

Obsah

1	Úvod do teórie strojového učenia	2
1.1	Matematický model	2
1.2	Analýza veľkostí chýb	3
1.2.1	Teoretické limity	3
1.2.2	Bias-variance tradeoff	4
1.2.3	Bias-variance tradeoff, verzia 2.	9
1.3	Ako sa vysporiadať s preučeníím/podučeníím?	10
1.3.1	Regularizácia	11
1.3.2	Holdout testing	11
1.4	Cvičenia	12
2	PAC učenie	14
2.1	Konečné množiny hypotéz	15
2.1.1	Základné výsledky	15
2.1.2	Problém konjunkcie	16
2.2	Nekonečné množiny hypotéz	19
2.2.1	Obdĺžniková hra	20
2.2.2	Vapnik-Červonenkisova dimenzia	21

Kapitola 1

Úvod do teórie strojového učenia

Chceme sa naučiť na základe nejakých vstupných dát x predikovať y . Môžeme si to predstaviť tak, že príroda vie poskytovať pozorovania, každé v tvare dvojice (x, y) . Dostali sme od nej sadu t pozorovaní, na základe ktorých chceme navrhnúť nejakú funkciu h , ktorá predpovedá y na základe x . Dobrá funkcia je taká, ktorá je schopná *zovšeobecňovať*, teda sa jej “dobré darí” aj na dátach mimo tréningovej množiny. Proces, ktorým h zostrojíme, si môžeme predstaviť ako algoritmus, ktorý berie ako vstupy tréningové dáta a vráti nám funkciu.

1.1 Matematický model

Z matematického hľadiska, prírodu vieme formalizovať ako pravdepodobnostnú distribúciu P . Množinu všetkých možných x označíme X , množinu možných y označíme Y .

V tejto časti sa nebudeme zaoberať výpočtovou stránkou strojového učenia, od detailov ako časová zložitosť, ..., abstrahujeme. Algoritmus si teda predstavíme iba ako niečo, čo vezme ako vstup tréningové dáta $(x_1, y_1), \dots, (x_t, y_t)$ a na výstup vráti funkciu $h : X \rightarrow Y$. Túto funkciu budeme volať *hypotéza*. Množinu všetkých možných funkcií, ktoré môže náš algoritmus vrátiť, budeme volať *množina hypotéz* a značiť ho H .

Chyba hypotézy. Ako vyjadriť mieru toho, že sa funkcií “dobré darí”? Správime tak pomocou *chybovej funkcie* $\text{err} : Y \times Y \rightarrow \mathbb{R}^+$, ktorej význam je nasledovný: $\text{err}(y, y')$ vyjadruje, ako veľmi sa od seba líšia y a y' . Pomocou tejto funkcie vieme odmerať priemernú chybu hypotézy h , ktorú budeme tiež označovať err , nasledovne:

$$\text{err}(h) = \mathbb{E}_{x,y} [\text{err}(h(x), y)]$$

Pod $\mathbb{E}_{x,y}$ sa rozumie stredná hodnota cez (x, y) z pravdepodobnostnej distribúcie P , teda $(x, y) \sim P$. Pri klasifikácii sa zvykne používať chybová funkcia

$$\text{err}(y, y') = \begin{cases} 0, & \text{ak } y = y' \\ 1, & \text{inak} \end{cases}$$

a potom zrejme

$$\mathbb{E}_{x,y} [\text{err}(h(x), y)] = \mathbb{P}_{x,y} (h(x) \neq y).$$

Pri regresii máme viacero možností, bežné voľby sú kvadratická chyba $(y - y')^2$ a absolútna chyba $|y - y'|$.

Chyba algoritmu. Ako vyjadriť chybu celého učiaceho algoritmu? Uvedomme si, že výstup algoritmu je závislý od tréningových dát $T = \{(x_1, y_1), \dots, (x_t, y_t)\}$, ktoré dostane. Takže výstupná funkcia je od nich závislá, budeme ju označovať \hat{h} . Potom priemerná chyba algoritmu (alebo inak *priemerná chyba priemernej hypotézy*), braná cez všetky možné vzorky tréningových dát, je rovná

$$\mathbb{E}_T [\text{err}(\hat{h})] = \mathbb{E}_T \left[\mathbb{E}_{x,y} [\text{err}(\hat{h}(x), y)] \right].$$

Pod \mathbb{E}_T sa rozumie stredná hodnota cez všetky možné t -tice tréningových dát T , brané nezávisle z pravdepodobnostnej distribúcie P .

Tréningová chyba. Pri vyššie uvedených chybách sme vždy merali vzhľadom na skutočnú distribúciu P . Môže nás ale zaujímať, aká je priemerná chyba hypotézy na tréningových dátach T . Túto chybu budeme označovať $\text{err}_T(h)$, a vypočítame ju ako

$$\text{err}_T(h) = \mathbb{E}_{x_i, y_i} [\text{err}(h(x_i), y_i)] = \frac{1}{t} \cdot \sum_{i=1}^t \text{err}(h(x_i), y_i).$$

Priemerná tréningová chyba z pohľadu algoritmu bude

$$\mathbb{E}_T [\text{err}_T(\hat{h})].$$

V nasledujúcom texte budeme vynechávať premenné, cez ktoré prebiehajú stredné hodnoty, všade tam, kde budú zrejmé z kontextu.

1.2 Analýza veľkostí chýb

V tejto časti sa podrobnejšie pozrieme na to, ako závisia vyššie uvedené štatistiky (tj. priemerná testovacia a tréningová chyba priemernej hypotézy) od veľkosti tréningovej množiny T a od veľkosti množiny hypotéz H .

V celej časti budeme predpokladať, že úloha je regresného charakteru a chyba sa meria ako kvadratická odchýlka, teda

$$\text{err}(y, y') = (y - y')^2.$$

1.2.1 Teoretické limity

Najprv sa ale pozrieme na teoretické limity toho, ako dobrá vôbec môže nejaká funkcia byť. Označme h^\square najlepšiu možnú funkciu, nemusí byť nutne z H . Teda

$$h^\square = \arg \min_h (\text{err}(h)) = \arg \min_h \left(\mathbb{E}_{x,y} [(h(x) - y)^2] \right).$$

Jediné obmedzenia kladené na h sú, že je to funkcia: pre každé x musí vrátiť vždy jednu a tú istú hodnotu. Distribúcia P ale nemusí pre dané x vždy vrátiť to isté y : môže byť zašumená, alebo jednoducho x neobsahuje dostatočnú informáciu. Napríklad, ak podľa plochy bytu určujeme jeho cenu, niektoré dva byty môžu mať rovnakú plochu a predsa rôznu cenu. Ako uvidíme, tento nedeterminizmus je jediný dôvod, prečo hypotéza h^\square nemusí mať nulovú chybu.

Chybu ľubovoľnej hypotézy h vieme upraviť nasledovne:

$$\text{err}(h) = \mathbb{E}_{x,y} [(h(x) - y)^2] \tag{1.1}$$

$$= \mathbb{E}_x \left[\mathbb{E}_{y|x} [(h(x) - y)^2] \right] \tag{1.2}$$

Pozrime sa na vnútornú strednú hodnotu. V nej je x konštanta, a teda aj $h(x) = c$ je konštanta. Aká konštanta minimalizuje danú strednú hodnotu? Nie je ťažké vidieť (napríklad zderivovaním), že minimum sa nadobúda pre $c = E[y]$. Takže

$$h^\square(x) = E_{y|x}[y],$$

a jeho priemerná chyba je

$$\text{err}(h^\square) = E_x \left[E_{y|x} [(y - E[y])^2] \right] = E_x \left[\text{Var}(y) \right].$$

Vidíme teda, že pokiaľ je y jednoznačne určené x -om, tak h^\square bude mať nulovú chybu.

1.2.2 Bias-variance tradeoff

V tomto odseku si ukážeme zaujímavý výsledok, ktorý nám za určitých predpokladov umožňuje vyjadriť chyby pomocou iných, jasnejších veličín: tzv. *výchylky* a *rozptylu*. Označme najlepšiu hypotézu z množiny H ako h^* , teda

$$h^* = \arg \min_h (\text{err}(h)).$$

Budeme upravovať výraz reprezentujúci priemernú chybu priemernej hypotézy \hat{h} .

$$\text{chyba algoritmu} = E_T [\text{err}(\hat{h})] \quad (1.3)$$

$$= E_T \left[E_{x,y} [(\hat{h}(x) - y)^2] \right] \quad (1.4)$$

$$= E_T \left[E_{x,y} \left[\left((\hat{h}(x) - h^*(x)) + (h^*(x) - y) \right)^2 \right] \right] \quad (1.5)$$

V tomto momente prichádza netriviálny technický krok, ktorý si vyžaduje dodatočné predpoklady. Tieto technické detaily prenecháme na koniec časti, sústreďme sa na to hlavné.

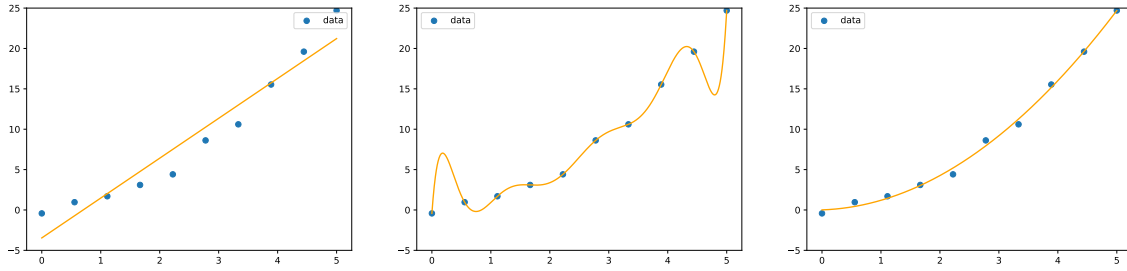
$$\text{chyba algoritmu} = E_T \left[E_{x,y} [(\hat{h}(x) - h^*(x))^2] \right] + E_T \left[E_{x,y} [(h^*(x) - y)^2] \right]$$

Druhý zo sčítancov sa dá ešte zjednodušiť. Keďže h^* ani y nezávisia od tréningových dát, môžeme sa zbaviť vonkajšej strednej hodnoty. Dostávame tak výslednú rovnosť

$$\text{chyba algoritmu} = \underbrace{E_T \left[E_{x,y} [(\hat{h}(x) - h^*(x))^2] \right]}_{\text{rozptyl}} + \underbrace{E_{x,y} [(h^*(x) - y)^2]}_{\text{výchylka}}.$$

Prvý zo sčítancov budeme volať *rozptyl*. Tréningový algoritmus s malým rozptylom vracia funkcie, ktoré sú blízko optima v množine H . Tým, že mu zväčšíme množinu tréningových dát, si veľmi neprilepíme. Naopak, algoritmus s veľkým rozptylom vracia funkcie ďaleko od optima, teoreticky by sme sa teda vedeli k optimu priblížiť tým, že zväčšíme množstvo tréningových dát.

Druhý zo sčítancov budeme volať *výchylka*. Vyjadruje chybu, ktorá je spôsobená tým, že sa náš algoritmus obmedzil na nejakú konkrétnu množinu hypotéz H . Čím väčšia množina hypotéz, tým menšia výchylka (nakolko h^* je najlepšia hypotéza v množine H , jej zväčšením si môžeme iba prilepšiť). Zložitejšia množina hypotéz ale ľahšie “napasuje” na ľubovoľné tréningové dáta. To zvyšuje riziko toho, že výsledná hypotéza bude špecifická pre obdržané dáta, a nebude schopná zovšeobecňovať mimo nich. Je teda potreba väčšieho množstva tréningových dát.



Obr. 1.1: Podučenie, preučenie, akurát.

Ak máme fixné trénovacie dáta T , pri voľbe množiny hypotéz H sa snažíme nájsť kompromis medzi malým rozptylom a malou výchyľkou. Zložitá H bude mať malú výchyľku ale veľký rozptyl, čo vedie k tzv. *preučeniu*. Jednoduchá H bude mať malý rozptyl, ale veľkú výchyľku, tzv. *podučenie*.

Na obrázku 1.1 ilustrujeme oba koncepty: úlohou je modelovať kvadratickú funkciu. Ak za množinu hypotéz zvolíme lineárne funkcie, ich chyby sa nebudú veľmi od seba líšiť, ale všetky budú zlé. Ak za množinu hypotéz zvolíme polynómy nejakého vysokého stupňa, ľahko nájdeme polynóm prechádzajúci cez trénovacie dáta, avšak mimo nich bude dávať výsledky úplne mimo.

Výchyľku vieme upraviť ďalej. Hypotéza h^* ani y nezávisia od trénovacej množiny T . Z ich pohľadu sú teda testovacie dáta x, y a trénovacie dáta x_i, y_i nerozoznatelné. Takže na meranie chyby h^* môžeme použiť trénovacie dáta (berúc v úvahu ich náhodný výber):

$$\text{výchyľka} = \mathbb{E}_T \left[\mathbb{E}_{x_i, y_i} [(h^*(x_i) - y_i)^2] \right] \quad (1.6)$$

$$= \mathbb{E}_T \left[\mathbb{E}_{x_i, y_i} \left[\left((h^*(x_i) - \hat{h}(x_i)) + (\hat{h}(x_i) - y_i) \right)^2 \right] \right] \quad (1.7)$$

Použitím ďalšieho technického kroku dostaneme:

$$\text{výchyľka} = \underbrace{\mathbb{E}_T \left[\mathbb{E}_{x_i, y_i} [(h^*(x_i) - \hat{h}(x_i))^2] \right]}_{\text{trénovací rozptyl}} + \underbrace{\mathbb{E}_T \left[\mathbb{E}_{x_i, y_i} [(\hat{h}(x_i) - y_i)^2] \right]}_{\text{priemerná trénovacia chyba}} \quad (1.8)$$

Prvý zo sčítancov budeme volať *trénovací rozptyl*. Uvedomme si, že pre ľubovoľné trénovacie dáta T platí

$$\text{err}_T(\hat{h}) \leq \text{err}_T(h^*),$$

nakoľko \hat{h} je optimálna hypotéza pre danú množinu trénovacích dát. Hypotéza h síce je najlepšia pre H , trénovacie dáta sú ale len malá vzorka z H . Trénovací rozptyl teda môžeme chápať ako mieru toho, ako veľmi reprezentatívnu vzorku trénovacích dát sme dostali. Čím menší je, tým viac reprezentatívna vzorka je.

Druhý zo sčítancov budeme volať *priemerná trénovacia chyba*. Je to priemerná chyba, ktorej sa dopustí výstu z algoritmu \hat{h} na tých istých dátach, pomocou ktorých sme \hat{h} zostrojili.

Platí

$$\text{priemerná trénovacia chyba} \leq \text{výchyľka} \leq \text{chyba algoritmu}.$$

Na konkrétnych trénovacích dátach ale nemusí platiť, že trénovacia chyba je menšia ako testovacia chyba: mohli sme si (síce s malou pravdepodobnosťou, ale predsa) vytiahnuť zlé trénovacie dáta, ktoré sa výrazne líšia od skutočných dát.

Na základe dosiaľ uvedeného vieme graficky znázorniť, ako sa zhruba správajú rozptyl, výchyľka, trénovací rozptyl a priemerná trénovacia chyba, v závislosti od veľkosti trénovacej množiny (obrázok ??) a od zložitosti množiny hypotéz (obrázok ??).

TODO obrázok kriviek učenia, vysvetlenie

Technické detaily. Nakoniec sa vyjadríme k spomínanému technickému kroku. Začneme jeho znením a potom uvedieme jeho predpoklady.

Veta 1.1. *Predpokladajme, že vstupom do hypotéz sú vektory reálnych čísel (tj. $X = \mathbb{R}^n$), cieľom je predpovedať jedno reálne číslo (tj. $Y = \mathbb{R}$), a že pravdepodobnostné rozdelenie P je spojité.*

Nech množina hypotéz H je uzavretá na lineárne kombinácie a na limity (teda ak postupnosť funkcií v H konverguje, jej limita je tiež v H).

Ďalej predpokladajme, že tréningový algoritmus vždy vráti takú funkciu $\hat{h} \in H$, ktorá minimalizuje tréningovú chybu. Inak zapísané,

$$\hat{h} = \arg \min_{h \in H} \left(\mathbb{E}_T [\text{err}_T(h)] \right).$$

Potom platí

$$\mathbb{E}_T \left[\mathbb{E}_{x,y} \left[\left((\hat{h}(x) - h^*(x)) + (h^*(x) - y) \right)^2 \right] \right] = \mathbb{E}_T \left[\mathbb{E}_{x,y} \left[(\hat{h}(x) - h^*(x))^2 \right] \right] + \mathbb{E}_T \left[\mathbb{E}_{x,y} \left[(h^*(x) - y)^2 \right] \right]$$

Poznámka 1.1. Dokazovaná rovnosť je ekvivalentná s nasledovnou, stručnejšou:

$$\mathbb{E}_T \left[\mathbb{E}_{x,y} \left[(\hat{h}(x) - h^*(x)) \cdot (h^*(x) - y) \right] \right] = 0.$$

Túto kratšiu verziu získame roznásobením a použitím linearity strednej hodnoty. V dôkaze budeme dokazovať túto rovnosť.

Poznámka 1.2. Všimnite si, že potrebujeme uzavretosť množiny H na limity na to, aby vôbec $\arg \min_{h \in H}(\dots)$ existovalo. Vo všeobecnosti nemusí existovať taká funkcia, ale môže existovať nekonečná postupnosť funkcií, každá ďalšia lepšia, ako tá predchádzajúca. (Inak povedané, neexistuje minimum, iba infimum.)

Poznámka 1.3. Je namieste otázka, či je $\arg \min_{h \in H}(\dots)$ dobre definované, teda či je taká funkcia h práve jedna. Za chvíľu uvidíme, že naše predpoklady to zaručujú.

Poznámka 1.4. Veta by sa dala rozšíriť aj na iné množiny X, Y , napríklad keď predpovedaná premenná je vektor ($Y = \mathbb{R}^m$), ... Možno ani P nemusí byť spojitá. Pre jednoduchosť argumentu ale budeme uvažovať vetu tak, ako je popísaná vyššie.

Poznámka 1.5. Predpoklady vety sú značne obmedzujúce. Napríklad si uvedomte, že ju nie je možné použiť na klasifikáciu, či dokonca ani na ľubovoľnú ohraničenú regresiu (kde rozumné hodnoty y sú ohraničené). Ale taká je teória.

Pri našom dôkaze využijeme niekoľko vlastností funkcií, ktoré uvádzame v nasledujúcom odseku. Skúsený čitateľ-matematik ho môže preskočiť.

Definícia 1. (Skalárny súčin.) Nech f, g sú funkcie z X do \mathbb{R} , z nejakej príjemne sa správajúcej množiny funkcií (tj. rovnomerne spojitá, ..., čokoľvek, aby nasledujúce argumenty prešli). Definujeme ich skalárny súčin $\langle \cdot, \cdot \rangle$ ako

$$\langle f, g \rangle = \int f(x) \cdot g(x) \, d\rho x \tag{1.9}$$

$$= \mathbb{E}_x [f(x) \cdot g(x)], \tag{1.10}$$

kde ρ je hustota pravdepodobnosti distribúcie P . Rozmyslite si, že takto definovaný skalárny súčin má všetky vlastnosti, ktoré sa bežne požadujú od skalárnych súčinov:

- Je symetrický od svojich argumentov, teda $\langle f, g \rangle = \langle g, f \rangle$.
- Je lineárny: $\langle f, g + h \rangle = \langle f, g \rangle + \langle f, h \rangle$ a tiež $\langle k \cdot f, g \rangle = k \cdot \langle f, g \rangle$.
- $\langle f, f \rangle \geq 0$ pre ľubovoľné f , pričom rovnosť nastáva práve vtedy, keď je f konštantne nulové.

Definícia 2. (Kolmost.) Dve funkcie f, g sú na seba kolmé, ak ich skalárny súčin je 0. Značíme $f \perp g$.

Definícia 3. (Norma.) Podľa skalárneho súčinu definujeme normu funkcie (jej “dĺžku”):

$$\|f\| = \sqrt{\langle f, f \rangle} = \sqrt{\mathbb{E}_x[f^2(x)]}$$

Splňa *trojuholníkovú nerovnosť*: pre ľubovoľné funkcie f, g platí

$$\|f\| + \|g\| \geq \|f + g\|.$$

Definuje nám teda (euklidovskú) metriku nad funkciami, podľa ktorej definujeme limity a konvergenciu.

Lemma 1.2. (Pytagorova veta.) Nech $f \perp g$. Potom platí:

$$\|f\|^2 + \|g\|^2 = \|f + g\|^2$$

Dôkaz. Pozrime sa na pravú stranu. Iba v nej zapíšeme normu ako skalárny súčin a využijeme jeho linearitu a symetriu:

$$\|f + g\|^2 = \langle f + g, f + g \rangle \tag{1.11}$$

$$= \langle f, f \rangle + \langle g, g \rangle + 2 \cdot \langle f, g \rangle \tag{1.12}$$

Pretože $f \perp g$, posledný sčítanec je nulový, čím dostávame dokazované tvrdenie. \square

Definícia 4. (Projekcia na množinu.) Projekciu funkcie f na množinu H budeme označovať f_H a budeme pod ňou rozumieť nasledovný výraz:

$$f_H = \arg \min_{h \in H} d(f, h)$$

Poznámka 1.6. Ako sa už spomínalo, nie je zrejmé, že projekcia je dobre definovaná. Preto v nasledujúcej lemme definujeme f_H trochu iným spôsobom, ako jednu z možno viacerých funkcií, ktoré minimalizujú vzdialenosť k H .

Lemma 1.3. (Kolmost' projekcie.) Pre ľubovoľnú funkciu $h \in H$ platí $h \perp f - f_H$.

Dôkaz. Sporom, predpokladajme, že $h \not\perp f - f_H$. Takže $\langle h, f - f_H \rangle \neq 0$. Ukážeme, že potom existuje v H funkcia, ktorá je k funkcii f bližšie, ako funkcia f_H . To bude hľadaný spor s definíciou f_H .

Pozrime sa na všetky funkcie, ktoré ležia na priamke $f_H + \Delta \cdot h$. Tieto funkcie sú v množine H , pretože $f_H, h \in H$ a množina H je uzavretá na lineárne kombinácie. Každú z týchto funkcií vieme asociovať s jedným reálnym číslom Δ . Pozrime sa na ich vzdialenosti od funkcie f , vyjadrené ako funkcia od Δ :

$$\text{dist}(\Delta) = d(f, f_H + \Delta \cdot h) \tag{1.13}$$

$$= \langle (f - f_H) + \Delta \cdot h, (f - f_H) + \Delta \cdot h \rangle \tag{1.14}$$

$$= \langle f - f_H, f - f_H \rangle + 2\Delta \cdot \langle h, f - f_H \rangle + \Delta^2 \cdot \langle h, h \rangle \tag{1.15}$$

Pozrime sa na deriváciu tejto funkcie. Podľa definície f_H by malo byť $f - f_H$ najkratšie možné, teda pre $\Delta = 0$ by mala funkcia dist nadobúdať minimum, a teda mať tam nulovú deriváciu. Uvidíme, že tomu tak nie je:

$$\frac{\partial \text{dist}}{\partial \Delta}(0) = \lim_{\Delta \rightarrow 0} \left(\frac{\text{dist}(\Delta) - \text{dist}(0)}{\Delta} \right) \quad (1.16)$$

$$= \lim_{\Delta \rightarrow 0} \left(\frac{2\Delta \cdot \langle h, f - f_H \rangle + \Delta^2 \cdot \langle h, h \rangle}{\Delta} \right) \quad (1.17)$$

$$= 2 \cdot \langle h, f - f_H \rangle \quad (1.18)$$

To je nenulové, nakoľko $h \not\perp f - f_H$. Čo je hľadaný spor. \square

Lemma 1.4. *Projekcia na množinu H je dobre definovaná, teda vždy existuje nanajvýš jedna funkcia $f_H \in H$, ktorá minimalizuje vzdialenosť k f .*

Dôkaz. Predpokladajme, že také funkcie sú dve, označme ich g, h . Ukážeme, že potom nutne $g = h$.

Podľa predchádzajúcej lemmy platí

$$f - g \perp g, \text{ odkiaľ } \langle f - g, g \rangle = 0 \quad (1.19)$$

$$\langle f - h, g \rangle = 0 \quad (1.20)$$

$$\langle f - g, h \rangle = 0 \quad (1.21)$$

$$\langle f - h, h \rangle = 0 \quad (1.22)$$

Z týchto rovností dostaneme

$$\langle g, g \rangle = \langle g, h \rangle = \langle h, g \rangle = \langle h, h \rangle.$$

Nakoniec, pozrime sa na normu funkcie $g - h$:

$$\|g - h\| = \sqrt{\langle g - h, g - h \rangle} \quad (1.23)$$

$$= \sqrt{\langle g, g \rangle - \langle g, h \rangle - \langle h, g \rangle + \langle h, h \rangle} \quad (1.24)$$

$$= 0 \quad (1.25)$$

To môže nastať jedine vtedy, keď $g = h$. \square

Lemma 1.5. *Hypotéza h^* je projekciou h^\square na H , teda $h^* = h_H^\square$.*

Dôkaz. Vychádzajme z definície h^* .

$$h^* = \arg \min_{h \in H} \mathbb{E}_{x,y} [(h(x) - y)^2] \quad (1.26)$$

$$= \arg \min_{h \in H} \mathbb{E}_{x,y} [((h(x) - h^\square(x)) + (h^\square(x) - y))^2] \quad (1.27)$$

$$= \arg \min_{h \in H} \left(\begin{array}{l} \mathbb{E}_{x,y} [(h(x) - h^\square(x))^2] \\ + \mathbb{E}_{x,y} [(h^\square(x) - y)^2] \\ + 2 \cdot \mathbb{E}_{x,y} [(h(x) - h^\square(x)) \cdot (h^\square(x) - y)] \end{array} \right) \quad (1.28)$$

Druhý sčítanec je konštanta, teda nám $\arg \min$ nijak neovplyvňuje. Tretí sčítanec vieme upraviť nasledovne:

$$\text{tretí sčítanec} = \mathbb{E}_x \left[\mathbb{E}_{y|x} [(h(x) - h^\square(x)) \cdot (h^\square(x) - y)] \right] \quad (1.29)$$

$$= \mathbb{E}_x \left[(h(x) - h^\square(x)) \cdot \mathbb{E}_{y|x} [h^\square(x) - y] \right] \quad (1.30)$$

$$= 0 \quad (1.31)$$

A teda je to tiež konštanta. Dostávame tak

$$h^* = \arg \min_{h \in H} \mathbb{E}_{x,y} [(h(x) - h^\square(x))^2],$$

čo je presne definícia projekcie h^\square na množinu H . \square

Vyzbrojení týmito znalosťami, môžeme sa vrhnúť na dôkaz vety 1.1. Pripomeňme si ešte pred tým dokazovanú rovnosť:

$$\mathbb{E}_T \left[\mathbb{E}_{x,y} [(\hat{h}(x) - h^*(x)) \cdot (h^*(x) - y)] \right] = 0.$$

Dôkaz. Ľavú stranu dokazovanej rovnosti vieme prepísať do nasledovného, ekvivalentného tvaru:

$$= \mathbb{E}_T \left[\mathbb{E}_{x,y} [(\hat{h}(x) - h^*(x)) \cdot ((h^*(x) - h^\square(x)) + (h^\square(x) - y))] \right]$$

Vieme, že $\varepsilon := h^\square(x) - y$ sa správa pre dané x ako náhodná premenná, ktorá má strednú hodnotu 0 a je nezávislá od ostatných premenných vystupujúcich vo výraze. Z výrazu ju teda môžeme vyhodíť, dostaneme tak

$$= \mathbb{E}_T \left[\mathbb{E}_{x,y} [(\hat{h}(x) - h^*(x)) \cdot (h^*(x) - h^\square(x))] \right]$$

Stačí nám teda dokázať $\hat{h} - h^* \perp h^* - h^\square$. To ale vyplýva z lemy o kolmosti projekcie (1.3). Overíme, že jej podmienky sú splnené: z uzavretosti na lineárne kombinácie platí $\hat{h} - h^* \in H$, a podľa lemy 1.5 platí $h^* = h_H^\square$, odkiaľ $h^* - h^\square = -(h^\square - h_H^\square)$. Záporné znamienko na kolmosti nič nemení. \square

Podobným spôsobom sa dá dokázať aj korektnosť druhého technického kroku. To prenechávame čitateľovi ako cvičenie.

1.2.3 Bias-variance tradeoff, verzia 2.

V literatúre pod názvom *bias-variance tradeoff* vystupuje aj podobný, ale predsa odlišný výsledok, ako bolo uvedené vyššie. Ukážeme a odvodíme si ho.

Veta 1.6. *Nech $y : X \rightarrow \mathbb{R}$ je funkcia, ktorú sa snažíme modelovať. Predpokladajme, že sa dá rozložiť na časti: $y = f(x) + \varepsilon$, kde ε hrá rolu šumu: je nezávislý od všetkého a $\mathbb{E}[\varepsilon] = 0$. Označíme jeho pravdepodobnostnú distribúciu E .*

Nech výstupom tréningového algoritmu je \hat{f} . Za chybovú funkciu zvolíme kvadratickú chybu. Chybu algoritmu vieme teda vypočítavať nasledovne:

$$\text{chyba algoritmu} = \mathbb{E}_{(x,y) \sim P, T \sim P^t, \varepsilon \sim E} [(\hat{f}(x) - y)^2].$$

Tvrdíme, že sa dá rozložiť na tri nasledovné časti:

$$\text{chyba algoritmu} = \underbrace{\text{Var}(\hat{f}(x) - f(x))}_{\text{rozptyl}} + \underbrace{(\mathbb{E}[\hat{f}(x)] - \mathbb{E}[f(x)])^2}_{\text{výchylka}^2} + \underbrace{\text{Var}(\varepsilon)}_{\text{šum}}$$

Poznámka 1.7. V poslednej rovnici sme kvôli stručnosti vynechali pri stredných hodnotách a rozptyloch premenné a distribúcie, z ktorých ich berieme. V dôkaze budeme vždy brať všetky premenné z ich príslušných distribúcií.

Poznámka 1.8. Funkcia f hrá v podstate tú istú rolu, čo najlepšia možná hypotéza spomedzi všetkých funkcií (nielen tých v množine hypotéz), h^\square .

Poznámka 1.9. V tomto znení bias-variance tradeoff-u názvy *rozptyl* a *výchylka* zodpovedajú príslušným štatistickým/pravdepodobnostným pojmom.

Poznámka 1.10. Na rozdiel od predchádzajúcej verzie bias-variance tradeoff-u, tu nebudeme potrebovať žiadne dodatočné predpoklady od algoritmu ani od jeho množiny hypotéz. (Nemusí teda vracieť hypotézu, ktorá je spomedzi hypotéz v H najlepšia na daných tréningových dátach. Takisto od množiny hypotéz nepožadujeme žiadne vlastnosti.)

Dôkaz. Upravujeme pôvodný výraz.

$$\text{chyba algoritmu} = \mathbb{E} \left[(\hat{f}(x) - y)^2 \right] \quad (1.32)$$

$$= \mathbb{E} \left[(\hat{f}(x) - f(x) - \varepsilon)^2 \right] \quad (1.33)$$

$$= \mathbb{E} \left[(\hat{f}(x) - f(x))^2 \right] + \mathbb{E} [\varepsilon^2] - 2 \cdot \mathbb{E} [\varepsilon \cdot (\hat{f}(x) - f(x))] \quad (1.34)$$

$$= \mathbb{E} \left[(\hat{f}(x) - f(x))^2 \right] + \mathbb{E} [\varepsilon^2] \quad (1.35)$$

Výraz sme upravili, roznásobili a využili linearitu strednej hodnoty. V poslednom kroku sme použili $\mathbb{E}[ab] = \mathbb{E}[a] \cdot \mathbb{E}[b]$, ktorý platí pre ľubovoľné nezávislé premenné, s $a := \varepsilon$, $b := \hat{f}(x) - f(x)$. Zamerajme sa ďalej na prvý sčítanec.

$$\text{prvý sčítanec} = \mathbb{E} \left[(\hat{f}(x) - f(x))^2 \right] \quad (1.36)$$

$$= \mathbb{E}[\hat{f}(x)^2] + \mathbb{E}[f(x)^2] - 2 \cdot \mathbb{E}[\hat{f}(x) \cdot f(x)] \quad (1.37)$$

$$= (\text{Var}(\hat{f}(x)) + \mathbb{E}[\hat{f}(x)]^2) + (\text{Var}(f(x)) + \mathbb{E}[f(x)]^2) - 2 \cdot \mathbb{E}[\hat{f}(x) \cdot f(x)] \quad (1.38)$$

V poslednom kroku sme využili vzťah $\text{Var}(a) = \mathbb{E}[a^2] - \mathbb{E}[a]^2$. Pokračujme ďalej v úpravách.

$$\begin{aligned} \text{prvý sčítanec} &= \text{Var}(\hat{f}(x)) + \text{Var}(f(x)) + (\mathbb{E}[\hat{f}(x)] - \mathbb{E}[f(x)])^2 \\ &\quad + 2 \cdot \mathbb{E}[\hat{f}(x)] \cdot \mathbb{E}[f(x)] - 2 \cdot \mathbb{E}[\hat{f}(x) \cdot f(x)] \end{aligned} \quad (1.39)$$

$$= \text{Var}(\hat{f}(x)) + \text{Var}(f(x)) + (\mathbb{E}[\hat{f}(x)] - \mathbb{E}[f(x)])^2 - 2 \cdot \text{Cov}(\hat{f}(x), f(x)) \quad (1.40)$$

$$= \text{Var}(\hat{f}(x) - f(x)) + (\mathbb{E}[\hat{f}(x)] - \mathbb{E}[f(x)])^2 \quad (1.41)$$

Využili sme najprv vzťah $\text{Cov}(a, b) = \mathbb{E}[ab] - \mathbb{E}[a] \cdot \mathbb{E}[b]$, a potom $\text{Var}(a - b) = \text{Var}(a) + \text{Var}(b) - 2 \cdot \text{Cov}(a, b)$. Keď to teda celé dáme do jednej rovnice, dostaneme

$$\text{chyba algoritmu} = \underbrace{\text{Var}(\hat{f}(x) - f(x))}_{\text{rozptyl}} + \underbrace{(\mathbb{E}[\hat{f}(x)] - \mathbb{E}[f(x)])^2}_{\text{výchylka}^2} + \underbrace{\text{Var}(\varepsilon)}_{\text{šum}}$$

□

1.3 Ako sa vysporiadať s preučeníím/podučením?

V tejto časti sa budeme zaoberať otázkou: “Ako zvoliť vhodne zložitú množinu hypotéz?” Ako sme videli, príliš jednoduché hypotézy vedú k síce malému rozptylu, ale veľkej výchylke, zatiaľ čo príliš zložené hypotézy vedú k malej výchylke, ale veľkému rozptylu.

Predstavme si, že máme na výber z viacerých množín hypotéz, čím ďalej tým zložitejších:

$$H_1 \subseteq H_2 \subseteq H_3 \subseteq \dots$$

Z ktorej množiny hypotéz chceme vybrať?

Pri tréningu sa snažíme nájsť hypotézu h , ktorá minimalizuje chybu na tréningových dátach $\text{err}_T(h)$. Táto chyba nám ale nehovorí nič o rozptyle. Ak by sme si graficky znázornili testovacie a tréningové chyby najlepších hypotéz z jednotlivých množín, vyzeralo by to zhruba ako na obrázku ??.

TODO obrázok

1.3.1 Regularizácia

V tomto prístupe do minimalizovaného výrazu umelo pridáme člen, ktorý aproximuje rozptyl: $\text{pokuta}(h)$, pričom z čím zložitejšej množiny hypotéza h je, tým väčšia pokuta. Takže výstupom algoritmu je

$$\hat{h} = \arg \min_{h \in H_1 \cup H_2 \cup \dots} (\text{err}_T(h) + \text{pokuta}(h)).$$

Uvedomte si, že vrámci jednej množiny H_i ostáva ako najlepšia hypotéza stále tá istá, ako pred zavedením pokuty. V jednej množine sú totiž všetky hypotézy penalizované rovnako, nerobí to teda rozdiel. Penalizácia nám ale umožňuje “férovejšie” porovnávať hypotézy z rôznych množín, nakoľko bez pokuty by na tom boli (neprávom) lepšie zložitejšie hypotézy.

Množiny H_i nemusia byť explicitné, môžu byť implicitne skryté v tom, aký tvar má výraz $\text{pokuta}(h)$. Do jednej množiny patria tie hypotézy, ktoré majú rovnakú penalizáciu.

Uvedieme si niekoľko príkladov výrazov, ktoré môžu byť použité ako pokuta. Vo všetkých prípadoch je pokuta je parametrizovaná reálnym parametrom λ hovoriacim, ako veľké pokuty chceme udeľovať. Budeme predpokladať, že celá množina hypotéz, z ktorej vyberáme (tj. $H_1 \cup H_2 \cup \dots$) je množina lineárnych funkcií $\mathbb{R}^n \rightarrow \mathbb{R}$. Hypotézy majú teda tvar

$$h(x) = a_1x_1 + a_2x_2 + \dots + a_nx_n.$$

- L_2 regularizácia (známa aj ako *ridge regression*). V nej penalizujeme veľké váhy: čím dôležitejší atribút, tým väčšie váhy si môže dovoliť mať.

$$\text{pokuta}(h) = \lambda \cdot \|(a_1, a_2, \dots, a_n)\|^2 = \lambda \cdot (a_1^2 + a_2^2 + \dots + a_n^2)$$

- L_1 regularizácia (známa aj ako *lasso*). Opäť penalizujeme veľké váhy, avšak pokuta je iná.

$$\text{pokuta}(h) = \lambda \cdot (|a_1| + |a_2| + \dots + |a_n|)$$

Táto pokuta “tlačí” nepotrebné atribúty do nuly, čo je výhodné: nulové atribúty vôbec nemusíme uvažovať, čo nám zníži výpočtové nároky. Na druhej strane sa táto pokuta neoptimalizuje ľahko (z optimalizačného hľadiska).

1.3.2 Holdout testing

V tomto prístupe si rozdelíme dostupné dáta na dve časti: trénovaciu množinu a *validačnú množinu*. Pomocou validačnej množiny budeme odhadovať testovacie chyby pre jednotlivé množiny hypotéz, na základe ktorých zistíme, ktorá množina hypotéz je pre náš problém najvhodnejšia. Konkrétnejšie:

1. Trénovaciu množinu použijeme na natrénovanie hypotéz z jednotlivých množín.
2. Ako odhad testovacej chyby jednotlivých hypotéz použijeme ich chybu na validačnej množine. Podľa týchto odhadov zistíme, ktorá množina hypotéz je pre náš problém najvhodnejšia.
3. Použijeme všetky dáta, ktoré máme k dispozícii (tj. z trénovacej aj validačnej množiny), na natrénovanie najlepšej možnej hypotézy. Berieme samozrejme v úvahu iba hypotézy z najlepšej množiny hypotéz. Výsledná hypotéza je výstupom.

V kroku 2 je dôležité, aby bola validačná množina nezávislá od trénovacej. Prečo je to dôležité? Môžeme uvažovať extrémny prípad, keď je validačná množina totožná s trénovacou. Potom ale ako náš “odhad” dostaneme trénovaciu chybu, ktorá rozhodne nie je dobrým odhadom testovacej chyby. Nezávislosť nám teda zaručuje, že odhad získaný na validačnej množine je dobrý.

Poznámka 1.11. Treba podotknúť, že chyba na validačnej množine je iba odhad. Ak by sme mali dostatočne veľa rôznych modelov, z ktorých vyberáme (tj. H_1, H_2, \dots), pre niektorý z nich by sa mohlo stať čistou náhodou, že jeho validačná chyba je nízka, napriek tomu, že jeho skutočná testovacia chyba je vysoká.

Toto je podobné, ako pri overovaní vedeckých hypotéz: napríklad si predstavme 10 hypotéz, každá s 10% šancou, že bude konzistentná s nazbieranými dátami. Potom môžeme očakávať, že jedna z nich bude konzistentná s nazbieranými dátami, napriek tomu, že všetky sú úplne náhodné.

V oboch prípadoch sa to dá samozrejme eliminovať jedným spôsobom: viac dát.

***k*-fold evaluation.** Pri tomto prístupe je dôležité mať dobrý odhad testovacej chyby pre jednotlivé množiny hypotéz. Dát ale môže byť málo, a v takom prípade môže byť odhad nestabilný/nepresný. Môžeme ale experiment zopakovať niekoľkokrát: v každej iterácii teda zvolíme inú tréningovú množinu a inú validačnú množinu, a dostaneme iný odhad testovacej chyby. Keď tieto odhady spriemerujeme, dostaneme oveľa presnejší odhad, ako keby sme vykonali iba jednu iteráciu.

V tomto konkrétnom prístupe je k iterácií, a množiny sa volia nasledovne: všetky dáta sa rozdelia na k zhruba rovnako veľkých a navzájom nezávislých množín K_1, K_2, \dots, K_k . Následne, v iterácii i sa ako validačné dáta použije množina K_i . Všetko ostatné budú tréningové dáta.

Testovacia množina. Ak chceme zmerať testovaciu chybu výstupnej hypotézy, musíme si na to rezervovať ďalšiu časť dát: *testovaciu množinu*. Tú nepoužívame ani pri tréningu, ani pri validácii. Iba úplne na konci celého procesu na nej vypočítame chybu našej hypotézy.

Poznámka 1.12. “Ak sa mi model trénuje príliš dobre, väčšinou to je veľmi zle!”

1.4 Cvičenia

V nasledujúcich dvoch cvičeniach môžete predpokladať, že tréningový algoritmus vždy vráti nejakú funkciu (nemusí byť len jedna) s minimálnou chybou na tréningových dátach.

1.1. Je rozumné predpokladať (a všade vyššie sme tak činili), že s väčším množstvom tréningových dát sa nám bude testovacia chyba znižovať. Sú ale zostrojiteľné situácie, kedy tomu tak nie je. Nájdite jednu takú situáciu.

Konkrétne, nájdite takú množinu hypotéz H funkciu $\mathbb{R}^n \rightarrow \mathbb{R}$ a pravdepodobnostné rozdelenie P , pre ktoré sa nám bude testovacia chyba so zvyšujúcim sa počtom tréningových chýb *zvyšovať*. Jediná podmienka je kladená na množinu hypotéz: pre každú možnú tréningovú množinu T musí existovať hypotéza v H , ktorá minimalizuje tréningovú chybu. (Teda vždy musí existovať minimum, vo všeobecnosti existuje iba infimum.)

1.2. Za určitých podmienok ale skutočne platí, že viac tréningových dát nám vo veľkom merítke neuškodí. Nech množina hypotéz H je konečná a všetky jej funkcie ($\mathbb{R}^n \rightarrow \mathbb{R}$) sú ohraničené. Dokážte, keď $t \rightarrow \infty$, tak chyba hypotézy \hat{h} sa bude blížiť k chybe najlepšej možnej hypotézy h^* . Konkrétnejšie, dokážte

$$\lim_{t \rightarrow \infty} \mathbb{E}_T [\text{err}(\hat{h}) - \text{err}(h^*)] = 0.$$

1.3. Dokážte korektnosť druhého technického kroku, v odvodení rozkladu výchyľky na tréningovú rozptyl a priemernú tréningovú chybu. Konkrétnejšie, dokážte

$$\mathbb{E}_T \left[\mathbb{E}_{x_i, y_i} \left[(h^*(x_i) - \hat{h}(x_i)) \cdot (\hat{h}(x_i) - y_i) \right] \right] = 0.$$

Predpoklady kladené na množinu hypotéz sú rovnaké: musí byť uzavretá na lineárne kombinácie a na limity.

1.4. Jednou výhodou L_2 regularizácie oproti L_1 regularizácie je, že sa ľahšie minimalizuje výsledný výraz. Ako príklad uvidíme lineárnu regresiu. V nej je hypotéza parametrizovaná stĺpcovým vektorom $\theta = (\theta_1, \dots, \theta_n)^T$. Výstupom pre vstup $x = (x_1, \dots, x_n)$ je $x \cdot \theta$.

Označme X maticu, ktorej riadkami sú vstupe jednotlivých tréningových príkladov. Ďalej nech y je stĺpcový vektor cieľových výstupov na jednotlivých príkladoch. Ako určite vieme, optimálnymi parametrami lineárnej hypotézy je taký stĺpcový vektor θ , ktorý je riešením rovnice

$$X^T X \cdot \theta = X^T y.$$

Dokáže, že keď k minimalizovanej hodnote pridáme pokutu vo forme $\lambda \cdot \|\theta\|^2$, tak sa optimálnymi parametrami stane θ riešiaci rovnicu

$$(X^T X + \lambda I) \cdot \theta = X^T y.$$

Rozmyslite si taktiež, že takýto explicitný “vzorec” nie je možné priamočiaro získať pre L_1 regularizáciu.

Kapitola 2

PAC učenie

V tejto kapitole sa budeme zaoberať otázkou toho, ako závisí chyba algoritmu od veľkosti trénovacej množiny. Konkrétne sa budeme zaoberať otázkami ako:

- “Pri danej veľkosti trénovacej množiny t , akú chybu algoritmu môžeme očakávať?”
- “Pri danom t , s akou pravdepodobnosťou nám algoritmus vráti hypotézu, ktorej chyba je menšia ako ε ?”

Na základe odpovedí na tieto dve otázky potom budeme schopní zodpovedať nasledovné, príbuzné otázky:

- “Akú veľkú trénovaciu množinu máme zvoliť, aby sme dosiahli dostatočne malú ($\leq \epsilon$) chybu algoritmu?”
- “Aké t máme zvoliť, aby sme s vysokou pravdepodobnosťou ($\geq \delta$) dostali dostatočne dobrú ($\leq \epsilon$) hypotézu?”

Odtiaľ sa odvíja názov *PAC učenie* (z anglického *probably approximately correct learning*).

Obe typy “chýb” sú potrebné, keď sa chceme rozprávať o tom, aký vplyv má veľkosť trénovacej množiny na trénovací algoritmus. Po prvé, ε je potrebné ako miera toho, čo je dostatočne dobrá hypotéza. Po druhé, δ je potrebné, nakoľko vo všeobecnosti nevieme garantovať, že dostaneme dobrú hypotézu: mohli sme si (s malou pravdepodobnosťou) vytiahnuť zlé trénovacie dáta.

Začneme definíciou PAC učenia, ktorá ešte nebude brať do úvahy výpočtovú stránku učenia. Zameriame sa na klasifikačné úlohy, v ktorých je cieľom rozlíšiť medzi reprezentantmi nejakého cieľového konceptu od nereprezentantov. Napríklad cieľovým konceptom môže byť “písmeno A”. Hypotéza dostane na vstupe obrázok 32×32 a má povedať, či tento obrázok vyobrazuje písmeno A alebo nie.

Máme teda množinu konceptov C , z ktorej pochádza cieľový koncept c . Náš trénovací algoritmus ale nevie, ktorý z nich to je. Jeho úlohou je nájsť dobrú aproximáciu, ktorú bude hľadať v množine hypotéz H . Budeme predpokladať, že $H \supseteq C$, aby bolo zaručená existencia dobrej hypotézy.

Algoritmus na vstupe dostane niekoľko trénovacích príkladov v tvare dvojíc $(x, c(x))$, pričom x pochádza z pravdepodobnostného rozdelenia P . Uvažovať y v rozdelení P nie je potrebné, nakoľko je jednoznačne určené cez x a c .

Definícia 5. Nech C je množina konceptov nad množinou vstupov X . Hovoríme, že C je *PAC naučiteľná* ak existuje algoritmus L s nasledujúcou vlastnosťou: pre každý (cieľový) koncept $c \in C$, pre každé $\varepsilon > 0$, $\delta > 0$ a pre každé možné pravdepodobnostné rozdelenie P , ak algoritmu L dáme na vstupe t trénovacích príkladov $(x_i, c(x_i))$ náhodne vybraných podľa P a t je dostatočne veľké, tak nám algoritmus s pravdepodobnosťou nanajvýš δ vráti hypotézu $\hat{h} \in C$ spĺňajúcu

$\text{err}(\hat{h}) \leq \varepsilon$. Táto pravdepodobnosť sa berie cez náhodnosť v ťahaní trénovacích príkladov a prípadnú náhodnosť algoritmu L .

Poznámka 2.1. Neskôr túto definíciu rozšírime tak, aby brala do úvahy aj výpočtovú stránku algoritmu. Vtedy sa budeme zaoberať aj tým, v akom čase náš algoritmus beží a koľko trénovacích príkladov potrebuje. Budeme hovoriť o *efektívnej PAC naučiteľnosti*.

Budeme hovoriť, že trénovací príklad (x, y) je *pozitívny*, ak $y = 1$. V opačnom prípade budeme hovoriť, že príklad je *negatívny*.

2.1 Konečné množiny hypotéz

V tejto časti sa zameriame na konečné množiny hypotéz.

Budeme predpokladať, že algoritmus vždy vráti *konzistentný klasifikátor*: takú hypotézu, ktorá je konzistentná s trénovacími príkladmi, teda že pre ľubovoľnú trénovaciu množinu T platí $\text{err}_T(\hat{h}) = 0$. Za predpokladu $H \supseteq C$ je to splniteľný predpoklad: jedným konzistentným klasifikátorom je samotný cieľový koncept c .

2.1.1 Základné výsledky

Veta 2.1. *Nech je dané $\varepsilon > 0$. Hypotézu nazveme zlú, ak jej chybovosť je väčšia ako ε . Potom vieme pomocou počtu trénovacích príkladov t odhadnúť pravdepodobnosť, že nám algoritmus vráti zlú hypotézu, nasledovne:*

$$P(\text{err}(\hat{h}) > \varepsilon) < |H| \cdot e^{-\varepsilon t}$$

Dôkaz. Ak žiadna zlá hypotéza nie je konzistentná s príkladmi, tak výstupom algoritmu nemôže byť zlá hypotéza. Budeme sa teda snažiť zhora odhadnúť pravdepodobnosť, že aspoň jedna zlá hypotéza “prežila”.

Nech h je ľubovoľná zlá hypotéza. Pravdepodobnosť, že je konzistentná s trénovacími príkladmi, je rovná $(1 - \text{err}(h))^t$, čo vieme odhadnúť nasledovne:

$$(1 - \text{err}(h))^t < (1 - \varepsilon)^t \leq e^{-\varepsilon t}$$

Zlých hypotéz je nanajvýš toľko, koľko je všetkých hypotéz, teda $|H|$. Pravdepodobnosť, že aspoň jedna z nich bude konzistentná s príkladmi, sa dá odhadnúť zhora ako súčet ich pravdepodobností:

$$P(\text{aspoň jedna zlá}) < |H| \cdot e^{-\varepsilon t}$$

Odkiaľ dostávame požadovanú nerovnosť. □

Poznámka 2.2. Vo vyššie uvedenom dôkaze sme nepredpokladali nič o pravdepodobnostnom rozdelení P ani o cieľovom koncepte c . Uvedená veta teda platí pre ľubovoľné P a c .

Čo ak nás zaujíma druhá otázka: “Ako závisí chyba algoritmu od počtu trénovacích príkladov?” Pri klasifikačných úlohách je táto otázka úzko spätá s predošlou otázkou, kde sme sa zaujímali o ε a δ .

Veta 2.2. *Platí*

$$\text{chyba algoritmu} \leq \frac{1}{t} \cdot (\ln |H| + \ln t + 1).$$

Dôkaz. Ak nám algoritmus vráti dobrú hypotézu (s chybou nanajvýš ε), vieme jej chybu odhadnúť zhora ako ε . Ak nám algoritmus vráti zlú hypotézu, jej chyba je nanajvýš 1. Z toho dostávame nasledovný horný odhad na celkovú chybu algoritmu:

$$\text{chyba algoritmu} = E[\text{err}(\hat{h})] \leq P(\text{err}(\hat{h}) \leq \varepsilon) \cdot \varepsilon + P(\text{err}(\hat{h}) > \varepsilon) \cdot 1$$

Pritom pravdepodobnosti na pravej strane vieme odhadnúť zhora: $P(\text{err}(\hat{h}) \leq \varepsilon) \leq 1$, a druhú vieme odhadnúť pomocou vety 2.1. Dostávame tak odhad

$$\text{chyba algoritmu} \leq 1 \cdot \varepsilon + |H| \cdot e^{-\varepsilon t}.$$

My sme si ale mohli zvoliť ε ľubovoľne. Ak teda chceme dostať čo najlepší odhad, nájdeme ε , pre ktoré je výraz na pravej strane čo najmenší. Zderivujeme a položíme rovné nule:

$$1 - |H| \cdot t \cdot e^{-\varepsilon t} = 0 \quad (2.1)$$

$$\varepsilon = \frac{1}{t} \cdot (\ln |H| + \ln t) \quad (2.2)$$

Odtiaľ dosadením dostaneme požadovaný odhad na chybu algoritmu. \square

Na základe týchto dvoch viet vieme sformulovať postačujúce podmienky na t také, aby boli príslušné chyby (ε, δ a chyba algoritmu) dostatočne malé. Sformulujeme a dokážeme jednu z nich.

Dôsledok 2.3. *Množina hypotéz H je PAC-naučiteľná: pre každé $\varepsilon > 0$, $\delta > 0$, ľubovoľné pravdepodobnostné rozdelenie P a ľubovoľný cieľový koncept $f \in H$ existuje počet tréningových príkladov t taký, že platí*

$$P(\text{err}(\hat{h}) \leq \varepsilon) \geq 1 - \delta.$$

Ekvivalentne,

$$P(\text{err}(\hat{h}) > \varepsilon) \leq \delta.$$

Dôkaz. Podľa vety 2.1 platí

$$P(\text{err}(\hat{h}) > \varepsilon) < |H| \cdot e^{-\varepsilon t}.$$

Stačí nám teda zvoliť také t , aby bol výraz na pravej strane menší rovný δ . Odtiaľ dostaneme postačujúci počet tréningových príkladov t :

$$|H| \cdot e^{-\varepsilon t} \leq \delta \quad (2.3)$$

$$\ln |H| - \varepsilon t \leq \ln \delta \quad (2.4)$$

$$\varepsilon t \geq \ln |H| - \ln \delta \quad (2.5)$$

$$t \geq \frac{1}{\varepsilon} \cdot \left(\ln |H| + \ln \frac{1}{\delta} \right) \quad (2.6)$$

\square

2.1.2 Problém konjunkcie

Jedným príkladom problému, kde je množina hypotéz konečná, je *problém konjunkcie pozitívnych literálov*. Na ňom si ukážeme, že (aspoň v niektorých problémoch) sú vyššie uvedené odhady relatívne tesné.

Nech je dané n . Množina vstupov sú všetky možné ohodnotenia boolovských premenných x_1, \dots, x_n . Napríklad $x_1 = 0, x_2 = 1, x_3 = 0$ je priradenie hodnôt. Priradenia vieme zapísať vo vektorovom tvare: vyššie uvedený príklad by sme zapísali ako $x = (0, 1, 0)$. Všetkých vstupov je zrejme 2^n .

Množina konceptov C sú všetky konjunkcie nad pozitívnymi literálmi x_1, \dots, x_n . Tieto konjunkcie sú chápané ako funkcie, ktoré vracajú 1 iba ak dané ohodnotenie premenných spĺňa túto konjunkciu. Napríklad $x_1 \wedge x_3 \wedge x_4$ vráti 1 na všetkých tých vstupoch, kde táto konjunkcia platí: musí platiť $x_1 = 1, x_3 = 1, x_4 = 1$, ale ostatné premenné už môžu mať ľubovoľnú hodnotu. Všetkých konceptov je zrejme tiež 2^n . Množina hypotéz H bude rovnaká, ako množina konceptov.

Aby sme videli, že konzistentnosť s tréningovými príkladmi je realistická požiadavka, ukážeme si, ako sa dá na základe tréningových príkladov zostrojiť nejaký konzistentný klasifikátor.

Veta 2.4. *Nech hypotéza \hat{h} je konjunkciou všetkých tých premenných, ktoré sa vyskytujú vo všetkých pozitívnych tréningových príkladoch. Potom \hat{h} je konzistentná so všetkými tréningovými príkladmi.*

Dôkaz. Ak (x, y) je pozitívny príklad, v cieľovej hypotéze c môžu byť jedine tie premenné, ktoré majú v x priradenú hodnotu 1. Keď túto úvahu zopakujeme pre všetky pozitívne príklady, dostaneme niekoľko množín “povolených premenných”. Cieľová premenná musí byť podmnožinou všetkých z nich, teda je podmnožinou ich prieniku. Pritom ich prienik je práve hypotéza \hat{h} .

Z toho vyplýva, že ak nie je splnené c , nemôže byť splnené ani h . Naša hypotéza totiž kladie ešte väčšie požiadavky na hodnoty premenných. Teda,

$$(c(x) = 0) \implies (h(x) = 0),$$

takže h je konzistentná s negatívnymi príkladmi.

Čo sa týka pozitívnych príkladov, v h sú všetky tie premenné, ktoré sme do nej mohli dať tak, aby bola konzistentná so všetkými pozitívnymi príkladmi. Teda jej konzistentnosť s pozitívnymi príkladmi vyplýva priamo z jej konštrukcie. \square

Uvedieme teraz niektoré výsledky z predchádzajúcej časti tak, ako platia pre problém konjunkcie.

Dôsledok 2.5. *Platí*

$$\text{chyba algoritmu} \leq \frac{1}{t} \cdot (n \ln 2 + \ln t + 1) = O\left(\frac{n + \ln t}{t}\right).$$

Dôsledok 2.6. *Aby sme mali zaručené (s pravdepodobnosťou aspoň $1 - \delta$), že dostaneme hypotézu s chybou najviac ε , stačí zvoliť veľkosť tréningovej množiny nasledovne:*

$$t \geq \frac{1}{\varepsilon} \cdot \left(n \ln 2 + \ln \frac{1}{\delta}\right) = \Omega\left(\frac{n + \ln \frac{1}{\delta}}{\varepsilon}\right)$$

Tieto odhady sú relatívne tesné. Vo všeobecnom prípade je ťažké dostať nejaký dolný odhad, nakoľko pravdepodobnostné rozdelenie P a cieľový koncept c môžu byť degenerované a “uľahčiť algoritmu robotu”. Uvidíme ale, že pre niektoré “ťažké” prípady vieme spraviť dolný odhad.

Veta 2.7. *Existuje pravdepodobnostné rozdelenie P a cieľový koncept $c \in C$ také, že nech je tréningový algoritmus ľubovoľný, pre jeho chybu platí nasledovný dolný odhad:*

$$\text{chyba algoritmu} \geq \frac{1}{2e} \cdot \frac{n - 1}{t + 1} = \Omega\left(\frac{n}{t}\right)$$

V znení vety je trochu obmedzujúce, že cieľový koncept musí byť pevne vybraný. Ukážeme teda najprv, že pokiaľ tvrdenie dokážeme pre náhodne vybrané c , bude z neho plynúť pôvodné tvrdenie.

Lemma 2.8. *Nech P_C je pravdepodobnostné rozdelenie nad množinou konceptov C . Cieľovú hypotézu vyberieme náhodne podľa P_C . Ak pre nejaké číslo k platí, že stredná hodnota chyby algoritmu je aspoň k , tj.*

$$\mathbb{E}_{c \sim P_C} [\text{chyba algoritmu}] \geq k,$$

tak existuje voľba cieľového konceptu $c \in C$ taká, že chyba algoritmu bude tiež aspoň k .

Dôkaz. Ak by na každom cieľovom koncepte bola chyba algoritmu menšia ako k , potom by aj ľubovoľný vážený priemer chýb (zodpovedajúci strednej hodnote s rozdelením P_C) bol menší ako k . To by bol spor. \square

Ďalej pokračujeme dôkazom vety 2.7. V ňom si už môžeme dovoliť vyberať cieľový koncept náhodne.

Dôkaz. Zvolíme pravdepodobnostné rozdelenie, ktoré priradí nenulovú pravdepodobnosť iba nasledovnej sade vstupov. Konkrétne pravdepodobnosti určíme neskôr v dôkaze tak, ako sa nám to bude hodiť.

$$x^{(1)} = (0, 1, 1, \dots, 1, 1) \quad (2.7)$$

$$x^{(2)} = (1, 0, 1, \dots, 1, 1) \quad (2.8)$$

$$\vdots \quad (2.9)$$

$$x^{(n)} = (1, 1, 1, \dots, 1, 0) \quad (2.10)$$

Tieto vstupy majú nasledujúcu vlastnosť. Pre ľubovoľný koncept c platí, že sa v ňom nachádza konjunkcia x_i vtedy a len vtedy, keď $f(x^{(i)}) = 0$. Každý z týchto vstupov sa dá teda chápať ako “test na niektorú premennú”.

Trénovacie príklady potom môžeme chápať tak, že algoritmu dávajú informáciu o jednotlivých premenných: “Je alebo nie je v cieľovej hypotéze?” Pokiaľ ale niektorý zo vstupov $x^{(i)}$ nie je medzi trénovacími príkladmi, nemáme žiadnu informáciu o tom, či sa tam príslušná premenná nachádza. Ak si potom \hat{h} vytiahne pri testovaní takýto vstup, môže jedine hádať, aká je správna odpoveď. Šanca úspechu pri hádaní bude $\frac{1}{2}$, pokiaľ vyberieme c rovnomerne náhodne z celej množiny konceptov.

Označme si pravdepodobnosti pridelené jednotlivým vstupom p_1, \dots, p_n . Aká je šanca, že pri trénovaní vstup $x^{(i)}$ nedostaneme a potom si ho pri testovaní vytiahneme?

$$p_i \cdot (1 - p_i)^t$$

Celková pravdepodobnosť, že si vytiahneme pri testovaní vstup mimo trénovacej množiny, je potom súčet jednotlivých pravdepodobností (nakoľko sú jednotlivé udalosti dizjunktné):

$$\sum_{i=1}^n p_i \cdot (1 - p_i)^t$$

V týchto prípadoch bude mať \hat{h} chybu $\frac{1}{2}$. V ostatných prípadoch sme si vytiahli počas testovania nejaký príklad, ktorý bol aj v trénovacej množine. Pokiaľ sme si ho zapamätali, tak budeme mať chybu 0, v každom prípade bude ale chyba aspoň 0. Dostávame tak nasledovný dolný odhad na chybu algoritmu:

$$\mathbb{E}_{c \in C} [\text{chyba algoritmu}] \geq \frac{1}{2} \cdot \left(\sum_{i=1}^n p_i \cdot (1 - p_i)^t \right)$$

Ako ale zvoliť p_1, \dots, p_n tak, aby sme dostali dobrý dolný odhad? Skúsime zvoliť p_i také, pre ktoré nadobúda výraz $p_i \cdot (1 - p_i)^t$ maximum. Zderivovaním a položením rovné 0 dostaneme

$$p_i = \frac{1}{t + 1}.$$

Pokiaľ ale $t \neq n - 1$, nemôžeme zvoliť všetky pravdepodobnosti takéto: pre $t < n - 1$ je súčet pravdepodobností priveľký a pre $t > n - 1$ prímalý. Prvý prípad nás nezaujíma, nakoľko je to “len konštanta” (v zmysle $t \rightarrow \infty$). V druhom prípade si vieme zvoliť jedného “obetného baránka”

p_1 , ktorému priradíme celú zvyšnú pravdepodobnosť.

$$p_1 = 1 - \frac{n-1}{t+1} \quad (2.11)$$

$$p_2 = \frac{1}{t+1} \quad (2.12)$$

$$\vdots \quad (2.13)$$

$$p_n = \frac{1}{t+1} \quad (2.14)$$

Dosadíme a dostaneme tak odhad:

$$\mathbb{E}_{c \in C} [\text{chyba algoritmu}] \geq \frac{1}{2} \cdot \left((n-1) \cdot \frac{1}{t+1} \cdot \left(1 - \frac{1}{t+1}\right)^t + \frac{n-1}{t+1} \cdot \left(1 - \frac{n-1}{t+1}\right)^t \right)$$

Odignorujeme druhý sčítanec, čím sa nám výraz na pravej strane môže len zmenšiť, takže nerovnosť sa zachová. Ďalej použijeme odhad

$$\left(1 - \frac{1}{t+1}\right)^t \geq \frac{1}{e},$$

ktorý platí pre všetky prirodzené t . (Tento odhad sa dá dokázať napríklad tak, že sa dokáže, že daný výraz je rastúci od t . Odhad z toho už plynie ľahko, nakoľko jeho limita pre $t \rightarrow \infty$ je práve $\frac{1}{e}$.) Dostaneme tak požadovanú nerovnosť:

$$\mathbb{E}_{c \in C} [\text{chyba algoritmu}] \geq \frac{1}{2e} \cdot \frac{n-1}{t+1}$$

□

Dolný odhad zodpovedajúci dôsledku 2.6 uvidíme bez dôkazu.

Veta 2.9. *Pre ľubovoľný tréningový algoritmus existuje pravdepodobnostné rozdelenie P a cieľová hypotéza $f \in H$, ktoré vynútiť, že algoritmus bude potrebovať aspoň*

$$t = \Omega \left(\frac{1}{\varepsilon} \cdot \left(n + \ln \frac{1}{\delta} \right) \right),$$

aby platilo

$$\mathbb{P}(\text{err}(\hat{h}) > \varepsilon) \leq \delta.$$

TODO dôkaz

2.2 Nekonečné množiny hypotéz

Vyššie uvedený postupy pre konečné množiny hypotéz zlyhávajú, ak $|H| = \infty$: nedostaneme z nich žiaden odhad. Aj pre nekonečné množiny hypotéz je ale možné odvodiť horné odhady, ktoré sa už nebudú odvíjať od veľkosti množiny H , ale od nejakej jej miery zložitosti. Touto mierou zložitosti bude tzv. *Vapnik-Červonenkisova dimenzia*. Predtým sa ale pozrieme na jednoduchý príklad.

2.2.1 Obdĺžniková hra

Ide o klasifikačnú úlohu. Množina vstupov sú všetky body v rovine. Množina konceptov sú všetky obdĺžniky, ktorých strany sú rovnobežné so súradnicovými osami. Pre každý bod v rovine sa teda pýtame, či je vo vnútri cieľového obdĺžnika c alebo nie. Body na okraji obdĺžnika považujeme, že sú vnútri.

Pre jednoduchosť argumentu použijeme konkrétny trénovací algoritmus. Dôkaz by sa dal zovšeobecniť na prípad, kedy jediné, čo o algoritme predpokladáme je, že nám vracia konzistentný klasifikátor. S nasledovným konzistentným klasifikátorom to ale bude jednoduchšie.

Z trénovacích príkladov vezmeme tie pozitívne. Ako hypotézu \hat{h} zvolíme najmenší (vzhľadom na inklúziu) osovorovnobezný obdĺžnik, ktorý obsahuje všetky body vo vybraných príkladoch. Na obrázku ?? ilustrujeme konštrukciu \hat{h} a porovnávame ho s c .

TODO obrázok

Lemma 2.10. *Každá množina bodov $(x_1, y_1), \dots, (x_t, y_t)$ má unikátny bounding box: najmenší (vzhľadom na inklúziu) osovorovnobezný obdĺžnik, ktorý ich všetky obsahuje.*

Dôkaz. Každý obdĺžnik je jednoznačne určený x -ovými súradnicami ľavého a pravého okraja a y -ovými súradnicami dolného a horného okraja. Označme ich $x_{\leftarrow}, x_{\rightarrow}$ a $y_{\downarrow}, y_{\uparrow}$. Vo vnútri sú práve tie body (x, y) , ktoré spĺňajú

$$x_{\leftarrow} \leq x \leq x_{\rightarrow}, \quad y_{\downarrow} \leq y \leq y_{\uparrow}.$$

Ak teda majú byť všetky spomínané body vo vnútri, musí platiť

$$x_{\leftarrow} \leq \min x_i, \quad x_{\rightarrow} \geq \max x_i, \quad y_{\downarrow} \leq \min y_i, \quad y_{\uparrow} \geq \max y_i$$

Najmenší osovorovnobezný obdĺžnik, ktorý toto spĺňa, je ten, pre ktorý nastávajú v jednotlivých nerovnostiach rovnosti. Teda $x_{\leftarrow} = \min x_i$, $x_{\rightarrow} = \max x_i$, $y_{\downarrow} = \min y_i$ a $y_{\uparrow} = \max y_i$. \square

Dôsledok 2.11. *Obdĺžnik \hat{h} je podmnožinou cieľového obdĺžnika c .*

Dôsledok 2.12. *Obdĺžnik \hat{h} je konzistentný klasifikátor.*

Dôkaz. Konzistentnosť na pozitívnych príkladoch vyplýva z konštrukcie \hat{h} . Každý negatívny príklad je mimo c a teda aj mimo \hat{h} , je teda konzistentný aj s negatívnymi príkladmi. \square

Veta 2.13. *Obdĺžniková hra je PAC naučiteľná: pre každý cieľový obdĺžnik $c \in C$, $\varepsilon > 0$, $\delta > 0$ a pravdepodobnostné rozdelenie P , ak vyššie popísanému algoritmu dáme t trénovacích príkladov z P a platí*

$$t \geq \frac{4}{\varepsilon} \cdot \ln \frac{4}{\delta},$$

tak nám algoritmus s vysokou pravdepodobnosťou vráti hypotézu s nízkou chybou:

$$P(\text{err}(\hat{h}) > \varepsilon) \leq \delta.$$

Dôkaz. Zle klasifikované body sú práve tie, ktoré sú vo vnútri c ale nie sú vo vnútri \hat{h} . Každý z týchto bodov padne do aspoň jedného z okrajových “pásikov” (nerátajúc jeden z okrajov), ako je zobrazené na obrázku ??.

Pod váhou množiny budeme rozumieť pravdepodobnosť, že náhodný bod z rozdelenia P padne do tejto množiny. Označme $p_{\leftarrow}, p_{\rightarrow}, p_{\uparrow}, p_{\downarrow}$ postupne váhy ľavého, pravého, horného, resp. dolného pásika. Chyba \hat{h} sa dá zhora odhadnúť ako súčet týchto váh:

$$\text{err}(\hat{h}) \leq p_{\leftarrow} + p_{\rightarrow} + p_{\uparrow} + p_{\downarrow}.$$

Ak by sme voľbou dostatočne veľkého t vedeli zaručiť (s pravdepodobnosťou aspoň $1 - \delta$), že dostaneme takú hypotézu \hat{h} , pre ktorú je výraz na pravej strane nanajvýš ε , vyhrali by sme. Ukážeme, že vieme zaručiť

$$p_{\leftarrow}, p_{\rightarrow}, p_{\uparrow}, p_{\downarrow} \leq \frac{\varepsilon}{4}.$$

Postup bude vo všetkých štyroch prípadoch ten istý, ukážeme si to teda len na ľavom pásiku.

TODO dôkaz (eww, grc)

□

2.2.2 Vapnik-Červonenkisova dimenzia

Vapnik-Červonenkisova dimenzia (alebo skrátene VC dimenzia) je miera zložitosti množiny hypotéz, pomocou ktorej sme schopní tvoriť tvrdenia v štýle PAC učenia. Na jej definíciu budeme potrebovať niekoľko pomocných pojmov.

Definícia 6. Majme konečnú podmnožinu vstupov $S \subseteq X$. Hovoríme, že množina hypotéz H *rozbíja* množinu S , ak platí: nech označíme vstupy v S ako pozitívne alebo negatívne akokoľvek, v množine H existuje hypotéza konzistentná s týmto označením.

Intuitívne, množina S je rozbitá hypotézami H , ak nie je pre hypotézy v H príliš náročné modelovať príklady v S : sú schopné ich modelovať akokoľvek.

Definícia 7. Vapnik-Červonenkisova dimenzia množiny hypotéz H , označovaná $VCD(H)$, je veľkosť najväčšej množiny $S \subseteq X$, ktorá sa dá rozbiť hypotézami v H . Ak sa dá rozbiť ľubovoľne veľká množina S , definujeme $VCD(H) = \infty$.

Teda na to, aby sme ukázali dolný odhad $VCD(H) \geq d$, stačí nám nájsť jednu množinu S veľkosti d , ktorá sa dá rozbiť. Aby sme ale ukázali horný odhad $VCD(H) \leq d$, musíme ukázať, že žiadna množina veľkosti $d + 1$ sa nedá rozbiť: teda že existuje také označenie vstupov, pre ktoré neexistuje v H konzistentný klasifikátor. Z tohto dôvodu je obvykle ťažšie dokázať horný odhad na VC dimenziu, ako dolný odhad.

Než uvedieme vety hovoriace o PAC naučiteľnosti, uvedieme príklady niektorých nekonečných množín hypotéz, a neformálne zdôvodníme ich VC dimenzie. Väčšinou budú geometrického charakteru (čo je prirodzené, ak máme mať nekonečnú množinu hypotéz).

Intervaly na reálnych číslach. Nie je ťažké vidieť, že ľubovoľná dvojica bodov sa dá rozbiť. Na druhej strane, ak máme tri body, tak ich vieme označiť ako na obrázku ??, čo nie je konzistentné so žiadnym intervalom.

Polroviny v rovine. Každá trojica bodov tvoriacich nedegenerovaný trojuholník sa dá rozbiť, ako je ilustrované na obrázku ??. Na druhej strane, žiadna štvorica bodov sa rozbiť nedá:

- Ak tvoria konvexný štvoruholník, tak jednu protiľahlú dvojicu označíme kladne a druhú záporne.
- Ak tvoria nekonvexný štvoruholník, jeden z bodov je vo vnútri trojuholníka tvoreného ostatnými tromi: tento bod označíme kladne a ostatné body záporne.

Osovorovnoběžné obdĺžniky. Štvorica bodov, ktorá sa dá rozbiť, je ilustrovaná na obrázku ??. Na druhej strane, žiadna päťica bodov sa nedá rozbiť. Rozoberieme dva prípady:

- Ak je aspoň jeden z bodov vo vnútri bounding boxu (nie na okraji), tak ho označíme záporne a všetky body na okraji kladne.
- Ak sú všetky body na obode bounding boxu, na jeden strane toho obdĺžnika musia ležať aspoň dva body. Jeden z nich označíme kladne a druhý záporne.