

Announcements

- Homework 5 is due Tuesday.
-

Today's Lecture

- Interpreting the VC generalization bound
- Implications in dataset methodology

VC GENERALIZATION BOUND AND ITS INTERPRETATION

Plug in: $m_H(2N) \leq (2N)^{d_{vc}} + 1$
to $\epsilon_{eff} \Rightarrow$

$$1. \quad \epsilon_{eff} \leq \sqrt{\frac{8}{N} \ln \frac{4[(2N)^{d_{vc}} + 1]}{\delta}} = \epsilon_{vc}$$

For: $(2N)^{d_{vc}} \gg 1$

$$\begin{aligned} \epsilon_{vc} &\approx \sqrt{\frac{8}{N} [\ln(4(2N)^{d_{vc}}) - \ln \delta]} \\ &= \sqrt{\frac{8 \ln 4}{N} + \frac{8 d_{vc}}{N} \ln(2N) - \frac{8}{N} \ln \delta} \end{aligned}$$

$$\lim_{N \rightarrow \infty} \epsilon_{vc} = 0.$$

\Rightarrow THIS PROVES THAT LEARNING IS FEASIBLE, EVEN WITH INFINITE HYPOTHESIS SETS, IF d_{vc} IS FINITE.

2. USING A TEST SET TO BOUND $E_{out}(h_g)$:
 IF D_{Test} HAS NOT BEEN USED TO PICK h_g OR H , THEN:

$$E_{out}(h_g) \leq E_{Test}(h_g) + \sqrt{\frac{1}{2N} \ln \frac{2M}{\delta}}$$

WITH PROBABILITY $1 - \delta$.

→ USE $M = 1$.

3. TH'N 2.5 AGAIN! (using E_{vc})

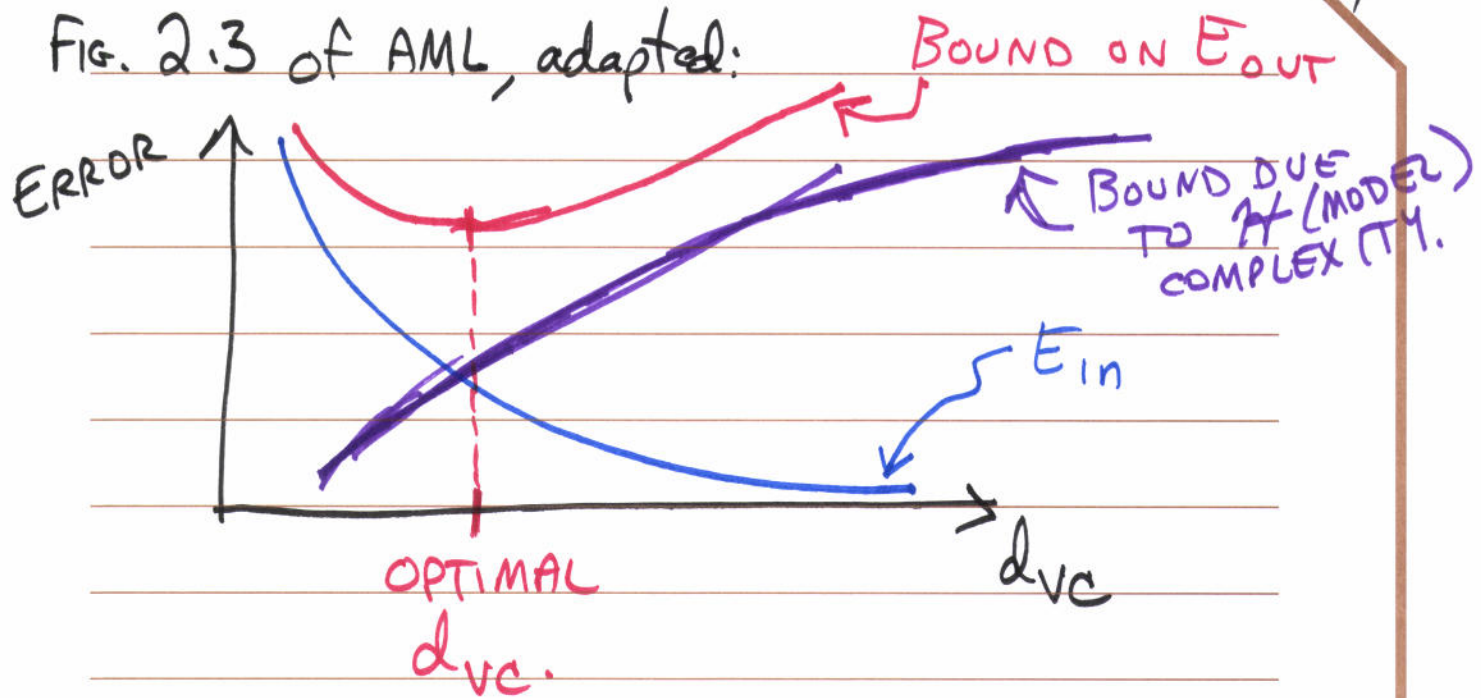
$$E_{out}(h_g) \leq E_{in}(h_g) + \sqrt{\frac{8}{N} \ln \frac{4[(2N)^{d_{vc}} + 1]}{\delta}}$$

$d_{vc} \downarrow$
 $d_{vc} \uparrow$
 $typ. \uparrow$
 $typ. \downarrow$

\downarrow
 \uparrow

BOUND DUE TO
 MODEL COMPLEXITY
 (FOR A GIVEN N, δ).

Fig. 2.3 of AML, adapted:



IMPLICATIONS IN DATASET METHODOLOGY

5

CONSIDER THIS SCENARIO;

I. 1. COLLECT DATA AND CONSTRUCT \mathcal{D} .

2. DO PRELIMINARY DATA ANALYSIS.

- VISUALIZE (OR LOOK AT) DATA,

- PLOT SOME HISTOGRAMS, LEARN
FROM THEM.

3. PREPROCESSING

~~RECHARGEZ~~

~~NORMALIZE~~ NORMALIZE

:

4. FEATURE EXTRACTION.

- ASSESS WHETHER EXTRACTED
FEATURES ARE USEFUL FOR
CLASSIFICATION.

5. DO SOME PRELIMINARY TRIALS ON \mathcal{D} .

- VARYING DIMENSION (#FEATURES)

- TRY A FEW HYPOTHESIS SETS AND
LEARNING ALGORITHMS.

6. SET UP AND RUN MODEL SELECTION AND LEARNING ALGORITHMS.

— DECIDE ON \mathcal{H}

— DIVIDE $\mathcal{D} \begin{cases} \mathcal{D}' \\ \mathcal{D}_{\text{Test}} \end{cases}$

— FOR CROSS-VALIDATION ON \mathcal{D}' TO CHOOSE PARAMETERS, TRAIN, AND FIND h_g .
(USING $\mathcal{D}' \begin{cases} \mathcal{D}_{\text{Tr}} \\ \mathcal{D}_{\text{Val}} \end{cases}$)

7. EVALUATE ITS PERFORMANCE.

(a) CALCULATE $E_{\text{in}}(h_g)$ USING \mathcal{D}_{Tr} OR \mathcal{D}_{Val}
CALCULATE E_{VC} USING d_{VC} OF \mathcal{H}
(ALSO N, δ).

GET "ERROR BAR" ON $E_{\text{out}}(h_g)$.

(b) CALCULATE $E_{\text{Test}}(h_g)$ USING $\mathcal{D}_{\text{Test}}$
CALCULATE E_M USING $M=1$
GET "ERROR BAR" ON $E_{\text{out}}(h_g)$.

ARE THESE $E_{out}(h_g)$ ERROR BARS VALID?

7(b) - No.

7(a) - No.

How to fix this?

For 7(b): ~~Do~~ MODIFY I TO SEPARATE

$\sigma \left(\begin{array}{l} \sigma' \\ \sigma_{Test} \end{array} \right)$ AFTER STEP 1. ~~These~~

THEN USE ONLY σ' IN STEPS 2-6.

~~7(b)~~ 7(b) WILL BE VALID.