

Thursday, 9/17/2020

EE 660

MACHINE LEARNING
FROM SIGNALS:
FOUNDATIONS AND METHODS

Prof. B. Keith Jenkins

Lecture 8

Announcements

- Homework 2 (Week 3) is due tomorrow (Friday).
- Homework 3 will be posted.
- Homework late submission policy has been posted
 - In Content > Syllabus and Course Information

Today's Lecture

- 2 aspects of learning feasibility
- Noisy targets
- Generalization error (revisited)
- Toward an effective number of hypotheses
 - Dichotomies
 - Growth function
 - Shattering
 - Break points
 - VC dimension (if time)

Noisy Target Functions [AML 1.4.2]

Sometimes the target function (class labels or regression output) is noisy.

Instead of $y = y(\underline{x}) = f(\underline{x}) = \text{deterministic}$, we let:

$$f(\underline{x}) = \text{deterministic} ; \quad y = y(\underline{x}) = \text{random} \sim p(y | \underline{x})$$

$$\text{Data points come from } p(\underline{x}, y) = p(y | \underline{x}) p(\underline{x})$$

Ex: $y(\underline{x}) = f(\underline{x}) + \epsilon_{\text{noise}}$ (regression)

$$y(\underline{x}) = \text{sgn} \{ g(\underline{x}) + \epsilon_{\text{noise}} \} \quad (\text{classification}) \quad f(\underline{x}) = \text{sgn} \{ g(\underline{x}) \}$$

$g(\underline{x}) = \text{discriminant function.}$

\underline{x}_n is drawn from $p(\underline{x}_n)$; $y_n = y(\underline{x}_n)$ is drawn from $p(y | \underline{x} = \underline{x}_n)$

$$\text{Then: } E_{\mathcal{D}}(h) = \frac{1}{N} \sum_{n=1}^N [h(\underline{x}_n) \neq y(\underline{x}_n)]$$

$$\text{and } E_{\text{out}}(h) = P[h(\underline{x}) \neq y(\underline{x})]$$

} for classification
(h, y are binary or integer)

We can reduce the noisy (probabilistic) case to the deterministic case:

$$\text{let } p(y|x) = \mathbb{I}[y=f(x)]$$

Thus, the noisy target is the more general case.

For the rest of the AML material, we will use the notation of the deterministic case most of the time, with the understanding that it can be generalized to the noisy target case as described above.

Generalization Error (revisited)

From before:

$$M = |\mathcal{H}| = \# \text{ of hypotheses in } \mathcal{H}$$

$$(2) \quad P[|E_{\mathcal{D}}(h) - E_{\text{out}}(h)| > \epsilon] \leq \underbrace{2Me^{-2\epsilon^2 N}}_{1-\delta}, \text{ for any } \epsilon > 0$$

True for any h in \mathcal{H} (including h_g) $\leftarrow \delta$

Re-arrange:

$$P[|E_{\text{out}}(h) - E_{\mathcal{D}}(h)| \leq \epsilon] \geq \underbrace{1 - 2Me^{-2\epsilon^2 N}}_{1-\delta}$$

δ is our tolerance. $\delta = 2Me^{-2\epsilon^2 N} \Rightarrow \epsilon = \sqrt{\frac{1}{2N} \ln \frac{2M}{\delta}} \triangleq \epsilon_m$

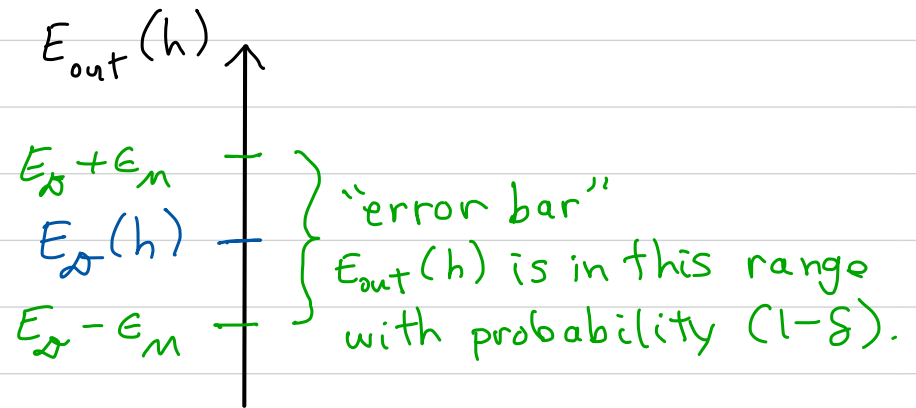
Note: $|E_{\text{out}} - E_{\mathcal{D}}| \leq \epsilon \Rightarrow$ [case $E_{\text{out}} \geq E_{\text{in}}$] (i) $E_{\text{out}} \leq E_{\mathcal{D}} + \epsilon$
 [case $E_{\text{out}} \leq E_{\text{in}}$] (ii) $E_{\text{out}} \geq E_{\mathcal{D}} - \epsilon$

$$(2.1) \therefore P[E_{\text{out}}(h_g) \leq E_{\mathcal{D}}(h_g) + \sqrt{\frac{1}{2N} \ln \frac{2M}{\delta}}] \geq 1 - \delta$$

$$(2.1)' \text{ and } P[E_{\text{out}}(h) \geq E_{\mathcal{D}}(h) - \underbrace{\sqrt{\frac{1}{2N} \ln \frac{2M}{\delta}}}_{\epsilon_m}] \geq 1 - \delta$$

(2.1) gives upper bound on $E_{\text{out}}(h_g)$ in terms of δ .

Note N, M dependence of ϵ_n



(2.1)' gives lower bound on $E_{\text{out}}(h)$, for any $h \in \mathcal{H}$.

If we had chosen a different hypothesis h_g' in \mathcal{H} , such that $E_{\delta}(h_g') > E_{\delta}(h_g)$,

$$\text{then } \underbrace{E_{\delta}(h_g') - \epsilon}_{\text{lower bound on } E_{\text{out}}(h_g')} > \underbrace{E_{\delta}(h_g) - \epsilon}_{\text{lower bound on } E_{\text{out}}(h_g)}$$

Two Aspects of Learning Feasibility [AML 1.3.3]

From last time: $\left\{ \begin{array}{l} P[|E_{\mathcal{S}}(h_g) - E_{\text{out}}(h_g)| > \epsilon] \leq 2Me^{-2\epsilon^2 N}, \text{ for any } \epsilon > 0 \\ M = |\mathcal{H}| = \# \text{ of hypotheses in } \mathcal{H} \end{array} \right.$

Feasibility of learning:
Can we get $E_{\text{out}}(h_g)$ small enough?

Can we get $|E_{\mathcal{S}}(h_g) - E_{\text{out}}(h_g)|$
or $E_{\mathcal{M}}$ small enough?

↑
Generalization error

Can we make $E_{\mathcal{S}}(h_g)$
small enough?

↑
Measured error of h_g on \mathcal{S}

Comment: if we use $\epsilon = \epsilon_n$ for generalization error with a linear perceptron classifier: ($\underline{w} = \underline{w}^{(*)}$)

$$\mathcal{H} = \{h_{\underline{w}}(\underline{x}) = \text{sgn}(\underline{w}^T \underline{x}) \mid \underline{w} \in \mathbb{R}^{D+1}\}$$

then $M = |\mathcal{H}| = \# \text{hypotheses} = ? = \infty$

$$\Rightarrow \text{generalization error} = \epsilon_n = \sqrt{\frac{1}{2N} \ln \frac{2M}{\delta}} = \infty$$

\Rightarrow We need a better measure than $M = |\mathcal{H}|$ for complexity of \mathcal{H} .

Toward an Effective Number of Hypotheses [AML 2.1.1]

For 2-class (binary) classification problems
 $\Rightarrow f(\underline{x}) \in \{-1, +1\}$ or $\{0, +1\}$

Consider how each $h_i(\underline{x}) \in \mathcal{H}$ behaves on N data points
 \underline{x}_n , $n=1, 2, 3, \dots, N$, drawn from \mathcal{X} .

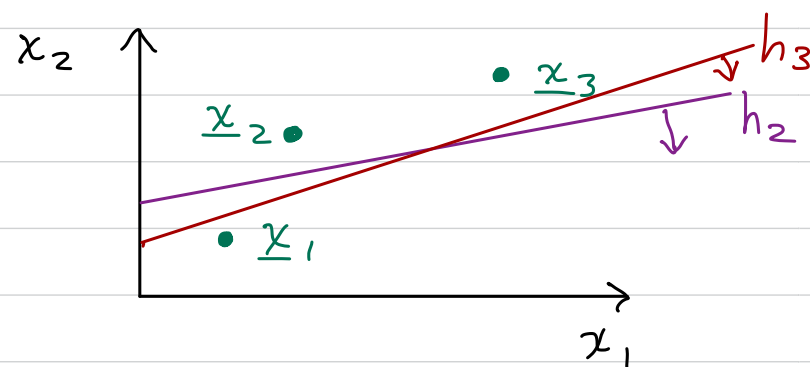
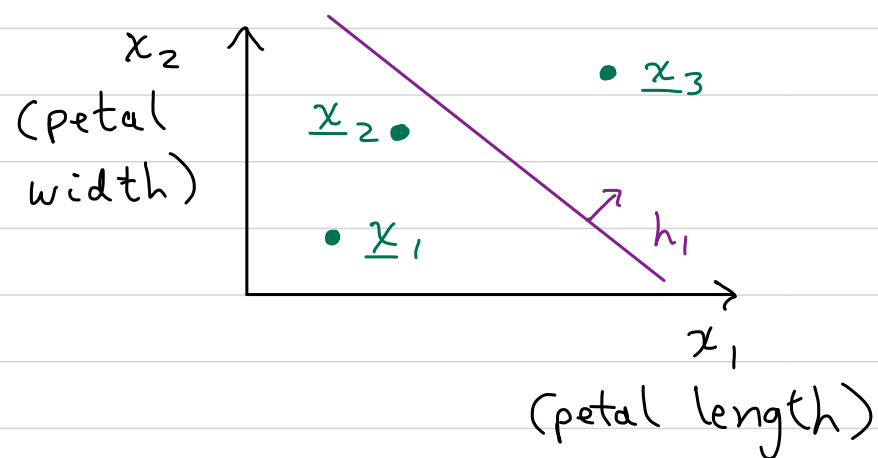
Def: The set of dichotomies generated by \mathcal{H} on $\{\underline{x}_n\}_{n=1}^N$ is

$$\mathcal{H}(\underline{x}_1, \underline{x}_2, \dots, \underline{x}_N) = \left\{ \underbrace{(h_i(\underline{x}_1), h_i(\underline{x}_2), \dots, h_i(\underline{x}_N))}_{N\text{-tuple}} \mid h_i \in \mathcal{H} \right\}.$$

(Duplicate N -tuples count as the same member of $\{\cdot\}$.)

Ex: Linear perceptron in 2D

Let $N=3$, $h=+1 \Rightarrow$ setosa, $h=-1 \Rightarrow$ virginica



$$h_1: (h_1(\underline{x}_1), h_1(\underline{x}_2), h_1(\underline{x}_3)) \\ = (-1, -1, +1)$$

$$h_2: (1, -1, -1) \\ h_3: (1, -1, -1)$$

Let $|\mathcal{H}(\underline{x}_1, \underline{x}_2, \dots, \underline{x}_N)| = \frac{\text{Cardinality of } \mathcal{H}(\underline{x}_1, \underline{x}_2, \dots, \underline{x}_N)}{\text{# of dichotomies (or # of unique N-tuples) that } \mathcal{H} \text{ can realize on the points } \{\underline{x}_n\}_{n=1}^N}$

Given N points, what is the maximum # of dichotomies? 2^N

Growth Function

Def: Growth function for a given \mathcal{H} is:

$$m_{\mathcal{H}}(N) \triangleq \max_{x_1, x_2, x_3, \dots, x_N \in \mathcal{X}} |\mathcal{H}(x_1, x_2, x_3, \dots, x_N)|$$

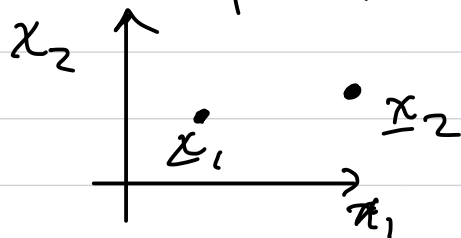
i.e., find the set of N points that maximizes the # of dichotomies that \mathcal{H} can realize. $m_{\mathcal{H}}(N) =$ this # of dichotomies.

$$m_{\mathcal{H}}(N) \leq 2^N \text{ always.}$$

[Example: in text and in discussion 4.]

Def: If \mathcal{H} can realize all possible dichotomies of a set of N points x_1, x_2, \dots, x_N , then \mathcal{H} can shatter x_1, x_2, \dots, x_N .

Ex: Can a linear perceptron in 2D shatter $N=2$ data points shown below?



→ Yes.


Break Points

Def: If there is no set of k distinct points that can be shattered by \mathcal{H} , then k is a break point for \mathcal{H} , and

$$m_{\mathcal{H}}(k) < 2^k$$

Ex: For $\mathcal{H}_L^{(1D)} = \{\text{linear perceptron in 1D}\}$,

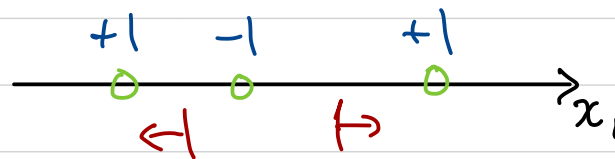
$N=2$:



$$\Rightarrow m_{\mathcal{H}_L^{(1D)}}(N=2) = 2^2 = 4$$

Can be shattered by $\mathcal{H}_L^{(1D)}$

$N=3$:



← $(+1, -1, +1)$ can't be realized by $\mathcal{H}_L^{(1D)}$.

$\Rightarrow k=3$ is a break point for $\mathcal{H}_L^{(1D)}$
Any $k \geq 3$ is a break point for $\mathcal{H}_L^{(1D)}$

True for all sets of $N=3$ distinct data points

VC Dimension [AML 2.1.3]

(Vapnik-Chervonenkis)

(binary classification problems)

VC dimension (or VCdim) is a measure of the "flexibility" or "complexity" of the hypothesis set (or decision boundary).

Def: The VC dimension of \mathcal{H} , $d_{VC}(\mathcal{H})$, is the largest value of N for which $m_{\mathcal{H}}(N) = 2^N$.

If $m_{\mathcal{H}}(N) = 2^N \quad \forall N$, then $d_{VC}(\mathcal{H}) = \infty$.