

1 Teória strojového učenia I

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 trénovacia chyba z pohľadu algoritmu bude

$$\frac{E}{T} \left[\text{err}_T(\hat{h}) \right].$$

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 trénovacia chyba priemernej hypotézy) od veľkosti trénovacej 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

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

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:

$$\text{err}(h) = E_{x,y} [(h(x) - y)^2] \tag{1}$$

$$= E_x \left[E_{y|x} [(h(x) - y)^2] \right] \tag{2}$$

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, verzia 1.

V tomto odseku si ukážeme zaujímavý výsledok, ktorý nám za určitých predpokladov umožňuje vyjadriť chyby pomocou iných, jasnejších veličín: tzv. *výchylky* a *rozptylu*.

Odvodenie. 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} .

$$\text{chyba algoritmu} = \mathbb{E}_T [\text{err}(\hat{h})] \quad (3)$$

$$= \mathbb{E}_T \left[\mathbb{E}_{x,y} [(\hat{h}(x) - y)^2] \right] \quad (4)$$

$$= \mathbb{E}_T \left[\mathbb{E}_{x,y} \left[\left((\hat{h}(x) - h^*(x)) + (h^*(x) - y) \right)^2 \right] \right] \quad (5)$$

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

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

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

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

Prvý zo sčítancov budeme volať *rozptyl*. Vyjadruje, ako ďaleko je naša funkcia od najlepšej možnej, v rámci množiny hypotéz H . Druhý zo sčítancov budeme volať *výchylka*. Vyjadruje chybu, ktorá je spôsobená výberom množiny hypotéz.

Výchylku vieme upraviť ďalej. Pretože hypotéza h^* ani y nezávisia od tréningovej množiny T , merať chybu na testovacích dátach x, y je to isté, ako merať ju na tréningových dátach x_i, y_i , berúc ich náhodný výber. Teda

$$\text{výchylka} = \mathbb{E}_T \left[\mathbb{E}_{x_i, y_i} [(h^*(x_i) - y_i)^2] \right] \quad (6)$$

$$= \mathbb{E}_T \left[\mathbb{E}_{x_i, y_i} \left[\left((h^*(x_i) - \hat{h}(x_i)) + (\hat{h}(x_i) - y_i) \right)^2 \right] \right] \quad (7)$$

Opäť, použitím toho istého technického kroku dostaneme:

$$\text{výchylka} = \underbrace{\mathbb{E}_T \left[\mathbb{E}_{x_i, y_i} [(h^*(x_i) - \hat{h}(x_i))^2] \right]}_{\text{tréningový rozptyl}} + \underbrace{\mathbb{E}_T \left[\mathbb{E}_{x_i, y_i} [(\hat{h}(x_i) - y_i)^2] \right]}_{\text{priemerná tréningová chyba}} \quad (8)$$

Tréningový rozptyl vyjadruje, ako ďaleko je naša hypotéza \hat{h} od najlepšej možnej h^* z H . Na rozdiel od rozptylu ale túto vzdialenosť meriame na tréningových dátach, nie na testovacích. To spraví rozdiel, nakoľko \hat{h} je závislé od tréningových dát. Priemerná tréningová chyba je priemerná chyba, ktorej sa dopustí výstup z algoritmu \hat{h} na tých istých dátach, pomocou ktorých sme \hat{h} zostrojili.

Záver. Podarilo sa nám teda rozložiť chybu algoritmu na dve, prípadne tri časti. Načo je to ale dobré? Ukážeme si, ako pomocou nich vieme získať intuíciu o tom, ako sa správa chyba algoritmu v závislosti od veľkosti tréningovej množiny a veľkosti (tj. zložitosti) množiny hypotéz.

TODO obrázok kriviek učenia, vysvetlenie

TODO podučenie, preučenie

	vektory	funkcie
súradnice	y_1, y_2, \dots, y_n	$f(x)$
dĺžka	$\ y\ = \sqrt{\sum_{i=1}^n y_i^2}$	$\ f\ = \sqrt{\int f^2(x) d\rho x} = \sqrt{\mathbb{E}_x[f^2(x)]}$
skalárny súčin	$\langle y, w \rangle = \sum_{i=1}^n y_i \cdot w_i$	$\langle f, g \rangle = \int f(x) \cdot g(x) d\rho x = \mathbb{E}_x[f(x) \cdot g(x)]$
vzdialenosť	$d(y, w) = \ y - w\ $	$d(f, g) = \ f - g\ $

Tabuľka 1: Vektory/funkcie.

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

Po prvé, pre jednoduchosť budeme predpokladať, že vstup je vektor reálnych čísel (tj. $X = \mathbb{R}^n$), snažíme sa predpovedať jedno reálne číslo (tj. $Y = \mathbb{R}$), a že pravdepodobnostné rozdelenie P je spojité. Jeho hustotu pravdepodobnosti označíme ρ .

Po druhé, predpokladáme, že trénovací algoritmus vráti takú funkciu $\hat{h} \in H$, ktorá minimalizuje trénovaciu chybu. Inak zapísané,

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

Po tretie, kladieme obmedzenia na množinu hypotéz H : musí byť uzavretá na lineárne kombinácie a na limity. Dá sa na ňu teda pozeráť ako na (nie nutne konečnorozmerný) vektorový priestor, ktorý je navyše uzavretý: ak postupnosť vektorov (v našom prípade funkcií) konverguje, tak jej limita je tiež v tom priestore.

V tabuľke 1 uvádzame na jednej strane vlastnosti konečnorozmerných vektorov, na druhej strane vlastnosti funkcií ako vektorov.

TODO dokončiť dôkaz

1.2.3 Bias-variance tradeoff, verzia 2.

Ak ste si dali vyhľadávať tento pojem, určite ste narazili na kopu materiálov, v ktorých sa výsledok vôbec nepodobal na to, čo sme ukazovali vyššie. Ide totiž o veľmi príbuzný, avšak predsa odlišný výsledok, ukážeme a odvodíme si ho.

Veta 1. *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 trénovacieho algoritmu je \hat{f} . Za chybovú funkciu zvolíme kvadratickú chybu. Chybu algoritmu vieme teda vypočítať nasledovne:

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

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

Poznámka 3. V tomto znení bias-variance tradeoff-u názvy *rozptyl* a *výchylka* zodpovedajú príslušným štatistickým/pravdepodobnostným pojmom.

Poznámka 4. Na rozdiel od predchádzajúcej verzie bias-variance tradeoff-u, tu nebudeme potrebovať žiadne dodatočné predpoklady od algoritmu ani od jeho množiny hypotéz. (Nemusí teda vracať hypotézu, ktorá je spomedzi hypotéz v H najlepšia na daných tréningových dátach. Takisto od množiny hypotéz nepožadujeme žiadne vlastnosti.)

Dôkaz. Upravujeme pôvodný výraz.

$$\text{chyba algoritmu} = \mathbb{E} \left[(\hat{f}(x) - y)^2 \right] \quad (9)$$

$$= \mathbb{E} \left[(\hat{f}(x) - f(x) - \varepsilon)^2 \right] \quad (10)$$

$$= \mathbb{E} \left[(\hat{f}(x) - f(x))^2 \right] + \mathbb{E} [\varepsilon^2] - 2 \cdot \mathbb{E} [\varepsilon \cdot (\hat{f}(x) - f(x))] \quad (11)$$

$$= \mathbb{E} \left[(\hat{f}(x) - f(x))^2 \right] + \mathbb{E} [\varepsilon^2] \quad (12)$$

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.

$$\text{prvý sčítanec} = \mathbb{E} \left[(\hat{f}(x) - f(x))^2 \right] \quad (13)$$

$$= \mathbb{E}[\hat{f}(x)^2] + \mathbb{E}[f(x)^2] - 2 \cdot \mathbb{E}[\hat{f}(x) \cdot f(x)] \quad (14)$$

$$= (\text{Var}(\hat{f}(x)) + \mathbb{E}[\hat{f}(x)]^2) + (\text{Var}(f(x)) + \mathbb{E}[f(x)]^2) - 2 \cdot \mathbb{E}[\hat{f}(x) \cdot f(x)] \quad (15)$$

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

$$\begin{aligned} \text{prvý sčítanec} &= \text{Var}(\hat{f}(x)) + \text{Var}(f(x)) + (\mathbb{E}[\hat{f}(x)] - \mathbb{E}[f(x)])^2 \\ &\quad + 2 \cdot \mathbb{E}[\hat{f}(x)] \cdot \mathbb{E}[f(x)] - 2 \cdot \mathbb{E}[\hat{f}(x) \cdot f(x)] \end{aligned} \quad (16)$$

$$= \text{Var}(\hat{f}(x)) + \text{Var}(f(x)) + (\mathbb{E}[\hat{f}(x)] - \mathbb{E}[f(x)])^2 - 2 \cdot \text{Cov}(\hat{f}(x), f(x)) \quad (17)$$

$$= \text{Var}(\hat{f}(x) - f(x)) + (\mathbb{E}[\hat{f}(x)] - \mathbb{E}[f(x)])^2 \quad (18)$$

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?

TODO regularizácia

TODO holdout testing

TODO k -fold cross validation

TODO best practices