

# Case study: stochastic simulation via Rademacher bootstrap

Maxim Raginsky

December 4, 2013

In this lecture, we will look at an application of statistical learning theory to the problem of efficient stochastic simulation, which arises frequently in engineering design. The basic question is as follows. Suppose we have a system with input space  $Z$ . The system has a tunable parameter  $\theta$  that lies in some set  $\Theta$ . We have a *performance index*  $\ell : Z \times \Theta \rightarrow [0, 1]$ , where we assume that the lower the value of  $\ell$ , the better the performance. Thus, if we use the parameter setting  $\theta \in \Theta$  and apply input  $z \in Z$ , the performance of the corresponding system is given by the scalar  $\ell(z, \theta) \in [0, 1]$ . Now let's suppose that the input to the system is actually a *random variable*  $Z \in Z$  with some distribution  $P \in \mathcal{P}(Z)$ . Then we can define the *operating characteristic*

$$L(\theta) \triangleq \mathbb{E}_P[\ell(Z, \theta)] \equiv \int_Z \ell(z, \theta) P_Z(dz), \quad \theta \in \Theta. \quad (1)$$

The goal is to find an *optimal operating point*  $\theta^* \in \Theta$  that achieves (or comes arbitrarily close to)  $\inf_{\theta \in \Theta} L(\theta)$ .

In practice, the problem of minimizing  $L(\theta)$  is quite difficult for large-scale systems. First of all, computing the integral in (1) may be a challenge. Secondly, we may not even know the distribution  $P_Z$ . Thirdly, there may be more than one distribution of the input, each corresponding to different operating regimes and/or environments. For this reason, engineers often resort to Monte Carlo simulation techniques: Assuming we can efficiently sample from  $P_Z$ , we draw a large number of independent samples  $Z_1, Z_2, \dots, Z_n$  and compute

$$\widehat{\theta}_n = \underset{\theta \in \Theta}{\operatorname{argmin}} L_n(\theta) \equiv \underset{\theta \in \Theta}{\operatorname{argmin}} \frac{1}{n} \sum_{i=1}^n \ell(Z_i, \theta),$$

where  $L_n(\cdot)$  denotes the empirical version of the operating characteristic (1). Given an accuracy parameter  $\varepsilon > 0$  and a confidence parameter  $\delta \in (0, 1)$ , we simply need to draw enough samples, so that

$$L(\widehat{\theta}_n) \leq \inf_{\theta \in \Theta} L(\theta) + \varepsilon$$

with probability at least  $1 - \delta$ , regardless of what the true distribution  $P_Z$  happens to be.

This is, of course, just another instance of the ERM algorithm we have been studying extensively. However, there are two issues. One is how many samples we need to guarantee that

the empirically optimal operating point will be good. The other is the complexity of actually computing an empirical minimizer.

The first issue has already come up in the course under the name of *sample complexity* of learning. The second issue is often handled by relaxing the problem a bit: We choose a probability distribution  $Q$  over  $\Theta$  (assuming it can be equipped with an appropriate  $\sigma$ -algebra) and, instead of minimizing  $L(\theta)$  over  $\theta \in \Theta$ , set some *level parameter*  $\alpha \in (0, 1)$ , and seek any  $\hat{\theta} \in \Theta$ , for which there exists some *exceptional set*  $\Lambda \subset \Theta$  with  $Q(\Lambda) \leq \alpha$ , such that

$$\inf_{\theta} L(\theta) - \varepsilon \leq L(\hat{\theta}) \leq \inf_{\theta \in \Theta \setminus \Lambda} L(\theta) + \varepsilon \quad (2)$$

with probability at least  $1 - \delta$ . Unless the actual optimal operating point  $\theta^*$  happens to lie in the exceptional set  $\Lambda$ , we will come to within  $\varepsilon$  of the optimum with confidence at least  $1 - \delta$ . Then we just need to draw a large enough number  $n$  of samples  $Z_1, \dots, Z_n$  from  $P_Z$  and a large enough number  $m$  of samples  $\theta_1, \dots, \theta_m$  from  $Q$ , and then compute

$$\hat{\theta} = \operatorname{argmin}_{\theta \in \{\theta_1, \dots, \theta_m\}} L_n(\theta).$$

In the next several lectures, we will see how statistical learning theory can be used to develop such simulation procedures. Moreover, we will learn how to use Rademacher averages<sup>1</sup> to determine how many samples we need *in the process of learning*. The use of statistical learning theory for simulation has been pioneered in the context of control by M. Vidyasagar [Vid98, Vid01]; the refinement of his techniques using Rademacher averages is due to Koltchinskii et al. [KAA<sup>+</sup>00a, KAA<sup>+</sup>00b]. We will essentially follow their presentation, but with slightly better constants.

We will follow the following plan. First, we will revisit the abstract ERM problem and its sample complexity. Then we will introduce a couple of refined tools pertaining to Rademacher averages. Next, we will look at *sequential* algorithms for empirical approximation, in which the sample complexity is not set *a priori*, but is rather determined by a data-driven *stopping rule*. And, finally, we will see how these sequential algorithms can be used to develop robust and efficient stochastic simulation strategies.

## 1 Empirical Risk Minimization: a quick review

Recall the *abstract Empirical Risk Minimization problem*: We have a space  $Z$ , a class  $\mathcal{P}$  of probability distributions over  $Z$ , and a class  $\mathcal{F}$  of measurable functions  $f : Z \rightarrow [0, 1]$ . Given an i.i.d. sample  $Z^n$  drawn according to some unknown  $P \in \mathcal{P}$ , we compute

$$\hat{f}_n \triangleq \operatorname{argmin}_{f \in \mathcal{F}} P_n(f) \equiv \operatorname{argmin}_{f \in \mathcal{F}} \frac{1}{n} \sum_{i=1}^n f(Z_i).$$

---

<sup>1</sup>More precisely, their stochastic counterpart, in which we do not take the expectation over the Rademacher sequence, but rather use it as a *resource* to aid the simulation.

We would like for  $P(\hat{f}_n)$  to be close to  $\inf_{f \in \mathcal{F}} P(f)$  with high probability. To that end, we have derived the bound

$$P(\hat{f}_n) - \inf_{f \in \mathcal{F}} P(f) \leq 2\|P_n - P\|_{\mathcal{F}},$$

where, as before, we have defined the *uniform deviation*

$$\|P_n - P\|_{\mathcal{F}} \triangleq \sup_{f \in \mathcal{F}} |P_n(f) - P(f)| = \sup_{f \in \mathcal{F}} \left| \frac{1}{n} \sum_{i=1}^n f(Z_i) - \mathbb{E}_P f(Z) \right|.$$

Hence, if  $n$  is sufficiently large so that, for every  $P \in \mathcal{P}$ ,  $\|P_n - P\|_{\mathcal{F}} \leq \varepsilon/2$  with  $P$ -probability at least  $1 - \delta$ , then  $P(\hat{f}_n)$  will be  $\varepsilon$ -close to  $\inf_{f \in \mathcal{F}} P(f)$  with probability at least  $1 - \delta$ . This motivates the following definition:

**Definition 1.** Given the pair  $(\mathcal{F}, \mathcal{P})$ , an accuracy parameter  $\varepsilon > 0$ , and a confidence parameter  $\delta \in (0, 1)$ , the sample complexity of empirical approximation is

$$N(\varepsilon; \delta) \triangleq \min \left\{ n \in \mathbb{N} : \sup_{P \in \mathcal{P}} \mathbb{P} \{ \|P_n - P\|_{\mathcal{F}} \geq \varepsilon \} \leq \delta \right\}. \quad (3)$$

In other words, for any  $\varepsilon > 0$  and any  $\delta \in (0, 1)$ ,  $N(\varepsilon/2; \delta)$  is an *upper bound* on the number of samples needed to guarantee that  $P(\hat{f}_n) \leq \inf_{f \in \mathcal{F}} P(f) + \varepsilon$  with probability (confidence) at least  $1 - \delta$ .

## 2 Empirical Rademacher averages

As before, let  $Z^n$  be an i.i.d. sample of length  $n$  from some  $P \in \mathcal{P}(Z)$ . On multiple occasions we have seen that the performance of the ERM algorithm is controlled by the *Rademacher average*

$$R_n(\mathcal{F}(Z^n)) \triangleq \frac{1}{n} \mathbb{E}_{\sigma^n} \left[ \sup_{f \in \mathcal{F}} \left| \sum_{i=1}^n \sigma_i f(Z_i) \right| \right], \quad (4)$$

where  $\sigma^n = (\sigma_1, \dots, \sigma_n)$  is an  $n$ -tuple of i.i.d. Rademacher random variables independent of  $Z^n$ . More precisely, we have established the fundamental *symmetrization inequality*

$$\mathbb{E} \|P_n - P\|_{\mathcal{F}} \leq 2\mathbb{E} R_n(\mathcal{F}(Z^n)), \quad (5)$$

as well as the concentration bounds

$$\mathbb{P} \{ \|P_n - P\|_{\mathcal{F}} \geq \mathbb{E} \|P_n - P\|_{\mathcal{F}} + \varepsilon \} \leq e^{-2n\varepsilon^2} \quad (6)$$

$$\mathbb{P} \{ \|P_n - P\|_{\mathcal{F}} \leq \mathbb{E} \|P_n - P\|_{\mathcal{F}} - \varepsilon \} \leq e^{-2n\varepsilon^2} \quad (7)$$

These results show two things:

1. The uniform deviation  $\|P_n - P\|_{\mathcal{F}}$  tightly concentrates around its expected value.

2. The expected value  $\mathbb{E}\|P_n - P\|_{\mathcal{F}}$  is bounded from above by  $\mathbb{E}R_n(\mathcal{F}(Z^n))$ .

It turns out that the expected Rademacher average  $\mathbb{E}R_n(\mathcal{F}(Z^n))$  also furnishes a *lower bound* on  $\mathbb{E}\|P_n - P\|_{\mathcal{F}}$ :

**Lemma 1** (Desymmetrization inequality). *For any class  $\mathcal{F}$  of measurable functions  $f : Z \rightarrow [0, 1]$ , we have*

$$\frac{1}{2}\mathbb{E}R_n(\mathcal{F}(Z^n)) - \frac{1}{2\sqrt{n}} \leq \frac{1}{2n} \mathbb{E} \left[ \sup_{f \in \mathcal{F}} \left| \sum_{i=1}^n \sigma_i [f(Z_i) - P(f)] \right| \right] \leq \mathbb{E}\|P_n - P\|_{\mathcal{F}}. \quad (8)$$

*Proof.* We will first prove the second inequality in (8). To that end, for each  $1 \leq i \leq n$  and each  $f \in \mathcal{F}$ , let us define  $U_i(f) \triangleq f(Z_i) - P(f)$ . Then  $\mathbb{E}U_i(f) = 0$ . Let  $\bar{Z}_1, \dots, \bar{Z}_n$  be an independent copy of  $Z_1, \dots, Z_n$ . Then we can define  $\bar{U}_i(f), 1 \leq i \leq n$ , similarly. Moreover, since  $\mathbb{E}U_i(f) = 0$ , we can write

$$\begin{aligned} \mathbb{E} \left[ \sup_{f \in \mathcal{F}} \left| \sum_{i=1}^n \sigma_i [f(Z_i) - P(f)] \right| \right] &= \mathbb{E} \left[ \sup_{f \in \mathcal{F}} \left| \sum_{i=1}^n \sigma_i U_i(f) \right| \right] \\ &= \mathbb{E} \left[ \sup_{f \in \mathcal{F}} \left| \sum_{i=1}^n \sigma_i [U_i(f) - \mathbb{E}\bar{U}_i(f)] \right| \right] \\ &\leq \mathbb{E} \left[ \sup_{f \in \mathcal{F}} \left| \sum_{i=1}^n \sigma_i [U_i(f) - \bar{U}_i(f)] \right| \right]. \end{aligned}$$

Since, for each  $i$ ,  $U_i(f)$  and  $\bar{U}_i(f)$  are i.i.d., the difference  $U_i(f) - \bar{U}_i(f)$  is a symmetric random variable. Therefore,

$$\left\{ \sigma_i [U_i(f) - \bar{U}_i(f)] : 1 \leq i \leq n \right\} \stackrel{(d)}{=} \left\{ U_i(f) - \bar{U}_i(f) : 1 \leq i \leq n \right\}.$$

Using this fact and the triangle inequality, we get

$$\begin{aligned} \mathbb{E} \left[ \sup_{f \in \mathcal{F}} \left| \sum_{i=1}^n \sigma_i [U_i(f) - \bar{U}_i(f)] \right| \right] &= \mathbb{E} \left[ \sup_{f \in \mathcal{F}} \left| \sum_{i=1}^n [U_i(f) - \bar{U}_i(f)] \right| \right] \\ &\leq 2 \mathbb{E} \left[ \sup_{f \in \mathcal{F}} \left| \sum_{i=1}^n U_i(f) \right| \right] \\ &= 2 \mathbb{E} \left[ \sup_{f \in \mathcal{F}} \left| \sum_{i=1}^n [f(Z_i) - P(f)] \right| \right] \\ &= 2n \cdot \mathbb{E}\|P_n - P\|_{\mathcal{F}}. \end{aligned}$$

To prove the first inequality in (8), we write

$$\begin{aligned}
\mathbb{E}R_n(\mathcal{F}(Z^n)) &= \frac{1}{n} \mathbb{E} \left[ \sup_{f \in \mathcal{F}} \left| \sum_{i=1}^n \sigma_i [f(Z_i) - P(f) + P(f)] \right| \right] \\
&\leq \frac{1}{n} \mathbb{E} \left[ \sup_{f \in \mathcal{F}} \left| \sum_{i=1}^n \sigma_i [f(Z_i) - P(f)] \right| \right] + \frac{1}{n} \mathbb{E} \left[ \sup_{f \in \mathcal{F}} P(f) \cdot \left| \sum_{i=1}^n \sigma_i \right| \right] \\
&= \frac{1}{n} \mathbb{E} \left[ \sup_{f \in \mathcal{F}} \left| \sum_{i=1}^n \sigma_i [f(Z_i) - P(f)] \right| \right] + \frac{1}{n} \mathbb{E} \left[ \sum_{i=1}^n \sigma_i \right] \\
&\leq \frac{1}{n} \mathbb{E} \left[ \sup_{f \in \mathcal{F}} \left| \sum_{i=1}^n \sigma_i [f(Z_i) - P(f)] \right| \right] + \frac{1}{\sqrt{n}}.
\end{aligned}$$

Rearranging, we get the desired inequality.  $\square$

In this section, we will see that we can get a lot of mileage out of the *stochastic version* of the Rademacher average. To that end, let us define

$$r_n(\mathcal{F}(Z^n)) \triangleq \frac{1}{n} \sup_{f \in \mathcal{F}} \left| \sum_{i=1}^n \sigma_i f(Z_i) \right|. \quad (9)$$

The key difference between (4) and (9) is that, in the latter, we do *not* take the expectation over the Rademacher sequence  $\sigma^n$ . In other words, both  $R_n(\mathcal{F}(Z^n))$  and  $r_n(\mathcal{F}(Z^n))$  are random variables, but the former depends only on the training data  $Z^n$ , while the latter also depends on the  $n$  Rademacher random variables  $\sigma_1, \dots, \sigma_n$ . We see immediately that  $R_n(\mathcal{F}(Z^n)) = \mathbb{E}[r_n(\mathcal{F}(Z^n)) | Z^n]$  and  $\mathbb{E}R_n(\mathcal{F}(Z^n)) = \mathbb{E}r_n(\mathcal{F}(Z^n))$ , where the expectation on the right-hand side is over both  $Z^n$  and  $\sigma^n$ . The following result will be useful:

**Lemma 2** (Concentration inequalities for Rademacher averages). *For any  $\varepsilon > 0$ ,*

$$\mathbb{P} \{ r_n(\mathcal{F}(Z^n)) \geq \mathbb{E}R_n(\mathcal{F}(Z^n)) + \varepsilon \} \leq e^{-n\varepsilon^2/2} \quad (10)$$

and

$$\mathbb{P} \{ r_n(\mathcal{F}(Z^n)) \leq \mathbb{E}R_n(\mathcal{F}(Z^n)) - \varepsilon \} \leq e^{-n\varepsilon^2/2}. \quad (11)$$

*Proof.* For each  $1 \leq i \leq n$ , let  $U_i \triangleq (Z_i, \sigma_i)$ . Then  $r_n(\mathcal{F}(Z^n))$  can be represented as a real-valued function  $g(U^n)$ . Moreover, it is easy to see that this function has bounded differences with  $c_1 = \dots = c_n = 2/n$ . Hence, McDiarmid's inequality tells us that for any  $\varepsilon > 0$

$$\mathbb{P} \{ g(U^n) \geq \mathbb{E}g(U^n) + \varepsilon \} \leq e^{-n\varepsilon^2/2},$$

and the same holds for the probability that  $g(U^n) \leq \mathbb{E}g(U^n) - \varepsilon$ . This completes the proof.  $\square$

### 3 Sequential learning algorithms

In a sequential learning algorithm, the sample complexity is a *random variable*. It is not known in advance, but rather is computed from data in the process of learning. In other words, instead of using a training sequence of fixed length, we keep drawing independent samples until we decide that we have acquired enough of them, and then compute an empirical risk minimizer.

To formalize this idea, we need the notion of a *stopping time*. Let  $\{U_n\}_{n=1}^\infty$  be a random process. A random variable  $\tau$  taking values in  $\mathbb{N}$  is called a *stopping time* if and only if, for each  $n \geq 1$ , the occurrence of the event  $\{\tau = n\}$  is determined by  $U^n = (U_1, \dots, U_n)$ . More precisely:

**Definition 2.** For each  $n$ , let  $\Sigma_n$  denote the  $\sigma$ -algebra generated by  $U^n$  (in other words,  $\Sigma_n$  consists of all events that occur by time  $n$ ). Then a random variable  $\tau$  taking values in  $\mathbb{N}$  is a *stopping time* if and only if, for each  $n \geq 1$ , the event  $\{\tau = n\} \in \Sigma_n$ .

In other words, denoting by  $U^\infty$  the entire sample path  $(U_1, U_2, \dots)$  of our random process, we can view  $\tau$  as a function that maps  $U^\infty$  into  $\mathbb{N}$ . For each  $n$ , the indicator function of the event  $\{\tau = n\}$  is a function of  $U^\infty$ :

$$\mathbf{1}_{\{\tau=n\}} \equiv \mathbf{1}_{\{\tau(U^\infty)=n\}}.$$

Then  $\tau$  is a stopping time if and only if, for each  $n$  and for all  $U^\infty, V^\infty$  with  $U^n = V^n$  we have

$$\mathbf{1}_{\{\tau(U^\infty)=n\}} = \mathbf{1}_{\{\tau(V^\infty)=n\}}.$$

Our sequential learning algorithms will work as follows. Given a desired accuracy parameter  $\varepsilon > 0$  and a confidence parameter  $\delta > 0$ , let  $\bar{n}(\varepsilon, \delta)$  be the initial sample size; we will assume that  $\bar{n}(\varepsilon, \delta)$  is a nonincreasing function of both  $\varepsilon$  and  $\delta$ . Let  $\mathcal{T}(\varepsilon, \delta)$  denote the set of all stopping times  $\tau$  such that

$$\sup_{P \in \mathcal{P}} \mathbb{P} \{ \|P_\tau - P\|_{\mathcal{F}} \leq \varepsilon \} \geq \delta.$$

Now if  $\tau \in \mathcal{T}(\varepsilon, \delta)$  and we let

$$\hat{f}_\tau \triangleq \operatorname{argmin}_{f \in \mathcal{F}} P_\tau(f) \equiv \operatorname{argmin}_{f \in \mathcal{F}} \frac{1}{\tau} \sum_{i=1}^{\tau} f(Z_i),$$

then we immediately see that

$$\sup_{P \in \mathcal{P}} \left\{ P(\hat{f}_\tau) \geq \inf_{f \in \mathcal{F}} P(f) + 2\varepsilon \right\} \leq \delta.$$

Of course, the whole question is how to construct an appropriate stopping time *without knowing*  $P$ .

**Definition 3.** A parametric family of stopping times  $\{v(\varepsilon, \delta) : \varepsilon > 0, \delta \in (0, 1)\}$  is called *strongly efficient (SE)* (w.r.t.  $\mathcal{F}$  and  $\mathcal{P}$ ) if there exist constants  $K_1, K_2, K_3 \geq 1$ , such that for all  $\varepsilon > 0, \delta \in (0, 1)$

$$v(\varepsilon, \delta) \in \mathcal{T}(K_1\varepsilon, \delta) \tag{12}$$

and for all  $\tau \in \mathcal{T}(\varepsilon, \delta)$

$$\sup_{P \in \mathcal{P}} \mathbb{P} \{v(K_2\varepsilon, \delta) > \tau\} \leq K_3\delta. \tag{13}$$

In other words, Eq. (12) says that any SE stopping time  $\{v(\varepsilon, \delta)\}$  guarantees that we can approximate statistical expectations by empirical expectations with accuracy  $K_1\varepsilon$  and confidence  $1 - \delta$ ; similarly, Eq. (13) says that, with probability at least  $1 - K_3\delta$ , we will require at most as many samples as would be needed by *any* sequential algorithm for empirical approximation with accuracy  $\varepsilon/K_2$  and confidence  $1 - \delta$ .

**Definition 4.** A family of stopping times  $\{v(\varepsilon, \delta) : \varepsilon > 0, \delta \in (0, 1)\}$  is weakly efficient (WE) for  $(\mathcal{F}, \mathcal{P})$  if there exist constants  $K_1, K_2, K_3 \geq 1$ , such that for all  $\varepsilon > 0, \delta \in (0, 1)$

$$v(\varepsilon, \delta) \in \mathcal{T}(K_1\varepsilon, \delta) \quad (14)$$

and

$$\sup_{P \in \mathcal{P}} \mathbb{P}\{v(K_2\varepsilon, \delta) > N(\varepsilon; \delta)\} \leq K_3\delta. \quad (15)$$

If  $v(\varepsilon, \delta)$  is a WE stopping time, then Eq. (14) says that we can solve the empirical approximation problem with accuracy  $K_1\varepsilon$  and confidence  $1 - \delta$ ; Eq. (15) says that, with probability at most  $1 - \delta$ , the sample complexity will be less than the sample complexity of empirical approximation with accuracy  $\varepsilon/K_2$  and confidence  $1 - \delta$ .

If  $N(\varepsilon; \delta) \geq \bar{n}(\varepsilon, \delta)$ , then  $N(\varepsilon, \delta) \in \mathcal{T}(\varepsilon, \delta)$ . Hence, any WE stopping time is also SE. The converse, however, is not true.

### 3.1 A strongly efficient sequential learning algorithm

Let  $\{Z_n\}_{n=1}^\infty$  be an infinite sequence of i.i.d. draws from some  $P \in \mathcal{P}$ ; let  $\{\sigma_n\}_{n=1}^\infty$  be an i.i.d. Rademacher sequence independent of  $\{Z_n\}$ . Choose

$$\bar{n}(\varepsilon, \delta) \geq \left\lceil \frac{2}{\varepsilon^2} \log \frac{2}{\delta(1 - e^{-\varepsilon^2/2})} \right\rceil + 1 \quad (16)$$

and let

$$v(\varepsilon, \delta) \triangleq \min\{n \geq \bar{n}(\varepsilon, \delta) : r_n(\mathcal{F}(Z^n)) \leq \varepsilon\}. \quad (17)$$

This is clearly a stopping time for each  $\varepsilon > 0$  and each  $\delta \in (0, 1)$ .

**Theorem 1.** The family  $\{v(\varepsilon, \delta) : \varepsilon > 0, \delta \in (0, 1)\}$  defined in (17) with  $\bar{n}(\varepsilon, \delta)$  set according to (16) is SE for any class  $\mathcal{F}$  of measurable functions  $f : Z \rightarrow [0, 1]$  and  $\mathcal{P} = \mathcal{P}(Z)$  with  $K_1 = 5, K_2 = 6, K_3 = 1$ .

*Proof.* Let  $\bar{n} = \bar{n}(\varepsilon, \delta)$ . We will first show that, for any  $P \in \mathcal{P}(Z)$ ,

$$\|P_n - P\|_{\mathcal{F}} \leq 2r_n(\mathcal{F}(Z^n)) + 3\varepsilon, \quad \forall n \geq \bar{n} \quad (18)$$

with probability at least  $1 - \delta$ . Since for  $n = v(\varepsilon, \delta) \geq \bar{n}$  we have  $r_n(\mathcal{F}(Z^n)) \leq \varepsilon$ , we will immediately be able to conclude that

$$\mathbb{P}\{\|P_{v(\varepsilon, \delta)} - P\|_{\mathcal{F}} \geq 5\varepsilon\} \leq \delta,$$

which will imply that  $v(\varepsilon, \delta) \in \mathcal{T}(5\varepsilon, \delta)$ . Now we prove (18). First of all, applying Lemma 2 and the union bound, we can write

$$\begin{aligned} \mathbb{P} \left\{ \bigcup_{n \geq \bar{n}} \{r_n(\mathcal{F}(Z^n)) \geq \mathbb{E}R_n(\mathcal{F}(Z^n)) + \varepsilon\} \right\} &\leq \sum_{n \geq \bar{n}} e^{-n\varepsilon^2/2} \\ &= e^{-\bar{n}\varepsilon^2/2} \sum_{n \geq 0} e^{-n\varepsilon^2/2} \\ &= \frac{e^{-\bar{n}\varepsilon^2/2}}{1 - e^{-\varepsilon^2/2}} \\ &\leq \delta/2. \end{aligned}$$

From the symmetrization inequality (5), we know that  $\mathbb{E}\|P_n - P\|_{\mathcal{F}} \leq 2\mathbb{E}R_n(\mathcal{F}(Z^n))$ . Moreover, using (6) and the union bound, we can write

$$\begin{aligned} \mathbb{P} \left\{ \bigcup_{n \geq \bar{n}} \{\|P_n - P\|_{\mathcal{F}} \geq \mathbb{E}\|P_n - P\|_{\mathcal{F}} + \varepsilon\} \right\} &\leq \sum_{n \geq \bar{n}} e^{-2n\varepsilon^2} \\ &\leq \sum_{n \geq \bar{n}} e^{-n\varepsilon^2/2} \\ &\leq \delta/2. \end{aligned}$$

Therefore, with probability at least  $1 - \delta$ ,

$$\|P_n - P\|_{\mathcal{F}} \leq \mathbb{E}\|P_n - P\|_{\mathcal{F}} + \varepsilon \leq 2\mathbb{E}R_n(\mathcal{F}(Z^n)) + \varepsilon \leq 2r_n(\mathcal{F}(Z^n)) + 3\varepsilon, \quad \forall n \geq \bar{n}$$

which is (18). This shows that (12) holds for  $v(\varepsilon, \delta)$  with  $K_1 = 5$ .

Next, we will prove that, for any  $P \in \mathcal{P}(Z)$ ,

$$\mathbb{P} \left\{ \min_{\bar{n} \leq n < v(6\varepsilon, \delta)} \|P_n - P\|_{\mathcal{F}} < \varepsilon \right\} \leq \delta. \quad (19)$$

In other words, (19) says that, with probability at least  $1 - \delta$ ,  $\|P_n - P\|_{\mathcal{F}} \geq \varepsilon$  for all  $\bar{n} \leq n < v(6\varepsilon, \delta)$ . This means that, for any  $\tau \in \mathcal{T}(\varepsilon, \delta)$ ,  $v(6\varepsilon, \delta) \leq \tau$  with probability at least  $1 - \delta$ , which will give us (13) with  $K_2 = 6$  and  $K_3 = 1$ .

To prove (19), we have by (7) and the union bound that

$$\mathbb{P} \left\{ \bigcup_{n \geq \bar{n}} \{\|P_n - P\|_{\mathcal{F}} \leq \mathbb{E}\|P_n - P\|_{\mathcal{F}} - \varepsilon\} \right\} \leq \delta/2.$$

By the desymmetrization inequality (8), we have

$$\mathbb{E}\|P_n - P\|_{\mathcal{F}} \geq \frac{1}{2}\mathbb{E}R_n(\mathcal{F}(Z^n)) - \frac{1}{2\sqrt{n}}, \quad \forall n.$$

Finally, by the concentration inequality (10) and the union bound,

$$\mathbb{P} \left\{ \bigcup_{n \geq \bar{n}} \{r_n(\mathcal{F}(Z^n)) \geq \mathbb{E}R_n(\mathcal{F}(Z^n)) + \varepsilon\} \right\} \leq \delta/2.$$



Therefore, with probability at least  $1 - \delta$ ,

$$\|P_n - P\|_{\mathcal{F}} \geq \frac{1}{2} r_n(\mathcal{F}(Z^n)) - \frac{1}{2\sqrt{n}} - \frac{3\varepsilon}{2}, \quad \forall n \geq \bar{n}.$$

If  $\bar{n} \leq n < v(6\varepsilon, \delta)$ , then  $r_n(\mathcal{F}(Z^n)) > 6\varepsilon$ . Therefore, using the fact that  $n \geq \bar{n}$  and  $\bar{n}(\varepsilon, \delta)^{-1/2} \leq \varepsilon$ , we see that, with probability at least  $1 - \delta$ ,

$$\|P_n - P\|_{\mathcal{F}} > \frac{3\varepsilon}{2} - \frac{1}{2\sqrt{n}} \geq \frac{3\varepsilon}{2} - \frac{1}{2\sqrt{\bar{n}}} \geq \varepsilon, \quad \bar{n} \leq n < v(6\varepsilon, \delta).$$

This proves (19), and we are done.  $\square$

### 3.2 A weakly efficient sequential learning algorithm

Now choose

$$\bar{n}(\varepsilon, \delta) \geq \left\lfloor \frac{2}{\varepsilon^2} \log \frac{4}{\delta} \right\rfloor + 1, \quad (20)$$

for each  $k = 0, 1, 2, \dots$  let  $n_k \triangleq 2^k \bar{n}(\varepsilon, \delta)$ , and let

$$v(\varepsilon, \delta) \triangleq \min \{n_k : r_{n_k}(\mathcal{F}(Z^{n_k})) \leq \varepsilon\}. \quad (21)$$

**Theorem 2.** *The family  $\{v(\varepsilon, \delta) : \varepsilon > 0, \delta \in (0, 1/2)\}$  defined in (21) with  $\bar{n}(\varepsilon, \delta)$  set according to (20) is WE for any class  $\mathcal{F}$  of measurable functions  $f : Z \rightarrow [0, 1]$  and  $\mathcal{P} = \mathcal{P}(Z)$  with  $K_1 = 5$ ,  $K_2 = 18$ ,  $K_3 = 3$ .*

*Proof.* As before, let  $\bar{n} = \bar{n}(\varepsilon, \delta)$ . The proof of (14) is similar to what we have done in the proof of Theorem 1, except we use the bounds

$$\begin{aligned} \mathbb{P} \left\{ \bigcup_{k=0}^{\infty} \{r_{n_k}(\mathcal{F}(Z^{n_k})) \geq \mathbb{E}R_{n_k}(\mathcal{F}(Z^{n_k})) + \varepsilon\} \right\} &\leq \sum_{k=0}^{\infty} e^{-2^k \bar{n} \varepsilon^2 / 2} \\ &= e^{-\bar{n} \varepsilon^2 / 2} + e^{-\bar{n} \varepsilon^2 / 2} \sum_{k=1}^{\infty} e^{-\frac{\bar{n} \varepsilon^2}{2} (2^k - 1)} \\ &\leq e^{-\bar{n} \varepsilon^2 / 2} + e^{-\bar{n} \varepsilon^2 / 2} \sum_{k=1}^{\infty} e^{-(2^k - 1)} \\ &\leq e^{-\bar{n} \varepsilon^2 / 2} + e^{-\bar{n} \varepsilon^2 / 2} \sum_{k=1}^{\infty} e^{-k} \\ &\leq 2e^{-\bar{n} \varepsilon^2 / 2} \\ &\leq \delta / 2, \end{aligned}$$

where in the third step we have used the fact that  $\bar{n} \varepsilon^2 / 2 \geq 1$ . Similarly,

$$\mathbb{P} \left\{ \bigcup_{k=0}^{\infty} \{\|P_{n_k} - P\|_{\mathcal{F}} \leq \mathbb{E}\|P_{n_k} - P\|_{\mathcal{F}} + \varepsilon\} \right\} \leq \delta^2.$$

Therefore,

$$\|P_{n_k} - P\|_{\mathcal{F}} \leq 2r_{n_k}(\mathcal{F}(Z^{n_k})) + 3\varepsilon, \quad \forall k = 0, 1, 2, \dots$$

and consequently

$$\mathbb{P}\{\|P_{v(\varepsilon, \delta)} - P\|_{\mathcal{F}} \geq 5\varepsilon\} \leq \delta,$$

which proves (14).

Now we prove (15). Let  $N = N(\varepsilon, \delta)$ , the sample complexity of empirical approximation that we have defined in (3). Let us choose  $k$  so that  $n_k \leq N < n_{k+1}$ , which is equivalent to  $2^k \bar{n} \leq N < 2^{k+1} \bar{n}$ . Then

$$\mathbb{P}\{v(18\varepsilon, \delta) > N\} \leq \mathbb{P}\{v(18\varepsilon, \delta) > n_k\}.$$

We will show that the probability on the right-hand side is less than  $3\delta$ . First of all, since  $N \geq \bar{n}$  (by hypothesis), we have  $n_k \geq \bar{n}/2 \geq 1/\varepsilon^2$ . Therefore, with probability at least  $1 - \delta$

$$\|P_{n_k} - P\|_{\mathcal{F}} \geq \frac{1}{2}r_{n_k}(\mathcal{F}(Z^{n_k})) - \frac{1}{2\sqrt{n_k}} - \frac{9\varepsilon}{2} \geq \frac{1}{2}r_{n_k}(\mathcal{F}(Z^{n_k})) - 5\varepsilon. \quad (22)$$

If  $v(18\varepsilon, \delta) > n_k$ , then by definition  $r_{n_k}(\mathcal{F}(Z^{n_k})) > 18\varepsilon$ . Writing  $r_{n_k} = r_{n_k}(\mathcal{F}(Z^{n_k}))$  for brevity, we see get

$$\begin{aligned} \mathbb{P}\{v(18\varepsilon, \delta) > n_k\} &\leq \mathbb{P}\{r_{n_k} > 18\varepsilon\} \\ &= \mathbb{P}\{r_{n_k} > 18\varepsilon, \|P_{n_k} - P\|_{\mathcal{F}} \geq 18\varepsilon\} + \mathbb{P}\{r_{n_k} > 18\varepsilon, \|P_{n_k} - P\|_{\mathcal{F}} < 4\varepsilon\} \\ &\leq \mathbb{P}\{\|P_{n_k} - P\|_{\mathcal{F}} \geq 4\varepsilon\} + \mathbb{P}\{r_{n_k} > 18\varepsilon, \|P_{n_k} - P\|_{\mathcal{F}} < 4\varepsilon\}. \end{aligned}$$

If  $r_{n_k} > 18\varepsilon$  but  $\|P_{n_k} - P\|_{\mathcal{F}} < 4\varepsilon$ , the event in (22) cannot occur. Indeed, suppose it does. Then it must be the case that  $4\varepsilon > 9\varepsilon - 5\varepsilon = 4\varepsilon$ , which is a contradiction. Therefore,

$$\mathbb{P}\{r_{n_k} > 18\varepsilon, \|P_{n_k} - P\|_{\mathcal{F}} < 4\varepsilon\} \leq \delta,$$

and hence

$$\mathbb{P}\{v(18\varepsilon, \delta) > n_k\} \leq \mathbb{P}\{\|P_{n_k} - P\|_{\mathcal{F}} \geq 4\varepsilon\} + \delta.$$

For each  $f \in \mathcal{F}$  and each  $n \in \mathbb{N}$  define

$$S_n(f) \triangleq \sum_{i=1}^n [f(Z_i) - P(f)]$$

and let  $\|S_n\|_{\mathcal{F}} \triangleq \sup_{f \in \mathcal{F}} |S_n(f)|$ . Then

$$\mathbb{P}\{\|P_{n_k} - P\|_{\mathcal{F}} \geq 4\varepsilon\} = \mathbb{P}\{\|S_{n_k}\|_{\mathcal{F}} \geq 4\varepsilon n_k\} \leq \mathbb{P}\{\|S_{n_k}\|_{\mathcal{F}} \geq 2\varepsilon N\}.$$

Since  $n_k \leq N$ , the  $\mathcal{F}$ -indexed stochastic processes  $S_{n_k}(f)$  and  $S_N(f) - S_{n_k}(f)$  are independent. Therefore, we use a technical result stated as Lemma 4 in the appendix with  $\xi_1 = S_{n_k}$  and  $\xi_2 = S_N(f) - S_{n_k}(f)$  to write

$$\mathbb{P}\{\|S_{n_k}\|_{\mathcal{F}} \geq 2\varepsilon N\} \leq \frac{\mathbb{P}\{\|S_N\|_{\mathcal{F}} \geq \varepsilon N\}}{\inf_{f \in \mathcal{F}} \mathbb{P}\{|S_N(f) - S_{n_k}(f)| \leq \varepsilon N\}}.$$

By definition of  $N = N(\varepsilon, \delta)$ , the probability in the numerator is at most  $\delta$ . To analyze the probability in the denominator, we use Hoeffding's inequality to get

$$\begin{aligned} \inf_{f \in \mathcal{F}} \mathbb{P} \{ |S_N(f) - S_{n_k}(f)| \leq \varepsilon N \} &= 1 - \sup_{f \in \mathcal{F}} \mathbb{P} \{ |S_N(f) - S_{n_k}(f)| > \varepsilon N \} \\ &\geq 1 - 2e^{-N\varepsilon^2/2} \\ &\geq 1 - \delta. \end{aligned}$$

Therefore,

$$\mathbb{P} \{ \nu(18\varepsilon, \delta) > n_k \} \leq \frac{\delta}{1 - \delta} + \delta \leq 3\delta$$

for  $\delta < 1/2$ . Therefore,  $\{\nu(\varepsilon, \delta) : \varepsilon \in (0, 1), \delta \in (0, 1/2)\}$  is WE with  $K_1 = 5, K_2 = 18, K_3 = 3$ .  $\square$

## 4 A sequential algorithm for stochastic simulation

Armed with these results on sequential learning algorithms, we can take up the question of constructing efficient simulation strategies. We fix an accuracy parameter  $\varepsilon > 0$ , a confidence parameter  $\delta \in (0, 1)$ , and a level parameter  $\alpha \in (0, 1)$ . Given two probability distributions,  $P$  on the input space  $Z$  and  $Q$  on the parameter space  $\Theta$ , we draw a large i.i.d. sample  $Z_1, \dots, Z_n$  from  $P$  and a large i.i.d. sample  $\theta_1, \dots, \theta_m$  from  $Q$ . We then compute

$$\hat{\theta} = \operatorname{argmin}_{\theta \in \{\theta_1, \dots, \theta_m\}} L_n(\theta),$$

where

$$L_n(\theta) \triangleq \frac{1}{n} \sum_{i=1}^n \ell(Z_i, \theta).$$

The goal is to pick  $n$  and  $m$  large enough so that, with probability at least  $1 - \delta$ ,  $\hat{\theta}$  is an  $\varepsilon$ -minimizer of  $L$  to level  $\alpha$ , i.e., with probability at least  $1 - \delta$  there exists some set  $\Lambda \subset \Theta$  with  $Q(\Lambda) \leq \alpha$ , such that Eq. (2) holds with probability at least  $1 - \delta$ .

To that end, consider the following algorithm based on Theorem 2, proposed by Koltchinskii et al. [KAA<sup>+</sup>00a, KAA<sup>+</sup>00b]:

**Algorithm 1**

choose positive integers  $m$  and  $n$  such that  
 $m \geq \frac{\log(2/\delta)}{\log[1/(1-\alpha)]}$  and  $n \geq \lfloor \frac{50}{\varepsilon^2} \log \frac{8}{\delta} \rfloor + 1$   
draw  $m$  independent samples  $\theta_1, \dots, \theta_m$  from  $Q$   
draw  $n$  independent samples  $Z_1, \dots, Z_n$  from  $P_Z$   
evaluate the stopping variable  
 $\gamma = \max_{1 \leq j \leq m} \left| \frac{1}{n} \sum_{i=1}^n \sigma_i \ell(Z_i, \theta_j) \right|$   
where  $\sigma_1, \dots, \sigma_n$  are i.i.d. Rademacher r.v.'s independent of  $\theta^m$  and  $Z^n$   
if  $\gamma > \varepsilon/5$ , then  
add  $n$  more i.i.d. samples from  $P_Z$  and repeat  
else stop and output  
 $\hat{\theta} = \arg \min_{\theta \in \{\theta_1, \dots, \theta_m\}} L_n(\theta)$

Then we claim that, with probability at least  $1 - \delta$ ,  $\hat{\theta}$  is an  $\varepsilon$ -minimizer of  $L$  to level  $\alpha$ . To see this, we need the following result [Vid03, Lemma 11.1]:

**Lemma 3.** *Let  $Q$  be a probability distribution on the parameter set  $\Theta$ , and let  $h : \Theta \rightarrow \mathbb{R}$  be a (measurable) real-valued function on  $\Theta$ , bounded from above, i.e.,  $h(\theta) < +\infty$  for all  $\theta \in \Theta$ . Let  $\theta_1, \dots, \theta_m$  be  $m$  i.i.d. samples from  $Q$ , and let*

$$\bar{h}(\theta^m) \triangleq \max_{1 \leq j \leq m} h(\theta_j).$$

Then for any  $\alpha \in (0, 1)$

$$Q(\{\theta \in \Theta : h(\theta) > \bar{h}(\theta^m)\}) \leq \alpha \quad (23)$$

with probability at least  $1 - (1 - \alpha)^m$ .

*Proof.* For each  $c \in \mathbb{R}$ , let

$$F(c) \triangleq \mathbb{P}(\{\theta \in \Theta : h(\theta) \leq c\}).$$

Note that  $F$  is the CDF of the random variable  $\xi = h(\theta)$  with  $\theta \sim Q$ . Therefore, it is right-continuous, i.e.,  $\lim_{c' \searrow c} F(c') = F(c)$ . Now define

$$c_\alpha \triangleq \inf\{c : F(c) \geq 1 - \alpha\}.$$

Since  $F$  is right-continuous,  $F(c_\alpha) \geq 1 - \alpha$ . Moreover, if  $c < c_\alpha$ , then  $F(c) < 1 - \alpha$ . Now let us suppose that  $\bar{h}(\theta^m) \geq c_\alpha$ . Then, since  $F$  is monotone nondecreasing,

$$\mathbb{P}(\{\theta \in \Theta : h(\theta) \leq \bar{h}(\theta^m)\}) = F(\bar{h}(\theta^m)) \geq F(c_\alpha) \geq 1 - \alpha,$$

or, equivalently, if  $\bar{h}(\theta^m) \geq c_\alpha$ , then

$$\mathbb{P}(\{\theta \in \Theta : h(\theta) > \bar{h}(\theta^m)\}) \leq \alpha.$$

Therefore, if  $\theta^m$  is such that

$$\mathbb{P}(\{\theta \in \Theta : h(\theta) > \bar{h}(\theta^m)\}) > \alpha,$$

then it must be the case that  $\bar{h}(\theta^m) < c_\alpha$ , which in turn implies that  $F(\bar{h}(\theta^m)) < 1 - \alpha$ , the complement of the event in (23). But  $\bar{h}(\theta^m) < c_\alpha$  means that  $h(\theta_j) < c_\alpha$  for every  $1 \leq j \leq m$ . Since the  $\theta_j$ 's are independent, the events  $\{h(\theta_j) < c_\alpha\}$  are independent, and each occurs with probability at most  $1 - \alpha$ . Therefore,

$$\mathbb{P}(\{\theta^m \in \Theta^m : Q(\{\theta \in \Theta : h(\theta) > \bar{h}(\theta^m)\})\}) \leq (1 - \alpha)^m,$$

which is what we intended to prove.  $\square$

We apply this lemma to the function  $h(\theta) = -L(\theta)$ . Then, provided  $m$  is chosen as described in Algorithm 1, we will have

$$Q\left(\left\{\theta \in \Theta : L(\theta) < \min_{1 \leq j \leq m} L(\theta_j)\right\}\right) \leq \delta/2.$$

Now consider the *finite* class of functions  $\mathcal{F} = \{f_j(z) = \ell(z, \theta_j) : 1 \leq j \leq m\}$ . By Theorem 2, the final output  $\hat{\theta} \in \{\theta_1, \dots, \theta_m\}$  will satisfy

$$\left|L(\hat{\theta}) - \min_{1 \leq j \leq m} L(\theta_j)\right| \leq \varepsilon$$

with probability at least  $1 - \delta/2$ . Hence, with probability at least  $1 - \delta$  there exists a set  $\Lambda \subset \Theta$  with  $Q(\Lambda) \leq \alpha$ , such that (2) holds. Moreover, the total number of samples used up by Algorithm 1 will be, with probability at least  $1 - 3\delta/2$ , no more than

$$N_{\mathcal{F}, P_Z}(\varepsilon/18, \delta/2) \equiv \min\{n \in \mathbb{N} : \mathbb{P}(\|P_n - P_Z\|_{\mathcal{F}} > \varepsilon/18) < \delta/2\}.$$

We can estimate  $N_{\mathcal{F}, P_Z}(\varepsilon/18, \delta/2)$  as follows. First of all, the function

$$\Delta(Z^n) \triangleq \|P_n - P_Z\|_{\mathcal{F}} \equiv \max_{1 \leq j \leq m} |P_n(f_j) - P_Z(f_j)|$$

has bounded differences with  $c_1 = \dots = c_n = 1/n$ . Therefore, by McDiarmid's inequality

$$\mathbb{P}(\Delta(Z^n) \geq \mathbb{E}\Delta(Z^n) + t) \leq e^{-2nt^2}, \quad \forall t > 0.$$

Secondly, since the class  $\mathcal{F}$  is finite with  $|\mathcal{F}| = m$ , the symmetrization inequality (5) and the Finite Class Lemma give the bound

$$\mathbb{E}\|P_n - P_Z\|_{\mathcal{F}} \leq 4\sqrt{\frac{\log m}{n}}.$$

Therefore, if we choose  $t = \varepsilon/18 - 4\sqrt{n^{-1}\log m}$  and  $n$  is large enough so that  $t > \varepsilon/20$  (say), then

$$\mathbb{P}(\|P_n - P\|_{\mathcal{F}} > \varepsilon/18) \leq e^{-n\varepsilon^2/200}.$$

Hence, a fairly conservative estimate is

$$N_{\mathcal{F}, P_Z}(\varepsilon/18, \delta/2) \leq \max\left\{\left\lceil \frac{200}{\varepsilon^2} \log \frac{2}{\delta} \right\rceil + 1, \left\lceil \left(\frac{720}{\varepsilon}\right)^2 \log m \right\rceil + 1\right\}$$

It is instructive to compare Algorithm 1 with a simple Monte Carlo strategy:

**Algorithm 0**

choose positive integers  $m$  and  $n$  such that  
 $m \geq \frac{\log(2/\delta)}{\log[1/(1-\alpha)]}$  and  $n \geq \frac{1}{2\varepsilon^2} \log \frac{4m}{\delta}$   
draw  $m$  independent samples  $\theta_1, \dots, \theta_m$  from  $Q$   
draw  $n$  independent samples  $Z_1, \dots, Z_n$  from  $P_Z$   
for  $j = 1$  to  $m$   
  compute  $L_n(\theta_j) = \frac{1}{n} \sum_{i=1}^n \ell(Z_i, \theta_j)$   
end for  
output  $\hat{\theta} = \operatorname{argmin}_{\theta \in \{\theta_1, \dots, \theta_m\}} L_n(\theta_j)$

The selection of  $m$  is guided by the same considerations as in Algorithm 1. Moreover, for each  $1 \leq j \leq m$ ,  $L_n(\theta_j)$  is an average of  $n$  independent random variables  $\ell(Z_i, \theta_j) \in [0, 1]$ , and  $L(\theta_j) = \mathbb{E}L_n(\theta_j)$ . Hence, Hoeffding's inequality says that

$$\mathbb{P}(\{Z^n \in Z^n : |L_n(\theta_j) - L(\theta_j)| > \varepsilon\}) \leq 2e^{-2n\varepsilon^2}.$$

If we choose  $n$  as described in Algorithm 0, then

$$\begin{aligned} \mathbb{P}\left(\left|L_n(\hat{\theta}) - \min_{1 \leq j \leq m} L(\theta_j)\right| > \varepsilon\right) &\leq \mathbb{P}\left(\bigcup_{j=1}^m |L_n(\theta_j) - L(\theta_j)| > \varepsilon\right) \\ &\leq \sum_{j=1}^m \mathbb{P}(|L_n(\theta_j) - L(\theta_j)| > \varepsilon) \\ &\leq \delta/2. \end{aligned}$$

Hence, with probability at least  $1 - \delta$  there exists a set  $\Lambda \subset \Theta$  with  $Q(\Lambda) \leq \alpha$ , so that (2) holds. It may seem at first glance that Algorithm 0 is more efficient than Algorithm 1. However, this is not the case in high-dimensional situations. There, one can actually show that, with probability practically equal to one, the empirical minimum of  $L$  can be *much larger* than the true minimum (cf. [KAA<sup>+</sup>00b] for a very vivid numerical illustration). This is an instance of the so-called *Curse of Dimensionality*, which adaptive schemes like Algorithm 1 can often avoid.

## A Technical lemma

**Lemma 4.** Let  $\{\xi_1(f) : f \in \mathcal{F}\}$  and  $\{\xi_2(f) : f \in \mathcal{F}\}$  be two independent  $\mathcal{F}$ -indexed stochastic processes with

$$\|\xi_j\|_{\mathcal{F}} \triangleq \sup_{f \in \mathcal{F}} |\xi_j(f)| < \infty, \quad j = 1, 2.$$

Then for all  $t > 0, c > 0$

$$\mathbb{P}\{\|\xi_1\|_{\mathcal{F}} \geq t + c\} \leq \frac{\mathbb{P}\{\|\xi_1 - \xi_2\|_{\mathcal{F}} \geq t\}}{\inf_{f \in \mathcal{F}} \mathbb{P}\{|\xi_2(f)| \leq c\}}. \quad (24)$$

*Proof.* If  $\|\xi_1\|_{\mathcal{F}} \geq t + c$ , then there exists some  $f \in \mathcal{F}$ , such that  $|\xi_1(f)| \geq t + c$ . Then for this particular  $f$  by the triangle inequality we see that

$$|\xi_2(f)| \leq c \quad \Rightarrow \quad |\xi_1(f) - \xi_2(f)| \geq t$$

Therefore,

$$\inf_{f \in \mathcal{F}} \mathbb{P}_{\xi_2} \left\{ |\xi_2(f)| \leq c \right\} \leq \mathbb{P}_{\xi_2} \left\{ |\xi_2(f)| \leq c \right\} \leq \mathbb{P}_{\xi_2} \left\{ |\xi_1(f) - \xi_2(f)| \geq t \right\} \leq \mathbb{P}_{\xi_2} \left\{ \|\xi_1 - \xi_2\|_{\mathcal{F}} \geq t \right\}.$$

The leftmost and the rightmost terms in the above inequality do not depend on the particular  $f$ , and the inequality between them is valid on the event  $\{\|\xi_1\|_{\mathcal{F}} \geq t + c\}$ . Therefore, integrating the two sides w.r.t.  $\xi_1$  on this event, we get

$$\inf_{f \in \mathcal{F}} \mathbb{P}_{\xi_2} \left\{ |\xi_2(f)| \leq c \right\} \cdot \mathbb{P}_{\xi_1} \left\{ \|\xi_1\|_{\mathcal{F}} \geq t + c \right\} \leq \mathbb{P}_{\xi_1, \xi_2} \left\{ \|\xi_1 - \xi_2\|_{\mathcal{F}} \geq t \right\}.$$

Rearranging, we get (24). □

## References

- [KAA<sup>+</sup>00a] V. Koltchinskii, C. T. Abdallah, M. Ariola, P. Dorato, and D. Panchenko. Improved sample complexity estimates for statistical learning control of uncertain systems. *IEEE Transactions on Automatic Control*, 45(12):2383–2388, 2000.
- [KAA<sup>+</sup>00b] V. Koltchinskii, C. T. Abdallah, M. Ariola, P. Dorato, and D. Panchenko. Statistical learning control of uncertain systems: it is better than it seems. Technical Report EECE-TR-00-001, University of New Mexico, April 2000.
- [Vid98] M. Vidyasagar. Statistical learning theory and randomized algorithms for control. *IEEE Control Magazine*, 18(6):162–190, 1998.
- [Vid01] M. Vidyasagar. Randomized algorithms for robust controller synthesis using statistical learning theory. *Automatica*, 37:1515–1528, 2001.
- [Vid03] M. Vidyasagar. *Learning and Generalization*. Springer, 2 edition, 2003.