# On the Stability of Bootstrap Estimators

#### A. Christmann and M. Salibían-Barrera and S. Van Aelst

<sup>1</sup> University of Bayreuth, Department of Mathematics, Bayreuth, GERMANY. e-mail: andreas.christmann@uni-bayreuth.de

<sup>2</sup> University of British Columbia, Department of Statistics, Vancouver, CANADA. e-mail: matias@stat.ubc.ca

<sup>3</sup> University of Ghent, Department of Applied Mathematics and Computer Science, Ghent, BELGIUM.

e-mail: Stefan.VanAelst@UGent.be

**Abstract:** It is shown that bootstrap approximations of an estimator which is based on a continuous operator from the set of Borel probability measures defined on a compact metric space into a complete separable metric space is stable in the sense of qualitative robustness. Support vector machines based on shifted loss functions are treated as special cases.

**Keywords and phrases:** bootstrap, statistical machine learning, stability, support vector machine, robustness.

#### 1. Introduction

The finite sample distribution of many nonparametric methods from statistical learning theory is unknown because the distribution P from which the data were generated is unknown and because there are often only asymptotical results on the behaviour of such methods known.

The goal of this paper is to show that bootstrap approximations of an estimator which is based on a continuous operator from the set of Borel probability distributions defined on a compact metric space into a complete separable metric space is stable in the sense of qualitative robustness. As a special case it is shown that bootstrap approximations for the support vector machine (SVM) are stable, both for the risk functional and for the SVM operator itself. The results can be interpreted as generalizations of theorems derived by [4].

The rest of the paper has the following structure. Section 2 gives the general

result and Section 3 contains the results for SVMs. All proofs are given in the appendix.

## 2. On Qualitative Robustness of Bootstrap Estimators

If not otherwise mentioned, we will use the Borel  $\sigma$ -algebra  $\mathcal{B}(A)$  on a set A and denote the Borel  $\sigma$ -algebra on  $\mathbb{R}$  by  $\mathcal{B}$ .

**Assumption 1.** Let  $(\Omega, \mathcal{A}, \mu)$  be a probability space, where  $\mu$  is unknown,  $(\mathcal{Z}, d_{\mathcal{Z}})$  be a compact metric space, and  $\mathcal{B}(\mathcal{Z})$  be the Borel  $\sigma$ -algebra on  $\mathcal{Z}$ . Denote the set of all Borel probability measures on  $(\mathcal{Z}, \mathcal{B}(\mathcal{Z}))$  by  $\mathcal{M}_1(\mathcal{Z}, \mathcal{B}(\mathcal{Z}))$ . On  $\mathcal{M}_1(\mathcal{Z}, \mathcal{B}(\mathcal{Z}))$  we use the Borel  $\sigma$ -algebra  $\mathcal{B}(\mathcal{M}_1(\mathcal{Z}, \mathcal{B}(\mathcal{Z})))$  and the bounded Lipschitz metric  $d_{BL}$ , see (4.11). Let S be a statistical operator defined on  $\mathcal{M}_1(\mathcal{Z},\mathcal{B}(\mathcal{Z}))$  with values in a complete, separable metric space  $(\mathcal{W},d_{\mathcal{W}})$  enclipped with its Borel  $\sigma$ -algebra  $\mathcal{B}(\mathcal{W})$ . Let  $Z, Z_n : (\Omega, \mathcal{A}, \mu) \to (\mathcal{Z}, \mathcal{B}(\mathcal{Z}))$ ,  $n \in \mathbb{N}$ , be independent and identically distributed random variables and denote the image measure by  $P := Z \circ \mu$ . Let  $S_n(Z_1, \ldots, Z_n)$  be a statistic with values in  $(W, \mathcal{B}(W))$ . Denote the empirical measure of  $(Z_1, \ldots, Z_n)$  by  $P_n := \frac{1}{n} \sum_{i=1}^n \delta_{Z_i}$ . The statistic  $S_n$  is defined via the operator

$$S: (\mathcal{M}_1(\mathcal{Z}, \mathcal{B}(\mathcal{Z})), \mathcal{B}(\mathcal{M}_1(\mathcal{Z}, \mathcal{B}(\mathcal{Z}))) \to (\mathcal{W}, \mathcal{B}(\mathcal{W}))$$

where  $S(P_n) = S_n(Z_1, \ldots, Z_n)$ . Denote the distribution of  $S_n(Z_1, \ldots, Z_n)$ when  $Z_i \stackrel{i.i.d.}{\sim} P$  by  $\mathfrak{L}_n(S; P) := \mathfrak{L}(S_n(Z_1, \ldots, Z_n))$ . Accordingly, we denote the distribution of  $S_n(Z_1, \ldots, Z_n)$  when  $Z_i \stackrel{i.i.d.}{\sim} P_n$  by  $\mathfrak{L}_n(S; P_n)$ .

Efron [9, 10] proposed the bootstrap, whose main idea is to approximate the unknown distribution  $\mathfrak{L}_n(S; P)$  by  $\mathfrak{L}_n(S; P_n)$ . Note that these bootstrap approximations  $\mathfrak{L}_n(S; P_n)$  are (probability measure-valued) random variables with values in  $\mathcal{M}_1(\mathcal{W}, \mathcal{B}(\mathcal{W}))$ .

Following [4] we call a sequence of bootstrap approximations  $\mathfrak{L}_n(S; P_n)$ qualitatively robust at  $P \in \mathcal{M}_1(\mathcal{Z}, \mathcal{B}(\mathcal{Z}))$  if the sequence of transformations

$$g_n: \mathcal{M}_1(\mathcal{Z}, \mathcal{B}(\mathcal{Z})) \to \mathcal{M}_1(\mathcal{W}, \mathcal{B}(\mathcal{W})), \quad g_n(Q) = \mathfrak{L}(\mathfrak{L}_n(S; Q_n)), \qquad n \in \mathbb{N},$$

$$(2.1)$$

is asymptotically equicontinuous at  $P \in \mathcal{M}_1(\mathcal{Z}, \mathcal{B}(\mathcal{Z}))$ , i.e. if

$$\forall \varepsilon > 0 \; \exists \, \delta > 0 \; \exists \, n_0 \in \mathbb{N} :$$

$$d_{\mathrm{BL}}(\mathbf{Q}, \mathbf{P}) < \delta \quad \Rightarrow \quad \sup_{n \ge n_0} d_{\mathrm{BL}} \big( \mathfrak{L}(\mathfrak{L}_n(S; \mathbf{Q}_n)), \mathfrak{L}(\mathfrak{L}_n(S; \mathbf{P}_n)) \big) < \varepsilon.$$

$$(2.2)$$

Following [4] again, we call a sequence of statistics  $(S_n)_{n\in\mathbb{N}}$  uniformly qualitatively robust in a neighborhood  $\mathcal{U}(P_0)$  of  $P_0 \in \mathcal{M}_1(\mathcal{Z}, \mathcal{B}(\mathcal{Z}))$  if

$$\exists n_0 \in \mathbb{N} \ \forall \varepsilon > 0 \ \forall n \ge n_0 \ \exists \delta > 0 \ \forall P \in \mathcal{U}(P_0) :$$

$$d_{BL}(Q, P) < \delta \quad \Rightarrow \quad d_{BL}(\mathfrak{L}_n(S; Q), \mathfrak{L}_n(S; P)) < \varepsilon.$$
(2.3)

The following two results and Theorem 8 in the next section are the main results of this paper.

**Theorem 2.** If Assumption 1 is valid and if S is uniformly continuous in a neighborhood  $\mathcal{U}(P_0)$  of  $P_0 \in \mathcal{M}_1(\mathcal{Z}, \mathcal{B}(\mathcal{Z}))$ , then  $(S_n(Z_1, \ldots, Z_n))_{n \in \mathbb{N}}$  is uniformly qualitatively robust in  $\mathcal{U}(P_0)$ .

**Theorem 3.** If Assumption 1 is valid and if  $(S_n(Z_1,\ldots,Z_n))_{n\in\mathbb{N}}$  is uniformly qualitatively robust in a neighborhood  $\mathcal{U}(P_0)$  of  $P_0 \in \mathcal{M}_1(\mathcal{Z}, \mathcal{B}(\mathcal{Z}))$ , then the sequence  $\mathfrak{L}_n(S; P_n)$  of bootstrap approximations of  $\mathfrak{L}_n(S; P)$  is qualitatively robust for  $P_0$ .

As an immediate consequence from both theorems given above we obtain

Corollary 4. If Assumption 1 is valid and if S is a continuous operator, then the sequence  $\mathfrak{L}_n(S; P_n)$  of bootstrap approximations of  $\mathfrak{L}_n(S; P)$  is qualitatively robust for all  $P \in \mathcal{M}_1(\mathcal{Z}, \mathcal{B}(\mathcal{Z}))$ .

**Remark 5.** The Theorems 2 and 3 can be considered as a generalization of [4, Thm. 2, Thm. 3], who considered the case  $W := A \subset \mathbb{R}$  being a finite interval and  $\mathcal{Z} := \mathbb{R}$ -valued random variables  $Z_1, \ldots, Z_n$ . In our case, the statistics  $S_n(Z_1,\ldots,Z_n)$  are W-valued statistics, where W is a complete separable metric space and its dimension can be infinite.

# 3. On Qualitative Robustness of Bootstrap SVMs

In this section we will apply the previous results to support vector machines which belong to the modern class of statistical machine learning methods. I.e., we will consider the special case that  $\mathcal{W}$  is a reproducing kernel Hilbert space H used by a support vector machine (SVM). Note that H typically has an infinite dimension, which is true, e.g., if the popular Gaussian RBF kernel  $k: \mathcal{X} \times \mathcal{X} \to \mathbb{R}, k(x, x') := \exp(-\gamma ||x - x'||_2^2)$  for  $\gamma > 0$  is used.

To state our result on the stability of bootstrap SVMs in Theorem 8 below, we need the following assumptions on the loss function and the kernel.

**Assumption 6.** Let  $\mathcal{Z} = \mathcal{X} \times \mathcal{Y}$  be a compact metric space with metric  $d_{\mathcal{Z}}$ , where  $\mathcal{Y} \subset \mathbb{R}$  is closed. Let  $L: \mathcal{X} \times \mathcal{Y} \times \mathbb{R} \to [0, \infty)$  be a loss function such that L is continuous and convex with respect to its third argument and that L is uniformly Lipschitz continuous with respect to its third argument with uniform Lipschitz constant  $|L|_1 > 0$ , i.e.  $|L|_1$  is the smallest constant c such that  $\sup_{(x,y)\in\mathcal{X}\times\mathcal{Y}}|L(x,y,t)-L(x,y,t')|\leq c|t-t'|$  for all  $t,t'\in\mathbb{R}$ . Denote the shifted loss function by  $L^*(x,y,t) := L(x,y,t) - L(x,y,0), (x,y,t) \in \mathcal{X} \times \mathcal{Y} \times \mathcal{Y}$  $\mathbb{R}$ . Let  $k: \mathcal{X} \times \mathcal{X} \to \mathbb{R}$  be a continuous kernel with reproducing kernel Hilbert space H and assume that k is bounded by  $||k||_{\infty} := (\sup_{x \in \mathcal{X}} k(x,x))^{1/2} \in$  $(0,\infty)$ . Let  $\lambda \in (0,\infty)$ .

These assumptions can be considered as standard assumptions for stable SVMs, see, e.g., [1] and [15, Chap. 10], .

In this paper the RKHS H, the penalyzing constant  $\lambda$ , and the loss function L and thus the shifted loss function  $L^*$  are fixed. Therefore, we write in the next definition just S and R instead of  $S_{L^*,H,\lambda}$  and  $R_{L^*,H,\lambda}$  to shorten the notation.

**Definition 7.** The SVM operator  $S: \mathcal{M}_1(\mathcal{Z}, \mathcal{B}(\mathcal{Z})) \to H$  is defined by

$$S(P) := f_{L^*,P,\lambda} := \arg\min_{f \in H} \mathbb{E}_P L^*(X,Y,f(X)) + \lambda \|f\|_H^2.$$
 (3.4)

The SVM risk functional  $R: \mathcal{M}_1(\mathcal{Z}, \mathcal{B}(\mathcal{Z})) \to \mathbb{R}$  is defined by

$$R(P) := \mathbb{E}_{P} L^{\star}(X, Y, S(P)(X)) = \mathbb{E}_{P} L^{\star}(X, Y, f_{L^{\star}, P, \lambda}(X)). \tag{3.5}$$

If Assumption 6 is valid, then S is well-defined because  $S(P) \in H$  exists and is unique, R is well-defined because  $R(P) \in \mathbb{R}$  exists and is unique, and it holds, for all  $P \in \mathcal{M}_1(\mathcal{X} \times \mathcal{Y})$ ,

$$||S(P)||_{\infty} \le \frac{1}{\lambda} |L|_1 ||k||_{\infty}^2 < \infty \text{ and } |R(P)| \le \frac{1}{\lambda} |L|_1^2 ||k||_{\infty}^2 < \infty, \quad (3.6)$$

see [2, Thm 5, Thm. 6, (17), (18)].

**Theorem 8.** If the general Assumption 1 and the Assumption 6 are valid, then the SVM operator S and the SVM risk functional R fulfill:

- (i) The sequence  $\mathfrak{L}_n(S; P_n)$  of bootstrap SVM estimators of  $\mathfrak{L}_n(S; P)$  is qualitatively robust for all  $P \in \mathcal{M}_1(\mathcal{Z}, \mathcal{B}(\mathcal{Z}))$ .
- (ii) The sequence  $\mathfrak{L}_n(R; P_n)$  of bootstrap SVM risk estimators of  $\mathfrak{L}_n(R; P)$ is qualitatively robust for all  $P \in \mathcal{M}_1(\mathcal{Z}, \mathcal{B}(\mathcal{Z}))$ .

#### 4. Proofs

## 4.1. Proofs of the results in Section 2

For the proofs we need Theorem 9 and Theorem 10, see below. To state Theorem 9 on uniform Glivenko-Cantelli classes, we need the following notation. For any metric space  $(\mathcal{S}, d)$  and real-valued function  $f: \mathcal{S} \to \mathbb{R}$ , we denote the bounded Lipschitz norm of f by

$$||f||_{\text{BL}} := \sup_{x \in \mathcal{S}} |f(x)| + \sup_{x,y \in \mathcal{S}, x \neq y} \frac{|f(x) - f(y)|}{d(x,y)}.$$
 (4.7)

Let  $\tilde{F}$  be a set of measurable functions from  $(\mathcal{S}, \mathcal{B}(\mathcal{S})) \to (\mathbb{R}, \mathcal{B})$ . For any function  $G: \tilde{F} \to \mathbb{R}$  (such as a signed measure) define

$$||G||_{\tilde{F}} := \sup\{|G(f)| : f \in \tilde{F}\}.$$
 (4.8)

**Theorem 9.** [8, Prop. 12] For any separable metric space (S, d) and  $M \in$  $(0,\infty),$ 

$$\tilde{\mathcal{F}}_M := \{ f : (\mathcal{S}, \mathcal{B}(\mathcal{S})) \to (\mathbb{R}, \mathcal{B}); \|f\|_{BL} \le M \}$$
(4.9)

is a universal Glivenko-Cantelli class. It is a uniform Glivenko-Cantelli class, i.e., for all  $\varepsilon > 0$ ,

$$\lim_{n \to \infty} \sup_{\nu \in \mathcal{M}_1(\mathcal{S}, \mathcal{B}(\mathcal{S}))} \Pr^* \left( \sup_{m \ge n} \|\nu_m - \nu\|_{\tilde{\mathcal{F}}_M} > \varepsilon \right) = 0, \tag{4.10}$$

if and only if (S, d) is totally bounded. Here,  $Pr^*$  denotes the outer probability.

Note that the term  $\|\nu_m - \nu\|_{\tilde{\mathcal{F}}_M}$  in (4.10) equals the bounded Lipschitz metric  $d_{\rm BL}$  of the probability measures  $\nu_m$  and  $\nu$  if M=1, i.e.

$$\|\nu_{m} - \nu\|_{\tilde{\mathcal{F}}_{1}} = \sup_{f \in \tilde{\mathcal{F}}_{1}} |(\nu_{m} - \nu)(f)| = \sup_{f : \|f\|_{\mathrm{BL}} \le 1} \left| \int f \, d\nu_{m} - \int f \, d\nu \right| =: d_{\mathrm{BL}}(\nu_{m}, \nu),$$
(4.11)

see [7, p. 394]. Hence, Theorem 9 can be interpreted as a generalization of [4, Lemma 1, p. 186], which says that if  $A \subset \mathbb{R}$  is a finite interval, then  $d_{\mathrm{BL}}(\mathrm{P}_m,\mathrm{P})$  converges almost surely to 0 uniformly in  $\mathrm{P}\in\mathcal{M}_1(A,\mathcal{B}(A))$ . For various characterizations of Glivenko-Cantelli classes, we refer to [16, Thm. |22| and |6|.

We next list the other main result we need for the proof of Theorem 8. This result is an analogon of the famous Strassen theorem for the bounded Lipschitz metric  $d_{\rm BL}$  instead of the Prohorov metric.

**Theorem 10.** [13, Thm. 4.2, p. 30] Let  $\mathcal{Z}$  be a Polish space with topology  $\tau_{\mathcal{Z}}$ . Let  $d_{\mathrm{BL}}$  be the bounded Lipschitz metric defined on the set  $\mathcal{M}_1(\mathcal{Z},\mathcal{B}(\mathcal{Z}))$ of all Borel probability measures on Z. Then the following two statements are equivalent:

- (i) There are random variables  $\xi_1$  with distribution  $\nu_1$  and  $\xi_2$  with distribution  $\nu_2$  such that  $\mathbb{E}[d_{\mathrm{BL}}(\xi_1, \xi_2)] \leq \varepsilon$ .
- (ii)  $d_{\rm BL}(\nu_1, \nu_2) \leq \varepsilon$ .

**Proof of Theorem 2.** We closely follow the proof by [4, Thm. 2]. However, we use Theorem 9 instead of their Lemma 1 and we use [3, Lem. 1] instead of [12, Lem. 1].

Let  $\mathcal{P}_n \subset \mathcal{M}_1(\mathcal{Z}, \mathcal{B}(\mathcal{Z}))$  be the set of empirical distributions of order  $n \in \mathbb{N}$ ,

$$\mathcal{P}_n := \left\{ P_n \in \mathcal{M}_1(\mathcal{Z}, \mathcal{B}(\mathcal{Z})); \exists (z_1, \dots, z_n) \in \mathcal{Z}^n \text{ such that } P_n = \frac{1}{n} \sum_{i=1}^n \delta_{z_i} \right\},$$

$$(4.12)$$

and let  $\mathcal{E}_n \subset \mathcal{P}_n$ . If misunderstandings are unlikely, we identify  $\mathcal{E}_n$  with the set  $\{z_1,\ldots,z_n\}$  of atoms.

It is enough to show that

$$\forall \varepsilon > 0 \; \exists \delta > 0 \; \forall P \in \mathcal{U}(P_0) \; \exists \text{ sequence } (\mathcal{E}_n)_{n \in \mathbb{N}} \subset \mathcal{P}_n$$
 (4.13)

such that  $P^n(\mathcal{E}_n) > 1 - \varepsilon$  and for all  $Q_n \in \mathcal{E}_n$  and for all  $\tilde{Q}_n \in \mathcal{P}_n$  we have

$$d_{\mathrm{BL}}(\mathbf{Q}_n, \tilde{\mathbf{Q}}_n) < \delta \quad \Rightarrow \quad d_{\mathcal{W}}(S(\mathbf{Q}_n), S(\tilde{\mathbf{Q}}_n)) < \varepsilon.$$
 (4.14)

From this we obtain that  $(S_n)_{n\in\mathbb{N}}$  is uniformly qualitatively robust by [3, Lem. 1].

Let  $\varepsilon > 0$ . Since the operator S is uniformly continuous in  $\mathcal{U}(P_0)$  we obtain

$$\exists \, \delta_0 > 0 \,\,\forall \, P \in \mathcal{U}(P_0) : \quad d_{BL}(P, Q) < \delta_0 \quad \Rightarrow \quad d_{\mathcal{W}}(S(P), S(Q)) < \varepsilon/2 \,. \tag{4.15}$$

Hence by Theorem 9 for the special case M=1 and by (4.11), we get

$$\exists n_0 \in \mathbb{N} : \sup_{P \in \mathcal{U}(P_0)} \Pr^* \left( \sup_{n \ge n_0} d_{BL}(P_n, P) < \delta_0 \right) > 1 - \varepsilon.$$
 (4.16)

For  $n \geq n_0$  and  $P \in \mathcal{U}(P_0)$ , define

$$\mathcal{E}_{n,P} := \{ Q_n \in \mathcal{P}_n : d_{BL}(Q_n, P) < \delta_0/2 \}.$$
 (4.17)

It follows, that  $P^n(\mathcal{E}_{n,P}) > 1 - \varepsilon$  together with  $Q_n \in \mathcal{E}_{n,P}$  and  $d_{BL}(Q_n, \tilde{Q}_n) < \delta_0/2$  implies that

$$d_{\mathrm{BL}}(\mathbf{Q}_n, \mathbf{P}) < \delta_0/2$$
 and  $d_{\mathrm{BL}}(\tilde{\mathbf{Q}}_n, \mathbf{P}) < \delta_0$ .

The triangle inequality thus yields due to (4.15)

$$d_{\mathcal{W}}(S(Q_n), S(\tilde{Q}_n)) \le d_{\mathcal{W}}(S(Q_n), S(P)) + d_{\mathcal{W}}(S(P), S(\tilde{Q}_n)) < \varepsilon, \quad (4.18)$$

from which the assertion follows.

**Proof of Theorem 3.** The proof mimics the proof of [4, Thm. 3], but uses Theorem 9 instead of [4, Lem. 1].

Fix  $P_0 \in \mathcal{M}_1(\mathcal{Z}, \mathcal{B}(\mathcal{Z}))$  and  $\varepsilon > 0$ . By the uniform qualitative robustness of  $(S_n)_{n \in \mathbb{N}}$  in  $\mathcal{U}(P_0)$ , there exists  $n \in \mathbb{N}$  such that for all  $\varepsilon > 0$  there exists  $\delta > 0$  such that

$$d_{\mathrm{BL}}(\mathbf{Q}, \mathbf{P}) < \delta \quad \Rightarrow \quad \sup_{m \ge n} \sup_{\mathbf{P} \in \mathcal{U}(\mathbf{P}_0)} d_{\mathrm{BL}}(\mathfrak{L}_m(S; \mathbf{Q}), \mathfrak{L}_m(S; \mathbf{P})) < \varepsilon.$$
 (4.19)

Define  $\delta_1 := \delta/2$ . Due to Theorem 9 for the special case M = 1 and by (4.11), we have, for all  $\varepsilon > 0$ ,

$$\lim_{n \to \infty} \sup_{P \in \mathcal{M}_1(\mathcal{Z}, \mathcal{B}(\mathcal{Z}))} \Pr^* \left( \sup_{m \ge n} d_{BL}(P_m, P) > \varepsilon \right) = 0.$$
 (4.20)

Hence (4.19) and Varadarajan's theorem on the almost sure convergence of empirical measures to a Borel probability measure defined on a separable metric space, see e.g. [7, Thm. 11.4.1, p. 399], yields for the empirical distributions  $Q_n$  from Q and  $P_{0,n}$  from  $P_0$  that,

$$\exists n_1 > n \,\forall n \geq n_1: \, d_{\mathrm{BL}}(\mathbf{Q}, \mathbf{P}_0) < \delta_1 \quad \Rightarrow \quad d_{\mathrm{BL}}(\mathbf{Q}_n, \mathbf{P}_{0,n}) < \delta \quad \text{almost surely}.$$

$$\tag{4.21}$$

It follows from the uniform qualitative robustness of  $(S_n)_{n\in\mathbb{N}}$ , see (4.19), that

$$\exists n_1 \in \mathbb{N} \ \forall \varepsilon > 0 \ \forall n \ge n_1 \ \exists \delta > 0 \ \forall P \in \mathcal{U}(P_0) :$$

$$d_{\mathrm{BL}}(Q, P) < \delta \quad \Rightarrow \quad d_{\mathrm{BL}}(\mathfrak{L}_n(S; Q_n), \mathfrak{L}_n(S; P_{0,n})) < \varepsilon \text{ almost surely.}$$

$$(4.22)$$

For notational convenience, we write for the sequences of bootstrap estimators

$$\xi_{1,n} := \mathfrak{L}_n(S; Q_n), \qquad \xi_{2,n} := \mathfrak{L}_n(S; P_{0,n}), \qquad n \in \mathbb{N}.$$
 (4.23)

Note that  $\xi_{1,n}$  and  $\xi_{2,n}$  are (measure-valued) random variables with values in the set  $\mathcal{M}_1(\mathcal{W}, \mathcal{B}(\mathcal{W}))$ . We denote the distribution of  $\xi_{j,n}$  by  $\mu_{j,n}$  for  $j \in \{1, 2\}$  and  $n \in \mathbb{N}$ . Hence (4.22) yields

$$d_{\mathrm{BL}}(\xi_{1,n}, \xi_{2,n}) < \varepsilon$$
 almost surely for all  $n \ge n_1$  (4.24)

and it follows

$$\mathbb{E}[d_{\mathrm{BL}}(\xi_{1,n},\xi_{2,n})] \le \varepsilon, \qquad \forall \, n \ge n_1. \tag{4.25}$$

Now an application of an analogon of Strassen's theorem, see Theorem 10, yields

$$\sup_{n \ge n_1} d_{\mathrm{BL}}(\mathfrak{L}(\xi_{1,n}), \mathfrak{L}(\xi_{2,n})) \le \varepsilon \qquad \forall n \ge n_1, \tag{4.26}$$

which completes the proof, because

$$\mathfrak{L}(\xi_{1,n}) = \mathfrak{L}(\mathfrak{L}_n(S; \mathbf{Q}_n)) \quad \text{and} \quad \mathfrak{L}(\xi_{2,n}) = \mathfrak{L}(\mathfrak{L}_n(S; \mathbf{P}_{0,n})). \tag{4.27}$$

## 4.2. Proofs of the results in Section 3

**Proof of Theorem 8.** Proof of part (i). By assumption,  $(\mathcal{Z}, d_{\mathcal{Z}})$  is a compact metric space, where  $\mathcal{Z} = \mathcal{X} \times \mathcal{Y}$ . Let  $\mathcal{B}(\mathcal{Z})$  be the Borel  $\sigma$ -algebra on  $\mathcal{Z}$ . It is well-known that the bounded Lipschitz metric  $d_{\mathrm{BL}}$  metrizes the weak topology on the space  $\mathcal{M}_1(\mathcal{Z}, \mathcal{B}(\mathcal{Z}))$ , see [7, Thm. 11.3.3], and that  $(\mathcal{M}_1(\mathcal{Z}, \mathcal{B}(\mathcal{Z})), d_{\mathrm{BL}})$  is a compact metric space if and only if  $(\mathcal{Z}, d_{\mathcal{Z}})$  is a compact metric space, see [14, p. 45, Thm. 6.4]. From the compactness of  $(\mathcal{M}_1(\mathcal{Z}, \mathcal{B}(\mathcal{Z})), d_{\mathrm{BL}})$ , it of course follows that this metric space is separable and totally bounded, see [5, Thm. 1.4.26].

Under the assumptions of the theorem we have, for all fixed  $\lambda \in (0, \infty)$ , that the SVM operator  $S: \mathcal{M}_1(\mathcal{Z}, \mathcal{B}(\mathcal{Z})) \to H$ ,  $S(P) = f_{L^*,P,\lambda}$ , is well-defined because it exists and is unique, see [2, Thm. 5, Thm. 6] and is continuous with respect to the combination of the weak topology on  $\mathcal{M}_1(\mathcal{Z}, \mathcal{B}(\mathcal{Z}))$  and the norm topology on H, see [11, Thm. 3.3, Cor. 3.4]. There it was also shown that the operator  $\tilde{S}: \mathcal{M}_1(\mathcal{Z}, \mathcal{B}(\mathcal{Z})) \to \mathcal{C}_b(\mathcal{Z})$ ,  $P \mapsto f_{L^*,P,\lambda}$ , is continuous with respect to the combination of weak topology on  $\mathcal{M}_1(\mathcal{Z}, \mathcal{B}(\mathcal{Z}))$  and the norm topology on  $\mathcal{C}_b(\mathcal{Z})$ . Because  $(\mathcal{M}_1(\mathcal{Z}, \mathcal{B}(\mathcal{Z})), d_{BL})$  is a compact metric space, the operators S and  $\tilde{S}$  are therefore even uniformly continuous on the whole space  $\mathcal{M}_1(\mathcal{Z}, \mathcal{B}(\mathcal{Z}))$  with respect to the mentioned topologies, see [5, Prop. 1.5.9].

Because the reproducing kernel Hilbert space  $\mathcal{W} := H$  is a Hilbert space, His complete. Furthermore, because the input space  $\mathcal{X}$  is separable and the kernel k is continuous, the RKHS H is also separable, see [15, Lem. 4.33]. Therefore, Theorem 2 yields that the sequence of H-valued statistics

$$S_n((X_1, Y_1), \dots, (X_n, Y_n)) = \arg\min_{f \in H} \frac{1}{n} \sum_{i=1}^n L^*(X_i, Y_i, f(X_i)) + \lambda \|f\|_H^2, \ n \in \mathbb{N},$$
(4.28)

is uniformly qualitatively robust in a neighborhood  $\mathcal{U}(P_0)$  for every probability measure  $P_0 \in \mathcal{M}_1(\mathcal{Z})$ . Now we apply Theorem 3, which yields that the sequence  $(\mathfrak{L}_n(S; P_n))_{n \in \mathbb{N}}$  of bootstrap SVM estimators of  $\mathfrak{L}_n(S; P)$  is qualitatively robust for all  $P_0 \in \mathcal{M}_1(\mathcal{Z}, \mathcal{B}(\mathcal{Z}))$ , which gives the first assertion of the theorem.

*Proof of part (ii)*. The proof consists of two steps. In Step 1 the continuity of the SVM risk functional R will be shown. In Step 2, the Theorems 2 and 3 will be used to show that the sequence  $(\mathfrak{L}_n(R; P_n))_{n \in \mathbb{N}}$ ,  $n \in \mathbb{N}$ , of bootstrap SVM risk estimators is qualitatively robust.

Step 1. We will first show that the SVM risk functional  $R: \mathcal{M}_1(\mathcal{Z}, \mathcal{B}(\mathcal{Z})) \to$ R is continuous with respect to the combination of the weak topology on  $\mathcal{M}_1(\mathcal{Z},\mathcal{B}(\mathcal{Z}))$  and the standard topology on  $\mathbb{R}$ .

As mentioned in part (i), the assumption that  $(\mathcal{Z}, d_{\mathcal{Z}})$  is a compact metric space implies that  $(\mathcal{M}_1(\mathcal{Z},\mathcal{B}(\mathcal{Z})),d_{\mathrm{BL}})$  is a compact metric space and hence this space is separable and totally bounded.

Under the assumptions of the theorem, the SVM operator  $S: \mathcal{M}_1(\mathcal{Z}, \mathcal{B}(\mathcal{Z})) \to$  $H, S(P) = f_{L^*,P,\lambda}$ , is well-defined because S(P) exists and is unique for all  $P \in \mathcal{M}_1(\mathcal{Z}, \mathcal{B}(\mathcal{Z}))$  and for all  $\lambda \in (0, \infty)$ , see [2, Thm. 5, Thm. 6]. Furthermore, S is continuous with respect to the combination of the weak topology on  $\mathcal{M}_1(\mathcal{Z},\mathcal{B}(\mathcal{Z}))$  and the norm topology on H, see [11, Thm. 3.3]. Hence the function

$$g_{\mathcal{P}}: \mathcal{X} \times \mathcal{Y} \to \mathbb{R}, \quad g_{\mathcal{P}}(x, y) := L^{\star}(x, y, S(\mathcal{P})(x)) = L^{\star}(x, y, f_{L^{\star}, \mathcal{P}, \lambda}(x))$$

$$(4.29)$$

is well-defined. Because the kernel k is bounded and continuous, all functions  $f \in H$ , and hence in particular  $S(P) = f_{L^*,P,\lambda} \in H$ , are continuous, see e.g. [15, Lem. 4.28, Lem. 4.29]. Hence the function  $g_P$  is continuous (with respect to (x,y), because the loss function L and hence the shifted loss function  $L^*(x,y,t) = L(x,y,t) - L(x,y,0), (x,y,t) \in \mathcal{X} \times \mathcal{Y} \times \mathbb{R}$ , are continuous. Furthermore, the function  $g_P$  is bounded, because  $(\mathcal{Z}, d_{\mathcal{Z}})$  with  $\mathcal{Z} := \mathcal{X} \times \mathcal{Y}$  is

by assumption a compact metric space, the Lipschitz continuous loss function L maps from  $\mathcal{X} \times \mathcal{Y} \times \mathbb{R}$  to  $[0, \infty)$ , and  $||S(P)||_{\infty} \leq \frac{1}{\lambda}|L|_1 ||k||_{\infty}^2 < \infty$ , see [2, p. 314, (17)]. Hence  $g_P \in \mathcal{C}_b(\mathcal{Z}, \mathbb{R})$ . Because the bounded Lipschitz metric  $d_{\text{BL}}$  metrizes the weak topology on  $\mathcal{M}_1(\mathcal{Z}, \mathcal{B}(\mathcal{Z}))$ , it follows that

$$\forall \varepsilon_1 > 0 \; \exists \, \delta_1 > 0 : \quad d_{\mathrm{BL}}(\mathbf{Q}, \mathbf{P}) < \delta_1 \quad \Longrightarrow \quad \left| \int g_{\mathbf{P}} \, d\mathbf{Q} - \int g_{\mathbf{P}} \, d\mathbf{P} \right| < \varepsilon_1.$$

$$(4.30)$$

Recall that  $S: \mathcal{M}_1(\mathcal{Z}, \mathcal{B}(\mathcal{Z})) \to H$  is continuous with respect to the combination of the weak topology on  $\mathcal{M}_1(\mathcal{Z}, \mathcal{B}(\mathcal{Z}))$  and the norm topology on H, see [11, Thm. 3.3]. Hence

$$\forall \varepsilon_2 > 0 \ \exists \delta_2 > 0 : d_{BL}(Q, P) < \delta_2 \implies \|S(Q) - S(P)\|_H < \varepsilon_2.$$
 (4.31)

Fix  $\varepsilon > 0$ . Define

$$\varepsilon_1 := \frac{\varepsilon}{3}$$
 and  $\varepsilon_2 := \frac{\varepsilon}{3|L|_1 ||k||_{\infty}}$ .

Using the triangle inequality in (4.33), the definition of the shifted loss function  $L^*$  in (4.34), the definition of the function  $g_P$  in (4.35), the Lipschitz continuity of L in (4.36), and the well-known formula

$$||f||_{\infty} \le ||k||_{\infty} ||f||_{H}, \qquad f \in H,$$
 (4.32)

see e.g. [15, p. 124] we obtain that  $d_{\rm BL}({\rm Q,P}) < \delta_2$  implies

$$|R(Q) - R(P)| = \left| \int L^{*}(x, y, S(Q)(x)) dQ(x, y) - \int L^{*}(x, y, S(P)(x)) dP(x, y) \right|$$

$$\leq \left| \int L^{*}(x, y, S(Q)(x)) dQ(x, y) - \int L^{*}(x, y, S(P)(x)) dQ(x, y) \right| (4.33)$$

$$+ \left| \int L^{*}(x, y, S(P)(x)) dQ(x, y) - \int L^{*}(x, y, S(P)(x)) dP(x, y) \right|$$

$$\leq \int |L(x, y, S(Q)(x)) - L(x, y, S(P)(x))| dQ(x, y)$$

$$+ \left| \int g_{P} dQ - \int g_{P} dP \right|$$

$$(4.35)$$

$$\stackrel{\text{(4.30)}}{\leq} |L|_1 \|S(Q) - S(P)\|_{\infty} + \varepsilon_1 \tag{4.36}$$

$$\stackrel{\text{(4.32)}}{\leq} |L|_1 ||k||_{\infty} ||S(\mathbf{Q}) - S(\mathbf{P})||_H + \varepsilon_1 \tag{4.37}$$

$$\stackrel{\text{(4.31)}}{\leq} |L|_1 ||k||_{\infty} \varepsilon_2 + \varepsilon_1 = \frac{2}{3} \varepsilon. \tag{4.38}$$

Hence, R is continuous with respect to the combination of the weak topology on  $\mathcal{M}_1(\mathcal{Z}, \mathcal{B}(\mathcal{Z}))$  and the standard topology on  $(\mathbb{R}, \mathcal{B})$ .

Step 2. Because  $(\mathcal{M}_1(\mathcal{Z}, \mathcal{B}(\mathcal{Z})), d_{BL})$  is a compact metric space and the risk functional  $R: \mathcal{M}_1(\mathcal{Z}, \mathcal{B}(\mathcal{Z})) \to \mathbb{R}$  is continuous, R is even uniformly continuous with respect to the mentioned topologies, see [5, Prop. 1.5.9]. Obviously  $(\mathcal{W}, d_{\mathcal{W}}) := (\mathbb{R}, |\cdot|)$  is a complete separable metric space. Therefore, Theorem 2 yields that the sequence of  $\mathbb{R}$ -valued statistics

$$R_n((X_1, Y_1), \dots, (X_n, Y_n)) = \frac{1}{n} \sum_{i=1}^n L^*(X_i, Y_i, f_{L^*, D, \lambda}(X_i)), \quad n \in \mathbb{N},$$

where  $f_{L^*,D,\lambda} := \arg\min_{f \in H} \frac{1}{n} \sum_{j=1}^n L^*(X_j, Y_j, f(X_j)) + \lambda \|f\|_H^2$ , is uniformly qualitatively robust in a neighborhood  $\mathcal{U}(P_0)$  for every probability measure  $P_0 \in \mathcal{M}_1(\mathcal{Z})$ . Now we apply Theorem 3, which yields that the sequence  $\mathfrak{L}_n(R; P_n)$  of bootstrap SVM estimators of  $\mathfrak{L}_n(R; P)$  is qualitatively robust for all  $P_0 \in \mathcal{M}_1(\mathcal{Z}, \mathcal{B}(\mathcal{Z}))$ , which completes the proof.

### References

- [1] A. Christmann and I. Steinwart. Consistency and robustness of kernel based regression. *Bernoulli*, 13:799–819, 2007.
- [2] A. Christmann, A. Van Messem, and I. Steinwart. On consistency and robustness properties of support vector machines for heavy-tailed distributions. *Statistics and Its Interface*, 2:311–327, 2009.
- [3] A. Cuevas. Qualitative robustness in abstract inference. *J. Statist. Plann. Inference*, 18:277–289, 1988.
- [4] A. Cuevas and R. Romo. On robustness properties of bootstrap approximations. *J. Statist. Plann. Inference*, 1993.
- [5] Z. Denkowski, S. Migórski, and N. Papageorgiou. An introduction to nonlinear analysis: Theory. Kluwer Academic Publishers, Boston, 2003.
- [6] R. M. Dudley. *Uniform Central Limit Theorems*. Cambridge University Press, Cambridge, 1999.
- [7] R. M. Dudley. *Real Analysis and Probability*. Cambridge University Press, Cambridge, 2002.
- [8] R. M. Dudley, E. Giné, and J. Zinn. Uniform and universal Glivenko-Cantelli classes. J. Theor. Prob., 4:485–510, 1991.
- [9] B. Efron. Bootstrap methods: Another look at the jackknife. *Annals of Statistics*, 7:1–26, 1979.

- [10] B. Efron. The Jackknife, the Bootstrap, and Other Resampling Plans, volume 38. CBMS Monograph, Society for Industrial and Applied Mathematics, Philadelphia, 1982.
- [11] R. Hable and A. Christmann. Qualitative robustness of support vector machines. *Journal of Multivariate Analysis*, 102:993–1007, 2011.
- [12] F. R. Hampel. A general qualitative definition of robustness. *Ann. Math. Statist.*, 42:1887–1896, 1971.
- [13] P. J. Huber. Robust Statistics. John Wiley & Sons, New York, 1981.
- [14] K. R. Parthasarathy. *Probability Measures on Metric Spaces*. Academic Press, New York, 1967.
- [15] I. Steinwart and A. Christmann. Support Vector Machines. Springer, New York, 2008.
- [16] M. Talagrand. The Glivenko-Cantelli problem. Ann. Probability, 15:837–870, 1987.