

# Obsah

<b>1</b>	<b>Teória strojového učenia I</b>	<b>2</b>
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

# Kapitola 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”? Spravíme 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  $\text{err}$ , 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.

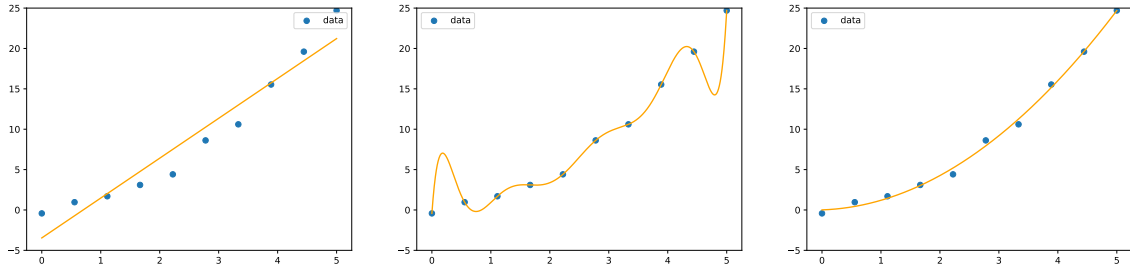
$$\text{chyba algoritmu} = E_T \left[ E_{x,y} [(\hat{h}(x) - h^*(x))^2] \right] + E_T \left[ E_{x,y} [(h^*(x) - y)^2] \right]$$

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 neprilepšíme. 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 prílepiť). 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 trérovacie dáta  $T$ , pri voľbe množiny hypotéz  $H$  sa snažíme nájsť kompromis medzi malým rozptylom a malou výchyľkou. Zložité  $H$  bude mať malú výchyľku ale veľký rozptyl, čo vedie k tzv. *preučeniu*. Jednoduché  $H$  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 trérovacie dáta, avšak mimo nich bude dávať výsledky úplne mimo.

Výchyľku vieme upraviť ďalej. Hypotéza  $h^*$  ani  $y$  nezávisia od trérovacej množiny  $T$ . Z ich pohľadu sú teda testovacie dáta  $x, y$  a trérovacie dáta  $x_i, y_i$  nerozoznatelné. Takže na meranie chyby  $h^*$  môžeme použiť trérovacie 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:

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

Prvý zo sčítancov budeme volať *trérovací rozptyl*. Uvedomme si, že pre ľubovoľné trérovacie dáta  $T$  platí

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

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

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

Platí

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

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

Na základe dosiaľ uvedeného vieme graficky znázorniť, ako sa zhruba správajú rozptyl, výchyľka, trérovací rozptyl a priemerná trérovacia chyba, v závislosti od veľkosti trérovacej 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  $\rho$  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  $\Delta$ . Pozrime sa na ich vzdialenosti od funkcie  $f$ , vyjadrené ako funkcia od  $\Delta$ :

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

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

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

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  $\text{dist}$  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 lemy 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  $\varepsilon$  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/podučeníím?

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  $\lambda$  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 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 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 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  $\theta$ , 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  $\theta$  riešiacie rovnicu

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