

Obsah

1	Úvod do teórie strojového učenia	2
1.1	Model	2
1.2	Definície a označenia	3
1.2.1	Meranie chyby hypotézy	3
1.2.2	Meranie chyby algoritmu	4
1.2.3	Príklady chybových funkcií	4
1.2.4	Ďalšie označenia	5
1.2.5	Diskusia o výpočtových detailoch	5
1.3	Analýza veľkostí chýb	5
1.3.1	Teoretické limity	5
1.3.2	Kompromis medzi výchylkou a rozptylom (bias-complexity tradeoff) . . .	6
1.3.3	Kompromis medzi výchylkou a rozptylom (bias-variance tradeoff)	8
1.4	Podučenie/preučenie	9
1.4.1	Regularizácia	10
1.4.2	Holdout testing	11
1.5	Cvičenia	12
1.6	Appendix: iný kompromis medzi inou výchylkou a iným rozptylom	12
1.7	Appendix: technické detaily k bias-variance tradeoff	14
2	PAC učenie	17
2.1	Definície, označenia a predpoklady	17
2.2	Konečné množiny hypotéz	18
2.2.1	Konečné množiny sú PAC★ naučiteľné	18
2.2.2	Problém konjunkcie	20
2.3	Nekonečné množiny hypotéz	23
2.3.1	Obdĺžniková hra	23
2.3.2	Vapnik-Chervonenkisova dimenzia	25
	Appendix	28
2.A	Dôkaz o VC dimenzii a PAC naučiteľnosti	28

Kapitola 1

Úvod do teórie strojového učenia

V oblasti teórie strojového učenia sa zaoberáme učením; formalizujeme ho, navrhujeme rôzne algoritmy schopné tohto učenia, analyzujeme ako rýchlo sa tieto algoritmy učia, ... Čo to ale vlastne je to “učenie”, a ako sa nad týmto konceptom dá zamýšľať?

Začnime tým, že si vybavíme, čo všetko sa nám spája s učením. Napríklad učenie sa na skúšky, ktoré pre niektorých vyzerá tak, že sa snažia si zapamätať všetkých 80 strán skript naspamäť. Alebo keď sa snažíme naučiť novú skladbu na klavíri: hráme ju znova a znova, až kým v tom nie sme dobrí. Alebo tréningovanie na dôležitý zápas vo futbale, ...

Vidíme teda, že učenie je zložitý koncept. Konkrétne detaily toho, ako presne sa učíme (resp. prečo to funguje), sú známe málo ľuďom. Avšak vidíme, že jednotlivé “učenia” mali niečo spoločné: nejaká aktivita sa opakovala veľa krát, a čím viackrát sa zopakovala, tým lepší bol výkon. Zároveň ale nechceme, aby sa zakaždým odohralo to isté; nová skúsenosť je dôležitou súčasťou učenia.

Ak by sme len tak vychrlili nejakú definíciu, ktorá nám príde rozumná, riskujeme, že nebude “sedieť s našou intuíciou”, alebo nebude dostatočne všeobecná. Kvôli tejto zložitosti teda namiesto toho, aby sme zachytili “učenie” v jednej definícii, študujeme veľa rôznych modelov učenia, ktoré sú aplikovateľné v rôznych kontextoch. Pod “modelom” rozumieme akúsi hračku, zjednodušený mentálny obraz, s ktorým sa jednoducho pracuje a zároveň “sedí s našou intuíciou”. Takto síce nijak neručíme, že pokrývame “všetko učenie”; čím viac viet a teorém v rôznych modeloch ale dokážeme, tým lepší budeme mať obraz o tom, čo to učenie je.

My sa v týchto skriptách budeme zaoberať hlavne oblasťou *učenia s učiteľom*. Uvedieme hlavný model (tzv. *statistical learning framework*), s ktorým budeme pracovať, a v prípade potreby rozširovať/dolaďovať.

1.1 Model

Predstavme si, že chceme vedieť na základe nejakých vstupných dát (ktoré budeme spravidla označovať x) predpovedať výstupné dáta (označované y). Napríklad chceme na základe rozlohy bytu vedieť predpovedať jeho cenu. Alebo vedieť z obrázku (vo formáte 32×32 čiernobielych pixelov) povedať, či sa v ňom nachádza mačka alebo nie.

Snažíme sa teda zachytiť nejaký vzťah medzi dátami x a y . To, akým spôsobom to budeme robiť, je nasledovné: nejakým procesom P získame t *trénovacích príkladov* $(x_1, y_1), \dots, (x_t, y_t)$. Na základe týchto príkladov budeme chcieť navrhnúť nejakú funkciu h , ktorá bude vedieť podľa vstupu x predpovedať výstup y s dostatočnou presnosťou.

Prírodný spôsob, akým môžeme merať úspešnosť funkcie h , je podľa toho, ako sa jej darí na našich t príkladoch. Prístupu, kde hľadáme funkciu, ktorej sa čo najlepšie darí na trénovacích príkladoch, hovoríme *empirical risk minimization* (skrátene ERM). Dobrá funkcia h ale bude schopná aj *generalizovať*: bude sa jej dariť na ďalších dátach, ktoré vieme získať tým istým procesom P .

Schopnosť generalizácie je jednou z najdôležitejších vlastností, ktoré od trénovania pomocou techník strojového učenia požadujeme. Predstavme si, že by sme funkciu h zostrojili tak, že si “zapamätáme” všetky trénovacie príklady: ak je na vstupe x , ktoré bolo medzi trénovacími príkladmi, vrátime príslušné y . Takáto funkcia h má zrejme nulovú chybu na trénovacích príkladoch (pokiaľ pre jedno x existuje vždy len jedno správne y), ale mimo nich sa jej vôbec nebude dariť.

Vidíme teda, že ERM vo svojej holej podstate nefunguje. Problém je v tom, že neobmedzujeme to, aké funkcie považujeme sa “kandidátov”: nič nezaručuje, že informácia z trénovacích príkladov sa preniesie aj na ostatné príklady. Obmedzíme teda množinu funkcií, ktoré uvažujeme, na nejakú užšiu skupinu, ktorú nazveme *množina hypotéz*.

Napríklad si predstavme, že hľadáme funkciu $f : \mathbb{R} \rightarrow \mathbb{R}$ a dostali sme trénovacie príklady $(0, 0)$ a $(1, 1)$. Ak o probléme nič ďalšie nevieme, tak nevieme povedať nič o tom, ako sa bude f správať pre $x \notin \{0, 1\}$. Ak ale vieme, že hľadáme neklesajúcu funkciu s oborom hodnôt $\{0, 1\}$, tak vieme, že $f(x) = 0$ pre $x \leq 0$, $f(x) = 1$ pre $x \geq 1$, a niekde v intervale $(0, 1)$ sa to “láme”.

Tým, že sme obmedzili množinu hypotéz, sme v podstate zaviedli *inductive bias*: predpoklady, ktoré využívame na to, aby sme odvodili výsledky aj na tých vstupoch, ktoré sme ešte nevideli.

1.2 Definície a označenia

Množinu všetkých vstupov označíme X , niekedy jej budeme hovoriť aj *množina inštancií*. Množinu všetkých výstupov označíme Y .

Množinu všetkých trénovacích príkladov budeme volať *trénovacia množina*, označíme ju T . (Z matematického hľadiska to ale nie je množina, skôr multimnožina, keďže jedna dvojica (x, y) sa medzi trénovacími príkladmi môže vyskytovať viackrát. Pre jednoduchosť to ale budeme nazývať množinou.) Počet trénovacích príkladov označíme t .

Proces, ktorým získavame dáta (x, y) , vieme formalizovať ako pravdepodobnostné rozdelenie P nad priestorom $X \times Y$, ktoré každému možnému pozorovaniu priradí nejakú pravdepodobnosť (resp. hustotu pravdepodobnosti v prípade spojitého rozdelenia). Trénovacie príklady získame ako t nezávislých vzoriek z rozdelenia P .

Postup, ktorým zostrojíme funkciu h , vieme formalizovať ako algoritmus L , ktorý na vstupe dostane niekoľko trénovacích príkladov $(x_1, y_1), \dots, (x_t, y_t)$, a na výstupe vráti funkciu $h : X \rightarrow Y$, ktorú budeme volať *hypotéza*. Množinu hypotéz označíme H .

1.2.1 Meranie chyby hypotézy

Ako vyjadriť mieru toho, že sa hypotéze “dobre darí”? Spravíme tak pomocou *chybovej funkcie* $\text{err} : Y \times Y \rightarrow \mathbb{R}_0^+$, ktorej význam je nasledovný: $\text{err}(y, y')$ vyjadruje, ako veľmi sa od seba líšia výstupy y a y' .

Chyba na jednom príklade. Keď dostaneme nejaký príklad (x, y) , vieme zmerať, ako veľmi sa naša hypotéza h na tomto vstupe pomýlila, ako $\text{err}(h(x), y)$.

Očakávaná testovacia chyba hypotézy (OTeChHyp). Celkovú chybu našej hypotézy potom vieme zmerať ako očakávanú chybu na náhodne vybranom príklade z rozdelenia P . Túto chybu budeme označovať Err , a vypočítame ju nasledovne:

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

Častokrát ale bude rozdelenie P jasné z kontextu, v takom prípade budeme skrátene písať $\text{Err}(h)$.

Trénovacia chyba. Niekedy nás bude zaujímať ale aj chyba, ktorú hypotéza nadobúda na trénovacej množine. Budeme ju tiež označovať Err , a vypočítame ju nasledovne:

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

Tento zápis sa dá chápať aj tak, že tréningovú množinu interpretujeme ako pravdepodobnostné rozdelenie, pričom pravdepodobnosť každého príkladu je úmerná tomu, koľkokrát sa v množine vyskytuje. Niekedy túto chybu budeme zapisovať aj $\text{Err}_T(h)$.

1.2.2 Meranie chyby algoritmu

Hypotéza je iba výstupom tréningového algoritmu, a závisí od toho, akú konkrétnu tréningovú množinu dostaneme a od prípadnej náhodnosti samotného algoritmu. Ak by sme chceli porovnať dve rôzne algoritmy, na úrovni výstupných hypotéz to nevieme spraviť: pre jednu tréningovú množinu môže byť lepší prvý algoritmus, pre inú množinu zas druhý. Môžeme ale hovoriť o tom, ako dobré budú tie algoritmy v priemere, ak vezmeme do úvahy všetky možné náhodné faktory.

Tréningový algoritmus označíme L (z anglického *learner*). Jeho výstup na tréningovej množine T označíme $L(T)$; prípadne \hat{h}_T , alebo len \hat{h} ak bude z kontextu jasné, aká je tréningová množina.

Očakávaná testovacia chyba algoritmu (*OTeChAlg*). Celkovú chybu algoritmu na t tréningových príkladoch zmeriame ako *OTeChHyp*, ktorú získame natréningovaním na náhodne zvolenej tréningovej množine veľkosti t . Túto chybu budeme opäť značiť Err , vypočítame ju ako

$$\text{Err}(L, P)[t] = \mathbb{E}_{T \sim P^t} [\text{Err}(L(T), P)] = \mathbb{E}_{T \sim P^t} [\text{Err}(\hat{h}_T)].$$

Očakávaná tréningová chyba algoritmu (*OTrChAlg*). Niekedy nás bude zaujímať aj to, akú môžeme očakávať tréningovú chybu výstupnej hypotézy. Uvedomme si, že táto hodnota sa bude líšiť od očakávanej testovacej chyby: totiž výstupná hypotéza priamo závisí od tréningovej množiny T , ale iba nepriamo závisí od pravdepodobnostného rozdelenia P . (Ako príklad uvedieme algoritmus “zapamätaj si tréningové dáta”.) Túto chybu budeme značiť Err_T , a vypočítame ju nasledovne:

$$\text{Err}_T(L, P)[t] = \mathbb{E}_{T \sim P^t} [\text{Err}(L(T), T)] = \mathbb{E}_{T \sim P^t} [\text{Err}_T(\hat{h}_T)].$$

Pri oboch týchto chybách, pokiaľ bude z kontextu jasné rozdelenie P alebo počet príkladov t , zo zápisu ich vynecháme. Teda v najkratšej forme budeme *OTeChAlg* označovať $\text{Err}(L)$ a *OTrChAlg* budeme značiť $\text{Err}_T(L)$. Rozlišovať medzi chybami hypotéz/algoritmov a tréningovými/testovacími chybami budeme podľa kontextu (t.j. či je v indexe T a či je v zátvorke hypotéza h alebo algoritmus L).

1.2.3 Príklady chybových funkcií

Uvedieme ešte príklady používaných chybových funkcií. Pri klasifikácii sa naša hypotéza vždy buď trafi, alebo netrafi do správnej odpovede. Neexistuje ale nejaká miera toho, ako veľmi sa netrafila, resp. ako blízko bola ku správnej odpovedi. Dáva teda zmysel každej správnej odpovedi priradiť chybu 0, a každej nesprávnej odpovedi chybu 1:

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

Potom sa nám viacero vzorcov pre chyby zjednoduší: testovacia chyba hypotézy je

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

a podobne sa dá zjednodušiť aj tréningová chyba.

Pri regresii naopak takáto miera existuje, dokonca by sa dalo povedať, že každá takáto miera zodpovedá jednej chybovej funkcii. Bežnými voľbami sú kvadratická chyba $(y - y')^2$ a absolútna chyba $|y - y'|$.

1.2.4 Ďalšie označenia

Pri zápise stredných hodnôt budeme vynechávať to, odkiaľ premenné pochádzajú, pokiaľ to bude z kontextu jasné. Konkrétne zavedieme nasledovné skrátene zápisy:

- ak sa stredná hodnota berie cez príklady z rozdelenia P :

$$\mathbb{E}_{(x,y) \sim P} \equiv \mathbb{E}_{x,y}$$

- ak sa stredná hodnota berie cez všetky tréningové príklady z množiny T :

$$\mathbb{E}_{(x,y) \in T} \equiv \mathbb{E}_{x,y_i}$$

- ak sa stredná hodnota berie cez všetky možné tréningové množiny veľkosti t :

$$\mathbb{E}_{T \in P^t} \equiv \mathbb{E}_T$$

Podobné skrátene zápisy budeme používať aj pri pravdepodobnostiach a pod.

1.2.5 Diskusia o výpočtových detailoch

V tejto kapitole sa nebudeme zaoberať výpočtovou stránkou strojového učenia, takže od detailov ako časová zložitosť, spôsob hľadania optimálneho riešenia a pod. abstrahujeme.

Namiesto konkrétneho tréningového algoritmu teda budeme predpokladať, že máme *všemo-*
húci algoritmus: taký, ktorý nám vždy vráti tú hypotézu z H , ktorá je z ERM hľadiska najlepšia (teda má najmenšiu chybu na tréningovej množine). Teda

$$\hat{h}_T = \arg \min_{h \in H} (\text{Err}_T(h)).$$

Vo všeobecnosti to nie je realistický predpoklad; nie je jednoduché navrhnúť efektívny algoritmus, ktorý vždy nájde hypotézu s najmenšou chybou. Dokonca pre niektoré dostatočne zložené množiny hypotéz H taký efektívny algoritmus ani nemusí existovať. Nájsť *nejaký* algoritmus ale väčšinou nie je problém, len potom beží v čase $O(t^{!t})$ (len pre ilustráciu, nič faktické za týmto výrazom nehľadajte). Ak teda chceme holistický pohľad na učenie, nemôžeme len tak zanedbať časovú zložitosť a iné výpočtové záležitosti.

1.3 Analýza veľkostí chýb

V tejto časti sa podrobnejšie pozrieme na to, ako závisia vyššie uvedené metriky (t.j. OTechAlg a OTrChAlg) 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.

1.3.1 Teoretické limity

Uvedomme si, že naša hypotéza musí byť funkciou, teda pre jedno x musí vždy vracať jednu a tú istú hodnotu $y = h(x)$. Rozdelenie P ale môže pre jedno x obsahovať viaceré hodnoty y , pre ktoré je pravdepodobnosť nenulová; napríklad P môže reprezentovať zašumené dáta.

Teda ani najlepšia možná hypotéza-funkcia (nie nutne z H) nemusí mať nulovú chybu. Označme túto hypotézu h^\square . Ak by sa nám podarilo vyjadriť jej testovaciu chybu, získali by sme akýsi ireducibilný komponent, ktorý má každá hypotéza; môžeme sa snažiť znížiť iba tú časť chyby, ktorá je “navyše”.

Z definície

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

Chybu ľubovoľnej hypotézy h vieme upraviť nasledovne:

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

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

Pozrime sa na vnútornú strednú hodnotu $(\mathbb{E}_{y|x})$. 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 $c = \mathbb{E}_{y|x=x}[y]$. Takže

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

a testovacia chyba najlepšej hypotézy-funkcie teda je

$$\text{Err}(h^\square) = \mathbb{E}_x \left[\mathbb{E}_{y|x=x} [(y - \mathbb{E}[y])^2] \right] = \mathbb{E}_x \left[\text{Var}(y) \right].$$

1.3.2 Kompromis medzi výchylkou a rozptylom (bias-complexity tradeoff)

V tomto odseku si ukážeme, ako sa dá rozložiť celková OTeChAlg na niekoľko komponentov. To nám dá lepší obraz o tom, čo treba spraviť, ak chceme znížiť OTeChAlg nášho tréningového algoritmu L (a dosiahnuť tak lepšie výsledky).

Pripomeňme si, že náš algoritmus je všemohúci, vracia teda tú hypotézu z H , ktorá má najmenšiu chybu na tréningovej množine:

$$\hat{h}_T = \arg \min_{h \in H} (\text{Err}_T(h))$$

Označme h^* objektívne najlepšiu hypotézu z H , teda takú, ktorá má najmenšiu testovaciu chybu.

$$h^* = \arg \min_{h \in H} (\text{Err}(h))$$

Uvedomme si, že tieto dve hypotézy sa líšia iba tým, na akom pravdepodobnostnom rozdelení sú optimálne. Hypotéza \hat{h}_T je optimálna na tréningovej množine, ktorá je iba konečnou aproximáciou skutočného rozdelenia P , zatiaľ čo h^* je optimálna na tomto skutočnom rozdelení. Platia teda nerovnosti

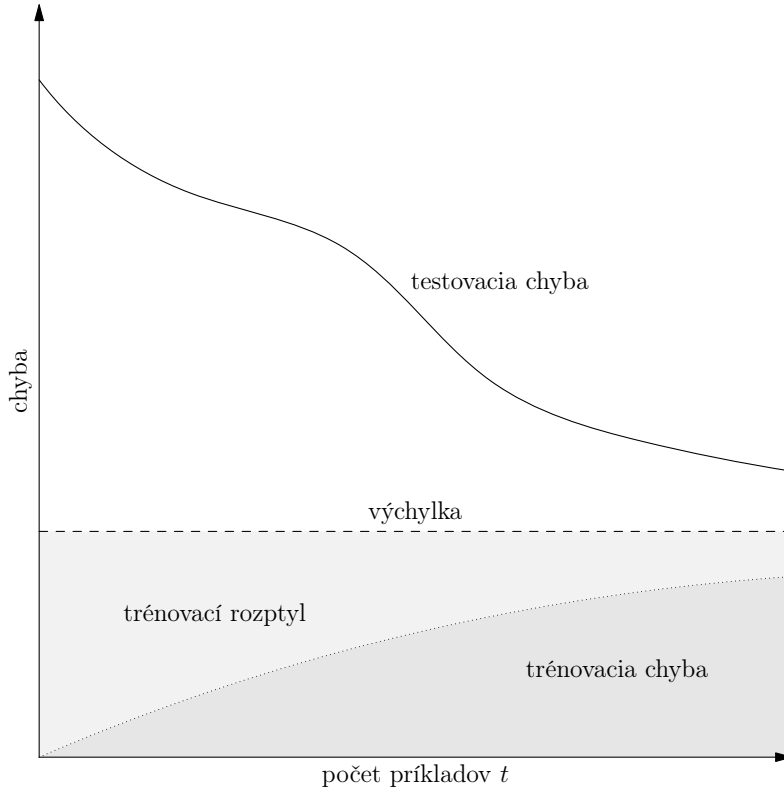
$$\begin{aligned} \text{Err}_T(\hat{h}_T) &\leq \text{Err}_T(h^*) \\ \text{Err}(h^*) &\leq \text{Err}(\hat{h}_T) \end{aligned}$$

Toto platí pre ľubovoľnú tréningovú množinu T . Keď to teda vezmeme dokopy cez všetky možné T , dostaneme “očakávanú” formu týchto nerovností:

$$\begin{aligned} \text{Err}_T(L) &\leq \mathbb{E}_T [\text{Err}_T(h^*)] \\ \text{Err}(h^*) &\leq \text{Err}(L) \end{aligned}$$

Pritom očakávaná tréningová chyba h^* je rovnaká, ako jej očakávaná testovacia chyba, nakoľko h^* je nezávislé od tréningovej množiny. (Akokoľvek je tréningová množina vybraná, z pohľadu h^* to vyzerá tak, akoby sme vybrali niekoľko náhodných príkladov z P .) Dostávame tak

$$\text{Err}_T(L) \leq \text{Err}(h^*) \leq \text{Err}(L).$$

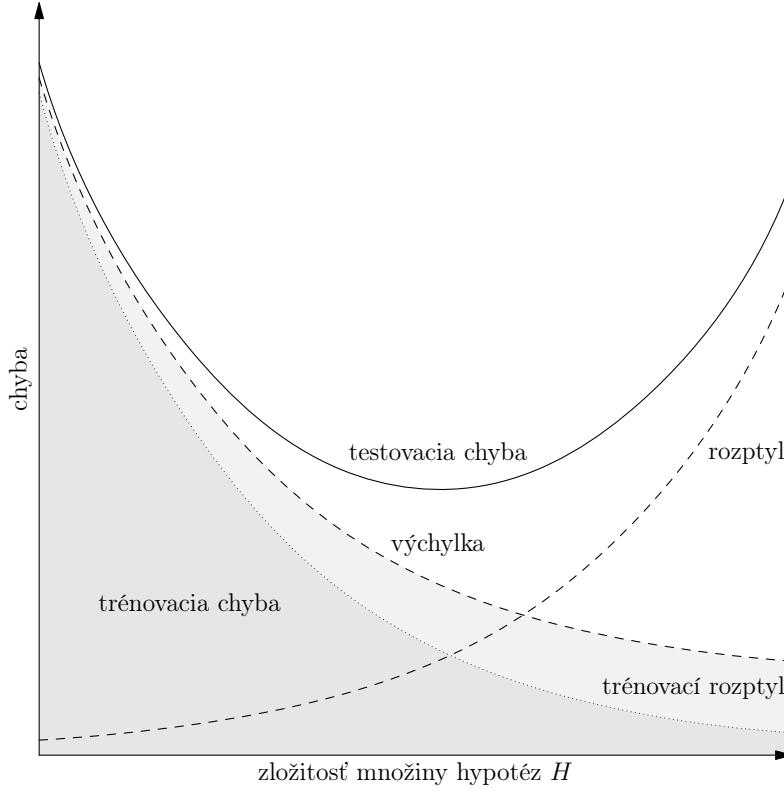
Obr. 1.1: Závislosť chyby algoritmu od počtu tréningových príkladov t .

OTeCh nášho algoritmu sa od najlepšej možnej chyby líši o hodnotu $\epsilon_{\text{est.}} = \text{Err}(L) - \epsilon_{\text{approx.}}$, túto hodnotu budeme volať *rozptyl* (po anglicky *estimation error*). Je spôsobený tým, že sme pri minimalizácii chyby neuvažovali skutočné rozdelenie P (ktoré nepoznáme), ale iba niekoľko tréningových príkladov. To, čo je najlepšie v T , nie je nutne najlepšie v P . Čím väčšia je ale tréningová množina, tým viac sa (pravdepodobne) podobá na skutočné rozdelenie P a tým nižší rozptyl.

Prostredný člen budeme volať *výchylka* a označovať $\epsilon_{\text{approx.}}$ (po anglicky *approximation error*). Vyjadruje chybu, ktorá je spôsobená tým, že sa náš algoritmus obmedzil na nejakú konkrétnu množinu hypotéz H . Čím väčšia množina hypotéz, tým menšia výchylka: keďže h^* je najlepšia hypotéza v množine H , jej zväčšením si môžeme iba prilepšiť. Zložitejšia množina hypotéz sa ale ľahšie “napasuje” na ľubovoľné tréningové dáta. To zvyšuje riziko toho, že výsledná hypotéza bude špecifická pre tréningové dáta a nebude schopná generalizácie. Teda čím väčšia množina hypotéz, tým väčší rozptyl. (Neplatí to úplne v každom prípade, dajú sa skonštruovať situácie, kedy zväčšením množiny hypotéz sa rozptyl zachová alebo dokonca zmenší. V prirodzených situáciách to ale do istej miery platí.)

OTrCh nášho algoritmu sa od očakávanej chyby h^* na tréningových dátach líši o $\epsilon_{\text{est.}}^T = \text{Err}(h^*) - \text{Err}_T(L)$, túto hodnotu budeme volať *trénovací rozptyl*. Na skutočnom rozdelení P sa nám nemôže dariť lepšie, ako hypotéze h^* ; fakt, že na tréningových dátach sa nám darí lepšie, je spôsobený tým, že to, čo je najlepšie v P , nemusí byť najlepšie v T . Opäť ale platí, že čím väčšia tréningová množina, tým viac sa (pravdepodobne) podobá na P a tým nižší trénovací rozptyl.

Na základe dosiaľ uvedeného vieme graficky znázorniť, ako sa zhruba správajú OTeChAlg, rozptyl, výchylka, trénovací rozptyl a OTrChAlg, v závislosti od veľkosti tréningovej množiny (obrázok 1.1) a od zložitosti množiny hypotéz (obrázok 1.2).

Obr. 1.2: Závislosť chyby algoritmu od zložitosti množiny hypotéz H .

1.3.3 Kompromis medzi výchylkou a rozptylom (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

$$h^* = \arg \min_{h \in H} (\text{Err}(h)).$$

Budeme upravovať výraz reprezentujúci OTechAlg.

$$\text{Err}(L) = \mathbb{E}_T [\text{Err}(\hat{h}_T)] \quad (1.3)$$

$$= \mathbb{E}_T \left[\mathbb{E}_{x,y} [(\hat{h}_T(x) - y)^2] \right] \quad (1.4)$$

$$= \mathbb{E}_T \left[\mathbb{E}_{x,y} \left[\left((\hat{h}_T(x) - h^*(x)) + (h^*(x) - y) \right)^2 \right] \right] \quad (1.5)$$

Posledná rovnosť vychádza z netriviálneho technického kroku, ktorý si vyžaduje dodatočné predpoklady. Tieto technické detaily však prenecháme na koniec kapitoly.

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

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

$$\text{Err}(L) = \underbrace{\mathbb{E}_T \left[\mathbb{E}_{x,y} [(\hat{h}_T(x) - h^*(x))^2] \right]}_{\text{rozptyl}} + \underbrace{\mathbb{E}_{x,y} [(h^*(x) - y)^2]}_{\text{výchylka}}.$$

Prvý zo sčítancov budeme volať *rozptyl*. Trénovací algoritmus s malým rozptylom vracia funkcie, ktoré sú blízko optima v množine H . Tým, že mu zväčšíme množinu trénovacích dát, si veľmi neprilepíme. Naopak, algoritmus s veľkým rozptylom vracia funkcie ďaleko od optima, vieme sa teda k optimu priblížiť tým, že zväčšíme množstvo trénovacích 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 (nakoľko h^* je najlepšia hypotéza v množine H , jej zväčšením si môžeme iba prilipsiť). Zložitejšia množina hypotéz sa ale ľahšie “napasuje” na ľubovoľné trénovacie dáta. To zvyšuje riziko toho, že výsledná hypotéza bude špecifická pre trénovacie dáta a nebude schopná generalizácie. Teda čím zložitejšia množina hypotéz, tým väčší rozptyl.

Výchylku 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 nerozlíšiteľné. Takže na meranie chyby h^* môžeme použiť trénovacie dáta, berúc v úvahu ich náhodný výber:

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

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

Použitím ďalšieho technického kroku dostaneme:

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

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

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

nakoľko \hat{h}_T 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 je naša vzorka.

Druhý zo sčítancov je OTrChAlg. Je to stredná hodnota chyby, ktorej sa dopustí výstup z algoritmu \hat{h}_T na tých istých dátach, pomocou ktorých sme \hat{h}_T zostrojili.

Platí

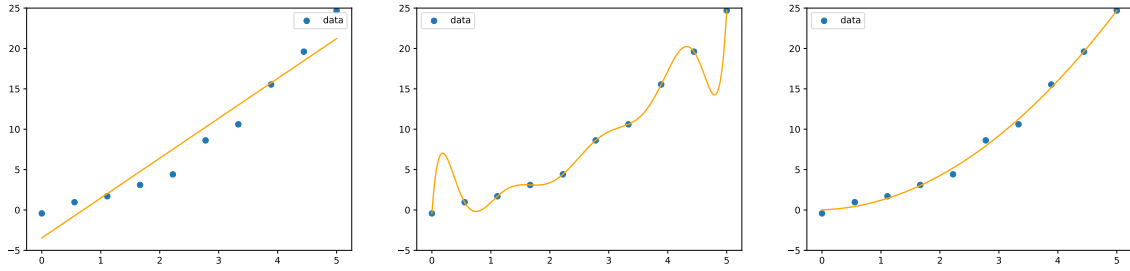
$$\text{Err}_T(L) \leq \text{výchylka} \leq \text{Err}(L).$$

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, na ktorých sa žiadnej hypotéze z H nedarí.

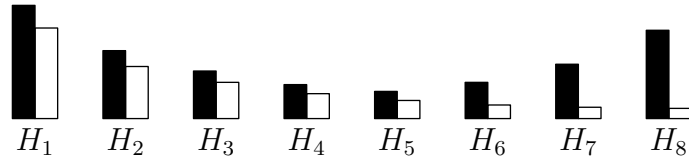
Na základe dosiaľ uvedeného vieme graficky znázorniť, ako sa zhruba správajú OTeChAlg, rozptyl, výchylka, trénovací rozptyl a OTrChAlg, v závislosti od veľkosti trénovacej množiny (obrázok 1.1) a od zložitosti množiny hypotéz (obrázok 1.2).

1.4 Podučenie/preučenie

Ak máme fixné trénovacie dáta T , pri voľbe množiny hypotéz H sa snažíme nájsť kompromis medzi malým rozptylom a malou výchylkou. Zložité H bude mať malú výchylku ale veľký rozptyl, čo vedie k tzv. *preučeniu* (angl. *overfitting*). Jednoduché H bude mať malý rozptyl, ale veľkú výchylku, tzv. *podučenie* (angl. *underfitting*).



Obr. 1.3: Podučenie, preučenie, akurát.

Obr. 1.4: Testovacie (čierne) a trénovacie (biele) chyby v čím ďalej, tým zložitejších množinách hypotéz H .

Na obrázku 1.3 ilustrujeme oba koncepty. Úlohou je modelovať kvadratickú funkciu, ku ktorej sme pridali malé množstvo šumu. Rozdelenie P teda vracia nejaké x (napríklad z intervalu $\langle 0, 10 \rangle$) a $y = x^2 + \varepsilon$, kde ε je zvolené náhodne z intervalu $\langle -1, 1 \rangle$. Ak za množinu hypotéz zvolíme lineárne funkcie, vďaka ich jednoduchosti už pri malých trénovacích množinách bude trénovacia chyba blízko očakávanej testovacej chyby (nízky rozptyl). Za to ale budú všetky chyby vysoké (vysoká výchylka). Ak za množinu hypotéz zvolíme polynómy nejakého vysokého stupňa, ľahko nájdeme polynóm presne prechádzajúci cez trénovacie dáta (nízka trénovacia chyba), avšak mimo nich bude dávať výsledky úplne mimo (vysoká testovacia chyba, a teda vysoký rozptyl).

1.4.1 Regularizácia

Predstavme si, že máme na výber z viacerých množín hypotéz H_1, H_2, \dots , čím ďalej tým zložitejších. Teda $H_1 \subseteq H_2 \subseteq H_3 \subseteq \dots$. Ak by sme si graficky znázornili testovacie a trénovacie chyby najlepších hypotéz z jednotlivých množín, vyzeralo by to zhruba ako na obrázku 1.4.

Z ktorej množiny hypotéz chceme vybrať? Pri trénovaní sa snažíme nájsť hypotézu h , ktorá minimalizuje chybu na trénovacích dátach. Táto chyba sa ale od očakávanej testovacej chyby líši zhruba o $\epsilon_{\text{est.}} + \epsilon_{\text{est.}}^T$ (v očakávanom prípade).

V prístupe zvanom *regularizácia* do minimalizovaného výrazu umelo pridáme *pokutu*, ktorá aproximuje $\epsilon_{\text{est.}} + \epsilon_{\text{est.}}^T$. Tento člen označíme $\Lambda(h)$. Z čím zložitejšej množiny hypotéza h je, tým väčšia bude pokuta. Výstupom algoritmu potom je

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

Uvedomme 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 $\Lambda(h)$. Do jednej množiny patria tie hypotézy, ktoré majú rovnakú penalizáciu.

Vo všetkých prípadoch je pokuta parametrizovaná reálnym parametrom λ hovoriacim, ako veľké pokuty chceme udeľovať. Veľké λ (v porovnaní s trénovacou chybou) nám hovorí, že sa snažíme hlavne dosiahnuť jednoduché hypotézy; s malým λ zas kladíme dôraz na hypotézy s menšou trénovacou chybou.

Uvedieme si niekoľko príkladov výrazov, ktoré sú bežne používané ako pokuta. 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*).

$$\Lambda(h) = \lambda \cdot \|(a_1, a_2, \dots, a_n)\|^2 = \lambda \cdot (a_1^2 + a_2^2 + \dots + a_n^2)$$

Táto pokuta “tlačí” váhy nepotrebných atribútov do nuly, pričom väčšie váhy tlačí viac. Takže čím dôležitejší atribút, tým väčšiu váhu si môže dovoliť mať.

- L_1 regularizácia (známa aj ako *lasso*).

$$\Lambda(h) = \lambda \cdot (|a_1| + |a_2| + \dots + |a_n|)$$

Opäť “tlačíme” nepotrebné atribúty do nuly, pričom ale všetky váhy tlačíme rovnako. To nám vie vynulovať nepotrebné atribúty, čo nám vie znížiť výpočtové nároky: nemusíme pri výpočte uvažovať tie členy, ktoré majú nulový koeficient.

1.4.2 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 do úvahy iba hypotézy z tej množiny hypotéz, ktorú sme identifikovali ako najvhodnejšiu. Výsledná hypotéza je výstupom algoritmu.

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.

k-fold evaluation. Pri holdout testovaní 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 trénovaciu 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 trénovacie dáta.

Leave-one-out. Ak chceme zmerať testovaciu chybu výstupnej hypotézy, musíme si na to rezervovať ďalšiu časť dát: *testovaciu množinu*. Tú nepoužívame ani pri trénovaní, ani pri validácii. Iba úplne na konci celého procesu na nej vypočítame chybu výslednej hypotézy.

1.5 Cvičenia

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

1.1. Je rozumné predpokladať, že s väčším množstvom trénovacích dát sa nám bude testovacia chyba zmenšovať. 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 rozdelenie P , pre ktoré sa bude $\text{Err}(\hat{h}_T)$ so zvyšujúcim sa počtom trénovacích chýb *zvyšovať*. Na množinu hypotéz je kladená jedna podmienka: pre každú možnú trénovaciu množinu T musí existovať hypotéza v H , ktorá minimalizuje trénovaciu chybu. (Teda vždy musí existovať minimum, môže ich byť ale viac. Pre všeobecné H vieme povedať iba to, že existuje infimum.)

1.2. Za určitých podmienok ale skutočne platí, že viac trénovacích 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, že pre $t \rightarrow \infty$ sa bude OTeCh hypotézy \hat{h}_T blížiť k OTeCh najlepšej novej hypotézy h^* . Inak zapísané, dokážte

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

1.3. Jednou výhodou L_2 regularizácie oproti L_1 regularizácie je, že sa ľahšie minimalizuje výsledný výraz. Ako príklad uvedieme 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 vstupy jednotlivých trénovacích príkladov. Ďalej nech y je stĺpcový vektor cieľových výstupov na jednotlivých príkladoch. Je známe, že 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 tiež, že takéto explicitné vyjadrenie nie je možné priamočiaro získať pre L_1 regularizáciu.

1.6 Appendix: iný kompromis medzi inou výchylkou a iným rozptylom

V literatúre pod názvom *bias-variance tradeoff* vystupuje odlišný výsledok, ako vyššie spomenutý bias-complexity tradeoff. Uvádzame ho, pretože sa často zamieňajú a je v tom zmätok. Pokúsime sa vyjasniť, kde sú rozdiely medzi týmito dvomi výsledkami.

Veta 1.1. Zamerajme sa na jeden konkrétny vstup x , a merajme chybu výslednej hypotézy \hat{h}_T iba na tomto vstupe. (Stále ale môžeme dostať rôzne y .) Na meranie chyby použijeme kvadratickú odchýlku. Označme očakávanú hodnotu tejto chyby (cez všetky možné trénovacie množiny T a výstupy y) ako $\text{Err}_{|x=x}(L)$. Tvrdíme, že sa dá vyjadriť nasledovne:

$$\text{Err}_{|x=x}(L) = \underbrace{\text{Var}_T(\hat{h}_T(x))}_{\text{rozptyl}} + \underbrace{\left(\mathbb{E}_T[\hat{h}_T(x)] - h^\square(x) \right)^2}_{\text{výchylka}^2} + \underbrace{\mathbb{E}_y[\varepsilon^2]}_{\text{šum}},$$

kde $\varepsilon := y - h^\square(x)$.

Zastavme sa najprv nad tým, ako toto tvrdenie interpretovať. Pre každú možnú vzorku tréningových dát T dostaneme niekú inú hypotézu \hat{h}_T . Rozptyl je nízky práve vtedy, keď budú všetky tieto hypotézy dávať podobné výsledky. Ak je vysoký, znamená to, že výsledná hypotéza je veľmi citlivá na tréningové dáta; oplatí sa preto zväčšiť ich množstvo. (Ak chceme byť veľmi skeptický, nič nezaručuje, že zväčšením tréningovej množiny sa rozptyl zníži. Možno existuje iný dôvod, kvôli ktorému je vysoký; takmer určite sa dajú skonštruovať takéto umelé protipríklady.)

Keď už uvažujeme všetky možné výsledné hypotézy \hat{h}_T , môžeme sa pozrieť na to, aký je ich “priemerný odhad”: ak by sme spriemerovali všetky ich výsledky, čo by sme dostali? Presne $E_T[\hat{h}_T(x)]$; výchylka-na-druhú potom meria, o koľko sa takáto priemerná hypotéza líši od najlepšej možnej funkcie h^\square . Ak je výchylka vysoká ale variancia nízka, znamená to, že takmer všetky \hat{h}_T majú problém na vstupe x , treba teda zväžiť voľbu zložitejšej množiny hypotéz.

Nakoniec, šum zodpovedá ireducibilnej chybe, ktorú bude mať každá hypotéza. Je to v dôsledku toho, že jednému x môže pripadať viacero rôznych y . Presne túto chybu nadobúda najlepšia hypotéza-funkcia h^\square (pre ktorú sú rozptyl aj výchylka-na-druhú nulové).

Dôkaz. Upravujeme pôvodný výraz.

$$\text{Err}_{|x=x}(L) = E_{T,y} \left[(\hat{h}_T(x) - y)^2 \right] \quad (1.9)$$

$$= E_{T,y} \left[((\hat{h}_T(x) - h^\square(x)) - \varepsilon)^2 \right] \quad (1.10)$$

$$= E_{T,y} \left[(\hat{h}_T(x) - h^\square(x))^2 + \varepsilon^2 - 2 \cdot \varepsilon \cdot (\hat{h}_T(x) - h^\square(x)) \right] \quad (1.11)$$

$$= E_T \left[(\hat{h}_T(x) - h^\square(x))^2 \right] + E_y [\varepsilon^2] - E_{T,y} \left[2 \cdot \varepsilon \cdot (\hat{h}_T(x) - h^\square(x)) \right] \quad (1.12)$$

$$= E_T \left[(\hat{h}_T(x) - h^\square(x))^2 \right] + E_y [\varepsilon^2] - 2 \cdot E_y [\varepsilon] \cdot E_T \left[(\hat{h}_T(x) - h^\square(x)) \right] \quad (1.13)$$

$$= E_T \left[(\hat{h}_T(x) - h^\square(x))^2 \right] + E_y [\varepsilon^2] \quad (1.14)$$

Výraz sme upravili, potom v kroku 1.11 roznásobili a v kroku 1.12 využili linearitu strednej hodnoty. Ďalej v kroku 1.13 sme využili, že stredná hodnota súčinu nezávislých premenných je súčinom stredných hodnôt tých premenných. Nakoniec v kroku 1.14 využívame $E[\varepsilon] = 0$. Zamerajme sa ďalej na prvý sčítanec.

$$= E_T \left[(\hat{h}_T(x) - h^\square(x))^2 \right] \quad (1.15)$$

$$= E_T \left[\hat{h}_T(x)^2 + h^\square(x)^2 - 2 \cdot \hat{h}_T(x) \cdot h^\square(x) \right] \quad (1.16)$$

$$= E_T \left[\hat{h}_T(x)^2 \right] + E_T \left[h^\square(x)^2 \right] - E_T \left[2 \cdot \hat{h}_T(x) \cdot h^\square(x) \right] \quad (1.17)$$

$$= \left(\text{Var}(\hat{h}_T(x)) + E_T \left[\hat{h}_T(x) \right]^2 \right) + \left(\text{Var}(h^\square(x)) + E_T \left[h^\square(x) \right]^2 \right) - E_T \left[2 \cdot \hat{h}_T(x) \cdot h^\square(x) \right] \quad (1.18)$$

$$(1.19)$$

V kroku 1.16 sme výraz roznásobili a potom v kroku 1.17 využili linearitu strednej hodnoty. V poslednom kroku (1.18) sme využili vzťah $\text{Var}(a) = E[a^2] - E[a]^2$. Pokračujme ďalej v úpravách.

$$= \left(\text{Var}(\hat{h}_T(x)) + E_T \left[\hat{h}_T(x) \right]^2 \right) + h^\square(x)^2 - 2 \cdot E_T \left[\hat{h}_T(x) \right] \cdot h^\square(x) \quad (1.20)$$

$$= \text{Var}(\hat{h}_T(x)) + \left(E_T \left[\hat{h}_T(x) \right]^2 + h^\square(x)^2 - 2 \cdot E_T \left[\hat{h}_T(x) \right] \cdot h^\square(x) \right) \quad (1.21)$$

$$= \text{Var}(\hat{h}_T(x)) + \left(E_T \left[\hat{h}_T(x) \right] - h^\square(x) \right)^2 \quad (1.22)$$

V kroku 1.20 sme využili, že x je vopred dané, a teda $h^\square(x)$ je konštanta. Má teda nulový rozptyl a jeho stredná hodnota je identická s jeho hodnotou. Ďalej sme už len upravovali. Keď to celé dáme do jednej rovnice, dostaneme

$$\text{Err}_{|x=x}(L) = \underbrace{\text{Var}_T(\hat{h}_T(x))}_{\text{rozptyl}} + \underbrace{\left(\mathbb{E}_T[\hat{h}_T(x)] - h^\square(x)\right)^2}_{\text{výchylka}^2} + \underbrace{\mathbb{E}_y[\varepsilon^2]}_{\text{šum}}.$$

□

1.7 Appendix: technické detaily k bias-variance tradeoff

V tejto časti dokážeme tvrdenie, ktoré sme použili pri odvodení vzťahu (1.3.3).

Pripomeňme si, že $h^\square(x)$ je cieľová funkcia, ktorá je definovaná ako $h^\square(x) = \mathbb{E}_{y|x=x}[y]$ a táto cieľová funkcia nemusí byť súčasťou triedy hypotéz H . Funkcia $h^*(x)$ je taká funkcia z H , ktorá má najmenšiu možnú očakávanú chybu, t.j. $h^* = \arg \min_{h \in H} (\mathbb{E}_{x,y}[(h(x) - y)^2])$. Túto definíciu vieme vziať ku funkcii h^\square pomocou nasledujúcej lemy.

Lemma 1.2. *Pre funkciu $h^*(x)$ platí:*

$$h^* = \arg \min_{h \in H} \left(\mathbb{E}_x[(h(x) - h^\square(x))^2] \right).$$

Dôkaz. Úpravou výrazu z definície h^* získame:

$$\begin{aligned} \mathbb{E}_{x,y}[(h(x) - y)^2] &= \mathbb{E}_{x,y}[(h(x) - h^\square(x)) + (h^\square(x) - y)]^2 \\ &= \mathbb{E}_{x,y}[\underbrace{(h(x) - h^\square(x))^2}_{\text{nezávisí od } y}] + \mathbb{E}_{x,y}[(h^\square(x) - y)^2] \\ &\quad + 2 \mathbb{E}_{x,y}[(h(x) - h^\square(x))(h^\square(x) - y)] \\ &= \mathbb{E}_x[(h(x) - h^\square(x))^2] + \mathbb{E}_{x,y}[(h^\square(x) - y)^2] \\ &\quad + 2 \mathbb{E}_x[(h(x) - h^\square(x))(h^\square(x) - \mathbb{E}_{y|x=x}[y])] \end{aligned}$$

Druhý sčítanec je konštanta vzhľadom na h , preto ho pri hľadaní minima nemusíme uvažovať. Z definície $h^\square(x) = \mathbb{E}_{y|x=x}[y]$, preto tretí sčítanec bude nulový. Zostáva nám teda:

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

čo bolo treba dokázať. □

Z predchádzajúcej lemy teda vyplýva, že funkcia h^* je priemetom funkcie h^\square do H . Pre zjednodušenie si v tejto kapitole zavedieme označenie $\langle f, g \rangle = \mathbb{E}_x[f(x)g(x)]$. Takisto si zavedieme pojem *vzdialenosti funkcií* $\text{dist}^2(f, g) = \mathbb{E}_x[(f(x) - g(x))^2] = \langle f(x) - g(x), f(x) - g(x) \rangle$. Funkcia $h^*(x)$ je tak vlastne definovaná ako funkcia z H s najmenšou vzdialenosťou od h^\square .

Lemma 1.3. *Nech množina hypotéz H je uzavretá na lineárne kombinácie. Potom pre ľubovoľnú funkciu $h \in H$ platí:*

$$\mathbb{E}_x[h(x)(h^*(x) - h^\square(x))] = \langle h, h^* - h^\square \rangle = 0.$$

Dôkaz. Nech h je ľubovoľná funkcia z H . Uvažujme triedu funkcií F , ktoré sú tvaru $f_\Delta = h^* + \Delta h$, kde Δ je ľubovoľné reálne číslo. Všimnime si, že $F \subseteq H$, lebo H je uzavretá na lineárne kombinácie.

Vzhľadom ku leme 1.2, ktorá charakterizuje h^* , musí byť zo všetkých funkcií z F najbližšou k h^\square funkcia $f_0 = h^*$, z čoho vyplýva, že derivácia funkcie $D(\Delta) = \text{dist}^2(f_\Delta, h^\square)$ musí byť v bode $\Delta = 0$ nulová.

Upravme funkciu $D(\Delta)$:

$$D(\Delta) = \text{dist}^2(f_\Delta, h^\square) \quad (1.23)$$

$$= \langle h^* + \Delta h - h^\square, h^* + \Delta h - h^\square \rangle \quad (1.24)$$

$$= \langle (h^* - h^\square) + \Delta h, (h^* - h^\square) + \Delta h \rangle \quad (1.25)$$

$$= \langle h^* - h^\square, h^* - h^\square \rangle + \Delta^2 \langle h, h \rangle + 2\Delta \langle h^* - h^\square, h \rangle \quad (1.26)$$

Vypočítame teraz deriváciu $D(\Delta)$ podľa Δ :

$$D'(\Delta) = \frac{d\langle h^* - h^\square, h^* - h^\square \rangle}{d\Delta} + \frac{d\Delta^2 \langle h, h \rangle}{d\Delta} + \frac{d(2\Delta \langle h^* - h^\square, h \rangle)}{d\Delta} \quad (1.27)$$

$$= 2\Delta \langle h, h \rangle + 2\langle h^* - h^\square, h \rangle \quad (1.28)$$

Teraz položíme $\Delta = 0$, a keďže derivácia musí byť v tomto bode nulová, dostávame rovnicu:

$$0 = d'(0) = 2\langle h^* - h^\square, h \rangle, \quad (1.29)$$

čo bolo treba dokázať. \square

Vráťme sa teraz k rovnici (1.3.3). Chceme ukázať, že

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

čiže potrebujeme ukázať, že

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

Platí:

$$\begin{aligned} \mathbb{E}_{x,y} \left[(\hat{h}_T(x) - h^*(x)) \cdot (h^*(x) - y) \right] &= \mathbb{E}_{x,y} \left[(\hat{h}_T(x) - h^*(x)) \cdot (h^*(x) - h^\square(x) + h^\square(x) - y) \right] \\ &= \mathbb{E}_{x,y} \left[\underbrace{(\hat{h}_T(x) - h^*(x)) \cdot (h^*(x) - h^\square(x))}_{\text{nezávisí od } y} \right] \\ &\quad + \mathbb{E}_{x,y} \left[(\hat{h}_T(x) - h^*(x)) \cdot (h^\square(x) - y) \right] \\ &= \mathbb{E}_x \left[(\hat{h}_T(x) - h^*(x)) \cdot (h^*(x) - h^\square(x)) \right] \\ &\quad + \mathbb{E}_x \left[(\hat{h}_T(x) - h^*(x)) \cdot (h^\square(x) - \mathbb{E}_{y|x=x} [y]) \right]. \end{aligned}$$

Keďže z definície $h^\square(x) = \mathbb{E}_{y|x=x} [y]$, druhý sčítanec sa vynuluje a dostávame:

$$\begin{aligned} \mathbb{E}_{x,y} \left[(\hat{h}_T(x) - h^*(x)) \cdot (h^*(x) - y) \right] &= \mathbb{E}_x \left[(\hat{h}_T(x) - h^*(x)) \cdot (h^*(x) - h^\square(x)) \right] \\ &= \langle \hat{h}_T - h^*, h^* - h^\square \rangle. \end{aligned}$$

Uvedomme si, že funkcia $h = \hat{h}_T - h^*$ je lineárnou kombináciou dvoch funkcií z H , a preto tiež patrí do H . Preto môžeme použiť Lemu 1.3 a dostávame pre ľubovoľnú tréningovú množinu T :

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

a teda aj:

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

čo sme chceli dokázať.

Poznámka 1.1. Označenie $\langle f, g \rangle$ sa vizuálne nepodobá na skalárny súčin vektorov náhodou. V skutočnosti je toto označenie priamočiarym rozšírením skalárneho súčinu vektorov na funkcie a množstvo vlastností skalárneho súčinu vektorov možno analogicky preniesť aj na takto definované skalárne súčiny funkcií. Analogicky ku kolmosti vektorov možno zdefinovať aj kolmosť funkcií $f \perp g$, ak $\langle f, g \rangle = 0$. Lema 1.3 je v takom prípade ekvivalentom tvrdenia o kolmosti v prípade projekcií vektorov do podpriestoru. Napríklad pre priamku p a bod x platí, že ak y je najbližší bod priamky p ku x , tak vektor $x - y$ je kolmý na p . V tomto texte je H analógom priamky p , h^* je analógom bodu y a h^\square je analógom bodu x .

Poznámka 1.2. Všetky tvrdenia doposiaľ predpokladali, že v množine hypotéz H existuje funkcia h^* najbližšia ku H . Pri niektorých množinách hypotéz však tento predpoklad nemusí byť splnený. Môže sa nám totiž stať, že vieme zostrojiť nekonečnú postupnosť funkcií, kde každá ďalšia funkcia je bližšie ku h^\square , no táto postupnosť funkcií konverguje ku funkcii mimo H . Aby sme zaručili, že takýto prípad nenastane, musíme do predpokladov zahrnúť aj požiadavku, že H je uzavreté na limity.

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:

1. “Akú OTeChAlg môžeme očakávať pri danej veľkosti trénovacej množiny t ?”
2. “Pri danom t , s akou pravdepodobnosťou nám algoritmus vráti hypotézu, ktorej OTeCh je menšia rovná ε ?”

Na základe odpovedí na tieto dve otázky potom budeme schopní zodpovedať nasledovné (dalo by sa povedať, že “inverzné”) otázky:

3. “Akú veľkú trénovaciu množinu máme zvoliť, aby sme dosiahli dostatočne malú ($\leq \varepsilon$) OTeChAlg?”
4. “Aké t máme zvoliť, aby sme iba s malou pravdepodobnosťou ($\leq \delta$) dostali zlú hypotézu (s OTeCh väčšou ako ε)?”

Odtiaľ sa odvíja názov *PAC učenie* (z anglického *probably approximately correct learning*). Presnejšie, PAC učenie sa zaoberá hlavne otázkami 2 a 4; my sa ale budeme zaoberať všeobecne tým, ako závisia rôzne metriky (OTeChAlg, ...) od počtu trénovacích príkladov t .

Uvedomme si rozdiel medzi otázkami (1, 3) a otázkami (2, 4). V prvom zmienenom sa pýtame na očakávaný prípad; ak ale chceme dosiahnuť dobrú hypotézu nie v očakávanom prípade, ale zaručene, vieme to zaručiť iba s určitou pravdepodobnosťou. Vždy sa totiž môže stať (s malou pravdepodobnosťou), že si vytiahneme nereprezentatívne trénovacie dáta. Je teda nutné hovoriť aj o veľkosti tejto záruky, takže hranica δ je potrebná.

2.1 Definície, označenia a predpoklady

Hypotézu h nazveme *zlou*, ak jej OTeCh je väčšia ako ε . Inak ju nazveme *dobrou*. Tieto pojmy budeme používať tak, aby bola hranica ε vždy jasná z kontextu.

Budeme riešiť hlavne binárne klasifikačné úlohy, v ktorých je cieľom pre každý vstup x povedať, či spĺňa určitú podmienku ($y = 1$) alebo nespĺňa ($y = 0$). Napríklad môžeme chcieť rozlíšiť obrázky 32×32 ktoré obsahujú aspoň jednu mačku od tých, ktoré mačky neobsahujú. Príklady, kde $y = 1$ nazveme *pozitívnymi*; ostatné príklady nazveme *negatívnymi*. Hypotézy v klasifikačných úlohách budeme nazývať *klasifikátormi*.

Pripomíname, že chyba klasifikátora h sa počíta ako

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

Budeme pracovať so štandardným modelom učenia (1.1) spolu s predpokladom všemohúceho trénovacieho algoritmu (vždy vráti hypotézu s najmenšou chybou na trénovacej množine). Okrem toho budeme predpokladať, že množina hypotéz H je *úplná*: existuje v nej “úplne správny

klasifikátor”, ktorý na každom vstupe dáva správny výsledok, má teda nulovú testovaciu chybu. Tento klasifikátor je zrejme najlepšou možnou hypotézou v H , mohli by sme ho preto značiť štandardne h^* . My ho ale budeme značiť c ako “cieľová hypotéza”. Treba si ale dať pozor, že správnych hypotéz v H môže byť viac, a to v prípade, keď niektorým vstupom majú v rozdelení P pravdepodobnosť 0. (V takom prípade je c ľubovoľná jedna z týchto hypotéz.)

Uvedomme si, že nech si vytiahneme ľubovoľnú trénovaciu množinu T , vďaka úplnosti množiny H v nej vždy vieme nájsť klasifikátor s nulovou chybou na trénovacej množine. (Napríklad h^* , môžu tam ale byť aj ďalšie klasifikátory, s odlišnými výstupmi na vstupoch mimo T .) Preto trénovací algoritmus vždy nájde hypotézu s nulovou trénovacou chybou. Takýto klasifikátor budeme volať *konzistentný klasifikátor*; názov pochádza od toho, že jeho výstupy sú konzistentné s trénovacími príkladmi.

Za týchto predpokladov si budeme klásť nasledovnú otázku: ktoré množiny hypotéz H sa dajú naučiť? Pod “naučením” sa rozumie to, že čím viac trénovacích príkladov algoritmu poskytneme, tým lepšiu hypotézu \hat{h}_T môžeme očakávať; v limite $t \rightarrow \infty$ si môžeme byť takmer určite istí, že dostaneme veľmi dobrú hypotézu. Toto sformalizujeme do nasledovnej definície:

Definícia 1. Množina hypotéz je *PAC \star naučiteľná*, ak pre ľubovoľné pravdepodobnostné rozdelenie P nad $X \times Y$, ľubovoľne malé hranice $\varepsilon > 0$ a $\delta > 0$ existuje hranica t_0 na počet trénovacích príkladov taká, že pre všetky počty príkladov $t \geq t_0$ platí: zlú hypotézu (s chybou aspoň ε) dostaneme iba s malou pravdepodobnosťou (nanajvýš δ).

$$\forall P : \forall \varepsilon > 0, \delta > 0 : \exists t_0 : \forall t \geq t_0 : \Pr_{T \sim P^t} (\text{Err}(\hat{h}_T) \geq \varepsilon) \leq \delta$$

Používame názov PAC \star naučiteľná preto, lebo skutočná (v literatúre zaužívaná) definícia PAC naučiteľnosti je o čosi komplikovanejšia: berie väčší ohľad na výpočtové detaily, ktorým sa my chceme zatiaľ vyhnúť.

Niekedy budeme uvažovať konkrétny trénovací algoritmus; v takých prípadoch budeme hovoriť o PAC Δ naučiteľnosti.

2.2 Konečné množiny hypotéz

V tejto časti sa zameriame na konečné množiny hypotéz. Uvedieme niekoľko príkladov takých množín: logické formuly obmedzenej veľkosti, konečné automaty obmedzenej veľkosti, ... V princípe všetko, v čom nevystupujú reálne čísla; s tým že uvažujeme iba malé hypotézy.

2.2.1 Konečné množiny sú PAC \star naučiteľné

Veta 2.1. Pomocou počtu trénovacích príkladov t vieme zhora odhadnúť pravdepodobnosť, že nám algoritmus vráti zlú hypotézu, nasledovne:

$$\Pr_T(\text{Err}(\hat{h}_T) \geq \varepsilon) \leq |H| \cdot e^{-\varepsilon t}$$

Dôkaz. Predstavme si, že na začiatku vidíme 0 trénovacích príkladov, a postupne sa nám odkrývajú. Sledujme jednotlivé hypotézy: vždy, keď sa odkryje nejaký príklad, môže sa stať, že je nekonzistentný s hypotézou. Každú takúto hypotézu zo súťaže vylúčime. Po odkrytí všetkých t trénovacích príkladov zostávajú už len tie hypotézy, ktoré sú s príkladmi konzistentné a môžu byť výsledkom algoritmu. Ak ani jedna z týchto hypotéz nie je zlá, výsledkom bude dobrá hypotéza. Aká je pravdepodobnosť, že aspoň jedna z týchto preživších hypotéz je zlá?

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

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

Zlých hypotéz je nanajvýš toľko, koľko je všetkých hypotéz, teda $|H|$. Pravdepodobnosť, že aspoň jedna z nich bude konzistentná s príkladmi, je nanajvýš súčet pravdepodobností jednotlivých zlých hypotéz:

$$\Pr_T(\text{prežije aspoň jedna zlá hypotéza}) \leq |H| \cdot e^{-\varepsilon t}$$

Odkiaľ dostávame požadovanú nerovnosť. \square

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

Veta 2.2. *Platí*

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

Dôkaz. Ak nám algoritmus vráti dobrú hypotézu, vieme jej OTeCh odhadnúť zhora ako ε . Ak nám algoritmus vráti zlú hypotézu, jej OTeCh je nanajvýš 1. Z toho dostávame nasledovný horný odhad na OTeChAlg:

$$\text{Err}(L) = \mathbb{E}_T [\text{Err}(\hat{h}_T)] \leq \Pr_T(\text{Err}(\hat{h}_T) < \varepsilon) \cdot \varepsilon + \Pr_T(\text{Err}(\hat{h}_T) \geq \varepsilon) \cdot 1$$

Pritom sčítance na pravej strane vieme odhadnúť zhora:

$$\Pr_T(\text{Err}(\hat{h}_T) < \varepsilon) \leq 1$$

a druhý sčítanec vieme odhadnúť pomocou predchádzajúcej vety (2.1). Dostávame tak odhad

$$\text{Err}(L) \leq 1 \cdot \varepsilon + |H| \cdot e^{-\varepsilon t}.$$

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

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

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

Odtiaľ dosadením dostaneme požadovaný odhad na chybu algoritmu. \square

Na základe týchto dvoch viet vieme sformulovať postačujúce podmienky na t také, aby boli hranice ε a δ dostatočne malé.

Dôsledok 2.3. *Ľubovoľná konečná množina hypotéz H je PAC \star naučiteľná: pre ľubovoľné rozdelenie P , ľubovoľne malé hranice $\varepsilon > 0$ a $\delta > 0$, pre počet tréningových príkladov t splňajúci*

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

platí, že dostaneme zlú hypotézu iba s malou pravdepodobnosťou.

Dôkaz. Podľa vety 2.1 platí

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

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

$$|H| \cdot e^{-\varepsilon t} \leq \delta \tag{2.3}$$

$$\ln |H| - \varepsilon t \leq \ln \delta \tag{2.4}$$

$$\varepsilon t \geq \ln |H| - \ln \delta \tag{2.5}$$

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

\square

2.2.2 Problém konjunkcie

Máme n boolovských premenných. Množinou inštancií X sú všetky možné ohodnotenia týchto premenných, napríklad pre $n = 3$ je jedným možným priradením $x_1 = 0, x_2 = 1, x_3 = 0$. Každé priradenie si vieme predstaviť ako postupnosť bitov dĺžky n . Je zrejmé, že rôznych priradení je 2^n .

Množinou hypotéz sú všetky formuly tvaru $x_{i_1} \wedge x_{i_2} \wedge \dots \wedge x_{i_k}$. Sú to teda konjunkcie premenných, s tým že premenné vo formule nesmú vystupovať negované. Tieto formuly sú chápané ako funkcie, ktoré vracajú 1 vtedy, keď dané ohodnotenie premenných spĺňa túto formulu. Napríklad $x_1 \wedge x_3 \wedge x_4$ vráti 1 na všetkých vstupoch, kde $x_1 = 1, x_3 = 1$ a $x_4 = 1$ (ostatné premenné môžu mať ľubovoľnú hodnotu). Rôznych hypotéz je zrejme tiež 2^n .

Metafora. Vyššie uvedený popis je možno príliš formálny, problém sa dá vysvetliť aj prirodzenejším spôsobom. Predstavme si, že náš kamarát má batoh, do ktorého dal niektoré z n vecí. My chceme zistiť, ktoré z týchto vecí v batohu sú a ktoré v ňom nie sú. Množina inštancií sú teda jednotlivé veci, a hypotézy sú všetky možné obsahy batoha. Čo ale zodpovedá tréningovým príkladom? Predstavme si, že sa môžeme spýtať niekoľko otázok tvaru: “Obsahuje batoh iba týchto k vecí?” (t.j. nič iné nesmie obsahovať) a kamarát nám na každú otázku odpovie áno/nie. Potom tréningové príklady zodpovedajú práve takýmto otázkam, spolu s odpoveďami na ne.

Zavedieme prirodzenú notáciu pre túto metaforu: namiesto $h(x) = 1$ budeme písať $h \subseteq x$, a namiesto $h(x) = 0$ píšeme $h \not\subseteq x$.

Všemožný tréningový algoritmus. Uvidíme, že konzistentnosť s tréningovými príkladmi je realistická požiadavka: ukážeme, ako sa dá na základe tréningových príkladov zostrojiť nejaký konzistentný klasifikátor.

Veta 2.4. *Nech hypotéza \hat{h}_T obsahuje všetky tie veci, ktoré sa vyskytujú v každom pozitívnom tréningovom príklade. Potom \hat{h} je konzistentná so všetkými tréningovými príkladmi.*

Dôkaz. Ak (x, y) je pozitívny príklad, v batohu môže byť iba nejaká podmnožina vecí z x . Keď túto úvahu zopakujeme pre všetky pozitívne príklady, dostaneme niekoľko množín “povolených vecí”. Skutočný obsah batoha musí byť podmnožinou všetkých z nich, teda je podmnožinou ich prieniku. Toto zároveň zachytáva všetku informáciu zo všetkých pozitívnych príkladov, teda ľubovoľná hypotéza, ktorá je podmnožinou tohto prieniku, je konzistentná s pozitívnymi príkladmi.

Čo ale s negatívnymi príkladmi? Ak by sme zobrali príliš malú množinu vecí ako hypotézu, mohol by nastať problém s negatívnymi príkladmi. (Napríklad ak máme príklady $(0110, 0)$ a $(1100, 1)$, nemôžeme ako hypotézu zvoliť 0100 ani 0000 , lebo by to bolo v rozpore s prvým príkladom.)

Ak ale zoberieme celý prienik ako našu hypotézu \hat{h}_T , nemôže sa to stať. Totiž, skutočný obsah batoha je podmnožinou oného prieniku, a teda ak $\hat{h}_T \subseteq x$, tak nutne aj $c \subseteq x$. Obmenou tejto implikácie dostaneme, že ak $c \not\subseteq x$, tak aj $\hat{h}_T \not\subseteq x$, takže naša hypotéza dáva na negatívnych príkladoch rovnaké výsledky. \square

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

Veta 2.5. *Pre problém konjunkcie platí*

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

Veta 2.6. *Pre problém konjunkcie platí, že aby sme zlí hypotézu (s chybou aspoň ε) dostali iba s malou pravdepodobnosťou (nanajvýš δ), stačí zvoliť veľkosť trénovacej množiny nasledovne:*

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

Dolné odhady. Tieto odhady sú relatívne tesné. Ako veľmi ich vieme zlepšiť? Vieme dostať aj nejaké dolné odhady, t.j. tvrdenia hovoriace, že ak sa chceme nejakú množinu hypotéz naučiť, potrebujeme aspoň toľkoto trénovacích príkladov?

Netriviálny dolný odhad na všetky prípady sa nedá dostať, nakoľko pravdepodobnostné rozdelenie P môže byť degenerované a “uľahčiť algoritmu robotu”. (Napríklad keď P priradí všetkých 100% jednej dvojici (x, y) .) Prirodzenejšie je teda hľadať nejaké jedno ťažké rozdelenie P . Toto rozdelenie ale môže závisieť od konkrétného trénovacieho algoritmu; nestačí len predpokladať, že algoritmus je všemohúci, nakoľko možných výsledných hypotéz môže byť viacero a ťažkosť rozdelenia P závisí od toho, ako sa hypotéze darí mimo trénovacích príkladov.

Budeme teda hľadať tvrdenia tvaru: pre každý algoritmus L existuje rozdelenie P , na ktorom potrebuje L aspoň nejaké množstvo trénovacích príkladov, aby sa mu dobre darilo.

Dokazovať to v takomto tvare sa ale ukazuje byť ťažké. Problém je v tom, že nájsť konkrétne ťažké rozdelenie P môže byť ťažké. V matematike nie vždy vieme ľahko nájsť konkrétny príklad s nejakou vlastnosťou; možno ale vieme ľahko dokázať, že taký objekt musí existovať. Tak tomu bude aj teraz; v dôkaze budeme využívať nasledovnú všeobecnú metavu:

Veta 2.7. *Ak chceme dokázať, že existuje objekt o spĺňajúci $f(o) \geq l$, kde f reprezentuje nejakú vlastnosť objektu o a hodnota l je nejaká ľubovoľná hranica, stačí dokázať, že pre o náhodne vybrané z nejakého rozdelenia O platí*

$$\mathbb{E}_{o \sim O} [f(o)] \geq l.$$

Dôkaz. Ak je hodnota $f(o)$ v priemere aspoň l , tak aspoň v jednom konkrétnom prípade (t.j. pre jedno konkrétne o) musí byť tiež väčšia rovná l . Formálne zapísaná obmena tejto implikácie:

$$(\forall o \in O : f(o) < l) \implies \mathbb{E}_{o \sim O} [f(o)] < l.$$

□

V našom prípade namiesto hľadania jedného konkrétného ťažkého rozdelenia P ho zvolíme náhodne. Potom ukážeme, že OTechAlg bude vysoká (pričom stredná hodnota sa berie aj cez náhodné voľby rozdelenia P).

Veta 2.8. *Pre problém konjunkcie platí, že pre ľubovoľný algoritmus L existuje rozdelenie P také, že pre $t \geq n - 1$ platí pre OTechAlg platí nasledovný dolný odhad:*

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

Pre $t < n - 1$ platí iný dolný odhad:

$$\text{Err}(L) \geq \frac{1}{2e} = \Omega(1).$$

Dôkaz. Predstavme si, že v rozdelení P vystupujú (t.j. majú nenulovú pravdepodobnosť) iba tie vstupy, ktoré obsahujú presne $n - 1$ vecí, teda práve jedna vec v množine nie je. Označme tieto vstupy $x^{(1)}, \dots, x^{(n)}$ podľa toho, ktorá vec chýba.

Ľubovoľná hypotéza h vráti na vstupe $x^{(i)}$ hodnotu 1 práve vtedy, keď vec i v hypotéze nie je; hodnotu 0 vráti vtedy, keď tam vec i je. Tieto vstupy teda zodpovedajú otázkam tvaru: “Je pravda, že vec v nie je v batohu?”

Takže každý trérovací príklad nám dáva informáciu iba o prítomnosti jednej veci. Počas testovania teda vieme správnu odpoveď iba na tých vstupoch, ktoré sme dostali ako trérovacie príklady. Na ostatných môžeme vo všeobecnosti iba hádať, či v batohu príslušná vec je alebo nie je.

Takže na vstupoch, ktoré sú v trérovacej množine, budeme mať chybu 0. Ak bol obsah batoha (t.j. výstupy y) zvolený rovnomerne náhodne spomedzi všetkých 2^n možností, tak na ostatných vstupoch budeme mať očakávanú chybu $\frac{1}{2}$.

Označme si pravdepodobnosti pridelené jednotlivým vstupom p_1, \dots, p_n . OTeChAlg vieme zmerať tak, že si náhodne vytiahneme trérovacie príklady a náhodne vytiahneme testovací príklad. Potom na trérovacích príkladoch natrénujeme hypotézu \hat{h}_T , zmeriame jej chybu na testovacím príklade, a spravíme strednú hodnotu tohto procesu. Aká je šanca, že pri trérovaní vstup $x^{(i)}$ nedostaneme a potom si ho pri testovaní vytiahneme?

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

Celková pravdepodobnosť, že si pri testovaní vytiahneme vstup mimo trérovacej množiny, je potom súčtom týchto pravdepodobností cez všetky vstupy $x^{(i)}$:

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

Tieto prípady prispievajú do chyby $1/2$, ostatné prípady prispievajú 0. Dostávame tak nasledovný dolný odhad na chybu algoritmu:

$$\text{Err}(L) \geq \frac{1}{2} \cdot \left(\sum_{i=1}^n p_i \cdot (1 - p_i)^t \right)$$

Ako zvoliť p_1, \dots, p_n tak, aby sme dostali dobrý dolný odhad? Skúsme 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 $t \neq n-1$, nemôžeme zvoliť všetky pravdepodobnosti takéto: pre $t > n-1$ by bol súčet pravdepodobností primálny a pre $t < n-1$ priveľký. Rozoberme tieto prípady osobitne.

Prvý prípad ($t \geq n-1$). Všetkým vstupom okrem jedného priradíme pravdepodobnosť $1/(t+1)$; posledný vstup dostane všetkú zvyšnú pravdepodobnosť.

$$p_1 = p_2 = \dots = p_{n-1} = \frac{1}{t+1}, \quad p_n = 1 - \frac{n-1}{t+1}$$

Dosadíme a dostaneme tak odhad

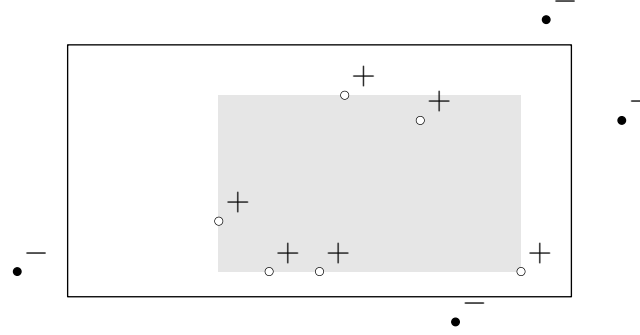
$$\text{Err}(L) \geq \frac{1}{2} \cdot \left((n-1) \cdot \frac{1}{t+1} \cdot \left(1 - \frac{1}{t+1} \right)^t + \frac{n-1}{t+1} \cdot \left(1 - \frac{n-1}{t+1} \right)^t \right).$$

Druhý sčítanec zdola odhadneme ako 0 a ďalej použijeme odhad

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

ktorý platí pre všetky prirodzené t . Dostaneme tak nerovnosť

$$\text{Err}(L) \geq \frac{1}{2e} \cdot \frac{n-1}{t+1}.$$



Obr. 2.1: Obdĺžniková hra: najmenší obdĺžnik konzistentný s príkladmi v porovnaní s cieľovým obdĺžnikom.

Druhý prípad ($t < n - 1$). Prvým $t + 1$ vstupom priradíme pravdepodobnosť $1/(t + 1)$, ostatné vstupy dostanú pravdepodobnosť 0. Dostaneme tak odhad

$$\text{Err}(L) \geq \frac{1}{2} \cdot (t + 1) \cdot \frac{1}{t + 1} \cdot \left(1 - \frac{1}{t + 1}\right)^t \quad (2.7)$$

$$\text{Err}(L) \geq \frac{1}{2e} \quad (2.8)$$

□

2.3 Nekonečné množiny hypotéz

Vyššie uvedené postupy pre konečné množiny hypotéz zlyhajú, ak $|H| = \infty$. Aj pre nekonečné množiny hypotéz je ale možné odvodiť horné odhady. Tie sa už nebudú odvíjať od veľkosti množiny H , ale od nejakej jej miery zložitosti. Touto mierou zložitosti bude tzv. *Vapnik-Chervonenkisova dimenzia*. Najprv sa ale pozrieme na jednoduchý príklad.

2.3.1 Obdĺžniková hra

Množinou inštancií sú všetky body roviny a množinou hypotéz sú všetky obdĺžniky, ktorých strany sú rovnobežné so súradnicovými osami. Obdĺžnik chápeme ako funkciu, ktorá vracia 1 pre tie body, ktoré sú vo vnútri obdĺžnika (vrátane jeho okrajov).

Pre jednoduchosť argumentu použijeme konkrétny trénovací algoritmus. Dôkaz by sa dal zovšeobecniť na prípad, kedy jediné, čo o algoritme predpokladáme je to, že je všemohúci. S nasledovným konkrétnym algoritmom ale bude analýza jednoduchšia.

Všemohúci trénovací algoritmus. Z trénovacích príkladov vezmeme tie pozitívne, a ako náš obdĺžnik \hat{h}_T zvolíme najmenší osovorovnobezný obdĺžnik obsahujúci body z týchto príkladov (tzv. *bounding box*). Túto konštrukciu ilustrujeme na obrázku 2.1.

Takýto obdĺžnik zrejme obsahuje všetky pozitívne trénovacie príklady, je teda s nimi konzistentný. Pretože je to najmenší taký obdĺžnik, musí byť celý vo vnútri cieľového obdĺžnika. Nemôže sa teda stať, že by nejaký bod nepatrí do c ale patrí do nášho obdĺžnika; je teda konzistentný aj s negatívnymi trénovacími príkladmi.

Veta 2.9. *Obdĺžniková hra je PAC Δ naučiteľná: pre ľubovoľné rozdelenie P , ľubovoľne malé hranice $\varepsilon > 0$ a $\delta > 0$, ak pre počet pozitívnych trénovacích príkladov t^+ platí*

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

tak nám algoritmus vráti zlú hypotézu (s OTeCh aspoň ε) iba s malou pravdepodobnosťou (nanajvýš δ).

Všimnime si, že v našom tvrdení vystupuje t^+ namiesto t . Ak rozdelenie P vráti pozitívny príklad s nenulovou pravdepodobnosťou, nie je to problém: s veľkým počtom tréningových príkladov t potom vieme dosiahnuť aj veľký počet pozitívnych tréningových príkladov t^+ .

Čo ak rozdelenie P pozostáva iba z negatívnych príkladov? V takom prípade náš algoritmus nepotrebuje žiaden príklad, hneď dosiahne OTeCh 100%. Vždy totiž vráti degenerovaný obdĺžnik, ktorý neobsahuje žiadne body (ani len na okraji).

V dôkaze budeme riešiť už len prvý prípad.

Dôkaz. Pod váhou množiny budeme rozumieť pravdepodobnosť, že náhodný bod z rozdelenia P padne do tejto množiny.

Zle klasifikované body roviny sú tie, ktoré sú vo vnútri cieľového obdĺžnika c ale nie sú vo vnútri nášho obdĺžnika \hat{h}_T . Každý z týchto zlých bodov padne do aspoň jedného z pásov na okraji c , ako je vyobrazené na obrázku 2.2a.

Zlá hypotéza je taká, v ktorej je celková váha zlých bodov aspoň ε . Keďže každý zlý bod padne do aspoň jedného z pásov, aspoň jeden z týchto pásov musí mať váhu aspoň $\varepsilon/4$. To ale znamená, že ak žiaden z pásov nebude mať takúto veľkú váhu, naša hypotéza bude zaručene dobrá. Inými slovami, pravdepodobnosť zlej hypotézy je zhora ohraničená pravdepodobnosťou, že aspoň jeden z pásov má váhu aspoň $\varepsilon/4$.

Pozrime sa na ľubovoľný jeden konkrétny pás, BUNV ten ľavý. Aká je pravdepodobnosť, že váha tohto pásu je aspoň $\varepsilon/4$? Aby sa toto stalo, museli by všetky body z tréningových príkladov padnúť ďaleko od ľavého okraja. Ako presne ďaleko?

Predstavme si súvislý úsek bodov od ľavého okraja, ktorého váha je presne $\varepsilon/4$. Bod je dostatočne ďaleko, ak nepadne do tohto úseku; ilustrujeme na obrázku 2.2b. Pravdepodobnosť, že náhodný príklad padne takto ďaleko, je $1 - \varepsilon/4$; pravdepodobnosť, že každý z t náhodne vybraných tréningových príkladov padol takto ďaleko, je $(1 - \varepsilon/4)^t$.

S touto predstavou je ale drobný technický problém: nie pre každé rozdelenie P platí, že vieme nájsť úsek s váhou presne $\varepsilon/4$. Napríklad si predstavme rozdelenie, v ktorom s pravdepodobnosťou 50% dostaneme bod v strede obdĺžnika, a inak si vyberieme rovnomerne náhodne bod v pravej polovici. Najmenší ľavý súvislý úsek s kladnou váhou bude mať hneď váhu 50%.

Musíme teda postupovať inak. Uvedomme si, že bod x je dostatočne ďaleko vtedy, keď váha bodov naľavo od x je aspoň $\varepsilon/4$. To ale znamená, že vieme vyznačiť súvislý úsek bodov pri ľavom okraji, ktorý má váhu aspoň $\varepsilon/4$ a ktorý je celý naľavo od všetkých takýchto dostatočne vzdialených bodov x . Z toho vyplýva, že váha všetkých dostatočne vzdialených bodov môže byť nanajvýš $1 - \varepsilon/4$. Pravdepodobnosť, že náhodný príklad padne takto ďaleko, je preto nanajvýš $1 - \varepsilon/4$; pravdepodobnosť, že každý z t náhodne vybraných tréningových príkladov padol takto ďaleko, je nanajvýš $(1 - \varepsilon/4)^t$. To ďalej zhora odhadneme ako $e^{-t \cdot \varepsilon/4}$.

Tento odhad spravíme pre každý z okrajových pásov. Dostaneme tak odhad na pravdepodobnosť zlej hypotézy:

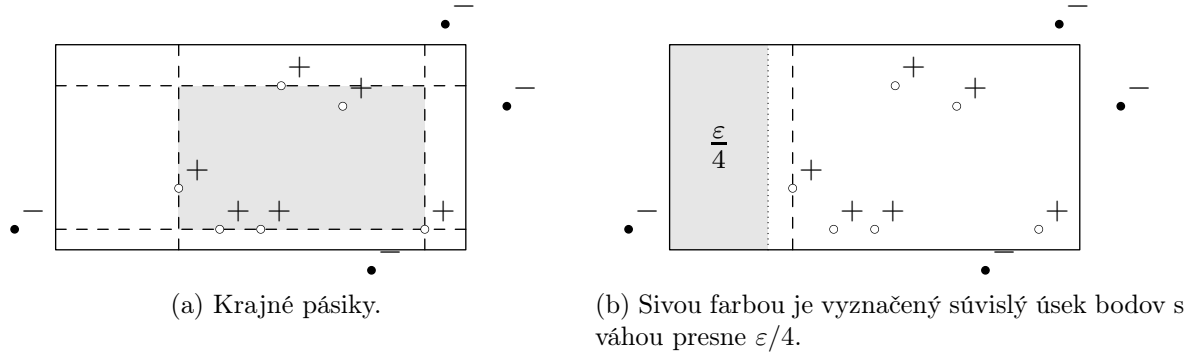
$$\Pr_T(\text{Err}(\hat{h}_T) \geq \varepsilon) \leq 4 \cdot e^{-\frac{t\varepsilon}{4}}$$

Ak chceme dostať zlú hypotézu iba s malou pravdepodobnosťou (nanajvýš δ), stačí zvoliť také t , že pravá strana je nanajvýš δ . Odtiaľ dostaneme odhad na postačujúci počet tréningových príkladov:

$$4 \cdot e^{-\frac{t\varepsilon}{4}} \leq \delta \quad (2.9)$$

$$-\frac{t\varepsilon}{4} \leq \ln\left(\frac{\delta}{4}\right) \quad (2.10)$$

$$t \geq \frac{4}{\varepsilon} \cdot \ln\left(\frac{4}{\delta}\right) \quad (2.11)$$



Obr. 2.2: Obdĺžniková hra.

□

2.3.2 Vapnik-Chervonenkisova dimenzia

Vapnik-Chervonenkisova dimenzia (alebo skráteno 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 2. Majme konečnú množinu vstupov $S \subseteq X$. Hovoríme, že množina hypotéz H rozbieja množinu S , ak platí: nech označíme vstupy v S ako pozitívne alebo negatívne akokoľvek, v množine H existuje hypotéza konzistentná s týmto označením.

Inými slovami, množina S je rozbitá množinou hypotéz H , ak sa hypotézy z tejto množiny vedú správať na vstupoch S akokoľvek; všetkých $2^{|S|}$ možných označení je konzistentných s niektorou hypotézou v H . Takéto označenia budeme nazývať *realizovateľné*.

Definícia 3. Vapnik-Chervonenkisova dimenzia množiny hypotéz H , označovaná $VC(H)$, je veľkosť najväčšej rozbiteľnej množiny vstupov $S \subseteq X$. Ak sa dá rozbiť ľubovoľne veľká množina S , definujeme $VC(H) = \infty$.

Pokiaľ bude z kontextu jasné, akú množinu hypotéz uvažujeme, budeme jej VC dimenziu označovať d .

Od miery zložitosti očakávame isté intuitívne vlastnosti, napríklad že zväčšením množiny hypotéz sa môže zložitosť iba zvýšiť. Toto skutočne platí pre VC dimenziu.

Veta 2.10. Nech pre množiny hypotéz H a H' platí $H \supseteq H'$. Potom $VC(H) \succeq VC(H')$, kde \succeq je relácia \geq prirodzene rozšírená na $\mathbb{N} \cup \{\infty\}$.

Dôkaz. Množina S sa dá rozbiť pomocou H' vtedy, keď každé označenie je realizovateľné niektorou hypotézou v H' . Pretože každá hypotéza v H' je aj v H , je každé označenie realizovateľné aj pomocou H . □

Súvis s PAC naučiteľnosťou

Konkrétny súvis medzi VC dimenziou a PAC naučiteľnosťou (a dôkazy týchto vzťahov) zatiaľ vynecháme, uvedieme tu iba znenia viet. Pre záujemcov odporúčame sekciu appendix.

Veta 2.11 (Blumer et. al., 1989 [1]). Ak je VC dimenzia konečná, tak je množina hypotéz H PAC* naučiteľná: pre ľubovoľné rozdelenie P , ľubovoľne malé hranice ε a δ , ak počet tréningových príkladov t spĺňa

$$t \geq c_0 \left(\frac{1}{\varepsilon} \cdot \ln \left(\frac{1}{\delta} \right) + \frac{d}{\varepsilon} \cdot \ln \left(\frac{1}{\varepsilon} \right) \right)$$



Obr. 2.3: Žiaden interval nevytvorí takéto označenie bodov.

pre nejakú konštantu $c_0 > 0$ závislú iba od H , tak dostaneme zlú hypotézu (s OTeCh aspoň ε) iba s malou pravdepodobnosťou (nanajvýš δ).

Veta 2.12 (Haussler et. al., 1994 [3]). OTeChAlg vieme zhora asymptoticky odhadnúť nasledovne:

$$\text{Err}(L) = O\left(\frac{d}{t} \cdot \ln\left(\frac{t}{d}\right)\right).$$

Dosiaľ sme všade predpokladali, že naša množina hypotéz H je úplná. Ukazuje sa ale, že aj v prípade neúplnej množiny hypotéz vieme dostať odhad na OTeCh hypotézy. Už ale nemôžeme požadovať ľubovoľne malú OTeCh hypotézu; namiesto toho ju budeme porovnávať s h^* (najlepšou možnou hypotézou v H).

Veta 2.13 (Vapnik & Chervonenkis, 1971 [2]). Pre ľubovoľné rozdelenie P , ľubovoľne malé hranice ε a δ , ak počet tréningových príkladov spĺňa

$$t \geq c_0 \left(\frac{1}{\varepsilon} \cdot \ln\left(\frac{1}{\delta}\right) + \frac{d}{\varepsilon} \cdot \ln\left(\frac{1}{\varepsilon}\right) \right),$$

pre nejakú konštantu c_0 závislú iba od množiny hypotéz H , tak sa iba s malou pravdepodobnosťou stane, že hypotéza \hat{h}_T je oveľa horšia od hypotézy h^* :

$$\Pr_T(\text{Err}(\hat{h}_T) \geq \text{Err}(h^*) + \varepsilon) \leq \delta.$$

Použitie

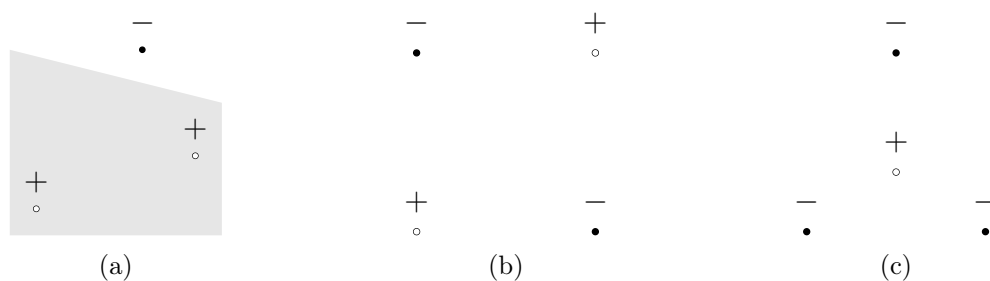
Na to, aby sme ukázali dolný odhad $\text{VC}(H) \geq d$, nám stačí nájsť jednu množinu S veľkosti d , ktorá sa dá rozbiť. Aby sme ale ukázali horný odhad $\text{VC}(H) \leq d$, musíme ukázať, že žiadna množina veľkosti $d+1$ sa nedá rozbiť; teda že pre každú $(d+1)$ -ticu vstupov existuje nerealizovateľné označenie. Z tohto dôvodu je obvykle ťažšie dokázať horný odhad na VC dimenziu, ako dolný odhad.

Dalej uvedieme príklady niektorých nekonečných množín hypotéz a zdôvodníme ich VC dimenzie.

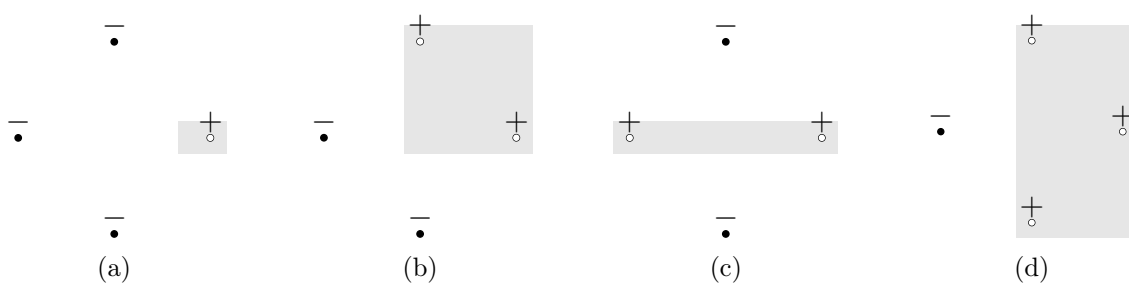
Intervaly na reálnych číslach. Nie je ťažké vidieť, že ľubovoľná dvojica bodov sa dá rozbiť. Na druhej strane, ak máme tri body, tak ich vieme označiť ako na obrázku 2.3; toto označenie nie je realizovateľné žiadnym intervalom. Takže VC dimenzia je 2.

Polroviny v rovine. Každá trojica bodov tvoriacich nede degenerovaný trojuholník sa dá rozbiť, ako je ilustrované na obrázku 2.4. Ukážeme, že žiadna štvorica bodov sa rozbiť nedá. Ak tvoria konvexný štvoruholník, tak jednu protiľahlú dvojicu označíme pozitívne a druhú negatívne. Ak tvoria nekonvexný štvoruholník, jeden z jeho vrcholov je vo vnútri trojuholníka tvoreného ostatnými tromi: tento bod označíme pozitívne a ostatné body negatívne. Takže VC dimenzia je 3.

Osovorovnoběžné obdĺžniky. Štvorica bodov, ktorá sa dá rozbiť, je ilustrovaná na obrázku 2.5 (nie každá štvorica bodov sa ale dá rozbiť). Ukážeme, že žiadna päťica bodov sa nedá rozbiť. Ak je aspoň jeden z bodov vo vnútri bounding boxu (nie na okraji), tak ho označíme negatívne a všetky body na okraji pozitívne. Ak sú všetky body na obode bounding boxu, na jednej strane toho boxu musia ležať aspoň dva body. Jeden z nich označíme pozitívne a druhý negatívne. Takže VC dimenzia je 4.



Obr. 2.4: Polroviny v rovine: v situácii (a) vieme nájsť deliacu priamku, v (b) a (c) nie.



Obr. 2.5: Osovorovnoběžné obdĺžniky: štvorica bodov, ktorá sa dá rozbiť. Jednotlivé obrázky zobrazujú všetky netriviálne označenia bodov a príslušné konzistentné klasifikátory.

Appendix

2.A Dôkaz o VC dimenzii a PAC naučiteľnosti

TODO všetko

Literatúra

- [1] Anselm Blumer, Andrzej Ehrenfeucht, David Haussler, and Manfred K Warmuth. Learnability and the Vapnik-Chervonenkis dimension. *Journal of the ACM (JACM)*, 36(4):929–965, 1989.
- [2] AIA Chervonenkis and VN Vapnik. Theory of uniform convergence of frequencies of events to their probabilities and problems of search for an optimal solution from empirical data(average risk minimization based on empirical data, showing relationship of problem to uniform convergence of averages toward expectation value). *Automation and Remote Control*, 32:207–217, 1971.
- [3] David Haussler, Nick Littlestone, and Manfred K Warmuth. Predicting $\{0, 1\}$ -functions on randomly drawn points. *Information and Computation*, 115(2):248–292, 1994.