

Concentration

for functions of bounded

interaction

Andreas Maurer

Which properties of a bounded function

$$f : \mathcal{X}^n \rightarrow \mathbb{R}$$

guarantee nice behaviour of the random variable  $f(\mathbf{X})$ ,  
with  $X = (X_1, \dots, X_n)$  a vector of independent random variables?

## Nice behaviour of sums

$$f(\mathbf{x}) = \sum_{i=1}^n g_i(x_i) \text{ with } g_i : \mathcal{X} \rightarrow [a, b].$$

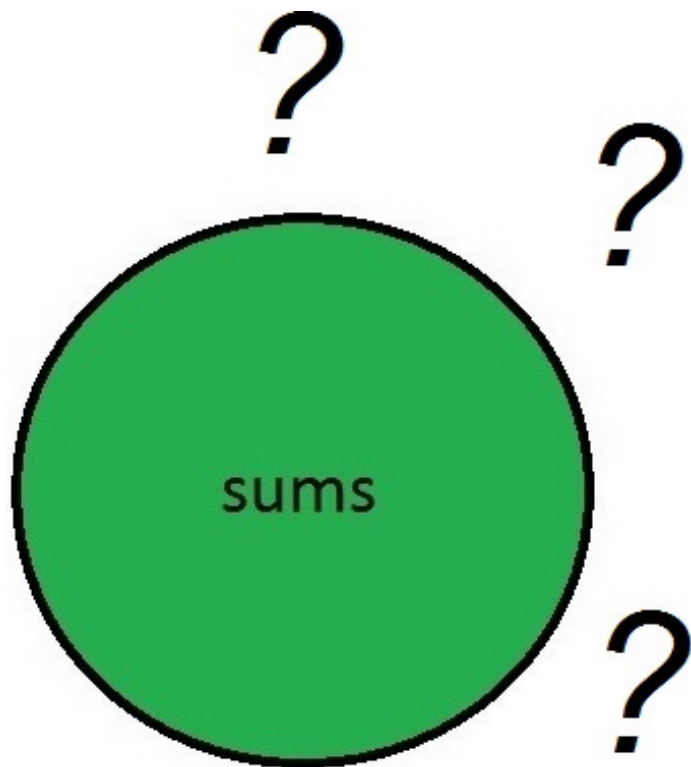
Additive variance  $\sigma^2(f) = \sum_{i=1}^n \sigma^2(g_i)$

normal approximation  $\frac{f - Ef}{\sigma(f)} \approx \mathcal{N}(0, 1)$  for large  $n$

Hoeffding inequality  $\Pr\{f - Ef > t\} \leq \exp\left(\frac{-2t^2}{n(b-a)^2}\right)$

Bernstein inequality  $\Pr\{f - Ef > t\} \leq \exp\left(\frac{-t^2}{2\sigma^2(f) + 2(b-a)t/3}\right)$

What about functions which are not sums?



# The bounded difference inequality

Partial difference operator

$$D_{y,y'}^k f(\mathbf{x}) := f(\dots, x_{k-1}, y, x_{k+1}, \dots) - f(\dots, x_{k-1}, y', x_{k+1}, \dots).$$

Define maximal variation in any argument

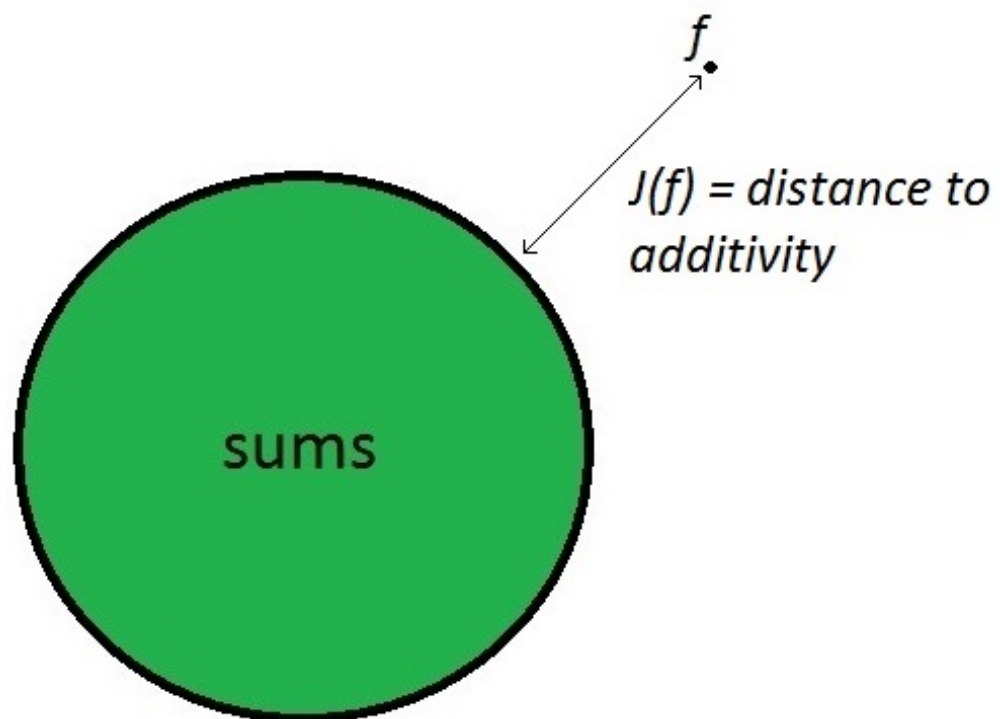
$$M(f) := \max_k \sup_{\mathbf{x}, y, y'} D_{y,y'}^k f(\mathbf{x}).$$

**Theorem** (Hoeffding, Azuma, McDiarmid):

$$\Pr\{f - Ef > t\} \leq \exp\left(\frac{-2t^2}{nM(f)^2}\right), \text{ for all } f : \mathcal{X}^n \rightarrow \mathbb{R}$$

Extends Hoeffding's inequality to general functions.

What about functions which are close to being sums?



# Interaction functional

For bounded  $f : \mathcal{X}^n \rightarrow \mathbb{R}$  define

$$J(f) := n \max_{k,l:k \neq l} \sup_{\mathbf{x}, y, y', z, z'} D_{z, z'}^l D_{y, y'}^k f(\mathbf{x}).$$

- $J$  is a seminorm which vanishes for sums
- Key property for "sum-like" behaviour:  
 $J(f)$  is at most of the same order in  $n$  as  $M(f)$   
or, put differently:  
maximal mixed second difference  $\leq O(1/n) \times$  maximal first difference

# Weak interactions

**Definition:**

$f : \mathcal{X}^n \rightarrow \mathbb{R}$  has  $(a, b)$ -weak interactions,  
if  $M(f) \leq a/n$  and  $J(f) \leq b/n$ .

A sequence  $(f_n)_{n \geq 2}$  of functions  $f_n : \mathcal{X}^n \rightarrow \mathbb{R}$  has  $(a, b)$ -weak interactions  
if every  $f_n$  has  $(a, b)$ -weak interactions.



# Outline

Examples of weak interactions:

- U- and V-statistics
- Lipschitz L-statistics
- Generalization error of the Gibbs algorithm
- Generalization error of  $\ell_2$ -regularized classification

Properties of weak interactions:

- Small bias in the Efron-Stein inequality
- Bernstein's inequality
- Variance estimation
- Normal approximation

Proof of Bernstein's inequality

# V- and U-statistics

Fix  $1 \leq m < n$ ,

for  $\mathbf{j} = (j_1, \dots, j_m) \in \{1, \dots, n\}^m$  let

$$\kappa_{\mathbf{j}} : \mathcal{X}^m \rightarrow \mathbb{R}, \quad |\kappa_{\mathbf{j}}| \leq 1$$

and define  $V, U : \mathcal{X}^m \rightarrow \mathbb{R}$ ,

$$V(\mathbf{x}) = n^{-m} \sum_{\mathbf{j} \in \{1, \dots, n\}^m} \kappa_{\mathbf{j}}(x_{j_1}, \dots, x_{j_m})$$
$$U(\mathbf{x}) = \binom{n}{m}^{-1} \sum_{1 \leq j_1 < \dots < j_m \leq m} \kappa_{\mathbf{j}}(x_{j_1}, \dots, x_{j_m})$$

$V$  = **V**on Mises statistic (1947)

$U$  = **U**nbiased statistic (Hoeffding, 1948)

V- and U-statistics have weak interactions

$$V(\mathbf{x}) = n^{-m} \sum_{\mathbf{j} \in \{1, \dots, n\}^m} \kappa_{\mathbf{j}}(x_{j_1}, \dots, x_{j_m})$$

$$\begin{aligned} D_{y, y'}^k V(\mathbf{x}) &\leq \frac{2}{n^m} |\{\mathbf{j} : k \in \mathbf{j}\}| = \frac{2}{n^m} \left| \bigcup_{r=1}^m \left\{ \mathbf{j} : r = \min_{j_i=k} i \right\} \right| \\ &= \frac{2mn^{m-1}}{n^m} = \frac{2m}{n} \end{aligned}$$

$$\begin{aligned} D_{z, z'}^l D_{y, y'}^{k: k \neq l} V(\mathbf{x}) &\leq \frac{4}{n^m} |\{\mathbf{j} : k, l \in \mathbf{j}\}| = \frac{4}{n^m} \left| \bigcup_{r, s: r \neq s} \left\{ \mathbf{j} : r = \min_{j_i=k} i \wedge s = \min_{j_i=l} i \right\} \right| \\ &= \frac{4m(m-1)n^{m-2}}{n^m} = \frac{4m(m-1)}{n^2}. \end{aligned}$$

So  $V$  has  $(2m, 4m(m-1))$ -weak interactions!

Similar argument and result for  $U$  (A.M,17a)

## Lipschitz L-statistics

$\mathcal{X} = [a, b]$  and  $(x_{(1)}, \dots, x_{(n)})$  = order statistic of  $\mathbf{x} \in \mathcal{X}^n$

$$f(\mathbf{x}) = \frac{1}{n} \sum_{i=1}^n F(i/n) x_{(i)}$$

where  $F : [0, 1] \rightarrow \mathbb{R}$  has Lipschitz constant  $L$ .

Examples: mean, smoothly trimmed mean, smoothed quantiles, etc

Then  $f$  as  $(L(b-a), L(b-a))$ -weak interactions

(bound on  $M(f)$  easy, bound on  $J(f)$  cumbersome - many cases)

## A chain rule

Extend definition of  $M$  and  $J$  to Banach space-valued functions  $f : \mathcal{X}^n \rightarrow B$

$$M(f) = \max_k \sup_{x,y,y'} \left\| D_{yy'}^k f(x) \right\| \quad \text{and} \quad J(f) = n \max_{k \neq l} \sup_{x,y,y',z,z'} \left\| D_{zz'}^l D_{yy'}^k f(x) \right\|.$$

**Lemma:**  $B$  be a Banach space,  $U \subseteq B$  convex,  $f : \mathcal{X}^n \rightarrow U$ ,  $F : U \rightarrow \mathbb{R}$  be twice Fréchet-differentiable. Then

$$\begin{aligned} M(F \circ f) &\leq \sup_{v \in U} \left\| F'(v) \right\| M(f) \quad \text{and} \\ J(F \circ f) &\leq n \sup_{v \in U} \left\| F''(v) \right\| M(f)^2 + \sup_{v \in U} \left\| F'(v) \right\| J(f). \end{aligned}$$

If  $f$  has weak interactions and  $\left\| F''(v) \right\|$  and  $\left\| F'(v) \right\|$  are bounded on  $U$ , then  $F \circ f$  also has weak interactions.

# Free energy and Gibbs distributions

$\Omega$  a mble space with finite, positive measure  $\rho$ .

$H : \Omega \times \mathcal{X}^n \rightarrow [-1, 1]$  a "Hamiltonian",

$\beta > 0$  an "inverse temperature".

For  $x \in \mathcal{X}^n$  define 
$$A_H(\mathbf{x}) = \ln Z_H(\mathbf{x}) := \ln \int_{\Omega} e^{-\beta H(\omega, \mathbf{x})} d\rho(\omega)$$

The chain rule with

$f : \mathbf{x} \in \mathcal{X}^n \mapsto H(\cdot, \mathbf{x}) \in L_{\infty}(\Omega)$  and

$F : G(\cdot) \in L_{\infty}(\Omega) \mapsto \ln \int_{\Omega} e^{-\beta G(\omega)} d\rho(\omega)$

gives 
$$M(A_H) \leq \beta M(H) \text{ and } J(A_H) \leq \beta J(H) + 2n\beta^2 M(H)^2$$

Also define Gibbs distribution on  $\Omega$

$$d\pi_H(\mathbf{x}) = Z_H^{-1}(\mathbf{x}) e^{-\beta H(\omega, \mathbf{x})} d\rho(\omega)$$

## "Generalization" of the Gibbs algorithm

loss of model  $\omega$  on datum  $x$  :  $\ell(\omega, x)$  where  $\ell : \Omega \times \mathcal{X} \rightarrow [0, 1]$

empirical loss on sample  $\mathbf{x}$  :  $H(\omega, \mathbf{x}) = \frac{1}{n} \sum_{i=1}^n \ell(\omega, x_i)$

true loss on r.v.  $X$  :  $\bar{H}(\omega) = E_X[\ell(\omega, X)]$ .

Gibbs measure for empirical loss :  $d\pi_H(\mathbf{x})$

Gibbs measure for true loss :  $d\pi_{\bar{H}}$

a measure of "generalization" :  $KL(d\pi_H(\mathbf{x}), d\pi_{\bar{H}}) =: f(\mathbf{x})$

By the chain rule

$$M(f) \leq \frac{6\beta(1+2\beta)}{n} \text{ and } J(f) \leq \frac{24\beta^2(1+2\beta)}{n},$$

So  $f$  has  $(6\beta(1+2\beta), 24\beta^2(1+2\beta))$ -weak interactions!

# Generalization of $\ell_2$ -regularized algorithms

$(H, \langle \cdot, \cdot \rangle, \|\cdot\|)$  a real Hilbert space with unit ball  $\mathcal{X}$   
define  $g : \mathcal{X}^n \rightarrow H$  by

$$\text{returned weight vector } g(\mathbf{x}) = \arg \min_{w \in H} \frac{1}{n} \sum_{i=1}^n \ell(\langle x_i, w \rangle) + \lambda \|w\|^2$$

$$\text{empirical loss } \hat{L}(\mathbf{x}) = \frac{1}{n} \sum_i \ell(\langle x_i, g(\mathbf{x}) \rangle),$$

$$\text{true expected loss } L(\mathbf{x}) = E[\ell(\langle X, g(\mathbf{x}) \rangle)],$$

$$\text{generalization error } \Delta(\mathbf{x}) = L(\mathbf{x}) - \hat{L}(\mathbf{x})$$

Then  $\Delta$  has  $(O(\lambda^{-3/2}) \|\ell''\|_\infty, O(\lambda^{-3}) \|\ell'''\|_\infty)$ -weak interactions!  
(A.M.17b, by implicit differentiation)



# Properties of functions with weak interactions

- Small bias in the Efron-Stein inequality
- Bernstein's inequality
- Variance estimation
- Normal approximation

## The bias of the Efron-Stein inequality

$$k\text{-th conditional variance} \quad : \quad \sigma_k^2(f)(\mathbf{x}) = \frac{1}{2} E_{(y,y') \sim \mu_k \times \mu_k} \left[ \left( D_{y,y'}^k f(\mathbf{x}) \right)^2 \right]$$

$$\text{sum of conditional variances} \quad : \quad \Sigma^2(f)(\mathbf{x}) = \sum_{k=1}^n \sigma_k^2(f)(\mathbf{x})$$

$$\text{Efron-Stein inequality} \quad : \quad \sigma^2(f) \leq E \left[ \Sigma^2(f) \right]$$

**Theorem** (Houdré, 1998):

$$E \left[ \Sigma^2(f) \right] \leq \sigma^2(f) + \frac{1}{4} \sum_{k,l:k \neq l} E_{\mathbf{x},z,z',y,y'} \left[ \left( D_{zz'}^l D_{yy'}^k f(\mathbf{x}) \right)^2 \right] \leq \sigma^2(f) + \frac{J(f)^2}{4}.$$

If  $f$  has weak interactions then  $\sigma^2(f) = E \left[ \Sigma^2(f) \right] + O(1/n^2)$ .

# Bernstein's inequality

**Theorem** (A.M.17a): For bounded mble  $f : \mathcal{X}^n \rightarrow \mathbb{R}$

$$\Pr \{f - E[f] > t\} \leq \exp \left( \frac{-t^2}{2E[\Sigma^2(f)] + (2M(f)/3 + J(f)) t} \right)$$

extends Bernstein's inequality from sums to general functions.

**Corollary:** If  $f$  has  $(a, b)$ -weak interactions then

(using  $E[\Sigma^2(f)] \leq \sigma^2(f) + J(f)/4$ )

$\forall \delta \in (0, 1/e)$  with probability at least  $1 - \delta$

$$f \leq E[f] + \sqrt{2\sigma^2(f) \ln(1/\delta)} + \frac{(2a/3 + 2b) \ln(1/\delta)}{n}.$$

## Estimating variance

**Theorem:** For any bounded  $f : \mathcal{X}^n \rightarrow \mathbb{R}$   
there exists  $g : \mathcal{X}^{n+1} \rightarrow \mathbb{R}$  such that  
for any iid sequence  $X_1, \dots, X_n, \dots$  with values in  $\mathcal{X}$   
and for  $0 < \delta \leq 1/e$  with probability at least  $1 - \delta$

$$\sqrt{g(\mathbf{X})} - K_1(f) \sqrt{\ln(2/\delta)} \leq \sqrt{\sigma^2(f)} \leq \sqrt{g(\mathbf{X})} + K_2(f) \sqrt{\ln(2/\delta)}$$

$$\text{with } K_1(f) = J(f)/2 + \sqrt{M(f)^2 + 8J(f)^2}$$

$$\text{and } K_2(f) = \sqrt{\max\{M(f)^2 + 8J(f)^2, M(f)(M(f) + 2J(f))\}}$$

Also:  $g$  is an unbiased estimator for the Efron Stein bound  $E[\Sigma^2(f)]$ .

# The variance estimator

For any  $n$  and  $\mathbf{x} \in \mathcal{X}^n$  define

$$\begin{aligned} \text{replacement operator} \quad S_y^k \mathbf{x} &= (x_1, \dots, x_{k-1}, y, x_{k+1}, \dots, x_n) \in \mathcal{X}^n \\ \text{deletion operator} \quad S_-^k \mathbf{x} &= (x_1, \dots, x_{k-1}, x_{k+1}, \dots, x_n) \in \mathcal{X}^{n-1}. \end{aligned}$$

The variance estimator  $g : \mathcal{X}^{n+1} \rightarrow \mathbb{R}$  is

$$g(\mathbf{x}) = \frac{1}{2(n+1)} \sum_{i=1}^{n+1} \sum_{j:j \neq i} \left( f(S_-^j \mathbf{x}) - f(S_-^j S_{x_j}^i \mathbf{x}) \right)^2.$$

Needs  $O(n^2)$  computations of  $f$ , but only a sample of  $O(n)$ !

So for weak interactions with high probability

$$\sqrt{\sigma^2(f)} = \sqrt{g(X)} + O\left(\frac{1}{n}\right).$$

## Normal approximation

$$\delta(W, V) = \{\sup |E[h(W)] - E[h(V)]| : h \text{ a Lipschitz-1 function}\}$$

**Theorem** : Let  $E[f] = 0$  and  $Z \sim \mathcal{N}(0, 1)$ . Then

$$\delta(f(\mathbf{X}), Z) \leq \frac{\sqrt{8n}M(f)(J(f) + M(f))}{\sigma^2(f)} + \frac{nM(f)^3}{2\sigma^3(f)}.$$

If  $(f_n)$  has  $(a, b)$ -weak interactions and  $\sigma(f_n) \geq Cn^{-p}$  for constant  $C$ , then

$$\delta(f_n(\mathbf{X}), Z) \leq \frac{\sqrt{8}Ca(a+b) + a^3}{C^3n^{2-3p}}.$$

$(1/2 \leq p < 2/3) \implies$  asymptotic normality.

$(p = 1/2) \implies$  rate is  $n^{-1/2}$ .

.

# Intermission

## A Bernstein-type inequality

**Theorem:** Let  $f : \mathcal{X}^n \rightarrow \mathbb{R}$  and for some  $b$  and all  $m \geq 2$

$$\sum_{k=1}^n E_k [(f - E_k f)^m] \leq \frac{\Sigma^2(f)}{2} m! b^{m-2}.$$

Then for  $t > 0$

$$\Pr \{f - Ef > t\} \leq \exp \left( \frac{-t^2}{2E [\Sigma^2(f)] + (2b + J(f)) t} \right).$$

The proof uses the *entropy method* (Boltzmann, Gibbs, Shannon, Nelson, Lieb, Ledoux, Bobkov, Massart, Boucheron, Lugosi, and many others)

We prove this for  $b = 1$ . The theorem then follows from rescaling



## Definitions

thermal measure	$\mu_{\beta f} : = e^{\beta f} d\mu / E [e^{\beta f}]$
thermal expectation	$E_{\beta f} [g] : = E [g e^{\beta f}] / E [e^{\beta f}]$
thermal variance	$\sigma_{\beta f}^2 [g] : = E_{\beta f} [(g - E_{\beta f} [g])^2]$
entropy	$\text{Ent} (\beta f) : = KL (d\mu_{\beta f}, d\mu) = \beta E_{\beta f} [f] - \ln E [e^{\beta f}]$
conditional expectation	$E_k [g] : = E [g   X_1, \dots, X_{k-1}, X_{k+1}, \dots, X_n]$
the conditional quantities	$E_{k, \beta f} [\cdot], \sigma_{k, \beta f}^2 [\cdot], \text{Ent}_k (\beta f)$ are all w.r.t. $E_k [\cdot]$

replacement operator	$S_y^k x = (x_1, \dots, x_{k-1}, y, x_{k+1}, \dots, x_n) \in \mathcal{X}^n$
	$(S_y^k f) (x) = f (S_y^k x)$ for $f : \mathcal{X}^n \rightarrow \mathcal{Y}$

## Facts used in the proof

Markov inequality	$\Pr \{f - Ef > t\} \leq \inf_{\beta \geq 0} \exp \left( \ln E \left[ e^{\beta(f - Ef)} \right] - \beta t \right)$
representation of mgf	$\ln E \left[ e^{\beta(f - Ef)} \right] = \beta \int_0^\beta \frac{\text{Ent}(\gamma f)}{\gamma^2} d\gamma$
subadditivity of entropy	$\text{Ent}(\beta f) \leq E_{\beta f} \left[ \sum_{k=1}^n \text{Ent}_k(\beta f) \right]$
fluctuation representation	$\text{Ent}(\beta f) = \int_0^\beta \int_t^\beta \sigma_{sf}^2[f] ds dt$
	$\text{Ent}(\beta f) \leq E_{\beta f} \left[ \sum_{k=1}^n \int_0^\beta \int_t^\beta \sigma_{k,sf}^2[f] ds dt \right]$
Fenchel-Young inequality	$E_{\beta f} [g] \leq \text{Ent}(\beta f) + \ln E [e^g]$

## The sum of conditional variances, I

**Lemma:** Let  $\beta \in (0, 1)$  and suppose that for all  $m \geq 2$

$$\sum_{k=1}^n E_k [(f - E_k f)^m] \leq \frac{\Sigma^2(f)}{2} m!.$$

$$\text{Then} \quad \text{Ent}(\beta f) \leq \frac{\beta^2}{2(1-\beta)^2} E_{\beta f} [\Sigma^2(f)].$$

**Proof:** For  $s \in (0, \beta)$

$$\begin{aligned} \sigma_{k,sf}^2(f) &\leq E_{k,sf} [(f - E_k(f))^2] \quad (\text{variational property of variance}) \\ &= \frac{E_k [(f - E_k(f))^2 e^{s(f - E_k f)}]}{E_k [e^{s(f - E_k f)}]} \end{aligned}$$

## The sum of conditional variances, II

**Lemma:** Let  $\beta \in (0, 1)$  and suppose that for all  $m \geq 2$

$$\sum_{k=1}^n E_k [(f - E_k f)^m] \leq \frac{\Sigma^2(f)}{2} m!.$$

$$\text{Then} \quad \text{Ent}(\beta f) \leq \frac{\beta^2}{2(1-\beta)^2} E_{\beta f} [\Sigma^2(f)].$$

**Proof:** For  $s \in (0, \beta)$

$$\begin{aligned} \sigma_{k, sf}^2(f) &\leq \frac{E_k [(f - E_k(f))^2 e^{s(f - E_k f)}]}{E_k [e^{s(f - E_k f)}]} \\ &\leq E_k [(f - E_k(f))^2 e^{s(f - E_k f)}] \quad (\text{Jensen's inequality}) \end{aligned}$$

## The sum of conditional variances, III

**Lemma:** Let  $\beta \in (0, 1)$  and suppose that for all  $m \geq 2$

$$\sum_{k=1}^n E_k [(f - E_k f)^m] \leq \frac{\Sigma^2(f)}{2} m!.$$

$$\text{Then} \quad \text{Ent}(\beta f) \leq \frac{\beta^2}{2(1-\beta)^2} E_{\beta f} [\Sigma^2(f)].$$

**Proof:** For  $s \in (0, \beta)$

$$\begin{aligned} \sum_{k=1}^n \sigma_{k, sf}^2(f) &\leq \sum_{k=1}^n E_k [(f - E_k f)^2 e^{s(f - E_k f)}] \\ &= \sum_{m=0}^{\infty} \frac{s^m}{m!} \sum_{k=1}^n E_k [(f - E_k f)^{m+2}] \end{aligned}$$

## The sum of conditional variances, IV

**Lemma:** Let  $\beta \in (0, 1)$  and suppose that for all  $m \geq 2$

$$\sum_{k=1}^n E_k [(f - E_k f)^m] \leq \frac{\Sigma^2(f)}{2} m!.$$

$$\text{Then} \quad \text{Ent}(\beta f) \leq \frac{\beta^2}{2(1-\beta)^2} E_{\beta f} [\Sigma^2(f)].$$

**Proof:** For  $s \in (0, \beta)$

$$\begin{aligned} \sum_{k=1}^n \sigma_{k, sf}^2(f) &\leq \sum_{m=0}^{\infty} \frac{s^m}{m!} \sum_{k=1}^n E_k [(f - E_k f)^{m+2}] \\ &\leq \frac{\Sigma^2(f)}{2} \sum_{m=0}^{\infty} (m+1)(m+2) s^m \quad (\text{by hypothesis}) \end{aligned}$$

## The sum of conditional variances, $V$

**Lemma:** Let  $\beta \in (0, 1)$  and suppose that for all  $m \geq 2$

$$\sum_{k=1}^n E_k [(f - E_k f)^m] \leq \frac{\Sigma^2(f)}{2} m!.$$

$$\text{Then} \quad \text{Ent}(\beta f) \leq \frac{\beta^2}{2(1-\beta)^2} E_{\beta f} [\Sigma^2(f)].$$

**Proof:**

$$\begin{aligned} \text{Ent}(\beta f) &\leq E_{\beta f} \left[ \int_0^\beta \int_t^\beta \sum_{k=1}^n \sigma_{k, sf}^2(f) ds dt \right] \text{ (subadditivity + fluctuation rep.)} \\ &\leq \frac{E_{\beta f} [\Sigma^2(f)]}{2} \sum_{m=0}^{\infty} (m+1)(m+2) \int_0^\beta \int_t^\beta s^m ds dt \end{aligned}$$

## The sum of conditional variances, VI

**Lemma:** Let  $\beta \in (0, 1)$  and suppose that for all  $m \geq 2$

$$\sum_{k=1}^n E_k [(f - E_k f)^m] \leq \frac{\Sigma^2(f)}{2} m!.$$

$$\text{Then} \quad \text{Ent}(\beta f) \leq \frac{\beta^2}{2(1-\beta)^2} E_{\beta f} [\Sigma^2(f)].$$

**Proof:**

$$\begin{aligned} \text{Ent}(\beta f) &\leq \frac{E_{\beta f} [\Sigma^2(f)]}{2} \sum_{m=0}^{\infty} (m+1)(m+2) \int_0^{\beta} \int_t^{\beta} s^m ds dt \\ &= \frac{E_{\beta f} [\Sigma^2(f)]}{2} \beta^2 \sum_{m=0}^{\infty} (m+1) \beta^m \quad \square \end{aligned}$$



## Spin-off: another version of Bernstein's inequality

**Theorem** (McDiarmid 1998): If  $f$  satisfies the conditions of the lemma then

$$\begin{aligned}\Pr\{f - Ef > t\} &\leq \inf_{\beta \in (0,1)} \exp\left(\beta \int_0^\beta \frac{\text{Ent}(\gamma f)}{\gamma^2} d\gamma - \beta t\right) \\ &\leq \inf_{\beta \in (0,1)} \exp\left(\beta \int_0^\beta \frac{1}{2(1-\gamma)^2} E_{\gamma f}[\Sigma^2(f)] d\gamma - \beta t\right) \\ &\leq \inf_{\beta \in (0,1)} \exp\left(\frac{\|\Sigma^2(f)\|_\infty}{2} \frac{\beta^2}{1-\beta} - \beta t\right) \\ &\leq \exp\left(\frac{-t^2}{2(\|\Sigma^2(f)\|_\infty + t)}\right)\end{aligned}$$

## Another bound on entropy, I

**Define**

$$D^2 f \quad : \quad = \sum_k \left( f - \inf_{y \in \mathcal{X}} S_y^k f \right)^2$$

**Lemma:**

$$\text{Ent}(\beta f) \leq \frac{\beta^2}{2} E_{\beta f} [D^2(f)]$$

**Proof:** For  $0 < s \leq \beta$ . Let  $h := f - \inf_{y \in \mathcal{X}} S_y^k f$ .

$$\begin{aligned} \frac{d}{ds} E_{k,sf} [h^2] &= \frac{d}{ds} E_{k,sh} [h^2] \\ &= E_{k,sh} [h^3] - E_{k,sh} [h^2] E_{k,sh} [h] \geq 0 \end{aligned}$$

so

$$\sigma_{k,sf}^2(f) \leq E_{k,sf} \left[ \left( f - \inf_{y \in \mathcal{X}} S_y^k f \right)^2 \right] \leq E_{k,\beta f} \left[ \left( f - \inf_{y \in \mathcal{X}} S_y^k f \right)^2 \right]$$

## Another bound on entropy, II

**Define**

$$D^2 f \quad : \quad = \sum_k \left( f - \inf_{y \in \mathcal{X}} S_y^k f \right)^2$$

**Lemma:**

$$\text{Ent}(\beta f) \leq \frac{\beta^2}{2} E_{\beta f} [D^2(f)]$$

**Proof:** For  $0 < s \leq \beta$ .

$$\begin{aligned} \text{Ent}(\beta f) &\leq E_{\beta f} \left[ \int_0^\beta \int_t^\beta \sum_{k=1}^n \sigma_{k,sf}^2(f) \, ds \, dt \right] \\ &\leq E_{\beta f} \left[ \int_0^\beta \int_t^\beta \sum_{k=1}^n E_{k,\beta f} \left[ \left( f - \inf_{y \in \mathcal{X}} S_y^k f \right)^2 \right] \, ds \, dt \right] \\ &= \frac{\beta^2}{2} E_{\beta f} [D^2(f)] \quad (\text{because } E_{\beta f} [E_{k,\beta f} [\cdot]] = E_{\beta f} [\cdot] ) \quad \square \end{aligned}$$

## Spin-off concentration inequality

**Theorem:** Let  $f : \mathcal{X}^n \rightarrow \mathbb{R}$

$$\begin{aligned} \Pr \{f - Ef > t\} &\leq \inf_{\beta > 0} \exp \left( \beta \int_0^\beta \frac{\text{Ent}(\gamma f)}{\gamma^2} d\gamma - \beta t \right) \\ &\leq \inf_{\beta > 0} \exp \left( \frac{\beta^2 \|D^2(f)\|_\infty}{2} - \beta t \right) \\ &\leq \exp \left( \frac{-t^2}{2 \|D^2(f)\|_\infty} \right) \end{aligned}$$

**Applications:** Concentration of convex Lipschitz functions, shortest T.S.P., largest eigenvalue of random symmetric matrix, and many more ...

## Self-bounded functions

**Lemma:** If  $D^2 f \leq a^2 f$  then for  $\beta \in (0, 2/a^2)$

$$\ln E [e^{\beta f}] \leq \frac{\beta}{1 - a^2 \beta / 2} E [f] .$$

**Proof:**

$$\begin{aligned} \ln E [e^{\beta(f - E[f])}] &= \beta \int_0^\beta \frac{\text{Ent}(\gamma f)}{\gamma^2} d\gamma \\ &\leq \frac{\beta}{2} \int_0^\beta E_{\gamma f} [D^2 f] d\gamma \\ &\leq \frac{a^2 \beta}{2} \int_0^\beta E_{\gamma f} [f] d\gamma = \frac{a^2 \beta}{2} \ln E [e^{\beta f}] \quad \square \end{aligned}$$

The sum of conditional variances is self-bounded, I

**Proposition:**

$$D^2 \left( \Sigma^2 (f) \right) \leq J (f)^2 \Sigma^2 (f)$$

**Proof:** Fix  $x \in \mathcal{X}^n$  and let  $z \in \mathcal{X}^n$ ,  $z_l := \arg \min_z S_z^l \Sigma^2 (f)$ .

$$\begin{aligned} D^2 \left( \Sigma^2 (f) \right) &= \sum_l \left( \Sigma^2 (f) - S_{z_l}^l \Sigma^2 (f) \right)^2 \\ &= \sum_l \left( \sum_{k:k \neq l} \left( \sigma_k^2 (f) - S_{z_l}^l \sigma_k^2 (f) \right) \right)^2. \end{aligned}$$

The sum of conditional variances is self-bounded, II

**Proposition:**

$$D^2 \left( \Sigma^2 (f) \right) \leq J (f)^2 \Sigma^2 (f)$$

**Proof:** Fix  $x \in \mathcal{X}^n$  and let  $z \in \mathcal{X}^n$ ,  $z_l := \arg \min_z S_z^l \Sigma^2 (f)$ .

$$\begin{aligned} 4D^2 \left( \Sigma^2 (f) \right) &= 4 \sum_l \left( \sum_{k:k \neq l} \left( \sigma_k^2 (f) - S_{z_l}^l \sigma_k^2 (f) \right) \right)^2 \\ &= \sum_l \left( \sum_{k \neq l} E_{(y,y') \sim \mu_k^2} \left[ \left( D_{y,y'}^k f \right)^2 - \left( D_{y,y'}^k S_{z_l}^l f \right)^2 \right] \right)^2 \end{aligned}$$

The sum of conditional variances is self-bounded, III

**Proposition:**

$$D^2 \left( \Sigma^2 (f) \right) \leq J (f)^2 \Sigma^2 (f)$$

**Proof:** Fix  $x \in \mathcal{X}^n$  and let  $z \in \mathcal{X}^n$ ,  $z_l := \arg \min_z S_z^l \Sigma^2 (f)$ .

$$\begin{aligned} & 4D^2 \left( \Sigma^2 (f) \right) \\ &= \sum_l \left( \sum_{k \neq l} E_{(y,y') \sim \mu_k^2} \left[ \left( D_{y,y'}^k f \right)^2 - \left( D_{y,y'}^k S_{z_l}^l f \right)^2 \right] \right)^2 \\ &= \sum_l \left( \sum_{k \neq l} E_{(y,y') \sim \mu_k^2} \left[ \left( D_{y,y'}^k f - D_{y,y'}^k S_{z_l}^l f \right) \left( D_{y,y'}^k f + D_{y,y'}^k S_{z_l}^l f \right) \right] \right)^2 \end{aligned}$$



The sum of conditional variances is self-bounded, IV

**Proposition:**

$$D^2 \left( \Sigma^2 (f) \right) \leq J (f)^2 \Sigma^2 (f)$$

**Proof:** Fix  $x \in \mathcal{X}^n$  and let  $z \in \mathcal{X}^n$ ,  $z_l := \arg \min_z S_z^l \Sigma^2 (f)$ .

$$\begin{aligned} & 4D^2 \left( \Sigma^2 (f) \right) \\ &= \sum_l \left( \sum_{k \neq l} E_{(y,y') \sim \mu_k^2} \left[ \left( D_{y,y'}^k f - D_{y,y'}^k S_{z_l}^l f \right) \left( D_{y,y'}^k f + D_{y,y'}^k S_{z_l}^l f \right) \right] \right)^2 \\ &\leq \sum_l \sum_{k:k \neq l} E_{(y,y') \sim \mu_k^2} \left[ D_{y,y'}^k \left( f - S_{z_l}^l f \right) \right]^2 \sum_{k:k \neq l} E_{(y,y') \sim \mu_k^2} \left[ D_{y,y'}^k f + D_{y,y'}^k S_{z_l}^l f \right]^2 \end{aligned}$$

The sum of conditional variances is self-bounded, V

**Proposition:**

$$D^2 \left( \Sigma^2 (f) \right) \leq J (f)^2 \Sigma^2 (f)$$

**Proof:** Fix  $x \in \mathcal{X}^n$  and let  $z \in \mathcal{X}^n$ ,  $z_l := \arg \min_z S_z^l \Sigma^2 (f)$ .

$$\begin{aligned} & 4D^2 \left( \Sigma^2 (f) \right) \\ & \leq \sum_l \sum_{k:k \neq l} E_{(y,y') \sim \mu_k^2} \left[ D_{y,y'}^k \left( f - S_{z_l}^l f \right) \right]^2 \sum_{k:k \neq l} E_{(y,y') \sim \mu_k^2} \left[ D_{y,y'}^k f + D_{y,y'}^k S_{z_l}^l f \right]^2 \\ & \leq 2 \sum_l \sum_{k:k \neq l} \sup_{z,z',y,y'} \left[ D_{z,z'}^l D_{y,y'}^k (f) \right]^2 \left( \Sigma^2 (f) + S_{z_l}^l \Sigma^2 (f) \right) \\ & \leq 4J (f)^2 \Sigma^2 (f) \end{aligned} \quad \square$$

# Proof of Bernstein's inequality, I

Let  $0 < \gamma \leq \beta < 1/(1 + J/2)$ ,

$$\theta := \gamma / (J(1 - \gamma)) \implies \gamma^2 / (2(1 - \gamma)^2) < \theta < 2/J^2.$$

$$\begin{aligned} \theta \operatorname{Ent}(\gamma f) &\leq \frac{\gamma^2}{2(1 - \gamma)^2} E_{\gamma f} [\theta \Sigma^2(f)] && \text{(1st Lemma)} \\ &\leq \frac{\gamma^2}{2(1 - \gamma)^2} \left( \operatorname{Ent}(\gamma f) + \ln E \left[ e^{\theta \Sigma^2(f)} \right] \right) && \text{(Fenchel-Young)} \end{aligned}$$

$$\begin{aligned} \operatorname{Ent}(\gamma f) \left( \theta - \frac{\gamma^2}{2(1 - \gamma)^2} \right) &\leq \frac{\gamma^2}{2(1 - \gamma)^2} \ln E \left[ e^{\theta \Sigma^2(f)} \right] \\ \operatorname{Ent}(\gamma f) &\leq \frac{\gamma J}{2(1 - (1 + (J/2))\gamma)} \ln E \left[ e^{\theta \Sigma^2(f)} \right] \end{aligned}$$

## Proof of Bernstein's inequality, II

Let  $0 < \gamma \leq \beta < 1/(1 + J/2)$ ,

$$\theta := \gamma / (J(1 - \gamma)) \implies \gamma^2 / (2(1 - \gamma)^2) < \theta < 2/J^2.$$

$$\begin{aligned} \text{Ent}(\gamma f) &\leq \frac{\gamma J}{2(1 - (1 + (J/2))\gamma)} \ln E \left[ e^{\theta \Sigma^2(f)} \right] \\ \ln E \left[ e^{\theta \Sigma^2(f)} \right] &\leq \frac{\theta}{1 - J^2\theta/2} E \left[ \Sigma^2(f) \right] \quad (\text{self bounded } \Sigma^2(f)) \\ &= \frac{\gamma/J}{1 - (1 + J/2)\gamma} E \left[ \Sigma^2(f) \right]. \\ \text{Ent}(\gamma f) &\leq \frac{\gamma^2}{2(1 - (1 + J/2)\gamma)^2} E \left[ \Sigma^2(f) \right] \end{aligned}$$

## Proof of Bernstein's inequality, III

Let  $0 < \gamma \leq \beta < 1/(1 + J/2)$ ,

$$\begin{aligned} \text{Ent}(\gamma f) &\leq \frac{\gamma^2}{2(1 - (1 + J/2)\gamma)^2} E[\Sigma^2(f)] \\ \Pr\{f - Ef > t\} &\leq \inf_{\beta \in (0, 1/(1+J/2))} \exp\left(\beta \int_0^\beta \frac{\text{Ent}(\gamma f)}{\gamma^2} d\gamma - \beta t\right) \\ &\leq \inf_{\beta \in (0, 1/(1+J/2))} \exp\left(\frac{E[\Sigma^2(f)]}{2} \frac{\beta^2}{1 - (1 + J/2)\beta} - \beta t\right) \\ &\leq \exp\left(\frac{-t^2}{2(E[\Sigma^2(f)] + (1 + J/2)t)}\right) \quad \square \end{aligned}$$

# References

- [1] S. Bernstein, Theory of Probability, Moscow, 1927.
- [2] S. Boucheron, G. Lugosi, P. Massart, Concentration Inequalities using the entropy method, *Annals of Probability* 31, Nr 3, 2003
- [3] S. Boucheron, G. Lugosi, P. Massart, On concentration of self-bounding functions, *Electronic Journal of Probability* Vol.14 (2009), Paper no. 64, 1884–1899, 2009
- [4] S. Boucheron, G. Lugosi, P. Massart. Concentration Inequalities, Oxford University Press (2013)

- [5] Efron, B., & Stein, C. (1981). The jackknife estimate of variance. *The Annals of Statistics*, 586-596.
  
- [6] M.Ledoux, *The Concentration of Measure Phenomenon*, AMS Surveys and Monographs 89, 2001.
  
- [7] A.Maurer, Thermodynamics and concentration. *Bernoulli* 18.2 (2012): 434-454.
  
- [8] Maurer, A. (2017). A Bernstein-type inequality for functions of bounded interaction. arXiv preprint arXiv:1701.06191.
  
- [9] C.McDiarmid, Concentration, in *Probabilistic Methods of Algorithmic Discrete Mathematics*, p. 195–248. Springer, Berlin, 1998.

- [10] J.M.Steele, An Efron-Stein inequality for nonsymmetric statistics, *Annals of Statistics* 14:753–758, 1986