

Learning Theory-Infinite Hypothesis sets

DINA ELREEDY

Generalization bound (M hypotheses)

- Hoeffding's Inequality of the selected g hypothesis out of M hypotheses using N training points.

$$\Pr[|E_{in}(g) - E_{out}(g)| > \varepsilon] \leq 2Me^{-2\varepsilon^2 N}$$

Then, with probability at least $1-\delta$:

$$E_{out}(g) \leq E_{in}(g) + \varepsilon$$

$$E_{out}(g) \leq E_{in}(g) + \sqrt{\frac{1}{2N} \ln\left(\frac{2M}{\delta}\right)}$$

$$\text{where } \delta = 2Me^{-2\varepsilon^2 N}$$

Generalization bound

- What if we have an infinite set of hypotheses?
- Example: the perceptron learning algorithm.
- Can we derive a generalization bound?

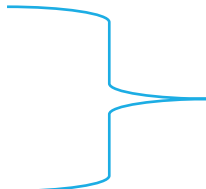
Review: Union bound of M hypotheses

Remember that M comes from applying a union bound.

Let the bad events B_i be:

$$\Pr[B_i] = P[|E_{in}(h_i) - E_{out}(h_i)| > \varepsilon]$$

The union bound for M hypotheses:

$$P[B_1 \text{ or } B_2 \text{ or } \dots \text{ or } B_M] \leq P[B_1] + P[B_2] + \dots + P[B_M] \leq \sum_{i=1}^M \Pr[B_i]$$


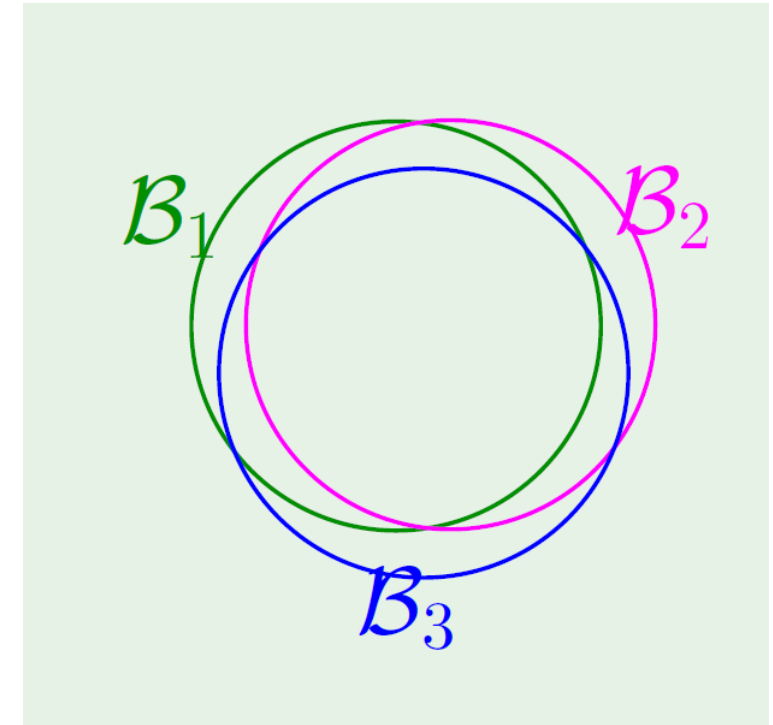
M terms

$$\text{Since } P(|E_{in}(h_i) - E_{out}(h_i)| > \varepsilon) \leq 2e^{-2\varepsilon^2 N}$$

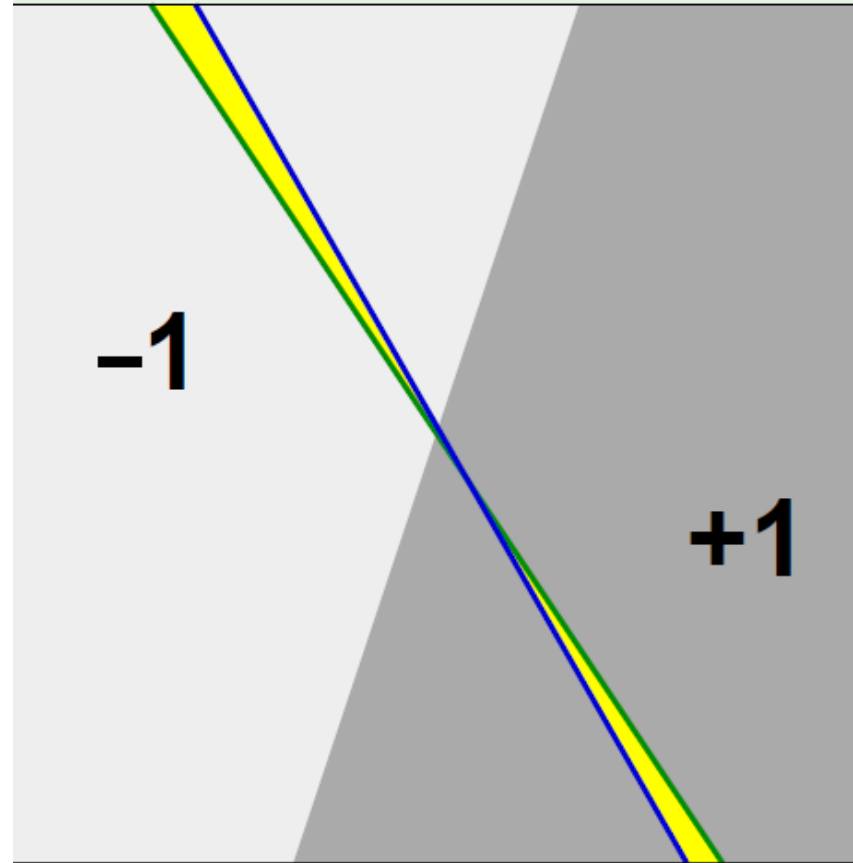
$$\text{Then, } P(|E_{in}(g) - E_{out}(g)| > \varepsilon) \leq 2Me^{-2\varepsilon^2 N}$$

Union bound of M hypotheses

- The good news is that the B events are overlapping!!
- Since many hypotheses are similar, so their corresponding B events $P[|E_{in}(h_i) - E_{out}(h_i)| > \varepsilon]$ are overlapping.
- Thus, the derived union bound was loose.
- We shall derive a tighter bound for infinite hypotheses.



Example- linear classifiers



Dichotomies

- For a binary classification problem:

A hypothesis $h: X \rightarrow \{+1, -1\}$.

- Apply a hypothesis h to a **finite sample of input points** not the whole input space.
- Assume we have N sample points $\{x_1, x_2, \dots, x_N\}$
- Apply $h \in H$ to the N points, we get a **dichotomy** which is an N -tuple of $(h(x_1), h(x_2), \dots, h(x_N))$.

Dichotomies (cont.)

- For a hypothesis set H , the dichotomies generated by H are defined as:

$$H(x_1, x_2, \dots, x_N) = \{(h(x_1), h(x_2), \dots, h(x_N)) | h \in H\}$$

- A hypothesis $h: X \rightarrow \{+1, -1\}$.
- A dichotomy: $\{x_1, x_2, \dots, x_N\} \rightarrow \{+1, -1\}$
- Number of hypotheses $|H|$ can be **infinite**.
- Maximum number of dichotomies $|H(x_1, x_2, \dots, x_N)|$ is **2^N** .
- The greater the number of dichotomies $|H(x_1, x_2, \dots, x_N)|$ is the more diverse and powerful the hypothesis set H .

Growth Function

- The growth function $m_H(N)$ for a hypothesis set H is defined as:
 - The **largest** number of dichotomies that can be generated by H on **any** N points.

$$m_H(N) = \max_{x_1, x_2, \dots, x_N \in X} |H(x_1, x_2, \dots, x_N)|$$

- Since the maximum number of dichotomies $|H(x_1, x_2, \dots, x_N)|$ is 2^N , then:

$$m_H(N) \leq 2^N$$

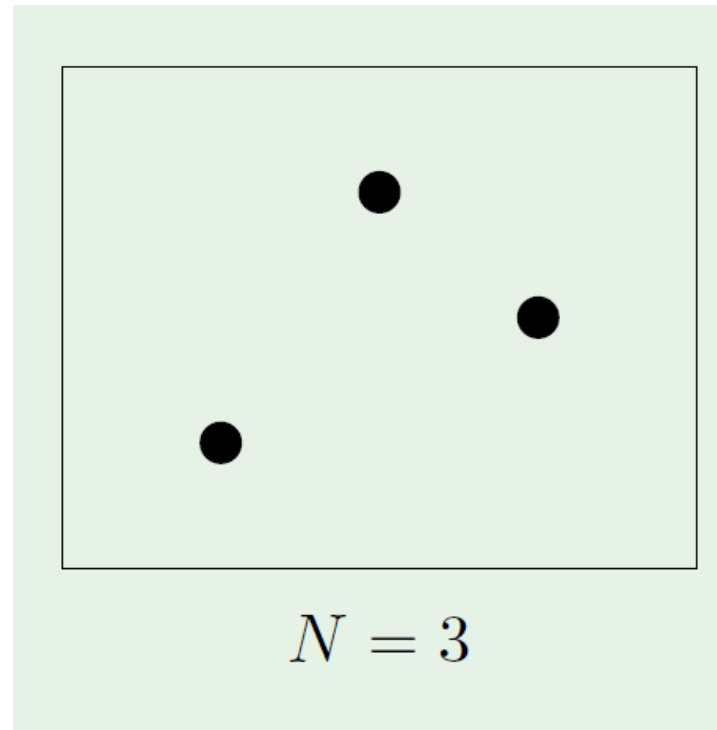
Growth Function (cont.)

- Given H , if H can generate all possible dichotomies on data points (x_1, x_2, \dots, x_N) such that $|H(x_1, x_2, \dots, x_N)| = 2^N$ then H **shatters** (x_1, x_2, \dots, x_N) .

Growth function-Perceptron

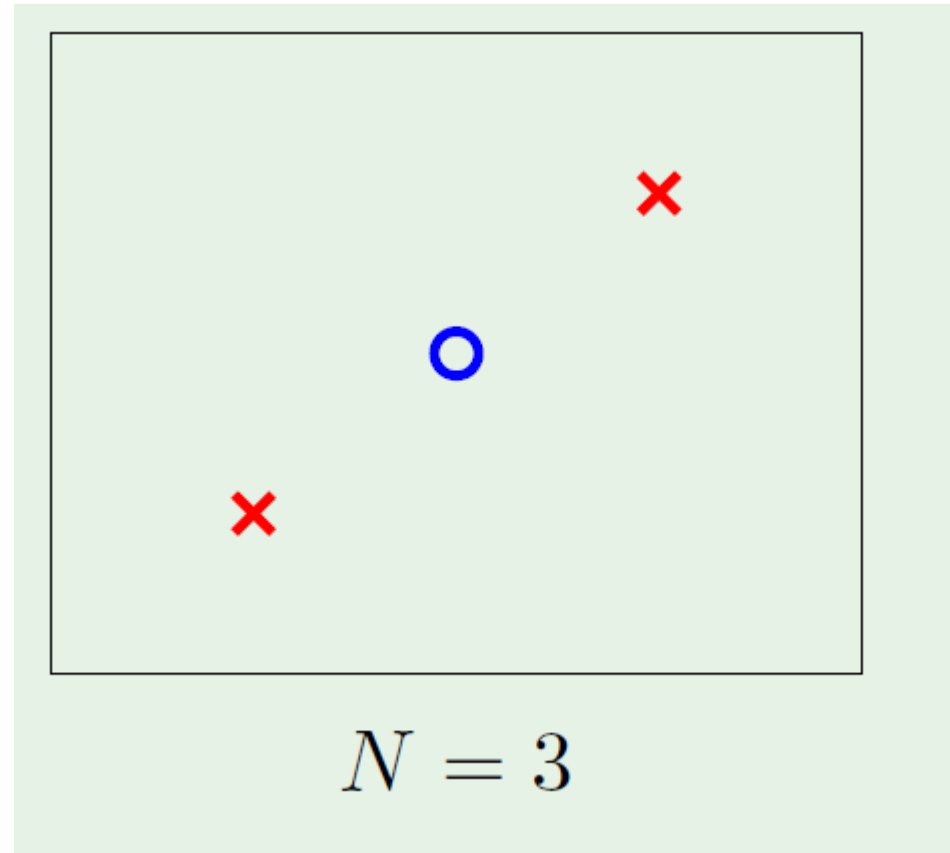
$N=3$

$$m_H(3) = 8$$



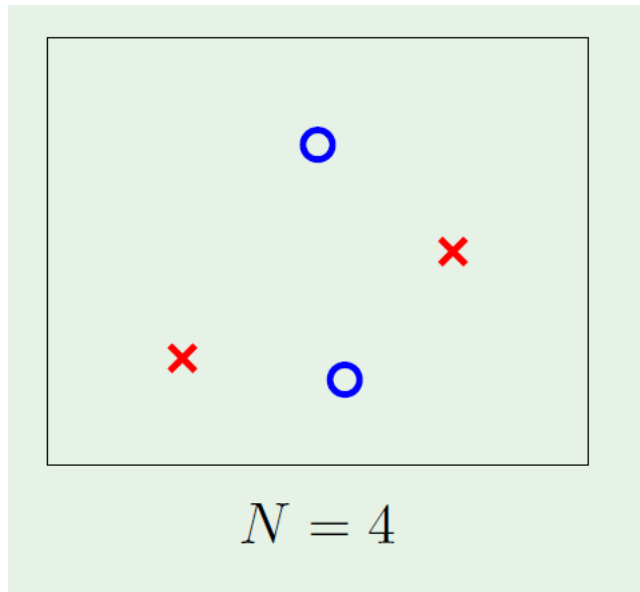
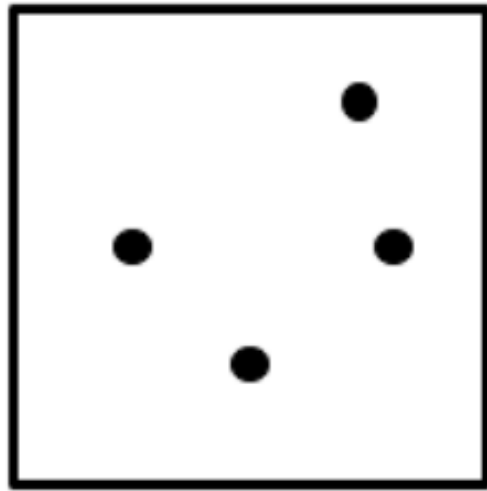
Growth function-Perceptron

$N=3$ (co-linear points)



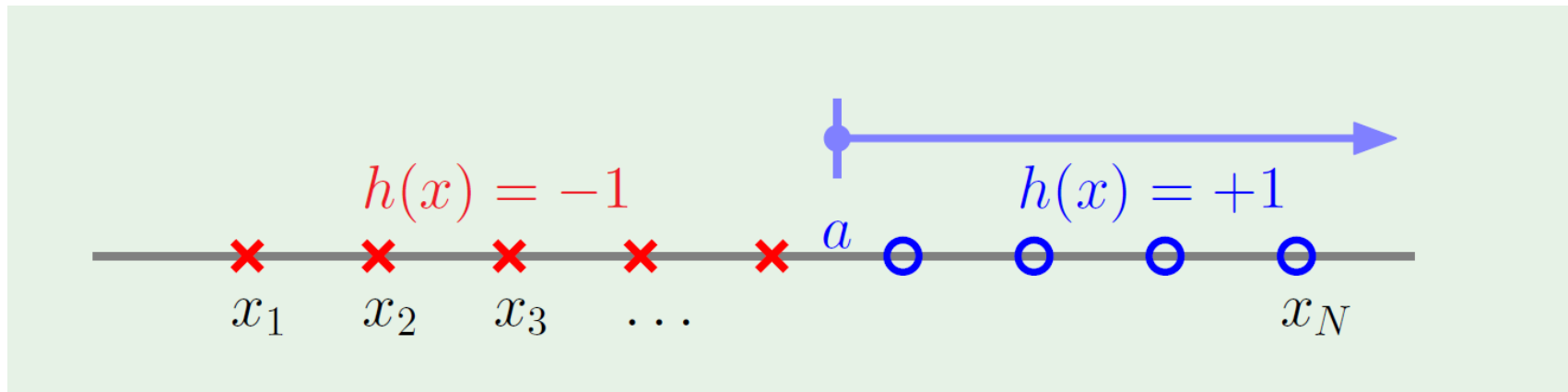
Growth function-Perceptron

$$m_H(4) = 14$$



Growth function- Positive rays

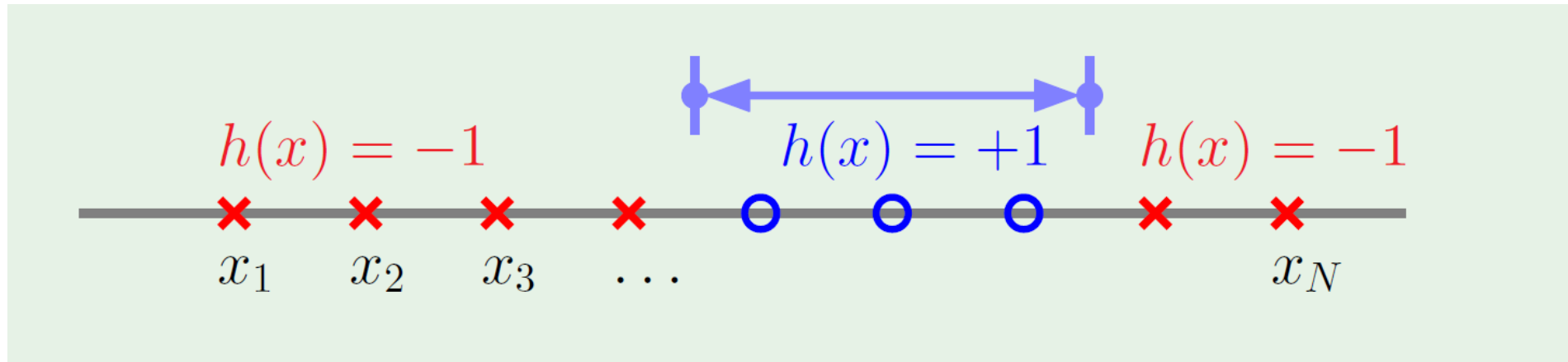
- H consists of all hypotheses h of the form: $h(x) = \text{sign}(x - a)$



- $m_{H(N)} = N + 1$

Growth function- Positive intervals

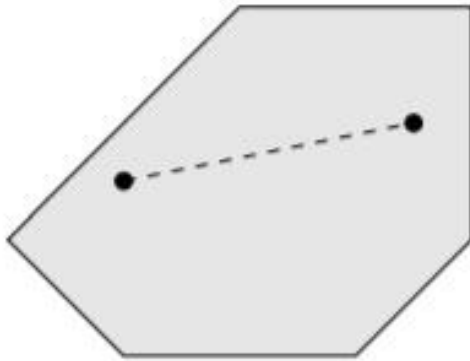
- H consists of all hypotheses in one dimension that returns +1 within some interval and -1 otherwise.
- Each hypothesis is specified by the two end values of the interval (a,b) .



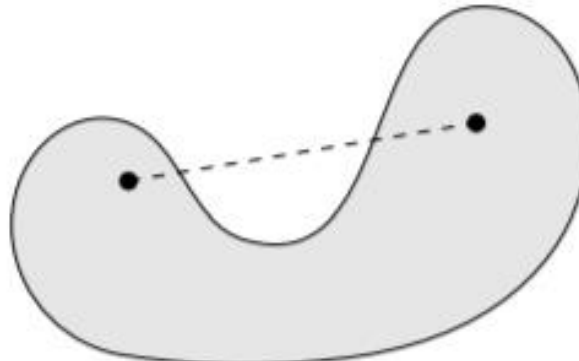
- $m_{H(N)} = \binom{N+1}{2} + 1 = \frac{1}{2}N^2 + \frac{1}{2}N + 1$

Convex sets

A set is convex if the line between any two points in the set entirely lies within the set.



CONVEX



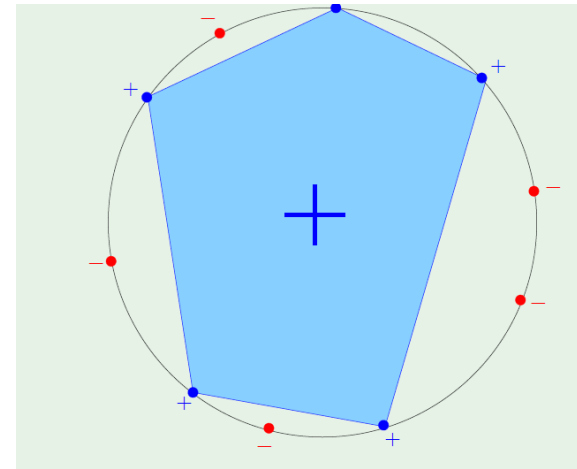
NOT CONVEX

Image source:

<https://faculty.math.illinois.edu/~mlavrov/docs/484-spring-2019/ch2lec1.pdf>

Growth function- Convex sets

- H consists of all the hypotheses in two dimensions that are positive inside a convex set and negative elsewhere.
 - Choose N points on the circumference of a circle.
 - Connect positive points with a polygon.
-
- $m_H(N) = 2^N$
 - H shatters these N points.



Breakpoint

- If **no data set** of size k **can be shattered** by H , then k is said to be a **breakpoint** for H .
- If k is a breakpoint then:

$$m_H(k) < 2^k$$

Breakpoint Examples

- For 2D perceptron:

$k=4$

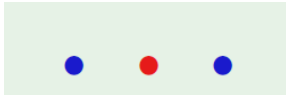
- For the 1D positive rays:

$k=2$



- For the 1D positive interval:

$k=3$



- For the convex sets:

$k=\infty$

The VC Dimension

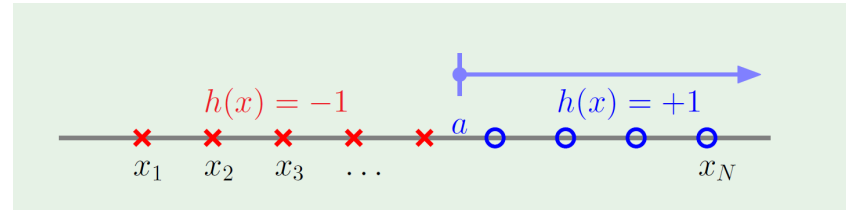
- The VC dimension of a hypothesis set H denoted by $d_{vc}(H)$ is the **largest value of N** for which **$m_H(N) = 2^N$** .
- $d_{vc}(H) = k - 1$
- If there is **no break point** for the hypothesis set H , **$m_H(N) = 2^N \quad \forall N$** , then $d_{vc}(H) = \infty$.

VC dimension-examples

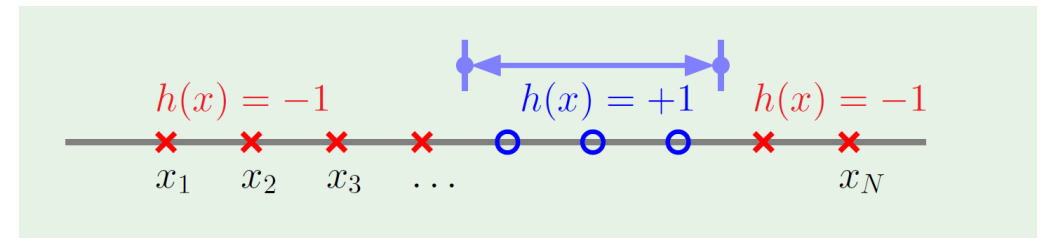
- The VC dimension for 2D perceptron is $d_{vc}(H) = 3$.
- The VC dimension for d-dimension perceptron is $d_{vc}(H) = d+1$.

VC dimension-examples (cont.)

- The VC dimension for positive rays $d_{vc}(H) = 1$.



- The VC dimension for positive intervals $d_{vc}(H) = 2$.



- VC dimension can be interpreted as the **effective** number of parameters (degrees of freedom).

Theorem

- If H has a break point, then $m_{H(N)}$ is bounded by a polynomial in N .

$$m_{H(N)} \leq \sum_{i=0}^{k-1} \binom{N}{i}$$

$$m_{H(N)} \leq \sum_{i=0}^{d_{vc}(H)} \binom{N}{i}$$

VC-Generalization bound for Infinite hypothesis set

The VC inequality using the growth function instead of M :

$$\Pr[|Ein(g) - Eout(g)| > \varepsilon] \leq 4m_H(2N)e^{-\frac{1}{8}\varepsilon^2 N}$$

VC Generalization bound

- The VC inequality:

$$\Pr[|E_{in}(g) - E_{out}(g)| > \varepsilon] \leq 4m_H(2N)e^{-\frac{1}{8}\varepsilon^2 N}$$

- Thus, with probability at least $1-\delta$ where $\delta = 4m_H(2N)e^{-\frac{1}{8}\varepsilon^2 N}$:

$$E_{out}(g) \leq E_{in}(g) + \varepsilon$$

- Accordingly, the VC generalization bound:

$$E_{out}(g) \leq E_{in}(g) + \sqrt{\frac{8}{N} \ln\left(\frac{4m_H(2N)}{\delta}\right)}$$

- If VC dimension is finite, then the growth function $m_H(2N)$ is polynomial, and the generalization error converges to zero as N increases.

Sample Complexity

- The sample complexity means the number of training examples N needed to achieve a certain generalization performance.
- The generalization performance is characterized by two parameters:
 - ϵ : Error tolerance defines the allowed generalization error
 - δ : Defines how often the error tolerance ϵ is violated.

Sample Complexity (cont.)

- To get a generalization error $\leq \varepsilon$:

$$\sqrt{\frac{8}{N} \ln\left(\frac{4m_{H(2N)}}{\delta}\right)} \leq \varepsilon$$

- Then, the sample complexity N to achieve that generalization error would be:

$$N \geq \frac{8}{\varepsilon^2} \ln\left(\frac{4m_{H(2N)}}{\delta}\right)$$

which is a function of N , solve it using numerical iterative methods.

As a rule of thumb $N \geq 10 d_{vc}(H)$

Model Complexity

- With probability at least $1 - \delta$:

$$E_{out}(g) \leq E_{in}(g) + \sqrt{\frac{8}{N} \ln\left(\frac{4m_{H(2N)}}{\delta}\right)}$$

- With fixed N (number of training samples), the term $\sqrt{\frac{8}{N} \ln\left(\frac{4m_{H(2N)}}{\delta}\right)}$ can be regarded as **model complexity**:

$$\Omega(N, H, \delta) = \sqrt{\frac{8}{N} \ln\left(\frac{4m_{H(2N)}}{\delta}\right)}$$

Thus, with probability at least $1 - \delta$:

$$E_{out}(g) \leq E_{in}(g) + \Omega(N, H, \delta)$$

Model complexity vs. generalization error Trade-off

Increases as dvc increases (complex models)

$$E_{out}(g) \leq E_{in}(g) +$$

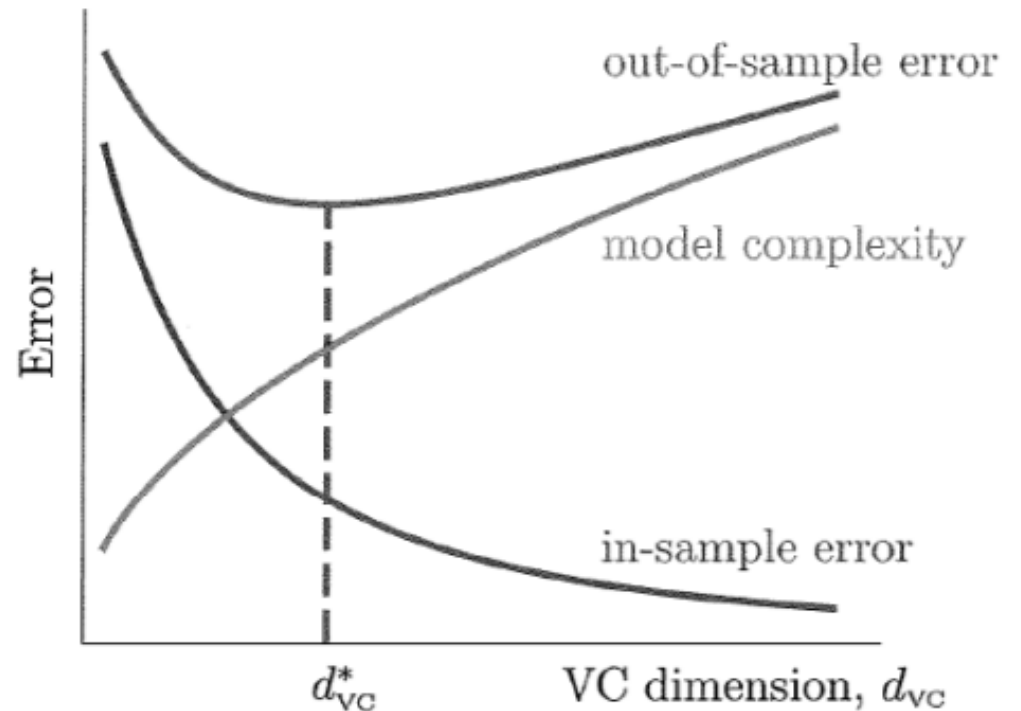
$$\sqrt{\frac{8}{N} \ln\left(\frac{4m_H(2N)}{\delta}\right)}$$

Decreases as dvc increases (complex models)

- The more complex H , the higher $m_H(2N)$ and the generalization error would be larger (worse).
- However, the more complex H , the less the in-sample error.

Generalization Error

- A very simple model with low d_{vc} , will not fit the training data well and will have a high in-sample error.
- A more complex learning model with higher d_{vc} would fit the training data better, resulting in a lower in-sample error, but the generalization will be worse.
- Some intermediate d_{vc}^* represents a trade-off between the two errors.



Ein vs. Eout

- **Ein** is the error on the training data. (in-sample error)
- **Eout** is the error over the entire input space X . (out of sample error)
- To estimate **Eout**, we should **use new test points that are never used for training.**

Estimating E_{out} in practice

- The VC bound is loose, we need a more accurate estimate of E_{out} for real-world applications.

Evaluate a sample estimate for E_{out} as follows:

- Use a fresh new test set of size K not involved in the training process.
 - Test the final hypothesis g on the test set and report E_{test} .
- However, what about the generalization error between E_{test} and E_{out} ?

We can apply Hoeffding's Inequality as g is not affected by test data.

$$P(|E_{test}(g) - E_{out}(g)| > \varepsilon) \leq 2e^{-2\varepsilon^2 K}$$

where K is the test set size.

Summary

- Dichotomies
- Growth function
- Breakpoint
- VC dimension
- VC generalization bound for infinite hypotheses
- Sample complexity
- Model Complexity and trade-off between in-sample and generalization error
- Estimate E_{out} in practice