



# Teoria systemów uczących się i wymiar Vapnika-Chervonenkisa

*Hung Son Nguyen*

Wydział Matematyki, Informatyki i Mechaniki  
Uniwersytet Warszawski  
email: `son@mimuw.edu.pl`

Grudzień 2009

# Plan wykładu

---



- 1 Wstęp do komputerowego uczenia się pojęć
- 2 Model PAC (probably approximately correct)
- 3 Wyuczalność klasy pojęć
- 4 Wymiar Vapnika Chervonenkisa (VC dimension)
- 5 Podstawowe twierdzenia teorii uczenia się
- 6 Appendix: „Nie ma nic za darmo” czyli “Non Free Lunch Theorem”



- Np. Pokazać, że dla każdego  $n \in \mathbb{N}$  zachodzi

$$\Psi(n) : \quad 1^2 + 2^2 + \dots + n^2 = \frac{n(n+1)(2n+1)}{6}$$



- Np. Pokazać, że dla każdego  $n \in \mathbb{N}$  zachodzi

$$\Psi(n) : \boxed{1^2 + 2^2 + \dots + n^2 = \frac{n(n+1)(2n+1)}{6}}$$

- Indukcja pełna:

$$\Psi(1) \quad \text{oraz} \quad \forall_{n \geq 1} [\Psi(n) \implies \Psi(n+1)]$$



- Np. Pokazać, że dla każdego  $n \in \mathbb{N}$  zachodzi

$$\Psi(n) : \boxed{1^2 + 2^2 + \dots + n^2 = \frac{n(n+1)(2n+1)}{6}}$$

- Indukcja pełna:

$$\Psi(1) \quad \text{oraz} \quad \forall_{n \geq 1} [\Psi(n) \implies \Psi(n+1)]$$

- Indukcja niepełna: czy wystarczy sprawdzić, np.

$$\Psi(1), \Psi(2), \Psi(3), \Psi(4)?$$

Podejście indukcyjne:

Wnioskowanie na podstawie skończonego zbioru obserwacji



## Podjęcie indukcyjne:

Wnioskowanie na podstawie skończonego zbioru obserwacji

**Jakie prawa rządzą procesem indukcyjnego uczenia się pojęć?**



Podjęcie indukcyjne:

Wnioskowanie na podstawie skończonego zbioru obserwacji

**Jakie prawa rządzą procesem indukcyjnego uczenia się pojęć?**

Szukamy teorii obejmującej zagadnienia:

- Szansy na skuteczne wyuczanie się pojęć;





## Podjęcie indukcyjne:

Wnioskowanie na podstawie skończonego zbioru obserwacji

**Jakie prawa rządzą procesem indukcyjnego uczenia się pojęć?**

Szukamy teorii obejmującej zagadnienia:

- Szansy na skuteczne wyuczanie się pojęć;
- Niezbędnej liczby przykładów treningowych;



## Podjęcie indukcyjne:

Wnioskowanie na podstawie skończonego zbioru obserwacji

**Jakie prawa rządzą procesem indukcyjnego uczenia się pojęć?**

Szukamy teorii obejmującej zagadnienia:

- Szansy na skuteczne wyuczanie się pojęć;
- Niezbędnej liczby przykładów treningowych;
- Złożoności przestrzeni hipotez;

## Podjęcie indukcyjne:

Wnioskowanie na podstawie skończonego zbioru obserwacji

**Jakie prawa rządzą procesem indukcyjnego uczenia się pojęć?**

Szukamy teorii obejmującej zagadnienia:

- Szansy na skuteczne wyuczanie się pojęć;
- Niezbędnej liczby przykładów treningowych;
- Złożoności przestrzeni hipotez;
- Jakości aproksymacji;

## Podjęcie indukcyjne:

Wnioskowanie na podstawie skończonego zbioru obserwacji

**Jakie prawa rządzą procesem indukcyjnego uczenia się pojęć?**

Szukamy teorii obejmującej zagadnienia:

- Szansy na skuteczne wyuczanie się pojęć;
- Niezbędnej liczby przykładów treningowych;
- Złożoności przestrzeni hipotez;
- Jakości aproksymacji;
- Metod reprezentacji danych treningowych;

# Ogólny model uczenia indukcyjnego

---

- Niech



# Ogólny model uczenia indukcyjnego

---

- Niech
  - $\mathcal{X}$  – (skończony lub nieskończony) zbiór obiektów;



# Ogólny model uczenia indukcyjnego

---

## ■ Niech

- $\mathcal{X}$  – (skończony lub nieskończony) zbiór obiektów;
- $\mathbb{C}$  – klasa pojęć w  $\mathcal{X}$ , tj.  $\mathbb{C} = \{f : \mathcal{X} \rightarrow \{0, 1\}\}$



## ■ Niech

- $\mathcal{X}$  – (skończony lub nieskończony) zbiór obiektów;
- $\mathbb{C}$  – klasa pojęć w  $\mathcal{X}$ , tj.  $\mathbb{C} = \{f : \mathcal{X} \rightarrow \{0, 1\}\}$
- $c \in \mathbb{C}$  – pojęcie docelowe lub funkcja celu;





# Ogólny model uczenia indukcyjnego

---

## ■ Niech

- $\mathcal{X}$  – (skończony lub nieskończony) zbiór obiektów;
- $\mathbb{C}$  – klasa pojęć w  $\mathcal{X}$ , tj.  $\mathbb{C} = \{f : \mathcal{X} \rightarrow \{0, 1\}\}$
- $c \in \mathbb{C}$  – pojęcie docelowe lub funkcja celu;

## ■ Dane są



## ■ Niech

- $\mathcal{X}$  – (skończony lub nieskończony) zbiór obiektów;
- $\mathbb{C}$  – klasa pojęć w  $\mathcal{X}$ , tj.  $\mathbb{C} = \{f : \mathcal{X} \rightarrow \{0, 1\}\}$
- $c \in \mathbb{C}$  – pojęcie docelowe lub funkcja celu;

## ■ Dane są

- skończona próbka etykietowanych obiektów:

$$D = \{\langle x_1, c(x_1) \rangle, \dots, \langle x_m, c(x_m) \rangle\} \in \mathcal{S}(m, c)$$

gdzie  $x_1, \dots, x_m \in \mathcal{X}$ .



## ■ Niech

- $\mathcal{X}$  – (skończony lub nieskończony) zbiór obiektów;
- $\mathbb{C}$  – klasa pojęć w  $\mathcal{X}$ , tj.  $\mathbb{C} = \{f : \mathcal{X} \rightarrow \{0, 1\}\}$
- $c \in \mathbb{C}$  – pojęcie docelowe lub funkcja celu;

## ■ Dane są

- skończona próbka etykietowanych obiektów:

$$D = \{\langle x_1, c(x_1) \rangle, \dots, \langle x_m, c(x_m) \rangle\} \in \mathcal{S}(m, c)$$

gdzie  $x_1, \dots, x_m \in \mathcal{X}$ .

- przestrzeń hipotez  $\mathbb{H} = \{h : \mathcal{X} \rightarrow \{0, 1\}\}$ ;

## ■ Niech

- $\mathcal{X}$  – (skończony lub nieskończony) zbiór obiektów;
- $\mathbb{C}$  – klasa pojęć w  $\mathcal{X}$ , tj.  $\mathbb{C} = \{f : \mathcal{X} \rightarrow \{0, 1\}\}$
- $c \in \mathbb{C}$  – pojęcie docelowe lub funkcja celu;

## ■ Dane są

- skończona próbka etykietowanych obiektów:

$$D = \{\langle x_1, c(x_1) \rangle, \dots, \langle x_m, c(x_m) \rangle\} \in \mathcal{S}(m, c)$$

gdzie  $x_1, \dots, x_m \in \mathcal{X}$ .

- przestrzeń hipotez  $\mathbb{H} = \{h : \mathcal{X} \rightarrow \{0, 1\}\}$ ;

## ■ Szukana



## ■ Niech

- $\mathcal{X}$  – (skończony lub nieskończony) zbiór obiektów;
- $\mathbb{C}$  – klasa pojęć w  $\mathcal{X}$ , tj.  $\mathbb{C} = \{f : \mathcal{X} \rightarrow \{0, 1\}\}$
- $c \in \mathbb{C}$  – pojęcie docelowe lub funkcja celu;

## ■ Dane są

- skończona próbka etykietowanych obiektów:

$$D = \{\langle x_1, c(x_1) \rangle, \dots, \langle x_m, c(x_m) \rangle\} \in \mathcal{S}(m, c)$$

gdzie  $x_1, \dots, x_m \in \mathcal{X}$ .

- przestrzeń hipotez  $\mathbb{H} = \{h : \mathcal{X} \rightarrow \{0, 1\}\}$ ;

## ■ Szukana

- hipoteza  $h \in \mathbb{H}$  będąca dobrą aproksymacją pojęcia  $c$ .



## ■ Niech

- $\mathcal{X}$  – (skończony lub nieskończony) zbiór obiektów;
- $\mathbb{C}$  – klasa pojęć w  $\mathcal{X}$ , tj.  $\mathbb{C} = \{f : \mathcal{X} \rightarrow \{0, 1\}\}$
- $c \in \mathbb{C}$  – pojęcie docelowe lub funkcja celu;

## ■ Dane są

- skończona próbka etykietowanych obiektów:

$$D = \{\langle x_1, c(x_1) \rangle, \dots, \langle x_m, c(x_m) \rangle\} \in \mathcal{S}(m, c)$$

gdzie  $x_1, \dots, x_m \in \mathcal{X}$ .

- przestrzeń hipotez  $\mathbb{H} = \{h : \mathcal{X} \rightarrow \{0, 1\}\}$ ;

## ■ Szukana

- hipoteza  $h \in \mathbb{H}$  będąca dobrą aproksymacją pojęcia  $c$ .

## ■ Wymagane



## ■ Niech

- $\mathcal{X}$  – (skończony lub nieskończony) zbiór obiektów;
- $\mathbb{C}$  – klasa pojęć w  $\mathcal{X}$ , tj.  $\mathbb{C} = \{f : \mathcal{X} \rightarrow \{0, 1\}\}$
- $c \in \mathbb{C}$  – pojęcie docelowe lub funkcja celu;

## ■ Dane są

- skończona próbka etykietowanych obiektów:

$$D = \{\langle x_1, c(x_1) \rangle, \dots, \langle x_m, c(x_m) \rangle\} \in \mathcal{S}(m, c)$$

gdzie  $x_1, \dots, x_m \in \mathcal{X}$ .

- przestrzeń hipotez  $\mathbb{H} = \{h : \mathcal{X} \rightarrow \{0, 1\}\}$ ;

## ■ Szukana

- hipoteza  $h \in \mathbb{H}$  będąca dobrą aproksymacją pojęcia  $c$ .

## ■ Wymagane

- dobra jakość aproksymacji



## ■ Niech

- $\mathcal{X}$  – (skończony lub nieskończony) zbiór obiektów;
- $\mathbb{C}$  – klasa pojęć w  $\mathcal{X}$ , tj.  $\mathbb{C} = \{f : \mathcal{X} \rightarrow \{0, 1\}\}$
- $c \in \mathbb{C}$  – pojęcie docelowe lub funkcja celu;

## ■ Dane są

- skończona próbka etykietowanych obiektów:

$$D = \{\langle x_1, c(x_1) \rangle, \dots, \langle x_m, c(x_m) \rangle\} \in \mathcal{S}(m, c)$$

gdzie  $x_1, \dots, x_m \in \mathcal{X}$ .

- przestrzeń hipotez  $\mathbb{H} = \{h : \mathcal{X} \rightarrow \{0, 1\}\}$ ;

## ■ Szukana

- hipoteza  $h \in \mathbb{H}$  będąca dobrą aproksymacją pojęcia  $c$ .

## ■ Wymagane

- dobra jakość aproksymacji
- szybki czas wyuczania.



# Przykład

---



# Przykład

---



# Przykład

---

- Pojęcie: "człowieka o średniej budowie ciała".



# Przykład

---

- Pojęcie: "człowieka o średniej budowie ciała".
- Dane – czyli osoby – są reprezentowane przez ich wagę( $kg$ ) i wzrost( $cm$ ) i są etykietowane przez  $+$  i  $-$ .



- Pojęcie: "człowieka o średniej budowie ciała".
- Dane – czyli osoby – są reprezentowane przez ich wagę( $kg$ ) i wzrost( $cm$ ) i są etykietowane przez  $+$  i  $-$ .
- Dodatkowa wiedza: szukane pojęcie można wyrazić za pomocą PROSTOKĄTA

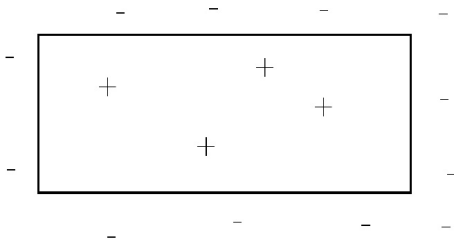


- Pojęcie: "człowieka o średniej budowie ciała".
- Dane – czyli osoby – są reprezentowane przez ich wagę( $kg$ ) i wzrost( $cm$ ) i są etykietowane przez  $+$  i  $-$ .
- Dodatkowa wiedza: szukane pojęcie można wyrazić za pomocą PROSTOKĄTA



# Przykład

- Pojęcie: "człowieka o średniej budowie ciała".
- Dane – czyli osoby – są reprezentowane przez ich wagę( $kg$ ) i wzrost( $cm$ ) i są etykietowane przez  $+$  i  $-$ .
- Dodatkowa wiedza: szukane pojęcie można wyrazić za pomocą PROSTOKĄTA

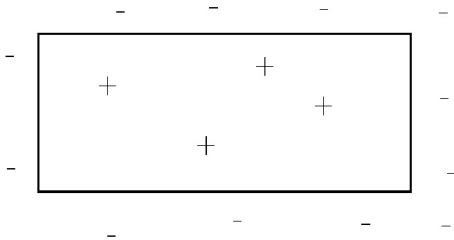


Uczenie prostokąta

# Przykład



- Pojęcie: "człowieka o średniej budowie ciała".
- Dane – czyli osoby – są reprezentowane przez ich wagę( $kg$ ) i wzrost( $cm$ ) i są etykietowane przez  $+$  i  $-$ .
- Dodatkowa wiedza: szukane pojęcie można wyrazić za pomocą PROSTOKĄTA



## Uczenie prostokąta

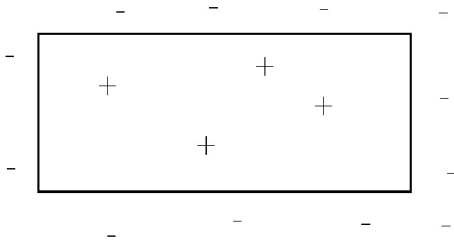
- $\mathcal{X} = \mathbb{R}^2$ ;



# Przykład



- Pojęcie: "człowieka o średniej budowie ciała".
- Dane – czyli osoby – są reprezentowane przez ich wagę( $kg$ ) i wzrost( $cm$ ) i są etykietowane przez  $+$  i  $-$ .
- Dodatkowa wiedza: szukane pojęcie można wyrazić za pomocą PROSTOKĄTA



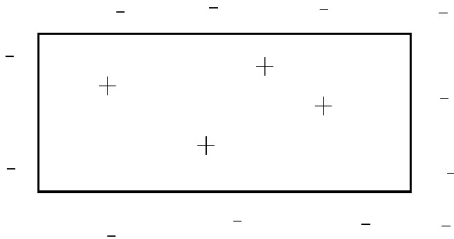
## Uczenie prostokąta

- $\mathcal{X} = \mathbb{R}^2$ ;
- $\mathbb{C} = \mathbb{H} =$  zbiór prostokątów;

# Przykład



- Pojęcie: "człowieka o średniej budowie ciała".
- Dane – czyli osoby – są reprezentowane przez ich wagę( $kg$ ) i wzrost( $cm$ ) i są etykietowane przez  $+$  i  $-$ .
- Dodatkowa wiedza: szukane pojęcie można wyrazić za pomocą PROSTOKĄTA

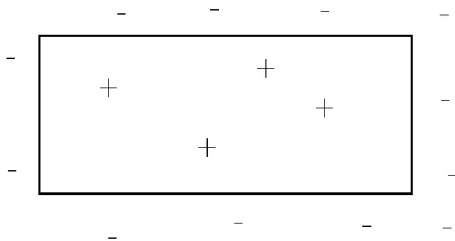


## Uczenie prostokąta

- $\mathcal{X} = \mathbb{R}^2$ ;
- $\mathbb{C} = \mathbb{H} =$  zbiór prostokątów;
- Przykład zbioru treningowego  
 $((84, 184), +)$ ,  
 $((70, 170), +)$ ,  
 $((75, 163), -)$ ,  
 $((80, 180), +)$ ,  
 $((81, 195), -)$ ,  
 $((63, 191), -)$ ,  
 $((77, 187), -)$ ,  
 $((68, 168), +)$

# Przykład

- Pojęcie: "człowieka o średniej budowie ciała".
- Dane – czyli osoby – są reprezentowane przez ich wagę(*kg*) i wzrost(*cm*) i są etykietowane przez + i -.
- Dodatkowa wiedza: szukane pojęcie można wyrazić za pomocą PROSTOKĄTA



## Uczenie prostokąta

- $\mathcal{X} = \mathbb{R}^2$ ;
- $\mathbb{C} = \mathbb{H} =$  zbiór prostokątów;
- Przykład zbioru treningowego  
((84, 184), +),  
((70, 170), +),  
((75, 163), -),  
((80, 180), +),  
((81, 195), -),  
((63, 191), -),  
((77, 187), -),  
((68, 168), +)  
■ ((79, 183, ?))

- **Uczenie półosi (lub dyskretyzacji):**

$$\mathcal{X} = \mathfrak{R}; \quad \mathbb{C} = \mathbb{H} = \{[\lambda, \infty) : \alpha \in \mathfrak{R}\}$$



- **Uczenie półosi (lub dyskretyzacji):**

$$\mathcal{X} = \mathbb{R}; \quad \mathbb{C} = \mathbb{H} = \{[\lambda, \infty) : \lambda \in \mathbb{R}\}$$

- **Uczenie hiperpłaszczyzny:**

$$\mathcal{X} = \mathbb{R}^n; \quad \mathbb{H} = \{f_{w_0, w_1, \dots, w_n} : \mathbb{R}^n \rightarrow \{0, 1\}\}$$

gdzie  $f_{w_0, \dots, w_n}(x_1, \dots, x_n) = \text{sgn}(w_0 + w_1 x_1 + \dots + w_n x_n)$ .



- **Uczenie półosi (lub dyskretyzacji):**

$$\mathcal{X} = \mathbb{R}; \quad \mathbb{C} = \mathbb{H} = \{[\lambda, \infty) : \lambda \in \mathbb{R}\}$$

- **Uczenie hiperpłaszczyzny:**

$$\mathcal{X} = \mathbb{R}^n; \quad \mathbb{H} = \{f_{w_0, w_1, \dots, w_n} : \mathbb{R}^n \rightarrow \{0, 1\}\}$$

gdzie  $f_{w_0, \dots, w_n}(x_1, \dots, x_n) = \text{sgn}(w_0 + w_1x_1 + \dots + w_nx_n)$ .

- **Uczenie jednomianów Boolowskich:**

$$\mathcal{X} = \{0, 1\}^n; \quad c : \{0, 1\}^n \rightarrow \{0, 1\};$$

$\mathbb{H} = \mathbf{M}_n =$  zbiór jednomianów Boolowskich o  $n$  zmiennych.

Błąd rzeczywisty



## Błąd rzeczywisty

- $\Omega = (\mathcal{X}, \mu)$  – przestrzeń probabilistyczna na  $\mathcal{X}$ ;





## Błąd rzeczywisty

- $\Omega = (\mathcal{X}, \mu)$  – przestrzeń probabilistyczna na  $\mathcal{X}$ ;
- Błąd hipotezy  $h \in \mathbb{H}$  względem funkcji celu  $c$ :

$$er_{\Omega}(h, c) = er_{\Omega}^c(h) = \mu(\mathcal{X}_{h \neq c})$$

gdzie  $\mathcal{X}_{h \neq c} = \{x \in \mathcal{X} : h(x) \neq c(x)\}$ .

## Błąd rzeczywisty

- $\Omega = (\mathcal{X}, \mu)$  – przestrzeń probabilistyczna na  $\mathcal{X}$ ;
- Błąd hipotezy  $h \in \mathbb{H}$  względem funkcji celu  $c$ :

$$er_{\Omega}(h, c) = er_{\Omega}^c(h) = \mu(\mathcal{X}_{h \neq c})$$

gdzie  $\mathcal{X}_{h \neq c} = \{x \in \mathcal{X} : h(x) \neq c(x)\}$ .

**Statystyka:** Jeśli przykłady z  $D$  są wybrane zgodnie z miarą prawdopodobieństwa  $\mu$  w sposób niezależny oraz  $|D| \geq 30$ , to

## Błąd rzeczywisty

- $\Omega = (\mathcal{X}, \mu)$  – przestrzeń probabilistyczna na  $\mathcal{X}$ ;
- Błąd hipotezy  $h \in \mathbb{H}$  względem funkcji celu  $c$ :

$$er_{\Omega}(h, c) = er_{\Omega}^c(h) = \mu(\mathcal{X}_{h \neq c})$$

gdzie  $\mathcal{X}_{h \neq c} = \{x \in \mathcal{X} : h(x) \neq c(x)\}$ .

**Statystyka:** Jeśli przykłady z  $D$  są wybrane zgodnie z miarą prawdopodobieństwa  $\mu$  w sposób niezależny oraz  $|D| \geq 30$ , to

- $er_{\Omega}^c(h) \approx er_D^c(h) = \frac{|D \cap \mathcal{X}_{h \neq c}|}{|D|}$ ,

## Błąd rzeczywisty

- $\Omega = (\mathcal{X}, \mu)$  – przestrzeń probabilistyczna na  $\mathcal{X}$ ;
- Błąd hipotezy  $h \in \mathbb{H}$  względem funkcji celu  $c$ :

$$er_{\Omega}(h, c) = er_{\Omega}^c(h) = \mu(\mathcal{X}_{h \neq c})$$

gdzie  $\mathcal{X}_{h \neq c} = \{x \in \mathcal{X} : h(x) \neq c(x)\}$ .

**Statystyka:** Jeśli przykłady z  $D$  są wybrane zgodnie z miarą prawdopodobieństwa  $\mu$  w sposób niezależny oraz  $|D| \geq 30$ , to

- $er_{\Omega}^c(h) \approx er_D^c(h) = \frac{|D \cap \mathcal{X}_{h \neq c}|}{|D|}$ ,
- z prawdopodobieństwem  $(1 - \varepsilon)$

$$|er_{\Omega}^c - er_D^c| \leq s_{\frac{\varepsilon}{2}} \cdot \sqrt{\frac{er_D^c(1 - er_D^c)}{|D|}}$$



- 1 Wstęp do komputerowego uczenia się pojęć
- 2 Model PAC (probably approximately correct)
- 3 Wyuczalność klasy pojęć
- 4 Wymiar Vapnika Chervonenkisa (VC dimension)
- 5 Podstawowe twierdzenia teorii uczenia się
- 6 Appendix: „Nie ma nic za darmo” czyli “Non Free Lunch Theorem”

# Model uczenia się PAC

---

Idea modelu PAC (Probably Approximately Correct):

Określenie warunków, przy których uczeń (algorytm uczenia się) z „dużym prawdopodobieństwem” znajdzie „dobrą hipotezę” na podstawie danych  $D$ .



# Model uczenia się PAC

Idea modelu PAC (Probably Approximately Correct):

Określenie warunków, przy których uczeń (algorytm uczenia się) z „dużym prawdopodobieństwem” znajdzie „dobrą hipotezę” na podstawie danych  $D$ .

PAC-owy uczeń

Niech  $\mathcal{L}$  będzie algorytmem uczenia się, jeśli

dla każdych  $0 < \varepsilon, \delta < 1$ , istnieje liczba  $m_0 = m_0(\varepsilon, \delta)$  taka, że dla dowolnego pojęcia  $c \in \mathbb{C}$ , dla dowolnego rozkładu  $\Omega$  na  $X$  i dla  $m > m_0$  mamy

$$\mu^m \{D \in \mathcal{S}(m, c) : er_{\Omega}(\mathcal{L}(D)) < \varepsilon\} > 1 - \delta$$

Wówczas mówimy w skrócie, że  $\mathcal{L}$  jest **PAC** dla klasy  $\mathbb{C}$  (“prawdopodobnie aproksymacyjnie poprawny”).

$\varepsilon$  = dopuszczalny poziom błędów;  $(1 - \delta)$  = poziom zaufania.

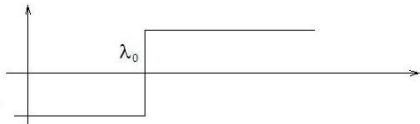
- $\mathbb{H} = \mathbb{C} = \{f_\lambda : \mathfrak{R} \rightarrow \{0, 1\} : f_\lambda(x) = 1 \Leftrightarrow x \geq \lambda\}$





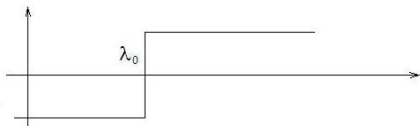
## Przykład problemu dyskretyzacji

- $\mathbb{H} = \mathbb{C} = \{f_\lambda : \mathfrak{R} \rightarrow \{0, 1\} : f_\lambda(x) = 1 \Leftrightarrow x \geq \lambda\}$
- $c = f_{\lambda_0}$



## Przykład problemu dyskretyzacji

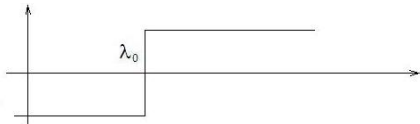
- $\mathbb{H} = \mathbb{C} = \{f_\lambda : \mathfrak{R} \rightarrow \{0, 1\} : f_\lambda(x) = 1 \Leftrightarrow x \geq \lambda\}$
- $c = f_{\lambda_0}$



- znaleźć  $\lambda_0$  na podstawie losowo wygenerowanych przykładów  $D = \{\langle x_1, f_{\lambda_0}(x_1) \rangle, \dots, \langle x_m, f_{\lambda_0}(x_m) \rangle\}$

## Przykład problemu dyskretyzacji

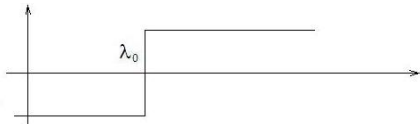
- $\mathbb{H} = \mathbb{C} = \{f_\lambda : \mathfrak{R} \rightarrow \{0, 1\} : f_\lambda(x) = 1 \Leftrightarrow x \geq \lambda\}$
- $c = f_{\lambda_0}$



- znaleźć  $\lambda_0$  na podstawie losowo wygenerowanych przykładów  $D = \{\langle x_1, f_{\lambda_0}(x_1) \rangle, \dots, \langle x_m, f_{\lambda_0}(x_m) \rangle\}$

### Algorytm:

- $\mathbb{H} = \mathbb{C} = \{f_\lambda : \mathfrak{R} \rightarrow \{0, 1\} : f_\lambda(x) = 1 \Leftrightarrow x \geq \lambda\}$
- $c = f_{\lambda_0}$

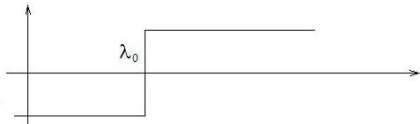


- znaleźć  $\lambda_0$  na podstawie losowo wygenerowanych przykładów  $D = \{\langle x_1, f_{\lambda_0}(x_1) \rangle, \dots, \langle x_m, f_{\lambda_0}(x_m) \rangle\}$

### Algorytm:

- 1 Set  $\lambda^* := \min_{i \in \{1, \dots, m\}} \{x_i : f_{\lambda_0}(x_i) = 1\}$ ;

- $\mathbb{H} = \mathbb{C} = \{f_\lambda : \mathfrak{R} \rightarrow \{0, 1\} : f_\lambda(x) = 1 \Leftrightarrow x \geq \lambda\}$
- $c = f_{\lambda_0}$

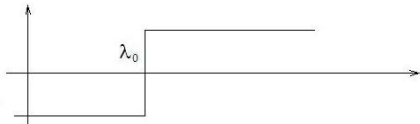


- znaleźć  $\lambda_0$  na podstawie losowo wygenerowanych przykładów  $D = \{\langle x_1, f_{\lambda_0}(x_1) \rangle, \dots, \langle x_m, f_{\lambda_0}(x_m) \rangle\}$

## Algorytm:

- 1 Set  $\lambda^* := \min_{i \in \{1, \dots, m\}} \{x_i : f_{\lambda_0}(x_i) = 1\}$ ;
- 2  $\mathcal{L}(D) := f_{\lambda^*}$ ;

- $\mathbb{H} = \mathbb{C} = \{f_\lambda : \mathfrak{R} \rightarrow \{0, 1\} : f_\lambda(x) = 1 \Leftrightarrow x \geq \lambda\}$
- $c = f_{\lambda_0}$



- znaleźć  $\lambda_0$  na podstawie losowo wygenerowanych przykładów  $D = \{\langle x_1, f_{\lambda_0}(x_1) \rangle, \dots, \langle x_m, f_{\lambda_0}(x_m) \rangle\}$

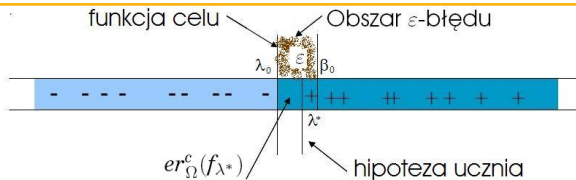
## Algorytm:

- 1 Set  $\lambda^* := \min_{i \in \{1, \dots, m\}} \{x_i : f_{\lambda_0}(x_i) = 1\}$ ;
- 2  $\mathcal{L}(D) := f_{\lambda^*}$ ;

## Twierdzenie:

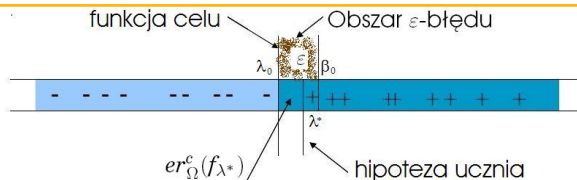
Powyższy algorytm jest PAC

# Dowód



■  $er_{\Omega}^c(f_{\lambda^*}) = \mu([\lambda_0, \lambda^*])$ .

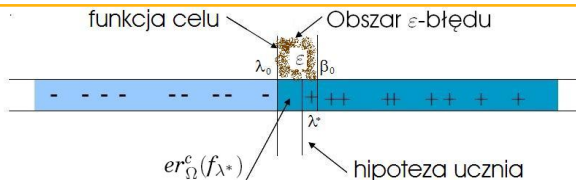
# Dowód



- $er_{\Omega}^c(f_{\lambda^*}) = \mu([\lambda_0, \lambda^*])$ .
- Niech  $\beta_0 = \sup\{\beta : \mu([\lambda_0, \beta]) < \varepsilon\}$ .



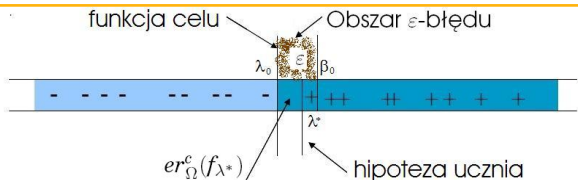
# Dowód



- $er_{\Omega}^c(f_{\lambda^*}) = \mu([\lambda_0, \lambda^*])$ .
- Niech  $\beta_0 = \sup\{\beta : \mu([\lambda_0, \beta]) < \varepsilon\}$ .
- Wówczas  $er_{\Omega}^c(f_{\lambda^*}) \geq \varepsilon \Leftrightarrow \forall x_i \in D : x_i \notin [\lambda_0, \beta_0]$ ;



# Dowód



- $er_{\Omega}^c(f_{\lambda^*}) = \mu([\lambda_0, \lambda^*])$ .
- Niech  $\beta_0 = \sup\{\beta : \mu([\lambda_0, \beta]) < \varepsilon\}$ .
- Wówczas  $er_{\Omega}^c(f_{\lambda^*}) \geq \varepsilon \Leftrightarrow \forall x_i \in D : x_i \notin [\lambda_0, \beta_0]$ ;
- Stąd

$$\mu^m \{(x_1, \dots, x_m) : \forall x_i \in D : x_i \notin [\lambda_0, \beta_0]\} \leq (1 - \varepsilon)^m$$
$$\mu^m \{D \in \mathcal{S}(m, f_{\lambda_0}) : er_{\Omega}(f_{\lambda^*}) \leq \varepsilon\} \geq 1 - (1 - \varepsilon)^m$$

- Aby to prawdopodobieństwo było  $> 1 - \delta$ , wystarczy przyjąć  $m \geq m_0 = \left\lceil \frac{1}{\varepsilon} \ln \frac{1}{\delta} \right\rceil$



- Niech  $\Omega$  będzie rozkładem dyskretnym zdefiniowanym przez  $\mu_1 = \mu(x_1), \dots, \mu_n = \mu(x_n)$  – dla pewnych  $x_1, \dots, x_n \in X$  – takich, że  $\mu_1 + \dots + \mu_n = 1$ . Niech  $\varepsilon_{\min} = \min_i \mu_i$ .



- Niech  $\Omega$  będzie rozkładem dyskretnym zdefiniowanym przez  $\mu_1 = \mu(x_1), \dots, \mu_n = \mu(x_n)$  – dla pewnych  $x_1, \dots, x_n \in X$  – takich, że  $\mu_1 + \dots + \mu_n = 1$ . Niech  $\varepsilon_{\min} = \min_i \mu_i$ .
- Jeśli  $\mathcal{L}$  jest PAC, i jeśli  $\varepsilon \leq \varepsilon_{\min}$  to warunek  $er_{\Omega}^c(L(D)) < \varepsilon$  jest równoważny z  $er_{\Omega}^c(L(D)) = 0$ . Stąd dla każdego  $\delta$ , istnieje  $m_0 = m_0(\varepsilon_{\min}, \delta)$  taka, że dla dowolnego  $c \in \mathbb{C}$  i  $\Omega$

$$m > m_0 \Rightarrow \mu^m \{D \in \mathcal{S}(m, t) | er_{\Omega}(L(D)) = 0\} > 1 - \delta$$



- Niech  $\Omega$  będzie rozkładem dyskretnym zdefiniowanym przez  $\mu_1 = \mu(x_1), \dots, \mu_n = \mu(x_n)$  – dla pewnych  $x_1, \dots, x_n \in X$  – takich, że  $\mu_1 + \dots + \mu_n = 1$ . Niech  $\varepsilon_{\min} = \min_i \mu_i$ .
- Jeśli  $\mathcal{L}$  jest PAC, i jeśli  $\varepsilon \leq \varepsilon_{\min}$  to warunek  $er_{\Omega}^c(L(D)) < \varepsilon$  jest równoważny z  $er_{\Omega}^c(L(D)) = 0$ . Stąd dla każdego  $\delta$ , istnieje  $m_0 = m_0(\varepsilon_{\min}, \delta)$  taka, że dla dowolnego  $c \in \mathbb{C}$  i  $\Omega$

$$m > m_0 \Rightarrow \mu^m \{D \in \mathcal{S}(m, t) | er_{\Omega}(L(D)) = 0\} > 1 - \delta$$

- Wówczas mówimy, że prawdopodobnie  $\mathcal{L}$  jest dokładnym algorytmem (*jest PEC – probably exactly correct*)

# Plan wykładu

---



- 1 Wstęp do komputerowego uczenia się pojęć
- 2 Model PAC (probably approximately correct)
- 3 Wyuczalność klasy pojęć**
- 4 Wymiar Vapnika Chervonenkisa (VC dimension)
- 5 Podstawowe twierdzenia teorii uczenia się
- 6 Appendix: „Nie ma nic za darmo” czyli “Non Free Lunch Theorem”

- Niech  $D = \{\langle x_1, c(x_1) \rangle, \dots, \langle x_m, c(x_m) \rangle\}$  i niech

$$\mathbb{H}^c(D) = \{h \in \mathbb{H} : h|D = c|D\}$$

zbiór hipotez zgodnych z  $c$  na próbce  $D$ .





- Niech  $D = \{\langle x_1, c(x_1) \rangle, \dots, \langle x_m, c(x_m) \rangle\}$  i niech

$$\mathbb{H}^c(D) = \{h \in \mathbb{H} : h|D = c|D\}$$

zbiór hipotez zgodnych z  $c$  na próbce  $D$ .

- $\mathbb{B}_\varepsilon^c = \{h \in \mathbb{H} : er_\Omega(h) \geq \varepsilon\}$  – zbiór  $\varepsilon$ -złych hipotez

- Niech  $D = \{\langle x_1, c(x_1) \rangle, \dots, \langle x_m, c(x_m) \rangle\}$  i niech

$$\mathbb{H}^c(D) = \{h \in \mathbb{H} : h|D = c|D\}$$

zbiór hipotez zgodnych z  $c$  na próbce  $D$ .

- $\mathbb{B}_\varepsilon^c = \{h \in \mathbb{H} : er_\Omega(h) \geq \varepsilon\}$  – zbiór  $\varepsilon$ -złych hipotez

- Niech  $D = \{\langle x_1, c(x_1) \rangle, \dots, \langle x_m, c(x_m) \rangle\}$  i niech

$$\mathbb{H}^c(D) = \{h \in \mathbb{H} : h|D = c|D\}$$

zbiór hipotez zgodnych z  $c$  na próbce  $D$ .

- $\mathbb{B}_\varepsilon^c = \{h \in \mathbb{H} : er_\Omega(h) \geq \varepsilon\}$  – zbiór  $\varepsilon$ -złych hipotez

## Definicja: Potencjalna wyuczalność

Mówimy, że  $\mathbb{C}$  jest potencjalnie wyuczalna za pomocą  $\mathbb{H}$ , jeśli dla każdego rozkładu  $\Omega$  na  $\mathcal{X}$  i dowolnego pojęcia  $c \in \mathbb{C}$  oraz dla dowolnych  $0 < \varepsilon, \delta < 1$  istnieje  $m_0 = m_0(\varepsilon, \delta)$  takie, że

$$m \geq m_0 \Rightarrow \mu^m \{D \in \mathcal{S}(m, c) : \mathbb{H}^c(D) \cap \mathbb{B}_\varepsilon^c = \emptyset\} > 1 - \delta$$

# Potencjalna wyuczalność

---

Algorytm  $\mathcal{L}$  nazywamy **niesprzecznym** jeśli  $\mathcal{L}(D) \in \mathbb{H}^c(D)$  dla każdego zbioru  $D$ .

## Twierdzenie

W przestrzeni potencjalnie wyuczalnej, każdy wzorowy uczeń (niespreczny algorytm) jest PAC-owy.



# Potencjalna wyuczalność

Algorytm  $\mathfrak{L}$  nazywamy **niesprzecznym** jeśli  $\mathfrak{L}(D) \in \mathbb{H}^c(D)$  dla każdego zbioru  $D$ .

## Twierdzenie

W przestrzeni potencjalnie wyuczalnej, każdy wzorowy uczeń (niespreczny algorytm) jest PAC-owy.

## Twierdzenie (Haussler, 1988)

Jeśli  $\mathbb{C} = \mathbb{H}$  i  $|\mathbb{C}| < \infty$ , to  $\mathbb{C}$  jest potencjalnie wyuczalna.

Dowód: Niech  $h \in \mathbb{B}_\varepsilon$  (tzn.  $er_\Omega(h) \geq \varepsilon$ ). Wówczas

$$\mu^m \{D \in \mathcal{S}(m, c) : er_D(h) = 0\} \leq (1 - \varepsilon)^m$$

$$\Rightarrow \mu^m \{D : \mathbb{H}^c(D) \cap \mathbb{B}_\varepsilon \neq \emptyset\} \leq |\mathbb{B}_\varepsilon| (1 - \varepsilon)^m \leq |\mathbb{H}| (1 - \varepsilon)^m$$

Aby  $|\mathbb{H}| (1 - \varepsilon)^m < \delta$  wystarczy wybrać  $m \geq m_0 = \left\lceil \frac{1}{\varepsilon} \ln \frac{|\mathbb{H}|}{\delta} \right\rceil$

# Plan wykładu

---



- 1 Wstęp do komputerowego uczenia się pojęć
- 2 Model PAC (probably approximately correct)
- 3 Wyuczalność klasy pojęć
- 4 Wymiar Vapnika Chervonenkisa (VC dimension)
- 5 Podstawowe twierdzenia teorii uczenia się
- 6 Appendix: „Nie ma nic za darmo” czyli “Non Free Lunch Theorem”

- Niech  $\vec{\mathbf{x}} = \langle x_1, \dots, x_m \rangle \in \mathcal{X}^m$ . Niech

$$\Pi_{\mathbb{H}}(\vec{\mathbf{x}}) = |\{\langle h(x_1), \dots, h(x_m) \rangle \in \{0, 1\}^m : h \in \mathbb{H}\}|$$



- Niech  $\vec{\mathbf{x}} = \langle x_1, \dots, x_m \rangle \in \mathcal{X}^m$ . Niech

$$\Pi_{\mathbb{H}}(\vec{\mathbf{x}}) = |\{\langle h(x_1), \dots, h(x_m) \rangle \in \{0, 1\}^m : h \in \mathbb{H}\}|$$

- $\Pi_{\mathbb{H}}(\vec{\mathbf{x}})$  jest liczbą podziałów zbioru elementów  $\vec{\mathbf{x}}$  wyznaczonych przez  $\mathbb{H}$ . Mamy  $\Pi_{\mathbb{H}}(\vec{\mathbf{x}}) \leq 2^m$ .

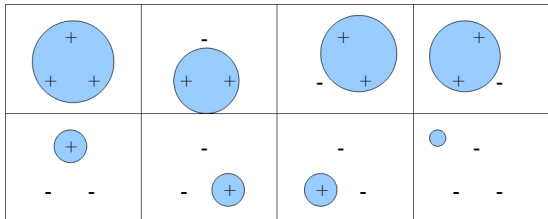




- Niech  $\vec{\mathbf{x}} = \langle x_1, \dots, x_m \rangle \in \mathcal{X}^m$ . Niech

$$\Pi_{\mathbb{H}}(\vec{\mathbf{x}}) = |\{\langle h(x_1), \dots, h(x_m) \rangle \in \{0, 1\}^m : h \in \mathbb{H}\}|$$

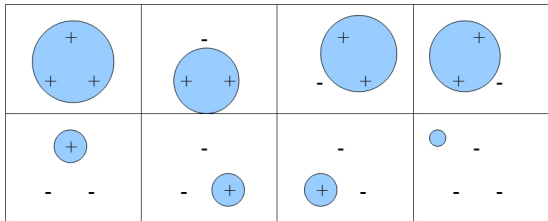
- $\Pi_{\mathbb{H}}(\vec{\mathbf{x}})$  jest liczbą podziałów zbioru elementów  $\vec{\mathbf{x}}$  wyznaczonych przez  $\mathbb{H}$ . Mamy  $\Pi_{\mathbb{H}}(\vec{\mathbf{x}}) \leq 2^m$ .
- Gdy  $\Pi_{\mathbb{H}}(\vec{\mathbf{x}}) = 2^m$ , mówimy, że  $\mathbb{H}$  **rozbija**  $\mathbf{x}$ .



- Niech  $\vec{\mathbf{x}} = \langle x_1, \dots, x_m \rangle \in \mathcal{X}^m$ . Niech

$$\Pi_{\mathbb{H}}(\vec{\mathbf{x}}) = |\{\langle h(x_1), \dots, h(x_m) \rangle \in \{0, 1\}^m : h \in \mathbb{H}\}|$$

- $\Pi_{\mathbb{H}}(\vec{\mathbf{x}})$  jest liczbą podziałów zbioru elementów  $\vec{\mathbf{x}}$  wyznaczonych przez  $\mathbb{H}$ . Mamy  $\Pi_{\mathbb{H}}(\vec{\mathbf{x}}) \leq 2^m$ .
- Gdy  $\Pi_{\mathbb{H}}(\vec{\mathbf{x}}) = 2^m$ , mówimy, że  $\mathbb{H}$  **rozbija**  $\mathbf{x}$ .

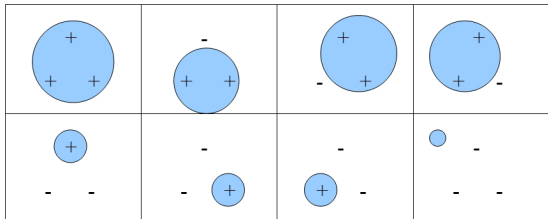


- Niech  $\Pi_{\mathbb{H}}(m) = \max_{\vec{\mathbf{x}} \in \mathcal{X}^m} \Pi_{\mathbb{H}}(\vec{\mathbf{x}})$

- Niech  $\vec{\mathbf{x}} = \langle x_1, \dots, x_m \rangle \in \mathcal{X}^m$ . Niech

$$\Pi_{\mathbb{H}}(\vec{\mathbf{x}}) = |\{\langle h(x_1), \dots, h(x_m) \rangle \in \{0, 1\}^m : h \in \mathbb{H}\}|$$

- $\Pi_{\mathbb{H}}(\vec{\mathbf{x}})$  jest liczbą podziałów zbioru elementów  $\vec{\mathbf{x}}$  wyznaczonych przez  $\mathbb{H}$ . Mamy  $\Pi_{\mathbb{H}}(\vec{\mathbf{x}}) \leq 2^m$ .
- Gdy  $\Pi_{\mathbb{H}}(\vec{\mathbf{x}}) = 2^m$ , mówimy, że  $\mathbb{H}$  **rozbija**  $\mathbf{x}$ .



- Niech  $\Pi_{\mathbb{H}}(m) = \max_{\vec{\mathbf{x}} \in \mathcal{X}^m} \Pi_{\mathbb{H}}(\vec{\mathbf{x}})$
- Na przykład:** W przypadku klasy pojęć "półosi" postaci  $[\alpha, \infty)$  mamy  $\Pi_{\mathbb{H}}(m) = m + 1$ .

## Uwagi:

- Jeśli  $\Pi_{\mathbb{H}}(m) = 2^m$ , to istnieje pewien zbiór o mocy  $m$  taki, że  $\mathbb{H}$  może definiować każdy jego podzbiór ( $\mathbb{H}$  rozbija ten zbiór).



## Uwagi:



- Jeśli  $\Pi_{\mathbb{H}}(m) = 2^m$ , to istnieje pewien zbiór o mocy  $m$  taki, że  $\mathbb{H}$  może definiować każdy jego podzbiór ( $\mathbb{H}$  rozbija ten zbiór).
- Maksymalna wartość  $m$ , dla której  $\Pi_{\mathbb{H}}(m) = 2^m$  można uważać za siłę wyrażalności przestrzeni  $\mathbb{H}$

## Uwagi:

- Jeśli  $\Pi_{\mathbb{H}}(m) = 2^m$ , to istnieje pewien zbiór o mocy  $m$  taki, że  $\mathbb{H}$  może definiować każdy jego podzbiór ( $\mathbb{H}$  rozbija ten zbiór).
- Maksymalna wartość  $m$ , dla której  $\Pi_{\mathbb{H}}(m) = 2^m$  można uważać za siłę wyrażalności przestrzeni  $\mathbb{H}$

### Definicja: wymiar $VCdim$

Wymiarem Vapnika-Chervonenkisa przestrzeni hipotez  $\mathbb{H}$  nazywamy liczbę

$$VCdim(\mathbb{H}) = \max\{m : \Pi_{\mathbb{H}}(m) = 2^m\}$$

gdzie maksimum wynosi  $\infty$  jeśli ten zbiór jest nieograniczony.

# Przykłady wymiarów $VCdim$

---



# Przykłady wymiarów $VCdim$

---





# Przykłady wymiarów $VCdim$

---

- $H = \{\text{okręgi ...}\} \implies VC(H) = 3$



# Przykłady wymiarów $VCdim$

---

- $H = \{\text{okręgi ...}\} \implies VC(H) = 3$
- $H = \{\text{prostokąty ...}\} \implies VC(H) = 4$



# Przykłady wymiarów $VCdim$

---



- $H = \{\text{okręgi ...}\} \implies VC(H) = 3$
- $H = \{\text{prostokąty ...}\} \implies VC(H) = 4$
- $H = \{\text{funkcje progowe ...}\} \implies$   
 $VC(H) = 1$  jeśli "+" są zawsze po  
prawej stronie;  
 $VC(H) = 2$  jeśli "+" mogą być po obu  
stronach

# Przykłady wymiarów $VCdim$



- $H = \{\text{okręgi ...}\} \implies VC(H) = 3$
- $H = \{\text{prostokąty ...}\} \implies VC(H) = 4$
- $H = \{\text{funkcje progowe ...}\} \implies$   
 $VC(H) = 1$  jeśli "+" są zawsze po  
prawej stronie;  
 $VC(H) = 2$  jeśli "+" mogą być po obu  
stronach
- $H = \{\text{przedziały ...}\} \implies$   
 $VC(H) = 2$  jeśli "+" są zawsze w środku  
 $VC(H) = 3$  jeśli w środku mogą być  
zarówno "+" i "-"

# Przykłady wymiarów $VCdim$



- $H = \{\text{okręgi ...}\} \implies VC(H) = 3$
- $H = \{\text{prostokąty ...}\} \implies VC(H) = 4$
- $H = \{\text{funkcje progowe ...}\} \implies$   
 $VC(H) = 1$  jeśli "+" są zawsze po  
prawej stronie;  
 $VC(H) = 2$  jeśli "+" mogą być po obu  
stronach
- $H = \{\text{przedziały ...}\} \implies$   
 $VC(H) = 2$  jeśli "+" są zawsze w środku  
 $VC(H) = 3$  jeśli w środku mogą być  
zarówno "+" i "-"
- $H = \{\text{półpłaszczyzny w } \mathbb{R}^2 \text{ ...}\}$   
 $\implies VC(H) = 3$

# Przykłady wymiarów $VCdim$

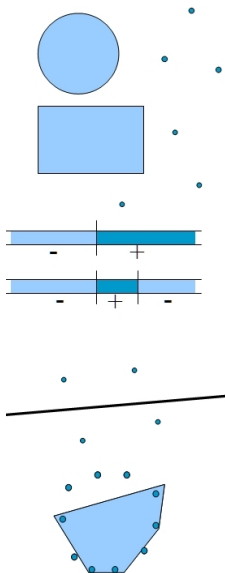


- $H = \{\text{okręgi ...}\} \implies VC(H) = 3$
- $H = \{\text{prostokąty ...}\} \implies VC(H) = 4$
- $H = \{\text{funkcje progowe ...}\} \implies$   
 $VC(H) = 1$  jeśli "+" są zawsze po  
prawej stronie;  
 $VC(H) = 2$  jeśli "+" mogą być po obu  
stronach
- $H = \{\text{przedziały ...}\} \implies$   
 $VC(H) = 2$  jeśli "+" są zawsze w środku  
 $VC(H) = 3$  jeśli w środku mogą być  
zarówno "+" i "-"
- $H = \{\text{półpłaszczyzny w } \mathbb{R}^2 \text{ ...}\}$   
 $\implies VC(H) = 3$
- czy istnieje  $H$  dla której  $VC(H) = \infty$ ?

# Przykłady wymiarów $VCdim$



- $H = \{\text{okręgi ...}\} \implies VC(H) = 3$
- $H = \{\text{prostokąty ...}\} \implies VC(H) = 4$
- $H = \{\text{funkcje progowe ...}\} \implies$   
 $VC(H) = 1$  jeśli "+" są zawsze po  
prawej stronie;  
 $VC(H) = 2$  jeśli "+" mogą być po obu  
stronach
- $H = \{\text{przedziały ...}\} \implies$   
 $VC(H) = 2$  jeśli "+" są zawsze w środku  
 $VC(H) = 3$  jeśli w środku mogą być  
zarówno "+" i "-"
- $H = \{\text{półpłaszczyzny w } \mathbb{R}^2 \text{ ...}\}$   
 $\implies VC(H) = 3$
- czy istnieje  $H$  dla której  $VC(H) = \infty$ ?



## Twierdzenie

Dla każdej liczby naturalnej  $n$ , niech  $P_n$  będzie perceptronem o  $n$  wejściach rzeczywistych. Wówczas

$$VCdim(P_n) = n + 1$$

**Dowód:**



## Twierdzenie

Dla każdej liczby naturalnej  $n$ , niech  $P_n$  będzie perceptronem o  $n$  wejściach rzeczywistych. Wówczas

$$VCdim(P_n) = n + 1$$

## Dowód:

- $VCdim(P_n) \leq n + 1$ : Wynika z Twierdzenia Radona: Dla dowolnego zbioru  $E$  zawierającego  $n + 2$  punktów w przestrzeni  $\mathbb{R}^n$  istnieje niepusty podzbiór  $S \subset E$  taki, że

$$\text{conv}(S) \cap \text{conv}(E \setminus S) \neq \emptyset$$

## Twierdzenie

Dla każdej liczby naturalnej  $n$ , niech  $P_n$  będzie perceptronem o  $n$  wejściach rzeczywistych. Wówczas

$$VCdim(P_n) = n + 1$$

## Dowód:

- $VCdim(P_n) \leq n + 1$ : Wynika z Twierdzenia Radona: Dla dowolnego zbioru  $E$  zawierającego  $n + 2$  punktów w przestrzeni  $\mathbb{R}^n$  istnieje niepusty podzbiór  $S \subset E$  taki, że

$$\text{conv}(S) \cap \text{conv}(E \setminus S) \neq \emptyset$$

- $VCdim(P_n) \geq n + 1$ : Wystarczy wybrać  $\mathbf{x} = \{0, e_1, \dots, e_n\}$  i pokazać, że każdy jego podzbiór jest definiowany przez jakiś perceptron.

## Twierdzenie



## Twierdzenie

1 Jeśli  $|\mathbb{H}| < \infty$  to  $VCdim(\mathbb{H}) \leq \log |\mathbb{H}|$ .

## Twierdzenie



1 Jeśli  $|\mathbb{H}| < \infty$  to  $VCdim(\mathbb{H}) \leq \log |\mathbb{H}|$ .

2 (**Lemat Sauer'a**) Jeśli  $VCdim(\mathbb{H}) = d \geq 0$  i  $m \geq 1$ , to

$$\Pi_{\mathbb{H}}(m) \leq 1 + \binom{m}{1} + \dots + \binom{m}{d} = \Phi(d, m)$$

## Twierdzenie



1 Jeśli  $|\mathbb{H}| < \infty$  to  $VCdim(\mathbb{H}) \leq \log |\mathbb{H}|$ .

2 (**Lemat Sauer'a**) Jeśli  $VCdim(\mathbb{H}) = d \geq 0$  i  $m \geq 1$ , to

$$\Pi_{\mathbb{H}}(m) \leq 1 + \binom{m}{1} + \dots + \binom{m}{d} = \Phi(d, m)$$

3 Wniosek:  $\Phi(d, m) \leq \left(\frac{em}{d}\right)^d \Rightarrow \Pi_{\mathbb{H}}(m) \leq \left(\frac{em}{d}\right)^d$

## Twierdzenie

1 Jeśli  $|\mathbb{H}| < \infty$  to  $VCdim(\mathbb{H}) \leq \log |\mathbb{H}|$ .

2 (**Lemat Sauer'a**) Jeśli  $VCdim(\mathbb{H}) = d \geq 0$  i  $m \geq 1$ , to

$$\Pi_{\mathbb{H}}(m) \leq 1 + \binom{m}{1} + \dots + \binom{m}{d} = \Phi(d, m)$$

3 Wniosek:  $\Phi(d, m) \leq \left(\frac{em}{d}\right)^d \Rightarrow \Pi_{\mathbb{H}}(m) \leq \left(\frac{em}{d}\right)^d$

4 Jeśli  $|\mathcal{X}| < \infty$ ,  $\mathbb{H} \subset 2^{\mathcal{X}}$  oraz  $\mathbb{H} > 1$

$$|\mathcal{X}| < \infty \implies VCdim(\mathbb{H}) > \frac{\ln |\mathbb{H}|}{1 + \ln |\mathcal{X}|}$$

# Plan wykładu

---



- 1 Wstęp do komputerowego uczenia się pojęć
- 2 Model PAC (probably approximately correct)
- 3 Wyuczalność klasy pojęć
- 4 Wymiar Vapnika Chervonenkisa (VC dimension)
- 5 Podstawowe twierdzenia teorii uczenia się
- 6 Appendix: „Nie ma nic za darmo” czyli “Non Free Lunch Theorem”





**Twierdzenie:** (Warunek konieczny)

Jeśli przestrzeń hipotez ma *nieskończony* wymiar  $VCdim$  to *nie jest potencjalnie wyuczalna*.



**Twierdzenie:** (Warunek konieczny)

Jeśli przestrzeń hipotez ma *nieskończony* wymiar  $VCdim$  to *nie jest potencjalnie wyuczalna*.

**Twierdzenie:** (fundamentalne)

Jeśli przestrzeń hipotez ma skończony wymiar  $VC$ , to jest ona potencjalnie wyuczalna.

1 Definiujemy

$$Q_m^\varepsilon = \{D \in \mathcal{S}(m, c) : H^c[D] \cap B_\varepsilon \neq \emptyset\}$$



1 Definiujemy

$$Q_m^\varepsilon = \{D \in \mathcal{S}(m, c) : H^c[D] \cap B_\varepsilon \neq \emptyset\}$$

2 Szukamy górnego ograniczenia  $f(m, \varepsilon)$  dla  $\mu^m(Q_m^\varepsilon)$ ,  
które powinny

- być niezależne od  $c \in \mathbb{C}$  i  $\mu$  (rozkład).
- dążyć do 0 przy  $m \rightarrow \infty$

1 Definiujemy

$$Q_m^\varepsilon = \{D \in \mathcal{S}(m, c) : H^c[D] \cap B_\varepsilon \neq \emptyset\}$$

2 Szukamy górnego ograniczenia  $f(m, \varepsilon)$  dla  $\mu^m(Q_m^\varepsilon)$ , które powinno

- być niezależne od  $c \in \mathbb{C}$  i  $\mu$  (rozkład).
- dążyć do 0 przy  $m \rightarrow \infty$

3 **Twierdzenie** Niech  $\mathbb{H}$  będzie przestrzenią hipotez określonych na  $X$ . Dla dowolnych  $c, \mu, \varepsilon$  (ale ustalonych) mamy

$$\mu^m(Q_m^\varepsilon) < 2\Pi_{\mathbb{H}}(2m)2^{-\varepsilon m/2}$$

o ile  $m \geq 8/\varepsilon$ .

1 Definiujemy

$$Q_m^\varepsilon = \{D \in \mathcal{S}(m, c) : H^c[D] \cap B_\varepsilon \neq \emptyset\}$$

2 Szukamy górnego ograniczenia  $f(m, \varepsilon)$  dla  $\mu^m(Q_m^\varepsilon)$ , które powinno

- być niezależne od  $c \in \mathbb{C}$  i  $\mu$  (rozkład).
- dążyć do 0 przy  $m \rightarrow \infty$

3 **Twierdzenie** Niech  $\mathbb{H}$  będzie przestrzenią hipotez określonych na  $X$ . Dla dowolnych  $c, \mu, \varepsilon$  (ale ustalonych) mamy

$$\mu^m(Q_m^\varepsilon) < 2\Pi_{\mathbb{H}}(2m)2^{-\varepsilon m/2}$$

o ile  $m \geq 8/\varepsilon$ .

4 Korzystamy z lematu Sauer'a, aby pokazać, że  $\mu^m(Q_m^\varepsilon) < \delta$  dla dostatecznie dużych  $m$ .

- Dla skończonych przestrzeni hipotez  $\mathbb{H}$  mamy

$$m_L(\mathbb{H}, \delta, \varepsilon) \leq \left\lceil \frac{1}{\varepsilon} \ln \frac{|\mathbb{H}|}{\delta} \right\rceil = \left\lceil \frac{1}{\varepsilon} (\ln |\mathbb{H}| + \ln(1/\delta)) \right\rceil$$



- Dla skończonych przestrzeni hipotez  $\mathbb{H}$  mamy

$$m_L(\mathbb{H}, \delta, \varepsilon) \leq \left\lceil \frac{1}{\varepsilon} \ln \frac{|\mathbb{H}|}{\delta} \right\rceil = \left\lceil \frac{1}{\varepsilon} (\ln |\mathbb{H}| + \ln(1/\delta)) \right\rceil$$

- **Twierdzenie** Niech  $VCdim(\mathbb{H}) = d \geq 1$ . Wówczas każdy algorytm niesprzeczny  $\mathfrak{L}$  jest PAC oraz wymagana liczba przykładów dla  $\mathfrak{L}$  wynosi

$$m_L(\mathbb{H}, \delta, \varepsilon) \leq \left\lceil \frac{4}{\varepsilon} \left( d \log \frac{12}{\varepsilon} + \log \frac{2}{\delta} \right) \right\rceil$$



- Dla skończonych przestrzeni hipotez  $\mathbb{H}$  mamy

$$m_L(\mathbb{H}, \delta, \varepsilon) \leq \left\lceil \frac{1}{\varepsilon} \ln \frac{|\mathbb{H}|}{\delta} \right\rceil = \left\lceil \frac{1}{\varepsilon} (\ln |\mathbb{H}| + \ln(1/\delta)) \right\rceil$$

- **Twierdzenie** Niech  $VCdim(\mathbb{H}) = d \geq 1$ . Wówczas każdy algorytm niesprzeczny  $\mathcal{L}$  jest PAC oraz wymagana liczba przykładów dla  $\mathcal{L}$  wynosi

$$m_L(\mathbb{H}, \delta, \varepsilon) \leq \left\lceil \frac{4}{\varepsilon} \left( d \log \frac{12}{\varepsilon} + \log \frac{2}{\delta} \right) \right\rceil$$

- Dolne ograniczenia:

- Dla skończonych przestrzeni hipotez  $\mathbb{H}$  mamy

$$m_L(\mathbb{H}, \delta, \varepsilon) \leq \left\lceil \frac{1}{\varepsilon} \ln \frac{|\mathbb{H}|}{\delta} \right\rceil = \left\lceil \frac{1}{\varepsilon} (\ln |\mathbb{H}| + \ln(1/\delta)) \right\rceil$$

- **Twierdzenie** Niech  $VCdim(\mathbb{H}) = d \geq 1$ . Wówczas każdy algorytm niesprzeczny  $\mathcal{L}$  jest PAC oraz wymagana liczba przykładów dla  $\mathcal{L}$  wynosi

$$m_L(\mathbb{H}, \delta, \varepsilon) \leq \left\lceil \frac{4}{\varepsilon} \left( d \log \frac{12}{\varepsilon} + \log \frac{2}{\delta} \right) \right\rceil$$

- Dolne ograniczenia:
  - $m_L(\mathbb{H}, \delta, \varepsilon) \geq d(1 - \varepsilon)$

- Dla skończonych przestrzeni hipotez  $\mathbb{H}$  mamy

$$m_L(\mathbb{H}, \delta, \varepsilon) \leq \left\lceil \frac{1}{\varepsilon} \ln \frac{|\mathbb{H}|}{\delta} \right\rceil = \left\lceil \frac{1}{\varepsilon} (\ln |\mathbb{H}| + \ln(1/\delta)) \right\rceil$$

- **Twierdzenie** Niech  $VCdim(\mathbb{H}) = d \geq 1$ . Wówczas każdy algorytm niesprzeczny  $\mathcal{L}$  jest PAC oraz wymagana liczba przykładów dla  $\mathcal{L}$  wynosi

$$m_L(\mathbb{H}, \delta, \varepsilon) \leq \left\lceil \frac{4}{\varepsilon} \left( d \log \frac{12}{\varepsilon} + \log \frac{2}{\delta} \right) \right\rceil$$

- Dolne ograniczenia:

- $m_L(\mathbb{H}, \delta, \varepsilon) \geq d(1 - \varepsilon)$
- Jeśli  $\delta \leq 1/100$  i  $\varepsilon \leq 1/8$ , to  $m_L(\mathbb{H}, \delta, \varepsilon) > \frac{d-1}{32\varepsilon}$

- Dla skończonych przestrzeni hipotez  $\mathbb{H}$  mamy

$$m_L(\mathbb{H}, \delta, \varepsilon) \leq \left\lceil \frac{1}{\varepsilon} \ln \frac{|\mathbb{H}|}{\delta} \right\rceil = \left\lceil \frac{1}{\varepsilon} (\ln |\mathbb{H}| + \ln(1/\delta)) \right\rceil$$

- **Twierdzenie** Niech  $VCdim(\mathbb{H}) = d \geq 1$ . Wówczas każdy algorytm niesprzeczny  $\mathfrak{L}$  jest PAC oraz wymagana liczba przykładów dla  $\mathfrak{L}$  wynosi

$$m_L(\mathbb{H}, \delta, \varepsilon) \leq \left\lceil \frac{4}{\varepsilon} \left( d \log \frac{12}{\varepsilon} + \log \frac{2}{\delta} \right) \right\rceil$$

- Dolne ograniczenia:

- $m_L(\mathbb{H}, \delta, \varepsilon) \geq d(1 - \varepsilon)$
- Jeśli  $\delta \leq 1/100$  i  $\varepsilon \leq 1/8$ , to  $m_L(\mathbb{H}, \delta, \varepsilon) > \frac{d-1}{32\varepsilon}$
- $m_L(\mathbb{H}, \delta, \varepsilon) > \frac{1-\varepsilon}{\varepsilon} \ln \frac{1}{\delta}$



## 1. Wyuczalność

Kiedy każdy „wzorowy uczeń” będzie PAC-owy?



## 1. Wyuczalność

Kiedy każdy „wzorowy uczeń” będzie PAC-owy?

## 2. Liczba przykładów

Ile przykładów musi mieć uczeń, by się nauczyć?



## 1. Wyuczalność

Kiedy każdy „wzorowy uczeń” będzie PAC-owy?

## 2. Liczba przykładów

Ile przykładów musi mieć uczeń, by się nauczyć?

## Skończoność wymiaru $VCdim()$

1  $VCdim(\mathbb{C}) = d < \infty \Leftrightarrow \mathbb{C}$  jest wyuczalna;

2 Wówczas  $L(\frac{1}{\varepsilon}, \frac{1}{\delta}, d) < m(\varepsilon, \delta) < U(\frac{1}{\varepsilon}, \frac{1}{\delta}, d)$

## 3. Ocena ucznia



$$R(\alpha) = \min_{\alpha \in A} \int Q_{\Omega}^c(h_{\alpha}) d\mu$$

na podstawie  $N$  losowych przykładów

$$R(\alpha_N) = \min_{\alpha_i \in D} \frac{1}{N} \sum_{i=1}^N Q^c(h_{\alpha_i})$$

Kiedy i jak szybko  $R(\alpha_N) \rightarrow R(\alpha)$ ?



## 3. Ocena ucznia



$$R(\alpha) = \min_{\alpha \in A} \int Q_{\Omega}^c(h_{\alpha}) d\mu$$

na podstawie  $N$  losowych przykładów

$$R(\alpha_N) = \min_{\alpha_i \in D} \frac{1}{N} \sum_{i=1}^N Q^c(h_{\alpha_i})$$

Kiedy i jak szybko  $R(\alpha_N) \rightarrow R(\alpha)$ ?

## Skończoność wymiaru $VCdim()$

3 Dla algorytmów typu ERM,  $R(\alpha_N) \rightarrow R(\alpha)$  szybko.

$$\text{Prob} \{ R(\alpha) - R(\alpha_N) > \varepsilon \} < e^{-2c\varepsilon^2 N}$$

# Plan wykładu

---



- 1 Wstęp do komputerowego uczenia się pojęć
- 2 Model PAC (probably approximately correct)
- 3 Wyuczalność klasy pojęć
- 4 Wymiar Vapnika Chervonenkisa (VC dimension)
- 5 Podstawowe twierdzenia teorii uczenia się
- 6 Appendix: „Nie ma nic za darmo” czyli “Non Free Lunch Theorem”

# O co chodzi w NFL?

---



- Znaleźć optimum nieznannej funkcji  $f : S \rightarrow W$  ( $f \in \mathcal{F}$ ), gdzie  $S, W$  są skończonymi zbiorami.



- Znaleźć optimum nieznannej funkcji  $f : S \rightarrow W$  ( $f \in \mathcal{F}$ ), gdzie  $S, W$  są skończonymi zbiorami.
- Działanie algorytmu przeszukiwania  $\mathcal{A}$  dla funkcji  $f$  jest identyfikowany z wektorem:

$$V_{\mathcal{A}}(f, t) = \langle (s_1, f(s_1)), (s_2, f(s_2)), \dots, (s_t, f(s_t)) \rangle$$

# O co chodzi w NFL?



- Znaleźć optimum nieznannej funkcji  $f : S \rightarrow W$  ( $f \in \mathcal{F}$ ), gdzie  $S, W$  są skończonymi zbiorami.
- Działanie algorytmu przeszukiwania  $\mathcal{A}$  dla funkcji  $f$  jest identyfikowany z wektorem:

$$V_{\mathcal{A}}(f, t) = \langle (s_1, f(s_1)), (s_2, f(s_2)), \dots, (s_t, f(s_t)) \rangle$$

- Ocena algorytmu:  $M : \{V_{\mathcal{A}}(f, t) | \mathcal{A}, f, t\} \rightarrow \mathbb{R}$ ;  
Np.  $M(V_{\mathcal{A}}(f, t)) = \min\{i | f(s_i) = f_{\max}\}$

- **Warunek NFL:** Dla dowolnej funkcji  $M$ , i dla dowolnych algorytmów  $\mathcal{A}, \mathcal{A}'$

$$\sum_{f \in \mathcal{F}} M(V_{\mathcal{A}}(f, |S|)) = \sum_{f \in \mathcal{F}} M(V_{\mathcal{A}'}(f, |S|))$$



- **Warunek NFL:** Dla dowolnej funkcji  $M$ , i dla dowolnych algorytmów  $\mathcal{A}, \mathcal{A}'$

$$\sum_{f \in \mathcal{F}} M(V_{\mathcal{A}}(f, |S|)) = \sum_{f \in \mathcal{F}} M(V_{\mathcal{A}'}(f, |S|))$$

- **$\mathcal{F}$  jest zamknięta wzg. permutacji:** dla dowolnej funkcji  $f \in \mathcal{F}$  i dowolnej permutacji  $\sigma \in \text{Perm}(S)$  mamy  $\sigma f \in \mathcal{F}$



- **Warunek NFL:** Dla dowolnej funkcji  $M$ , i dla dowolnych algorytmów  $\mathcal{A}, \mathcal{A}'$

$$\sum_{f \in \mathcal{F}} M(V_{\mathcal{A}}(f, |S|)) = \sum_{f \in \mathcal{F}} M(V_{\mathcal{A}'}(f, |S|))$$

- **$\mathcal{F}$  jest zamknięta wzg. permutacji:** dla dowolnej funkcji  $f \in \mathcal{F}$  i dowolnej permutacji  $\sigma \in \text{Perm}(S)$  mamy  $\sigma f \in \mathcal{F}$

## Twierdzenie o NFL





- **Warunek NFL:** Dla dowolnej funkcji  $M$ , i dla dowolnych algorytmów  $\mathcal{A}, \mathcal{A}'$

$$\sum_{f \in \mathcal{F}} M(V_{\mathcal{A}}(f, |S|)) = \sum_{f \in \mathcal{F}} M(V_{\mathcal{A}'}(f, |S|))$$

- **$\mathcal{F}$  jest zamknięta wzg. permutacji:** dla dowolnej funkcji  $f \in \mathcal{F}$  i dowolnej permutacji  $\sigma \in \text{Perm}(S)$  mamy  $\sigma f \in \mathcal{F}$

### Twierdzenie o NFL

- zachodzi równoważność

$$NFL \Leftrightarrow \mathcal{F} \text{ jest zamknięta wzg. permutacji}$$

- **Warunek NFL:** Dla dowolnej funkcji  $M$ , i dla dowolnych algorytmów  $\mathcal{A}, \mathcal{A}'$

$$\sum_{f \in \mathcal{F}} M(V_{\mathcal{A}}(f, |S|)) = \sum_{f \in \mathcal{F}} M(V_{\mathcal{A}'}(f, |S|))$$

- **$\mathcal{F}$  jest zamknięta wzg. permutacji:** dla dowolnej funkcji  $f \in \mathcal{F}$  i dowolnej permutacji  $\sigma \in \text{Perm}(S)$  mamy  $\sigma f \in \mathcal{F}$

### Twierdzenie o NFL

- zachodzi równoważność

$NFL \Leftrightarrow \mathcal{F}$  jest zamknięta wzg. permutacji

- Prawdopodobieństwo wylosowania niepustej klasy funkcji zamkniętej wzg. permutacji wynosi:

$$\frac{2^{\binom{|S|+|W|-1}{|S|}} - 1}{2^{|S||W|} - 1}$$

# The No Free Lunch Theorem for learning

---

- Algorytm  $\mathcal{L}$  dobrze się uczy pojęcia  $c$  jeśli  $er_{\Omega}^c$  jest mały.



# The No Free Lunch Theorem for learning

---



- Algorytm  $\mathcal{L}$  dobrze się uczy pojęcia  $c$  jeśli  $er_{\Omega}^c$  jest mały.
- Niech  $\mathbb{P}(X) = \{c : X \rightarrow \{0, 1\}\}$ .  
Czy można stwierdzić wiedzieć, że  $\mathcal{L}_1$  uczy się wszystkich pojęć z  $\mathbb{P}(X)$  lepiej od  $\mathcal{L}_2$ ?

# The No Free Lunch Theorem for learning

---



- Algorytm  $\mathcal{L}$  dobrze się uczy pojęcia  $c$  jeśli  $er_{\Omega}^c$  jest mały.
- Niech  $\mathbb{P}(X) = \{c : X \rightarrow \{0, 1\}\}$ .  
Czy można stwierdzić wiedzieć, że  $\mathcal{L}_1$  uczy się wszystkich pojęć z  $\mathbb{P}(X)$  lepiej od  $\mathcal{L}_2$ ?
- “No Free Lunch theorem” (Wolpert, Schaffer) w wersji problemów uczenia się głosi, że:

# The No Free Lunch Theorem for learning

---



- Algorytm  $\mathcal{L}$  dobrze się uczy pojęcia  $c$  jeśli  $er_{\Omega}^c$  jest mały.
- Niech  $\mathbb{P}(X) = \{c : X \rightarrow \{0, 1\}\}$ .  
Czy można stwierdzić wiedzieć, że  $\mathcal{L}_1$  uczy się wszystkich pojęć z  $\mathbb{P}(X)$  lepiej od  $\mathcal{L}_2$ ?
- “No Free Lunch theorem” (Wolpert, Schaffer) w wersji problemów uczenia się głosi, że:
  - Żaden algorytm nie może być najlepszy w uczeniu wszystkich pojęć.

# The No Free Lunch Theorem for learning

---



- Algorytm  $\mathcal{L}$  dobrze się uczy pojęcia  $c$  jeśli  $er_{\Omega}^c$  jest mały.
- Niech  $\mathbb{P}(X) = \{c : X \rightarrow \{0, 1\}\}$ .  
Czy można stwierdzić wiedzieć, że  $\mathcal{L}_1$  uczy się wszystkich pojęć z  $\mathbb{P}(X)$  lepiej od  $\mathcal{L}_2$ ?
- “No Free Lunch theorem” (Wolpert, Schaffer) w wersji problemów uczenia się głosi, że:
  - Żaden algorytm nie może być najlepszy w uczeniu wszystkich pojęć.
  - Każdy algorytm jest najlepszy dla takiej samej liczby pojęć

# The No Free Lunch Theorem for learning



- Algorytm  $\mathcal{L}$  dobrze się uczy pojęcia  $c$  jeśli  $er_{\Omega}^c$  jest mały.
- Niech  $\mathbb{P}(X) = \{c : X \rightarrow \{0, 1\}\}$ .  
Czy można stwierdzić wiedzieć, że  $\mathcal{L}_1$  uczy się wszystkich pojęć z  $\mathbb{P}(X)$  lepiej od  $\mathcal{L}_2$ ?
- “No Free Lunch theorem” (Wolpert, Schaffer) w wersji problemów uczenia się głosi, że:
  - Żaden algorytm nie może być najlepszy w uczeniu wszystkich pojęć.
  - Każdy algorytm jest najlepszy dla takiej samej liczby pojęć
  - Ale interesuje nas tylko pewna klasa problemów czyli klasa pojęć  $\mathbb{C} \subset \mathbb{P}(X)$



# The No Free Lunch Theorem for learning



- Algorytm  $\mathcal{L}$  dobrze się uczy pojęcia  $c$  jeśli  $er_{\Omega}^c$  jest mały.
- Niech  $\mathbb{P}(X) = \{c : X \rightarrow \{0, 1\}\}$ .  
Czy można stwierdzić wiedzieć, że  $\mathcal{L}_1$  uczy się wszystkich pojęć z  $\mathbb{P}(X)$  lepiej od  $\mathcal{L}_2$ ?
- “No Free Lunch theorem” (Wolpert, Schaffer) w wersji problemów uczenia się głosi, że:
  - Żaden algorytm nie może być najlepszy w uczeniu wszystkich pojęć.
  - Każdy algorytm jest najlepszy dla takiej samej liczby pojęć
  - Ale interesuje nas tylko pewna klasa problemów czyli klasa pojęć  $C \subset \mathbb{P}(X)$
  - Wniosek: Należy znaleźć odp. algorytm do każdego problemu.