

Supplementary Material for: Optimization by Gradient Boosting

G erard Biau and Beno t Cadre

Proof (Theorem 1) Assume that, for some $t_0 \geq 0$, $\sup_{f \in \mathcal{F}} -\mathbb{E}\xi(F_{t_0}(X), Y)f(X) = 0$. Then, by the symmetry of the class \mathcal{F} , for all $f \in \mathcal{F}$, $\mathbb{E}\xi(F_{t_0}(X), Y)f(X) = 0$. We conclude by technical Lemma 2 that

$$C(F_t) = \inf_{F \in \text{lin}(\mathcal{F})} C(F) \quad \text{for all } t \geq t_0,$$

and the result is proved. Thus, in the following, it is assumed that

$$\sup_{f \in \mathcal{F}} -\mathbb{E}\xi(F_t(X), Y)f(X) > 0 \quad \text{for all } t \geq 0.$$

Consequently, $-\mathbb{E}\xi(F_t(X), Y)f_{t+1}(X) > 0$ and $w_t > 0$ for all t . Since $w_t \rightarrow 0$ (by Lemma 1 of the Main Document), there exists a subsequence $(w_{t'})_{t'}$ such that

$$\begin{aligned} w_{t'+1} &= -(2L)^{-1} \mathbb{E}\xi(F_{t'}(X), Y)f_{t'+1}(X) \\ &= (2L)^{-1} \sup_{f \in \mathcal{F}} -\mathbb{E}\xi(F_{t'}(X), Y)f(X). \end{aligned} \quad (1)$$

Let $\varepsilon > 0$. For all t' large enough and all $f \in \mathcal{F}$, by the symmetry of \mathcal{F} ,

$$-\mathbb{E}\xi(F_{t'}(X), Y)f(X) \leq \varepsilon \quad \text{and} \quad \mathbb{E}\xi(F_{t'}(X), Y)f(X) \leq \varepsilon,$$

and thus $\lim_{t' \rightarrow \infty} \mathbb{E}\xi(F_{t'}(X), Y)f(X) = 0$ for all $f \in \mathcal{F}$. We conclude that, for all $G \in \text{lin}(\mathcal{F})$,

$$\lim_{t' \rightarrow \infty} \mathbb{E}\xi(F_{t'}(X), Y)G(X) = 0. \quad (2)$$

G erard Biau
Sorbonne Universit , CNRS, LPSM, 4 place Jussieu, 75005 Paris, France e-mail: gerard.biau@sorbonne-universite.fr

Beno t Cadre
Univ Rennes, CNRS, IRMAR - UMR 6625, F-35000 Rennes, France e-mail: benoit.cadre@univ-rennes2.fr

Assume, without loss of generality, that $F_0 = 0$, and observe that $F_t = \sum_{k=1}^t w_k f_k$. Thus, we may write

$$\begin{aligned} \mathbb{E}\xi(F_{t'}(X), Y)F_{t'}(X) &= \sum_{k=1}^{t'} w_k \mathbb{E}\xi(F_{t'}(X), Y)f_k(X) \\ &\leq \sup_{f \in \mathcal{F}} \mathbb{E}\xi(F_{t'}(X), Y)f(X) \sum_{k=1}^{t'} w_k \\ &= \sup_{f \in \mathcal{F}} -\mathbb{E}\xi(F_{t'}(X), Y)f(X) \sum_{k=1}^{t'} w_k \\ &\quad \text{(by the symmetry of } \mathcal{F} \text{)} \\ &= 2Lw_{t'+1} \sum_{k=1}^{t'} w_k, \end{aligned}$$

by definition of $w_{t'+1}$ —see (1). So,

$$\mathbb{E}\xi(F_{t'}(X), Y)F_{t'}(X) \leq 2Lw_{t'} \sum_{k=1}^{t'} w_k = 2Lw_{t'} \sum_{k=1}^{t'} w_k^{-1} w_k^2$$

(because $w_{t'+1} \leq w_{t'}$).

Since $\sum_{k \geq 1} w_k^2 < \infty$, and since the sequence $(w_t)_t$ is nonincreasing, positive, and tends to 0 as $t \rightarrow \infty$, Kronecker's lemma reveals that $w_{t'} \sum_{k=1}^{t'} w_k^{-1} w_k^2 \rightarrow 0$ as $t' \rightarrow \infty$. Therefore,

$$\limsup_{t' \rightarrow \infty} \mathbb{E}\xi(F_{t'}(X), Y)F_{t'}(X) \leq 0. \quad (3)$$

Let $\varepsilon > 0$ and let $F_\varepsilon^* \in \text{lin}(\mathcal{F})$ be such that

$$\inf_{F \in \text{lin}(\mathcal{F})} C(F) \geq C(F_\varepsilon^*) - \varepsilon.$$

By the convexity of C , we have, for all t' ,

$$\begin{aligned} \inf_{F \in \text{lin}(\mathcal{F})} C(F) &\geq C(F_\varepsilon^*) - \varepsilon \\ &\geq C(F_{t'}) + \mathbb{E}\xi(F_{t'}(X), Y)(F_\varepsilon^*(X) - F_{t'}(X)) - \varepsilon \\ &\geq \inf_k C(F_k) + \mathbb{E}\xi(F_{t'}(X), Y)F_\varepsilon^*(X) - \mathbb{E}\xi(F_{t'}(X), Y)F_{t'}(X) - \varepsilon. \end{aligned}$$

Combining (2) and (3), we conclude that $\inf_{F \in \text{lin}(\mathcal{F})} C(F) \geq \inf_k C(F_k) - \varepsilon$ for all $\varepsilon > 0$, so that

$$\lim_{t \rightarrow \infty} C(F_t) = \inf_k C(F_k) = \inf_{F \in \text{lin}(\mathcal{F})} C(F),$$

which is the desired result. \square

Proof (Theorem 2) The first step is to establish that there exists a subsequence $(F_{t''})_{t''}$ such that $\lim_{t'' \rightarrow \infty} \mathbb{E} \xi(F_{t''}(X), Y) G(X) \rightarrow 0$ for all $G \in \text{lin}(\mathcal{P})$. We start by observing that, by Lemma 2 of the Main Document, $C(F_t) \leq C(F_0)$. Thus, by technical Lemma 3, $\sup_t \|F_t\|_{\mu_X} \leq B$ for some positive constant B . Now,

$$\begin{aligned} & |\mathbb{E} \xi(F_t(X), Y) f_{t+1}(X)| \\ &= |\mathbb{E} \mathbb{E}(\xi(F_t(X), Y) | X) f_{t+1}(X)| \\ &\leq \mathbb{E} |\mathbb{E}(\xi(F_t(X), Y) - \xi(0, Y) | X)| \cdot |f_{t+1}(X)| + \mathbb{E} |\xi(0, Y) f_{t+1}(X)| \\ &\leq L \mathbb{E} |F_t(X) f_{t+1}(X)| + \mathbb{E} |\xi(0, Y) f_{t+1}(X)| \\ &\quad \text{(by Assumption A}_3\text{)}. \end{aligned}$$

So,

$$\begin{aligned} |\mathbb{E} \xi(F_t(X), Y) f_{t+1}(X)| &\leq L \|F_t\|_{\mu_X} \|f_{t+1}\|_{\mu_X} + (\mathbb{E} \xi(0, Y)^2)^{1/2} \|f_{t+1}\|_{\mu_X} \\ &\quad \text{(by the Cauchy-Schwarz inequality)} \\ &\leq (LB + (\mathbb{E} \xi(0, Y)^2)^{1/2}) \|f_{t+1}\|_{\mu_X}. \end{aligned}$$

Consequently, since $\lim_{t \rightarrow \infty} \|f_{t+1}\|_{\mu_X} = 0$ (by Lemma 2 of the Main Document),

$$\begin{aligned} \inf_{f \in \mathcal{P}} (2\mathbb{E} \xi(F_t(X), Y) f(X) + \|f\|_{\mu_X}^2) &= 2\mathbb{E} \xi(F_t(X), Y) f_{t+1}(X) + \|f_{t+1}\|_{\mu_X}^2 \\ &\rightarrow 0 \text{ as } t \rightarrow \infty. \end{aligned}$$

Accordingly, by the symmetry of \mathcal{P} , for all $\varepsilon > 0$ and all t large enough, we have, for all $f \in \mathcal{P}$,

$$2\mathbb{E} \xi(F_t(X), Y) f(X) + \|f\|_{\mu_X}^2 \geq -\varepsilon \quad \text{and} \quad -2\mathbb{E} \xi(F_t(X), Y) f(X) + \|f\|_{\mu_X}^2 \geq -\varepsilon.$$

So, for all t large enough and all $f \in \mathcal{P}$,

$$|2\mathbb{E} \xi(F_t(X), Y) f(X)| \leq \varepsilon + \|f\|_{\mu_X}^2.$$

Since ε was arbitrary, we conclude that, for all $f \in \mathcal{P}$,

$$2 \limsup_{t \rightarrow \infty} |\mathbb{E} \xi(F_t(X), Y) f(X)| \leq \|f\|_{\mu_X}^2. \quad (4)$$

On the other hand, by Assumption A₃,

$$|\mathbb{E}(\xi(F_t(X), Y) | X)| \leq \mathbb{E}(|\xi(0, Y)| | X) + L|F_t(X)|.$$

Since $\sup_t \|F_t\|_{\mu_X} < \infty$, we deduce that

$$\sup_t \|\mathbb{E}(\xi(F_t(X), Y) | X = \cdot)\|_{\mu_X} < \infty.$$

Next, since $\sum_{k \geq 1} \|f_k\|_{\mu_X}^2 < \infty$, there exists a subsequence $(f_{t'})_{t'}$ satisfying $t' \|f_{t'+1}\|_{\mu_X}^2 \rightarrow 0$. Besides, recalling that the unit ball of $L^2(\mu_X)$ is weakly compact, there exists a subsequence $(F_{t''})_{t''}$ of $(F_{t'})_{t'}$ and $\tilde{F} \in L^2(\mu_X)$ such that, for all $G \in \text{lin}(\mathcal{P})$,

$$\mathbb{E}\xi(F_{t''}(X), Y)G(X) = \mathbb{E}\mathbb{E}(\xi(F_{t''}(X), Y) | X)G(X) \rightarrow \mathbb{E}\tilde{F}(X)G(X).$$

Combining this identity with (4) reveals that $2|\mathbb{E}\tilde{F}(X)f(X)| \leq \|f\|_{\mu_X}^2$ for all $f \in \mathcal{P}$. In particular, for all $\varepsilon > 0$ and all $f \in \mathcal{P}$, $2|\mathbb{E}\tilde{F}(X)\varepsilon f(X)| \leq \varepsilon^2 \|f\|_{\mu_X}^2$, and thus, letting $\varepsilon \downarrow 0$, we find that $\mathbb{E}\tilde{F}(X)f(X) = 0$ for all $f \in \mathcal{P}$. By a linearity argument, we conclude that $\mathbb{E}\tilde{F}(X)G(X) = 0$ for all $G \in \text{lin}(\mathcal{P})$. Therefore, for all $G \in \text{lin}(\mathcal{P})$,

$$\lim_{t'' \rightarrow \infty} \mathbb{E}\xi(F_{t''}(X), Y)G(X) = 0, \quad (5)$$

which was our first objective.

The next step is to prove that $\limsup_{t'' \rightarrow \infty} \mathbb{E}\xi(F_{t''}(X), Y)F_{t''}(X) \leq 0$. To simplify the notation, we assume, without loss of generality, that $F_0 = 0$. Fix $\varepsilon > 0$. Since $\sum_{k \geq 1} \|f_k\|_{\mu_X}^2 < \infty$, there exists $T \geq 0$ such that $\sum_{k \geq T+1} \|f_k\|_{\mu_X}^2 \leq \varepsilon$. In addition, for all $t > T$, $F_t = F_T + \nu \sum_{k=T+1}^t f_k$, so that

$$\mathbb{E}\xi(F_t(X), Y)F_t(X) = \mathbb{E}\xi(F_t(X), Y)F_T(X) + \nu \sum_{k=T+1}^t \mathbb{E}\xi(F_t(X), Y)f_k(X). \quad (6)$$

Also, by the very definition of f_{t+1} and the symmetry of \mathcal{P} , we have, for all $f \in \mathcal{P}$,

$$2\mathbb{E}\xi(F_t(X), Y)f_{t+1}(X) + \|f_{t+1}\|_{\mu_X}^2 \leq -2\mathbb{E}\xi(F_t(X), Y)f(X) + \|f\|_{\mu_X}^2, \quad (7)$$

i.e., for all $f \in \mathcal{P}$,

$$2\mathbb{E}\xi(F_t(X), Y)f(X) \leq -2\mathbb{E}\xi(F_t(X), Y)f_{t+1}(X) - \|f_{t+1}\|_{\mu_X}^2 + \|f\|_{\mu_X}^2.$$

Using (6), this leads to

$$\begin{aligned} & \mathbb{E}\xi(F_t(X), Y)F_t(X) \\ & \leq \mathbb{E}\xi(F_t(X), Y)F_T(X) \\ & \quad + \frac{\nu}{2} \left(t \left(-2\mathbb{E}\xi(F_t(X), Y)f_{t+1}(X) - \|f_{t+1}\|_{\mu_X}^2 \right) + \sum_{k \geq T+1} \|f_k\|_{\mu_X}^2 \right) \\ & \leq \frac{\varepsilon\nu}{2} + \mathbb{E}\xi(F_t(X), Y)F_T(X) + \frac{\nu t}{2} \left(-2\mathbb{E}\xi(F_t(X), Y)f_{t+1}(X) - \|f_{t+1}\|_{\mu_X}^2 \right). \end{aligned} \quad (8)$$

But, according to inequality (7) applied with $f = -2f_{t+1}$ (which belongs to \mathcal{P} by assumption),

$$2\mathbb{E}\xi(F_t(X), Y)f_{t+1}(X) + \|f_{t+1}\|_{\mu_X}^2 \leq 4\mathbb{E}\xi(F_t(X), Y)f_{t+1}(X) + 4\|f_{t+1}\|_{\mu_X}^2,$$

i.e.,

$$-2\mathbb{E}\xi(F_t(X), Y)f_{t+1}(X) \leq 3\|f_{t+1}\|_{\mu_X}^2.$$

Combining this inequality with (8) shows that

$$\mathbb{E}\xi(F_t(X), Y)F_t(X) \leq \frac{\varepsilon\nu}{2} + \mathbb{E}\xi(F_t(X), Y)F_T(X) + \nu t\|f_{t+1}\|_{\mu_X}^2.$$

Since $F_T \in \text{lin}(\mathcal{F})$, we know from (5) that $\mathbb{E}\xi(F_{t''}(X), Y)F_T(X) \rightarrow 0$. Therefore, recalling that $t''\|f_{t''+1}\|_{\mu_X}^2 \rightarrow 0$, for all $\varepsilon > 0$,

$$\limsup_{t'' \rightarrow \infty} \mathbb{E}\xi(F_{t''}(X), Y)F_{t''}(X) \leq \frac{\varepsilon\nu}{2}.$$

Since ε is arbitrary, we have just shown that

$$\limsup_{t'' \rightarrow \infty} \mathbb{E}\xi(F_{t''}(X), Y)F_{t''}(X) \leq 0, \quad (9)$$

as desired.

Let $\varepsilon > 0$ and let $F_\varepsilon^* \in \text{lin}(\mathcal{F})$ be such that

$$\inf_{F \in \text{lin}(\mathcal{F})} C(F) \geq C(F_\varepsilon^*) - \varepsilon.$$

By the convexity of C , along t'' ,

$$\begin{aligned} \inf_{F \in \text{lin}(\mathcal{F})} C(F) &\geq C(F_\varepsilon^*) - \varepsilon \\ &\geq \inf_k C(F_k) + \mathbb{E}\xi(F_{t''}(X), Y)F_\varepsilon^*(X) - \mathbb{E}\xi(F_{t''}(X), Y)F_{t''}(X) - \varepsilon. \end{aligned}$$

Putting (5) and (9) together, we conclude that

$$\lim_{t \rightarrow \infty} C(F_t) = \inf_k C(F_k) = \inf_{F \in \text{lin}(\mathcal{F})} C(F).$$

□

Proof (Theorem 3) For $\beta \in \mathbb{R}^N$, we let $F_\beta = \sum_{j=1}^N \beta_j \mathbb{1}_{A_j^n}$ and notice that $\bar{F}_n = F_\alpha$ for some (data-dependent) $\alpha \in \mathbb{R}^N$. Let the event S be defined by

$$S = \{\forall j = 1, \dots, N : P_n(A_j^n) \geq P(A_j^n)/2\}.$$

Observe that

$$\|\bar{F}_n\|_{P_n}^2 \leq \frac{\frac{1}{n} \sum_{i=1}^n \phi(0, Y_i)}{\gamma_n} \leq \frac{\bar{\phi}}{\gamma_n},$$

and, similarly, that

$$\|\bar{F}_n\|_{P_n}^2 = \sum_{j=1}^N \alpha_j^2 P_n(A_j^n).$$

Therefore, on S ,

$$\frac{1}{2} \sum_{j=1}^N \alpha_j^2 P(A_j^n) \leq \frac{\bar{\phi}}{\gamma_n},$$

and so

$$\frac{\inf_{\mathcal{X}} g}{2} \cdot \nu_n \sum_{j=1}^N \alpha_j^2 \leq \frac{\bar{\phi}}{\gamma_n}.$$

We have just shown that, on the event S , $\alpha \in T$, where

$$T = \left\{ \beta \in \mathbb{R}^N : \sum_{j=1}^N \beta_j^2 \leq \frac{2\bar{\phi}}{\inf_{\mathcal{X}} g} \cdot \frac{1}{\nu_n \gamma_n} \right\}.$$

Now, observe that

$$\begin{aligned} \mathbb{E} C_n(\bar{F}_n) &= \mathbb{E} \inf_{F \in \text{lin}(\mathcal{F}_n)} C_n(F) \\ &= \mathbb{E} \inf_{F \in \text{lin}(\mathcal{F}_n)} C_n(F) \mathbb{1}_S + \mathbb{E} \inf_{F \in \text{lin}(\mathcal{F}_n)} C_n(F) \mathbb{1}_{S^c} \\ &\leq \mathbb{E} \inf_{F \in \text{lin}(\mathcal{F}_n)} C_n(F) \mathbb{1}_S + \mathbb{E} C_n(0) \mathbb{1}_{S^c} \\ &= \mathbb{E} \inf_{\beta \in T} C_n(F_\beta) \mathbb{1}_S + \mathbb{E} A_n(0) \mathbb{1}_{S^c} \\ &\leq \mathbb{E} \inf_{\beta \in T} C_n(F_\beta) + \bar{\phi} \mathbb{P}(S^c). \end{aligned}$$

Define

$$D_n(F) = A(F) + \gamma_n \|F\|_{P_n}^2.$$

Since $C_n(F) - D_n(F) = A_n(F) - A(F)$, we deduce from Lemma 4 and Lemma 6 that whenever

$$\frac{\log N}{n \nu_n} \rightarrow 0 \quad \text{and} \quad \frac{1}{\sqrt{n \nu_n \gamma_n}} \zeta \left(\sqrt{\frac{2\bar{\phi}}{\nu_n \gamma_n \inf_{\mathcal{X}} g}} \right) \rightarrow 0,$$

we have

$$\begin{aligned} \limsup_{n \rightarrow \infty} \mathbb{E} C_n(\bar{F}_n) &\leq \limsup_{n \rightarrow \infty} \mathbb{E} \inf_{\beta \in T} C_n(F_\beta) \\ &\leq \limsup_{n \rightarrow \infty} \mathbb{E} \inf_{\beta \in T} D_n(F_\beta) + \limsup_{n \rightarrow \infty} (\mathbb{E} \sup_{\beta \in T} |A_n(F_\beta) - A(F_\beta)|) \\ &= \limsup_{n \rightarrow \infty} \mathbb{E} \inf_{\beta \in T} D_n(F_\beta). \end{aligned} \tag{10}$$

Let $\varepsilon > 0$. By Lemma 5, there exists $(\beta_1^\varepsilon, \dots, \beta_N^\varepsilon) \in T$ such that

$$\left\| F^\star - \sum_{j=1}^N \beta_j^\varepsilon \mathbb{1}_{A_j^n} \right\|_P \leq \varepsilon.$$

Define $F_\varepsilon^\star = \sum_{j=1}^N \beta_j^\varepsilon \mathbb{1}_{A_j^n}$. Then, according to (10),

$$\begin{aligned}
\limsup_{n \rightarrow \infty} \mathbb{E}C_n(\bar{F}_n) &\leq \limsup_{n \rightarrow \infty} (A(F_\varepsilon^*) + \gamma_n \mathbb{E}\|F_\varepsilon^*\|_{P_n}^2) \\
&= \limsup_{n \rightarrow \infty} (A(F_\varepsilon^*) + \gamma_n \|F_\varepsilon^*\|_P^2) \\
&\leq A(F_\varepsilon^*).
\end{aligned} \tag{11}$$

Since A is continuous, we conclude that $\limsup_{n \rightarrow \infty} \mathbb{E}C_n(\bar{F}_n) \leq A(F^*)$.

On the other hand, $C_n(\bar{F}_n) \geq A_n(\bar{F}_n)$, and, by Lemma 4 and Lemma 6,

$$\begin{aligned}
\mathbb{E}|A_n(\bar{F}_n) - A(\bar{F}_n)| &\leq \mathbb{E} \sup_{\beta \in T} |A_n(F_\beta) - A(F_\beta)| + \bar{\phi} \mathbb{P}(S^c) \\
&\rightarrow 0 \quad \text{as } n \rightarrow \infty.
\end{aligned}$$

Therefore,

$$\limsup_{n \rightarrow \infty} \mathbb{E}A(\bar{F}_n) \leq \limsup_{n \rightarrow \infty} C_n(\bar{F}_n).$$

So, with (11),

$$\limsup_{n \rightarrow \infty} \mathbb{E}A(\bar{F}_n) \leq A(F^*),$$

which is the desired result. \square

Some technical lemmas

Lemma 1 Assume that Assumptions \mathbf{A}_1 and \mathbf{A}_3 are satisfied. Then, for all $a > 0$ and all $F, G \in L^2(\mu_X)$,

$$C(F) - C(F + aG) \geq -a^2 L \|G\|_{\mu_X}^2 - a \mathbb{E} \xi(F(X), Y) G(X).$$

Proof By inequality (1) of the Main Document,

$$\begin{aligned}
C(F) &\geq C(F + aG) - a \mathbb{E} \xi(F(X) + aG(X), Y) G(X) \\
&= C(F + aG) - a \mathbb{E} (\xi(F(X) + aG(X), Y) - \xi(F(X), Y)) G(X) \\
&\quad - a \mathbb{E} \xi(F(X), Y) G(X) \\
&= C(F + aG) - a \mathbb{E} \mathbb{E} (\xi(F(X) + aG(X), Y) - \xi(F(X), Y) | X) G(X) \\
&\quad - a \mathbb{E} \xi(F(X), Y) G(X) \\
&\geq C(F + aG) - a (\mathbb{E} \mathbb{E}^2 (\xi(F(X) + aG(X), Y) - \xi(F(X), Y) | X))^{1/2} \|G\|_{\mu_X} \\
&\quad - a \mathbb{E} \xi(F(X), Y) G(X) \\
&\quad \text{(by the Cauchy-Schwarz inequality)}.
\end{aligned}$$

Thus, by Assumption \mathbf{A}_3 ,

$$C(F) \geq C(F + aG) - a^2L\|G\|_{\mu_X}^2 - a\mathbb{E}\xi(F(X), Y)G(X).$$

□

Lemma 2 Assume that Assumption **A**₁ is satisfied, and let $(F_t)_t$ be defined by [Algorithm 1](#) with $(w_t)_t$ as in (8) of the Main Document. If, for some $t_0 \geq 0$,

$$\mathbb{E}\xi(F_{t_0}(X), Y)f(X) = 0 \quad \text{for all } f \in \mathcal{F},$$

then $C(F_{t_0}) = \inf_{F \in \text{lin}(\mathcal{F})} C(F)$.

Proof Fix $t_0 \geq 0$ and assume that $\mathbb{E}\xi(F_{t_0}(X), Y)f(X) = 0$ for all $f \in \mathcal{F}$. By linearity, $\mathbb{E}\xi(F_{t_0}(X), Y)G(X) = 0$ for all $G \in \text{lin}(\mathcal{F})$. Let $\varepsilon > 0$ and let $F_\varepsilon^* \in \text{lin}(\mathcal{F})$ be such that

$$\inf_{F \in \text{lin}(\mathcal{F})} C(F) \geq C(F_\varepsilon^*) - \varepsilon.$$

By the convexity inequality (1) of the Main Document,

$$C(F_\varepsilon^*) \geq C(F_0) + \mathbb{E}\xi(F_0(X), Y)(F_\varepsilon^*(X) - F_0(X)) = C(F_0).$$

Thus,

$$\inf_{F \in \text{lin}(\mathcal{F})} C(F) \geq C(F_0) - \varepsilon.$$

Since ε is arbitrary, the result follows. □

Lemma 3 Assume that Assumptions **A**₁ and **A**₂ are satisfied. Then, for all $F \in L^2(\mu_X)$,

$$\|F\|_{\mu_X} \leq \frac{2}{\alpha} (\mathbb{E}\xi(0, Y)^2)^{1/2} + \sqrt{\frac{2C(F)}{\alpha}}.$$

Proof By inequality (2) of the Main Document and the Cauchy-Schwarz inequality,

$$\begin{aligned} C(F) &\geq C(0) + \mathbb{E}\xi(0, Y)F(X) + \frac{\alpha}{2} \|F\|_{\mu_X}^2 \\ &\geq C(0) - (\mathbb{E}\xi(0, Y)^2)^{1/2} \|F\|_{\mu_X} + \frac{\alpha}{2} \|F\|_{\mu_X}^2. \end{aligned}$$

Let $\kappa = (\mathbb{E}\xi(0, Y)^2)^{1/2}$. Since $C(0) \geq 0$,

$$C(F) + \kappa \|F\|_{\mu_X} - \frac{\alpha}{2} \|F\|_{\mu_X}^2 \geq 0.$$

Therefore,

$$\|F\|_{\mu_X} \leq \frac{\kappa + \sqrt{\kappa^2 + 2\alpha C(F)}}{\alpha} \leq \frac{2\kappa}{\alpha} + \sqrt{\frac{2C(F)}{\alpha}}.$$

□

Lemma 4 Let the event S be defined by

$$S = \{\forall j = 1, \dots, N : P_n(A_j^n) \geq P(A_j^n)/2\}.$$

If $\frac{\log N}{nv_n} \rightarrow 0$, then $\lim_{n \rightarrow \infty} \mathbb{P}(S^c) = 0$.

Proof We have

$$\begin{aligned}
\mathbb{P}(S^c) &= \mathbb{P}(\exists j \leq N : P_n(A_j^n) < P(A_j^n)/2) \\
&= \mathbb{P}(\exists j \leq N : P_n(A_j^n) - P(A_j^n) < -P(A_j^n)/2) \\
&= \mathbb{P}\left(\exists j \leq N : \frac{P(A_j^n) - P_n(A_j^n)}{\sqrt{P(A_j^n)}} > \sqrt{P(A_j^n)/2}\right) \\
&\leq \mathbb{P}\left(\max_{1 \leq j \leq N} \frac{P(A_j^n) - P_n(A_j^n)}{\sqrt{P(A_j^n)}} > \sqrt{v_n \inf_{\mathcal{X}} g/2}\right) \\
&\leq c_1 N e^{-nv_n \inf_{\mathcal{X}} g/c_2},
\end{aligned}$$

where c_1 and c_2 are positive constants. In the last inequality, we used a Vapnik-Chervonenkis inequality [Vapnik, 1988] for relative deviations. \square

In the sequel, we let

$$T = \left\{ \beta \in \mathbb{R}^N : \sum_{j=1}^N \beta_j^2 \leq \frac{2\bar{\phi}}{\inf_{\mathcal{X}} g} \cdot \frac{1}{v_n \gamma_n} \right\},$$

where $\bar{\phi} = \sup_{\mathcal{X}} \phi(0, y) < \infty$. We recall that $A^n(x) := A_j^n$ whenever $x \in A_j^n$.

Lemma 5 Assume that $\text{diam}(A^n(X)) \rightarrow 0$ in probability and that $\gamma_n \rightarrow 0$ as $n \rightarrow \infty$. For all $\varepsilon > 0$ and all n large enough, there exists $(\beta_1^\varepsilon, \dots, \beta_N^\varepsilon) \in T$ such that

$$\left\| F^\star - \sum_{j=1}^N \beta_j^\varepsilon \mathbb{1}_{A_j^n} \right\|_p \leq \varepsilon.$$

Proof Let K be a bounded and uniformly continuous function on \mathbb{R}^d , with $\int K d\lambda = 1$. Let, for $p > 0$,

$$K_p(x) = p^d K\left(\frac{x}{p}\right), \quad x \in \mathbb{R}^d.$$

With a slight abuse of notation, we consider F^\star as a function defined on the whole space \mathbb{R}^d (instead of \mathcal{X}) by implicitly assuming that $F^\star = 0$ on \mathcal{X}^c . We also define $F_p^\star = F^\star \star K_p$, i.e.,

$$F_p^\star(x) = \int_{\mathbb{R}^d} K_p(z) F^\star(x - z) dz, \quad x \in \mathbb{R}^d.$$

Let $(L^2(\lambda), \|\cdot\|_\lambda)$ be the vector space of all real-valued square integrable functions on \mathbb{R}^d . For all p large enough, we have

$$\|F_p^\star - F^\star\|_\lambda \leq \frac{\varepsilon}{2\sqrt{\sup_{\mathcal{X}} g}}$$

[see, e.g., [Wheeden and Zygmund, 1977](#), Theorem 9.6]. Therefore, for all p large enough,

$$\|F_p^\star - F^\star\|_P \leq \varepsilon/2. \quad (12)$$

In addition, F_p^\star is uniformly continuous on \mathcal{X} [[Wheeden and Zygmund, 1977](#), Theorem 9.4]. Thus, there exists $\eta = \eta(\varepsilon, p) > 0$ such that, for all $(x, x') \in \mathcal{X}^2$ with $\|x - x'\| \leq \eta$,

$$|F_p^\star(x) - F_p^\star(x')| \leq \varepsilon/\sqrt{8}.$$

For each $j \in \{1, \dots, N\}$, choose an arbitrary $a_j^n \in A_j^n$ and set $G_p^\star = \sum_{j=1}^N F_p^\star(a_j^n) \mathbb{1}_{A_j^n}$. Then

$$\begin{aligned} \|G_p^\star - F_p^\star\|_P^2 &= \sum_{j=1}^N \mathbb{E}(G_p^\star(X) - F_p^\star(X))^2 \mathbb{1}_{[X \in A_j^n, \text{diam}(A^n(X)) \leq \eta]} \\ &\quad + \sum_{j=1}^N \mathbb{E}(G_p^\star(X) - F_p^\star(X))^2 \mathbb{1}_{[X \in A_j^n, \text{diam}(A^n(X)) > \eta]} \\ &= \sum_{j=1}^N \mathbb{E}(F_p^\star(a_j^n) - F_p^\star(X))^2 \mathbb{1}_{[X \in A_j^n, \text{diam}(A^n(X)) \leq \eta]} \\ &\quad + \sum_{j=1}^N \mathbb{E}(G_p^\star(X) - F_p^\star(X))^2 \mathbb{1}_{[X \in A_j^n, \text{diam}(A^n(X)) > \eta]} \\ &\leq \frac{\varepsilon^2}{8} \sum_{j=1}^N \mathbb{P}(X \in A_j^n) + 4 \sup_{\mathcal{X}} (F^\star)^2 \sum_{j=1}^N \mathbb{P}(X \in A_j^n, \text{diam}(A^n(X)) > \eta) \\ &\quad (\text{since } \sup_{\mathcal{X}} |F_p^\star| \leq \sup_{\mathcal{X}} |F^\star| < \infty \text{ and } \sup_{\mathcal{X}} |G_p^\star| \leq \sup_{\mathcal{X}} |G^\star| < \infty) \\ &\leq \frac{\varepsilon^2}{8} + 4 \sup_{\mathcal{X}} (F^\star)^2 \mathbb{P}(\text{diam}(A^n(X)) > \eta), \end{aligned}$$

because the $(A_j^n)_{1 \leq j \leq N}$ form a partition of \mathcal{X} . Since $\text{diam}(A^n(X)) \rightarrow 0$ in probability, we see that for all n large enough (depending on ε and p),

$$\|G_p^\star - F_p^\star\|_P \leq \varepsilon/2.$$

Letting $\beta_j^\varepsilon = F_p^\star(a_j^n)$, $1 \leq j \leq N$, and combining this inequality and (12), we conclude that for every fixed $\varepsilon > 0$ and all n large enough, there exists $(\beta_1^\varepsilon, \dots, \beta_N^\varepsilon) \in \mathbb{R}^N$ such that

$$\left\| F^\star - \sum_{j=1}^N \beta_j^\varepsilon \mathbb{1}_{A_j^n} \right\|_P \leq \varepsilon.$$

To complete the proof, it remains to show that $(\beta_1^\varepsilon, \dots, \beta_N^\varepsilon) \in T$. Observe that

$$\sum_{j=1}^N (\beta_j^\varepsilon)^2 \leq \sup_{\mathcal{X}} (F^*)^2 N.$$

The right-hand side is bounded by $\frac{2\bar{\phi}}{\inf_{\mathcal{X}} g} \cdot \frac{1}{v_n \gamma_n}$ for all n large enough. To see this, just note that

$$N v_n \leq \sum_{j=1}^N \lambda(A_j^n) = \lambda(\mathcal{X}) < \infty.$$

Therefore, $N v_n \gamma_n \leq \lambda(\mathcal{X}) \gamma_n \rightarrow 0$ as $n \rightarrow \infty$. This concludes the proof of the lemma. \square

Lemma 6 For $\beta \in \mathbb{R}^N$, let $F_\beta = \sum_{j=1}^N \beta_j \mathbb{1}_{A_j^n}$. Assume that Assumption **A4** is satisfied. If

$$\frac{1}{\sqrt{n v_n \gamma_n}} \zeta \left(\sqrt{\frac{2\bar{\phi}}{v_n \gamma_n \inf_{\mathcal{X}} g}} \right) \rightarrow 0,$$

then

$$\lim_{n \rightarrow \infty} \mathbb{E} \sup_{\beta \in T} |A_n(F_\beta) - A(F_\beta)| = 0.$$

Proof Let

$$s_n = \sqrt{\frac{2\bar{\phi}}{v_n \gamma_n \inf_{\mathcal{X}} g}},$$

and let $\|\beta\|_\infty = \max_{1 \leq j \leq N} |\beta_j|$ be the supremum norm of $\beta = (\beta_1, \dots, \beta_N) \in \mathbb{R}^N$. By definition of T , we have, for all $\beta \in T$,

$$\sup_{\mathcal{X}} |F_\beta| = \sup_{\mathcal{X}} \left| \sum_{j=1}^N \beta_j \mathbb{1}_{A_j^n} \right| \leq \|\beta\|_\infty \leq s_n.$$

In addition, according to Assumption **A4**, we may write, for β_1 and $\beta_2 \in T$,

$$|\phi(F_{\beta_1}(x), y) - \phi(F_{\beta_2}(x), y)| \leq \zeta(s_n) |F_{\beta_1}(x) - F_{\beta_2}(x)| \leq \zeta(s_n) \|\beta_1 - \beta_2\|_\infty.$$

This shows that the process

$$\left(\frac{A_n(F_\beta) - A(F_\beta)}{\zeta(s_n)} \right)_{\beta \in T}$$

is subgaussian [e.g., [van Handel, 2016](#), Chapter 5] for the distance $d(\beta_1, \beta_2) = \frac{1}{\sqrt{n}} \|\beta_1 - \beta_2\|_\infty$. Now, let $N(T, d, \varepsilon)$ denote the ε -covering number of T for the distance d . Then, by Dudley's inequality [[van Handel, 2016](#), Corollary 5.25], one has

$$\begin{aligned}\mathbb{E} \sup_{\beta \in T} (A_n(F_\beta) - A(F_\beta)) &\leq 12\zeta(s_n) \int_0^\infty \sqrt{\log \left(N(T, \frac{1}{\sqrt{n}} \|\cdot\|_\infty, \varepsilon) \right)} d\varepsilon \\ &= 12\zeta(s_n) \cdot \frac{1}{\sqrt{n}} \int_0^\infty \sqrt{\log(N(T, \|\cdot\|_\infty, \varepsilon))} d\varepsilon.\end{aligned}$$

Let $B_2(0, 1)$ denote the unit Euclidean ball in $(\mathbb{R}^N, \|\cdot\|_2)$. Since $T = s_n B_2(0, 1)$, we see that

$$\mathbb{E} \sup_{\beta \in T} (A_n(F_\beta) - A(F_\beta)) \leq 12\zeta(s_n) \cdot \frac{s_n}{\sqrt{n}} \int_0^\infty \sqrt{\log(B_2(0, 1), \|\cdot\|_\infty, \varepsilon)} d\varepsilon.$$

But $\|\cdot\|_2 \leq \sqrt{N} \|\cdot\|_\infty$, and so

$$\begin{aligned}\mathbb{E} \sup_{\beta \in T} (A_n(F_\beta) - A(F_\beta)) &\leq 12\zeta(s_n) \cdot \frac{s_n}{\sqrt{n}} \int_0^\infty \sqrt{\log \left(B_2(0, 1), \frac{1}{\sqrt{N}} \|\cdot\|_2, \varepsilon \right)} d\varepsilon \\ &= 12\zeta(s_n) \cdot \frac{s_n}{\sqrt{n}} \cdot \frac{1}{\sqrt{N}} \int_0^\infty \sqrt{\log(3/\varepsilon)^N} d\varepsilon \\ &= 12 \frac{s_n \zeta(s_n)}{\sqrt{n}} \int_0^\infty \sqrt{\log(3/\varepsilon)} d\varepsilon.\end{aligned}$$

In the last equality, we used the fact that $N(B_2(0, 1), \|\cdot\|_2, \varepsilon)$ equals 1 for $\varepsilon \geq 1$ and is not larger than $(3/\varepsilon)^N$ for $\varepsilon < 1$ [e.g., [van Handel, 2016](#), Chapter 5]. The same conclusion holds for $\mathbb{E} \sup_{\beta \in T} (A(F_\beta) - A_n(F_\beta))$, and this proves the result. \square

References

- R. van Handel. *Probability in High Dimension*. APC 550 Lecture Notes, Princeton University, 2016.
- V.N. Vapnik. *Statistical Learning Theory*. John Wiley & Sons, New York, 1988.
- R.L. Wheeden and A. Zygmund. *Measure and Integral*. Marcel Dekker, New York, 1977.