

## A Appendix

The appendix is devoted to the proof of Theorem 3.1 and Theorem 3.2. Although Theorem 3.2 uses Theorem 3.1, we still present the proof of Theorem 3.2 first, as it is a more standard proof. After this proof, we will present that of Theorem 3.1. However, first we present some results that we will need later multiple times.

**Lemma A.1** (Elementary Properties of Subgaussian Random Variables). *Let  $U$  be a subgaussian random variable with parameters  $(\beta, \Gamma)$ . Then,*

$$(i) \quad \mathbb{E}[|U|] \leq \frac{1}{\beta} \sqrt{\Gamma - 1};$$

(ii) for any constant  $c \geq 0$ ,  $U + c$  is subgaussian.

*Proof.* (i) follows from

$$\Gamma \geq \mathbb{E}[\exp(|\beta U|^2)] \geq \exp(\mathbb{E}[|\beta U|^2]) \geq \exp(\mathbb{E}[|\beta U|]^2) \geq 1 + \mathbb{E}[|\beta U|]^2.$$

(ii) follows from

$$\mathbb{E}[\exp(\frac{1}{2}\beta^2(U+c)^2)] \leq \mathbb{E}[\exp(\beta^2 U^2 + \beta^2 c^2)] \leq e^{\beta^2 c^2} \mathbb{E}[\exp(|\beta U|^2)] \leq e^{\beta^2 c^2} \Gamma.$$

□

We will also need the following result:

**Lemma A.2.** *Let  $U > 0$ ,  $\mathcal{C} = \mathcal{H} + \mathcal{G}(U)$ . Then, a.s.*

$$\int_0^1 \sqrt{H(u, \mathcal{C}, \|\cdot\|_n)} du \leq 2C_H + 2C_G(U),$$

where  $C_G(U) = \rho^{1/2} \int_0^1 \log^{1/2}(\frac{4U+u}{u}) du (= O(\sqrt{\rho(\log(U))_+})$ .

*Proof of Lemma A.2.* Since  $\mathcal{C} = \mathcal{H} + \mathcal{G}(U)$ , a standard argument shows that

$$H(u; \sigma) \leq H(u/2; \mathcal{H}; \|\cdot\|_n) + H(u/2; \mathcal{G}(U), \|\cdot\|_n). \quad (6)$$

Now, note that  $\|\cdot\|_n \leq \|\cdot\|_{\infty, n}$ . Thus,

$$\begin{aligned} \int_0^1 H^{1/2}(u/2, \mathcal{H}, \|\cdot\|_n) du &= 2 \int_0^{1/2} H^{1/2}(u, \mathcal{H}, \|\cdot\|_n) du \leq 2 \int_0^1 H^{1/2}(u, \mathcal{H}, \|\cdot\|_n) du \\ &\leq 2 \int_0^1 H^{1/2}(u, \mathcal{H}, \|\cdot\|_{\infty, n}) du \leq 2C_H, \end{aligned}$$

where the last inequality is by Assumption 3.3. Moreover, since  $\|g\|_n \leq \|g\|_{\infty}$ ,  $\mathcal{G}(U)$  is a subset of the ball  $B_{\mathcal{G}, \|\cdot\|_n}(0, U)$ . Thus,

$$\begin{aligned} \int_0^1 H^{1/2}(u/2, \mathcal{G}(U), \|\cdot\|_n) du &\leq 2 \int_0^1 H^{1/2}(u, \mathcal{G}(U), \|\cdot\|_n) du \leq 2 \int_0^1 H^{1/2}(u, B_{\mathcal{G}, \|\cdot\|_n}(0, U), \|\cdot\|_n) du \\ &\leq 2\rho^{1/2} \int_0^1 \log^{1/2}\left(\frac{4U+u}{u}\right) du = 2C_G(U), \end{aligned}$$

where the second inequality is by Corollary 2.6 of [van de Geer, 2000], which states that  $H(\varepsilon, B_{\mathcal{G}, \|\cdot\|_n}(0, \sigma)) \leq \rho \log(\frac{4\sigma+\varepsilon}{\varepsilon})$ . Using (6) and  $\sqrt{a+b} \leq \sqrt{a} + \sqrt{b}$  which holds for  $a, b \geq 0$ , we conclude that

$$\int_0^1 \sqrt{H(u; \sigma)} du \leq 2C_H + 2C_G(U),$$

finishing the proof of the claim. □

## B The Proof of Theorem 3.2

In this section we prove Theorem 3.2 assuming that Theorem 3.1 holds.

Let  $U$  be as in Theorem 3.1 and let  $E$  denote the event when

$$\sup_{h \in \mathcal{H}} \|g_{h,n}\|_\infty \leq U.$$

For any  $z \geq 0$ ,

$$\begin{aligned} \mathbb{P}(L(f_n) - L(f^*) > z) &= \mathbb{P}(L(f_n) - L(f^*) > z, E^c) + \mathbb{P}(L(f_n) - L(f^*) > z, E) \\ &\leq \mathbb{P}(E^c) + \mathbb{P}(L(f_n) - L(f^*) > z, E). \end{aligned} \quad (7)$$

Thus, to study the tail probabilities of  $L(f_n) - L(f^*)$ , it suffices to study  $L(f_n) - L(f^*)$  on the event  $E$ .

Define  $\mathcal{G}(U) = \{g \in \mathcal{G} : \|g\|_\infty \leq U\}$  and  $\mathcal{C} = \mathcal{H} + \mathcal{G}(U)$ . On  $E$ , we claim that  $f_n \in \mathcal{C}$ . We have  $f_n = h_n + g_n$  and since  $h_n \in \mathcal{H}$  by definition, it remains to show that  $g_n \in \mathcal{G}(U)$ . By appropriately selecting  $g_{h,n}$ , we can arrange for  $g_n = g_{h_n,n}$ . Hence,  $\|g_n\|_\infty \leq \sup_{h \in \mathcal{H}} \|g_{h,n}\|_\infty \leq U$ , showing that  $f_n \in \mathcal{C}$  indeed holds.

Now, by increasing  $U$  if necessary, we can always arrange for  $f^* = h^* + g^* \in \mathcal{C}$  (for this we may need to increase  $U$  so that  $\|g^*\|_\infty \leq U$ ). Hence, in what follows, we will assume this.<sup>4</sup> By (3), on  $E$  it holds almost surely that

$$\begin{aligned} L(f_n) - L(f^*) &\leq L_n(f^*) - L(f^*) - (L_n(f_n) - L(f_n)) \\ &= (\tilde{\Delta}_n(f_n) - \tilde{\Delta}_n(f^*)) + (\bar{\Delta}_n(f_n) - \bar{\Delta}_n(f^*)) \\ &\leq \underbrace{\sup_{f \in \mathcal{C}} \tilde{\Delta}_n(f) - \tilde{\Delta}_n(f^*)}_{\tilde{\Delta}_n^*(\mathcal{C})} + \underbrace{\sup_{f \in \mathcal{C}} |\bar{\Delta}_n(f) - \bar{\Delta}_n(f^*)|}_{\bar{\Delta}_n^*(\mathcal{C})}, \end{aligned} \quad (8)$$

where we introduced  $\bar{\Delta}_n(f) = L_n(f) - \bar{L}_n(f)$  and  $\tilde{\Delta}_n(f) = L(f) - \bar{L}_n(f)$  with  $\bar{L}_n(f) = \frac{1}{n} \sum_{k=1}^n \mathbb{E}[\ell(Y_k, f(X_k)) | X_k]$ . Note that the first term does not depend on the (unbounded) responses  $Y_1, \dots, Y_n$ . Furthermore, by our assumptions,  $\tilde{\Delta}_n(f)$  is bounded for  $f$  bounded. Hence, we can analyze these terms using tools developed for bounded random variables and empirical processes. Now, while the last term involves  $Y_1, \dots, Y_n$ ,  $\bar{\Delta}_n$  compares average losses over the sample  $X_1, \dots, X_n$ , this last term concerns *in-sample* generalization. Hence, as we will show it below, it can be analyzed using tools developed for the so-called “fixed design” setting. In fact, the following result gives tail bounds for this part:

**Lemma B.1.** *Let Assumptions 3.1 to 3.4 hold and WLOG assume that  $U \geq \max(1, \|g^*\|_\infty)$ . Then, there exist constants  $c, \alpha > 0$  such that for any  $0 < \delta < 1$  satisfying  $\log \frac{1}{\delta} \geq c$  with probability at least  $1 - \delta$ ,*

$$\bar{\Delta}_n^*(\mathcal{C}) \leq 2(r + U) \sqrt{\frac{\log \frac{2}{\delta}}{\alpha n}}. \quad (9)$$

The proof is based on Theorem 3.3 of van de Geer [1990], which we quote below for completeness. Let  $(\Lambda, d)$  be a pseudo-metric space and for  $u > 0$  let  $B_{\Lambda,d}(\lambda, u)$  be the  $d$ -ball in  $\Lambda$  that has radius  $u$  and is centered at  $\lambda$ . We will allow  $d$  to be replaced with a pseudo-norm meaning the ball where the pseudo-metric is defined by the chosen pseudo-norm. The theorem of van de Geer bounds the tails of the suprema of centered, Lipschitz empirical processes of  $\Lambda$  over balls of  $\Lambda$ :

**Theorem B.2** (Theorem 3.3 of van de Geer [1990]). *Let  $(\Lambda, d)$  be a pseudo-metric space with  $d^2 = (1/n) \sum_{k=1}^n d_k^2$  where  $d_1, \dots, d_n$  pseudo-metrics on  $\Lambda$ . Let  $U_1, \dots, U_n$  be real-valued, independent, centered process on  $\Lambda$  such that for  $Z_n = \frac{1}{\sqrt{n}} \sum U_k$ ,  $Z_n(\lambda_0) = 0$  for some  $\lambda_0 \in \Lambda$ . For  $u > 0$ , denote by  $H(u; \sigma) = H(u, B_{\Lambda,d}(\lambda_0, \sigma), d)$ , the  $u$ -entropy of the ball  $B_{\Lambda,d}(\lambda_0, \sigma)$ . Assume further that  $|U_k(\lambda) - U_k(\lambda')| \leq M_k d_k(\lambda, \lambda')$  with  $M_k \geq 0$  random such that  $\mathbb{E}[\exp(|\beta M_k|^2)] \leq \Gamma < \infty$  for some positive constants  $\beta$  and  $\Gamma$ . Then, there exist  $\alpha, \eta, C_1, C_2 > 0$  depending only on  $\beta$  and  $\Gamma$  such that*

$$\mathbb{P} \left( \sup_{\lambda \in B_{\Lambda,d}(\lambda_0, \sigma)} |Z_n(\lambda)| \geq t \right) \leq 2 \exp \left( -\frac{\alpha t^2}{\sigma^2} \right)$$

<sup>4</sup>Note that this is assumed only to simplify the presentation.

holds for any  $t > 0$  and  $\sigma > 0$  that satisfies  $t/\sigma > C_1$  and  $t > C_2 \int_0^{t_0} \sqrt{H(u; \sigma)} du$  where  $t_0 \geq \inf\{u : H(u; \sigma) \leq \eta t^2 / \sigma^2\}$ .

Let us now turn to the proof of Lemma B.1.

*Proof of Lemma B.1.* Let  $(W, \mathcal{W}, \mathbb{P})$  be the probability space that holds  $(X_1, Y_1), \dots, (X_n, Y_n)$ . Note that with no loss of generality, we can assume that  $(W, \mathcal{W})$  is a Borel-space (this is because all our random variables leave in complete, separable metric spaces). For  $x_1, \dots, x_n \in \mathcal{X}$ , let  $x_{1:n} = (x_1, \dots, x_n)$ . Similarly, let  $X_{1:n} = (X_1, \dots, X_n)$ . Define  $(\mathbb{P}_{x_{1:n}})_{x_{1:n} \in \mathcal{X}^n}$  to be the disintegration of the probability measure  $\mathbb{P}$  with respect to  $X_{1:n}$ , also known as the regular conditional probability measure obtained from  $\mathbb{P}$  by conditioning on  $X_{1:n}$ .<sup>5</sup> The expectation operator corresponding to  $\mathbb{P}_{x_{1:n}}$  will be denoted by  $\mathbb{E}_{x_{1:n}}$ . Note that, by the definition of  $\mathbb{P}_{x_{1:n}}$ , for any random variable  $Z$  on  $(W, \mathcal{W}, \mathbb{P})$ ,  $\mathbb{E}_{X_{1:n}}[Z] = \mathbb{E}[Z|X_{1:n}]$  holds everywhere and in particular, for a measurable function  $s : \mathcal{Y} \rightarrow \mathbb{R}$ ,  $\mathbb{E}_{X_{1:n}}[s(Y_k)] = \mathbb{E}[s(Y_k)|X_{1:n}] = \mathbb{E}[s(Y_k)|X_k]$ , where the last equality holds since by assumption  $(X_t, Y_t)$  is i.i.d.

Let

$$U_{k,x_{1:n}}(f) = \Delta_{k,x_{1:n}}(f) - \Delta_{k,x_{1:n}}(f^*)$$

where

$$\Delta_{k,x_{1:n}}(f) = \ell(Y_k, f(x_k)) - \mathbb{E}_{x_{1:n}}[\ell(Y_k, f(x_k))], \quad f \in \Lambda.$$

Note that  $U_{k,x_{1:n}}$ 's are independent, centered processes over  $W_{x_{1:n}}$ . Let  $Z_{n,x_{1:n}} = \frac{1}{\sqrt{n}} \sum_{k=1}^n U_{k,x_{1:n}}$ . By construction,  $Z_{n,x_{1:n}}(f^*) = 0$ . We now show that it is enough to study the deviations of the suprema of  $Z_{n,x_{1:n}}(f)$  over the probability spaces  $W_{x_{1:n}}$ .

We have

$$\bar{\Delta}_n(f) = L_n(f) - \bar{L}_n(f) = \frac{1}{n} \sum_{k=1}^n \Delta_{k,X_{1:n}}(f)$$

and so

$$\sqrt{n}(\bar{\Delta}_n(f) - \bar{\Delta}_n(f^*)) = Z_{n,X_{1:n}}(f).$$

By the construction of  $\mathbb{P}_{x_{1:n}}$ , for  $z \geq 0$ ,

$$\mathbb{P}(\bar{\Delta}_n^*(\mathcal{C}) \geq z) = \int \mathbb{P}_{x_{1:n}}(\bar{\Delta}_n^*(\mathcal{C}) \geq z) P_{X_{1:n}}(dx_{1:n}), \quad (10)$$

and

$$\mathbb{P}_{x_{1:n}}(\bar{\Delta}_n^*(\mathcal{C}) \geq z/\sqrt{n}) = \mathbb{P}_{x_{1:n}}(\sqrt{n} \sup_{f \in \mathcal{C}} |\bar{\Delta}_n(f) - \bar{\Delta}_n(f^*)| \geq z) = \mathbb{P}_{x_{1:n}}(\sup_{f \in \mathcal{C}} Z_{n,x_{1:n}}(f) \geq z). \quad (11)$$

Let  $\Lambda = \mathcal{C}$ ,  $d_{k,x_{1:n}}(f, f') = |f(x_k) - f'(x_k)|$  and  $d_{x_{1:n}}^2(f, f') = \frac{1}{n} \sum_{k=1}^n d_{k,x_{1:n}}^2(f, f')$ . By construction,  $d_{x_{1:n}}(f, f') = \|f - f'\|_n$ . Since for any  $f = h + g \in \mathcal{C}$ ,

$$\|f - f^*\|_n = \|h - h^*\|_n + \|g - g^*\|_n \leq \|h - h^*\|_\infty + \|g - g^*\|_\infty \leq 2(r + U) =: \sigma,$$

thus,  $\mathcal{C} \subset B_{\Lambda, d_{x_{1:n}}}(f^*, \sigma) \subset \Lambda = \mathcal{C}$  and

$$\mathbb{P}_{x_{1:n}}(\bar{\Delta}_n^*(\mathcal{C}) > z/\sqrt{n}) = \mathbb{P}_{x_{1:n}} \left( \sup_{f \in B_{\Lambda, d_{x_{1:n}}}(f^*, \sigma)} Z_{n,x_{1:n}}(f) > z \right). \quad (12)$$

Thus, it remains to bound this latter probability. Fix  $x_{1:n} \in \mathcal{X}^n$  such that

$$\mathbb{E}_{x_{1:n}}[\exp((\beta K_\ell(Y, c)^2)] \leq \Gamma_c, \quad \text{for all } c > 0. \quad (13)$$

---

<sup>5</sup> The defining properties of  $(\mathbb{P}_{x_{1:n}})$  are that for each  $x_{1:n} \in \mathcal{X}^n$ ,  $\mathbb{P}_{x_{1:n}}$  is a probability measure on  $(W, \mathcal{W})$  concentrated on  $\{X_{1:n} = x_{1:n}\}$ ,  $x_{1:n} \mapsto \mathbb{P}_{x_{1:n}}$  is measurable and for any  $f : (W, \mathcal{W}) \rightarrow [0, \infty)$  measurable function  $\int f(w) \mathbb{P}(dw) = \int (\int f(w) \mathbb{P}_{x_{1:n}}(dw)) P_{X_{1:n}}(dx_{1:n})$ . The existence of  $(\mathbb{P}_{x_{1:n}})$ , which is also called a regular conditional probability distribution is ensured thanks to the assumption that  $(W, \mathcal{W})$  is Borel. Moreover,  $(\mathbb{P}_{x_{1:n}})$  is unique up to an almost sure equivalence in the sense that if  $(\hat{\mathbb{P}}_{x_{1:n}})$  is another disintegration of  $\mathbb{P}$  w.r.t.  $X_{1:n}$  then  $P_X(\{x_{1:n} : \mathbb{P}_{ux} \neq \hat{\mathbb{P}}_{x_{1:n}}\}) = 0$ . For background on disintegration and conditioning, the reader is referred to [Chang and Pollard \[1997\]](#).

Let us now apply Theorem B.2 to  $W_{x_{1:n}} = (W, \mathcal{W}, \mathbb{P}_{x_{1:n}})$  with  $\Lambda$ ,  $(d_{k,x_{1:n}})$  and  $(U_{k,x_{1:n}})$  ( $k = 1, \dots, n$ ), as defined above. To verify the uniform subgaussian property of the Lipschitz coefficient of  $U_{k,x_{1:n}}$ , note that for  $f, f' \in \mathcal{C}$ , by Assumption 3.1(iii),

$$\begin{aligned} |U_{k,x_{1:n}}(f) - U_{k,x_{1:n}}(f')| &= |\Delta_{k,x_{1:n}}(f) - \Delta_{k,x_{1:n}}(f')| \\ &\leq |\ell(Y_k, f(x_k)) - \ell(Y_k, f'(x_k))| + |\mathbb{E}_{x_{1:n}}[\ell(Y_k, f(x_k)) - \ell(Y_k, f'(x_k))]| \\ &\leq K_l(Y_k, r + U)|f(x_k) - f'(x_k)| + \mathbb{E}_{x_{1:n}}[K_l(Y_k, r + U)]|f(x_k) - f'(x_k)|. \end{aligned}$$

By Lemma A.1(i),  $\mathbb{E}_{x_{1:n}}[K_l(Y_k, r + U)] \leq \frac{1}{\beta} \sqrt{\Gamma_{r+U} - 1}$  and so by part (ii) of the same lemma,  $K_l(Y_k, r + U) + \mathbb{E}_{x_{1:n}}[K_l(Y_k, r + U)]$  is subgaussian, with parameters  $\beta'$  and  $\Gamma'$  only depending on  $r + U$ .

Therefore, from Theorem B.2 we conclude that there exists  $C_1, C_2, \eta > 0$  such that for any  $t > 0$  satisfying  $\eta t^2 / \sigma^2 \geq H(1; \sigma)$ ,  $t > C_1 \sigma$  and  $t > C_2 \int_0^1 \sqrt{H(u; \sigma)} du$ , it holds that

$$\mathbb{P}_{x_{1:n}} \left( \sup_{f \in \mathcal{C}} |Z_{n,x_{1:n}}(f)| \geq t \right) = \mathbb{P}_{x_{1:n}} \left( \sup_{f \in B_{\Lambda, d_{x_{1:n}}}(f^*, \sigma)} |Z_{n,x_{1:n}}(f)| \geq t \right) \leq 2 \exp \left( -\frac{\alpha t^2}{\sigma^2} \right). \quad (14)$$

It still remains to check that  $H(1, \sigma)$  and  $\int_0^1 \sqrt{H(u; \sigma)} du$  are finite (otherwise the result is vacuous). By definition,  $H(u; \sigma) = H(u, B_{\Lambda, d_{x_{1:n}}}(f^*, \sigma), d_{x_{1:n}}) = H(u, \mathcal{C}, d_{x_{1:n}}) = H(u, \mathcal{C}, \|\cdot\|_n)$ . Hence, by Lemma A.2,  $\int_0^1 H^{1/2}(u; \sigma) du \leq 2C_H + 2C_G(U)$ . Noting that  $H(u; \sigma)$  is monotonically decreasing in  $u$ , we calculate  $H^{1/2}(1; \sigma) \leq \int_0^1 H^{1/2}(u; \sigma) du \leq 2C_H + 2C_G(U)$  and so  $H(1; \sigma) \leq (2C_H + 2C_G(U))^2 < \infty$ . We conclude that (14) holds for any  $t \geq t_{\min} := \max\{C_1 \sigma, C_2(2C_H + 2C_G(U)), (2C_H + 2C_G(U))\sigma \eta^{-1/2}\}$ .

Since by Assumption 3.1(iii), (13) holds  $[P_X]$ -almost surely, combining (10), (12) and (14), we get

$$\mathbb{P} \left( \bar{\Delta}_n^*(\mathcal{C}) \geq t / \sqrt{n} \right) \leq 2 \exp \left( -\frac{\alpha t^2}{\sigma^2} \right). \quad (15)$$

Inverting this inequality, we see that for any  $0 < \delta < 1$  such that  $\log(2/\delta) \geq t_{\min}^2 \alpha / \sigma^2$ , with probability at least  $1 - \delta$ ,

$$\bar{\Delta}_n^*(\mathcal{C}) \leq \sigma \sqrt{\frac{\log \frac{2}{\delta}}{\alpha n}},$$

finishing the proof.  $\square$

It remains to bound  $\tilde{\Delta}_n^*(\mathcal{C}) = \sup_{f \in \mathcal{C}} \tilde{\Delta}_n(f) - \tilde{\Delta}_n(f^*)$ . For this, define

$$\bar{\ell}(x, p) = \mathbb{E}[\ell(Y, p) | X = x].$$

With a slight abuse of notation, we also introduce  $\bar{\ell}(x, f) = \bar{\ell}(x, f(x))$ . Let

$$B(\bar{\ell}, U) = \left\| \sup_{p \in [-r-U, r+U]} \bar{\ell}(X, p) \right\|_{L^\infty},$$

where  $\|\cdot\|_{L^\infty}$  denotes the essential supremum of its argument. We also let  $\bar{L}$  be the Lipschitz constant of  $\bar{\ell}$  when  $p \in [-r-U, r+U]$ :

$$\text{Lip}(\bar{\ell}, U) = \left\| \sup_{p, p' \in [-r-U, r+U], p \neq p'} \frac{\bar{\ell}(X, p) - \bar{\ell}(X, p')}{|p - p'|} \right\|_{L^\infty}.$$

The next lemma shows that both quantities are finite:

**Lemma B.3.** *Let  $r' = \max(r + U, \|\hat{h}\|_\infty)$ . Then,  $B(\bar{\ell}, U) \leq Q + \frac{2r'}{\beta} \sqrt{\Gamma_{r'} - 1} < +\infty$  and  $\text{Lip}(\bar{\ell}, U) < \frac{\sqrt{\Gamma_{r+U} - 1}}{\beta} < +\infty$ .*

*Proof.* For the second statement, for any  $t, s \in [-b, b]$  we have

$$\bar{\ell}(X, t) - \bar{\ell}(X, s) \leq \mathbb{E}[\ell(Y, t) - \ell(Y, s) | X] \leq \mathbb{E}[K_\ell(Y, b)|t - s| | X] \leq \frac{\sqrt{\Gamma_b - 1}}{\beta} |t - s|,$$

where we used Assumption 3.1(iii) and Lemma A.1(i). Thus,  $\text{Lip}(\bar{\ell}, U) \leq \frac{\sqrt{\Gamma_{r+U-1}}}{\beta} < +\infty$ .

For the first statement take some  $|p| \leq r + U$  and write

$$\begin{aligned} \bar{\ell}(X, p) &\leq \bar{\ell}(X, \hat{h}(X)) + |\bar{\ell}(X, p) - \bar{\ell}(X, \hat{h}(X))| \leq Q + \text{Lip}(\bar{\ell}, r')|p - \hat{h}(X)| \leq Q + \text{Lip}(\bar{\ell}, r')(r + U + \|\hat{h}\|_\infty) \\ &\leq Q + \frac{\Gamma_{r'-1}}{\beta}(2r'), \end{aligned}$$

where in the second inequality we used Assumption 3.1(ii), while in the last one we used the bound on the Lipschitz coefficient.  $\square$

As it is well known, the Rademacher complexity of  $\mathcal{C}$ , defined next, captures *exactly* the behavior of  $\mathbb{E} \left[ \tilde{\Delta}_n^*(\mathcal{C}) \right]$  (e.g., [Tewari and Bartlett \[2013\]](#)).

**Definition 3** (Rademacher Complexity of Subsets of  $\mathbb{R}^n$ ). *Let  $A \subset \mathbb{R}^n$ ,  $(\sigma_1, \dots, \sigma_n) \in \{-1, +1\}^n$  be independent Rademacher random variables (i.e.,  $\mathbb{P}(\sigma_k = 1) = 1/2$ ). The Rademacher complexity of  $A$ ,  $\mathfrak{R}(A)$  is*

$$\mathfrak{R}(A) = \frac{1}{n} \mathbb{E} \left[ \sup_{a \in A} \sum_{i=1}^n \sigma_i a_i \right].$$

**Definition 4** (Rademacher Complexity of Function Sets). *Let  $\mathcal{F} \subset \{f : \mathcal{X} \rightarrow \mathbb{R}\}$  and  $P$  be a measure on  $\mathcal{X}$ . Then, the  $n$ th Rademacher number of  $\mathcal{F}$  induced by  $P$  is*

$$\mathfrak{R}_n(\mathcal{F}) = \mathbb{E} [\mathfrak{R}(\mathcal{F}(X_{1:n}))],$$

where  $\mathcal{F}(X_{1:n}) = \{(f(X_1), \dots, f(X_n)) : f \in \mathcal{F}\}$  is the projection of  $\mathcal{F}$  to an i.i.d. sample  $X_{1:n} = (X_1, \dots, X_n)$  from  $P$ . When  $n$  and  $P$  are uniquely identified from the context, we also call  $\mathfrak{R}_n(\mathcal{F})$  the Rademacher-complexity of  $\mathcal{F}$ .

The Rademacher complexity enjoys a number of useful properties, amongst which we need the following contraction property:

**Theorem B.4.** *For  $A \subset \mathbb{R}^n$  and  $\phi = (\phi_1, \dots, \phi_n) : \mathbb{R}^n \rightarrow \mathbb{R}^n$ , define  $\phi \circ A = \{(\phi_1(a_1), \dots, \phi_n(a_n)) : a \in A\}$ . Assume that all the component functions  $\phi_i$  are  $L$ -Lipschitz over  $A$ . Then,  $\mathfrak{R}(\phi \circ A) \leq L\mathfrak{R}(A)$ .*

Note that this theorem is usually stated for the case when  $\phi_1 = \dots = \phi_n$ . The simpler form is sufficient for “margin based losses” (used in classification) that have the form  $\ell(y, p) = g(y p)$  with some  $g$ . As we will see, here we need this more general form as our losses are less constrained. However, the proof of this more general result still follows the standard reasoning.

*Proof.* We follow the proof of Theorem 11.9 in [\[Rakhlin and Sridharan, 2014\]](#) and write

$$\begin{aligned} n\mathfrak{R}(\phi \circ A) &= \mathbb{E} \left[ \sup_{a \in A} \sum_{i=1}^n \sigma_i \phi_i(a_i) \right] \\ &= \frac{1}{2} \left\{ \mathbb{E} \left[ \sup_{a \in A} \sum_{i=1}^{n-1} \sigma_i \phi_i(a_i) + \phi_n(a_n) \mid \sigma_n = 1 \right] + \mathbb{E} \left[ \sup_{b \in A} \sum_{i=1}^{n-1} \sigma_i \phi_i(b_i) - \phi_n(b_n) \mid \sigma_n = -1 \right] \right\} \\ &= \frac{1}{2} \left\{ \mathbb{E} \left[ \sup_{a, b \in A} \sum_{i=1}^{n-1} \sigma_i (\phi_i(a_i) + \phi_i(b_i)) + (\phi_n(a_n) - \phi_n(b_n)) \right] \right\} \\ &\leq \frac{1}{2} \left\{ \mathbb{E} \left[ \sup_{a, b \in A} \sum_{i=1}^{n-1} \sigma_i (\phi_i(a_i) + \phi_i(b_i)) + L|a_n - b_n| \right] \right\}. \end{aligned}$$

Now assume that some  $(a^*, b^*)$  achieves the supremum (the proof when the supremum is not achieved is easy once we know how to prove the statement for the case when the supremums involved are all achieved). If  $a_n^* \geq b_n^*$ ,

the absolute value can be removed. Otherwise,  $(b^*, a^*)$  will achieve the same supremum, and again the absolute value can be removed. Thus, the last expression is bounded by

$$\begin{aligned} & \frac{1}{2} \left\{ \mathbb{E} \left[ \sup_{a, b \in A} \sum_{i=1}^{n-1} \sigma_i (\phi_i(a_i) + \phi_i(b_i)) + L(a_n - b_n) \right] \right\} \\ &= \frac{1}{2} \left\{ \mathbb{E} \left[ \sup_{a \in A} \sum_{i=1}^{n-1} \sigma_i \phi_i(a_i) + L a_n \mid \sigma_n = 1 \right] + \mathbb{E} \left[ \sup_{b \in A} \sum_{i=1}^{n-1} \sigma_i \phi_i(b_i) - L b_n \mid \sigma_n = -1 \right] \right\} \\ &= \mathbb{E} \left[ \sup_{a \in A} \sum_{i=1}^{n-1} \sigma_i \phi_i(a_i) + L \sigma_n a_n \right]. \end{aligned}$$

Continuing this way,

$$\mathbb{E} \left[ \sup_{a \in A} \sum_{i=1}^{n-1} \sigma_i \phi_i(a_i) + L \sigma_n a_n \right] \leq \mathbb{E} \left[ \sup_{a \in A} \sum_{i=1}^{n-2} \sigma_i \phi_i(a_i) + L(\sigma_{n-1} a_{n-1} + \sigma_n a_n) \right] \leq L \mathbb{E} \left[ \sup_{a \in A} \sum_{i=1}^n \sigma_i a_i \right],$$

thus finishing the proof.  $\square$

Let  $\mathcal{L} = \{s_f : \mathcal{X} \rightarrow \mathbb{R} : s_f(x) = \bar{\ell}(x, f) - \bar{\ell}(x, f^*), f \in \mathcal{C}, x \in \mathcal{X}\}$ . Note that  $\tilde{\Delta}_n(f) - \tilde{\Delta}_n(f^*) = (L(f) - \bar{L}_n(f)) - (L(f^*) - \bar{L}_n(f^*)) = \mathbb{E} [\bar{\ell}(X, f) - \bar{\ell}(X, f^*)] - \frac{1}{n} \sum_{k=1}^n (\bar{\ell}(X_k, f) - \bar{\ell}(X_k, f^*)) = \mathbb{E} [s_f(X)] - \frac{1}{n} \sum_{k=1}^n s_f(X_k)$ . Following the standard argument, since the range of functions in  $\mathcal{L}$  is bounded by  $B(\bar{\ell}, U)$ , by McDiarmid's inequality, for any  $0 < \delta < 1$ , with probability at least  $1 - \delta$ ,

$$\tilde{\Delta}_n^*(\mathcal{C}) = \sup_{s \in \mathcal{L}} \mathbb{E} [s(X)] - \frac{1}{n} \sum_{k=1}^n s(X_k) \leq \mathbb{E} \left[ \sup_{s \in \mathcal{L}} \mathbb{E} [s(X)] - \frac{1}{n} \sum_{k=1}^n s(X_k) \right] + B(\bar{\ell}, U) \sqrt{\frac{2 \log \frac{1}{\delta}}{n}}.$$

Following the calculation before Theorem 7 in Section 3.2 of [Tewari and Bartlett \[2013\]](#),

$$\mathbb{E} \left[ \sup_{s \in \mathcal{L}} \mathbb{E} [s(X)] - \frac{1}{n} \sum_{k=1}^n s(X_k) \right] \leq 2\mathfrak{R}_n(\mathcal{L}).$$

Let us now bound  $\mathfrak{R}_n(\mathcal{L}) = \mathbb{E} [\mathfrak{R}(\mathcal{L}(X_{1:n}))]$ . We can write

$$\mathcal{L}(X_{1:n}) = \{s_f(X_{1:n}) : f \in \mathcal{C}\} = \phi \circ \mathcal{C}(X_{1:n}),$$

where  $\phi = (\phi_1, \dots, \phi_n) : \mathbb{R}^n \rightarrow \mathbb{R}^n$  is defined by  $\phi_k(t) = \bar{\ell}(X_k, t) - \bar{\ell}(X_k, f^*(X_k))$  (note that  $\phi$  is random). By definition, each component of  $\phi$  is almost surely Lipschitz over any bounded interval  $[-b, b]$  with the same Lipschitz constant (depending on  $b$ ). Indeed, for any  $t, s \in [-b, b]$ ,

$$\begin{aligned} \|\phi_k(t) - \phi_k(s)\|_{L^\infty} &= \|\bar{\ell}(X_k, t) - \bar{\ell}(X_k, s)\|_{L^\infty} \\ &= \inf \{a \in \mathbb{R} : \mathbb{P} (|\bar{\ell}(X_k, t) - \bar{\ell}(X_k, s)| > a)\} \\ &= \inf \{a \in \mathbb{R} : \mathbb{P} (|\bar{\ell}(X, t) - \bar{\ell}(X, s)| > a)\} \\ &= \|\bar{\ell}(X, t) - \bar{\ell}(X, s)\|_{L^\infty} \\ &\leq \text{Lip}(\bar{\ell}, b) |t - s|, \end{aligned}$$

where the second and fourth equalities used the definition of  $\|\cdot\|_{L^\infty}$  and the third used that  $X_k$  and  $X$  are identically distributed. Now, since  $\mathcal{C}$  contains functions bounded by  $r + U$ , by Theorem [B.4](#),

$$\mathfrak{R}(\phi \circ \mathcal{C}(X_{1:n})) \leq \text{Lip}(\bar{\ell}, U) \mathfrak{R}(\mathcal{C}(X_{1:n})) \quad \text{a.s.}$$

and hence

$$\mathfrak{R}_n(\mathcal{L}) = \mathbb{E} \mathfrak{R}(\mathcal{L}(X_{1:n})) = \mathbb{E} \mathfrak{R}(\phi \circ \mathcal{C}(X_{1:n})) \leq \text{Lip}(\bar{\ell}, U) \mathbb{E} \mathfrak{R}(\mathcal{C}(X_{1:n})) = \text{Lip}(\bar{\ell}, U) \mathfrak{R}_n(\mathcal{C}).$$

Our next goal is to bound  $\mathfrak{R}_n(\mathcal{C})$ . By Dudley’s entropy integral bound [Dudley, 1967] (e.g., Theorem 10 of Tewari and Bartlett [2013], for a statement with a proof see Theorem 11.4 of Rakhlin and Sridharan [2014]),

$$\mathfrak{R}_n(\mathcal{C}) \leq \frac{12}{\sqrt{n}} \mathbb{E} \int_0^1 H^{1/2}(u, \mathcal{C}, \|\cdot\|_n) du \leq \frac{12}{\sqrt{n}} (2C_H + 2C_G(U)),$$

where the second inequality holds thanks to Lemma A.2 and we also used that Dudley’s bound holds regardless the scale of the range of functions in  $\mathcal{C}$ , which is not hard to check by inspecting the proof of the bound. Combining all the inequalities we get that with probability at least  $1 - \delta$ ,

$$\tilde{\Delta}_n^*(\mathcal{C}) \leq \frac{48(C_H + C_G(U)) \text{Lip}(\bar{\ell}, U)}{\sqrt{n}} + B(\bar{\ell}, U) \sqrt{\frac{\log \frac{1}{\delta}}{2n}}. \quad (16)$$

Combining Equations (7) and (8), we have for any  $z \geq 0$ ,

$$\mathbb{P}(L(f_n) - L(f^*) > z) \leq \mathbb{P}(E^c) + \mathbb{P}\left(\tilde{\Delta}_n^*(\mathcal{C}) + \bar{\Delta}_n^*(\mathcal{C}) > z\right). \quad (17)$$

Now, by Lemma B.1 and (16), for any  $0 < \delta < 1$  such that  $\log(1/\delta) \geq c$ , with probability at least  $1 - 2\delta$ ,

$$\tilde{\Delta}_n^*(\mathcal{C}) + \bar{\Delta}_n^*(\mathcal{C}) \leq \frac{48(C_H + C_G(U)) \text{Lip}(\bar{\ell}, U)}{\sqrt{n}} + 2(r + U) \sqrt{\frac{\log \frac{2}{\delta}}{\alpha n}} + B(\bar{\ell}, U) \sqrt{\frac{\log \frac{1}{\delta}}{2n}} =: \pi(\delta).$$

Together with (17) and Theorem 3.1, we thus get that with probability  $1 - 3\delta$ , provided that  $\log(1/\delta) \geq c$  and  $n \geq c_1 + c_2 \frac{\log(\frac{2\rho}{\delta})}{\lambda_{\min}}$ ,

$$L(f_n) - L(f^*) \leq \pi(\delta),$$

thus finishing the proof.

## C The Proof of Theorem 3.1

In this section we present the proof of Theorem 3.1, which calls for a bound of

$$\sup_{h \in \mathcal{H}} \|g_{h,n}\|_\infty$$

that holds with high probability. Fix  $h \in \mathcal{H}$ . Then,  $g_{h,n}(x) = \langle \theta, \phi(x) \rangle \leq \|\theta_{h,n}\|_2 \|\phi(x)\|_2$ , where  $\theta_{h,n}$  is the parameter vector of  $g_{h,n}$ . Since  $\|\phi(x)\|_2 \leq 1$ , it suffices to bound  $\|\theta_{h,n}\|_2$ . On  $G_{\lambda_{\min}}$ , which is defined as the event  $\{\hat{\lambda}_{\min} \geq \lambda_{\min}/2\}$ , we have

$$g_{h,n}^2(x) \leq \|\theta_{h,n}\|_2^2 \leq \frac{\theta_{h,n}^\top \hat{G} \theta_{h,n}}{\hat{\lambda}_{\min}} = \frac{2 \|g_{h,n}\|_n}{\lambda_{\min}}. \quad (18)$$

Hence, the problem is reduced to proving a uniform ( $h$ -independent) upper bound on the empirical norm of  $g_{h,n}$  and showing that  $G_{\lambda_{\min}}$  happens with “large probability”.

For the latter, we use a result of Gittens and Tropp [2011]. This is summarized in the lemma which also includes some observations that will prove to be useful later:

**Lemma C.1.** *The following hold:*

- (i) *With probability one, for any  $\theta \in \mathbb{R}^d$ ,  $\theta^\top \hat{G} \theta \leq \frac{\theta^\top G \theta}{\lambda_{\min}}$ .*
- (ii) *Assuming that  $n \in \mathbb{N}$  and  $\delta \in (0, 1)$  are such that*

$$n \geq \frac{2}{\lambda_{\min} \log\left(\frac{\rho}{\delta}\right)} \log\left(\frac{\rho}{\delta}\right), \quad (19)$$

where  $\rho$  and  $\lambda_{\min}$  are respectively the rank and the smallest positive eigenvalue of  $G$ , with probability at least  $1 - \delta$ , it holds that  $\hat{\lambda}_{\min} \geq \frac{\lambda_{\min}}{2} > 0$ .

(iii) For any  $n, \delta$  satisfying (19), with probability  $1 - \delta$  it holds that for any  $\theta \in \mathbb{R}^d$  and  $[P_X]$  almost every  $x \in \mathcal{X}$ ,

$$|\langle \theta, \phi(x) \rangle| \leq \sqrt{\frac{2\theta^\top \hat{G}\theta}{\lambda_{\min}}}.$$

The (easy) proof of the lemma is deferred to Appendix C.2.

To get an upper bound on the empirical norm of  $g_{h,n}$ , we will use

$$\|g_{h,n}\|_n \leq \|g_{h,n} - \bar{g}_{h,n}\|_n + \|\bar{g}_{h,n}\|_n \quad (20)$$

and develop uniform bound on the two terms on the r.h.s..

**Lemma C.2.** *It holds almost surely that*

$$\sup_{h \in \mathcal{H}} \|\bar{g}_{h,n}\|_n \leq \bar{R},$$

where  $\bar{R} = R_{C_0} + r$ ,  $C_0 = \frac{2\hat{r}}{\beta} \sqrt{\Gamma_{\hat{r}} - 1} + Q$ ,  $\hat{r} = \max(r, \|\hat{h}\|_\infty)$  and  $\hat{h}$  is the function from Assumption 3.1(ii).

The constant  $R_{C_0}$  that appears in the statement is defined in our “level-set assumption” (cf. Assumption 3.1(iv)).

*Proof.* Fix some  $h \in \mathcal{H}$ . We have  $\|\bar{g}_{h,n}\|_n = \|h + \bar{g}_{h,n} + (-h)\|_n \leq \|h + \bar{g}_{h,n}\|_n + \|-h\|_n \leq \|h + \bar{g}_{h,n}\|_n + r$  thanks to  $\|h\|_\infty \leq r$ . Hence, it remains to bound  $\|h + \bar{g}_{h,n}\|_n$ .

By Assumption 3.1(iv), for this it suffices if we show a bound on  $\bar{L}_n(h + \bar{g}_{h,n})$  since by this assumption if  $\bar{L}_n(h + \bar{g}_{h,n}) \leq c$  then  $\|h + \bar{g}_{h,n}\|_n \leq R_c$ . By the optimizing property of  $\bar{g}_{h,n}$ , we have  $\bar{L}_n(h + \bar{g}_{h,n}) = \bar{L}_{n,h}(\bar{g}_{h,n}) \leq \bar{L}_{n,h}(0) = \bar{L}_n(h)$ . Now, by definition

$$\bar{L}_n(h) = \mathbb{E} \left[ \frac{1}{n} \sum_i \ell(Y_i, h(X_i)) \middle| X_{1:n} \right],$$

hence, it suffices to bound  $\mathbb{E}[\ell(Y_i, h(X_i)) | X_i]$ . For this, we have

$$\mathbb{E}[\ell(Y_i, h(X_i)) | X_i] \leq \mathbb{E} \left[ |\ell(Y_i, h(X_i)) - \ell(Y_i, \hat{h}(X_i))| \middle| X_i \right] + \mathbb{E} \left[ \ell(Y_i, \hat{h}(X_i)) \middle| X_i \right],$$

where we used that by assumption the loss is nonnegative. By Assumption 3.1(ii),

$$\mathbb{E} \left[ \ell(Y_i, \hat{h}(X_i)) \middle| X_i \right] \leq Q.$$

Therefore it is sufficient to bound

$$\mathbb{E} \left[ |\ell(Y_i, h(X_i)) - \ell(Y_i, \hat{h}(X_i))| \middle| X_i \right].$$

Note that by Assumption 3.1(iii), almost surely  $\mathbb{E}[\exp(|\beta K_\ell(Y, r)|^2) | X] \leq \Gamma_r$ . So, by Lemma A.1 (i),  $\mathbb{E}[K_\ell(Y, r) | X] \leq \frac{1}{\beta} \sqrt{\Gamma_r - 1}$  a.s.. Thus, with  $\hat{r} = \max(r, \|\hat{h}\|_\infty)$ ,

$$\mathbb{E} \left[ |\ell(Y_i, h(X_i)) - \ell(Y_i, \hat{h}(X_i))| \middle| X_i \right] \leq \mathbb{E} [2\hat{r} K_\ell(Y_i, \hat{r}) | X_i] \leq \frac{2\hat{r}}{\beta} \sqrt{\Gamma_{\hat{r}} - 1}.$$

Putting together the inequalities, we obtain that  $\bar{L}_n(h + \bar{g}_{h,n}) \leq \frac{2\hat{r}}{\beta} \sqrt{\Gamma_{\hat{r}} - 1} + Q =: C_0$  and thus  $\|h + \bar{g}_{h,n}\|_n \leq R_{C_0}$ .  $\square$

Let us now consider bounding  $\|g_{h,n} - \bar{g}_{h,n}\|_n$ . In fact, we will only bound this on the event  $G_{\lambda_{\min}}$  when  $\hat{\lambda}_{\min} \geq \lambda_{\min}/2$ . Since we use this event to upper bound  $1/\hat{\lambda}_{\min}$  by  $2/\lambda_{\min}$ , there is no loss in bounding  $\|g_{h,n} - \bar{g}_{h,n}\|_n$  on this event only. Note that by Lemma C.1 (ii),  $G_{\lambda_{\min}}$  holds with probability at least  $1 - \delta$ .

**Lemma C.3.** *There exist problem-dependent positive constants  $C_0$  and  $L_0 \geq 1$  such that for any  $n \geq 16L_0^4$ , it holds that*

$$\mathbb{P} \left( \sup_{h \in \mathcal{H}} \|g_{h,n} - \bar{g}_{h,n}\|_n \geq 1, G_{\lambda_{\min}} \right) \leq \exp \left( -\frac{C_0 n}{4} \right). \quad (21)$$



The proof of this lemma follows the proofs in the paper of [van de Geer \[1990\]](#), who studied the deviations  $\|g_{h,n} - \bar{g}_{h,n}\|_n$  for  $h = 0$  (see also [van de Geer 2000](#)). It turns out the techniques of the mentioned paper are just strong enough to bound the uniform deviation  $\sup_{h \in \mathcal{H}} \|g_{h,n} - \bar{g}_{h,n}\|_n$ . As the proof is lengthy and technical, it is developed in a separate section.

Now, combining (18), (20) and Lemma C.2 we get that on  $G_{\lambda_{\min}}$ ,

$$G_{n,\infty} \doteq \sup_{h \in \mathcal{H}} \|g_{h,n}\|_\infty \leq \frac{2}{\lambda_{\min}} \sup_{h \in \mathcal{H}} \|g_{h,n}\|_n \leq \frac{2}{\lambda_{\min}} \left( \bar{R} + \sup_{h \in \mathcal{H}} \|g_{h,n} - \bar{g}_{h,n}\|_n \right). \quad (22)$$

Since for any  $A > 0$ ,

$$\mathbb{P}(G_{n,\infty} > A) \leq \mathbb{P}(G_{\lambda_{\min}}^c) + \mathbb{P}(G_{n,\infty} > A, G_{\lambda_{\min}})$$

and by (22),

$$\mathbb{P}(G_{n,\infty} > A, G_{\lambda_{\min}}) \leq \mathbb{P}\left(\frac{2}{\lambda_{\min}} \left( \bar{R} + \sup_{h \in \mathcal{H}} \|g_{h,n} - \bar{g}_{h,n}\|_n \right) > A, G_{\lambda_{\min}}\right),$$

choosing  $A = \frac{2}{\lambda_{\min}} (\bar{R} + 1)$ , we see that

$$\mathbb{P}\left(G_{n,\infty} > \frac{2}{\lambda_{\min}} (\bar{R} + 1)\right) \leq \mathbb{P}(G_{\lambda_{\min}}^c) + \mathbb{P}\left(\sup_{h \in \mathcal{H}} \|g_{h,n} - \bar{g}_{h,n}\|_n \geq 1, G_{\lambda_{\min}}\right).$$

By Eq. (19) and Lemma C.3, provided that  $n \geq \frac{2}{\lambda_{\min} \log(\frac{6}{5})} \log(\frac{2\rho}{\delta})$ ,  $n \geq 16L_0^4$  and  $n \geq \frac{4 \log(\frac{2}{\delta})}{C_0}$  we get that

$$\mathbb{P}\left(G_{n,\infty} > \frac{2}{\lambda_{\min}} (\bar{R} + 1)\right) \leq \delta,$$

which is the desired statement. In particular, we can choose  $U = \frac{2}{\lambda_{\min}} (\bar{R} + 1)$ .

### C.1 The Proof of Lemma C.3

The proof follows the ideas from the paper of [van de Geer \[1990\]](#). Lemma C.3 calls for a uniform (in  $h \in \mathcal{H}$ ) bound for  $\|g_{h,n} - \bar{g}_{h,n}\|_n$ . Fix  $h \in \mathcal{H}$ . We consider a self-normalized “version” of the differences  $g_{h,n} - \bar{g}_{h,n}$ , which are easier to deal with. This is done as follows: For  $g \in \mathcal{G}$ , define

$$\omega_{g,h} = \frac{g - \bar{g}_{h,n}}{1 + K \|g - \bar{g}_{h,n}\|_n} \quad \text{and} \quad \Omega_{h,n} = \{\omega_{g,h} : g \in \mathcal{G}\},$$

where  $K > 0$  is to be chosen later. Then, for any  $\omega \in \Omega_{h,n}$ ,  $\|\omega\|_n < \frac{1}{K}$  and

$$\begin{aligned} \|g - \bar{g}_{h,n}\|_n &= \frac{\|g - \bar{g}_{h,n}\|_n}{1 + K \|g - \bar{g}_{h,n}\|_n} \left(1 + K \|g - \bar{g}_{h,n}\|_n\right) = \|\omega_{g,h}\|_n \left(1 + K \|g - \bar{g}_{h,n}\|_n\right) \\ &= \frac{\|\omega_{g,h}\|_n}{1 - K \|\omega_{g,h}\|_n}. \end{aligned} \quad (23)$$

Thus, we see that is enough to control the empirical norm of

$$\hat{\omega}_{h,n} = \omega_{g_{h,n},h} = \frac{g_{h,n} - \bar{g}_{h,n}}{1 + K \|g_{h,n} - \bar{g}_{h,n}\|_n}.$$

The first step is to bound this norm in terms of the increments of the empirical process

$$\Delta_{h,n}(g) := L_{h,n}(g) - \bar{L}_{h,n}(g).$$

**Lemma C.4** (“Basic Inequality”). *Let Assumption 3.2 hold. There exists a constant  $\eta$ , such that on the event  $G_{\lambda_{\min}}$ , for any  $h \in \mathcal{H}$ ,*

$$\eta \|\hat{\omega}_{h,n}\|_n^2 \leq \Delta_{h,n}(\bar{g}_{h,n}) - \Delta_{h,n}(\bar{g}_{h,n} + \hat{\omega}_{h,n}).$$

The proof, which is stated in Appendix C.3, follows standard arguments. Based on this, we can reduce the study of the supremum of the empirical norm of  $\hat{\omega}_{h,n}$  to that of the supremum of the increments  $\mathcal{V}_{h,n}(\omega) = \sqrt{n}(\Delta_{h,n}(\bar{g}_{h,n}) - \Delta_{h,n}(\bar{g}_{h,n} + \omega))$  normalized by  $\omega$ . In particular, it follows from Lemma C.4 that for  $L, \sigma > 0$ ,

$$\begin{aligned} & \mathbb{P}\left(\sup_{h \in \mathcal{H}} \|\hat{\omega}_{h,n}\|_n \geq L\sigma, G_{\lambda_{\min}}\right) \\ &= \mathbb{P}\left(\exists h \in \mathcal{H} : \|\hat{\omega}_{h,n}\|_n \geq L\sigma, \frac{\mathcal{V}_{h,n}(\hat{\omega}_{h,n})}{\|\hat{\omega}_{h,n}\|_n^2} \geq \eta\sqrt{n}, G_{\lambda_{\min}}\right) \\ &\leq \mathbb{P}\left(\sup_{(g,h) \in \mathcal{G} \times \mathcal{H} : \|\omega_{g,h}\|_n \geq L\sigma} \frac{\mathcal{V}_{h,n}(\omega_{g,h})}{\|\omega_{g,h}\|_n^2} \geq \eta\sqrt{n}, G_{\lambda_{\min}}\right). \end{aligned} \quad (24)$$

The supremum of normalized increments similar to the one appearing above was studied by van de Geer [1990]. In fact, we will adapt Lemma 3.4 of this paper to our purposes. The lemma requires minimal modifications: In our case, the empirical process is indexed with elements of  $\{\omega_{g,h} : g \in \mathcal{G}, h \in \mathcal{H}\}$ , the product set  $\mathcal{G} \times \mathcal{H}$ , whereas van de Geer [1990] considers a similar result for  $h = 0$ . As a result, whereas van de Geer [1990] reduces the study of this probability to bounding the “size” of balls in the the index space, we will reduce it to bounding the size of “tubes”.

To state the generalization of Lemma 3.4 of van de Geer [1990], we introduce the following abstract setting: Let  $(V, d_{V,k}), (\Lambda, d_{\Lambda,k})$  be pseudo-metric spaces ( $k = 1, \dots, n$ ),  $d_k^2$  be the pseudo-metric on  $V \times \Lambda$ , which for  $\gamma = (\nu, \lambda), \tilde{\gamma} = (\tilde{\nu}, \tilde{\lambda})$  in  $V \times \Lambda$  is defined by  $d_k^2(\gamma, \tilde{\gamma}) = d_{V,k}^2(\nu, \tilde{\nu}) + d_{\Lambda,k}^2(\lambda, \tilde{\lambda})$ . Further, let  $d^2$  be the pseudo-metric on  $V \times \Lambda$  defined by  $d^2 = \frac{1}{n} \sum_{k=1}^n d_k^2$ . Consider the real-valued processes  $U_1, U_2, \dots, U_n$  on  $V \times \Lambda$  and the process

$$Z_n = \frac{1}{\sqrt{n}} \sum_{k=1}^n U_k.$$

For  $\sigma > 0$ , denote by  $H(\varepsilon, \sigma) \doteq H(\varepsilon, T(\sigma), d)$ , the metric entropy of the  $\sigma$ -“tube”

$$T(\sigma) = \cup_{\nu \in V} \{\nu\} \times \{\lambda \in \Lambda : d_{\Lambda}(\lambda_{\nu}, \lambda) \leq \sigma\} \subset V \times \Lambda,$$

where for  $\nu \in V, \Lambda_{\nu} \subset \Lambda$  and  $d_{\Lambda}$  (defining the “tube”) is the a pseudo-metric on  $\Lambda$  defined by  $d_{\Lambda}^2(\lambda, \tilde{\lambda}) = \frac{1}{n} \sum_k d_{\Lambda,k}^2(\lambda, \tilde{\lambda})$ . For  $L > 0$ , define

$$\alpha_n(L, \sigma) = \frac{\int_0^1 \sqrt{H(uL\sigma, L\sigma)} du}{\sqrt{n}L\sigma}.$$

With this, we are ready to state our generalization of Lemma 3.4 of van de Geer [1990]:

**Lemma C.5.** *Assume that the following conditions hold:*

(i)  $U_1, U_2, \dots, U_n$  are independent, centered; for all  $\nu \in V, Z_n(\nu, \lambda_{\nu}) = 0$  for some  $\lambda_{\nu} \in \Lambda$ , and

$$|U_k(\gamma) - U_k(\tilde{\gamma})| \leq M_k d_k(\gamma, \tilde{\gamma}), \quad \gamma, \tilde{\gamma} \in V \times \Lambda,$$

where  $M_1, M_2, \dots, M_n$  are uniformly subgaussian, i.e., for some positive  $\beta$  and  $\Gamma$ ,

$$\mathbb{E}[\exp(|\beta M_k|^2)] \leq \Gamma < \infty, k = 1, 2, \dots, n.$$

(ii) Assume that  $\sigma > 0$  is such that  $\sqrt{n}\sigma \geq 1$  and suppose

$$\lim_{L \rightarrow \infty} \alpha_n(L, \sigma) = 0.$$

Then, there exist constants  $L_0 \geq 1$  and  $C_0$ , depending only on  $(\beta, \Gamma)$  and the map  $L \mapsto \alpha_n(L, \sigma)$ , such that for all  $L \geq L_0$ ,

$$\mathbb{P}\left(\sup_{\nu \in V} \sup_{\substack{\lambda \in \Lambda_\nu: \\ d_\Lambda(\lambda_\nu, \lambda) > L\sigma}} \frac{|Z_n(\nu, \lambda)|}{d_\Lambda^2(\lambda_\nu, \lambda)} \geq \sqrt{n}\right) \leq \exp(-C_0 L^2 \sigma^2 n).$$

*Remark C.1.* The proof is obtained by modifying the proof of [van de Geer \[1990\]](#)'s Lemma 3.4 in a straightforward manner and hence it is omitted. A careful investigation of the original proof will find that the result also holds if we find  $L_0$  and  $C_0$  depending on an *upper bound*  $\tilde{\alpha}_n(L, \sigma)$  for  $\alpha_n(L, \sigma)$  provided that  $\lim_{L \rightarrow \infty} \tilde{\alpha}_n(L, \sigma) = 0$  still holds. Moreover, if the upper bound is selected such that it does not depend on  $n$  and  $\sigma$  but only on  $L$  and the "size" of the spaces  $V$ ,  $(\Lambda_\nu)_{\nu \in V}$ , then  $L_0$  and  $C_0$  will depend only on  $(\beta, \Gamma)$  and the mentioned "size".

To apply Lemma C.5 to our problem, we choose the spaces to be  $V = \mathcal{H}$ ,  $\Lambda = \cup_{h \in \mathcal{H}} \Lambda_h$ , where  $\Lambda_h = \Omega_{h,n}$ . Further, we choose the pseudo-metrics to be  $d_{V,k}^2(h, \tilde{h}) = |h(X_k) - \tilde{h}(X_k)|^2 + \|h - \tilde{h}\|_{\infty,n}^2$  ( $h, \tilde{h} \in V$ ), and  $d_{\Lambda,k}(\omega, \tilde{\omega}) = |\omega(X_k) - \tilde{\omega}(X_k)|$  ( $\omega, \tilde{\omega} \in \Lambda$ ). We also choose  $\Lambda_h = \Omega_{h,n} \subset \Lambda$ . Since these pseudo-metrics are random (they depend on  $X_{1:n}$ ), for a proper use of Lemma C.5 we again need to "condition" on  $X_{1:n}$  when using this lemma. Making this argument formal has been discussed in Appendix B.

For  $f \in L^1(\mathcal{X}, P_X)$ ,  $\omega \in \Lambda$ ,  $h \in \mathcal{H}$  set

$$\begin{aligned} \Delta_k(f) &= \frac{1}{\eta} (\ell(Z_k, f) - \mathbb{E}_{x_{1:n}}[\ell(Z_k, f)]), \\ U_k(h, \omega) &= \Delta_k(h + \bar{g}_{h,n}) - \Delta_k(h + \bar{g}_{h,n} + \omega). \end{aligned}$$

(We remind the reader that, although not shown to minimize clutter,  $\Delta_k$  and  $U_k$  do depend on  $x_{1:n}$ .)

Now, for  $h \in \mathcal{H}$ , we set  $\lambda_h = 0$ . Thus,  $U_k(h, \lambda_h) = U_k(h, 0) = 0$ . Furthermore, for  $Z_n(h, \omega) = \frac{1}{\sqrt{n}} \sum_{k=1}^n U_k(h, \omega)$  we have  $Z_n(h, \omega) = \frac{1}{\eta} \mathcal{V}_{h,n}(\omega)$  and therefore (using that  $\lambda_h = 0$  and  $d_\Lambda(\omega, \tilde{\omega}) = \|\omega - \tilde{\omega}\|_n$ )

$$\sup_{h \in \mathcal{H}} \sup_{\substack{\omega \in \Lambda_h: \\ d_\Lambda(\lambda_h, \omega) > L\sigma}} \frac{Z_n(h, \omega)}{d_\Lambda^2(\lambda_h, \omega)} = \sup_{h \in \mathcal{H}} \sup_{\substack{\omega \in \Omega_{h,n}: \\ \|\omega\|_n > L\sigma}} \frac{\mathcal{V}_{h,n}(\omega)}{\eta \|\omega\|_n^2} =: Q_n(L\sigma), \quad (25)$$

showing that the conclusion of the lemma suffices to bound the quantity of interest appearing in (24).

We claim that the condition of Lemma C.5 are satisfied for  $[P_X]$  almost every  $x_{1:n} \in \mathcal{X}^n$  such that  $\lambda_{\min}(x_{1:n}) \doteq \lambda_{\min}(\Phi(x_{1:n})^\top \Phi(x_{1:n})) \geq \lambda_{\min}/2$ . Let  $\mathcal{N} \subset \mathcal{X}^n$  be the  $[P_X]$  null-set where the claim is not required to hold (we will construct  $\mathcal{N}$  in the proof). That  $U_k$  are centered and  $Z_n(h, \lambda_h) = 0$  for any  $h \in \mathcal{H}$  holds by construction. As far as the remaining conditions are concerned, we have:

*Condition (i), the independence of  $(U_k)$ :* This follows from the definition of  $\mathbb{P}_{x_{1:n}}$  and the independence of  $(X_k, Y_k)$ .

*Condition (i), the Lipschitzness of  $U_k$ :* Our goal is to show (for later use) that the Lipschitz coefficients  $M_k$  can be chosen independently of  $n$  and  $x_{1:n}$  as long  $\lambda_{\min}(x_{1:n}) \geq \lambda_{\min}/2$ . For this, we will assume that

$$K \geq 1. \quad (26)$$

Since  $U_k$  is defined as a function of  $\Delta_k$ , we consider the Lipschitzness of  $\Delta_k$  first. Using the definition of  $\Delta_k$  and the Lipschitzness of  $\ell$  (cf. Assumption 3.1(iii)), for any  $f, f' \in L^1(\mathcal{X}, P_X)$  we have

$$|\Delta_k(f) - \Delta_k(f')| \leq \frac{1}{\eta} \left( \frac{|\ell(Z_k, f) - \ell(Z_k, f')|}{|f(X_k) - f'(X_k)|} + \frac{\mathbb{E}[|\ell(Z_k, f) - \ell(Z_k, f')| | X_k]}{|f(X_k) - f'(X_k)|} \right) |f(X_k) - f'(X_k)|.$$

Denote  $\frac{|\ell(Z_k, f) - \ell(Z_k, f')|}{|f(X_k) - f'(X_k)|} + \frac{\mathbb{E}[|\ell(Z_k, f) - \ell(Z_k, f')| | X_k]}{|f(X_k) - f'(X_k)|}$  by  $N_k(f, f')$ . Thus, for  $h, \tilde{h} \in \mathcal{H}$ ,  $\omega, \tilde{\omega} \in \Lambda$ , letting  $f = h + \bar{g}_{h,n}$ ,  $\tilde{f} = \tilde{h} + \bar{g}_{\tilde{h},n}$ ,

$$\begin{aligned} &|U_k(h, \omega) - U_k(\tilde{h}, \tilde{\omega})| \\ &\leq \frac{1}{\eta} N_k(f, \tilde{f}) \left| f(X_k) - \tilde{f}(X_k) \right| \\ &\quad + \frac{1}{\eta} N_k(f + \omega, \tilde{f} + \tilde{\omega}) \left\{ \left| f(X_k) - \tilde{f}(X_k) \right| + |\omega(X_k) - \tilde{\omega}(X_k)| \right\} \end{aligned}$$

Now, by assumption  $|h(x_k)|, |\tilde{h}(x_k)| \leq r$ . From  $\lambda_{\min}(x_{1:n}) \geq \lambda_{\min}/2$ , (18) and Lemma C.2 it follows that  $|\bar{g}_{h,n}(x_k)|, |\bar{g}_{\tilde{h},n}(x_k)| \leq \frac{2\bar{R}}{\lambda_{\min}}$ . Also, by the same argument as in Lemma C.7, again thanks to  $\lambda_{\min}(x_{1:n}) \geq \lambda_{\min}/2$ ,  $|\omega(x_k)|, |\tilde{\omega}(x_k)| \leq \frac{1}{K(\lambda_{\min}/2)^{1/2}} \leq \frac{1}{(\lambda_{\min}/2)^{1/2}}$ , where we used (26). Hence,

$$N_k(f, \tilde{f}) \leq K_\ell \left( Y_k, r + \frac{2\bar{R}}{\lambda_{\min}} \right) + \mathbb{E} \left[ K_\ell \left( Y_k, r + \frac{2\bar{R}}{\lambda_{\min}} \right) \mid X_k \right]$$

and similarly,

$$N_k(f + \omega, \tilde{f} + \tilde{\omega}) \leq K_\ell \left( Y_k, r + \frac{2\bar{R}}{\lambda_{\min}} + \frac{1}{(\lambda_{\min}/2)^{1/2}} \right) + \mathbb{E} \left[ K_\ell \left( Y_k, r + \frac{2\bar{R}}{\lambda_{\min}} + \frac{1}{(\lambda_{\min}/2)^{1/2}} \right) \mid X_k \right]$$

Now,

$$\begin{aligned} |f(x_k) - \tilde{f}(x_k)| &\leq |h(x_k) - \tilde{h}(x_k)| + |\bar{g}_{h,n}(x_k) - \bar{g}_{\tilde{h},n}(x_k)| \\ &\leq |h(x_k) - \tilde{h}(x_k)| + K_h \|h - \tilde{h}\|_{\infty, n}, \end{aligned}$$

where the second inequality follows since by Assumption 3.4,  $h \mapsto \bar{g}_{h,n}(x_k)$  is  $K_h$ -Lipschitz. Therefore, by the choice of  $d_{V,k}$  and  $d_{\Lambda,k}$ ,

$$|U_k(h, \omega) - U_k(\tilde{h}, \tilde{\omega})| \leq \frac{2M_k}{\eta} \left( d_{V,k}(h, \tilde{h}) + d_{\Lambda,k}(\omega, \tilde{\omega}) \right) \leq M'_k d_k \left( (h, \omega), (\tilde{h}, \tilde{\omega}) \right)$$

where  $M_k = 2K_\ell \left( Y_k, r + \frac{2\bar{R}}{\lambda_{\min}} + \frac{1}{(\lambda_{\min}/2)^{1/2}} \right) + 2\mathbb{E} \left[ K_\ell \left( Y_k, r + \frac{2\bar{R}}{\lambda_{\min}} + \frac{1}{(\lambda_{\min}/2)^{1/2}} \right) \mid X_k \right]$ . Note that by Lemma A.1(i), it holds almost surely that  $\mathbb{E} \left[ K_\ell \left( Y_k, r + \frac{2\bar{R}}{\lambda_{\min}} + \frac{1}{(\lambda_{\min}/2)^{1/2}} \right) \mid X_k \right] \leq \frac{1}{\beta} \sqrt{\Gamma_{r+2\bar{R}/\lambda_{\min}+1/(\lambda_{\min}/2)^{1/2}} - 1}$ . Then by Lemma A.1(ii),  $M_k$  is uniformly subgaussian, so is  $M'_k$ .

*Condition (ii):* We want to verify that  $\alpha_n(L, \sigma) \rightarrow 0$  as  $L \rightarrow \infty$  and show that in fact an upper bound  $\tilde{\alpha}(L)$  on  $\alpha_n(L, \sigma)$  which is independent of  $x_{1:n}$ ,  $n$ ,  $K$  and  $\sigma$  exists such that  $\tilde{\alpha}(L) \rightarrow 0$  still holds. Since  $\alpha_n(L, \sigma)$  depends on the entropy numbers  $H(\varepsilon, T(\sigma), d)$  of the tube w.r.t.  $d^2 = \frac{1}{n} \sum_k d_k^2$ , first we need to estimate these entropy numbers. For  $\gamma = (h, \omega)$ ,  $\tilde{\gamma} = (\tilde{h}, \tilde{\omega})$ , we have

$$\begin{aligned} d^2(\gamma, \tilde{\gamma}) &= \frac{1}{n} \sum_k d_{V,k}^2(h, \tilde{h}) + \frac{1}{n} \sum_k d_{\Lambda,k}^2(\omega, \tilde{\omega}) \\ &= \|h - \tilde{h}\|_n^2 + \|h - \tilde{h}\|_{\infty, n}^2 + \|\omega - \tilde{\omega}\|_n^2 \leq 2 \left( \|h - \tilde{h}\|_{\infty, n}^2 + \|\omega - \tilde{\omega}\|_n^2 \right). \end{aligned}$$

Further,  $d_{\Lambda}^2(\omega, \tilde{\omega}) = \|\omega - \tilde{\omega}\|_n^2$  and therefore by the choice  $\Lambda_h = \Omega_{h,n}$  and  $\lambda_h = 0$ ,

$$T(\sigma) = \{(h, \omega) : h \in \mathcal{H}, \omega \in \Omega_{h,n} \text{ s.t. } \|\omega\|_n \leq \sigma\}.$$

Therefore, it suffices to estimate the metric entropy of  $T(\sigma)$  at different scales w.r.t. the pseudo-norm  $\|\cdot\|_T$  defined by  $\|(h, \omega_{g,h})\|_T = \|h\|_{\infty, n} + \|\omega_{g,h}\|_n$ . This is done in the following proposition, which also shows that the integrability assumption is satisfied (the proof is presented in the appendix):

**Proposition C.6.** *Let Assumptions 3.1 to 3.4 hold. Take  $n \geq 1$ ,  $K > 0$ ,  $\varepsilon > 0$ ,  $1 \geq \sigma \geq \varepsilon$  such that  $K\sigma \leq 1/2$ . Then on  $G_{\lambda_{\min}}$ ,*

$$H(\varepsilon, T(\sigma), \|\cdot\|_T) \leq \rho \log(\sigma/\varepsilon) + \rho \log(241) + AH\left(\frac{\varepsilon}{A}, \mathcal{H}, \|\cdot\|_{\infty, n}\right)$$

holds a.s. for some positive (non-random) constant  $A$  that depends only on  $K_h$ .

Furthermore, on  $G_{\lambda_{\min}}$ ,

$$\int_0^1 H^{1/2}(u\sigma, T(\sigma), \|\cdot\|_T) du \leq A' \sqrt{\rho} + \frac{A''}{\sigma},$$

holds a.s. for some universal constant  $A' > 0$  and some non-random constant  $A''$  that depends on  $C_H$  and  $K_h$  only.

Now,  $H(\varepsilon, \sigma) = H(\varepsilon, T(\sigma), d) \leq CH(\varepsilon, T(\sigma), \|\cdot\|_T)$  with some universal constant  $C$ , hence  $H(uL\sigma, L\sigma) \leq CH(uL\sigma, T(L\sigma), \|\cdot\|_T)$  and by the previous result,

$$\int_0^1 H^{1/2}(uL\sigma, L\sigma) du \leq C^{1/2} \int_0^1 H^{1/2}(uL\sigma, T(L\sigma), \|\cdot\|_T) du \leq C' \left(1 + \frac{1}{\sigma}\right) \leq \frac{2C'}{\sigma}$$

where  $C'$  is a constant that is independent of  $L, n, K, \sigma$  and we assumed that  $\sigma \leq 1$ . Hence,

$$\alpha_n(L, \sigma) \leq \frac{2C'}{\sqrt{n}L\sigma^2} \leq \frac{2C'}{L}$$

provided that  $\sqrt{n}\sigma^2 \geq 1$ . Thus, under this condition,  $\alpha_n(L, \sigma) \rightarrow 0$  as  $L \rightarrow \infty$ , as required. Furthermore, the upper bound on  $\alpha_n(L, \sigma)$  is independent of  $x_{1:n}, K, n$  and  $\sigma$ . Therefore,  $L_0$  and  $C_0$  can be selected independently of  $x_{1:n}, K, n$  and  $\sigma$ , finishing the verification of the conditions of Lemma C.5.

Therefore, using (25) we conclude that for any  $L \geq L_0, K, n, \sigma$  such that  $\sqrt{n}\sigma^2 \geq 1$  and  $K\sigma \leq 1/2$  and  $K \geq 1$ , for  $[P_X]$  almost all  $x_{1:n}$  such that  $\lambda_{\min}(x_{1:n}) \geq \lambda_{\min}/2$ ,

$$\mathbb{P}_{x_{1:n}}(Q_n(L\sigma) \geq \sqrt{n}) \leq \exp(-C_0L^2\sigma^2n).$$

Now, by the definition of  $\mathbb{P}_{x_{1:n}}$ ,

$$\begin{aligned} \mathbb{P}(Q_n(L\sigma) \geq \sqrt{n}, G_{\lambda_{\min}}) &= \int \mathbb{P}_{x_{1:n}}(Q_n(L\sigma) \geq \sqrt{n}, G_{\lambda_{\min}}) P_X(dx_{1:n}) \\ &= \int_{\lambda_{\min}(x_{1:n}) \geq \lambda_{\min}/2} \mathbb{P}_{x_{1:n}}(Q_n(L\sigma) \geq \sqrt{n}) P_X(dx_{1:n}) \\ &\leq \int_{\lambda_{\min}(x_{1:n}) \geq \lambda_{\min}/2} \exp(-C_0L^2\sigma^2n) P_X(dx_{1:n}) \leq \exp(-C_0L^2\sigma^2n), \end{aligned}$$

where the second equality follows since  $G_{\lambda_{\min}}$  is  $X_{1:n}$ -measurable.

Hence, by combining (23) and (24), using the definition of  $Q_n(L\sigma)$  in (25) and choosing  $L = L_0$ ,

$$\mathbb{P}\left(\sup_{h \in \mathcal{H}} \|g_{h,n} - \bar{g}_{h,n}\|_n \geq \frac{L_0\sigma}{1 - KL_0\sigma}, G_{\lambda_{\min}}\right) \leq \mathbb{P}(Q_n(L\sigma) \geq \sqrt{n}, G_{\lambda_{\min}}) \leq \exp(-C_0L_0^2\sigma^2n).$$

Choosing  $\sigma = 1/(2L_0)$  and  $K = 1$ , noting that  $n \geq \sigma^{-4}$  then translates into  $n \geq 16L_0^4$  gives that

$$\mathbb{P}\left(\sup_{h \in \mathcal{H}} \|g_{h,n} - \bar{g}_{h,n}\|_n \geq 1, G_{\lambda_{\min}}\right) \leq \exp(-C_0n/4),$$

which is the desired result (we also used that  $L_0 \geq 1$  by assumption and hence  $\sigma \leq 1$  which gives that  $\sqrt{n}\sigma \geq \sqrt{n}\sigma^2 \geq 1$ ).

## C.2 Eigenvalue Bound

**Lemma C.1.** *The following hold:*

(i) *With probability one, for any  $\theta \in \mathbb{R}^d$ ,  $\theta^\top \hat{G}\theta \leq \frac{\theta^\top G\theta}{\lambda_{\min}}$ .*

(ii) *Assuming that  $n \in \mathbb{N}$  and  $\delta \in (0, 1)$  are such that*

$$n \geq \frac{2}{\lambda_{\min} \log\left(\frac{\rho}{\delta}\right)} \log\left(\frac{\rho}{\delta}\right), \tag{19}$$

*where  $\rho$  and  $\lambda_{\min}$  are respectively the rank and the smallest positive eigenvalue of  $G$ , with probability at least  $1 - \delta$ , it holds that  $\hat{\lambda}_{\min} \geq \frac{\lambda_{\min}}{2} > 0$ .*

(iii) *For any  $n, \delta$  satisfying (19), with probability  $1 - \delta$  it holds that for any  $\theta \in \mathbb{R}^d$  and  $[P_X]$  almost every  $x \in \mathcal{X}$ ,*

$$|\langle \theta, \phi(x) \rangle| \leq \sqrt{\frac{2\theta^\top \hat{G}\theta}{\lambda_{\min}}}.$$

*Proof. Part (i):* We first show that  $\text{Ker}(G) \subseteq \text{Ker}(\hat{G})$  holds almost surely: In particular, this can be seen by proving that  $G\theta = 0$  for some  $\theta \in \mathbb{R}^d$  then with probability one,  $\hat{G}\theta = 0$  also holds. Indeed, if the latter did not hold with probability one, then for some  $\varepsilon > 0$ ,  $\mathbb{P}(\theta^\top \hat{G}\theta \geq \varepsilon) > 0$  would hold. Then,  $\theta^\top G\theta = \mathbb{E}[\theta^\top \hat{G}\theta] \geq \varepsilon \mathbb{P}(\theta^\top \hat{G}\theta \geq \varepsilon) > 0$ , which means that  $\theta \notin \text{Ker}(G)$ . Now, if we take a set of vectors  $\{\theta_1, \dots, \theta_m\}$  spanning  $\text{Ker}(G)$ , then on some event  $E$  with  $\mathbb{P}(E) = 1$ ,  $\hat{G}\theta_i = 0$  holds for all  $1 \leq i \leq m$ . Now, on  $E$ ,  $\text{Ker}(G) \subset \text{Ker}(\hat{G})$ . Indeed, take an arbitrary  $\theta \in \text{Ker}(G)$  and expand it using  $\{\theta_i\}$ :  $\theta = \sum_{i=1}^m \lambda_i \theta_i$ . Then,  $\hat{G}\theta = \sum_i \lambda_i \hat{G}\theta_i$  and since  $\hat{G}\theta_i = 0$  simultaneously for all  $i$ , the statement follows.

Now, for proving Part (i), consider the event  $E$  where  $\text{Ker}(G) \subset \text{Ker}(\hat{G})$ . We prove the result on  $E$ : Pick any  $\theta \in \mathbb{R}^d$  and decompose it into  $\theta = \theta_\perp + \theta_\parallel$  such that  $\theta_\perp \perp \text{Im}(G)$  and  $\theta_\parallel \in \text{Im}(G)$ . Hence,  $\theta^\top G\theta = \theta_\parallel^\top G\theta_\parallel$ . Since  $\theta_\perp \in \text{Ker}(G)$  and  $\text{Ker}(G) \subset \text{Ker}(\hat{G})$ , we have  $\hat{G}\theta_\perp = 0$ . Hence,  $\theta^\top \hat{G}\theta = \theta_\parallel^\top \hat{G}\theta_\parallel$ . Now, since  $\|\phi(x)\|_2 \leq 1$  it holds that  $\hat{\lambda}_{\max} \leq 1$ , where  $\hat{\lambda}_{\max}$  denotes the largest eigenvalue of  $\hat{G}$ . Therefore, on  $E$ ,

$$\theta^\top \hat{G}\theta = \theta_\parallel^\top \hat{G}\theta_\parallel \leq \|\theta_\parallel\|_2^2 \leq \frac{\theta_\parallel^\top G\theta_\parallel}{\lambda_{\min}} = \frac{\theta^\top G\theta}{\lambda_{\min}}.$$

Since  $\mathbb{P}(E) = 1$ , the result follows.

*Part (ii):* By the ‘‘Eigenvalue Chernoff Bound’’ (Theorem 4.1) of [Gittens and Tropp \[2011\]](#), with probability at least  $1 - \rho \exp(-n\lambda_{\min}(\varepsilon + (1 - \varepsilon)\log(1 - \varepsilon)))$ ,  $\hat{\lambda}_{\min} \geq (1 - \varepsilon)\lambda_{\min}$ . Choosing  $\varepsilon = 1/2$  gives the result.

*Part (iii):* Fix  $n, \delta$  as required. Let  $E$  be the event where  $\text{Ker}(G) \subset \text{Ker}(\hat{G})$  and let  $F_\delta$  be the event where the inequality of Part (ii) holds. Take the set  $S$  of those  $x \in \text{supp}(P_X)$  where  $\text{Ker}(G) \subset \text{Ker}(\phi(x)\phi(x)^\top)$  holds. It follows from the argument presented in Part (i) that  $P_X(\mathcal{X} \setminus S) = 0$ .

Since  $\mathbb{P}(E \cap F_\delta) \geq 1 - \delta$ , it suffices to prove the statement on  $E \cap F_\delta$ . Hence, in what follows all statements are meant to hold on this event. Pick any  $\theta \in \mathbb{R}^d$ ,  $x \in S$  and decompose  $\theta$  as before. Then, thanks to  $x \in S$  it holds that  $\theta_\perp \in \text{Ker}(\phi(x)\phi(x)^\top)$ . Hence,  $\langle \theta, \phi(x) \rangle^2 = \theta^\top \phi(x)\phi(x)^\top \theta = \theta_\parallel^\top \phi(x)\phi(x)^\top \theta_\parallel = \langle \theta_\parallel, \phi(x) \rangle^2$ . Now, owing to  $\|\phi(x)\|_2 \leq 1$ ,

$$\langle \theta_\parallel, \phi(x) \rangle^2 \leq \|\theta_\parallel\|_2^2 \leq \frac{\theta_\parallel^\top \hat{G}\theta_\parallel}{\hat{\lambda}_{\min}} \leq \frac{2\theta_\parallel^\top \hat{G}\theta_\parallel}{\lambda_{\min}} = \frac{2\theta^\top \hat{G}\theta}{\lambda_{\min}},$$

where the last inequality follows from Part (ii). □

### C.3 Proof of the ‘‘Basic Inequality’’ (Lemma C.4)

We start with a uniform bound for the infinity norm of elements in  $\Omega_{h,n}$ . Let

$$K_\infty = \frac{1}{K(\lambda_{\min}/2)^{1/2}}.$$

Recall that  $G_{\lambda_{\min}}$  is the event when  $\hat{\lambda}_{\min} \geq \lambda_{\min}/2$ .

**Lemma C.7.** *On the event  $G_{\lambda_{\min}}$ ,*

$$\sup_{\omega \in \Omega_{h,n}} \|\omega\|_\infty < K_\infty.$$

*Proof.* Introduce  $\|x\|_M^2 = x^\top Mx$  for  $M$  positive definite. Let  $\bar{\theta}_{h,n}$  be the parameter of  $\bar{g}_{h,n}$ . Thus,

$$|\omega(x)| = \frac{|\langle \phi(x), \theta - \bar{\theta}_{h,n} \rangle|}{1 + K\sqrt{\|\theta - \bar{\theta}_{h,n}\|_{\hat{G}}^2}} \leq \frac{\|\theta - \bar{\theta}_{h,n}\|_2}{1 + K\hat{\lambda}_{\min}^{1/2}\|\theta - \bar{\theta}_{h,n}\|_2} < \frac{1}{K\hat{\lambda}_{\min}^{1/2}} \leq K_\infty.$$

□

With this, we can state the proof of Lemma C.4:

**Lemma C.4** (‘‘Basic Inequality’’). *Let Assumption 3.2 hold. There exists a constant  $\eta$ , such that on the event  $G_{\lambda_{\min}}$ , for any  $h \in \mathcal{H}$ ,*

$$\eta \|\hat{\omega}_{h,n}\|_n^2 \leq \Delta_{h,n}(\bar{g}_{h,n}) - \Delta_{h,n}(\bar{g}_{h,n} + \hat{\omega}_{h,n}).$$

*Proof.* The proof follows the ideas underlying the proof of Lemma 12.2 of the book of van de Geer [2000].

First, we will prove that  $L_{h,n}(\bar{g}_{h,n}) - L_{h,n}(\bar{g}_{h,n} + \hat{\omega}_{h,n}) \geq 0$ . Note that by the definition of  $g_{h,n}$ ,  $L_{h,n}(\bar{g}_{h,n}) - L_{h,n}(g_{h,n}) \geq 0$ . Thus,

$$0 \leq L_{h,n}(\bar{g}_{h,n}) - L_{h,n}(g_{h,n} - \bar{g}_{h,n} + \bar{g}_{h,n}) \leq \frac{1}{\alpha} (L_{h,n}(\bar{g}_{h,n}) - L_{h,n}((1-\alpha)\bar{g}_{h,n} + \alpha g_{h,n}))$$

for any  $0 < \alpha \leq 1$ , because of the convexity of  $L_{h,n}$ . Taking  $\alpha = \frac{1}{1+K\|\bar{g}_{h,n} - g_{h,n}\|_n}$ , the previous inequality gives

$$\frac{1}{\alpha} (L_{h,n}(\bar{g}_{h,n}) - L_{h,n}(\bar{g}_{h,n} + \hat{\omega}_{h,n})) \geq 0. \quad (27)$$

Now take  $\varepsilon > 0$  small enough so that it satisfies Assumption 3.2 and also  $\frac{\varepsilon}{K_\infty} \leq 1$ . Then we have  $\bar{L}_{h,n}(\bar{g}_{h,n} + \hat{\omega}_{h,n}) - \bar{L}_{h,n}(\bar{g}_{h,n}) \geq \frac{\varepsilon}{K_\infty} (\bar{L}_{h,n}(\bar{g}_{h,n} + \hat{\omega}_{h,n}) - \bar{L}_{h,n}(\bar{g}_{h,n}))$  because  $\bar{g}_{h,n}$  is a minimizer of  $\bar{L}_{h,n}$  (thus  $\bar{L}_{h,n}(\bar{g}_{h,n} + \hat{\omega}_{h,n}) - \bar{L}_{h,n}(\bar{g}_{h,n}) > 0$ ) and thus

$$\bar{L}_{h,n}(\bar{g}_{h,n} + \hat{\omega}_{h,n}) - \bar{L}_{h,n}(\bar{g}_{h,n}) \geq \bar{L}_{h,n}(\bar{g}_{h,n} + \frac{\varepsilon}{K_\infty} \hat{\omega}_{h,n}) - \bar{L}_{h,n}(\bar{g}_{h,n}) \geq \frac{\varepsilon^3}{K_\infty^2} \|\hat{\omega}_{h,n}\|_n^2. \quad (28)$$

Here, the first inequality holds by the convexity of  $\bar{L}_{h,n}$ . The second inequality follows from Assumption 3.2 used with  $a = \hat{\omega}_{h,n}|_{X_{1:n}}$ , once we verify that its conditions. That  $a \in [-\varepsilon, \varepsilon]^n$  follows from Lemma C.7, while  $a \in \text{Im}(\Phi)$  follows since both  $g_{h,n}|_{X_{1:n}}$  and  $\bar{g}_{h,n}|_{X_{1:n}}$  satisfy this, by construction. Combining (27) and (28) gives the desired result.  $\square$

#### C.4 Proof of Proposition C.6

The result we want to prove is as follows:

**Proposition C.6.** *Let Assumptions 3.1 to 3.4 hold. Take  $n \geq 1$ ,  $K > 0$ ,  $\varepsilon > 0$ ,  $1 \geq \sigma \geq \varepsilon$  such that  $K\sigma \leq 1/2$ . Then on  $G_{\lambda_{\min}}$ ,*

$$H(\varepsilon, T(\sigma), \|\cdot\|_T) \leq \rho \log(\sigma/\varepsilon) + \rho \log(241) + AH(\frac{\varepsilon}{A}, \mathcal{H}, \|\cdot\|_{\infty, n})$$

holds a.s. for some positive (non-random) constant  $A$  that depends only on  $K_h$ .

Furthermore, on  $G_{\lambda_{\min}}$ ,

$$\int_0^1 H^{1/2}(u\sigma, T(\sigma), \|\cdot\|_T) du \leq A' \sqrt{\rho} + \frac{A''}{\sigma},$$

holds a.s. for some universal constant  $A' > 0$  and some non-random constant  $A''$  that depends on  $C_H$  and  $K_h$  only.

We start by showing that the mapping  $g, h \mapsto (h, \omega_{g,h})$  is Lipschitz w.r.t  $\|\cdot\|_T$  ( $g \in \mathcal{G}$ ,  $h \in \mathcal{H}$ ) as this will allow us to bound the entropy of  $T(\sigma)$  in terms of the entropy of  $\mathcal{H}$  and the entropy of the union of balls in  $\cup_{h \in \mathcal{H}} \Omega_{h,n}$ , in particular  $\cup_{h \in \mathcal{H}} \Omega_{h,n}(\sigma)$ .

**Proposition C.8.** *Let Assumption 3.4 hold. Then, for any  $K, \sigma > 0$  satisfying  $K\sigma \leq 1/2$  and any  $(g_1, h_1), (g_2, h_2) \in \mathcal{G} \times \mathcal{H}$  s.t.  $\|\omega_{g_i, h_i}\|_n \leq \sigma$ ,*

$$\|\omega_{g_1, h_1} - \omega_{g_2, h_2}\|_n \leq K_g \|g_1 - g_2\|_n + K_g K_h \|h_1 - h_2\|_{\infty, n} \quad (29)$$

holds a.s. on the event  $G_{\lambda_{\min}}$ , where  $K_g = 4\sqrt{2}$ .

The constant  $K_h$  appearing in the bound is the Lipschitz constant defined in Assumption 3.4.

*Proof.* Take any  $(g_1, h_1), (g_2, h_2) \in \mathcal{G} \times \mathcal{H}$  with the required property. By the triangle inequality, we have

$$\|\omega_{g_1, h_1} - \omega_{g_2, h_2}\|_T \leq \|\omega_{g_1, h_1} - \omega_{g_2, h_1}\|_T + \|\omega_{g_2, h_1} - \omega_{g_2, h_2}\|_T.$$

Let us consider bounding  $\|\omega_{g_1, h_1} - \omega_{g_2, h_1}\|_T$  as the first step. To minimize clutter, introduce  $h = h_1$ ,  $\omega_i = \omega_{g_i, h}$ ,  $i = 1, 2$ . With this, our goal is to bound  $\|\omega_1 - \omega_2\|_T$ .

We have

$$\begin{aligned} |\omega_1(x) - \omega_2(x)| &= \left| \frac{(g_1 - \bar{g}_{h,n})(x)}{1 + K \|g_1 - \bar{g}_{h,n}\|_n} - \frac{(g_2 - \bar{g}_{h,n})(x)}{1 + K \|g_2 - \bar{g}_{h,n}\|_n} \right| \\ &= \left| g_1(x) \left( \frac{1}{1 + K \|g_1 - \bar{g}_{h,n}\|_n} - \frac{1}{1 + K \|g_2 - \bar{g}_{h,n}\|_n} \right) \right. \\ &\quad \left. + \frac{1}{1 + K \|g_2 - \bar{g}_{h,n}\|_n} (g_1 - g_2)(x) \right. \\ &\quad \left. - \bar{g}_{h,n}(x) \left( \frac{1}{1 + K \|g_1 - \bar{g}_{h,n}\|_n} - \frac{1}{1 + K \|g_2 - \bar{g}_{h,n}\|_n} \right) \right|. \end{aligned}$$

By the triangle inequality,

$$\left| \frac{1}{1 + K \|g_1 - \bar{g}_{h,n}\|_n} - \frac{1}{1 + K \|g_2 - \bar{g}_{h,n}\|_n} \right| \leq K \|g_1 - g_2\|_n.$$

Thus,

$$|\omega_1(x) - \omega_2(x)| \leq K |g_1(x) - \bar{g}_{h,n}(x)| \|g_1 - g_2\|_n + |(g_1 - g_2)(x)|$$

and therefore,

$$\begin{aligned} n \|\omega_1 - \omega_2\|_n^2 &\leq \sum_{i=1}^n \left\{ K |g_1(X_i) - \bar{g}_{h,n}(X_i)| \|g_1 - g_2\|_n + |(g_1 - g_2)(X_i)| \right\}^2 \\ &\leq 2 \sum_{i=1}^n \left\{ K^2 |g_1(X_i) - \bar{g}_{h,n}(X_i)|^2 \|g_1 - g_2\|_n^2 + |(g_1 - g_2)(X_i)|^2 \right\} \\ &\leq 2n(K^2 \|g_1 - \bar{g}_{h,n}\|_n^2 + 1) \|g_1 - g_2\|_n^2. \end{aligned}$$

By Equation (23),

$$\|g_1 - \bar{g}_{h,n}\|_n = \frac{\|\omega_1\|_n}{1 - K \|\omega_1\|_n}.$$

Since  $\omega_1 \in \Omega_{h,n}(\sigma)$ ,  $\|\omega_1\|_n \leq \sigma$  and  $K\sigma < 1$  by assumption,  $\|g_1 - \bar{g}_{h,n}\|_n \leq \frac{\sigma}{1 - K\sigma}$ . Combining this with the bound on  $n \|\omega_1 - \omega_2\|_n^2$ , after simplification we get

$$\begin{aligned} \|\omega_1 - \omega_2\|_n &\leq \sqrt{2 + 2 \left( \frac{K\sigma}{1 - K\sigma} \right)^2} \|g_1 - g_2\|_n \leq \sqrt{2} \left( 1 + \frac{K\sigma}{1 - K\sigma} \right) \|g_1 - g_2\|_n \\ &= \frac{2\sqrt{2}}{1 - K\sigma} \|g_1 - g_2\|_n \leq K_g \|g_1 - g_2\|_n, \end{aligned} \tag{30}$$

where  $K_g = 4\sqrt{2}$  in the last two steps we used that by assumption  $K\sigma \leq 1/2$ .

Let us now consider bounding

$$\|\omega_{g_2, h_1} - \omega_{g_2, h_2}\|_n.$$

Noticing that apart from a sign,  $\bar{g}_{h,n}$  and  $g$  play a symmetric role in the definition of  $\omega_{g,h}$ , following the derivation in the first part we get that, similarly to (30),

$$\|\omega_{g_2, h_1} - \omega_{g_2, h_2}\|_n \leq \frac{2\sqrt{2}}{1 - K\sigma} \|\bar{g}_{h_1, n} - \bar{g}_{h_2, n}\|_n \leq K_g \|\bar{g}_{h_1, n} - \bar{g}_{h_2, n}\|_n.$$

Since by Assumption 3.4,  $\|\bar{g}_{h_1, n} - \bar{g}_{h_2, n}\|_n \leq K_h \|h_1 - h_2\|_{\infty, n}$  holds a.s. on  $G_{\lambda_{\min}}$ , we get

$$\|\omega_{g_2, h_1} - \omega_{g_2, h_2}\|_n \leq K_g K_h \|h_1 - h_2\|_{\infty, n}.$$



Putting together the bounds obtained, we get that on the event  $G_{\lambda_{\min}}$ ,

$$\|\omega_{g_1, h_1} - \omega_{g_2, h_2}\|_n \leq K_g \|g_1 - g_2\|_n + K_g K_h \|h_1 - h_2\|_{\infty, n}$$

as required. □

With this, we can state the proof of Proposition C.6.

*Proof of Proposition C.6.* We can write

$$T(\sigma) = \cup_{h \in \mathcal{H}} \{h\} \times \Omega_{h, n}(\sigma),$$

where

$$\Omega_{h, n}(\sigma) = \{\omega \in \Omega_{h, n} : \|\omega\|_n \leq \sigma\}.$$

We first show that

$$H(\varepsilon, T(\sigma), \|\cdot\|_T) \leq H(\frac{\varepsilon}{2}, \mathcal{H}, \|\cdot\|_{\infty, n}) + H(\frac{\varepsilon}{2}, \Omega_n(\sigma), \|\cdot\|_n), \quad (31)$$

where  $\Omega_n(\sigma) = \cup_{h \in \mathcal{H}} \Omega_{h, n}(\sigma)$ . In short, this follows since  $T(\sigma) \subset \mathcal{H} \times \Omega_n(\sigma)$  and since, by definition,  $\|\cdot\|_T$  is obtained by “summing”  $\|\cdot\|_{\infty, n}$  and  $\|\cdot\|_n$ .

In details, we have: Let  $C$  be an integer s.t.  $C \geq \exp(H(\varepsilon/2, \mathcal{H}, \|\cdot\|_{\infty, n}))$ . Then, there exists  $\{h_1, \dots, h_C\} \subset \mathcal{H}$  such that for any  $h \in \mathcal{H}$ ,  $\|h - h_i\|_{\infty, n} \leq \varepsilon/2$  for some  $i \in \{1, \dots, C\}$ . Similarly, let  $D$  be an integer s.t.

$$D \geq \exp(H(\frac{\varepsilon}{2}, \Omega_n(\sigma), \|\cdot\|_n)) \geq \max_{1 \leq i \leq C} \exp(H(\frac{\varepsilon}{2}, \Omega_{h_i, n}(\sigma), \|\cdot\|_n))$$

and  $\{\omega_1, \dots, \omega_D\} \subset \Omega_n(\sigma)$  be an  $\varepsilon/2$ -net of  $\Omega_n(\sigma)$  w.r.t.  $\|\cdot\|_n$ . Then,

$$\{(h_i, \omega_j) : 1 \leq i \leq C, 1 \leq j \leq D\}$$

is an  $\varepsilon$ -net of  $T(\sigma)$ : To show this pick any  $(h, \omega) \in T(\sigma)$ . Then, take the index  $i$  such that  $\|h - h_i\|_{\infty, n} \leq \varepsilon/2$  and take the index  $j$  such that  $\|\omega - \omega_j\|_n \leq \varepsilon/2$ . Then,  $\|(h, \omega) - (h_i, \omega_j)\|_T = \|h - h_i\|_{\infty, n} + \|\omega - \omega_j\|_n \leq \varepsilon$  as required. This shows that (31) indeed holds.

Next, we bound  $H(\varepsilon, \Omega_n(\sigma), \|\cdot\|_n)$ . We have

$$\begin{aligned} \Omega_n(\sigma) &= \{\omega_{g, h} : h \in \mathcal{H}, g \in \mathcal{G}, \|\omega_{g, h}\|_n \leq \sigma\} \\ &\subset \left\{ \omega_{g, h} : h \in \mathcal{H}, g \in \mathcal{G}, \|g - \bar{g}_{h, n}\|_n \leq \frac{\sigma}{1 - K\sigma} \right\}, \end{aligned} \quad (32)$$

where the containment follows since by Equation (23),  $\|g - \bar{g}_{h, n}\|_n = \frac{\|\omega_{g, h}\|_n}{1 - K\|\omega_{g, h}\|_n}$ . For  $s \geq 0$  define

$$\mathcal{G}_h(s) = \left\{ g \in \mathcal{G} : \|g - \bar{g}_{h, n}\|_n \leq s \right\}.$$

Pick  $\hat{\mathcal{H}} \subset \mathcal{H}$  and an arbitrary “discretization” map  $N : \mathcal{H} \rightarrow \hat{\mathcal{H}}$ . We claim that on  $G_{\lambda_{\min}}$ ,

$$\Omega_n(\sigma) \subset \left\{ \omega_{g, h} : h \in \mathcal{H}, g \in \mathcal{G}_{N(h)} \left( \frac{\sigma}{1 - K\sigma} + K_h \|N(h) - h\|_{\infty, n} \right) \right\} \text{ a.s.} \quad (33)$$

By (32) it suffices to show that for any  $h \in \mathcal{H}$  and  $g \in \mathcal{G}_h \left( \frac{\sigma}{1 - K\sigma} \right)$ ,

$$g \in \mathcal{G}_{N(h)} \left( \frac{\sigma}{1 - K\sigma} + K_h \|N(h) - h\|_{\infty, n} \right) \quad (34)$$

also holds true. For brevity introduce  $h' = N(h)$ . Thanks to the choice  $g$  and Assumption 3.4,

$$\|g - \bar{g}_{h', n}\|_n \leq \|g - \bar{g}_{h, n}\|_n + \|\bar{g}_{h, n} - \bar{g}_{h', n}\|_n \leq \frac{\sigma}{1 - K\sigma} + K_h \|h - h'\|_{\infty, n}$$

holds a.s. on  $G_{\lambda_{\min}}$ , which shows that (34) indeed holds.

The following statements holds a.s. on  $G_{\lambda_{\min}}$  – hence we will not mention this condition to minimize clutter. If  $\hat{\mathcal{H}}$  is an  $\varepsilon/(2K_g K_h)$ -net of  $\mathcal{H}$  w.r.t.  $\|\cdot\|_{\infty, n}$  and  $N(h) = \arg \min_{h' \in \mathcal{H}} \|h - h'\|_{\infty, n}$  then  $K_h \|N(h) - h\|_{\infty, n} \leq \varepsilon/(2K_g) \leq \varepsilon/2$  and therefore for any  $h' \in \hat{\mathcal{H}}$ ,

$$\mathcal{G}_{h'} \left( \frac{\sigma}{1-K\sigma} + K_h \|N(h) - h\|_{\infty, n} \right) \subset \mathcal{G}_{h'} \left( 2\sigma + \frac{\varepsilon}{2} \right).$$

For each  $h \in \hat{\mathcal{H}}$ , let  $\hat{\mathcal{G}}_{h'}$  be an  $\varepsilon/2K_g$ -net of  $\mathcal{G}_{h'}(2\sigma + \frac{\varepsilon}{2})$ . We claim that

$$S = \left\{ \omega_{g', h'} : h' \in \hat{\mathcal{H}}, g' \in \hat{\mathcal{G}}_{h'} \right\}$$

is an  $\varepsilon$ -net of  $\Omega_n(\sigma)$  w.r.t.  $\|\cdot\|_n$ . Indeed, let  $\omega = \omega_{g, h} \in \Omega_n(\sigma)$  arbitrary. Let  $h'$  be the nearest neighbor of  $h$  in  $\hat{\mathcal{H}}$  w.r.t.  $\|\cdot\|_{\infty, n}$  and let  $g'$  be the nearest neighbor of  $g$  in  $\hat{\mathcal{G}}_{h'}$  w.r.t.  $\|\cdot\|_n$ . Note that  $g \in \mathcal{G}_{h'}(2\sigma + \varepsilon/2)$ . Then, by Proposition C.8,

$$\|\omega_{g, h} - \omega_{g', h'}\|_n \leq K_g \|g - g'\|_n + K_g K_h \|h - h'\|_{\infty, n}.$$

Now, because  $g \in \mathcal{G}_{h'}(2\sigma + \varepsilon/2)$  and  $\hat{\mathcal{G}}_{h'}$  is an  $\varepsilon/(2K_g)$ -net of this set,  $K_g \|g - g'\|_n \leq \varepsilon/2$ . Similarly, by the choice of  $\mathcal{H}$ ,  $\|h - h'\|_{\infty, n} \leq \varepsilon/2$ , showing that  $S$  is indeed an  $\varepsilon$ -net of  $\Omega_n(\sigma)$ . Note that the cardinality of  $S$  can be bounded by

$$|S| \leq |\hat{\mathcal{H}}| \max_{h' \in \hat{\mathcal{H}}} |\hat{\mathcal{G}}_{h'}|.$$

Hence,

$$H(\varepsilon, \Omega_n(\sigma), \|\cdot\|_n) \leq H\left(\frac{\varepsilon}{2K_g}, \mathcal{G}_{h_0}(2\sigma + \varepsilon/2), \|\cdot\|_n\right) + H\left(\frac{\varepsilon}{2K_g K_h}, \mathcal{H}, \|\cdot\|_{\infty, n}\right),$$

for an arbitrary  $h_0 \in \mathcal{H}$ , where we used that  $\|\cdot\|_n$  is translation invariant.

Combining this with (31), we get

$$\begin{aligned} H(\varepsilon, T(\sigma), \|\cdot\|_T) &\leq H\left(\frac{\varepsilon}{2}, \Omega_n(\sigma), \|\cdot\|_n\right) + H\left(\frac{\varepsilon}{2}, \mathcal{H}, \|\cdot\|_{\infty, n}\right) \\ &\leq H\left(\frac{\varepsilon}{4K_g}, \mathcal{G}_{h_0}(2\sigma + \varepsilon/2), \|\cdot\|_n\right) + H\left(\frac{\varepsilon}{4K_g K_h}, \mathcal{H}, \|\cdot\|_{\infty, n}\right) + H\left(\frac{\varepsilon}{2}, \mathcal{H}, \|\cdot\|_{\infty, n}\right) \\ &\leq H\left(\frac{\varepsilon}{4K_g}, \mathcal{G}_{h_0}(2\sigma + \varepsilon/2), \|\cdot\|_n\right) + AH\left(\frac{\varepsilon}{A}, \mathcal{H}, \|\cdot\|_{\infty, n}\right) \end{aligned} \quad (35)$$

for  $A$  large enough.

By Corollary 2.6 in the book of van de Geer [2000],  $H(\varepsilon, \mathcal{G}_{h_0}(\sigma), \|\cdot\|_n) \leq \rho \log(\frac{4\sigma + \varepsilon}{\varepsilon})$  a.s.. Hence,

$$H\left(\frac{\varepsilon}{4K_g}, \mathcal{G}_{h_0}(2\sigma + \varepsilon/2), \|\cdot\|_n\right) \leq \rho \log\left(\frac{32\sigma K_g + 8K_g \varepsilon + \varepsilon}{\varepsilon}\right) \leq \rho \log(241) + \rho \log(\sigma/\varepsilon), \text{ a.s..}$$

Here the second inequality follows from bounding the  $\varepsilon$  in numerator by  $\sigma$  (since  $\sigma \geq \varepsilon$ ), and  $K_g < 6$ . Combining this with (35) finishes the proof of the first statement.

To prove the second part, note that  $\int_0^1 (-\log(x))^{1/2} dx = \sqrt{\pi}/2$ . Thus, with  $\mathcal{I}(\sigma, \mathcal{H}) = \int_0^1 H^{1/2}(u\sigma, \mathcal{H}, \|\cdot\|_{\infty, n}) du$ ,

$$\begin{aligned} \int_0^1 H^{1/2}(u\sigma, T(\sigma), \|\cdot\|_T) du &\leq \rho^{1/2} \int_0^1 \log^{1/2}\left(\frac{1}{u}\right) du + \rho^{1/2} \log^{1/2}(241) + A^{1/2} \mathcal{I}\left(\frac{\sigma}{A}, \mathcal{H}\right) \\ &\leq \rho^{1/2} \sqrt{\pi}/2 + \rho^{1/2} \log^{1/2}(241) + A^{1/2} \mathcal{I}\left(\frac{\sigma}{A}, \mathcal{H}\right). \end{aligned}$$

Now, note that

$$\int_0^1 H^{1/2}(u\sigma, \mathcal{H}, \|\cdot\|_{\infty, n}) du = \frac{1}{\sigma} \int_0^\sigma H^{1/2}(v, \mathcal{H}, \|\cdot\|_{\infty, n}) dv \leq \frac{1}{\sigma} \int_0^1 H^{1/2}(v, \mathcal{H}, \|\cdot\|_{\infty, n}) dv \leq C_H/\sigma,$$

where the first inequality follows since  $\sigma \leq 1$ , while the last inequality follows by Assumption 3.3. The desired result follows by choosing  $A' = \sqrt{\pi}/2 + \log^{1/2}(241)$  and  $A'' = A^{3/2} C_H$ .  $\square$