

A Proof of Theorem 2

Before we begin the proof, let us introduce more notation. Since we only consider a single critical point, for simplicity of notation we denote the critical point as $\boldsymbol{\theta} = (\boldsymbol{w}, \mathbf{V}_1, \mathbf{V}_2, \mathbf{U}_2, \dots, \mathbf{V}_L, \mathbf{U}_L, \mathbf{z})$, without $*$. For $l \in [2 : L]$, let $J_l(x) := \nabla \phi_{\mathbf{z}}^l(\mathbf{U}_l h_{l-1}(x)) \in \mathbb{R}^{n_l \times m_l}$, i.e., $J^l(x)$ is the Jacobian matrix of $\phi_{\mathbf{z}}^l(\cdot)$ evaluated at $\mathbf{U}_l h_{l-1}(x)$, whenever it exists. Also, let $\mathcal{U} := \text{col}([\mathbf{U}_2^T \ \dots \ \mathbf{U}_L^T]) \subsetneq \mathbb{R}^{d_x}$.

The proof is divided into two cases: 1) if $\boldsymbol{w} \notin \mathcal{U}$, and 2) if $\boldsymbol{w} \in \mathcal{U}$. For Case 1, we will show that $\mathfrak{R}(\boldsymbol{\theta}^*) \leq \mathfrak{R}_{\text{lin}}$; we also note that our representation coverage condition $\text{rank}(\mathbb{E}_{(x,y) \sim \mathcal{P}} [\ell''(f_{\boldsymbol{\theta}}(x); y) h_L(x) h_L(x)^T]) = d_x$ is not required for Case 1. For Case 2, we will show that at least one of $\mathfrak{R}(\boldsymbol{\theta}^*) \leq \mathfrak{R}_{\text{lin}}$ or $\lambda_{\min}(\nabla^2 \mathfrak{R}(\boldsymbol{\theta}^*)) < 0$ has to hold.

Case 1: If $\boldsymbol{w} \notin \mathcal{U}$. From standard matrix calculus, we can calculate the partial derivatives of \mathfrak{R} with respect to \boldsymbol{w} and \mathbf{V}_l 's. Since $\boldsymbol{\theta}$ is a critical point we have

$$\begin{aligned} \frac{\partial \mathfrak{R}}{\partial \boldsymbol{w}}(\boldsymbol{\theta}) &= \mathbb{E}[\ell'(f_{\boldsymbol{\theta}}(x); y) h_L(x)] = \mathbf{0}, \\ \frac{\partial \mathfrak{R}}{\partial \mathbf{V}_l}(\boldsymbol{\theta}) &= \mathbb{E} \left[\ell'(f_{\boldsymbol{\theta}}(x); y) \prod_{k=l+1}^L (I + \mathbf{U}_k^T J_k(x)^T \mathbf{V}_k^T) \boldsymbol{w} \phi_{\mathbf{z}}^l(\mathbf{U}_l h_{l-1}(x))^T \right] = \mathbf{0}, \quad l = 2, \dots, L, \\ \frac{\partial \mathfrak{R}}{\partial \mathbf{V}_1}(\boldsymbol{\theta}) &= \mathbb{E} \left[\ell'(f_{\boldsymbol{\theta}}(x); y) \prod_{k=2}^L (I + \mathbf{U}_k^T J_k(x)^T \mathbf{V}_k^T) \boldsymbol{w} \phi_{\mathbf{z}}^1(x)^T \right] = \mathbf{0}. \end{aligned}$$

For $\mathbf{V}_2, \dots, \mathbf{V}_L$, note that we can arrange terms and express the partial derivatives as

$$\frac{\partial \mathfrak{R}}{\partial \mathbf{V}_l}(\boldsymbol{\theta}) = \boldsymbol{w} \mathbb{E} [\ell'(f_{\boldsymbol{\theta}}(x); y) \phi_{\mathbf{z}}^l(\mathbf{U}_l h_{l-1}(x))]^T + \sum_{k=l+1}^L \mathbf{U}_k^T E_k = \mathbf{0}, \quad (1)$$

where $E_k \in \mathbb{R}^{m_l \times m_l}$ are appropriately defined matrices. Note that any column of $\sum_{k=l+1}^L \mathbf{U}_k^T E_k$ is in \mathcal{U} . Since $\boldsymbol{w} \notin \mathcal{U}$, the sum being zero (1) implies that $\mathbb{E} [\ell'(f_{\boldsymbol{\theta}}(x); y) \phi_{\mathbf{z}}^l(\mathbf{U}_l h_{l-1}(x))] = \mathbf{0}$ (because $\boldsymbol{w} \notin \mathcal{U}$ already implies that $\boldsymbol{w} \neq \mathbf{0}$), for all $l \in [2 : L]$. Similarly, we have $\mathbb{E} [\ell'(f_{\boldsymbol{\theta}}(x); y) \phi_{\mathbf{z}}^1(x)] = \mathbf{0}$.

Now, from $\mathbb{E} [\ell'(f_{\boldsymbol{\theta}}(x); y) h_L(x)] = \mathbf{0}$,

$$\begin{aligned} \mathbf{0} &= \mathbb{E} [\ell'(f_{\boldsymbol{\theta}}(x); y) h_L(x)] \\ &= \mathbb{E} [\ell'(f_{\boldsymbol{\theta}}(x); y) (h_{L-1}(x) + \mathbf{V}_L \phi_{\mathbf{z}}^L(\mathbf{U}_L h_{L-1}(x)))] \\ &= \mathbb{E} [\ell'(f_{\boldsymbol{\theta}}(x); y) h_{L-1}(x)] + \mathbf{V}_L \mathbb{E} [\ell'(f_{\boldsymbol{\theta}}(x); y) \phi_{\mathbf{z}}^L(\mathbf{U}_L h_{L-1}(x))] \\ &= \mathbb{E} [\ell'(f_{\boldsymbol{\theta}}(x); y) h_{L-1}(x)] = \dots = \mathbb{E} [\ell'(f_{\boldsymbol{\theta}}(x); y) x]. \end{aligned}$$

Recall that by convexity, $\ell(p; y) - \ell(q; y) \leq \ell'(p; y)(p - q)$. Now for any $t \in \mathbb{R}^{d_x}$, we can apply this inequality for $p = f_{\boldsymbol{\theta}}(x) = \boldsymbol{w}^T h_L(x)$ and $q = t^T x$:

$$\begin{aligned} \mathbb{E} [\ell(f_{\boldsymbol{\theta}}(x); y)] - \mathbb{E} [\ell(t^T x; y)] &\leq \mathbb{E} [\ell'(f_{\boldsymbol{\theta}}(x); y) (\boldsymbol{w}^T h_L(x) - t^T x)] \\ &= \boldsymbol{w}^T \mathbb{E} [\ell'(f_{\boldsymbol{\theta}}(x); y) h_L(x)] - t^T \mathbb{E} [\ell'(f_{\boldsymbol{\theta}}(x); y) x] = 0. \end{aligned}$$

Thus, $\mathbb{E} [\ell(f_{\boldsymbol{\theta}}(x); y)] \leq \mathbb{E} [\ell(t^T x; y)]$ for all t , so taking infimum over t gives $\mathfrak{R}(\boldsymbol{\theta}^*) \leq \mathfrak{R}_{\text{lin}}$.

Case 2: If $\boldsymbol{w} \in \mathcal{U}$. For this case, we will consider the Hessian of \mathfrak{R} with respect to \boldsymbol{w} and \mathbf{V}_l , for each $l \in [L]$. We will show that if $\mathbb{E} [\ell'(f_{\boldsymbol{\theta}}(x); y) \phi_{\mathbf{z}}^l(\mathbf{U}_l h_{l-1}(x))] \neq \mathbf{0}$, then $\lambda_{\min}(\nabla^2 \mathfrak{R}(\boldsymbol{\theta})) < 0$. This implies that if $\mathbb{E} [\ell'(f_{\boldsymbol{\theta}}(x); y) \phi_{\mathbf{z}}^l(\mathbf{U}_l h_{l-1}(x))] = \mathbf{0}$ for all $l \in [L]$, then by the same argument as in Case 1 we have $\mathfrak{R}(\boldsymbol{\theta}^*) \leq \mathfrak{R}_{\text{lin}}$; otherwise, we have $\lambda_{\min}(\nabla^2 \mathfrak{R}(\boldsymbol{\theta})) < 0$.

Because $\boldsymbol{\theta}$ is a twice-differentiable critical point of $\mathfrak{R}(\cdot)$, if we apply perturbation $\boldsymbol{\delta}$ to $\boldsymbol{\theta}$ and do Taylor expansions, what we get is

$$\mathfrak{R}(\boldsymbol{\theta} + \boldsymbol{\delta}) = \mathfrak{R}(\boldsymbol{\theta}) + \frac{1}{2} \boldsymbol{\delta}^T \nabla^2 \mathfrak{R}(\boldsymbol{\theta}) \boldsymbol{\delta} + o(\|\boldsymbol{\delta}\|^2). \quad (2)$$

So, if we apply a particular form of perturbation δ , calculate $\mathfrak{R}(\theta + \delta)$, and then show that the sum of all second-order perturbation terms are negative for such a δ , it is equivalent to showing $\frac{1}{2}\delta^T \nabla^2 \mathfrak{R}(\theta) \delta < 0$, hence $\lambda_{\min}(\nabla^2 \mathfrak{R}(\theta)) < 0$.

Now fix any $l \in [2 : L]$, and consider perturbing w by ϵ and V_l by Δ , while leaving all other parameters unchanged. We will choose $\Delta = \alpha\beta^T$, where $\alpha \in \mathbb{R}^{d_x}$ is chosen from $\alpha \in \mathcal{U}^\perp$, the orthogonal complement of \mathcal{U} , and $\beta \in \mathbb{R}^{n_l}$ will be chosen later. We will now compute $\mathfrak{R}(\theta + \delta)$ directly from the network architecture. The residual block output $h_1(x), \dots, h_{l-1}(x)$ stays invariant after perturbation because their parameters didn't change. For l -th residual block, the output after perturbation, denoted as $\tilde{h}_l(x)$, becomes

$$\tilde{h}_l(x) = h_l(x) + \Delta \phi_z^l(\mathbf{U}_l h_{l-1}(x)).$$

The next residual block output is

$$\begin{aligned} \tilde{h}_{l+1}(x) &= \tilde{h}_l(x) + \mathbf{V}_{l+1} \phi_z^{l+1}(\mathbf{U}_{l+1} \tilde{h}_l(x)) \\ &= h_l(x) + \Delta \phi_z^l(\mathbf{U}_l h_{l-1}(x)) + \mathbf{V}_{l+1} \phi_z^{l+1}(\mathbf{U}_{l+1} h_l(x) + \mathbf{U}_{l+1} \Delta \phi_z^l(\mathbf{U}_l h_{l-1}(x))) \\ &\stackrel{(a)}{=} h_l(x) + \Delta \phi_z^l(\mathbf{U}_l h_{l-1}(x)) + \mathbf{V}_{l+1} \phi_z^{l+1}(\mathbf{U}_{l+1} h_l(x)) \\ &= h_{l+1}(x) + \Delta \phi_z^l(\mathbf{U}_l h_{l-1}(x)), \end{aligned}$$

where (a) used the fact that $\mathbf{U}_{l+1} \Delta = \mathbf{U}_{l+1} \alpha \beta^T = 0$ because $\alpha \in \mathcal{U}^\perp$. We can propagate this up to $\tilde{h}_L(x)$ and similarly show $\tilde{h}_L(x) = h_L(x) + \Delta \phi_z^l(\mathbf{U}_l h_{l-1}(x))$. Using this, the network output after perturbation, denoted as $f_{\theta+\delta}(\cdot)$, is

$$\begin{aligned} f_{\theta+\delta}(x) &= (\mathbf{w} + \epsilon)^T (h_L(x) + \Delta \phi_z^l(\mathbf{U}_l h_{l-1}(x))) \\ &= f_\theta(x) + \epsilon^T h_L(x) + \mathbf{w}^T \Delta \phi_z^l(\mathbf{U}_l h_{l-1}(x)) + \epsilon^T \Delta \phi_z^l(\mathbf{U}_l h_{l-1}(x)) \\ &\stackrel{(b)}{=} f_\theta(x) + \epsilon^T h_L(x) + \epsilon^T \Delta \phi_z^l(\mathbf{U}_l h_{l-1}(x)), \end{aligned}$$

where (b) used $\mathbf{w}^T \Delta = \mathbf{w}^T \alpha \beta^T = 0$ because $\mathbf{w} \in \mathcal{U}$ and $\alpha \in \mathcal{U}^\perp$. Using this, the risk function value after perturbation is

$$\begin{aligned} \mathfrak{R}(\theta + \delta) &= \mathbb{E}[\ell(f_{\theta+\delta}(x); y)] \\ &= \mathbb{E}[\ell(f_\theta(x) + \epsilon^T h_L(x) + \epsilon^T \Delta \phi_z^l(\mathbf{U}_l h_{l-1}(x)); y)] \\ &\stackrel{(c)}{=} \mathbb{E}\left[\ell(f_\theta(x); y) + \ell'(f_\theta(x); y)(\epsilon^T h_L(x) + \epsilon^T \Delta \phi_z^l(\mathbf{U}_l h_{l-1}(x)))\right. \\ &\quad \left. + \frac{1}{2} \ell''(f_\theta(x); y)(\epsilon^T h_L(x))^2 + o(\|\delta\|^2)\right] \\ &\stackrel{(d)}{=} \mathfrak{R}(\theta) + \mathbb{E}\left[\ell'(f_\theta(x); y) \epsilon^T \Delta \phi_z^l(\mathbf{U}_l h_{l-1}(x)) + \frac{1}{2} \ell''(f_\theta(x); y)(\epsilon^T h_L(x))^2\right] + o(\|\delta\|^2), \end{aligned}$$

where (c) used Taylor expansion of $\ell(\cdot; y)$ and (d) used that $\mathbb{E}[\ell'(f_\theta(x); y) h_L(x)] = \frac{\partial \mathfrak{R}}{\partial \mathbf{w}}(\theta) = \mathbf{0}$. Comparing with the expansion (2), the second term in the RHS corresponds to the second-order perturbation $\frac{1}{2} \delta^T \nabla^2 \mathfrak{R}(\theta) \delta$.

Now note that

$$\begin{aligned} &\mathbb{E}\left[\ell'(f_\theta(x); y) \epsilon^T \Delta \phi_z^l(\mathbf{U}_l h_{l-1}(x)) + \frac{1}{2} \ell''(f_\theta(x); y)(\epsilon^T h_L(x))^2\right] \\ &= \epsilon^T \Delta \mathbb{E}\left[\ell'(f_\theta(x); y) \phi_z^l(\mathbf{U}_l h_{l-1}(x))\right] + \frac{1}{2} \epsilon^T \mathbb{E}\left[\ell''(f_\theta(x); y) h_L(x) h_L(x)^T\right] \epsilon. \end{aligned}$$

Let $A := \mathbb{E}[\ell''(f_\theta(x); y) h_L(x) h_L(x)^T]$ and $b := \mathbb{E}[\ell'(f_\theta(x); y) \phi_z^l(\mathbf{U}_l h_{l-1}(x))]$ for simplicity. By the representation coverage condition of the theorem A is full-rank, hence invertible. We can choose $\epsilon = -A^{-1} \Delta b$ to minimize the expression above, then the minimum value we get is $-\frac{1}{2} b^T \Delta^T A^{-1} \Delta b$.

First, note that A is positive definite, and so is A^{-1} . If $b \neq 0$, we can choose $\beta = b$, so $\Delta b = \alpha \beta^T b = \|b\|^2 \alpha \neq \mathbf{0}$, so $-\frac{1}{2} b^T \Delta^T A^{-1} \Delta b < 0$. This proves that $\lambda_{\min}(\nabla^2 \mathfrak{R}(\theta)) < 0$ if $\mathbb{E}[\ell'(f_\theta(x); y) \phi_z^l(\mathbf{U}_l h_{l-1}(x))] \neq 0$, as desired.

The case when $l = 1$ can be done similarly, by perturbing w and V_1 . This finishes the proof.

B Proof of Theorem 4

Since we only consider a single critical point, we denote the critical point as $\theta = (\mathbf{w}, \mathbf{z})$, without $*$. By the same argument as in Case 1 of Proof of Theorem 2, we can use convexity of ℓ to get the following bound:

$$\begin{aligned} \mathbb{E}[\ell(f_{\theta}(x); y)] - \mathbb{E}[\ell(\hat{t}^T x; y)] &\leq \mathbb{E}[\ell'(f_{\theta}(x); y)(\mathbf{w}^T h_L(x) - \hat{t}^T x)] \\ &= (\mathbf{w} - \hat{t})^T \mathbb{E}[\ell'(f_{\theta}(x); y)h_L(x)] + \hat{t}^T \mathbb{E}[\ell'(f_{\theta}(x); y)(h_L(x) - x)] \\ &\stackrel{(a)}{=} \hat{t}^T \mathbb{E}\left[\ell'(f_{\theta}(x); y) \sum_{l=1}^L \phi_{\mathbf{z}}^l(h_{l-1}(x))\right] \\ &\leq \mu \|\hat{t}\| \sum_{l=1}^L \mathbb{E}[\|\phi_{\mathbf{z}}^l(h_{l-1}(x))\|], \end{aligned}$$

where (a) used the fact that $\mathbb{E}[\ell'(f_{\theta}(x); y)h_L(x)] = \frac{\partial \mathfrak{R}}{\partial \mathbf{w}} = 0$. Now, for any fixed $l \in [L]$, using Assumption 5.1 we have

$$\begin{aligned} \|\phi_{\mathbf{z}}^l(h_{l-1}(x))\| &\leq \rho_l \|h_{l-1}(x)\| \\ &\leq \rho_l (\|h_{l-2}(x)\| + \|\phi_{\mathbf{z}}^{l-1}(h_{l-2}(x))\|) \\ &\leq \rho_l (1 + \rho_{l-1}) \|h_{l-2}(x)\| \\ &\leq \dots \leq \rho_l \prod_{k=1}^{l-1} (1 + \rho_k) \|x\|. \end{aligned}$$

Substituting this bound to the one above, we get

$$\mathfrak{R}(\theta) - \mathfrak{R}_{\text{in}} \leq \mu \|\hat{t}\| \mathbb{E}[\|x\|] \sum_{l=1}^L \rho_l \prod_{k=1}^{l-1} (1 + \rho_k) = \mu \|\hat{t}\| \mathbb{E}[\|x\|] \left(\prod_{k=1}^L (1 + \rho_k) - 1 \right).$$

C Proof of Theorem 5

First, we collect the symbols used in this section. Given a real number p , define $[p]_+ := \max\{p, 0\}$ and $[p]_- := \max\{-p, 0\}$. Notice that $|p| = [p]_+ + [p]_-$. Recall that given a vector x , let $\|x\|$ denotes its Euclidean norm. Recall also that given a matrix M , let $\|M\|$ denote its spectral norm, and $\|M\|_F$ denote its Frobenius norm.

The proof is done by a simple induction argument using the ‘‘peeling-off’’ technique used for Rademacher complexity bounds for neural networks. Before we start, let us define the function class of hidden layer representations, for $0 \leq l \leq L$:

$$\mathcal{H}_l := \{h_l : \mathbb{R}^{d_x} \mapsto \mathbb{R}^{d_x} \mid \|\mathbf{V}_j\|_F, \|\mathbf{U}_j\|_F \leq M_j \text{ for all } j \in [l]\},$$

defined with the same bounds as used in \mathcal{F}_L . Note that \mathcal{H}_0 is a singleton with the identity mapping $x \mapsto x$. Also, define \mathcal{F}_l to be the class of functions represented by a l -block ResNet ($0 \leq l \leq L$):

$$\mathcal{F}_l := \{x \mapsto w^T h_l(x) \mid \|w\| \leq 1, h_l \in \mathcal{H}_l\}.$$

Note that if $l = L$, this recovers the definition of \mathcal{F}_L in the theorem statement. Since

$$\mathcal{F}_0 := \{x \mapsto w^T x \mid \|w\| \leq 1\},$$

it is well-known that $\widehat{\mathcal{R}}_n(\mathcal{F}_0|_S) \leq \frac{B}{\sqrt{n}}$. The rest of the proof is done by proving the following:

$$\widehat{\mathcal{R}}_n(\mathcal{F}_l|_S) \leq (1 + 2M_l^2) \widehat{\mathcal{R}}_n(\mathcal{F}_{l-1}|_S),$$

for $l \in [L]$.

Fix any $l \in [L]$. Then, by the definition of Rademacher complexity,

$$\begin{aligned}
n\widehat{\mathcal{R}}_n(\mathcal{F}_l|_S) &= \mathbb{E}_{\epsilon_{1:n}} \left[\sup_{\substack{\|w\| \leq 1, \\ h_l \in \mathcal{H}_l}} \sum_{i=1}^n \epsilon_i w^T h_l(x_i) \right] \\
&= \mathbb{E}_{\epsilon_{1:n}} \left[\sup_{\substack{\|w\| \leq 1, \\ h_{l-1} \in \mathcal{H}_{l-1}}} \sup_{\substack{\|\mathbf{V}_l\|_F \leq M_l, \\ \|\mathbf{U}_l\|_F \leq M_l}} \sum_{i=1}^n \epsilon_i w^T (h_{l-1}(x_i) + \mathbf{V}_l \sigma(\mathbf{U}_l h_{l-1}(x_i))) \right] \\
&\leq \mathbb{E}_{\epsilon_{1:n}} \left[\sup_{\substack{\|w\| \leq 1, \\ h_{l-1} \in \mathcal{H}_{l-1}}} \sum_{i=1}^n \epsilon_i w^T h_{l-1}(x_i) \right] + \underbrace{\mathbb{E}_{\epsilon_{1:n}} \left[\sup_{\substack{\|w\| \leq 1, \\ h_{l-1} \in \mathcal{H}_{l-1}}} \sup_{\substack{\|\mathbf{V}_l\|_F \leq M_l, \\ \|\mathbf{U}_l\|_F \leq M_l}} \sum_{i=1}^n \epsilon_i w^T \mathbf{V}_l \sigma(\mathbf{U}_l h_{l-1}(x_i)) \right]}_{=: \mathcal{A}}.
\end{aligned}$$

The first term in RHS is $n\widehat{\mathcal{R}}_n(\mathcal{F}_{l-1}|_S)$ by definition. It is left to show an upper bound for the second term in RHS, which we will call \mathcal{A} .

First, because $\|w\| \leq 1$ and $\|\mathbf{V}_l\| \leq \|\mathbf{V}_l\|_F \leq M_l$, we have $\|\mathbf{V}_l^T w\| \leq M_l$. So, by using dual norm,

$$\mathcal{A} = \mathbb{E} \left[\sup_{\substack{\|v\| \leq M_l, \\ \|\mathbf{U}_l\|_F \leq M_l, \\ h_{l-1} \in \mathcal{H}_{l-1}}} v^T \sum_{i=1}^n \epsilon_i \sigma(\mathbf{U}_l h_{l-1}(x_i)) \right] = M_l \mathbb{E} \left[\sup_{\substack{\|\mathbf{U}_l\|_F \leq M_l, \\ h_{l-1} \in \mathcal{H}_{l-1}}} \left\| \sum_{i=1}^n \epsilon_i \sigma(\mathbf{U}_l h_{l-1}(x_i)) \right\| \right].$$

Let $u_1^T, u_2^T, \dots, u_k^T$ be the rows of \mathbf{U}_l . Then, by positive homogeneity of ReLU σ , we have

$$\left\| \sum_{i=1}^n \epsilon_i \sigma(\mathbf{U}_l h_{l-1}(x_i)) \right\|^2 = \sum_{j=1}^k \|u_j\|^2 \left(\sum_{i=1}^n \epsilon_i \sigma \left(\frac{u_j^T h_{l-1}(x_i)}{\|u_j\|} \right) \right)^2.$$

The supremum of this quantity over u_1, \dots, u_k under the constraint that $\|\mathbf{U}_l\|_F^2 = \sum_{j=1}^k \|u_j\|^2 \leq M_l^2$ is attained when $\|u_j\| = M_l$ for some j and $\|u_{j'}\| = 0$ for all other $j' \neq j$. This means that

$$\begin{aligned}
\frac{\mathcal{A}}{M_l} &= \mathbb{E} \left[\sup_{\substack{\|\mathbf{U}_l\|_F \leq M_l, \\ h_{l-1} \in \mathcal{H}_{l-1}}} \left\| \sum_{i=1}^n \epsilon_i \sigma(\mathbf{U}_l h_{l-1}(x_i)) \right\| \right] = \mathbb{E} \left[\sup_{\substack{\|u\| \leq M_l, \\ h_{l-1} \in \mathcal{H}_{l-1}}} \left| \sum_{i=1}^n \epsilon_i \sigma(u^T h_{l-1}(x_i)) \right| \right] \\
&= \mathbb{E} \left[\sup_{\substack{\|u\| \leq M_l, \\ h_{l-1} \in \mathcal{H}_{l-1}}} \left[\sum_{i=1}^n \epsilon_i \sigma(u^T h_{l-1}(x_i)) \right]_+ + \left[\sum_{i=1}^n \epsilon_i \sigma(u^T h_{l-1}(x_i)) \right]_- \right] \\
&\leq \mathbb{E} \left[\sup_{\substack{\|u\| \leq M_l, \\ h_{l-1} \in \mathcal{H}_{l-1}}} \left[\sum_{i=1}^n \epsilon_i \sigma(u^T h_{l-1}(x_i)) \right]_+ \right] + \mathbb{E} \left[\sup_{\substack{\|u\| \leq M_l, \\ h_{l-1} \in \mathcal{H}_{l-1}}} \left[\sum_{i=1}^n \epsilon_i \sigma(u^T h_{l-1}(x_i)) \right]_- \right] \\
&\stackrel{(a)}{=} 2\mathbb{E} \left[\sup_{\substack{\|u\| \leq M_l, \\ h_{l-1} \in \mathcal{H}_{l-1}}} \left[\sum_{i=1}^n \epsilon_i \sigma(u^T h_{l-1}(x_i)) \right]_+ \right] \stackrel{(b)}{=} 2\mathbb{E} \left[\left[\sup_{\substack{\|u\| \leq M_l, \\ h_{l-1} \in \mathcal{H}_{l-1}}} \sum_{i=1}^n \epsilon_i \sigma(u^T h_{l-1}(x_i)) \right]_+ \right] \\
&\stackrel{(c)}{=} 2\mathbb{E} \left[\sup_{\substack{\|u\| \leq M_l, \\ h_{l-1} \in \mathcal{H}_{l-1}}} \sum_{i=1}^n \epsilon_i \sigma(u^T h_{l-1}(x_i)) \right] \stackrel{(d)}{\leq} 2\mathbb{E} \left[\sup_{\substack{\|u\| \leq M_l, \\ h_{l-1} \in \mathcal{H}_{l-1}}} \sum_{i=1}^n \epsilon_i u^T h_{l-1}(x_i) \right],
\end{aligned}$$

where equality (a) is due to symmetry of Rademacher random variables and (b) uses $\sup [t]_+ = [\sup t]_+$. Equality (c) uses the fact that the supremum is nonnegative, because setting $u = 0$

already gives $\sum_{i=1}^n \epsilon_i \sigma(u^T h_{l-1}(x_i)) = 0$. Inequality (d) uses contraction property of Rademacher complexity.

Lastly, one can notice that

$$\mathbb{E} \left[\sup_{\substack{\|u\| \leq M_l, \\ h_{l-1} \in \mathcal{H}_{l-1}}} \sum_{i=1}^n \epsilon_i u^T h_{l-1}(x_i) \right] = M_l \mathbb{E} \left[\sup_{\substack{\|w\| \leq 1, \\ h_{l-1} \in \mathcal{H}_{l-1}}} \sum_{i=1}^n \epsilon_i w^T h_{l-1}(x_i) \right] = M_l n \widehat{\mathcal{R}}_n(\mathcal{F}_{l-1}|_S).$$

This establishes

$$\mathcal{A} \leq 2M_l^2 n \widehat{\mathcal{R}}_n(\mathcal{F}_{l-1}|_S),$$

which leads to the conclusion that

$$\widehat{\mathcal{R}}_n(\mathcal{F}_l|_S) \leq (1 + 2M_l^2) \widehat{\mathcal{R}}_n(\mathcal{F}_{l-1}|_S),$$

as desired.