
Generalization Error Bounds for Deep Unfolding RNNs (Supplementary Material)

Boris Joukovsky ^{1,2}

Tanmoy Mukherjee ^{1,2}

Huynh Van Luong ^{1,2}

Nikos Deligiannis ^{1,2}

¹Department of Electronics and Informatics, Vrije Universiteit Brussel, Pleinlaan 2, B-1050 Brussels, Belgium

²imec, Kapeldreef 75, B-3001 Leuven, Belgium

1 PROXIMAL OPERATOR OF REWEIGHTED-RNN

The proximal operator $\phi_{\frac{\lambda_1}{c}, \frac{\lambda_2}{c}, \hbar}(u)$ is defined as

$$\phi_{\frac{\lambda_1}{c}, \frac{\lambda_2}{c}, \hbar}(u) = \arg \min_{v \in \mathbb{R}^h} \left\{ \frac{1}{c} g(v) + \frac{1}{2} \|v - u\|_2^2 \right\}, \quad (1)$$

where $g(v) = \lambda_1 g|v| + \lambda_2 g|v - \hbar|$.

Following the proof in Appendix of [Luong et al., 2021], $\Phi_{\frac{\lambda_1}{c} \mathbf{g}_l, \frac{\lambda_2}{c} \mathbf{g}_l, \hbar}(u)$ is given by

$$\Phi_{\frac{\lambda_1}{c} \mathbf{g}_l, \frac{\lambda_2}{c} \mathbf{g}_l, \hbar}(u) = \begin{cases} u - \frac{\lambda_1}{c} \mathbf{g}_l - \frac{\lambda_2}{c} \mathbf{g}_l, & \hbar + \frac{\lambda_1}{c} \mathbf{g}_l + \frac{\lambda_2}{c} \mathbf{g}_l < u < \infty \\ \hbar, & \hbar + \frac{\lambda_1}{c} \mathbf{g}_l - \frac{\lambda_2}{c} \mathbf{g}_l \leq u \leq \hbar + \frac{\lambda_1}{c} \mathbf{g}_l + \frac{\lambda_2}{c} \mathbf{g}_l \\ u - \frac{\lambda_1}{c} \mathbf{g}_l + \frac{\lambda_2}{c} \mathbf{g}_l, & \frac{\lambda_1}{c} \mathbf{g}_l - \frac{\lambda_2}{c} \mathbf{g}_l < u < \hbar + \frac{\lambda_1}{c} \mathbf{g}_l - \frac{\lambda_2}{c} \mathbf{g}_l \\ 0, & -\frac{\lambda_1}{c} \mathbf{g}_l - \frac{\lambda_2}{c} \mathbf{g}_l \leq u \leq \frac{\lambda_1}{c} \mathbf{g}_l - \frac{\lambda_2}{c} \mathbf{g}_l \\ u + \frac{\lambda_1}{c} \mathbf{g}_l + \frac{\lambda_2}{c} \mathbf{g}_l, & -\infty < u < -\frac{\lambda_1}{c} \mathbf{g}_l - \frac{\lambda_2}{c} \mathbf{g}_l, \end{cases} \quad (2)$$

for $\hbar \geq 0$ and

$$\Phi_{\frac{\lambda_1}{c} \mathbf{g}_l, \frac{\lambda_2}{c} \mathbf{g}_l, \hbar}(u) = \begin{cases} u - \frac{\lambda_1}{c} \mathbf{g}_l - \frac{\lambda_2}{c} \mathbf{g}_l, & \frac{\lambda_1}{c} \mathbf{g}_l + \frac{\lambda_2}{c} \mathbf{g}_l < u < \infty \\ 0, & -\frac{\lambda_1}{c} \mathbf{g}_l + \frac{\lambda_2}{c} \mathbf{g}_l \leq u \leq \frac{\lambda_1}{c} \mathbf{g}_l + \frac{\lambda_2}{c} \mathbf{g}_l \\ u + \frac{\lambda_1}{c} \mathbf{g}_l - \frac{\lambda_2}{c} \mathbf{g}_l, & \hbar - \frac{\lambda_1}{c} \mathbf{g}_l + \frac{\lambda_2}{c} \mathbf{g}_l < u < -\frac{\lambda_1}{c} \mathbf{g}_l + \frac{\lambda_2}{c} \mathbf{g}_l \\ \hbar, & \hbar - \frac{\lambda_1}{c} \mathbf{g}_l - \frac{\lambda_2}{c} \mathbf{g}_l \leq u \leq \hbar - \frac{\lambda_1}{c} \mathbf{g}_l + \frac{\lambda_2}{c} \mathbf{g}_l \\ u + \frac{\lambda_1}{c} \mathbf{g}_l + \frac{\lambda_2}{c} \mathbf{g}_l, & -\infty < u < \hbar - \frac{\lambda_1}{c} \mathbf{g}_l - \frac{\lambda_2}{c} \mathbf{g}_l \end{cases} \quad (3)$$

for $\hbar < 0$.

2 SUPPORTS FOR RADEMACHER COMPLEXITY CALCULUS

The contraction lemma in Shalev-Shwartz and Ben-David [2014] shows the Rademacher complexity of the composition of a class of functions with ρ -Lipschitz functions.

Proposition 2.1. [Shalev-Shwartz and Ben-David, 2014, Lemma 26.9—Contraction lemma]

Let \mathcal{F} be a set of functions, $\mathcal{F} = \{f : \mathcal{X} \mapsto \mathbb{R}\}$, and Φ_1, \dots, Φ_m , ρ -Lipschitz functions, namely, $|\Phi_i(\alpha) - \Phi_i(\beta)| \leq \rho |\alpha - \beta|$

for all $\alpha, \beta \in \mathbb{R}$ for some $\rho > 0$. For any sample set S of m points $\mathbf{x}_1, \dots, \mathbf{x}_m \in \mathcal{X}$, let $(\Phi \circ f)(\mathbf{x}_i) = \Phi(f(\mathbf{x}_i))$. Then,

$$\frac{1}{m} \mathbb{E}_{\epsilon \in \{\pm 1\}^m} \left[\sup_{f \in \mathcal{F}} \sum_{i=1}^m \epsilon_i (\Phi \circ f)(\mathbf{x}_i) \right] \leq \frac{\rho}{m} \mathbb{E}_{\epsilon \in \{\pm 1\}^m} \left[\sup_{f \in \mathcal{F}} \sum_{i=1}^m \epsilon_i f(\mathbf{x}_i) \right], \quad (4)$$

alternatively, $\mathfrak{R}_S(\Phi \circ \mathcal{F}) \leq \rho \mathfrak{R}_S(\mathcal{F})$, where Φ denotes $\Phi_1(\mathbf{x}_1), \dots, \Phi_m(\mathbf{x}_m)$ for S .

Proposition 2.2. [Mohri et al., 2018, Proposition A.1—Hölder's inequality]

Let $p, q \geq 1$ be conjugate: $\frac{1}{p} + \frac{1}{q} = 1$. Then, for all $\mathbf{x}, \mathbf{y} \in \mathbb{R}^n$,

$$\|\mathbf{x} \cdot \mathbf{y}\|_1 \leq \|\mathbf{x}\|_p \|\mathbf{y}\|_q, \quad (5)$$

with the equality when $|\mathbf{y}_i| = |\mathbf{x}_i|^{p-1}$ for all $i \in [1, n]$.

3 PROOF FOR GENERALIZATION ERROR BOUNDS FOR DEEP UNFOLDING RNNs

Proof. We consider the real-valued family of functions $\mathcal{F}_{d,T} : \mathbb{R}^h \times \mathbb{R}^n \mapsto \mathbb{R}$ for the functions $f_{\mathbf{W}, \mathbf{U}}^{(d)}$ to update $\mathbf{h}_T^{(d)}$ in layer d , time step T , defined as

$$\mathcal{F}_{d,T} = \left\{ (\mathbf{h}_{T-1}^{(d)}, \mathbf{x}_T) \mapsto \Phi(\mathbf{w}_d^T f_{\mathbf{W}, \mathbf{U}}^{(d-1)}(\mathbf{h}_{T-1}^{(d)}, \mathbf{x}_T) + \mathbf{u}_d^T \mathbf{x}_T) : \|\mathbf{W}_d\|_{1,\infty} \leq \alpha_d, \|\mathbf{U}_d\|_{1,\infty} \leq \beta_d \right\}, \quad (6)$$

where $\mathbf{w}_d, \mathbf{u}_d$ are the corresponding rows from $\mathbf{W}_d, \mathbf{U}_d$, respectively, and α_l, β_l , with $1 < l \leq d$, are nonnegative hyper-parameters. For the first layer and the first time step, i.e., $l = 1, t = 1$, the real-valued family of functions, $\mathcal{F}_{1,1} : \mathbb{R}^h \times \mathbb{R}^n \mapsto \mathbb{R}$, for the functions $f_{\mathbf{W}, \mathbf{U}}^{(1)}$ is defined by:

$$\mathcal{F}_{1,1} = \left\{ (\mathbf{h}_0, \mathbf{x}_1) \mapsto \Phi(\mathbf{w}_1^T \mathbf{h}_0 + \mathbf{u}_1^T \mathbf{x}_1) : \|\mathbf{W}_1\|_{1,\infty} \leq \alpha_1, \|\mathbf{U}_1\|_{1,\infty} \leq \beta_1 \right\}, \quad (7)$$

where α_1, β_1 are nonnegative hyper-parameters. We denote the input layer as $f_{\mathbf{W}, \mathbf{U}}^{(0)} = \mathbf{h}_0$ at the first time step. From the definition of Rademacher complexity [Definition 3.1] and the family of functions in (6) and (7), we obtain:

$$\begin{aligned} m \mathfrak{R}_S(\mathcal{F}_{d,T}) &\leq \mathbb{E}_{\epsilon \in \{\pm 1\}^m} \left[\sup_{\substack{\mathbf{W}, \mathbf{U} \\ \|\mathbf{w}_d\|_1 \leq \alpha_d \\ \|\mathbf{u}_d\|_1 \leq \beta_d}} \sum_{i=1}^m \epsilon_i \Phi(\mathbf{w}_d^T f_{\mathbf{W}, \mathbf{U}}^{(d-1)}(\mathbf{h}_{T-1,i}, \mathbf{x}_{T,i}) + \mathbf{u}_d^T \mathbf{x}_{T,i}) \right] \\ &\leq \frac{1}{\lambda} \log \exp \left(\mathbb{E}_{\epsilon \in \{\pm 1\}^m} \left[\sup_{\substack{\mathbf{W}, \mathbf{U} \\ \|\mathbf{w}_d\|_1 \leq \alpha_d \\ \|\mathbf{u}_d\|_1 \leq \beta_d}} \lambda \sum_{i=1}^m \epsilon_i (\mathbf{w}_d^T f_{\mathbf{W}, \mathbf{U}}^{(d-1)}(\mathbf{h}_{T-1,i}, \mathbf{x}_{T,i}) + \mathbf{u}_d^T \mathbf{x}_{T,i}) \right] \right) \end{aligned} \quad (8a)$$

$$\leq \frac{1}{\lambda} \log \mathbb{E}_{\epsilon \in \{\pm 1\}^m} \left[\sup_{\substack{\mathbf{W}, \mathbf{U} \\ \|\mathbf{w}_d\|_1 \leq \alpha_d \\ \|\mathbf{u}_d\|_1 \leq \beta_d}} \exp \left(\lambda \sum_{i=1}^m \epsilon_i (\mathbf{w}_d^T f_{\mathbf{W}, \mathbf{U}}^{(d-1)}(\mathbf{h}_{T-1,i}, \mathbf{x}_{T,i})) + \lambda \sum_{i=1}^m \epsilon_i \mathbf{u}_d^T \mathbf{x}_{T,i} \right) \right] \quad (8a)$$

$$\leq \frac{1}{\lambda} \log \mathbb{E}_{\epsilon \in \{\pm 1\}^m} \left[\sup_{\substack{\mathbf{W}, \mathbf{U} \\ \|\mathbf{w}_d\|_1 \leq \alpha_d}} \exp \left(\lambda \sum_{i=1}^m \epsilon_i (\mathbf{w}_d^T f_{\mathbf{W}, \mathbf{U}}^{(d-1)}(\mathbf{h}_{T-1,i}, \mathbf{x}_{T,i})) \right) \sup_{\|\mathbf{u}_d\|_1 \leq \beta_d} \exp \left(\lambda \sum_{i=1}^m \epsilon_i \mathbf{u}_d^T \mathbf{x}_{T,i} \right) \right], \quad (8b)$$

where $\lambda > 0$ is an arbitrary parameter, Eq. (8a) follows Lemma 2.1 for 1-Lipschitz Φ along with Inequality (20), and (8b) holds by Inequality (17).

For layer $1 \leq l \leq d$ and time step t , let us denote:

$$\Delta_{\mathbf{h}_{t-1}, \mathbf{x}_t}^{(l)} = \sup_{\substack{\mathbf{w}, \mathbf{u} \\ \|\mathbf{w}_l\|_1 \leq \alpha_l}} \exp \left(\lambda \Lambda_l \sum_{i=1}^m \epsilon_i \left(\mathbf{w}_l^T f_{\mathbf{w}, \mathbf{u}}^{(l-1)}(\mathbf{h}_{t-1,i}, \mathbf{x}_{t,i}) \right) \right), \quad (9)$$

$$\Delta_{\mathbf{x}_t}^{(l)} = \sup_{\|\mathbf{u}_l\|_1 \leq \beta_l} \exp \left(\lambda \Lambda_l \sum_{i=1}^m \epsilon_i \left(\mathbf{u}_l^T \mathbf{x}_{t,i} \right) \right), \quad (10)$$

where Λ_l is defined as follows: $\Lambda_d = 1$, $\Lambda_l = \prod_{k=l+1}^d \alpha_k$ with $1 \leq l \leq d-1$, and $\Lambda_0 = \prod_{k=1}^d \alpha_k$.

Following the Hölder's inequality in (5) in case of $p = 1$ and $q = \infty$ applied to \mathbf{w}_l^T and $f_{\mathbf{w}, \mathbf{u}}^{(l-1)}(\mathbf{h}_{t-1,i}, \mathbf{x}_{t,i})$ in (9), respectively, we get:

$$\begin{aligned} \Delta_{\mathbf{h}_{t-1}, \mathbf{x}_t}^{(d)} &\leq \sup_{\substack{\mathbf{w}, \mathbf{u} \\ \|\mathbf{w}_{d-1}\|_{1,\infty} \leq \alpha_{d-1} \\ \|\mathbf{u}_{d-1}\|_{1,\infty} \leq \beta_{d-1}}} \exp \left(\lambda \alpha_d \left\| \sum_{i=1}^m \epsilon_i \Phi \left(\mathbf{W}_{d-1} f_{\mathbf{w}, \mathbf{u}}^{(d-2)}(\mathbf{h}_{t-1,i}, \mathbf{x}_{t,i}) + \mathbf{U}_{d-1} \mathbf{x}_{t,i} \right) \right\|_{\infty} \right) \\ &\leq \sup_{\substack{\mathbf{w}, \mathbf{u} \\ \|\mathbf{w}_{d-1,k}\|_1 \leq \alpha_{d-1} \\ \|\mathbf{u}_{d-1,k}\|_1 \leq \beta_{d-1}}} \exp \left(\lambda \alpha_d \max_{k \in \{1, \dots, h\}} \left| \sum_{i=1}^m \epsilon_i \Phi \left(\mathbf{w}_{d-1,k}^T f_{\mathbf{w}, \mathbf{u}}^{(d-2)}(\mathbf{h}_{t-1,i}, \mathbf{x}_{t,i}) + \mathbf{u}_{d-1,k}^T \mathbf{x}_{t,i} \right) \right| \right) \\ &\leq \sup_{\substack{\mathbf{w}, \mathbf{u} \\ \|\mathbf{w}_{d-1,k}\|_1 \leq \alpha_{d-1} \\ \|\mathbf{u}_{d-1,k}\|_1 \leq \beta_{d-1}}} \exp \left(\lambda \alpha_d \left| \sum_{i=1}^m \epsilon_i \Phi \left(\mathbf{w}_{d-1,k}^T f_{\mathbf{w}, \mathbf{u}}^{(d-2)}(\mathbf{h}_{t-1,i}, \mathbf{x}_{t,i}) + \mathbf{u}_{d-1,k}^T \mathbf{x}_{t,i} \right) \right| \right). \end{aligned} \quad (11)$$

Similarly, from (10), we obtain:

$$\Delta_{\mathbf{x}_t}^{(d)} \leq \sup_{\|\mathbf{u}_d\|_1 \leq \beta_d} \exp \left(\lambda \sum_{i=1}^m \epsilon_i \mathbf{u}_d^T \mathbf{x}_{t,i} \right) \leq \exp \left(\lambda \beta_d \left\| \sum_{i=1}^m \epsilon_i \mathbf{x}_{t,i} \right\|_{\infty} \right) \leq \exp \left(\lambda \beta_d \left| \sum_{i=1}^m \epsilon_i \mathbf{x}_{\tau,i,\kappa} \right| \right), \quad (12)$$

where $\{\tau, \kappa\} = \arg\max_{t \in \{1, \dots, T\}, j \in \{1, \dots, n\}} \left| \sum_{i=1}^m \epsilon_i \mathbf{x}_{t,i,j} \right|$.

From (8b), (11), and (12), we get:

$$\begin{aligned} m \mathfrak{R}_S(\mathcal{F}_{d,T}) &\leq \frac{1}{\lambda} \log \left(\mathbb{E}_{\epsilon \in \{\pm 1\}^m} \left[\sup_{\substack{\mathbf{w}, \mathbf{u} \\ \|\mathbf{w}_{d-1,k}\|_1 \leq \alpha_{d-1} \\ \|\mathbf{u}_{d-1,k}\|_1 \leq \beta_{d-1}}} \exp \left(\lambda \alpha_d \left| \sum_{i=1}^m \epsilon_i \Phi \left(\mathbf{w}_{d-1,k}^T f_{\mathbf{w}, \mathbf{u}}^{(d-2)}(\mathbf{h}_{T-1,i}, \mathbf{x}_{T,i}) + \mathbf{u}_{d-1,k}^T \mathbf{x}_{T,i} \right) \right| + \lambda \beta_d \left| \sum_{i=1}^m \epsilon_i \mathbf{x}_{\tau,i,\kappa} \right| \right| \right] \right) \\ &\leq \frac{1}{\lambda} \log \left(\mathbb{E}_{\epsilon \in \{\pm 1\}^m} \left[\sup_{\substack{\mathbf{w}, \mathbf{u} \\ \|\mathbf{w}_{d-1,k}\|_1 \leq \alpha_{d-1} \\ \|\mathbf{u}_{d-1,k}\|_1 \leq \beta_{d-1}}} \left(\exp \left(\lambda \alpha_d \sum_{i=1}^m \epsilon_i \Phi \left(\mathbf{w}_{d-1,k}^T f_{\mathbf{w}, \mathbf{u}}^{(d-2)}(\mathbf{h}_{T-1,i}, \mathbf{x}_{T,i}) + \mathbf{u}_{d-1,k}^T \mathbf{x}_{T,i} \right) \right) + \lambda \beta_d \sum_{i=1}^m \epsilon_i \mathbf{x}_{\tau,i,\kappa} \right) \right. \right. \\ &\quad \left. \left. + \exp \left(\lambda \alpha_d \sum_{i=1}^m \epsilon_i \Phi \left(\mathbf{w}_{d-1,k}^T f_{\mathbf{w}, \mathbf{u}}^{(d-2)}(\mathbf{h}_{T-1,i}, \mathbf{x}_{T,i}) + \mathbf{u}_{d-1,k}^T \mathbf{x}_{T,i} \right) - \lambda \beta_d \sum_{i=1}^m \epsilon_i \mathbf{x}_{\tau,i,\kappa} \right) \right. \right. \\ &\quad \left. \left. + \exp \left(- \lambda \alpha_d \sum_{i=1}^m \epsilon_i \Phi \left(\mathbf{w}_{d-1,k}^T f_{\mathbf{w}, \mathbf{u}}^{(d-2)}(\mathbf{h}_{T-1,i}, \mathbf{x}_{T,i}) + \mathbf{u}_{d-1,k}^T \mathbf{x}_{T,i} \right) + \lambda \beta_d \sum_{i=1}^m \epsilon_i \mathbf{x}_{\tau,i,\kappa} \right) \right. \right. \\ &\quad \left. \left. + \exp \left(- \lambda \alpha_d \sum_{i=1}^m \epsilon_i \Phi \left(\mathbf{w}_{d-1,k}^T f_{\mathbf{w}, \mathbf{u}}^{(d-2)}(\mathbf{h}_{T-1,i}, \mathbf{x}_{T,i}) + \mathbf{u}_{d-1,k}^T \mathbf{x}_{T,i} \right) - \lambda \beta_d \sum_{i=1}^m \epsilon_i \mathbf{x}_{\tau,i,\kappa} \right) \right) \right] \right) \end{aligned}$$

$$\leq \frac{1}{\lambda} \log \left(4 \mathbb{E}_{\epsilon \in \{\pm 1\}^m} \left[\Delta_{\mathbf{h}_{T-1}, \mathbf{x}_T}^{(d-1)} \Delta_{\mathbf{x}_T}^{(d-1)} \exp \left(\beta_d \lambda \sum_{i=1}^m \epsilon_i \mathbf{x}_{\tau, i, \kappa} \right) \right] \right) \quad (13a)$$

$$\leq \frac{1}{\lambda} \log \left(4^{d-1} \mathbb{E}_{\epsilon \in \{\pm 1\}^m} \left[\Delta_{\mathbf{h}_{T-1}, \mathbf{x}_T}^{(1)} \Delta_{\mathbf{x}_T}^{(1)} \exp \left(\lambda \left(\sum_{l=2}^d \beta_l \Lambda_l \right) \sum_{i=1}^m \epsilon_i \mathbf{x}_{\tau, i, \kappa} \right) \right] \right) \quad (13b)$$

$$\begin{aligned} &\leq \frac{1}{\lambda} \log \left(4^{d-1} \mathbb{E}_{\epsilon \in \{\pm 1\}^m} \left[\exp \left(\lambda \left(\sum_{l=2}^d \beta_l \Lambda_l \right) \sum_{i=1}^m \epsilon_i \mathbf{x}_{\tau, i, \kappa} \right) \right] \sup_{\|\mathbf{w}_1\|_1 \leq \alpha_1} \exp \left(\lambda \Lambda_1 \sum_{i=1}^m \epsilon_i (\mathbf{w}_1^\top \mathbf{h}_{T-1, i}) \right) \right. \\ &\quad \cdot \left. \sup_{\|\mathbf{u}_1\|_1 \leq \beta_1} \exp \left(\lambda \Lambda_1 \sum_{i=1}^m \epsilon_i (\mathbf{u}_1^\top \mathbf{x}_{T, i}) \right) \right] \end{aligned} \quad (13c)$$

$$\begin{aligned} &\leq \frac{1}{\lambda} \log \left(4^{d-1} \mathbb{E}_{\epsilon \in \{\pm 1\}^m} \left[\exp \left(\lambda \left(\sum_{l=2}^d \beta_l \Lambda_l \right) \sum_{i=1}^m \epsilon_i \mathbf{x}_{\tau, i, \kappa} \right) \right] \sup_{\|\mathbf{w}_d\|_1 \leq \alpha_d} \exp \left(\lambda \Lambda_0 \left\| \sum_{i=1}^m \epsilon_i \mathbf{h}_{T-1, i} \right\|_\infty \right) \right. \\ &\quad \cdot \left. \exp \left(\lambda \beta_1 \Lambda_1 \left\| \sum_{i=1}^m \epsilon_i \mathbf{x}_{T, i} \right\|_\infty \right) \right] \end{aligned} \quad (13d)$$

$$\begin{aligned} &\leq \frac{1}{\lambda} \log \left(4^d \mathbb{E}_{\epsilon \in \{\pm 1\}^m} \left[\exp \left(\lambda \left(\sum_{l=1}^d \beta_l \Lambda_l \right) \sum_{i=1}^m \epsilon_i \mathbf{x}_{\tau, i, \kappa} \right) \right. \right. \\ &\quad \cdot \left. \left. \sup_{\substack{\|\mathbf{w}_d\|_1 \leq \alpha_d \\ \|\mathbf{u}_d\|_1 \leq \beta_d}} \exp \left(\lambda \Lambda_0 \sum_{i=1}^m \epsilon_i \Phi(\mathbf{w}_d^\top f_{\mathbf{w}, \mathbf{u}}^{(d-1)}(\mathbf{h}_{T-2, i}, \mathbf{x}_{T-1, i}) + \mathbf{u}_d^\top \mathbf{x}_{T-1, i}) \right) \right] \right), \end{aligned} \quad (13e)$$

where (13a) holds by Inequality (17) and (13b) follows by repeating the process from layer $d-1$ to layer 1 for time step T . Furthermore, (13c) is returned as the beginning of the process for time step $T-1$ and (13d) follows Inequality (5).

Proceeding by repeating the above procedure in (13e) from time step $T-1$ to time step 1, we get:

$$m\mathfrak{R}_S(\mathcal{F}_{d, T}) \leq \frac{1}{\lambda} \log \left(4^{dT} \mathbb{E}_{\epsilon \in \{\pm 1\}^m} \left[\exp \left(\lambda \left(\sum_{l=1}^d \beta_l \Lambda_l \right) \left(\frac{\Lambda_0^T - 1}{\Lambda_0 - 1} \right) \sum_{i=1}^m \epsilon_i \mathbf{x}_{\tau, i, \kappa} \right) \exp \left(\lambda \Lambda_0^T \left\| \sum_{i=1}^m \epsilon_i \mathbf{h}_0 \right\|_\infty \right) \right] \right). \quad (14)$$

Let us denote $\mu = \operatorname{argmax}_{j \in \{1, \dots, h\}} \left| \sum_{i=1}^m \epsilon_i \mathbf{h}_{0, j} \right|$, from (14), we have:

$$\begin{aligned} m\mathfrak{R}_S(\mathcal{F}_{d, T}) &\leq \frac{1}{\lambda} \log \left(4^{dT} \mathbb{E}_{\epsilon \in \{\pm 1\}^m} \left[\exp \left(\lambda \left(\sum_{l=1}^d \beta_l \Lambda_l \right) \left(\frac{\Lambda_0^T - 1}{\Lambda_0 - 1} \right) \sum_{i=1}^m \epsilon_i \mathbf{x}_{\tau, i, \kappa} \right) \exp \left(\lambda \Lambda_0^T \sum_{i=1}^m \epsilon_i \mathbf{h}_{0, \mu} \right) \right] \right) \\ &\leq \frac{2dT \log 2}{\lambda} + \frac{1}{2\lambda} \log \left(\mathbb{E}_{\epsilon \in \{\pm 1\}^m} \left[\exp \left(\lambda \left(\sum_{l=1}^d \beta_l \Lambda_l \right) \left(\frac{\Lambda_0^T - 1}{\Lambda_0 - 1} \right) \sum_{i=1}^m \epsilon_i \mathbf{x}_{\tau, i, \kappa} \right) \exp \left(\lambda \Lambda_0^T \sum_{i=1}^m \epsilon_i \mathbf{h}_{0, \mu} \right) \right] \right)^2 \\ &\leq \frac{2dT \log 2}{\lambda} + \frac{1}{2\lambda} \log \mathbb{E}_{\epsilon \in \{\pm 1\}^m} \left[\exp \left(2\lambda \left(\sum_{l=1}^d \beta_l \Lambda_l \right) \left(\frac{\Lambda_0^T - 1}{\Lambda_0 - 1} \right) \sum_{i=1}^m \epsilon_i \mathbf{x}_{\tau, i, \kappa} \right) \right] + \frac{1}{2\lambda} \log \mathbb{E}_{\epsilon \in \{\pm 1\}^m} \left[\exp \left(2\lambda \Lambda_0^T \sum_{i=1}^m \epsilon_i \mathbf{h}_{0, \mu} \right) \right] \end{aligned} \quad (15a)$$

$$\begin{aligned} &\leq \frac{2dT \log 2}{\lambda} + \frac{1}{2\lambda} \log \sum_{j=1}^n \mathbb{E}_{\epsilon \in \{\pm 1\}^m} \left[\exp \left(2\lambda \left(\sum_{l=1}^d \beta_l \Lambda_l \right) \left(\frac{\Lambda_0^T - 1}{\Lambda_0 - 1} \right) \sum_{i=1}^m \epsilon_i \mathbf{x}_{\tau, i, j} \right) \right] \\ &\quad + \frac{1}{2\lambda} \log \sum_{j=1}^h \mathbb{E}_{\epsilon \in \{\pm 1\}^m} \left[\exp \left(2\lambda \Lambda_0^T \sum_{i=1}^m \epsilon_i \mathbf{h}_{0, j} \right) \right] \end{aligned} \quad (15b)$$

$$\leq \frac{2dT \log 2}{\lambda} + \frac{1}{2\lambda} \log \sum_{j=1}^n \prod_{i=1}^m \mathbb{E}_{\epsilon \in \{\pm 1\}^m} \left[\exp \left(2\lambda \left(\sum_{l=1}^d \beta_l \Lambda_l \right) \left(\frac{\Lambda_0^T - 1}{\Lambda_0 - 1} \right) \epsilon_i \mathbf{x}_{\tau, i, j} \right) \right]$$

$$\begin{aligned}
& + \frac{1}{2\lambda} \log \sum_{j=1}^h \prod_{i=1}^m \mathbb{E}_{\epsilon \in \{\pm 1\}^m} \left[\exp \left(2\lambda \Lambda_0^T \epsilon_i h_{0,j} \right) \right] \\
& \leq \frac{2dT \log 2}{\lambda} + \frac{1}{2\lambda} \log \sum_{j=1}^n \prod_{i=1}^m \left[\frac{1}{2} \exp \left(2\lambda \left(\sum_{l=1}^d \beta_l \Lambda_l \right) \left(\frac{\Lambda_0^T - 1}{\Lambda_0 - 1} \right) \mathbf{x}_{\tau,i,j} \right) + \frac{1}{2} \exp \left(- 2\lambda \left(\sum_{l=1}^d \beta_l \Lambda_l \right) \left(\frac{\Lambda_0^T - 1}{\Lambda_0 - 1} \right) \mathbf{x}_{\tau,i,j} \right) \right] \\
& + \frac{1}{2\lambda} \log \sum_{j=1}^h \prod_{i=1}^m \left[\frac{1}{2} \exp \left(2\lambda \Lambda_0^T h_{0,j} \right) + \frac{1}{2} \exp \left(- 2\lambda \Lambda_0^T h_{0,j} \right) \right] \\
& \leq \frac{2dT \log 2}{\lambda} + \frac{1}{2\lambda} \log \sum_{j=1}^n \left[\exp \left(2\lambda^2 \left(\sum_{l=1}^d \beta_l \Lambda_l \right)^2 \left(\frac{\Lambda_0^T - 1}{\Lambda_0 - 1} \right)^2 \sum_{i=1}^m x_{\tau,i,j}^2 \right) \right] + \frac{1}{2\lambda} \log \sum_{j=1}^h \left[\exp \left(2\lambda^2 \Lambda_0^{2T} \sum_{i=1}^m h_{0,j}^2 \right) \right]
\end{aligned} \tag{15c}$$

$$\begin{aligned}
& \leq \frac{2dT \log 2}{\lambda} + \frac{\log n}{2\lambda} + \lambda \left(\sum_{l=1}^d \beta_l \Lambda_l \right)^2 \left(\frac{\Lambda_0^T - 1}{\Lambda_0 - 1} \right)^2 m B_{\mathbf{x}}^2 + \frac{\log h}{2\lambda} + \lambda \Lambda_0^{2T} m \|\mathbf{h}_0\|_{\infty}^2 \\
& \leq \frac{2dT \log 2 + \log \sqrt{n} + \log \sqrt{h}}{\lambda} + \lambda \left(\left(\sum_{l=1}^d \beta_l \Lambda_l \right)^2 \left(\frac{\Lambda_0^T - 1}{\Lambda_0 - 1} \right)^2 m B_{\mathbf{x}}^2 + \Lambda_0^{2T} m \|\mathbf{h}_0\|_{\infty}^2 \right),
\end{aligned} \tag{15d}$$

where (15a) follows Inequality (19), (15b) holds by replacing with $\sum_{j=1}^n$ and $\sum_{j=1}^h$, respectively. In addition, (15c) follows (18) and (15d) is received by the following definition: At time step t , we define $\mathbf{X}_t \in \mathbb{R}^{n \times m}$, a matrix composed of m columns from the m input vectors $\{\mathbf{x}_{t,i}\}_{i=1}^m$; we also define $\|\mathbf{X}_t\|_{2,\infty} = \sqrt{\max_{k \in \{1, \dots, n\}} \sum_{i=1}^m \mathbf{x}_{t,i,k}^2} \leq \sqrt{m} B_{\mathbf{x}}$, representing the maximum of the ℓ_2 -norms of the rows of matrix \mathbf{X}_t , and $\|\mathbf{h}_0\|_{\infty} = \max_j |h_{0,j}|$.

Choosing $\lambda = \sqrt{\frac{2dT \log 2 + \log \sqrt{n} + \log \sqrt{h}}{\left(\sum_{l=1}^d \beta_l \Lambda_l \right)^2 \left(\frac{\Lambda_0^T - 1}{\Lambda_0 - 1} \right)^2 m B_{\mathbf{x}}^2 + \Lambda_0^{2T} m \|\mathbf{h}_0\|_{\infty}^2}}$, we achieve the upper bound:

$$\mathfrak{R}_S(\mathcal{F}_{d,T}) \leq \sqrt{\frac{2(4dT \log 2 + \log n + \log h)}{m} \left(\left(\sum_{l=1}^d \beta_l \Lambda_l \right)^2 \left(\frac{\Lambda_0^T - 1}{\Lambda_0 - 1} \right)^2 B_{\mathbf{x}}^2 + \Lambda_0^{2T} \|\mathbf{h}_0\|_{\infty}^2 \right)}. \tag{16}$$

It can be noted that $\mathfrak{R}_S(\mathcal{F}_{d,T})$ in (16) is derived for the real-valued functions $\mathcal{F}_{d,T}$. For the vector-valued functions $\mathcal{F}_{d,T} : \mathbb{R}^h \times \mathbb{R}^n \mapsto \mathbb{R}^h$ [in Theorem 3.3] we apply the contraction lemma [Lemma 2.1] to a Lipschitz loss to obtain the complexity of such vector-valued functions by means of the complexity of the real-valued functions. Specifically, in Theorem 3.3 under the assumption of the 1-Lipschitz loss function and from Theorem 3.2, Lemma 2.1, we complete the proof.

□

Supporting inequalities:

(i) If A, B are sets of positive real numbers, then:

$$\sup(AB) = \sup(A) \cdot \sup(B). \tag{17}$$

(ii) Given $x \in \mathbb{R}$, we have:

$$\frac{\exp(x) + \exp(-x)}{2} \leq \exp(x^2/2). \tag{18}$$

(iii) Let X and Y be random variables, the Cauchy–Bunyakovsky–Schwarz inequality gives:

$$(\mathbb{E}[XY])^2 \leq \mathbb{E}[X^2] \cdot \mathbb{E}[Y^2]. \tag{19}$$

(iv) If ψ is a convex function, the Jensen’s inequality gives:

$$\psi(\mathbb{E}[X]) \leq \mathbb{E}[\psi(X)]. \tag{20}$$

4 EXTENSION OF GENERALIZATION ERROR BOUND FOR CLASSIFICATION

Proof. Let $\mathbf{y} = \mathbf{Y}\mathbf{h}_t^{(d)} \equiv y(\mathbf{h}_t^{(d)})$ be a linear classifier with $\mathbf{Y} \in \mathbb{R}^{c \times h}$. Let \mathbf{Y}_i denote the i^{th} row of \mathbf{Y} . Below, we show that each entry $y_i(\mathbf{h}_t^{(d)})$ is ρ -Lipschitz on its input with $\rho = \min(\max_i \|\mathbf{Y}_i\|_2, \max_i \|\mathbf{Y}_i\|_1)$:

$$\forall \mathbf{h}, \mathbf{h}' \in \mathbb{R}^h, \forall i \in \{1, \dots, c\} : |y_i - y'_i| = |\mathbf{Y}_i^T \mathbf{h} - \mathbf{Y}_i^T \mathbf{h}'| = |\mathbf{Y}_i^T (\mathbf{h} - \mathbf{h}')| \quad (21)$$

$$\leq \|\mathbf{Y}_i\|_2 \|\mathbf{h} - \mathbf{h}'\|_2 \quad (22)$$

$$\leq \max_{j \in \{1, \dots, c\}} \|\mathbf{Y}_j\|_2 \|\mathbf{h} - \mathbf{h}'\|_2 \quad (23)$$

In the development above, line 22 was obtained by applying the triangular inequality. Moreover, in line 23, we have identified a unique Lipschitz constant that is valid for all i . Alternatively, we can also write that:

$$\forall \mathbf{h}, \mathbf{h}' \in \mathbb{R}^h, \forall i \in \{1, \dots, c\} : |y_i - y'_i| \leq \|\mathbf{Y}_i\|_1 \|\mathbf{h} - \mathbf{h}'\|_\infty \quad (24)$$

$$\leq \max_{j \in \{1, \dots, c\}} \|\mathbf{Y}_j\|_1 \|\mathbf{h} - \mathbf{h}'\|_\infty \quad (25)$$

$$\leq \max_{j \in \{1, \dots, c\}} \|\mathbf{Y}_j\|_1 \|\mathbf{h} - \mathbf{h}'\|_2 \quad (26)$$

In the development above, line 25 was obtained using Hölder's inequality [see Proposition 2.2] and line 26 was obtained considering that the ℓ_2 norm is an upper bound of the ℓ_∞ norm. Setting $\rho = \min\{\max_i \|\mathbf{Y}_i\|_2, \max_i \|\mathbf{Y}_i\|_1\}$ completes the proof and shows that ρ is a Lipschitz constant for each entry $y_i(\mathbf{h}_t^{(d)})$. To obtain the generalization upper bound proposed in section 4.3 using the ramp loss evaluated on the classification margin $\ell_\gamma(\mathbf{y})$ (which is $\frac{1}{\gamma}$ -Lipschitz), it suffies to apply the contraction lemma twice [see Proposition 2.1], first for the composition with multivariate linear classifier function and secondly with the ℓ_γ loss function, leading to:

$$L_{\mathcal{D}, \gamma}(f) - L_{S, \gamma}(f) \leq \frac{2\rho}{\gamma} \mathfrak{R}_S(\mathcal{F}_{d, T}) + 4\sqrt{\frac{2\log(4/\delta)}{m}} \quad (27)$$

□

References

- Huynh Van Luong, Boris Joukovsky, and Nikos Deligiannis. Designing interpretable recurrent neural networks for video reconstruction via deep unfolding. *IEEE Transactions on Image Processing*, 30:4099–4113, 2021.
- Mehryar Mohri, Afshin Rostamizadeh, and Ameet Talwalkar. *Foundations of Machine Learning, Second edition*. MIT Press, Cambridge, Massachusetts, USA, 2018. ISBN 9780262039406.
- Shai Shalev-Shwartz and Shai Ben-David. *Understanding Machine Learning: From Theory to Algorithms*. Cambridge University Press, New York, NY, USA, 2014.