

Summary Part 3: Applications of Concentration Inequalities

Tianpei Xie

Jan. 26th., 2023

Contents

1	Applications	2
1.1	U-Statistics	2
1.2	Jackknife Estimation and Bootstrapping	2
1.3	Kernel-Density Estimation	2
1.4	Random Graph	3
1.5	Minimum Weight Spanning Tree	4
1.6	Rademacher Complexity	4
1.7	Dimensionality Reduction	5
2	Self-Bounding Functions	5
2.1	Definitions, Variance Bounds and Concentration	5
2.2	Configuration Function	6
2.3	VC-Dimension and Growth Function	7
2.4	Longest Increasing Subsequence	7
2.5	Weakly Self-Bounding Functions	7
3	Random Matrices	7
3.1	Definitions	7
3.2	Concentration Inequalities of Random Vectors	7
3.3	Concentration of Norm of Gaussian Vectors	7
3.4	Spectral Distribution of Hermitian Matrix: Semi-Circular Law	7
3.5	Largest Eigenvalue of Hermitian Random Matrix	7
4	Empirical Process	7
4.1	Definition	7
4.2	Tail Bounds for Empirical Processes	7
4.3	Uniform Law of Large Numbers	7
4.4	Suprema of Empirical Processes	7
4.5	Covering Number, Packing Number and Metric Entropy	7
4.6	Chaining	7
4.7	VC-Dimension	7
4.8	Variance Bounds	7

1 Applications

1.1 U-Statistics

- **Example (*U-Statistics*)** [Wainwright, 2019]

Let $g : \mathbb{R}^2 \rightarrow \mathbb{R}$ be a **symmetric function** of its arguments. Given an i.i.d. sequence $\{X_k, k \geq 1\}$, of random variables, the quantity

$$U := \frac{1}{\binom{n}{2}} \sum_{j < k} g(X_j, X_k) \quad (1)$$

is known as a **pairwise U-statistic**. For instance, if $g(s, t) = |s - t|$, then U is an *unbiased estimator of the mean absolute pairwise deviation* $\mathbb{E}[|X_1 - X_2|]$. Note that, while U is **not** a sum of independent random variables, the dependence is relatively weak, and this fact can be revealed by a martingale analysis.

If g is bounded (say $\|g\|_\infty \leq b$), then the *Bounded Difference Inequality* can be used to establish the concentration of U around its mean. Viewing U as a function $f(x) = f(x_1, \dots, x_n)$, for any given coordinate k , we have

$$\begin{aligned} \left| f(x) - f(x^{(-k)}) \right| &\leq \frac{1}{\binom{n}{2}} \sum_{j \neq k} |g(x_j, x_k) - g(x_j, x'_k)| \\ &\leq \frac{(n-1)2b}{\binom{n}{2}} = \frac{4b}{n}, \end{aligned}$$

so that the bounded differences property holds with parameter $L_k = \frac{4b}{n}$ in each coordinate. Thus, we conclude that

$$\mathbb{P}\{|U - \mathbb{E}[U]| \geq t\} \leq 2 \exp\left(-\frac{nt^2}{8b^2}\right),$$

This tail inequality implies that U is a consistent estimate of $\mathbb{E}[U]$, and also yields *finite sample bounds* on its quality as an estimator. Similar techniques can be used to obtain *tail bounds on U-statistics of higher order*, involving sums over k -tuples of variables. ■

1.2 Jackknife Estimation and Bootstrapping

1.3 Kernel-Density Estimation

- **Example (*Kernel Density Estimation*)**

Let Z_1, \dots, Z_n be i.i.d. samples drawn according to some (unknown) density ϕ on the real line. The density is estimated by the kernel estimate

$$\phi_n(z) = \frac{1}{n h_n} \sum_{i=1}^n K\left(\frac{z - Z_i}{h_n}\right),$$

where $h_n > 0$ is a *smoothing parameter*, and K is a nonnegative function with $\int K(z) dz = 1$. The performance of the estimate is typically measured by **the L_1 error**:

$$X(n) := f(Z_1, \dots, Z_n) = \int |\phi(z) - \phi_n(z)| dz.$$

It is easy to see that

$$\begin{aligned} |f(z_1, \dots, z_n) - f_i(z_1, \dots, z_{i-1}, z'_i, z_{i+1}, \dots, z_n)| &\leq \frac{1}{nh_n} \int \left| K\left(\frac{z - z_i}{h_n}\right) - K\left(\frac{z - z'_i}{h_n}\right) \right| dz \\ &\leq \frac{2}{n}, \end{aligned}$$

so without further work we obtain

$$\text{Var}(X(n)) \leq \frac{1}{n}$$

It is known that for every ϕ , $\sqrt{n}\mathbb{E}[X(n)] \rightarrow \infty$, which implies, by *Chebyshev's inequality*, that for every $\epsilon > 0$

$$\mathbb{P}\left\{\left|\frac{X(n)}{\mathbb{E}[X(n)]} - 1\right| > \epsilon\right\} = \mathbb{P}\{|X(n) - \mathbb{E}[X(n)]| > \epsilon\mathbb{E}[X(n)]\} \leq \frac{\text{Var}(X(n))}{\epsilon^2(\mathbb{E}[X(n)])^2} \rightarrow 0$$

as $n \rightarrow \infty$. That is, $\frac{X(n)}{\mathbb{E}[X(n)]} \rightarrow 1$ **in probability**, or in other words, $X(n)$ is **relatively stable**. This means that **the random L_1 -error** essentially *behaves like its expected value*.

By bounded difference inequality, we have

$$\mathbb{P}\{|X(n) - \mathbb{E}[X(n)]| \geq t\} \leq 2 \exp\left(-\frac{nt^2}{2}\right) \quad \blacksquare$$

1.4 Random Graph

- **Example (*Clique Number in Erdős-Rényi Random Graphs*)** [Wainwright, 2019]
Let $\mathcal{G} = (\mathcal{V}, \mathcal{E})$ be an *undirected graph*, where $\mathcal{V} = \{1, \dots, d\}$ is the vertex set and $\mathcal{E} = \{(i, j), i, j \in \mathcal{V}\}$ is the undirected edge set. A **graph clique** C is a subset of vertices such that $(i, j) \in \mathcal{E}$ for all $i, j \in C$. **The clique number** $C(\mathcal{G})$ of the graph is **the cardinality of the largest clique**. Note that $C(\mathcal{G}) \in [1, d]$. When the edges \mathcal{E} of the graph are drawn according to some random process, then *the clique number* $C(\mathcal{G})$ is a *random variable*, and we can study its concentration around its mean $\mathbb{E}[C(\mathcal{G})]$.

The Erdős-Rényi ensemble of random graphs is one of the most well-studied models: it is defined by a parameter $p \in (0, 1)$ that specifies the probability with which each edge (i, j) is *included* in the graph, **independently across all** $\binom{d}{2}$ edges. More formally, for each $i < j$, let us introduce a **Bernoulli edge-indicator variable** $X_{i,j}$ with parameter p , where $X_{i,j} = 1$ means that edge (i, j) is *included* in the graph, and $X_{i,j} = 0$ means that it is *not included*.

Note that the $\binom{d}{2}$ -dimensional random vector $Z := \{X_{i,j}\}_{i < j}$ specifies the edge set; thus, we may view *the clique number* $C(\mathcal{G})$ as a function $Z \rightarrow f(Z)$. Based on definition in Section ??, we see that $f(Z)$ is a **configuration function** with property of “*being in a clique*”.

Let Z' denote a vector in which a *single coordinate* of Z has been changed, and let \mathcal{G}' and \mathcal{G} be the associated graphs. It is easy to see that $C(\mathcal{G}')$ can differ from $C(\mathcal{G})$ by at most 1, so that

$$|f(Z) - f(Z')| \leq 1,$$

Thus, the function $C(\mathcal{G}) = f(Z)$ satisfies *the bounded difference property* in each coordinate with parameter $L = 1$, so that

$$\mathbb{P} \left\{ \frac{1}{n} |C(\mathcal{G}) - \mathbb{E}[C(\mathcal{G})]| \geq \delta \right\} \leq 2 \exp(-2n\delta^2).$$

Consequently, we see that the clique number of an *Erdős-Rényi random graph* is *very sharply concentrated around its expectation*. ■

1.5 Minimum Weight Spanning Tree

1.6 Rademacher Complexity

- **Example (*Rademacher Complexity*)** [Wainwright, 2019]

Let $\{\epsilon_k\}_{k=1}^n$ be an i.i.d. sequence of *Rademacher variables* (i.e., taking the values $\{-1, +1\}$ *equiprobably*). Given a collection of vectors $\mathcal{A} \subset \mathbb{R}^n$, define the random variable

$$Z := \sup_{a \in \mathcal{A}} \sum_{k=1}^n \epsilon_k a_k = \sup_{a \in \mathcal{A}} \langle a, \epsilon \rangle. \quad (2)$$

The random variable Z measures the **size** of \mathcal{A} in a certain sense, and its expectation

$$\mathfrak{R}(\mathcal{A}) := \mathbb{E}[Z(\mathcal{A})] \quad (3)$$

is known as **the Rademacher complexity** of the set \mathcal{A} .

Let us now show how the bounded difference inequality can be used to establish that $Z(\mathcal{A})$ is **sub-Gaussian**. Viewing $Z(\mathcal{A})$ as a function $(\epsilon_1, \dots, \epsilon_n) \rightarrow f(\epsilon_1, \dots, \epsilon_n)$, we need to *bound the maximum change* when coordinate k is changed. Given two Rademacher vectors $\epsilon, \epsilon' \in \{-1, +1\}^n$, recall our definition of the modified vector $\epsilon^{(-k)}$. Since

$$f(\epsilon^{(-k)}) \geq \langle a, \epsilon^{(-k)} \rangle, \quad \text{for any } a \in \mathcal{A},$$

we have

$$\langle a, \epsilon \rangle - f(\epsilon^{(-k)}) \leq \langle a, \epsilon - \epsilon^{(-k)} \rangle = a_k(\epsilon_k - \epsilon'_k) \leq 2|a_k|.$$

Taking *the supremum over \mathcal{A} on both sides*, we obtain the inequality

$$f(\epsilon) - f(\epsilon^{(-k)}) \leq 2 \sup_{a \in \mathcal{A}} |a_k|.$$

Since the same argument applies with the roles of ϵ and $\epsilon^{(-k)}$ *reversed*, we conclude that f satisfies *the bounded difference inequality* in coordinate k with parameter $L_k := 2 \sup_{a \in \mathcal{A}} |a_k|$.

Consequently, the bounded difference inequality implies that the random variable $Z(\mathcal{A})$ is *sub-Gaussian* with parameter at most $2 \sqrt{\sum_{k=1}^n \sup_{a \in \mathcal{A}} a_k^2}$. This sub-Gaussian parameter can be reduced to the (potentially much) smaller quantity $\sqrt{\sup_{a \in \mathcal{A}} \sum_{k=1}^n a_k^2}$ using alternative techniques. ■

1.7 Dimensionality Reduction

2 Self-Bounding Functions

2.1 Definitions, Variance Bounds and Concentration

- Another simple property which is satisfied for many important examples is the so-called *self-bounding property*.

Definition (*Self-Bounding Property*)

A *nonnegative* function $f : \mathcal{X}^n \rightarrow [0, \infty)$ has the *self-bounding property* if there exist functions $f_i : \mathcal{X}^{n-1} \rightarrow \mathbb{R}$ such that for all $z_1, \dots, z_n \in \mathcal{X}$ and all $i = 1, \dots, n$,

$$0 \leq f(z_1, \dots, z_n) - f_i(z_1, \dots, z_{i-1}, z_{i+1}, \dots, z_n) \leq 1 \quad (4)$$

and also

$$\sum_{i=1}^n (f(z_1, \dots, z_n) - f_i(z_1, \dots, z_{i-1}, z_{i+1}, \dots, z_n)) \leq f(z_1, \dots, z_n). \quad (5)$$

- **Remark** Clearly if f has the *self-bounding property*,

$$\sum_{i=1}^n (f(z_1, \dots, z_n) - f_i(z_1, \dots, z_{i-1}, z_{i+1}, \dots, z_n))^2 \leq f(z_1, \dots, z_n) \quad (6)$$

Taking expectation on both sides, we have the following inequality

- **Corollary 2.1** [Boucheron et al., 2013]
If f has the *self-bounding property*, then

$$\text{Var}(f(Z)) \leq \mathbb{E}[f(Z)].$$

- **Remark (*Relative Stability*)** [Boucheron et al., 2013]
A sequence of nonnegative random variables $(Z_n)_{n \in \mathbb{N}}$ is said to be *relatively stable* if

$$\frac{Z_n}{\mathbb{E}[Z_n]} \xrightarrow{\mathbb{P}} 1.$$

This property guarantees that *the random fluctuations of Z_n around its expectation are of negligible size when compared to the expectation*, and therefore *most information about the size of Z_n is given by $\mathbb{E}[Z_n]$* .

Bounding the variance of Z_n by its expected value implies, in many cases, the relative stability of $(Z_n)_{n \in \mathbb{N}}$. If Z_n has the *self-bounding property*, then, by Chebyshev's inequality, for all $\epsilon > 0$,

$$\mathbb{P} \left\{ \left| \frac{Z_n}{\mathbb{E}[Z_n]} - 1 \right| > \epsilon \right\} \leq \frac{\text{Var}(Z_n)}{\epsilon^2 (\mathbb{E}[Z_n])^2} \leq \frac{1}{\epsilon^2 \mathbb{E}[Z_n]}.$$

Thus, for relative stability, it suffices to have $\mathbb{E}[Z_n] \rightarrow \infty$.

2.2 Configuration Function

- An important class of functions satisfying *the self-bounding property* consists of the so-called **configuration functions**.

Definition (Configuration Function)

Assume that we have a property Π *defined over the union of finite products* of a set \mathcal{X} , that is, a sequence of sets

$$\Pi_1 \subset \mathcal{X}, \Pi_2 \subset \mathcal{X} \times \mathcal{X}, \dots, \Pi_n \subset \mathcal{X}^n.$$

We say that $(z_1, \dots, z_m) \in \mathcal{X}^m$ *satisfies the property* Π if $(z_1, \dots, z_m) \in \Pi_m$.

We assume that Π is **hereditary** in the sense that if (z_1, \dots, z_m) satisfies Π then so does *any sub-sequence* $\{z_{i_1}, \dots, z_{i_k}\}$ of (z_1, \dots, z_m) .

The function f that maps any vector $z = (z_1, \dots, z_n)$ to *the size of a largest sub-sequence satisfying* Π is the configuration function associated with property Π .

- **Corollary 2.2** [Boucheron et al., 2013]

Let f be a **configuration function**, and let $X = f(Z_1, \dots, Z_n)$, where Z_1, \dots, Z_n are *independent* random variables. Then

$$\text{Var}(f(Z)) \leq \mathbb{E}[f(Z)].$$

2.3 VC-Dimension and Growth Function

- **Example (VC Dimension)**

Let \mathcal{H} be an arbitrary collection of subsets of \mathcal{X} , and let $x = (x_1, \dots, x_n)$ be a vector of n points of \mathcal{X} . Define the **trace** of \mathcal{H} on x by

$$\text{tr}(x) = \{A \cap \{x_1, \dots, x_n\} : A \in \mathcal{H}\}.$$

The shatter coefficient, (or *Vapnik-Chervonenkis growth function*) of \mathcal{H} in x is $\tau_{\mathcal{H}}(x) = |\text{tr}(x)|$, *the size of the trace*. $\tau_{\mathcal{H}}(x)$ is the number of different subsets of the n -point set $\{x_1, \dots, x_n\}$ generated by intersecting it with elements of \mathcal{H} . A subset $\{x_{i_1}, \dots, x_{i_k}\}$ of $\{x_1, \dots, x_n\}$ is said to be **shattered** if $2^k = T(x_{i_1}, \dots, x_{i_k})$.

The VC dimension $D(x)$ of \mathcal{H} (with respect to x) is the *cardinality* k of the *largest shattered subset* of x . From the definition it is obvious that $f(x) = D(x)$ is a **configuration function** (associated with the property of “**shatteredness**”) and therefore if X_1, \dots, X_n are *independent random variables*, then

$$\text{Var}(D(X)) \leq \mathbb{E}[D(X)].$$

2.4 Longest Increasing Subsequence

2.5 Weakly Self-Bounding Functions

3 Random Matrices

3.1 Definitions

3.2 Concentration Inequalities of Random Vectors

3.3 Concentration of Norm of Gaussian Vectors

3.4 Spectral Distribution of Hermitian Matrix: Semi-Circular Law

3.5 Largest Eigenvalue of Hermitian Random Matrix

4 Empirical Process

4.1 Definition

4.2 Tail Bounds for Empirical Processes

4.3 Uniform Law of Large Numbers

4.4 Suprema of Empirical Processes

4.5 Covering Number, Packing Number and Metric Entropy

4.6 Chaining

4.7 VC-Dimension

4.8 Variance Bounds

References

- Stéphane Boucheron, Gábor Lugosi, and Pascal Massart. *Concentration inequalities: A nonasymptotic theory of independence*. Oxford university press, 2013.
- Martin J Wainwright. *High-dimensional statistics: A non-asymptotic viewpoint*, volume 48. Cambridge University Press, 2019.