

# Czego mogą się nauczyć komputery?

Andrzej Skowron, Hung Son Nguyen

`son@mimuw.edu.pl`; `skowron@mimuw.edu.pl`

Instytut Matematyki,  
Wydział MIM, UW

# Spis treści

- Wprowadzenie do uczenia maszynowego;
- Co to znaczy, że komputer się wyuczy jakichś pojęć?
  - Model PAC
  - Efektywne uczenie się
- Jakich klas pojęć jest się w stanie wyuczyć komputer?
  - Wymiar Vapnika-Chervonenkisa
  - Ciekawe oszacowania

# Kto się uczy?

Ograniczymy się do programów komputerowych zwanych ”*algorytmami uczącymi się*”. Tak jak zwykli uczniowie, są na nie nałożone pewne rygorystyczne wymagania: np.

- *muszą być skuteczne,*
- *mogą korzystać z ograniczonej pamięci,*
- *mogą używać pewnej specyficznej reprezentacji wiedzy.*

# Dziedzina: “czego” się uczy?

Przykłady uczenia się:

- *pojęć*: uczeń próbuje znaleźć reguł odróżniających pozytywne przykłady od negatywnych. (np. pojęcie “krzesło”).
- *nieznanych urządzeń* – np. używanie VCR
- *nieznanych środowisk* – (np. nowe miasto)
- *procesów* (np. pieczenie ciasta)
- *nieznanych rodzin podobnych wzorców* (np. rozpoznawanie mowy, twarzy lub pisma)
- *funkcji*: (np. funkcje boolowskie)

# Model uczenia

Każdy “model uczenia” powinien uwzględnić proces uczenia się różnych przedmiotów.

*Np. jeśli uczymy się funkcji, to ważne jest aby “algorytm uczenia się” nie ograniczał się do jednej konkretnej funkcji. Żądamy aby “modele uczenia” działały skutecznie na klasach funkcji.*

# Źródło informacji:

Uczeń może pozyskać informacje o dziedzinie poprzez:

1. **Przykłady:** Uczeń dostaje pozytywne i/lub negatywne przykłady. Przykłady mogą być wybrane
  - (a) losowo według pewnego znanego lub nieznanego rozkładu;
  - (b) arbitralnie;
  - (c) złośliwie (np. przez kontrolera, który chciałby poznać najgorsze zachowanie algorytmu uczenia się);
  - (d) specjalnie przez życzliwego nauczyciela (aby ułatwić proces uczenia się)
2. **Zapytania:**
3. **Eksperymentowanie:** (aktywne uczenie się)

# Kryteria oceny jakości:

## Skąd wiemy, czy uczeń się nauczył lub jak dobrze się nauczył?

- Miary off-line (batch) vs. on-line (interactive).
- Jakość opisu vs. jakość predykcji
- Skuteczność: obliczona na podstawie błędu klasyfikacji, dokładności opisu ...
- Efektywność uczenia: wymagana jest wielomianowa złożoność obliczeniowa.

# Przykład

- Załóżmy, że chcemy nauczyć się pojęcia "człowieka o średniej budowie ciała". Dane – czyli osoby – są reprezentowane przez punkty  $(wzrost(c^m), waga(Kg))$  i są etykietowane przez  $+$  dla pozytywnych przykładów i  $-$  dla negatywnych.
- Dodatkowa wiedza: szukane pojęcie można wyrazić za pomocą PROSTOKĄTA
- Na przykład dany jest etykietowany zbiór:  
 $((84, 184), +)$ ,  $((70, 170), +)$ ,  $((75, 163), -)$ ,  
 $((80, 180), +)$ ,  $((81, 195), -)$ ,  $((63, 191), -)$ ,  
 $((77, 187), -)$ ,  $((68, 168), +)$
- Znajdź etykietę  $((79, 183, ?))$



# Problem uczenia się prostokąta

Możemy definiować problem jak następująco:

- Cel: Znaleźć w  $\mathbb{R}^2$  prostokąt  $R$  o bokach równoległych do osi.
- Wejście: Zbiór zawierający przykłady w postaci punktów  $(x, y), +/ -$  . Te punkty zostały wygenerowane losowo.
- Wyjście: Znaleźć hipotetyczny prostokąt  $R'$  będący ”dobrą aproksymacją”  $R$ .
- Dodatkowe wymagania: Algorytm powinien być efektywny (czasowo) używając do uczenia najmniejszej liczby przykładów.

# Ogólny model uczenia się

- Dane są
  - zbiór wszystkich obiektów  $X$  (skończony lub nie);
  - pojęcie  $c \in \mathbb{C}$  (funkcja celu);
  - skończona próbka  $D$  obiektów  $x_1, \dots, x_m \in X$  wraz z wartością funkcji  $c$  na tych obiektach;
  - przestrzeń hipotez  $\mathbb{H}$ ;
- Szukane
  - hipoteza  $h \in \mathbb{H}$  będąca dobrą aproksymacją pojęcia  $c$ .
- Wymagane
  - dobra jakość aproksymacji
  - szybki czas działania.

# Inne przykłady

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

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

- **Uczenie hiperpłaszczyzny:**

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

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 jednomianów Boolowskich:**

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

$\mathbb{H} = M_n =$  zbiór jednomianów Boolowskich o  $n$  zmiennych, t.j. zbiór koniunkcji literałów będących albo zmiennymi bądź negacjami zmiennych.

# Błąd hipotezy

Niech

- $X$  – zbiór wszystkich obiektów.
- $\Omega = (X, \mu)$  – przestrzeń probabilistyczna określona na  $X$ .

Błąd hipotezy  $h \in \mathbb{H}$  względem pojęcia  $c$  (funkcji docelowej):

$$er_{\Omega}(h, c) = er_{\Omega}^c(h) = \mu\{x \in X | h(x) \neq c(x)\}$$

Z prawdopodobieństwem  $(1 - \varepsilon)$  możemy oszacować  $er_{\Omega}^c$ :

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

# The No Free Lunch Theorem

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

Niech  $\mathbb{C}(X) = \{c : X \rightarrow \{0, 1\}\}$ . Czy można powiedzieć, że  $L_1$  uczy się wszystkich pojęć z  $\mathbb{C}(X)$  lepiej od  $L_2$ ?

”No Free Lunch theorem” (Wolpert, Schaffer) głosi, że:

- Żaden algorytm nie może być najlepszy w uczeniu wszystkich pojęć.
- Każdy algorytm jest najlepszy dla takiej samej liczby pojęć
- Ale interesują nam tylko konkretne problemy czyli klasy pojęć  $\mathbb{C} \subset \mathbb{C}(X)$
- Wniosek: Dopasuj algorytm do problemu.

# Model uczenia się PAC

Dane: dziedzina  $X$ , klasa pojęć  $\mathbb{C}$  i przestrzeń hipotez ucznia  $\mathbb{H}$ .

# Model uczenia się PAC

Dane: dziedzina  $X$ , klasa pojęć  $\mathbb{C}$  i przestrzeń hipotez ucznia  $\mathbb{H}$ .

- Uczeń uczy się pojęcia  $c \in \mathbb{C}$  na podstawie przykładów wygenerowanych przez zm.1.  $EX(\Omega, c)$  (zwaną *wyrocznią*).

rodzina wszystkich zbiorów zawierających  $m$  przykładów

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

# Model uczenia się PAC

Dane: dziedzina  $X$ , klasa pojęć  $\mathbb{C}$  i przestrzeń hipotez ucznia  $\mathbb{H}$ .

- Uczeń uczy się pojęcia  $c \in \mathbb{C}$  na podstawie przykładów wygenerowanych przez zm.1.  $EX(\Omega, c)$  (zwaną *wyrocznią*).  
rodzina wszystkich zbiorów zawierających  $m$  przykładów

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

- Celem do wyuczenia przez ucznia jest znalezienie hipotezy minimalizującej błąd rzeczywisty względem  $c$  dla rozkładu  $\Omega$ , czyli  $er_{\Omega}^c$ .



# Model uczenia się PAC

Dane: dziedzina  $X$ , klasa pojęć  $\mathbb{C}$  i przestrzeń hipotez ucznia  $\mathbb{H}$ .

- Uczeń uczy się pojęcia  $c \in \mathbb{C}$  na podstawie przykładów wygenerowanych przez zm.1.  $EX(\Omega, c)$  (zwaną *wyrocznią*).  
rodzina wszystkich zbiorów zawierających  $m$  przykładów

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

- Celem do wyuczenia przez ucznia jest znalezienie hipotezy minimalizującej błąd rzeczywisty względem  $c$  dla rozkładu  $\Omega$ , czyli  $er_{\Omega}^c$ .
- Zasadnicza idea modelu PAC = określenie warunków, pod jakimi uczeń (lub algorytm uczenia się) znajdzie “dobrą hipotezę” z “dużym prawdopodobieństwem”.

o ograniczonym błędzie rzeczywistym

powyżej określonego progu

# PAC uczenie się

Definicja: Mówimy, że *prawdopodobnie algorytm uczenia się*  $L$  jest *aproxymacyjnie poprawny*  $\Leftrightarrow$  dla każdego  $0 < \varepsilon, \delta < 1$ , istnieje liczba  $m_0 = m_0(\varepsilon, \delta)$  taka, że dla dowolnego pojęcia  $c \in \mathbb{C}$  i dla dowolnego rozkładu  $\Omega$  na  $X$ , jeśli  $m > m_0$

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

Wówczas mówimy w skrócie, że  **$L$  jest PAC**.

*(Probably Approximately Correct)*

- $\varepsilon$  – dopuszczalny poziom błędu
- $(1 - \delta)$  – poziom zaufania

# Przykład problemu dyskretyzacji

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

# Przykład problemu dyskretyzacji

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

# Przykład problemu dyskretyzacji

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

- $X = \mathfrak{R}; \mathbb{H} = \mathbb{C} = \{f_\lambda : \mathfrak{R} \rightarrow \{0, 1\} | f_\lambda(x) = 1 \Leftrightarrow x \geq \lambda\}$
- $c = f_{\lambda_0}$
- 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.  $L(D) := f_{\lambda^*};$

# Przykład problemu dyskretyzacji

- $X = \mathfrak{R}; \mathbb{H} = \mathbb{C} = \{f_\lambda : \mathfrak{R} \rightarrow \{0, 1\} \mid f_\lambda(x) = 1 \Leftrightarrow x \geq \lambda\}$
- $c = f_{\lambda_0}$
- 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.  $L(D) := f_{\lambda^*}$ ;

**Twierdzenie:** Powyższy algorytm jest PAC

# Przykład (c.d.)

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



# Przykład (c.d.)

- $er_{\Omega}^{\varepsilon} = \mu([\lambda_0, \lambda^*))$
- Niech  $\beta_0 = \sup\{\beta \mid \mu([\lambda_0, \beta)) < \varepsilon\}$ . Wówczas  $er_{\Omega}^{\varepsilon}(f_{\lambda^*}) \leq \varepsilon \Leftrightarrow \lambda^* \leq \beta_0 \Leftrightarrow$  jeden z przykładów  $x_i$  znajduje się w przedziale  $[\lambda_0, \beta_0]$ ;

# Przykład (c.d.)

- $er_{\Omega}^c = \mu([\lambda_0, \lambda^*))$
- Niech  $\beta_0 = \sup\{\beta | \mu([\lambda_0, \beta)) < \varepsilon\}$ . Wówczas  $er_{\Omega}^c(f_{\lambda^*}) \leq \varepsilon \Leftrightarrow \lambda^* \leq \beta_0 \Leftrightarrow$  jeden z przykładów  $x_i$  znajduje się w przedziale  $[\lambda_0, \beta_0]$ ;
- Prawdopodobieństwo tego, że żaden spośród  $m$  przykładów nie należy do  $[\lambda_0, \beta_0]$  jest  $\leq (1 - \varepsilon)^m$ . Stąd

$$\mu^m \{D \in \mathcal{S}(m, f_{\lambda_0}) | er_{\Omega}(L(D)) \leq \varepsilon\} \geq 1 - (1 - \varepsilon)^m$$

# Przykład (c.d.)

- $er_{\Omega}^c = \mu([\lambda_0, \lambda^*))$
- Niech  $\beta_0 = \sup\{\beta | \mu([\lambda_0, \beta)) < \varepsilon\}$ . Wówczas  $er_{\Omega}^c(f_{\lambda^*}) \leq \varepsilon \Leftrightarrow \lambda^* \leq \beta_0 \Leftrightarrow$  jeden z przykładów  $x_i$  znajduje się w przedziale  $[\lambda_0, \beta_0]$ ;
- Prawdopodobieństwo tego, że żaden spośród  $m$  przykładów nie należy do  $[\lambda_0, \beta_0]$  jest  $\leq (1 - \varepsilon)^m$ . Stąd

$$\mu^m\{D \in \mathcal{S}(m, f_{\lambda_0}) | er_{\Omega}(L(D)) \leq \varepsilon\} \geq 1 - (1 - \varepsilon)^m$$

- Aby to prawdopodobieństwo było  $> 1 - \delta$ , wystarczy wybrać

$$m \geq m_0 = \left\lceil \frac{1}{\varepsilon} \ln \frac{1}{\delta} \right\rceil$$

# Przykład (c.d.)

- $er_{\Omega}^c = \mu([\lambda_0, \lambda^*))$
- Niech  $\beta_0 = \sup\{\beta | \mu([\lambda_0, \beta)) < \varepsilon\}$ . Wówczas  $er_{\Omega}^c(f_{\lambda^*}) \leq \varepsilon \Leftrightarrow \lambda^* \leq \beta_0 \Leftrightarrow$  jeden z przykładów  $x_i$  znajduje się w przedziale  $[\lambda_0, \beta_0]$ ;
- Prawdopodobieństwo tego, że żaden spośród  $m$  przykładów nie należy do  $[\lambda_0, \beta_0]$  jest  $\leq (1 - \varepsilon)^m$ . Stąd

$$\mu^m\{D \in \mathcal{S}(m, f_{\lambda_0}) | er_{\Omega}(L(D)) \leq \varepsilon\} \geq 1 - (1 - \varepsilon)^m$$

- Aby to prawdopodobieństwo było  $> 1 - \delta$ , wystarczy wybrać

$$m \geq m_0 = \left\lceil \frac{1}{\varepsilon} \ln \frac{1}{\delta} \right\rceil$$

# Dokładne uczenie się

- 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  $L$  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  $L$  jest dokładnym algorytmem (*jest PEC – probably exactly correct*)

# Potencjalna wyuczalność

- Algorytm  $L$  nazywamy **niesprzecznym** jeśli dla każdego pojęcia  $c$  i każdego zbioru  $D$  mamy  $er_D^c(L(D)) = 0$  (tzn.  $L(D)(x_i) = c(x_i)$  dla dowolnego przykładowego  $(x_i, c(x_i)) \in D$ ).
- $\mathbb{H}^c(D) = \{h \in \mathbb{H} \mid h(x_i) = c(x_i) (i = 1, \dots, m)\}$ .  
 $L$  jest niespreczny jeśli  $L(D) \in \mathbb{H}^c(D)$  dla każdego  $D$ .
- $\mathbb{B}_\varepsilon^c = \{h \in \mathbb{H} \mid er_\Omega(h) \geq \varepsilon\}$

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

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

# Potencjalna wyuczalność

**Twierdzenie:** Jeśli

1.  $\mathbb{C}$  jest potencjalnie wyuczalne za pomocą  $\mathbb{H}$
2.  $L$  jest algorytmem niesprzecznym dla  $\mathbb{C}$

Wówczas  $L$  jest PAC

**Twierdzenie:** (Haussler, 1988) Jeśli  $\mathbb{C} = \mathbb{H}$  i  $|\mathbb{C}| < \infty$ , to  $\mathbb{C}$  jest potencjalnie wyuczalne.

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

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

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

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

# Wymiar Vapnika-Chervonenkisa

- Niech  $\mathbf{x} = \{x_1, x_2, \dots, x_m\} \subset X$ . Oznaczmy przez  $\Pi_{\mathbb{H}}(\mathbf{x})$  liczbę podziałów zbioru  $\mathbf{x}$  dokonanych przez  $\mathbb{H}$ , t.j. liczbę różnych wektorów postaci

$$(h(x_1), \dots, h(x_m)) \in \{0, 1\}^m$$

po wszystkich  $h \in H$ .

- $\Pi_{\mathbb{H}}(\mathbf{x}) \leq 2^m$ . Jeśli zachodzi równość, mówimy, że  $\mathbb{H}$  **rozbija**  $\mathbf{x}$ .
- Niech  $\Pi_{\mathbb{H}}(m) = \max_{\mathbf{x} \in X^m} \Pi_{\mathbb{H}}(\mathbf{x})$

Na przykład: W przypadku przestrzeni półosi postaci  $[\alpha, \infty)$  mamy  $\Pi_{\mathbb{H}}(m) = m + 1$ .

- Na ogół trudno znaleźć wzór na  $\Pi_{\mathbb{H}}(m)$ !!!



# Wymiar Vapnika-Chervonenkisa (c.d.)

Uwagi:

- Jeśli  $\Pi_{\mathbb{H}}(m) = 2^m$ , to  $\mathbb{H}$  rozbija prawie każdy zbiór o mocy  $m$  (prawie zawsze można znaleźć niesprzeczną hipotezę).
- Maksymalna wartość  $m$ , dla której  $\Pi_{\mathbb{H}}(m) = 2^m$  można uważać za moc wyrażania przestrzeni  $\mathbb{H}$

**Definicja:** 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

- Niech  $\mathbb{H}_1$  będzie przestrzenią półosi. Wówczas  $VCdim(\mathbb{H}_1) = 1$
- Niech  $X = \mathbb{R}^2$  i  $\mathbb{H} =$  zbiór wszystkich półpłaszczyzn. Wówczas  $VCdim(\mathbb{H}) = 3$
- Niech  $\mathbb{H} =$  zbiór wszystkich półprzestrzeni wyznaczonych przez hiperpłaszczyzny w przestrzeni  $\mathbb{R}^m$ . Wówczas  $VCdim(\mathbb{H}) = m + 1$
- **Twierdzenie** Jeśli  $\mathbb{H}$  jest skończoną przestrzenią, to  $VCdim(\mathbb{H}) \leq \log |\mathbb{H}|$
- Niech  $M_n$  będzie przestrzenią wszystkich jednomianów Boolowskich o  $n$  zmiennych. Ponieważ  $|M_n| = 3^n$ , mamy

$$VCdim(M_n) \leq n \log 3$$

# Lemat Sauer'a

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

$$\Pi_{\mathbb{H}}(m) \leq 1 + \underbrace{\binom{m}{1} + \binom{m}{2} + \dots + \binom{m}{d}}_{\Phi(d,m)}$$

- Wnioski

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

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

# VC vs. potencjalna wyuczalność

**Twierdzenie** Każda przestrzeń hipotez o nieskończonym wymiarze VC nie jest potencjalnie wyuczalna.

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

$$\mu^m \{D : \mathbb{H}[D] \cap \mathbb{B}_\varepsilon \neq \emptyset\} < 2\Pi_{\mathbb{H}}(2m)2^{-\varepsilon m/2}$$

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

**Twierdzenie (Fundamentalne twierdzenie)** Jeśli przestrzeń hipotez ma skończony wymiar VC, to jest ona potencjalnie wyuczalna.

# Złożoność zbioru treningowego

- Z Fundamentalnego Twierdzenia wynika, że jeśli  $VCdim(\mathbb{H}) < \infty$ , to dla danych  $\delta$  i  $\varepsilon$ , istnieje  $m_0 = m_0(\mathbb{H}, \delta, \varepsilon)$  takie, że

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

- Wówczas każdy niesprzeczny algorytm  $L$  jest PAC oraz wymagana liczba przykładów  $m_L(\mathbb{H}, \delta, \varepsilon)$  dla  $L$  jest ograniczona z góry przez  $m_0(\mathbb{H}, \delta, \varepsilon)$ .
- Dla skończonych przestrzeni hipotez  $\mathbb{H}$  mamy

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

# Złożoność zbioru treningowego (c.d.)

- **Twierdzenie** Niech  $VCdim(\mathbb{H}) = d \geq 1$ . Wówczas każdy algorytm niesprzeczny  $L$  jest PAC oraz wymagana liczba przykładów dla  $L$  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  $m_L(\mathbb{H}, \delta, \varepsilon) > \frac{d-1}{32\varepsilon}$
  - $m_L(\mathbb{H}, \delta, \varepsilon) > \frac{1-\varepsilon}{\varepsilon} \ln \frac{1}{\delta}$