

ECE595 / STAT598: Machine Learning I

Lecture 26 Growth Function

Spring 2020

Stanley Chan

School of Electrical and Computer Engineering
Purdue University



Outline

- Lecture 25 Generalization
- **Lecture 26 Growth Function**
- Lecture 27 VC Dimension

Today's Lecture:

- **Overcoming the M Factor**
 - **Decisions based on Training Samples**
 - **Dichotomy**
- Examples of $m_{\mathcal{H}}(N)$
 - Finite 2D Set
 - Positive ray
 - Interval
 - Convex set

Probably Approximately Correct

- **Probably:** Quantify error using probability:

$$\mathbb{P}[|E_{\text{in}}(h) - E_{\text{out}}(h)| \leq \epsilon] \geq 1 - \delta$$

- **Approximately Correct:** In-sample error is an approximation of the out-sample error:

$$\mathbb{P}[|E_{\text{in}}(h) - E_{\text{out}}(h)| \leq \epsilon] \geq 1 - \delta$$

- If you can find an algorithm \mathcal{A} such that for any ϵ and δ , there exists an N which can make the above inequality holds, then we say that the target function is **PAC-learnable**.

The Factor “ M ”

- Testing

$$\mathbb{P}\left\{ |E_{\text{in}}(h) - E_{\text{out}}(h)| > \epsilon \right\} \leq 2e^{-2\epsilon^2 N},$$

- Training

$$\mathbb{P}\left\{ |E_{\text{in}}(g) - E_{\text{out}}(g)| > \epsilon \right\} \leq 2Me^{-2\epsilon^2 N}.$$

- So what? M is a constant.
- Bad news: M can be large, or even ∞ .
- A linear regression has $M = \infty$.
- Good news: It is possible to bound M .
- We will do it later.

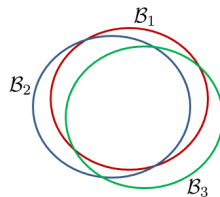
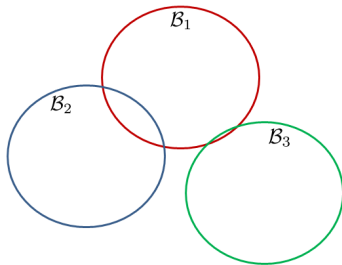
Overcoming the M Factor

- The *Bad* events \mathcal{B}_m are

$$\mathcal{B}_m = \{|E_{\text{in}}(h_m) - E_{\text{out}}(h_m)| > \epsilon\}$$

- The factor M is here because of the Union bound:

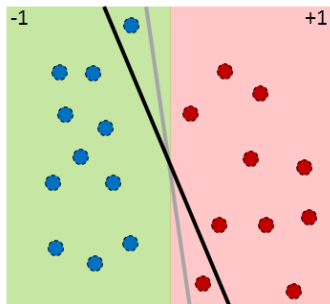
$$\mathbb{P}[\mathcal{B}_1 \text{ or } \dots \text{ or } \mathcal{B}_M] \leq \mathbb{P}[\mathcal{B}_1] + \dots + \mathbb{P}[\mathcal{B}_M].$$



Counting the Overlapping Area

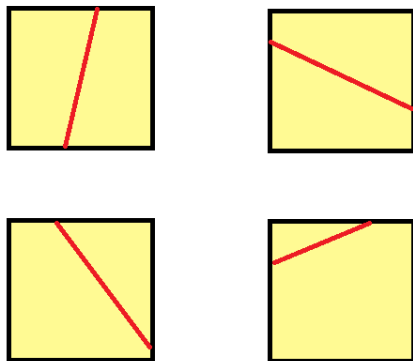
- ΔE_{out} = change in the +1 and -1 area
- Example below: Change a little bit
- ΔE_{in} = change in labels of the training samples
- Example below: Change a little bit, too
- So we should expect the probabilities

$$\mathbb{P}[|E_{\text{in}}(h_1) - E_{\text{out}}(h_1)| > \epsilon] \approx \mathbb{P}[|E_{\text{in}}(h_2) - E_{\text{out}}(h_2)| > \epsilon].$$



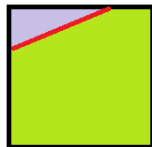
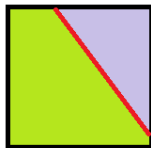
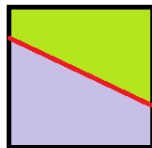
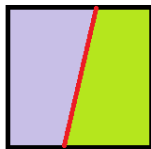
Looking at the Training Samples Only

- Here is our goal: Find something to replace M .
- But M is big because the whole input space is big.
- Let us look at the input space.



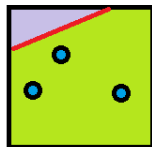
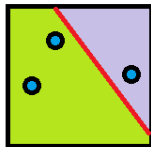
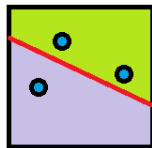
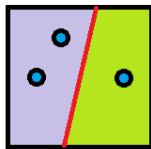
Looking at the Training Samples Only

- If you move the hypothesis a little, you get a different partition
- Literally there are infinitely many hypotheses
- This is M



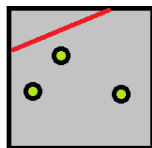
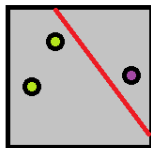
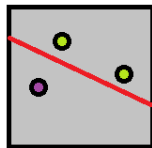
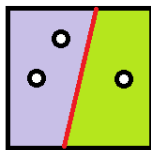
Looking at the Training Samples Only

- Here is our goal: Find something to replace M
- But M is big because the whole input space is big
- Can we restrict ourselves to just the training sets?



Looking at the Training Samples Only

- The idea is: Just look at the training samples
- Put a mask on your dataset
- Don't care until a training sample flips its sign

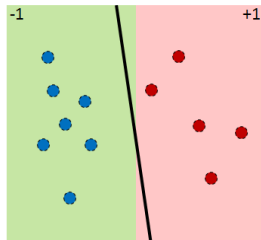
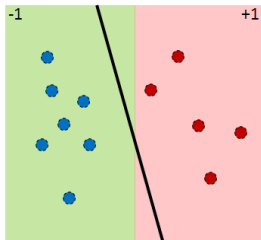


Dichotomies

- We need a new name: dichotomy.
- Dichotomy = mini-hypothesis.

Hypothesis	Dichotomy
$h : \mathcal{X} \rightarrow \{+1, -1\}$ for all population samples number can be infinite	$h : \{\mathbf{x}_1, \dots, \mathbf{x}_N\} \rightarrow \{+1, -1\}$ for training samples only number is at most 2^N

- Different hypothesis, same dichotomy.

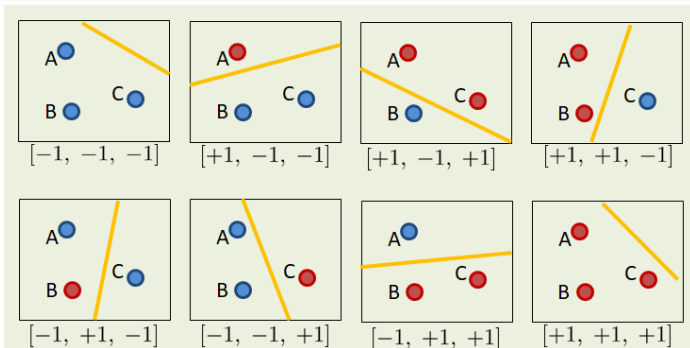


Dichotomy

Definition

Let $\mathbf{x}_1, \dots, \mathbf{x}_N \in \mathcal{X}$. The **dichotomies** generated by \mathcal{H} on these points are

$$\mathcal{H}(\mathbf{x}_1, \dots, \mathbf{x}_N) = \{(h(\mathbf{x}_1), \dots, h(\mathbf{x}_N)) \mid h \in \mathcal{H}\}.$$

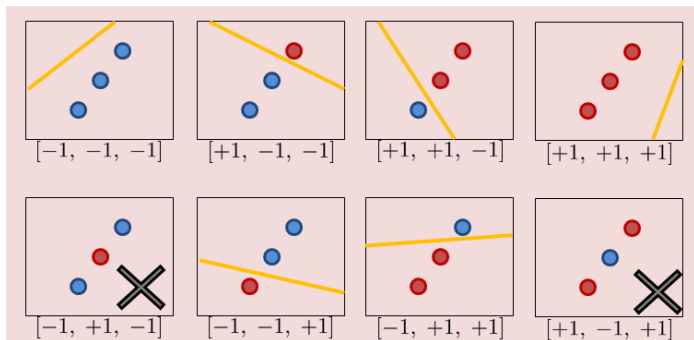


Dichotomy

Definition

Let $\mathbf{x}_1, \dots, \mathbf{x}_N \in \mathcal{X}$. The **dichotomies** generated by \mathcal{H} on these points are

$$\mathcal{H}(\mathbf{x}_1, \dots, \mathbf{x}_N) = \{(h(\mathbf{x}_1), \dots, h(\mathbf{x}_N)) \mid h \in \mathcal{H}\}.$$



Candidate to Replace M

- So here is our candidate replacement for M .
- Define **Growth Function**

$$m_{\mathcal{H}}(N) = \max_{\mathbf{x}_1, \dots, \mathbf{x}_N \in \mathcal{X}} |\mathcal{H}(\mathbf{x}_1, \dots, \mathbf{x}_N)|$$

- You give me a hypothesis set \mathcal{H}
- You tell me there are N training samples
- My job: Do whatever I can, by allocating $\mathbf{x}_1, \dots, \mathbf{x}_N$, so that the number of dichotomies is maximized
- Maximum number of dichotomy = the best I can do with your \mathcal{H}
- $m_{\mathcal{H}}(N)$: How expressive your hypothesis set \mathcal{H} is
- Large $m_{\mathcal{H}}(N)$ = more expressive \mathcal{H} = more complicated \mathcal{H}
- $m_{\mathcal{H}}(N)$ only depends on \mathcal{H} and N
- Doesn't depend on the learning algorithm \mathcal{A}
- Doesn't depend on the distribution $p(\mathbf{x})$ (because I'm giving you the max.)

Outline

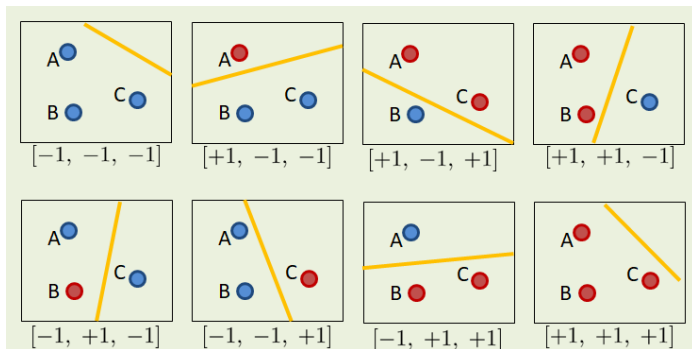
- Lecture 25 Generalization
- **Lecture 26 Growth Function**
- Lecture 27 VC Dimension

Today's Lecture:

- Overcoming the M Factor
 - Decisions based on Training Samples
 - Dichotomy
- **Examples of $m_{\mathcal{H}}(N)$**
 - **Finite 2D Set**
 - **Positive ray**
 - **Interval**
 - **Convex set**

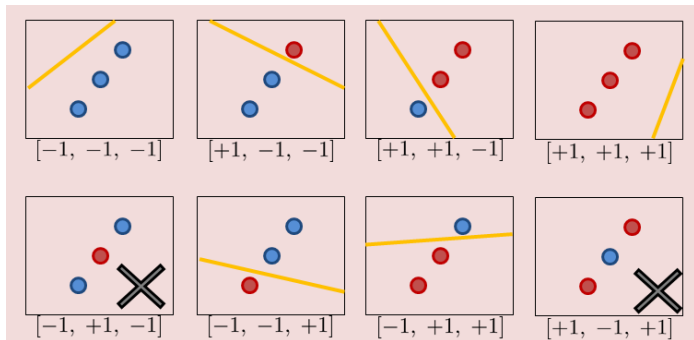
Examples of $m_{\mathcal{H}}(N)$

- \mathcal{H} = linear models in 2D
- $N = 3$
- How many dichotomies can I generate by moving the three points?
- This gives you 8. Are we the best?

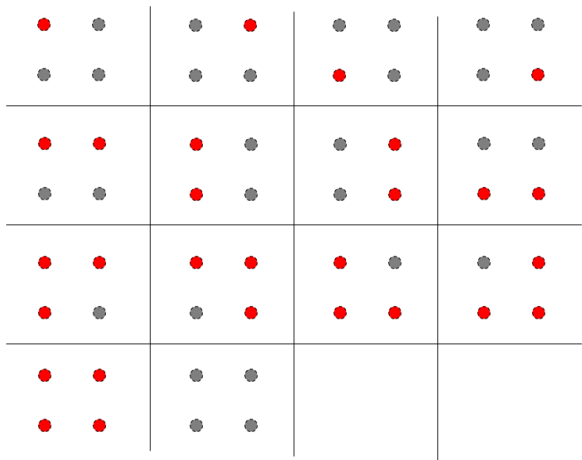


Examples of $m_{\mathcal{H}}(N)$

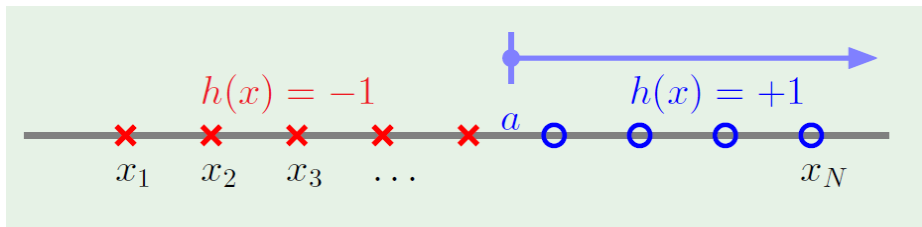
- \mathcal{H} = linear models in 2D
- $N = 3$
- How many dichotomies can I generate by moving the three points?
- This gives you 6. The previous is the best. So $m_{\mathcal{H}}(3) = 8$.



What about $m_{\mathcal{H}}(4)$? Ans: 14.

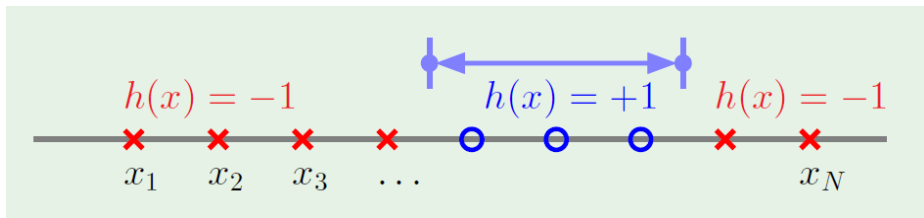


Another Example



- \mathcal{H} = set of $h: \mathbb{R} \rightarrow \{+1, -1\}$
- $h(x) = \text{sign}(x - a)$
- Cut the line into two halves
- You can only move along the line
- $m_{\mathcal{H}}(N) = N + 1$
- The N comes from the N points
- The $+1$ comes from the two ends

Another Example



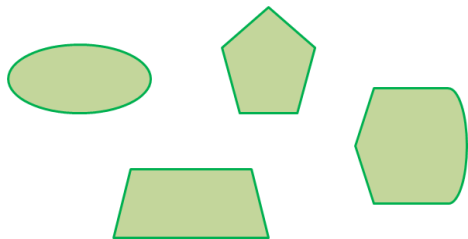
- $\mathcal{H} = \text{set of } h: \mathbb{R} \rightarrow \{+1, -1\}$
- Put an interval
- Length of the interval is N points

$$m_{\mathcal{H}}(N) = \binom{N+1}{2} + 1 = \frac{N^2}{2} + \frac{N}{2} + 1$$

- Think of $N + 1$ balls, pick 2.

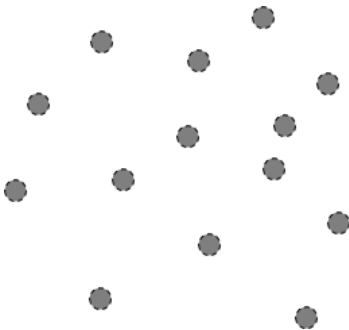
Another Example

- $\mathcal{H} = \text{set of } h: \mathbb{R}^2 \rightarrow \{+1, -1\}$
- $h(\mathbf{x}) = +1$ is convex
- Here are some examples

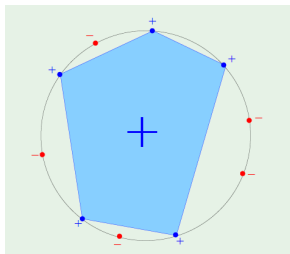


Another Example

- How about this collection of data points?
- Can you find an h such that you get a convex set?
- Yes. Do convex hull.
- Does it give you the maximum number of dichotomies?
- No. All interior points do not contribute.



Another Example



- The best you can do is this.
- Put all the points on a circle.
- Then you can get at most 2^N different dichotomies
- So

$$m_{\mathcal{H}}(N) = 2^N$$

- That is the best you can ever get with N points

Summary of the Examples

- \mathcal{H} is positive ray:

$$m_{\mathcal{H}}(N) = N + 1$$

- \mathcal{H} is positive interval:

$$m_{\mathcal{H}}(N) = \binom{N+1}{2} + 1 = \frac{N^2}{2} + \frac{N}{2} + 1$$

- \mathcal{H} is convex set:

$$m_{\mathcal{H}}(N) = 2^N$$

- So if we can replace M by $m_{\mathcal{H}}(N)$
- And if $m_{\mathcal{H}}(N)$ is a polynomial
- Then we are good.

Reading List

- Yasar Abu-Mostafa, Learning from Data, chapter 2.1
- Mehrya Mohri, Foundations of Machine Learning, Chapter 3.2
- Stanford Note <http://cs229.stanford.edu/notes/cs229-notes4.pdf>