

# Linear classifiers

# Perceptron

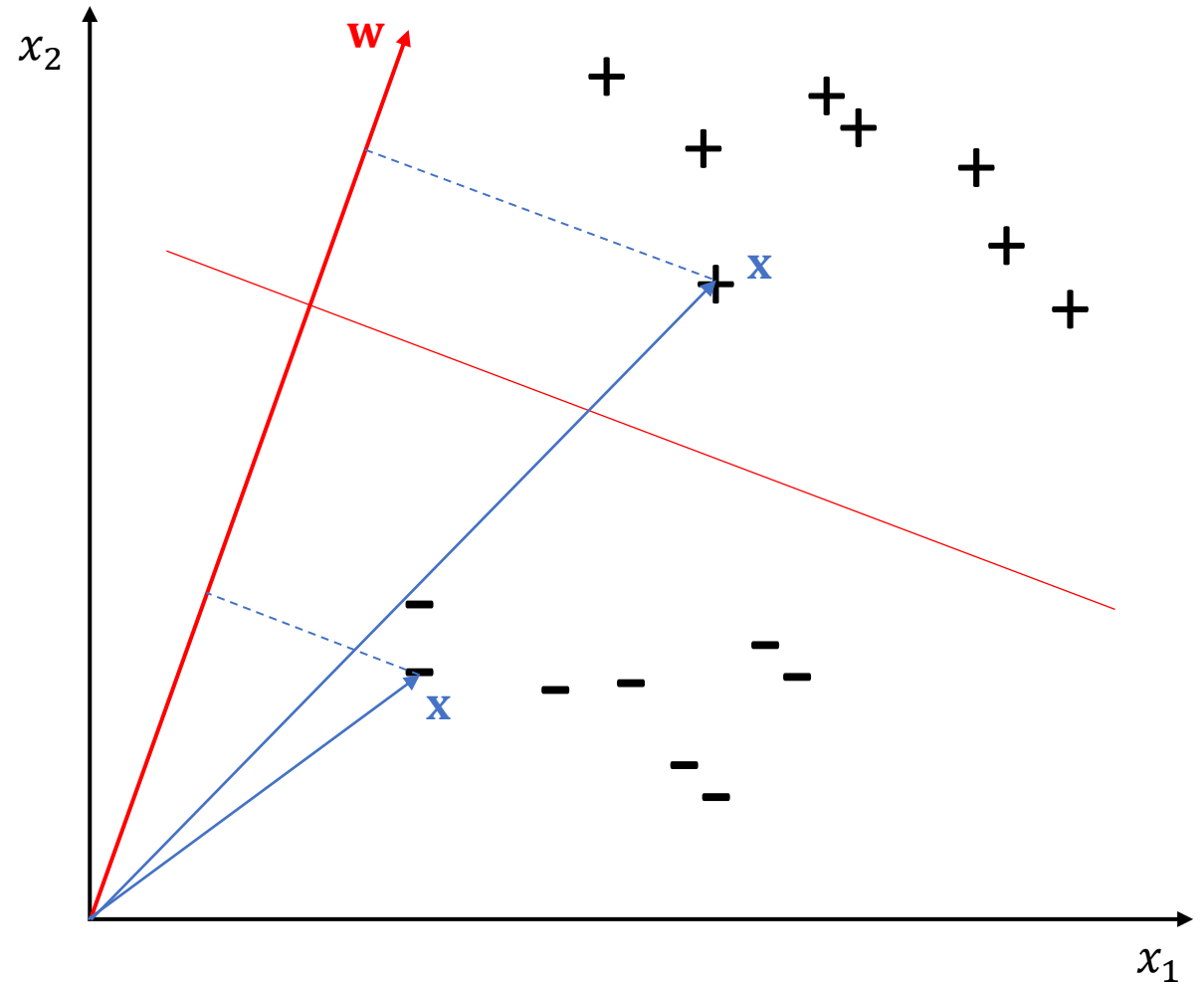
$$y \in \{-1, 1\}$$

$$h(\mathbf{x}) = \text{sign} \left( \left( \sum_{i=1}^d w_i x_i \right) - \text{threshold} \right)$$

$$h(\mathbf{x}) = \text{sign} \left( \left( \sum_{i=1}^d w_i x_i \right) - w_0 \right)$$

$$h(\mathbf{x}) = \text{sign} \left( \sum_{i=0}^d w_i x_i \right)$$

$$h(\mathbf{x}) = \text{sign}(\mathbf{w}^T \mathbf{x})$$



# Perceptron training

$$h(\mathbf{x}) = \text{sign}(\mathbf{w}^T \mathbf{x})$$

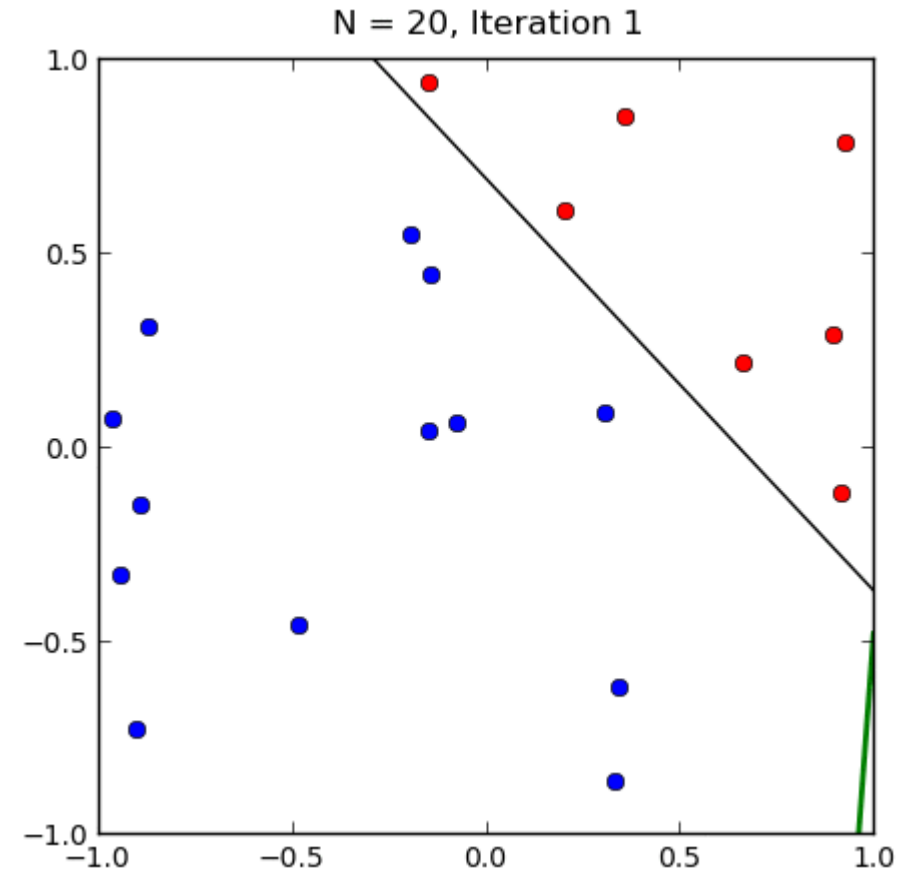
Data:  $(\mathbf{x}_1, y_1), (\mathbf{x}_2, y_2), \dots, (\mathbf{x}_N, y_N)$

Algorithm:

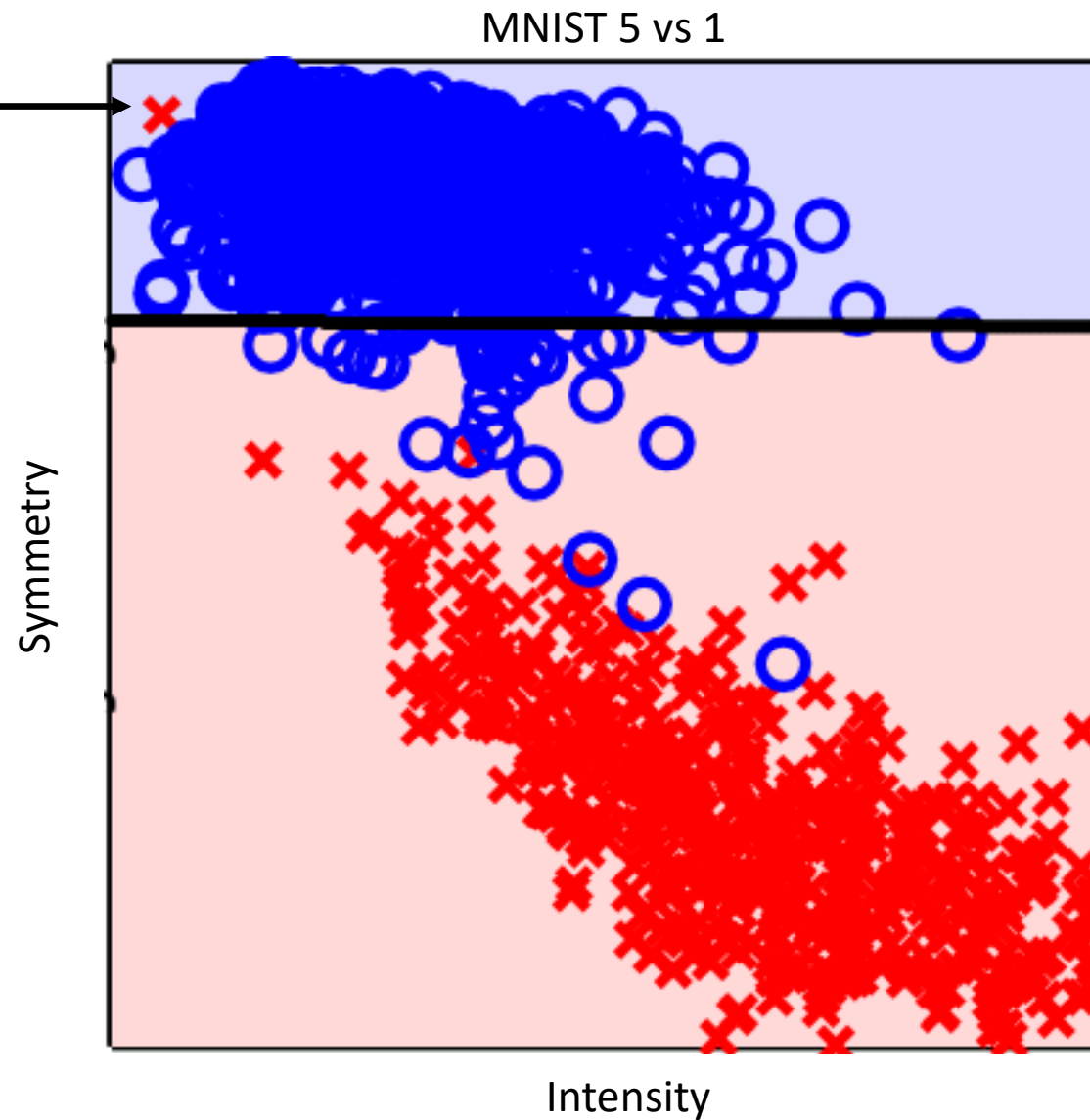
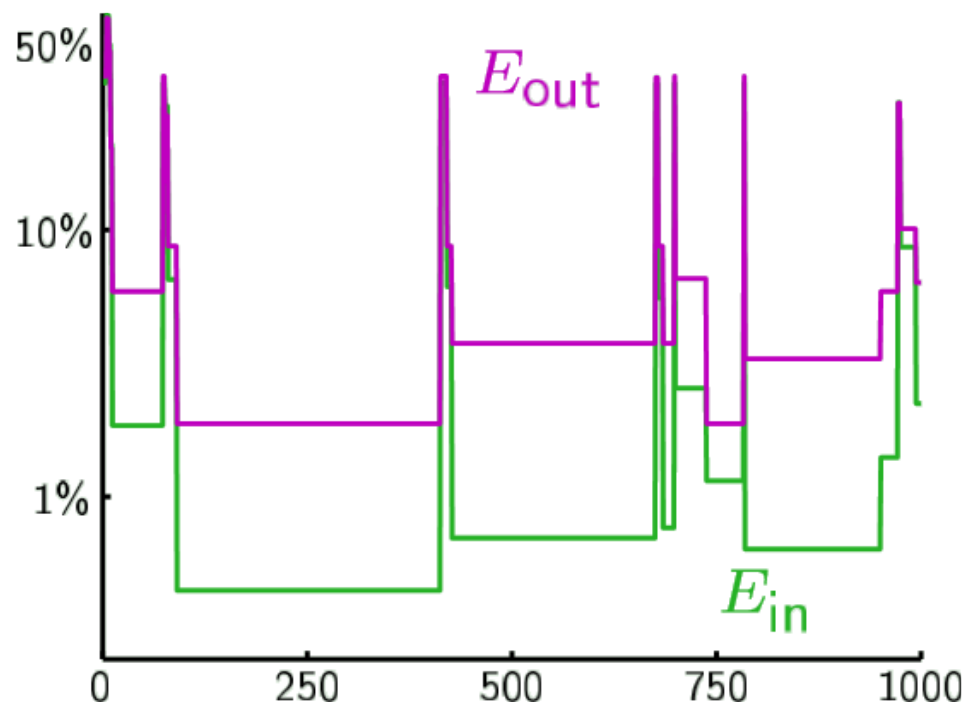
Start with the random  $\mathbf{w}$

Find such  $\mathbf{x}_i$ , that  $h(\mathbf{x}_i) \neq y_i$

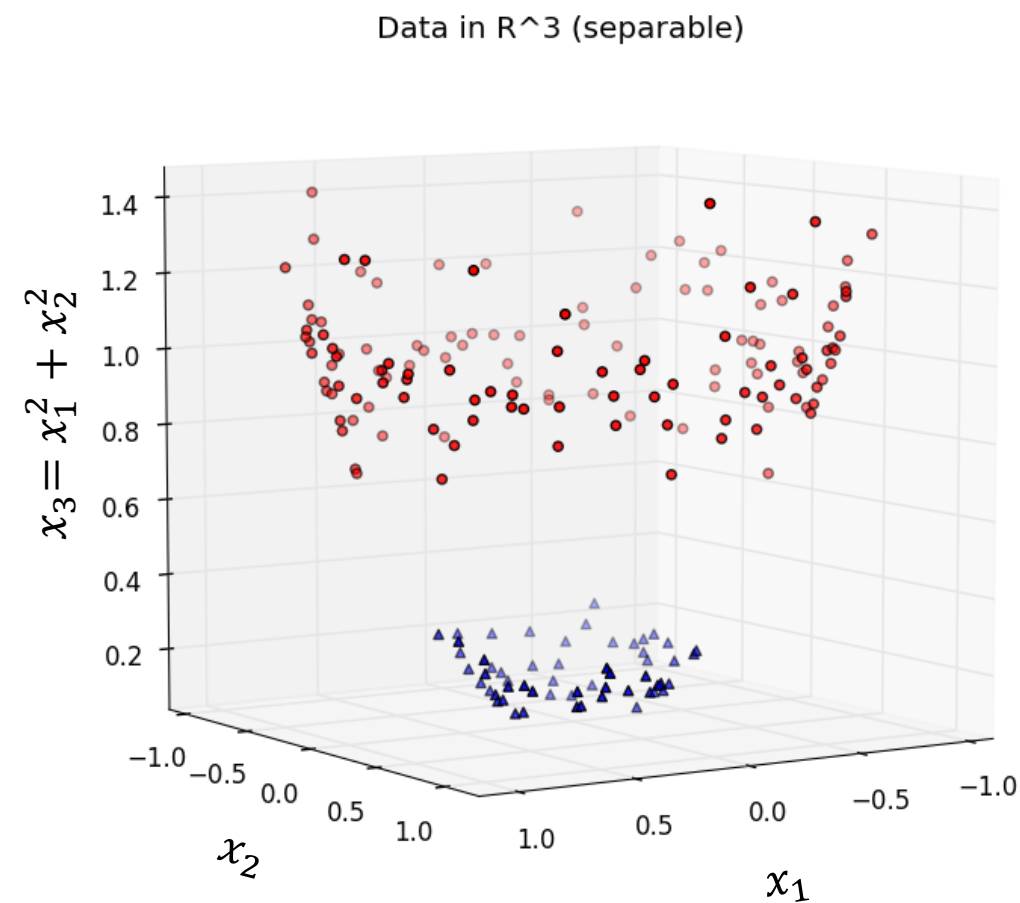
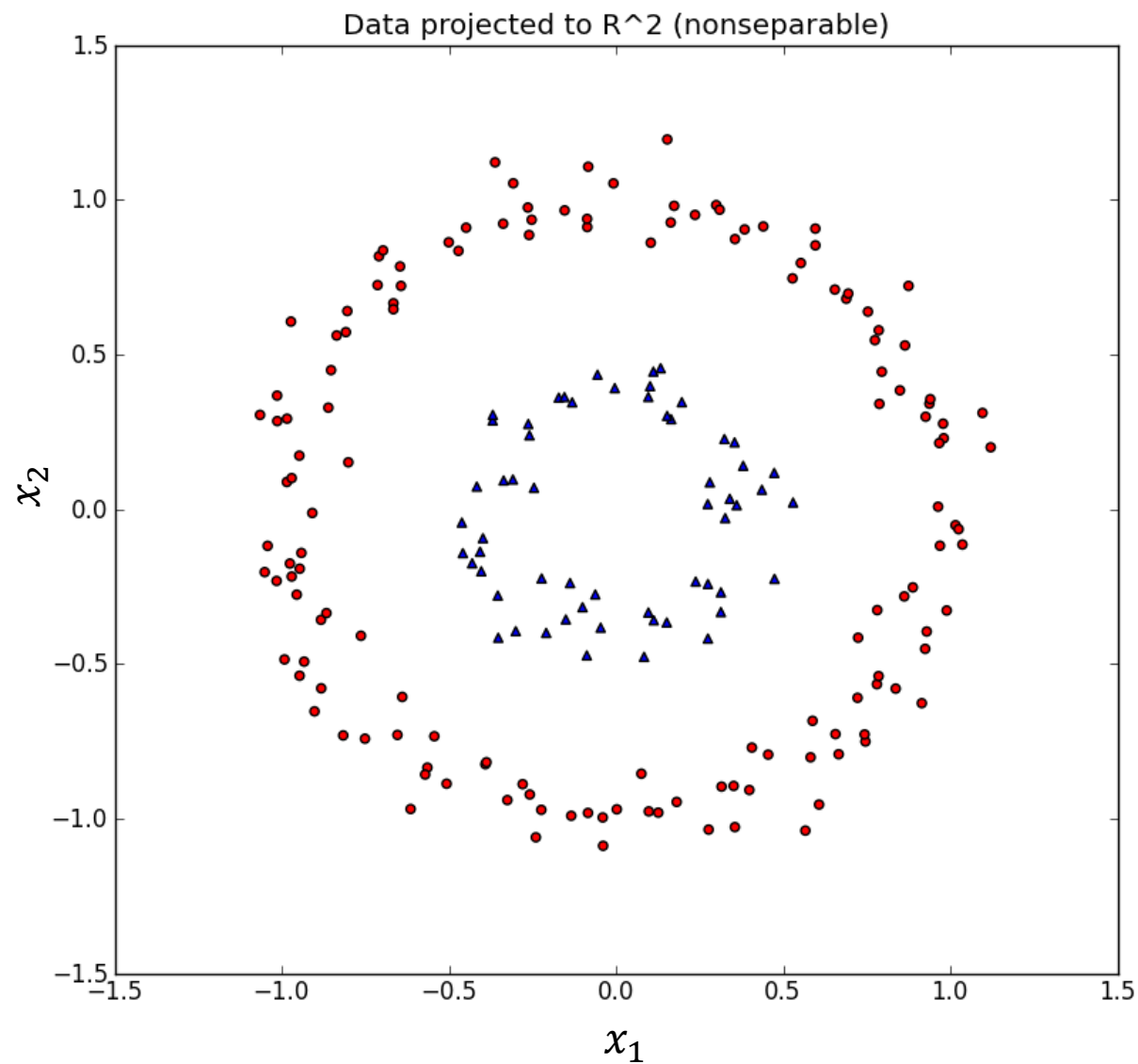
$$\mathbf{w} \leftarrow \mathbf{w} + y_i \mathbf{x}_i$$



# Pocket algorithm



# Feature engineering



# Theory of error

$$E_{in}(h) = \frac{1}{N} \sum_1^N e(h(\mathbf{x}_n), f(\mathbf{x}_n)) \quad \text{in sample}$$

$$E_{out}(h) = E_{\mathbf{x}}[e(h(\mathbf{x}), f(\mathbf{x}))] \quad \text{out of sample}$$

# Hoeffding's inequality

$$P[|E_{in}(h) - E_{out}(h)| > \epsilon] \leq 2e^{-\epsilon^2 N}$$

Generalization error

Sometimes  $E_{out}$  is denoted as generalization error

# Hoeffding's inequality

$$P[|E_{in}(h) - E_{out}(h)| > \epsilon] \leq 2e^{-\epsilon^2 N}$$

$$P[|E_{in}(h) - E_{out}(h)| > \epsilon] \leq M2e^{-\epsilon^2 N}$$

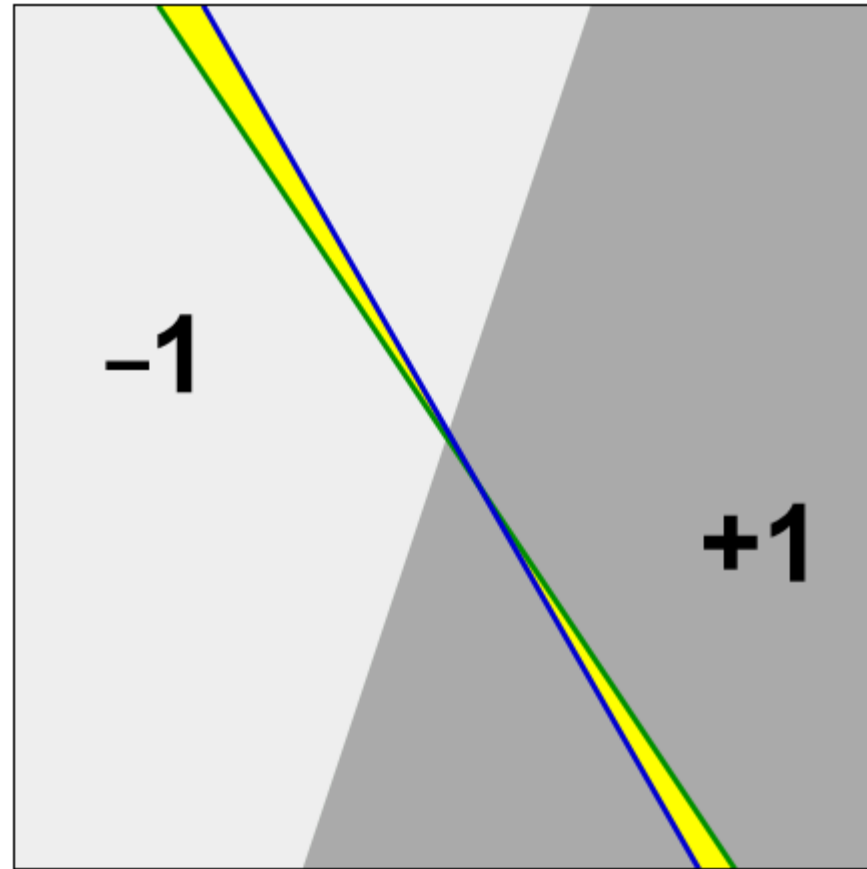
M is the number of hypothesis



# Close hypotheses

$$|E_{\text{in}}(h_1) - E_{\text{out}}(h_1)| \approx |E_{\text{in}}(h_2) - E_{\text{out}}(h_2)|$$

up



down

# From hypotheses to dichotomies

- Hypothesis:  $h: X \rightarrow \{-1, +1\}$
- Dichotomy:  $h: \{\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_N\} \rightarrow \{-1, +1\}$
- Maximum number of dichotomies:  $2^N$

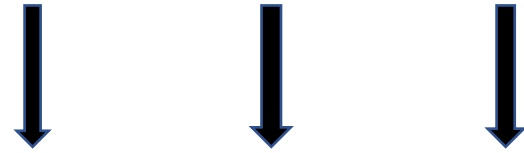
## Growth function

$$m_H(N) = \max_{\mathbf{x}_1, \dots, \mathbf{x}_N} |H(\mathbf{x}_1, \dots, \mathbf{x}_N)|$$

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

# Vapnik–Chervonenkis Inequality

$$P[|E_{in}(h) - E_{out}(h)| > \epsilon] \leq M2e^{-\epsilon^2 N}$$

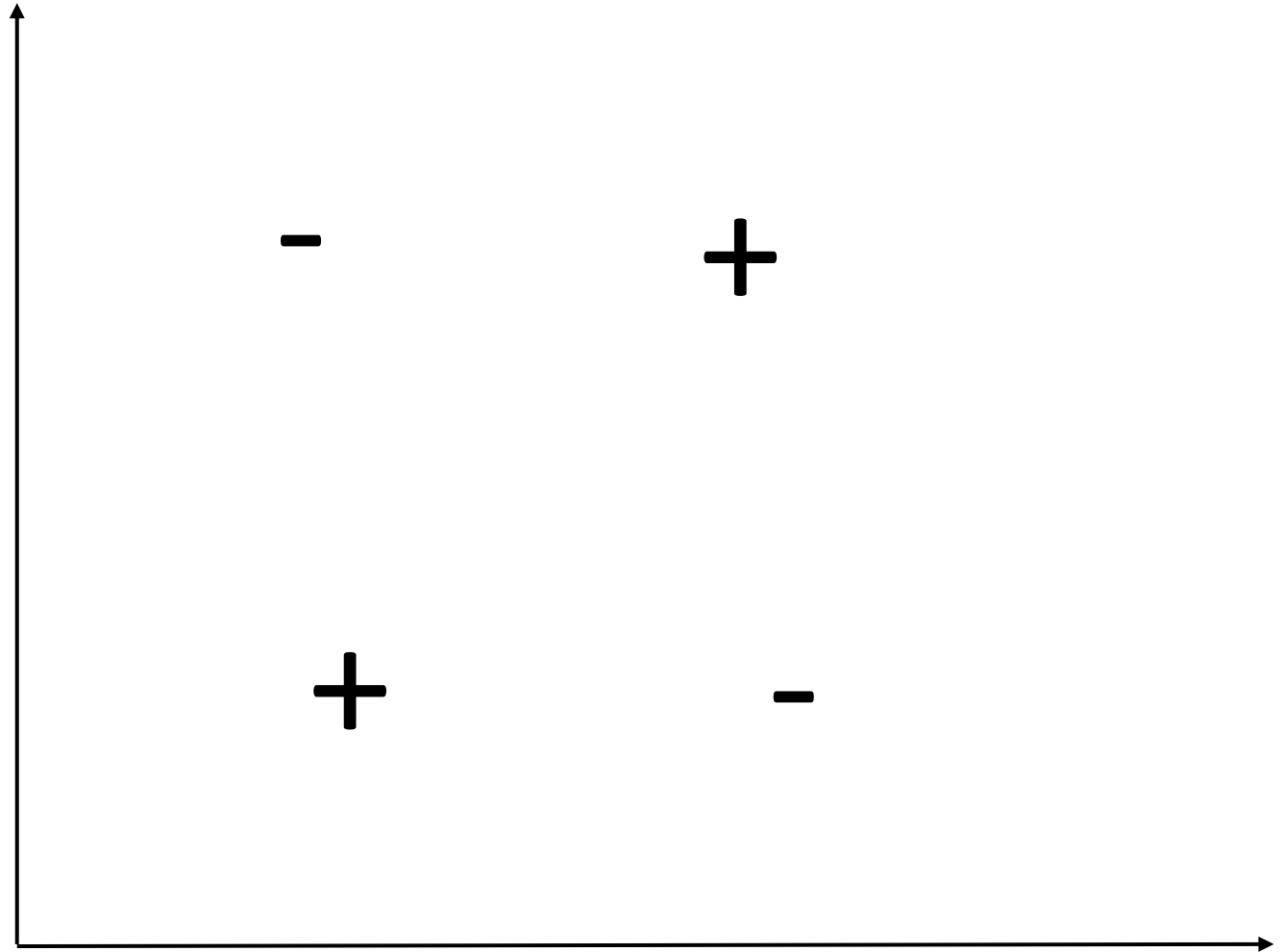


$$P[|E_{in}(h) - E_{out}(h)| > \epsilon] \leq m_H(\textcolor{red}{2}N)\textcolor{red}{4}e^{-\epsilon^{\textcolor{red}{1}\over\textcolor{red}{8}}N}$$

# Breakpoint

$$\min(k: m_H(k) < 2^k)$$

For 2D perceptron,  **$k = 4$** .



# Proof of polynomiality of a growth function in the presence of a breakpoint

- $B(N, k) = m_H(N)$  with breakpoint  $k$
- $B(N, k) = \alpha + 2\beta$
- $\alpha + \beta \leq B(N - 1, k)$
- $\beta \leq B(N - 1, k - 1)$
- $B(N, k) \leq B(N - 1, k) + B(N - 1, k - 1)$
- Let's prove that  $B(N, k) \leq \sum_{i=0}^{k-1} C_N^i$

	$\mathbf{x}_1$	$\mathbf{x}_2$	$\dots$	$\mathbf{x}_{N-1}$	$\mathbf{x}_N$	
$\alpha$	+1	+1	$\dots$	+1	+1	1 class for $\mathbf{x}_N$
	-1	+1	$\dots$	+1	-1	
	$\vdots$	$\vdots$	$\vdots$	$\vdots$	$\vdots$	
	+1	-1	$\dots$	-1	-1	
	-1	+1	$\dots$	-1	+1	
$\beta$	+1	-1	$\dots$	+1	+1	2 classes for $\mathbf{x}_N$
	-1	-1	$\dots$	+1	+1	
	$\vdots$	$\vdots$	$\vdots$	$\vdots$	$\vdots$	
	+1	-1	$\dots$	+1	+1	
	-1	-1	$\dots$	-1	+1	
$\beta$	+1	-1	$\dots$	+1	-1	2 classes for $\mathbf{x}_N$
	-1	-1	$\dots$	+1	-1	
	$\vdots$	$\vdots$	$\vdots$	$\vdots$	$\vdots$	
	+1	-1	$\dots$	+1	-1	
	-1	-1	$\dots$	-1	-1	

# Proof

$$B(N, k) \leq \sum_{i=0}^{k-1} C_N^i$$

$$B(N, 1) = 1, \quad B(1, k > 1) = 2$$

$$B(N, k) \leq B(N-1, k) + B(N-1, k-1) \leq \sum_{i=0}^{k-1} C_{N-1}^i + \sum_{i=0}^{k-2} C_{N-1}^i =$$

$$= 1 + \sum_{i=1}^{k-1} C_{N-1}^i + \sum_{i=1}^{k-1} C_{N-1}^{i-1} = 1 + \sum_{i=1}^{k-1} (C_{N-1}^i + C_{N-1}^{i-1}) = 1 + \sum_{i=1}^{k-1} C_N^i = \sum_{i=0}^{k-1} C_N^i$$

# VC-dimension

$d_{VC}(H)$  for hypotheses type  $H$ , is the maximum number  $N$ ,  
such that  $m_H(N) = 2^N$ .

$d_{VC}(H) = k - 1$ , where  $k$  – is the breakpoint.

# Growth function and VC-dimension

$$m_H(N) \leq \sum_{i=0}^{k-1} C_N^i$$

$$m_H(N) \leq \sum_{i=0}^{d_{VC}} C_N^i \leq N^{d_{VC}} + 1$$



# VC-dimension of perceptron

If  $d$  is the space dimensionality,  $d_{VC} = d + 1$

$$\begin{aligned}
 &\mathbf{x}_1, \mathbf{x}_2, \mathbf{x}_3 \dots, \mathbf{x}_{d+1} \\
 &\text{sign}(\mathbf{X}\mathbf{w}) = \mathbf{y} \\
 &\mathbf{X}\mathbf{w} = \mathbf{y} \\
 &\mathbf{w} = \mathbf{X}^{-1}\mathbf{y}
 \end{aligned}
 \quad
 \mathbf{X} = \begin{bmatrix} -\mathbf{x}_1^T & - \\ -\mathbf{x}_2^T & - \\ -\mathbf{x}_3^T & - \\ \vdots & \\ -\mathbf{x}_k^T & - \end{bmatrix} = \begin{bmatrix} 1 & 0 & 0 & 0 & \dots & 0 \\ 1 & 1 & 0 & 0 & \dots & 0 \\ 1 & 0 & 1 & 0 & \dots & 0 \\ & & & \vdots & & \\ 1 & 0 & 0 & 0 & \dots & 1 \end{bmatrix}$$

$$\begin{aligned}
 &\mathbf{x}_1, \mathbf{x}_2, \mathbf{x}_3 \dots, \mathbf{x}_{d+1}, \mathbf{x}_{d+2} \\
 &\mathbf{x}_j = \sum_{i \neq j} \mathbf{x}_i a_i \quad \mathbf{w}^T \mathbf{x}_j = \sum_{i \neq j} \mathbf{w}^T \mathbf{x}_i a_i \quad y_i = \text{sign}(a_i) \quad y_j = -1
 \end{aligned}$$

Sufficient data

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

$$\approx N^{d_{VC}} e^{-N}$$

$$N \geq 10 d_{VC}$$

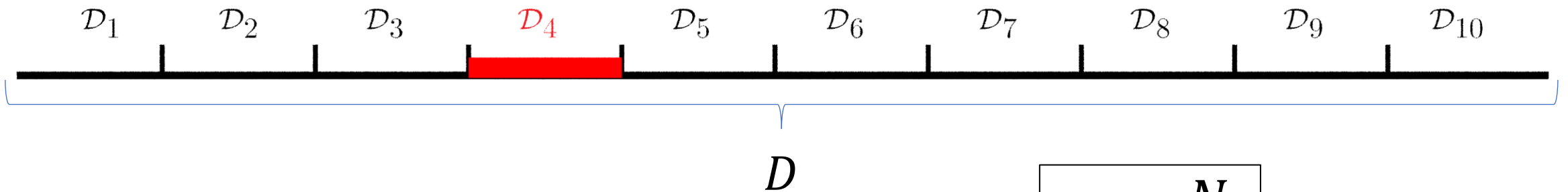
# Validation

$$(\mathbf{x}_1, \mathbf{y}_1), (\mathbf{x}_2, \mathbf{y}_2), \dots, (\mathbf{x}_K, \mathbf{y}_K) \in \mathcal{D}_{val} \quad E_{val}(h) = \frac{1}{K} \sum_1^K e(h(\mathbf{x}_k), f(\mathbf{x}_k))$$

$$P[|E_{val}(h) - E_{out}(h)| > \epsilon] \leq 2e^{-\epsilon^2 N}$$

$$K = \frac{N}{5}$$

# Cross-validation



$$K = \frac{N}{10}$$

# Train-Val-Test

- Train the algorithm on the **train** dataset.
- Optimize hyperparameters on **val (cross-val)**.
- Check the final performance on **test**.