

UNIVERZITA PAVLA JOZEFA ŠAFÁRIKA V KOŠICIACH
PRÍRODOVEDECKÁ FAKULTA

Gabriela Andrejková
Ľubomír Antoni

Strojové učenie

Košice 2020, Univerzita Pavla Jozefa Šafárika v Košiciach



EURÓPSKA ÚNIA

Európsky sociálny fond
Európsky fond regionálneho rozvoja



OPERAČNÝ PROGRAM
ĽUDSKÉ ZDROJE



MINISTERSTVO
ŠKOLSTVA, VEDY,
VÝSKUMU A ŠPORTU
SLOVENSKEJ REPUBLIKY



*Tento projekt sa realizuje vďaka podpore z Európskeho sociálneho fondu
a Európskeho fondu regionálneho rozvoja v rámci Operačného programu Ľudské zdroje*

www.minedu.sk www.employment.gov.sk/sk/esf/ www.itakademia.sk

UNIVERZITA PAVLA JOZEFA ŠAFÁRIKA V KOŠICIACH
PRÍRODOVEDECKÁ FAKULTA



STROJOVÉ UČENIE

Spracované v rámci národného projektu
IT Akadémia – vzdelávanie pre 21. storočie

Gabriela Andrejková, Ľubomír Antoni

Košice 2020

Tento projekt sa realizuje vďaka podpore z Európskeho sociálneho fondu v rámci Operačného programu Ľudské zdroje.

*Spracované s finančnou podporou národného projektu
IT Akadémia – vzdelávanie pre 21. storočie*

Strojové učenie

Vysokoškolský učebný text

Autori:

doc. RNDr. Gabriela Andrejková, CSc.

Ústav informatiky, Prírodovedecká fakulta UPJŠ v Košiciach

RNDr. Ľubomír Antoni, PhD.

Ústav informatiky, Prírodovedecká fakulta UPJŠ v Košiciach

Recenzenti:

doc. PaedDr. Jozef Kapusta, PhD.

Katedra informatiky, Fakulta prírodných vied UKF v Nitre

RNDr. Dávid Hudák, PhD.

VSL Software, a.s.

Za odbornú a jazykovú stránku publikácie zodpovedajú autori. Rukopis neprešiel redakčnou ani jazykovou úpravou.

Tento text je publikovaný pod licenciou Creative Commons 4.0 - Attribution CC BY
Creative Commons Attribution 4.0 („Uveďte pôvod“).



Umiestnenie: www.unibook.upjs.sk

Dostupné od: 20. 10. 2020

ISBN 978-80-8152-912-2 (e-publikácia)

1 Úvod	6
2 Koncepty, hypotézy, učiace sa algoritmy	8
2.1 Koncepty	9
2.2 Trénovanie a učenie sa	10
2.3 Učenie sa pomocou konštrukcie	13
2.4 Učenie sa očíslovaním	14
2.5 Úlohy	15
3 Booleovské formuly (hypotézy) a ich reprezentácie	16
3.1 Učiaci sa algoritmus pre monočleny	16
3.2 Učenie disjunkcií malých monočlenov	18
3.3 Reprezentácia hypotézového priestoru	20
3.4 Čas behu učiacich algoritmov	20
3.5 Konzistencia tréningovej vzorky	21
3.6 Úlohy	23
4 Pravdepodobnostné učenie	25
4.1 Algoritmus pre učenie líčov	25
4.2 Pravdepodobnostné aproximačne správne učenie, PAC učenie	26
4.3 Učenie líčov je PAC	30
4.4 Exaktné učenie	31
4.5 Úlohy	33
5 Lineárna regresia, logistická regresia	34
5.1 Lineárne modelovanie	34
5.1.1 Definícia modelu	35
5.1.2 Odvodenie vzťahov pre výpočet k a q	37

5.1.3	Predikcia rekordu	38
5.2	Zovšeobecnené lineárne modely	39
5.3	Úlohy	39
6	Lineárne modelovanie, maximálna dôveryhodnosť	40
6.1	Chyby ako šum	40
6.2	Dôveryhodnosť (hodnovernosť)	41
6.3	Hodnovernosť datasetu	43
6.4	Maximálna dôveryhodnosť	44
6.5	Charakteristika max. dôveryhodnosti	45
6.6	Úlohy	46
7	Klasifikácia	47
7.1	Úvod	47
7.2	Klasifikácia do dvoch tried	48
7.3	Pásový perceptrón	51
7.4	Klasifikácia s podpornými vektormi - SVM (Support Vector Machines)	52
7.5	Nie celkom presná klasifikácia	53
7.6	Nelineárne SVM	54
7.7	Úlohy	56
8	Rozhodovacie stromy	57
8.1	Miery výberu atribútov	58
8.2	Algoritmus pre generovanie rozhodovacích stromov	60
8.3	Preučenie rozhodovacieho stromu a jeho prerezávanie	61
9	Zhlukovanie	63
9.1	Typy údajov pri zhlukovaní	64
9.2	Segmentačné metódy zhlukovania	66
9.3	Hierarchické metódy zhlukovania	69
9.4	Metódy zhlukovania založené na hustote	71
9.5	Požiadavky na zhlukovanie	73
10	Bayesovský prístup ku strojovému učeniu	75

10.1 Úvod	75
10.2 Hra Hod mincou	75
10.3 Bayesovský prístup k hre	77
10.3.1 Bayesova veta - všeobecný tvar	77
10.3.2 Bayesova veta v hre s mincou	78
10.3.3 Pravdepodobnosť vs. dôveryhodnosť (likelihood)	78
10.3.4 Dvojice dôveryhodnosť-apriórna pravdepodobnosť	81
10.3.5 Exaktné posteriórne rozdelenie	82
10.4 Grafické modely	83
10.5 Bayesov prístup k problému predikcie rekordov na OH	83
10.5.1 Bayesov prístup k predikcii rekordov	84
10.5.2 Prípád lineárneho modelu $t_n = w_0 + w_1 x_n$	85
10.5.3 Ako urobiť predikciu	87
10.5.4 Úlohy	89

Kapitola 1

Úvod

Aj keď je strojové učenie často spájané so všeobecnejšou oblasťou umelá inteligencia, jeho realistickejší popis je, že sa zaoberá otázkami ako počítačové programy automaticky rozpoznávajú zložité vzory a robia rozhodnutia na základe dát.

Algoritmus strojového učenia sa¹ je teda algoritmus, ktorý sa učí z dát. Ale čo máme na mysli tým, že "sa učí"? Mitchell (1997) uvádza definíciu: "Počítačový program sa učí zo skúsenosti **E** s ohľadom na určitú triedu úloh **T** a meranie výkonnosti **P**, ak sa jeho výkon pri úlohách v triede **T**, meraný pomocou **P**, zlepšuje použitím skúseností z **E**." V rámci tejto definície si môžeme si predstaviť veľmi širokú škálu skúseností **E**, úloh **T** a meraní výkonnosti **P**.

V týchto vysokoškolských učebných textoch prezentujeme formálnu definíciu toho, čo môže byť použité pre každú z týchto entít. Ale tiež budeme pracovať intuitívne, pretože zavádzanie formalizmu v niektorých prípadoch by značne zväčšilo rozsah tejto publikácie. Budeme sa zaoberať viacerými triedami úloh, skúsenosti budú vyjadrené v dátach a na "učenie sa" budeme používať rôzne algoritmy.

Strojové učenie nám umožňuje riešiť úlohy, ktoré sú veľmi náročné na riešenie použitím programov napísaných a navrhnutých ľuďmi. Z vedeckého a filozofického hľadiska je strojové učenie sa zaujímavé, pretože rozvíjanie nášho chápania strojového učenia znamená rozvíjať naše chápanie princípov, ktoré tvoria základ umelej inteligencie.

Úlohy

V tomto relatívne formálnom vymedzení slova „úloha“, samotný proces „učenia sa“ nie je úlohou. Učenie sa je náš prostriedok na dosiahnutie schopnosti splniť úlohu. Napríklad, ak chceme, aby bol robot schopný chodiť, potom chôdza je úloha. Mohli by sme robota naprogramovať, aby sme ho naučili chodiť, alebo by sme sa mohli pokúsiť priamo napísať ručne program, ktorý špecifikuje, ako má chodiť.

¹V slovenčine sa používajú slovesá „učiť“ a „učiť sa“, preto je vhodné uvedomiť si, o ktorú aktivitu ide. Algoritmy, ktoré adaptujú svoje parametre na základe nejakých poznatkov „sa učia“.

Typickými príkladmi úloh sú klasifikačné, regresné a aproximačné úlohy.

Aby sme mohli vyhodnotiť schopnosti algoritmu strojového učenia, musíme navrhnuť kvantitatívne meranie jeho výkonnosti. Zvyčajne toto meranie P je špecifické pre úlohu T vykonávanú systémom.

**Meranie
výkonnosti**

V prípade úloh, ako je klasifikácia, často meriame presnosť modelu. Presnosť predstavuje pomer príkladov, pre ktoré model produkuje správny výstup. Ekvivalentné informácie môžeme získať meraním miery chybovosti, podielu príkladov, pri ktorých model produkuje nesprávny výstup.

Voľba typu merania výkonnosti sa môže zdať jednoduchá a objektívna, ale často je veľmi dôležité vybrať si tak, aby meranie zodpovedalo požadovanému správaniu sa systému. V niektorých prípadoch je preto dôležité rozhodnúť, čo sa má merať. Napríklad, pri vykonávaní regresnej úlohy by sme mali penalizovať systém viac, ak často robí stredne veľké chyby alebo ak zriedka robí veľmi veľké chyby? Tieto možnosti výberu závisia od aplikácie.

Učiacie sa algoritmy môžu byť charakterizované podľa typu učenia sa: bez dozoru alebo pod dohľadom. Ale aké skúsenosti môžu mať počas procesu učenia sa? Väčšina učiacich sa algoritmov v našom prípade bude používať celý súbor údajov. Súbor údajov je súbor mnohých príkladov. Niekedy ich budeme tiež nazývať príklady dátových bodov. Skúsenosť sa teda bude nachádzať v príkladoch. Napríklad, pomocou príkladu sa dá vyjadriť, že to, čo má štyri nohy, operadlo a sedaciu časť, je stolička.

Skúsenosti

V týchto vysokoškolských učebných textoch venujeme prvé tri kapitoly teoretickým poznatkom z oblasti strojového učenia a výpočtového učenia, v ktorých ukážeme tvorbu konceptov pre objekty, s ktorými budeme pracovať, prípravu konceptov – hypotéz, ktoré budú predstavovať naučený algoritmus a proces učenia sa z príkladov. Táto časť je motivovaná knihou Anthony [1].

Ďašie tri kapitoly sú venované lineárnemu modelovaniu a klasifikácii predovšetkým do dvoch tried. Lineárnemu modelovaniu je venovaná pozornosť aj z hľadiska dôveryhodnosti dosahovaných výsledkov. Klasifikácia je analyzovaná pomocou perceptrónov a pomocou metódy podporných vektorov (Support Vector Machines). Motivácia pre tieto kapitoly a desiatu kapitolu bola nájdená v knihe Rogers & Girolami [6].

Rozhodovacie stromy a metódy zhlukovania sú predmetom ďalších dvoch kapitol, ktoré sú prezentované pomocou pojmov definovaných v knihe Han a kol. [3]. Posledná kapitola sa zaoberá Bayesovským prístupom k strojovému učeniu, v ktorej predovšetkým na príklade vysvetľujeme pravdepodobnosti vystupujúce v Bayesovom vzorci.

Vysokoškolské učebné texty sú určené pre študentov bakalárskych, magisterských, inžinierskych a doktorandských študijných programov informatických odborov.

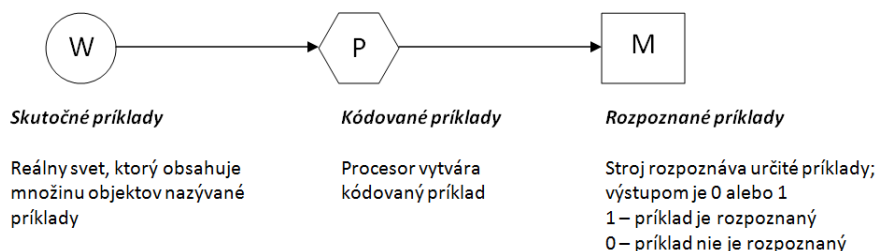
Kapitola 2

Koncepty, hypotézy, učiace sa algoritmy

Existuje viacero typov aktivít označovaných ako "učenie sa". V úvode budeme študovať matematický model takého procesu. Tento model sa zdá byť použiteľný, pretože zachytáva základ určitých aktivít, ktoré boli popísané predtým pomocou nepresných výrazov, a zároveň umožňuje vytvoriť netriviálne matematické tvrdenia, ktoré môžu byť dokázané. Terminológia a text tejto kapitoly a nasledujúcich dvoch kapitol je silne motivovaný učebnicou Anthony (1997).

V reálnom svete máme množstvo objektov, ktoré používame v rôznych situáciách, napríklad stoly, stoličky, a pod. Množina rôznych stoličiek tvorí množinu príkladov stoličiek, na základe ktorých sme si vytvorili akýsi model (vzor) stoličky. Napríklad, môžeme brať do úvahy nasledujúce atribúty pre stoličku: má 4 nohy, má operadlo, má sedáciu časť, nemá chvost, na farbe nezáleží. Takýto abstraktný vzor budeme nazývať koncept. Na základe definovaného konceptu budeme vedieť podobné objekty zaradiť alebo nezariť k stoličkám.

Na obrázku [2.1](#) je uvedený postup rozpoznávania objektu, ktorý je použiteľný pre nejaký robot alebo stroj. Stroj dostane zakódované príklady a musí mať algoritmus, na základe ktorého vytvorí svoju odpoveď.



Obr. 2.1: Spracovanie príkladov reálneho sveta pomocou natrénovaného stroja

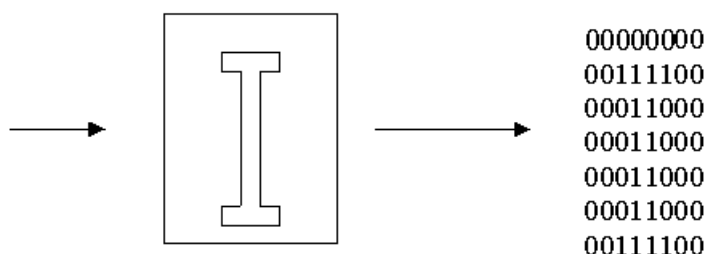
Kvalita odpovedí stroja M závisí od kvality algoritmu a bolo by vhodné mať algoritmus, ktorý sa dokáže modifikovať na základe spracovaných správne (aj nesprávne) príkladov. Proces adaptácie algoritmu z príkladov budeme nazývať *učenie sa*.

2.1 Koncepty

Sformalizujeme pojem koncept, ktorý budeme popisovať pomocou množiny príkladov. \square

Nech množina Σ je abeceda na popis príkladov. Napríklad, $\Sigma = \{0, 1\}$, $\Sigma = \mathcal{R}$, množina reálnych čísel.

Množina $\Sigma^n, n \geq 1$, predstavuje množinu všetkých reťazcov dĺžky n v abecede Σ . Do množiny Σ^* patrí prázdny reťazec a všetky reťazce konečnej dĺžky.



Obr. 2.2: Popis písmena I pomocou bitmapy v abecede $\Sigma = \{0, 1\}$.

**Koncept,
priestor
príkladov**

Definícia 2.1 *Nech $X \subseteq \Sigma^*$. Množina X sa nazýva priestor príkladov. Koncept v abecede Σ je funkcia $c, c : X \rightarrow \{0, 1\}$. Prvok $x, x \in X$ sa nazýva príklad. Ak pre $x \in X$ platí $c(x) = 1$, tak x je **pozitívny príklad**, ak platí $c(x) = 0$, tak x je **negatívny príklad**.*

Zjednotenie množiny kladných a záporných príkladov je definičný obor funkcie c . Teda za predpokladu, že definičný obor je známy, c určuje a je určené množinou svojich pozitívnych príkladov.

Príklad 2.1 Príklady konceptov

1. Koncept PARITA pre reťazce vytvorené v abecede $\Sigma = \{0, 1\}$

¹Koncept (z lat.) môže byť:

- prvé predbežné spracovanie, návrh, náčrt (najmä textu);
- nákres, osnova, koncepcia, poňatie
- vo filozofii: formulácia, rozumový obraz, všeobecná myšlienka, pojem
- v logickej sémantike: zmysel mena

$\Sigma = \{0, 1\}$, príkladový priestor $X = \Sigma^*$, $c : \Sigma^* \rightarrow \{0, 1\}$.

c definujeme takto: Nech $y \in \Sigma^*$, $y = y_1 \dots y_n$, potom

$$c(y) = \begin{cases} 1 & \text{ak } y \text{ je nepárny počet jedničiek} \\ 0 & \text{ak } y \text{ je párny počet jedničiek} \end{cases}$$

1011101 predstavuje kladný príklad, 1000001 záporný príklad.

2. Koncept PALINDROM

$\Sigma = \{0, 1\}$, príkladový priestor $X = \Sigma^*$, $c : \Sigma^* \rightarrow \{0, 1\}$.

c definujeme takto: Nech $y \in \Sigma^*$, $y = y_1 \dots y_n$, potom

$$c(y) = \begin{cases} 1 & \text{ak } y_i = y_{n-i+1} \quad i = 1, 2, \dots, \frac{n}{2} \\ 0 & \text{inak} \end{cases}$$

1011100 predstavuje záporný príklad, 1000001 kladný príklad.

3. Koncept n-ROZMERNÁ JEDNOTKOVÁ GUĽA

$\Sigma = R$, príkladový priestor $X = \Sigma^n$, $c : \Sigma^n \rightarrow \{0, 1\}$.

c definujeme takto: Nech $y \in X$, $y = y_1 \dots y_n$, potom

$$c(y) = \begin{cases} 1 & \text{ak } y_1^2 + y_2^2 + \dots + y_n^2 \leq 1 \\ 0 & \text{inak} \end{cases}$$

V prípade $n = 3$, $(0.2, 0.3, 0.1)$ predstavuje kladný príklad a $(2, 3, 1)$ záporný príklad.

2.2 Trénovanie a učenie sa

Sú dve množiny konceptov používaných v rámci učenia sa popísaného na obr. 2.1.

Prvou množinou je množina konceptov odvodených z reálneho sveta, ktorá je predkladaná na rozpoznanie (napríklad, je to písmeno A). Táto množina môže obsahovať koncepty ako "písmeno A", "písmeno B", ...; každé písmeno môže byť zakódované. Každý koncept má svoje množiny kladných a záporných príkladov. Keď je množina konceptov určovaná týmto spôsobom, budeme pre ňu používať výraz *konceptový priestor* a označovať **K**.

Konceptový priestor

Druhá množina konceptov obsiahnutých v rámci učenia na tom istom obrázku je množina, ktorú stroj M je schopný rozpoznať. Budeme predpokladať, že M je možné nastaviť do rôznych stavov a v danom stave bude klasifikovať niektoré vstupy ako kladné (výstup 1) a zvyšok ako záporné (výstup 0). Teda stav stroja M určuje koncept, ktorý môžeme chápať ako *hypotézu*. Množina všetkých konceptov, ktoré M určuje, bude nazývaná *hypotézový priestor* a označovaná **H**.

Hypotézový priestor

Ak sa vrátíme k množine stoličiek, tak vidíme, že na základe konkrétnych príkladov je vytvorený koncept stoličky. Takých konceptov môže byť viac a tvoria konceptový priestor (napríklad konceptový priestor pre nábytok). Hypotézový priestor obsahuje hypotézy, na základe ktorých stroj dáva na výstupe zaradenie vstupného objektu medzi stoličky alebo nie.

Príklad 2.2 *Stoličku je možné popísať pomocou atribútov: 4 nohy, chvost, sedací priestor, zafarbenie, žije. Príklady reprezentácie:*

$$\begin{array}{l}
 (4 \text{ nohy}, \text{ chvost}, \text{ sedací priestor}, \text{ zafarbenie}, \text{ žije}), b \\
 (1 \quad 0 \quad 1 \quad 0 \quad 0), 1 \\
 (1 \quad 0 \quad 1 \quad 1 \quad 0), 1
 \end{array}$$

Cieľom procesu učenia sa je vytvoriť hypotézu, ktorá v nejakom zmysle zodpovedá konceptu z konceptového priestoru vzhľadom na vyššie uvedené úvahy. V ďalšom texte uvedieme detaily, kedy a akým spôsobom je možné toto urobiť.

Cieľ procesu učenia sa

Máme teda 2 množiny konceptov:

- K - konceptový priestor – vychádzajúci z reality (koncept sa vytvára z príkladov),
- H - hypotézový priestor – množina hypotéz, ktoré popisujú koncepty zovšeobecneným spôsobom

a **problémom** je ku každému konceptu $k, k \in K$, nájsť nejaké $h, h \in H$, ktoré je jeho dobrou aproximáciou.

V reálnych situáciách sú hypotézy tvorené na základe určitých informácií, ktoré nedávajú explicitnú definíciu k . My budeme predpokladať, že táto informácia je poskytovaná postupnosťou kladných a záporných príkladov patriacich do X . Chceme nájsť hypotézu $h \in H$, ktorá bude popisovať koncept k čo najlepšie.

V praxi sú kladené obmedzenia na výpočtové zdroje, a preto sa musíme uspokojiť s hypotézou h , ktorá "pravdepodobne" reprezentuje k (aproximuje k) v nejakom vopred definovanom zmysle.

Nech $X \subseteq \Sigma^*$ je príkladový priestor. $\Sigma = \{0, 1\}$ alebo $\Sigma = R$. Vzorka dĺžky m je postupnosť m príkladov, t.j. m -tica $\bar{x} = (x_1, x_2, \dots, x_m) \in X^m$, kde x_i sú príklady. Postupnosť môže obsahovať rovnaké príklady viackrát. Niekedy budeme predpokladať, že sú navzájom rôzne (bez újmy na všeobecnosti).

Tréningová vzorka \bar{s} je postupnosť príkladov patriaca do množiny $(X \times \{0, 1\})^m$, t.j.

Tréningová vzorka príkladov

$$\bar{s} = ((x_1, b_1), (x_2, b_2), \dots, (x_m, b_m)),$$

kde x_i sú príklady a b_i vyjadruje, či príklad je kladný alebo záporný.

Budeme predpokladať, že vo vzorke nie sú žiadne **sporné príklady**, t.j. ak

$$x_i = x_j \quad \Rightarrow \quad b_i = b_j.$$

To teda znamená, že existuje funkcia f , definovaná ako $f(x_i) = b_i$ ($1 \leq i \leq m$).

Definícia 2.2 Budeme hovoriť, že \bar{s} je **tréningová vzorka vhodná pre cieľový koncept** $t \in H$, ak platí $b_i = t(x_i)$, pre $1 \leq i \leq m$.

Príklad 2.3 Príklad tréningovej vzorky pre koncept "PALINDROM"

$$((0010, 0), (1001001001, 1), (111, 1), (010101, 0), (111101, 0))$$

Cieľový koncept $t \in H$:

Nech $x = (x^1 \dots x^n)$, $x \in X$ je príklad, potom

$$t(x) = \begin{cases} 1 & \text{ak } x^i = x^{n-i+1} \quad \text{pre } 1 \leq i \leq n \\ 0 & \text{inak} \end{cases}$$

Tento cieľový sme vedeli vyjadriť presne jednoduchou funkciou. V praxi pre zložitejšie problémy to nie je také jednoduché.

Uvažujme teraz o povahe učiaceho procesu, ktorý tu chceme študovať. Majme dané a príklady v abecede Σ , K - konceptový priestor a H - hypotézový priestor.

Učiaci algoritmus pre (K, H)

Definícia 2.3 Učiaci algoritmus pre (K, H) , niekedy nazývaný (K, H) -učiaci algoritmus je procedúra, ktorá akceptuje vstupné tréningové vzorky pre koncepty v K a výstupy sú v tvare hypotéz v H .

Aby táto procedúra mohla byť považovaná za algoritmus, musí skončiť a musí byť efektívna. Ak ignorujeme problém efektívnosti, tak učiaci algoritmus pre (K, H) je teda funkcia L , ktorá k ľubovoľnej tréningovej vzorke \bar{s} pre cieľový koncept $t \in C$ priradí funkciu (hypotézu) $h \in H$. Píšeme $h = L(\bar{s})$.

Poznamenajme, že L by mala byť definovaná na celom príkladovom priestore X (mala by byť všeobecnejšia), aj keď $L(\bar{s})$ bola konštruovaná na konečnej podmnožine X (lebo zahŕňa len príklady vzorky (x_1, \dots, x_m)).

Hypotéza $h \in H$ je **konzistentná** s tréningovou vzorkou \bar{s} alebo súhlasí s \bar{s} , ak platí $h(x_i) = b_i$, pre $1 \leq i \leq m$.

Konzistentná hypotéza, konzistentný algoritmus

Vo všeobecnosti nepredpokladáme, že $L(\bar{s})$ je konzistentná s \bar{s} , ale keď táto podmienka platí pre všetky tréningové vzorky \bar{s} , hovoríme že L je **konzistentný algoritmus**. V tomto prípade je funkcia $L(\bar{s})$ rozširujúca funkcia na celom príkladovom priestore X .

Vo všeobecnosti platí, nie každé rozšírenie tréningovej vzorky bude vhodným zovšeobecnením, pretože cieľový koncept je len parciálne definovaný príkladmi vzorky. Okrem toho tréningová vzorka môže byť nereprezentatívna, alebo zavádzajúca. Napríklad, ak vhodne zakódujeme všetky zvieratá a cieľový koncept je "mačka", tak sa môže stať, že tréningová vzorka pozostáva z bezchvostových mačiek.

Budeme predpokladať, že nereprezentatívne vzorky sú nepravdepodobné a že väčšina vzoriek je dostatočne reprezentatívna, takže rozšírenia funkcií sú vyhovujúce.

Príklad 2.4 *Ak uvažujeme koncept lineárnej separácie bodov v rovine (dve triedy oddelené priamkou), tak konceptovým priestorom sú všetky možné dvojice disjunktných množín bodov.*

Hypotézovým priestorom budú hypotézy, ktoré správne klasifikujú body zadaných množín, napríklad priamky, ktoré oddelia tieto body do dvoch polrovín.

Učiacim algoritmom pre (K, H) tu môže byť učiaci algoritmus pre Rosenblattov perceptrón.

V ďalšom uvedieme dva veľmi jednoduché a veľmi všeobecné algoritmy, ktoré sú ale neefektívne. Neskôr budeme venovať pozornosť efektívnejším algoritmom.

2.3 Učenie sa pomocou konštrukcie

Nech X je príkladový priestor, t cieľový koncept, $X^+, X^- \subseteq X$ množina kladných príkladov.

Jeden spôsob učenia sa pre cieľový koncept t je skonštruovať množinu X^+ explicitne. Môžeme začať s prázdnu množinou prechodom cez tréningovú vzorku pridať každý pozitívny príklad. Formálne to môžeme vyjadriť pomocou nasledujúceho algoritmu.

Algoritmus 2.1 *Učenie sa pomocou konštrukcie*

begin

for all $x \in X$ set $h(x) = 0$;

for $i := 1$ to m do if $b_i = 1$ then set $h(x_i) = 1$;

$L(\bar{s}) = h$;

end

Je zrejmé, že pre záporné príklady a príklady, ktoré neboli v tréningovej vzorke algoritmus bude dávať odpoveď 0. Niektoré otázky, týkajúce sa algoritmu:

1. Čo ak X je nekonečný priestor príkladov? Odpoveď: Tréningová vzorka je vždy konečná, teda zostane nekonečne veľa záporných príkladov.
2. Ako vhodne vyjadriť hypotézový priestor tak, aby hypotézy boli vhodne vyjadriteľné?

Ak odsunieme otázku efektívnosti, zaujímavé sú nasledujúce poznámky. Zrejme, výstupná hypotéza $L(\bar{s})$ je rovná cieľovému konceptu $t \iff$ keď \bar{s} obsahuje všetky kladné príklady pre t . Pretože \bar{s} je konečná postupnosť, to znamená, že len koncepty s konečným počtom kladných príkladov môžu byť naučené úplne.

Napríklad, koncept "PARITA" je definovaný nad celým $\{0,1\}^*$, teda algoritmus nemôže skonštruovať celú množinu kladných príkladov. Ak sa obmedzíme na paritu reťazcov dĺžky n , tak koncept "PARITA" nad $\{0,1\}^n$ je naučiteľný, počet kladných príkladov je 2^{n-1} , preto musíme voľiť počet príkladov tréningovej vzorky aspoň tak veľký, t. j. $m = 2^{n-1}$.

Tento algoritmus má aj dobré vlastnosti:

1. je **konzistentný** t.j. výstupná hypotéza $L(\bar{s})$ klasifikuje všetky príklady vyskytujúce sa v \bar{s} korektne.
2. každý príklad tréningovej vzorky sa použije práve 1x. Toto je veľmi silná vlastnosť on-line vlastnosť. V praxi to znamená, že príklady môžu byť prezentované učiacemu algoritmu, bez nutnosti mať extra "pamäť", ktorá ich uloží pre ďalšie použitie.

Definícia 2.4 *Hovoríme, že algoritmus je bezpamäťový (on-line) algoritmus, ak pre danú tréningovú vzorku \bar{s} vytvára postupnosť hypotéz h_0, h_1, \dots, h_m , takých, že h_{i+1} závisí len od h_i a od priebežne spracovávaného príkladu vzorky (x_i, b_i) .*

2.4 Učenie sa očíslovaním

Nasledujúca metóda učenia určite nie je bezpamäťový on-line algoritmus. Predpokladáme, že hypotézový priestor H je spočítateľný a má explicitné očíslovanie hypotéz, $H = \{h^{(1)}, h^{(2)}, \dots\}$

Predpokladajme, že \bar{s} je tréningová vzorka pre cieľový koncept t .

Metóda: Porovnať každú hypotézu s každým príkladom v \bar{s} , odmietnuť hypotézu, ak nesúhlasí s hodnotou niektorého príkladu. Po odmietnutí hypotézy je tréningová vzorka testovaná tým istým spôsobom pre ďalšiu hypotézu. Proces sa zastaví, keď je nájdená hypotéza, ktorá vyhovuje všetkým príkladom tréningovej vzorky. Formálny zápis je vyjadrený v Algoritme [2.2](#)

Algoritmus 2.2 Učenie sa očíslovaním

Nech r - poradové číslo hypotézy, i - poradové číslo vzorky.

```
begin
   $r := 1; i := 1;$ 
  repeat
    if  $h^{(r)}(x_i) \neq b_i$  then
      begin  $r := r + 1; i := 1$  end
    else  $i := i + 1;$ 
  until  $i = m + 1;$ 
   $L(\bar{s}) := h^{(r)};$ 
end
```

Množina H môže byť konečná, a teda môže sa stať, že sa vhodná hypotéza nenájde. Modifikáciu algoritmu vieme ľahko urobiť. V praxi sa musíme vyhnúť používaniu neprimeraných veľkých hypotézových priestorov. Počet všetkých hypotéz $h : \{0, 1\}^n \rightarrow \{0, 1\}$ je 2^{2^n} . Ak $n = 10$, $2^{2^n} = 2^{1024} = 4^{512} = 8^{256} = 16^{128}$

Z poznámok vyplýva, že na to, aby sa táto metóda stala vhodnou metódou učenia, je potrebné urobiť určité obmedzenia na hypotézový priestor H a jeho vzťah k priestoru konceptov K . Toto vedie k pojmu ”**inductive bias**”. Je to predpoklad, že „učiaci sa“ má nejakú vopred nastavenú ideu o tom, akú metódu klasifikácie učenie sa používa, t.j. učiaci sa vie, alebo má nejaké informácie o konceptovom priestore.

Najjednoduchší spôsob modelovania takého predpokladu je stanoviť $H = K$ a v tomto prípade hovoríme o učiacom algoritme pre H , čo znamená, že pracujeme s (H, H) . Väčšina preberaných algoritmov v ďalšom bude tohto typu.

2.5 Úlohy

1. Aký je počet kladných príkladov konceptu ”palindrom”, keď príkladový priestor je $\{0, 1\}^n$?
2. Nech w je nasledujúci koncept: $\{0, 1\}^n, y \in \{0, 1\}^n, y = y_1 \dots y_n$:

$$w(y) = \begin{cases} 1 & \text{ak } y \text{ obsahuje najviac 2 jedničky} \\ 0 & \text{inak} \end{cases}$$

Ukážte, že počet kladných príkladov v tomto koncepte je kvadratickou funkciou n .

3. Predpokladajme, že v konečnom ”učení sa očíslovaním” sme si istí, že hypotézy sú očíslované tak, že tá, ktorú chceme, je v prvej polovici. Ak môžeme spracovať 1 milión hypotéz za sekundu a príkladový priestor je $\{0, 1\}^9$, koľko to bude trvať v najhoršom prípade?

Kapitola 3

Booleovské formuly (hypotézy) a ich reprezentácie

Booleovské formuly tvorené logickými spojkami, "AND (\wedge), OR (\vee), NOT (\neg)", nám poskytujú zjednodušený pohľad na koncepty a učiace algoritmy. Keďže sa tu pracuje v abecede $\{0, 1\}$ je jednoduché odhadnúť potrebné počty príkladov tréningovej vzorky, je ľahké vidieť nekonzistentnosť tréningových vzoriek. Budeme sa tu zaoberať dvoma priestormi konceptov, a síce množinou booleovských formúl, ktoré obsahujú len konjunkcie a negácie ("AND, NOT", t.j. monočlenmi) a zložitejším konceptom, v ktorom sú do formúl pridané aj disjunkcie ("AND, NOT, OR").

3.1 Učiaci sa algoritmus pre monočleny

Medzi najjednoduchšie priestory konceptov patrí množina *monočlenov*.

Monočlen je booleovská funkcia vyjadrená ako konjunkcia literálov, t. j. premenných alebo ich negácií. Konjunkcie vo formulách nebudeme vyjadrovať žiadnou spojkou (v tejto kapitole), negáciu premennej vyjadríme čiarou nad ňou. Napríklad formula x_1x_2 vyjadruje x_1 AND x_2 .

Monočlen

Literály sú premenné a ich negácie. Ak má formula n premenných, môže mať najviac $2n$ rôznych literálov.

Idea algoritmu pochádza od Valianta (1984). Konjunkcia nadobúda hodnotu 1 (true), ak sú všetky literály v nej pravdivé. Odstránením literálov, ktoré to kazia, dostaneme pravdivú formulu.

Začíname bez informácií, t.j. predpokladáme výskyt všetkých $2n$ rôznych literálov vo výslednej formule.

$$h_u : \quad u_1\bar{u}_1u_2\bar{u}_2 \dots u_n\bar{u}_n$$

Každý pozitívny príklad $y = y_1 \dots y_n$ umožňuje odstránenie tých literálov u_j , pre ktoré $y_j = 0$ a tých literálov \bar{u}_j , pre ktoré $y_j = 1$.

Algoritmus 3.1 Učenie monočlenov

Predpokladajme, že \bar{s} je tréningová vzorka

$$\bar{s} = ((x_1, b_1), \dots, (x_m, b_m))$$

$$x_i = ((x_i)_1(x_i)_2 \dots (x_i)_n), \quad 1 \leq i \leq m$$

$h_U \dots$ monočlenná funkcia obsahujúca literály v množine U .

begin

set $U = \{u_1, \bar{u}_1, \dots, u_n, \bar{u}_n\}$;

for $i:=1$ to m do

if $b_i=1$ then

for $j:=1$ to n do

if $(x_i)_j = 1$ then delete \bar{u}_j from U

else delete u_j from U ;

$L(\bar{s}) = h_U$;

end

Tento algoritmus sa nazýva *štandardný učiaci algoritmus pre monočleny*.

Redukovaná tréningová vzorka vznikne z pôvodnej tréningovej vzorky vynechaním tých položiek v tréningových príkladoch, ktoré neovplyvnia hodnotu h_U , t.j. im odpovedajúce premenné (ani ich negácie) sa vo formule h_U nevyskytujú.

Redukovaná tréningová vzorka

Veta 3.1 *Štandardný učiaci algoritmus pre monočleny je konzistentný s redukovanou tréningovou vzorkou, ak je bezosporná.*

Dôkaz: To znamená, že výsledná neprázdna formula je konzistentná s redukovanou tréningovou vzorkou, každý príklad je konzistentne klasifikovaný.

V každom kroku algoritmu je odstránených niekoľko (možno žiadny) literálov. Nech V je množina literálov, z ktorých je vytvorená výsledná formula.

Ak $h_U(x) = 1$ a $V \subset U$, tak $h_V(x) = 1$. Po každej prezentácii kladného príkladu x , mazacia procedúra zaručí, že $h_U(x) = 1$ a teda klasifikácia x je korektná. Teda koncová hypotéza $L(\bar{s})$ korektne klasifikuje všetky pozitívne príklady.

Ľubovoľný negatívny príklad pre cieľový koncept t (t.j. výsledkom t pre tento príklad je 0) je založený na niektorom literále v t , ktorý nebol odstránený.

Redukovaná vzorka je bezosporná, takže kladný a záporný príklad sa líšia v aspoň jednej položke. Nech je to i -tá položka.

Nech u_i je odpovedajúci literál v t . Ak u_i je premenná, tak všetky pozitívne príklady mali v tejto premennej nastavenú hodnotu 1. Záporný príklad bude mať v tejto položke 0. A teda výsledná hodnota t pre tento záporný príklad bude 0.

Ak u_i je negácia premennej, tak všetky pozitívne príklady mali v tejto premennej nastavenú hodnotu 0. Záporný príklad bude mať v tejto položke 1. A teda výsledná hodnota t pre tento záporný príklad bude 0.

Teda všetky negatívne príklady redukovanej tréningovej vzorky pre t (a čiastočne tie vo vzorke) sú korektne klasifikované pomocou $L(\bar{s})$.

Príklad 3.1 *Predpokladajme, že počet premenných vo formule je $n = 8$ a je daná tréningová vzorka*

$$(10101010, 1), (00001111, 0), (10101111, 1), (10100011, 1)$$

Na začiatku predpokladáme, že vo formule je všetkých 16 literálov. Kroky úpravy po použití jednotlivých príkladov:

1. *Pred 1. príkladom: $u_1\bar{u}_1u_2\bar{u}_2u_3\bar{u}_3u_4\bar{u}_4u_5\bar{u}_5u_6\bar{u}_6u_7\bar{u}_7u_8\bar{u}_8$*
2. *Po 1. príklade: $u_1\bar{u}_2u_3\bar{u}_4u_5\bar{u}_6u_7\bar{u}_8$*
3. *Po 2. príklade: $u_1\bar{u}_2u_3\bar{u}_4u_5\bar{u}_6u_7\bar{u}_8$. Druhý príklad je negatívny, teda nič neovplyvní.*
4. *Po 3. príklade: $u_1\bar{u}_2u_3\bar{u}_4u_5u_7$*
5. *Po 4. príklade: $u_1\bar{u}_2u_3\bar{u}_4u_7$*

Algoritmus pre učenie monočlenov má časovú zložitosť $O(m.n)$, kde m je počet príkladov tréningovej vzorky a n je počet premenných monočlena. Počet príkladov je najviac 2^n . Počet možných formúl je 3^n (premenná, jej negácia alebo ani jedna z nich).

3.2 Učenie disjunkcií malých monočlenov

Pri práci s booleovskými premennými sa stretáme s ich štandardnými tvarmi, s ktorými sa dobre pracuje. Sú to

- **Disjunktívna normálna forma (DNF)**

$$\mu_1 \vee \mu_2 \vee \cdots \vee \mu_n$$

kde μ_i je monočlen, $1 \leq i \leq r$.

- **Konjunktívna normálna forma (KNF)**

$$\gamma_1 \wedge \gamma_2 \wedge \cdots \wedge \gamma_n$$

kde γ_i je, $1 \leq i \leq n$ je klauzula, t.j. disjunkcia literálov.

Tvary booleovských funkcií vieme ovplyvniť napríklad tým, že určíme ohraničenia na počty premenných v klauzulách alebo v monočlenoch.

Označenie:

- M_n -množina monočlenov nad $\{0, 1\}^n$,
- $M_{n,k}$ -množina monočlenov nad $\{0, 1\}^n$, z ktorých každý má najviac k literálov,
- $D_{n,k}$ -množina disjunkcií členov z množiny $M_{n,k}$.

Nasledujúci algoritmus navrhol Valiant (1984). Algoritmus začína počiatočnou hypotézou, ktorá je disjunkciou všetkých monočlenov dĺžky najviac k . Každý ďalší krok aplikuje jednoduchú logickú dedukciu, a síce: pre záporný príklad je potrebné, aby všetky monočleny mali hodnotu 0, preto odstránime tie monočleny, ktoré to kazia (nadobúdajú hodnotu true).

Algoritmus 3.2 *Učenie disjunkcií malých monočlenov*

Predpokladáme, že dĺžka tréningovej vzorky je m , μ je monočlen a $\mu(x_i)$ predstavuje booleovskú hodnotu monočlena μ pre príklad x_i .

begin

*$h :=$ disjunkcia všetkých monočlenov dĺžky najviac k ;
for $i := 1$ to m do
if $(b_i = 0) \wedge (h(x_i) = 1)$ then
vymazať monočleny μ , pre ktoré $\mu(x_i) = 1$;
 $L(\bar{s}) := h$;*

end

Je zrejmé, že pre negatívne príklady naučený algoritmus bude dávať správne odpovede. V prípade pozitívnych príkladov je potrebné opäť uvažovať bezospornú redukovanú vzorku príkladov, na ktorej sa bude chovať konzistentne.

Časová zložitosť algoritmu je $O(m \cdot d)$, kde d je počet všetkých uvažovaných monočlenov.

Príklad 3.2 *Použitie algoritmu pre učenie disjunkcií malých monočlenov. Budeme pracovať s formulami v $D_{3,2}$.*

Zoznam všetkých relevantných monočlenov je:

$$U = \{u_1, u_2, u_3, \bar{u}_1, \bar{u}_2, \bar{u}_3, u_1u_2, u_1u_3, u_2u_3, \bar{u}_1u_2, \bar{u}_1u_3, \bar{u}_2u_3, u_1\bar{u}_2, u_1\bar{u}_3, u_2\bar{u}_3, \bar{u}_1\bar{u}_2, \bar{u}_1\bar{u}_3, \bar{u}_2\bar{u}_3\}.$$

Majme tréningovú vzorku:

(000, 1), (001, 1), (010, 1), (011, 1), (100, 0), (101, 0), (110, 1), (111, 0)

Algoritmus nepoužije prvé štyri príklady, lebo sú kladné. V ďalších príkladoch jeho chovanie je nasledujúce:

1. z U budú odstránené monočleny $u_1, \bar{u}_2, \bar{u}_3, u_1\bar{u}_2, u_1\bar{u}_3, \bar{u}_2\bar{u}_3$;
2. z U budú odstránené monočleny $u_3, u_1u_3, \bar{u}_2u_3, \bar{u}_1\bar{u}_3$;
3. U zostane nezmenené;
4. z U budú odstránené monočleny $\bar{u}_3, u_1u_2, u_1u_3, u_2u_3$;

Výsledná formula má tvar $h = \bar{u}_1 \vee \bar{u}_1u_2 \vee \bar{u}_1u_3 \vee u_2\bar{u}_3 \vee \bar{u}_1\bar{u}_2 \vee \bar{u}_1\bar{u}_3$.

Túto formulu je možné zjednodušiť na tvar: $\bar{u}_1 \vee u_2\bar{u}_3$.

3.3 Reprezentácia hypotézového priestoru

Učenie prediskutované v tejto časti malo zjednodušené predpoklady, a síce že konceptový priestor je ten istý ako hypotézový. V skutočnosti sme uvažovali o tom, že cieľové koncepty majú nejaký popis pomocou logických formlí alebo strojov. Hoci tento predpoklad sa môže zdať reštriktívny, je prirodzený pri hľadaní riešení viacerých problémov.

3.4 Čas behu učiacich algoritmov

Učiace algoritmy popísané v predchádzajúcom texte sa zaoberajú booleovskými konceptami. V týchto prípadoch príkladový priestor je $\{0, 1\}^n$ pre nejaké pevné n a hypotézový priestor je množina funkcií definovaných na príkladovom priestore. Pre každý z týchto algoritmov parameter n je ľubovoľný v zmysle, že algoritmus je definovaný pre ľubovoľné n a navyše operuje v podstate tým istým spôsobom pre každú hodnotu n .

Napríklad, štandardný učiaci algoritmus pre priestor M_n monočlenov je definovaný celkom všeobecne, hoci výpočet vždy prebieha pre danú hodnotu n . Chceme kvantifikovať chovanie sa učiacich algoritmov vzhľadom na n , a je obvykle používať nasledujúce definície.

Definícia 3.1 *Hovoríme, že zjednotenie hypotézových priestorov $H = \bigcup H_n$ je odstupňované (graded) príkladmi veľkosti n , ak H_n označuje hypotézový priestor definovaný len na príkladoch z X^n .*

Definícia 3.2 *Nech je daný učiaci algoritmus pre $H = \bigcup H_n$, t. j. $L : X^* \rightarrow H$, taký, že ak \bar{s} je tréningová vzorka pre $h \in H_n$, tak $L(\bar{s}) \in H_n$. Hovoríme, že L podporuje stupňovanie (grading) priestoru H .*

Uvažujme učiaci algoritmus L pre booleovský hypotézový priestor $H = \bigcup H_n$, odstupňovaný veľkosťou príkladov. Vstup do L je tréningová vzorka, ktorá pozostáva z m n -bitových vektorov spolu s m 1-bitovými označeniami. Celkový počet bitov na vstupe je $m(n+1)$ a bolo by možné použiť toto jediné číslo ako mieru veľkosti vstupu. Avšak je výhodné sledovať aj m aj n oddelene.

Použijeme $R_L(m, n)$ **na označenie najhoršieho času behu** L na tréningovej vzorke m n -bitových vektorov.

Príklad 3.3 *Nech L je učiaci algoritmus pre monočleny popísaný vyššie. Hypotézový priestor je zjednotenie $\bigcup M_n$. Hlavný krok algoritmu vyžaduje kontrolu každého bitu u každého pozitívneho príkladu a možno vymazanie niektorých literálov. V najhoršom prípade, každý príklad v tréningovej vzorke môže byť pozitívny príklad, a tak by sme mali ošetriť tento krok m -krát, každý krok obsahuje kontrolu n bitov. Iné časti výpočtu vyžadujú porovnateľne toľko operácií, takže môžeme hovoriť, že čas behu $R_L(m, n)$ je v tomto prípade $O(m * n)$.*

3.5 Konzistencia tréningovej vzorky

Nech $H = \bigcup H_n$ je hypotézový priestor booleovských funkcií odstupňovaný príkladmi veľkosti n . Problém konzistencie tréningovej vzorky pre H môže byť stanovený nasledovne:

H - konzistencia:

- Inštancia: Tréningová vzorka \bar{s} vyjadrená n -bitovými ohodnotenými príkladmi.
- Otázka: Existuje hypotéza v H_n konzistentná s \bar{s} ?

Ukážeme, že v niektorých netriviálnych prípadoch je tento problém NP-ťažký. Na to, aby sme vyjadrili praktický dosah tohto výsledku, potrebujeme urobiť niekoľko všeobecných komentárov.

Označme C_n^k klauzulu, t. j. disjunkciu vytvorenú výberom najviac k literálov z n premenných. Ukážeme, že pre pevné $k \geq 3$, problém konzistencie pre $C^k = \bigcup C_n^k$ je NP-ťažký. Teda nie je pravda, že existuje polynomiálny učiaci algoritmus pre C_n^k , ktorý produkuje konzistentnú hypotézu.

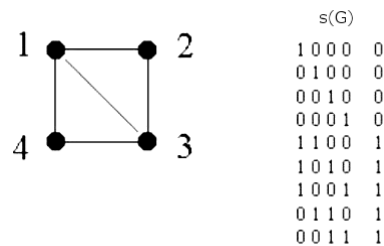
Dôkaz spočíva v prevedení problému farbenia grafu s n vrcholmi pomocou k farieb na tento problém, čo je NP-úplný problém pre $k \geq 3$ (Garey, Johnson 1979).

$G = (V, E)$, k -farbenie je funkcia $\chi : V \rightarrow \{1, 2, \dots, k\}$ s vlastnosťou: ak $\langle v_i, v_j \rangle \in E$, potom $\chi(v_i) \neq \chi(v_j)$.

Problém farbenia grafu:

Predpokladajme, že máme $G = (V, E)$, $V = \{1, 2, \dots, k\}$. Skonstruujeme tréningovú vzorku $s(G)$ nasledovne: Pre každý vrchol $i \in V$ určíme záporný príklad vektor v_i , ktorý má 1 v pozícii i -tej súradnice a 0 inde. Pre každú hranu $\langle i, j \rangle \in E$ vezmeme ako pozitívny príklad vektor $v_i + v_j$.

Príklad 3.4 Na obrázku je znázornený graf $G = (\{1, 2, 3, 4\}, \{12, 13, 14, 23, 34\})$. tréningová vzorka je vytvorená príkladmi kvrcholom a hranám.



Obr. 3.1: Graf spolu s pripravenou tréningovou vzorkou

Veta 3.2 Existuje funkcia $h \in C_n^k$, ktorá je konzistentná so vzorkou $s(G) \iff$ graf G je k -zafarbiteľný.

Dôkaz:

\Rightarrow

Predpokladajme, že $h \in C_n^k$ a je konzistentná s tréningovou vzorkou. Podľa definície h je konjunkcia $h = h_1 \wedge h_2 \wedge \dots \wedge h_k$ klauzúl. Pre každý vrchol $i \in V$, $h(v_i) = 0$ a teda musí ex. aspoň 1 klauzula h_f , pre ktorú $h_f(v_i) = 0$. Funkcia $\chi : V \rightarrow \{1, 2, \dots, k\}$ taká, že:

$$\chi(i) = \min\{f \mid h_f(v_i) = 0\}$$

Zostáva ukázať, že χ je farbenie grafu G ; inak povedané, ak i a j sú dva vrcholy, pre ktoré $\chi(i) = \chi(j)$, potom $\langle i, j \rangle \notin E$. Predp., že $\chi(i) = \chi(j) = f$ a teda $h_f(v_i) = h_f(v_j) = 0$. Pretože h_f je klauzula, každý literál, ktorý sa v nej vyskytuje musí byť 0 na v_i a na v_j . Teraz v_i má 1 len v i -tej pozícii a tak $h_f(v_i) = 0$ implikuje, že len jeden negovaný

literál, ktorý sa môže vyskytnúť v h_f je \bar{u}_i . Pretože to isté platí pre \bar{u}_j , dostávame, že h_f obsahuje len niektoré literály u_z , pre ktoré $z \neq i, j$. Teda $h_f(v_i + v_j) = 0$ a $h(v_i + v_j) = 0$. Ak by $\langle i, j \rangle$ bola hranou v G , potom by platilo $h(v_i + v_j) = 1$, pretože sme predpokladali, že h je konzistentná s $s(G)$. Teda $\langle i, j \rangle$ nie je hranou v G a χ je farbenie.

⇐

Predpokladajme, že je dané farbenie $\chi : V \rightarrow \{1, 2, \dots, k\}$. Pre $1 \leq f \leq k$ definujme h_f ako klauzulu

$$\langle \bigvee u_{i\chi(i) \neq f} \rangle$$

a definujme $h = h_1 \wedge h_2 \wedge \dots \wedge h_k$. Tvrdíme, že h je konzistentná s $s(G)$.

Najprv, predpokladajme, že pre vrchol i $\chi(i) = g$. Klauzula h_g je definovaná tak, že obsahuje len tie (nie negované) literály, ktoré zodpovedajú vrcholom nezafarbeným farbou g , a teda u_i sa nenachádza v h_g . Teda $h_g(v_i) = 0$ a $h(v_i) = 0$.

Ďalej, nech $\langle i, j \rangle$ je hrana v G . Pre každú farbu f aspoň jedno v_i alebo v_j nemá farbu f ; označme vhodný výber $i(f)$. Potom h_f obsahuje literál $u_{i(f)}$, ktorý je 1 na $v_i + v_j$. Teda klauzula h_f je 1 na $v_i + v_j$ a $h(v_i + v_j) = 1$, ako sme požadovali.

Príklad 3.5 Daný graf je 3-zafarbitelný. Odpovedajúce formuly sú: $h_1 = v_2 \vee v_3 \vee v_4$, $h_2 = v_1 \vee v_3$, $h_3 = v_1 \vee v_2 \vee v_4$, $h = h_1 \wedge h_2 \wedge h_3$



Obr. 3.2: Graf s ofarbením vrcholov

3.6 Úlohy

1. Napíšte postupnosť hypotéz generovaných algoritmom učenia monochlenov, keď na vstupe je prezentovaná tréningová vzorka

$$(11100101, 1), (00100011, 0), (11001001, 1)$$

Ak cieľový koncept je $\langle u_2 \bar{u}_4 u_8 \rangle$, doplňte príklady do vzorky, ktoré sú pre to nutné.

2. Napíšte *DNF* formuly pre koncepty *parita* a *palindrom* na príkladových priestoroch $\{0, 1\}^5$.

3. Uveďte príklad booleovskej funkcie 3 premenných, ktorá nie je $D_{3,2}$.
4. Prečo je obvyčajne vhodné uvažovať veľkosť vstupu v tvare $\lg n$ namiesto n , keď uvažujeme o otázkach efektívnosti?
5. Nasledujúci algoritmus je rýchlym algoritmom pre výpočet m -tej mocniny daného čísla u . (výstupom je finálna hodnota uložená v *bot*.)

```

bot:=1; top:=u; q:=m;
while q>0 do
begin
  if q mod 2 =1 then bot:=top*bot;
  top:=sqr(top);
  q:=q div 2;
end;

```

Ukážte, že efektívnosť algoritmu je $O(s)$, kde s je miera veľkosti m , ako to bolo povedané v predchádzajúcom príklade.

6. Navrhните algoritmus, ktorý rozhodne, či daný n -bitový reťazec je palindrom a odhadnite jeho efektívnosť vzhľadom na veľkosť vstupu n .
7. Nech $G = (V, E)$ je graf s množinou vrcholov $V = \{1, 2, 3, 4, 5\}$ a množinou hrán $E = \{12, 13, 15, 23, 25, 34, 35\}$. Vytvorte odpovedajúcu tréningovú vzorku $s(G)$ podľa postupu v dôkaze. Nájdite najmenšiu možnú hodnotu k , pre ktorú je funkcia $h \in C_5^k$ konzistentná s $s(G)$ a danú formulu vyjadrite explicitne.
8. Nech $C_n = C_n^1$, čo je priestor booleovských funkcií na $\{0, 1\}^n$, ktorý môže byť reprezentovaný jednou klauzulou. Sformulujte konzistentný učiaci algoritmus pre C_n , ktorý je "duálny" k štandardnému algoritmu pre monočleny a vyhodnoďte tvrdenie, že jeho čas behu je polynomiálny v m a n .

Kapitola 4

Pravdepodobnostné učenie

4.1 Algoritmus pre učenie lúčov

V tejto kapitole sa budeme zaoberať veľmi jednoduchým algoritmom pre učenie sa v reálnom hypotézovom priestore. Budeme sa snažiť odhadnúť, aký presný je dosiahnutý výsledok a určiť do akej miery je dôveryhodný.

Pôjde o učenie sa jednej reálnej meranej hodnoty (napríklad objemu telesa), ktorá existuje, z postupnosti už nameraných hodnôt, pričom vieme, že namerané hodnoty nemôžu byť menšie než skutočná hodnota (skutočný kladný objem).

Koncept, ktorým sa budeme zaoberať je **zlava uzavretý a sprava otvorený interval (lúč)** v množine reálnych čísel \mathbb{R} , t.j. $\langle \Theta, \infty \rangle$ pre ľubovoľné reálne číslo Θ . Budeme ho označovať r_Θ . Tento koncept je definovaný na príkladovom priestore \mathbb{R} funkciou

$$r_\Theta(y) = y \iff y \geq \Theta$$

Algoritmus pre učenie v hypotézovom priestore $H = \{r_\Theta | \Theta \in \mathbb{R}\}$ je založený na idee, že za aktuálnu hypotézu vezmeme "najmenší" lúč obsahujúci všetky pozitívne príklady v tréningovej vzorke.

V prípade, že neexistujú kladné príklady, vtedy budeme hovoriť o **prázdnom lúči**. Bude označovaný $r_{+\infty}$.

Pre danú tréningovú vzorku

$$\bar{s} = ((x_1, b_1), (x_2, b_2), \dots, (x_m, b_m))$$

výstupná hypotéza $L(\bar{s})$ by mala byť r_λ , kde

$$\lambda = \lambda(\bar{s}) = \min_{1 \leq i \leq m} \{x_i | b_i = 1\}$$

$\lambda = +\infty$, ak vzorka neobsahuje kladné príklady. Jednoduchá modifikácia algoritmu, ktorý počíta minimum konečnej množiny je postačujúca pre naše účely. Toto poskytuje nasledujúci bezpamäťový on-line algoritmus:

Algoritmus 4.1 Učenie lúčov

Je daná tréningová vzorka $\bar{s} = ((x_1, b_1), (x_2, b_2), \dots, (x_m, b_m))$

Učenie lúčov

begin

 set $\lambda = +\infty$;

 for $i:=1$ to m do

 if $(b_i = 1)$ and $(x_i < \lambda)$ then set $\lambda = x_i$;

$L(\bar{s}) := r_\lambda$;

end

Je ľahké vidieť, že ak tréningová vzorka je pre cieľovú hypotézu r_Θ , potom $L(\bar{s})$ bude lúč r_λ s $\lambda = \lambda(s) \geq \Theta$. Pretože je len konečný počet príkladov v tréningovej vzorke a celý príkladový priestor je nespočítateľný, nemôžeme očakávať, že $\lambda = \Theta$. Avšak, zdá sa, že ak dĺžka tréningovej vzorky rastie, tak by sa mala zvyšovať pravdepodobnosť toho, že chyba medzi r_λ a r_Θ je veľmi malá (prípadne klesá).

Prakticky táto vlastnosť môže byť charakterizovaná nasledovne:

Predpokladajme, že spustíme algoritmus s veľkou tréningovou vzorkou a potom sa rozhodneme použiť výstupnú hypotézu r_λ pre cieľovú (neznámu) hypotézu r_Θ . Inak povedané, uspokojíme sa s tým, že "učiaci sa" bol adekvátne trénovaný. Ak λ nie je blízke k Θ , toto indikuje, že pozitívne príklady, ktoré by boli blízke Θ sú relatívne nepravdepodobné a nevyskytovali sa v tréningovej vzorke. Z toho vyplýva, keď klasifikujeme niektoré ďalšie príklady, ktoré sú síce prezentované podľa toho istého rozloženia, aj tak môžeme urobiť niekoľko chýb ako dôsledok použitia r_λ namiesto r_Θ .

Zaujíma nás, akej chyby sa dopustíme a do akej miery je výsledok dôveryhodný. Je zrejmé, že výsledok bude závisieť od tréningovej vzorky, resp. od pravdepodobnostného rozdelenia príkladov.

4.2 Pravdepodobnostné aproximačne správne učenie, PAC učenie

□

Uvažujme model, v ktorom tréningová vzorka \bar{s} pre cieľový koncept t je generovaná výberom príkladov x_1, x_2, \dots, x_m z X "náhodne" podľa nejakého známeho, pevne daného, pravdepodobnostného rozdelenia. Učiaci algoritmus L produkuje hypotézu $L(\bar{s})$, ktorá je očakávaná ako dobrá aproximácia pre t . Vyžadujeme nasledujúce: ak počet príkladov m v trénujúcej vzorke vzrastie, tak z pravdepodobnosti vyplynie, že chyba, ktorá je

¹PAC - Probably Approximately Correct learning

výsledkom použitia $L(\bar{s})$ namiesto t , je menšia.

Základné pojmy:

X - pravdepodobnostný priestor,

\mathbf{A} - trieda podmnožín množiny X ,

μ - pravdepodobnostné rozdelenie, miera pravdepodobnosti, z \mathbf{A} do intervalu $\langle 0, 1 \rangle$.

$$\mathbf{A} \rightarrow \langle 0, 1 \rangle.$$

Od triedy \mathbf{A} sa vyžaduje, aby bola uzavretá vzhľadom na operácie komplementu, konečného prieniku a spočítateľného zjednotenia.

Prvok $A \in \mathbf{A}$, sa nazýva *udalosť* a $\mu(A)$ je pravdepodobnosť udalosti A .

Od μ sa vyžaduje, aby spĺňovala nasledujúce podmienky:

$$\mu(\emptyset) = 0, \mu(X) = 1,$$

a pre ľubovoľné po dvoch disjunktné množiny $A_1, A_2, \dots \in \mathbf{A}$

$$\mu\left(\bigcup_{i=1}^{\infty} A_i\right) = \sum_{i=1}^{\infty} \mu(A_i).$$

Pre nás X je príkladový priestor, a príklady sú buď booleovské alebo reálne. V prvom prípade je X konečná alebo spočítateľná množina, a teda \mathbf{A} môže byť triedou všetkých podmnožín X . V reálnom prípade môžeme zobrať za \mathbf{A} ľubovoľnú dostatočne veľkú triedu obsahujúcu množiny, ktoré potrebujeme uvažovať; postačí uvažovať triedu borelovských množín v R^n .

V oboch prípadoch budeme používať vhodnú triedu bez explicitného vyjadrovania detailov.

Budeme jednoducho hovoriť *pravdepodobnostné rozdelenie μ na X* , čím mienime funkciu μ definovanú na vhodnej triede \mathbf{A} a splňujúcej axiomy uvedené vyššie. Treba zdôrazniť, že v aplikáciách, o ktorých sme sa zmieňovali, nerobíme žiadne predpoklady o μ , okrem podmienok uvedených v definícii. Situácia, ktorú sme modelovali, je vyjadrená množinou príkladov prezentovaných učiacemu sa a chovácej sa (modeluje) podľa nejakého pevného, ale neznámeho rozdelenia. Učiteľovi je povolené klasifikovať príklady ako pozitívne a negatívne, ale nemôže riadiť postupnosť, v ktorej príklady budú prezentované.

Budeme pokračovať s predpokladom, že cieľový koncept patrí do hypotézového priestoru H , ktorý je dostupný učiacemu sa. K danému cieľovému konceptu $t \in H$ definujeme chybu ľubovoľnej hypotézy $h \in H$ vzhľadom na t a bude to pravdepodobnosť udalosti $h(x) \neq t(x)$, t. j.

$$err_{\mu}(h, t) = \mu\{x \in X \mid h(x) \neq t(x)\}.$$

v kučeravých zátvorkách je error set - chybová množina a predpokladáme, že existuje udalosť taká, že tejto udalosti môže byť priradená pravdepodobnosť. Keď pôjde o t známe z konceptu, budeme tiež používať označenie $err_{\mu}(h)$.

Príklad 4.1 Pravdepodobnosť chyby

Nech $X = \{0, 1\}^3$, predpokladajme, že cieľový koncept je $\langle u_1 \rangle$. Chybová množina pre hypotézu $\langle u_1 \bar{u}_2 \rangle$ obsahuje dva príklady, 110 a 111. Tak

$$err_{\langle u_1 \bar{u}_2 \rangle} = \mu\{110, 111\}.$$

Napríklad, ak μ - rovnomerné rozdelenie na X - $\frac{1}{8}$ je pravdepodobnosť každého príkladu, potom

$$err_{\langle u_1 \bar{u}_2 \rangle} = \frac{1}{4}.$$

Ak z nejakých dôvodov príklady, u ktorých je $y_2 = 1$ sú málo pravdepodobné, potom err_{μ} bude o niečo menšia.

Keď je daná množina X so štruktúrou pravdepodobnostného priestoru, karteziánsky súčin množín X^m preberá pravdepodobnostnú štruktúru X . Detaily nebudeme rozoberať, je postačujúce poznamenať, že konštrukcia nám umožňuje považovať komponenty m -tice (x_1, x_2, \dots, x_m) za nezávislé premenné, rozdelenie každej z nich je podľa pravdepodobnostného rozdelenia μ na X . Odpovedajúce pravdepodobnostné rozdelenie na X^m je označované μ^m .

Neformálne, pre dané $Y \subseteq X^m$ budeme interpretovať hodnotu $\mu^m(Y)$ ako **”pravdepodobnosť, že náhodná vzorka m príkladov vybratých z X podľa rozdelenia μ patrí do Y ”**.

Nech $S(m, t)$ označuje množinu tréningových vzoriek dĺžky m pre daný cieľový koncept t , kde príklady sú vybrané z príkladového priestoru X . Ľubovoľná postupnosť príkladov $x \in X$ determinuje a je determinovaná tréningovou vzorkou $\bar{s} \in S(m, t)$:

ak $x = (x_1, x_2, \dots, x_m)$, potom $\bar{s} = ((x_1, t(x_1)), (x_2, t(x_2)), \dots)$.

Inak povedané, existuje zobrazenie Φ

$$\Phi : X^m \rightarrow S(m, t), \quad \text{pre ktorú} \quad \Phi(x) = \bar{s}$$

Teda môžeme interpretovať pravdepodobnosť, že $\bar{s} \in S(m, t)$ má nejakú danú vlastnosť P , nasledujúcim spôsobom. Definujeme

$$\mu^m\{\bar{s} \in S(m, t) \mid \bar{s} \text{ má vlastnosť } P\}$$

to znamená

$$\mu^m\{x \in X^m \mid \Phi(x) \in S(m, t) \text{ má vlastnosť } P\}$$

Z toho vyplýva, že keď príkladový priestor X je vybavený pravdepodobnostným rozdelením μ , môžeme zaviesť precíznejšiu interpretáciu pre

- (i) chybu hypotézy, ktorá vznikne, keď učiaci algoritmus L pracuje s \bar{s} ; Táto veličina pracuje s $err_\mu(L(\bar{s}, t))$ alebo zjednodušene $err_\mu(L(\bar{s}))$.
- (ii) a pre pravdepodobnosť, že táto chyba je menšia než vopred zvolené $0 < \epsilon < 1$.

Druhá je pravdepodobnosť vzhľadom na μ^m , teda že \bar{s} má vlastnosť $err_\mu(L(\bar{s})) < \epsilon$

**PAC -
algoritmus**

Definícia 4.1 PAC - algoritmus

Hovoríme, že algoritmus L je **probably approximately correct** (pravdepodobnostne aproximáčne správny) učiaci algoritmus, pre hypotézový priestor H , ak

- k ľubovoľnému reálnemu číslu δ , ... dôveryhodnosť, $0 \leq \delta \leq 1$
- k ľubovoľnému reálnemu číslu ϵ , ... pravdepodobnosť chyby učenia sa, $0 \leq \epsilon \leq 1$
- existuje kladné celé číslo $m_0 = m_0(\delta, \epsilon)$ také, že
- pre ľub. cieľový koncept $t \in H$,
pre ľub. pravdepodobnostné rozloženie μ na X pre všetky $m \geq m_0$ platí

$$\mu^m \{ \bar{s} \in S(m, t) \mid err_\mu(L(\bar{s}, t)) < \epsilon \} > 1 - \delta$$

t.j.

$$\forall_{(0 \leq \delta \leq 1)} \forall_{(0 \leq \epsilon \leq 1)} \exists_{(m_0 = m_0(\epsilon, \delta))} \forall_{(t \in H)} \forall_{(\mu \text{ na } X)} \forall_{(m \geq m_0)} \mu^m \{ \bar{s} \in S(m, t) \mid err_\mu(L(\bar{s}, t)) < \epsilon \} > 1 - \delta$$

Skutočnosť, že m_0 závisí od δ a ϵ , ale nie od t a μ odráža to, že učiaci sa môže byť schopný špecifikovať predpokladanú úroveň dôvery a presnosti, aj keď cieľový koncept a rozloženie príkladov sú neznáme. Dôvodom k tomu, že je možné splniť podmienku pre ľubovoľné μ je, že vyjadruje vzťah medzi dvoma veličinami, ktoré obsahujú μ : chyba err_μ a pravdepodobnosť vzhľadom na μ^m určitej množiny.

PAC učenie je, v istom zmysle, najlepšie, v čo môžeme dúfať pri tomto pravdepodobnostnom pohľade. Nereprezentatívne tréningové vzorky, hoci málo pravdepodobné, budú príležitostne prezentované učiacemu algoritmu, ale môžeme očakávať, že je viac pravdepodobné, že je prezentovaná použiteľná tréningová vzorka. Naopak, aj keď máme reprezentatívnu tréningovú vzorku, rozšírenie tréningovej vzorky nebude vo všeobecnosti koincidovať s cieľovým konceptom, takže aj tak výstupná hypotéza bude len aproximáčne správna.

4.3 Učenie lúčov je PAC

Veta 4.1 *Algoritmus L pre učenie lúčov je PAC.*

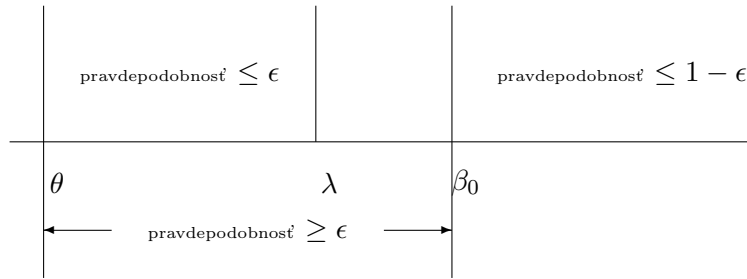
Proof. Predpokladajme, že $\delta, \epsilon, r_\Theta$ a μ sú zvolené, ale sú pevné. Nech \bar{s} je tréningová vzorka dĺžky m pre r_Θ a nech $L(\bar{s}) = r_\lambda$. Zrejme, chybový interval je $\langle \Theta, \lambda \rangle$. Pre danú hodnotu ϵ a dané μ definujme

$$\beta_0 = \beta_0(\epsilon, \mu) = \sup\{\beta \mid \mu\langle \Theta, \beta \rangle < \epsilon\}$$

Ak uvažujeme $\lambda \leq \beta_0$ tak máme

$$\text{err}_\mu(L(\bar{s}, t)) = \mu\langle \Theta, \lambda \rangle \leq \mu\langle \Theta, \beta_0 \rangle \leq \epsilon$$

Udalosť, že \bar{s} má vlastnosť $\lambda \leq \beta_0$ je práve udalosť, že aspoň jeden príklad v \bar{s} je v intervale $\langle \Theta, \beta_0 \rangle$.



Obr. 4.1: V tejto situácii chyba výstupu je aspoň ϵ

Pretože $\mu\langle \Theta, \beta_0 \rangle \geq \epsilon$, pravdepodobnosť, že jeden príklad nie je v tomto intervale, je najviac $1 - \epsilon$. Preto pravdepodobnosť, že žiadny z m príkladov vzorky \bar{s} nie je v tomto intervale je najviac $(1 - \epsilon)^m$. Keď budeme uvažovať komplementárnu udalosť (existuje príklad, ktorý je z tohto intervalu), z toho vyplýva, že pravdepodobnosť, že $\lambda \leq \beta_0$ je aspoň $1 - (1 - \epsilon)^m$.

Ako sme už poznamenali vyššie, že udalosť $\lambda \leq \beta_0$ implikuje udalosť $\text{err}_\mu(L(\bar{s}, t)) \leq \epsilon$ a tak

$$\mu^m\{s \in S(m, r_\Theta) \mid \text{err}_\mu(L(\bar{s}, t)) \leq \epsilon\} \geq 1 - (1 - \epsilon)^m$$

Položíme

$$m \geq m_0 = \frac{1}{\epsilon} * \ln \frac{1}{\delta}$$

$$(1 - \epsilon)^m \leq (1 - \epsilon)^{m_0} < e^{-\epsilon m_0} < e^{\ln \delta} = \delta$$

tento výpočet ukazuje, že algoritmus je PAC.

Dôkaz korektnosti poskytuje explicitnú formulu pre dĺžku vzorky postačujúcu na to, aby boli splnené predpísané hodnoty presnosti a dôveryhodnosti. Predpokladajme, že $\delta = 0.001$, $\epsilon = 0.01$

$$m_0 = \frac{1}{0.01} * \ln \frac{1}{0.001} = 100 * \ln 1000 = 691.$$

Takže aspoň 691 príkladov je treba, aby sme si boli istí na 99,9%, že najviac 1% príkladov bude klasifikovaných nesprávne, za predpokladu, že sú z toho istého zdroja ako tréningová vzorka.

Odvodenie vzťahu:

$$\delta = (1 - \epsilon)^m$$

$$\ln \delta = m * \ln(1 - \epsilon) = m * \frac{-\epsilon}{1-\epsilon} \quad \ln \delta \leq m * \frac{-\epsilon}{1-\epsilon}$$

$$\ln(1 - \epsilon) \leq \ln(1) + f'(0) * \frac{\epsilon}{1!} \leq 0 + \frac{1}{1-\epsilon} * (-1) * \frac{\epsilon}{1!}$$

$$-\ln \delta = m * \frac{\epsilon}{1-\epsilon} \quad \frac{1-\epsilon}{\epsilon} \ln \frac{1}{\delta} \leq m \Rightarrow \left(\frac{1}{\epsilon} - 1\right) \ln \frac{1}{\delta} \leq m$$

$$\frac{1}{\epsilon} * \ln \frac{1}{\delta} + \ln \delta \leq m$$

4.4 Exaktné učenie

Keď príkladový priestor X je konečný, pojem PAC - učenie má ďalšie dodatočné obmedzenia. Začneme tým, že ľubovoľné pravdepodobnostné rozdelenie na konečnej množine X je determinované hodnotami na jej 1-prvkových množinách $\{x\}$, použitím axiómy o aditivite. Budeme písať $\mu(x)$ namiesto $\mu(\{x\})$. Ak budú nejaké príklady, pre ktoré $\mu(x) = 0$, s pravdepodobnosťou 1 sa nebudú vyskytovať v konečnej náhodnej vzorke a môžu byť ignorované. Inými slovami, môžeme, ak je to nutné, predefinovať X tak, že $\mu(x) > 0$ pre všetky $x \in X$. Pretože X je konečný, veličina

$$\epsilon_\mu = \min_{\{x \in X\}} \mu(x) > 0$$

je dobre definovaná.

Predpokladajme, že máme algoritmus L , ktorý je PAC pre hypotézový priestor H definovaný na X . Vo význame definície PAC algoritmu máme dané δ, ϵ, μ , a t v ich obvyklom význame

$$m \geq m_0 \Rightarrow \mu^m \{ \bar{s} \in S(m, t) \mid \text{err}_\mu(L(\bar{s})) < \epsilon \} > 1 - \delta$$

Predpokladajme, že presnosť ϵ je vybraná tak, aby nebola väčšia než ϵ_μ . Potom podmienka $\text{err}_\mu(L(\bar{s})) < \epsilon$ implikuje, že chybová množina pre $L(\bar{s})$ je prázdna, pretože neexistujú žiadne príklady, ktoré majú pravdepodobnosť menšiu než ϵ . Teda podmienka implikuje, že $L(\bar{s}) = t$, t.j. výstupná hypotéza je presne rovná cieľovému konceptu t .

Záverom týchto úvah je, že učenie na konečnom priestore je presné PEC- „probably exactly correct“. Ale je v tom háčik. Jednoduchá vzorka dĺžky m_0 v definícii PAC-učenia závisí od parametrov δ, ϵ ale nezávisí od μ (a t).

Argument uvedený vyššie obsahuje výber ϵ pomocou ϵ_μ , a tak hodnota m_0 vyžadovaná pre exaktné učenie bude závisieť od δ a μ . Toto je v spore s našim originálnym cieľom dokazovania výsledkov, ktoré nie sú závislé od rozdelenia μ (možno neznáme rozdelenie príkladov).

Príklad 4.2 Štandardný učiaci algoritmus pre monočleny na $\{0, 1\}^n$ pre pevné n . Uvedieme "PEC" vlastnosť. Kľúčovým zdelením tu je, že algoritmus poskytuje správne hypotézy, poskytované všetkými kladnými príkladmi, ktoré boli zahrnuté do tréningovej vzorky. Dĺžka tréningovej vzorky rastie a rastie tiež pravdepodobnosť, že vzorka obsahuje všetky kladné príklady; dôsledkom toho pravdepodobnosť spôsobí, že výstup je korektný. Presnejšie, nech ϵ_μ bude najmenšia hodnota $\mu(x)$, ktorú uvažujeme nad množinou príkladov x z $\{0, 1\}^n$ s nenulovou pravdepodobnosťou. Potom pravdepodobnosť, že tréningová vzorka dĺžky m neobsahuje daný príklad je najviac $(1 - \epsilon_\mu)^m$. Pravdepodobnosť, že existuje jeden z danej množiny p príkladov, ktoré nie sú v tréningovej vzorke je preto $p \cdot (1 - \epsilon_\mu)^m$. Ak X^+ je množina kladných príkladov pre daný cieľový koncept t , pravdepodobnosť, že existuje príklad v X^+ , ktorý nie je vo vzorke je najviac

$$|X^+|(1 - \epsilon_\mu)^m$$

Potrebuje vyjadriť m , teda

$$\begin{aligned} |X^+|(1 - \epsilon_\mu)^m &< \delta \\ m \lg(1 - \epsilon_\mu) + \lg |X^+| &< \lg \delta \\ \lg |X^+| - \lg \delta &< -m \lg(1 - \epsilon_\mu) < m\epsilon_\mu \end{aligned}$$

Použijeme nejaké známe skutočnosti:

$$\begin{aligned} |X^t| \leq |X| \leq 2^n \quad a \quad 1 - \epsilon_\mu &< \exp(-\epsilon_\mu) \\ \log |X^+| \leq n \quad \log(1 - \epsilon_\mu) &< -\epsilon_\mu \\ m &\geq \left\lceil \frac{n}{\epsilon_\mu} \ln 2 + \frac{1}{\epsilon_\mu} \ln \frac{1}{\delta} \right\rceil \end{aligned}$$

Poznamenajme, že dĺžka vzorky je nezávislá od t , ale závisí od rozloženia cez parameter ϵ_μ .

Poznámka: Niekoľko variantov PAC - učenia bolo dosiahnutých tak, že bolo umožnené, že učiaci algoritmus a vzorka dostatočnej dĺžky m_0 záviseli nejakým spôsobom buď na rozdelení pravdepodobnosti μ alebo na cieľovom koncepte t . Toto nie je umelé: v mnohých učiacich problémoch je niečo známe - rozdelenie alebo cieľ. Výsledné definície naučiteľnosti sú menej atraktívne než bezkonceptové a bez-rozdelenia PAC definície, ale sú veľmi často ľahko splniteľné. Je možné nájsť veľa publikovaných prác skúmajúcich takého "neuniformné" PAC učenie.

4.5 Úlohy

1. Napíšte postupnosť hypotéz generovaných algoritmom uvedeným vyššie pre učenie lúčov, ak tréningová vzorka je nasledujúca: (6.1436, 1), (1.5987, 0), (4.2381, 1), (5.7462, 1), (4.3964, 1), (4.2167, 1).
2. Aká je dĺžka tréningovej vzorky pre algoritmus učenia lúčov, aby bola dôveryhodnosť 99.5% a najviac 25% príkladov vybraných podľa nejakého rozdelenia do vzorky bolo zle klasifikovaných?
3. Modifikujte algoritmus pre učenie lúčov tak, že namiesto prázdneho lúča na začiatku vezmeme nejaké veľké číslo. Je tento algoritmus konzistentný?
4. Dané sú dve reálne čísla a, b , a je menšie alebo rovné b . Intervalový koncept $C_{a,b}$ je definovaný takto:

$C_{a,b}(y) = 1$, ak y je v intervale $[a, b]$, inak je rovné 0.

Nech H je hypotézový priestor všetkých intervalov, do ktorého patrí aj prázdny interval.

Nasleduje učiaci algoritmus pre H :

```
empty:=true;
for i:-1 to m do
if bi=1 then
  if empty then begin a:=xi ; b:=xi;
                    empty:= false;
                    end else begin
                        if xi > b then b:=xi;
                        if xi < a then a:=xi end;
  if empty then L(s):=prazdny interval
  else          L(s):=C<a,b>
```

Dokážte, že L je PAC s vhodne zvolenou dĺžkou vzorky. Akou?

5. Modifikujte algoritmus z predchádzajúceho cvičenia pre priestor intervalových konceptov, ktorý nepoužíva funkciu identicky rovnú nule.

Kapitola 5

Lineárna regresia, logistická regresia

Dôležitým problémom v strojovom učení, ktorý má široké použitie, je odvodenie funkcionálnych vzťahov (učenie sa) medzi atribútovými premennými a s nimi asociovanou odozvou alebo cieľovými premennými tak, že je možné predpovedať odozvu pre ľubovoľnú množinu hodnôt atribútových premenných. Napríklad chceme vytvoriť „model“, ktorý bude predpovedať príchod geomagnetickej búrky o niekoľko hodín.¹ Na to je potrebné mať k dispozícii namerané údaje ďalších veličín, od ktorých je geomagnetická búrka závislá a poznať ich vlastnosti. Iným príkladom by mohla byť predikcia transakcií na burze.

5.1 Lineárne modelovanie

Najjednoduchším príkladom učiacich sa problémov je lineárne modelovanie, t. j. učenie sa lineárnych vzťahov medzi atribútmi a odozvami. V tabuľke 5.1² sú rekordné výsledky v skoku do diaľky u mužov. Na obrázku 5.1² sú tieto výsledky² znázornené graficky v paneli A. Rekordy tvoria rastúcu funkciu. Náš cieľ je použiť tieto údaje na učenie pre model funkcionálnej závislosti (ak existuje) medzi rokmi rekordov a dosiahnutými výsledkami. Je nám jasné, že rok nie je jediným faktorom, ktorý ovplyvňuje rekord. Keď by sme chceli predpovedať seriózne, budeme musieť brať do úvahy aj iné faktory, napríklad formu súťažiacich. Avšak na prvý kontakt s lineárnym modelom nám tieto údaje postačia.

Nakoniec by sme chceli vedieť, v ktorom roku by sme mohli očakávať ďalší rekord.

¹Modelom v strojovom učení bude vlastne program, ktorý rieši daný problém.

²wikipedia: https://sk.wikipedia.org/wiki/Skok_do_diaľky

Tab. 5.1: Muži, skok do diaľky

Výkon (m)	Atlét	Miesto	Dátum
7,61	Sp. kráľovstvo Peter O'Connor	Dublin	5. 8. 1901
7,69	USA Edwin Gourdin	Cambridge	23. 7. 1923
7,76	USA Robert LeGendre	Paríž	7. 7. 1924
7,89	USA William DeHart Hubbard	Chicago	13. 6. 1925
7,90	USA Edward Hamm	Cambridge	7. 7. 1928
7,93	Haiti Sylvio Cator	Paríž	9. 9. 1928
7,98	Japonsko Chuhei Nambu	Tokio	27. 10. 1931
8,13	USA Jesse Owens	Ann Arbor	25. 5. 1935
8,21	USA Ralph Boston	Walnut	12. 8. 1960
8,24	USA Ralph Boston	Modesto	27. 5. 1961
8,28	USA Ralph Boston	Moskva	16. 7. 1961
8,31	ZSSR Igor Ter-Ovanesyan	Jerevan	10. 6. 1962
8,31	USA Ralph Boston	Kingston	15. 8. 1964
8,34	USA Ralph Boston	Los Angeles	12. 9. 1964
8,35	USA Ralph Boston	Modesto	29. 5. 1965
8,35	ZSSR Igor Ter-Ovanesyan	Mexiko	19. 10. 1967
8,90	USA Bob Beamon	Mexiko	18. 10. 1968
8,95	USA Mike Powell	Tokio	30. 8. 1991

5.1.1 Definícia modelu

Začneme definíciou nášho modelu ako funkcie, ktorá mapuje naše *vstupné* atribúty, v tomto prípade roky rekordov v skoku do diaľky na *výstupné* alebo cieľové hodnoty - dĺžky skokov.

Je mnoho funkcií, ktoré by sme mohli použiť na toto mapovanie. Ak označíme x vstupný atribút (rok) a t výstupný atribút (dĺžku skoku), matematicky to budeme zapisovať $t = f(x)$. Okrem vstupných atribútov (premenných) funkcia spravidla obsahuje ďalšie parametre, z ktorých niektoré sú modifikovateľné (napríklad učenie). Parametre učiaceho sa modelu sú ústrednou témou strojového učenia. Budeme používať zápis $t = f(x; a)$ pre vyjadrenie toho, že funkcia pracuje na atribútoch x a má parametre a . Je zrejmé, že aj premenných aj parametrov môže byť viac.

Náš prvý model bude lineárny a bude mať jednu premennú a dva parametre,

$$y = f(x) = kx + q,$$

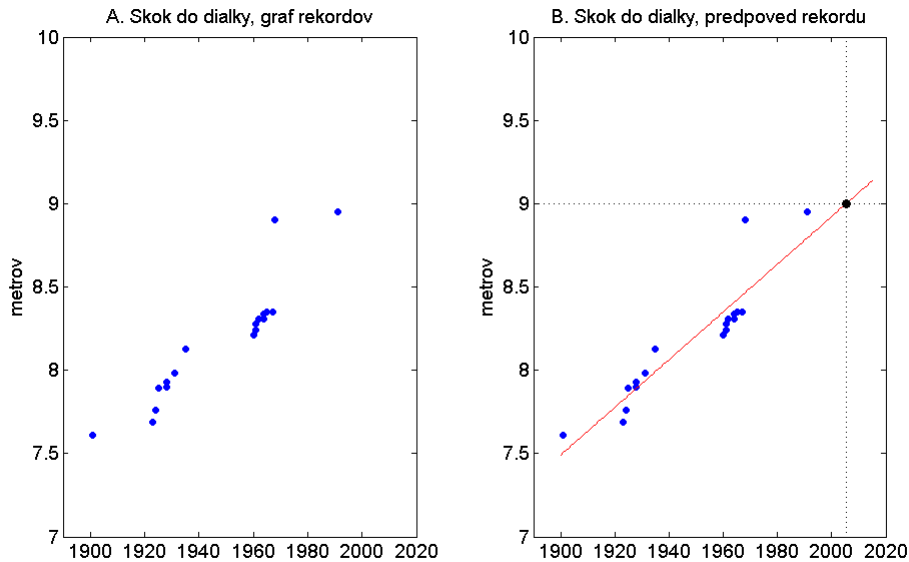
kde k, q sú konštanty a x je premenná a z geometrického hľadiska táto funkcia predstavuje priamku.

Predpoklady pre modelovanie Aby sme vedeli vybrať model, potrebujeme určiť predpoklady charakterizujúce vzťah medzi x a t .

Pre prvý model sme stanovili, že vzťah bude lineárny. Dáta na obrázku [5.1](#) v paneli A by mohli byť modelované pomocou priamky. Teda

$$y = f(x; k, q) = kx + q, \quad (5.1)$$

Prvý lineárny model



Obr. 5.1: Graf rekordov skoku do diaľky mužov a predpoveď ďalšieho rekordu

K tomu, aby sme pre k a q vybrali najlepšie hodnoty, potrebujeme definovať, čo to znamená *najlepšie*. Naším cieľom bude, aby model pre všetky vstupné údaje dával výstupy, ktoré sa budú od skutočných výstupných hodnôt líšiť čo najmenej (chyba modelu bude čo najmenšia). Predpokladajme, že sú dané dvojice vstupných a výstupných hodnôt $(x_1, t_1), (x_2, t_2), \dots, (x_N, t_N)$.

Ohodnotenie modelu

Chyba modelu pre jednu dvojicu je

$$E_n(t_n; f(x_n; k, q)) = (t_n - f(x_n; k, q))^2 \quad (5.2)$$

Pre všetky dvojice budeme používať kvadratickú chybovú funkciu

$$E = \frac{1}{N} \sum_{n=1}^N E_n(t_n; f(x_n; k, q)) \quad (5.3)$$

Najlepším modelom bude taký, pre ktorý bude chyba 5.3 najmenšia. Matematicky to zapíšeme

$$\arg \min_{k,q} \frac{1}{N} \sum_{n=1}^N E_n(t_n; f(x_n; k, q)), \quad (5.4)$$

čo znamená *nájsť argumenty, ktoré minimalizujú súčet chýb*.

5.1.2 Odvodenie vzťahov pre výpočet k a q

Pre dané dvojice vstupných a výstupných hodnôt $(x_1, t_1), (x_2, t_2), \dots, (x_N, t_N)$ máme chybu

$$\begin{aligned} E &= \frac{1}{N} \sum_{n=1}^N E_n(t_n; f(x_n; k, q)) = \frac{1}{N} \sum_{n=1}^N (t_n - (kx_n + q))^2 \\ &= \frac{1}{N} \sum_{n=1}^N (t_n^2 - 2t_n(kx_n + q) + k^2x_n^2 + 2kqx_n + q^2) \\ &= \frac{1}{N} \sum_{n=1}^N (k^2x_n^2 + 2kx_n(q - t_n) + q^2 - 2t_nq + t_n^2) \end{aligned} \quad (5.5)$$

K výpočtu minima vypočítame najprv parciálne derivácie podľa k a q . Parciálna derivácia podľa k je

$$\frac{\partial E}{\partial k} = \frac{2k}{N} \sum_{n=1}^N x_n^2 + \frac{2}{N} \sum_{n=1}^N (x_n(q - t_n)) \quad (5.6)$$

Parciálna derivácia podľa q je

$$\frac{\partial E}{\partial q} = \frac{2k}{N} \sum_{n=1}^N x_n + 2q - \frac{2}{N} \sum_{n=1}^N t_n \quad (5.7)$$

Označme \bar{x} priemer hodnôt x_1, \dots, x_n a \bar{t} priemer hodnôt t_1, \dots, t_n , t. j.

$$\bar{x} = \frac{1}{N} \sum_{n=1}^N x_n, \quad \bar{t} = \frac{1}{N} \sum_{n=1}^N t_n$$

Položíme (5.7) rovné 0

$$\frac{2k}{N} \sum_{n=1}^N x_n + 2q - \frac{2}{N} \sum_{n=1}^N t_n = 2k\bar{x} + 2q - 2\bar{t} = 0$$

teda $q = \bar{t} - k\bar{x}$.

Označme $\overline{x^2}$ priemer hodnôt x_1^2, \dots, x_n^2 a \overline{xt} priemer hodnôt x_1t_1, \dots, x_nt_n , t. j.

$$\overline{x^2} = \frac{1}{N} \sum_{n=1}^N x_n^2, \quad \overline{xt} = \frac{1}{N} \sum_{n=1}^N x_nt_n$$

Položíme (5.6) rovné 0

$$\frac{2k}{N} \sum_{n=1}^N x_n^2 + \frac{2}{N} \sum_{n=1}^N (x_n(q - t_n)) = 2k\overline{x^2} + 2q\bar{x} - 2\overline{xt} = 0$$

Dosadením $q = \bar{t} - k\bar{x}$ dostaneme

$$2k\bar{x}^2 + 2(\bar{t} - k\bar{x})\bar{x} - 2\bar{x}\bar{t} = 0 \quad (5.8)$$

Úpravou dostaneme výrazy pre k a q

$$k = \frac{\bar{x}\bar{t} - \bar{x}\bar{t}}{\bar{x}^2 - \bar{x}^2}, \quad q = \bar{t} - k\bar{x} \quad (5.9)$$

Výpočtom druhej derivácie sa presvedčíme, že ide o minimum.

Pre náš príklad dostávame výsledky uvedené v tabuľke [5.1.2](#).

Tab. 5.2: Muži, skok do diaľky a vypočítané hodnoty

Por. č	Rok	Skok	Rok*Skok	Rok*Rok
1	1901	7,61	14466,61	3613801
2	1923	7,69	14787,87	3697929
3	1924	7,76	14930,24	3701776
4	1925	7,89	15188,25	3705625
5	1928	7,90	15231,20	3717184
6	1928	7,93	15289,04	3717184
7	1931	7,98	15409,38	3728761
8	1935	8,13	15731,55	3744225
9	1960	8,21	16091,60	3841600
10	1961	8,24	16158,64	3845521
11	1961	8,28	16237,08	3845521
12	1962	8,31	16304,22	3849444
13	1964	8,31	16320,84	3857296
14	1964	8,34	16379,76	3857296
15	1965	8,35	16407,75	3861225
16	1967	8,35	16424,45	3869089
17	1968	8,90	17515,20	3873024
18	1991	8,95	17819,45	3964081
$(1/N) \sum_{n=1}^N$	1947,7	8,1739	15927,40	3793921

Hodnoty pre k , q a E sú

$$k = 0.0143; \quad q = -19.7139; \quad E = 0.0193$$

V paneli B na obrázku [5.1](#) sa nachádza vypočítaná priamka aj nakreslená a vidíme, aký trend mali rekordy v skoku do diaľky.

5.1.3 Predikcia rekordu

Trendovú priamku môžeme použiť na predikciu, aký rekord by sme mohli očakávať v niektorom z nasledujúcich rokov. Očakávaný rekord v roku 2018 je $y = 0.0143 * 2018 - 19.7139 = 9,14 m.$ a v roku 2019 je $y = 0.0143 * 2019 - 19.7139 = 9,16 m.$

Treba tu poznamenať, že tieto presné predikcie sú len orientačné. Toto sa prejavuje aj v predchádzajúci rokoch.

5.2 Zovšeobecnené lineárne modely

Doteraz sme uvažovali lineárne modely, lineárne v parametroch aj v premenných. Ak zavedieme indexované označenie parametrov $\mathbf{w} = (w_1, w_0) = (k, q)$, tak (5.1) môžeme zapísať

$$y = f(x; k, q) = kx + q = w_1x + w_0 = f(x; \mathbf{w}), \quad (5.10)$$

Model je stále lineárny v parametroch, ale kvadratický v dátach.

Všeobecnejší model je polynomiálny

$$f(x; \mathbf{w}) = \sum_{k=0}^K w_k x^k, x^0 = 1; \quad (5.11)$$

Matica dát má tvar

$$\mathbf{X} = \begin{pmatrix} x_1^0 & x_1^1 & x_1^2 & \dots & x_1^K \\ x_2^0 & x_2^1 & x_2^2 & \dots & x_2^K \\ \vdots & \vdots & \vdots & \vdots & \vdots \\ x_N^0 & x_N^1 & x_N^2 & \dots & x_N^K \end{pmatrix}$$

Nie je nutné používať len polynomiálne funkcie, je možné použiť K iných funkcií, $h_k(x)$.

Matica dát má potom tvar

$$\mathbf{X} = \begin{pmatrix} h_1(x_1) & h_2(x_1) & h_3(x_1) & \dots & h_K(x_1) \\ h_1(x_2) & h_2(x_2) & h_3(x_2) & \dots & h_K(x_2) \\ \vdots & \vdots & \vdots & \vdots & \vdots \\ h_1(x_N) & h_2(x_N) & h_3(x_N) & \dots & h_K(x_N) \end{pmatrix}$$

5.3 Úlohy

1. Analyzujte skoky do diaľky u žien.
2. Na predchádzajúcich dátach vyskúšajte použitie parametricky lineárneho modelu s funkciami:

$$h_1(x) = 1; \quad h_2(x) = x; \quad h_3(x) = \sin\left(\frac{x-a}{b}\right)$$

Zvoľte rôzne a a b .

Kapitola 6

Lineárne modelovanie, maximálna dôveryhodnosť

Lineárna regresia je jedným z najzákladnejších algoritmov strojového učenia a jednoducho sa snaží analyzovať naše údaje pomocou preloženej priamky s najmenšou chybou.

V tejto kapitole uvidíme lineárnu regresiu z inej perspektívy, a síce z hľadiska pravdepodobnosti, najmä pomocou odhadu maximálnej dôveryhodnosti - Maximum Likelihood (ML). Budeme explicitne modelovať **šum** (chyby medzi modelom a pozorovaniami) v dátach pomocou **náhodnej premennej**.

6.1 Chyby ako šum

Pri vytváraní lineárneho modelu sme predpokladali lineárny vzťah medzi nezávislými premennými a závislou premennou. Tento model zachytáva všeobecný trend v dátach, ale ignoruje to, že v niektorých bodoch sa dopúšťa veľkej chyby. Z hľadiska modelu je ťažké tieto chyby obhájiť. Vieme, že chyby v dátach budú a mali by sme ich pri tvorbe modelu zohľadniť. Je možné s modelom experimentovať, modifikovať jeho parametre, ale je tu ešte iný prístup. Budeme brať náš modelovací problém ako **generatívny**, to znamená, že môžeme vytvárať model, ktorý bude použitý na vytvorenie (generovanie) datasetu, ktorý vyzerá podobne ako náš skutočný dataset.

Poznamenajme si niektoré vlastnosti chýb (majme na mysli náš príklad rekordov):

- Chyby sú rôzne v každom roku (bode x). Niekdy sú pozitívne, niekedy negatívne a majú rôzne hodnoty.
- Nie je žiadny vzťah medzi veľkosťou chýb v rokoch (bodoch x). Chyby nie sú funkciou roka (x).

Použijeme náhodnú premennú ϵ_n na vyjadrenie chyby v roku \mathbf{x}_n a náš model teraz bude vyzerat takto:

$$t_n = \mathbf{w}^T \mathbf{x}_n + \epsilon_n \quad (6.1)$$

Na doplnenie definície modelu potrebujeme rozhodnúť o štatistickom rozdelení ϵ_n . Najprv si vyjasníme, že rozdiel medzi modelom a dĺžkami skokov je spojitá veličina. Preto ϵ_n je spojitá náhodná premenná. Tiež si uvedomíme, že nemáme len jednu náhodnú premennú, ale máme pre každý rok jednu. Môžeme predpokladať, že tieto hodnoty sú nezávislé preto pre pravdepodobnosť p všetkých n skokov platí

$$p(\epsilon_1, \dots, \epsilon_n) = \prod_{n=1}^N p_n(\epsilon_n) \quad (6.2)$$

Ďalší predpoklad sa týka tvaru rozdelenia $p_n(\epsilon_n)$. Budeme predpokladať, že pre všetky n je to Gaussovo (normálne) rozdelenie s priemerom 0 a varianciou σ^2 . Vieme, že tak ϵ_n môže nadobúdať kladné aj záporné hodnoty a má zaujímavé modelovacie vlastnosti.

Použitím normálneho rozdelenia pre ϵ , t.j.

$$p(y|\mu, \sigma^2) = \frac{1}{\sigma\sqrt{2\pi}} \exp\left(-\frac{1}{2\sigma^2}(y - \mu)^2\right) \quad (6.3)$$

s $\mu = 0$ a $\sigma^2 = 0.05$ dostaneme realistickejší dataset.

Náš model teraz pozostáva z dvoch komponentov:

Trend a šum

- deterministický komponent ($\mathbf{w}^T \mathbf{x}_n$), ktorý sa nazýva **trend** alebo **drift**,
- náhodný komponent (ϵ_n), ktorý sa nazýva **šum**.

Je pravda, že šum môže mať iné štatistické rozdelenie než Gaussovo a tiež nemusí byť aditívny, ale multiplikatívny ($t_n = \mathbf{w}^T \mathbf{x}_n \epsilon_n$), ktorý sa často využíva pri zašumení obrázkov.

6.2 Dôveryhodnosť (hodnovernosť)

Náš model je teraz v nasledujúcom tvare

$$t_n = f(\mathbf{x}_n; \mathbf{w}) + \epsilon_n, \epsilon_n \sim \mathcal{N}(0, \sigma^2), \quad (6.4)$$

kde $\mathcal{N}(0, \sigma^2)$ predstavuje Gaussovo (normálne) rozdelenie s priemerom μ a varianciou σ^2 .

Potrebujeme nájsť optimálnu hodnotu pre \mathbf{w} , ktorú označíme $\hat{\mathbf{w}}$. Tiež tu máme ďalší parameter σ^2 , ktorý je potrebné nastaviť. V regresnom

modeli chyba bola meraná pomocou rozdielu pozorovanej hodnoty a hodnoty, ktorú predikoval model.

Efekt pridania náhodnej premennej do modelu je, že výstup modelu, t , je teraz sám náhodná premenná. Inak povedané, neexistuje len jedna hodnota t_n pre partikulárne \mathbf{x}_n . Teda nie je možné použiť chybu ako prostriedok pre optimalizáciu \mathbf{w} a σ^2 .

Pridanie konštanty ($\mathbf{w}^T \mathbf{x}_n$) do Gaussovej náhodnej premennej je ekvivalentné inej Gaussovej náhodnej premennej s priemerom posunutým o uvedenú konštantu:

$$\begin{aligned} y &= a + z \\ p(z) &= \mathcal{N}(m, s) \\ p(y) &= \mathcal{N}(m + a, s) \end{aligned}$$

Preto náhodná premenná t_n má nasledujúcu funkciu hustoty

$$p(t_n | \mathbf{x}_n, \mathbf{w}, \sigma^2) = \mathcal{N}(\mathbf{w}^T \mathbf{x}_n, \sigma^2) \quad (6.5)$$

t_n teda závisí od partikulárnych hodnôt \mathbf{x}_n a \mathbf{w} (určujú priemer) a od σ^2 (určuje varianciu).

Aby sme videli, ako môžeme toto použiť na hľadanie optimálnej hodnoty \mathbf{w} a σ^2 , vyberme jeden rok z datasetu - 1980. Na základe modelu predchádzajúceho (základného) lineárneho modelu majme určené (w_0, w_1) a predpokladajme, že $\sigma^2 = 0.05$. Na obrázku 6.2 je znázornená $p(t_n | x_n = 1980, \mathbf{w}, \sigma^2)$ ako funkcia t_n . Pre tento rok dostaneme

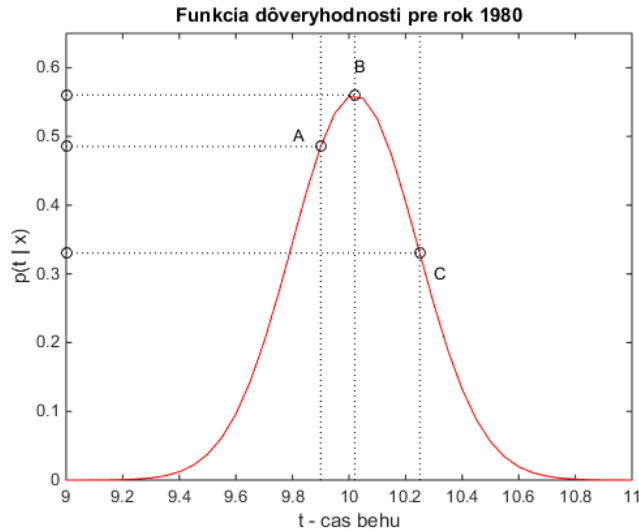
$$p(t_n | \mathbf{x}_n = [1, 1980]^T, \mathbf{w} = [36.416, -0.0133]^T, \sigma^2 = 0.05), \quad (6.6)$$

pri Gaussovom rozdelení $\mu = 36.416, -0.0133 * 1980 = 10.02$ a $\sigma^2 = 0.05$.

Pripomeňme, že pre spojitú premennú $t, p(t)$ nemôže byť interpretovaná ako pravdepodobnosť. Výška krivky v partikulárnej hodnote pre t môže byť interpretovaná takto: nakoľko je hodnoverné (dôveryhodné), že budeme pozorovať túto partikulárnu hodnotu t v danom roku $x = 1980$. Najhodnovernejšia hodnota v roku 1980 by bola 10.02 sekúnd (pre Gaussovo rozdelenie najhodnovernejší - najvyšší bod zodpovedá priemeru). Na obrázku máme 3 príklady časov - A, B, C. Z nich B je najhodnovernejší, a C najmenej hodnoverný.

Skutočný nameraný čas v roku 1980 je C (10.25 sekúnd). Hustota pravdepodobnosti $p(t_n | \mathbf{x}_n, \mathbf{w}, \sigma^2)$ vyhodnotená v $t_n = 10.25$ je dôležitá veličina, ktorá je známa ako **dôveryhodnosť (hodnovernosť)** n -tého dátového bodu. Nemôžeme zmeniť $t_n = 10.25$ (máme to v dátach), ale môžeme zmeniť \mathbf{w} a σ^2 a skúšať, posúvať a hľadať čo najväčšiu dôveryhodnosť.

Idea hľadania parametrov, ktoré maximalizujú dôveryhodnosť týmto spôsobom, je kľúčový koncept v strojovom učení.



Obr. 6.1: Graf pre 1980, x-ová os t , y-ová os $p(t|x)$

6.3 Hodnovernosť datasetu

V skutočnosti sa sa nezaujíname o dôveryhodnosť jedného dátového bodu, ale o dôveryhodnosť celého datasetu. Označme $\mathbf{X} = [\mathbf{x}_1, \dots, \mathbf{x}_N]$, $\mathbf{t} = [t_1, \dots, t_N]$. Keď máme N bodov, zaujíma nás podmienená pravdepodobnosť

$$p(t_1, \dots, t_N | \mathbf{x}_1, \dots, \mathbf{x}_N, \mathbf{w}, \sigma^2) = p(\mathbf{t} | \mathbf{X}, \mathbf{w}, \sigma^2).$$

Vyhodnotenie tejto hustoty v pozorovaných dátových bodoch dáva jednu hodnotu dôveryhodnosti pre celý dataset, ktorú môžeme optimalizovať zmenou \mathbf{w} a σ^2 .

Predpoklad, že šum v každom dátovom bode je nezávislá veličina, znamená $p(\epsilon_1, \dots, \epsilon_N) = \prod_n p(\epsilon_n)$ a umožňuje nám použiť zaujímavé odvodenia. Táto spojená (joint) podmienená hustota môže byť rozdelená do N separálnych výrazov, jeden pre každý dátový objekt:

$$L = p(\mathbf{t} | \mathbf{X}, \mathbf{w}, \sigma^2) = \prod_{n=1}^N p(t_n | \mathbf{x}_n, \mathbf{w}, \sigma^2) = \prod_{n=1}^N \mathcal{N}(\mathbf{w}^T \mathbf{x}_n, \sigma^2) \quad (6.7)$$

Nepovedali sme, že t_n hodnoty sú navzájom úplne nezávislé. Sú v priemere v čase klesajúce, a teda ponúkajú jasnú štatistickú závislosť. Ak by boli úplne nezávislé, nemohli by sme vytvárať model. V skutočnosti sú **podmienene nezávislé** – vzhľadom k daným hodnotám \mathbf{w} (deterministická časť modelu) t_n sú nezávislé (model ponúka výpočet v ľubovoľnom bode); bez nich nie sú celkom nezávislé.

Predstavme si, že nám chýbajú dáta z jedného roku, napríklad 1960. Ostatné roky sú v poriadku. \mathbf{X} , \mathbf{t} neobsahujú údaje o roku 1960. Z daných dát by sme chceli vedieť niečo o roku 1960. Teda nás zaujíma

$$p(t_{1960} | \mathbf{x}_{1960}, \mathbf{X}, \mathbf{t})$$

Z definície podmienenej pravdepodobnosti dostaneme

$$p(t_{1960}|\mathbf{x}_{1960}, \mathbf{X}, \mathbf{t}) = \frac{p(t_{1960}, \mathbf{t}|\mathbf{x}_{1960}, \mathbf{X})}{p(\mathbf{t}|\mathbf{X})}$$

Za prepokladu, že \mathbf{t} sú nezávislé, výsledky v t_{1960} závisia len od \mathbf{x}_{1960} :

$$p(t_{1960}|\mathbf{x}_{1960}, \mathbf{X}, \mathbf{t}) = \frac{p(t_{1960}|\mathbf{x}_{1960})\prod_n p(t_n|\mathbf{x}_n)}{\prod_n p(t_n|\mathbf{x}_n)} = p(t_{1960}|\mathbf{x}_{1960})$$

Avšak pre náš model toto nie je použiteľné, t_{1960} musí byť v nejakom zmysle závislé od ostatných dát. Táto závislosť je zapúzdrená v parametroch \mathbf{w} . Ak poznáme \mathbf{w} , tak to, čo zostáva neznáme, je chyba medzi pozorovanými dátami a $\mathbf{w}^T \mathbf{x}_n$. Predpokladáme, že tieto chyby sú nezávislé. Bez modelu pozorovania sú nezávislé, ale keď berieme do úvahy model (\mathbf{w}), pretože je tu chyba, pozorovania nie sú nezávislé.

6.4 Maximálna dôveryhodnosť

Budeme hľadať parametre modelu \mathbf{w} a σ^2 také, ktoré urobia náš model najdôveryhodnejší. Z analytických dôvodov budeme maximalizovať prirodzený algoritmus dôveryhodnosti L (logaritmus je rastúca funkcia a teda je to možné urobiť).

$$\begin{aligned} \log L &= \sum_{n=1}^N \log \left(\frac{1}{\sqrt{2\pi\sigma^2}} \exp\left\{-\frac{1}{2\sigma^2}(t_n - f(\mathbf{x}_n; \mathbf{w}))^2\right\} \right) \\ &= \sum_{n=1}^N \left(-\frac{1}{2} \log(2\pi) - \log \sigma - \frac{1}{2\sigma^2}(t_n - f(\mathbf{x}_n; \mathbf{w}))^2 \right) \\ &= -\frac{N}{2} \log 2\pi - N \log \sigma - \frac{1}{2\sigma^2} \sum_{n=1}^N (t_n - f(\mathbf{x}_n; \mathbf{w}))^2 \end{aligned}$$

Substitúciou $f(\mathbf{x}_n; \mathbf{w}) = \mathbf{w}^T \mathbf{x}_n$ dostaneme

$$L = -\frac{N}{2} \log 2\pi - N \log \sigma - \frac{1}{2\sigma^2} \sum_{n=1}^N (t_n - \mathbf{w}^T \mathbf{x}_n)^2 \quad (6.8)$$

Parciálne derivácie podľa \mathbf{w}

$$\begin{aligned} \frac{\partial \log L}{\partial \mathbf{w}} &= \frac{1}{\sigma^2} \sum_{n=1}^N \mathbf{x}_n (t_n - \mathbf{w}^T \mathbf{x}_n) \\ &= \frac{1}{\sigma^2} \sum_{n=1}^N (\mathbf{x}_n t_n - \mathbf{x}_n \mathbf{x}_n^T \mathbf{w}) = 0 \end{aligned}$$

Označme

$$\mathbf{X} = [\mathbf{x}_1^T \mathbf{x}_2^T \dots \mathbf{x}_N^T]^T, \mathbf{t} = [t_1 t_2 \dots t_N]^T \quad (6.9)$$

$$\mathbf{X} = \begin{pmatrix} 1 & x_1 \\ 1 & x_2 \\ \vdots & \vdots \\ 1 & x_N \end{pmatrix}$$

Maximálna dôveryhodnosť z hľadiska \mathbf{w} je v maticovom tvare

$$\hat{\mathbf{w}} = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{t} \quad (6.10)$$

Pre σ^2 dostávame

$$\frac{\partial L}{\partial \sigma} = -\frac{N}{\sigma} + \frac{1}{\sigma^3} \sum_{n=1}^N (t_n - \mathbf{x}_n^T \hat{\mathbf{w}})^2 = 0 \quad (6.11)$$

$$\hat{\sigma}^2 = \frac{1}{N} \sum_{n=1}^N (t_n - \mathbf{x}_n^T \hat{\mathbf{w}})^2 \quad (6.12)$$

Tento výraz je perfektný - variancia je jednoducho priemer kvadratickej chyby. V maticovom označení

$$\hat{\sigma}^2 = \frac{1}{N} (\mathbf{t} - \mathbf{X} \hat{\mathbf{w}})^T (\mathbf{t} - \mathbf{X} \hat{\mathbf{w}}) = \frac{1}{N} (\mathbf{t}^T \mathbf{t} - 2\mathbf{t}^T \mathbf{X} \hat{\mathbf{w}} + \hat{\mathbf{w}}^T \mathbf{X}^T \mathbf{X} \hat{\mathbf{w}}) \quad (6.13)$$

Toto zjednodušíme substitúciou (6.10), navyše platí $\hat{\mathbf{w}}^T = \mathbf{t}^T \mathbf{X} (\mathbf{X}^T \mathbf{X})^{-1}$ a $(\mathbf{X}^T \mathbf{X})^{-1}$ je symetrická, a preto sa rovná svojej transponovanej matici.

Pre rekordy v behu dostávame

$$\hat{\mathbf{w}} = [36.4165, -0.0133]^T, \sigma^2 = 0.0503$$

\mathbf{w} je to isté ako pri pri lineárnej regresii, σ^2 predstavuje varianciu Gaussovho šumu, ktorý sme predpokladali v dátach.

6.5 Charakteristika max. dôveryhodnosti

Aby sme dokázali, že je to maximum, mali by sme vypočítať druhú deriváciu vo vypočítanom bode. Hodnota $\log L$ v maxime je

$$\begin{aligned} \log L &= -\frac{N}{2} \log 2\pi - N \log \hat{\sigma}^2 - \frac{1}{2\hat{\sigma}^2} N \hat{\sigma}^2 \\ &= -\frac{N}{2} (1 + \log 2\pi) - \frac{N}{2} \log \hat{\sigma}^2 \end{aligned} \quad (6.14)$$

L bude rásť, keď $\hat{\sigma}^2$ bude klesať. Toto ale ovplyvní \mathbf{w} , teda deterministický model (výber komplikovanejšieho modelu).

Vzťah medzi zovšeobecnením a pretrénovaním sa niekedy popisuje ako bias – variancia tradeoff.

Ak by sme mali prístup k rozdeleniu, podľa ktorého tréningová vzorka vznikla, $p(x|t)$, mohli by sme vypočítať očakávanú hodnotu kvadratickej chyby medzi odhadnutými a skutočnými hodnotami.

Chceme, aby táto hodnota $\bar{\mathcal{M}}$ bola čo najmenšia. Môže byť dekomponovaná do dvoch výrazov **bias**, \mathcal{B} , a **variancia**, \mathcal{V} :

$$\bar{\mathcal{M}} = \mathcal{B}^2 + \mathcal{V}$$

Bias opisuje systematické chyby medzi naším modelom a generovanými dátami.

Ak je model príliš jednoduchý, bias je veľký (podtrénovanie). Ak vezmeme zložitejší model, bias vieme znížiť. Avšak zložitejšie modely majú vysokú varianciu.

Dôležité je nájsť správne vyváženie oboch.

6.6 Úlohy

1. Predpokladajme, že dataset N binárnych hodnôt x_1, \dots, x_N bol generovaný Bernoulliho rozdelením. Vypočítajte odhad maximálnej dôveryhodnosti pre Bernoulliho parameter.

$$\text{Bernoulliho rozdelenie: } P(X = x) = q^x(1 - x)^{1-x}$$

2. Analyzujte skoky do diaľky u žien použitím šumu. Vygenerujte podobnú tréningovú vzorku.

Kapitola 7

Klasifikácia

7.1 Úvod

Dve najčastejšie objavujúce sa úlohy v strojovom učení sú klasifikácia a predikcia. Hlavný rozdiel spočíva v type cieľového atribútu, ktorého hodnotu predpovedáme. Klasifikácia je skupina metód strojového učenia, ktoré slúžia na predpovedanie hodnôt nominálnych atribútov (tried), napr. pridelenie alebo nepridelenie úveru.

Predikcia je skupina metód strojového učenia, ktoré slúžia na predpovedanie hodnôt numerických atribútov, napr. predikcia spotreby vody na určitom území.

Niektorí autori používajú namiesto pojmu predikcia pojem regresia a následne spoločným pojmom predikcia označujú klasifikáciu a regresiu. My budeme používať pôvodné označenie.

Klasifikácia – predpovedanie hodnôt nominálnych atribútov (tried).

1. fáza: vybudovať (naučiť sa) model správania dát na základe tréningovej množiny – obsahuje záznamy o objektoch, ktorých je hodnota cieľového atribútu už známa.
2. fáza: zostavený model sa použije na predikciu hodnoty cieľového atribútu u nových objektov (na základe hodnôt ostatných atribútov, ktoré sú o objekte k dispozícii).

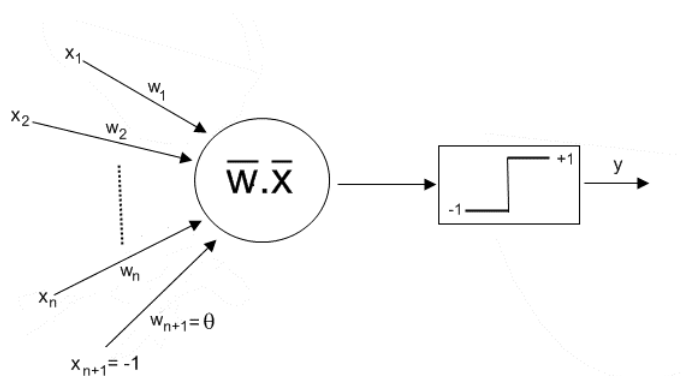
Príklad cieľových atribútov: pohlavie, áno/nie, typ školy a podobne.

Príklad: pridelenie úveru – zaradenie nových žiadateľov do kategórie prideliť/neprideliť na základe podobných údajov minulých žiadateľov o úver.

7.2 Klasifikácia do dvoch tried

Perceptrón

Perceptrón je základný prvok umelej neurónovej siete rovnako ako je neurón základným prvkom organických neurónových sietí. Existuje viac modelov perceptrónu, ktoré sa navzájom líšia svojimi funkciami a tiež tým, do akej miery sa blížia ku funkcii skutočného neurónu. V tejto kapitole sa budeme zaoberať modelom perceptrónu s bipolárnou binárnou aktivačnou funkciou, ktorý navrhol americký psychológ Frank Rosenblatt v roku 1958. Model je znázornený na obrázku 7.2. Na vstupe perceptrónu je vstupný vektor $\mathbf{x} = (x_1, x_2, \dots, x_N)$. Nastaviteľnými parametrami sú váhy a prah perceptrónu, t. j. vektor $\mathbf{w} = (w_1, w_2, \dots, w_N)$ a θ .



Obr. 7.1: Model perceptrónu

Funkciu f budeme nazývať *aktivačnou funkciou*.

$$f(x) = \begin{cases} 1 & \text{ak } x \geq 1, \\ -1 & \text{inak.} \end{cases} \quad (7.1)$$

Pre vstupný príklad \mathbf{x} perceptrón vypočíta hodnotu

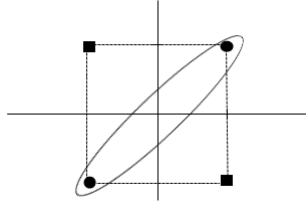
$$y = f(\mathbf{w} \cdot \mathbf{x} - \theta) = f\left(\sum_{n=1}^{N+1} w_n x_n\right) = f\left(\sum_{n=1}^N w_n x_n - \theta\right), \quad (7.2)$$

kde $w_{N+1} = -1, x_{N+1} = \theta$.

Jediné čo perceptrón dokáže, je rozdelenie priestoru vzorov na dva polpriestory (označíme si tieto polpriestory PL1 a PL2), ktoré sú oddelené lineárnou deliacou hranicou - nadrovinou. Deliacá hranica v N -rozmernom priestore je vyjadrená rovnicou (7.3)

$$w_1 x_1 + w_2 x_2 + \dots + w_N x_N - \theta = 0 \quad (7.3)$$

Ak budeme uvažovať 2-rozmerný priestor, t. j. perceptrón s dvoma vstupmi, tak ten nedokáže riešiť situáciu na obrázku 7.2, pretože neexistuje nadroviná (priamka), ktorá oddelí štvorcové body od krúžkových.



Obr. 7.2: Usporiadanie bodov v rovine, ktoré perceptrón nedokáže oddeliť priamkou.

Riešime nasledujúci problém

Formulácia problému

Je daná tréningová vzorka

$$\bar{s} = ((\mathbf{x}_1, y_1), (\mathbf{x}_2, y_2), \dots, (\mathbf{x}_M, y_M)),$$

kde $\mathbf{x}_i = (x_{i,1}, x_{i,2}, \dots, x_{i,N})$ a $y_i \in \{-1, +1\}$. Je potrebné nájsť nastavenie váh a prahu tak, aby klasifikácia do dvoch tried podľa perceptrónu bola konzistentná s tréningovou vzorkou.

Rossenblatt navrhol základný učiaci algoritmus pre učenie perceptrónu, avšak tu budeme hľadať riešenie iným spôsobom.

Pre p -tý príklad, $p = 1, \dots, M$ platí, ak je správne klasifikovaný

$$\mathbf{w} \cdot \mathbf{x}_p - \theta > 0, \quad \text{ak } y_p = +1$$

$$\mathbf{w} \cdot \mathbf{x}_p - \theta < 0, \quad \text{ak } y_p = -1$$

Uvedené nerovnosti je možné zapísať pomocou jednej, ak použijeme $-y_p$, teda

$$-y_p(\mathbf{w} \cdot \mathbf{x}_p - \theta) < 0 \tag{7.4}$$

Výraz $\max\{0, -y_p(\mathbf{w} \cdot \mathbf{x}_p - \theta)\}$ nadobúda hodnotu 0, ak je príklad \mathbf{x}_p správne klasifikovaný a nadobúda kladnú hodnotu, ak je klasifikovaný nesprávne.

Chceme, aby všetky príklady boli klasifikované správne, preto uvažujeme nasledujúcu funkciu

$$f_1(\mathbf{w}, \theta) = \sum_{p=1}^N \max\{0, -y_p(\mathbf{w} \cdot \mathbf{x}_p - \theta)\} \tag{7.5}$$

a tiež chceme, aby chyba klasifikácie bola minimálna, preto riešime problém

$$\min_{\mathbf{w}, \theta} \sum_{p=1}^N \max\{0, -y_p(\mathbf{w} \cdot \mathbf{x}_p - \theta)\}, \tag{7.6}$$

ktorý nájde optimálne parametre pre separujúcu nadrovinu. Avšak sú tu dva technické problémy, s ktorými sa treba vysporiadať pri hľadaní riešenia:

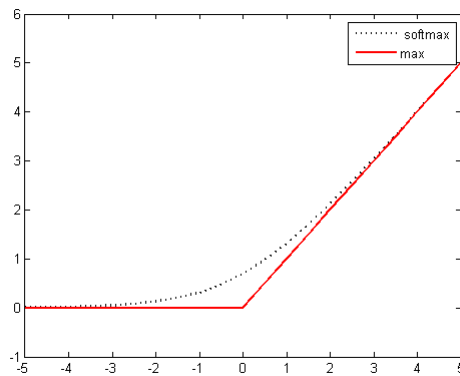
- jedným z optimálnych riešení je $\theta = 0, \mathbf{w} = (0, 0, \dots, 0)$; toto riešenie je potrebné ignorovať vo väčšine prípadov;
- funkcia f_1 je spojitá, ale nie je všade diferencovateľná, preto sa tu nedá použiť metóda poklesu gradientu.

Jedným z používaných riešení je nahradenie nediferencovateľnej funkcie \max vo výraze (7.5), kde sa jedná o maximum z dvoch hodnôt, funkciou **Logistická funkcia**

$$\text{softmax}(e_1, e_2) = \log(e^{e_1} + e^{e_2}); \quad (7.7)$$

V našom prípade

$$\text{softmax}(e_1, e_2) = \text{softmax}(0, -y_p(\mathbf{w} \cdot \mathbf{x}_p - \theta)) = \log(1 + e^{-y_p(\mathbf{w} \cdot \mathbf{x}_p - \theta)}) \quad (7.8)$$



Obr. 7.3: Porovnanie priebehov funkcií \max a softmax .

Funkcia softmax je dobrou aproximáciou funkcie \max a je diferencovateľná. Takže funkciu f_1 nahradíme funkciou f_2 (7.9), ktorú budeme nazývať *softmax cenová funkcia*.

$$f_2(\mathbf{w}, \theta) = \sum_{p=1}^N \log(1 + e^{-y_p(\mathbf{w} \cdot \mathbf{x}_p - \theta)}) \quad (7.9)$$

Táto funkcia už nemá vyššie uvedené dva problémy funkcie f_1 . Formálne sa náš problém redukoval na problém (7.10).

$$\min_{\mathbf{w}, \theta} \sum_{p=1}^N \log(1 + e^{-y_p(\mathbf{w} \cdot \mathbf{x}_p - \theta)}) \quad (7.10)$$

Použijeme gradientovú metódu na hľadanie minima funkcie $f_2(\mathbf{w}, \theta)$. Zaved' me označenie $\tilde{\mathbf{x}}_p = \begin{bmatrix} \mathbf{x}_p \\ -1 \end{bmatrix}$ a $\tilde{\mathbf{w}} = \begin{bmatrix} \mathbf{w} \\ \theta \end{bmatrix}$. Funkcia f_2 sa zmení na

$$f_2(\tilde{\mathbf{w}}) = \sum_{p=1}^N \log(1 + e^{-y_p(\tilde{\mathbf{w}} \cdot \tilde{\mathbf{x}}_p)}) \quad (7.11)$$

Deriváciou funkcie f_2 dostaneme vzťah (7.12).

$$\nabla f_2(\tilde{\mathbf{w}}) = \sum_{p=1}^N \sigma(-y_p \tilde{\mathbf{w}} \cdot \tilde{\mathbf{x}}_p) (1 - \sigma(-y_p \tilde{\mathbf{w}} \cdot \tilde{\mathbf{x}}_p)) \tilde{\mathbf{x}}_p \cdot \tilde{\mathbf{x}}_p \quad (7.12)$$

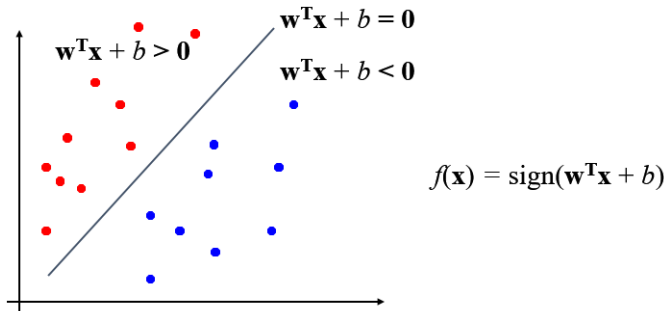
kde $\sigma(-t) = \frac{1}{1+e^t}$ je logistická sigmoidálna funkcia. Riešenie pre nastavenie váh získame, keď položíme $\nabla f_2(\tilde{\mathbf{w}}) = \mathbf{0}$.

7.3 Pásový perceptrón

Nadrovina, ktorá oddeľuje body dvoch tried, je daná rovnicou, ktorú zapíšeme vo vektorovom tvare $\mathbf{w} \cdot \mathbf{x} = 0$. Ak nechceme, aby v blízkosti nadroviny sa nachádzali nejaké body, môžeme stanoviť nasledujúce podmienky

$$\begin{aligned} \mathbf{w} \cdot \mathbf{x}_p - \theta &\geq d, \text{ ak } y_p = +1 \\ \mathbf{w} \cdot \mathbf{x}_p - \theta &\leq -d, \text{ ak } y_p = -1, \end{aligned} \quad (7.13)$$

kde d určuje nadrovinu posunutú do kladného alebo záporného polpriestoru.



Obr. 7.4: Lineárne separovateľné dáta je možné oddeliť nadrovinou.

Analogickým spôsobom ako pri obyčajnom perceptróne vieme určiť funkciu, ktorú budeme minimalizovať (7.14).

$$\min_{\mathbf{w}, \theta} \sum_{p=1}^M \max(0, d - y_p(\mathbf{w} \cdot \mathbf{x}_p - \theta)) \quad (7.14)$$

Aj v tomto prípade je potrebné riešiť to, že funkcia (7.14) je v aspoň jednom bode nediferencovateľná. Je tu tiež použiteľná funkcia *softmax*. Podobne je potom možné použiť gradientovú metódu, ktorá vedie k nastaveniu vhodných váh a prahu.

Zavedieme základnú funkciu, ktorá presne počíta počet bodov z tréningovej vzorky, ktoré boli klasifikované nekorektne.

**Presnosť
naučeného
klasifikátora**

$$f_0(\mathbf{w}, \theta) = \sum_{p=1}^M \max(0, \text{sign}(-y_p(\mathbf{w} \cdot \mathbf{x}_p - \theta))) \quad (7.15)$$

Funkcia je rastúca pri každom nesprávne klasifikovanom príklade. Funkcia môže byť použitá pri určovaní presnosti klasifikátora, ktorú je možné vyjadriť vzťahom

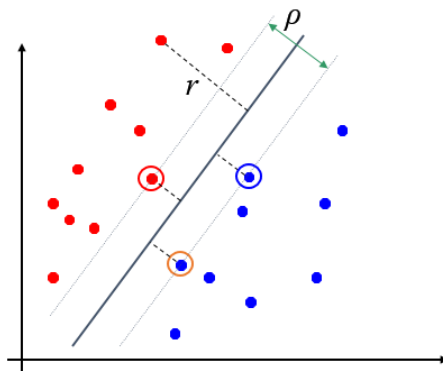
$$\text{presnosť} = 1 - \frac{f_0(\mathbf{w}^*, \theta^*)}{M}. \quad (7.16)$$

kde \mathbf{w}^*, θ^* predstavujú naučené hodnoty parametrov \mathbf{w} a θ .

Táto metrika nadobúda hodnoty $\langle 0, 1 \rangle$ a môže byť použitá aj na percentuálne vyjadrenie.

7.4 Klasifikácia s podpornými vektormi - SVM (Support Vector Machines)

Ideu pásového perceptróna budeme ďalej rozvíjať v tom smere, že budeme hľadať najväčšie možné d (najširší pás), ktoré bude triedy oddeľovať.



Obr. 7.5: Hranice pásu oddeľujúceho obe triedy je možné vypočítať pomocou vzdialenosti dvoch bodov normály na separujúcu priamku, $\frac{2}{\|\mathbf{w}\|}$.

Klasifikácia s maximálnym ρ

- Implikuje, že len podporné vektory sú dôležité; ostatné tréningové príklady sú ignorovateľné.

- Intuícia nám hovorí, že maximálne ρ je dobré a je to tiež v súlade s PAC teóriou.

Základný SVM problém

$$\text{minimalizovať}_{b, \mathbf{w}} \quad \|\mathbf{w}\|^2$$

$$\text{vzhľadom na} \quad \max (0, 1 - y_p(b + \mathbf{x}_p^T \mathbf{w})) = 0, p = 1, \dots, N$$

7.5 Nie celkom presná klasifikácia

Stará formulácia

- Nájsť \mathbf{w} a b také, že $\Phi(\mathbf{w}) = \frac{1}{2}\|\mathbf{w}\|^2 = \frac{1}{2}\|\mathbf{w}\|^T\|\mathbf{w}\|$ je minimálne a pre všetky $(\mathbf{x}_i, y_i), i = 1, \dots, n$ platí $y_i(\mathbf{w}^T \mathbf{x}_i + b) \geq 1$

Modifikovaná formulácia obsahuje voľné premenné

- Nájsť \mathbf{w} a b také, že $\Phi(\mathbf{w}) = \frac{1}{2}\|\mathbf{w}\|^T\|\mathbf{w}\| + C \sum_{i=1}^N \xi_i$ je minimálne a pre všetky $(\mathbf{x}_i, y_i), i = 1, \dots, n$ platí $y_i(\mathbf{w}^T \mathbf{x}_i + b) \geq 1 - \xi_i, \xi_i \geq 0$

Použitie nie celkom presných hraníc nám dáva voľný parameter C , ktorý musí byť fixovaný. Na parameter C sa môžeme dívať ako na prostriedok kontrolujúci "pretrénovanie".

C zohľadňuje maximálny rozsah separátora a tiež fitovanie tréningových dát.

Podobnými úpravami ako v presnej klasifikácii dostaneme nasledujúci problém kvadratického programovania:

$$\text{argmax}_{\mathbf{w}} \quad \sum_{i=1}^N \alpha_i - \frac{1}{2} \sum_{i,j=1}^N \alpha_i \alpha_j y_i y_j \mathbf{x}_i^T \mathbf{x}_j, \quad (7.17)$$

s podmienkami $\sum_{i=1}^N \alpha_i y_i = 0$ a $0 \leq \alpha_i \leq C$ pre všetky i .

Hodnoty α_i v tomto prípade majú ohraničenie zhora C .

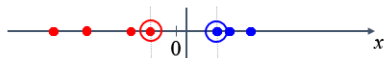
Vplyv každého tréningového bodu v našej rozhodovacej funkcii je úmerný α_i .

C tento vplyv ohraničí.

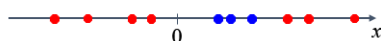
Ak C klesá, maximálny vplyv každého tréningového bodu je znížený, a tak viac z nich sa stáva aktívnymi v rozhodovacej funkcii. C by malo byť fixované.

7.6 Nelineárne SVM

V nasledujúcich dvoch obrázkoch vidíme dve situácie rozloženia bodov na priamke. V prvom prípade sú separovateľné v druhom nie sú.

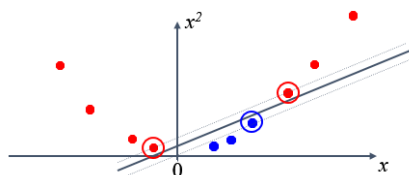


Obr. 7.6: Separovateľné body na priamke pomocou SVM.



Obr. 7.7: Neseparovateľné body na priamke pomocou SVM.

Keby sme tieto body pretransformovali pomocou kvadratickej funkcie do roviny, mohli by sme už lineárny separátor nájsť tak, ako to je na obrázku [7.6](#).



Obr. 7.8: Transformácia dát do roviny umožňuje nájsť lineárny separátor.

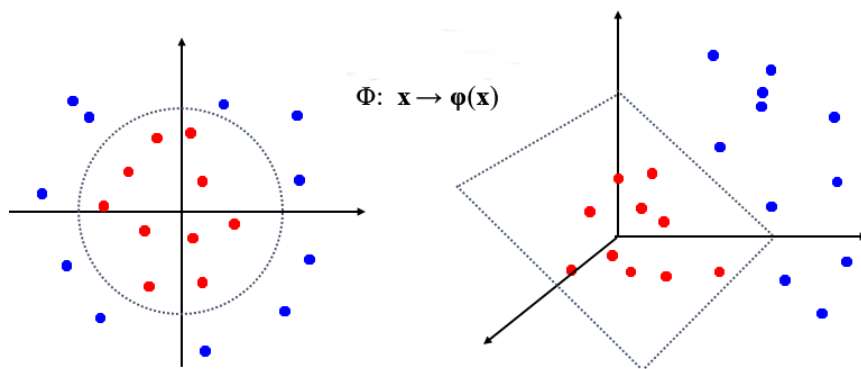
Namapujeme teda dáta do priestoru s vyššou dimenziou. V priestore s vyššou dimenziou je jednoduchšie nájsť oddeľujúcu nadrovinu. Napríklad, transformácia z roviny do 3-rozmerného priestoru je na obrázku [7.6](#).

Trik na použitie jadrovej funkcie pri transformácii dát do priestoru a vyššou dimenziou - Kernel Trick

- Lineárne separátory pracujú so skalárnym súčinom medzi vektormi
 $K(\mathbf{x}_i, \mathbf{x}_j) = \mathbf{x}_i^T \mathbf{x}_j$
- Ak každý dátový bod je namapovaný do priestoru vyššej dimenzie pomocou transformácie $\phi : \mathbf{x} \rightarrow \phi(\mathbf{x})$, dostávame skalárny súčin

$$K(\mathbf{x}_i, \mathbf{x}_j) = \phi(\mathbf{x}_i)^T \phi(\mathbf{x}_j)$$

- **Jadrová funkcia** je funkcia, ktorá je ekvivalentná skalárnemu súčinu v základnom priestore.



Obr. 7.9: Transformácia z roviny do 3-rozmerného priestoru.

Príklad na použitie jadrovej funkcie:

- Uvažujme 2- dimenzionálny priestor, vektory sú v tvare $\mathbf{x} = [x_1, x_2]^T$;
Nech jadrová funkcia má tvar $K(\mathbf{x}_i, \mathbf{x}_j) = (1 + \mathbf{x}_i^T \mathbf{x}_j)^2$
- $K(\mathbf{x}_i, \mathbf{x}_j) = (1 + \mathbf{x}_i^T \mathbf{x}_j)^2$
 $= (1 + x_{i1}x_{j1} + x_{i2}x_{j2})^2 = 1 + x_{i1}^2x_{j1}^2 + 2x_{i1}x_{i2}x_{j1}x_{j2} + x_{i2}^2x_{j2}^2 + 2x_{i1}x_{j1} + 2x_{i2}x_{j2}$
- Uvažujme transformáciu
 $\phi(\mathbf{x})^T = [1, x_1^2, \sqrt{2}x_1x_2, x_2^2, \sqrt{2}x_1, \sqrt{2}x_2]$
- Potrebujeme dokázať $K(\mathbf{x}_i, \mathbf{x}_j) = \phi(\mathbf{x}_i)^T \phi(\mathbf{x}_j)$
- $\phi(x_i)^T \phi(x_j) =$
 $[1, x_{i1}^2, \sqrt{2}x_{i1}x_{i2}, x_{i2}^2, \sqrt{2}x_{i1}, \sqrt{2}x_{i2}]^T * [1, x_{j1}^2, \sqrt{2}x_{j1}x_{j2}, x_{j2}^2, \sqrt{2}x_{j1}, \sqrt{2}x_{j2}] = 1 + x_{i1}^2x_{j1}^2 + 2x_{i1}x_{i2}x_{j1}x_{j2} + x_{i2}^2x_{j2}^2 + 2x_{i1}x_{j1} + 2x_{i2}x_{j2}$
- Požadovaná rovnosť je dokázaná. Teda jadrová funkcia implicitne mapuje dáta do priestoru vyššej dimenzie.

K jadrovej funkcii $K(\mathbf{x}_i, \mathbf{x}_j)$ a dátam je možné vytvoriť maticu v tvare:

$$K = \begin{matrix} & \begin{matrix} K(\mathbf{x}_1, \mathbf{x}_1) & K(\mathbf{x}_1, \mathbf{x}_2) & K(\mathbf{x}_1, \mathbf{x}_3) & \dots & K(\mathbf{x}_1, \mathbf{x}_n) \end{matrix} \\ \begin{matrix} K(\mathbf{x}_2, \mathbf{x}_1) \\ K(\mathbf{x}_2, \mathbf{x}_2) \\ K(\mathbf{x}_2, \mathbf{x}_3) \\ \dots \\ K(\mathbf{x}_n, \mathbf{x}_1) \\ K(\mathbf{x}_n, \mathbf{x}_2) \\ K(\mathbf{x}_n, \mathbf{x}_3) \\ \dots \\ K(\mathbf{x}_n, \mathbf{x}_n) \end{matrix} & \begin{matrix} & & & & \\ & & & & \\ & & & & \\ & & & & \\ & & & & \\ & & & & \\ & & & & \\ & & & & \\ & & & & \\ & & & & \end{matrix} \end{matrix}$$

Poznámka: Ak vytvorená matica je semi-pozitívne definitná, tak funkcia $K(\mathbf{x}_i, \mathbf{x}_j)$ je použiteľná ako jadrová funkcia (Mercerova veta).

Semi-pozitívne definitná symetrická matica zodpovedá semi-pozitívne definitnej symetrickej Gramovej matici.

Príklady jadrových (kernel) funkcií

- Lineárna jadrová funkcia

$$K(\mathbf{x}_i, \mathbf{x}_j) = \mathbf{x}_i^T \mathbf{x}_j$$

- Polynomiálna jadrová funkcia

$$K(\mathbf{x}_i, \mathbf{x}_j) = (1 + \mathbf{x}_i^T \mathbf{x}_j)^p$$

- Gaussova jadrová funkcia

$$K(\mathbf{x}_i, \mathbf{x}_j) = \exp\left(-\frac{\|\mathbf{x}_i - \mathbf{x}_j\|^2}{2\sigma^2}\right)$$

- Sigmoidálna jadrová funkcia

$$K(\mathbf{x}_i, \mathbf{x}_j) = \tanh(\beta_0 \mathbf{x}_i^T \mathbf{x}_j + \beta_1)$$

Keď pracujeme s jadrovými funkciami máme optimalizačný problém v nasledujúcom tvare:

$$\operatorname{argmax}_w \sum_{i=1}^N \alpha_i - \frac{1}{2} \sum_{i=1}^N \sum_{j=1}^N \alpha_i \alpha_j y_i y_j \kappa(\mathbf{x}_i, \mathbf{x}_j)$$

$$\text{vzhľadom na } \sum_{i=1}^N \alpha_i y_i = 0, 0 \leq \alpha_n \leq C, i = 1, \dots, N$$

$$y_{nov} = \operatorname{sign} \left(\sum_{i=1}^N \alpha_i y_i \kappa(\mathbf{x}_i, \mathbf{x}_{nov}) + b \right)$$

7.7 Úlohy

1. Pre niektorú vyššie uvedenú jadrovú funkciu ukážte, že je jadrová.
2. Navrhните množinu dát v dvojrozmernom priestore, ktorá nie je lineárne separovateľná, aplikujte na ňu použitie niektorej jadrovej funkcie.

Kapitola 8

Rozhodovacie stromy

Rozhodovacie stromy predstavujú jednu z často používaných metód na riešenie klasifikačných úloh. Ich výhody spočívajú v prehľadnosti a dobrej interpretácii. Pripomeňme si základné pojmy, ktoré využijeme pri popise rozhodovacích stromov.

Príklady (pozorovania, objekty) budeme označovať ako n -tice

$$\mathbf{x}_i \in \mathcal{R}^n, \quad i \in \{1, 2, \dots, m\},$$

pričom n je prirodzené číslo.

Pre hodnoty cieľového atribútu y_i pre $i \in \{1, 2, \dots, m\}$ platí, že

$$y_i \in C, \quad C = \{1, 2, \dots, c\},$$

ak cieľový atribút môže nadobúdať c rôznych hodnôt.

Tréningová množina predstavuje množinu dvojíc

$$S = \{(\mathbf{x}_1, y_1), \dots, (\mathbf{x}_k, y_k)\} \subseteq \mathcal{R}^n \times \{1, 2, \dots, c\}, \quad k < m,$$

Testovacia množina má nasledujúci tvar

$$T = \{\mathbf{x}_{k+1}, \dots, \mathbf{x}_m\} \subseteq \mathcal{R}^n.$$

Klasifikátorom budeme nazývať funkciu K , ktorá každému príkladu z testovacej množiny priradí hodnotu, triedu, cieľového atribútu, formálne

$$K : T \rightarrow C.$$

Klasifikačná metóda rozhodovacích stromov je založená na rekurzívnom delení tréningovej množiny na dve alebo viac častí – vytvára sa strom. V každom kroku sa použije jedna vysvetľujúca premenná, teda vyberie

sa atribút, pomocou ktorého vetvíme rozhodovací strom. Deliaci bod sa určí tak, aby sa čo najviac zvýšila homogenita vzniknutých podmnožín príkladov (pomocou entropie, Giniho indexu alebo iné kritériá, napríklad korelačné koeficienty, či štatistické testy). Vytváranie stromu je ukončené, ak v koncových vrchoch (listoch) stromu ostali len pozorovania z rovnakej triedy alebo je splnená nejaká obmedzujúca podmienka na minimálny počet pozorovaní v listovom vrchole, hĺbku stromu a podobne. Každému listovému vrcholu (ktorý vlastne zodpovedá zreťazeniu logických podmienok) sa priradí trieda, ktorá sa najčastejšie vyskytuje v pozorovaniach tréningovej množiny zodpovedajúcich danému listovému vrcholu.

Niektoré z kritérií pre výber vysvetľujúcich atribútov na delenie rozhodovacieho stromu si predstavíme v nasledujúcej podkapitole.

8.1 Miery výberu atribútov

Miera výberu atribútov slúži na určenie deliaceho kritéria, ktoré oddelí danú množinu príkladov do tried. Ak sme rozdelili množinu príkladov na základe istého deliaceho kritéria, každá podmnožina by mala byť v ideálnom prípade homogénna, teda obsahuje iba príklady jednej triedy. Najlepšie deliace kritérium je teda také, ktoré nám zabezpečí takýto výsledok.

Miery výberu atribútov sa zvyknú nazývať aj **deliace pravidlá**, pretože určujú, akým spôsobom sú príklady v danom uzle rozhodovacieho stromu rozdelené. Tieto miery poskytujú zoradenie atribútov vzhľadom na tzv. kvalitu rozdelenia uvažovanej množiny príkladov do tried. Atribút, ktorý dosiahne najvyššie skóre, je vybraný za deliaci atribút. Ak deliaci atribút obsahuje spojité hodnoty a uvažujeme binárne stromy, je potrebné určiť ako deliace kritérium tzv. deliaci bod, resp. deliacu podmnožinu.

Ak je k uzlu stromu pre množinu príkladov S priradené deliace kritérium, hrany vychádzajúce z tohto uzla zodpovedajú jednotlivým výstupom tohto kritéria.

Entropia

Entropia je miera homogenity v skupine príkladov, miera pomocou ktorej vetvíme rozhodovací strom. Hodnoty nadobúda z jednotkového intervalu. Čím vyššia hodnota, tým je homogenita nižšia.

Nech $x_i(A) = a$ vyjadruje hodnotu a atribútu A , ktorá prislúcha príkladu x_i pre $i \in \{1, 2, \dots, k\}$, pričom k je prirodzené číslo, ktoré vyjadruje počet príkladov v tréningovej množine.

Entropiu v jednej hrane uzla rozhodovacieho stromu môžeme vypočítať nasledovne:

$$H(x_i(A) = a) = - \sum_{y_j \in \{1, \dots, c\}} p(y_j) \cdot \log_2(p(y_j)),$$

pričom $p(y_j)$ je pravdepodobnosť, že príklad patriaci hrane $x_i(A) = a$ bude klasifikovaný do triedy y_j .

Celková entropia v uzle má tvar

$$H(A) = \sum_{a \in A} p(x_i(A) = a) \cdot H(x_i(A) = a).$$

Informačný zisk

Informačný zisk je definovaný ako rozdiel medzi entropiou v jednej hrane uzla a celkovou entropiou v uzle. Informačný zisk vyjadrujeme

$$I(x_i(A) = a) = H(x_i(A) = a) - H(A),$$

pričom $H(x_i(A) = a)$ je entropia v jednej hrane uzle a $H(A)$ je celková entropia v uzle.

Informačný zisk sa pri výbere atribútu snažíme maximalizovať. Táto miera výberu atribútov má nevýhodu, že väčšinou uprednostní také atribúty, ktoré majú viac vetiev. Uvedený nedostatok je možné odstrániť normalizovaním informačného zisku, zavedením tzv. **pomerového informačného zisku** pomocou tzv. pomerovej entropie.

Pomerovú entropiu definujeme

$$H_P(A) = - \sum_{a \in A} p(x_i(A) = a) \cdot \log_2(p(x_i(A) = a)),$$

pričom $p(x_i(A) = a)$ je pravdepodobnosť, že príklad x_i nadobúda v atribúte A hodnotu a .

Následne môžeme pomocou pomerovej entropie definovať pomerový informačný zisk nasledujúcim spôsobom.

Pomerový informačný zisk

Pomerový informačný zisk môžeme zapísať v tvare

$$I_P(x_i(A) = a) = \frac{I(x_i(A) = a)}{H_P(A)},$$

pričom $I(x_i(A) = a)$ je informačný zisk.

Giniho index je možné použiť pre kategoriálne atribúty, resp. numerické atribúty po ich diskretizácii. Nižšie hodnoty Giniho indexu indikujú lepšiu schopnosť atribútu rozdeliť príklady do tried.

Giniho index v jednej hrane uzla môžeme vyjadriť v tvare

$$G(x_i(A) = a) = 1 - \sum_{y_j \in \{1, \dots, c\}} p(y_j) * p(y_j),$$

pričom $p(y_j)$ je pravdepodobnosť, že príklad patriaci hrane $x_i(A) = a$ bude klasifikovaný do triedy y_j .

Celkový Giniho index v uzle vypočítame

$$G(A) = \sum_{a \in A} p(x_i(A) = a) \cdot G(x_i(A) = a).$$

V tejto podkapitole sme predstavili niekoľko mier na výber atribútov pre delenie rozhodovacieho stromu. Bolo navrhnutých viacero ďalších mier založených napríklad na chí-kvadrát testovaní nezávislosti. Niektoré postupy používajú aj viacatribútové delenie, pri ktorom je delenie stromu založené na kombinácii atribútov. Napriek viacerým štúdiám nebolo preukázané, aby niektorá z týchto mier bola nadradená nad všetkými ostatnými. Z tohto dôvodu závisí výber vhodnej miery od použitej aplikáčnej domény, typu atribútov a ďalších faktorov.

8.2 Algoritmus pre generovanie rozhodovacích stromov

Ak už poznáme spôsob, akým ohodnotiť kvalitu rozdelenia rozhodovacieho stromu v jednotlivých uzloch, môžeme definovať všeobecný algoritmus na generovanie rozhodovacieho stromu:

Algoritmus 8.1 *Konštrukcia rozhodovacieho stromu*

KonštruujRS(tréningsová množina S , $\text{minPodiel} \in [0, 1]$)

ak minPodiel objektov z S patrí do triedy $y_j \in C$

potom vytvor listový uzol zaradzujúci objekty do y_j ;

inak

pre každý atribút A

pre každé rozdelenie hodnôt $a \in A$

Algoritmus na generovanie rozhodovacieho stromu

ohodnot kvalitu rozdelenia;
vykonaj najlepšie rozdelenie;
vzniknú množiny S_1, S_2, \dots, S_p ;
KonštruujRS(S_1 , *minPodiel*);
KonštruujRS(S_2 , *minPodiel*);
 ⋮
KonštruujRS(S_p , *minPodiel*).

Ak sa pozrieme na konkrétne implementácie tohto všeobecného algoritmu zistíme, že algoritmus **ID3**, ktorého autorom je Quinlan, využíva **entropiu** a od nej odvodený **informačný zisk** ako kritérium delenia rozhodovacieho stromu. Quinlan navrhol aj algoritmus C4.5, ktorý je vylepšeným a doplneným algoritmom ID3, využíva **entropiu** a od nej odvodený **pomernový informačný zisk**. Algoritmus **CART** využíva **Giniho index**.

8.3 Preučenie rozhodovacieho stromu a jeho prerezávanie

Preučenie (tzv. overfitting) rozhodovacieho stromu znamená, že rozhodovací strom klasifikuje správne príklady v tréningovej údajovej sade, ale pri klasifikácii príkladov z testovacej údajovej sady poskytuje nesprávne výsledky.

Po vytvorení rozhodovacieho stromu môžu byť niektoré vetvy ovplyvnené anomáliami v tréningovej sade údajov z dôvodu extrémnych hodnôt alebo chybných údajov. **Metódy prerezávania rozhodovacích stromov** slúžia na riešenie takýchto situácií. Zvyčajne sa používajú štatistické metódy na odstránenie tých najmenej spoľahlivých vetiev stromu. Orezané stromy majú tendenciu byť jednoduchšie, a teda aj ľahšie interpretovateľné.

Prerezávanie rozhodovacieho stromu

Existujú dva základné prístupy k prerezávaniu rozhodovacích stromov:

- prerezávanie počas generovania rozhodovacieho stromu (tzv. prepruning),
- prerezávanie po vygenerovaní rozhodovacieho stromu (tzv. postpruning),

Pri prvom prístupe sa rozhodovací strom prerezáva predčasným zastavením jeho konštrukcie (napríklad rozhodnutím nedeliť ďalej podmnožinu tréningových n -tíc v danom uzle. Po zastavení sa uzol stane listom.

Príkladom v liste priradíme triedu, ktorá zodpovedá najčastejšej hodnote medzi podmnožinou n -tíc alebo na základe rozdelenia pravdepodobnosti týchto n -tíc do jednotlivých tried.

Na posúdenie správnosti rozdelenia v uzle používame pri konštrukcii rozhodovacieho stromu miery výberu atribútov, akými sú napríklad informačný zisk, Giniho index, štatistická významnosť a iné miery. Ak pri rozdeľovaní príkladov v uzle by vzniklo rozdelenie, ktorého kvalita klesne pod vopred stanovenú hranicu, rozdelenie danej podmnožiny zastavíme. Niekedy je však náročné určiť primeranú hranicu, teda vhodný prah pre zastavenie delenia. Vysoké prahy môžu mať za následok príliš zjednodušené stromy, zatiaľ čo nízke prahové hodnoty by mohli viesť k veľmi malému zjednodušeniu.

Druhým a bežnejším prístupom je prerezávanie po vygenerovaní rozhodovacieho stromu, pri ktorom sa odstraňujú podstromy z úplne vytvoreného rozhodovacieho stromu. Podstrom v danom uzle sa orezáva odstránením jeho vetiev a jeho nahradením listom. Príkladom v liste je priradená trieda, ktorá je najčastejšia v príslušnom podstrome.

Bagging a boosting

Na zvyšovanie presnosti klasifikácie s použitím rozhodovacích stromov sa používa aj metóda tzv. **baggingu**, pri ktorej sa postupne vytvára určitý počet rozhodovacích stromov (odporúča sa veľký počet), pričom ich vytváranie je nezávislé (môže byť realizované paralelne). Každý rozhodovací strom sa tak vytvára nad inou náhodnou vzorkou príkladov z tréningovej množiny. Náhoda sa zapája aj do výberu vysvetľujúcej premennej pri každom delení, pričom sa nevyberá zo všetkých premenných, ale len z menšej náhodne vybranej podmnožiny premenných. Výsledná klasifikácia sa získa takým spôsobom, že sa vytvoria všetky rozhodovacie stromy, pričom trieda pre konkrétny príklad sa určí ako modus zo všetkých tried, ktoré boli pridelené danému príkladu v rámci všetkých vygenerovaných rozhodovacích stromov. V takomto prípade hovoríme o tzv. **náhodných lesoch**. Je možné matematicky ukázať, že zvyšovaním počtu rozhodovacích stromov nedochádza k preučeniu, teda je vhodné vytvoriť čo najviac rozhodovacích stromov, pokiaľ to dovoľuje výpočtový výkon. Praktické skúsenosti z mnohých oblastí ukazujú, že náhodné lesy dokázali mnohokrát výrazným spôsobom zvýšiť presnosť klasifikácie v porovnaní s jedným rozhodovacím stromom.

Ďalším spôsobom zvyšovania presnosti klasifikácie pri použití rozhodovacích stromov je metóda tzv. **boostingu**. Na rozdiel od náhodných lesov pri boostingu nevytvárame jednotlivé rozhodovacie stromy nezávisle, ale postupne ich budujeme tak, aby každý ďalší rozhodovací strom čo najviac prispel k znižovaniu chyby klasifikácie pomocou doteraz vytvorených rozhodovacích stromov. Pri veľkom počte rozhodovacích stromov alebo pri nevhodnej veľkosti rozhodovacích stromov môže dôjsť k preučeniu, preto je pri tejto metóde náročnejšia inicializácia vstupných hyperparametrov.

Kapitola 9

Zhlukovanie

Zhlukovanie (zhluková analýza, klastering) je metóda, ktorej cieľom je zoskupovanie množiny navzájom podobných objektov do tried, tzv. zhlukov.

Zhluk je množina navzájom podobných objektov, ktoré sa od objektov nepatriacich do zhluku odlišujú. Podobnosť objektov popisujeme na základe hodnôt ich atribútov vzhľadom na použité metrické miery (tzv. metriky).

Aj klasifikácia je metódou, ktorá rozdeľuje množinu objektov do tried alebo skupín, ale pri zhlukovaní nie sú cieľové triedy vopred dané, označenie cieľových tried vopred nepoznáme. Klasifikácia vyžaduje poznať označenia cieľových tried pre objekty v tréningovej množine, na základe čoho sa vybuduje klasifikátor na modelovanie každej skupiny. Zhlukovanie predstavuje v podstate opačný proces. Najprv rozdelíme množinu objektov do skupín na základe ich podobnosti a následne priradíme označenie pre relatívne malý počet zhlukov, ktoré takýmto spôsobom vzniknú.

Zhlukovanie má dôležité postavenie v strojovom učení, rozpoznávaní vzorov, analýze obrazu, dolovaní údajov, bioinformatike a iných odvetviach.

Jednotlivé algoritmy zhlukovania sa líšia spôsobom hľadania zhlukov, efektívnou vyhľadávania a obsahom jednotlivých zhlukov. Môžeme uvažovať:

- **klasický prístup:** každý objekt patrí do určitého zhluku, alebo do neho nepatrí,
- **fuzzy prístup:** každý objekt patrí do každého zhluku s určitým stupňom príslušnosti.

Vo všeobecnosti, rozlišujeme niekoľko typov techník zhlukovania:

- segmentačné metódy,
- hierarchické metódy,
- metódy založené na hustote,

**Typy
zhlukovania**

- metódy založené na mriežke,
- metódy založené na modeloch.

V nasledujúcich častiach si predstavíme niektoré príklady týchto metód.

9.1 Typy údajov pri zhlukovaní

Predpokladajme, že množina údajov obsahuje n objektov, ktoré môžu reprezentovať osoby, krajiny, dokumenty a iné príklady. Pri zhlukovaní najčastejšie používame údajové štruktúry **matice údajov** (štruktúra typu objekt-atribút) a **matice podobnosti, resp. odlišnosti** (štruktúra typu objekt-objekt).

Matica údajov X je matica v tvare $n \times p$, pomocou ktorej je možné reprezentovať hodnoty n objektov $\{o_1, o_2, \dots, o_n\}$ v p atribútoch $\{a_1, a_2, \dots, a_p\}$:

Matica údajov

$$\begin{pmatrix} x_{11} & \dots & x_{1j} & \dots & x_{1p} \\ \dots & \dots & \dots & \dots & \dots \\ x_{i1} & \dots & x_{ij} & \dots & x_{ip} \\ \dots & \dots & \dots & \dots & \dots \\ x_{n1} & \dots & x_{nj} & \dots & x_{np} \end{pmatrix} \quad (9.1)$$

Matica podobnosti, resp. odlišnosti reprezentuje vzťahy medzi dvojicami objektov pre všetkých n objektov. Štruktúra je vo forme matice typu $n \times n$, pričom hodnota $d_{i,l}$ reprezentuje zistenú mieru podobnosti, resp. odlišnosti medzi objektmi i a j pre $i \in \{1, 2, \dots, n\}$ a $l \in \{1, 2, \dots, n\}$:

Matica podobnosti

$$\begin{pmatrix} 1 & & & & \\ d(2,1) & 1 & & & \\ d(3,1) & d(3,2) & 1 & & \\ \vdots & \vdots & \vdots & & \\ d(n,1) & d(n,2) & \dots & \dots & 1 \end{pmatrix} \quad (9.2)$$

V matici podobnosti alebo matici odlišnosti platí $d(i, l) = d(l, i)$ pre $i \in \{1, 2, \dots, n\}$ a $l \in \{1, 2, \dots, n\}$. Z matice údajov je možné vytvoriť maticu podobnosti, resp. maticu odlišnosti podľa podmienok použitého algoritmu zhlukovania.

V prípade intervalových premenných (napríklad výška osoby, hmotnosť osoby, teplota) používame **standardizáciu hodnôt atribútu**.

Pri štandardizácii transformujeme hodnoty daného atribútu na hodnoty bez jednotiek (bez ohľadu na to, či meriame v metroch alebo palcoch). Napríklad pre atribút a_j , $j \in \{1, 2, \dots, p\}$ je postup nasledujúci:

1. Nech $x_{1j}, x_{2j}, \dots, x_{nj}$ sú hodnoty všetkých objektov v atribúte a_j pre $j \in \{1, 2, \dots, p\}$. Nech m_j je aritmetický priemer hodnôt v atribúte a_j pre $j \in \{1, 2, \dots, p\}$. Potom priemernú absolútnu odchýlku s_{a_j} vypočítame:

$$s_{a_j} = \frac{1}{n}(|x_{1j} - m_j| + |x_{2j} - m_j| + \dots + |x_{nj} - m_j|).$$

2. Z-skóre (štandardizované skóre) vypočítame ako:

$$z_{ij} = \frac{x_{ij} - m_j}{s_j}$$

pre $i \in \{1, 2, \dots, n\}$ a $j \in \{1, 2, \dots, p\}$. Z-skóre môže byť užitočne použité na detekciu anomálií v množine údajov.

Po štandardizácii údajov alebo v prípade možnosti aj bez nej nasleduje vytvorenie matice odlišnosti, resp. matice podobnosti.

Odlišnosť, resp. podobnosť medzi objektmi, ktoré sú vyjadrené intervalovými premennými definujeme ako vzdialenosť medzi objektmi. Používa sa napríklad **Euklidovská vzdialenosť**

$$d(i, l) = \sqrt{(x_{i1} - x_{l1})^2 + (x_{i2} - x_{l2})^2 + \dots + (x_{ip} - x_{lp}^2)},$$

pričom x_{ij} a x_{lj} predstavuje hodnoty objektov o_i a o_l v atribúte a_j pre $i \in \{1, 2, \dots, n\}$, $l \in \{1, 2, \dots, n\}$ a $j \in \{1, 2, \dots, p\}$.

Vo všeobecnosti, každému atribútu môžeme priradiť váhu.

Nech $\{w_1, w_2, \dots, w_p\}$ sú váhy jednotlivých atribútov pre $j \in \{1, 2, \dots, p\}$. V takom prípade môžeme vyjadriť **váženú Euklidovskú vzdialenosť** nasledujúcim spôsobom:

$$d(i, k) = \sqrt{w_1(x_{i1} - x_{k1})^2 + w_2(x_{i2} - x_{k2})^2 + \dots + w_p(x_{ip} - x_{kp}^2)},$$

pričom x_{ij} a x_{kj} predstavuje hodnoty objektov o_i a o_k v atribúte a_j pre $i \in \{1, 2, \dots, n\}$, $k \in \{1, 2, \dots, n\}$ a $j \in \{1, 2, \dots, p\}$.

Iným príkladom vzdialenostnej metriky je **Manhattanská vzdialenosť**, ktorá je definovaná nasledujúcim spôsobom:

$$d(i, j) = |x_{i1} - x_{j1}| + |x_{i2} - x_{j2}| + \dots + |x_{ip} - x_{jp}|,$$

pričom x_{ij} a x_{kj} predstavuje hodnoty objektov o_i a o_k v atribúte a_j pre $i \in \{1, 2, \dots, n\}$, $k \in \{1, 2, \dots, n\}$ a $j \in \{1, 2, \dots, p\}$.

Minkowského vzdialenosť je zovšeobecnením Euklidovej vzdialenosti alebo Manhattanskej vzdialenosti. Je definovaná takýmto spôsobom:

$$d(i, j) = \left(|x_{i1} - x_{j1}|^m + |x_{i2} - x_{j2}|^m + \dots + |x_{ip} - x_{jp}|^m \right)^{\frac{1}{m}},$$

pričom m je prirodzené číslo väčšie ako nula, x_{ij} a x_{kj} predstavuje hodnoty objektov o_i a o_k v atribúte a_j pre $i \in \{1, 2, \dots, n\}$, $k \in \{1, 2, \dots, n\}$ a $j \in \{1, 2, \dots, p\}$. Takáto vzdialenosť sa označuje aj ako L_m norma, L_1 predstavuje Manhattanskú vzdialenosť, L_2 predstavuje Euklidovskú vzdialenosť.

V prípade binárnych atribútov môžeme použiť tabuľku početností, v ktorej budeme pre každú dvojicu objektov reprezentovať, v koľkých prípadoch sa hodnoty ich atribútov zhodujú. Podobný princíp môžeme zvoliť aj pri kategoriálnych premenných s rôznym počtom hodnôt. Pri ordinálnych premenných môžeme navyše reprezentovať poradie hodnôt atribútov a maticu podobnosti generovať na základe zachovávaní, resp. nezachovávaní poradia hodnôt medzi dvojicou objektov.

9.2 Segmentačné metódy zhlukovania

Tieto metódy zhlukovania sú založené na vytvorení vopred známeho počtu zhlukov, pričom každý zhluk musí obsahovať aspoň jeden objekt a každý objekt musí patriť práve do jedného zhlukov. V niektorých fuzzy algoritmoch môže byť druhá podmienka zjemnená. Tieto metódy zhlukovania sú založené na iteratívnom prístupe, pričom našim cieľom je zlepšiť zhlukovanie postupným presúvaním objektov medzi zhlukmi.

Najznámejšími a najpoužívanějšími segmentačnými metódami zhlukovania je **k -means zhlukovanie**, **zhlukovanie k -medoidov** a ich rôzne varianty.

Nech X je matica údajov reprezentujúca hodnoty n objektov $\{o_1, o_2, \dots, o_n\}$ v p atribútoch $\{a_1, a_2, \dots, a_p\}$, napríklad x_{ij} vyjadruje hodnotu objektu o_i v atribúte a_j pre $i \in \{1, 2, \dots, n\}$ a $j \in \{1, 2, \dots, p\}$. Nech k je počet zhlukov ($k \leq n$).

Metóda **k -means zhlukovania** má nasledujúci postup:

1. Náhodne vyberieme k objektov, každý z nich reprezentuje tzv. **centrum zhluku**, ktoré označíme C_1, C_2, \dots, C_k .
2. Každý objekt priradíme do zhluku, ktorý mu je najpodobnejší vzhľadom na vzdialenosť medzi objektom a centrom zhluku.
3. Ak je každý objekt priradený k nejakému zhluku, vypočítame nové centrá zhlukov tak, že vzdialenosť medzi všetkými objektmi zhluku a novým centrom zhluku bude minimalizovaná, spriemerovaná. Napríklad ak budeme uvažovať v 2-rozmernom priestore zhluk s dvoma objektmi o_1, o_2 , tak centrum zhluku C_1 bude nadobúdať hodnoty

$$x_{C_11} = \frac{x_{11} + x_{21}}{2} \quad \text{a} \quad x_{C_12} = \frac{x_{12} + x_{22}}{2}.$$

**K-means
zhlukovanie**

4. Ak po prepočítaní centier zhlukov nastala zmena aspoň v jednom centre zhluku, pokračujeme krokom 2.

Na určenie toho, či vytvorené zhluky sú kompaktné a kvalitne oddelené, môžeme použiť **kritérium štvorcovej chyby**:

$$E = \sum_{i=1}^k \left(\sum_{o_l \in C_i} |x_{l1} - x_{C_i1}|^2 + |x_{l2} - x_{C_i2}|^2 + \dots + |x_{lp} - x_{C_i p}|^2 \right),$$

pričom l označuje indexy objektov patriacich do zhluku C_i pre $i \in \{1, 2, \dots, k\}$. Čím je hodnota kritéria štvorcovej chyby menšia, tým sú zhluky kompaktnéjšie a kvalitnejšie oddelené.

Algoritmus pre k -means zhlukovanie môžeme vo všeobecnosti zapísať nasledovne:

Algoritmus 9.1 *K-means zhlukovanie*

Nech X je matica údajov, ktorá reprezentuje hodnoty n objektov $\{o_1, o_2, \dots, o_n\}$ v p atribútoch $\{a_1, a_2, \dots, a_p\}$. Nech k je počet zhlukov, $k \leq n$.

begin

náhodne vyber k objektov z množiny objektov $\{o_1, o_2, \dots, o_n\}$;

vybrané objekty označ za centrá zhlukov C_i pre $i \in \{1, 2, \dots, k\}$;

repeat

pre každý objekt $\{o_1, o_2, \dots, o_n\}$ nájdí najbližšie centrum zhluku C_i pre $i \in \{1, 2, \dots, k\}$;

pre každé centrum zhluku C_i , $i \in \{1, 2, \dots, k\}$

aktualizuj $x_{C_i1}, \dots, x_{C_i p}$ tak, aby vzdialenosť od objektov zhluku bola minimálna;

until *centrá zhlukov C_i , $i \in \{1, 2, \dots, k\}$ sú po aktualizácii nezmenené;*

end

Metóda k -means zhlukovania môže byť použitá len v prípade, že je definované centrum zhluku pomocou priemerných hodnôt. To nie je možné v prípade kategoriálnych premenných. Potreba vopred určiť počet zhlukov môže byť tiež nevýhodou. Existuje niekoľko ďalších variánt tejto metódy, ktoré sa líšia v spôsobe výberu počiatočných zhlukov, výpočte podobnosti medzi objektmi, a v spôsobe aktualizácie centier zhlukov.

Metóda k -means je citlivá na extrémne hodnoty, keďže používame aritmetické priemery a kritérium štvorcovej chyby. Tieto nevýhody môžu byť odstránené použitím metódy zhlukovania k -medoidov, ktorú si predstavíme v nasledujúcej časti.

Metóda k-medoidov zhlukovania používa namiesto výpočtu priemerných hodnôt pri aktualizácii centier zhlukov ako centrá zhlukov samotné objekty z matice údajov. Konkrétne za centrum zhluku vyberieme objekt, ktorý je k stredu zhluku najbližšie. Rozdelenie objektov do príslušných zhlukov je založené na princípe minimalizácie súčtu odlišností medzi objektmi a centrami zhlukov.

Kritérium absolútnej chyby pri tejto metóde je definované

$$E = \sum_{i=1}^k \left(\sum_{o_l \in C_i} |x_{l1} - x_{C_i1}| + |x_{l2} - x_{C_i2}| + \dots + |x_{lp} - x_{C_ip}| \right),$$

pričom l označuje indexy objektov patriacich do zhluku C_i pre $i \in \{1, 2, \dots, k\}$.

Kritérium absolútnej chyby môže byť použité na vyjadrenie tzv. **nákladovej funkcie**, ktorá určuje hodnotu, o ktorú sa zmenila absolútna chyba po nahradení aktuálneho centra zhluku iným objektom. Celková nákladová funkcia nahradenia predstavuje sumu nákladov za všetky objekty, ktoré nie sú priradené k centrá zhlukov. Ak sú celkové náklady záporné (absolútna chyba klesla), tak aktuálne centrum zhluku je nahradené príslušným objektom. Ak sú celkové náklady pozitívne, centrum zhluku sa nemení.

Algoritmus zhlukovania k-medoidov môžeme vo všeobecnosti zapísať nasledovne:

Algoritmus 9.2 Zhlukovanie k-medoidov

Nech X je matica údajov, ktorá reprezentuje hodnoty n objektov $\{o_1, o_2, \dots, o_n\}$ v p atribútoch $\{a_1, a_2, \dots, a_p\}$. Nech k je počet zhlukov, $k \leq n$.

begin

náhodne vyber k objektov z množiny objektov $\{o_1, o_2, \dots, o_n\}$;

vybrané objekty označ za centrá zhlukov C_i pre $i \in \{1, 2, \dots, k\}$;

repeat

pre každý objekt $\{o_1, o_2, \dots, o_n\}$ nájdí najbližšie centrum zhluku C_i pre $i \in \{1, 2, \dots, k\}$;

vypočítaj chybu E ;

pre každé centrum zhluku C_i , $i \in \{1, 2, \dots, k\}$

náhodne vyber objekt o_l z množiny $\{o_1, o_2, \dots, o_n\}$;

vypočítaj zmenu chyby, ak aktuálne centrum zhluku C_i nahradíme o_l ;

ak chyba klesla, nahrad' centrum zhluku C_i objektom o_l ;

until centrá zhlukov C_i , $i \in \{1, 2, \dots, k\}$ sú po aktualizácii nezmenené;

end

Metóda zhlukovania k -medoidov je robustnejšia ako metóda k -means zhlukovania s ohľadom na zašumené a anomálne údaje, pretože reprezentácia pomocou medoidov (reprezentujúcich objektov) je menej ovplyvnená extrémnymi hodnotami ako metóda k -means s centrami počítanými pomocou priemerných hodnôt. Spracovanie metódou k -medoidov je však výpočtovo náročnejšie. Obe metódy vyžadujú, aby používateľ určil počet zhlukov.

V prípade veľkého množstva objektov môžeme pre výber vhodných medoidov vybrať namiesto všetkých údajov len ich časť. Na takejto myšlienke sú založené napríklad algoritmy CLARANS (zhlukovanie vo veľkých aplikáciách založené na náhodnom vyhľadávaní).

Okrem použitia priemeru alebo medoidu ako miery stredu zhluku, iné varianty segmentačných metód využívajú pri rozdelení objektov do zhlukov medián alebo modus. Výsledkom pri takomto postupe sú **metódy zhlukovania k -mediánov alebo metóda k -modusov**.

9.3 Hierarchické metódy zhlukovania

Hierarchické metódy zhlukovania sú založené na zoskupovaní objektov do stromu zhlukov. Hierarchická dekompozícia môže byť formovaná buď zdola nahor (aglomeratívne hierarchické zhlukovanie) alebo zhora nadol (rozkladajúce hierarchické zhlukovanie). Integrácia hierarchických a segmentačných metód zhlukovania sa ukázala ako vhodné riešenie pri mnohých aplikáciách.

Agglomeratívne hierarchické zhlukovanie je založené na spájaní, pričom v nulte iterácii je každému objektu priradený vlastný zhluk. Tieto jednoprvkové zhluky sú následne spájané do stále väčších zhlukov, až kým nie je splnená podmienka ukončenia, napríklad počet zhlukov (v teoretickom prípade môže nastať situáciu, že po ukončení zhlukovania vznikne jeden zhluk, ktorý bude obsahovať všetky objekty). Väčšina metód hierarchického zhlukovania patrí do tejto skupiny, pričom sa líšia len v definícii a spôsobe určovania podobnosti zhlukov.

Rozkladajúce hierarchické zhlukovanie je založené na rozdeľovaní, pričom v nulte iterácii patria všetky objekty jednému zhluku. Tento zhluk je ďalej rozdeľovaný do stále menších zhlukov, až kým nie je splnená podmienka ukončenia (v teoretickom prípade môže nastať aj taká situácia, že po ukončení zhlukovania vzniknú jednoprvkové zhluky). Väčšina metód hierarchického zhlukovania sú rozkladajúce, líšia sa definíciou a spôsobom určovania podobnosti zhlukov.

Dendrogram je stromová štruktúra, ktorú môžeme použiť na znázornenie hierarchického zhlukovania. Ilustruje, akým spôsobom sú objekty v jednotlivých iteráciách spájané, resp. rozdeľované.

Nech C_i a C_j sú zhluky pre $i \in \{1, 2, \dots, k\}$, $j \in \{1, 2, \dots, k\}$, $i \neq j$ a nech $o \in C_i$ a $o' \in C_j$. Nech d vyjadruje funkciu vzdialenosti medzi dvoma objektmi a funkcia d^* vyjadruje funkciu vzdialenosti medzi dvoma zhlukmi.

Pri spájaní, resp. rozdeľovaní zhlukov sa najčastejšie používajú nasledujúce miery na určovanie vzdialenosti medzi zhlukmi:

1. **Minimálna vzdialenosť** – vzdialenosť medzi dvoma objektmi z dvoch rôznych zhlukov, ktorá je najmenšia:

$$d_{\min}^*(C_i, C_j) = \min\{d(o, o') : o \in C_i, o' \in C_j\}.$$

V tomto prípade hovoríme o zhlukovacom algoritme **najbližších susedov**. Ak je navyše daná podmienka ukončenia zhlukovania, ak vzdialenosť medzi najbližšími zhlukmi presiahne určitý prah, hovoríme o **algoritme jednoduchého spojenia**. Ak si predstavíme objekty ako vrcholy v grafe a hrany, ktoré predstavujú cestu medzi vrcholmi v zhluku, potom spojenie dvoch zhlukov C_i a C_j zodpovedá pridaniu hrany medzi najbližšou dvojicou vrcholov v C_i a C_j . Výsledný graf tvorí strom – aglomeratívne hierarchické zhlukovanie, ktoré používa minimálnu vzdialenosť, nazývame aj **algoritmus minimálnej kostry grafu**.

2. **Maximálna vzdialenosť** – vzdialenosť medzi dvoma objektmi z dvoch rôznych zhlukov, ktorá je najväčšia:

$$d_{\max}^*(C_i, C_j) = \max\{d(o, o') : o \in C_i, o' \in C_j\}.$$

V tomto prípade hovoríme o zhlukovacom algoritme **najvzdialenejších susedov**. Ak je navyše daná podmienka ukončenia zhlukovania, ak vzdialenosť medzi najbližšími zhlukmi presiahne určitý prah, hovoríme o **algoritme úplného spojenia**. Ak si predstavíme objekty ako vrcholy v grafe a hrany, ktoré spájajú vrcholy, tak každý zhluk predstavuje úplný podgraf, teda hrany spájajú všetky vrcholy v zhluku. Ak sú zhluky približne rovnakej veľkosti, metóda produkuje zhluky vysokej kvality.

3. **Priemerná vzdialenosť** – priemerná vzdialenosť medzi všetkými dvojicami objektov z dvoch rôznych zhlukov:

$$d_{\text{average}}^*(C_i, C_j) = \frac{1}{n_i n_j} \sum_{o \in C_i} \sum_{o' \in C_j} d(o, o'),$$

pričom n_i a n_j predstavujú počet objektov v zhlukoch C_i a C_j pre $i \in \{1, 2, \dots, k\}$, $j \in \{1, 2, \dots, k\}$.

4. **Vzdialenosť priemerov** – vzdialenosť medzi centrami dvoch zhlukov.

Miery minimálnej a maximálnej vzdialenosti sú citlivé na extrémne hodnoty. V takýchto prípadoch je vhodné používať priemernú vzdialenosť alebo vzdialenosť priemerov. Kým vzdialenosť priemerov má jednoduchý výpočet, priemerná vzdialenosť dokáže spracovať aj kategoriálne a numerické údaje.

Algoritmus aglomeratívneho hierarchického zhlukovania pre jednoduché spojenie môžeme vo všeobecnosti zapísať nasledovne:

Algoritmus 9.3 *Aglomeratívne hierarchické zhlukovanie pre jednoduché spojenie*

Nech X je matica údajov, ktorá reprezentuje hodnoty n objektov $\{o_1, o_2, \dots, o_n\}$ v p atribútoch $\{a_1, a_2, \dots, a_p\}$. Nech θ je prah spájania zhlukov.

begin

pre každý objekt o_i , $i \in \{1, 2, \dots, n\}$ vytvor zhluk C_i , $i \in \{1, 2, \dots, n\}$;

$k := n$;

repeat

$\min \text{Vzdialenosť} := \min \{d_{\min}^*(C_i, C_j) : i \in \{1, 2, \dots, k\}, j \in \{1, 2, \dots, k\}, i \neq j\}$;

ak $d_{\text{extrmmin}}^*(C_i, C_j) = \min \text{Vzdialenosť}$

$C_i := C_i \cup C_j$;

zmeň označenia zhlukov na C_1, C_2, \dots, C_{k-1} ;

$k := k-1$;

until $\min \text{Vzdialenosť} > \theta$;

end

Ak v jednej iterácii dôjde k spojeniu dvoch zhlukov, nasledujúci krok už pracuje iba s jedným novým zhlukom, ktorý už nie je možné znovu rozdeliť. Túto nevýhodu je možné odstrániť použitím kombinovaných metód zhlukovania.

9.4 Metódy zhlukovania založené na hustote

Metódy zhlukovania založené na hustote definujú zhluk ako maximálnu množinu husto spojených objektov. Táto metóda si vyžaduje použitie nasledujúcich pojmov:

- ε -**okolie objektu** je okolie objektu s polomerom ε .
- **Vnútorňý objekt** je objekt, ktorého ε -okolie obsahuje aspoň m objektov.
- Objekt o je **priamo dosiahnuteľný vzhľadom na hustotu** z objektu o' , ak o je v ε -okolí o' a o' je vnútorňým objektom.
- Objekt o je **dosiahnuteľný vzhľadom na hustotu** z objektu o' pre dané ε a m , ak existuje postupnosť objektov o_1, o_2, \dots, o_t taká, že $o_1 = o'$ a $o_t = o$, pričom o_{i+1} je priamo dosiahnuteľný vzhľadom na hustotu z objektu o_i pre $i \in \{1, 2, \dots, t\}$, ε a m .
- Objekt o je **husto spojený** s objektom o' s ohľadom na ε a m , ak existuje objekt o^* taký, že o a o' sú dosiahnuteľné vzhľadom na hustotu s ohľadom na ε a m .

Iba vnútorné objekty sú navzájom dosiahnuteľné vzhľadom na hustotu. Husté spojenie je symetrická relácia. **Zhluk založený na hustote** je množina husto spojených objektov, ktoré sú maximálne z pohľadu dosiahnuteľnosti vzhľadom na hustotu.

Algoritmus zhlukovania založený na hustote môžeme vo všeobecnosti zapísať nasledovne:

Algoritmus 9.4 *Algoritmus zhlukovania založený na hustote*

Nech X je matica údajov, ktorá reprezentuje hodnoty n objektov $\{o_1, o_2, \dots, o_n\}$ v p atribútoch $\{a_1, a_2, \dots, a_p\}$. Nech ε je okolie objektov a m je minimálny počet objektov v ε okolí vnútorného objektu.

begin

pre každý objekt o_i , $i \in \{1, 2, \dots, n\}$:

$m_i :=$ počet objektov v ε -okolí objektu o_i ;

ak $m_i > m$:

vytvor nový zhluk s vnútorným objektom o_i ;

repeat

pre každý vnútorný objekt o_i

pridaj do zhluku C_i všetky priamo dosiahnuteľné objekty o z o_i ;

until *žiadne objekty o nie sú pridané do niektorého zhluku;*

end

Vhodným zvolením okolia ε a minimálneho počtu objektov v okolí vnútorného bodu m , je možné pomocou metód zhlukovania založených na hustote nájsť tvary zhlukov, ktoré sú odlišné od tvarov zhlukov generovaných segmentačnými a hierarchickými metódami. Prezentovaný algoritmus je možné nájsť v literatúre aj pod názvom DBSCAN (Density Based Spatial Clustering of Applications with Noise). Výhodou algoritmu je aj schopnosť pracovať so zašumenými údajmi.

9.5 Požiadavky na zhlukovanie

Potenciálne aplikácie zhlukovania v praxi si vyžadujú špeciálne požiadavky. Napríklad v oblasti strojové učenia a dolovania údajov môžeme uvažovať:

- **Schopnosť pracovať s rôznym typom atribútov:** Mnoho algoritmov je navrhnutých na zhlukovanie objektov s numerickými atribútmi. Aplikácie môžu vyžadovať vytváranie zhlukov aj pre iné typy údajov, akými sú napríklad binárne, kategorické (nominálne) a ordinálne údaje, alebo kombinácia týchto typov údajov.
- **Škálovateľnosť:** Väčšina algoritmov zhlukovania má dobré výsledky na malých množinách údajov s niekoľkými stovkami objektov. Veľká databáza môže obsahovať milióny objektov. Zhlukovanie na vzorke daného veľkého súboru údajov tak môže viesť k skresleným výsledkom. Je potrebné teda využívať vysoko škálovateľné algoritmy zhlukovania.
- **Znalosť aplikačnej domény na určenie vstupných parametrov:** Niektoré algoritmy zhlukovania vyžadujú, aby používatelia zadali určité vstupné parametre zhlukovania (napríklad počet požadovaných zhlukov). Výsledky zhlukovania môžu byť citlivé na vstupné parametre. Parametre je často ťažké určiť, najmä pre množiny údajov obsahujúce vysokodimenzionálne objekty.
- **Nájdenie zhlukov ľubovoľného tvaru:** Viacero algoritmov zhlukovania je založených na euklidovských alebo manhattanských mierach. Metódy s meraním vzdialenosti majú tendenciu nájsť zhluky s podobnou veľkosťou a hustotou. Zhluk však môže mať akýkoľvek tvar, preto je vhodné poznať aj algoritmy na detekciu zhlukov ľubovoľného tvaru.
- **Schopnosť práce so zašumenými údajmi:** Databázy v reálnom svete obsahujú anomálne hodnoty alebo chýbajúce, neznáme, či chybné údaje. Niektoré algoritmy zhlukovania sú na tieto údaje citlivé a môžu viesť k zoskupeniam nízkej kvality.
- **Dimenzionalita údajov:** Databáza alebo dátový sklad môže obsahovať niekoľko dimenzií. Mnoho algoritmov je vhodných pri spracovávaní nízkodimenzionálnych údajov, ktoré obsahujú dva až tri rozmery. Ľudské oči dokážu posúdiť kvalitu zhlukovania až do troch dimenzií. Hľadanie zhlukov pre objekty vo vysoko rozmernom priestore je náročné, prípadne takýchto objektov môže byť málo a výsledok môže byť skreslený.
- **Zhlukovanie založené na obmedzeniach:** Je možné, že aplikácie v reálnom svete budú musieť vykonať zhlukovanie vzhľadom na rôzne obmedzenia. Napríklad, ak úlohou je výber vhodných miest pre daný počet nových bankomatov v meste, pri riešení môžeme zoskupovať domácnosti vzhľadom na rieky, diaľničné siete, typ a počet zákazníkov pre jeden zhluk. Nájsť skupiny s požadovaným správaním v zhlukoch, ktoré vyhovujú týmto obmedzeniam, môže byť náročnou úlohou.

- **Interpretovateľnosť a použiteľnosť:** Používatelia očakávajú, že ich zhľuky budú interpretovateľné, zrozumiteľné a použiteľné. To znamená, že vytváranie zhľukov bude potrebné spájať s konkrétnym významom interpretácie a aplikácie.
- **Prírastkové zhľukovanie, či necitlivosť na poradie vstupných údajov:** Niektoré zhľukovacie algoritmy nemôžu obsahovať novo-vložené údaje (t. j. aktualizáciu databázy) do existujúcich výsledkov. Je potrebné použiť nové zhľukovanie od začiatku. Niektoré algoritmy zhľukovania sú citlivé na poradie vstupných údajov. Vzhľadom na vstupnú množinu objektov môže takýto algoritmus určiť rôzne zoskupenia v závislosti od poradia prezentácie vstupných objektov. Je dôležité zaoberať sa aj prírastkovými algoritmi zhľukovania a algoritmi, ktoré sú necitlivé na poradie vstupu.

Kapitola 10

Bayesovský prístup ku strojovému učeniu

10.1 Úvod

- Pridanie šumu do dát umožňuje dosiahnuť lepšie výsledky v predikcii.
- Sme schopní kvantifikovať neurčitost', ktorá je prítomná v odhadoch našich parametrov, a z toho vyplývajúcej predikcie.
- Od idey, že máme neurčitost' v odhadoch našich parametrov, je len krok k uvažovaniu daných parametrov ako náhodných premenných.
- Bayesovské metódy sa tu stávajú veľmi dôležitými. Urobíme úvod do ich použitia podľa [6].

10.2 Hra Hod mincou

Idete cez trh a narazíte na stánok, kde sa zákazníci zúčastňujú hry *Hod mincou*. Pravidlá hry sú jednoduché:

- Zákazník hádže mincou 10 krát
- Ak padne hlava 6x alebo menej krát, zákazník dostáva naspäť dvojnásobok vložených peňazí (1 EUR stávka). Inak zákazník o peniaze príde.

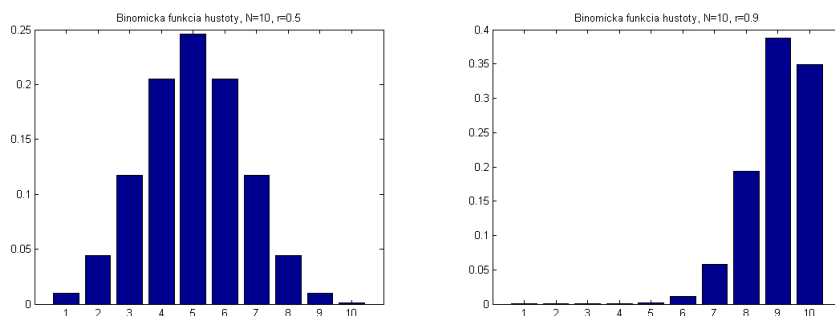
Úvahy o možnosti vyhrať:

Predpokladajme, že pravdepodobnosť, že padne hlava je r . Potom, pravdepodobnosť, že hlava padne y krát z N hodov, je

$$P(Y = y) = \binom{N}{y} r^y (1 - r)^{N-y} \quad (10.1)$$

Pre $P(Y \leq 6)$ dostávame

$$P(Y \leq 6) = 1 - P(Y > 6) = 1 - (P(Y = 7) + P(Y = 8) + P(Y = 9) + P(Y = 10))$$



Obr. 10.1: Prvý panel predstavuje funkciu hustoty, ak minca je spravodlivá, t. j. $r = 0.5$, druhý panel predstavuje prípad, keď viac padá hlava, t. j. $r = 0.9$ (pri 10 hodoch).

$$= 1 - (0.1172 + 0.0439 + 0.0098 + 0.001) = 0.8281$$

čo vyzerá celkom zaujímavo. Je možné tiež vypočítať očakávanú finančnú sumu z hry. Očakávaná hodnota funkcie $f(X)$ náhodnej premennej X je počítaná

$$\mathbf{E}_{P(x)}\{f(X)\} = \sum_x f(x)P(x)$$

suma je cez všetky možné hodnoty náhodnej premennej X .

Nech X bude náhodná premenná, ktorá nadobúda hodnotu 1 pri výhre (padne hlava 6 a menej krát) a 0 inak: $P(X = 1) = P(Y \leq 6)$.

Ak $X = 1$, je to výhra a zákazník dostáva 2 EUR, teda $f(1) = 2$. Ak prehrá, nedostáva nič, teda $f(0) = 0$.

Očakávaný zisk je

$$f(1)P(X = 1) + f(0)P(X = 0) = 2 * P(Y \leq 6) + 0 * P(Y > 6) = 1.6562$$

Teda v priemere vyhráva $1.6562 - 1 = 0.6562$ centov na hre.

Zdá sa, že hrať je OK. Ale je v tom háčik. V skutočnosti výhier nie je až tak veľa. Ďalším problémovým parametrom v hre môže byť r .

Pozorovanie

Možno predpoklady podliehajúce výpočtom sú nesprávne. Ide o nasledujúce:

- Počet hláv môže byť modelovaný ako náhodná premenná s binomickým rozdelením a pravdepodobnosť, že padne hlava v ľubovoľnom hode je r .
- Minca je spravodlivá - pravdepodobnosť oboch strán je $r = 0.5$.

Mohli by sme r uvažovať ako parameter a použiť ho pri fitovaní dát.

Úvahy o
postupnosti
hodov

Sú traja ľudia, ktorí hrajú túto hru. Predpokladajme, že prvý hráč dostane postupnosť: H, O, H, H, H, H, H, H, H.

Je možné vypočítať maximálnu hodnotu vierohodnosti (likelihood) na základe binomického rozdelenia pre r takto:

$$P(Y = y|r, N) = \binom{N}{y} r^y (1-r)^{N-y} \quad (10.2)$$

Zlogaritmovaním

$$L = \log(P(Y = y|r, N)) = \log \binom{N}{y} + y \log r + (N-y) \log(1-r) \quad (10.3)$$

Zderivujeme podľa r

$$\begin{aligned} \frac{\partial L}{\partial r} &= \frac{y}{r} - \frac{N-y}{1-r} = 0 \\ r &= \frac{y}{N} \end{aligned}$$

Ak $y = 9$ a $N = 10$ dostávame $r = 0.9$ a $P(Y \leq 6) = 0.0128$.

Očakávaný zisk $2 * P(Y \leq 6) + 0 * P(Y > 6) = 0.0256$, $0.0256 - 1 = -0.9744$ na jednu hru. Teda strata 97 centov.

$P(Y \leq 6) = 0.0128$ znamená, že len jeden zákazník zo 100 vyhrá. Toto by mohlo odradiť ďalších hráčov.

10.3 Bayesovský prístup k hre

10.3.1 Bayesova veta - všeobecný tvar

Predpokladajme, že náhodné javy B_i , $1 \leq i \leq k$ sú navzájom nezávislé javy a v každom pokuse nastáva práve jeden z nich, takže súčet ich pravdepodobností je rovný 1, t. j. $\sum_{i=1}^k P(B_i) = 1$

Ak poznáme podmienené pravdepodobnosti $P(A|B_i)$ javu A pri podmienke B_i pre všetky i a *apriórne* pravdepodobnosti javov $P(B_i)$, potom je možné pravdepodobnosť javu $P(B_j|A)$ vypočítať podľa Bayesovho vzorca:

$$P(B_j|A) = \frac{P(A|B_j) * P(B_j)}{P(A)} = \frac{P(A|B_j) * P(B_j)}{\sum_{i=1}^k P(A|B_i) * P(B_i)} \quad (10.4)$$

- $P(B_j|A)$ je posteriórna pravdepodobnosť javu B_j ;
- $P(A|B_j)$ je dôveryhodnosť javu B_j ;
- $P(B_j)$ je apriórna pravdepodobnosť javu B_j ;
- $\sum_{i=1}^k P(A|B_i) * P(B_i)$ je marginálna pravdepodobnosť dát.

V prípade spojitých hodnôt premenných pravdepodobosť je nahradená funkciou hustoty pravdepodobnosti. Potom pre hypotézu B_j dostávame

$$f(B_j|A) = \frac{f(A|B_j) * f(B_j)}{P(A)} = \frac{f(A|B_j) * f(B_j)}{\int_{B_i} f(A|B_i) * f(B_i) dB_i} \quad (10.5)$$

- $f(B_j|A)$ je posteriórna hustota pravdepodobnosti pre B_j ;
- $f(A|B_j)$ je dôveryhodnosť hypotézy B_j ;
- $f(B_j)$ je apriórna hustota pravdepodobnosti pre B_j ;
- $\int_{B_i} f(A|B_i) * f(B_i) dB_i$ je marginálna pravdepodobnosť dát.

10.3.2 Bayesova veta v hre s mincou

Hodnoty pravdepodobnosti r (že padne hlava) sa menia v čase a mohli by sme sa na ňu dívať ako na náhodnú premennú R .

Nezávisle od dĺžky postupnosti hodov vždy bude nejaká neurčitnosť v r - budeme ju brať ako náhodnú premennú s asociovaným rozdelením. Toto nám pomôže merať a porozumieť túto neurčitnosť.

Nech Y_N je náhodná premenná vyjadrujúca počet padnutých hláv v N hodoch. Chceme poznať pravdepodobnostné rozdelenie pre r podmienené hodnotami Y_N , teda $p(r|y_N)$.

Pri danom rozdelení bude možné vypočítať očakávanú pravdepodobnosť výhry tak, že vezmeme očakávania $P(Y_{nov} \leq 6|r)$ vzhľadom na $p(r|y_N)$:

$$P(Y_{nov} \leq 6|y_N) = \int P(Y_{nov} \leq 6|r)p(r|y_N)dr,$$

kde Y_{nov} je náhodná premenná opisujúca počet padnutých hláv v budúcej množine desiatich hodov.

Podľa Bayesovho pravidla dostaneme:

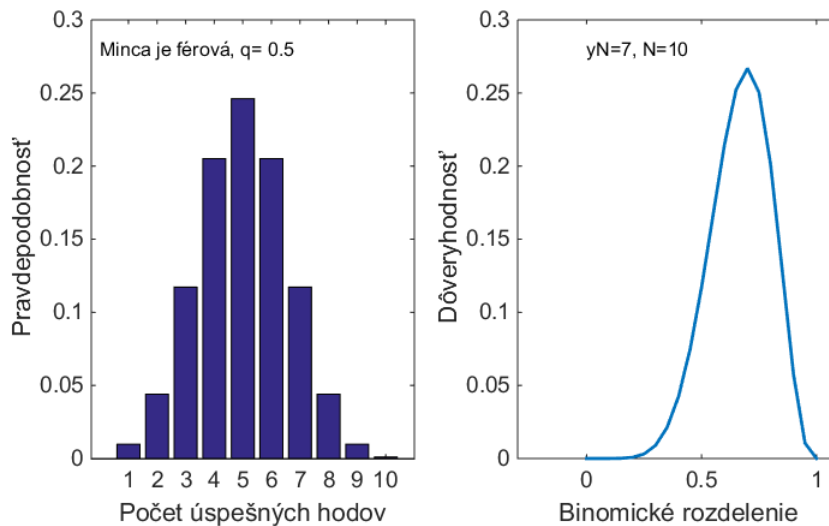
$$p(r|y_N) = \frac{P(y_N|r)p(r)}{P(y_N)}.$$

$p(y_N|r)$ – pravdepodobnostné rozdelenie nad počtom hláv v N nezávislých hodoch za predpokladu, že pravdepodobnosť, že padne hlava v jednom hode, je r .

Príklad: Binomické rozdelenie pravdepodobnosti a dôveryhodnosť:

10.3.3 Pravdepodobnosť vs. dôveryhodnosť (likelihood)

- Binomické rozdelenie pravdepodobnosti je diskrétné (11 možných experimentálnych výsledkov podľa počtu hodených hláv).
- Dôveryhodnosť je spojitá, parameter r nadobúda hodnoty medzi 0 a 1.



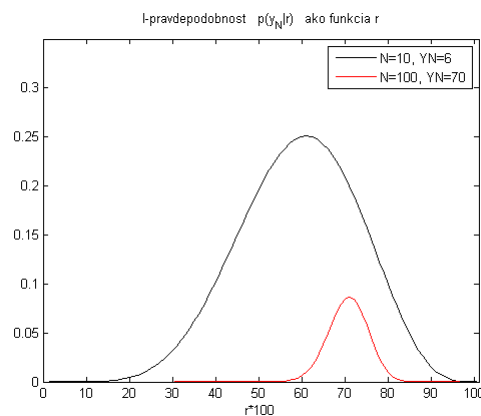
Obr. 10.2: Pri férovej minci a 7 úspešných hodoch z 10.

- Pravdepodobnosti v prvom paneli dávajú súčet 1.
- Dôveryhodnosť po zintegrovaní bude menej ako 1.

Dôveryhodnosť (likelihood), význam $P(y_N|r)$ - ako je dôveryhodné, že by sme mohli pozorovať naše dáta (y_N) pri určitej hodnote r (náš model).

V našom prípade to je binomické rozdelenie. Táto hodnota bude vysoká, ak r by mohlo nadobudnúť dosiahnuteľný výsledok y_N a nízka, ak je výsledok veľmi nepravdepodobný.

Dva scenáre pre rôzne nastavenia N a y_N



Obr. 10.3: Príklad dvoch scenárov pre vyjadrenie $p(y_N|r)$ ako funkcie r .

Dve dôležité vlastnosti dôveryhodnosti:

- Nie je to hustota pravdepodobnosti. Ak by bola, obsah plochy pod krivkou by bol rovný 1.
- Uvedené dva príklady sa líšia v tom, koľko nám povedia o r . V prvom prípade dôveryhodnosť je nenulová pre veľký interval pre r . V druhom prípade je tento interval redukovaný. Máme viac dát, a teda viac vieme (presnejšie) o r .

Apriórne rozdelenie (prior distribution), $p(r)$ - umožňuje vyjadriť nejakú vieru v r predtým, než niečo vieme o dátach.

Uvažujme nasledujúce 3 prípady:

1. Nevieme nič o hádzanej minci ani o vlastníčkovi stánku.
2. Myslíme si, že minca je fér.
3. Myslíme si, že minca je vytvorená tak, že viac padá hlava.

Každý z týchto predpokladov vieme zakódovať iným apriórnym rozdelením.

r môže nadobúdať hodnoty od 0 po 1, preto musí byť modelované ako spojitá náhodná premenná.

Nevyberieme žiadny scenár, jednoducho budeme zatiaľ uvažovať $p(r|y_N)$.

Pre tento prípad je vhodná **beta funkcia hustoty** definovaná pre náhodnú premennú R ako:

$$p(r) = \frac{\Gamma(\alpha + \beta)}{\Gamma(\alpha)\Gamma(\beta)} r^{\alpha-1} (1-r)^{\beta-1} \quad (10.6)$$

Γ je známa funkcia, ktorá zaručí, že rozdelenie je normalizované. Parametre $\alpha, \beta \geq 0$ riadia tvar výsledného rozdelenia.

$$\Gamma(x) = \int_0^\infty t^{x-1} e^{-t} dt, \quad \text{pre } x > 1$$

$$\frac{\Gamma(\alpha + \beta)}{\Gamma(\alpha)\Gamma(\beta)} = \int_{r=0}^{r=1} r^{\alpha-1} (1-r)^{\beta-1} dr \quad (10.7)$$

zaručuje, že

$$\int_{r=0}^{r=1} \frac{\Gamma(\alpha + \beta)}{\Gamma(\alpha)\Gamma(\beta)} r^{\alpha-1} (1-r)^{\beta-1} dr = 1$$

V uvedených troch prípadoch zvolíme:

- Nevieme nič: $\alpha = 1, \beta = 1$
- Férová minca: $\alpha = 50, \beta = 50$
- Neférová minca: $\alpha = 5, \beta = 1$

Marginálne rozdelenie (marginal distribution) pre $y_N, P(y_N)$ – pôsobí ako normalizujúca konštanta, ktorá zaručí, že $p(r|y_N)$ je vlastne definovaná hustota. Počíta sa ako združená hustota $p(y_N|r)$ cez r :

$$P(y_N) = \int_{r=0}^{r=1} p(y_N, r) dr \quad (10.8)$$

čo je možné vyjadriť

$$P(y_N) = \int_{r=0}^{r=1} P(y_N|r)p(r) dr \quad (10.9)$$

$P(y_N)$ je známa ako marginálna dôveryhodnosť dát, y_N , spriemernená cez všetky hodnoty parametrov.

Posteriórne rozdelenie (posterior distribution), $p(r|y_N)$ - o toto sa zaujíname.

Je to výsledok úpravy apriórnej pravdepodobnosti $p(r)$ vo svetle nových poznatkov o y_N .

- Tvar rozdelenia je zaujímavý - vyjadruje koľko informácie máme o r po kombinácii s tým, čo sme vedeli predtým a čo sme videli (dôveryhodnosť).
- Je to rozdelenie, preto nám poskytuje indikáciu o úrovni neurčitosti, stále máme v r informáciu o nejakých pozorovaných dátach.
- Posteriórne rozdelenie je možné použiť na výpočet očakávaní (výhry).

Posteriórnu hustotu vieme použiť na výpočet očakávanej hodnoty. Napríklad

$$\mathbf{E}_{p(r|y_N)}\{P(Y_{10} \leq 6)\} = \int_{r=0}^{r=1} P(Y_{10} \leq 6|r) p(r|y_N) dr$$

je očakávaná hodnota pravdepodobnosti, že vyhráme. Je potrebné brať do úvahy pozorované dáta, naša apriórna viera a neurčitosť zostáva. Treba sa rozhodnúť, či budeme hrať.

10.3.4 Dvojice dôveryhodnosť-apriórna pravdepodobnosť

Konjugované dvojice **apriórna – dôveryhodnosť** umožňujú použiť posteriórne rozdelenie v rovnakom tvare ako apriórne.

Apriórna pravdepodobnosť – dôveryhodnosť, niektoré konjugované dvojice:

- Gauss – Gauss
- Beta – Binomické
- Gamma – Gauss
- Dirichlet – Multinomial

10.3.5 Exaktné posteriórne rozdelenie

Beta rozdelenie je vhodný výber pre apriórnu pravdepodobnosť, ak dôveryhodnosť je binomické rozdelenie. Vypočítame ho presne.

Beta rozdelenie je známe ako konjugované rozdelenie k binomickej dôveryhodnosti. Použitie konjugovaného vzťahu uľahčí výpočty.

Ak v Bayesovom vzťahu vynecháme $P(y_N)$, dostaneme

$$p(r|y_N) \propto P(y_N|r)p(r)$$

Nahradením výrazov na pravej strane binomickým a beta rozd.

$$p(r|y_N) \propto \left[\binom{N}{y_N} r^{y_N} (1-r)^{N-y_N} \right] \times \left[\frac{\Gamma(\alpha+\beta)}{\Gamma(\alpha)\Gamma(\beta)} r^{\alpha-1} (1-r)^{\beta-1} \right] \quad (10.10)$$

Pretože apriórne a beta rozd. sú konjugované, vieme, že $p(r|y_N)$ má byť beta hustota. Beta hustota s parametrami δ a γ a konštantou K má nasledujúci tvar

$$p(r) = Kr^{\delta-1}(1-r)^{\gamma-1}$$

Vhodnými úpravami dostaneme

$$\begin{aligned} p(r|y_N) &\propto \left[\binom{N}{y_N} \frac{\Gamma(\alpha+\beta)}{\Gamma(\alpha)\Gamma(\beta)} \right] \times [r^{y_N} r^{\alpha-1} (1-r)^{N-y_N} (1-r)^{\beta-1}] \\ p(r|y_N) &\propto r^{y_N+\alpha-1} (1-r)^{N-y_N+\beta-1} \\ p(r|y_N) &\propto r^{\delta-1} (1-r)^{\gamma-1}, \quad \delta = y_N + \alpha, \gamma = N - y_N + \beta \end{aligned}$$

Preto

$$p(r|y_N) = \frac{\Gamma(\alpha+\beta+N)}{\Gamma(\alpha+y_N)\Gamma(\beta+N-y_N)} r^{\alpha+y_N-1} (1-r)^{\beta+N-y_N} \quad (10.11)$$

Dostali sme posteriórnu hustotu pre r na základe apriórnej $p(r)$ a dát y_N .

Pripomeňme si, ako boli vypočítané - sčítaním počtov hodených hláv (y_n) do prvého apriórneho parametra α a počtov opaku mince ($N - y_N$) do β .

Toto nám umožní získať intuíciu ohľadom α a β . Ak uvažujeme scenár korektnej mince, zo 100 hodov je 50 hlava a 50 opak mince. Mohli by sme použiť $\alpha = \beta = 50$,

**Poznámka k
trom
scenárom**

Dôslednú analýzu všetkých scenárov je možné nájsť v [\[6\]](#).

Chceme tu pripomenúť, že vo všetkých troch prípadoch dostaneme iné hodnoty očakávanej pravdepodobnosti pre výhru. Naše úvahy sú potvrdené tým, že najväčšia pravdepodobnosť je pri spravodlivej minci a najmenšia, ak hlava padá častejšie.

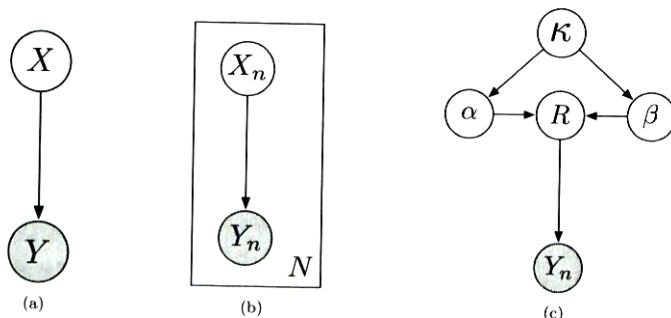
Dôležité je to, že Baeysovský prístup nám umožňuje kombinovať pozorované dáta (hody mincou) s nejakými predchádzajúcimi znalosťami (jeden zo scenárov) principiálnym spôsobom. Posteriórna hustota explicitne modeluje neurčitost', ktorá zostáva v r v každom stave a môže byť použitá na predikcie.

Pridávanie ďalších údajov (hodov) mení scenár 1, lebo viac vieme o hodoch.

10.4 Grafické modely

Grafický model je graf, v ktorom vrcholy zodpovedajú náhodným premenným a orientované hrany závislostiam medzi náhodnými premennými.

Napríklad, jedna náhodná premenná reprezentuje hod mincou (X) a druhá predstavuje to, čo je výsledkom hodu (Y). Model je definovaný pomocou podmieneného rozdelenia $P(Y = y|X = x)$. Podmienená premenná je v šedom vrchole a k nej smeruje orientácia hrany. Nasledujúci obr. (a). Šedý vrchol vyjadruje pozorovanú premennú.



Obr. 10.4: Grafické modely vyjadrujúce vzťahy medzi náhodnými premennými, [6]

Ak si predstavíme, že hod mincou sa opakuje N -krát, máme $2N$ náhodných premenných X_1, \dots, X_N a Y_1, \dots, Y_N . Toto zakreslíme pomocou obdĺžnika, v ktorom sa nachádza to, čo sa opakuje N -krát, na obr. (b). A je evidentné, že Y_n je podmienená od X_n .

Grafická reprezentácia nášho modelu pre hod mincou je na obr. (c). Má jednu pozorovanú náhodnú premennú, ktorá reprezentuje počet hodených hláv v N hodoch y_N . Toto je podmienené náhodnou premennou R , ktorá závisí od náhodných premenných α a β . A nakoniec α a β závisia od nejakých hyperparametrov κ .

10.5 Bayesov prístup k problému predikcie rekordov na OH

Cieľ je predikovať výsledok na nasledujúcich OH z daných dát. K tomu budeme potrebovať:

- definovať apriórnu pravdepodobnosť a dôveryhodnosť
- vypočítať posteriórnu hustotu nad parametrami nášho modelu (r v prípade mince)
- keď máme posteriórnu hustotu, vieme predikovať.

Bayesov prístup - Model

Budeme používať polynomiálny model s Gaussovým šumom:

$$t_n = w_0 + w_1 x_n + w_2 x_n^2 + \dots + w_K x_n^K + \epsilon_n, \quad \epsilon \sim \mathcal{N}(0, \sigma^2).$$

Vo vektorovom tvare

$$t_n = \mathbf{w}^T \mathbf{x}_n + \epsilon_n, \quad \mathbf{w} = [w_0, \dots, w_K]^T, \quad \mathbf{x}_n = [1, x_n, x_n^2, \dots, x_n^K]^T.$$

V maticovom tvare pre všetky dané dáta

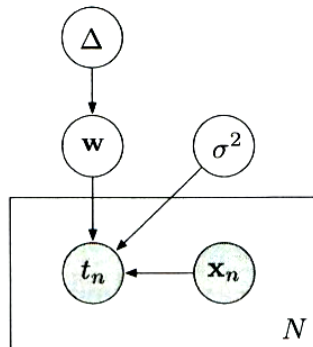
$$\mathbf{t} = \mathbf{X}\mathbf{w} + \epsilon, \quad \mathbf{X} = [\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_N]^T, \quad \mathbf{t} = [t_1, t_2, \dots, t_N]^T$$

Kvôli zjednodušeniu výpočtov budeme predpokladať, že poznáme σ^2 .

Bayesovo pravidlo

$$p(\mathbf{w}|\mathbf{t}, \mathbf{X}, \sigma^2, \Delta) = \frac{p(\mathbf{t}|\mathbf{w}, \mathbf{X}, \sigma^2, \Delta)p(\mathbf{w}|\Delta)}{p(\mathbf{t}|\mathbf{X}, \sigma^2, \Delta)} = \frac{p(\mathbf{t}|\mathbf{w}, \mathbf{X}, \sigma^2)p(\mathbf{w}|\Delta)}{p(\mathbf{t}|\mathbf{X}, \sigma^2, \Delta)}$$

kde Δ predstavuje parametre potrebné pri definovaní apriórnej pravdepodobnosti pre \mathbf{w} (uvedieme neskôr).



Obr. 10.5: Grafický model pre Bayesov prístup - rekordy na OH, [6].

10.5.1 Bayesov prístup k predikcii rekordov

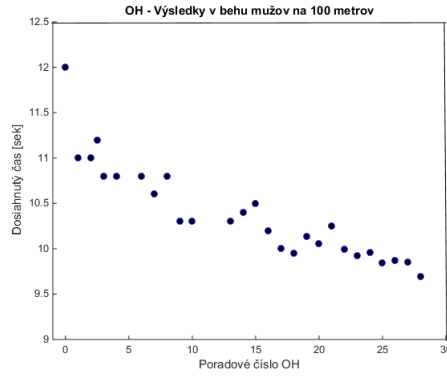
Potrebuje určiť

- **apriórnu pravdepodobnosť** – $p(\mathbf{w}|\Delta)$

Budeme používať Gaussovú apriórnu hustotu so zvolenými parametrami, ktoré označíme μ_0, Σ_0

$$p(\mathbf{w}|\mu_0, \Sigma_0) = \mathcal{N}(\mu_0, \Sigma_0)$$

- **dôveryhodnosť** – $p(\mathbf{t}|\mathbf{w}, \mathbf{X}, \sigma^2)$



Obr. 10.6: Výsledky v behu mužov na 100 m na OH

Náš model je $\mathbf{t} = \mathbf{X}\mathbf{w} + \epsilon$, kde $\epsilon \sim \mathcal{N}(\mathbf{0}, \sigma^2 \mathbf{I}_N)$. Teda

$$p(\mathbf{t}|\mathbf{w}, \mathbf{X}, \sigma^2) = \mathcal{N}(\mathbf{X}\mathbf{w}, \sigma^2 \mathbf{I}_n)$$

N -dimenzionálna Gaussova hustota s priemerom $\mathbf{X}\mathbf{w}$ a varianciou $\sigma^2 \mathbf{I}_N$

- **posteriórnu pravdepodobnosť** Pretože vieme, že posteriórna pravdepodobnosť bude Gaussova, zatiaľ si nebudeme šímať marginálnu pravdepodobnosť, preto

$$\begin{aligned} p(\mathbf{w}|\mathbf{t}, \mathbf{X}, \sigma^2) &\propto p(\mathbf{t}|\mathbf{w}, \mathbf{X}, \sigma^2)p(\mathbf{w}|\mu_0, \Sigma_0) & (10.12) \\ &= \frac{1}{(2\pi)^{N/2}|\sigma^2 \mathbf{I}|^{1/2}} \exp\left(-\frac{1}{2}(\mathbf{t} - \mathbf{X}\mathbf{w})^T (\sigma^2 \mathbf{I})^{-1} (\mathbf{t} - \mathbf{X}\mathbf{w})\right) \end{aligned}$$

Tiež tu použijeme Gaussovo rozdelenie

$$\begin{aligned} p(\mathbf{w}|\mathbf{t}, \mathbf{X}, \sigma^2) &= \mathcal{N}(\mu_{\mathbf{w}}, \Sigma_{\mathbf{w}}) & (10.13) \\ &\propto \exp\left(-\frac{1}{2}(\mathbf{w} - \mu_{\mathbf{w}})^T \Sigma_{\mathbf{w}}^{-1} (\mathbf{w} - \mu_{\mathbf{w}})\right) \end{aligned}$$

a na základe úprav a porovnania (10.12) a (10.13) dostaneme

$$p(\mathbf{w}|\mathbf{t}, \mathbf{X}, \sigma^2) = \mathcal{N}(\mu_{\mathbf{w}}, \Sigma_{\mathbf{w}}) \quad (10.14)$$

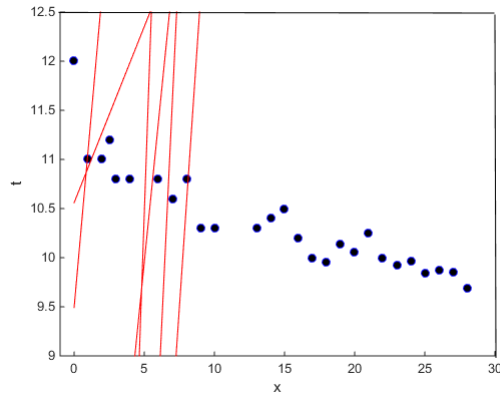
$$\Sigma_{\mathbf{w}} = \left(\frac{1}{\sigma^2} \mathbf{X}^T \mathbf{X} + \Sigma_0^{-1}\right)^{-1}, \quad (10.15)$$

$$\mu_{\mathbf{w}} = \Sigma_{\mathbf{w}} \left(\frac{1}{\sigma^2} \mathbf{X}^T \mathbf{t} + \Sigma_0^{-1} \mu_0\right) \quad (10.16)$$

10.5.2 Prípád lineárneho modelu $t_n = w_0 + w_1 x_n$

Predpokladajme, že časy x_n v behu na 100 m mužov na OH sú znázornené na obrázku 10.6

Vstupy sú teda $\mathbf{x}_n = [1, x_n]^T$. Predpokladáme, že reálne nič nevieme o tom, ako vybrať parametre, tak zvolíme $\mu_0 = [0, 0]^T$. Pre kovarianciu zvolíme $\Sigma_0 =$

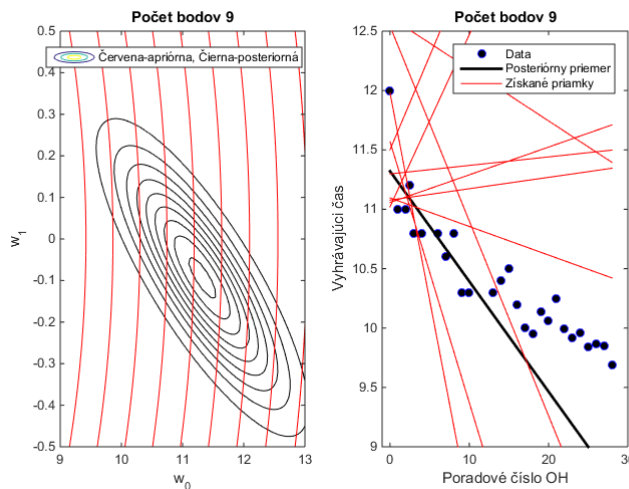


Obr. 10.7: Situácia po spracovaní prvého bodu.

$((100, 0), (0, 5))$, $\sigma^2 = 10$, pretože sme už videli odhad pre w_0 , ktorý bol väčší ako pre w_1 .

Na základe vzťahov (10.14), (10.15) a (10.16) pre $\sigma^2 = 10$ a prvý bod dát $\mathbf{X} = [1, 0]$, $\mathbf{T} = 12$ vieme vypočítať posteriórne rozdelenie. Jeden bod nám niečo povie o hodnote w_0 , ale takmer nič o hodnote w_1 . Vieme vygenerovať nejaké priamky, ale tie sú od dát veľmi vzdialené, vidíme to na obrázku 10.7.

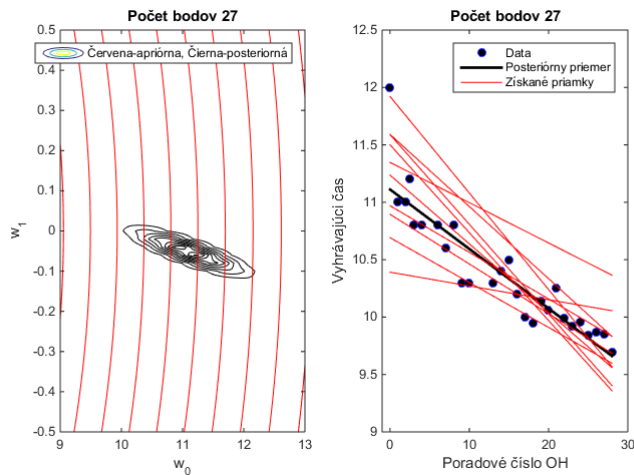
Ak vezmeme 9 bodov, situácia je lepšia tak, ako je to na obrázku ... Všetky body nám dajú výsledky znázornené na obrázku 10.8.



Obr. 10.8: Apriórna a posteriórna pravdepodobnosť po 9 OH

Pri použití všetkých bodov dostaneme výsledky znázornené na obrázku 10.9.

V obrázkoch 10.8 a 10.9 pozorujeme, že posteriórna pravdepodobnosť sa stala viac kondenzovaná, čo znamená, že máme viac poznatkov pre \mathbf{w} . Pozorujeme napríklad, že keď zvýšime w_0 , musíme znížiť w_1 . V apriórnej pravdepodobnosti sme predpokladali, že sú nezávislé.

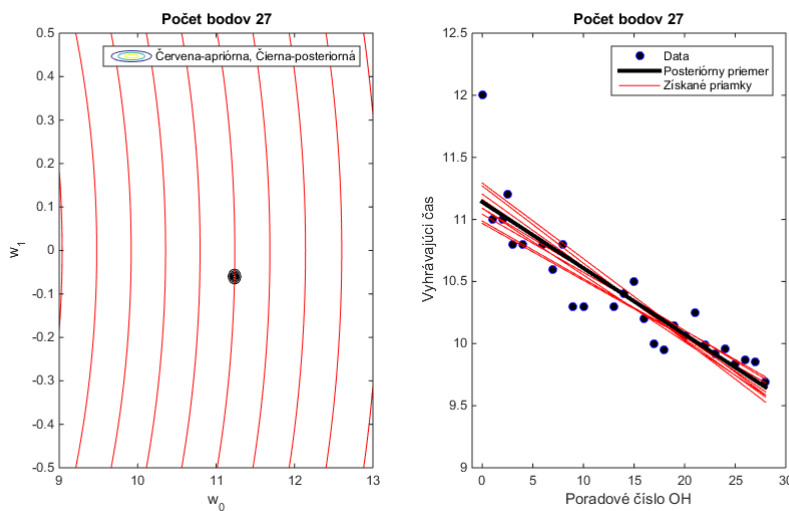


Obr. 10.9: Apriórna a posteriórna pravdepodobnosť po 27 OH

V obrázku 10.9 pozorujeme ešte dosť veľkú variabilitu, čo je spôsobené tým, že $\sigma^2 = 2$. Zmenšením σ je možné docieľiť menšiu variabilitu pre w_0 a w_1 .

10.5.3 Ako urobiť predikciu

Ak zvolíme $\sigma^2 = 0.05$, dostaneme výsledky znázornené na obrázku 10.10



Obr. 10.10: Apriórna a posteriórna pravdepodobnosť po 27 OH

Predpokladáme nové pozorovanie t_{nove} , teda nás zaujíma hustota

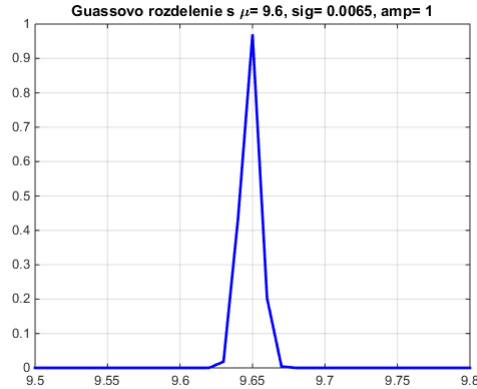
$$p(t_{nove} | \mathbf{x}_{nove}, \mathbf{X}, \mathbf{t}, \sigma^2). \quad (10.17)$$

Výraz neobsahuje \mathbf{w} , ale naše očakávania sú zamerané na posteriórnu pravdepodobnosť $p(\mathbf{w} | \mathbf{t}, \mathbf{X}, \sigma^2)$. Potrebujeme vypočítať

$$\begin{aligned}
p(t_{nove}|\mathbf{x}_{nove}, \mathbf{X}, \mathbf{t}, \sigma^2) &= \mathbf{E}_{p(\mathbf{w}|\mathbf{t}, \mathbf{X}, \sigma^2)}\{p(t_{nove}|\mathbf{x}_{nove}, \mathbf{X}, \sigma^2)\} \\
&= \int p(t_{nove}|\mathbf{x}_{nove}, \mathbf{X}, \sigma^2)p(\mathbf{w}|\mathbf{t}, \mathbf{X}, \sigma^2) d\mathbf{w}.
\end{aligned} \tag{10.18}$$

$p(t_{nove}|\mathbf{x}_{nove}, \mathbf{X}, \sigma^2)$ je v našom modeli definované ako súčin \mathbf{x}_{nove} a \mathbf{w} s pridaním nejakého gausssovského šumu.

$$p(t_{nove}|\mathbf{x}_{nove}, \mathbf{X}, \sigma^2) = \mathcal{N}(\mathbf{x}_{nove}^T \mathbf{w}, \sigma^2).$$



Obr. 10.11: Predikcia na 28 OH

Výsledok teda očakávame v tvare Gaussovej funkcie. Teda

$$p(t_{nove}|\mathbf{x}_{nove}, \mathbf{X}, \mathbf{t}, \sigma^2) = \mathcal{N}(\mathbf{x}_{nove}^T \mu_w, \sigma^2 + \mathbf{x}_{nove}^T \Sigma_w \mathbf{x}_{nove})$$

Pre situáciu na obrázku [10.10](#) dostávame

$$p(t_{nove}|\mathbf{x}_{nove}, \mathbf{X}, \mathbf{t}, \sigma^2) = \mathcal{N}(9.6484, 0.0065)$$

Prediktívne rozdelenie je na obrázku [10.11](#). Najviac očávaný je čas 9.65, ale je nenulová pravdepodobnosť aj horších časov.

$$f(x, y) = A \exp \left(- \left(\frac{(x - x_o)^2}{2\sigma_X^2} + \frac{(y - y_o)^2}{2\sigma_Y^2} \right) \right).$$

Koeficient A je amplitúda, (x_o, y_o) je centrum.

**Príklad 2-
dimensionálnej
Gaussovej
funkcie:**

10.5.4 Úlohy

1. Odvodili sme pre výraz pre Gaussovu posteriórnu pravdepodobnosť pre lineárny model (rekordy na OH). Ak dosadíme za $\mu_0 = [0, 0, \dots, 0]^T$, vidíme určitú podobnosť priemerom posteriórnej pravdepodobnosti

$$\mu_w = \frac{1}{\sigma^2} \left(\frac{1}{\sigma^2} \mathbf{X}^T \mathbf{X} + \Sigma_0^{-1} \right)^{-1} \mathbf{X}^T \mathbf{t} \quad (10.19)$$

regularizovaným riešením pomocou metódy najmenších štvorcov

$$\hat{\mathbf{w}} = (\mathbf{X}^T \mathbf{X} + N\lambda \mathbf{I})^{-1} \mathbf{X}^T \mathbf{t} \quad (10.20)$$

Nájdite kovariančnú maticu Σ_0 , ktorá urobí tieto výrazy identickými, t. j. Σ_0 pomocou λ .

2. V príklade OH analyzujte efekt redukcie σ^2 .

Literatúra

- [1] M. Anthony, N. Biggs: Computational Learning Theory, Cambridge University Press, 1991, 1997
- [2] I. Goodfellow, Y. Bengio, A. Courville: Deep Learning. MIT Press book in preparation, <http://www.deeplearningbook.org/version-2016-03-08/>
- [3] J. Han, M. Kamber, J. Pei: Data Mining – Concepts and Techniques, 3rd edition. Morgan Kaufmann Publishers, 2011
- [4] T. Mitchell: Machine Learning, McGraw Hill, 1997
- [5] J. D. Riggs, T. L. Lalonde: Non-Gaussian and correlated data. Handbook for applied modeling. Cambridge University Press, 2017
- [6] S. Rogers, M. Girolami: A First Course in Machine learning. CRC Press, Taylor & Francis Group, 2012, 2017
- [7] J. Watt, R. Borhani, A. K. Katsaggelos: Machine Learning Refined. Foundations, Algorithms, and applications. Cambridge University Press, 2016
- [8] L. G. Valiant: Theory of the Learnable. Comm. of the ACM, Vol. 27, No. 1, November 1984

Strojové učenie

Vysokoškolský učebný text

Autori: doc. RNDr. Gabriela Andrejková, CSc.
RNDr. Ľubomír Antoni, PhD.

Vydavateľ: Univerzita Pavla Jozefa Šafárika v Košiciach
Vydavateľstvo ŠafárikPress

Rok vydania: 2020
Počet strán: 91
Rozsah: 4,5 AH
Vydanie: prvé



ISBN 978-80-8152-912-2 (e-publikácia)