

## Compression

$$H = \left\{ \mathbb{R}^d \ni x \mapsto y = W_N \sigma(\dots W_2 \sigma(W_1 x) \dots) \in \mathbb{R}^d : W_n \in \mathbb{R}^{d \times d}, n \in [N] \right\}$$

Assume  $\sigma$  is  $\delta$ -Lipschitz and  $\sigma(0) = 0$ ,  $X = \{x \in \mathbb{R}^d, \|x\| \leq 1\}$

For  $r \in [d]$  Let  $H_r$  to be the hypotheses space corresponding to the same network above when its weight matrices are constrained to have rank  $r$  or less, i.e.:

$$H_r = \left\{ x \mapsto y = W_N V_N^T \sigma(U_{N-1} V_{N-1}^T \dots \sigma(U_1 V_1^T x) \dots) : U_n, V_n \in \mathbb{R}^{d \times r}, n \in [N] \right\}$$

Assume that each of the  $2Ndr$  Parameters representing  $H_r$  is stored in memory using  $b$  bits.

Given a loss  $\ell: \mathbb{R}^d \times \mathbb{R}^d \rightarrow [0, 1]$  s.t.  $\forall y, \hat{y}, y' \quad |\ell(y, \hat{y}) - \ell(y, y')| \leq \rho / \|y' - y\|$

we would like to derive generalization bounds for  $H$ .

(a) Fix  $r \in [d]$  and derive a generalization bound for  $H$  by compressing it into  $H_r$ .

(b) Derive a generalization bound for  $H$  by simultaneously compressing it into  $H_r$  for all  $r \in [d]$

## Solution (a)

Let  $r \in [d]$ .

By assumption each of the  $2Ndr$  parameters representing  $H_r$  using

$$b \text{ bits, thus } |H_r| \leq \underbrace{2^b}_{\substack{\# \text{ options} \\ \text{for the} \\ \text{first param}}} \dots \underbrace{2^b}_{\substack{\# \text{ options} \\ \text{for the} \\ 2Ndr \text{ param}}} = 2^{b \cdot 2Ndr}$$

In addition,  $\ell(\cdot, \cdot)$  is  $\delta$ -Lipschitz, therefore we can use the theorem we proved in class in order to get:

For  $\hat{h}$  the returned hypothesis from the algorithm  $\forall \delta \in (0, 1)$  w.p.  $\geq 1 - \delta$  over  $S \sim \mathcal{D}^m$ :

$$L_D(\hat{h}) - L_S(\hat{h}) \leq \sqrt{\frac{(b \cdot 2Ndr + 1) \ln(2) + \ln(\frac{1}{\delta})}{2m}} + 2 \cdot \rho \cdot d(\hat{h}, H_r)$$

We would like to bound the term  $d(\hat{h}, H_r)$ .

Note that we proved in class the following (under the same assumptions except the rank of the approximation matrices, in class it was  $\text{rank}=1$  and now  $\text{rank}=r$ ):

$$d(\hat{h}, h_r) \leq \gamma^{N-1} \cdot \sum_{n=1}^N \prod_{j \in [N] \setminus n} \|W_j\|_{\text{spectral}} \cdot \|W_n - W'_n\|_{\text{spectral}}$$

For  $\hat{h}$  with weighted matrices  $W_N, \dots, W_1$

and  $h_r$  with weighted matrices  $W'_N, \dots, W'_1$  with  $\text{rank} \leq r$

the closest approximations to  $W_N, \dots, W_1$ .

In the proof in class the only place we used that the approximations were to 1-rank matrix was when explaining  $\|W'_n\|_{\text{spectral}} = \|W_n\|_{\text{spectral}} \forall n$

Note, that this equation holds for approximations to rank  $r \geq 1$  too,

since if we denote  $W_n$  as spectral decomposition  $W_n = \sum_{i=1}^d \sigma_i u_i v_i^T$

where  $|\sigma_1| \geq |\sigma_2| \geq \dots \geq |\sigma_d|$  then the best rank- $r$  approximation

to  $W_n$  is  $W'_n = \sum_{i=1}^r \sigma_i u_i v_i^T$ , thus we can see clearly

that  $\|W'_n\|_{\text{spectral}} = |\sigma_1| = \|W_n\|_{\text{spectral}}$  and thus the

bound denote above holds for each  $n \in [d]$ .

In addition  $\|W_n - W'_n\|_{\text{spectral}} = \left\| \sum_{i=1}^d \sigma_i u_i v_i^T - \sum_{i=1}^r \sigma_i u_i v_i^T \right\|_F = \left\| \sum_{i=r+1}^d \sigma_i u_i v_i^T \right\|_F = |\sigma_{r+1}|$

$$= |\sigma_{r+1}(W_n)|$$

$\Rightarrow$  In total we'll get  $\forall \epsilon \in (0, 1)$  w.p.  $\geq 1 - \delta$  over  $S \sim D^m$ :

$$L_D(\hat{h}) - L_S(\hat{h}) \leq \sqrt{\frac{(b(2M/r+1)(n(2) + (n(\frac{1}{\delta})))}{2m}} + 2 \cdot \gamma \cdot \left[ \gamma^{N-1} \sum_{n=1}^N \prod_{j \in [N] \setminus n} \|W_j\|_{\text{spec}} \cdot |\sigma_{r+1}(W_n)| \right]$$



### Solution (b):

We'll use the bound from (a) with union bound in order to get bound for H.

From (a) we've got that  $\forall \delta \in (0,1)$  w.p  $\geq 1 - \frac{\delta}{d}$  over  $S \sim D^m$  and fixed  $1 \leq r \leq d$ :

$$\begin{aligned} L_0(\hat{h}) - L_S(\hat{h}) &\leq \sqrt{\frac{(b \cdot 2Nd + 1) \ln(2) + \ln\left(\frac{1}{\delta}\right)}{2(m-1)}} + 2\beta \cdot \left[ \gamma^{N-1} \sum_{n=1}^N \prod_{j \neq n} \|W_j\|_{\text{FP}} \cdot |\sigma_{r+1}(W_n)| \right] \\ &= \sqrt{\frac{(b \cdot 2Nd + 1) \ln(2) + \ln\left(\frac{1}{\delta}\right)}{2(m-1)}} + 2\beta \cdot \left[ \gamma^{N-1} \sum_{n=1}^N \prod_{j \neq n} \|W_j\|_{\text{FP}} \cdot |\sigma_{r+1}(W_n)| \right] \end{aligned}$$

We will sign  $A_r = \sqrt{\frac{(b \cdot 2Nd + 1) \ln(2) + \ln\left(\frac{1}{\delta}\right)}{2(m-1)}}$ ,  $B_r = 2\beta \cdot \left[ \gamma^{N-1} \sum_{n=1}^N \prod_{j \neq n} \|W_j\|_{\text{FP}} \cdot |\sigma_{r+1}(W_n)| \right]$

Thus  $\Pr(L_0(\hat{h}) - L_S(\hat{h}) \geq A_r + B_r) \leq \frac{\delta}{d} \quad \forall 1 \leq r \leq d$

$\Rightarrow \Pr(\exists r (L_0(\hat{h}) - L_S(\hat{h}) \geq A_r + B_r)) \leq \sum_{r=1}^d \Pr(L_0(\hat{h}) - L_S(\hat{h}) \geq A_r + B_r) \leq \sum_{r=1}^d \frac{\delta}{d} = \delta$   
union bound

$\Rightarrow$  w.p  $\geq 1 - \delta$  over  $S \sim D^m$ :

$$L_0(\hat{h}) - L_S(\hat{h}) \leq \min_{1 \leq r \leq d} (A_r + B_r) = \min_{1 \leq r \leq d} \left\{ \sqrt{\frac{(b \cdot 2Nd + 1) \ln(2) + \ln\left(\frac{1}{\delta}\right)}{2(m-1)}} + 2\beta \left[ \gamma^{N-1} \sum_{n=1}^N \prod_{j \neq n} \|W_j\|_{\text{FP}} \cdot |\sigma_{r+1}(W_n)| \right] \right\}$$

Note that as  $r$  gets bigger:  $A_r$  gets bigger and  $B_r$  gets lower (better approximation to the matrices but another case of more weighted parameters)

