

# Lecture 14: VC Dimension, Validation, Neural Networks, Kernels (workshop)

Nisha Chandramoorthy

October 5, 2023

# VC dimension

- ▶ A hypothesis class restricted to some  $C$  *shatters*  $C$  if any binary function on  $C$  is in the class

# VC dimension

- ▶ A hypothesis class restricted to some  $C$  *shatters*  $C$  if any binary function on  $C$  is in the class
- ▶ VC dimension,  $\text{VCdim}(\mathcal{H})$ : maximal size of  $C \subseteq \mathcal{X}$  that is shattered by  $\mathcal{H}$ .

# VC dimension

- ▶ A hypothesis class restricted to some  $C$  *shatters*  $C$  if any binary function on  $C$  is in the class
- ▶ VC dimension,  $\text{VCdim}(\mathcal{H})$ : maximal size of  $C \subseteq \mathcal{X}$  that is shattered by  $\mathcal{H}$ .
- ▶ Eg 1.  $\text{VCdim}$  of threshold functions on  $\mathbb{R}$  is 1.

# VC dimension

- ▶ A hypothesis class restricted to some  $C$  *shatters*  $C$  if any binary function on  $C$  is in the class
- ▶ VC dimension,  $\text{VCdim}(\mathcal{H})$ : maximal size of  $C \subseteq \mathcal{X}$  that is shattered by  $\mathcal{H}$ .
- ▶ Eg 1.  $\text{VCdim}$  of threshold functions on  $\mathbb{R}$  is 1.
- ▶ Eg 2:  $\text{VCdim}$  of indicator functions on intervals of  $\mathbb{R}$  is 2.

# VC dimension

- ▶ A hypothesis class restricted to some  $C$  *shatters*  $C$  if any binary function on  $C$  is in the class
- ▶ VC dimension,  $\text{VCdim}(\mathcal{H})$ : maximal size of  $C \subseteq \mathcal{X}$  that is shattered by  $\mathcal{H}$ .
- ▶ Eg 1.  $\text{VCdim}$  of threshold functions on  $\mathbb{R}$  is 1.
- ▶ Eg 2:  $\text{VCdim}$  of indicator functions on intervals of  $\mathbb{R}$  is 2.
- ▶ Eg 3:  $\text{VCdim}$  of a finite class  $\mathcal{H} \leq \log_2 |\mathcal{H}|$

# Generalization bounds based on VC dimension

►  $\mathcal{H} = \{h_\theta(x) = \sin(\theta x) : \theta \in \mathbb{R}\}.$

# Generalization bounds based on VC dimension

- ▶  $\mathcal{H} = \{h_\theta(x) = \sin(\theta x) : \theta \in \mathbb{R}\}.$
- ▶ VCdim is  $\infty$



# Generalization bounds based on VC dimension

- ▶  $\mathcal{H} = \{h_\theta(x) = \sin(\theta x) : \theta \in \mathbb{R}\}$ .
- ▶ VCdim is  $\infty$
- ▶ Binary classification generalization for 0-1 loss over class  $\mathcal{H}$  with VCdim =  $d$ : there exist constants  $C_1, C_2 > 0$  such that

$$C_1 \frac{d + \log(1/\delta)}{\epsilon^2} \leq m_{\mathcal{H}}(\epsilon, \delta) \leq C_2 \frac{d + \log(1/\delta)}{\epsilon}$$

- ▶ Consistent with bound from lecture 2.

- ▶ Consistent with bound from lecture 2.
- ▶  $m_{\mathcal{H}} = \frac{\log(|\mathcal{H}|/\delta)}{\epsilon}$
- ▶  $\mathcal{D}^m(S : R(h_S) \geq \epsilon) \leq |\mathcal{H}_b| e^{-\epsilon m} \leq |\mathcal{H}| e^{-\epsilon m}.$

- ▶ Consistent with bound from lecture 2.
- ▶  $m_{\mathcal{H}} = \frac{\log(|\mathcal{H}|/\delta)}{\epsilon}$
- ▶  $\mathcal{D}^m(S : R(h_S) \geq \epsilon) \leq |\mathcal{H}_b| e^{-\epsilon m} \leq |\mathcal{H}| e^{-\epsilon m}.$

# VC dimension contd

# Cross Validation

- ▶ Recall Hoeffding's inequality:  $X_1, \dots, X_m$  iid sequence with  $P(a \leq X \leq b) = 1$ . Then, with probability  $\geq 1 - \delta$ ,

$$\left| \frac{1}{m} \sum_{i=1}^m X_i - EX \right| < (b - a) \sqrt{\frac{\log(2/\delta)}{2m}}$$

# Cross Validation

- ▶ Recall Hoeffding's inequality:  $X_1, \dots, X_m$  iid sequence with  $P(a \leq X \leq b) = 1$ . Then, with probability  $\geq 1 - \delta$ ,

$$\left| \frac{1}{m} \sum_{i=1}^m X_i - EX \right| < (b - a) \sqrt{\frac{\log(2/\delta)}{2m}}$$

- ▶ For any set  $V$  of size  $m$ , when  $\text{loss} \in (0, 1)$ , by Hoeffding's inequality,

$$|R_V(h) - R(h)| \leq \sqrt{\frac{\log(2/\delta)}{2m}}$$

- ▶  $k$ -fold cross validation to choose models (e.g., regularization parameters):
  - ▶ Divide given set  $S$  into  $k$  subsets (folds)
  - ▶ For each parameter, each fold: run learning algorithm on union of all folds except one; calculate test/validation loss on fold
  - ▶ Run algorithm on  $S$  using parameter with minimum total test/validation loss.



# Models/representations, algorithms, statistical principles

- ▶ how to make and test conjectures about how large language models (LLMs) learn?

# Models/representations, algorithms, statistical principles

- ▶ how to make and test conjectures about how large language models (LLMs) learn?
- ▶ “how to train them better (more efficiently)” – number of practical questions perhaps benefit from theory

# Models/representations, algorithms, statistical principles

- ▶ how to make and test conjectures about how large language models (LLMs) learn?
- ▶ “how to train them better (more efficiently)” – number of practical questions perhaps benefit from theory
- ▶ AI safety, fair and ethical ethical use, combining with other domain knowledge (e.g., physics, chemistry etc).... and many more!

# Models/representations, algorithms, statistical principles

- ▶ how to make and test conjectures about how large language models (LLMs) learn?
- ▶ “how to train them better (more efficiently)” – number of practical questions perhaps benefit from theory
- ▶ AI safety, fair and ethical use, combining with other domain knowledge (e.g., physics, chemistry etc).... and many more!
- ▶ Perhaps biggest contribution advance to LLMs: transformers and their training.

# (partial) History - trace back from transformers (source:Wikipedia)

- ▶ Transformer architecture: 2017, Google Brain [Vaswani et al]
- ▶ Deep learning, unsupervised learning 2010s (e.g., GANs 2014)...
- ▶ ImageNet: 2009, Fei Fei Li
- ▶ Long-short term memory (LSTM) architecture: 1997, [Hochreiter and Schmidhuber]
- ▶ Convolutional NNs: (inspired from) 1979 work by [Fukushima]; Recurrent neural networks: 1982 [Hopfield]
- ▶ ...
- ▶ Automatic Differentiation: 1970 [Linnainmaa]
- ▶ ...
- ▶ First neural networks: 1950s [Minsky and others]

# Fully connected Neural Networks

- ▶ Neuron: input  $\sum_j w_j h_j$ ; output  $\sigma(\sum_j w_j h_j)$
- ▶ Organized into layers of depth  $l$  and width  $n$
- ▶ Graph:  $V, E, \sigma, w$ ; weight function.

# Fully connected Neural Networks

- ▶ Neuron: input  $\sum_j w_j h_j$ ; output  $\sigma(\sum_j w_j h_j)$
- ▶ Organized into layers of depth  $l$  and width  $n$
- ▶ Graph:  $V, E, \sigma, w$ ; weight function.
- ▶ VC dim of  $\mathcal{H}_{V,E,\text{sign}} \leq C|E| \log |E|$

# Fully connected Neural Networks

- ▶ Neuron: input  $\sum_j w_j h_j$ ; output  $\sigma(\sum_j w_j h_j)$
- ▶ Organized into layers of depth  $l$  and width  $n$
- ▶ Graph:  $V, E, \sigma, w$ ; weight function.
- ▶ VC dim of  $\mathcal{H}_{V,E,\text{sign}} \leq C|E| \log |E|$ 
  - ▶ Proof: Growth function  $\tau_{\mathcal{H}}(m) = \max_{C, |C|=m} |\mathcal{H}|_C|$



# Fully connected Neural Networks

- ▶ Neuron: input  $\sum_j w_j h_j$ ; output  $\sigma(\sum_j w_j h_j)$
- ▶ Organized into layers of depth  $l$  and width  $n$
- ▶ Graph:  $V, E, \sigma, w$ ; weight function.
- ▶ VC dim of  $\mathcal{H}_{V,E,\text{sign}} \leq C|E| \log |E|$ 
  - ▶ Proof: Growth function  $\tau_{\mathcal{H}}(m) = \max_{C, |C|=m} |\mathcal{H}|_C|$  Fact 1: Let  $\mathcal{H} = \mathcal{H}_l \circ \dots \circ \mathcal{H}_1$ . Then,  $\tau_{\mathcal{H}}(m) \leq \prod_{t=1}^l \tau_{\mathcal{H}_t}(m)$ .

# Fully connected Neural Networks

- ▶ Neuron: input  $\sum_j w_j h_j$ ; output  $\sigma(\sum_j w_j h_j)$
- ▶ Organized into layers of depth  $l$  and width  $n$
- ▶ Graph:  $V, E, \sigma, w$ ; weight function.
- ▶ VC dim of  $\mathcal{H}_{V,E,\text{sign}} \leq C|E| \log |E|$ 
  - ▶ Proof: Growth function  $\tau_{\mathcal{H}}(m) = \max_{C, |C|=m} |\mathcal{H}|_C|$  Fact 1: Let  $\mathcal{H} = \mathcal{H}_l \circ \dots \circ \mathcal{H}_1$ . Then,  $\tau_{\mathcal{H}}(m) \leq \prod_{t=1}^l \tau_{\mathcal{H}_t}(m)$ .
  - ▶ Fact 2: Let  $\mathcal{H} = \mathcal{H}^{(1)} \dots \circ \mathcal{H}^{(n)}$ . Then,  $\tau_{\mathcal{H}}(m) \leq \prod_{t=1}^l \tau_{\mathcal{H}^{(t)}}(m)$ .

# Fully connected Neural Networks

- ▶ Neuron: input  $\sum_j w_j h_j$ ; output  $\sigma(\sum_j w_j h_j)$
- ▶ Organized into layers of depth  $l$  and width  $n$
- ▶ Graph:  $V, E, \sigma, w$ ; weight function.
- ▶ VC dim of  $\mathcal{H}_{V,E,\text{sign}} \leq C|E| \log |E|$ 
  - ▶ Proof: Growth function  $\tau_{\mathcal{H}}(m) = \max_{C, |C|=m} |\mathcal{H}|_C|$  Fact 1: Let  $\mathcal{H} = \mathcal{H}_l \circ \dots \circ \mathcal{H}_1$ . Then,  $\tau_{\mathcal{H}}(m) \leq \prod_{t=1}^l \tau_{\mathcal{H}_t}(m)$ .
  - ▶ Fact 2: Let  $\mathcal{H} = \mathcal{H}^{(1)} \dots \circ \mathcal{H}^{(n)}$ . Then,  $\tau_{\mathcal{H}}(m) \leq \prod_{t=1}^l \tau_{\mathcal{H}^{(t)}}(m)$ .
  - ▶ Fact 3: Sauer's Lemma:  $\tau_{\mathcal{H}}(m) = (em/d)^d$ , where  $d \geq \text{VCdim}(\mathcal{H})$

Restriction:

$$h \in \mathcal{H}$$

$C$

$$h|_C(x_i) = h(x_i) \quad x_i \in C$$

$$\mathcal{H}|_C = \{h|_C : h \in \mathcal{H}\}$$

Restriction of  $\mathcal{H}$  to a set  $C$ .

$$\rightarrow C = \{x_1, x_2, x_3\}$$

$h|_C$  should take values

$$\mathcal{H} \quad \left\{ \begin{bmatrix} 0 \\ 0 \\ 0 \end{bmatrix}, \begin{bmatrix} 0 \\ 0 \\ 1 \end{bmatrix}, \dots, \begin{bmatrix} 1 \\ 1 \\ 1 \end{bmatrix} \right\}$$

$|\mathcal{H}|_C$  should have at least 8 elements if  $\mathcal{H}$  can shatter  $C$ .

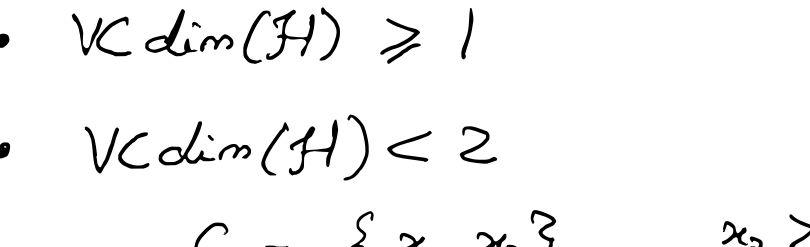
$$\rightarrow \text{If } \text{VCdim}(\mathcal{H}) = m$$

- $\exists$  some set  $C \subseteq \mathcal{X}$  that is shattered by  $\mathcal{H}$  and  $|C| = m$
- $\nexists$  a set  $C \subseteq \mathcal{X}$  of size  $|C| = m+1$  that is shattered by  $\mathcal{H}$ .

$$\rightarrow \text{VCdim}(\mathcal{H}) = 1$$

$$\mathcal{H} = \{x \rightarrow h_\theta(x) : \theta \in \mathbb{R}\}$$

$$h_\theta(x) = 1 \text{ if } x < \theta \text{ and } h_\theta(x) = 0 \text{ otherwise}$$



$$\bullet \text{VCdim}(\mathcal{H}) \geq 1$$

$$\bullet \text{VCdim}(\mathcal{H}) < 2$$

$$C = \{x_1, x_2\} \quad x_2 > x_1 \text{ wlog}$$

$$h|_C(x_1) = 0 \text{ and } h|_C(x_2) = 1$$

Not possible if  $x_2 > x_1$ .

$$\rightarrow \text{VCdim}(\mathcal{H}) = 2$$

$$\mathcal{H} = \{h_{\theta_1, \theta_2}(x) = \begin{cases} 1 & \text{if } \theta_1 < x < \theta_2 \\ 0 & \text{otherwise} \end{cases} : \theta_1, \theta_2 \in \mathbb{R}\}$$



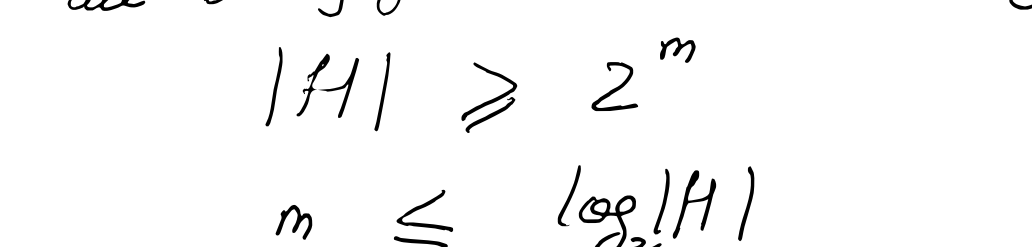
$$0, 1, 0$$

$$1, 0, 1 \rightarrow \text{not possible}$$

x

.

x



$$\rightarrow |\mathcal{H}| < \infty$$

$$\text{VCdim}(\mathcal{H}) < \log_2 |\mathcal{H}|$$

Proof

$$\text{Say } \text{VCdim}(\mathcal{H}) = m.$$

$$\rightarrow \exists C \text{ with } |C| = m \text{ s.t. all binary functions on } C \text{ are in } \mathcal{H}|_C$$

$$|\mathcal{H}| \geq 2^m$$

$$m \leq \log_2 |\mathcal{H}|$$

$$\rightarrow \text{VCdim}(\mathcal{H}) = d$$

(parameter dim)

$$\bullet \mathcal{H} = \{x \rightarrow \langle w, x \rangle : w \in \mathbb{R}^d, x \in \mathbb{R}^d\}$$

$$\bullet \text{Suppose } \text{VCdim}(\mathcal{H}) = d$$

$$\exists x_1, x_2, \dots, x_d \text{ is shattered by } \mathcal{H}.$$

and

$$\nexists x_1, x_2, \dots, x_{d+1} \text{ that are shattered by } \mathcal{H}.$$

$$\bullet \text{Assume that } \exists x_1, x_2, \dots, x_{d+1} \text{ that is shattered by } \mathcal{H} \rightarrow \text{show this is not possible.}$$

$$\bullet \begin{bmatrix} x_{11} & \dots & x_{1d} \\ \vdots & & \vdots \\ x_{d+1,1} & \dots & x_{d+1,d} \end{bmatrix} \begin{bmatrix} w_1 \\ w_2 \\ \vdots \\ w_d \end{bmatrix} = \begin{bmatrix} \vdots \\ \vdots \\ \vdots \end{bmatrix} \in \mathbb{R}^d$$

$h|_C$

if  $x_1, \dots, x_d$  are linearly independent,

$$\text{yes!} \Rightarrow \text{VCdim}(\mathcal{H}) \geq d.$$

$$\bullet \begin{bmatrix} x_{11} & \dots & x_{1d} \\ \vdots & & \vdots \\ x_{d+1,1} & \dots & x_{d+1,d} \end{bmatrix} \begin{bmatrix} w_1 \\ \vdots \\ w_d \end{bmatrix} = \begin{bmatrix} \vdots \\ \vdots \\ \vdots \end{bmatrix} \in \mathbb{R}^{d+1}$$

$$x_{d+1} = \sum_{i=1}^d a_i x_i \quad (\because x_i \in \mathbb{R}^d)$$

$$x_{d+1} = a_d x_d$$

$$\langle x_{d+1}, w \rangle = a_d \langle x_d, w \rangle$$

$$\text{sgn}(a_d) < 0$$

$$\bullet x_{d+1} = \sum_{i=1}^d a_i x_i \quad \exists \text{ some } a_i \in \mathbb{R}^{d+1}$$

$$\sum_{i=1}^{d+1} a_i x_i = 0 \rightarrow \textcircled{A}$$

$$P = \{i : a_i > 0\} \quad N = \{i : a_i < 0\}$$

$$\sum_{i \in P} a_i x_i = \sum_{i \in N} |a_i| x_i \quad \textcircled{A}$$

$$\text{if } x_1, x_2, \dots, x_{d+1} \text{ is shattered by } \mathcal{H}, \exists w \text{ s.t.}$$

$$\text{for any } i \in P, \langle x_i, w \rangle > 0$$

$$\text{and } i \in N, \langle x_i, w \rangle < 0 \checkmark$$

$$\rightarrow$$

$$\mathcal{H} = \{x \rightarrow \sin \theta x : \theta \in \mathbb{R}\}$$



with  $P_{\mathcal{S}}$  (over  $S \sim D^m$ )  $> 1 - \delta$ ,

$$R(h) \leq R_S(h) + \frac{C(\mathcal{H})}{m} + \sqrt{\frac{\log(1/\delta)}{2m}}$$

Sample Complexity:  $m_{\mathcal{H}}(\epsilon, \delta)$   
 as the minimal number of  
 samples s.t.  $|R(h) - R_S(h)| < \epsilon$   
 with  $P_{\mathcal{S}} = 1 - \delta$ .

$$m_H(\varepsilon, \delta) = O\left(\frac{\log|H| + \log(1/\delta)}{\varepsilon}\right)$$

Recap of Proof from Lec 2

$$\rightarrow H_\varepsilon = \{h \in H : R(h) > \varepsilon\}$$

$\rightarrow$  Realizability assumption

$$\min_{h \in H} R_S(h) = 0$$

$h_S$  is ERM for  $H, S$ .

$$\Pr(R_S(h_S) = 0 \text{ and } R(h_S) > \varepsilon)$$

$$\leq (1 - \varepsilon)^m \leq e^{-\varepsilon m}$$

$\rightarrow$  Next time:

$$VC \dim(\text{FCNN with } |E|) \leq C|E| \log |E|$$