

Appendix

477 **Organization of the Appendix** Section A contains the proofs of the results of the main paper.
 478 Section B contains the details of the numerical illustrations presented in Section 4.3.

479 A Proofs

480 A.1 Proof of Proposition 1

481 The function

$$(t, h) \mapsto \sum_{i=1}^m \theta_i(t) f_i(h)$$

482 is locally Lipschitz-continuous with respect to its first variable and globally Lipschitz-continuous
 483 with respect to its second variable. Therefore the existence and uniqueness of the solution of the
 484 initial value problem (5) for $t \geq 0$ comes as a consequence of the Picard-Lindelöf theorem (see, e.g.,
 485 Luk, 2017 for a self-contained presentation and Arnold, 1992 for a textbook).

486 A.2 Proof of Proposition 2

487 For $x \in \mathcal{X}$, let H be the solution of the initial value problem (5) with parameter θ and with the initial
 488 condition $H_0 = x$. Let us first upper-bound $\|f_i(H_t)\|$ for all $i \in \{1, \dots, m\}$ and $t > 0$. To this aim,
 489 for $t \geq 0$, we have

$$\begin{aligned} \|H_t - H_0\| &= \left\| \int_0^t \sum_{i=1}^m \theta_i(s) f_i(H_s) ds \right\| \\ &\leq \int_0^t \sum_{i=1}^m |\theta_i(s)| \|f_i(H_0)\| ds + \int_0^t \sum_{i=1}^m |\theta_i(s)| \|f_i(H_s) - f_i(H_0)\| ds \\ &\leq M \int_0^t \sum_{i=1}^m |\theta_i(s)| ds + K_f \int_0^t \left(\|H_s - H_0\| \sum_{i=1}^m |\theta_i(s)| \right) ds \\ &\leq tMR_\Theta + K_f R_\Theta \int_0^t \|H_s - H_0\| ds. \end{aligned}$$

490 Next, Grönwall's inequality yields, for $t \in [0, 1]$,

$$\|H_t - H_0\| \leq tMR_\Theta \exp(tK_f R_\Theta) \leq MR_\Theta \exp(K_f R_\Theta).$$

491 Hence

$$\|H_t\| \leq \|H_0\| + \|H_t - H_0\| \leq R_{\mathcal{X}} + MR_\Theta \exp(K_f R_\Theta),$$

492 yielding the first result of the proposition. Furthermore, for any $i \in \{1, \dots, m\}$,

$$\|f_i(H_t)\| \leq \|f_i(H_t) - f_i(H_0)\| + \|f_i(H_0)\| \leq M(K_f R_\Theta \exp(K_f R_\Theta) + 1) =: C.$$

493 Now, let \tilde{H} be the solution of the initial value problem (5) with another parameter $\tilde{\theta}$ and with the
 494 same initial condition $\tilde{H}_0 = x$. Then, for any $t \geq 0$,

$$H_t - \tilde{H}_t = \int_0^t \sum_{i=1}^m \theta_i(s) f_i(H_s) ds - \int_0^t \sum_{i=1}^m \tilde{\theta}_i(s) f_i(\tilde{H}_s) ds.$$

495 Hence

$$\begin{aligned}
\|H_t - \tilde{H}_t\| &= \left\| \int_0^t \sum_{i=1}^m (\theta_i(s) - \tilde{\theta}_i(s)) f_i(H_s) ds + \int_0^t \sum_{i=1}^m \tilde{\theta}_i(s) (f_i(H_s) - f_i(\tilde{H}_s)) ds \right\| \\
&\leq \int_0^t \sum_{i=1}^m |\theta_i(s) - \tilde{\theta}_i(s)| \|f_i(H_s)\| ds + \int_0^t \sum_{i=1}^m |\tilde{\theta}_i(s)| \|f_i(H_s) - f_i(\tilde{H}_s)\| ds \\
&\leq \int_0^t \sum_{i=1}^m |\theta_i(s) - \tilde{\theta}_i(s)| \|f_i(H_s)\| ds + K_f \int_0^t \left(\|H_s - \tilde{H}_s\| \sum_{i=1}^m |\tilde{\theta}_i(s)| \right) ds \\
&\leq tC \|\theta - \tilde{\theta}\|_{1,\infty} + K_f R_\Theta \int_0^t \|H_s - \tilde{H}_s\| ds.
\end{aligned}$$

496 Then Grönwall's inequality implies that, for $t \in [0, 1]$,

$$\begin{aligned}
\|H_t - \tilde{H}_t\| &\leq tC \|\theta - \tilde{\theta}\|_{1,\infty} \exp(tK_f R_\Theta) \\
&\leq M(K_f R_\Theta \exp(K_f R_\Theta) + 1) \exp(K_f R_\Theta) \|\theta - \tilde{\theta}\|_{1,\infty} \\
&\leq 2MK_f R_\Theta \exp(2K_f R_\Theta) \|\theta - \tilde{\theta}\|_{1,\infty}
\end{aligned}$$

497 since $1 \leq K_f R_\Theta \exp(K_f R_\Theta)$ because $K_f \geq 1, R_\Theta \geq 1$.

498 A.3 Proof of Proposition 3

499 We first prove the result for $m = 1$. Let G_x be an $\varepsilon/2K_\Theta$ -grid of $[0, 1]$ and G_y an $\varepsilon/2$ -grid of
500 $[-R_\Theta, R_\Theta]$. Formally, we can take

$$G_x = \left\{ \frac{k\varepsilon}{2K_\Theta}, 0 \leq k \leq \left\lceil \frac{2K_\Theta}{\varepsilon} \right\rceil \right\} \quad \text{and} \quad G_y = \left\{ -R_\Theta + \frac{k\varepsilon}{2}, 1 \leq k \leq \left\lfloor \frac{4R_\Theta}{\varepsilon} \right\rfloor \right\}$$

501 Our cover consists of all functions that start at a point of G_y , are piecewise linear with kinks in G_x ,
502 where each piece has slope $+K_\Theta$ or $-K_\Theta$. Hence our cover is of size

$$\mathcal{N}_1(\varepsilon) = |G_y| 2^{|G_x|} \leq \frac{4R_\Theta}{\varepsilon} 2^{\frac{2K_\Theta}{\varepsilon} + 2} = \frac{16R_\Theta}{\varepsilon} 4^{\frac{K_\Theta}{\varepsilon}}.$$

503 Now take a function $f : [0, 1] \rightarrow \mathbb{R}$ that is uniformly bounded by R_Θ and K_Θ -Lipschitz. We
504 construct a cover member at distance ε from f as follows. Choose a point y_0 in G_y at distance at
505 most $\varepsilon/2$ from $f(0)$. Since $f(0) \in [-R_\Theta, R_\Theta]$, this is clearly possible, except perhaps at the end of
506 the interval. To verify that it is possible at the end of the interval, note that R_Θ is at a distance less
507 than $\varepsilon/2$ of the last element of the grid, since

$$R_\Theta - \left(-R_\Theta + \left\lfloor \frac{4R_\Theta}{\varepsilon} \right\rfloor \frac{\varepsilon}{2} \right) = 2R_\Theta - \left\lfloor \frac{4R_\Theta}{\varepsilon} \right\rfloor \frac{\varepsilon}{2} \in \left[2R_\Theta - \frac{4R_\Theta}{\varepsilon} \frac{\varepsilon}{2}, 2R_\Theta - \left(\frac{4R_\Theta}{\varepsilon} - 1 \right) \frac{\varepsilon}{2} \right] = \left[0, \frac{\varepsilon}{2} \right].$$

508 Then, among the cover members that start at y_0 , choose the one which is closest to f at each point
509 of G_x (in case of equality, pick any one). Let us denote this cover member as \tilde{f} . Let us show
510 recursively that f is at ℓ_∞ -distance at most ε from \tilde{f} . More precisely, let us first show by induction
511 on k that for all $k \in \{0, \dots, \lceil \frac{2K_\Theta}{\varepsilon} \rceil\}$,

$$\left| f\left(\frac{k\varepsilon}{2K_\Theta}\right) - \tilde{f}\left(\frac{k\varepsilon}{2K_\Theta}\right) \right| \leq \frac{\varepsilon}{2}. \quad (15)$$

512 First, $|f(0) - \tilde{f}(0)| \leq \frac{\varepsilon}{2}$. Then, assume that (15) holds for some k . Then we have the following
513 inequalities:

$$\begin{aligned}
\tilde{f}\left(\frac{k\varepsilon}{2K_\Theta}\right) - \varepsilon &\leq f\left(\frac{k\varepsilon}{2K_\Theta}\right) - \frac{\varepsilon}{2} && \text{(by induction)} \\
&\leq f\left(\frac{(k+1)\varepsilon}{2K_\Theta}\right) && (f \text{ is } K_\Theta\text{-Lipschitz)} \\
&\leq f\left(\frac{k\varepsilon}{2K_\Theta}\right) + \frac{\varepsilon}{2} && (f \text{ is } K_\Theta\text{-Lipschitz)} \\
&\leq \tilde{f}\left(\frac{k\varepsilon}{2K_\Theta}\right) + \varepsilon && \text{(by induction)}.
\end{aligned}$$

514 Moreover, by definition, $\tilde{f}(\frac{(k+1)\varepsilon}{K_\Theta})$ is the closest point to $f(\frac{(k+1)\varepsilon}{K_\Theta})$ among

$$\left\{ \tilde{f}\left(\frac{k\varepsilon}{K_\Theta}\right) - \frac{\varepsilon}{2}, \tilde{f}\left(\frac{k\varepsilon}{K_\Theta}\right) + \frac{\varepsilon}{2} \right\}.$$

515 The bounds above show that, among those two points, at least one is at distance no more than $\varepsilon/2$
516 from $f(\frac{(k+1)\varepsilon}{K_\Theta})$. This shows (15) at rank $k + 1$.

517 To conclude, take now $x \in [0, 1]$. There exists $k \in \{0, \dots, \lceil \frac{2K_\Theta}{\varepsilon} \rceil\}$ such that x is at distance at most
518 $\varepsilon/4K_\Theta$ from $\frac{k\varepsilon}{2K_\Theta}$. Again, this is clear except perhaps at the end of the interval, where it is also true
519 since

$$1 - \left\lceil \frac{2K_\Theta}{\varepsilon} \right\rceil \frac{\varepsilon}{2K_\Theta} \leq 1 - \frac{2K_\Theta}{\varepsilon} \frac{\varepsilon}{2K_\Theta} = 0,$$

520 meaning that 1 is located between two elements of the grid G_x , showing that it is at distance at
521 most $\varepsilon/4K_\Theta$ from one element of the grid. Then, we have

$$\begin{aligned} |f(x) - \tilde{f}(x)| &\leq \left| f(x) - f\left(\frac{k\varepsilon}{2K_\Theta}\right) \right| + \left| f\left(\frac{k\varepsilon}{2K_\Theta}\right) - \tilde{f}\left(\frac{k\varepsilon}{2K_\Theta}\right) \right| + \left| \tilde{f}\left(\frac{k\varepsilon}{2K_\Theta}\right) - \tilde{f}(x) \right| \\ &\leq \frac{\varepsilon}{4} + \frac{\varepsilon}{2} + \frac{\varepsilon}{4}, \end{aligned}$$

522 where the first and third terms are upper-bounded because f and \tilde{f} are K_Θ -Lip, while the second
523 term is upper bounded by (15). Hence $\|f - \tilde{f}\|_\infty \leq \varepsilon$, proving the result for $m = 1$.

524 Finally, to prove the result for a general m , note that the Cartesian product of ε/m -covers for each
525 coordinate of θ gives a ε -cover for θ . Indeed, consider such covers and take $\theta \in \Theta$. Since each
526 coordinate of θ is uniformly bounded by R_Θ and K_Θ -Lipschitz, the proof above shows the existence
527 of a cover member $\tilde{\theta}$ such that, for all $i \in \{1, \dots, m\}$, $\|\theta_i - \tilde{\theta}_i\|_\infty \leq \varepsilon/m$. Then

$$\|\theta - \tilde{\theta}\|_{1,\infty} = \sup_{0 \leq t \leq 1} \sum_{i=1}^m |\theta_i(t) - \tilde{\theta}_i(t)| \leq \sup_{0 \leq t \leq 1} \sum_{i=1}^m \|\theta_i - \tilde{\theta}_i\|_\infty \leq \varepsilon.$$

528 As a consequence, we conclude that

$$\mathcal{N}(\varepsilon) \leq \left(\mathcal{N}_1\left(\frac{\varepsilon}{m}\right) \right)^m = \left(\frac{16mR_\Theta}{\varepsilon} \right)^m 4^{\frac{m^2 K_\Theta}{\varepsilon}}.$$

529 Taking the logarithm yields the result.

530 A.4 Proof of Theorem 1

531 First note that, for any $\theta \in \Theta$, $x \in \mathcal{X}$ and $y \in \mathcal{Y}$,

$$|\ell(F_\theta(x), y)| \leq |\ell(F_\theta(x), y) - \ell(y, y)| + |\ell(y, y)| \leq K_\ell \|F_\theta(x) - y\|.$$

532 since, by assumption, ℓ is K_ℓ -Lipschitz with respect to its first variable and $\ell(y, y) = 0$. Thus

$$|\ell(F_\theta(x), y)| \leq K_\ell (\|F_\theta(x)\| + \|y\|) \leq K_\ell (R_\mathcal{X} + MR_\Theta \exp(K_f R_\Theta) + R_\mathcal{Y}) =: \overline{M}.$$

533 by Proposition 2.

534 Now, taking $\delta > 0$, a classical computation involving McDiarmid's inequality (see, e.g., [Wainwright, 2019](#), proof of thm 4.10) yields that, with probability at least $1 - \delta$,

$$\mathcal{R}(\hat{\theta}_n) \leq \widehat{\mathcal{R}}_n(\hat{\theta}_n) + \mathbb{E} \left[\sup_{\theta \in \Theta} |\mathcal{R}(\theta) - \widehat{\mathcal{R}}_n(\theta)| \right] + \frac{\overline{M}\sqrt{2}}{\sqrt{n}} \sqrt{\log \frac{1}{\delta}}.$$

536 Denote $C = 2MK_f R_\Theta \exp(2K_f R_\Theta)$. Then we show that \mathcal{R} and $\widehat{\mathcal{R}}_n$ are CK_ℓ -Lipschitz with
537 respect to $(\theta, \|\cdot\|_{1,\infty})$: for $\theta, \tilde{\theta} \in \Theta$,

$$\begin{aligned} |\mathcal{R}(\theta) - \mathcal{R}(\tilde{\theta})| &\leq \mathbb{E} [|\ell(F_\theta(x), y) - \ell(F_{\tilde{\theta}}(x), y)|] \\ &\leq K_\ell \mathbb{E} [\|F_\theta(x) - F_{\tilde{\theta}}(x)\|] \\ &\leq CK_\ell \|\theta - \tilde{\theta}\|_{1,\infty}, \end{aligned}$$

538 according to Proposition 2. The proof for the empirical risk is very similar.

539 Let now $\varepsilon > 0$ and $\mathcal{N}(\varepsilon)$ be the covering number of Θ endowed with the $(1, \infty)$ -norm. By Proposi-
540 tion 3,

$$\log \mathcal{N}(\varepsilon) \leq m \log \left(\frac{16mR_\Theta}{\varepsilon} \right) + \frac{m^2 K_\Theta \log(4)}{\varepsilon}.$$

541 Take $\theta^{(1)}, \dots, \theta^{(\mathcal{N}(\varepsilon))}$ the associated cover elements. Then, for any $\theta \in \Theta$, denoting $\theta^{(i)}$ the cover
542 element at distance at most ε from θ ,

$$\begin{aligned} |\mathcal{R}(\theta) - \widehat{\mathcal{R}}_n(\theta)| &\leq |\mathcal{R}(\theta) - \mathcal{R}(\theta^{(i)})| + |\mathcal{R}(\theta^{(i)}) - \widehat{\mathcal{R}}_n(\theta^{(i)})| + |\widehat{\mathcal{R}}_n(\theta^{(i)}) - \widehat{\mathcal{R}}_n(\theta)| \\ &\leq 2CK_\ell\varepsilon + \sup_{i \in \{1, \dots, \mathcal{N}(\varepsilon)\}} |\mathcal{R}(\theta^{(i)}) - \widehat{\mathcal{R}}_n(\theta^{(i)})|. \end{aligned}$$

543 Hence

$$\mathbb{E} \left[\sup_{\theta \in \Theta} |\mathcal{R}(\theta) - \widehat{\mathcal{R}}_n(\theta)| \right] \leq 2CK_\ell\varepsilon + \mathbb{E} \left[\sup_{i \in \{1, \dots, \mathcal{N}(\varepsilon)\}} |\mathcal{R}(\theta^{(i)}) - \widehat{\mathcal{R}}_n(\theta^{(i)})| \right].$$

544 Since $\widehat{\mathcal{R}}_n(\theta)$ is the average of n independent random variables, which are each almost surely bounded
545 by \overline{M} , it is \overline{M}/\sqrt{n} sub-Gaussian, hence we have the classical inequality on the expectation of the
546 maximum of sub-Gaussian random variables (see, e.g., [Rigollet and Hütter, 2017](#), Theorem 1.14)

$$\mathbb{E} \left[\sup_{i \in \{1, \dots, \mathcal{N}(\varepsilon)\}} |\mathcal{R}(\theta^{(i)}) - \widehat{\mathcal{R}}_n(\theta^{(i)})| \right] \leq \overline{M} \sqrt{\frac{2 \log(2\mathcal{N}(\varepsilon))}{n}}.$$

547 The remainder of the proof consists in computations to put the result in the required format. More
548 precisely, we have

$$\begin{aligned} \mathbb{E} \left[\sup_{\theta \in \Theta} |\mathcal{R}(\theta) - \widehat{\mathcal{R}}_n(\theta)| \right] &\leq 2CK_\ell\varepsilon + \overline{M} \sqrt{\frac{2 \log(2\mathcal{N}(\varepsilon))}{n}} \\ &\leq 2CK_\ell\varepsilon + \overline{M} \sqrt{\frac{2 \log(2) + 2m \log \left(\frac{16mR_\Theta}{\varepsilon} \right) + \frac{2m^2 K_\Theta \log(4)}{\varepsilon}}{n}} \\ &\leq 2CK_\ell\varepsilon + \overline{M} \sqrt{\frac{2(m+1) \log \left(\frac{16mR_\Theta}{\varepsilon} \right) + \frac{2m^2 K_\Theta \log(4)}{\varepsilon}}{n}}. \end{aligned}$$

549 The third step is valid if $\frac{16mR_\Theta}{\varepsilon} \geq 2$. We will shortly take ε to be equal to $\frac{1}{\sqrt{n}}$, thus this condition
550 holds true under the assumption from the Theorem that $mR_\Theta\sqrt{n} \geq 3$. Hence we obtain

$$\mathcal{R}(\widehat{\theta}_n) \leq \widehat{\mathcal{R}}_n(\widehat{\theta}_n) + 2CK_\ell\varepsilon + \overline{M} \sqrt{\frac{2(m+1) \log \left(\frac{16mR_\Theta}{\varepsilon} \right) + \frac{2m^2 K_\Theta \log(4)}{\varepsilon}}{n}} + \frac{\overline{M}\sqrt{2}}{\sqrt{n}} \sqrt{\log \frac{1}{\delta}}. \quad (16)$$

551 Now denote $\tilde{B} = 2\overline{M}K_f \exp(K_f R_\Theta)$. Then $CK_\ell \leq \tilde{B}$ and $2\overline{M} \leq \tilde{B}$. Taking $\varepsilon = \frac{1}{\sqrt{n}}$, we obtain

$$\begin{aligned} \mathcal{R}(\widehat{\theta}_n) &\leq \widehat{\mathcal{R}}_n(\widehat{\theta}_n) + \frac{2\tilde{B}}{\sqrt{n}} + \frac{\tilde{B}}{2} \sqrt{\frac{2(m+1) \log(16mR_\Theta\sqrt{n})}{n} + \frac{2m^2 K_\Theta \log(4)}{\sqrt{n}}} + \frac{\tilde{B}}{\sqrt{n}} \sqrt{\log \frac{1}{\delta}} \\ &\leq \widehat{\mathcal{R}}_n(\widehat{\theta}_n) + \frac{2\tilde{B}}{\sqrt{n}} + \frac{\tilde{B}}{2} \sqrt{\frac{2(m+1) \log(16mR_\Theta\sqrt{n})}{n}} + \frac{\tilde{B}}{2} \frac{m\sqrt{2K_\Theta \log(4)}}{n^{1/4}} + \frac{\tilde{B}}{\sqrt{n}} \sqrt{\log \frac{1}{\delta}} \\ &\leq \widehat{\mathcal{R}}_n(\widehat{\theta}_n) + \frac{3\tilde{B}}{2} \sqrt{\frac{2(m+1) \log(16mR_\Theta\sqrt{n})}{n}} + \tilde{B} \frac{m\sqrt{K_\Theta}}{n^{1/4}} + \frac{\tilde{B}}{\sqrt{n}} \sqrt{\log \frac{1}{\delta}}, \end{aligned}$$

552 since $2 \leq 2\sqrt{\log(2)} \leq \sqrt{2(m+1) \log(16mR_\Theta\sqrt{n})}$ since $16mR_\Theta\sqrt{n} \geq 2$ by the Theorem's
553 assumptions, and $\sqrt{2 \log(4)} \leq 2$. We finally obtain that

$$\mathcal{R}(\widehat{\theta}_n) \leq \widehat{\mathcal{R}}_n(\widehat{\theta}_n) + 3\tilde{B} \sqrt{\frac{(m+1) \log(mR_\Theta n)}{n}} + \tilde{B} \frac{m\sqrt{K_\Theta}}{n^{1/4}} + \frac{\tilde{B}}{\sqrt{n}} \sqrt{\log \frac{1}{\delta}},$$

554 by noting that $n \geq 9 \max(m^{-2}R_\Theta^{-2}, 1)$ implies that

$$\log(16mR_\Theta\sqrt{n}) \leq 2 \log(mR_\Theta n).$$

555 The result unfolds since the constant B in the Theorem is equal to $3\tilde{B}$.

556 **A.5 Proof of Corollary 1**

557 The corollary is an immediate consequence of Theorem 1. To obtain the result, note that $m = d^2$,
 558 thus in particular $\sqrt{m+1} = \sqrt{d^2+1} \leq d+1$, and besides $\log(R_{\mathcal{W}}d^2n) \leq 2\log(R_{\mathcal{W}}dn)$ since
 559 $R_{\mathcal{W}}n \leq R_{\mathcal{W}}^2n^2$ by assumption on n .

560 **A.6 Proof of Proposition 4**

561 For $x \in \mathcal{X}$, let $(H_k)_{0 \leq k \leq L}$ be the values of the layers defined by the recurrence (10) with the
 562 weights \mathbf{W} and the input $H_0 = x$. We denote by $\|\cdot\|$ the ℓ_2 -norm for vectors and the spectral norm
 563 for matrices. Then, for $k \in \{0, \dots, L-1\}$, we have

$$\|H_{k+1}\| \leq \|H_k\| + \frac{1}{L}\|W_k\sigma(H_k)\| \leq \|H_k\| + \frac{1}{L}\|W_k\|\|\sigma(H_k)\| \leq \left(1 + \frac{K_\sigma R_{\mathcal{W}}}{L}\right)\|H_k\|,$$

564 where the last inequality uses that the spectral norm of a matrix is upper-bounded by its $(1, 1)$ -norm
 565 and that $\sigma(0) = 0$. As a consequence, for any $k \in \{0, \dots, L\}$,

$$\|H_k\| \leq \left(1 + \frac{K_\sigma R_{\mathcal{W}}}{L}\right)^k \|H_0\| \leq \exp(K_\sigma R_{\mathcal{W}})R_{\mathcal{X}} =: C,$$

566 yielding the first claim of the Proposition.

567 Now, let \tilde{H} be the values of the layers (10) with another parameter $\tilde{\mathbf{W}}$ and with the same input
 568 $\tilde{H}_0 = x$. Then, for any $k \in \{0, \dots, L-1\}$,

$$H_{k+1} - \tilde{H}_{k+1} = H_k - \tilde{H}_k + \frac{1}{L}(W_k\sigma(H_k) - \tilde{W}_k\sigma(\tilde{H}_k)).$$

569 Hence, using again that the spectral norm of a matrix is upper-bounded by its $(1, 1)$ -norm and that
 570 $\sigma(0) = 0$,

$$\begin{aligned} \|H_{k+1} - \tilde{H}_{k+1}\| &\leq \|H_k - \tilde{H}_k\| + \frac{1}{L}\|W_k(\sigma(H_k) - \sigma(\tilde{H}_k))\| + \frac{1}{L}\|(W_k - \tilde{W}_k)\sigma(\tilde{H}_k)\| \\ &\leq \left(1 + K_\sigma \frac{R_{\mathcal{W}}}{L}\right)\|H_k - \tilde{H}_k\| + \frac{K_\sigma}{L}\|W_k - \tilde{W}_k\|\|\tilde{H}_k\| \\ &\leq \left(1 + K_\sigma \frac{R_{\mathcal{W}}}{L}\right)\|H_k - \tilde{H}_k\| + \frac{CK_\sigma}{L}\|W_k - \tilde{W}_k\|. \end{aligned}$$

571 Then, dividing by $(1 + K_\sigma \frac{R_{\mathcal{W}}}{L})^{k+1}$ and using the method of differences, we obtain that

$$\begin{aligned} \frac{\|H_k - \tilde{H}_k\|}{(1 + K_\sigma \frac{R_{\mathcal{W}}}{L})^k} &\leq \|H_0 - \tilde{H}_0\| + \frac{CK_\sigma}{L} \sum_{j=0}^{k-1} \frac{\|W_j - \tilde{W}_j\|}{(1 + K_\sigma \frac{R_{\mathcal{W}}}{L})^{j+1}} \\ &\leq \frac{CK_\sigma}{L} \|\mathbf{W} - \tilde{\mathbf{W}}\|_{1,1,\infty} \sum_{j=0}^{k-1} \frac{1}{(1 + K_\sigma \frac{R_{\mathcal{W}}}{L})^{j+1}}. \end{aligned}$$

572 Finally note that

$$\begin{aligned} \sum_{j=0}^{k-1} \frac{(1 + K_\sigma \frac{R_{\mathcal{W}}}{L})^k}{(1 + K_\sigma \frac{R_{\mathcal{W}}}{L})^{j+1}} &= \sum_{j=0}^{k-1} (1 + K_\sigma \frac{R_{\mathcal{W}}}{L})^j \\ &= \frac{L}{K_\sigma R_{\mathcal{W}}} \left((1 + K_\sigma \frac{R_{\mathcal{W}}}{L})^k - 1 \right) \\ &\leq \frac{L}{K_\sigma R_{\mathcal{W}}} (\exp(K_\sigma R_{\mathcal{W}}) - 1). \end{aligned}$$

573 We conclude that

$$\|H_k - \tilde{H}_k\| \leq \frac{C}{R_{\mathcal{W}}} (\exp(K_\sigma R_{\mathcal{W}}) - 1) \|\mathbf{W} - \tilde{\mathbf{W}}\|_{1,1,\infty} \leq \frac{R_{\mathcal{X}}}{R_{\mathcal{W}}} \exp(2K_\sigma R_{\mathcal{W}}) \|\mathbf{W} - \tilde{\mathbf{W}}\|_{1,1,\infty}.$$

574 **A.7 Proof of Proposition 5**

575 For two integers a and b , denote respectively $a//b$ and $a\%b$ the quotient and the remainder of the
576 Euclidian division of a by b . Then, for $\mathbf{W} \in \mathbb{R}^{L \times d \times d}$, let $\phi(\mathbf{W}) : [0, 1] \rightarrow \mathbb{R}^{d^2}$ the piecewise-affine
577 function defined as follows: $\phi(\mathbf{W})$ is affine on every interval $\left[\frac{k}{L}, \frac{k+1}{L}\right]$ for $k \in \{0, \dots, L-1\}$; for
578 $k \in \{1, \dots, L\}$ and $i \in \{1, \dots, d^2\}$,

$$\phi(\mathbf{W})_i\left(\frac{k}{L}\right) = \mathbf{W}_{\frac{k}{L}, (i//d)+1, (i\%d)+1},$$

579 and $\phi(\mathbf{W})_i(0) = \phi(\mathbf{W})_i(1/L)$. Then $\phi(\mathbf{W})$ satisfies two properties. First, it is a linear function
580 of \mathbf{W} . Second, for $\mathbf{W} \in \mathbb{R}^{L \times d \times d}$,

$$\|\phi(\mathbf{W})\|_{1, \infty} = \|\mathbf{W}\|_{1, 1, \infty},$$

581 because, for $x \in [0, 1]$, $\phi(\mathbf{W})(x)$ is a convex combination of two vectors that are bounded in ℓ_1 -
582 norm by $\|\mathbf{W}\|_{1, 1, \infty}$, so it is itself bounded in ℓ_1 -norm by $\|\mathbf{W}\|_{1, 1, \infty}$, implying that $\|\phi(\mathbf{W})\|_{1, \infty} \leq$
583 $\|\mathbf{W}\|_{1, 1, \infty}$. Reciprocally,

$$\|\phi(\mathbf{W})\|_{1, \infty} = \sup_{0 \leq t \leq 1} \|\phi(\mathbf{W})(x)\|_1 \geq \sup_{1 \leq k \leq L} \left\| \phi(\mathbf{W})\left(\frac{k}{L}\right) \right\|_1 = \|\mathbf{W}\|_{1, 1, \infty}.$$

584 Now, take $\mathbf{W} \in \mathcal{W}$. The second property of ϕ implies that $\|\phi(\mathbf{W})\|_{1, \infty} \leq R_{\mathcal{W}}$. Moreover, each
585 coordinate of $\phi(\mathbf{W})$ is $K_{\mathcal{W}}$ -Lipschitz, since the slope of each piece of $\phi(\mathbf{W})_i$ is at most $K_{\mathcal{W}}$. As a
586 consequence, $\phi(\mathcal{W})$ belongs to

$$\Theta_{\mathcal{W}} = \{\theta : [0, 1] \rightarrow \mathbb{R}^{d^2}, \|\theta\|_{1, \infty} \leq R_{\mathcal{W}} \text{ and } \theta_i \text{ is } K_{\mathcal{W}}\text{-Lipschitz for } i \in \{1, \dots, d^2\}\}.$$

587 Therefore $\phi(\mathcal{W})$ is a subset of $\Theta_{\mathcal{W}}$, thus its covering number is less than the one of $\Theta_{\mathcal{W}}$. Moreover,
588 ϕ is clearly injective, thus we can define ϕ^{-1} on its image. Consider an ε -cover $(\theta_1, \dots, \theta_N)$ of
589 $(\phi(\mathcal{W}), \|\cdot\|_{1, \infty})$. Let us show that $(\phi^{-1}(\theta_1), \dots, \phi^{-1}(\theta_N))$ is an ε -cover of $(\mathcal{W}, \|\cdot\|_{1, 1, \infty})$: take
590 $\mathbf{W} \in \mathcal{W}$ and consider θ_i a cover member at distance less than ε from $\phi(\mathbf{W})$. Then

$$\|\mathbf{W} - \phi^{-1}(\theta_i)\|_{1, 1, \infty} = \|\phi(\mathbf{W} - \phi^{-1}(\theta_i))\|_{1, \infty} = \|\phi(\mathbf{W}) - \theta_i\|_{1, \infty} \leq \varepsilon,$$

591 where the second equality holds by linearity of ϕ . Therefore, the covering number of $(\mathcal{W}, \|\cdot\|_{1, 1, \infty})$
592 is upper bounded by the one of $(\phi(\mathcal{W}), \|\cdot\|_{1, \infty})$, which itself is upper bounded by the one of
593 $(\Theta_{\mathcal{W}}, \|\cdot\|_{1, \infty})$, yielding the result by Proposition 3.

594 **A.8 Proof of Theorem 2**

595 The proof structure is the same as the one of Theorem 1, but some constants change. Similarly
596 to (16), we obtain that, if $\frac{16d^2 R_{\mathcal{W}}}{\varepsilon} \geq 2$ (which holds true for $\varepsilon = 1/\sqrt{n}$ and under the assumption of
597 the Theorem),

$$\mathcal{R}(\widehat{\mathbf{W}}_n) \leq \widehat{\mathcal{R}}_n(\widehat{\mathbf{W}}_n) + 2CK_\ell\varepsilon + \overline{M} \sqrt{\frac{2(d^2 + 1) \log\left(\frac{16d^2 R_{\mathcal{W}}}{\varepsilon}\right) + \frac{2d^4 K_{\mathcal{W}}}{\varepsilon} \log(4)}{n}} + \frac{\overline{M}\sqrt{2}}{\sqrt{n}} \sqrt{\log \frac{1}{\delta}},$$

598 with

$$\overline{M} = K_\ell(R_{\mathcal{X}} \exp(K_\sigma R_{\mathcal{W}}) + R_{\mathcal{Y}})$$

599 and

$$C = \frac{R_{\mathcal{X}}}{R_{\mathcal{W}}} \exp(2K_\sigma R_{\mathcal{W}}).$$

600 Finally denote

$$\tilde{B} = 2\overline{M} \max\left(\frac{\exp(K_\sigma R_{\mathcal{W}})}{R_{\mathcal{W}}}, 1\right).$$

601 Then $CK_\ell \leq \tilde{B}$ and $2\overline{M} \leq \tilde{B}$. Taking $\varepsilon = \frac{1}{\sqrt{n}}$, we obtain as in the proof of Theorem 1 that

$$\mathcal{R}(\widehat{\mathbf{W}}_n) \leq \widehat{\mathcal{R}}_n(\widehat{\mathbf{W}}_n) + 3\tilde{B} \sqrt{\frac{(d^2 + 1) \log(d^2 R_{\mathcal{W}} n)}{n}} + \tilde{B} \frac{d^2 \sqrt{K_{\mathcal{W}}}}{n^{1/4}} + \frac{\tilde{B}}{\sqrt{n}} \sqrt{\log \frac{1}{\delta}}.$$

602 for $n \geq 9R_{\mathcal{W}}^{-1} \max(d^{-4}R_{\mathcal{W}}^{-1}, 1)$. Thus

$$\mathcal{R}(\widehat{\mathbf{W}}_n) \leq \widehat{\mathcal{R}}_n(\widehat{\mathbf{W}}_n) + 3\sqrt{2}\tilde{B}(d+1)\sqrt{\frac{\log(dR_{\mathcal{W}}n)}{n}} + \tilde{B}\frac{d^2\sqrt{K_{\mathcal{W}}}}{n^{1/4}} + \frac{\tilde{B}}{\sqrt{n}}\sqrt{\log\frac{1}{\delta}},$$

603 since $\sqrt{d^2+1} \leq d+1$ and $R_{\mathcal{W}}n \leq R_{\mathcal{W}}^2n^2$ by assumption on n . The result unfolds since the
604 constant B in the Theorem is equal to $3\sqrt{2}\tilde{B}$.

605 A.9 Proof of Corollary 2

606 Let

$$A(\mathbf{W}) = \left(\prod_{k=1}^L \left\| I + \frac{1}{L}W_k \right\| \right) \left(\sum_{k=1}^L \frac{\|W_k^T\|_{2,1}^{2/3}}{L^{2/3}\|I + \frac{1}{L}W_k\|^{2/3}} \right)^{3/2},$$

607 where $\|\cdot\|_{2,1}$ denotes the $(2, 1)$ -norm defined as the ℓ_1 -norm of the ℓ_2 -norms of the columns, and I
608 is the identity matrix (and we recall that $\|\cdot\|$ denotes the spectral norm). We apply Theorem 1.1 from
609 [Bartlett et al. \(2017\)](#) by taking as reference matrices the identity matrix. The theorem shows that,
610 under the assumptions of the corollary,

$$\mathbb{P}\left(\arg\max_{1 \leq j \leq d} F_{\mathbf{W}}(x)_j \neq y\right) \leq \widehat{\mathcal{R}}_n(\mathbf{W}) + C\frac{R_{\mathcal{X}}A(\mathbf{W})\log(d)}{\gamma\sqrt{n}} + \frac{C}{\sqrt{n}}\sqrt{\log\frac{1}{\delta}},$$

611 where, as in the corollary, $\widehat{\mathcal{R}}_n(\mathbf{W}) \leq n^{-1} \sum_{i=1}^n \mathbf{1}_{F_{\mathbf{W}}(x_i)_{y_i} \leq \gamma + \max_{j \neq y_i} f(x_i)_j}$ and C is a universal
612 constant. Let us upper bound $A(\mathbf{W})$ to conclude. On the one hand, we have

$$\begin{aligned} \prod_{k=1}^L \left\| I + \frac{1}{L}W_k \right\| &\leq \prod_{k=1}^L \left(\|I\| + \frac{1}{L}\|W_k\| \right) \\ &\leq \prod_{k=1}^L \left(1 + \frac{1}{L}\|W_k\|_{1,1} \right) \\ &\leq \prod_{k=1}^L \left(1 + \frac{1}{L}R_{\mathcal{W}} \right) \\ &\leq \exp(R_{\mathcal{W}}) \end{aligned}$$

613 On the other hand, for any $k \in \{1, \dots, L\}$,

$$\|W_k^T\|_{2,1} \leq \|W_k^T\|_{1,1} \leq R_{\mathcal{W}},$$

614 while

$$\left\| I + \frac{1}{L}W_k \right\| \geq 1 - \frac{1}{L}\|W_k\| \geq 1 - \frac{R_{\mathcal{W}}}{L} \geq \frac{1}{2},$$

615 under the assumption that $L \geq R_{\mathcal{W}}$. All in all, we obtain that

$$A(\mathbf{W}) \leq \exp(R_{\mathcal{W}}) \left(2^{2/3} L^{1/3} R_{\mathcal{W}}^{2/3} \right)^{3/2} = 2R_{\mathcal{W}} \exp(R_{\mathcal{W}}) \sqrt{L},$$

616 which yields the result.

617 B Experimental details

618 Our code is available at [\[XXX\]](#).

619 We use the following model, corresponding to model (10) with additional projections at the beginning
620 and at the end:

$$\begin{aligned} H_0 &= Ax \\ H_{k+1} &= H_k + \frac{1}{L}W_{k+1}\sigma(H_k), \quad 0 \leq k \leq L-1 \\ F_{\mathbf{W}}(x) &= BH_L, \end{aligned}$$

Name	Value
d	30
L	1000
σ	ReLU

Table 1: Values of the model hyper-parameters.

621 where $x \in \mathbb{R}^{768}$ is a vectorized MNIST image, $A \in \mathbb{R}^{d \times 768}$, and $B \in \mathbb{R}^{10 \times d}$. Table 1 gives the
622 value of the hyper-parameters.

623 We use the initialization scheme outlined in Section 4.1: we initialize, for $k \in \{1, \dots, L\}$ and
624 $i, j \in \{1, \dots, d\}$,

$$\mathbf{W}_{k,i,j} = \frac{1}{\sqrt{d}} f_{i,j} \left(\frac{k}{L} \right),$$

625 where $f_{i,j}$ are independent Gaussian processes with a RBF kernel (with bandwidth equal to 0.1).
626 We refer to Marion et al. (2022) and Sander et al. (2022) for further discussion on this initialization
627 scheme. However, A and B are initialized with a more usual scheme, namely with i.i.d. $\mathcal{N}(0, 1/c)$
628 random variables, where c denotes the number of columns of A (resp. B).

629 In Figure 1a, we repeat training 10 times independently. Each time, we perform 30 epochs, and
630 compute after each epoch both the Lipschitz constant of the weights and the generalization gap. This
631 gives 300 pairs (Lipschitz constant, generalization gap), which each corresponds to one dot in the
632 figure. Furthermore, we report results for two setups: when A and B are trained or when they are
633 fixed random matrices.

634 In Figure 1b, A and B are not trained. The reason is to assess the effect of the penalization on \mathbf{W}
635 for a fixed scale of A and B . If we allow A and B to vary, then it is possible that the effect of the
636 penalization might be neutralized by a scale increase of A and B during training.

637 For all experiments, we use the standard MNIST datasplit (60k training samples and 10k testing
638 samples). We train using the cross entropy loss, mini-batches of size 128, and the optimizer Adam
639 (Kingma and Ba, 2015) with default parameters and a learning rate of 0.02.

640 We use PyTorch (Paszke et al., 2019) and PyTorch Lightning for our experiments.

641 The code takes about 60 hours to run on a standard laptop (no GPU).