = \min_{x} \sup_{y} f(x,y) \,.$ 

 $(\widehat{\boldsymbol{\nu}}, \widehat{\boldsymbol{\alpha}})$  : a saddle point for  $X_1^n$ .

$$Z^{(i)} = \inf_{\nu \in \mathcal{M}(A)} \sup_{\alpha} \sum_{j} \alpha_{j} \mathbb{E}_{\nu} [\mathbb{1}_{X_{j}^{(i)} \neq Y_{j}}] \geq \inf_{\nu \in \mathcal{M}(A)} \sum_{j} \widehat{\alpha}_{j} \mathbb{E}_{\nu} [\mathbb{1}_{X_{j}^{(i)} \neq Y_{j}}].$$

 $\tilde{\nu}$ : distribution on A that achieves the infimum.

$$Z = \inf_{\nu} \sum_{j} \widehat{\alpha}_{j} \mathbb{E}_{\nu} [\mathbb{1}_{X_{j} \neq Y_{j}}] \leq \sum_{j} \widehat{\alpha}_{j} \mathbb{E}_{\widetilde{\nu}} [\mathbb{1}_{X_{j} \neq Y_{j}}].$$

$$Z - Z^{(i)} \leq \sum_{j} \widehat{\alpha}_{j} \mathbb{E}_{\widetilde{\nu}} [\mathbb{1}_{X_{j} \neq Y_{j}} - \mathbb{1}_{X_{j}^{(i)} \neq Y_{j}}] = \widehat{\alpha}_{i} \mathbb{E}_{\widetilde{\nu}} [\mathbb{1}_{X_{i} \neq Y_{i}} - \mathbb{1}_{X_{i}^{(i)} \neq Y_{i}}] \leq \widehat{\alpha}_{i} .$$

$$V_{+} \leq \sum_{i} \widehat{\alpha}_{i}^{2} = 1.$$

 $\operatorname{Var}[d_T(X_1^n, A)] \leq 1$  Efron-Stein inequality !

 $\mathbb{P}\left[d_T(X_1^n, A) - \mathbb{E}d_T(X_1^n, A) \ge t\right] \le e^{-t^2/4} \quad \text{exponential Efron-Stein inequality !.}$ 

 $\operatorname{Var}[d_T(X_1^n, A)] \leq 1$  Efron-Stein inequality !

 $\mathbb{P}\left[d_T(X_1^n, A) - \mathbb{E}d_T(X_1^n, A) \ge t\right] \le e^{-t^2/4} \quad \text{exponential Efron-Stein inequality !.}$ 

$$\mathbb{P}[d_T(X_1^n, A) - \mathbb{E}d_T(X_1^n, A) \le -t] \le \frac{\mathsf{Var}[d_T(X_1^n, A)]}{t^2} \le \frac{1}{t^2} \quad \text{Chebyshev inequality}$$

$$\mathbb{P}[A] \le \frac{1}{\left(\mathbb{E}[d_T(\cdot, A)]\right)^2}$$

 $\operatorname{Var}[d_T(X_1^n, A)] \leq 1$  Efron-Stein inequality !

 $\mathbb{P}\left[d_T(X_1^n, A) - \mathbb{E}d_T(X_1^n, A) \ge t\right] \le e^{-t^2/4} \quad \text{exponential Efron-Stein inequality !.}$ 

 $\mathbb{P}[d_T(X_1^n, A) - \mathbb{E}d_T(X_1^n, A) \le -t] \le \frac{\mathsf{Var}[d_T(X_1^n, A)]}{t^2} \le \frac{1}{t^2} \quad \text{Chebyshev inequality}$ 

$$\mathbb{E}d_T(X_1^n, A) \le \frac{1}{\sqrt{\mathbb{P}[A]}}$$

 $\mathbb{P}[d_T(X_1^n, A) \ge t + \sqrt{2}] \le e^{-t^2/4} \text{ if } \mathbb{P}[A] \ge 1/2$ 

 $\operatorname{Var}[d_T(X_1^n, A)] \leq 1$  Efron-Stein inequality !

 $\mathbb{P}\left[d_T(X_1^n, A) - \mathbb{E}d_T(X_1^n, A) \ge t\right] \le e^{-t^2/4} \quad \text{exponential Efron-Stein inequality !.}$ 

 $\mathbb{P}[d_T(X_1^n, A) - \mathbb{E}d_T(X_1^n, A) \le -t] \le \frac{\mathsf{Var}[d_T(X_1^n, A)]}{t^2} \le \frac{1}{t^2} \quad \text{Chebyshev inequality}$ 

$$\mathbb{E}d_T(X_1^n, A) \le \frac{1}{\sqrt{\mathbb{P}[A]}}$$

 $t > \sqrt{2} {>} (t - \sqrt{2})^2 \ge t^2/2 - 4\log 2$ 

 $\mathbb{P}[d_T(X_1^n, A) \ge t] \le 2e^{-t^2/8} \quad \text{for } t > \sqrt{2} \text{ and } \mathbb{P}[A] \ge 1/2$ 

# **Plan III : Statistical learning applications**

- From exponential inequalities to the concentration of measure phenomenon
- Concentration inequalities using the entropy method
- Learning-theoretical applications
  - Fat-shattering VC-dimension
  - VC entropy
  - Conditional Rademacher averages
  - Supremum of empirical processes
  - **9** ...
- Moment inequalities using the generalized entropy method

## Where are we?

We hope (and have some evidence) that we have developped valuable tools for deriving tail bounds for general functions of independent variables.

$$Z = f(X_1, ..., X_n)$$
  

$$Z^{(i)} = f(X_1, ..., X'_i, ..., X_n)$$
  

$$Z_i = f_i(X_1, ..., X_{i-1}, X_{i+1}, ..., X_n)$$
  

$$V^+ = \mathbb{E}' \Big[ \sum_i (Z - Z^{(i)})^2 \mathbb{1}_{Z > Z^{(i)}} \mid X_1^n \Big]$$
  

$$V = \sum_i (Z - Z_i)^2$$

$$\operatorname{Var}[Z] \leq \mathbb{E}[V^+] \leq \mathbb{E}[V]$$
$$\log \mathbb{E}\left[e^{\lambda(Z - \mathbb{E}[Z])}\right] \leq \frac{\theta\lambda}{1 - \theta\lambda} \log \mathbb{E}\left[e^{\lambda V^+/\theta}\right]$$

We need to check that for some learning problems, either V or  $V^+$  is manageable.

# Lecture III: Roadmap

Concentration inequalities have proved helpful in statistical learning theory because they are key ingredients in the derivation of risk bounds for empirical risk minimizers in classification and bounded regression.

Concentration inequalities for self-bounded functionals allow to prove that many of the quantities that have been considered in order to quantify the complexity of a class of functions, like the empirical VC-dimension, the empirical VC-entropy and conditional Rademacher averages are sharply concentrated.

This paves the way to data-dependent estimation of the complexity of function classes, which is of great importance in model selection.

The most important consequence of concentration inequalities concerns suprema of empirical processes. The most refined versions of Talagrand's concentration inequality for suprema of empirical processes may be considered as process versions of Bernstein inequality. They provide new insights on Vapnik-Chervonenkis inequalities.

## **Fat-shattering dimension**

Solution Fat-shattering dimension  $fat(X_1^n, \gamma, \mathcal{F})$ . Given  $\gamma > 0$ , and  $\mathcal{F}$ , and sample  $X_1, \ldots, X_n$ ,  $fat(X_1^n, \gamma, \mathcal{F})$  is the largest d such that there exists  $\{i_1, \ldots, i_d\} \subseteq \{1, \ldots, n\}$  with

$$\forall (\epsilon_{i_j})_{j \le d} \in \{-1, 1\}^d, \quad \exists f \in \mathcal{F}, \quad \epsilon_{i_j} f(X_{i_j}) \ge \gamma \,.$$

 $Z = \mathsf{fat}(X_1^n, \gamma, \mathcal{F})$ 

Fat-shattering dimension captures the complexity of a function class at certain scale on a certain sample.

If a set of linear classifiers  $\mathcal{F}$  separates a sample  $X_1, \ldots, X_n$  with margin  $\gamma > 0$ :

 $\forall (\epsilon_i)_{i \leq n} \in \{-1, 1\}^n, \quad \exists f \in \mathcal{F}, \quad \epsilon_i f(X_i) \geq \gamma \quad \text{then fat}(X_1^n, \gamma, \mathcal{F}) = n., ,$ 

Depending on the choice of kernels, on underlying distribution,  $fat(X_1^n, \gamma, \mathcal{F})$  may scale very differently with *n*. Is the fat-shattering dimension relatively stable?

$$\frac{\operatorname{fat}(X_1^n, \gamma, \mathcal{F})}{\mathbb{E}\left[\operatorname{fat}(X_1^n, \gamma, \mathcal{F})\right]} \to 1 ?$$

### **Fat-shattering dimension**

Set Fat-shattering dimension  $fat(X_1^n, \gamma, \mathcal{F})$ . Given  $\gamma > 0$ , and  $\mathcal{F}$ , and sample  $X_1, \ldots, X_n$ ,  $fat(X_1^n, \gamma, \mathcal{F})$  is the largest d such that there exists  $\{i_1, \ldots, i_d\} \subseteq \{1, \ldots, n\}$  with

$$\begin{aligned} \forall (\epsilon_{i_j})_{j \leq d} \in \{-1, 1\}^d, \quad \exists f \in \mathcal{F}, \quad \epsilon_{i_j} f(X_{i_j}) \geq \gamma \,. \end{aligned}$$
$$\begin{aligned} Z = \mathrm{fat}(X_1^n, \gamma, \mathcal{F}) \end{aligned}$$

Z is a configuration function: a  $\gamma$ -shattered sub-sample witnesses the value of Z

With 
$$Z_i = fat(X_1^{i-1}, X_{i+1}^n, \gamma, \mathcal{F}), \quad 0 \le Z - Z_i \le 1$$
 &  $\sum_i Z - Z_i \le Z.$ 

Z is sub-Poissonian !!!  $\log \mathbb{E} \left[ e^{\lambda(Z - \mathbb{E}[Z])} \right] \leq \mathbb{E}[Z] \tau^*(\lambda)$ .

$$\mathbb{P}\left[Z \ge \mathbb{E}Z + t\right] \le \exp\left[-\frac{t^2}{2\mathbb{E}Z + 2t/3}\right]$$

For every  $0 < t \le \mathbb{E}Z$ ,  $\mathbb{P}[Z \le \mathbb{E}Z - t] \le \exp\left[-\frac{t^2}{2\mathbb{E}Z}\right]$ .

 $\mathcal{F}$ : a class of  $\{-1, 1\}$ -valued functions.  $X_1^n$  a sample.

Trace of 
$$\mathcal{F}$$
 on  $X_1^n$ :  $\left\{ (b_i)_{i \le n}; (b_i) \in \{-1, 1\}^n, \exists f \in \mathcal{F}, b_i = f(X_i) \right\}$ 

The VC-entropy of  $\mathcal{F}$  in  $X_1^n$  is defined as:

$$Z = \log_2 \left| \left\{ (b_i)_{i \le n}; (b_i) \in \{-1, 1\}^n, \ \exists f \in \mathcal{F}, b_i = f(X_i) \right\} \right|$$

 $Z_i$  is defined as

$$Z_{i} = \log_{2} \left| \left\{ (b_{j})_{j \le n, j \ne i}; (b_{j}) \in \{-1, 1\}^{n-1}, \ \exists f \in \mathcal{F}, b_{j} = f(X_{j}) \ \forall j \right\} \right|$$

Obvious:  $0 \leq Z - Z_i \leq 1$ ,

P a probability on a finite set.

Shannon entropy of 
$$P$$
  $H(P) = \sum_{i} -P(i) \log_2 P(i)$ 

Shannon entropy is positive and maximal when P is uniform.

 $\begin{array}{ll} \mbox{Notation} & H(X) = H(\mbox{dist}(X)) \\ \mbox{Conditional entropy} & H(X \mid Y) = \mathbb{E} \Big[ H(\mbox{dist}(X \mid Y)) \Big] \\ \mbox{Chain rule.} & H(X,Y) = H(Y) + H(X \mid Y) \\ \mbox{Conditionning} & H(X \mid Y) \leq H(X) & ...\mbox{decreases Shannon entropy.} \end{array}$ 

Han inequality.

$$\sum_{i} H(X_{1}^{n}) - H(X_{1}^{i-1}, X_{i+1}^{n}) = \sum_{i=1}^{n} H(X_{i} \mid X_{1}^{i-1}, X_{i+1}^{n})$$

$$\leq \sum_{i=1}^{n} H(X_{i} \mid X_{1}^{i-1})$$

$$\leq H(X_{1}^{n}).$$

The uniform distribution on the trace of  $\mathcal{F}$  on  $X_1^n$  defines a random element  $Y_1^n$  of  $\{-1,1\}^n$ 

$$Z = H(Y_1^n)$$

 $Z_i$  is the Shannon-entropy of the uniform distribution on the trace of  $\mathcal{F}$  on  $X_1^{i-1}, X_{i+1}^n$ .

$$Z_i \ge H(Y_1^{i-1}, Y_{i+1}^n)$$
.

realize Han inequality entails  $\sum_i Z - Z_i \leq Z$ .

The VC-entropy is a self bounded functional. It enjoys a sub-Poissonian behavior.

#### $V \leq Z \rightarrow \operatorname{Var}[Z] \leq \mathbb{E}[Z]$

INFIGURE VC-entropy may be almost-surely constant (half-spaces when samples are in general position with probability one.)

INF VC-entropy may be approximately Gaussian: samples of size n from the uniform distribution on a shattered set of size  $\alpha n$  where  $\alpha$  is fixed and n tends to infinity...

If we just know that the Z is a VC-entropy, the inequality is tight.

 $\checkmark$   $\mathcal{F}$  a VC-class of classifiers.

$$\mathbb{E}\Big[\sup_{f\in\mathcal{F}}\Big|\sum_{i}\mathbbm{1}_{f(X_{i})\neq Y_{i}}-n\mathbb{E}\big[\mathbbm{1}_{f(X)\neq Y}\big]\Big|\Big]$$

$$=\mathbb{E}\Big[\sup_{f\in\mathcal{F}}\Big|\sum_{i}\mathbbm{1}_{f(X_{i})\neq Y_{i}}-\sum_{i}\mathbb{E}'\big[\mathbbm{1}_{f(X'_{i})\neq Y'_{i}}\big]\Big|\Big]$$

$$\leq\mathbb{E}\mathbb{E}'\Big[\sup_{f\in\mathcal{F}}\Big|\sum_{i}\mathbbm{1}_{f(X_{i})\neq Y_{i}}-\mathbbm{1}_{f(X'_{i})\neq Y'_{i}}\Big|\Big]$$

$$=\mathbb{E}\mathbb{E}'\mathbb{E}_{\epsilon}\Big[\sup_{f\in\mathcal{F}}\Big|\sum_{i}\epsilon_{i}\big(\mathbbm{1}_{f(X_{i})\neq Y_{i}}-\mathbbm{1}_{f(X'_{i})\neq Y'_{i}}\big)\Big|\Big]$$

$$\leq\mathbb{E}\mathbb{E}_{\epsilon}\Big[\sup_{f\in\mathcal{F}}\Big|\sum_{i}\epsilon_{i}\mathbbm{1}_{f(X_{i})\neq Y_{i}}\Big|+\mathbb{E}'\mathbb{E}_{\epsilon}\Big[\sup_{f\in\mathcal{F}}\Big|\sum_{i}\epsilon_{i}\mathbbm{1}_{f(X'_{i})\neq Y'_{i}}\Big|\Big]$$

$$=\mathbb{E}\Big[\mathbb{E}_{\epsilon}\Big[\sup_{f\in\mathcal{F}}\Big|\sum_{i}\epsilon_{i}\mathbbm{1}_{f(X_{i})\neq Y_{i}}\Big|+\mathbb{E}'\mathbb{E}_{\epsilon}\Big[\sup_{f\in\mathcal{F}}\Big|\sum_{i}\epsilon_{i}\mathbbm{1}_{f(X'_{i})\neq Y'_{i}}\Big|\Big]$$

The conditional Rademacher average only depends on the trace of  $\mathcal{F}$  on  $X_1^n$  ....

 $\mathbb{E}_{\epsilon} \left[ \sup_{f \in \mathcal{F}} \left| \sum_{i} \epsilon_{i} \mathbb{1}_{f(X_{i}) \neq Y_{i}} \right| \mid X_{1}^{n} \right] \leq \sqrt{2 \log |\operatorname{trace}(\mathcal{F}, X_{1}^{n})|} \times \sup_{f \in \mathcal{F}} \sqrt{\sum_{i} \mathbb{1}_{f(X_{i}) \neq Y_{i}}} \right]$  $\textbf{So For a fixed } f \in \mathcal{F} \mathbb{E}_{\epsilon} \left[ e^{\lambda \left| \sum_{i} \epsilon_{i} \mathbb{1}_{f(X_{i}) \neq Y_{i}} \right|} \right] \leq 2 \exp \left( \frac{\lambda^{2}}{2} \sum_{i} \mathbb{1}_{f(X_{i}) \neq Y_{i}} \right)$  Hoeffding! Revisiting the union bound:  $\exp\left(\lambda \mathbb{E}\left[\max_{f \in \text{trace}(\mathcal{F}, X_1^n)} \left|\sum_{i} \epsilon_i \mathbb{1}_{f(X_i) \neq Y_i}\right|\right]\right)$  $\leq \mathbb{E}\Big[\exp\Big(\lambda \max_{f \in \mathsf{trace}(\mathcal{F}, X_1^n)} \big| \sum_{i} \epsilon_i \mathbb{1}_{f(X_i) \neq Y_i} \big| \Big) \Big]$  $\leq \mathbb{E}\Big[\sum_{f \in \mathsf{trace}(\mathcal{F}, X_{i}^{n})} \exp\left(\lambda \Big|\sum_{i} \epsilon_{i} \mathbb{1}_{f(X_{i}) \neq Y_{i}}\Big|\Big)\Big]$  $\leq \sum_{f \in \operatorname{trace}(\mathcal{F}, X_1^n)} \mathbb{E} \Big[ \exp \left( \lambda \Big| \sum_i \epsilon_i \mathbb{1}_{f(X_i) \neq Y_i} \Big| \right) \Big]$  $\leq 2|\operatorname{trace}(\mathcal{F}, X_1^n)| \max_{f \in \operatorname{trace}(\mathcal{F}, X_1^n)} \exp\left(\frac{\lambda^2}{2} \sum_{i} \mathbb{1}_{f(X_i) \neq Y_i}\right)$ 

The bound follows by optimizing with respect to  $\lambda$ .

$$\begin{split} \mathbb{E}[Z] &= \mathbb{E}\Big[\sup_{f\in\mathcal{F}}\Big|\sum_{i}\mathbbm{1}_{f(X_{i})\neq Y_{i}} - n\mathbb{E}\big[\mathbbm{1}_{f(X)\neq Y}\big]\Big|\Big] \\ &\leq \mathbb{E}\Big[\sqrt{2\log|\operatorname{trace}(\mathcal{F},X_{1}^{n})|} \times \sup_{f\in\mathcal{F}}\sqrt{\sum_{i}\mathbbm{1}_{f(X_{i})\neq Y_{i}}}\Big] \\ &\leq \sqrt{\mathbb{E}\Big[2\log|\operatorname{trace}(\mathcal{F},X_{1}^{n})|\Big]} \times \sqrt{\mathbb{E}\Big[\sup_{f\in\mathcal{F}}\sum_{i}\mathbbm{1}_{f(X_{i})\neq Y_{i}}\Big]} \\ &\leq \sqrt{\mathbb{E}\Big[2\log|\operatorname{trace}(\mathcal{F},X_{1}^{n})|\Big]} \times \sqrt{\sup_{f\in\mathcal{F}}n\mathbb{E}\big[\mathbbm{1}_{f(X)\neq Y}\big] + \mathbb{E}[Z]} \\ &\mathbb{E}[Z] \leq \mathbb{E}\Big[2\log|\operatorname{trace}(\mathcal{F},X_{1}^{n})|\Big] + \sup_{f\in\mathcal{F}}n\mathbb{E}\big[\mathbbm{1}_{f(X)\neq Y}\big] \end{split}$$

Solution May be highly relevant if  $\mathcal{F}$  is constituted by classifiers with law error rate (order 1/n).

Line of reasoning is at the origin of (some) localization procedures.

## **Rademacher complexity**

The empirical fat-shattering dimension, and the empirical VC-entropy, are two capacity concepts in statistical learning theory. Their relevance stems from their relationship to the average value of  $\sup_{f \in \mathcal{F}} \left| \sum_{i} \mathbb{1}_{f(X_i) \neq Y_i} - n \mathbb{E} \left[ \mathbb{1}_{f(X) \neq Y} \right] \right|$ .

This relationship is proved using conditional Rademacher averages (symmetrization). Why shouldn't we prove that conditional Rademacher averages are concentrated around their mean value.

 ${\cal F}$  : countable class of measurable centered real-valued functions of [-1,1] -valued functions.

$$Z = \mathbb{E}\Big[\sup_{f \in \mathcal{F}} \left| \sum_{i} \epsilon_{i} f(X_{i}) \right| \mid X_{1}^{n} \Big]$$

 $\epsilon_i$ : independent centered  $\{-1, 1\}$ -valued Random variables.

$$Z_{i} = \mathbb{E} \Big[ \sup_{f \in \mathcal{F}} \left| \sum_{j \neq i} \epsilon_{j} f(X_{j}) \right| \mid X_{1}^{i-1}, X_{i+1}^{n} \Big]$$

# **Rademacher complexity**

$$Z_{i} = \mathbb{E} \Big[ \sup_{f \in \mathcal{F}} \left| \sum_{j \neq i} \epsilon_{j} f(X_{j}) + \mathbb{E}_{\epsilon_{i}} \epsilon_{i} f(X_{j}) \right| | X_{1}^{n} \Big]$$

$$\leq \mathbb{E} \Big[ \sup_{f \in \mathcal{F}} \left| \sum_{j \neq i} \epsilon_{j} f(X_{j}) + \epsilon_{i} f(X_{j}) \right| | X_{1}^{n} \Big] \text{ Jensen}$$

$$= Z$$

$$Z \leq \mathbb{E}_{\epsilon} \Big[ \sup_{f \in \mathcal{F}} \big| \sum_{j \neq i} \epsilon_j f(X_j) \big| + 1 \mid X_1^n \Big] = Z_i + 1.$$

$$\sum_{i} (Z - Z_i) \leq Z \text{ triangle inequality.}$$

# **Rademacher complexity**

Rademacher complexities are self-bounded functionals. They also enjoy sub-Poissonian behavior.

 $V \leq Z \quad \textbf{\rightarrow} \quad \mathsf{Var}[Z] \leq \mathbb{E}[Z]$ 

For VC-classes,  $\mathbb{E}[Z] \approx C\sqrt{\operatorname{vc} n}$ .

E-S inequality implies that the typical fluctuations of Rademacher complexities. are of order at most  $(dn)^{1/4}$ .

Again this may be too conservative.

#### Suprema of positive bounded empirical processes

$$Z = \sup_{f \in \mathcal{F}} \sum_{i} f(X_i)$$

where  $\mathcal{F}$  is a set of positive functions, ... Assumption  $\sup_f \sup_x f(x) \leq 1$ 

$$Z_i = \sup_{f \in \mathcal{F}} \sum_{j \neq i} f(X_j)$$

$$0 \le Z - Z_i \le 1$$
 and  $\sum_i Z - Z_i \le Z$ .

Z is sub-Poissonian !!!  $\log \mathbb{E}\left[e^{\lambda(Z-\mathbb{E}[Z])}\right] \leq \mathbb{E}[Z]\tau^*(\lambda)$ .

$$\mathbb{P}\left[Z \ge \mathbb{E}Z + t\right] \le \exp\left[-\frac{t^2}{2\mathbb{E}Z + 2t/3}\right]$$

For every  $0 < t \leq \mathbb{E}Z$ ,  $\mathbb{P}\left[Z \leq \mathbb{E}Z - t\right] \leq \exp\left[-\frac{t^2}{2\mathbb{E}Z}\right]$ .

# Suprema of positive bounded empirical processes

Gram matrix (Kernel-machines)  $G_{i,j} = \langle X_i, X_j \rangle$ 

Let  $\lambda_1 \geq \lambda_2 \geq \ldots \lambda_n \geq 0$  denote the ordered sequence of eigenvalues of G.

What can we say about  $\sum_{i=1}^{r} \lambda_i$ ?

It may be representeed as a supremum of a positive bounded empirical processes  $\Pi_V$ : projection on subspace *V*.

$$\sum_{i=1}^{r} \lambda_i = \sup_{V:\dim(V)=r} \|\Pi_V(X_i)\|^2$$

If the distribution of Xi has bounded support, this is the supremum of a bounded positive empirical process...

Suprema of empirical processes

$$\sup_{f \in \mathcal{F}} \left| \sum_{i} \left( f(X_i) - \mathbb{E}[f(X_i)] \right) \right|$$

where  ${\mathcal F}$  may be a set of classifiers, regression functions, ... Assumption  $\sup_f \sup_x |f(x)| \le 1$ 

The analysis of empirical processes are at the root of the Vapnik-Chervonenkis theory of learning. More generally, they prove to be central in the annalysis of M-estimators.

The supremum is one among many quantities associated with an empirical process. Another is the modulus of continuity.

As the supremum of an empirical process is the sup-norm of a sum of random vectors, a natural question is: are there extensions of Bernstein inequalities for suprema of empirical processes?

Using specifically modified logarithmic Sobolev inequalities, Rio and Bousquet have proved that Bernstein inequality scales up to the vector-valued setting.

Suprema of bounded empirical processes provide an exemple of a two-steps approach of concentration inequalities:

- Compute a tractable upper-bound for Efron-Stein estimates.
- Analyse the concentration properties of this Efron-Stein estimate.

 $V^+ \leq \sup_{f \in \mathcal{F}} \sum_i \mathbb{E}' \left[ (f(X_i) - f(X'_i))^2 \right] \text{ If } \mathbb{E}[f] = 0 \text{ for all } f \in \mathcal{F}$ 

$$V^{+} \leq \sup_{f \in \mathcal{F}} \sum_{i} f^{2}(X_{i}) + \sup_{f \in \mathcal{F}} \mathbb{E}\left[f^{2}(X_{i})\right]$$

Solution The stochastic part of the upper-bound on  $V^+$  is a supremum of a bounded positive process !

$$\log \mathbb{E}\Big[e^{\lambda V^+}\Big] \leq \lambda \Big(\mathbb{E}[\sup_{f \in \mathcal{F}} \sum_i f^2(X_i)] + \sup_{f \in \mathcal{F}} \mathbb{E}\Big[f^2(X_i)\Big]\Big) + \mathbb{E}[\sup_{f \in \mathcal{F}} \sum_i f^2(X_i)]\tau^*(\lambda).$$

$$\mathbb{E}\Big[\sup_{f\in\mathcal{F}}\sum_{i}f^{2}(X_{i})\Big] \leq \mathbb{E}\Big[\sup_{f\in\mathcal{F}}\sum_{i}f^{2}(X_{i}) - \mathbb{E}[f^{2}(X_{i})]\Big] + \sup_{f\in\mathcal{F}}\sum_{i}\mathbb{E}\Big[f^{2}(X_{i})\Big] \\
\leq 2\mathbb{E}\Big[\sup_{f\in\mathcal{F}}\sum_{i}\epsilon_{i}f^{2}(X_{i})\Big] + \sup_{f\in\mathcal{F}}\sum_{i}\mathbb{E}\Big[f^{2}(X_{i})\Big] \quad \text{symmetrization} \\
\leq 4\mathbb{E}\Big[\sup_{f\in\mathcal{F}}\sum_{i}\epsilon_{i}f(X_{i})\Big] + \sup_{f\in\mathcal{F}}\mathbb{E}\Big[f^{2}(X_{i})\Big] \quad \text{contraction} \\
\leq 4\mathbb{E}\Big[\sup_{f\in\mathcal{F}}\Big|\sum_{i}f(X_{i}) - \mathbb{E}[f(X_{i})]\Big|\Big] + \sup_{f\in\mathcal{F}}\sum_{i}\mathbb{E}\Big[f^{2}(X_{i})\Big]$$

$$\log \mathbb{E}\left[e^{\lambda V^{+}}\right] \leq (2\lambda + \tau^{*}(\lambda)) \sum_{i} \left(\sup_{f \in \mathcal{F}} \mathbb{E}\left[f^{2}(X_{i})\right]\right) + 4(\tau^{*}(\lambda) + \lambda) \mathbb{E}\left[\sup_{f \in \mathcal{F}} \sum_{i} |f(X_{i}) - \mathbb{E}[f(X_{i})]|\right]$$

Plugging

$$\log \mathbb{E}\left[e^{\lambda V^{+}}\right] \leq (2\lambda + \tau^{*}(\lambda)) \Big(\sup_{f \in \mathcal{F}} \sum_{i} \mathbb{E}\left[f^{2}(X_{i})\right]\Big) + 4(\tau^{*}(\lambda) + \lambda) \mathbb{E}[\sup_{f \in \mathcal{F}} \sum_{i} |f(X_{i}) - \mathbb{E}[f(X_{i})]|]$$

in exponential Efron-Stein inequality.

$$\log \mathbb{E} \left[ e^{\lambda (Z - \mathbb{E}[Z])_{+}} \right] \\ \frac{\theta \lambda}{1 - \theta \lambda} \left( \left( \frac{2\lambda}{\theta} + \tau^{*} \left( \frac{\lambda}{\theta} \right) \right) \sup_{f \in \mathcal{F}} \sum_{i} \mathbb{E} \left[ f^{2}(X_{i}) \right] + 4 \left( \frac{\lambda}{\theta} + \tau^{*} \left( \frac{\lambda}{\theta} \right) \right) \mathbb{E}[Z] \right)$$

The official concentration inequality for suprema of bounded empirical processes. Talagrand (94,96), Leddoux (97), Massart(2000), Rio(2001), Bousquet (2002)

$$v = \sup_{f \in \mathcal{F}} \sum_{i} \mathbb{E} \Big[ f^2(X_i) \Big] + 2\mathbb{E}[Z]$$

$$\mathbb{P}\Big\{Z \ge \mathbb{E}[Z] + \sqrt{2xv} + \frac{x}{3}\Big\} \le e^{-x}$$

This holds for very small classes as well as for large classes ... What distinguishes large and small classes is  $\mathbb{E}[Z]$  not the concentration phenomenon.

 $\mathbb{E}[Z]$  may be analyzed using different techniques (chaining).

Moment inequalities
# Organization

- From exponential inequalities to the concentration of measure phenomenon
- Concentration inequalities using the entropy method
- Learning-theoretical applications
- Moment inequalities using the generalized entropy method

# Plan IV : when increments are not bounded

- From exponential inequalities to the concentration of measure phenomenon
- Concentration inequalities using the entropy method
- Learning-theoretical applications
- Moment inequalities
  - Beyond exponential inequalities
  - Main moment inequalities
  - $\phi$ -Sobolev inequalities
  - **Solution** From  $\phi$ -Sobolev inequalities to moment inequalities
  - **Solution** Rosenthal inequalities from  $\phi$ -Sobolev inequalities
  - Conditional Rademacher averages for general processes

## Where are we?

The entropy method (lecture II) has allowed us to scale up Efron-Stein inequalities (lecture I) to the level of exponential moments. In lecture III, we have illustrated the power of exponential Efron-Stein inequalities by showing that some quantities that play a central role in Vapnik-Chervonenkis theory of learning are indeed concentrated around their mean value.

Such results proved vvery helpful when understanding classification problems. However, the concentration inequalities described in lecture II, took full advantage of some boundedness assumptions.

In some settings (regression, SVM-soft-margin classification), assuming that everything is bounded seems very restrictive. Hence there is a need to interpolate between Efron-Stein and exponential Efron-Stein inequalities.

# Lecture IV: Roadmap

Exponential inequalities prove powerless when dealing with sums of poorly integrable independent random variables. In order to get tail inequalities for such sums, we traditionally resort to moment inequalities known as Rosenthal-Pinelis inequalities. The aim of this lecture is to describe the generalized entropy method. This generalized entropy method provides upper-bound on the *q*th norm of functions of many independent random variables.

The approach is to relate the *q*th norm of some functional of independent random variables, with the *q*th norm of the square root of  $V^+$  (the quantity which shows up in the Efron-Stein upper-bound), and then to upper-bound  $||\sqrt{V^+}||_q$ .

The relationship between  $||(Z - \mathbb{E}[Z])_+||_q$  and  $||\sqrt{V^+}||_q$ , interpolates between Efron-Stein inequality and modified logarithmic Sobolev inequalities. It is proved by resorting to functionals called  $\phi$ -entropies that generalize both entropy and variance. Those functionals called  $\phi$ -entropies enjoy a duality property that warrants that they also enjoy the tensorization property. Using tensorization and optimization, it is thus possible to show that product probability distributions enjoy  $\phi$ -Sobolev inequalities that may be used to derive moment inequalities.

The Martingale approach is well-adapted to derive moment inequalities. Recall: ( $\mathcal{F}_i$ ) =  $\sigma(X_1^i)$ -algebra.  $\left(M_i = \mathbb{E}[Z|\mathcal{F}_i]\right)_i$  is an  $\mathcal{F}_i$ -adapted martingale.

$$\langle Z \rangle = \sum_{i=1}^{n} (M_i - M_{i-1})^2 .$$

**Burkholder** inequalities

$$\left\|Z - \mathbb{E}[Z]\right\|_{q} \le (q-1)\sqrt{\left\|\langle Z \rangle\right\|_{q/2}} = (q-1)\left\|\sqrt{\langle Z \rangle}\right\|_{q} .$$

Burkholder inequalities

$$|Z - \mathbb{E}[Z]||_q \le (q-1)\sqrt{\|\langle Z \rangle\|_{q/2}} = (q-1) \left\|\sqrt{\langle Z \rangle}\right\|_q .$$

$$V_{+} \stackrel{\Delta}{=} \sum_{i} \mathbb{E}\left[ \left( Z - Z^{(i)} \right)^{2} \mathbb{1}_{Z > Z^{(i)}} \mid X_{1}^{n} \right]$$

 $\operatorname{Var}(Z) \leq \mathbb{E}\left[V^+\right].$ 

What about other moments?

Burkholder inequalities

$$|Z - \mathbb{E}[Z]||_q \le (q-1)\sqrt{||\langle Z \rangle||_{q/2}} = (q-1) \left\|\sqrt{\langle Z \rangle}\right\|_q .$$

$$\mathsf{Var}(Z) \le \mathbb{E}\left[V^+\right]$$

What about other moments?

Exponential Efron-Stein inequality.

$$\log \mathbb{E}\left[e^{\lambda(Z-\mathbb{E}[Z])}\right] \leq \frac{\lambda\theta}{1-\lambda\theta} \log \mathbb{E}\left[e^{\frac{\lambda V^+}{\theta}}\right] \quad \text{for} \qquad \lambda \quad \in \qquad (0,1/\theta).$$

Requires exponential integrability of  $V^+$ .

Burkholder inequalities

$$|Z - \mathbb{E}[Z]||_q \le (q-1)\sqrt{||\langle Z \rangle||_{q/2}} = (q-1) \left\|\sqrt{\langle Z \rangle}\right\|_q .$$

$$\mathsf{Var}(Z) \le \mathbb{E}\left[V^+\right]$$

What about other moments?

Exponential Efron-Stein inequality.

$$\log \mathbb{E}\left[e^{\lambda(Z-\mathbb{E}[Z])}\right] \leq \frac{\lambda\theta}{1-\lambda\theta} \log \mathbb{E}\left[e^{\frac{\lambda V^+}{\theta}}\right] \quad \text{for} \qquad \lambda \quad \in \quad (0,1/\theta).$$

Requires exponential integrability of  $V^+$ .

Goal: Burkholder-type inequalities relating the qth norm of Z with the q/2th norm V or  $V^+$ 

Goal: relate the integrability of Z with the integrability of  $V^+$ . Viewing Burkholder inequalities, it is reasonable to assume that Z is as integrable as  $\sqrt{V^+}$ .

It seems highly desirable to get a sharp dependence on q.

Goal: relate the integrability of Z with the integrability of  $V^+$ . Viewing Burkholder inequalities, it is reasonable to assume that Z is as integrable as  $\sqrt{V^+}$ .

It seems highly desirable to get a sharp dependence on q.

Moment inequalities might be interesting per se:

$$\mathbb{P}\left\{Z \ge t\right\} \le \min_{q} \left(\frac{\|Z\|_{q}}{t}\right)^{q} \le \inf_{\lambda} \frac{\mathbb{E}\left[e^{\lambda Z}\right]}{e^{\lambda t}}$$

Goal: relate the integrability of Z with the integrability of  $V^+$ . Viewing Burkholder inequalities, it is reasonable to assume that Z is as integrable as  $\sqrt{V^+}$ .

It seems highly desirable to get a sharp dependence on q.

Moment inequalities might be interesting per se:

$$\mathbb{P}\left\{Z \ge t\right\} \le \min_{q} \left(\frac{\|Z\|_{q}}{t}\right)^{q} \le \inf_{\lambda} \frac{\mathbb{E}\left[e^{\lambda Z}\right]}{e^{\lambda t}}$$

The relationship between  $||Z||_q$  and q reflects the tail behavior of Z.
"Equivalence of moments" principle asserts that  $||Z||_q$  grows slower than  $q^d$  iff  $Z^{1/d}$  is exponentially integrable.

Goal: relate the integrability of Z with the integrability of  $V^+$ . Viewing Burkholder inequalities, it is reasonable to assume that Z is as integrable as  $\sqrt{V^+}$ .

It seems highly desirable to get a sharp dependence on q.

Moment inequalities might be interesting per se:

$$\mathbb{P}\left\{Z \ge t\right\} \le \min_{q} \left(\frac{\|Z\|_{q}}{t}\right)^{q} \le \inf_{\lambda} \frac{\mathbb{E}\left[e^{\lambda Z}\right]}{e^{\lambda t}}$$

- The relationship between  $||Z||_q$  and q reflects the tail behavior of Z.
  "Equivalence of moments" principle asserts that  $||Z||_q$  grows slower than  $q^d$  iff  $Z^{1/d}$  is exponentially integrable.
- Moment inequalities with a tight dependence on q provide Bernstein-like inequalities.

$$\text{if } \|(Z - \mathbb{E}\left[Z\right])_{+}\|_{q} \leq \sum_{i=1}^{d} A_{i}q^{i/2} \quad \mathbb{P}\{Z > \mathbb{E}\left[Z\right] + t\} \leq \exp\left(-\log 2\min_{i \leq d} \left(\frac{dt}{2A_{i}}\right)^{2/i}\right)$$

# Main inequalities

 $q\in\mathbb{N},q\geq2.$ 

\_

Gaussian type behavior, bounded-difference inequality.

if  $V^+ \leq c$ , then  $q \geq 2$ ,  $\|(Z - \mathbb{E}[Z])_+\|_q \leq \sqrt{qc}$ .

# Main inequalities

 $q\in\mathbb{N},q\geq2.$ 

\_

Gaussian type behavior, bounded-difference inequality.

if  $V^+ \leq c$ , then  $q \geq 2$ ,  $\|(Z - \mathbb{E}[Z])_+\|_q \leq \sqrt{qc}$ .

Burkholder-like inequality

$$\left\| (Z - \mathbb{E}[Z])_+ \right\|_q \le \sqrt{3q} \left\| \sqrt{V^+} \right\|_q,$$

# Main inequalities

 $q \in \mathbb{N}, q \geq 2.$ 

Gaussian type behavior, bounded-difference inequality.

if  $V^+ \leq c$ , then  $q \geq 2$ ,  $||(Z - \mathbb{E}[Z])_+||_q \leq \sqrt{qc}$ .

Burkholder-like inequality

$$\left\| (Z - \mathbb{E} \left[ Z \right])_+ \right\|_q \le \sqrt{3q} \left\| \sqrt{V^+} \right\|_q,$$

Burkholder-like inequality (II)
If  $Z_i \leq Z$  for all *i*,

$$\left| (Z - \mathbb{E}[Z])_+ \right\|_q \le \sqrt{3q} \left\| \sqrt{V} \right\|_q$$

Modified  $\phi$ -Sobolev inequalities.

The milestone in the proof of the Burkholder-like inequalities is a relationship with the following flavor:

 $q \geq 2$  and  $\alpha$  satisfies  $q/2 \leq \alpha \leq q-1$ . Then

$$\mathbb{E}\left[\left(Z - \mathbb{E}\left[Z\right]\right)_{+}^{q}\right] - \mathbb{E}\left[\left(Z - \mathbb{E}\left[Z\right]\right)_{+}^{\alpha}\right]^{q/\alpha} \leq \frac{q\left(q - \alpha\right)}{2}\mathbb{E}\left[V\left(Z - \mathbb{E}\left[Z\right]\right)_{+}^{q-2}\right],$$

Letting  $\phi(x) = x^{q/\alpha}$ , this translates into

$$\mathbb{E}\left[\phi\left(\left(Z-\mathbb{E}\left[Z\right]\right)_{+}^{\alpha}\right)\right] - \phi\left(\mathbb{E}\left[\left(Z-\mathbb{E}\left[Z\right]\right)_{+}^{\alpha}\right]\right) \leq \frac{q\left(q-\alpha\right)}{2}\mathbb{E}\left[V\left(Z-\mathbb{E}\left[Z\right]\right)_{+}^{q-2}\right],$$

For q = 2 and  $\alpha = 1$ , this is (almost) Efron-Stein !!!

In order to establish Burkholder-like inequality, we proceed by induction on q. At each step we take  $\alpha = q - 1$ , that is  $\phi(x) = x^{q/(q-1)}$ . when dealing with moments of high order, we consider functions  $\phi$  that get closer and closer to  $x \log x$ .

$$\mathbb{E}\left[\left(Z - \mathbb{E}\left[Z\right]\right)_{+}^{q}\right] \leq \mathbb{E}\left[\left(Z - \mathbb{E}\left[Z\right]\right)_{+}^{q-1}\right]^{q/(q-1)} + \frac{q}{2}\mathbb{E}\left[V\left(Z - \mathbb{E}\left[Z\right]\right)_{+}^{q-2}\right]$$

In order to establish Burkholder-like inequality, we proceed by induction on q. At each step we take  $\alpha = q - 1$ , that is  $\phi(x) = x^{q/(q-1)}$ . when dealing with moments of high order, we consider functions  $\phi$  that get closer and closer to  $x \log x$ .

$$\mathbb{E}\left[\left(Z - \mathbb{E}\left[Z\right]\right)_{+}^{q}\right] \leq \mathbb{E}\left[\left(Z - \mathbb{E}\left[Z\right]\right)_{+}^{q-1}\right]^{q/(q-1)} + \frac{q}{2}\mathbb{E}\left[V\left(Z - \mathbb{E}\left[Z\right]\right)_{+}^{q-2}\right],$$

Hölder

$$\mathbb{E}\left[\left(Z - \mathbb{E}\left[Z\right]\right)_{+}^{q}\right] \leq \mathbb{E}\left[\left(Z - \mathbb{E}\left[Z\right]\right)_{+}^{q-1}\right]^{q/(q-1)} + \frac{q}{2}\mathbb{E}\left[V^{q/2}\right]^{2/q}\mathbb{E}\left[\left(Z - \mathbb{E}\left[Z\right]\right)_{+}^{q}\right]^{(q-2)/q},$$

Hölder

$$\mathbb{E}\left[\left(Z - \mathbb{E}\left[Z\right]\right)_{+}^{q}\right] \leq \mathbb{E}\left[\left(Z - \mathbb{E}\left[Z\right]\right)_{+}^{q-1}\right]^{q/(q-1)} + \frac{q}{2}\mathbb{E}\left[V^{q/2}\right]^{2/q}\mathbb{E}\left[\left(Z - \mathbb{E}\left[Z\right]\right)_{+}^{q}\right]^{(q-2)/q},$$

$$m_q \stackrel{\Delta}{=} \| (Z - \mathbb{E} [Z])_+ \|_q \text{ and } v_q \stackrel{\Delta}{=} \left\| \sqrt{V} \right\|_q$$

- Base case : Efron-Stein:  $m_2 \leq v_2$ .
- Assume  $m_{q-1} \leq \sqrt{3(q-1)}v_{q-1}$ ,

$$m_{q}^{q} \leq m_{q-1}^{q} + \frac{q}{2}v_{q}^{2}m_{q}^{q-2}$$

$$\leq \sqrt{3(q-1)}^{q}v_{q-1}^{q} + \frac{q}{2}v_{q}^{2}m_{q}^{q-2}$$

$$\dots$$

$$\leq \sqrt{3q}v_{q}.$$

 $\Phi$ : functions  $\phi$  on  $\mathbb{R}^+$  (Latala and Oleskiewicz)

- convex and continuous,
- twice differentiable on  $(0, +\infty)$ ,
- $\phi^{\prime\prime}$  is positive ,
- $1/\phi''$  is concave.

$$\phi$$
-entropy.  $H_{\phi}(Z) = \mathbb{E}[\phi(Z)] - \phi(\mathbb{E}[Z])$ .

Examples :

$$H_{x \log x}(Z) = \mathbb{E}[Z \log Z] - \mathbb{E}[Z] \log \mathbb{E}[Z]$$
$$H_{x^{p}}(Z) = \mathbb{E}[Z^{p}] - (\mathbb{E}[Z])^{p} \quad p \in (1, 2].$$

Tensorization of  $\phi$ -entropy.

$$H_{\phi}(Z) \leq \sum_{i=1}^{n} \mathbb{E}\left[\mathbb{E}\left[\phi(Z) \mid X_{1}^{i-1}, X_{i+1}^{n}\right] - \phi\left(\mathbb{E}\left[Z \mid X_{1}^{i-1}, X_{i+1}^{n}\right]\right)\right].$$

$$H_{\phi}(Z) \leq \sum_{i=1}^{n} H_{\phi}(Z \mid X_{1}^{i-1}, X_{i+1}^{n}).$$

Tensorization of  $\phi$ -entropy.

$$H_{\phi}(Z) \leq \sum_{i=1}^{n} \mathbb{E}\left[\mathbb{E}\left[\phi(Z) \mid X_{1}^{i-1}, X_{i+1}^{n}\right] - \phi\left(\mathbb{E}\left[Z \mid X_{1}^{i-1}, X_{i+1}^{n}\right]\right)\right].$$

**Duality** and  $\phi$ -entropy.

 $H_{\phi}$  is convex and continuous in  $\mathbb{L}_{1}^{+}$ . Obvious when the base space is discrete.

$$H_{\phi}(Z) = \sup_{T \in \mathbb{L}_{1}^{+}, T \neq 0} \left\{ \mathbb{E}\left[ \left( \phi'(T) - \phi'(\mathbb{E}[T]) \right) (Z - T) + \phi(T) \right] - \phi(\mathbb{E}[T]) \right\} .$$

The duality representation implies that  $H_{\phi}$  is convex and lower-semi-continuous. This is the basis of a Jensen-like property:

 $H(\mathbb{E}[Z \mid X_1]) \le H(Z \mid X_2)$ 

**Duality** and  $\phi$ -entropy.

 $H_{\phi}$  is convex and continuous in  $\mathbb{L}_1^+$ . Obvious when the base space is discrete.

$$H_{\phi}(Z) = \sup_{T \in \mathbb{L}_{1}^{+}, T \neq 0} \left\{ \mathbb{E}\left[ \left( \phi'(T) - \phi'(\mathbb{E}[T]) \right) (Z - T) + \phi(T) \right] - \phi(\mathbb{E}[T]) \right\} .$$

The duality representation implies that  $H_{\phi}$  is convex and lower-semi-continuous. This is the basis of a Jensen-like property:

 $H(\mathbb{E}[Z \mid X_1]) \le H(Z \mid X_2)$ 

Examples.

$$H_{\phi}(Z) = \sup_{T} \left\{ \mathbb{E}\left[ \left( \log\left(T\right) - \log\left(\mathbb{E}\left[T\right]\right) \right) Z \right] \right\} \text{ for } \phi(x) = x \log x$$
  
$$H_{\phi}(Z) = \sup_{T} \left\{ p \mathbb{E}\left[ Z\left(T^{p-1} - \left(\mathbb{E}\left[T\right]\right)^{p-1}\right) \right] - (p-1) H_{\phi}(T) \right\}, \text{ for } \phi(x) = x^{p}$$

Optimization.

If  $\phi \in \Phi$ , then both  $\phi'$  and  $x \to (\phi(x) - \phi(0)) / x$  are concave functions on  $\mathbb{R}^*_+$ .

Optimization.

If  $\phi \in \Phi$ , then both  $\phi'$  and  $x \to (\phi(x) - \phi(0)) / x$  are concave functions on  $\mathbb{R}^*_+$ .

Optimization.

If  $\phi \in \Phi$ , then both  $\phi'$  and  $x \to (\phi(x) - \phi(0)) / x$  are concave functions on  $\mathbb{R}^*_+$ . Let  $\psi$  denote the function  $x \to (\phi(x) - \phi(0)) / x$ .

 $H_{\phi}\left(f\left(Z\right)\right) \leq \mathbb{E}\left[V^{+}f^{\prime 2}\left(Z\right)\psi^{\prime}\left(f\left(Z\right)\right)\right] \quad \text{if}\psi \circ f \text{ is convex}.$ 

Optimization.

If  $\phi \in \Phi$ , then both  $\phi'$  and  $x \to (\phi(x) - \phi(0)) / x$  are concave functions on  $\mathbb{R}^*_+$ . Application :  $f(z) = \exp(\lambda z)$  and  $\phi(x) = x \log(x)$  :

$$H_{x \log x}\left(e^{\lambda Z}\right) \leq \lambda^2 \mathbb{E}\left[V^+ \exp\left(\lambda Z\right)\right],$$

Let  $q \geq 2$  and let  $\alpha$  satisfy  $q/2 \leq \alpha \leq q-1$ . Then

$$\mathbb{E}\left[\left(Z - \mathbb{E}\left[Z\right]\right)_{+}^{q}\right] \leq \mathbb{E}\left[\left(Z - \mathbb{E}\left[Z\right]\right)_{+}^{\alpha}\right]^{q/\alpha} + \frac{q\left(q - \alpha\right)}{2}\mathbb{E}\left[V\left(Z - \mathbb{E}\left[Z\right]\right)_{+}^{q-2}\right],$$

#### Khinchine inequalities (warm-up)

 $a_1, \ldots, a_n$ : non-negative constants,  $X_1, \ldots, X_n$  independent Rademacher variables. If  $Z = \sum_{i=1}^n a_i X_i$ then for any integer  $q \ge 2$ ,

$$||(Z)_+||_q = ||(Z)_-||_q \le \sqrt{2q} \sqrt{\sum_{i=1}^n a_i^2}$$

re Optimal constants can be obtained using hypercontractivity arguments.

If  $V^+ \le c$ ,  $||(Z - \mathbb{E}[Z])_+||_q \le \sqrt{qc}$ .

$$V^{+} = \sum_{i=1}^{n} \mathbb{E} \left[ (a_{i}(X_{i} - X_{i}'))_{+}^{2} \mid X_{i} \right] = 2 \sum_{i=1}^{n} a_{i}^{2} \mathbb{1}_{a_{i}X_{i} > 0} \le 2 \sum_{i=1}^{n} a_{i}^{2} ,$$

Rosenthal-Pinelis inequalities deal with suprema of empirical processes in the absence of uniform boundedness assumption.

Let  $\mathcal{F}$ : a countable class of measurable functions from  $\mathcal{X} \to \mathbb{R}$ .

 $X_1, \ldots, X_n$ : independent  $\mathcal{X}$ -valued random variables such that for all  $f \in \mathcal{F}$ ,  $\mathbb{E}f(X_i) = 0$ .

$$\Sigma^{2} = \mathbb{E} \Big[ \sup_{f} \sum_{i} f^{2}(X_{i}) \Big]$$
  
$$\sigma^{2} = \sup_{f} \mathbb{E} \Big[ \sum_{i} f^{2}(X_{i}) \Big]$$
  
$$M = \sup_{f,i} |f(X_{i})|.$$

In Talagrand's inequality, M is assumed to be bounded. Here, integrability of M is related wit integrability of Z.

$$Z = \sup_{f \in \mathcal{F}} \left| \sum_{i=1}^{n} f(X_i) \right|.$$

$$Z = \sup_{f \in \mathcal{F}} \left| \sum_{i=1}^{n} f(X_i) \right|.$$

The application of Burkholder like inequality leads to:

$$\left\| (Z - \mathbb{E}[Z])_+ \right\|_q \leq \sqrt{6q} \left( \Sigma + \sigma \right) + 6q \left( \left\| M \right\|_q + \sup_{i, f \in \mathcal{F}} \left\| f(X_i) \right\|_2 \right),$$

 $||Z||_q \leq 2\mathbb{E}Z + 2\sigma\sqrt{6q} + 60q ||M||_q + 4\sqrt{3q} ||M||_2$ .

$$Z = \sup_{f \in \mathcal{F}} \left| \sum_{i=1}^{n} f(X_i) \right|.$$

$$\|(Z - \mathbb{E}[Z])_+\|_q \leq \sqrt{6q} \left(\Sigma + \sigma\right) + 6q \left(\|M\|_q + \sup_{i, f \in \mathcal{F}} \|f(X_i)\|_2\right),$$

As  $\mathbb{E}[f] = 0$  for all  $f \in \mathcal{F}$ :

$$V^{+} \leq \sup_{f \in \mathcal{F}} \sum_{i} f^{2}(X_{i}) + \sup_{f \in \mathcal{F}} \mathbb{E}\left[f^{2}(X_{i})\right] = \sup_{f \in \mathcal{F}} \sum_{i} f^{2}(X_{i}) + \sigma^{2}$$

$$\sqrt{V^+} \leq \sup_{f \in \mathcal{F}} \sup_{\alpha: \sum_i \alpha_i^2 \leq 1} \sum_i \alpha_i f(X_i) + \sigma.$$

Applying the second Burkholder-like inequality again:

$$\left\| \sup_{f \in \mathcal{F}} \sup_{\alpha: \sum_{i} \alpha_{i}^{2} \leq 1} \sum_{i} \alpha_{i} f(X_{i}) \right\|_{q} \leq \Sigma + \sqrt{3q} \|M\|_{q}$$

$$Z = \sup_{f \in \mathcal{F}} \left| \sum_{i=1}^{n} f(X_i) \right|.$$

 $||Z||_q \leq 2\mathbb{E}Z + 2\sigma\sqrt{6q} + 60q ||M||_q + 4\sqrt{3q} ||M||_2$ 

To get this inequality, need to relate  $\Sigma, \mathbb{E}[Z]$  and  $\sigma$ 

In the uniformly bounded case, symmetrization and the contraction principle were enough. We need a new ingredient.

Hoffman-Jorgensen inequality. Let  $\epsilon_1, \ldots, \epsilon_n$  denote independent Rademacher variables. Let  $\lambda > 4$  and define  $t_0 = \sqrt{\lambda \mathbb{E}[M^2]}$ .

$$\mathbb{E}\left[\sup_{f}\left|\sum_{i}\epsilon_{i}f^{2}(X_{i})\mathbb{1}_{\sup_{f}|f(X_{i})|>t_{0}}\right|\right] \leq \frac{1}{(1-2/\sqrt{\lambda})^{2}}\mathbb{E}\left[M^{2}\right] .$$

$$Z = \sup_{f \in \mathcal{F}} \left| \sum_{i=1}^{n} f(X_i) \right|.$$

$$\begin{split} \Sigma^{2} &\leq \mathbb{E} \left[ \sup_{f} \left| \sum_{i} f^{2}(X_{i}) - \mathbb{E}[f^{2}(X_{i})] \right| \right] + \sup_{f} \mathbb{E} \left[ \sum_{i} f^{2}(X_{i}) \right] \\ &\leq \sigma^{2} + 2 \mathbb{E} \left[ \sup_{f} \left| \sum_{i} \epsilon_{i} f^{2}(X_{i}) \mathbb{1}_{\sup_{f} |f(X_{i})| \leq t_{0}|} \right] \\ &+ 2 \mathbb{E} \left[ \sup_{f} \left| \sum_{i} \epsilon_{i} f^{2}(X_{i}) \mathbb{1}_{\sup_{f} |f(X_{i})| > t_{0}|} \right] \\ &\leq \sigma^{2} + 4t_{0} \mathbb{E} \left[ \sup_{f} \left| \sum_{i} \epsilon_{i} f(X_{i}) \right| \right] + 2 \mathbb{E} \left[ \sup_{f} \left| \sum_{i} \epsilon_{i} f^{2}(X_{i}) \mathbb{1}_{\sup_{f} |f(X_{i})| > t_{0}|} \right] \\ &\leq \sigma^{2} + 4t_{0} \mathbb{E} \left[ \sup_{f} \left| \sum_{i} \epsilon_{i} f(X_{i}) \right| \right] + 2 \mathbb{E} \left[ \sup_{f} \left| \sum_{i} \epsilon_{i} f^{2}(X_{i}) \mathbb{1}_{\sup_{f} |f(X_{i})| > t_{0}|} \right] \\ &\leq \sigma^{2} + 4t_{0} \mathbb{E} \left[ \sup_{f} \left| \sum_{i} \epsilon_{i} f(X_{i}) \right| \right] + \frac{2}{(1 - \lambda/2)^{2}} \mathbb{E} [M^{2}] \\ &\leq \sigma^{2} + 8 \sqrt{\lambda \mathbb{E} [M^{2}]} ||Z||_{1} + \frac{2}{(1 - 2/\sqrt{\lambda})^{2}} \mathbb{E} \left[ M^{2} \right]_{\text{S. Boucheron Concentration Inequalities, August, 2003 - p.444} \end{split}$$

#### **Rademacher complexity**

 $\mathcal F$  : countable class of measurable centered real-valued functions.

$$Z = \mathbb{E} \Big[ \sup_{f \in \mathcal{F}} \left| \sum_{i} \epsilon_{i} f(X_{i}) \right| \mid X_{1}^{n} \Big]$$

$$M = \sup_{i,f} \left| f(X_i) \right|.$$

Tools derived from Burkholder-like inequalities:

 $\| {\rm f} \, V \leq WZ, \quad \| (Z - \mathbb{E}[Z])_+ \|_q \leq 2 \sqrt{q \, \|W\|_q \, \mathbb{E}[Z]} + 2q \, \|W\|_q \, \, .$ 

#### **Rademacher complexity**

Tools derived from Burkholder-like inequalities:

If 
$$V \le WZ$$
,  $||(Z - \mathbb{E}[Z])_+||_q \le 2\sqrt{q ||W||_q \mathbb{E}[Z]} + 2q ||W||_q$ .

$$Z_{i} = \mathbb{E} \left[ \sup_{f \in \mathcal{F}} \left| \sum_{j \neq i} \epsilon_{j} f(X_{j}) \right| \mid X^{(i)} \right]$$

$$\forall i, \qquad 0 \leq Z - Z_i \leq M \& \sum_i Z - Z_i \leq Z \lor V \leq M Z.$$

#### **Rademacher complexity**

Tools derived from Burkholder-like inequalities:

If 
$$V \le WZ$$
,  $||(Z - \mathbb{E}[Z])_+||_q \le 2\sqrt{q ||W||_q \mathbb{E}[Z]} + 2q ||W||_q$ .

Hence:

$$||(Z - \mathbb{E}[Z])_+||_q \le 2\sqrt{q ||M||_q \mathbb{E}[Z]} + 2q ||M||_q$$
.

If  $\mathcal{F}$  is uniformly-bounded, conditional Rademacher averages have sub-Poissonian tails. If  $\mathcal{F}$  is a VC-class of VC-dimension d,  $\mathbb{E}[Z] \leq C\sqrt{dn}$ , and

$$\mathbb{P}\left\{Z \ge \mathbb{E}\left[Z\right] + t\right\} \le e^{-\frac{t^2}{2(\mathbb{E}[Z] + t/3)}}.$$

If the maximum of n envelops of  $\mathcal{F}$  is q-integrable, so is Z.
## References

## References

Concentration inequalities Massart. About the constants in Talagrand's concentration inequalities for empirical processes, Ann. Probab. 28, 863-884, 2000. Boucheron, Lugosi and Massart. A sharp concentration inequality and applications, Rand. Struct. Alg., 16, 277-292, 2000. Boucheron, Lugosi and Massart. Concentration inequalities by the entropy method, Ann. Probab. 31, No. 3, 2003. Bousquet. Concentration inequalities for sub-additive funnctions. . 2002. Moment inequalities. Boucheron, Bousquet, Lugosi and Massart. Moment inequalities for functions of independent random variables, Manuscript, 2003.