

Ch 20 Neural Networks: sample complexity

Friday, March 27, 2020

3:05 PM

i.e. II Estimation / generalization error

In direct competition w, approximation error, now we want d_f to be small (low-complexity)

We'll analyze by bounding VC dimension.

We learn $|E|$ parameters (weights), so via discretization trick, we can think of

$\mathcal{H}_{V,E,\sigma}$ as having finite size $|E|$ -constant, so expect $\text{VCdim}(\mathcal{H}_{V,E,\sigma}) \approx O(|E|)$
 \uparrow but large, like 64

Case 1: $\sigma = \text{sign}$

Thm 20.6 $\text{VCdim}(\mathcal{H}_{V,E,\text{sign}}) = O(|E| \cdot \log(|E|))$

Proof Recall growth fun $T_{\mathcal{H}}(m) = \max_{\substack{C \subseteq X \\ |C|=m}} |\mathcal{H}|_C| \leq 2^m$
 i.e. $\#$ dichotomies.

Can extend to any finite Y , not just $Y = \{0,1\}$, if want to

Exercise 20.4: showed that if $\mathcal{H} = \{f \circ g, f \in \mathcal{H}_1, g \in \mathcal{H}_2\}$ then

$$T_{\mathcal{H}}(m) \leq T_{\mathcal{H}_1}(m) \cdot T_{\mathcal{H}_2}(m)$$

Using the ANN's layered structure, $\mathcal{H} = \mathcal{H}^{(1)} \circ \dots \circ \mathcal{H}^{(z)} \circ \mathcal{H}^{(n)}$

and each layer maps $\mathbb{R}^{|V_{t-1}|} \rightarrow \{\pm 1\}^{|V_t|}$
 \leftarrow since $\text{sign}(\dots)$

Also, each layer $\mathcal{H}^{(t)}$ is a direct product of its individual neurons (all independent)

$$\mathcal{H}^{(t)} = \mathcal{H}^{(t,1)} \times \mathcal{H}^{(t,2)} \times \dots \times \mathcal{H}^{(t,|V_t|)}$$

and using exercise 20.3,

$$T_{\mathcal{H}^{(t)}}(m) \leq \prod_{i \in [V_t]} T_{\mathcal{H}^{(t,i)}}(m)$$

\hookrightarrow halfspace classifier, so $\text{VCdim}(\mathcal{H}^{(t,i)}) = \text{dimension of input}$
 $= d_{t,i}$ ($\#$ edges arriving to i^{th} neuron in layer t)

and Sauer's Lemma:

$$T_{\mathcal{H}^{(t,i)}}(m) \leq \left(\frac{e \cdot m}{d_{t,i}} \right)^{d_{t,i}} \leq (e \cdot m)^{d_{t,i}}$$

So, combining everything,

$$T_{\mathcal{H}}(m) \leq \prod_t \prod_{i \in [V_t]} (e \cdot m)^{d_{t,i}} = (e \cdot m)^{\sum_t \sum_i d_{t,i}} = (e \cdot m)^{|E|}$$

$\#$ edges

now "reverse Sauer's lemma" vschg Lemma A.2 :

$$\text{If } \text{VCdim}(\mathcal{H}) = m \Rightarrow \mathcal{I}_{\mathcal{H}}(m) \geq 2^m, \text{ i.e., } 2^m \leq (em)^{|E|} \\ \text{i.e. } m \leq |E| \cdot \log_2(em)$$

$$\text{So Lemma A.2} \Rightarrow \underline{m \leq 4 \cdot |E| \cdot \log_2(2|E|) + \text{const.}} \Rightarrow \text{VCdim}(\mathcal{H}_{V,E,\text{sign}}) \leq O(|E| \cdot \log(|E|)) \quad \square$$

Case 2: $\sigma = \text{sigmoid}$

$$\text{Exercise 20.5: } \text{VCdim}(\mathcal{H}_{V,E,\sigma}) = \Omega(|E|^2) \\ \uparrow \text{i.e. at least. BAD!}$$

$$\text{and (not proven) } \text{VCdim}(\mathcal{H}_{V,E,\sigma}) \leq O(|E|^2 \cdot |V|^2) \text{ also bad}$$

$$\text{though in practice } \text{VCdim}(\mathcal{H}_{V,E,\sigma}) \approx 64 \cdot O(|E|) \\ \text{via discretization trick w, 64 bits}$$

Case 3: $\sigma = \text{ReLU}$ (not in our book)

Nick Harvey, Chris Liaw, and Abbas Mehrabian. Nearly-tight VC-dimension bounds for piece-wise linear neural networks. COLT 2017

let $\mathcal{H}_{V,E,\sigma}$ have $|E|$ edges (weights), T layers

$$\underline{\text{Thm (lower bound)}} \quad \text{VCdim}(\mathcal{H}_{V,E,\sigma}) \geq \Omega(|E| \cdot T \cdot \log(\frac{|E|}{T}))$$

$$\underline{\text{Thm (upper bound)}} \quad \text{VCdim}(\mathcal{H}_{V,E,\sigma}) \leq O(|E| \cdot T \cdot \log(|E|))$$

Since T is mild (even for a "deep network", $T = O(10)$),

this is saying $\text{VCdim} \approx |E| \cdot \log(|E|)$, same result we derived for $\sigma = \text{sign}$ case