

# Výpočtové učenie

21. marca 2006

# Kapitola 1

## Efektívne učenie II

### 1.1 Účinnosť versus dôveryhodnosť a presnosť

Diskusia v predchádzajúcej kapitole sa sústredila na chovanie sa učiaceho algoritmu v čase behu, ktorý závisel (bol funkciou) len od veľkosti príkladov  $n$ . Zrejme sú aj **iné faktory, ktoré určujú čas** behu učiaceho algoritmu, a mohli by sme zaviesť pojem efektívnosti vzhľadom na **úroveň dôveryhodnosti a úroveň presnosti**. Potom prediskutujeme efektívnosť vzhľadom na veľkosť reprezentácie cieľového konceptu. Tieto úvahy sú relevantné pre ľubovoľný hypotézový priestor a môžu byť kombinované s ideami nasledujúcej kapitoly, aby bola zavedená celkom všeobecná definícia pojmu **efektívny PAC učiaci algoritmus**.

V predchádzajúcej kapitole sme zaviedli pevné, ale ľubovoľné parametre, a síce parameter dôveryhodnosti  $\delta$  a parameter presnosti  $\epsilon$ . Je zřejmé, že zníženie aspoň jednej hodnoty z týchto veličín urobí učenie ťažším, a preto čas behu efektívneho PAC učiaceho algoritmu by mohol byť ohraničovaný nejakým vhodným spôsobom pomocou rastúcich  $\delta^{-1}$  a  $\epsilon^{-1}$ . Mohli by sme sa jednoducho pýtať, či čas behu rastie polynomiálne vzhľadom na  $\delta^{-1}$  a  $\epsilon^{-1}$ , ale táto závislosť na  $\delta^{-1}$  nie je celkom vhodná z nasledujúcich dôvodov. Ak dĺžka tréningovej vzorky, ktorá vstupuje do efektívneho učiaceho algoritmu je zdvojnásobená, môžeme očakávať, že pravdepodobnosť, že výstupná hypotéza je zlá, je približne kvadratická. Inak povedané, **vzťah medzi zložitou vzorkou a  $\delta^{-1}$  je logaritmický**. Motivovaní týmto, mohli by sme povedať, že učiaci algoritmus  $L$  je **efektívny vzhľadom na dôveryhodnosť**, ak jeho čas behu je polynomiálny v  $m$  a zložitost vzorky  $m_L(H, \delta, \epsilon)$  závisí polynomiálne od veličiny  $\ln(\delta^{-1})$ , čo budeme označovať  $\delta^*$ . V prípade parametra presnosti budeme hovoriť, že  $L$  je **efektívny vzhľadom na presnosť**, ak jeho čas behu je polynomiálny v  $m$  a zložitost vzorky závisí polynomiálne od  $\epsilon^{-1}$ . Ak obe tieto podmienky platia, potom čas behu potrebný na vytvorenie PAC výstupnej hypotézy je polynomiálny v  $\delta^*$  a  $\epsilon^{-1}$ .

Napríklad, ak  $H$  je ľubovoľný konečný hypotézový priestor a  $L$  je konzistentný učiaci algoritmus pre  $H$ , potom na základe predchádzajúcich výsledkov máme, že dolná hranica pre zložitost vzorky je  $m_0(H, \delta, \epsilon) = \lceil \epsilon^{-1} \cdot \ln(|H|/\delta) \rceil$ .

V tomto prípade  $m_0$  je zrejme ohraničené polynomiálnou funkciou vzhľadom na  $\delta^*$  a  $\epsilon^{-1}$ . Ak čas behu  $L$  je polynóm v  $m$ , potom  $L$  je PAC učiaci algoritmus pre  $H$ , ktorý beží v polynomiálnom čase vzhľadom na  $\delta^*$  a  $\epsilon^{-1}$ . Ten istý argument platí v odstupňovanom prípade. Ak  $H = \bigcup H_n$  je hypotézový priestor booleovských funkcií odstupňovaný veľkosťou príkladov, potom dolná hranica pre zložitost vzorky je

$$m_0(H_n, \delta, \epsilon) = \left\lceil \frac{1}{\epsilon} \cdot \ln \left( \frac{|H_n|}{\delta} \right) \right\rceil.$$

V tomto prípade, ak čas behu  $R_L(m, n)$  je polynóm v  $m$  a  $n$  a ak  $\ln |H_n|$  je polynóm v  $n$ , potom  $L$  PAC učí  $H_n$  v polynomiálnom čase behu nielen v  $n$ , ale tiež v  $\delta^*$  a  $\epsilon^{-1}$ .

### 1.2 PAC učenie a problém konzistencie

Analýza urobená na konci predchádzajúcej sekcie je motivovaná doteraz známymi vzťahmi medzi konzistenciou a PAC učením. Obrátením našej pozornosti na neodstupňovaný prípad výsledok je jednoduchý,

a síce ak existuje konzistentný učiaci algoritmus  $L$  pre konečný hypotézový priestor  $H$ , ktorý beží v polynomiálnom čase vzhľadom na dĺžku vzorky  $m$ , potom  $L$  sa PAC učí  $H$  v polynomiálnom čase vzhľadom na  $\delta^*$  a  $\epsilon^{-1}$ . Zhruba povedané, môžeme hovoriť, že efektívny "konzistentný-hľadač-hypotéz" je efektívny "PAC-učiaci sa". V tejto sekcii uvedieme, k čomu vedie opačná úvaha.

Toto by znamenalo, že efektívne PAC učenie implikuje efektívne hľadanie konzistentných hypotéz za predpokladu, že sme pripravení akceptovať **náhodný algoritmus**.

Predpokladajme, že máme daný nejaký **generátor náhodných čísel**, ktorý pre dané ľubovoľné celé číslo  $I \geq 2$  produkuje náhodné čísla  $i$  v intervale  $1 \leq i \leq I$ , pričom každá hodnota je rovnako pravdepodobná.

Náhodný algoritmus  $A$  má povolené brať tieto čísla ako časť svojho vstupu. **Výpočet algoritmu  $A$  je riadený jeho vstupom** tak, že závisí od partikulárnej postupnosti produkovanej generátorom náhodných čísel. Z toho vyplýva, že môžeme hovoriť o pravdepodobnosti, že  $A$  má daný výsledok. Tým je mienená relatívna frekvencia postupností, ktoré produkuje tento výsledok vzhľadom na celkový počet možných postupností.

Hovoríme, že náhodný algoritmus  $A$  "rieši" problém vyhľadávania  $\Pi$ , ak sa chová nasledujúcim spôsobom: Algoritmus vždy zastaví a produkuje výstup. Ak  $A$  padol pri hľadaní riešenia pre  $\Pi$ , dá jednoducho výstup *nie*. Ale s pravdepodobnosťou aspoň  $\frac{1}{2}$  (v zmysle vyjadrenom vyššie),  $A$  je úspešný pri hľadaní riešenia pre  $\pi$  a výstupom je jeho riešenie.

Praktická použiteľnosť náhodného algoritmu vyplýva z faktu, že opakovaním algoritmu niekoľkokrát veľmi rýchlo rastie pravdepodobnosť úspechu. Ak algoritmus padne pri prvom pokuse, čo sa stane s pravdepodobnosťou najviac  $\frac{1}{2}$ , potom jednoducho skúsime ďalej. Pravdepodobnosť, že padne dvakrát, je najviac  $\frac{1}{4}$ , že padne  $k$ -krát je najviac  $(\frac{1}{2})^k \rightarrow 0$ . Teda v praxi náhodný algoritmus je takmer tak dobrý ako obyčajný - samozrejme za predpokladu, že má polynomiálny čas behu. Máme nasledujúcu vetu Pitta a Valianta (1988) (tiež Matarjan (1989) a Haussler et. al. (1988)).

**Veta 1** *Nech  $H$  je hypotézový priestor a predpokladajme, že existuje PAC učiaci algoritmus pre  $H$  s časom behu polynomiálnym v  $\epsilon^{-1}$ . Potom existuje náhodný algoritmus, ktorý rieši problém hľadania hypotézy v  $H$  konzistentnej s danou tréningovou vzorkou a ktorý má čas behu polynomiálny v  $m$  (dĺžka tréningovej vzorky).*

**Dôkaz.** Predpokladajme, že  $\bar{s}$  je tréningová vzorka pre cieľovú hypotézu  $t \in H$  a  $\bar{s}$  obsahuje  $m$  rôzne označených príkladov. Ukážeme, že je možné nájsť hypotézu konzistentnú s  $\bar{s}$  spustením daného PAC algoritmu  $L$  na odpovedajúcej tréningovej vzorke. Definujme pravdepodobnostné rozdelenie  $\mu$  na príkladovom priestore  $X$  takto:

$$\mu(x) = \begin{cases} \frac{1}{m}, & \text{ak } x \text{ sa vyskytuje v } \bar{s} \\ 0, & \text{inak.} \end{cases}$$

Môžeme použiť generátor náhodných čísel s výstupnými hodnotami  $\in \langle 1, m \rangle$  na výber príkladov z  $X$  podľa tohoto rozdelenia: každé náhodné číslo pridáme ako návestie 1.. $m$  rovnako pravdepodobného príkladu. Teda výber tréningovej vzorky dĺžky  $m$  pre  $t$  na základe rozdelenia  $\mu$  môže byť simulovaný generovaním postupnosti  $m$  náhodných čísel v danom intervale.

Nech  $L$  je PAC učiaci algoritmus, ako bolo uvedené vyššie. Potom ak máme dané 4 veličiny  $\delta, \epsilon, \mu, t$ , tak môžeme nájsť celé číslo  $m_0(\delta, \epsilon)$

$$\epsilon > 0 \quad \delta > 0 \quad m_0(\delta, \epsilon) \quad \forall \mu \quad \forall t \quad \mu^m \{s \in S(m, t) : h(s, t) < \epsilon\} > 1 - \delta$$

Predpokladajme, že špecifikujeme dôveryhodnosť  $\delta = \frac{1}{2}$  a presnosť  $\epsilon = \frac{1}{m^*}$ .

Ak spustíme učiaci algoritmus  $L$  na tréningovej vzorke dĺžky  $m_0(\frac{1}{2}, \frac{1}{m^*})$ , získanej náhodne podľa rozdelenia  $\mu$ , vlastnosť PAC algoritmu zaručí, že pravdepodobnosť, že chyba výstupu je menšia než  $\frac{1}{m^*}$ , je väčšia než  $1 - \frac{1}{2} = \frac{1}{2}$ . Pretože nie sú žiadne príklady s pravdepodobnosťou striktno medzi 0 a  $\frac{1}{m^*}$ , z toho vyplýva, že pravdepodobnosť, že výstup súhlasí presne s tréningovou vzorkou je väčší než  $\frac{1}{2}$ .

□

Procedúra uvedená vo vyššie je základom pre náhodný algoritmus  $L^*$  pre hľadanie hypotézy, ktorá je konzistentá s danou tréningovou vzorkou  $\bar{s}^*$ .

Zhrnutie v krokoch pre  $L^*$ :

- Vyhodnotiť  $m_0 = m_0(\frac{1}{2}, \frac{1}{m^*})$ .

- Použitím gnč na skonštrukciu tréningovej vzorky  $\bar{s}$  dĺžky  $m_0$  podľa pravdepodobnostného rozdelenia  $\mu$ .
- Spustiť daný PAC učiaci algoritmus  $L$  na  $\bar{s}$ .
- Skontrolovať výslednú hypotézu  $L(\bar{s})$ , či je konzistentná s  $\bar{s}^*$ .
- Ak hypotéza nie je konzistentná s  $\bar{s}^*$ , výstup "nie". Ak hypotéza je konzistentná s  $\bar{s}^*$ , výstupom je táto hypotéza.

Ako sme uviedli, PAC vlastnosť algoritmu  $L$  zaručí, že  $L^*$  je úspešný s pravdepodobnosťou viac ako  $\frac{1}{2}$ . Nakoniec je zjavné, že ak čas behu algoritmu  $L$  je polynomiálny v  $\epsilon^{-1}$ , potom čas behu  $L^*$  je polynomiálny v  $m^* = \epsilon^{-1}$ .

□

Veta 1 nám umožňuje rozšíriť zložitosť výsledkov pre problém konzistencie, ako bolo dokázané vyššie, na PAC učenie. Pripomeňme, že oba problémy - problém rozhodnutia o konzistencii a problém hľadania konzistentnej hypotézy sú NP-ťažké v niektorých prípadoch, takých ako hypotézový priestor  $C^k = \bigcup C_n^k$ . Veta hovorí, že ak by sme mohli PAC naučiť  $C_n^k$  v polynomiálnom čase vzhľadom na  $\epsilon^{-1}$  a  $n$ , potom by sme mohli nájsť konzistentnú hypotézu použitím náhodného algoritmu s časom behu polynomiálnym v  $m$  a  $n$ . V jazyku teórie zložitosti by to mohlo znamenať, že uvedený problém je v triede RP, čo je trieda problémov, ktoré môžu byť riešené v "pravdepodobne polynomiálnom čase". Teraz sa predpokladá, že RP neobsahuje žiadny NP-ťažký problém - teda, že  $RP \neq NP$ , čo niektorí považujú za podobne ťažké ako  $NP=P$ . Takže keď toto akceptujeme, z toho vyplýva, že neexistuje žiadny PAC polynomiálny algoritmus pre odstupňovaný priestor  $C^k$ , ak  $k \geq 2$ .

Vyššie uvedená diskusia ukazuje, že pre  $k \geq 2$ ,  $C^k$  nie je efektívne PAC naučiteľný vzhľadom na veľkosť príkladu. Avšak pre ľub.  $n$ ,  $C_n^k$  je obsiahnuté v  $D_{n,k}$ , priestore disjunkcií jednočlenov s najviac  $k$  literálmi. Valiantov učiaci algoritmus pre  $D_k = \bigcup D_{n,k}$  popísaný vyššie je konzistentný algoritmus s časom behu  $R_L(m, n) = O(m \cdot n^k)$ , polynóm v  $n$  aj  $m$ . Preto pre ľub. tréningovú vzorku  $\bar{s}$  pre hypotézu v  $C_n^k$  tento algoritmus bude produkovať v polynomiálnom čase hypotézu v  $D_{n,k}$  konzistentnú s  $\bar{s}$ . V ďalších poznámkach predchádzajúcej kapitoly sme definovali, čo znamená učiaci algoritmus  $L$  pre odstupňovaný priestor  $\bigcup C_n$  iným odstupňovaným priestorom  $\bigcup H_n$ ;  $k$  danej tréningovej vzorky  $\bar{s}$  pre hypotézu z  $C_n$   $L$  vráti hypotézu  $L(s) \in H_n$ . Použitím tejto terminológie, štandardný učiaci algoritmus pre  $D_k$  je PAC učiaci algoritmus pre  $(C^k, D_{n,k})$ , efektívny vzhľadom na veľkosť príkladov. Teda v protiklade k negatívnemu výsledku vyššie,  $C^k$  je efektívne naučiteľný pomocou "väčšieho priestoru". Nie je tu žiadny spor. Zhruba povedané, je ťažké nájsť formulu pre konzistentnú hypotézu v  $C_n^k$ , pretože tento priestor je príliš "ohraničený". Daná je väčšia flexibilita v práci v "bohatšom" priestore  $D_{n,k}$ , v ktorom algoritmus môže vyjadriť svoje hypotézy vo výrazoch  $D_{n,k}$  formúl; môže byť dosiahnuté rýchlejšie učenie. Výsledok o **nenučiteľnosti** je preto uvažovaný aj v zmysle závislosti od reprezentácie.

### 1.3 Veľkosť reprezentácie

Už sme sa zmienili o tom, že výstup realistického učiaceho algoritmu nie je abstraktná funkcia ale skôr reprezentácia tejto funkcie cez formulu alebo stav stroja. Pretože booleovská funkcia, ktorá môže byť reprezentovaná krátkou booleovskou formulou je zrejme "jednoduchšia" než taká, ktorá vyžaduje dlhšiu formulu, môžeme očakávať, že je ju ťažšie naučiť než krátku formulu.

Rámec potrebný pre starostlivú diskusiu o takých veciach je poskytnutý dojmom reprezentácie  $\Omega \rightarrow H$  ako bolo uvedené vyššie. Množina  $\Omega$  môže byť myslená ako množina formúl alebo množina stavov stroja tak, že pre každé  $\omega \in \Omega$  existuje odpovedajúca hypotéza  $h_\omega$ . V nasledujúcich niekoľkých sekciách uvedieme ako reprezentácia hypotéz ovplyvňuje čas behu učiacich algoritmov.

Aby sme to mohli urobiť, potrebujeme nejakú mieru pre "veľkosť" reprezentácie hypotézy. Samozrejme, že neexistuje žiadna absolútna miera a tak musíme skonštruovať jednu, ktorá sa zdá byť rozumná pre skúmané problémy. Booleovský prípad je najpriamočiarejší.

Štandardná metóda reprezentujúca bool. funkciu pomocou formúl bola popísaná v 2.3. Formálne používame abecedu  $4 + 2n$  symbolov:  $() \wedge \vee u_1 \bar{u}_1 \dots u_n \bar{u}_n$ , ktoré sú skombinované podľa určitých pravidiel. Táto abeceda môže byť zakódovaná použitím  $3 + \lceil \log n \rceil$  bitov pre každý symbol ako je možné

ukázať v nasledujúcej tabuľke:

Symbol	Kód
(	110 0000...0
)	101 000...0
∨	101 000...0
∧	111 000...0
$u_1$	001 000...01
$\overline{u_1}$	000 000...01
$u_2$	<u>001</u> <u>000...10</u>
⋮	

Idea je v tom, že prvé 3 bity sú použité na kódovanie povahy symbolu a zvyšných  $\lceil \log_2 n \rceil$  bitov je použitých na reprezentáciu symbolov  $u_i, \overline{u_i}$ . Nech  $\omega$  je správna formula, ktorá je dosiahnutá z abecedy uvedenej vyššie. Ak  $\omega$  má  $s$  symbolov, potom môže byť zakódovaná pomocou  $s \cdot (3 + \lceil \log_2 n \rceil)$  bitov a tak môžeme definovať veľkosť  $\omega$

$$\|\omega\| = s \cdot (3 + \lceil \log_2 n \rceil).$$

Napríklad, ak  $n = 3$ ,  $\omega = (u_1 \wedge u_2) \vee u_3$  (7 symbolov)

$$\|\omega\| = 7 \cdot (3 + \lceil \log_2 3 \rceil) = 35.$$

Ako sme už uviedli, výstup učiaceho algoritmu nie je abstraktná funkcia alebo hypotéza, ale reprezentácia hypotézy pomocou formuly alebo stroja. Z tohto hľadiska je rozumné porovnať veľkosť takého výstupu veľkosťou vstupu do algoritmu, čo je vlastne tréningová vzorka označených príkladov.

**Príklad 1** *Predpokladajme, že máme 20 príkladov 30 bitových do učiaceho algoritmu pre monočleny. Celkový počet bitov vstupu je  $20 \cdot (30 + 1) = 620$ . Výstupom je hypotéza-monočlen, ktorý môže byť akceptovaný pomocou 30 literálov a 29 konjunkcií. Použitím vyššie uvedenej kódovacej schémy dostávame počet bitov výstupu*

$$(30 + 29) \cdot (3 + \lceil \log_2 30 \rceil) = 59 \cdot 8 = 472$$

*Pretože toto číslo je menšie než počet bitov na vstupe, je celkom rozumné hovoriť, že výstup je (v nejakom zmysle) kompresovaná forma vstupu. □*

Tento príklad ilustruje, že môžeme očakávať, že učiaci algoritmus dáva na výstup reprezentáciu  $\omega$  hypotézy  $h_\omega$  takú, že  $h_\omega$  nie je len rozšírením tréningovej vzorky, ale  $\omega$  je kompresovaným tvarom vstupu. To je v zmysle, že  $\omega$  obsahuje tak veľa informácie ako bolo vo vstupnej tréningovej vzorke, definuje rozšírenú funkciu a vyžaduje menej bitov než tréningová vzorka. V ďalšom uvidíme, že ak učiaci algoritmus  $L$  dá na výstupe reprezentáciu hypotézy, ktorá nie je príliš dlhá a je signifikantne kompresovaná vzhľadom na vstup, potom  $L$  má určité pravdepodobnostné aproximatívne vlastnosti.

## 1.4 Hľadanie najmenej konzistentnej hypotézy

Predpokladajme, že je daná tréningová vzorka pre monočlen. Štandardný učiaci algoritmus skonštruje v polynomiálnom čase hypotézu konzistentnú s tréningovou vzorkou. Avšak môžeme očakávať viac a pýtať sa na najmenší monočlen konzistentný so vzorkou. V tomto kontexte môžeme ignorovať symboly konjunkcie v množine a budeme uvažovať aproximáciu  $k \log n$  pre dĺžku monočlena tvoreného  $k$  literálmi a definovaného na  $\{0, 1\}^n$ . Teda v ľubovoľnej danej podmnožine  $M_n$  najmenší monočlen je taký, ktorý má najmenší počet literálov. Ukážeme, že problém nájdenia najmenšieho monočlena konzistentného s tréningovou vzorkou je NP-ťažký problém.

Náš cieľ bude dosiahnutý pomocou známeho NP-úplného problému. Predpokladajme, že  $U$  je konečná množina a  $\mathbf{S}$  je konečný systém podmnožín množiny  $U$ , ktorého zjednotenie pokrýva celú množinu  $U$ . Hovoríme, že podsystém  $\mathbf{S}'$  systému  $\mathbf{S}$  je *podpokrytie*, ak zjednotenie množín v  $\mathbf{S}'$  pokrýva celú množinu  $U$ . Nasledujúci problém je jeden z prvých problémov, o ktorom bolo dokázané, že je NP-úplný (Karp, 1972).

### PODPOKRYTIE

Inštancia: Dvojica  $(U, \mathbf{S})$  podľa vyššie uvedenej definície a kladné celé číslo  $k \leq |\mathbf{S}|$

Otázka: Existuje podpokrytie pokrytia  $\mathbf{S}$  obsahujúce najviac  $k$  množín?

Je treba poznamenať, že veľkosť inštancie závisí od  $|U| = u$  aj od  $|\mathbf{S}| = n$ . V skutočnosti môžeme popísať  $(U, \mathbf{S})$  pomocou matice veľkosti  $u \times n$ , v ktorej ku každému prvku (riadok) je vyjadrená príslušnosť k množine v systéme  $\mathbf{S}$ . Hodnota  $u \cdot n$  môže byť považovaná za veľkosť inštancie  $(U, \mathbf{S})$  a toto je parameter, ktorého sa týka otázka polynomiálneho algoritmu.

Predchádzajúci problém je zapísaný v existenčnom tvare. Je možné sformulovať optimalizačný problém v nasledujúcom tvare:

PODPOKRYTIE

Inštancia: Dvojica  $(U, \mathbf{S})$  podľa vyššie uvedenej definície a kladné celé číslo  $k \leq |\mathbf{S}|$

Otázka: Aká je veľkosť minimálneho podpokrytia pre  $(U, \mathbf{S})$ ?

Je zrejmé, keď máme odpoveď na druhý problém čas behu je polynomiálny v  $u, n$ , tak je zodpovedaný aj prvý problém v polynomiálnom čase.

## 1.5 OCCAM algoritmy

Nech  $\Omega \rightarrow H$  je reprezentácia booleovských funkcií. Nech  $\|\omega\|$  je miera veľkosti reprezentácie definovanej pre každé  $\omega \in \Omega$ . Pre každé  $r \geq 1$  definujeme

$$\Omega_r = \{\omega \in \Omega \mid \|\omega\| = r\}$$

a nech  $H_r$  označuje podmnožinu  $H$  obsahujúcu tie hypotézy  $h_\omega$ , ktorých minimálna reprezentácia má veľkosť  $r$ .

Hovoríme, že také hypotézy majú veľkosť reprezentácie  $r$ .

Potom  $H$  môže byť odstupňované pomocou veľkosti reprezentácie  $H = \bigcup H_r$ . Učiaci algoritmus  $L$  pre  $H$  má vstup - tréningovú vzorku pre nejakú cieľovú funkciu  $t \in H$ . Predpokladajme, že  $t \in H_r$ ; inak povedané najmenšia reprezentácia  $t$  má veľkosť  $r$ . Výstup z  $L$  bude špecifikovaný reprezentáciou  $\omega \in \Omega_q$ . Potrebujeme uvažovať vzťah medzi  $q$  a  $r$ . Na základe výsledkov predchádzajúceho odseku môže byť ťažké nájsť najmenšiu možnú hodnotu  $q$ , ale môžeme požadovať nájdienie nie najkratšej novej reprezentáciu, ale dostatočne krátku. Táto idea je presnejšie vyjadrená v nasledujúcej definícii Blumera (1987).

**Definícia 1.5.1** *Hovoríme, že učiaci algoritmus  $L$  pre  $H$  je **Occam** vzhľadom na reprezentáciu  $\Omega \rightarrow H$ , ak*

- $L$  je konzistentný
- $k$  danej tréningovej vzorky  $\bar{s}$  dĺžky  $m$  pre cieľovú hypotézu  $t \in H_r$  výstupná hypotéza  $L(\bar{s}) = h_\omega$  je taká, že  $\|\omega\| \leq m^\alpha \cdot r^\beta$ , kde  $0 < \alpha < 1$  a  $\beta \geq 1$  sú konštanty.

Hranica pre  $\|\omega\|$  hovorí, že výstup je komprimovaný vzhľadom na dĺžku tréningovej vzorky a rastie len polynomicne s veľkosťou minimálnej reprezentácie dĺžky cieľovej hypotézy. Podmienka  $\alpha < 1$  znamená, že výstup je vlastne komprimovaný tvar vstupu; ak povolíme  $\alpha = 1$ , potom výstup by bol porovnateľný s veľkosťou tréningovej vzorky, ktorej bitová dĺžka je lineárna v  $m$  a žiadna signifikantná komprimácia by nebola dosiahnutá. Nasledujúca veta ukazuje, že výstup krátkej reprezentácie v tomto zmysle, postačí pre PAC učenie.

Aby sme sformulovali túto vetu, vráťme sa k našej originálnej definícii učiaceho algoritmu s konceptovým a hypotézovým priestorom ako rôznymi. Tento rozdiel je tu použiteľný, pretože sme sa zaujímali o hypotézy v  $H_r$  použitím plného zdroja  $H$ .

**Veta 2** *Nech  $H$  je priestor booleovských funkcií s reprezentáciou  $\Omega \rightarrow H$ , nech  $H = \bigcup H_r$  je odstupňovaný veľkosťou reprezentácie. Ak  $L$  je Occam učiaci algoritmus vzhľadom na danú reprezentáciu, potom pre každé  $r, L$  existuje PAC učiaci algoritmus pre  $(H_r, H)$  so zložitou vzorky  $m_L(H_r, \delta, \epsilon)$  polynomiálnou v  $r, \delta^*$  a  $\epsilon^{-1}$ .*

**Dôkaz:** Predpokladajme, že sú dané  $\delta, \epsilon, \mu$  a  $t \in H_r$ . Pre každé dané  $m$  nech  $L(m, t)$  označuje množinu hypotéz  $h \in H$  takú, že  $h$  je výstup  $L(\bar{s})$  algoritmu  $L$  pre nejakú tréningovú vzorku  $\bar{s}$  dĺžky  $m$  a cieľový koncept  $t$ . Inak povedané,  $L(m, t)$  je efektívny hypotézový priestor pre  $t$ . Podľa 2. podmienky Occam

algoritmu členmi  $L(m, t)$  sú hypotézy  $h_\omega$ , pre ktoré  $\omega$  má najviac  $M = \lfloor m^\alpha \cdot r^\beta \rfloor$  bitov, a celkový počet takých  $\omega$  je najviac  $2^{M+1}$ . Odtiaľ

$$|L(m, t)| \leq 2^{m^\alpha \cdot r^\beta + 1}.$$

Poznamenajme, že hranica závisí len od  $r$  nie od  $t$  samotnej; inak povedané, platí uniformne pre všetky  $t \in H_r$ . Teraz zopakujeme argument daný v predchádzajúcom odseku. Pravdepodobnosť, že ľubovoľná daná  $\epsilon$ -zlá hypotéza z  $H$  súhlasí s  $t$  na tréningovej vzorke dĺžky  $m$  je  $(1 - \epsilon)^m$ . Pretože  $L$  je konzistentný, jeho výstupné hypotézy súhlasia s tréningovou vzorkou a teda pravdepodobnosť, že výstupná hypotéza je  $\epsilon$ -zlá je najviac

$$|L(m, t)| \cdot (1 - \epsilon)^m \leq 2^{m^\alpha \cdot r^\beta + 1} (1 - \epsilon)^m.$$

Zostáva dokázať, že toto môže byť  $< \delta$ , ak vezmeme dostatočne veľké  $m$  a že  $m$  je polynóm v  $r$ ,  $\delta$  a  $\epsilon^{-1}$ . Použitím nerovnosti  $(1 - \epsilon)^m < e^{-\epsilon \cdot m}$  a preusporiadaním dostávame, že  $\epsilon \cdot m \geq A \cdot m^\alpha + B$ , kde  $A = r^\beta \cdot \ln 2$  a  $B = \ln(\frac{2}{\delta})$ .

Pretože  $\alpha > 1$ , podmienka platí ak

$$m^{1-\alpha} \geq (A + B)/\epsilon, \text{ t.j. } m \geq m_0 = \left\lceil \left( \frac{A + B}{\epsilon} \right)^{1/(1-\alpha)} \right\rceil.$$

Inak povedané, výraz pre  $m_0$  je horná hranica pre zložitosť vzorky. Zrejme  $m_0$  je polynóm v  $r$ , pretože  $A$  je  $O(r^\beta)$  a tak  $m_0$  je  $O(r^{\beta/(1-\alpha)})$ ; ďalej  $m_0$  je tiež polynóm v  $\delta^{-1}$  a  $\epsilon^{-1}$ .

□

Zdôrazňujeme znovu podmienku  $\alpha < 1$ ; je zrejmé, že podmienka  $\epsilon \cdot m > A \cdot m^\alpha + B$  môže byť splnená, ak  $\alpha = 1$ .

Existuje bezprostredný dôsledok tejto vety, že ak čas behu Occam algoritmu je polynomiálny v  $m$ , potom čas behu PAC učiaceho algoritmu je polynomiálny v  $r$ ,  $\delta^{-1}$  a  $\epsilon^{-1}$ . Inak povedané, Occam algoritmus  $L$  pre  $H$  PAC učí každý  $H_r$  podľa  $H$  a urobí to efektívne vzhľadom na veľkosť reprezentácie a  $\delta$  a  $\epsilon$ . Poznamenajme, že nemusí nutne platiť, že  $H$  samotný je PAC naučiteľný aj keď by to mohlo tak byť, ak existuje horná hranica na veľkosť reprezentácie hypotézy v  $H$ .

## 1.6 Príklady Occam algoritmov

Predpokladajme, že je daný systém  $\mathbf{S} = \{S_1, S_2, \dots, S_n\}$  konečných množín  $U = \bigcup_{i=1}^n S_i$ . Chceme nájsť najmenšie podpokrytie  $(U, \mathbf{S})$ ; t. j. najmenší podsystem  $\bar{\mathbf{S}}$ , ktorého zjednotením je  $U$ . Videli sme, že tento problém je NP-ťažký. Neznamena to, že neexistujú efektívne prostriedky na dosiahnutie aproximatívneho riešenia pre tento problém. Existuje jednoduchá intuitívna metóda na nájdenie aproximatívneho riešenia, založená na "greedy" metóde, ktorá sa zdá byť veľmi účinnou. Najprv vyberieme množinu  $S_{j_1}$ , ktorá obsahuje najväčší počet prvkov z  $U$  a odstránime ju z  $\bar{\mathbf{S}}$ . Potom vyberieme  $S_{j_2}$ , ktorá obsahuje najväčší počet zvyšných prvkov, atď. Pokračujeme týmto spôsobom, v každom kroku vyberieme množinu, ktorá obsahuje najväčší počet zvyšných prvkov.

### Greedy algoritmus pre minimálne pokrytie

```

set X=U;
while X≠∅ do
begin choose  $S_j$  such that  $|S_j \cap X|$  is maximal;
set X=X- $S_j$ ;
end;
```

Pretože  $\bar{\mathbf{S}}$  pokrýva  $U$ , proces musí skončiť s podpokrytím  $S' = \{S_{j_1}, S_{j_2}, \dots, S_{j_k}\}$ . Samozrejme, veľkosť  $k$  výsledného podpokrytia nebude vo všeobecnosti najmenším možným podpokrytím, ale bolo ukázané (Nigmatullin (1969), Yohansson (1974)), že platí nasledujúci vzťah:  $k \leq l \cdot (\ln |U| + 1)$ , kde  $l$  je veľkosť minimálneho podpokrytia.

Toto poskytuje dobrú hornú hranicu pre **pomer výkonnosti**  $\frac{k}{l}$  a v tomto zmysle greedy algoritmus je dobrá aproximácia pre problém.

Čas behu závisí od  $u = |U|$  a  $n = |\mathbf{S}|$ ,  $u \leq n$ ,  $k \leq \min(u, n) \Rightarrow k \leq n \cdot (\ln u + 1)$ .

Každý výberový krok obsahuje nájdenie maxima z najviac  $n$  celých čísel a vymazanie najviac  $u$  prvkov z každej z najviac  $n$  množín. Počet operácií  $O(u \cdot n)$ . Celkový čas behu je  $O(u \cdot n \cdot \min(u, n))$ .

**Greedy** metóda môže byť použitá na odvodenie algoritmov pre určité triedy booleovských formúl. Podľa Hausslera (1988) ukážeme technicky ako greedy algoritmus pre pokrytie môže byť použitý na priestor  $M_n$ -monočlenov, ukážeme, že výsledný učiaci algoritmus je Occam.

Počiatočná hypotéza je jednočlen bez literálov, identicky 1-ková funkcia. V každom kroku je pridaný 1 literál do priebežnej konjunkcie literálov podľa pravidla založeného na greedy algoritme pre pokrytie. Budeme hovoriť, že literál  $\lambda$  *eliminuje* negatívny príklad  $x$ , ak  $\langle \lambda \rangle(x) = 0$ . Vezmeme prvky, ktoré budú pokryté ako množinu záporných príkladov v danej tréningovej vzorke a pokrývajúce množiny ako množiny záporných príkladov eliminovaných literálmi určitého druhu. V každom stave vyberieme literál, ktorý eliminuje najväčší počet záporných príkladov zo vzorky, pridáme tento literál do formuly a vymažeme príklady, ktoré eliminuje.

Prečo toto pracuje?

Nech  $\bar{s}$  je tréningová vzorka pre monočlen a  $E$  nech je množina príkladov v  $\bar{s}$  takých, že  $E = E^+ \cup E^-$ , Pre ľubovoľný literál  $\lambda$  položíme  $S_\lambda = \{x \in E^- \mid \langle \lambda \rangle(x) = 0\}$

Nakoniec, nech

$$\Lambda = \{\lambda \mid \langle \lambda \rangle(x) = 1 \quad \forall x \in E^+\}$$

**Lema 1** *Systém množín  $\mathbf{S} = \{S_\lambda \mid \lambda \in \Lambda\}$  pokrýva  $E^-$ .*

**Dôkaz:** Pretože  $\bar{s}$  je tréningová vzorka pre monočleny, vieme, že existuje monočlen  $t = \langle \lambda_1 \wedge \dots \wedge \lambda_l \rangle$  taký, že pre  $x \in E$ ,  $t(x)$  je 1 alebo 0 podľa toho či  $x$  je v  $E^+$  alebo v  $E^-$ . Toto implikuje, že  $\lambda_1, \dots, \lambda_l$  všetky patria do  $\Lambda$ . A tiež, že pre ľubovoľné  $x \in E^-$  aspoň jeden z literálov  $\lambda_j$  vyskytujúcich sa v  $t$  je taký, že  $\langle \lambda_j \rangle(x) = 0$ . Inak povedané  $x \in S_{\lambda_j} \in \mathbf{S}$ .

□

**Lema 2** *Ak  $\mathbf{S}' = \{S_{\lambda_1}, \dots, S_{\lambda_k}\}$  je ľubovoľné podpokrytie  $(E^-, \mathbf{S})$ , potom monočlen  $h = \langle \lambda_1 \wedge \dots \wedge \lambda_k \rangle$  je konzistentný s  $\bar{s}$ .*

**Dôkaz:** Predpokladajme, že  $x \in E^+$ . Pretože  $\lambda_1, \dots, \lambda_k$  sú členmi  $\Lambda$ , na  $x$  dávajú hodnotu 1, teda  $h(x) = 1$ . Ak  $x \in E^-$ , potom pretože  $\mathbf{S}'$  je podpokrytie, existuje  $j$ ,  $1 \leq j \leq k$  také, že  $x \in S_{\lambda_j}$ . Teda  $\langle \lambda_j \rangle(x) = 0$ , a teda  $h(x) = 0$ .

□

Tieto lemy ukazujú, že greedy algoritmus pre problém pokrytia môže byť transformovaný na algoritmus pre hľadanie monočlenov konzistentných s danou tréningovou vzorkou.

Aby sme zreteľne videli, že je to Occam algoritmus, uvažujme jeho chovanie sa na tréningovej vzorke pre monočlen  $t$ , ktorého najmenšia reprezentácia je pomocou formuly, ktorá obsahuje  $l$  literálov. Minimálna reprezentácia  $t$  je veľkosti  $r = \lceil l \cdot \log n \rceil$ . Výsledok práce greedy algoritmu pre problém pokrytia implikuje, že počet  $k$  literálov vo výstupnej formule je taký, že  $k \leq l \cdot (\ln |E^-| + 1)$ . Preto veľkosť výstupnej formuly splňuje

$$\|\omega\| = \lceil k \cdot \log n \rceil \leq \lceil l \cdot (\ln |E^-| + 1) \cdot \log n \rceil \leq r \cdot (\ln |E^-| + 1)$$

Platí, že  $|E^-| \leq m$ .

$\|\omega\| \leq r \cdot (\ln m + 1)$ , čo triviálne implikuje Occam kompresnú podmienku.

$\|\omega\| \leq m^\alpha \cdot r^\beta \quad \alpha = \frac{1}{2}, \quad \beta = 1$ .

□

Greedy algoritmus sa líši evidentne od štandardného algoritmu pre monočleny. Namiesto začiatku s identicky nulovou funkciou (konjunkcia  $2n$  literálov) a nasledujúcim vymazaním literálov použitím kladných príkladov greedy algoritmus štartuje s funkciou identicky rovnou 1 (prázdna konjunkcia literálov) a potom pridáva literály použitím záporných príkladov. Teda, pokým štandardný algoritmus je bezpamäťový on-line algoritmus, greedy algoritmus určite taký nie je. Avšak greedy algoritmus ako Occam algoritmus má dôležitú výhodu v tom, že jeho výstupom sú konzistentné hypotézy, ktoré sú relatívne jednoduché.



## 1.7 Epac učenie

Predpokladajme, že  $H = \bigcup H_n$  je hypotézový veľkosťou príkladov odstupňovaný priestor a  $\Omega \rightarrow H$  je reprezentácia pre  $H$ . Potom by sme mohli odstupňovať každé  $H_n$  pomocou veľkosti reprezentácie takto  $H_n = \bigcup H_{n,r}$ , kde  $H_{n,r}$  pozostáva z tých hypotéz  $H_n$ , ktoré majú minimálnu veľkosť reprezentácie  $r$ . Teda

$$H = \bigcup_n \bigcup_r H_{n,r}$$

je dvojite odstupňovaný.

Nech  $L$  je učiaci algoritmus pre  $H$  v obvyklom zmysle, že  $L(\bar{s})$  je v  $H_n$  vždy, keď  $\bar{s}$  je tréningová vzorka pre hypotézy v  $H_n$ . Hovoríme, že  $L$  je *efektívny PAC* alebo *ePAC* ak (Valiant, 1991)

- čas behu  $R_L(m, n)$  je polynomiálny v  $m$  aj v  $n$ ;
- zložitosť vzorky  $m_L(H_{n,r}, \delta, \epsilon)$  je polynomiálna v  $n, r, \delta^*$  a  $\epsilon^{-1}$ .

Teda *ePAC* učiaci algoritmus zaručuje, že vydá pravdepodobnostne aproximovaný správny výstup s časom behu polynomiálnym v  $n, r, \delta^*$  a  $\epsilon^{-1}$ .

Jeden spôsob na zaručenie druhej podmienky je použitie nejakej verzie Occam podmienok. V tomto kontexte hovoríme, že  $L$  je Occam ak podmienky stanovené v definícii pre Occam algoritmus platia pre každé  $H_n$  s konštantami  $\alpha$  a  $\beta$  nezávislými od  $n$ . Potom máme nasledujúci výsledok.

**Veta 3** *Predpokladajme, že hypotézový priestor je  $H = \bigcup H_{n,r}$  ako bolo uvedené vyššie a  $L$  je Occam algoritmus pre učenie  $H_{n,r}$  pomocou  $H_n$  s polynomiálnym časom behu  $R_L(m, n)$ . Potom  $L$  je ePAC.*

**Dôkaz:** Z dôkazu vety vyššie uvedenej máme hornú hranicu

$$m_0(H_r, \delta, \epsilon) = \lceil \left( \frac{A+B}{\epsilon} \right)^{1/(1-\alpha)} \rceil$$

pre zložitosť vzorky algoritmu  $L$  na  $H_{r,n}$ , kde  $A = r^\beta \ln 2$  a  $B = \ln(2/\delta)$ . Ako sme už poznamenali, toto je polynóm v  $r, \delta^*$  a  $\epsilon^{-1}$ . Pretože  $\alpha$  a  $\beta$  sú nezávislé od  $n$ , aj  $m_0(H_r, \delta, \epsilon)$  je nezávislé od  $n$ . Výsledok vyplýva z toho, že horná hranica na čas behu algoritmu  $L$  v *PAC* učení  $H_{n,r}$  je

$$R_L(m_0(H_r, \delta, \epsilon), n)$$

čo je polynóm v  $n, r, \delta^*$  a  $\epsilon^{-1}$ .

□

**Príklad:** Odstupňovaný hypotézový priestor  $M = \bigcup M_n$  monočlenov môže byť odstupňovaný dvojite ako  $M = \bigcup M_{n,r}$ , kde  $M_{n,r}$  pozostáva z tých monočlenov  $n$  premenných, ktoré majú veľkosť reprezentácie  $r$ . V predchádzajúcom odseku bol popísaný algoritmus pre učenie  $M_{n,r}$  pomocou  $M_n$  založený na greedy metóde a ukázali sme, že má Occam vlastnosť pri  $\alpha = 1/2$  a  $\beta = 1$ . Čas behu  $R_L(m, n)$  je  $O(m.n.\min(m, n))$ , čo je určite polynóm v  $m$  a  $n$ . Z toho nám vyplýva, že greedy algoritmus pre  $M$  je *ePAC*.

□

## 1.8 Ďalšie poznámky

Ako sme už uviedli, fakt, že  $C^k$  nie je efektívne naučiteľný vzhľadom na veľkosť príkladov, je výsledok, ktorý závisí od reprezentácie. Keby výstupné hypotézy mohli byť reprezentované iným spôsobom než konjunkcie najviac  $k$  klauzúl, tak generovanie pravdepodobnostne aproximatívne korektných hypotéz by mohlo byť jednoduchšie. Kearns a Valiant (1989) v tomto smere dosiahli veľmi silný výsledok založený na kryptografických predpokladoch.

## 1.9 Úlohy

1. Ukážte, že  $C_n^k \subseteq D_{n,k}$  pre všetky  $k$  a  $n$  a že inklúzia je striktná pre niektoré hodnoty  $n$  a  $k$ .

2. Sformulujte problém pokrytia množiny, ktorý zodpovedá nájdeniu najkratšieho monočlena konzistentného s nasledujúcimi príkladmi.

$$E^+ = \{1110011, 1111011, 1011001, 1011011, 1110001\};$$

$$E^- = \{1010100, 0111011, 0001111, 1001010, 0101111, 1100000\}.$$

Riešte problém pokrytia a napíšte najkratší monočlen.

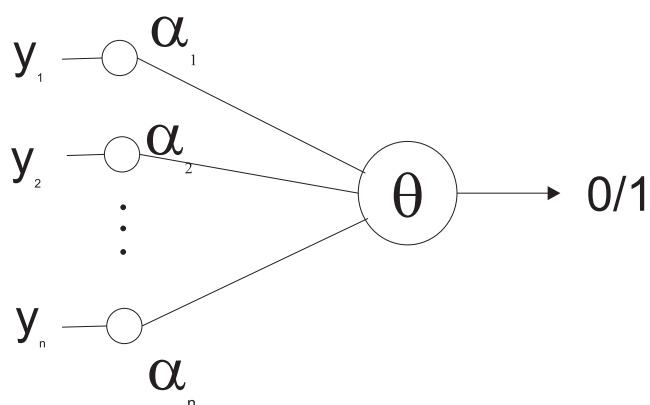
# Kapitola 2

## VC Dimenzia

### 2.1 Rastová funkcia

Vapnik, Chervonenkis - 1971

Najvýznamnejší pojem pre teóriu výpočtového učenia.



Obrázok 2.1: Perceptrón

#### Perceptrón

Pre stav  $\omega = (\alpha_1, \alpha_2, \dots, \alpha_n, \Theta)$  funkcia  $h_\omega \in H$  z  $X = \mathbb{R}^n$  do  $\{0, 1\}$  je daná

$$h_\omega(y) = \begin{cases} 1, & \text{ak } \sum_{i=1}^n \alpha_i y_i \geq \Theta \\ 0, & \text{inak} \end{cases}$$

$\omega \vdash h_\omega$  nie je injekcia;

pre ľubovoľné  $\lambda$ , stav  $\lambda_\omega$  definuje tú istú funkciu.

$\Pi_H(x)$  = počet klasifikácií podľa  $H$ , t.j. počet rôznych vektorov tvaru

$$(h(x_1), \dots, h(x_m))$$

kde  $h$  prebieha všetky hypotézy v  $H$ .

$H$  môže byť nekonečný ...  $H|E_x$ ,  $E_x = \{x_1 \dots x_m\}$  je konečný a má kardinalitu  $\Pi_H(x)$ .

$$\Pi_H(x) \leq 2^m$$

**Definícia 2.1.1** Rastová funkcia je

$$\Pi_H(m) = \max\{\Pi_H(x) : x \in X^m\}$$

**Definícia 2.1.2** Hovoríme, že vzorka  $x$  dĺžky  $m$  je rozbitá podľa  $H$  alebo  $H$  rozbieja  $x$ , ak počet možných klasifikácií podľa  $H$  je  $2^m$ , t.j.  $H$  poskytuje všetky možné klasifikácie.

- Ak príklady v  $x$  nie sú rôzne, tak  $x$  nemôže byť rozbitá.
- Ak príklady v  $x$  sú rôzne,  $x$  je rozbitá podľa  $H \Leftrightarrow$  ak pre ľubovoľné  $S \subseteq E_x$  existuje nejaká hypotéza  $h \in H$  taká, že pre ľubovoľné  $1 \leq i \leq m$  platí

$$h(x_i) = 1 \Leftrightarrow x_i \in S.$$

$S$  je potom podmnožina sústreďujúca kladné príklady.

**Definícia 2.1.3** VC dimenzia  $H$  je maximálna dĺžka vzorky, ktorú  $H$  rozbieja. Ak neexistuje, hovoríme, že VC je  $\infty$ .

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

**Príklady:**

1.  $X = \mathbb{R}$ ,  $H \subseteq$  lúče;  $x = (x_1, \dots, x_m)$  je tréningová vzorka,  $x_1 < x_2 < \dots < x_m$ .  
Pre  $\Theta \in \mathbb{R}$ ,  $r_\Theta = 1 \Leftrightarrow x_i \geq \Theta$ .  
Množina klasifikovaných vektorov ...  $m + 1$ .  
(0...00), (0...01), ... (1...11)  
ľubovoľná vzorka s rôznymi príkladmi má jeden z týchto klas. vektorov (spermutovaný).  $\Rightarrow$   
 $\Pi_H(m) = m + 1$   
V prípade rovnakých príkladov počet klasifikácií je menší.
2.  $r_\Theta$ ,  $H$  je priestor lúčov. Nech je daná tréningová vzorka  $(x_1, x_2)$  dĺžky 2,  $x_1 < x_2$ .  
 $\Rightarrow$  neexistuje lúč taký, že  $h(x_1) = 1$  a  $h(x_2) = 0$ , pretože by muselo platiť  $x_2 < \Theta \leq x_1$ .  
 $\Rightarrow H$  nerozbieja žiadnu vzorku dĺžky 2.  $H$  rozbieja ľubovoľnú vzorku dĺžky 1  $\Rightarrow VCdim(H) = 1$ .
3. Nech  $X = \mathbb{R}^2$ .  $H$  je hypotézový priestor perceptrónu  $P_2$ ,  $x = (x_1, x_2, x_3)$  sú tri nekolineárne rôzne body.  $X$  je rozbitá pomocou  $H \Leftrightarrow$  ak pre ľubovoľnú podmnožinu  $S \subseteq E_x = \{x_1, x_2, x_3\}$  platí, že  $S$  a  $E_x$  sú separovateľné.  
 $\Rightarrow VCdim(H) \geq 3$  ( $P_2$  rozbieja 3 nekolineárne body.)  
Rovnosť dokážeme tým, že nájdeme vzorku, ktorej príklady sa nedajú separovať  $\Rightarrow$  existuje vzorka dĺžky 4, ktorú  $P_2$  nerozbieja.

**Tvrdenie 1** Ak  $H$  je konečný hypotézový priestor, tak  $VCdim(H) \leq \ln |H|$ .

**Dôkaz:**

Počet klasifikácií podľa konečného  $H$  vzorky ľubovoľnej dĺžky je najviac počet rôznych hypotéz v  $H$   
 $\Rightarrow \Pi_H(m) \leq |H|$

VCdim je najväčšie  $d$ , pre ktoré  $\Pi_H(d) = 2^d$ .

$$2^d = \Pi_H(d) \leq |H| \Rightarrow d = \Pi_H(d) \leq \ln |H|.$$

□

**Príklad:**  $VCdim(M_n) \dots |M_n| = 3^n$ ,  $M_n$  sú monočleny.

$VCdim(M_n) \leq (\ln 3) + n$  je horná hranica. čo dolná hranica? Pokúsime sa dokázať, že je rovná  $n$ .

Tvrdíme, že  $M_n$  rozbieja každú vzorku dĺžky  $n$ :

$(e_1, e_2, \dots, e_n)$ , kde  $e_i = (0 \dots 010 \dots 0)$ , kde 1 je na  $i$ -tom mieste,  $1 \leq i \leq n$ .

Predpokladajme, že  $(q_1, \dots, q_n) = q \in \{0, 1\}^n$ . Ukážeme, že existuje  $h \in M_n$  také, že  $h(e_i) = q_i$  pre  $1 \leq i \leq n$ .

Ak  $q = (1 \dots 1)$  vezmeme za  $h$  prázdnu hypotézu.

Ak  $q = (0 \dots 00 \dots 01)$  tak  $h$  vytvoríme ako konjunkciu negácií literálov, pre ktoré  $q_j = 0$ .

$\Rightarrow VCdim(M_n) \geq n$  pre ľubovoľné  $n$ .

## 2.2 VC dimenzia reálnych perceptrónov

**Veta 4** Pre ľubovoľné  $n$  nech  $P_n$  je reálny perceptrón s  $n$  vstupmi. Potom

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

**Dôkaz:**  $P_n$  je v stave  $\omega = (\alpha_1 \alpha_2 \dots \alpha_n \Theta)$   
 $h_\omega$  - funkcia, ktorú perceptrón počíta

$$h_\omega(y) = 1 \Leftrightarrow \alpha_1 y_1 + \dots + \alpha_n y_n \geq \Theta.$$

Označme

$$l_\omega^+ = \{y \in \mathbb{R}^n : \sum_{i=1}^n \alpha_i y_i \geq \Theta\}$$

$$l_\omega^- = \{y \in \mathbb{R}^n : \sum_{i=1}^n \alpha_i y_i < \Theta\}$$

$$l_\omega = \{y \in \mathbb{R}^n : \sum_{i=1}^n \alpha_i y_i = \Theta\}$$

$C \subseteq \mathbb{R}^n$  je **konvexná**, ak pre ľubovoľné  $x, y \in C$  a ľubovoľné reálne číslo  $\lambda$ ,  $0 \leq \lambda \leq 1$ , bod  $\lambda x + (1 - \lambda)y \in C$ .

Prienik ľubovoľných 2 konvexných množín v  $\mathbb{R}^n$  je konvexná množina. Pre ľubovoľnú množinu bodov  $S \subseteq \mathbb{R}^n$  existuje najmenšia konvexná množina obsahujúca  $S$ ;

$conv(S)$  ... konvexný obal  $S$  je prienik všetkých konvexných množín obsahujúcich  $S$ ;

**Pripomeňme Radonovu vetu:**

**Veta 5** Nech  $n$  je kladné celé číslo,  $E$  je ľubovoľná množina  $n+2$  bodov z  $\mathbb{R}^n$ . Potom existuje  $\emptyset \neq S \subseteq E$  taká, že

$$conv(S) \cap conv(E \setminus S) \neq \emptyset.$$

Nech  $x = \underbrace{(x_1 \dots x_{n+2})}_{rozne}$  je ľubovoľná vzorka príkladov z  $\mathbb{R}^n$  dĺžky  $n+2$ .

$E_x$  je množina príkladov v  $x \dots |E_x| = n + 2$ . Podľa Radonovej vety existuje  $S \neq \emptyset$ ,  $S \subseteq E_x$

$$conv(S) \cap conv(E_x \setminus S) \neq \emptyset.$$

Predpokladajme, že existuje  $h_{\omega}$  v  $P_n$  také, že  $S$  je množina pozitívnych príkladov  $h_\omega$  v  $E_x$ .

$$\Rightarrow S \subseteq l_\omega^+, E_x \setminus S \subseteq l_\omega^-$$

Pretože uzavretý polpriestor  $l_\omega^+$  a otvorený  $l_\omega^-$  sú disjunktné a sú konvexné v  $\mathbb{R}^n \Rightarrow conv(S) \subseteq l_\omega^+$ ,  $conv(E_x \setminus S) \subseteq l_\omega^-$

$$conv(S) \cap conv(E_x \setminus S) \subseteq conv(S) \cap conv(E_x \setminus S) = \emptyset$$

$\Rightarrow$  neexistuje taká  $h_\omega$ , a preto  $x$  nie je rozbitá  $p_n$ .

$\Rightarrow$  žiadna vzorka dĺžky  $n+2$  nie je rozbitá pomocou  $p_n$ .

$\Rightarrow VCdim(P_n) \leq n + 1$ .

Opačná nerovnosť:  $o \in \mathbb{R}^n$ ,  $o = (0, 0 \dots 0)$ ,  $e_i = (0 \dots 10 \dots 0)$ ,  $1 \leq i \leq n$ . Ukážeme, že  $P_n$  rozbíja  $x = (o, e_1, \dots, e_n)$  dĺžky  $n + 1$ .

Nech  $S \subseteq E_x = \{o, e_1 \dots e_n\}$ . Nech

$$\alpha_i = \begin{cases} 1, & \text{ak } e_i \in S \\ -1, & \text{ak } e_i \notin S \end{cases}$$

$$\Theta = \begin{cases} -\frac{1}{2}, & \text{ak } o \in S \\ +\frac{1}{2}, & \text{ak } o \notin S \end{cases}$$

Priamou verifikáciou, ak  $\omega = (\alpha_1 \dots \alpha_n \Theta)$  je stav  $P_n$ , potom množina pozitívnych príkladov  $h_\omega$  je práve  $S \Rightarrow VCdim(P_n) \geq n + 1$ .  $\square$

## 2.3 Sauerova lema

Rastová funkcia - miera počtu rôznych klasifikácií vzorky dĺžky  $m$  na pozitívne a negatívne príklady podľa  $H$ , pokiaľ  $VCDim H$  je max. hodnota  $m$ , pre ktorú platí  $\Pi_H(m) = 2^m$ .

**Veta 6 (Sauerova lema):** Nech  $d \geq 0$  a  $m \geq 1$  sú celé kladné čísla, nech  $H$  je hypotézový priestor s  $VCDim(H) = d$ . Potom

$$\Pi_H(m) \leq \underbrace{1 + \binom{m}{1} + \binom{m}{2} + \dots + \binom{m}{d}}_{\Phi(d,m)}.$$

Binomické koeficienty spĺňajú

$$\binom{a}{b} = \binom{a-1}{b} + \binom{a-1}{b-1}$$

Zavedme funkciu:

$$\Phi(0, m) = 1 \quad (m \geq 1);$$

$$\Phi(d, 1) = 2 \quad (d \geq 1)$$

$$\Phi(d, m) = \Phi(d, m-1) + \Phi(d-1, m-1), \quad (d \geq 1, m \geq 2)$$

**Dôkaz:** Ak  $VCDim(H) = d = 0$ , potom príklad x-ľub,  $h(x)$  je rovnaké (buď 0 alebo 1) pre ľub. hypotézu  $h \in H$ .  $\Rightarrow \Pi_H(x) = 1$  pre ľub. vzorku dĺžky  $m \Rightarrow \Pi_H(m) = 1 = \Phi(0, m) \Rightarrow$  veta platí pre  $d = 0$ .

Ak  $m = 1$  a  $d \geq 1 \Rightarrow \Pi_H(1) \leq 2 = \Phi(d, 1)$

Indukciou na  $d + m$ :

- Prípád  $d + m = 2$  je dokázaný.

- Predpokladajme, že veta platí pre  $d + m \leq k$ , kde  $k \geq 2$  a nech  $H$  je hypotézový priestor s  $VCDim = d$ ,  $x$  je tréningová vzorka dĺžky  $m$ , kde  $d + m = k + 1$ .

Prípady  $(d, m) = (0, k + 1)$  a  $(d, m) = (k, 1)$  sú už dokázané.

Nech  $d \geq 1, m \geq 2$ ,  $x$  obsahuje rôzne príklady,  $E$  je množina príkladov v  $x$ ,  $H_E = H|E$  je obmedzenie hypotéz  $H$  na  $E$ .

$\Rightarrow H_E$  je konečný a  $\Pi_H(x) = |H_E|$ .

Potrebuje ukázať, že  $|H_E| \leq \Phi(d, m)$ .

Nech  $F = E \setminus \{x_m\}$ ,  $H_F = H|F$ .

Dve rôzne hypotézy  $h, g \in H_E$  pri obmedzení na  $F$  dávajú tú istú hypotézu z  $H_F$  práve vtedy, keď sa zhodujú na  $F$  a nezhodujú na  $x_m$ .

$H_*$  je množina všetkých hypotéz, ktoré vzniknú takto:

Ak  $h^* \in H_*$ , tak sú možné obe rozšírenia  $h^*$  na funkciu na  $E \dots$  hypotézu z  $H_E$ .  $h^*$  je rozšírenie

$\Rightarrow |H_E| = |H_F| + |H_*|$ .

Nech  $x' = (x_1 \dots x_{m-1})$ . Potom

$$|H_F| = \Pi_H(x') \leq \Pi_H(m-1) \leq \Phi(d, m-1)$$

pretože  $d + (m-1) \leq k$ .

Tvrdíme, že  $VCDim(H_*)$  je najviac  $d-1$ . Ak by  $VCDim(H_*) = d \Rightarrow h^*$  rozbíja nejakú vzorku  $z = (z_1 \dots z_d)$  dĺžky  $d$  príkladov z  $F$ .

Pre

$$h^* \in H_* \begin{cases} h_1 \in H_E \dots & h_1(x_m) = 0 \\ h_2 \in H_E \dots & h_2(x_m) = 1 \end{cases}$$

$\Rightarrow H_E$  a teda  $H$  rozbíja vzorku  $(z_1 \dots z_d x_m)$  dĺžky  $d+1$ , čo je v spore s  $VCDim(H) \leq d$ .

$\Rightarrow VCDim(H_*) \leq d-1$ .

Použitím indukčnej hypotézy

$$|H_*| = \Pi_{H_*}(x'0) \leq \Pi_{H_*}(m-1) \leq \Phi(d-1, m-1)$$

pretože  $(d-1) + (m-1) \leq k$ . Kombináciou oboch výsledkov dostaneme

$$\Pi_H(x) = |H_E| = |H_F| + |H_*| \leq \Phi(d, m-1) + \Phi(d-1, m-1) = \Phi(d, m).$$

□

**Tvrdenie 2** Pre všetky  $m \geq d \geq 1$  platí  $\Phi(d, m) < \left(\frac{e \cdot m}{d}\right)^d$ .

**Tvrdenie 3** Nech  $H$  je ľubovoľný hypotézový priestor obsahujúci aspoň 2 hypotézy a definovaný na konečnom príkladovom priestore  $X$ , potom  $VCdim(H) > \frac{\ln |H|}{1 - \ln |X|}$ .

## 2.4 Úlohy k VC-dim

1. Ukážte, že ak  $X = \mathbb{R}$  a  $H$  je množina všetkých uzavretých intervalov, tak  $\Pi_H(m) = 1 + m + \frac{1}{2}m(m-1)$ .
2. Popíšte expl. hypotézový priestor  $P_1$  a ukážte, že  $VCdim(P_1) = 2$ .
3. Ukážte, že ak  $H$  je hypotézový priestor reálneho perceptoru  $P_2$ , tak  $\Pi_H(4) = 14$ .
4. Nech  $H$  má konečnú  $VCdim$ . Pre  $h \in H$  definujme  $\bar{h}$

$$\bar{h} = 1 \quad \Leftrightarrow \quad h(x) = 0$$

a nech komplement  $H$  je priestor  $\{\bar{h} : h \in H\}$ . Dokážte, že majú oba priestory rovnakú VC dimenziu.

5. Dokážte
  - (a)  $\Phi(d, m) = \Phi(d, m-1) + \Phi(d-1, m-1)$ ,  $d \geq 1, m \geq 2$
  - (b)  $\Phi(d, m) \leq m^d$ ,  $m \geq d > 1$ .
6. Monočlen je monotónny, ak neobsahuje žiadne negované literály. Dokážte, že priestor monotónnych monočlenov definovaný na  $\{0, 1\}^*$  má VC dimenziu práve  $n$ .
7. Hypotézový priestor  $H$  je lineárne usporiadaný, ak má aspoň 2 hypotézy a ak pre ľubovoľné dve  $h, g \in H$  platí  
buď  $h(x) = 1 \Rightarrow g(x) = 1$   
alebo  $g(x) = 1 \Rightarrow h(x) = 1$   
Dokážte, že ak  $H$  je lineárne usporiadaný,  $VCdim(H) = 1$ .
8. Nech  $G_n$  je množina hypotéz z  $P_n$ , pre ktoré nulový vektor  $o$  je negatívny príklad. Predpokladajme, že vzorka  $x = (x_1 \dots x_m)$  je rozbitá pomocou  $G_n$ . Prečo sa žiadne  $x_i$  nesmie rovnať 0? Dokážte, že vzorka  $(x_1 \dots x_m, 0)$  je rozbitá pomocou  $P_n$ . Dokážte, že  $VCdim(G_n) = n$ .