

České vysoké učení technické v Praze

Fakulta jaderná a fyzikálně inženýrská

Akademie věd České republiky

Ústav informatiky a výpočetní techniky

Úvod do teorie neuronových sítí

Ing. František Hakl, CSc.

Ing. RNDr. Martin Holeňa, CSc.



1997

Ediční středisko ČVUT Praha

Předkládaná skripta shrnují látku prezentovanou v rámci přednášky „Úvod do teorie neuronových sítí“, která je určena pro studenty 5. ročníku matematického inženýrství na FJFI ČVUT Praha. Cílem této přednášky je seznámit posluchače se základními teoretickými poznatky a postupy, které jsou nezbytné k rigoróznímu studiu a efektivnímu použití paraelních výpočtů, označovaných jako „umělé neuronové sítě“.

Skriptum je rozděleno do dvou samostatných částí, z nichž každá byla psána s rozdílným přístupem k prezentaci učební látky. Přestože tyto dvě části na sebe navazují pouze volně, vzájemně jedna druhou vhodně doplňují a poskytují tak čtenáři možnost jak detailního studia, tak i možnost získání širokého přehledu.

První část skript (autor František Hakl) je psána pokud možno autonomně, tedy obsahuje důkazy veškerých lemm a tvrzení, s výjimkou několika široce známých tvrzení, se kterými byli posluchači 5. ročníku MI FJFI seznámeni již v průběhu svého studia, nebo jsou uvedeny v běžně dostupné literatuře.

Výklad v první části není zaměřen na detailní sumarizaci veškerých modelů neuronových sítí a jejich derivátů, ani na popis oblastí jejich aplikovatelnosti. Místo toho autor klade důraz na obecný popis paraelního zpracování informace neuronovými sítěmi, s cílem naznačit teoretické principy, na kterých je paraelismus neuronových sítí založen, se snahou zdůraznit souvislosti této teorie s poznatky standardních matematických disciplín. Důkladné zvládnutí těchto principů v prezentovaném rozsahu by mělo být postačující pro analýzu i některých jiných modelů neuronových sítí, které nejsou zde zmiňovány. Jednotlivé konkrétní modely neuronových sítí jsou uvedeny především na podporu a ilustraci vykládaných teoretických principů. Pokud bude mít student zájem o získání hlubšího přehledu vykládané látky, může využít literatury uvedené na konci skript.

Z důvodů omezeného rozsahu skript zcela opomíjíme v první části přístupy k neuronovým sítím založené na statistických metodách, teorii fuzzy množin, atd. Stejně tak vynecháváme výklad cyklických neuronových sítí, neuronových sítí odvozených od genetických algoritmů, a jiných výpočetních modelů, které s pojmem neuronových sítí nějak souvisí.

Druhá část skript (autor Martin Holeňa) má za cíl stručně sumarizovat současné poznatky o různých architekturách neuronových sítí a jejich dynamickém chování. Vzhledem k rozsahu předkládaných informací neobsahuje tato část důkazy jednotlivých tvrzení, eventuálně je obsahuje pouze v náznakové formě. Oproti předcházející první části se čtenář může seznámit s více modely neuronových sítí, pozornost je věnována souvislosti mezi neuronovými sítěmi a teorií fuzzy množin, dále pak rozvedení statistického pohledu na problematiku neuronových sítí.

Dovolte mi závěrem na tomto místě poděkovat slečně Ivaně Baranovské za přepis některých částí textu a celkovou jazykovou korekturu první části. Dále bych chtěl ocenit ochotu pana Arnošta Štedrého, který mi pomohl vyřešit mnohá úskalí a taje systémů \TeX a \LaTeX 2 ϵ , na která jsem při práci na těchto skriptech narazil.

Případné připomínky k textu, podněty či doporučení, směrujte prosím na níže uvedenou e-mailovou adresu.

Klasifikace moderních věd:

JE-LI TO ZELENÉ, NEBO SE TO HÝBE, PATŘÍ TO DO BIOLOGIE.
PÁCHNE-LI TO, PATŘÍ TO DO CHEMIE.
NEFUNGUJE-LI TO, PATŘÍ TO DO FYZIKY.
ZNÍ-LI TO JAKO NESMYSL, PATŘÍ TO DO EKONOMIE.
JE-LI TO NESROZUMITELNÉ, PATŘÍ TO DO MATEMATIKY.

Obsah

I Teoretické základy umělých neuronových sítí	
(AUTOR FRANTIŠEK HAKL)	9
1 Úvodní kapitola	13
2 Analýza booleovských sítí	9
3 Aproximační možnosti neuronových sítí	33
4 Vapnik-Chervonenkova dimenze	51
5 Teorie učení a neuronové sítě	61
6 Numerická analýza učicích algoritmů	77
II Relevantní výsledky související s neuronovými sítěmi	
(AUTOR MARTIN HOLEŇA)	87

Seznam obrázků

Část I

Teoretické základy umělých neuronových sítí

(AUTOR FRANTIŠEK HAKL)

Věnováno památce mých rodičů.

Kapitola 1

Úvodní kapitola

Edingtonova teorie

POČET NEJRŮZNĚJŠÍCH HYPOTÉZ,
SNAŽÍCÍCH SE OSVĚTLIT URČITÝ BIOLOGICKÝ JEV,
JE NEPŘÍMO ÚMĚRNÝ DOSTUPNÝM VĚDOMOSTEM.

V této úvodní kapitole nastíníme základní předmět, kterému jsou tato skripta věnována. Velmi krátce pohovoříme obecně o základním popisu umělých neuronových sítí, uvedeme některé vybrané modely a jejich základní charakteristiku.

1.1 Použité značení

V celém textu budeme používat následujícího značení a používání fontů:

$\stackrel{\text{def}}{=}$	defice nově zavedeného pojmu, množiny či čísla
$K = \{1, \dots, m\}$	množina přirozených čísel od 1 do m
B_α^N	koule v prostoru \mathfrak{R}^N se středem v počátku a poloměrem α
\mathfrak{R}^N	reálný N -rozměrný prostor
Ω^m	m -násobný kartézský součin množiny Ω
\bar{z}	prvek v kartézském (vícenásobném) součinu dané množiny
\bar{X}	obecná (pod)množina v dané množině
X	system podmnožin dané množiny
\mathcal{X}	množina systémů množin nad nějakou obecně danou množinou
$2^{\bar{X}}$	potenční množina množiny \bar{X} , $2^{\bar{X}} \stackrel{\text{def}}{=} \{\bar{Z} \mid \bar{Z} \subset \bar{X}\}$
$[\Omega]_\lambda$	lineární obal množiny Ω vektorů vektorového prostoru (=množina všech lineárních kombinací konečného počtu prvků z Ω)
$[\bar{B}]_\kappa$	konvexní obal množiny \bar{B}
\mathcal{H}_n^*	vektorový prostor (se specifikovanými parametry)
$C_{\bar{A}}$	prostor funkcí spojitých na dané množině \bar{A}
\mathbf{NP}_k^α	množina objektů s danými vlastnostmi, zpravidla množina funkcí, zobrazení, atd., dle kontextu
B^*	obecný učící algoritmus

$\bar{A} \triangle \bar{B}$	symetrická diference množin \bar{A} a \bar{B}
$\bar{A} - \bar{B}$	rozdíl množin \bar{A} a \bar{B}
$o(\tilde{f})$	$\tilde{g} = o(\tilde{f})$ vyjadřuje skutečnost, že funkce \tilde{g} je na asymptotickém okolí nuly zhora omezena v absolutní hodnotě funkcí \tilde{f}
$O(\tilde{f})$	$\tilde{g} = O(\tilde{f})$ vyjadřuje skutečnost, že funkce \tilde{g} je na asymptotickém okolí nekonečna zhora omezena v absolutní hodnotě funkcí \tilde{f}
\mathbf{A}	matice reálných čísel, sloupce či řádky budeme označovat společným názvem ŘADY MATICE
\vec{x}	vektor daného vektorového prostoru (vektory vždy budeme považovat za sloupce)
$\mathbf{A}^{(n)}$	čtvercová matice řádu 2^n
$\vec{x}^{(n)}$	vektor dimenze 2^n
$\langle \vec{x} \vec{y} \rangle$	skalární součin vektorů, $\langle \vec{x} \vec{y} \rangle = \sum_{i=1}^n \vec{x}_i \vec{y}_i$
$\vec{x} \odot \vec{y}$	tensorový součin vektorů či matic
$\overset{\wedge}{\text{AND}}$	binární logická funkce AND
$\overset{\vee}{\text{OR}}$	binární logická funkce OR
$\overset{\sqcup}{\text{XOR}}$	binární logická funkce XOR
$\overset{\neg}{\text{NOT}}$	unární logická funkce NOT
$\lfloor M \rfloor$	největší celé číslo menší než M (=celá dolní část)
$\lceil M \rceil$	nejmenší celé číslo větší než M (=celá horní část)
j^{\flat}	j^{\flat} představuje vektor z $\{-1, +1\}^n$, (resp. $\{0, 1\}^n$, podle kontextu), jehož souřadnice korespondují s binárním zápisem čísla j (souřadnice vektoru j^{\flat} odpovídající pozici 1 v binárním zápisu j , jsou rovny 1).
\vec{v}^{\sharp}	celé číslo, odvozené od vektoru $\vec{v} \in \{-1, +1\}^n$, (resp. $\vec{v} \in \{0, 1\}^n$, podle kontextu), kde jedničkové složky vektoru \vec{v}^{\sharp} odpovídají 1 v binárním zápisu čísla \vec{v}^{\sharp} , ostatní složky odpovídají pozicím 0 v tomto zápisu
$\binom{a}{i}$	binomický koeficient
\tilde{f}	obecné zobrazení mezi dvěma množinami
$(\tilde{f} * \tilde{g})$	konvoluce funkcí \tilde{f} a \tilde{g}
$\tilde{A}(b)$	hodnota funkce (zobrazení) \tilde{A} v argumentu b
$\log_e(x)$	přirozený logaritmus
$\log_2(x)$	logaritmus o základu 2
$x \equiv y \pmod{r}$	$x \equiv y \pmod{r}$ vyjadřuje, že číslo x je dělitelné číslem y se zbytkem r (děleno moduló r)
$Prob_{\tilde{\Pi}}(\tilde{A})$	pravděpodobnost množiny \tilde{A} při hustotě pravděpodobnosti $\tilde{\Pi}$
$ \tilde{V} $	mohutnost množiny \tilde{V}
$ \delta $	absolutní hodnota z čísla δ
$\ \vec{z}\ _{max}$	maximová norma vektoru
$\ \vec{z}\ _E$	Eukleidovská norma vektoru
$\ \vec{z}\ _{\mathcal{L}_2^{\bar{K}}}$	specifikovaná norma vektoru

Uvedme několik poznámek k použité terminologii. Nově definované pojmy jsou zvý-

razněny v textu použitím jiného fontu, např. **NOVĚ ZAVEDENÝ POJEM**. Co se týče české terminologie, v mnohých případech se jedná pouze o autorův překlad původního anglického termínu. Přestože snahou autora bylo těmito překlady vystihnout co nejpřesněji daný pojem (při co možná nejmenším pokřivení českého jazyka), je možné předpokládat nesoulad námi užívané terminologie s terminologií používanou jinými českými autory. V případech, kdy je tento nesoulad velmi pravděpodobný, uvádíme v příslušných definicích i původní anglický termín, např. (anglický termín: PAC learning), mimo jiné i proto, že ani anglická terminologie není ještě plně ustálena.

1.2 Motivace pro výzkum umělých neuronových sítí

Obecně lze chápat každý počítač jako zařízení, které zobrazuje množinu vstupních dat (např. koeficienty lineární soustavy) na množinu dat jiných (např. řešení této soustavy). Nebudeme-li na takovýto počítač klást další zužující podmínky, lze konstatovat, že i nervový systém individuálního živočicha je jakýmsi prototypem (každý je jiný) počítače, a mnohou duševní činnost či vegetativní aktivitu živočicha lze považovat za realizaci jakéhosi výpočtu (např. více či méně smysluplná odpověď na více či méně smysluplnou otázku, reakce na daný podnět prostředí, atd.). Mezi těmito dvěma paralelními světy „živých počítačů“ a klasických počítačů elektronických existují však zásadní rozdíly. Pokusme se sumarizovat ty z nich, o nichž se domníváme, že mají vztah k principu jejich činnosti:

klasický počítač	nervový systém
<ul style="list-style-type: none"> • existence centrálního procesoru provádějícího veškerou činnost • nízká paralelizace, vyjímecně více procesorů (≤ 1000) komunikujících mezi sebou • přesná znalost způsobu vyhodnocení a zpracování informace (stroje typu RAM - Random Acces Machine, RASP - Random Acces Stored Program, Turingův stroj) • nutnost detailní algoritmizace výpočtů, podložená hlubokým teoretickým pozadím (např. nutnost existence poznatků lineární algebry pro řešení soustav lineárních rovnic) • velmi vysoká rychlost početních elementárních operací (dnes $\simeq 300$ MHz) • přesně definovaná architektura procesoru 	<ul style="list-style-type: none"> • takovýto centrální procesor nebyl v nervových tkáních dosud objeven a nic nepotvrzuje jeho existenci • velmi vysoký počet jednodušších (z hlediska přenosu informace) procesorů-neuronů ($\simeq 10^{12}$), velmi vysoká hustota propojení mezi nimi (10^3 pro každý neuron) • velmi matná představa o činnosti jednotlivých elementů, téměř žádná představa o způsobu komunikace mezi elementy • zdá se, že schopnost vyhodnocování různých situací je integrální vlastností nervových systémů bez potřeby pochopení a uvědomění si způsobu zpracování informace • pomalý přenos informace (řádu milisekund) • velký počet lokálních elementů vzájemně propojených, s velkou variabilitou hustoty a způsobu spojení.

Výše uvedené rozdíly ospravedlňují názor, že schopnosti živých organismů efektivně zpracovávat informace ze svého okolí (v mnohých ohledech a případech výrazně efektivněji než nejvýkonnější počítače) jsou dány především způsobem komunikace mezi jednotlivými neurony, a že tyto neurony lze považovat za jakési elementární, principem shodné, stavební kameny nervových center.

Dá se očekávat, že analýza struktur, které jsou velmi jednoduchým modelem nervových center živočichů (na základě dosud známých poznatků z oblasti anatomie, biochemie a neurologie), může být přínosem jak pro výzkum živých organismů, tak i opačně – pro výzkum a implementaci nových typů vysoce paralelních počítačových architektur, které umožní řešit problémy, pro které na počítačích odvozených od principu Thuringova stroje neexistuje efektivní algoritmizace.

V následující podkapitole zavedeme základní modely umělých neuronových sítí a popíšeme stručně jejich dynamické chování.

1.2.1 Základní popis architektury a dynamiky neuronových sítí

Nyní popíšeme obecnou architekturu a dynamiku neuronových sítí za účelem definování předmětu analýzy pro následující kapitoly těchto skript.

Neuronovou síť budeme chápat jako systém sestávající se z jakési *architektury*, *dynamiky* a *úlohy*.

Úlohou budeme rozumět zobrazení, které by měla daná neuronová síť realizovat mezi množinou svých vstupů (= vstupní prostor) a množinou svých výstupů (= výstupní prostor). Interpretace této úlohy může být rozličná, nejčastěji hovoříme o separaci tříd, kdy vstupní prostor je tvořen množinou $\bar{A} \in \mathbb{R}^n$, výstupní prostor je množina $\{-1, +1\}$, a naším cílem je, aby výstupní hodnota neuronové sítě byla rovna charakteristické funkci množiny \bar{A} . Další častá úloha je dána jako problém aproximace. Pro dané zobrazení $\tilde{Z} : \mathbb{R}^n \rightarrow \mathbb{R}^m$ požadujeme, aby výstupy z neuronové sítě co nejlépe aproximovaly zobrazení \tilde{Z} v nějaké metrice prostoru \mathbb{R}^m . Poslední úlohu, kterou zde zmíníme, je problém interpolace a/nebo extrapolace (= též predikce), kdy pro zadané zobrazení $\tilde{Z} : \mathbb{R}^n \rightarrow \mathbb{R}^m$ a množiny $\bar{A} \in \mathbb{R}^n$, $\bar{B} \in \mathbb{R}^n$ požadujeme, aby síť na základě množiny $\tilde{Z}(\bar{A})$ interpolovala či extrapolovala množinu $\tilde{Z}(\bar{B})$.

Poznamenejme, že nejčastější způsob zadání množin \bar{A} a \bar{B} je konečná množina jejich prvků a nejčastější způsob zadání zobrazení je dán jako konečná množina uspořádaných dvojic $\{\vec{x}, \tilde{Z}(\vec{x})\}$.

Architektura neuronové sítě je daná orientovaným grafem, jehož uzly jsou tvořeny neurony a který popisuje způsob propojení jednotlivých neuronů, dále pak popisem tzv. *přenosové funkce* neuronů. Množiny uzlů (= neuronů), do nichž nevede hrana, nazveme vstupními neurony, ty uzly, z nichž nevede hrana, nazveme výstupními neurony, ostatní uzly nazveme vnitřní (někdy též skryté) neurony. Podle druhu použitého grafu je zvykem rozlišovat následující základní typy architektur:

- dvouvrstvá neuronová síť tvořená bipartitním grafem
- neuronová síť tvořená úplným grafem
- vícevrstvá neuronová síť, kdy množina uzlů je rozdělena do disjunktivních množin U_1, \dots, U_n a hrany grafu jsou pouze mezi množinami U_i a U_{i+1} , $i \in \{1, \dots, n-1\}$.

Přenosová funkce neuronu je zobrazení $\tilde{P} : \bar{W} \times \bar{X} \rightarrow \mathfrak{R}$, kde \bar{X} je množina všech možných vstupních vektorů (odpovídajících vstupním vrcholům) a \bar{W} je parametrický prostor daného neuronu v grafu sítě. Každému vektoru z tohoto prostoru budeme říkat váhový vektor neuronu. Nejčastější přenosová funkce neuronu je rovna funkci

$$\tilde{y}(\vec{x}) \stackrel{\text{def}}{=} \tilde{\sigma}(\langle \vec{w} | \vec{x} \rangle - t),$$

kde \vec{w} je vektor tvořený jednotlivými váhami, t je prahová hodnota daného neuronu, $\tilde{y}(\vec{x})$ je výstupní hodnota neuronu, \vec{x} je vstupní vektor neuronu. $\tilde{\sigma}$ bývá nejčastěji buď binární (tzv. tvrdá nelinearita):

$$\tilde{\sigma}(z) \stackrel{\text{def}}{=} \begin{cases} 1 & \text{pro } x \geq 0 \\ -1 & x < 0 \end{cases},$$

nebo tzv. sigmoidální funkce.

Definice 1.2.1 *Nechť $\tilde{\sigma}$ je reálná funkce reálné proměnné, monotónně rostoucí v celém definičním oboru \mathfrak{R} , vyhovující podmínkám*

$$\lim_{z \rightarrow +\infty} \tilde{\sigma}(z) = 1, \quad \lim_{z \rightarrow -\infty} \tilde{\sigma}(z) = -1.$$

Potom takovouto funkci nazveme SIGMOIDÁLNÍ FUNKCE.

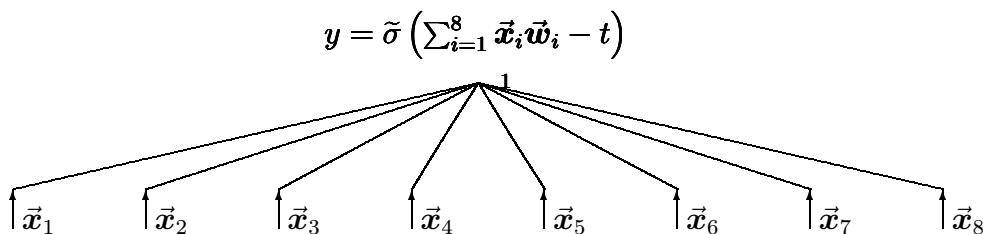
Dynamikou neuronové sítě budeme rozumět změny parametrů přenosových funkcí neuronů (eventuelně i samotného grafu sítě) v čase. Dynamika sestává ze dvou částí

- aktivní dynamika (proces vybavování), kdy na vstupní neurony vložíme vstupní hodnoty sítě, tyto vstupní hodnoty po hranách probíhají k dalším neuronům, přes ně přechází modifikované podle jejich přenosových funkcí k dalším neuronům, až se modifikované hodnoty objeví na výstupních neuronech, kde jsou odečítány a interpretovány jako výstupní hodnoty sítě.
- adaptivní dynamika (proces učení), při které se hledají a nastavují váhy spojnic (= hran grafu) jednotlivých neuronů a ostatní parametry přenosových funkcí jednotlivých neuronů, eventuelně se upravuje celý graf tak, aby síť řešila danou konkrétní úlohu.

Nejdůležitější přístupy k adaptivní dynamice jsou založeny na minimalizaci hodnoty celkové chybové funkce

$$\min \left\{ \sum_{i=1}^t \left\| \tilde{N}(\langle \vec{x}_i | \vec{w} \rangle) - y_i \right\|^2 \mid \vec{w} \in \bar{W} \right\},$$

kde $\{(\vec{x}_1, y_1), \dots, (\vec{x}_t, y_t)\}$ je úloha, $\tilde{N}(\langle \vec{x} | \vec{w} \rangle)$ je zobrazení realizované neuronovou sítí, \vec{w} je vektor všech vah v síti. Tato minimalizace se provádí nejčastěji gradientními metodami, např. učení metodou zpětného šíření (anglický termín: back-propagation), v mnoha případech ale také různými heuristickými algoritmy, ze kterých zmíníme významnější v následujících příkladech.



Obrázek 1.1: Schéma architektury perceptronu.

1.2.2 Příklady základních neuronových sítí

Jako prvotní zdroj informací o základních modelech neuronových sítí, jejich dynamickém chování a aplikacích lze využít např. [ŠJ96]. Mezi nejznámější modely neuronových sítí patří následující dva, které jsou patrně nejlépe prozkoumány z teoretického hlediska a z hlediska praktických aplikací nejvíce rozšířené.

Perceptron

Perceptron je dvouvrstvá síť, kde výstupní vrstva je tvořena právě jedním neuronem, spojeným se všemi neurony vstupní vrstvy. Výstupní hodnota perceptronu je rovna

$$y \stackrel{\text{def}}{=} \tilde{f} \left(\sum_{i=1}^n \vec{x}_i \vec{w}_i - t \right),$$

kde dvojice \vec{w}, t představuje parametry přenosové funkce výstupního neuronu a \tilde{f} je buď tzv. tvrdá nelinearita nebo libovolná sigmoidální funkce. Pro danou úlohu

$$\{(\vec{x}_1, y_1), \dots, (\vec{x}_t, y_t)\}$$

definujeme proces učení jako iterační algoritmus popsany vztahem

$$\vec{w}_{i+1} \stackrel{\text{def}}{=} \vec{w}_i + y \left(y - \vec{w}_i^T \vec{x} \right) \vec{x}, \quad (1.1)$$

pro nějaký výchozí vektor vah \vec{w}_0 . Vlastnosti této posloupnosti budeme studovat v části 6.1.

Vrstevnaté sítě, učení metodou back propagation

Velmi často používaným typem neuronové sítě je vrstevnatý model, kdy neurony jsou rozděleny do vrstev a každý neuron v dané vrstvě je spojen se všemi neurony ve vrstvě předcházející a ve vrstvě následující (tak, že příslušný graf nemá smyčky, hrany vedou pouze mezi vrstvami).

Předpokládejme přechodovou funkci ve tvaru 1.1, kde ale funkce \tilde{f} je spojitá, dvakrát spojitě diferencovatelná, sigmoidální funkce.

Předpokládejme, že máme opět zadanou úlohu $\{(\vec{x}_1, \vec{y}_1), \dots, (\vec{x}_t, \vec{y}_t)\}$, kde $\vec{x}_i \in \mathbb{R}^n$, $\vec{y}_i \in \mathbb{R}^m$, $i \in \{1, \dots, t\}$, a že pomocí gradientní metody největšího spádu minimalizujeme hodnotu chybové funkce

$$\tilde{E}(\vec{w}) \stackrel{\text{def}}{=} \frac{1}{2} \sum_{i=1}^t \left\| \tilde{N}(\vec{x}_i, \vec{w}) - \vec{y}_i \right\|_E^2.$$

Abychom pochopili princip metody zpětného šíření, vyslovíme následující lemmu, která dává explicitní vyjádření derivace vícenásobně zanořené složené funkce.

Lemma 1.2.1 Předpokládejme, že pro pevné přirozené k je dána posloupnost reálných čísel w_1, \dots, w_k , dále posloupnost $\tilde{\sigma}_1, \dots, \tilde{\sigma}_k$ reálných diferencovatelných funkcí jedné reálné proměnné, a necht' reálná posloupnost y_1, \dots, y_{k+1} splňuje rekurzivní vztah

$$y_{i+1} \stackrel{\text{def}}{=} \tilde{\sigma}_i(w_i y_i), \quad i \in \{1, \dots, k\}.$$

Potom platí, že

$$\frac{\partial y_{k+1}}{\partial w_k} = \tilde{\sigma}'_k(w_k y_k) \cdot y_k$$

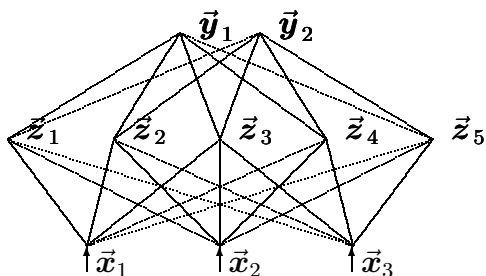
a

$$\frac{\partial y_{k+1}}{\partial w_j} = \left(\prod_{m=j+1}^k \tilde{\sigma}'_m(w_m y_m) \cdot w_m \right) \tilde{\sigma}'_j(w_j y_j) \cdot y_j,$$

pro $j \in \{1, \dots, k-1\}$.

■ Důkaz:

DŮKAZ JE UVEDEN V ORIGINALNÍ TIŠTĚNÉ VERZI. ■



Obrázek 1.2: Schéma architektury vrstvenaté sítě.

Uvědomíme-li si skutečnost, že formální pravidla pro derivaci složené funkce a pro derivaci součinu platí i v případě, kdy jednotlivé součinitele jsou matice a argumenty funkcí jsou vektory (zde předpokládáme, že $\tilde{\sigma}$ je vektorová funkce, aplikovaná na jednotlivé složky svého argumentu stejně jako v jednorozměrném případě), je zřejmé, že lemma 1.2.1 popisuje i formální zápis derivace zobrazení, které je realizováno vícevrstvou neuronovou sítí. Nyní je ale zřejmý význam termínu zpětné šíření. Chceme-li spočítat hodnotu derivace výstupu z nějakého neuronu

podle váhy mezi j -tou a $(j+1)$ -ní vrstvou, postupujeme tak, že nejdříve spočítáme hodnotu všech neuronů pro daný vstupní vektor \vec{x} , a potom podle analogických vzorců jako v lemmě 1.2.1 spočítáme derivace výstupních hodnot nejdříve pro váhy mezi poslední a předposlední vrstvou, dále pak pro váhy mezi vrstvami o jedničku níže a takto postupujeme až do konce pro váhy mezi první a druhou vrstvou.

Různé způsoby definování architektury sítě a její dynamiky umožňují postulovat nesčetné množství různých typů neuronových sítí. V současné době se hovoří o řádově několika desítkách mnohdy dosti odlišných neuronových sítí. Pro ilustraci a podání základní informace uvedeme pouze následující modely neuronových sítí, a to ve velmi sumarizované formě.

Hopfieldova síť

Tato síť je tvořena úplným grafem na n vrcholech. Každá váha mezi dvěma vrcholy má hodnotu ± 1 a pro danou posloupnost $\vec{x}_1, \dots, \vec{x}_t \in \{-1, +1\}^n$ je váha mezi i -tým a j -tým neuronem dána vzorcem

$$W_{ij} \stackrel{\text{def}}{=} \sum_{s=1}^t (\vec{x}_i)_s (\vec{x}_j)_s,$$

tedy výsledná matice všech vah \mathbf{W} má tvar

$$\mathbf{W} = \sum_{i=1}^t \vec{x}_i \vec{x}_i^T. \quad (1.2)$$

Adaptace vah podle vzorce 1.2 umožňuje postupné předkládání jednotlivých vzorů po sobě. Aktivní dynamika, tedy vybavování, je popsána iteračním předpisem

$$\vec{y}_{k+1} \stackrel{\text{def}}{=} \tilde{f}(\mathbf{W} \vec{y}_k), \quad \vec{y}_k, \vec{y}_{k+1} \in \{-1, +1\}^n, \quad (1.3)$$

kde hodnoty vektoru \vec{y} můžeme chápat jako výstupní hodnoty jednotlivých neuronů a \tilde{f} je již výše zmíněná tvrdá nelinearita. Iterace podle vzorce 1.3 se provádí, dokud není dosažen stabilní stav, kdy $\vec{y}_{k+1} = \vec{y}_k$.

Smyslem této neuronové sítě je zapamatovat si předkládané vzory. Dá se ukázat, že iterační posloupnost \vec{y}_k snižuje hodnoty energetické funkce definované jako

$$\tilde{H}(\vec{y}) \stackrel{\text{def}}{=} -\frac{1}{2} \sum_{i=1}^n \sum_{j=1}^n \mathbf{W}_{ij} \vec{y}_i \vec{y}_j \quad (1.4)$$

a že vede vždy do lokálního minima této funkce.

Kohonenova síť

Poměrně dosti odlišným typem neuronové sítě je tzv. Kohonenova síť (někdy též Kohonenova mapa). Graf této sítě má tu vlastnost, že počet vstupních hran každého neuronu je roven počtu vstupních neuronů celé sítě a k zakódování vstupních vektorů slouží váhy neuronů. Neuronu nemají přenosovou funkci, ale pouze počítají v dané metrice vstupního prostoru vzdálenost předloženého vzoru od vektoru vstupních vah, např. Eukleidovskou vzdálenost (\vec{w} je vektor vah j -tého neuronu)

$$d \stackrel{\text{def}}{=} \sqrt{\sum_{i=1}^n (\vec{x}_i - \vec{w}_i)^2}.$$

Tyto hodnoty lze považovat za výstupní hodnoty Kohonenovy sítě. Učení Kohonenovy sítě probíhá podle následujícího schématu.

- Nejdříve zvolíme parametr rychlosti učení $0 \leq \tilde{\eta}(k) \leq 1$, který pro k -tý vzor rozhoduje o míře změny vah. Dále pro každý neuron zvolíme jeho okolí v grafu a zvolíme strategii zmenšování tohoto okolí v závislosti na výstupní hodnotě neuronu. Pro předložený vstupní vektor $\vec{x} \in \mathfrak{R}^n$ spočteme pro každý neuron sítě hodnotu

$$d_j \stackrel{\text{def}}{=} \|\vec{x} - \vec{w}_j^k\|,$$

kde \vec{w}_j^k je vektor vah pro k -tý neuron v čase k .

- Pro každý neuron j , pro který platí $d_j = \min_i \{d_i\}$, změníme váhy tohoto neuronu a všech neuronů v jeho aktuálním okolí podle formule

$$\vec{w}_{n+1}^{k+1} \stackrel{\text{def}}{=} \vec{w}_n^k + \tilde{\eta}(k) (\vec{x} - \vec{w}_n^k).$$

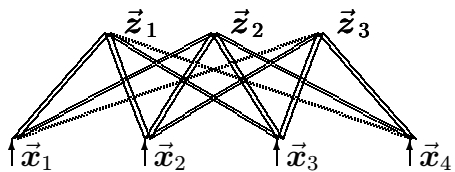
- Podle hodnoty d_j zredukujeme podle zvolené strategie okolí j -tého neuronu.
- Pokračujeme s tímto procesem pro každý další vstupní vektor.

Aktivní dynamika Kohonenovy sítě spočívá ve výpočtu hodnot d_j pro předložený vstupní vektor a všechny neurony v síti. Odpovědí sítě na tento vstup je pak pořadí neuronu s minimální hodnotou d_j .

Princip této sítě spočívá ve skutečnosti, že v průběhu učení se v síti vytvářejí jakési shluky neuronů, jejichž váhy jsou blízké předkládaným vzorům. Tyto shluky odpovídají třídám vstupních vektorů, které se pomocí Kohonenovy sítě vzájemně separují od sebe.

ART – Carpenter–Grossbergův klasifikátor

Návrhem této sítě se její autoři snažili potlačit nežádoucí jev projevující se u různých typů neuronových sítí, a to tzv. „problém proměnné stability“. Pod tímto pojmem se myslí skutečnost, kdy při předložení vzoru, který se má síť naučit, se poruší již dříve uložené informace, získané na základě dosavadního učení. Síť pak zpravidla již není konzistentní s dříve předloženými vzory pro učení. Proto byl rozpracován následující algoritmus Adaptivní Rezonanční Teorie.



Obrázek 1.3: Schéma architektury ART sítě ($N = 4$, $M = 3$).

Graf ART sítě je podobný dvouvrstvému modelu, ve kterém jsou spojeny všechny neurony vstupní vrstvy se všemi neurony vrstvy druhé, výstupní. Na rozdíl od vrstvenatých sítí, všechny hrany v grafu jsou obousměrné a každému směru je přiřazena váha. Hlavním rozdílem vzhledem k vrstvenatým sítím je opakovaná výměna vstupního vektoru mezi vstupní a výstupní vrstvou, tedy se dá říci, že vstupní vektor v této síti „rezonuje“. Dynamické chování sítě typu ART je

popsáno následujícím schématem:

1. **Inicializace** Počáteční nastavení dopředných a zpětných vah:

$$\widetilde{W}_{ij}(0) \stackrel{\text{def}}{=} \frac{1}{1+N}, \quad \widetilde{T}_{ij}(0) \stackrel{\text{def}}{=} 1,$$

kde $1 \leq i \leq N$, $1 \leq j \leq M$, a nastavení tzv. PRAHU OSTRÁŽITOSTI ρ v rozmezí $0 \leq \rho \leq 1$ (symbol $\widetilde{W}_{ij}(0)$ představuje dopřednou váhu, vedoucí z i -tého vstupu na j -tý výstup a symbol $\widetilde{T}_{ij}(0)$ představuje zpětnou váhu, vedoucí opačným směrem, tj. z výstupů ke vstupům). Ve zpětných vahách jsou zakódovány všechny naučené vzory. Práh ostrážitosti ρ určuje, jak moc je vstupní vzor vzdálen od nějakého již naučeného vzoru.

2. **Předložení nového vzoru na vstupy sítě** Přiložíme nový vzor $\vec{x} \in \{-1, +1\}^N$ na vstup sítě.
3. **Výpočet výstupních hodnot**

$$\mu_j \stackrel{\text{def}}{=} \sum_{i=1}^N \widetilde{W}_{ij}(t) \vec{x}_i, \quad 1 \leq j \leq M$$

(v této rovnici představuje μ_j výstup z j -tého neuronu a \vec{x}_i je i -tý element vstupního vzoru \vec{x} , jehož hodnota může být pouze -1 nebo 1).

4. Výběr nejlepšího výstupu

$$\mu_{j^*} \stackrel{\text{def}}{=} \max_j \mu_j.$$

Tento vztah představuje postranní inhibici ve výstupní vrstvě sítě.

5. **Test ostrážitosti** Porovnáme vstupní a vybraný vektor s prahem ostrážitosti. Jestliže je

$$S \stackrel{\text{def}}{=} \frac{\sum_{i=1}^N \mathbf{T}_{ij^*} \vec{x}_i}{\sum_{i=1}^N \vec{x}_i} > \rho,$$

potom přejdeme ke Kroku 7, jinak pokračujeme Krokem 6 (S zřejmě představuje podobnost vstupního vzoru s vektorem zpětných vah pro j -tý neuron výstupní vrstvy).

6. **Potlačení vybraného vzoru** Výstup neuronu, který doposavad představoval nejlepší vybraný vzor nastavíme na nulu a tím se nebude dále uplatňovat při výběru dalšího výstupu (v Kroku 4).

Pokud jsme nevyčerpali všechny neurony ve výstupní vrstvě, tak pokračujeme Krokem 4. V opačném případě, máme-li dostatek volných výstupních neuronů, přiřadíme vstupní vzor k nějakému volnému neuronu a vytvoříme tak novou třídu.

7. Přizpůsobení vah podle vybraného vzoru

$$\widetilde{\mathbf{T}}_{ij^*}(t+1) \stackrel{\text{def}}{=} \widetilde{\mathbf{T}}_{ij^*}(t) \vec{x}_i,$$

$$\widetilde{\mathbf{W}}_{ij^*}(t+1) \stackrel{\text{def}}{=} \frac{\widetilde{\mathbf{T}}_{ij^*}(t) \vec{x}_i}{\frac{1}{2} + \sum_{i=1}^N \widetilde{\mathbf{T}}_{ij^*}(t) \vec{x}_i}.$$

8. **Opakování algoritmu pro další vzor** Pokud chceme pokračovat v klasifikaci, povolíme všechny výstupy, které byly potlačeny v Kroku 6 a přejdeme ke Kroku 2. V opačném případě skončíme.

Předchozím modelem uzavřeme náš ilustrativní přehled nejužívanějších neuronových sítí a dále se budeme věnovat teoretické analýze vlastností zejména vrstevnatých sítí založených na myšlence perceptronu, důraz položíme na vlastnosti sítí s binárními vstupy a výstupy. Pro spojitě sítě obdobného typu ukážeme základní aproximační vlastnosti a v dalších kapitolách naznačíme souvislost neuronových sítí a obecné teorie učení. Dále pak ukážeme některé numerické aspekty gradientních algoritmů učení.

Kapitola 2

Analýza booleovských sítí

Cooperův zákon

POKUD V MATEMATICKÉM TEXTU NEROZUMÍTE NĚJAKÉMU SLOVU,
NEVĚŠTE HLAVU.
TEXT DÁVÁ SMYSL I BEZ NĚHO.
(JESTLI-ŽE NEDÁVÁ, NEDÁVAL BY SMYSL ANI S NÍM.)

V celé této kapitole se budeme věnovat analýze booleovských obvodů, to jest neuronových sítí, jejichž vstupy i výstupy, stejně jako stavy všech neuronů, jsou binární. Budeme zkoumat základní otázku, jaké booleovské funkce definované v definičním oboru $\{-1, +1\}^n$ je daná booleovská síť schopna spočítat. Vyjdeme z nejjednodušší booleovské sítě, booleovského perceptronu, který počítá tzv. prahový vektor svých vstupů, jak bude vidět dále.

V analýze booleovských funkcí hrají velmi důležitou roli vlastnosti Hadamardových matic, analyzované v následující části. Tyto vlastnosti nám umožní ostře separovat množiny booleovských funkcí, realizovaných obvody různé hloubky. Význam analýzy booleovských obvodů má význam i z hlediska praktického, v technické praxi je problém realizace dané logické funkce běžný, zejména v aplikacích řízení a optimalizace procesů. Kromě základní charakterizace funkcí počítaných různými obvody, se budeme zajímat i o horní a dolní odhad počtu lineárně separovatelných booleovských funkcí a horními a dolními odhady velikosti vah nezbytných k výpočtu dané funkce. V celé kapitole budeme budovat teorii pro binární vstupy ± 1 , ale zřejmě všechna tvrzení zůstávají v platnosti pro libovolné binární vstupy $a, b \in \mathbb{R}, a \neq b$.

2.1 Sylvestrova konstrukce Hadamardových matic

Tuto úvodní část věnujeme zavedení a popisu Hadamardových matic, hrajících klíčovou roli při analýze vlastností booleovských obvodů hloubky 1 a 2, které budou studovány v kapitole 2. Pro studium Hadamardových matic je vhodné zavést pojem tenzorového součinu podle následující definice.

Definice 2.1.1 *Nechť A je reálná matice typu $r \times s$ a B je reálná matice typu $m \times n$. Potom TENZOROVÝ SOUČIN MATIC A a B je matice $A \odot B$ typu $mr \times ns$,*

definovaná jako

$$\mathbf{A} \odot \mathbf{B} \stackrel{\text{def}}{=} \begin{pmatrix} \mathbf{A}b_{1,1} & \mathbf{A}b_{1,2} & \dots & \mathbf{A}b_{1,n} \\ \mathbf{A}b_{2,1} & \mathbf{A}b_{2,2} & \dots & \mathbf{A}b_{2,n} \\ \vdots & \vdots & \ddots & \vdots \\ \mathbf{A}b_{m,1} & \mathbf{A}b_{m,2} & \dots & \mathbf{A}b_{m,n} \end{pmatrix}.$$

Pro operaci tenzorového součinu platí asociativní zákon, což se dá velmi snadno ověřit přímo z definice této operace, a neplatí zákon komutativní. Definujme nyní Hadamardovy matice.

Definice 2.1.2 *Nechť \mathbf{A} je čtvercová matice řádu n , jejíž prvky jsou z množiny $\{-1, +1\}$. Potom \mathbf{A} je HADAMARDOVA MATICE, jestli-že platí*

$$\mathbf{A} \cdot \mathbf{A}^T \stackrel{\text{def}}{=} n\mathbf{I}.$$

Souvislost Hadamardových matic a tenzorového součinu ozřejmuje následující lemma.

Lemma 2.1.1 *Tenzorový součin Hadamardových matic je opět Hadamardova matice.*

■ *Důkaz:*

DŮKAZ JE UVEDEN V ORIGINALNÍ TIŠTĚNÉ VERZI. ■

Pro analýzu booleovských neuronových sítí jsou relevantní matice, vzniklé následující rekurzí z Hadamardovy matice řádu 2.

Definice 2.1.3 *Nechť $\mathbf{A} \stackrel{\text{def}}{=} \begin{pmatrix} 1 & 1 \\ 1 & -1 \end{pmatrix}$. Potom matici $\mathbf{B}^{(n)}$ řádu $2^n \times 2^n$ vzniklou n -násobným tenzorovým násobením matice \mathbf{A} nazveme MATICE PARITY ŘÁDU n .*

Způsob vytvoření matice $\mathbf{B}^{(n)}$ se v literatuře označuje jako Sylvestrova konstrukce. Matice $\mathbf{B}^{(n)}$ mají zajímavou algebraickou strukturu. Nejdůležitější vlastností z hlediska neuronových sítí je fakt, že sloupce matice $\mathbf{B}^{(n)}$ jsou tvořeny hodnotami booleovských funkcí $\{-1, +1\}^n \times \{-1, +1\}^n \rightarrow \{-1, +1\}$, které se nazývají FUNKCEMI PARITY, tj. funkcemi jejichž hodnota v argumentu (\vec{a}, \vec{b}) , $\vec{a}, \vec{b} \in \{-1, +1\}^n$ je rovna 1, právě když počet jedniček v \vec{a} na pozicích, kde v \vec{b} je jednička, se od počtu -1 na těchže pozicích liší o sudé číslo. Základní vlastnosti matic parity shrnuje následující pomocné tvrzení.

Lemma 2.1.2 *Pro matici parity $\mathbf{B}^{(n)}$ platí následující:*

1. $\mathbf{B}^{(n)}$ je symetrická Hadamardova matice,
2. $\mathbf{B}^{(n-1)}$ je hlavní vedoucí submaticí matice $\mathbf{B}^{(n)}$,
3. vynecháním sudých řádků a sloupců matice $\mathbf{B}^{(n)}$ dostaneme matici $\mathbf{B}^{(n-1)}$,
4. vynecháním lichých řádků a sloupců matice $\mathbf{B}^{(n)}$ dostaneme matici $-\mathbf{B}^{(n-1)}$,
5. počet jedniček v matici $\mathbf{B}^{(n)}$ se od počtu -1 v téže matici liší o hodnotu 2^n ,

6. nechť $\vec{a}, \vec{x} \in \{0, 1\}^j$, $\vec{b}, \vec{y} \in \{0, 1\}^{n-j}$, $0 < j < n$ a nechť dále zápis (\vec{z}, \vec{y}) představuje blokový vektor s bloky \vec{z} a \vec{y} . Potom platí následující vztah:

$$B^{(n)}_{(\vec{a}, \vec{b})^\#, (\vec{x}, \vec{y})^\#} = B^{(j)}_{\vec{a}^\#, \vec{x}^\#} \cdot B^{(n-j)}_{\vec{b}^\#, \vec{y}^\#},$$

(řady matice $B^{(n)}$ číslyme od 0 do $2^n - 1$),

7. pro všechna $\vec{i}, \vec{j} \in \{0, 1\}^n$ platí rovnost

$$B^{(n)}_{\vec{i}^\#, \vec{j}^\#} = (-1)^{\sum_{\alpha=1}^n \vec{i}_\alpha \cdot \vec{j}_\alpha},$$

(řady matice $B^{(n)}$ číslyme od 0 do $2^n - 1$).

■ *Důkaz:*

DŮKAZ JE UVEDEN V ORIGINALNÍ TIŠTĚNÉ VERZI.

■

Dalším nosným pojmem pro analýzu booleovských neuronových sítí je pojem rostoucí matice.

Definice 2.1.4 Čtvercová matice A je rostoucí, tvoří-li její řádky a sloupce neklesající posloupnosti. Čtvercová matice A je potenciálně rostoucí, existují-li permutační matice P a Q tak, že matice $P \cdot A \cdot Q$ je rostoucí.

Rostoucí matice musí mít hlavní diagonálu a všechny „diagonály“ rovnoběžné s hlavní diagonálou nerostoucí, což triviálně plyne z definice rostoucí matice. Stejně tak z definice triviálně plyne, že pro potenciálně rostoucí matici A je matice $-A$ také potenciálně rostoucí a součet matice A s libovolnou konstantní maticí je opět potenciálně rostoucí. Existuje jednoduchá ekvivalentní podmínka pro to, aby matice byla potenciálně rostoucí.

Lemma 2.1.3 Čtvercová matice A s prvky ± 1 je potenciálně rostoucí, právě když neobsahuje jako svojí podmatici matici $\begin{pmatrix} 1 & -1 \\ -1 & 1 \end{pmatrix}$ nebo matici $\begin{pmatrix} -1 & 1 \\ 1 & -1 \end{pmatrix}$.

■ *Důkaz:*

DŮKAZ JE UVEDEN V ORIGINALNÍ TIŠTĚNÉ VERZI.

■

Počet submatic řádu 2 matice řádu n je zřejmě $\binom{n}{2}^2$ a tedy ověření, zda matice řádu n je potenciálně rostoucí, lze provést v čase $o(n^4)$. Další dvě pomocná tvrzení 2.1.4, 2.1.5 zavádí pomocné vektory \vec{p} , \vec{v} a posloupnosti $\{z_i\}_1^\omega$ a $\{w_i\}_1^\omega$, jejichž algebraická struktura umožňuje dokázat stěžejní tvrzení této kapitoly.

Lemma 2.1.4 Nechť $\vec{p}^{(0)} \stackrel{\text{def}}{=} \vec{v}^{(0)} \stackrel{\text{def}}{=} 1$ a nechť pro $n \geq 1$ platí

$$\vec{p}^{(n)} \stackrel{\text{def}}{=} \vec{p}^{(n-1)} \odot \begin{pmatrix} 1 \\ -1 \end{pmatrix} \quad \text{a} \quad \vec{v}^{(n)} \stackrel{\text{def}}{=} \vec{v}^{(n-1)} \odot \begin{pmatrix} -1 \\ 1 \end{pmatrix}.$$

Nechť s_k (resp. t_k) je součet prvních k členů číselné posloupnosti $\{(\vec{p}^{(n)})_i\}_1^{2^n}$ (resp. $\{(\vec{v}^{(n)})_i\}_1^{2^n}$). Potom platí:

1. $\forall l \in \{1, \dots, 2^{n-1}\}$ je $s_{2l} = 0$
2. $\forall l \in \{1, \dots, 2^{n-1}\}$ je $s_{2l-1} = (\vec{p}^{(n)})_{2l-1} = (-1)^{\kappa(2l-1)}$, kde symbol $\kappa(i)$ označuje počet jedniček v binárním rozvoji čísla i .
3. $\forall l \in \{1, \dots, 2^n - 1\}$ platí $s_l = t_{2^n-l}$.

■ *Důkaz:*

DŮKAZ JE UVEDEN V ORIGINÁLNÍ TIŠTĚNÉ VERZI.

■

Vektor $\vec{p}^{(n)}$ definovaný v lemmatu je roven poslednímu sloupci matice $\mathbf{B}^{(n)}$, což přímo plyne z definice $\mathbf{B}^{(n)}$.

Lemma 2.1.5 *Nechť d_2, \dots, d_k jsou přirozená čísla a d_1 je celé nezáporné číslo. Definiujme posloupnost vektorů \vec{q}_j a \vec{u}_j rekurzivním předpisem:*

$$\begin{array}{ll}
\vec{q}_1 \stackrel{\text{def}}{=} \vec{p}^{(d_1)} & \vec{u}_1 \stackrel{\text{def}}{=} \vec{v}^{(d_1)} \\
\vec{q}_2 \stackrel{\text{def}}{=} -(\vec{q}_1^T, \vec{q}_1^T, \dots, \vec{q}_1^T)^T & d_2 \times \quad \vec{u}_2 \stackrel{\text{def}}{=} -(\vec{q}_1^T, \vec{q}_1^T, \dots, \vec{q}_1^T)^T & d_2 \times \\
\vec{q}_3 \stackrel{\text{def}}{=} \vec{q}_2 \odot \vec{p}^{(d_3)} & \vec{u}_3 \stackrel{\text{def}}{=} \vec{q}_2 \odot \vec{v}^{(d_3)} \\
\vec{q}_4 \stackrel{\text{def}}{=} -(\vec{q}_3^T, \vec{q}_3^T, \dots, \vec{q}_3^T)^T & d_4 \times \quad a \quad \vec{u}_4 \stackrel{\text{def}}{=} -(\vec{q}_3^T, \vec{q}_3^T, \dots, \vec{q}_3^T)^T & d_4 \times \\
\vec{q}_5 \stackrel{\text{def}}{=} \vec{q}_4 \odot \vec{p}^{(d_5)} & \vec{u}_5 \stackrel{\text{def}}{=} \vec{q}_4 \odot \vec{v}^{(d_5)} \\
\vec{q}_6 \stackrel{\text{def}}{=} \dots & \vec{u}_6 \stackrel{\text{def}}{=} \dots \\
\dots & \dots
\end{array} \tag{2.1}$$

Nechť posloupnost $\{z_i\}_1$ je definována jako zřetězení

$$\{z_i\}_1 \stackrel{\text{def}}{=} (\vec{q}_1^T, \vec{q}_2^T, \dots, \vec{q}_k^T),$$

a posloupnost $\{w_i\}_1$ jako zřetězení

$$\{w_i\}_1 \stackrel{\text{def}}{=} (\vec{u}_1^T, \vec{u}_2^T, \dots, \vec{u}_k^T),$$

Potom platí:

1. částečné součty posloupnosti $\{z_i\}_1$ jsou zhora omezeny číslem 1.
2. částečné součty posloupnosti $\{w_i\}_1$ jsou zdola omezeny číslem 0.

■ *Důkaz:*

DŮKAZ JE UVEDEN V ORIGINÁLNÍ TIŠTĚNÉ VERZI.

■

Pouze z důvodů stručného zápisu definujeme skalární součin matic jako běžný skalární součin, kdy matici chápeme jako vektor obsahující její prvky.

Definice 2.1.5 *Nechť \mathbf{A} a \mathbf{B} jsou libovolné matice řádu $m \times n$. Potom pod pojmem SKALÁRNÍ SOUČIN MATIC \mathbf{A} a \mathbf{B} rozumíme číslo*

$$\langle \mathbf{A} | \mathbf{B} \rangle \stackrel{\text{def}}{=} \sum_{i=1}^m \sum_{j=1}^n \mathbf{A}_{i,j} \cdot \mathbf{B}_{i,j}.$$

Závěrem této kapitoly uveďme stěžejní tvrzení, umožňující zhora odhadnout skalární součin matice parity a libovolné rostoucí matice tvořené prvky ± 1 . Jak uvidíme v kapitole o inkluzi tříd booleovských obvodů, tento odhad umožní jednotlivé třídy od sebe separovat.

Tvrzení 2.1.6 *Nechť \mathbf{C} je čtvercová matice dimenze 2^n , $n \geq 1$ taková, že $\mathbf{C}_{i,j} \stackrel{\text{def}}{=} 1$ pro $i + j \leq 2^n + 1$ a $\mathbf{C}_{i,j} \stackrel{\text{def}}{=} -1$ pro $i + j > 2^n + 1$ (tedy matice \mathbf{C} má na vedlejší diagonále a nad ní samé jedničky, pod vedlejší diagonálou jsou samé -1). Dále předpokládejme, že \mathbf{D} je libovolná rostoucí matice dimenze 2^n s prvky z množiny $\{-1, +1\}$. Potom platí, že*

$$2^n(n+1) = \langle \mathbf{B}^{(n)} | \mathbf{C} \rangle \geq \langle \mathbf{B}^{(n)} | \mathbf{D} \rangle.$$

■ *Důkaz:*

DŮKAZ JE UVEDEN V ORIGINÁLNÍ TIŠTĚNÉ VERZI.

■

Poznamenejme, že maximální vlastnost matice \mathbf{C} vyslovená v předešlé větě platí pro širší třídu matic s prvky ± 1 , které mají monotónní pouze sloupce, nebo mají monotónní pouze řádky (monotónní v tomto kontextu znamená, že daný sloupec (řádka) tvoří neklesající posloupnost).

Cvičení 2.1.0 Hadamardovy matice

1. Dokažte, že je-li n dimenze libovolné Hadamardovy matice, potom n je 1, 2 nebo násobek 4. (Stanovte podmínky, které musí splňovat první tři řádky.)
2. Dokažte následující Linsdeyovo lemma:

Lemma 2.1.7 (Linsdey) *Nechť \mathbf{X} je matice řádu $n \times m$, s prvky $\mathbf{X}_{i,j} = \pm 1$, jejíž sloupce jsou ortogonální. Nechť \bar{R} je podmnožina řádkových indexů, \bar{S} je podmnožina sloupcových indexů a $\gamma \stackrel{\text{def}}{=} \sum_{s \in \bar{S}} \sum_{r \in \bar{R}} \mathbf{X}_{r,s}$. Potom platí:*

$$|\gamma| \leq \sqrt{n |\bar{S}| |\bar{R}|}.$$

(Návod: Spočítejte hodnotu $\mathbf{T} \stackrel{\text{def}}{=} \sum_{j=1}^m (\sum_{r \in \bar{R}} \mathbf{X}_{r,j})^2$. Potom porovnejte tuto hodnotu s hodnotou γ . Využijte Cauchyho nerovnosti $\sqrt{\sum_{k=1}^q 1} \sqrt{\sum_{k=1}^q \alpha_k^2} \geq \sum_{k=1}^q \alpha_k$.)

3. V matici $\mathbf{B}^{(n)}$ nalezněte podmatici, která realizuje rovnost v Linsdeyově nerovnosti (tedy ukažte, že odhad daný Linsdeyovou nerovností nelze pro matici $\mathbf{B}^{(n)}$ zlepšit).

2.2 Práh množiny vektorů

Výchozím pojmem pro analýzu booleovských obvodů bude pojem prahového vektoru (viz. např. [RSO94]). Tento pojem nám umožní zúžit naši pozornost nejenom pouze na booleovské funkce definované na celé krychli $\{-1, +1\}^n$, ale i na booleovské funkce definované na libovolné podmnožině diskrétní množiny $\{-1, +1\}^n$. Každou booleovskou funkci můžeme popsat klasicky jako zobrazení $\tilde{b} : \{-1, +1\}^n \rightarrow \{-1, +1\}$. Toto zobrazení lze

popsat buď nějakým předpisem, nebo jako posloupnost uspořádaných dvojic argumentů $\vec{z} \in \{-1, +1\}^n$ a funkčních hodnot $\tilde{b}(\vec{z})$, což máme k dispozici i v případě definice \tilde{b} pomocí nějakého předpisu. Booleovskou funkci \tilde{b}' definovanou pouze na nějaké podmnožině krychle $\{-1, +1\}^n$ můžeme opět popsat posloupností uspořádaných dvojic argumentů, v nichž je \tilde{b}' definována, a odpovídajících funkčních hodnot. Druhou možností, jak popsat \tilde{b}' , je uvést všechny vektory $\vec{x}_1, \dots, \vec{x}_S$ (v našem případě S vektorů), z nichž vektor \vec{x}_1 je tvořen prvními složkami argumentů, ve kterých je \tilde{b}' definována, \vec{x}_2 je tvořen druhými složkami argumentů, v nichž je \tilde{b}' definována (v pevně zvoleném pořadí argumentů), \vec{x}_3 třetími složkami, atd. K těmto vektorům $\vec{x}_1, \dots, \vec{x}_S$ ještě přidáme vektor \vec{y} , který obsahuje odpovídající funkční hodnoty \tilde{b}' (opět ve stejném pořadí). Potom ale tedy systém vektorů $\vec{x}_1, \dots, \vec{x}_S$ a \vec{y} popisuje booleovskou funkci \tilde{b}' .

Dimenze těchto vektorů je zřejmě rovna počtu vrcholů krychle $\{-1, +1\}^n$, ve kterých je \tilde{b}' definována. Jak uvidíme dále, tento – dalo by se říci – „transponovaný“ pohled na způsob definice booleovských funkcí umožňuje prostřednictvím geometrických vlastností vektorů $\vec{x}_1, \dots, \vec{x}_S$ analyzovat vektory \vec{y} , které lze vyjádřit jako vážený součet $\vec{y} = \sum_{i=1}^S w_i \vec{x}_i$, tedy vektory \vec{y} spočitatelné pomocí booleovského perceptronu. Na rozdíl od klasické definice perceptronu nebudeme brát v úvahu existenci prahu, neboť ten lze zavést v našem pojetí přidáním vektoru \vec{x}_{S+1} , který je konstantní a je tvořen buď $+1$ či -1 . Výhodou tohoto postupu je kompaktnější a přehlednější zápis všech tvrzení a důkazů.

Definice 2.2.1 *Nechť funkce $\widetilde{sgn} : \mathbb{R}^n \rightarrow \{-1, +1\}$ je definována jako*

$$\widetilde{sgn}(z) \stackrel{\text{def}}{=} \begin{cases} 1 & \text{pro } x > 0 \\ -1 & \text{pro } x \leq 0 \end{cases}.$$

Definice 2.2.2 *Nechť $\vec{x}_1, \dots, \vec{x}_S$ jsou vektory z prostoru $\{-1, +1\}^n$. Pak vektor \vec{y} je PRAH VEKTORŮ $\vec{x}_1, \dots, \vec{x}_S$, jestliže existují čísla w_1, \dots, w_S tak, že vektor $\sum_{i=1}^S w_i \cdot \vec{x}_i$ má pouze nenulové složky, a platí:*

$$\vec{y} \stackrel{\text{def}}{=} \widetilde{sgn} \left(\sum_{i=1}^S w_i \cdot \vec{x}_i \right).$$

(Funkce \widetilde{sgn} je aplikována na vektor po složkách).

Základní vlastností prahového vektoru je fakt, že není kolmý na žádný vektor \vec{x}_i .

Lemma 2.2.1 *Nechť $\vec{y}^T (\vec{x}_1, \dots, \vec{x}_S) = \vec{0}^T$. Potom \vec{y} není prahem $(\vec{x}_1, \dots, \vec{x}_S)$.*

■ *Důkaz:*

DŮKAZ JE UVEDEN V ORIGINÁLNÍ TIŠTĚNÉ VERZI.

■

Následující tvrzení ukazuje, že pro libovolné různé prahy \vec{y} a \vec{z} systému $\vec{x}_1, \dots, \vec{x}_S$ jsou vektory $\vec{y}^T (\vec{x}_1, \dots, \vec{x}_S)$ a $\vec{z}^T (\vec{x}_1, \dots, \vec{x}_S)$ nutně od sebe různé, což umožňuje horní odhad počtu prahových vektorů.

Lemma 2.2.2 *Nechť \vec{y}^T je prah systému $(\vec{x}_1, \dots, \vec{x}_S)$, $\vec{x}_i \neq \vec{0}$, $i \in \{1, \dots, S\}$ a matice \mathbf{X} má za své sloupce právě vektory \vec{x}_i . Potom pro všechny vektory $\vec{z} \in \{-1, +1\}^n$, $\vec{y} \neq \vec{z}$, platí, že*

$$\vec{y}^T \cdot \mathbf{X} \neq \vec{z}^T \cdot \mathbf{X}.$$

■ *Důkaz:*

DŮKAZ JE UVEDEN V ORIGINÁLNÍ TIŠTĚNÉ VERZI.

■

Opačná implikace však neplatí, což lze dokumentovat na následujícím příkladě.

Příklad 2.2.1 *Zaveďme následující blokové rozdělení matice parity $\mathbf{B}^{(n)}$ definované jako*

$$\mathbf{B}^{(n)} \stackrel{\text{def}}{=} \begin{pmatrix} 1 & \vec{j}^T \\ \vec{j} & \mathbf{A} \end{pmatrix},$$

kde \mathbf{A} je čtvercová matice řádu $2^n - 1$ a \vec{j} je vektor z jedniček dimenze takéž $2^n - 1$. Z vlastností paritních matic bezprostředně vyplývá

$$\begin{pmatrix} 1 & \vec{j}^T \\ \vec{j} & \mathbf{A} \end{pmatrix} \cdot \begin{pmatrix} 1 & \vec{j}^T \\ \vec{j} & \mathbf{A} \end{pmatrix} = \begin{pmatrix} 2^n & \vec{0}^T \\ \vec{0} & 2^n \mathbf{I} \end{pmatrix}.$$

Současně ale z vlastností Schurova doplňku matice $\mathbf{1}$ v matici $\mathbf{B}^{(n)}$ plyne, (viz. [Fie81], str. 31), že tento Schurův doplněk je inverzní maticí k matici \mathbf{A} , tedy

$$\mathbf{A}^{-1} = \frac{1}{2^n} \left(\mathbf{A} - \vec{j} \cdot \vec{j}^T \right).$$

Odtud ale dostáváme, že \mathbf{A}^{-1} je nekladná matice (obsahuje buď nuly, nebo čísla $-\frac{2}{2^n}$).

Předpokládejme nyní, že pro nějaký vektor \vec{w} platí $\mathbf{A}\vec{w} = \vec{r} > \vec{0}$. Potom ale $\vec{w} = \mathbf{A}^{-1} \cdot \vec{r} > \vec{0}$. Poslední nerovnost zřejmě platí na základě nekladnosti matice \mathbf{A}^{-1} . Definujme dále matici

$$\mathbf{V}^{(n)} \stackrel{\text{def}}{=} \begin{pmatrix} -1 & \vec{j}^T \\ \vec{j} & \mathbf{A} \end{pmatrix},$$

která se od matice $\mathbf{B}^{(n)}$ liší pouze ve znaménku prvku vlevo nahoře. Z výše uvedeného rozboru vlastností matice \mathbf{A} vyplývá, že první sloupec matice $\mathbf{V}^{(n)}$ nemůže být prahovým vektorem zbylých sloupců, protože kdyby tomu tak bylo, pak by všechny koeficienty \vec{w}_i byly kladné, ale tím pádem se dostáváme do sporu, protože první indexy sloupců matice $\mathbf{V}^{(n)}$ (vyjma prvního) jsou kladné, a nelze z nich tedy nakombinovat záporné číslo.

Současně ale sloupce matice $\mathbf{V}^{(n)}$ tvoří bazi prostoru \mathbb{R}^{2^n} , přičemž první sloupec neleží v ortogonálním doplňku ostatních. Proto pro libovolné různé vektory $\vec{z}, \vec{y} \in \{-1, +1\}^{2^n}$ se vektory

$$\vec{z}^T \mathbf{V}^{(n)} \quad \text{a} \quad \vec{y}^T \mathbf{V}^{(n)}$$

musí lišit alespoň v jedné ze složek s indexem větším než 1. To dokazuje, že tvrzení lemma 2.2.2 nelze obrátit.

Zřejmě na základě lemma 2.2.2 lze vyslovit následující tvrzení.

Lemma 2.2.3 *Pro systém vektorů $\vec{x}_1, \dots, \vec{x}_s \in \{-1, +1\}^n$ existuje nanejvýše $(n+1)^S$ různých prahů.*

■ *Důkaz:*

DŮKAZ JE UVEDEN V ORIGINÁLNÍ TIŠTĚNÉ VERZI.

Významnou roli v analýze aproximačních možností booleovských obvodů hrají souřadnice prahového vektoru \vec{y} vzhledem k systému $\vec{x}_1, \dots, \vec{x}_S$.

Definice 2.2.3 Předpokládejme, že vektory $\vec{x}_1, \dots, \vec{x}_S$ jsou lineárně nezávislé a nechť

$$\vec{y} \stackrel{\text{def}}{=} \sum_{i=1}^S \beta_i \vec{x}_i + \vec{y}^\perp,$$

kde \vec{y}^\perp je z ortogonálního doplňku $[\vec{x}_1, \dots, \vec{x}_S]_\lambda$. Potom vektor $\vec{\beta} \stackrel{\text{def}}{=} (\beta_1, \dots, \beta_S)$ nazveme ZOBECNĚNÉ SPEKTRUM VEKTORU \vec{y} a číslo

$$\beta_{\max} \stackrel{\text{def}}{=} \max \{ |\beta_i| \mid i \in \{1, \dots, S\} \}$$

nazveme MAXIMUM SPEKTRA vektoru \vec{y} .

Čísla β_i lze spočítat na základě pojmu pseudoinverzní matice takto: Označme si $\mathbf{X} \stackrel{\text{def}}{=} (\vec{x}_1, \dots, \vec{x}_S)^T$. Potom platí

$$\vec{\beta} = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \vec{y}.$$

Základní vlastností vektoru $\vec{\beta}$ je následující odhad součtu absolutních hodnot jeho složek zdola.

Lemma 2.2.4 Nechť \vec{y} je práh systému $(\vec{x}_1, \dots, \vec{x}_S)$. Potom platí

$$\sum_{i=1}^S |\vec{\beta}_i| \geq 1.$$

■ *Důkaz:*

DŮKAZ JE UVEDEN V ORIGINÁLNÍ TIŠTĚNÉ VERZI.

Předešlé tvrzení poskytuje základní vztah mezi hodnotami zobecněného spektra a počtem S vektorů \vec{x}_i . Toto tvrzení, které je hlavním závěrem této části výkladu, lze využít k důkazu neexistence polynomiálního počtu vektorů $\vec{x}_1, \dots, \vec{x}_S$ takových, že \vec{y} je jejich práhem. Princip tohoto důkazu je ten, že shora odhadneme hodnoty $|\vec{\beta}_i|$, a budou-li polynomiálně malé vzhledem k dimenzi \vec{y} , pak S nemůže být omezeno polynomem vzhledem k dimenzi \vec{y} .

Tvrzení 2.2.5 Nechť \vec{y} je práh ortogonálních vektorů $\vec{x}_1, \dots, \vec{x}_S \in \{-1, +1\}^n$. Pak

$$S \geq \frac{n}{\max_i \langle \vec{y} | \vec{x}_i \rangle}. \quad (2.2)$$

■ *Důkaz:*

DŮKAZ JE UVEDEN V ORIGINÁLNÍ TIŠTĚNÉ VERZI.

Ukážeme ještě jednu formu tvrzení předešlého, zahrnujícího v sobě i odhad na velikost maximální váhy prahového vektoru.

Lemma 2.2.6 *Předpokládejme, že \vec{y} je práh systému $(\vec{x}_1, \dots, \vec{x}_S)$ vektorů z $\{-1, +1\}^n$ a že příslušné hodnoty w_1, \dots, w_s jsou celočíselné. Nechť $\tilde{w} \stackrel{\text{def}}{=} \max_i |w_i|$ a $\tilde{\tau} \stackrel{\text{def}}{=} \max_i |\langle \vec{y} | \vec{x}_i \rangle|$. Potom platí*

$$S \geq \frac{n}{\tilde{w}\tilde{\tau}}.$$

■ *Důkaz:*

DŮKAZ JE UVEDEN V ORIGINÁLNÍ TIŠTĚNÉ VERZI.

■

2.2.1 Vlastnosti prahových vektorů

Nyní se pokusíme analyzovat základní vlastnosti prahových vektorů obecně, bez ohledu na konkrétní systém $\vec{x}_1, \dots, \vec{x}_S$. Následující příklad ilustruje dobře známou skutečnost, že práh dvou vektorů $\vec{x}_1 \stackrel{\text{def}}{=} (1, -1, 1, -1)^T$, $\vec{x}_2 \stackrel{\text{def}}{=} (1, 1, -1, -1)^T$ nemůže realizovat logickou funkci XOR , definovanou na sobě si odpovídajících složkách vektorů \vec{x}_1 a \vec{x}_2 .

Příklad 2.2.2 *Jesliže přijmeme předpoklad, že platí rovnost*

$$\begin{pmatrix} -1 \\ 1 \\ 1 \\ -1 \end{pmatrix} = \widetilde{sgn} \left(\beta \begin{pmatrix} 1 \\ 1 \\ -1 \\ -1 \end{pmatrix} + \gamma \begin{pmatrix} 1 \\ -1 \\ 1 \\ -1 \end{pmatrix} \right),$$

potom musí být splněno

$$\left. \begin{array}{l} \beta + \gamma < 0 \\ \beta - \gamma > 0 \\ \beta + \gamma > 0 \\ \beta - \gamma < 0 \end{array} \right\} \rightarrow \left. \begin{array}{l} \beta + \gamma = 0 \\ \beta - \gamma = 0 \end{array} \right\} \rightarrow \left. \begin{array}{l} \beta = 0 \\ \gamma = 0 \end{array} \right\},$$

což je ve sporu s definicí prahového vektoru. Alternativně lze tento závěr obdržet na základě skutečnosti, že všechny tři uvažované vektory jsou vzájemně ortogonální, tedy žádný z nich nemůže být prahem zbývajících dvou (stejná argumentace platí i pro případ, kdy přidáme navíc konstantní vektor).

Jinými slovy booleovský perceptron se dvěma vstupními neurony a jedním neuronem výstupním nemůže počítat logickou funkci XOR , definovanou na svých vstupech. Bohužel tento negativní závěr byl v šedesátých letech interpretován jako charakterizující vlastnost sítí tvořených více booleovskými perceptrony, čímž došlo k výraznému zpoždění vývoje v tomto oboru.

V dalším nyní ukážeme, že tato skutečnost, popsaná v příkladě 2.2.2, má mnohem hlubší podstatu a lze ji velmi zobecnit. Definujme proto následující matici $\mathbf{M}_{\vec{y}}$.

Definice 2.2.4 *Nechť $n \stackrel{\text{def}}{=} 2k$ a $\vec{y} \in \{-1, +1\}^{2^n}$. Nechť $\mathbf{M}_{\vec{y}}$ je čtvercová matice řádu $2^{\frac{n}{2}}$, jejíž sloupce jsou tvořeny po sobě následujícími subvektory vektoru \vec{y} délky $2^{\frac{n}{2}}$. Matici $\mathbf{M}_{\vec{y}}$ nazveme SDRUŽENOU MATICÍ pro vektor \vec{y} .*

Například pro vektor $\vec{y} \stackrel{\text{def}}{=} (a, b, c, d)^T$ je $\mathbf{M}_{\vec{y}} = \begin{pmatrix} a & c \\ b & d \end{pmatrix}$. Při analýze sdrúžených matic bude užitečná lemma o vlastnostech rostoucích matic.

Lemma 2.2.7 *Platí následující tvrzení:*

1. *Je-li matice \mathbf{A} rostoucí (potenciálně rostoucí), je i matice $\widetilde{\text{sgn}}(\mathbf{A})$ rostoucí (potenciálně rostoucí).*
2. *Nechť prvky matice \mathbf{A} jsou tvořeny čísly ± 1 . Potom je-li $\widetilde{\text{sgn}}(\mathbf{A})$ rostoucí (potenciálně rostoucí), je i matice \mathbf{A} rostoucí (potenciálně rostoucí).*

■ *Důkaz:*

DŮKAZ JE UVEDEN V ORIGINALNÍ TIŠTĚNÉ VERZI. ■

Jak již bylo řečeno, při našem výkladu využíváme popisu booleovských funkcí pomocí pojmu prahových vektorů, které kromě jiného umožňují popis a analýzu booleovských funkcí definovaných pouze v nějaké podmnožině $\{-1, +1\}^n$. Pro popis booleovské funkce definované v celém oboru $\{-1, +1\}^n$, musíme používat takový systém vektorů $\vec{x}_1, \dots, \vec{x}_S$, který má tu vlastnost, že libovolná s -členná posloupnost tvořená z ± 1 je obsažena mezi řádky matice $\mathbf{X} \stackrel{\text{def}}{=} (\vec{x}_1, \dots, \vec{x}_S)$. Jedním (a nejpřirozenějším) z těchto systémů je systém základních vektorů parity podle následující definice.

Definice 2.2.5 *Předpokládejme, že sloupce paritní matice $\mathbf{B}^{(n)}$ jsou očíslovány počínaje nulou. Potom vektor tvořící sloupec matice $\mathbf{B}^{(n)}$, jehož index v binárním zápisu obsahuje pouze jednu jedničku na i -té pozici zleva, nazveme i -tý ZÁKLADNÍ VEKTOR PARITY a budeme ho značit $\tilde{\mathbf{p}}^{(n)}_i$.*

Příklad 2.2.3 *Pro lepší čitelnost předešlé definice uveďme následující příklad základních vektorů parity pro $n = 4$:*

$$\tilde{\mathbf{p}}^{(4)}_1 = (1, -1, 1, -1, 1, -1, 1, -1, 1, -1, 1, -1, 1, -1, 1, -1)^T,$$

$$\tilde{\mathbf{p}}^{(4)}_2 = (1, 1, -1, -1, 1, 1, -1, -1, 1, 1, -1, -1, 1, 1, -1, -1)^T,$$

$$\tilde{\mathbf{p}}^{(4)}_3 = (1, 1, 1, 1, -1, -1, -1, -1, 1, 1, 1, 1, -1, -1, -1, -1)^T,$$

$$\tilde{\mathbf{p}}^{(4)}_4 = (1, 1, 1, 1, 1, 1, 1, 1, -1, -1, -1, -1, -1, -1, -1, -1)^T,$$

(Srovnej s obrázkem 2.1).

Nechť \mathbf{Z} je matice z čísel 0 a 1 o n sloupcích a 2^n řádcích a nechť j -tý řádek této matice je roven binárnímu zápisu čísla $j - 1$. Dále předpokládejme, že matice \mathbf{A} vznikla z matice \mathbf{Z} záměnou 1 za -1 a 0 za 1. Potom $\mathbf{A} = (\tilde{\mathbf{p}}^{(n)}_n, \dots, \tilde{\mathbf{p}}^{(n)}_1)$.

Definice 2.2.6 *Označme symbolem $\mathbf{E}^{(n)}$ čtvercovou matici řádu 2^n obsahující samé jedničky. Symbolem $\mathbf{I}^{(n)}$ označme jednotkovou matici řádu 2^n .*

Algebraickou strukturu matic sdružených k základním vektorům parity popisuje následující lemma.

Lemma 2.2.8 *Předpokládejme, že $n \stackrel{\text{def}}{=} 2k$ a $\tilde{\mathbf{p}}^{(n)}_1, \dots, \tilde{\mathbf{p}}^{(n)}_n$ jsou odpovídající základní vektory parity. Potom platí:*

1. pro $1 \leq j \leq k$ je

$$M_{\tilde{\mathbf{p}}_j^{(n)}} = \mathbf{E}^{(j-1)} \odot \begin{pmatrix} 1 & 1 \\ -1 & -1 \end{pmatrix} \odot \mathbf{E}^{(k-j)}.$$

a

$$M_{\tilde{\mathbf{p}}_{k+j}^{(n)}} = \mathbf{E}^{(j-1)} \odot \begin{pmatrix} 1 & -1 \\ 1 & -1 \end{pmatrix} \odot \mathbf{E}^{(k-j)}.$$

2. Pro $1 \leq j \leq k$ platí, že

$$M_{\tilde{\mathbf{p}}_{k+j}^{(n)}} = M_{\tilde{\mathbf{p}}_j^{(n)}}^T.$$

■ *Důkaz:*

DŮKAZ JE UVEDEN V ORIGINALNÍ TIŠTĚNÉ VERZI.

■

Pro aplikaci předchozího algebraického popisu sružených matic základních vektorů parity bude užitečná následující identita, platná pro tenzorový součin matic.

Lemma 2.2.9 *Nechť matice A , B a C , D jdou spolu násobit v tomto pořadí. Potom platí identita*

$$(\mathbf{A} \odot \mathbf{B}) \cdot (\mathbf{C} \odot \mathbf{D}) = (\mathbf{A} \cdot \mathbf{C}) \odot (\mathbf{B} \cdot \mathbf{D}).$$

■ *Důkaz:*

DŮKAZ JE UVEDEN V ORIGINALNÍ TIŠTĚNÉ VERZI.

■

Za hlavní závěr této podkapitoly lze považovat tvrzení, jehož obsahem je skutečnost, že každá lineární kombinace sružených matic základních vektorů parity je potenciálně rostoucí.

Tvrzení 2.2.10 *Nechť $n \stackrel{\text{def}}{=} 2k$ a $\tilde{\mathbf{y}}$ je libovolná lineární kombinace základních vektorů parity $\tilde{\mathbf{p}}_1^{(n)}, \dots, \tilde{\mathbf{p}}_n^{(n)}$. Potom matice $M_{\tilde{\mathbf{y}}}$ je potenciálně rostoucí.*

■ *Důkaz:*

DŮKAZ JE UVEDEN V ORIGINALNÍ TIŠTĚNÉ VERZI.

■

Evidentním důsledkem je následující tvrzení:

Tvrzení 2.2.11 *Nechť $n \stackrel{\text{def}}{=} 2k$, $\tilde{\mathbf{y}} \in \{-1, +1\}^{2n}$ je práh systému základních vektorů parity a konstantního vektoru. Potom matice $M_{\tilde{\mathbf{y}}}$ je potenciálně rostoucí.*

■ *Důkaz:*

DŮKAZ JE UVEDEN V ORIGINALNÍ TIŠTĚNÉ VERZI.

■

Vezměme v úvahu booleovskou funkci $\overset{\sqcup}{\text{XOR}}$ definovanou na krychli $\{-1, +1\}^2$. Potom odpovídající základní vektory parity jsou $(1, -1, 1, -1)^T$ a $(1, 1, -1, -1)^T$. Vektor $\tilde{\mathbf{y}}$ funkčních hodnot $\overset{\sqcup}{\text{XOR}}$ je přitom roven $\tilde{\mathbf{y}} = (-1, 1, 1, -1)^T$. Potom ale $M_{\tilde{\mathbf{y}}} = \begin{pmatrix} -1 & 1 \\ 1 & -1 \end{pmatrix}$

není potenciálně rostoucí matice, tedy \vec{y} nemůže být prahem základních vektorů parity a konstantního vektoru.

Algebraická struktura základních vektorů parity nám umožní přesně spočítat počet od sebe různých prahových vektorů tohoto systému. Zdůrazněme však, že tento počet není roven počtu booleovských funkcí realizovaných binárním perceptronem v klasickém slova smyslu (ten není znám, v dalším uvedeme pouze některé známé odhady), nýbrž počtu booleovských funkcí realizovaných perceptronem s nulovým prahem. Bohužel zavedení nenulového prahu podstatným způsobem tuto analýzu počtu prahů zkomplikuje. Vezmeme-li ale v úvahu pouze nulový prah, platí následující.

Tvrzení 2.2.12 *Pro $n \stackrel{\text{def}}{=} 2k$ je počet různých prahů systému $\vec{p}_1^{(n)}, \dots, \vec{p}_n^{(n)}$ základních vektorů parity roven číslu $2^{k(k+1)}$.*

■ *Důkaz:*

DŮKAZ JE UVEDEN V ORIGINALNÍ TIŠTĚNÉ VERZI.

■

2.3 Inkluze základních tříd booleovských obvodů

V následujícím obrátíme pozornost na množiny booleovských funkcí, které lze vyjádřit booleovskými neuronovými sítěmi (obvody) různé hloubky. Zavedeme třídy booleovských funkcí, n -proměnných, které lze spočítat obvody, jejichž počet uzlů je polynomiálně omezen v n .

Uvidíme, že důležitou roli hrají tzv. α -symetrické booleovské vektory. Na základě výkladu o prahových vektorech ukážeme způsob separace jednotlivých tříd prahových vektorů počítaných obvody různé hloubky. Začneme s definicí booleovského obvodu.

Definice 2.3.1 **BOOLEOVSKÝ OBVOD** je acyklický orientovaný graf s hranovým a vrcholovým ohodnocením a jedinným vrcholem, ze kterého nevede žádná hrana. Vrchol tohoto grafu nazýváme **VSTUPNÍ VRCHOL**, nevede-li do něj žádná hrana, **VÝSTUPNÍ VRCHOL**, nevede-li z něj žádná hrana, a **VNITŘNÍ VRCHOL**, není-li ani vstupní, ani výstupní. Hodnoty hranového ohodnocení jsou obecně reálná čísla, ohodnocení hrany z vrcholu i do vrcholu j nazveme **VÁHOU HRANY $i \rightarrow j$** a budeme značit $w_{i \rightarrow j}$.

Ohodnocení vrcholů jsou uspořádané dvojice $(v, t) \in \{-1, +1\} \times \mathbb{R}$, které pro vrchol i budeme budeme značit jako (v_i, t_i) . Hodnotu v_i u vstupních vrcholů nazýváme **ustupními hodnotami**, ohodnocení v_i u výstupních vrcholů, **výstupními hodnotami**. Dále, ohodnocení vnitřního a výstupního vrcholu i splňuje rovnost

$$v_i \stackrel{\text{def}}{=} \widetilde{\text{sgn}} \left(\sum_{\{j|j \rightarrow i\}} w_{j \rightarrow i} \cdot v_j - t_i \right).$$

Definice 2.3.2 **VELIKOST BOOLEOVSKÉHO OBVODU** je počet vnitřních vrcholů. **HLOUBKA VRCHOLU v** je definována jako délka nejdelší cesty mezi libovolným vstupním vrcholem a vrcholem v . **HLOUBKA BOOLEOVSKÉHO OBVODU** je rovna maximu hloubek všech vnitřních vrcholů. **DIMENZE BOOLEOVSKÉHO OBVODU** je rovna počtu vstupních vrcholů.

V této definici booleovského obvodu připouštíme existenci prahových hodnot t_i jednotlivých uzlů (neuronů) booleovského obvodu. Toto rozšíření ve srovnání s definicí prahového vektoru systému $\vec{x}_1, \dots, \vec{x}_S$ ale stále umožňuje používat všech vlastností prahových vektorů při analýze booleovských obvodů, protože existenci prahových hodnot t_i v uzlech booleovského obvodu lze simulovat existencí konstantního vektoru v systému $\vec{x}_1, \dots, \vec{x}_S$.

Definice 2.3.3 *Booleovský obvod dimenze n , $n \geq 1$, počítá vektor $\vec{y} \in \{-1, +1\}^{2^n}$, právě když existuje takové očíslování vstupních vrcholů i_1, \dots, i_n a takové hranové ohodnocení, že pro každé ohodnocení vstupních vrcholů čísla ± 1 , je ohodnocení výstupního vrcholu rovno číslu \vec{y}_γ , kde index γ má v binárním zápisu na pozici k nulu, je-li $i_k = -1$, a jedničku, je-li $i_k = +1$.*

Pro lepší pochopení přesného významu předešlé definice uveďme jednoduchý příklad. Vektor $\vec{y} \stackrel{\text{def}}{=} (1, 1, -1, 1, -1, -1, 1, -1)^T$ je počítán booleovským obvodem B , právě když pro zobrazení $\tilde{B} : \{-1, +1\}^3 \rightarrow \{-1, +1\}$, které obvod realizuje platí

$$\begin{aligned} \tilde{B} \left(\begin{pmatrix} 1 & 1 & 1 \end{pmatrix}^T \right) &= 1, \\ \tilde{B} \left(\begin{pmatrix} 1 & 1 & -1 \end{pmatrix}^T \right) &= 1, \\ \tilde{B} \left(\begin{pmatrix} 1 & -1 & 1 \end{pmatrix}^T \right) &= -1, \\ \tilde{B} \left(\begin{pmatrix} 1 & -1 & -1 \end{pmatrix}^T \right) &= 1, \\ \tilde{B} \left(\begin{pmatrix} -1 & 1 & 1 \end{pmatrix}^T \right) &= -1, \\ \tilde{B} \left(\begin{pmatrix} -1 & 1 & -1 \end{pmatrix}^T \right) &= -1, \\ \tilde{B} \left(\begin{pmatrix} -1 & -1 & 1 \end{pmatrix}^T \right) &= 1, \\ \tilde{B} \left(\begin{pmatrix} -1 & -1 & -1 \end{pmatrix}^T \right) &= -1. \end{aligned}$$

Z tohoto příkladu je souvislost mezi hodnotami vstupních uzlů a hodnotami složek vektoru \vec{y} již zřejmá.

Dále definujeme třídy prahových vektorů, počítaných polynomiálně velkými booleovskými obvody hloubky d .

Definice 2.3.4 *Předpokládejme, že máme zadánu posloupnost*

$$\bar{A} \stackrel{\text{def}}{=} \left\{ \vec{y}_n \in \{-1, +1\}^{2^n} \mid n \in \{1, \dots, \infty\} \right\},$$

a nechtě $\tilde{p}(x)$ je daný reálný polynom. Potom \bar{A} leží ve třídě $\mathbf{LT}_{d, \tilde{p}}$, jestliže pro každé $n \geq 1$ existuje booleovský obvod dimenze n , jehož hloubka je nejvýše d a velikost nejvýše $\tilde{p}(n)$, takový, že počítá vektor \vec{y}_n . Dále definujeme třídu

$$\mathbf{LT}_d \stackrel{\text{def}}{=} \bigcup_{\{\tilde{p} \mid \tilde{p} \text{ je polynom}\}} \mathbf{LT}_{d, \tilde{p}}.$$

Nechtě navíc pro příslušné booleovské obvody platí, že jejich váhy jsou omezeny zhora hodnotou $\tilde{p}(n)$. Potom příslušnou třídu budeme značit $\widehat{\mathbf{LT}}_d$.

Definice 2.3.5 *Vektory tvořící sloupce paritní matice $\mathbf{B}^{(n)}$ budeme nazývat vektory parity.*

Definice 2.3.6 *Nechť $\bar{A} \stackrel{\text{def}}{=} \{\vec{y}_n \in \{-1, +1\}^{2^n} \mid n \in \{1, \dots, \infty\}\}$ a necht $\tilde{p}(x)$ je daný reálný polynom. Potom řekneme, že \bar{A} leží ve třídě $\mathbf{PT}_{1, \tilde{p}}$, existuje-li přirozené číslo $k \leq \tilde{p}(n)$ a systém vektorů parity $\vec{p}^{(n)}_1, \dots, \vec{p}^{(n)}_k \in \{-1, +1\}^{2^n}$, tak, že vektor \vec{y}_n je jejich prahem. Dále definujme třídu*

$$\mathbf{PT}_1 \stackrel{\text{def}}{=} \bigcup_{\{\tilde{p} \mid \tilde{p} \text{ je polynom}\}} \mathbf{PT}_{1, \tilde{p}}.$$

Nechť navíc příslušné koeficienty w_i jsou omezeny zhora hodnotou $\tilde{p}(n)$. Potom příslušnou třídu budeme značit $\widehat{\mathbf{PT}}_1$.

Jak již bylo řečeno v předešlé kapitole (viz. tvrzení 2.2.5), pro odhad počtu vektorů $\vec{x}_1, \dots, \vec{x}_S$, pro které je dané \vec{y} prahovým vektorem, lze zdola odhadnout na základě znalosti hodnot standardního skalárního součinu vektoru \vec{y} s vektory $\vec{x}_1, \dots, \vec{x}_S$.

Vlastnosti vektoru \vec{g} , definovaného v následující lemmě, nám umožní od sebe odseparovat třídy \mathbf{LT}_2 a \mathbf{PT}_1 .

Lemma 2.3.1 *Nechť $n = 2k$, i^b je binární zápis čísla i a definujme vektor $\vec{g}^{(n)}$ dimenze 2^n předpisem*

$$\vec{g}_i^{(n)} \stackrel{\text{def}}{=} \begin{cases} -1 & \sum_{l=1}^n (i^b)_l \equiv 4 \pmod{0, 1} \\ \text{pro} & \\ 1 & \sum_{l=1}^n (i^b)_l \equiv 4 \pmod{2, 3}. \end{cases}$$

Potom pro libovolné vektory parity $\vec{p}^{(n)}$ (t.j. sloupce matice $\mathbf{B}^{(n)}$) platí

$$\left\langle \vec{g}^{(n)} \mid \vec{p}^{(n)} \right\rangle \leq 2^{\frac{n}{2}}.$$

■ *Důkaz:*

DŮKAZ JE UVEDEN V ORIGINÁLNÍ TIŠTĚNÉ VERZI.

■

Nyní zavedeme další třídu prahových vektorů, kterou jak uvidíme, lze spočítat velmi snadno booleovským obvodem hloubky 2 lineární velikosti.

Definice 2.3.7 *Nechť n je přirozené číslo, $\vec{y} \in \{-1, +1\}^{2^n}$ a $\alpha \in \{0, \dots, 2^n - 1\}$. Potom vektor \vec{y} je α -SYMETRICKÝ VEKTOR, právě když hodnota \vec{y}_i závisí pouze na čísle*

$$\sum_{k=1}^n (i^b)_k \cdot (\alpha^b)_k$$

(tedy pouze na součtu jedniček v binárním zápise i , které jsou na těch pozicích, kde je v binárním zápise čísla α hodnota 1). Dále budeme symbolem \mathbf{SYM}_2 označovat množinu všech α -symetrických vektorů pro $\alpha \in \{0, \dots, 2^n - 1\}$.

Tvrzení 2.3.2 *Nechť n je přirozené číslo, $\vec{y} \in \{-1, +1\}^{2^n}$ je $(2^n - 1)$ -symetrický vektor a necht čísla d_1, \dots, d_m jsou délky po sobě jdoucích konstantních úseků v posloupnosti*

$$\vec{y}_{2^{d_1}-1}, \vec{y}_{2^{d_2}-1}, \dots, \vec{y}_{2^{d_m}-1}. \quad (2.3)$$

Dále pro $k \in \{1, \dots, m\}$ definujeme funkce $\widetilde{v}_k : \{-1, +1\}^n \rightarrow \{-1, +1\}$ jako

$$\widetilde{v}_k(\vec{x}) \stackrel{\text{def}}{=} \text{sgn} \left(\sum_{j=1}^n \vec{x}_j + n - 2 \sum_{j=1}^k d_j - \frac{1}{2} \right). \quad (2.4)$$

Potom platí, že

$$-\vec{y}_1 \left(1 - \sum_{j=1}^m (-1)^j (\widetilde{v}_j(\vec{x}) - 1) \right) = \vec{y}_{\vec{x}^\#}. \quad (2.5)$$

■ *Důkaz:*

DŮKAZ JE UVEDEN V ORIGINALNÍ TIŠTĚNÉ VERZI.

■

Toto tvrzení v podstatě popisuje booleovský obvod pro výpočet libovolného α -symetrického vektoru. Architektura tohoto vektoru je popsána na obrázku 2.3. Obsah předešlého tvrzení lze formulovat také následovně.

Důsledek 2.3.3 *Každý α -symetrický vektor je možno spočítat booleovským obvodem hloubky dva, který má lineárně omezené váhy a lineárně omezený počet vrcholů.*

Protože pro každé α je α -tý vektor parity α -symetrickým vektorem, lze vektory parity realizovat ve třídě \mathbf{LT}_2 . Uvažujme booleovský obvod zobrazený na obrázku 2.3. Nechť \bar{A} je množina těch vstupních vrcholů, pro které je $(\alpha^b)_i = 0$. Dále ohodnoťme veškeré hrany jdoucí z vrcholů v množině \bar{A} hodnotou 0. Nechť všechny zbylé hrany jsou ohodnoceny způsobem popsaným v tvrzení 2.3.2. Potom takovýto obvod hloubky 2 zřejmě počítá α -tý vektor parity.

Následující dvě tvrzení popisují základní inkluze mezi třídami \mathbf{LT}_2 , \mathbf{PT}_1 , \mathbf{LT}_3 , \mathbf{SYM}_2 . Důkazy obou dvou tvrzení jsou analogické, založené na odhadu ve výrazu 2.2.

Jednotlivé inkluze jsou vesměs zřejmé z předešlého výkladu. Ostrost inkluzí se obecně dokáže tak, že pro vybraný vektor \vec{y} spočítáme jeho standardní skalární součin se základními vektory parity a tyto hodnoty použijeme v odhadu pro počet vektorů \vec{x}_i nezbytných k tomu, aby \vec{y} byl jejich prahem. Tímto způsobem vždy dokážeme, že počet vektorů $\vec{x}_1, \dots, \vec{x}_S$ nemůže být polynomiálně pro daný vektor \vec{y} omezen.

Tvrzení 2.3.4 *Platí inkluze:*

1. $\mathbf{PT}_1 \subset \mathbf{LT}_3$ a $\widehat{\mathbf{PT}}_1 \subset \widehat{\mathbf{LT}}_3$,
2. $\mathbf{SYM}_2 \not\subset \mathbf{PT}_1$ a $\widehat{\mathbf{SYM}}_2 \not\subset \widehat{\mathbf{PT}}_1$.

■ *Důkaz:*

DŮKAZ JE UVEDEN V ORIGINALNÍ TIŠTĚNÉ VERZI.

■

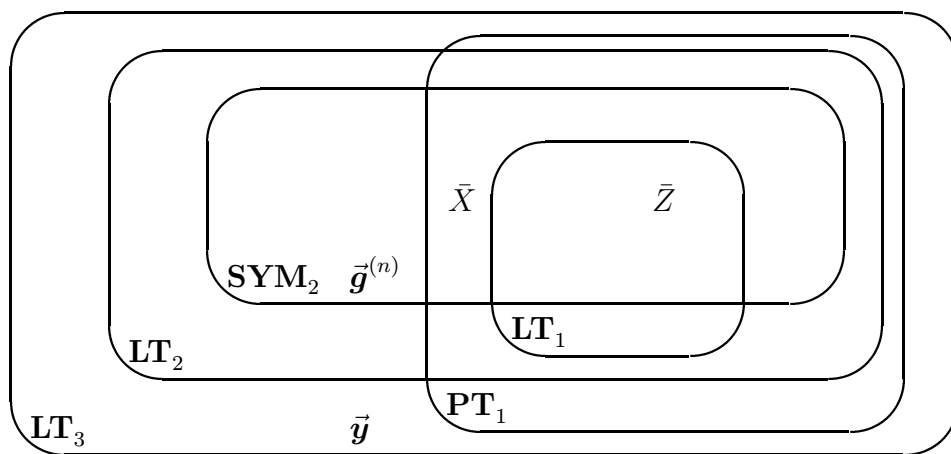
Tvrzení 2.3.5 $\mathbf{LT}_2 \subset \mathbf{LT}_3$ a $\mathbf{LT}_2 \neq \mathbf{LT}_3$.

■ *Důkaz:*

DŮKAZ JE UVEDEN V ORIGINÁLNÍ TIŠTĚNÉ VERZI.

■

Nyní shrneme dosavadní poznatky o inkluzi tříd booleovských obvodů. Necht' \vec{y} je takový, že $M_{\vec{y}} = B^{(n)}$, $n = 2k$, \bar{Z} označuje množinu všech základních vektorů parity a \bar{X} označuje množinu ostatních vektorů parity. Potom platí inkluze popsané následujícím diagramem:



Obrázek 2.3: Inkluze základních tříd booleovských obvodů.

Cvičení 2.3.0 Booleovské obvody

1. Navrhněte booleovský obvod hloubky 2, počítající vektor $\vec{g}^{(n)}$ definovaný v lemmě 2.3.1.
2. Navrhněte booleovský obvod hloubky 3, počítající vektor, jehož sdružená matice je rovna matici parity $B^{(n)}$.

2.4 Odhady velikosti vah prahového vektoru

Důležitou informací o booleovských obvodech je také odhad velikosti vah, nezbytných pro výpočet daného vektoru \vec{y} . Dosud jsme se tímto problémem nezabývali, ukážeme ale, že úzce souvisí s horním a dolním odhadem počtu prahů systému vektoru $\vec{x}_1, \dots, \vec{x}_S$, kde alespoň jeden z vektorů \vec{x}_i je konstantní.

Protože pro daný váhový vektor prahu \vec{y} je jeho libovolný nenulový násobek opět váhovým vektorem pro práh \vec{y} , mohou být velikosti vah libovolné. Nás ale zajímají velikosti vah především z hlediska posouzení toho, zda jsme schopni pro daný booleovský obvod jeho váhy uložit v paměti – zajímá nás tedy de fakto odhad na počet bitů nutných k uložení vah booleovského obvodu. Zřejmě vzhledem k hustotě racionálních čísel na reálné ose stačí brát v úvahu pouze racionální váhy. Proto můžeme bez újmy na obecnosti brát na zřetel pouze takové váhové vektory, které mají celočíselné složky, a zajímat se o

horní odhad velikosti těchto složek. Navíc většina tvrzení umožňuje zhora odhadnout tzv. celkovou váhu, tedy číslo $\sum_{i=1}^S |w_i|$, což je ale k verifikaci polynomiálního odhadu počtu bitů nutných k uložení vah postačující.

2.4.1 Dolní odhad – diskriminační lemma

Nejdříve vyslovíme tzv. diskriminační lemmu a na následujícím příkladě ukážeme její užitečnost pro odhad velikosti celkové váhy.

Lemma 2.4.1 *Nechť $\vec{x}_1, \dots, \vec{x}_S$ jsou z $\{-1, +1\}^n$ a vektor \vec{y} je prahem tohoto systému s celočíselnými koeficienty lineární kombinace w_1, \dots, w_S a $w \stackrel{\text{def}}{=} \sum_{j=1}^n |w_j|$. Nechť dále \mathbf{B} je čtvercová matice řádu n , pro kterou platí, že*

$$(\forall k \in \{1, \dots, n\}) ((\mathbf{B}\vec{y})_k \cdot \vec{y}_k \geq 0)$$

Potom platí:

$$w \cdot \max_{i \in \{1, \dots, S\}} \{|\langle \vec{x}_i | \mathbf{B}\vec{y} \rangle|\} \geq \sum_{k=1}^n |(\mathbf{B}\vec{y})_k|.$$

■ *Důkaz:*

DŮKAZ JE UVEDEN V ORIGINÁLNÍ TIŠTĚNÉ VERZI.

■

Význam diskriminační lemmy pro odhad velikosti vah spočívá v možnosti široké volby matice \mathbf{B} . Při vhodné volbě této matice můžeme pro daný systém vektorů dostat na základě tvrzení diskriminační lemmy dolní odhad na součet absolutních hodnot vah. Ilustrujme tuto skutečnost na následujících příkladech.

Příklad 2.4.1 *Nechť $\mathbf{T}^{(k)}$ je čtvercová matice řádu 2^k taková, že na hlavní diagonále a všude nad ní jsou jedničky a všude pod hlavní diagonálou jsou pouze -1 . Potom lze velmi snadno ověřit (neboť $\mathbf{M}_{\vec{p}_{\tilde{(2k)}_j}} = \mathbf{E}^{(j-1)} \odot \begin{pmatrix} 1 & 1 \\ -1 & -1 \end{pmatrix} \odot \mathbf{E}^{(k-j)}$ a $\mathbf{M}_{\vec{p}_{\tilde{(2k)}_{k+j}}} = \mathbf{M}_{\vec{p}_{\tilde{(2k)}_j}}^T$), že pro i -tý základní vektor parity $\vec{p}_{\tilde{(2k)}_i}$, $i \in \{1, \dots, k\}$, platí*

$$\left\langle \mathbf{T}^{(k)} \left| \mathbf{M}_{\vec{p}_{\tilde{(2k)}_i}} \right. \right\rangle = 2^{i+1}, \quad \left\langle \mathbf{T}^{(k)} \left| \mathbf{M}_{\vec{p}_{\tilde{(2k)}_{k+i}}} \right. \right\rangle = -2^{i+1}, \quad \left\langle \mathbf{T}^{(k)} \left| \mathbf{E}^{(k)} \right. \right\rangle = 2^k.$$

Nechť $\vec{y} \in \{-1, +1\}^{2^{2k}}$ je takový vektor, že matice $\mathbf{T}^{(k)}$ je k němu sdružená. Na základě diskriminační lemmy a z předchozích rovností dostáváme při volbě matice \mathbf{B} rovné identické matici, že celková váha \vec{y} , jakožto prahu základních vektorů parity, musí být alespoň rovna 2^{k-1} .

Tento odhad ale zdaleka není optimální, ukážeme však v následujícím příkladě, že vhodnou volbou matice \mathbf{B} ho lze podstatně zlepšit.

Příklad 2.4.2 *Nyní pro stejný vektor \vec{y} jako v příkladě předešlém provedeme odhad celkové váhy s využitím matice \mathbf{B} , zkonstruované následovně. Nechť pomocná matice \mathbf{G} je čtvercová řádu 2^k , mající těsně pod a nad hlavní diagonálou jedničky a na ostatních místech nuly s výjimkou prvku na pozici $(1, 1)$ (při číslování řad od jedné), který je opět jednička. Nechť \vec{z} je vektor, pro který je matice \mathbf{G} maticí sdruženou. Potom*

definujme matici \mathbf{B} jako diagonální matici, s vektorem \vec{z} na hlavní diagonále. Nyní spočítejme skalární součiny vektorů $\vec{p}_{\tilde{p}_1}^{(n)}, \dots, \vec{p}_{\tilde{p}_n}^{(n)}$, $n \stackrel{\text{def}}{=} 2k$, s vektorem $\mathbf{B}\vec{y}$, neboli spočteme skalární součiny matic $\mathbf{M}_{\vec{p}_{\tilde{p}_i}^{(n)}}$ s maticí $\mathbf{M}_{\mathbf{B}\vec{y}}$. Zřejmě ale matice $\mathbf{M}_{\mathbf{B}\vec{y}}$ má těsně nad hlavní diagonálou 1 a těsně pod hlavní diagonálou -1 a na pozici $(1, 1)$ jedničku. Odtud plyne, že všechny uvažované skalární součiny jsou rovny -1 , protože matice sdružené k základním vektorům parity jsou po sloupcích (řádcích) konstantní, tedy skalární součiny jednotlivých řad těchto matic s odpovídajícími řadami matice $\mathbf{M}_{\mathbf{B}\vec{y}}$ jsou nulové (v každém sloupci matice $\mathbf{M}_{\mathbf{B}\vec{y}}$ je pouze jedna 1 a jedna -1) s výjimkou prvního a posledního sloupce (řádku). Skalární součin těchto dvou řad je roven -1 . Stejně tak je $\langle \mathbf{E}^{(k)} | \mathbf{M}_{\mathbf{B}\vec{y}} \rangle = 1$. Navíc je zřejmé, že $\sum_{k=1}^{2^n} |(\mathbf{B}\vec{y})_k| = 2 \cdot 2^k - 1$, tedy dostáváme z diskriminační lemmy nerovnost

$$w \geq 2 \cdot 2^k - 1 = 2^{k+1} - 1.$$

Na složky vektoru \vec{y} z příkladu se můžeme dívat jako na výsledek porovnání velikosti indexů v matici $\vec{T}^{(k)} = \mathbf{M}_{\vec{y}}$. Lze tedy vyvodit závěr, že každý booleovský obvod hloubky 1, porovnávající dvě libovolná čísla s binárním zápisem délky nepřesahující hodnotu k , musí mít váhy se součtem absolutních hodnot alespoň $2^{k+1} - 1$. Tento odhad je optimální, neboť takovýto obvod lze snadno zkonstruovat. Vidíme tedy, že výsledky získané použitím diskriminační lemmy mohou být velmi rozdílné a jsou silně závislé na matici \mathbf{B} .

Cvičení 2.4.1 Diskriminační lemma

1. Ověřte, že pro $i \in \{1, \dots, k\}$, platí

$$\langle \mathbf{T}^{(k)} | \mathbf{M}_{\vec{p}_{\tilde{p}_i}^{(2k)}} \rangle = 2^{i+1}.$$

2. Dokažte, že vektor \vec{y} definovaný v příkladě 2.4.1 není prahem systému základních vektorů parity (jinými slovy, k jeho výpočtu je nutno mít k dispozici nenulové prahy).
3. Navrhněte booleovský obvod hloubky 1, počítající vektor \vec{y} jako práh systému základních vektorů parity a konstantního vektoru, který má součet absolutních hodnot vah roven $2^{k+1} - 1$.

2.4.2 Horní odhad velikosti vah a počtu prahových vektorů

Nyní obraťme pozornost na odhady počtu prahových vektorů klasického binárního perceptronu, t.j. perceptronu, jehož výstupem je hodnota $\text{sgn}(\sum_{i=1}^S \vec{w}_i \vec{z}_i - t)$, kde \vec{z}_i jsou vstupní hodnoty a t je prahová hodnota. Tento model perceptronu počítá prahové vektory systému základních vektorů parity a konstantního vektoru, jehož prostřednictvím je realizován práh. Výchozím pojmem k odhadu počtu takovýchto prahů je tzv. jádro prahu ([Hås94]), což je ve své podstatě v korespondenci s diskriminační lemmou minimální systém nerovnic,

$$\sum_{i=1}^n \vec{w}_i (\vec{x}_i)_j \leq \vec{y}_j, \quad \text{pro } \vec{y}_j = -1,$$

$$\sum_{i=1}^n \vec{w}_i (\vec{x}_i)_j \geq \vec{y}_j, \quad \text{pro } \vec{y}_j = 1,$$

který zaručuje, že jeho řešením jsou váhy pro práh \vec{y} . Uvidíme, že pro práh \vec{y} vektorů $\vec{x}_1, \dots, \vec{x}_S$ takovýto systém vždy existuje, nerovnice v něm nabývají rovnosti a počet

těchto nerovnic je roven právě číslu S . (Jinými slovy, matice \mathbf{B} realizující přesný odhad velikosti vah v diskriminační lemmě musí obsahovat pro systém $\vec{x}_1, \dots, \vec{x}_S$ alespoň S jedniček a na druhé straně takováto matice s právě S jedničkami existuje).

Definice 2.4.1 *Nechť \vec{y} je prahem systému vektorů $\vec{x}_1, \dots, \vec{x}_S$ z $\{-1, +1\}^n$, \mathbf{X} je matice, jejíž sloupce jsou tvořeny vektory \vec{x}_i a Ω je matice o S sloupcích a 2^S řádcích, jejíž i -tý sloupec je roven i -tému základnímu vektoru parity $\vec{p}^{(S)}_i$. Nechť $\vec{\Theta} \in \mathbb{R}^S$ splňuje následující tři kritéria (zde funkce \widetilde{sgn} je brána po složkách, $|\vec{z}|$ je vektor tvořený absolutními hodnotami složek vektoru \vec{z} a $\vec{1}$ označuje vektor tvořený jedničkami):*

1. $\widetilde{sgn}(\mathbf{X}\vec{\Theta}) = \vec{y}$
2. $|\Omega\vec{\Theta}| \geq \vec{1}$
3. počet složek vektoru $|\Omega\vec{\Theta}|$ rovných jedné je maximální přes všechny vektory $\vec{\Theta}^*$, splňující předešlé dvě podmínky (je-li takovýchto $\vec{\Theta}$ více, bereme v úvahu libovolný z nich).

Každé $\vec{\Theta}$ vyhovující předchozím podmínkám nazveme JÁDRO PRAHU \vec{y} .

Existence vektoru $\vec{\Theta}$ plyne z předpokladu, že \vec{y} je prahem systému $\vec{x}_1, \dots, \vec{x}_S$. Dále připomeňme, že i -tý řádek matice Ω odpovídá binárnímu zápisu čísla i . Protože počet těchto řádků je roven 2^S , obsahují řádky matice Ω všechny možné S -členné posloupnosti 1 a -1 .

Lemma 2.4.2 *Předpokládejme, že $\vec{\Theta}$ je jádrem prahového vektoru \vec{y} systému $\vec{x}_1, \dots, \vec{x}_S$ a že matice \mathbf{A} je podmatice matice Ω , tvořená těmi řádky \vec{r}^T matice Ω , pro které platí $|\vec{r}^T\vec{\Theta}| = 1$. Potom hodnota matice \mathbf{A} je rovna číslu S .*

■ *Důkaz:*

DŮKAZ JE UVEDEN V ORIGINÁLNÍ TIŠTĚNÉ VERZI.

■

K odhadu počtu různých prahů systému základních vektorů parity a konstantního vektoru bude užitečná následující elementární lemma o determinantu matic obsahujících pouze ± 1 .

Lemma 2.4.3 *Nechť \mathbf{A} je čtvercová matice řádu m s prvky z množiny $\{-1, +1\}$. Potom $\det \mathbf{A}$ je dělitelný číslem 2^{m-1} .*

■ *Důkaz:*

DŮKAZ JE UVEDEN V ORIGINÁLNÍ TIŠTĚNÉ VERZI.

■

Tvrzení 2.4.4 *Nechť \vec{y} je prahem systému vektorů $\vec{x}_1, \dots, \vec{x}_S$ z $\{-1, +1\}^n$. Potom existují celočíselné váhy w_i pro \vec{y} tak, že pro všechna $i \in \{1, \dots, n\}$ platí odhad*

$$|w_i| \leq \frac{S^{\frac{S}{2}}}{2^{S-1}}. \quad (2.6)$$

■ *Důkaz:*

DŮKAZ JE UVEDEN V ORIGINÁLNÍ TIŠTĚNÉ VERZI.

Skutečnost, že váhy jsou dány jako řešení soustavy rovnic s koeficienty ± 1 , umožňuje horní odhad počtu prahů.

Tvrzení 2.4.5 *Počet prahových vektorů systému $\vec{x}_1, \dots, \vec{x}_S$ je zhora omezen číslem 2^{S^2+S} .*

■ *Důkaz:*

DŮKAZ JE UVEDEN V ORIGINÁLNÍ TIŠTĚNÉ VERZI.

Ještě poznamenejme, že počet různých prahů pro systém základních vektorů parity o dimenzi 2^S a konstantního vektoru je roven počtu regulárních soustav $(S+1)$ -lineárních rovnic s koeficienty ± 1 a pravou stranou tvořenou také ± 1 .

2.4.3 Dolní odhad velikosti vah a počtu prahových vektorů

Nyní se krátce budeme věnovat dolnímu odhadu velikosti vah a počtu prahových vektorů. Uvidíme úzkou souvislost těchto dvou odhadů. K důkazu dolního odhadu počtu prahů je nutné tvrzení lemmy následující.

Lemma 2.4.6 *Nechť \vec{y} je prahem systému vektorů $\vec{x}_1, \dots, \vec{x}_S$ dimenze n , takových, že matice jejíž sloupce jsou rovny vektorům $\vec{x}_1, \dots, \vec{x}_S$, má všechny řádky po dvou od sebe různé. Potom existují váhy \vec{w} takové, že*

$$(\forall i, j \in \{1, \dots, n\}) \left(i \neq j \Rightarrow \left(\sum_{k=1}^n \vec{x}_k \vec{w}_k \right)_i \neq \left(\sum_{k=1}^n \vec{x}_k \vec{w}_k \right)_j \right).$$

■ *Důkaz:*

DŮKAZ JE UVEDEN V ORIGINÁLNÍ TIŠTĚNÉ VERZI.

Na podkladě této lemmy lze provést dolní odhad počtu prahových vektorů ([YI65], [Smi65]). Tento dolní odhad provedeme indukcí podle počtu prahových vektorů, a to tak, že pro libovolný práh systému k -základních vektorů parity a konstantního vektoru (vše dimenze 2^k) zkonstruujeme $2^k + 1$ od sebe různých prahových vektorů systému základních vektorů parity a konstantního vektoru dimenze 2^{k+1} . Navíc pro různé prahové vektory dimenze 2^k budou takto zkonstruované prahy dimenze 2^{k+1} od sebe navzájem různé.

Z geometrického pohledu je podstata zmíněné konstrukce následující (berme v úvahu $k = 2$). Představme si v rovině čtverec (např. $\langle -1, 1 \rangle \times \langle -1, 1 \rangle$) a libovolnou přímku p . Dále předpokládejme situaci, kdy s přímkou p pohybujeme v jejím normálovém směru. Potom v každé poloze tato přímka rozděluje všechny vrcholy čtverce na dvě množiny. Zřejmě můžeme bez újmy na obecnosti předpokládat (viz. lemma 2.4.6), že nikdy nenastane situace, aby přímka p protínala současně dva vrcholy čtverce. Tedy pohybem takové přímky p ve směru jejího normálového vektoru dostaneme celkem $2^k + 1 = 5$ různých rozdělení

vrcholů čtverce. Nechť uvažovaný čtverec leží v pevně zvolené rovině trojrozměrného prostoru a tvoří jednu stěnu trojrozměrné krychle. Ponechme přímkou p fixovanou v nějaké poloze. Rovinou obsahující opačnou stěnu krychle (označme tuto stěnu z) rovnoběžnou s původním čtvercem vedme přímkou q , která je rovnoběžná s přímkou p a pohybujeme přímkou q ve směru jejího normálového vektoru, který je rovnoběžný se stěnou z .

Zřejmě přímkou q rozdělujeme vrcholy ležící na stěně z stejným způsobem jako přímkou p v původním čtverci před svojí fixací. Evidentně rovina, obsahující obě přímkou p a q rozdělují vrcholy krychle do dvou množin, a to právě $2^k + 1$ způsoby podle aktuální polohy přímkou q . Zřejmě jsme pro pevně zvolenou přímkou p dostali $2^k + 1$ různých rozdělení vrcholů krychle do množin \bar{A} , \bar{B} , $\bar{A} \cap \bar{B} = \emptyset$. Přitom je evidentní, že \bar{A} a \bar{B} jsou lineárně separovatelné právě rovinou obsahující přímkou p a q .

Na druhé straně použijeme-li přímkou p' takovou, že množiny vrcholů původního čtverce vymezené přímkou p jsou rozdílné od množin vrcholů původního čtverce vymezené přímkou p' , dostáváme, že pro žádné dvě takto vzniklé množiny \bar{A}' a \bar{B}' nemůže platit $\bar{A} = \bar{A}'$, neboť se musí lišit ve vrcholech krychle ležících v původním čtverci. Tato konstrukce je obsahem důkazu následujícího tvrzení a je ilustrována na obrázku 2.4.

Tvrzení 2.4.7 *Počet prahových vektorů systému $\vec{p}_1^{(S)}, \dots, \vec{p}_S^{(S)}$, základních vektorů parity a jednotkového vektoru $\vec{e}^{(S)}$ je zdola omezen číslem $2^{\frac{S(S-1)}{2}}$.*

■ *Důkaz:*

DŮKAZ JE UVEDEN V ORIGINÁLNÍ TIŠTĚNÉ VERZI.

■

Tohoto tvrzení lze přímo využít pro následující tvrzení, které říká, že existuje prahový vektor, jehož celková váha je exponenciální v S , což znamená, že k zápisu jeho vah je potřeba lineárního počtu bitů.

Tvrzení 2.4.8 *Existuje práh systému $\vec{p}_1^{(S)}, \dots, \vec{p}_S^{(S)}$, základních vektorů parity a jednotkového vektoru $\vec{e}^{(S)}$, pro jehož všechny celočíselné váhové vektory \vec{w} platí odhad*

$$2^{\frac{S-2}{2}} \leq \sum_{k=1}^S |\vec{w}_k|.$$

Navíc \vec{w} lze volit tak, že platí

$$\sum_{k=1}^S |\vec{w}_k| \leq \frac{S^{\frac{S}{2}}}{2^S}.$$

■ *Důkaz:*

DŮKAZ JE UVEDEN V ORIGINÁLNÍ TIŠTĚNÉ VERZI.

■

2.5 Otázka existence prahového vektoru

V této podkapitole budeme analyzovat problém, zda pro daný vektor \vec{y} a pro daný systém vektorů $\vec{x}_1, \dots, \vec{x}_S$ lze efektivně rozhodnout otázku, zda \vec{y} je prahem zmíněného systému vektorů.

Pro tuto analýzu budeme potřebovat tvrzení Farkašovy věty, které zde nebudeme dokazovat a které vyslovíme v následujícím tvaru (důkaz lze nalézt v každé učebnici lineárního programování, pro studenty MI FJFI ČVUT doporučuji přednášku J. Pytlíček–Lineární programování):

Lemma 2.5.1 (Farkašovo) *Z následujících dvou soustav lineárních nerovnic*

1. $A\vec{w} \geq \vec{b}$, $\vec{b} \neq \vec{0}$
2. $\vec{z}^T A < 0$, $\vec{z}^T \vec{b} > 0$,

má právě jedna nezáporné řešení.

Nyní vyslovíme nutnou a postačující podmínku pro to, aby vektor \vec{y} nebyl prahovým vektorem systému vektorů $\vec{x}_1, \dots, \vec{x}_S$ ([Mur71]).

Tvrzení 2.5.2 *Vektor \vec{y} není prahem systému vektorů $\vec{x}_1, \dots, \vec{x}_S$ právě když soustava nerovnic*

$$\vec{z}^T \text{diag}(\vec{y}) \mathbf{X} < \vec{0}, \quad \vec{z}^T \vec{1} > 0, \quad (2.7)$$

(kde $\text{diag}(\vec{y})$ je diagonální matice s vektorem \vec{y} na hlavní diagonále) má nezáporné řešení.

■ *Důkaz:*

DŮKAZ JE UVEDEN V ORIGINÁLNÍ TIŠTĚNÉ VERZI.

■

Ekvivalence dokázaná v předchozím tvrzení má zásadní důsledek pro aplikaci booleovských obvodů, bohužel velmi negativní, protože následující problém

INSTANCE:LINEÁRNÍ PROGRAMOVÁNÍ

Celočíselná matice A , celočíselné vektory \vec{c} , \vec{b} , odpovídajících dimenzí.

PROBLÉM:

Existuje racionální vektor \vec{z} tak, že platí $\vec{z}^T A \leq \vec{c}$, $\vec{z}^T \vec{b} \geq 0$?

leží ve třídě NP–úplných problémů (viz [GJ79], str. 287–288). Tuto skutečnost lze na základě tvrzení 2.5.2 přeformulovat následovně:

Tvrzení 2.5.3 *Následující problém*

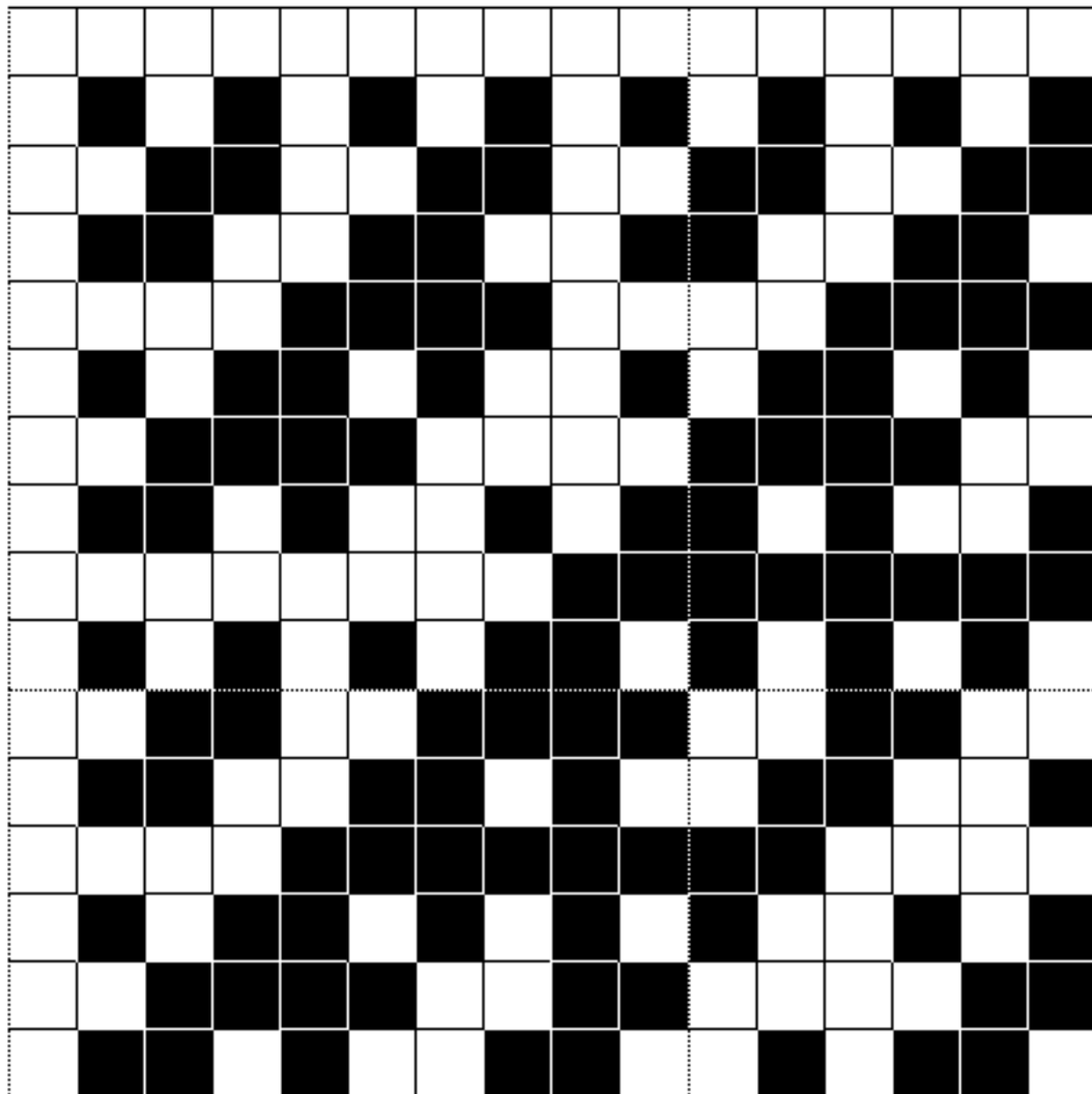
INSTANCE:PRÁH SYSTÉMU VEKTORŮ

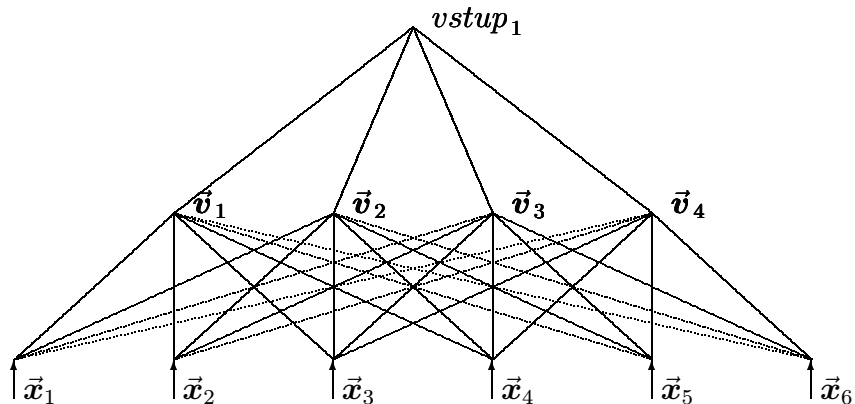
$\vec{x}_1, \dots, \vec{x}_S, \vec{y} \in \{-1, +1\}^n$.

PROBLÉM:

Je vektor \vec{y} prahem systému $\vec{x}_1, \dots, \vec{x}_S$?

leží ve třídě NP–úplných problémů.





Obrázek 2.2: Booleovský obvod pro výpočet α -symetrických vektorů (pro $n = 5$ a $m = 3$). Uzly \vec{x}_6 a \vec{v}_4 slouží pouze k přenosu konstantního vstupu (prahové váhy). Váhy hran mezi vstupní a prostřední vrstvou jsou popsány vzorci 2.4, váhy mezi uzly prostřední vrstvy a výstupního uzlu jsou popsány vzorci 2.5.

Obrázek 2.4: Konstrukce $2^k + 1$ prahových vektorů odvozených z daného prahového vektoru pro $k = 2$.

Kapitola 3

Aproximační možnosti neuronových sítí

Finmanovo pravidlo

NIKOMU SE NECHCE PRACOVAT SE VZORCI,
KTERÉ VYMYSLEL NĚKDO JINÝ.

Pod aproximačními vlastnostmi umělých neuronových sítí budeme rozumět především schopnost neuronových sítí odpovídat na zadané vstupní hodnoty takovou výstupní hodnotou, která bude blízko hodnotě předem dané funkce v argumentu, který obsahuje vstupní hodnoty sítě.

Ukážeme, že je-li \bar{K} kompaktní množina v \mathbb{R}^n , pak neuronová síť s jednou skrytou vrstvou může uniformě aproximovat libovolnou spojitou funkci definovanou na \bar{K} . Tento problém je zajímavý zejména z toho důvodu, že nejčastější způsob aplikování neuronových sítí, t.j. oddělení (separace) více množin vzájemně od sebe, lze formulovat jako problém aproximace charakteristických funkcí těchto množin.

Jak bylo popsáno v úvodní kapitole, jediný perceptron dokáže rozdělit vstupní prostor \mathbb{R}^n do dvou poloprostorů. Zřejmě tedy síť sestávající z více perceptronů, které jsou paralelně vedle sebe a jejichž výstupy jsou vstupními hodnotami jediného výstupního neuronu (tzv. třívrstvá neuronová síť), dokáže z prostoru \mathbb{R}^n vydělit průnik m poloprostorů, kde m je počet neuronů v prostřední vrstvě.

Budeme-li mít disjunktní konečné množiny $\bar{A}, \bar{B} \in \mathbb{R}^n$, pak zřejmě existují vzájemně po dvou disjunktní konvexní množiny $\bar{a}_1, \dots, \bar{a}_l$, $\bar{b}_1, \dots, \bar{b}_k$, z nichž každá je tvořena průnikem konečného počtu poloprostorů (je tedy polyedrem) a pro které platí, že $\bar{A} \subset \bigcup_1^l \bar{a}_i$ a $\bar{B} \subset \bigcup_1^k \bar{b}_j$.

Předpokládejme, že máme vícevrstvou neuronovou síť, tvořenou perceptrony, a že počet neuronů ve druhé vrstvě je roven počtu poloprostorů, nezbytných k vytvoření systému množin $\bar{a}_1, \dots, \bar{a}_l$ a $\bar{b}_1, \dots, \bar{b}_k$, a že každý z neuronů druhé vrstvy nabývá hodnoty ± 1 podle toho, na jaké straně daného poloprostoru je aktuální vstupní vektor. Chceme-li vytvořit síť, která pro vektory z $\bigcup_i^l \bar{a}_i$ dává hodnoty odlišné od vektorů z $\bigcup_j^k \bar{b}_j$, můžeme se na tento problém dívat jako na problém vytvoření konkrétní booleovské funkce, jejíž proměnné jsou hodnoty neuronů druhé vrstvy. Zřejmě funkce NOT, AND a OR jsou všechny lineárně separabilní a tudíž vyjádřitelné pomocí perceptronu. Současně ale trojice

funkcí $\overline{\text{NOT}}$, $\widehat{\text{AND}}$ a $\overline{\text{OR}}$ tvoří tzv. úplný systém výrokové logiky, což znamená, že jejich skládáním a uzávorkováním lze vyjádřit jakoukoli booleovskou funkci libovolného počtu proměnných. Z toho ale vyplývá, že vícevrstvá síť perceptronů, separující od sebe dvě dané konečné disjunktční množiny \bar{A} , \bar{B} (obecněji dva disjunktční systémy disjunktčních konvexních polyedrů) vždy existuje.

Otázkou zůstává horní a dolní odhad počtu neuronů ve vrstvách a počet těchto vrstev. V obecném případě aproximace funkce stejnoměrně spojitě na kompaktu \bar{K} můžeme vytvořit konečné ϵ -disjunktční pokrytí \bar{D} tohoto kompaktu a aproximovanou funkci nahradit funkcí $\tilde{\psi}$, po částech konstantní na tomto pokrytí. Zřejmě funkci $\tilde{\psi}$ můžeme přesně aproximovat neuronovou sítí (disjunktční pokrytí \bar{D} lze volit jako systém polyedrů) a ze stejnoměrné spojitosti plyne, že pro $\epsilon \rightarrow 0$ lze stejnoměrně spojitou funkci aproximovat libovolně přesně sítí sestávající se pouze z perceptronů. Otázkou však zůstává odhad na velikost takovéto sítě.

Obecně lze problém analýzy aproximačních vlastností neuronových sítí se třemi vrstvami formulovat jako problém vyjádření dané funkce $\tilde{f} : \mathfrak{R}^n \rightarrow \mathfrak{R}$ ve tvaru součtu

$$\tilde{f}(\vec{x}) \stackrel{\text{def}}{=} \tilde{g} \left(\sum_{j=1}^l \tilde{\sigma}(\vec{w}_j^T \vec{x} + t_j) \right),$$

kde \tilde{g} a $\tilde{\sigma}$ jsou funkce jedné reálné proměnné.

V následujících podkapitolách se tímto problémem budeme zabývat a ukážeme čtyři přístupy k analýze aproximačních vlastností neuronových sítí. Tyto přístupy jsou primárně založeny na Stone-Weierstrassově větě, vlastnostech duálních prostorů, využití Kolmogorovy věty a konečně na aplikaci konvoluce funkcí.

3.1 Důsledky Stone-Weierstrassovy věty

Nyní ukážeme přímý přístup k analýze aproximačních vlastností ([Ell94]). Základní přístup je založen na následujících principech. Nejdříve aproximujeme spojitou funkci jedné reálné proměnné funkcí po částech konstantní, a to s libovolnou přesností. Následně se pokusíme zobecnit tento typ aproximace na funkce více proměnných. Následující pojem modulu spojitosti bude užitečný v aproximačních odhadech.

Definice 3.1.1 *Nechť $\tilde{f} : \bar{K} \rightarrow \mathfrak{R}$, kde $\bar{K} \subset \mathfrak{R}^n$ je kompaktní množina. Potom definujeme MODUL SPOJITOSTI (anglický termín: modulus of continuity) jako číslo*

$$\omega(\tilde{f}, \delta) \stackrel{\text{def}}{=} \sup_{\substack{\vec{x}, \vec{y} \in \bar{K} \\ \|\vec{x} - \vec{y}\| < \delta}} |\tilde{f}(\vec{x}) - \tilde{f}(\vec{y})|.$$

Je-li funkce \tilde{f} spojitá, pak hodnota $\omega(\tilde{f}, \delta)$ je konečná a pro $\delta \rightarrow 0$ se blíží k nule také. Modul spojitosti v sobě zahrnuje různé způsoby vyjádření hladkosti funkce \tilde{f} . Například je-li \tilde{f} Lipstichovská, tedy $|\tilde{f}(\vec{x}) - \tilde{f}(\vec{y})| \leq L \cdot \|\vec{x} - \vec{y}\|$, $L > 0$, $\vec{x}, \vec{y} \in \bar{K}$, tak zřejmě platí $\omega(\tilde{f}, \delta) \leq L\delta$. Modul spojitosti dává základní odhad, jak přesně lze danou funkci \tilde{f} aproximovat funkcí, která je po částech konstantní.

Následující tvrzení je založené na zřejmém faktu, že pro každou sigmoidální funkci $\tilde{\sigma}$ (viz 1.2.1) platí, že pro $a \rightarrow +\infty$ funkce $\tilde{\sigma}(ax)$ konverguje bodově k 1 pro $x > 0$ a k -1

pro $x < 0$. Tedy $\tilde{\sigma}(ax)$ pro a rostoucí nade všechny meze konverguje bodově k prahové funkci \widetilde{sgn} až na vyjimku v bodě 0. Z toho plyne, že funkce $\frac{1}{2}(\tilde{\sigma}(ax) - \tilde{\sigma}(a(x-1)))$ konverguje k charakteristické funkci intervalu $(0, 1)$.

Tvrzení 3.1.1 *Nechť $\tilde{\sigma}$ je sigmoidální funkce. Pak existuje konstanta c tak, že pro libovolné $\tilde{f} \in C_{(0,1)}$ platí nerovnost*

$$\|\tilde{f} - \tilde{q}_n\|_{\infty} \leq c \cdot \omega\left(\tilde{f}, \frac{1}{n}\right),$$

kde posloupnost funkcí $\{\tilde{q}_n\}_1^{\infty}$ je definována předpisem

$$\tilde{q}_n(x) \stackrel{\text{def}}{=} \tilde{f}(0) + \sum_{v=1}^n \left(\tilde{f}\left(\frac{v}{n}\right) - \tilde{f}\left(\frac{v-1}{n}\right) \right) \tilde{\sigma}(A_n(nx - v))$$

a pro číselnou posloupnost $\{A_i\}_1^{\infty}$ platí, že A_n je nejmenší kladné celé číslo, pro které platí

$$|\tilde{\sigma}(x)| \leq \frac{1}{n} \quad \text{pro } x \leq -A_n \quad a \quad 1 - \frac{1}{n} \leq \tilde{\sigma}(x) \leq 1 + \frac{1}{n} \quad \text{pro } x \geq A_n.$$

■ *Důkaz:*

DŮKAZ JE UVEDEN V ORIGINÁLNÍ TIŠTĚNÉ VERZI.

■

Následující větu, jejíž důkaz zde neuvеdeme, (viz. např. [Miš89], str. 220), budeme potřebovat pro důkaz úplnosti všech sigmoidálních funkcí v prostoru všech funkcí spojitých na kompaktní množině.

Tvrzení 3.1.2 (Stone-Weierstrass) *Nechť \mathcal{B} je lineární podprostor v prostoru $C_{\mathbb{R}^n}$ takový, že obsahuje konstantní nenulovou funkci a v prostoru \mathbb{R}^n odděluje libovolné dva body (t.j. pro libovolné $\vec{x}, \vec{y} \in \mathbb{R}^n$, $\vec{x} \neq \vec{y}$ existuje $\tilde{f} \in \mathcal{B}$ tak, že $\tilde{f}(\vec{x}) \neq \tilde{f}(\vec{y})$). Potom \mathcal{B} je hustá v prostoru $C_{\mathbb{R}^n}$.*

Před vlastním důkazem hustoty sigmoidálních funkcí v prostoru funkcí spojitých na kompaktní množině \bar{K} dokážeme tvrzení pomocné, ukazující hustotu exponenciálních funkcí v prostoru $C_{\bar{K}}$.

Tvrzení 3.1.3 *Předpokládejme, že \bar{K} je kompaktní podmnožina prostoru \mathbb{R}^n . Potom lineární obal množiny funkcí $\bar{\Theta} \stackrel{\text{def}}{=} \left\{ \tilde{\rho}(\vec{x}) = e^{\langle \vec{a} | \vec{x} \rangle} \mid \vec{a} \in \mathbb{R}^n \right\}$ je hustý v prostoru $C_{\bar{K}}$.*

■ *Důkaz:*

DŮKAZ JE UVEDEN V ORIGINÁLNÍ TIŠTĚNÉ VERZI.

■

Následující tvrzení ve své podstatě říká, že neuronová síť typu perceptronu s libovolnou sigmoidální funkcí $\tilde{\sigma}$ může libovolně přesně aproximovat jakoukoli spojitou funkci na kompaktní množině.

Tvrzení 3.1.4 *Nechť \bar{K} je kompaktní podmnožina prostoru \mathbb{R}^n . Potom množina funkcí tvaru*

$$\tilde{g}(\vec{x}) \stackrel{\text{def}}{=} \sum_{j=1}^k a_j \tilde{\sigma}(\langle \vec{w} | \vec{x} \rangle + c_j),$$

kde $\tilde{\sigma}$ je spojitá sigmoidální funkce, je hustá v prostoru $C_{\bar{K}}$.

■ *Důkaz:*

DŮKAZ JE UVEDEN V ORIGINALNÍ TIŠTĚNÉ VERZI.

■

3.2 Aplikace duálních prostorů

Jeden z dalších způsobů důkazu hustoty daného systému funkcí v množině funkcí spojitých na kompaktní množině je založen na funkcionálních vlastnostech duálních prostorů. Připomeňme stručně, že duální prostor normovaného vektorového prostoru \mathcal{X} je lineární prostor, tvořený všemi omezenými lineárními funkcionály definovanými na prostoru \mathcal{X} . Každý duální prostor s normou definovanou jako

$$\|\tilde{f}\| \stackrel{\text{def}}{=} \sup_{\substack{\vec{x} \in \mathcal{X} \\ \|\vec{x}\|=1}} |\tilde{f}(\vec{x})|,$$

je Banachův prostor, což znamená, že každá Cauchyovská posloupnost (vzhledem ke zmíněné normě) funkcionálů duálního prostoru má v tomto prostoru limitní prvek. Relevance duálních prostorů k hustotě dané množiny funkcí v prostoru všech spojitých funkcí je předmětem následující věty.

Tvrzení 3.2.1 *Nechť \mathcal{X} je normovaný vektorový prostor nad tělesem \mathbb{R} , \bar{V} je jeho podprostor. Potom \bar{V} je hustý v \mathcal{X} , právě když jediný lineární funkcionál na \mathcal{X} , v jehož nulovém prostoru je \bar{V} obsažen, je nulový funkcionál.*

■ *Důkaz:*

DŮKAZ JE UVEDEN V ORIGINALNÍ TIŠTĚNÉ VERZI.

■

Chceme-li tedy pro daný podprostor \bar{V} prostoru \mathcal{X} dokázat hustotu tohoto podprostoru v celém prostoru, stačí ukázat, že jediný funkcionál, který anihiluje celý podprostor \bar{V} je nulový funkcionál.

Jako příklad aplikace duálních prostorů uvedeme následující analýzu v prostoru $\mathcal{L}_p^{\bar{K}}$ a $C_{\bar{K}}$.

3.2.1 Hustota sigmoidálních funkcí v prostoru $\mathcal{L}_p^{\bar{K}}$

Nyní využijeme předchozí tvrzení pro analýzu hustoty množiny všech sigmoidálních funkcí v prostoru funkcí, jejichž p -tá mocnina, $1 < p$, je Lesbagueovsky integrabilní na kompaktu $\bar{K} \subset \mathbb{R}^n$.

Je všeobecně známým poznatkem funkcionální analýzy, že každý lineární funkcionál na prostoru $\mathcal{L}_p^{\bar{K}}$ lze vyjádřit jednoznačně ve tvaru integrálu

$$\int_{\bar{K}} \tilde{f}(\vec{x}) \tilde{g}(\vec{x}) d(\vec{x}), \quad (3.1)$$

přičemž \tilde{g} je prvek z prostoru $\mathcal{L}_q^{\bar{K}}$, kde čísla p a q jsou spolu svázána vztahem $\frac{1}{p} + \frac{1}{q} = 1$. Abychom ukázali, že množina funkcí tvaru

$$\tilde{\sigma}(\langle \vec{w} | \vec{x} \rangle + c) \quad (3.2)$$

je hustá v $\mathcal{L}_p^{\bar{K}}$, dokážeme, že jediný anihilátor těchto funkcí je nulový funkcionál. Máme tedy dokázat následující implikaci:

$$(\forall \tilde{\sigma}(\langle \vec{x} | \vec{w} \rangle + c)) \left(\int_{\bar{K}} \tilde{\sigma}(\langle \vec{x} | \vec{w} \rangle + c) \tilde{g}(\vec{x}) d(\vec{x}) = 0 \right) \Rightarrow \tilde{g}(\vec{x}) = 0 \quad \text{skoro všude na } \bar{K}.$$

Důkaz provedeme sporem, předpokládejme tedy, že \tilde{g} je anihilátor všech funkcí tvaru 3.2 a $\|\tilde{g}\| > 0$. Nejdříve z důvodů jednoduššího zápisu dodefinujeme funkci \tilde{g} nulou vně kompaktu \bar{K} . Protože funkcionál 3.1 anihiluje všechny funkce tvaru 3.2, musí anihilovat i funkce

$$\tilde{\sigma}_n(\vec{x}) \stackrel{\text{def}}{=} \tilde{\sigma}(n(\langle \vec{a} | \vec{x} \rangle - \langle \vec{a} | \vec{y} \rangle)),$$

kde \vec{y} je libovolný vektor z \bar{K} a \vec{a} je libovolný jednotkový vektor, $\|\vec{a}\| = 1$. Zřejmě nadrovina $\langle \vec{a} | \vec{x} \rangle - \langle \vec{a} | \vec{y} \rangle = 0$ protíná množinu \bar{K} v bodě \vec{y} . Protože ale

$$\int_{\bar{K}} \tilde{\sigma}_n(\vec{x}) \tilde{g}(\vec{x}) d(\vec{x}) = 0, \quad (3.3)$$

pro libovolné n dostáváme, že $(\tilde{\sigma}_n)$ konverguje bodově k -1 na množině $\langle \vec{a} | \vec{x} \rangle < \langle \vec{a} | \vec{y} \rangle$ a k $+1$ na množině $\langle \vec{a} | \vec{x} \rangle > \langle \vec{a} | \vec{y} \rangle$ platí rovnost

$$\int_{\langle \vec{a} | \vec{x} \rangle < \langle \vec{a} | \vec{y} \rangle} \tilde{g}(\vec{x}) d(\vec{x}) = \int_{\langle \vec{a} | \vec{x} \rangle > \langle \vec{a} | \vec{y} \rangle} \tilde{g}(\vec{x}) d(\vec{x}),$$

protože jinak bychom na základě spojitosti Lebesgueova integrálu dostali spor s rovnicemi 3.3. Funkce s takovými vlastnostmi popisuje následující lemma:

Lemma 3.2.2 *Nechť \tilde{g} je funkce definovaná na kompaktní množině $\bar{K} \subset \mathbb{R}^n$ splňující následující podmínku:*

$$(\forall \vec{a} \in \mathbb{R}^n, \|\vec{a}\| = 1) (\forall \vec{y} \in \bar{K}) \left(\int_{\langle \vec{a} | \vec{x} \rangle < \langle \vec{a} | \vec{y} \rangle} \tilde{g}(\vec{x}) d(\vec{x}) = \int_{\langle \vec{a} | \vec{x} \rangle > \langle \vec{a} | \vec{y} \rangle} \tilde{g}(\vec{x}) d(\vec{x}). \right)$$

Potom \tilde{g} je skoro všude na \bar{K} rovna nule.

■ *Důkaz:*

DŮKAZ JE UVEDEN V ORIGINÁLNÍ TIŠTĚNÉ VERZI. ■

Tedy jediný anihilátor sigmoidálních funkcí v prostoru $\mathcal{L}_p^{\bar{K}}$ je nulový funkcionál a proto lineární obal množiny sigmoidálních funkcí 3.2 je hustý v $\mathcal{L}_p^{\bar{K}}$.

3.2.2 Hustota v RBF funkcí v prostoru $C_{\bar{K}}$

Nyní budeme zkoumat hustotu sigmoidální funkce v prostoru $C_{\bar{K}}$, na kterém je definovaná supremová norma $\|\tilde{f} - \tilde{g}\|_{\infty} \stackrel{\text{def}}{=} \sup_{\vec{x} \in \bar{K}} |\tilde{f}(\vec{x}) - \tilde{g}(\vec{x})|$.

Hustota v tomto funkcionálním prostoru umožňuje aproximovat spojité funkce na kompaktní množině lokálním způsobem, který zaručuje bodovou konvergenci, narozdíl od prostoru \mathcal{L}_p^K , kde je zaručena pouze konvergence podle středu. Stejně jako v předchozím případě (pro důkaz hustoty sigmoidálních funkcí v $C_{\bar{K}}$) potřebujeme znát obecný tvar lineárního funkcionálu na prostoru $C_{\bar{K}}$. Za tímto účelem definujeme množinu $\Phi_{\bar{K}}$ jako množinu všech reálných regulárně spočetně aditivních funkcí $\tilde{\phi}$, zadaných na σ -algebře \mathcal{B} , všech Borelovských množin na \bar{K} , majících konečnou totální variaci. Nyní vyslovíme Riesz-Markov-Katukaniho větu ([KA77]), popisující obecný tvar lineárního funkcionálu v $C_{\bar{K}}$.

Tvrzení 3.2.3 *Obecný tvar lineárního spojitého funkcionálu \tilde{g} na prostoru $C_{\bar{K}}$ lze zapsat ve tvaru*

$$\tilde{g}(\tilde{f}) = \int_{\bar{K}} \tilde{f}(\vec{x}) d(\tilde{\phi}(\vec{x})), \quad \tilde{f} \in C_{\bar{K}},$$

kde $\tilde{\phi}$ je libovolný prvek z množiny $\Phi_{\bar{K}}$.

Jedněmi z dosti rozšířených přechodových funkcí používaných v různých modelech neuronových sítí jsou tzv. RBF funkce (anglický termín: Radial Basis Function). Tyto funkce umožňují lokální aproximaci spojitých funkcí. Jejich obecný tvar lze zavést například následující definicí.

Definice 3.2.1 *Nechť \tilde{f} je libovolná spojitá funkce taková, že $\tilde{f} : \mathbb{R} \rightarrow \langle 0, 1 \rangle$, $\lim_{\alpha \rightarrow \pm\infty} \tilde{f}(\alpha) = 0$ a $\lim_{\alpha \rightarrow 0} \tilde{f}(\alpha) = 1$. Potom každou funkci $\tilde{g} : \mathbb{R}^n \rightarrow \mathbb{R}$ tvaru*

$$\tilde{g}(\vec{x}) \stackrel{\text{def}}{=} \tilde{f}(\beta \|\vec{x} - \vec{c}\|),$$

kde \vec{c} je libovolný vektor a $\beta \in \mathbb{R}$, nazveme RBF FUNKCÍ odvozenou od funkce \tilde{f} .

Nechť $\tilde{\phi}$ je nenulová funkce z množiny $\Phi_{\bar{K}}$. Jelikož $\tilde{\phi}$ má nenulovou konečnou variaci, existuje bod \vec{x} takový, že, bez újmy na obecnosti, funkce $\tilde{\phi}$ je kladná na nějakém kruhovém okolí bodu \vec{x} . Pak ale existuje pro dané \tilde{f} z definice RBF vektor \vec{c} a číslo β tak, že

$$\int_{\mathbb{R}^n} \tilde{f}(\beta \|\vec{x} - \vec{c}\|) \tilde{\phi}(\vec{x}) d(\vec{x}) > 0,$$

(vně daného okolí $\tilde{f}(\beta \|\vec{x} - \vec{c}\|)$ konverguje k nule, na okolí konverguje k 1 a $\tilde{\phi}$ má konečnou variaci, tedy i konečný integrál) což ale znamená, že $\tilde{\phi}$ není anihilátorem RBF funkcí odvozených od funkce \tilde{f} . Proto je lineární obal RBF funkcí odvozených od funkce \tilde{f} hustý v prostoru $C_{\bar{K}}$.

3.2.3 Hustota sigmoidálních funkcí v prostoru $C_{\langle a,b \rangle}$

Pro prostor $C_{\langle a,b \rangle}$ platí Riesz-Markov-Katukaniho věta 3.2.3, obecný tvar funkcionálu je tedy opět zadán integrálem. Jedním ze způsobů důkazu hustoty sigmoidálních funkcí v prostoru $C_{\langle a,b \rangle}$ je adaptace důkazu hustoty těchto funkcí v prostoru \mathcal{L}_p^K . Zde však ukážeme přímý přístup, který je ve své podstatě totožný s důkazem pro prostor \mathcal{L}_p^K , ale

díky jednorozměrnosti daného definičního oboru nejsme nuceni použít vlastností Fourierovy transformace.

Mějme funkci \tilde{g} z množiny $\Phi_{\langle a,b \rangle}^-$, která je anihilátorem všech sigmoidálních funkcí. Zvolme libovolná přirozená čísla p, q a n tak, aby platilo, že $\frac{p}{q} \in \langle a, b \rangle$ (bez újmy na obecnosti předpokládejme, že $0 < a < b$). Protože \tilde{g} je anihilátorem sigmoidálních funkcí, platí rovnost

$$\int_a^b \tilde{\sigma} \left(n \left(x - \frac{p}{q} \right) \right) d(\tilde{g}(x)) = 0.$$

Protože ale funkce $\tilde{\sigma} \left(n \left(x - \frac{p}{q} \right) \right)$ konverguje pro n rostoucí nade všechny meze k funkci $\widetilde{sgn} \left(x - \frac{p}{q} \right)$, platí, že

$$\begin{aligned} 0 &= \lim_{n \rightarrow +\infty} \left(\int_a^b \tilde{\sigma} \left(n \left(x - \frac{p}{q} \right) \right) d(\tilde{g}(x)) \right) = \int_a^b \widetilde{sgn} \left(x - \frac{p}{q} \right) d(\tilde{g}(x)) = \\ &= - \int_a^{\frac{p}{q}} d(\tilde{g}(x)) + \int_{\frac{p}{q}}^b d(\tilde{g}(x)). \end{aligned}$$

Opět tedy dostáváme, že pro libovolný bod c intervalu $\langle a, b \rangle$ platí rovnost

$$\int_a^c d(\tilde{g}(x)) = \int_c^b d(\tilde{g}(x)).$$

Odtud pro libovolná $d, e \in \langle a, b \rangle$ platí

$$\int_d^e d(\tilde{g}(x)) = 0,$$

a proto funkce \tilde{g} je identicky rovna nule skoro všude na celém intervalu $\langle a, b \rangle$ (stále používáme teorii Lebesgueova integrálu). Z toho ale plyne, že anihilátor sigmoidálních funkcí v prostoru $C_{\langle a,b \rangle}$ je pouze nulový funkcionál, a tedy lineární obal sigmoidálních funkcí je hustý v $C_{\langle a,b \rangle}$.

Cvičení 3.2.3 Aplikace duálních prostorů

1. **Definice 3.2.2** Necht' \tilde{f} je libovolná spojitá funkce taková, že $\tilde{f} : \mathfrak{R} \rightarrow \langle 0, 1 \rangle$, $\lim_{\alpha \rightarrow \pm\infty} \tilde{f}(\alpha) = 0$ a $\lim_{\alpha \rightarrow 0} \tilde{f}(\alpha) = 1$. Potom každou funkci $\tilde{g} : \mathfrak{R}^n \rightarrow \mathfrak{R}$ tvaru

$$\tilde{g}(\vec{x}) \stackrel{\text{def}}{=} \tilde{f}(\beta \langle \vec{x} | \vec{c} \rangle),$$

kde \vec{c} je libovolný vektor a $\beta \in \mathfrak{R}$, nazveme SEMILOKÁLNÍ FUNKCÍ odvozenou od funkce \tilde{f} .

Rozhodněte, zda lineární obal semilokálních funkcí odvozených od nějaké funkce \tilde{f} je hustý v prostorech $C_{\bar{K}}$, $C_{\mathfrak{R}^n}$, $\mathcal{L}_p^{\bar{K}}$ a $\mathcal{L}_p^{\mathfrak{R}^n}$.

2. Rozhodněte, zda lineární obal RBF funkcí odvozených od funkce \tilde{f} je hustý v $\mathcal{L}_p^{\bar{K}}$ a $\mathcal{L}_p^{\mathfrak{R}^n}$.

3.3 Analýza na základě Kolmogorovy věty

V této sekci ukážeme přímý přístup k analýze aproximačních možností třívrstvé perceptronové sítě (vstupní a výstupní vrstva plus tzv. skrytá vrstva mezi nimi) s jedním výstupním neuronem, zobrazující \mathfrak{R}^n do \mathfrak{R} . K tomu využijeme tzv. Kolmogorovu větu (1957), která vyjadřuje každou reálnou funkci n reálných proměnných, pomocí superpozice konečného počtu monotónních rostoucích funkcí jedné reálné proměnné. Zmíněnou Kolmogorovu větu uvedeme v následující podobě ([Spr93]):

Tvrzení 3.3.1 *Pro každé přirozené $n \geq 2$ existuje $n \times (2n + 1)$ spojitých funkcí $\tilde{\psi}_{pq}$ takových, že pro libovolnou reálnou funkci $\tilde{f} : \mathfrak{R}^n \rightarrow \mathfrak{R}$ existuje $2n$ spojitých funkcí $\tilde{\Phi}_q$ tak, že pro libovolné $\vec{x} \in \mathfrak{R}^n$ platí rovnice*

$$\tilde{f}(\vec{x}) = \sum_{q=0}^{2n} \tilde{\Phi}_q \left(\sum_{p=1}^n \tilde{\psi}_{pq}(\vec{x}_p) \right).$$

Důkaz této věty neuvеdeme pro jeho příliš velkou technickou náročnost. Původní důkaz navíc není relevantní k tematice, které je tato kapitola věnovaná. Kolmogorova věta v podobě právě zmíněné nemá přímý vztah k aproximačním vlastnostem vrstevnatých architektur neuronových sítí.

3.3.1 Konstrukce aproximačního jádra $\tilde{\psi}$

V původním znění Kolmogorovy věty je spojitá funkce vyjádřena jako součet funkcí, přičemž jejich argumenty jsou hodnoty dané součtem funkcí $\tilde{\psi}_{pq}$. Funkce $\tilde{\psi}_{pq}$ mohou být obecně vzájemně různé, což by v korespondenci s funkcemi implementovanými prostřednictvím neuronových sítí znamenalo, že každý neuron má jinou přechodovou funkci $\tilde{\psi}_{pq}$. Ukážeme ale, že lze zkonstruovat funkci $\tilde{\phi}$, která umožní vyjádřit spojitou funkci ve tvaru, odpovídajícímu zápisu funkcí, realizovaných třívrstvou neuronovou sítí s předchodovou funkcí $\tilde{\psi}$.

Definice 3.3.1 *Definujme následující intervaly*

$$\bar{D}_0 \stackrel{\text{def}}{=} \left\langle 0, \frac{1}{5!} \right\rangle \quad \text{a} \quad \bar{D}_1 \stackrel{\text{def}}{=} \left\langle 0, 1 + \frac{1}{5!} \right\rangle.$$

Nejdříve ukážeme konstrukci pomocné funkce $\tilde{\psi}^* : \bar{D}_0 \rightarrow \mathfrak{R}$, indukci podle k , kde za počáteční hodnotu k vezmeme hodnotu 5. Předpokládejme, že máme danu posloupnost racionálních čísel $\{\rho_k\}_5^\infty$, která vyhovuje podmínkám

$$\rho_5 < \frac{1}{2 \cdot 5!} \quad \text{a} \quad 0 < \rho_{k+1} < \frac{\rho_k}{k} \quad \text{pro} \quad k \geq 5. \quad (3.4)$$

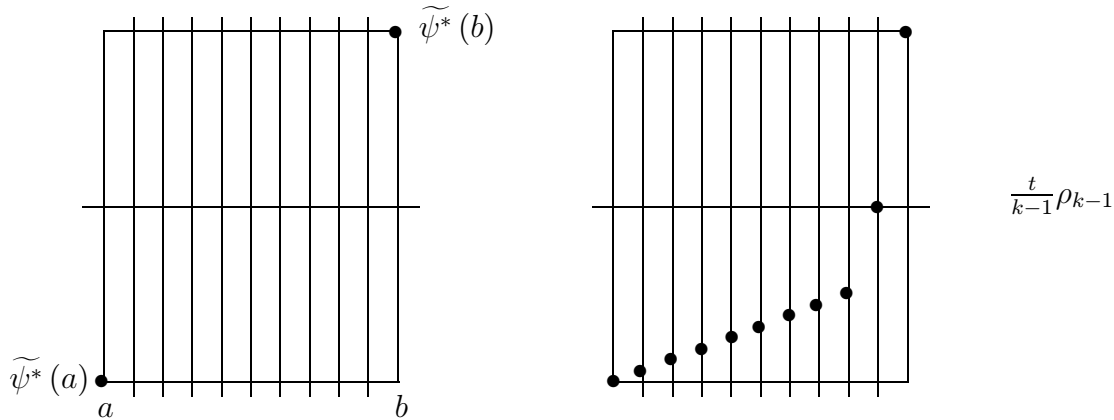
Nejdříve definujme hodnoty funkce $\tilde{\psi}^*$ v bodech $\frac{j}{6!}$, $j \in \{0, \dots, 6\}$ jako

$$\begin{aligned} \tilde{\psi}^*(0) &\stackrel{\text{def}}{=} 0, \quad \tilde{\psi}^*\left(\frac{1}{6!}\right) \stackrel{\text{def}}{=} \frac{1}{5}\rho_5, \quad \tilde{\psi}^*\left(\frac{2}{6!}\right) \stackrel{\text{def}}{=} \frac{2}{5}\rho_5, \quad \tilde{\psi}^*\left(\frac{3}{6!}\right) \stackrel{\text{def}}{=} \frac{3}{5}\rho_5, \\ \tilde{\psi}^*\left(\frac{4}{6!}\right) &\stackrel{\text{def}}{=} \frac{4}{5}\rho_5, \quad \tilde{\psi}^*\left(\frac{5}{6!}\right) \stackrel{\text{def}}{=} \frac{1}{2 \cdot 5!}, \quad \tilde{\psi}^*\left(\frac{1}{5!}\right) \stackrel{\text{def}}{=} \frac{1}{5!}. \end{aligned} \quad (3.5)$$

Předpokládejme, že hodnoty $\widetilde{\psi}^* \left(\frac{i}{(k-1)!} \right)$ jsou již definovány. Položme $j \stackrel{\text{def}}{=} ki + t$, kde $t \in \{0, \dots, k-1\}$, je-li $i \in \{0, \dots, \frac{k!-1}{5!}\}$ a $t \stackrel{\text{def}}{=} 0$, je-li $i = \frac{k!}{5!}$, a definujme

$$\widetilde{\psi}^* \left(\frac{j}{k!} \right) = \widetilde{\psi}^* \left(\frac{i}{(k-1)!} + \frac{t}{k!} \right) \stackrel{\text{def}}{=} \begin{cases} \widetilde{\psi}^* \left(\frac{i}{(k-1)!} \right) + \frac{t}{k-1} \rho_{k-1} & \text{pro } t \in \{0, \dots, k-2\}, \\ \frac{1}{2} \left\{ \widetilde{\psi}^* \left(\frac{i}{(k-1)!} \right) + \widetilde{\psi}^* \left(\frac{i+1}{(k-1)!} \right) \right\} & \text{pro } t = k-1. \end{cases} \quad (3.6)$$

Tedy funkci $\widetilde{\psi}^*$ definujeme rekurzivně na intervalech $\left\langle \frac{i}{(k-1)!}, \frac{i+1}{(k-1)!} \right\rangle$ a to tak, že hodnoty uvnitř tohoto intervalu závisí pouze na funkčních hodnotách $\widetilde{\psi}^*$ v krajních bodech tohoto intervalu. Interval $\left\langle \frac{i}{(k-1)!}, \frac{i+1}{(k-1)!} \right\rangle$ ekvidistantně rozdělíme na k částí a v bodech, které interval $\left\langle \frac{i}{(k-1)!}, \frac{i+1}{(k-1)!} \right\rangle$ rozdělují, definujeme nové funkční hodnoty vždy o $\frac{1}{k-1} \rho_{k-1}$ větší než již definovaná hodnota v levém sousedním bodě, s výjimkou bodu nejvíce napravo, jemuž přiřadíme funkční hodnotu rovnu aritmetickému průměru hodnot $\widetilde{\psi}^*$ v krajních bodech intervalu $\left\langle \frac{i}{(k-1)!}, \frac{i+1}{(k-1)!} \right\rangle$. Tuto konstrukci ilustruje obrázek 3.1.



Obrázek 3.1: Konstrukce funkce $\widetilde{\psi}^*$.

Z vlastností číselné posloupnosti $\{\rho_k\}_5^\infty$ plyne, že posloupnost takto konstruovaných hodnot funkce $\widetilde{\psi}^*$ je monotonně rostoucí vzhledem ke svým argumentům. Zřejmě množina

$$\left\{ \frac{i}{k!} \mid k \in \{1, \dots, \infty\}, i \in \{0, \dots, k!\} \right\}$$

je hustá v intervalu $\left\langle 0, \frac{1}{5!} \right\rangle$, máme tedy definovanou funkci $\widetilde{\psi}^*$ na husté podmnožině $\left\langle 0, \frac{1}{5!} \right\rangle$. Na doplňku této podmnožiny do intervalu $\left\langle 0, \frac{1}{5!} \right\rangle$ lze zřejmě funkci $\widetilde{\psi}^*$ spojitě dodefinovat.

S využitím funkce $\widetilde{\psi}^*$ nyní definujeme funkci $\widetilde{\psi} : \bar{D}_1 \rightarrow \mathfrak{R}$ podle následujícího předpisu:

$$\widetilde{\psi}(x) \stackrel{\text{def}}{=} \begin{cases} \widetilde{\psi}^*(x) & \text{pro } x \in \bar{D}_0, \\ \widetilde{\psi}^* \left(x - \frac{i}{5!} \right) + \frac{i}{5!}, & \text{pro } x \in \left\langle \frac{i}{5!}, \frac{i+1}{5!} \right\rangle, \quad i \in \{0, \dots, 5!\}. \end{cases} \quad (3.7)$$

Evidentně $\widetilde{\psi}$ je monotonně rostoucí, zobrazující interval \bar{D}_1 na sebe sama. Základní vlastnost funkce $\widetilde{\psi}$ vystihuje následující lemma.

Lemma 3.3.2 *Nechť platí*

$$\delta_k \stackrel{\text{def}}{=} \frac{1}{k!} - \sum_{r=k+1}^{\infty} \frac{1}{r!} \quad a \quad \nu_k \stackrel{\text{def}}{=} \sum_{r=k}^{\infty} \frac{r-1}{r} \rho_r.$$

Potom platí

$$\tilde{\psi} \left(\frac{i}{k!} + \delta_k \right) = \tilde{\psi} \left(\frac{i}{k!} \right) + \nu_k.$$

■ *Důkaz:*

DŮKAZ JE UVEDEN V ORIGINÁLNÍ TIŠTĚNÉ VERZI. ■

Ještě uvedme lemmu, která bude nutná pro důkaz tvrzení o hustotě.

Definice 3.3.2 *Posloupnost reálných čísel $\{\lambda_i\}_1^\infty$ je CELOČÍSELNĚ NEZÁVISLÁ, jestliže pro libovolnou konečnou posloupnost celých čísel t_1, \dots, t_n splňujících podmínku*

$$\sum_{p=1}^n |t_p| \neq 0$$

platí, že $\sum_{p=1}^n t_p \lambda_p \neq 0$. Dále nechť pro každé $k \geq 5$ symbol \bar{I}_k označuje množinu celých čísel

$$\bar{I}_k \stackrel{\text{def}}{=} \left\{ i \in \{0, \dots, +\infty\} \mid 0 \leq i \leq \left(1 + \frac{1}{5!}\right) k! \right\}.$$

Lemma 3.3.3 *Nechť $\{\lambda_i\}_1^\infty$, $\lambda_i > 0$ je celočíselně nezávislá posloupnost a nechť pro každé $k \geq 5$ je*

$$\mu_k \sigma_k \stackrel{\text{def}}{=} \min_{i_p, j_p \in \bar{I}_k} \left\{ \left| \sum_{p=1}^{k-3} \lambda_p \left[\tilde{\psi} \left(\frac{i_p}{k!} \right) - \tilde{\psi} \left(\frac{j_p}{k!} \right) \right] \right| \mid \sum_{p=1}^{k-3} |i_p - j_p| \neq 0 \right\}, \quad (3.8)$$

kde je

$$\sigma_k \stackrel{\text{def}}{=} 1 + \sum_{p=1}^{k-3} \lambda_p.$$

Potom je možno volit výchozí racionální hodnoty $\{\rho_k\}_5^\infty$ ve výrazu 3.4 tak, aby platilo $\nu_k < \mu_k$.

■ *Důkaz:*

DŮKAZ JE UVEDEN V ORIGINÁLNÍ TIŠTĚNÉ VERZI. ■

3.3.2 Konstrukce rekurzivního pokrytí jednotkové krychle

Na základě předchozího výkladu definujme systémy intervalů $\bar{E}_{k,i}$ a $\bar{E}_{k,i}^q$ podle následujících definic. Základní vlastnosti těchto intervalů jsou obsaženy v následujících dvou lemmách.

Definice 3.3.3 Pro každé přirozené $k \geq 5$ definujme systém uzavřených intervalů

$$\bar{E}_{k,i} \stackrel{\text{def}}{=} \left\langle \frac{i}{k!}, \frac{i}{k!} + \delta_k \right\rangle, \quad i \in \bar{I}_k.$$

Zřejmě tyto intervaly jsou pro každou pevnou hodnotu $k \geq 5$ od sebe vzdáleny mezerou délky

$$e_k \stackrel{\text{def}}{=} \frac{1}{k!} - \delta_k = \sum_{r=k+1}^{\infty} \frac{1}{r!} \quad (3.9)$$

a dá se přímým výpočtem ověřit, že mezi délkou těchto mezer a intervalů platí vztah

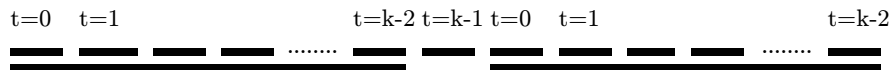
$$\delta_k > (k-1)e_k. \quad (3.10)$$

S výjimkou těchto mezer intervaly $\bar{E}_{k,i}$ pokrývají interval \bar{D}_1 . Navíc se snadno ověří, že platnost inkluze

$$\bar{E}_{k,j} \subset \bar{E}_{k-1,i}$$

je ekvivalentní s podmínkou, že $j = ki + t$, $t \in \{0, \dots, k-2\}$. Jestliže tato inkluze platí, pak počáteční body $\bar{E}_{k,j}$ a $\bar{E}_{k-1,i}$ jsou totožné pro $t = 0$ a jejich koncové body jsou totožné pro $t = k-2$. Pro $t = k-1$ interval $\bar{E}_{k,j}$ leží v mezeře, oddělující intervaly $\bar{E}_{k,j}$ a $\bar{E}_{k-1,i}$. Navíc z monotonie funkce $\tilde{\psi}$ a z lemy 3.3.2 plyne

$$\tilde{\psi} \left(\left\langle \frac{i}{k!}, \frac{i}{k!} + \delta_k \right\rangle \right) = \left\langle \tilde{\psi} \left(\frac{i}{k!} \right), \tilde{\psi} \left(\frac{i}{k!} \right) + \nu_k \right\rangle.$$



Obrázek 3.2: Vlastnosti intervalů $\bar{E}_{k,j} \subset \bar{E}_{k-1,i}$ ($\bar{E}_{k-1,i}$ dole, $\bar{E}_{k,j}$ nahoře).

Definice 3.3.4 Pro libovolná přirozená $k \geq 5$, $n \geq 2$ a $q \in \{0, \dots, 2n\}$ definujme systém uzavřených intervalů

$$\bar{E}_{k,i}^q \stackrel{\text{def}}{=} \left\langle \frac{i}{k!} - qe_{2n+1}, \frac{i}{k!} + \delta_k - qe_{2n+1} \right\rangle, \quad i \in \bar{I}_k,$$

kde $i \in \bar{I}_k$.

Intervaly $\bar{E}_{k,i}^q$ jsou de facto rovny intervalům $\bar{E}_{k,i}$, posunutým o hodnotu $-qe_{2n-1}$.

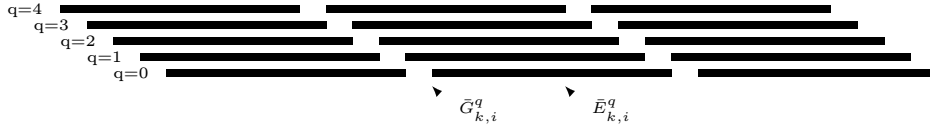
Lemma 3.3.4 Nechť je dáno celé číslo $n \geq 2$. Potom pro každé celé číslo $k \geq 2n+1$ a pro každé $q \in \{0, \dots, 2n\}$ platí

1. délka intervalu $\bar{E}_{k,i}^q$ se blíží k nule pro $k \rightarrow \infty$.

2. Kdykoli je $i \neq j$, pak

$$\bar{E}_{k,i}^q \cap \bar{E}_{k,j}^q = \emptyset. \quad (3.11)$$

Dále pro každý bod x z intervalu \bar{D}_1 existuje alespoň $2n$ hodnot q takových, že $x \in \bar{E}_{k,i}^q$.

Obrázek 3.3: Intervaly $\bar{E}_{k,i}^q$.

■ *Důkaz:*

DŮKAZ JE UVEDEN V ORIGINÁLNÍ TIŠTĚNÉ VERZI. ■

Lemma 3.3.5 *Nechť jsou dána celá čísla $n \geq 2$, $q \in \{0, \dots, 2n\}$, $k \geq 2n + 1$ a pro každé $p \in \{1, \dots, n\}$ nechť je $i_p \in \bar{I}_k$. Dále nechť posloupnost čísel $\{\lambda_i\}_1^n$ je celočíselně nezávislá. Dále definujme následující intervaly*

$$\bar{T}_{k,i_1,\dots,i_n}^q \stackrel{\text{def}}{=} \left\langle \sum_{p=1}^n \lambda_p \tilde{\psi} \left(\frac{i_p}{k!} + qe_{2n+1} \right), \sum_{p=1}^n \lambda_p \left[\tilde{\psi} \left(\frac{i_p}{k!} + qe_{2n+1} \right) + \nu_k \right] \right\rangle. \quad (3.12)$$

Potom

1. Velikost intervalů $\bar{T}_{k,i_1,\dots,i_n}^q$ se blíží k nule pro $k \rightarrow \infty$.
2. Kdykoli je $(i_1, \dots, i_n) \neq (j_1, \dots, j_n)$, tak

$$\bar{T}_{k,i_1,\dots,i_n}^q \cap \bar{T}_{k,j_1,\dots,j_n}^q = \emptyset. \quad (3.13)$$

Platí-li ještě navíc

$$\xi_q \stackrel{\text{def}}{=} \sum_{p=1}^n \lambda_p \tilde{\psi}(\bar{\mathbf{x}}_p + qe_{2n+1}), \quad (3.14)$$

tak pro libovolný bod $\bar{\mathbf{x}} \in \langle 0, 1 \rangle^n$ existuje alespoň $n+1$ hodnot q , pro které je $\xi_q \in \bar{T}_{k,i_1,\dots,i_n}^q$.

■ *Důkaz:*

DŮKAZ JE UVEDEN V ORIGINÁLNÍ TIŠTĚNÉ VERZI. ■

3.3.3 Konstrukce přenosových funkcí a neuronové sítě

Vlastnosti intervalů $\bar{E}_{k,i}^q$ uvedené v lemmě 3.3.5 umožňují bezprostřední důkaz modifikované Kolmogorovy věty ve tvaru vhodně korespondujícím s vrstevnatou architekturou neuronových sítí. Tato konstrukce je odvozena z původního Kolmogorova důkazu, uvedeného například v [Kol57].

Definice 3.3.5 *Pro každé $p \in \{1, \dots, n\}$ a $k \geq 2n + 1$ nechť je $i_p \in \bar{I}_k$. Potom definujme množinu*

$$\bar{\Omega}_{k,i_1,\dots,i_n}^q \stackrel{\text{def}}{=} \prod_{p=1}^n \bar{E}_{k,i_p}^q,$$

($\bar{\Omega}_{k,i_1,\dots,i_n}^q$ je tedy kartézským součinem intervalů \bar{E}_{k,i_p}^q).

Lemma 3.3.6 Pro libovolnou spojitou funkci $\tilde{g} : \langle 0, 1 \rangle \rightarrow \mathfrak{R}$ a pro libovolnou krychli $\bar{\Omega}_{k,i_1,\dots,i_n}^q$ označme

$$\Gamma_{k,i_1,\dots,i_n}^q \stackrel{\text{def}}{=} \max \left\{ \left| \tilde{g}(\vec{x}) - \alpha_{k,i_1,\dots,i_n}^q \right| \mid \vec{x} \in \bar{\Omega}_{k,i_1,\dots,i_n}^q \right\},$$

kde

$$\alpha_{k,i_1,\dots,i_n}^q \stackrel{\text{def}}{=} \min \left\{ \tilde{g}(\vec{x}) \mid \vec{x} \in \bar{\Omega}_{k,i_1,\dots,i_n}^q \right\}.$$

Nechť $\Gamma_{\tilde{g},k}$ je maximum čísel $\Gamma_{k,i_1,\dots,i_n}^q$ přes všechny intervaly $\bar{\Omega}_{k,i_1,\dots,i_n}^q$. Potom pro libovolné číslo $\epsilon > 0$ lze volit přirozené číslo k tak, že $\Gamma_{\tilde{g},k} \leq \epsilon$.

■ *Důkaz:*

DŮKAZ JE UVEDEN V ORIGINALNÍ TIŠTĚNÉ VERZI.

Definice 3.3.6 Nechť je dáno přirozené $n \geq 2$ a funkce $\tilde{f} : \langle 0, 1 \rangle^n \rightarrow \mathfrak{R}$. Pro každé $q \in \{0, \dots, 2n\}$ definujme funkce $\widetilde{\Phi}_{q,0} : \mathfrak{R} \rightarrow \mathfrak{R}$ jako nulové funkce a $\widetilde{f}_0 : \langle 0, 1 \rangle^n \rightarrow \mathfrak{R}$ také jako nulovou funkci. Dále definujme

$$M_0 \stackrel{\text{def}}{=} \sup_{\vec{x} \in \langle 0, 1 \rangle^n} \left| \tilde{f}(\vec{x}) \right|$$

a číslo k_0 tak, aby $\Gamma_{\tilde{f},k_0} \leq \frac{1}{2n+2} M_0$.

Předpokládejme, že pro nějaké přirozené $r \geq 0$ jsou již definovány funkce \widetilde{f}_r a $\widetilde{\Phi}_{q,r}$ a čísla M_r, k_r . Potom definujme

$$\widetilde{f}_{r+1}(\vec{x}) \stackrel{\text{def}}{=} \sum_{q=0}^{2n} \widetilde{\Phi}_{q,r} \left(\sum_{p=1}^n \lambda_p \tilde{\psi}(\vec{x}_p + qe_{2n+1}) \right)$$

a pro všechna $y \in \bar{T}_{k,i_1,\dots,i_n}^q$

$$\widetilde{\Phi}_{q,r+1}(y) \stackrel{\text{def}}{=} \widetilde{\Phi}_{q,r}(y) + \frac{1}{n+1} \left| \tilde{f}(\vec{z}) - \widetilde{f}_r(\vec{z}) \right|, \quad (3.15)$$

kde \vec{z} je libovolný bod ležící v krychli $\bar{\Omega}_{k,i_1,\dots,i_n}^q$. Dále definujme

$$M_{r+1} \stackrel{\text{def}}{=} \sup_{\vec{x} \in \langle 0, 1 \rangle^n} \left| \tilde{f}(\vec{x}) - \widetilde{f}_r(\vec{x}) \right|$$

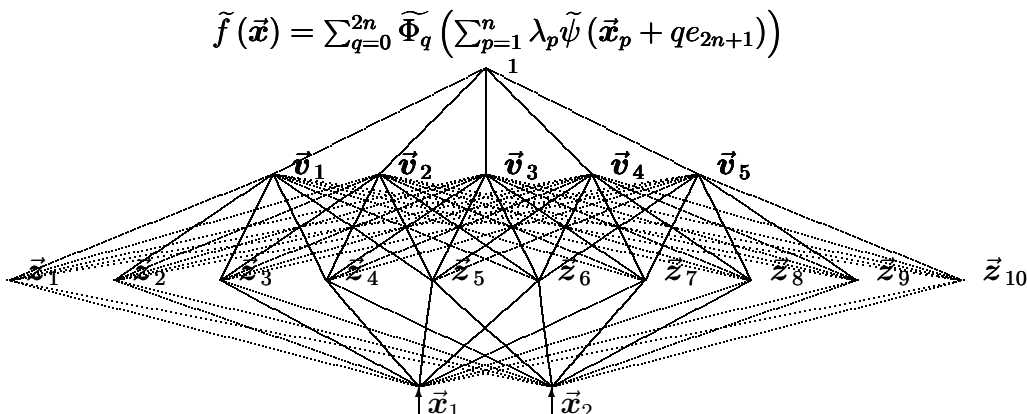
a číslo k_{r+1} tak, aby platilo $\Gamma_{(\tilde{f}-\widetilde{f}_r),k_{r+1}} \leq \frac{1}{2n+2} M_r$.

Lemma 3.3.7 Platí $\lim_{r \rightarrow \infty} M_r = 0$.

■ *Důkaz:*

DŮKAZ JE UVEDEN V ORIGINALNÍ TIŠTĚNÉ VERZI.

Poznamenejme, $\frac{2n+1}{2n+2}$ se pro velká n blíží k jedné, na základě čehož se dá usuzovat na obecně pomalou rychlost konvergence posloupnosti $\{M_r\}_0^\infty$, což se projeví nutností volit hodnoty k_r velmi vysoké, chceme-li zaručit dostatečnou rychlost konvergence $\{M_r\}_0^\infty$. Na základě konvergence posloupnosti $\{M_r\}_0^\infty$ můžeme vyslovit tvrzení, které je obdobné původní Kolmogorově větě, ale formálně popisuje zobrazení, které lze realizovat prostřednictvím neuronové sítě sestávající ze vstupní vrstvy, dvou následujících vrstev a jednoho výstupního neuronu.



Obrázek 3.4: Schéma architektury sítě odvozené od Kolmogorovy věty pro $n = 2$. ($\tilde{z}_{(p-1)n+q} = \tilde{\psi}(\vec{x}_p + qe_{2n+1})$, $\tilde{v}_q = \tilde{\Phi}_q \left(\sum_{p=1}^n \lambda_p \tilde{\psi}(\vec{x}_p + qe_{2n+1}) \right)$.)

Tvrzení 3.3.8 *Nechť $\lambda_1, \dots, \lambda_n$ je kladná, celočíselně nezávislá posloupnost. Potom existuje spojitá, monotónně rostoucí funkce $\tilde{\psi} : \bar{D}_1 \rightarrow \bar{D}_1$, taková, že pro každou reálnou spojitou funkci $\tilde{f} : \langle 0, 1 \rangle^n \rightarrow \mathfrak{R}$, $n \geq 2$, existují spojitě funkce $\tilde{\Phi}_q$, $q \in \{0, \dots, 2n\}$, takové, že platí rovnice*

$$\tilde{f}(\vec{x}) = \sum_{q=0}^{2n} \tilde{\Phi}_q \left(\sum_{p=1}^n \lambda_p \tilde{\psi}(\vec{x}_p + qe_{2n+1}) \right).$$

■ *Důkaz:*

DŮKAZ JE UVEDEN V ORIGINÁLNÍ TIŠTĚNÉ VERZI.

■

Zřejmý důsledek předchozího tvrzení udává architekturu a základní vlastnosti neuronů sítě, která přesně počítá spojitou funkci definovanou na jednotkové krychli prostoru \mathfrak{R}^n .

Důsledek 3.3.9 *Každá spojitá funkce $\tilde{f} : \langle 0, 1 \rangle^n \rightarrow \mathfrak{R}$, $n \geq 2$, může být přesně spočítána neuronovou sítí s n neurony ve vstupní vrstvě, $(2n + 1)n$ neurony ve druhé vrstvě, $2n + 1$ neurony ve třetí vrstvě a s jedním výstupním neuronem (=čtvrtá vrstva). Navíc platí, že váhy mezi vstupní vrstvou a druhou vrstvou jsou buď 1 nebo 0, váhy mezi druhou a třetí vrstvou jsou buď 0 nebo λ_i , $i \in \{1, \dots, n\}$, váhy mezi neurony třetí vrstvy a výstupním neuronem jsou rovny jedničce a žádná váha není závislá na \tilde{f} . Dále neurony ve druhé vrstvě a vstupní neuron nezávisí na implementované funkci \tilde{f} .*

3.4 Konvoluční přístup

Základní myšlenka konvoluční metody spočívá v konstrukci jádra, které je odvozené ze sigmoidálních funkcí a které může aproximovat Diracovu delta funkci (t.j. distribuci, jejíž nosič je pouze jednobodová množina). Samotná konvoluce a současně jádro jsou aproximovány kvadraturou, na základě které dostáváme požadovanou aproximaci dané spojitě

funkce. Začneme tvrzením o základních vlastnostech uniformní konvergence konvolucí. (Připomeňme, že $\|\vec{x}\|_\infty \stackrel{\text{def}}{=} \max_{i \in \{1, \dots, n\}} |\vec{x}_i|$ a $\|\vec{x}\|_1 \stackrel{\text{def}}{=} \sum_{i=1}^n |\vec{x}_i|$).

Tvrzení 3.4.1 Předpokládejme, že \tilde{f} je omezená, stejnoměrně spojitá funkce na \mathfrak{R}^n a že pro $\tilde{g} \in \mathcal{L}_1^{\mathfrak{R}^n}$ platí:

$$\int_{\mathfrak{R}^n} \tilde{g}(\vec{x}) d(\vec{x}) = 1.$$

Dále definujeme funkce $\tilde{g}_m(\vec{x}) \stackrel{\text{def}}{=} m^n \tilde{g}(m\vec{x})$. Potom platí:

1. $(\tilde{f} * \tilde{g}_m)$ konverguje stejnoměrně k funkci \tilde{f} pro $m \rightarrow \infty$.
2. pro libovolné číslo $R > 0$ platí

$$\|(\tilde{f} * \tilde{g}_m) - \tilde{f}\|_\infty \leq \omega\left(\tilde{f}, \frac{2R}{m}\right) \|\tilde{g}\|_1 + 2\|\tilde{f}\|_\infty \int_{\|\vec{s}\| > R} |\tilde{g}(\vec{s})| d\vec{s},$$

kde norma $\|\cdot\|_\infty$ je brána vzhledem k celému \mathfrak{R}^n .

3. Nechť \tilde{f} Lipstichovská s konstantou λ a

$$M \stackrel{\text{def}}{=} \int_{\mathfrak{R}^n} \|\vec{x}\| |\tilde{g}(\vec{x})| d\vec{x} < \infty.$$

Potom platí odhad $\|(\tilde{f} * \tilde{g}_m) - \tilde{f}\|_\infty \leq \frac{M\lambda}{m}$.

■ *Důkaz:*

DŮKAZ JE UVEDEN V ORIGINALNÍ TIŠTĚNÉ VERZI. ■

Definice 3.4.1 Pro libovolnou spojitou funkci $\tilde{\phi} \in C_{\mathfrak{R}^n}$ definujeme KONVOLUČNÍ JÁDRO jako funkci

$$\tilde{G}_{\tilde{\phi}}(\vec{x}) \stackrel{\text{def}}{=} \frac{1}{\alpha_n} \int_{\|\vec{u}\|=1} \tilde{\phi}(\langle \vec{x} | \vec{u} \rangle) d(\vec{u}), \quad \text{kde } \alpha_n \stackrel{\text{def}}{=} \int_{\|\vec{u}\|=1} d(\vec{u}), \quad (3.16)$$

(α_n je povrch jednotkové koule v prostoru \mathfrak{R}^n).

Toto jádro použijeme k důkazu následující stěžejní věty, která je založena na obsahu tvrzení 3.4.1.

Tvrzení 3.4.2 Nechť pro nějaké $a > 0$ je $\bar{K} \stackrel{\text{def}}{=} \langle a, a \rangle^n$ a $\tilde{\phi} \in C_{\mathfrak{R}}$ je stejnoměrně spojitá. Předpokládejme, že $\tilde{G}_{\tilde{\phi}}(\vec{x})$ je definována podle 3.16. Potom jestliže $\tilde{G}_{\tilde{\phi}}(\cdot) \in \mathcal{L}_1^{\mathfrak{R}^n}$ a současně

$$\int_{\mathfrak{R}^n} \tilde{G}_{\tilde{\phi}}(\vec{x}) d(\vec{x}) \neq 0,$$

tak množina funkcí ve tvaru $\tilde{\phi}(\langle \vec{x} | \vec{a} \rangle - c)$, $\vec{a} \in \mathfrak{R}^n$, $c \in \mathfrak{R}$, je hustá v prostoru $C_{\bar{K}}$.

■ *Důkaz:*

DŮKAZ JE UVEDEN V ORIGINALNÍ TIŠTĚNÉ VERZI. ■

3.4.1 Volba konvolučního jádra

Abychom ukázali, že množina všech funkcí tvaru $\tilde{\phi}(\langle \vec{x} | \vec{a} \rangle - c)$, $\vec{a} \in \mathfrak{R}^n$, $c \in \mathfrak{R}$ je hustá v $C_{\bar{K}}$, musíme jednak dokázat, že konvoluční jádro $\tilde{G}_{\tilde{\phi}}(\cdot)$ definované v 3.16 patří do prostoru $\mathcal{L}_1^{\mathfrak{R}^n}$, a dále, že integrál tohoto jádra přes celý definiční obor je od nuly různý. Následující lemma nám pomůže stanovit takovou mocninu proměnné $r \stackrel{\text{def}}{=} \|\vec{x}\|$, aby jádro $\tilde{G}_{\tilde{\phi}}(\vec{x})$, omezené pro velká r touto mocninou, bylo integrovatelné a tedy v $\mathcal{L}_1^{\mathfrak{R}^n}$.

Lemma 3.4.3 *Nechť $p, R \in \mathfrak{R}$, $R > 0$. Potom platí, že*

$$\int_{\|\vec{x}\|=1} \frac{1}{\|\vec{x}\|^p} d(\vec{x}) < \infty$$

právě když $p > n$.

■ *Důkaz:*

DŮKAZ JE UVEDEN V ORIGINÁLNÍ TIŠTĚNÉ VERZI.

Abychom dokázali, že $\tilde{G}_{\tilde{\phi}}(\cdot) \in \mathcal{L}_1^{\mathfrak{R}^n}$, stačí ukázat, že $\tilde{G}_{\tilde{\phi}}(\vec{x}) = O\left(\frac{1}{\|\vec{x}\|^n}\right)$ pro $\|\vec{x}\| \rightarrow \infty$. V tomto důkazu nám pomůže vyjádření jádra $\tilde{G}_{\tilde{\phi}}(\vec{x})$ v následujícím tvaru:

Lemma 3.4.4 *Nechť $\tilde{G}_{\tilde{\phi}}(\vec{x})$ je definováno dle 3.16. Potom $\tilde{G}_{\tilde{\phi}}(\vec{x}) = \tilde{g}_0(\tilde{\phi}, r)$, kde $r \stackrel{\text{def}}{=} \|\vec{x}\|$ a*

$$\begin{aligned} \tilde{g}_0(\tilde{\phi}, r) &\stackrel{\text{def}}{=} \frac{\alpha_{n-2}}{\alpha_{n-1}} \int_{-1}^1 \tilde{\phi}(rs) (1-s^2)^{\frac{n-3}{2}} d(s) \\ &= \frac{\alpha_{n-2}}{\alpha_{n-1}} \int_{-r}^r r^{2-n} \tilde{\phi}(t) (r^2-t^2)^{\frac{n-3}{2}} d(t), \quad r \neq 0. \end{aligned} \tag{3.17}$$

■ *Důkaz:*

DŮKAZ JE UVEDEN V ORIGINÁLNÍ TIŠTĚNÉ VERZI.

Ukážeme, že pro důkaz hustoty sigmoidálních funkcí v množině funkcí spojitých na kompaktní množině je vhodné zkoumat ne přímo vlastnosti dané sigmoidální funkce, ale jisté funkce od ní odvozené. Zkoumejme otázku, jakou podmínku musí splňovat funkce $\tilde{\phi}$, aby jádro $\tilde{G}_{\tilde{\phi}}(\cdot)$ bylo v $\mathcal{L}_1^{\mathfrak{R}^n}$. Základním požadavkem je, aby $\tilde{\phi}$ šlo dostatečně rychle k nule pro velké absolutní hodnoty argumentů. Zřejmě hodnota jádra $\tilde{G}_{\tilde{\phi}}(\vec{x})$ je nulová, je-li $\tilde{\phi}$ lichá funkce. Zvolme tedy $\tilde{\phi}$ tak, aby byla sudá (není to ale nutné). Pro sigmoidální funkci $\tilde{\sigma}$ definujme funkci

$$\tilde{\psi}(t) \stackrel{\text{def}}{=} \tilde{\sigma}(1+t) + \tilde{\sigma}(1-t).$$

Zřejmě $\tilde{\psi}$ je sudá funkce, která jde v $\pm\infty$ k nule. Předpokládejme, že $\tilde{\sigma}$ je spojitá funkce a že

$$|\tilde{\psi}(t)| < \frac{\eta}{|t|^p}, \quad p > n-2,$$

pro nějakou kladnou hodnotu η . Rozvineme-li výraz 3.17 do mocninné řady, dostaneme

$$\tilde{g}_0(\tilde{\phi}, r) = \frac{1}{r^{n-2}} \sum_{j=0}^{\lambda} \tilde{\beta}_j(r) r^{2\lambda-2j}, \tag{3.18}$$

kde

$$\widetilde{\beta}_j(r) \stackrel{\text{def}}{=} \frac{\alpha_{n-2}}{\alpha_{n-1}} \binom{\lambda}{j} \int_{-r}^{+r} \widetilde{\phi}(t) t^{2j} d(t).$$

Tvaru tohoto rozkladu využijeme pro důkaz následující lemma.

Lemma 3.4.5 *Nechť n je liché a pro $\widetilde{\psi}(t) \stackrel{\text{def}}{=} \widetilde{\phi}(1+t) + \widetilde{\phi}(1-t)$ a pro nějaké $\eta > 0$ platí $|\widetilde{\psi}(t)| \leq \frac{\eta}{|t|^p}$, $p > n - 2$. Potom $\widetilde{G}_{\widetilde{\phi}}(\cdot) \in \mathcal{L}_1^{\mathfrak{R}^n}$ tehdy a jen tehdy, je-li*

$$\int_0^\infty \widetilde{\psi}(t) t^{2j} d(t) = 0, \quad j \in \{0, \dots, \frac{n-3}{2}\}. \quad (3.19)$$

■ *Důkaz:*

DŮKAZ JE UVEDEN V ORIGINÁLNÍ TIŠTĚNÉ VERZI.

■

Lemma 3.4.6 *Předpokládejme, že n je liché a že funkce $\widetilde{\psi}$ splňuje předpoklady a tvrzení lemma 3.4.5. Dále předpokládejme, že $\widetilde{\psi}$ je sudá funkce. Potom platí*

$$\int_{\mathfrak{R}^n} \widetilde{G}_{\widetilde{\phi}}(\vec{x}) d(\vec{x}) = -2\alpha_{n-2}\tau_n \int_0^\infty \widetilde{\psi}(t) t^{n-1} d(t), \quad (3.20)$$

kde je

$$\tau_n \stackrel{\text{def}}{=} \int_0^1 r (1-r^2)^{\frac{n-3}{2}} d(r) > 0.$$

■ *Důkaz:*

DŮKAZ JE UVEDEN V ORIGINÁLNÍ TIŠTĚNÉ VERZI.

■

Spojíme-li všechny předchozí lemma dohromady, můžeme dokázat následující tvrzení, kterým uzavřeme kapitolu, týkající se aproximačních vlastností neuronových sítí.

Tvrzení 3.4.7 *Nechť $\widetilde{\sigma} \in C_{\mathfrak{R}}$ je stejnoměrně spojitá na \mathfrak{R} a $\widetilde{\psi}(t) \stackrel{\text{def}}{=} \widetilde{\sigma}(1+t) + \widetilde{\sigma}(1-t)$. Předpokládejme, že n je liché a $\bar{K} = \langle -a, a \rangle^n$, $a > 0$ a pro nějaké $\eta > 0$ platí $|\widetilde{\psi}(t)| \leq \frac{\eta}{|t|^p}$, $p > n - 2$. Dále nechť*

$$\int_0^\infty \widetilde{\psi}(t) t^{n-1} d(t) \neq 0.$$

Potom množina funkcí tvaru $\widetilde{\sigma}(\langle \vec{x} | \vec{w} \rangle + c)$, kde $\vec{w} \in \mathfrak{R}^n$ a $c \in \mathfrak{R}$, je hustá v prostoru $C_{\bar{K}}$.

■ *Důkaz:*

DŮKAZ JE UVEDEN V ORIGINÁLNÍ TIŠTĚNÉ VERZI.

■

Kapitola 4

Vapnik-Chervonenkova dimenze

Einsteinův axiom

POKUD VYCHÁZEJÍ MATEMATICKÉ POUČKY ZE SKUTEČNOSTI,
NEJSOU SPOLEHLIVÉ.
POKUD JSOU SPOLEHLIVÉ,
NEVYCHÁZEJÍ ZE SKUTEČNOSTI.

Nedílnou součástí při výkladu teorie umělých neuronových sítí je i problematika spojená se schopnostmi neuronových sítí efektivně přiřazovat svým vnitřním parametrům (váhy jednotlivých spojů, prahy uzlů atd.) takové hodnoty, aby odezva neuronové sítě na dané vstupní hodnoty byla rovna hodnotám požadovaným, nebo je alespoň uspokojivě aproximovala. Procesu nastavování těchto parametrů na základě extrakce vlastností vstupních hodnot se říká učení. Teorie učení byla rozpracována obecně, nejenom pro neuronové sítě. K rigoróznímu výkladu té části teorie učení, která bezprostředně souvisí s neuronovými sítěmi je zapotřebí nejdříve zavést a detailně analyzovat jistou míru složitosti množin, která se označuje jako VC-dimenze. Této analýze věnujeme tuto část.

4.1 Pojem konceptu a třídy konceptů

Základním předmětem, ke kterému s váže pojem VC-dimenze, je jakákoli podmnožina potenční množiny libovolné, pevně dané množiny \bar{X} . V literatuře relevantní k našemu dalšímu výkladu jsou pro takovéto systémy podmnožin používány termíny koncept, hypotéza, třída konceptů a třída hypotéz. Proto této terminologie budeme v následujících definicích a výkladu používat také.

Definice 4.1.1 *Nechť \bar{X} je libovolná množina. Potom $\bar{c} \subset \bar{X}$ nazveme KONCEPTEM (nad množinou \bar{X}). Neprázdnou množinu $C \subset 2^{\bar{X}}$ nazveme TŘÍDOU KONCEPTŮ (nad množinou X). Prvky této množiny $\bar{c} \in C$, budeme nazývat koncepty třídy C .*

Jedním ze způsobů popisu konceptu je zadání konceptu jako vzor dané množiny při nějakém zobrazení. Formálně tedy definujeme pro booleovskou funkci \tilde{f} koncept $C_{\tilde{f}}$.

Definice 4.1.2 *Nechť \tilde{f} je zobrazení z množiny \bar{X} do dvouprvkové množiny $\{-1, +1\}$. Potom definujeme koncept $\bar{c}_{\tilde{f}}$ jako množinu*

$$\bar{c}_{\tilde{f}} \stackrel{\text{def}}{=} \{x \in \bar{X} \mid \tilde{f}(x) = 1\}.$$

Předchozí definici lze zřejmě rozšířit z definice konceptu na definici celé třídy konceptů, odpovídající dané množině zobrazení.

Definice 4.1.3 Předpokládejme, že \bar{F} je libovolná množina funkcí definovaných na množině \bar{X} s oborem hodnot $\{-1, +1\}$. Pak říkáme, že \bar{F} REPREZENTUJE TŘÍDU KONCEPTŮ

$$\mathbf{C}_{\bar{F}} \stackrel{\text{def}}{=} \left\{ \bar{A} \subset \bar{X} \mid (\exists \tilde{f} \in \bar{F})(\bar{A} = \bar{c}_{\tilde{f}}) \right\}.$$

V souladu s touto definicí lze definovat třídu konceptů, které jsou generovány neuronovou sítí typu perceptron.

Definice 4.1.4 Necht n je přirozené číslo a

$$\bar{F} \stackrel{\text{def}}{=} \left\{ \tilde{f} : \mathfrak{R}^n \rightarrow \{-1, +1\} \mid \tilde{f} = \widehat{sgn}(\langle \vec{x} \mid \vec{w} \rangle - t), \quad t \in \mathfrak{R}, \quad \vec{w}, \vec{x} \in \mathfrak{R}^n \right\}.$$

Potom definujme třídu konceptů

$$\mathbf{HALFSPACE}_n \stackrel{\text{def}}{=} \mathbf{C}_{\bar{F}}.$$

Zřejmě třída konceptů $\mathbf{HALFSPACE}_n$ obsahuje veškeré poloprostory v Eukleidovském prostoru dimenze n .

4.2 VC-dimenze třídy konceptů

Nyní uvedeme stěžejní pojem vyjadřující schopnost pevně daného systému množin vymezit v jiné množině libovolnou podmnožinu. Zřejmě vlastnosti každé neuronové sítě velmi silně závisí na vlastnostech třídy konceptů, která je reprezentována všemi možnými zobrazeními, která je daná architektura neuronové sítě schopna realizovat s ohledem na možné hodnoty vnitřních parametrů sítě. Zřejmě čím složitější geometrické vlastnosti konceptů, které lze pomocí neuronové sítě generovat, tím potenciálně vyšší možnost existence parametrů (vah atd.), které jsou vhodné pro řešení požadované konkrétní úlohy (buď aproximace funkce či separace množin). Intuitivně se dá očekávat, že taková třída konceptů, jejíž pomocí lze v dané množině vymezit co největší počet všech možných podmnožin dané velikosti je vhodnější, pro úlohy založené na separaci množin. Mírou této schopnosti je z jistého úhlu pohledu právě Vapnik-Chervonenkova dimenze.

Definice 4.2.1 Necht \bar{X} je daná množina, $\bar{Y} \subset \bar{X}$, $\mathbf{C} \subset 2^{\bar{X}}$. Potom říkáme, že \mathbf{C} ROZDĚLUJE \bar{Y} , právě když platí:

$$(\forall \bar{z} \subset \bar{Y})(\exists \bar{c} \in \mathbf{C})(\bar{c} \cap \bar{Y} = \bar{z}).$$

Dále definujme pojem VAPNIK-CHERVONENKOVA DIMENZE systému množin \mathbf{C} jako

$$\mathbf{VC}_{dim}(\mathbf{C}) \stackrel{\text{def}}{=} \max \left\{ |\bar{Y}| \mid \mathbf{C} \text{ rozděluje } \bar{Y} \right\},$$

nebo jako $+\infty$ v případě, že toto maximum neexistuje.

Poznamenejme, že Vapnik-Chervonenkova dimenze je rovna mohutnosti maximální množiny $\bar{A} \in \mathbf{C}$, pro kterou platí

$$\left\{ \bar{B} \mid (\exists \bar{Z} \in \mathbf{C})(\bar{B} = \bar{A} \cap \bar{Z}) \right\} = 2^{\bar{A}}.$$

Uvedme několik příkladů na výpočet VC-dimenze ([AB92]).

Příklad 4.2.1

1. Necht' $\bar{X} = \mathfrak{R}$ a \mathbf{C} je tvořeno všemi intervaly tvaru $(a, +\infty)$, kde $a \in \mathfrak{R}$. Potom zřejmě pro libovolnou $\{b, c \in \mathfrak{R} \mid b < c\}$ neexistuje žádný interval $(a, +\infty)$, který obsahuje bod b a neobsahuje bod c . Zřejmě tedy $\mathbf{VC}_{\dim}(\mathbf{C}) = 1$.
2. Necht' $\bar{X} = \mathfrak{R}$, s je přirozené číslo, $s > 1$ a \mathbf{C} je množina obsahující sjednocení právě s intervalů. Necht' $\bar{S} = \{x_1, \dots, x_{2s} \mid x_i < x_{i+1}\}$. Zřejmě \bar{S} je rozdělena třídou konceptů \mathbf{C} . Je-li ale $\bar{S} = \{x_1, \dots, x_{2s+1} \mid x_i < x_{i+1}\}$, pak každý koncept obsahující body $x_1, x_3, \dots, x_{2s+1}$ nutně obsahuje i body x_2, x_4, \dots, x_{2s} . Tedy takováto \bar{S} není rozdělena třídou konceptů \mathbf{C} , tedy $\mathbf{VC}_{\dim}(\mathbf{C}) = 2s$.
3. Necht' $\bar{X} = \mathfrak{R}^n$ a \mathbf{C} je tvořena všemi intervaly v prostoru \mathfrak{R}^n . Evidentně množina $2n$ bodů ležících ve středu stěn jednotkové krychle je třídou \mathbf{C} rozdělena. Vezmeme-li ale libovolnou množinu \bar{A} $2n + 1$ bodů, lze k ní přiřadit minimální interval obsahující celou \bar{A} . Pak ale musí existovat prvek $\bar{x} \in \bar{A}$, který je buď uvnitř minimálního intervalu, nebo na nějaké jeho stěně, a zřejmě množinu $\bar{A} - \{\bar{x}\}$ nelze intervalem oddělit od množiny $\{\bar{x}\}$. Proto $\mathbf{VC}_{\dim}(\mathbf{C}) = 2n$.

Příklad 4.2.2 Necht' \mathbf{C} je konečná třída konceptů. Jelikož k rozdělení množiny mohutnosti d je třeba minimálně 2^d různých konceptů, množina o mohutnosti větší než $\log_2(|\mathbf{C}|)$ nemůže být třídou \mathbf{C} rozdělena. Proto $\mathbf{VC}_{\dim}(\mathbf{C}) \leq \log_2(|\mathbf{C}|)$.

Příklad 4.2.3 Pro úplnost ukážeme příklad třídy konceptů s nekonečnou VC-dimenzí. Položme

$$\mathbf{C} \stackrel{\text{def}}{=} \left\{ \bar{A}_\alpha \mid (\exists \alpha \in \mathfrak{R}^n) \left(\bar{A}_\alpha = \left\{ x \in \mathfrak{R} \mid \widetilde{\sin}(\alpha x) \geq 0 \right\} \right) \right\}.$$

Potom pro libovolné přirozené l vezměme posloupnost čísel $\bar{Z}_l \stackrel{\text{def}}{=} \{z_i\}_1^l$, $z_i = \frac{1}{10^i}$. Dále každou podmnožinu množiny \bar{Z}_l , lze charakterizovat posoupností $\delta_1, \dots, \delta_l$, $\delta_i \in \{0, 1\}$. Zvolíme-li

$$\alpha \stackrel{\text{def}}{=} \pi \left(\sum_{i=1}^l (1 - \delta_i) 10^i + 1 \right),$$

pak platí rovnost

$$\alpha \frac{1}{10^j} = \pi \left(\sum_{i=1}^l (1 - \delta_i) 10^{i-j} + \frac{1}{10^j} \right).$$

Odtud

$$\alpha \frac{1}{10^j} = \pi \left(\sum_{i=1}^{j-1} \frac{1 - \delta_i}{10^{j-i}} + (1 - \delta_j) + \sum_{i=j+1}^l (1 - \delta_i) 10^{i-j} + \frac{1}{10^j} \right).$$

Proto je $\widetilde{\sin}(\alpha z_j) < 0$ pro $\delta_j = 0$ a $\widetilde{\sin}(\alpha z_j) > 0$ pro $\delta_j = 1$. Zřejmě tedy množina \bar{A}_α separuje každou podmnožinu \bar{Z}_l definovanou prostřednictvím $\delta_1, \dots, \delta_l$. Tedy pro libovolné l existuje množina mohutnosti l , která je rozdělena třídou konceptů \mathbf{C} , tedy $\mathbf{VC}_{\dim}(\mathbf{C}) = +\infty$.

Velmi užitečným nástrojem pro horní odhad VC-dimenze systému poloprostorů Eukleidovského prostoru \mathfrak{R}^n je následující Radonovo lemma (viz. [Lej85]). Předtím, než toto lemma vyslovíme, připomeňme, že lineární kombinaci $\sum_{i=1}^n \alpha_i \vec{x}_i$ libovolného konečného počtu vektorů z daného vektorového prostoru nazveme afinní kombinací vektorů, právě když $\sum_{i=1}^n \alpha_i = 1$. Dále připomeňme, že v n dimenzionálním prostoru je každá $n+2$ členná posloupnost vektorů afinně závislá, což znamená, že platnost rovnic

$$\sum_{i=1}^k \alpha_i \vec{x}_i = \vec{0} \quad \text{a} \quad \sum_{i=1}^k \alpha_i = 0$$

implikuje, že všechna α_i jsou rovna 0, $i \in \{1, \dots, k\}$.

Lemma 4.2.1 (Radonova) *Nechť $\bar{S} \stackrel{\text{def}}{=} \{\vec{x}_1, \dots, \vec{x}_k\} \subset \mathfrak{R}^n$, $k \geq n+2$, \vec{x}_i jsou po dvou od sebe různé. Potom existují množiny \bar{S}_1 a \bar{S}_2 tak, že $\bar{S}_1 \cup \bar{S}_2 = \bar{S}$, $\bar{S}_1 \cap \bar{S}_2 = \emptyset$ a*

$$[\bar{S}_1]_\kappa \cap [\bar{S}_2]_\kappa \neq \emptyset,$$

($[\bar{S}]_\kappa$ značí konvexní obal množiny \bar{S}).

■ *Důkaz:*

DŮKAZ JE UVEDEN V ORIGINÁLNÍ TIŠTĚNÉ VERZI.

■

Aplikací Radonova lemmatu dokážeme následující tvrzení pro VC-dimenzi poloprostorů.

Tvrzení 4.2.2 $\text{vc}_{\dim}(\text{HALFSPACE}_n) = n + 1$.

■ *Důkaz:*

DŮKAZ JE UVEDEN V ORIGINÁLNÍ TIŠTĚNÉ VERZI.

■

Vezměme v úvahu pouze takové třídy konceptů, která jsou definovány nad konečnými množinami \bar{X} . Potom následující Sauerovo lemma dává horní odhad mohutnosti každé třídy konceptů se zadanou VC-dimenzí. Zdůrazněme, jak plyne z důkazu lemmatu ([Sau72]), že tento odhad je čistě kombinatorickou vlastností konečných množin.

Lemma 4.2.3 (Sauer) *Nechť \bar{X} je konečná množina, $\mathcal{C} \subset 2^{\bar{X}}$. Potom*

$$|\mathcal{C}| \leq \sum_{i=0}^{\text{vc}_{\dim}(\mathcal{C})} \binom{|\bar{X}|}{i}.$$

Navíc existuje $\mathcal{C} \subset 2^{\bar{X}}$ tak, že platí rovnost.

■ *Důkaz:*

DŮKAZ JE UVEDEN V ORIGINÁLNÍ TIŠTĚNÉ VERZI.

■

Pojem zavedený v následující definici má velmi úzký vztah k VC-dimenzi a také k tvrzení Sauerovy lemmy.

Definice 4.2.2 *Nechť \mathcal{C} je neprázdňá třída konceptů nad \bar{X} a $\bar{S} \subset \bar{X}$. Potom definujeme*

$$\Pi_{\mathcal{C}}(\bar{S}) \stackrel{\text{def}}{=} \{\bar{S} \cap \bar{c} \mid \bar{c} \in \mathcal{C}\}.$$

Dále pro pevné $m \geq 0$ definujeme

$$\Pi_{\mathcal{C}}(m) \stackrel{\text{def}}{=} \max \{|\Pi_{\mathcal{C}}(\bar{S})| \mid \bar{S} \subset \bar{X}, |\bar{S}| = m\},$$

(maximum se bere přes všechny množiny $\bar{S} \subset \bar{X}$ mohutnosti m).

Zřejmě $\Pi_{\mathcal{C}}(\bar{S})$ je systém všech podmnožin množiny \bar{S} , které lze oddělit od jejich doplňku v \bar{S} pomocí konceptu ležícího ve třídě konceptů \mathcal{C} . Číslo $\Pi_{\mathcal{C}}(m)$ vyjadřuje mohutnost maximálního takového systému $\Pi_{\mathcal{C}}(\bar{S})$, kde množina \bar{S} je konečná m prvková podmnožina \bar{X} . Mezi těmito čísly a VC-dimenzí platí evidentní vztah

$$\text{VC}_{\dim}(\mathcal{C}) = \max \{m \mid \Pi_{\mathcal{C}}(m) = 2^m\}.$$

Definice 4.2.3 *Pro všechna $d \geq 0$ a $m \geq 0$ položme $\Phi_{d,m} \stackrel{\text{def}}{=} \sum_{i=0}^d \binom{m}{i}$, je-li $m \geq d$ a $\Phi_{d,m} \stackrel{\text{def}}{=} 2^m$ je-li $m < d$.*

Tvrzení obsažená v následující lemmě budou nezbytná k důkazu horních odhadů délky vzorků, nutných pro učení algoritmů s danou přesností (ve smyslu nalezení hypotézy mající "malou" symetrickou diferencii s daným konceptem při požadované pravděpodobnosti nalezení takovéto hypotézy), které je analyzované v následujících oddílech.

Lemma 4.2.4 *Je-li \mathcal{C} libovolná třída konceptů nad konečnou množinou \bar{X} , tak platí*

1. *Je-li $\text{VC}_{\dim}(\mathcal{C}) = d$, tak $\Pi_{\mathcal{C}}(m) \leq \Phi_{d,m}$ pro všechna $m \geq 0$.*
2. *$\Phi_{d,m} \leq m^d + 1$ pro $d \geq 0$ a $\Phi_{d,m} \leq m^d$ pro $d \geq 0$ a $m \geq 2$.*
3. *$\Phi_{d,m} \leq 2 \frac{m^d}{d!} m + 1 \leq \left(\frac{em}{d}\right)^d$ pro $m \geq d \geq 1$.*

■ *Důkaz:*

DŮKAZ JE UVEDEN V ORIGINALNÍ TIŠTĚNÉ VERZI.

■

Důsledek 4.2.5 *Nechť \bar{X} je konečná množina, $\mathcal{C} \subset 2^{\bar{X}}$ a $\text{VC}_{\dim}(\mathcal{C}) > 0$. Potom*

$$\text{VC}_{\dim}(\mathcal{C}) > \frac{\log_e(|\bar{\mathcal{C}}|)}{1 + \log_e(|\bar{X}|)}.$$

■ *Důkaz:*

DŮKAZ JE UVEDEN V ORIGINALNÍ TIŠTĚNÉ VERZI.

■

1. Dokažte druhou část lemmy 4.2.4.
2. Nechť $\vec{x} \in \mathbb{R}^n$ a $G_{n,\vec{x}}$ je systém všech poloprostorů v \mathbb{R}^n , které neobsahují bod \vec{x} ,

$$G_{n,\vec{x}} \stackrel{\text{def}}{=} \{ \bar{c} \subset \mathbb{R}^n \mid (\exists \vec{w} \in \mathbb{R}^n, \alpha \in \mathbb{R}) (\bar{c} = \{ \vec{z} \in \mathbb{R}^n \mid \langle \vec{w} \mid \vec{z} \rangle \leq \alpha \text{ a } \vec{x} \notin \bar{c} \}) \} \quad (4.1)$$

Dokažte, že $\text{VC}_{\dim}(G_{n,\vec{x}}) = n$.

3. Ukažte, že $\Pi_{G_{n,\vec{x}}}(m) = \Phi_{n,m}$.
4. Nechť A je třída konceptů definovaná nad \bar{X} a $\text{VC}_{\dim}(A)$ je konečná. Nechť B je třída konceptů obsahující doplňky konceptů z A do množiny \bar{X} . Dokažte, že potom platí $\text{VC}_{\dim}(A) = \text{VC}_{\dim}(B)$.
5. Nechť C je systém všech uzavřených intervalů na \mathbb{R} . Ukažte, že

$$\Pi_C(m) = 1 + m + \frac{1}{2}m(m-1).$$

6. Řekneme, že třída konceptů C je LINEÁRNĚ USPOŘÁDANÁ, právě když obsahuje alespoň dva koncepty a pro libovolné $\bar{a}, \bar{b} \in C$ je buď $\bar{a} \subset \bar{b}$ nebo $\bar{b} \subset \bar{a}$. Potom platí, že $\text{VC}_{\dim}(C) = 1$. Dokažte.
7. Ukažte, že systém všech přímek v \mathbb{R}^n má VC-dimenzi rovnu 1.

4.3 Vapnik-Chervonenkova dimenze neuronových sítí

Pro libovolný systém funkcí můžeme definovat VC-dimenzi těchto funkcí jako VC-dimenzi všech množin, které jsou vzorem pro intervaly $(-\infty, \alpha)$ při zobrazení nějakou funkcí z uvažovaného systému funkcí. Obdobně lze definovat VC-dimenzi obecných neuronových sítí či booleovských obvodů, budeme-li na ně pohlížet jako na funkce zobrazující hodnoty na vstupních uzlech do množiny reálných čísel. Pro účely další analýzy VC-dimenze neuronových sítí zdefinujeme obecnější model neuronových sítí, než je booleovský obvod.

Definice 4.3.1 DOPŘEDNÝ OBVOD je acyklický orientovaný graf s hranovým a vrcholovým ohodnocením a s pouze jediným vrcholem, ze kterého nevede hrana. Vrchol tohoto grafu nazýváme VSTUPNÍ VRCHOL, nevede-li do něj žádná hrana, VÝSTUPNÍ VRCHOL, nevede-li z něj žádná hrana, a VNITŘNÍ VRCHOL, není-li ani vstupní ani výstupní. Hodnoty hranového ohodnocení jsou obecně reálná čísla, ohodnocení hrany z vrcholu i do vrcholu j nazveme VÁHOU HRANY $i \rightarrow j$ a budeme značit $w_{i \rightarrow j}$. Pro daný vrchol j nechť n_j je počet hran vedoucích do vrcholu j a $\vec{w}_j \in \mathbb{R}^{n_j}$ je vektor, jehož složky jsou rovny ohodnocení těchto hran. Ohodnocení vrcholů jsou reálná čísla, která pro vrchol j budeme značit jako v_j . Hodnoty v_j u vstupních vrcholů nazýváme vstupními hodnotami, ohodnocení v_j u výstupního vrcholu výstupní hodnotou. Nechť $\vec{s}_j \in \mathbb{R}^{n_j}$ je vektor tvořený hodnotami ohodnocení vrcholů, z nichž vede hrana do vrcholu j . Dále, ohodnocení vnitřního a výstupního vrcholu j splňuje rovnost

$$v_j \stackrel{\text{def}}{=} \widetilde{f}_j(\vec{w}_j, \vec{s}_j),$$

kde \widetilde{f}_j jsou libovolné pevně dané funkce, $\widetilde{f}_j: \mathbb{R}^{n_j} \times \mathbb{R}^{n_j} \rightarrow \mathbb{R}$.

V následujícím textu odvodíme základní horní odhad VC-dimenze dopředného obvodu využívající vlastostí funkcí \tilde{f}_j v tomto obvodu. K tomu bude zapotřebí pracovat s částmi dopředného obvodu, které budeme definovat prostřednictvím tzv. vlastního očíslování vrcholů obvodu.

Definice 4.3.2 Předpokládejme, že G je acyklický orientovaný graf a že uzly grafu G , které mají alespoň jednu vstupní hranu, jsou očíslovány přirozenými čísly $i \in \{1, \dots, z\}$. Potom toto očíslování nazveme VLASTNÍ OČÍSLOVÁNÍ, právě když pro všechny očíslované uzly platí, že hrany z daného uzlu vedou pouze do uzlů, jejichž očíslování je vyšší.

Vlastní očíslování acyklického orientovaného grafu dostaneme například takovým způsobem, že po odstranění všech uzlů nemajících vstupní hranu, v takto vzniklém grafu nejdříve očísloujeme postupně všechny uzly nemající vstupní hranu (to lze, neboť graf je acyklický), tyto uzly z grafu odstraníme, a tento postup zopakujeme na takto vzniklém grafu, přičemž číslujeme postupně dalšími čísly.

Definice 4.3.3 Necht D je dopředný obvod, který má $(z - 1)$ vnitřních vrcholů. Potom pro každé $l \in \{1, \dots, z\}$ definujme třídu konceptů

$$C_l \stackrel{\text{def}}{=} \left\{ \bar{c} \subset \mathfrak{R}^{n_l} \mid (\exists \bar{\omega} \in \mathfrak{R}^{n_l}) \left(\bar{c} = \left\{ \bar{s} \in \mathfrak{R}^{n_l} \mid \tilde{f}_l(\bar{\omega}, \bar{s}) \leq 0 \right\} \right) \right\}.$$

Každou třídu konceptů C_l nazveme LOKÁLNÍ TŘÍDOU KONCEPTŮ.

Necht dále $\bar{\Delta}$ je množina všech možných hranových ohodnocení obvodu D a hodnotu výstupního vrcholu z označme v závislosti na $\bar{\mathbf{h}} \in \bar{\Delta}$ a vstupních hodnotách $\bar{\mathbf{x}} \in \mathfrak{R}^n$ obvodu D jako $v_{z, \bar{\mathbf{h}}, \bar{\mathbf{x}}}$. Potom definujme TŘÍDU KONCEPTŮ DOPŘEDNÉHO OBVODU jako systém množin

$$C_D \stackrel{\text{def}}{=} \left\{ \bar{c} \subset \mathfrak{R}^n \mid (\exists \bar{\mathbf{h}} \in \bar{\Delta}) \left(\bar{c} = \left\{ \bar{\mathbf{x}} \in \mathfrak{R}^n \mid v_{z, \bar{\mathbf{h}}, \bar{\mathbf{x}}} \leq 0 \right\} \right) \right\}.$$

VC-DIMENZÍ DOPŘEDNÉHO OBVODU D budeme rozumět číslo $\text{VC}_{\dim}(C_D)$.

Jinými slovy, pro dané l je C_l je systém všech podmnožin $\bar{c} \in \mathfrak{R}^{n_l}$ takových, že existuje váhový vektor $v \in \mathfrak{R}^{n_l}$ tak, že množina \bar{c} je vzorem intervalu $(-\infty, 0)$ při zobrazení $\tilde{f}_l(\bar{\omega}, \bar{s})$. Obdobným způsobem je definována třída konceptů C_D .

Definice 4.3.4 Necht $\bar{V} \subset \mathfrak{R}^n$, $\bar{V} \stackrel{\text{def}}{=} \{\bar{\mathbf{x}}_1, \dots, \bar{\mathbf{x}}_m\}$ a D je dopředný obvod dimenze n s $n+z$ vrcholy a s vlastním očíslováním vrcholů. Potom řekneme, že dvě hranová ohodnocení ω_1 a ω_2 jsou l -ROZLIŠITELNÁ POMOCÍ \bar{V} , $l \in \{1, \dots, z\}$, právě když existuje $\bar{\mathbf{x}}_i \in \bar{V}$ a existuje $j \in \{1, \dots, l\}$ taková, že hodnota uzlu j se liší, je-li hranové ohodnocení obvodu rovno hodnotám ω_1 nebo hodnotám ω_2 . Dále označme symbolem $S_{l, \bar{V}}$ maximální počet všech hranových ohodnocení, která jsou vzájemně po dvou l -rozlišitelná pomocí \bar{V} .

Tvrzení 4.3.1 Necht m je libovolné kladné celé číslo, $\bar{V} \subset \mathfrak{R}^n$, $\bar{V} \stackrel{\text{def}}{=} \{\bar{\mathbf{x}}_1, \dots, \bar{\mathbf{x}}_m\}$. Potom pro všechna $l \in \{1, \dots, z\}$ platí odhad

$$S_{l, \bar{V}} \leq \Pi_{C_1}(m) \cdot \Pi_{C_2}(m) \cdot \dots \cdot \Pi_{C_l}(m).$$

Navíc platí

$$\Pi_{C_D}(m) \leq S_{z, \bar{V}} \leq \Pi_{C_1}(m) \cdot \Pi_{C_2}(m) \cdot \dots \cdot \Pi_{C_z}(m).$$

■ *Důkaz:*

DŮKAZ JE UVEDEN V ORIGINÁLNÍ TIŠTĚNÉ VERZI.

Až do tohoto místa výkladu jsme nikde nevyužili vlastností přechodových funkcí \tilde{f}_i dopředného obvodu, podstatná byla pouze skutečnost, že graf dopředného obvodu je orientovaný a acyklický. Tvzení 4.3.1 tedy zřejmě platí pro dosti širokou třídu neuronových sítí a dává nám návod, jak postupovat při odvozování horního odhadu VC-dimenze takovýchto obecnějších neuronových sítí, neboť umožňuje rozložit nalezení takového odhadu na odhady VC-dimenze sítí obsahujících pouze jednotlivé uzly sítě původní. Budeme ilustrovat tento postup na jednoduchém příkladě tzv. lineárního sigmoidálního obvodu, který v sobě zahrnuje i booleovské obvody.

Definice 4.3.5 *Dopředný obvod L nazveme LINEÁRNÍM SIGMOIDÁLNÍM OBVODEM jestliže pro každý nevstupní uzel obvodu L jemu příslušná funkce \tilde{f}_j je tvaru*

$$\tilde{f}_j(\vec{s}_j) \stackrel{\text{def}}{=} \tilde{\sigma}_j(\langle \vec{s}_j | \vec{w}_j \rangle - t_j), \quad (4.2)$$

kde $\tilde{\sigma}_j$ jsou libovolné, pevně dané sigmoidální funkce, $t_j \in \mathfrak{R}$, a \vec{s}_j, \vec{w}_j mají stejný význam jako v definici 4.3.1.

Určíme nyní horní odhad VC-dimenze lineárního sigmoidálního obvodu. Jelikož každý takovýto obvod je konkrétní realizací dopředného obvodu zmíněného v předchozím, všechny potřebné definice a tvrzení ohledně dopředného obvodu jsou platné i pro tyto obvody. Přístupme proto přímo k tvrzení, jež se týká lineárních sigmoidálních obvodů.

Tvrzení 4.3.2 *Nechť L je lineární sigmoidální obvod, w je počet jeho hran, z je počet jeho vrcholů, které nejsou vstupní, $q \stackrel{\text{def}}{=} w + z$. Potom pro každé $m \geq \max\{n_i | i \in \{1, \dots, z\}\}$ platí*

$$\Pi_{C_L}(m) \leq \left(\frac{ezm}{q}\right)^q \quad (4.3)$$

a dále platí odhad

$$\text{VC}_{\dim}(C_L) \leq 2q \log_2(ez). \quad (4.4)$$

■ *Důkaz:*

DŮKAZ JE UVEDEN V ORIGINÁLNÍ TIŠTĚNÉ VERZI.

Cvičení 4.3.0 VC-dimenze neuronových sítí

1. Předpokládejme, že v definičním vztahu 4.2 jsou všechna t_j rovna nule. Pro takto modifikovaný lineární sigmoidální obvod odvoďte odhad 4.4 na VC-dimenzi takového obvodu. Vysvětlete souvislost tohoto odhadu s vlastnostmi třídy konceptů $\mathbf{G}_{n, \vec{x}}$ definované v 4.1. Objasněte souvislost s existencí či neexistencí prahových hodnot t_j v lineárním sigmoidálním obvodu.

2. Mějme dopředný obvod, jehož přechodové funkce \tilde{f}_j jsou tvaru

$$\tilde{f}_j(\vec{s}_j) \stackrel{\text{def}}{=} \max\left\{ |(\vec{s}_j)_k - (\vec{w}_j)_k| \mid k \in \{1, \dots, n_j\} \right\} - t_j.$$

Spočtěte VC-dimenzi tohoto obvodu.

3. Mějme dopředný obvod, jehož přechodové funkce \widetilde{f}_j jsou tvaru

$$\widetilde{f}_j(\vec{\mathbf{s}}_j) \stackrel{\text{def}}{=} \min \left\{ \widetilde{\text{sgn}} \left((\vec{\mathbf{s}}_j)_k - (\vec{\mathbf{w}}_j)_k \right) \mid k \in \{1, \dots, n_j\} \right\}.$$

Spočtěte VC-dimenzi tohoto obvodu.

Kapitola 5

Teorie učení a neuronové sítě

Saganovská variace

TVRDIT, ŽE MATEMATICKÉ VZORCE JSOU JEN ČÁRY NA PAPIŘE,
JE STEJNÉ, JAKO TVRDIT,
ŽE SHAKESPEAROVY HRY JSOU JEN SLOVA.

Tuto kapitolu věnujeme výkladu základů tzv. teorie učení, která analyzuje možnosti extrakce vlastností množin na základě znalosti konečného počtu jejich reprezentantů. Tato teorie poskytuje základní pohled na možnost generalizace vjemů a vstupů, přicházejících do jakéhokoliv systému, který je schopen na základě dřívějších znalostí generovat hypotézy týkající se vlastností vjemů a vstupů následujících. Tato kapitola se tedy netýká pouze neuronových sítí, ale obecně analyzuje problém učení. Náš výklad bude založen zejména na teorii rozpracované v citaci ([BEHW89]). Ukážeme, že VC-dimenze třídy konceptů a VC-dimenze neuronových sítí jsou klíčové pojmy pro studium obecného problému učení. Dále zavedeme různá kritéria efektivity učících algoritmů, co se týče počtu reprezentantů nutných k realizaci daného algoritmu s požadovanou kvalitou výstupu. Ukážeme základní odhady počtu reprezentantů vzhledem ke složitosti konceptu a vzhledem k dimenzi reprezentantů konceptu pro polynomiálně složitě algoritmy učení. Vyložené pojmy plynule navazují na problematiku neuronových sítí, prezentovanou v předchozích kapitolách.

5.1 Deskripce učících algoritmů

Uvažujme problém, kdy máme zadánu jakousi množinu $\bar{A} \in \mathfrak{R}^n$, o které apriori víme, že \bar{A} je z geometrického hlediska koule. Předpokládejme, že máme danu posloupnost bodů $\vec{x}_i \in \mathfrak{R}^n$ a o každém jednom z nich je známo, zda leží v množině \bar{A} nebo ne. Naším úkolem je rozhodnout, jak vypadá množina \bar{A} , v tomto konkrétním případě najít střed a poloměr příslušné koule. Je zřejmé, že naše rozhodnutí je velmi závislé na vlastnostech posloupnosti \vec{x}_i . Předpokládejme, že na základě nějakého algoritmu (deterministického či nedeterministického) jsme se rozhodli pro nějakou kouli v \mathfrak{R}^n , o které budeme tvrdit, že je rovna množině \bar{A} . Je přirozené očekávat, že v této množině by měly být obsaženy všechny body $\vec{x}_j \in \bar{A}$ a naopak by v ní nemělo být žádné \vec{x}_j , které v \bar{A} neleží. Tomuto požadavku budeme říkat konzistentnost a o množině, o které se domníváme, že je rovna kouli \bar{A} , budeme hovořit jako o hypotéze (rovnost námi zvolené koule a \bar{A} je pouze naše hypotéza, založená na vlastnostech použitého algoritmu a na posloupnosti bodů \vec{x}_i).

Označme námi zvolenou kouli jako \bar{K} . Pak je zřejmé, že podmínka konzistentnosti musí být splněna v případě, kdy symetrická diference množin \bar{K} a \bar{A} bude prázdná množina, je tedy přirozené požadovat, aby toto platilo pro generovanou hypotézu, nebo aby to platilo alespoň s vysokou mírou pravděpodobnosti. Výše diskutované pojmy jsou obsahem následujících definic.

Definice 5.1.1 *Nechť \mathcal{C} je třída konceptů a pro systém množin \mathbf{H} platí $\mathcal{C} \subset \mathbf{H} \subset 2^{\bar{X}}$. Potom \mathbf{H} nazveme TŘÍDOU HYPOTÉZ pro třídu konceptů \mathcal{C} .*

Definice 5.1.2 *Nechť $\tilde{x} \stackrel{\text{def}}{=} \{x_1, \dots, x_m\}$, $x_i \in \bar{X}$, $i \in \{1, \dots, m\}$, $\tilde{z} \in \{-1, +1\}^m$ a necht' $\bar{c} \subset \bar{X}$. Potom uspořádanou dvojici*

$$(\tilde{x}, \tilde{z})$$

nazveme VZOREK KONCEPTU \bar{c} DÉLKY m právě když

$$(\forall i \in \{1, \dots, m\}) ((x_i \in \bar{c}) \Leftrightarrow (\tilde{z}_i = 1)).$$

Pro třídu konceptů \mathcal{C} definujme PROSTOR VZORŮ TŘÍDY KONCEPTŮ jako

$$\bar{S}_{\mathcal{C}} \stackrel{\text{def}}{=} \bigcup_{m \geq 1} \left\{ \bigcup_{\bar{c} \in \mathcal{C}} \{(\tilde{x}, \tilde{z}) \mid (\tilde{x}, \tilde{z}) \text{ je vzorek délky } m \text{ konceptu } \bar{c}\} \right\}.$$

Množina $\bar{b} \subset \bar{X}$ je KONZISTENTNÍ MNOŽINA se vzorkem (\tilde{x}, \tilde{z}) jestliže pro všechna $i \in \{1, \dots, m\}$, platí $(x_i \in \bar{b} \Leftrightarrow \tilde{z}_i = 1)$.

Pro další výklad musíme mít k dispozici nějakou míru odlišnosti dvou množin. Tuto míru odlišnosti definujeme takto:

Definice 5.1.3 *Nechť \bar{X} je pravděpodobnostní prostor s hustotou pravděpodobnosti \tilde{P} . Necht' \mathcal{C} je třída konceptů definovaná nad prostorem \bar{X} a necht' je \mathbf{H} třída hypotéz pro \mathcal{C} . Potom pro každé $\bar{c} \in \mathcal{C}$ a $\bar{h} \in \mathbf{H}$ číslo*

$$e_{\tilde{P}}(\bar{c}, \bar{h}) \stackrel{\text{def}}{=} \text{Prob}_{\tilde{P}}(\bar{h} \Delta \bar{c})$$

nazveme CHYBA HYPOTÉZY \bar{h} VZHLEDEM KE KONCEPTU \bar{c} a hustotě pravděpodobnosti \tilde{P} (symbol $\bar{h} \Delta \bar{c}$ označuje symetrickou diferenci množin \bar{h} a \bar{c}).

Protože množina \bar{A} byla popsána pouze posloupností svých prvků (plus apriorní informací o tvaru) a pouze na základě těchto údajů jsme se měli rozhodnout, která koule je rovna množině \bar{A} , je náš uvažovaný učící algoritmus ve své podstatě zobrazení mezi množinou všech vzorků množiny \bar{A} do množiny všech koulí v \mathbb{R}^n . Od smysluplného učícího algoritmu nebudeme požadovat přesnou hypotézu pro všechny možné úlohy tohoto typu, nýbrž budeme pouze požadovat, aby učící algoritmus ve většině případů produkoval "uspokojivou" hypotézu (např. hypotézu konzistentní se zadaným vzorkem). Kvantifikace tohoto požadavku je obsahem následující definice.

Definice 5.1.4 Mějme definovanou funkci $\tilde{m}(\epsilon, \delta)$ zobrazující množinu $(0, 1) \times (0, 1)$ do množiny přirozených čísel a \mathbf{H} nechť je třída hypotéz pro třídu konceptů \mathbf{C} definovanou nad množinou \bar{X} . Pak \tilde{P} -UČICÍ ALGORITMUS SLOŽITOSTI $\tilde{m}(\epsilon, \delta)$ třídy konceptů \mathbf{C} je každé zobrazení $\widetilde{A}^* : \bar{S}_{\mathbf{C}} \rightarrow \mathbf{H}$ takové, že pro všechna $\bar{c} \in \mathbf{C}$, pro všechna $0 < \epsilon < 1$ a $0 < \delta < 1$ a pro libovolnou pravděpodobnost \tilde{P} definovanou na \bar{X} je pravděpodobnost množiny

$$\{\tilde{x} \in \bar{X}^m \mid (\tilde{x}, \tilde{z}) \text{ je vzorek } \bar{c} \text{ a } e_{\tilde{P}}(\bar{c}, \widetilde{A}^*((\tilde{x}, \tilde{z}))) \geq \epsilon\}$$

menší než číslo δ . Jestliže takovýto učicí algoritmus existuje, říkáme, že \mathbf{C} JE UNIFORMĚ NAUČITELNÁ podle třídy hypotéz \mathbf{H} . Každý takovýto učicí algoritmus nazveme (ϵ, δ) -UČICÍM ALGORITMEM.

POZNÁMKA:

Tento model výpočetní složitosti se v literatuře označuje jako PAC učení (anglický termín: Probably Approximately Correct).

Je zřejmé, že každou množinu lze velmi přesně popsat velkým počtem jejích prvků. Algoritmy využívající vzorků velké délky mohou zřejmě produkovat mnohem věrohodnější hypotézy. Spokojíme-li se s předem danou přesností generování hypotéz, popsanou čísly ϵ a δ , je přirozené se ptát, jak dlouhý vzorek pro danou třídu konceptů musíme použít, abychom měli tuto přesnost zaručenu. Proto zavedeme pojem vzorové složitosti učících algoritmů následovně:

Definice 5.1.5 Minimální hodnotu $\tilde{m}(\epsilon, \delta)$, pro kterou je A^* učicím algoritmem, nazveme VZOROVOU SLOŽITOSTÍ učícího algoritmu A^* .

Tato složitost je definována pouze vzhledem k délce vzorku použitého učícím algoritmem, otázka vnitřní výpočetní složitosti vlastního algoritmu zde není brána v úvahu! Nejhrubší odhad vzorové složitosti dává následující tvrzení, založené pouze na vlastnostech pravděpodobnosti (poskytuje základní odhad počtu vzorů, které musíme pro učicí algoritmus produkující konzistentní hypotézu použít).

Tvrzení 5.1.1 Nechť \mathbf{C} je třída konceptů nad konečnou množinou \bar{X} a $\mathbf{H} = \mathbf{C}$. Potom pro každý učicí algoritmus vyžadující

$$\frac{1}{\epsilon} \log_e \left(\frac{|\mathbf{C}|}{\delta} \right) \quad (5.1)$$

dotazů a produkující pro daný koncept $\bar{c} \in \mathbf{C}$ konzistentní hypotézu a pro každou hustotu pravděpodobnosti \tilde{P} definovanou na \bar{X} platí

$$\text{Prob}_{\tilde{P}} \left(e_{\tilde{P}}(\bar{c}, \widetilde{A}^*((\tilde{x}, \tilde{z}))) \geq \epsilon \right) < \delta.$$

■ *Důkaz:*

DŮKAZ JE UVEDEN V ORIGINALNÍ TIŠTĚNÉ VERZI.

■

Poznamenejme, že odhad 5.1 není závislý na konkrétní míře pravděpodobnosti \tilde{P} . To lze vysvětlit tím, že vliv různé pravděpodobnosti prvků v \bar{X} je kompenzován pravděpodobností výběru tohoto prvku podle hustoty pravděpodobnosti \tilde{P} . Jinak řečeno, prvky z \bar{X} s velkou pravděpodobností, které v případě příslušnosti k symetrické diferencii $\bar{c} \Delta \bar{h}$

přispívají významně k velikosti chyby, jsou v průměrném vzorku častější, což vzhledem ke skutečnosti, že algoritmus produkuje vždy konzistentní hypotézu, eliminuje jejich vliv. Proto je odhad 5.1 nezávislý na pravděpodobnosti \tilde{P} .

V následujícím výkladu ukážeme, že tento odhad lze velmi výrazně zlepšit, právě na základě aplikace pojmu VC-dimenze.

5.2 Odhady složitosti učení a VC-dimenze

V této podkapitole se budeme zabývat dolními a horními odhady počtu vzorů, potřebných pro učící algoritmy složitosti $\tilde{m}(\epsilon, \delta)$.

Jeden takovýto odhad na postačující množství dotazů je dán v tvrzení 5.1.1. V této kapitole ukážeme, že dominantní roli při těchto odhadech hraje právě dříve studovaný pojem VC-dimenze. Na základě tohoto pojmu lze vyslovit lepší odhady na nezbytnou délku vzorku použitého pro generování hypotézy, aproximující daný koncept.

Nejdříve se budeme zabývat zmíněnými odhady pro pevně daný systém konceptů a hypotéz. V dalším budeme uvažovat případy, kdy dané koncepty jsou rozlišeny podle dimenze prostoru, jehož podmnožinami jsou, a budeme analyzovat vzorovou složitost učících algoritmů s ohledem na dimenzi vektorových prostorů, obsahujících jako své podmnožiny dané koncepty. Přirozeným požadavkem je postulovat kritérium, které zaručí, že délka vzorku $\tilde{m}(\epsilon, \delta)$ učícího algoritmu bude zhora omezena polynomem v dimenzi prostoru.

Závěrem provedeme obdobnou analýzu vzhledem k délce slov, která popisují koncepty z daného systému konceptů a která ve své podstatě vyjadřují svou délkou popisnou složitost konceptů.

5.2.1 Odhad počtu vzorů

V této části dokážeme sérii lemm a tvrzení, které budou nutné k důkazu tvrzení 5.2.6, udávajícího dolní a horní odhad vzorové složitosti pro (ϵ, δ) -algoritmy.

V průběhu celé této kapitoly a kapitol následujících budeme předpokládat, že každá uvažovaná třída konceptů, stejně jako každá uvažovaná třída hypotéz, bude tvořena pouze Borelovskými množinami. Pro úplnost připomeňme pojem Borelovských množin ([Jar55]). Předpokládejme, že Ω je neprázdný systém množin, který je σ -aditivní, to jest pro libovolnou posloupnost množin $A_n \in \Omega$ platí, že $\bigcup_{n=1}^{\infty} A_n \in \Omega$. Platí-li navíc, že pro libovolné $\bar{A}, \bar{B} \in \Omega$ je též i $\bar{A} - \bar{B} \in \Omega$, nazveme systém množin Ω σ -aditivním okruhem. Dá se ukázat, že pro každý neprázdný systém množin Ω existuje právě jeden nejmenší σ -aditivní okruh obsahující Ω (je to právě průnik všech σ -aditivních okruhů, obsahujících systém Ω). Tento minimální σ -aditivní okruh nazveme Borelovským okruhem nad systémem Ω . Budiž \mathcal{M} libovolný metrický prostor a Θ nechť je Borelovský okruh nad systémem všech otevřených množin v \mathcal{M} . Potom prvky tohoto okruhu nazveme BORELOVSKÝMI MNOŽINAMI prostoru \mathcal{M} .

Kromě požadavku ohledně Borelovskosti množin, které tvoří jednotlivé koncepty, bude nezbytné pro další důkazy detailněji charakterizovat uvažované koncepty a hypotézy. Tuto nutnou charakterizaci lze vhodně postihnout pojmem ϵ -transversála.

Definice 5.2.1 Pro libovolné $\mathbb{R} \subset 2^{\bar{X}}$ a pro libovolnou hustotu pravděpodobnosti \tilde{P} definovanou na \bar{X} a pro libovolné $\epsilon > 0$ definujme $\mathbb{R}_{\tilde{P}, \epsilon} \stackrel{\text{def}}{=} \{\bar{r} \in \mathbb{R} \mid \text{Prob}_{\tilde{P}}(\bar{r}) > \epsilon\}$. $\mathbb{R}_{\tilde{P}, \epsilon}$ Potom $\mathbb{N} \subset 2^{\bar{X}}$ nazveme ϵ -TRANSVERSÁLA, jestliže \mathbb{N} obsahuje bod z každé množiny patřící do systému $\mathbb{R}_{\tilde{P}, \epsilon}$.

Příklad 5.2.1 *Nechť $\bar{a} \stackrel{\text{def}}{=} \langle 0, 1 \rangle$ a \mathcal{C} je systém všech uzavřených intervalů na \bar{a} . Potom množina všech bodů tvaru ϵk , $1 \leq k \leq \frac{1}{\epsilon}$ tvoří ϵ -transversálu \bar{a} pro libovolnou hodnotu $\epsilon > 0$ a pro libovolnou hustotu pravděpodobnosti definovanou na \bar{a} . Naopak je-li \mathcal{C} tvořeno všemi otevřenými podmnožinami intervalu \bar{a} , pak pro rovnoměrnou hustotu pravděpodobnosti neexistuje pro žádné ϵ konečná ϵ -transversála.*

Budeme se zajímat o pravděpodobnost výběru ϵ -transversály systému \mathbf{R} na základě náhodného výběru bodů z \bar{X} . Konkrétně nás bude zajímat pravděpodobnost události popsané v následující definici.

Definice 5.2.2 *Pro každé $m \geq 1$ a $\epsilon > 0$ nechť $\bar{Q}_{m,\epsilon}$ označuje množinu všech $\bar{x} \in \bar{X}^m$ takovou, že navzájem různé prvky \bar{x} netvoří ϵ -transversálu pro \mathbf{R} s ohledem na hustotu \tilde{P} . Nechť potom \bar{J}_ϵ^{2m} označuje množinu všech $\bar{x}\bar{y} \in \bar{X}^{2m}$, kde $\bar{x}, \bar{y} \in \bar{X}^m$ tak, že existuje $\bar{r} \in \mathcal{R}_{\tilde{P},\epsilon}$, pro kterou je $\bar{x} \cap \bar{r} = \emptyset$ a $\left| \left\{ i \mid \bar{y}_i \in \bar{r}, i \in \{1, \dots, m\} \right\} \right| \geq \frac{\epsilon m}{2}$.*

V důkazech pravděpodobnostních vlastností učících algoritmů nás budou zajímat zejména vlastnosti systému množin $\mathbf{R} \stackrel{\text{def}}{=} \{ \bar{h} \Delta \bar{c} \mid \bar{h} \in \mathbf{H} \}$, kde \bar{c} je nějaký pevný koncept $\bar{c} \subset \bar{X}$ a \mathbf{H} je třída hypotéz pro tento koncept. Význam pojmu ϵ -transversála je ten, že jestliže vzorek konceptu \bar{c} je současně ϵ -transversálou pro \mathbf{R} , pak obsahuje protipříklad pro každou hypotézu, jejíž chyba vzhledem k cílovému konceptu \bar{c} je větší než ϵ . Abychom mohli zkoumat vlastnosti systému \mathbf{R} , vezmeme v úvahu třídy hypotéz splňující následující kritérium uvedené v definici.

Definice 5.2.3 *Třída hypotéz \mathbf{H} je DOBRĚ UTVOŘENÁ (anglický termín: well-behaved) jestliže množiny $\bar{Q}_{m,\epsilon}$ a \bar{J}_ϵ^{2m} jsou měřitelné pro libovolnou hustotu pravděpodobnosti \tilde{P} , libovolné $m \geq 1$, $\epsilon > 0$ a libovolný systém množin $\mathbf{R} \stackrel{\text{def}}{=} \{ \bar{h} \Delta \bar{c} \mid \bar{h} \in \mathbf{H} \}$, kde \bar{c} je libovolná Borelovská množina.*

Výrazná většina tříd konceptů a tříd hypotéz běžně používaných je dobře utvořená. Jednou z možností verifikace této vlastnosti je ověření univerzální separability.

Definice 5.2.4 *Třída hypotéz $\mathbf{H} \subset 2^{\bar{X}}$ je UNIVERZÁLNĚ SEPARABILNÍ, existuje-li spočetná podmnožina \mathbf{T} třídy \mathbf{H} tak, že pro všechny $\bar{h} \in \mathbf{H}$ existuje posloupnost $\{\bar{h}_i\}_1^\infty$ množin z \mathbf{T} taková, že*

$$\left(\forall x \in \bar{X} \right) \left(\exists n \geq 1 \right) \left(\left(\forall i \geq n \right) \left(x \in \bar{h}_i \text{ právě když } x \in \bar{h} \right) \right).$$

Platí implikace vyjádřená v následujícím tvrzení.

Tvrzení 5.2.1 *Je-li \mathbf{H} univerzálně separabilní, pak \mathbf{H} je dobře utvořená.*

■ *Důkaz:*

DŮKAZ JE UVEDEN V ORIGINALNÍ TIŠTĚNÉ VERZI.

■

Dokážeme nyní následující technickou lemmu poskytující základní nerovnosti nezbytné při odhadech vzorové složitosti učících algoritmů (zavedme konvenci, že \tilde{P}^m je hustota pravděpodobnosti definovaná na kartézském součinu \bar{X}^m , odvozená od hustoty pravděpodobnosti \tilde{P} definované na \bar{X}).

Lemma 5.2.2 *Nechť \mathbf{R} je neprázdná třída konceptů na \bar{X} a \tilde{P} je hustota pravděpodobnosti na \bar{X} , pro kterou $\bar{Q}_{m,\epsilon}$ a \bar{J}_ϵ^{2m} jsou měřitelné pro libovolné $m \geq 1$ a $\epsilon > 0$. Potom*

1. *pro každé $\epsilon > 0$ a $m \geq \frac{m}{\epsilon}$ platí*

$$\text{Prob}_{\tilde{P}^m}(\bar{Q}_{m,\epsilon}) < 2\text{Prob}_{\tilde{P}^{2m}}(\bar{J}_\epsilon^{2m}),$$

2. *pro každé $\epsilon > 0$ a $m \geq 1$ platí*

$$\text{Prob}_{\tilde{P}^{2m}}(\bar{J}_\epsilon^{2m}) \leq \Pi_{\mathbf{R}}(2m) 2^{-\frac{\epsilon m}{2}}.$$

■ *Důkaz:*

DŮKAZ JE UVEDEN V ORIGINÁLNÍ TIŠTĚNÉ VERZI.

Předešlá lemma umožňuje získat horní odhad pravděpodobnosti události, kdy na základě náhodného vzorku nezískáme ϵ -transversálu pro třídu hypotéz \mathbf{R} . K získání odhadu budeme potřebovat lemmu následující.

Lemma 5.2.3 *Pro každé $m \geq 1$, $\bar{c} \subset \bar{X}$ a $\mathbf{H} \subset 2^{\bar{X}}$ je $\Pi_{\mathbf{H}}(m) = \Pi_{\mathbf{R}}(m)$, kde je $\mathbf{R} \stackrel{\text{def}}{=} \{\bar{h} \Delta \bar{c} \mid \bar{h} \in \mathbf{H}\}$.*

■ *Důkaz:*

DŮKAZ JE UVEDEN V ORIGINÁLNÍ TIŠTĚNÉ VERZI.

Nyní vyslovme tvrzení, zhora odhadující pravděpodobnost existence hypotézy konzistentní se vzorkem konceptu \bar{c} a s chybou menší než ϵ .

Tvrzení 5.2.4 *Nechť \mathbf{H} je neprázdná dobře utvořená třída hypotéz nad \bar{X} , \tilde{P} je hustota pravděpodobnosti na \bar{X} a $\bar{c} \subset \bar{X}$ je libovolná Borelovská množina. Potom pro libovolnou pevně zvolenou $\epsilon > 0$, $m \geq 1$ a pro libovolný vzorek (\tilde{x}, \tilde{z}) konceptu \bar{c} délky m nechť Γ označuje pravděpodobnost existence hypotézy $\bar{h} \in \mathbf{H}$, která je konzistentní s (\tilde{x}, \tilde{z}) a pro kterou je $e_{\tilde{P}}(\bar{h}, \bar{c}) > \epsilon$. Potom platí odhad*

$$\Gamma \leq 2\Pi_{\mathbf{H}}(2m) 2^{-\frac{\epsilon m}{2}}.$$

■ *Důkaz:*

DŮKAZ JE UVEDEN V ORIGINÁLNÍ TIŠTĚNÉ VERZI.

Na základě lemmy 4.2.4 je zřejmé, že je-li VC-dimenze třídy hypotéz \mathbf{H} konečná, potom číslo $\Pi_{\mathbf{H}}(m)$ je polynomiálně omezené vzhledem k m . Proto asymptotické chování výrazu $\Pi_{\mathbf{H}}(2m) 2^{-\frac{\epsilon m}{2}}$ v okolí nekonečna je dáno členem exponenciálním a zřejmě se tato hodnota velmi rychle blíží nule pro velké hodnoty m . Nyní odhadněme velikost m (tedy délku vzorku) tak, aby výraz $2\Pi_{\mathbf{H}}(2m) 2^{-\frac{\epsilon m}{2}}$ byl menší než číslo δ .

Lemma 5.2.5 *Je-li*

$$m \geq \max \left(\frac{4}{\epsilon} \log_2 \left(\frac{2}{\delta} \right), \frac{8d}{\epsilon} \log_2 \left(\frac{13}{\epsilon} \right) \right),$$

tak platí $2\Phi_{d,2m} 2^{-\frac{\epsilon m}{2}} \leq \delta$.

■ *Důkaz:*

DŮKAZ JE UVEDEN V ORIGINALNÍ TIŠTĚNÉ VERZI.

■

Lemmy a tvrzení vyslovená v předchozím umožňují již přímo dokázat základní tvrzení této podkapitoly. Vynechejme ale z našich úvah triviální případy, popsané v následující definici.

Definice 5.2.5 *Třídou konceptů* \mathbf{C} *definovanou nad* \bar{X} *nazveme* TRIVIÁLNÍ TŘÍDA KONCEPTŮ, *obsahuje-li pouze jeden koncept, nebo obsahuje-li koncepty dva, které spolu tvoří disjunktní pokrytí množiny* \bar{X} .

Nyní již vyslovme stěžejní tvrzení této části.

Tvrzení 5.2.6 *Nechť* \mathbf{C} *je netriviální, dobře utvořená třída konceptů. Potom*

1. *Je-li* $\text{VC}_{\dim}(\mathbf{C}) = d$ *a* $d < \infty$, *tak*

(a) *pro každé* $0 < \epsilon < \frac{1}{2}$ *žádný učicí algoritmus vyžadující méně než*

$$\max \left(\frac{1-\epsilon}{\epsilon} \log_e \left(\frac{1}{\delta} \right), d(1 - 2(\epsilon(1-\delta) + \delta)) \right) \quad (5.2)$$

dotazů není pro žádný prostor hypotéz \mathbf{H} *(* ϵ, δ *)-učicím algoritmem.*

(b) *pro každé* $0 < \epsilon < 1$ *je každý učicí algoritmus vyžadující alespoň*

$$\max \left(\frac{4}{\epsilon} \log_2 \left(\frac{2}{\delta} \right), \frac{8d}{\epsilon} \log_2 \left(\frac{13}{\epsilon} \right) \right) \quad (5.3)$$

dotazů a produkující konzistentní hypotézu (ϵ, δ) -*učicí algoritmus.*

2. \mathbf{C} *je uniformě naučitelná, právě když* $\text{VC}_{\dim}(\mathbf{C}) < \infty$.

■ *Důkaz:*

DŮKAZ JE UVEDEN V ORIGINALNÍ TIŠTĚNÉ VERZI.

■

Zřejmě toto tvrzení je mnohem silnější než tvrzení 5.1.1, protože hodnota $|\mathbf{C}|$ v horním odhadu je nyní nahrazena hodnotou $\text{VC}_{\dim}(\mathbf{C})$, která může být podstatně menší.

5.2.2 Polynomiální učení vztažené k dimenzi

Vezměme v úvahu výše zmíněný problém nalézt v n dimenzionálním prostoru koncept (konkrétně nějakou kouli), který bude konzistentní s nějakým zadaným vzorkem. Předpokládejme, že máme k dispozici nějaký algoritmus, generující na základě vzorku odpovídající hypotézu a že tento algoritmus pracuje s přesností ϵ , δ ve smyslu uvedeném v předešlých kapitolách. Tvzení 5.2.6 nám dává odhad na nutnou délku vzorku, který zajistí, že námi každý použitý algoritmus produkující konzistentní hypotézy bude mít přesnost popsanou čísly ϵ a δ . Nás tato délka zajímá především proto, abychom si pro danou třídu konceptů byli jisti, že snaha hledat efektivní učící algoritmus je smysluplná v tom smyslu, že počet příkladů nutných k naučení algoritmu bude polynomiální v $\frac{1}{\epsilon}$ a $\frac{1}{\delta}$. Jsou ale obvyklé případy, kdy koncepty v dané třídě konceptů jsou navíc parametrizovány dalšími parametry, např. dimenzí konceptu (koule v prostoru dimenze n), a nás budou pochopitelně zajímat takové podmínky, které zaručí, že pro koncepty s vyšší dimenzí nevrzoste neúměrně délka vzorku nezbytně nutná ke garanci (ϵ, δ) přesnosti učení. Touto problematikou se budeme zabývat nyní.

Definice 5.2.6 *Nechť \bar{U} je nějaká univerzální množina, $n \geq 1$ a $\bar{X}_n \subset \bar{U}^n$. Nechť $C_n \subset 2^{\bar{X}_n}$ je třída konceptů definovaná nad \bar{X}_n a $H_n \subset 2^{\bar{X}_n}$ je prostor hypotéz pro C_n . Předpokládejme dále, že každá H_n je dobře utvořená a obsahuje prázdnou množinu a celý prostor \bar{X}_n . Množinu $\mathcal{C} \stackrel{\text{def}}{=} \bigcup_{i=1}^{\infty} C_n$ nazveme SYSTÉM KONCEPTŮ nad posloupností $\{\bar{X}_i\}_1^{\infty}$. Stejně tak množinu $\mathcal{H} \stackrel{\text{def}}{=} \bigcup_{i=1}^{\infty} H_n$ nazveme SYSTÉMEM HYPOTÉZ.*

Ještě definujme prostor vzorů pro systém konceptů analogicky, jako pro prostor vzorů pro třídu konceptů.

Definice 5.2.7 *Množinu*

$$\bar{S}_{\mathcal{C}} \stackrel{\text{def}}{=} \bigcup_{i=1}^{\infty} \bar{S}_{C_n}$$

nazveme PROSTOR VZORŮ SYSTÉMU KONCEPTŮ \mathcal{C} .

Z formálního hlediska zavedeme reprezentaci konceptů pomocí slov nějakého jazyka a budeme předpokládat, že tato reprezentace umožňuje pracovat s koncepty a jejich prvky v čase, který je polynomiální vzhledem k dimenzi jednotlivých prvků v daném konceptu, tedy vzhledem k n . To nám v dalším umožní odhlédnout od problému, jakou výpočetní složitost má samotný přístup ke konceptům a jejich prvkům, a zaměříme se pouze na problém složitosti vzhledem k počtu vzorů pro danou dimenzi konceptů. V následující definici symbol $\{0, 1\}^*$ označuje množinu všech konečných posloupností tvořených nulami a jedničkami.

Definice 5.2.8 *Řekneme, že systém konceptů \mathcal{C} (resp. třída konceptů \mathcal{C}) je POLYNOMIÁLNĚ DEFINOVANÝ SYSTÉM KONCEPTŮ, (resp. POLYNOMIÁLNĚ DEFINOVANÁ TŘÍDA KONCEPTŮ) existuje-li jazyk $R \subset \{0, 1\}^*$ takový, že:*

1. *Mezi jazykem R a množinou \mathcal{C} (resp. \mathcal{C}) existuje vzájemně jednoznačné zobrazení.*
2. *Existuje polynomiální algoritmus (polynomiální vzhledem k n), který rozhodne, zda slovo $z \in \{0, 1\}^*$ délky n je v jazyce R .*

3. *Existuje polynomiální algoritmus (polynomiální vzhledem k n), který pro dané slovo $r \in R$, n a $x \in \bar{X}_n$ (resp. $x \in \bar{X}$) rozhodne, zda x leží v konceptu, který je reprezentován slovem r .*

Stejně tak definujeme POLYNOMIÁLNĚ DEFINOVANÝ SYSTÉM HYPOTÉZ (resp. POLYNOMIÁLNĚ DEFINOVANOU TŘIDU HYPOTÉZ).

Zřejmě například HALFSPACE_n je polynomiálně definovaná třída konceptů. Většina tříd konceptů přirozeně vzniklých je polynomiálně definovaná třída konceptů.

V dalším textu se budeme zabývat již pouze takovými systémy konceptů a hypotéz, které budou polynomiálně definované. Dále rozšíříme pojem naučitelnosti tříd konceptů na naučitelnost systému konceptů podle následující definice.

Definice 5.2.9 *Předpokládejme, že \mathcal{C} je polynomiálně definovaný systém konceptů a \mathcal{H} je polynomiálně definovaný systém hypotéz, $\mathcal{C} \subset \mathcal{H}$ a \bar{N} je množina přirozených čísel. Mějme definovanou funkci $\bar{m}(\epsilon, \delta, n)$ zobrazující množinu $(0, 1) \times (0, 1) \times \bar{N}$ do množiny \bar{N} . Potom řekneme, že \mathcal{C} JE D-POLYNOMIÁLNĚ NAUČITELNÝ SYSTÉM KONCEPTŮ PODLE \mathcal{H} , existuje-li polynomiální algoritmus (vzhledem k času) A^* zobrazující $\bar{S}_{\mathcal{C}}$ do \mathcal{H} tak, že pro všechna $0 < \epsilon < 1$, $0 < \delta < 1$ a $n \geq 1$ existuje přirozené číslo $\bar{m}(\epsilon, \delta, n)$ polynomiálně omezené vzhledem k $\frac{1}{\epsilon}$, $\frac{1}{\delta}$ a n tak, že pro všechna $n \geq 1$, pro všechny koncepty $\bar{c} \in \mathcal{C}_n$ a všechny hustoty pravděpodobnosti \tilde{P} definované na \bar{X}_n je pravděpodobnost množiny*

$$\left\{ \bar{x} \in \bar{X}_n^{\bar{m}(\epsilon, \delta, n)} \mid (\bar{x}, \bar{z}) \text{ je vzorek } \bar{c} \text{ a } e_{\tilde{P}}(\bar{c}, A^*((\bar{x}, \bar{z}))) \geq \epsilon \right\}$$

menší než δ . Je-li \mathcal{C} d-polynomiálně naučitelný systém konceptů podle \mathcal{C} (tedy podle sebe sama), říkáme, že \mathcal{C} je TOTÁLNĚ NAUČITELNÝ SYSTÉM KONCEPTŮ. Minimální hodnotu $\bar{m}(\epsilon, \delta, n)$ s touto vlastností nazveme VZOROVOU SLOŽITOSTÍ učicího algoritmu A^ .*

Následující tvrzení postuluje ekvivalentní podmínku pro to, aby systém konceptů byl totálně naučitelný. Ukazuje se, že nutně musí být $\text{VC}_{\dim}(\mathcal{C}_n)$ polynomiálně omezena v n a musí existovat možnost, jak lze na základě nějakého nedeterministického algoritmu generovat hypotézu konzistentní s libovolným vzorkem z $\bar{S}_{\mathcal{C}}$ s pravděpodobností alespoň γ , kde $\gamma > 0$ je nezávislé na n .

Definice 5.2.10 *Nechť \mathcal{C} je polynomiálně definovaný systém konceptů. Nechť pro dané $\gamma > 0$ existuje nedeterministický polynomiální algoritmus A^* zobrazující $\bar{S}_{\mathcal{C}}$ do \mathcal{C} tak, že pro všechna $(\bar{x}, \bar{z}) \in \bar{S}_{\mathcal{C}}$ je koncept $A^*((\bar{x}, \bar{z}))$ konzistentní s (\bar{x}, \bar{z}) s pravděpodobností alespoň γ . Potom takovýto algoritmus nazveme NEDETERMINISTICKÝ POLYNOMIÁLNÍ GENERÁTOR HYPOTÉZ pro systém konceptů \mathcal{C} .*

Tvrzení 5.2.7 *Libovolný systém konceptů \mathcal{C} je totálně naučitelný, právě když existuje nedeterministický polynomiální generátor hypotéz pro \mathcal{C} a $\text{VC}_{\dim}(\mathcal{C}_n)$ je omezena polynomem v n .*

■ *Důkaz:*

DŮKAZ JE UVEDEN V ORIGINÁLNÍ TIŠTĚNÉ VERZI.

■

V případě, kdy systém konceptů je konstruován nad binární množinou $\{-1, +1\}$, existuje velmi užitečný vztah mezi čísly $\text{VC}_{\dim}(\mathbf{C}_n)$ a čísly $\log_2(|\mathbf{C}_n|)$, který je popsán v následující lemmě.

Lemma 5.2.8 *Předpokládejme, že $\bar{X}_n \stackrel{\text{def}}{=} \{-1, +1\}^n$ a $\mathbf{C}_n \subset 2^{\bar{X}_n}$, $n \geq 1$. Potom $\text{VC}_{\dim}(\mathbf{C}_n)$ roste polynomiálně vzhledem k n , právě když $\log_2(|\mathbf{C}_n|)$ roste polynomiálně v n .*

■ *Důkaz:*

DŮKAZ JE UVEDEN V ORIGINÁLNÍ TIŠTĚNÉ VERZI.

■

Cvičení 5.2.2 Učení systému konceptů

1. Předpokládejme, že třídy \mathbf{P} a \mathbf{NP} nejsou totožné (což se obecně předpokládá, že platí). Potom na základě nějakého vybraného NP-úplného problému zkonstruujte nepolynomiálně definovaný systém konceptů.
2. Předpokládejme, že $\bar{X}_n \stackrel{\text{def}}{=} \{-1, +1\}^n$. Zkonstruujte systém konceptů \mathbf{C} , pro který posloupnost $\log_2(|\mathbf{C}_n|)$, $\mathbf{C}_n \subset 2^{\bar{X}_n}$, $n \geq 1$ není omezena polynomem v n .

5.2.3 Přibližné řešení problému pokrytí množin

Při dokazování některých vlastností učicích algoritmů vyvstává problém vybrat pro daný konečný systém konečných množin \bar{S}_i , $i \in \{1, \dots, n\}$, a danou množinu $\bar{A} \subset \bigcup_{i=1}^n \bar{S}_i$, pokud možno co nejmenší počet množin \bar{S}_i , které mají tu vlastnost, že jejich sjednocení obsahuje množinu \bar{A} . Obecně se dá tento problém formulovat jako problém pokrytí konečného sjednocení konečných množin, kdy hledáme minimální počet množin tvaru $\bar{S}_i \cap \bar{A}$, jejichž sjednocení je právě rovno množině \bar{A} . Problém nalezení takového optimálního pokrytí (pokrytí sestávající z minimálního počtu množin) leží ve třídě NP-úplných problémů. Existuje však hladový algoritmus, který sice nenalezne optimální pokrytí, ale v čase úměrném $m \log_e(m)$ nalezne takové pokrytí množiny \bar{A} mohutnosti m , které obsahuje nejvýše $(\log_e(m) + 1)$ -krát více množin nežli optimální pokrytí. V dalším výkladu uvidíme, že tento odhad je dostatečný pro analýzu některých vlastností učicích algoritmů. Věnujme se nyní výkladu zmíněného hladového algoritmu (viz. [Joh74]).

Definice 5.2.11 *Nechť $\bar{S}_1, \dots, \bar{S}_n$ je libovolná konečná posloupnost konečných množin. Pro každé $i \in \{1, \dots, n\}$ položme $\bar{T}_{i,0} \stackrel{\text{def}}{=} \bar{S}_i$, číslo α_0 nechť splňuje podmínku*

$$|\bar{T}_{\alpha_0,0}| \stackrel{\text{def}}{=} \max_{i \in \{1, \dots, n\}} |\bar{T}_{i,0}|$$

a definujme

$$\bar{\Omega}_0 \stackrel{\text{def}}{=} \bar{S}_{\alpha_0}.$$

Potom pro každé $j \geq 1$ a $i \in \{1, \dots, n\}$ rekurzivně definujme množiny

$$\bar{T}_{i,j} \stackrel{\text{def}}{=} \bar{T}_{i,j-1} \setminus \bar{T}_{\alpha_{j-1},j-1}, \quad (5.4)$$

číslo α_j tak, aby platilo

$$|\bar{T}_{\alpha_j, j}| \stackrel{\text{def}}{=} \max_{i \in \{1, \dots, n\}} |\bar{T}_{i, j}|$$

a množinu

$$\bar{\Omega}_j \stackrel{\text{def}}{=} \bar{\Omega}_{j-1} \cup \bar{S}_{\alpha_j} .$$

Základní vlastnosti právě zavedených množin popisuje následující lemma.

Lemma 5.2.9 *Nechť $\bar{S}_1, \dots, \bar{S}_n$ je konečná posloupnost konečných množin a $\bar{\Omega}_j, j \geq 0$, je posloupnost množin definovaná v definici 5.2.11. Potom platí následující výroky:*

1. Existuje l takové, že $\bar{\Omega}_l = \bigcup_{i=1}^n \bar{S}_i$.
2. $(\forall j \geq 0) (\forall k > j) (\forall i \in \{1, \dots, n\})$ platí, že $\bar{S}_{\alpha_j} \cap \bar{T}_{i, k} = \emptyset$.
3. $(\forall r \geq 0)$ je

$$\bigcup_{j > r} \bar{T}_{\alpha_j, r} = \left(\bigcup_{i=1}^n \bar{S}_i \right) \dot{-} \bar{\Omega}_r. \quad (5.5)$$

4. Nechť $r \geq 0$. Potom platí

$$(\forall j > r) (\forall i \in \{1, \dots, n\}) \left(|\bar{T}_{i, j}| \leq |\bar{T}_{\alpha_r, r}| \right).$$

5. Nechť $q, r \geq 0$ jsou přirozená čísla a $\bar{M} \subset \{1, \dots, n\}$. Navíc nechť platí

$$\bigcup_{i \in \bar{M}} \bar{T}_{i, q+1} = \left(\bigcup_{i=1}^n \bar{S}_i \right) \dot{-} \bar{\Omega}_q. \quad (5.6)$$

Potom platí

$$\bigcup_{i \in \bar{M}} \left(\bar{T}_{i, q+1} \dot{-} \bar{T}_{i, r+1} \right) = \bar{\Omega}_r \dot{-} \bar{\Omega}_q.$$

■ *Důkaz:*

DŮKAZ JE UVEDEN V ORIGINÁLNÍ TIŠTĚNÉ VERZI.

■

Bod 1 předchozí lemmy nás opravňuje k zavedení následujícího pojmu.

Definice 5.2.12 *Nejmenší takové číslo l , pro které platí rovnost*

$$\bar{\Omega}_l = \bigcup_{i=1}^n \bar{S}_i ,$$

nazveme DÉLKOU APROXIMACE OPTIMÁLNÍHO POKRYTÍ množiny $\bigcup_{i=1}^n \bar{S}_i$.

Pro porovnání velikosti pokrytí, nalezeného algoritmem popsáným v definici 5.2.11, a velikosti optimálního pokrytí pro konkrétní zadanou úlohu slouží odhad vyslovený v lemmě 5.2.10.

Lemma 5.2.10 *Nechť l je délka aproximace optimálního pokrytí množiny $\bigcup_{i=1}^n \bar{S}_i$. Nechť dále pro nějaké $q \geq 0$ a $\bar{M} \subset \{1, \dots, n\}$ platí rovnost*

$$\bigcup_{i \in \bar{M}} \bar{T}_{i,q+1} = \left(\bigcup_{i=1}^n \bar{S}_i \right) \dot{-} \bar{\Omega}_q. \quad (5.7)$$

Potom platí nerovnost

$$l - q \leq \sum_{i \in \bar{M}} \left(\sum_{z=1}^{|\bar{T}_{i,q+1}|} \frac{1}{z} \right).$$

■ *Důkaz:*

DŮKAZ JE UVEDEN V ORIGINALNÍ TIŠTĚNÉ VERZI. ■

Přímým důsledkem vyslovené lemy je odhad, který říká, že pokrytí nalezené analyzovaným algoritmem je v porovnání s pokrytím optimálním v nejhorším případě

$$(\log_e(m) + 1)$$

krátě větší (co do počtu množin realizujících pokrytí), kde m je maximální mohutnost množin tvořících sjednocení.

Tvrzení 5.2.11 *Nechť $\bar{S}_1, \dots, \bar{S}_n$ je konečná posloupnost konečných množin splňujících pro nějaké přirozené m podmínku $|\bar{S}_i| \leq m$, $i \in \{1, \dots, n\}$. Dále nechť $\bar{M} \subset \{1, \dots, n\}$ je taková množina, že systém množin $\{\bar{S}_i \mid i \in \bar{M}\}$ tvoří minimální pokrytí množiny $\bigcup_{i=1}^n \bar{S}_i$, to jest*

$$|\bar{M}| \stackrel{\text{def}}{=} \min \left\{ |\bar{Q}| \mid \bar{Q} \subset \{1, \dots, n\} \text{ a } \bigcup_{i \in \bar{Q}} \bar{S}_i = \bigcup_{i=1}^n \bar{S}_i \right\}.$$

Nechť l je jakákoliv délka aproximace optimálního pokrytí množiny $\bigcup_{i=1}^n \bar{S}_i$. Potom platí odhad

$$\frac{l}{|\bar{M}|} \leq \sum_{z=1}^m \frac{1}{z} \leq \log_e(m) + 1.$$

■ *Důkaz:*

DŮKAZ JE UVEDEN V ORIGINALNÍ TIŠTĚNÉ VERZI. ■

Závěrem této malé exkurze do oblasti algoritmů přibližně řešících NP-úplné problémy vyslovme tvrzení o aproximaci optimálního pokrytí množin ve tvaru, který je vhodný z hlediska analýzy vlastností učících algoritmů pro třídy konceptů, kde jednotlivé koncepty jsou dány jako sjednocení (či průnik) konceptů s omezenou popisnou složitostí.

Tvrzení 5.2.12 *Nechť $\bar{S}_1, \dots, \bar{S}_n$ je konečná posloupnost konečných množin a \bar{A} je libovolná podmnožina $\bigcup_{i=1}^n \bar{S}_i$ mohutnosti m , nechť V označuje počet množin \bar{S}_i , tvořících minimální pokrytí množiny \bar{A} . Potom lze nalézt pokrytí množiny \bar{A} obsahující méně než $V(\log_e(m) + 1)$ množin \bar{S}_i v čase úměrném hodnotě $nmV\log_e(m)$.*

■ *Důkaz:*

DŮKAZ JE UVEDEN V ORIGINALNÍ TIŠTĚNÉ VERZI. ■

5.2.4 Polynomiální učení vztažené ke složitosti konceptu.

Obraťme nyní pozornost na další aspekt problematiky učení. Předpokládejme, že máme danou třídu konceptů \mathcal{C} a že koncepty z této třídy mají různou složitost. Přesněji, za tuto složitost budeme považovat hodnotu míry složitosti, což bude funkce definovaná pro každý koncept a vyjadující míru jeho složitosti. Budeme řešit otázku, jak odhadnout vzorovou složitost učících algoritmů vzhledem k míře složitosti cílového konceptu. Pro jednoduchost ztotožníme v této kapitole třídu hypotéz \mathcal{H} s třídou konceptů \mathcal{C} , což však není újma na obecnosti, protože všechny závěry této podkapitoly se dají elementárně zobecnit i pro obecný případ $\mathcal{C} \subset \mathcal{H}$. Nejdříve definujme zmíněnou míru složitosti.

Definice 5.2.13 *Nechť \mathcal{C} je třída konceptů nad \bar{X} a nechť je definována funkce \widetilde{cm} zobrazující třídu konceptů \mathcal{C} do množiny přirozených čísel. Pak tuto funkci nazveme KONCEPTUÁLNÍ MÍROU SLOŽITOSTI třídy konceptů \mathcal{C} .*

Analogicky jako v případě dimenzionální složitosti definujme konceptuální složitost algoritmu A^* a s-polynomiální naučitelnost třídy konceptů \mathcal{C} .

Definice 5.2.14 *Předpokládejme, že \mathcal{C} je polynomiálně definovaná třída konceptů, \bar{N} je množina přirozených čísel, a $\mathcal{H} = \mathcal{C}$ je polynomiálně definovaná třída hypotéz. Mějme definovanou funkci $\widetilde{m}(\epsilon, \delta, s)$ zobrazující množinu $(0, 1) \times (0, 1) \times \bar{N}$ do množiny přirozených čísel. Potom řekneme, že \mathcal{C} JE S-POLYNOMIÁLNĚ NAUČITELNÁ TŘÍDA KONCEPTŮ, existuje-li polynomiální algoritmus (vzhledem k času) A^* zobrazující $\bar{S}_{\mathcal{C}}$ do \mathcal{H} tak, že pro všechna $0 < \epsilon < 1$, $0 < \delta < 1$ a $s \geq 1$ existuje přirozené číslo $\widetilde{m}(\epsilon, \delta, s)$ polynomiálně omezené vzhledem k $\frac{1}{\epsilon}$, $\frac{1}{\delta}$ a s tak, že pro všechny koncepty $\bar{c} \in \mathcal{C}$ pro které je $\widetilde{cm}(\bar{c}) \leq s$ a všechny hustoty pravděpodobnosti \tilde{P} definované na \bar{X} , je pravděpodobnost množiny*

$$\left\{ \bar{x} \in \bar{X}^{\widetilde{m}(\epsilon, \delta, s)} \mid (\bar{x}, \bar{z}) \text{ je vzorek } \bar{c} \text{ a } e_{\tilde{P}}(\bar{c}, \widetilde{A}^*((\bar{x}, \bar{z}))) \geq \epsilon \right\}$$

menší než δ . Minimální hodnotu $\widetilde{m}(\epsilon, \delta, s)$ s touto vlastností nazveme KONCEPTUÁLNÍ SLOŽITOSTÍ učícího algoritmu A^ .*

Naše další úvahy velmi výrazně zjednoduší zavedení pojmu efektivní prostor hypotéz, neboť nám umožní pro daný algoritmus pracovat pouze s takovou podmnožinou třídy konceptů \mathcal{C} , pro jejíž každý koncept \bar{a} existuje při zobrazení daným algoritmem vzor splňující předem zadaná omezení na míru složitosti konceptu \bar{a} a na délku vzorku, nezbytného ke generování konceptu \bar{a} uvažovaným algoritmem.

Definice 5.2.15 *Předpokládejme, že \mathcal{C} je polynomiálně definovaná třída konceptů a že existuje polynomiální algoritmus (vzhledem k času) A^* zobrazující $\bar{S}_{\mathcal{C}}$ do \mathcal{C} tak, že pro všechna $(\bar{x}, \bar{z}) \in \bar{S}_{\mathcal{C}}$ je $\widetilde{A}^*((\bar{x}, \bar{z}))$ konzistentní s (\bar{x}, \bar{z}) . Pro pevně zvolená přirozená čísla m, s definujme následující množiny:*

$$\mathcal{Z}_s \stackrel{\text{def}}{=} \{ \bar{c} \in \mathcal{C} \mid \overline{ssc}(\bar{c}) \leq s \}$$

a

$$\mathcal{Y}_{m,s} \stackrel{\text{def}}{=} \left\{ (\bar{x}, \bar{z}) \mid \bar{x} \in \bar{X}^m \text{ a } (\exists \bar{c} \in \mathcal{Z}_s) \left((\bar{x}, \bar{z}) \text{ je konzistentní s } \bar{c} \right) \right\}.$$

Potom množinu

$$\text{EF}_{A^*, m, s} \stackrel{\text{def}}{=} \left\{ \widetilde{A}^*((\bar{x}, \bar{z})) \mid (\bar{x}, \bar{z}) \in \mathcal{Y}_{m,s} \right\}$$

nazveme EFEKTIVNÍ PROSTOR HYPOTÉZ ALGORITMU A^ PRO KONCEPTY SLOŽITOSTI s A VZORKY DÉLKY m (stručněji s, m -EFEKTIVNÍ PROSTOR).*

Definice 5.2.16 Předpokládejme, že pro nějaký algoritmus A^* je definován efektivní prostor hypotéz $EF_{A^*,m,s}$. Potom algoritmus A^* je OCCAMŮV ALGORITMUS, právě když existuje číslo α , $0 \leq \alpha < 1$ a polynom $\tilde{p}(s)$ tak, že pro všechna $m, s \geq 1$ platí, že

$$VC_{dim}(EF_{A^*,m,s}) \leq \tilde{p}(s)m^\alpha.$$

Nyní dokážeme jednu technickou lemmu sloužící k důkazu tvrzení následujícího.

Lemma 5.2.13

1. Je-li $0 < \epsilon, \delta < 1, k \geq 1, 0 \leq \alpha < 1$ a

$$m \stackrel{def}{=} \max \left(\frac{4}{\epsilon} \log_2 \left(\frac{2}{\delta} \right), \left(\frac{8k}{\epsilon} \log_2 \left(\frac{13}{\epsilon} \right) \right)^{\frac{1}{1-\alpha}} \right), \quad (5.8)$$

potom platí

$$\frac{2 \left(\frac{2em}{km^\alpha} \right)^{km^\alpha}}{2^{\frac{m\epsilon}{2}}} \leq \delta.$$

2. Je-li $l \geq 1$ a druhá část omezující m ve výrazu 5.8 je nahrazena výrazem

$$\frac{2^{l+4}k}{\epsilon} \left(\log_2 \left(\frac{8(2l+2)^{l+1}k}{\epsilon} \right) \right)^{l+1}, \quad (5.9)$$

pak

$$\frac{2(2m)^{k(\log_2(m))^l}}{2^{\frac{m\epsilon}{2}}} \leq \delta.$$

■ *Důkaz:*

DŮKAZ JE UVEDEN V ORIGINALNÍ TIŠTĚNÉ VERZI. ■

Následující tvrzení dává do vzájemné souvislosti Vapnik-Chervonenkovu dimenzi efektivního prostoru hypotéz a délku vzorku, nezbytnou pro příslušný učící algoritmus A^* .

Tvrzení 5.2.14 Nechť C je třída konceptů s definovanou konceptuální mírou složitosti. Potom platí následující tvrzení.

1. (a) Je-li $VC_{dim}(EF_{A^*,m,s}) \leq \tilde{p}(s)m^\alpha$ pro nějaký polynom $\tilde{p}(s) \geq 1$ a $0 \leq \alpha < 1$, potom A^* je (ϵ, δ) -učící algoritmus vyžaduje-li alespoň

$$m \stackrel{def}{=} \max \left(\frac{4}{\epsilon} \log_2 \left(\frac{2}{\delta} \right), \left(\frac{8\tilde{p}(s)}{\epsilon} \log_2 \left(\frac{13}{\epsilon} \right) \right)^{\frac{1}{1-\alpha}} \right) \quad (5.10)$$

dotazů.

- (b) Je-li $VC_{dim}(EF_{A^*,m,s}) \leq \tilde{p}(s)(\log_2(m))^l$ pro nějaký polynom $\tilde{p}(s) \geq 2$ a $l \geq 1$, potom předchozí odhad platí s tím, že druhá část je rovna výrazu

$$\frac{2^{l+4}\tilde{p}(s)}{\epsilon} \left(\log_2 \left(\frac{8(2l+2)^{l+1}\tilde{p}(s)}{\epsilon} \right) \right)^{l+1}. \quad (5.11)$$

2. *Existuje-li Occamův algoritmus pro \mathcal{C} , pak \mathcal{C} je s -polynomiálně naučitelná třída konceptů.*

■ *Důkaz:*

DŮKAZ JE UVEDEN V ORIGINÁLNÍ TIŠTĚNÉ VERZI.

Definice 5.2.17 *Nechť \mathcal{C} je polynomiálně definovaná třída konceptů a nechť existuje polynomiální algoritmus A^* zobrazující $\bar{S}_{\mathcal{C}}$ do \mathcal{C} tak, že pro všechna $(\tilde{x}, \tilde{z}) \in \bar{S}_{\mathcal{C}}$ je koncept $\widetilde{A^*}((\tilde{x}, \tilde{z}))$ konzistentní s (\tilde{x}, \tilde{z}) . Potom takovýto algoritmus nazveme POLYNOMIÁLNÍ GENERÁTOR HYPOTÉZ pro třídu konceptů \mathcal{C} . Dále definujeme PROBLÉM KONZISTENCE PRO \mathcal{C} jako rozhodovací problém, zda pro daný vzorek z $\bar{S}_{\mathcal{C}}$ existuje koncept $\bar{c} \in \mathcal{C}$, konzistentní s tímto vzorkem.*

Důležitou vlastností polynomiálního generátoru hypotéz je to, že existence takového generátoru zaručuje, že problém konzistence pro třídu konceptů \mathcal{C} leží ve třídě \mathbf{P} polynomiálně řešitelných problémů. Poznamenejme, že výše zavedený polynomiální generátor hypotéz je deterministický algoritmus.

Obdobně jako v analýze vlastností učicích algoritmů vzhledem k dimenzi platí i zde obdobné tvrzení, které říká, že za předpokladu existence polynomiálního generátoru hypotéz má třída konceptů \mathcal{C} z jistého úhlu pohledu "rozumné vlastnosti" ve smyslu tvrzení následující lemma.

Lemma 5.2.15 *Nechť třída konceptů má konečnou VC-dimenzi a problém konzistence pro třídu konceptů \mathcal{C} je ve třídě \mathbf{P} . Potom pro každou konečnou množinu $\bar{Z} \subset \bar{X}$ množina $\Pi_{\mathcal{C}}(\bar{Z})$ všech podmnožin množiny \bar{Z} , které jsou rozděleny třídou konceptů \mathcal{C} , může být proběhnuta v čase polynomiálním vzhledem k mohutnosti množiny \bar{Z} .*

■ *Důkaz:*

DŮKAZ JE UVEDEN V ORIGINÁLNÍ TIŠTĚNÉ VERZI.

Vlastnost třídy konceptů popsaná v předešlé lemmě nám umožní zabývat se otázkou, zda pro danou třídu konceptů \mathcal{C} (pro kterou například může existovat efektivní učicí algoritmus) nelze nalézt třídu konceptů, která je složitější (např. s vyšší VC-dimenzí), ale která současně od třídy \mathcal{C} převzala vlastnosti, zaručující efektivitu učicích algoritmů pro tuto novou odvozenou třídu konceptů. Uvidíme dále, že dva způsoby odvození takovýchto složitějších tříd lze realizovat prostřednictvím konečných průniků a sjednocení konceptů z původní třídy konceptů \mathcal{C} .

Definice 5.2.18 *Nechť $\mathcal{C} \subset 2^{\bar{X}}$ je třída konceptů. Označme $U_{\mathcal{C}}$ množinu všech konečných sjednocení množin z \mathcal{C} a $I_{\mathcal{C}}$ množinu všech konečných průniků množin z \mathcal{C} . Dále definujme pro každé přirozené číslo k množiny*

$$U_{k,\mathcal{C}} \stackrel{\text{def}}{=} \left\{ \bigcup_{i=1}^k \bar{c}_i \mid (\forall i \in \{1, \dots, k\}) (\bar{c}_i \in \mathcal{C}) \right\}$$

a

$$I_{k,\mathcal{C}} \stackrel{\text{def}}{=} \left\{ \bigcap_{i=1}^k \bar{c}_i \mid (\forall i \in \{1, \dots, k\}) (\bar{c}_i \in \mathcal{C}) \right\}.$$

Lemma 5.2.16 *Nechť $C \subset 2^{\bar{X}}$ je třída konceptů, $VC_{dim}(C) = d \geq 1$ je konečná. Potom pro všechna $k \geq 1$ je $VC_{dim}(U_{k,C}) \leq 2dk \log_2(3k)$ a $VC_{dim}(I_{k,C}) \leq 2dk \log_2(3k)$.*

■ *Důkaz:*

DŮKAZ JE UVEDEN V ORIGINALNÍ TIŠTĚNÉ VERZI.

■

Vyslovená lemma umožňuje důkaz následujícího tvrzení o s -polynomiální naučitelnosti tříd konceptů U_C a I_C .

Tvrzení 5.2.17 *Předpokládejme, že C je třída konceptů pro kterou existuje polynomiální generátor hypotéz, a nechť navíc $VC_{dim}(C)$ je konečná. Potom třídy konceptů U_C a I_C jsou s -polynomiálně naučitelné.*

■ *Důkaz:*

DŮKAZ JE UVEDEN V ORIGINALNÍ TIŠTĚNÉ VERZI.

■

Kapitola 6

Numerická analýza učicích algoritmů

Finaglův numerický postulát

POSLEDNÍ ITERACE ŘEŠÍ ZADANOU ÚLOHU.
NENECHTE SE ZMÁST FAKTY.

6.1 Analýza gradientních metod

Nyní se zabýváme numerickou analýzou vybraných učicích algoritmů, které využívají při nalezení cílových vah gradientních metod pro minimalizaci chybové (účelové, energetické) funkce. Tento přístup je široce používán v rozličných aplikacích neuronových sítí. Zmíníme se o tzv. delta pravidlu a různých vlastnostech jeho modifikací, dále pak stručně zmíníme nelineární zobecnění delta pravidla, jehož speciálním případem je i často používaná metoda zpětného šíření (anglický termín: back-propagation).

6.1.1 Delta pravidlo

Nejdříve vezměme v úvahu gradientní algoritmus pro učení jednoho perceptronu, jehož váhy budeme chápat jako vektor a označovat symbolem \vec{w} . Učící algoritmus po předložení vstupního vektoru \vec{x} a příslušné výstupní hodnoty y adaptuje hodnoty vah podle formule ($\eta > 0$)

$$\delta \vec{w} \stackrel{\text{def}}{=} \eta (y - \langle \vec{w} | \vec{x} \rangle) \vec{x}, \quad (6.1)$$

což vede k iterační posloupnosti definované jako

$$\vec{w}_{k+1} \stackrel{\text{def}}{=} \vec{w}_k + \delta \vec{w}_k = \vec{w}_k + \eta (y - \vec{w}_k^T \vec{x}) \vec{x} = (\mathbf{I} - \eta \vec{x} \vec{x}^T) \vec{w}_k + \eta y \vec{x}$$

(v dalším textu bude symbol \vec{w}_k označovat k -tý člen iterační posloupnosti, \vec{x}_i i -tý vstupní vektor). Adaptaci vah podle 6.1 budeme říkat učení podle delta pravidla. Význam tohoto učení spočívá v tom, že přírůstek váhy mezi i -tým neuronem vstupní vrstvy a neuronem výstupním je přímo úměrný jak chybě aproximace hodnoty y , tak i i -té souřadnici vektoru \vec{x} . Je-li chyba aproximace y kladná, pak zvyšujeme velikost těch vah, pro které je \vec{x}_i kladné (a tedy chyba se snižuje), a snižujeme hodnoty vah, pro které je \vec{x}_i záporné (a

tedy chyba se opět snižuje). Obdobně tato argumentace platí pro případ, kdy chyba aproximace y je záporná (všechny tyto úvahy platí samozřejmě pro dostatečně malou hodnotu parametru η). Z geometrického pohledu, vzorec 6.1 vyjadřuje tu skutečnost, že při adaptaci váhy \vec{w}_i se nadrovina definovaná váhovým vektorem \vec{w} natočí do takového směru, že hodnota lokální chyby $\|y - \langle \vec{w} | \vec{x} \rangle\|$ se sníží. Problémem je, aby toto snížení nebylo na straně druhé zapláceno příliš vysokým přírůstkem lokální chyby pro jiné vektory, na které chceme daný perceptron naučit (kritériem celkové kvality naučení se nejčastěji bere hodnota "energetické" funkce, kterou zavedeme níže).

V následujících lemmách a tvrzeních ukážeme, že pro malé hodnoty parametru η je učicí algoritmus založený na delta pravidle v jistém kontextu konvergentní a také numericky stabilní.

Lemma 6.1.1 *Nechť $\mathbf{B} \stackrel{\text{def}}{=} (\mathbf{I} - \eta \vec{x} \vec{x}^T)$. Potom \mathbf{B} má pouze vlastní číslo $1 - \eta \|\vec{x}\|^2$, které odpovídá vlastnímu vektoru \vec{x} , a vlastní číslo 1 odpovídající vlastním vektorům v ortogonálním doplňku vektoru \vec{x} .*

■ *Důkaz:*

DŮKAZ JE UVEDEN V ORIGINÁLNÍ TIŠTĚNÉ VERZI.

■

Lemma 6.1.2 *Nechť $\mathbf{B} \stackrel{\text{def}}{=} (\mathbf{I} - \eta \vec{x} \vec{x}^T)$ a $0 \leq \eta \leq \frac{2}{\|\vec{x}\|^2}$. Potom je $\|\mathbf{B}\| = \rho(\mathbf{B}) = 1$, kde $\rho(\mathbf{B})$ označuje spektrální poloměr matice \mathbf{B} .*

■ *Důkaz:*

DŮKAZ JE UVEDEN V ORIGINÁLNÍ TIŠTĚNÉ VERZI.

■

Matice \mathbf{A} definovaná v následující definici bude hrát klíčivou roli při analýze zmíněné iterační posloupnosti pro váhový vektor \vec{w} .

Definice 6.1.1 *Předpokládejme, že $\vec{x}_1, \dots, \vec{x}_t \in \mathbb{R}^n$. Pro každý vektor \vec{x}_p definujme matici $\mathbf{B}_p = (\mathbf{I} - \eta \vec{x}_p \vec{x}_p^T)$. Potom pro každou permutaci π množiny $(1, \dots, t)$ definujme matici*

$$\mathbf{A}_\pi \stackrel{\text{def}}{=} \mathbf{B}_{\pi(t)} \mathbf{B}_{\pi(t-1)} \cdots \mathbf{B}_{\pi(1)}$$

a vektor

$$\vec{h} \stackrel{\text{def}}{=} y_{\pi(1)} (\mathbf{B}_{\pi(t)} \mathbf{B}_{\pi(t-1)} \cdots \mathbf{B}_{\pi(2)}) \vec{x}_{\pi(1)} + \cdots + y_{\pi(t-1)} \mathbf{B}_{\pi(t)} \vec{x}_{\pi(t-1)} + y_{\pi(t)} \vec{x}_{\pi(t)}.$$

Lemma 6.1.3 *Nechť pro všechna $p \in \{1, \dots, t\}$ platí $0 \leq \eta \leq \frac{2}{\|\vec{x}_p\|^2}$ a nechť soubor vektorů $\vec{x}_1, \dots, \vec{x}_t$ generuje celý prostor \mathbb{R}^n . Potom pro libovolnou permutaci π je $\|\mathbf{A}_\pi\| < 1$.*

■ *Důkaz:*

DŮKAZ JE UVEDEN V ORIGINÁLNÍ TIŠTĚNÉ VERZI.

■

Nyní rigorózním způsobem zdefinujeme posloupnost váhových parametrů vzniklou aplikací delta pravidla, se kterou budeme pracovat dále.

Definice 6.1.2 Předpokládejme, že $(\vec{x}_1, y_1), \dots, (\vec{x}_t, y_t)$ je zadaná posloupnost uspořádaných dvojic z $\mathbb{R}^n \times \mathbb{R}$, $t \geq 1$. Nechť pro dané $\eta > 0$ je posloupnost vektorů $\{\vec{w}_i\}_1^\infty$ definovaná rekurzivní předpisem

$$\begin{aligned} \vec{w}_1 &\stackrel{\text{def}}{=} \vec{w}_0 && \vec{w}_0 \in \mathbb{R}^n, \text{ libovolné,} \\ \vec{w}_{k+1} &\stackrel{\text{def}}{=} (\mathbf{I} - \eta \vec{x}_k \vec{x}_k^T) \vec{w}_k + \eta y_k \vec{x}_k, && \text{pro } k \in \{1, \dots, t\}, \\ \vec{w}_{k+t} &\stackrel{\text{def}}{=} \Lambda_e \vec{w}_k + \eta \vec{h}_e && k \geq 1 \end{aligned}$$

kde e je identická permutace. Potom řekneme, že tato posloupnost vznikla iterací

$$(\vec{x}_1, y_1), \dots, (\vec{x}_t, y_t)$$

podle DELTA PRAVIDLA .

Předchozí definice popisuje proces učení perceptronu na základě předložení posloupnosti $\vec{x}_1, \dots, \vec{x}_t$ vstupních vektorů a odpovídající posloupnosti y_1, \dots, y_t výstupních hodnot. Učení probíhá tak, že postupně cyklicky nastavujeme váhové parametry na základě vektoru \vec{x}_j a hodnoty y_j . Dále budeme předpokládat, že pořadí ve kterém předkládáme jednotlivé vektory \vec{x}_j je pevné, a v průběhu učení se nemění.

Tvrzení 6.1.4 Nechť posloupnost $\{\vec{w}_i\}_1^\infty$ vznikla iterací $(\vec{x}_1, y_1), \dots, (\vec{x}_t, y_t)$ podle delta pravidla. Předpokládejme, že $\vec{x}_1, \dots, \vec{x}_t$ generuje prostor \mathbb{R}^n a že vektor \vec{w}^* minimalizuje hodnotu funkce

$$\tilde{E}(\vec{w}) \stackrel{\text{def}}{=} \sum_{p=1}^t (y_p - \langle \vec{w} | \vec{x}_p \rangle)^2. \quad (6.2)$$

Potom jestliže pro všechna $p \in \{1, \dots, t\}$ platí $0 \leq \eta \leq \frac{2}{\|\vec{x}_p\|^2}$, konverguje posloupnost $\{\vec{w}_i\}_1^\infty$ ke konečnému cyklu $\{\vec{w}_i\}_1^t$ a platí, že \vec{w}_i je jediný pevný bod kontrahujícího zobrazení

$$\tilde{F}_i(\vec{w}) \stackrel{\text{def}}{=} \Lambda_{\pi_i} \vec{w} + \eta \vec{h}_{\pi_i},$$

kde permutace π_i je definována jako $(i+1, \dots, t, 1, \dots, i)$ (tedy π_0 je identická permutace).

Nechť $\vec{w}(\eta)$ je libovolný člen tohoto cyklu pro pevně zvolenou hodnotu parametru $\eta > 0$. Potom platí:

1. $\|\vec{w}(\eta) - \vec{w}^*\| = o(\eta)$
2. $|\tilde{E}(\vec{w}(\eta)) - \tilde{E}(\vec{w}^*)| = o(\eta^2)$.

■ *Důkaz:*

DŮKAZ JE UVEDEN V ORIGINÁLNÍ TIŠTĚNÉ VERZI.

■

Poznamenejme, že důkaz předešlých lemm a tvrzení lze modifikovat pro případ, kdy pořadí předkládání vektorů \vec{x}_i se v průběhu adaptace vektoru \vec{w} mění. Chování takové posloupnosti $\{\vec{w}_i\}_1^\infty$ se oproti původní posloupnosti změní v tom, že není zajištěna konvergence k nějakému cyklu. Obecně může tato posloupnost oscilovat. Nezmění se však charakter konvergence vzhledem k parametru η , jelikož i nadále bude platit

$$\limsup_{k \rightarrow \infty} \|\vec{w}_k(\eta) - \vec{w}^*\| = o(\eta).$$

Stejně tak bude i nadále platit

$$\limsup_{k \rightarrow \infty} \left| \tilde{E}(\vec{w}_k(\eta)) - \tilde{E}(\vec{w}^*) \right| = o(\eta^2).$$

Zřejmě i v případě, že vektory \vec{x}_i vybíráme pro nastavení vektoru \vec{w} podle rovnoměrného rozložení pravděpodobnosti, kdy každý z vektorů \vec{x}_i má stejnou pravděpodobnost výběru, je s vysokou mírou pravděpodobnosti zachována podmínka, že vektory použité pro učení jsou předkládány v po sobě následujících permutacích, a tedy výše popsané chování iterační posloupnosti bude s vysokou mírou pravděpodobnosti platné i v tomto případě.

Závěry odvozené v této kapitole lze přenést bezprostředně i na případ, kdy síť je tvořena více perceptronů, které jsou zařazeny vedle sebe, protože v průběhu učení se změny vah u jednoho perceptronu nijak neprojeví na váhách jiných perceptronů. Tvrzení 6.1.4 tedy popisuje i konvergenci vah pro dvouvrstvou architekturu sítí.

6.1.2 Diskrétní delta pravidlo

Nyní zkoumejme diskrétní případ předešlého algoritmu. I když se v tomto případě již nejedná o gradientní algoritmus, uvádíme tento algoritmus právě zde z důvodů podobnosti s původním delta algoritmem.

Definice 6.1.3 *Předpokládejme, že $(\vec{x}_1, y_1), \dots, (\vec{x}_t, y_t)$ je zadaná posloupnost uspořádaných dvojic z $\mathcal{R}^n \times \{-1, +1\}$, $t \geq 1$. Nechť pro dané $\eta > 0$ je posloupnost vektorů $\{\vec{w}_i\}_1^\infty$ definovaná následujícím rekurzivním předpisem*

1. položme $\vec{w}_1 \stackrel{\text{def}}{=} \vec{0}$

2. zvolme index $j_k \in \{1, \dots, t\}$ tak, aby platilo $\widehat{\text{sgn}}(\langle \vec{w}_k | \vec{x}_{j_k} \rangle) \neq y_{j_k}$ a dále:

- (a) neexistuje-li takovéto j_k , položme $\vec{w}_{k+1} = \vec{w}_k$,

- (b) existuje-li takovéto j_k , položme

$$\vec{w}_{k+1} \stackrel{\text{def}}{=} \vec{w}_k + \eta \vec{x}_{j_k} (\widehat{\text{sgn}}(\langle \vec{w}_k | \vec{x}_{j_k} \rangle) - y_{j_k}).$$

Potom řekneme, že tato posloupnost vznikla iterací $(\vec{x}_1, y_1), \dots, (\vec{x}_t, y_t)$ podle DISKRÉTNÍHO DELTA PRAVIDLA .

Základní vlastnost diskrétního delta pravidla popisuje následující tvrzení ([ŠN96]):

Tvrzení 6.1.5 *Nechť posloupnost $\{\vec{w}_i\}_1^\infty$ vznikla iterací podle diskrétního delta pravidla a nechť dále existuje $\widehat{\vec{w}}$ takové, že pro všechna $i \in \{1, \dots, t\}$ platí $\widehat{\text{sgn}}(\langle \widehat{\vec{w}} | \vec{x}_i \rangle) = y_i$. Potom existuje přirozené číslo $z > 0$, pro které $\vec{w}_{z+1} = \vec{w}_z$, a navíc platí odhad*

$$z \leq \frac{\|\widehat{\vec{w}}\|_\alpha^2}{4\eta^2\beta^2} + 1 ,$$

kde je

$$\alpha \stackrel{\text{def}}{=} \max_{i \in \{1, \dots, t\}} \left\{ \|\vec{x}_i\|^2 \right\}$$

a

$$\beta \stackrel{\text{def}}{=} \min_{i \in \{1, \dots, t\}} \left\{ \left| \langle \widehat{\vec{w}} | \vec{x}_i \rangle \right| \right\}.$$

■ *Důkaz:*

DŮKAZ JE UVEDEN V ORIGINALNÍ TIŠTĚNÉ VERZI. ■

6.1.3 Rozšířené delta pravidlo a SVD rozklad

V obou případech diskrétního i obyčejného delta pravidla se hodnota váhového vektoru měnila vždy v závislosti na chybě aproximace pro jeden vstupní vektor. Existují ale i jiné přístupy ke způsobu nastavování váhového vektoru v průběhu učení. Jedním z nich, který se často zmiňuje, je algoritmus, při kterém se hodnoty vah nastavují v závislosti na všech chybách přes všechny vstupní vektory. Tento algoritmus adaptuje váhový vektor podle následující iterační formule:

$$\vec{w}_{k+1} \stackrel{\text{def}}{=} \vec{w}_k - \eta \sum_{p=1}^t (\vec{x}_p \vec{x}_p^T) \vec{w}_k + \eta \sum_{p=1}^t y_p \vec{x}_p \stackrel{\text{def}}{=} (\mathbf{I} - \eta \mathbf{L}) \vec{w}_k - \eta \sum_{p=1}^t y_p \vec{x}_p. \quad (6.3)$$

Zavedeme-li označení $\Omega \stackrel{\text{def}}{=} \mathbf{I} - \eta \mathbf{L}$, lze předchozí rovnici psát ve tvaru

$$\vec{w}_{k+1} = \Omega \vec{w}_k - \eta \sum_{p=1}^t (y_p \vec{x}_p). \quad (6.4)$$

Proveďme si nyní analýzu takovéto iterační posloupnosti. Dosadíme-li do rovnice 6.3 $\vec{w}_{k+1} = \vec{w}_k = \vec{w}$, vidíme (viz. ??), že řešení této rovnice je rovno vektoru, který minimalizuje hodnotu chybové funkce 6.2. Rovnice 6.3 vlastně představuje iterační posloupnost pro metodu nejmenších čtverců. Matice \mathbf{L} je symetrická a pozitivně semidefinitní. Budeme-li navíc předpokládat, že $\{\vec{x}_p\}_1^t$ generuje prostor \mathfrak{R}^n , je matice \mathbf{L} pozitivně definitní, viz. lemma 6.1.6. Všechna vlastní čísla matice Ω jsou ve tvaru $1 - \eta\lambda$, kde λ je vlastní číslo matice \mathbf{L} . Tedy pro dostatečně malé η je spektrální poloměr a odtud i indukovaná norma (Ω je symetrická) menší než 1 a iterační posloupnost 6.4 bude tedy konvergentní na základě věty o kontrahujícím zobrazení.

Ukážeme ale, že tato konvergence obecně neprobíhá dostatečně rychle, a to i v případech separace konvexních množin, které jsou od sebe velmi vzdáleny v porovnání s jejich rozměry. Abychom demonstrovali tento fakt, uveďme následující odhady velikosti maximálního vlastního čísla matice \mathbf{L} .

Lemma 6.1.6 *Budiž $\vec{x}_1, \dots, \vec{x}_t \in \mathfrak{R}^n$ a necht' λ je maximální vlastní číslo matice*

$$\mathbf{L} \stackrel{\text{def}}{=} \sum_{p=1}^t \vec{x}_p \vec{x}_p^T.$$

Potom platí odhad

$$\max \left\{ \|\vec{x}_p\|^2 \mid p \in \{1, \dots, t\} \right\} \leq \lambda \leq \sum_{p=1}^t \|\vec{x}_p\|^2.$$

■ *Důkaz:*

DŮKAZ JE UVEDEN V ORIGINÁLNÍ TIŠTĚNÉ VERZI. ■

Předpokládejme, že \vec{u} a \vec{v} jsou dva ortonormální vektory v \mathfrak{R}^n a že pro $\epsilon > 0$ a $z \geq 1$ máme dány libovolné množiny \bar{A} a \bar{B} ,

$$\bar{A} \stackrel{\text{def}}{=} \left\{ \vec{x}_p \in \mathfrak{R}^n, p \in \{1, \dots, z\} \mid \|\vec{x}_p - \vec{u}\| \leq \epsilon \right\}$$

a

$$\bar{B} \stackrel{\text{def}}{=} \left\{ \vec{y}_p \in \mathfrak{R}^n, p \in \{1, \dots, z\} \mid \|\vec{y}_p - \vec{v}\| \leq \epsilon \right\}.$$

Zřejmě pro malé hodnoty ϵ není problém množiny \bar{A} a \bar{B} od sebe lineárně odseparovat. Označme vlastní čísla matice \mathbf{L} jako $0 \leq \lambda_n \leq \lambda_{n-1} \leq \dots \leq \lambda_1$. Podle odhadů z lemma 6.1.6 je ale zřejmé, že $\lambda_1 \geq (1 - \epsilon)^2$ a současně matice \mathbf{L} se pro $\epsilon \rightarrow 0$ blíží matici, jejíž hodnota je rovna jedné. Tedy pro $j \geq 2$ je $\lim_{\epsilon \rightarrow 0} \lambda_j = 0$. Jinými slovy, matice \mathbf{L} je špatně podmíněná a to tak, že všechna vlastní čísla s výjimkou maximálního (které se blíží k 1, $\epsilon \rightarrow 0$) se stále více přibližují k nule. Protože $\mathbf{\Omega} = \mathbf{I} - \eta\mathbf{L}$, vlastní čísla matice $\mathbf{\Omega}$ jsou tvaru $1 - \eta\lambda_n \geq 1 - \eta\lambda_{n-1} \geq \dots \geq 1 - \eta\lambda_1$. To je ale pro konvergenci ten nejhorší možný případ, protože pro libovolnou počáteční iteraci se každá iterační posloupnost velmi rychle přiblíží k lineárnímu obalu vlastního vektoru příslušejícímu maximálnímu vlastnímu číslu (rychlost tohoto přiblížení je daná velikostí ostatních vlastních čísel, které jsou velmi malé, lineární operátor $\mathbf{\Omega}$ zúžený na ortogonální doplněk vlastního vektoru pro $1 - \eta\lambda_n$ je kontrahující s velmi malým kontrahujícím faktorem) a výsledná rychlost konvergence de fakto závisí pouze na maximálním vlastním čísle, které je poblíž hodnoty $1 - \eta$.

Zřejmě tedy individuální korekce vah podle jednotlivých vektorů použitých k učení, jako je tomu v případě původního delta pravidla, je dominantní pro efektivitu učícího algoritmu.

SVD rozklad a jev přeučení

Přestože rozšířené delta pravidlo je z numerického hlediska mnohem méně vhodné pro učení nežli původní delta pravidlo, použijeme nyní SVD (anglický termín: singular value decomposition) rozkladu k analýze jevu přeučení, který při aplikaci tohoto pravidla nastává. Důvod, proč tuto analýzu neprovedeme pro původní delta pravidlo, je ten, že tato analýza by byla technicky mnohem obtížnější. Jev přeučení spočívá v tom, že schopnost sítě generalizovat data se mnohdy s pokračující minimalizací chybové funkce snižuje, přestože vektor vah konverguje k nějaké limitní množině.

Vezměme v úvahu iterační posloupnost 6.4. Aplikujme na matici \mathbf{X} SVD rozklad, tedy vyjádříme matici \mathbf{X} ve tvaru

$$\mathbf{X} = \mathbf{P}\mathbf{S}\mathbf{Q}^T,$$

kde \mathbf{P} a \mathbf{Q} jsou ortogonální matice (mohou být různého řádu) a matice \mathbf{S} je obecně obdélníková diagonální matice, na jejíž diagonále jsou singulární čísla matice \mathbf{X} (což jsou vlastní čísla matice $\mathbf{X}^T\mathbf{X}$). Dosadíme-li tuto dekompozici do vztahu 6.4, obdržíme

$$\begin{aligned} \vec{w}_{k+1} &= (\mathbf{I} - \eta\mathbf{P}\mathbf{S}\mathbf{S}^T\mathbf{P}^T) \vec{w}_k - \eta\mathbf{X}\vec{y} \\ &= \mathbf{P}(\mathbf{I} - \eta\mathbf{S}\mathbf{S}^T) \mathbf{P}^T \vec{w}_k - \eta\mathbf{P}\mathbf{S}\mathbf{Q}^T \vec{y}. \end{aligned}$$

Po zavedení substituce $\bar{\mathbf{z}}_k \stackrel{\text{def}}{=} \mathbf{P}^T \bar{\mathbf{w}}_k$ a $\bar{\mathbf{u}} \stackrel{\text{def}}{=} \mathbf{P}^T \bar{\mathbf{g}}$ lze tento výraz přepsat do tvaru

$$\bar{\mathbf{z}}_{k+1} = (\mathbf{I} - \eta \mathbf{S} \mathbf{S}^T) \bar{\mathbf{z}}_k - \eta \mathbf{S} \bar{\mathbf{u}}.$$

Předpokládejme, že \mathbf{X} má r nenulových singulárních čísel, pro která platí $\nu_1 \geq \nu_2 \geq \dots \geq \nu_r$. Potom lze složky vektoru $\bar{\mathbf{z}}_{k+1}$ zapsat ve tvaru

$$(\bar{\mathbf{z}}_{k+1})_i = (1 - \eta \nu_i^2) (\bar{\mathbf{z}}_k)_i - \eta \nu_i \bar{\mathbf{u}}_i$$

pro všechny indexy $i \in \{1, \dots, r\}$ a

$$(\bar{\mathbf{z}}_{k+1})_i = (\bar{\mathbf{z}}_k)_i$$

pro indexy $r + 1 \leq i \leq n$.

Je-li η dostatečně malé na to, aby zaručovalo konvergenci (t.j. aby $0 < 1 - \eta \nu^2 < 1$), je zřejmé, že konvergence bude mnohem rychlejší pro taková $\bar{\mathbf{z}}_i$, která odpovídají velkým singulárním hodnotám. Takovéto složky vektorů $\bar{\mathbf{z}}$ jsou pro charakterizaci dat významné a v průběhu iterování zkonvergují nejdříve. Složky vektoru $\bar{\mathbf{z}}$, které korespondují s malými singulárními čísly, naopak konvergují pomalu, ale mohou významně přispívat k velikosti chybové funkce. Současně nejsou pro dané vektory $\bar{\mathbf{x}}_1, \dots, \bar{\mathbf{x}}_t$ signifikantní. V průběhu iterování se váhy nastaví nejdříve tak, aby daná síť separovala dominantní vlastnosti učicích vzorů, při dalším průběhu se síť snaží separovat nevýznamné vlastnosti, eventuelně šumy.

6.1.4 Zobecnění na nelineární systémy

Efektivní aplikovatelnost různých obměn delta pravidla je omezena pouze na případy, kdy separované množiny jsou lineárně separovatelné. To vyplývá z toho, že výchozím bodem pro návrhy učicích gradientních algoritmů založených na delta pravidle, je vztah 6.1. Tento výraz ale vyjadřuje změnu vah takovým způsobem, aby byla minimalizována chyba parametrizovaná polohou nějaké nadroviny, oddělující body patřící do separovaných množin. Budeme-li brát v úvahu takové neuronové sítě, jejichž dělicí plochy nejsou nadroviny ale plochy vyšších řádů, je potřeba adaptovat pravidlo 6.1 tak, aby změna vah byla úměrná chybě daného vzoru vzhledem k dělicím plochám, které odpovídají uvažovanému typu neuronových sítí. Proto nyní popíšeme způsob zobecnění delta pravidla na nelineární případ a ukážeme, že lze obdobným způsobem jako v lineárním případě analyzovat konvergentní vlastnosti iterační posloupnosti v okolí stabilních bodů, kterých ale v nelineárním případě může být (a zpravidla i je) více. Nechť M označuje celkový počet vah dané sítě, n je dimenze vstupního vektoru. Obecně síť realizuje zobrazení $\tilde{g} : \mathbb{R}^M \times \mathbb{R}^n \rightarrow \mathbb{R}^m$, kde m je dimenze výstupních vektorů. Pro daný vstupní vektor $\bar{\mathbf{x}}_p$ označme výstupní vektor $\bar{\mathbf{v}}_p = \tilde{g}(\bar{\mathbf{x}}_p, \bar{\mathbf{w}})$. Budeme předpokládat, že \tilde{g} je diferencovatelné a nechť jeho derivace je označena jako $\widetilde{\mathbf{D}}(\bar{\mathbf{x}}, \bar{\mathbf{w}})$. Ze spojitosti \tilde{g} plyne po aplikaci Tylerova rozvoje, že

$$\tilde{g}(\bar{\mathbf{w}} + \delta \bar{\mathbf{w}}, \bar{\mathbf{x}}) = \tilde{g}(\bar{\mathbf{w}}, \bar{\mathbf{x}}) + \widetilde{\mathbf{D}}(\bar{\mathbf{w}}, \bar{\mathbf{x}}) \delta \bar{\mathbf{w}} + o(\|\delta \bar{\mathbf{w}}\|).$$

Pro pevný vstupní vektor $\bar{\mathbf{x}}_p$ nechť $\bar{\mathbf{y}}_p$ označuje požadovaný výstupní vektor. Potom označme chybu aproximace tohoto výstupního vektoru jako

$$\epsilon_p^2 \stackrel{\text{def}}{=} (\bar{\mathbf{y}}_p - \bar{\mathbf{v}}_p)^T (\bar{\mathbf{y}}_p - \bar{\mathbf{v}}_p) \stackrel{\text{def}}{=} \bar{\mathbf{q}}_p^T \bar{\mathbf{q}}_p$$

a celkovou chybovou funkcí naučení vzorů $\vec{x}_1, \dots, \vec{x}_t$ jako

$$\tilde{E}(\vec{w}) \stackrel{\text{def}}{=} \sum_{i=1}^t \epsilon_i^2.$$

Zřejmě při změně váhového vektoru o hodnotu $\delta\vec{w}$ se změní chyba pro \vec{x}_p o hodnotu

$$\delta\epsilon_p^2 = (\vec{q}_p + \delta\vec{q}_p)^T (\vec{q}_p + \delta\vec{q}_p) - \vec{q}_p^T \vec{q}_p = 2 (\delta\vec{q}_p)^T \vec{q}_p + (\delta\vec{q}_p)^T \delta\vec{q}_p.$$

Protože \vec{y}_p je konstantní, platí

$$\delta\vec{q}_p = -\delta\vec{v}_p = -\tilde{D}(\vec{w}, \vec{x}_p)\delta\vec{w} + o(\|\delta\vec{w}\|).$$

Tedy

$$\begin{aligned} \delta\epsilon_p^2 &= -2 (\tilde{D}(\vec{w}, \vec{x}_p)\delta\vec{w})^T (\vec{y}_p - \tilde{g}(\vec{w}, \vec{x}_p)) + o(\|\delta\vec{w}\|) \\ &= -2 (\delta\vec{w})^T (\tilde{D}(\vec{w}, \vec{x}_p))^T (\vec{y}_p - \tilde{g}(\vec{w}, \vec{x}_p)) + o(\|\delta\vec{w}\|). \end{aligned} \quad (6.5)$$

Člen $o(\|\delta\vec{w}\|)$ můžeme vynechat a potom je podle 6.5 pro fixovanou hodnotu $\|\delta\vec{w}\|$ dosaženo maximálního úbytku hodnoty $\delta\epsilon_p^2$, právě když bude platit

$$\delta\vec{w} = \eta (\tilde{D}(\vec{w}, \vec{x}_p))^T (\vec{y}_p - \tilde{g}(\vec{w}, \vec{x}_p)). \quad (6.6)$$

Tato formule představuje zobecněné delta pravidlo, čehož lze snadno nahlédnout, položíme-li $\tilde{g}(\vec{w}, \vec{x}) = \vec{w}^T \vec{x}$. Po dosazení do 6.6 dostaneme původní delta pravidlo popsané výrazem 6.1.

Na základě vzorce 6.6 dostáváme iterační posloupnost pro nelineární delta pravidlo ve tvaru

$$\begin{aligned} \vec{w}_{k+1} &= \vec{w}_k + \delta\vec{w}_k \\ &= \vec{w}_k + \eta \tilde{D}(\vec{w}_k, \vec{x}_p)^T (\vec{y}_p - \tilde{g}(\vec{w}_k, \vec{x}_p)). \end{aligned}$$

Oproti lineárnímu případu, nemá chybová funkce $\tilde{E}(\vec{w})$ obecně pouze jedno globální minimum, ale může mít lokální minima, což je obecně potvrzeno (např. při učení metodou zpětného šíření). Zaměřme se proto na chování posloupnosti \vec{w}_k poblíž některého lokálního minima. Budeme-li předpokládat spojitost a omezenost derivace \tilde{D} na nějakém okolí lokálního minima $\widehat{\vec{w}}$ (což je zaručeno např. požadavkem, aby zobrazení \tilde{g} mělo spojitě derivace do druhého řádu včetně), lze psát

$$\begin{aligned} \vec{w}_{k+1} &= \vec{w}_k + \eta (\tilde{D}(\vec{w}_k, \vec{x}_p))^T (\vec{y}_p - \tilde{g}(\widehat{\vec{w}}, \vec{x}_p) - \tilde{D}(\widehat{\vec{w}}, \vec{x}_p) (\vec{w}_k - \widehat{\vec{w}})) + o(\|\vec{w}_k - \widehat{\vec{w}}\|) \\ &= \left(\mathbf{I} - \eta (\tilde{D}(\vec{w}_k, \vec{x}_p))^T \tilde{D}(\widehat{\vec{w}}, \vec{x}_p) \right) \vec{w}_k + \eta (\tilde{D}(\vec{w}_k, \vec{x}_p))^T \times \\ &\quad \times (\vec{y}_p - \tilde{g}(\widehat{\vec{w}}, \vec{x}_p) + \tilde{D}(\widehat{\vec{w}}, \vec{x}_p)\widehat{\vec{w}}) + o(\|\vec{w}_k - \widehat{\vec{w}}\|). \end{aligned}$$

Iterační matice $\left(\mathbf{I} - \eta (\tilde{D}(\vec{w}_k, \vec{x}_p))^T \tilde{D}(\widehat{\vec{w}}, \vec{x}_p) \right)$ není obecně symetrická, ale pro $\vec{w}_k \rightarrow \widehat{\vec{w}}$ se symetrické matici blíží. Dále budeme předpokládat, že matice $\tilde{D}(\vec{w}_k, \vec{x}_p)$ je Lipstichovská vzhledem k proměnné \vec{w} a stejnoměrně spojitá vzhledem k proměnné \vec{x} na nějakém okolí bodu $\widehat{\vec{w}}$. Za těchto předpokladů platí

$$\begin{aligned} \vec{w}_{k+1} &= \left(\mathbf{I} - \eta (\tilde{D}(\widehat{\vec{w}}, \vec{x}_p))^T \tilde{D}(\widehat{\vec{w}}, \vec{x}_p) \right) \vec{w}_k + \eta (\tilde{D}(\widehat{\vec{w}}, \vec{x}_p))^T \times \\ &\quad \times (\vec{y}_p - \tilde{g}(\widehat{\vec{w}}, \vec{x}_p) + \tilde{D}(\widehat{\vec{w}}, \vec{x}_p)\widehat{\vec{w}}) + o(\|\vec{w}_k - \widehat{\vec{w}}\|). \end{aligned}$$

Předpokládejme, že obdobně jako v lineárním případě i zde iterujeme cyklicky přes vektory $\vec{x}_1, \dots, \vec{x}_t$. Ukážeme, že i v tomto případě pro malé hodnoty η je lineární část (bez členu $o(\|\vec{w}_k - \widehat{w}\|)$) zobrazení převádějící \vec{w}_k na \vec{w}_{k+t} kontrahující. Předně pro $p \in \{1, \dots, t\}$ je matice

$$\widetilde{D}(\widehat{w}, \vec{x}_p)^T \widetilde{D}(\widehat{w}, \vec{x}_p) \quad (6.7)$$

pozitivně semidefinitní. Lze tedy volit η dostatečně malé tak, aby platilo

$$\left\| I - \eta \widetilde{D}(\widehat{w}, \vec{x}_p)^T \widetilde{D}(\widehat{w}, \vec{x}_p) \right\| \leq 1.$$

Protože matice 6.7 je pozitivně semidefinitní, můžeme rozložit prostor \mathfrak{R}^n vstupních vektorů na direktní součet nulového prostoru operátoru daného maticí 6.7 a jeho ortogonálního doplňku. Složka vektoru \vec{w}_{k+1} ležící v nulovém prostoru je totožná s odpovídající složkou vektoru \vec{w}_k . Stačí nám tedy ukázat, že operátor $I - \eta \widetilde{D}(\widehat{w}, \vec{x}_p)^T \widetilde{D}(\widehat{w}, \vec{x}_p)$ zúžený na zmíněný ortogonální doplněk je kontrahující. K tomu stačí zvolit η tak, aby platilo

$$\eta < \rho \left(\widetilde{D}(\widehat{w}, \vec{x}_p)^T \widetilde{D}(\widehat{w}, \vec{x}_p) \right)^{-1}.$$

Tato volba zaručí, že v okolí lokálního minima chybové funkce $\widetilde{E}(\vec{w})$ se bude iterační posloupnost chovat obdobným způsobem, jako v lineárním případě, který je popsán na začátku této kapitoly.

Část II

Relevantní výsledky související s neuronovými sítěmi

(AUTOR MARTIN HOLEŇA)

Celá druhá část je uvedena v originální
tištěné verzi.

Literatura

- [AB92] Martin Anthony and Norman Biggs. *Computational Learning Theory*. Press Syndicate of the University of Cambridge, 1992.
- [BEHW89] A. Blumer, A. Ehrenfeucht, D. Haussler, and M. K. Warmuth. Learnability and the Vapnik-Chervonenkis Dimension. *Journal of the Association for Computing Machinery*, 36:929–965, oct 1989.
- [Blo93] Arthur Bloch. *Murphyho zákon*. Nakladatelství svoboda, Libertas, Praha, 1993.
- [Ell94] S. W. Ellacott. Aspects of the numerical analysis of neural networks. *Acta Numerica*, pages 145–202, 1994.
- [Fie81] M. Fiedler. *Speciální matice a jejich použití v numerické matematice*. Edice Teoretická knižnice inženýra. Státní nakladatelství technické literatury, 1981.
- [GJ79] M. R. Garey and D. S. Johnson. *Computers and Intractability*. W.H.Freeman & Company, 1979.
- [Hås94] Johan Håstad. On the size of weights for threshold gates. *Siam Journal Discrete Mathematics*, 7(3):484–492, aug 1994.
- [Jar55] V. Jarník. *Integrální počet II*. NČAV Praha, 1955.
- [Joh74] David S. Johnson. Approximation algorithms for combinatorial problems. *Journal of Computer and System Sciences*, 9(3):256–278, dec 1974.
- [KA77] L. V. Kantorovič and G. P. Akilov. *Funkcionalnyj analiz*. Nauka Moskva, 1977.
- [Kol57] A. N. Kolmogorov. O predstavlenii nepreryvnykh funkciy neskol'kikh peremennykh v vide superpozicij nepreryvnykh funkciy odnogo peremennogo i slozenija. *Doklady akademii nauk SSSR*, 114(5):953–956, 1957.
- [Lej85] K. Lejchtvějs. *Vypuklité množestva*. Moskva Nauka, 1985.
- [Miš89] Ladislav Mišík. *Funkcionálna analýza*. ALFA Bratislava, 1989.
- [Mur71] Saburo Muroga. *Threshold Logic and Its Applications*. John Wiley & Sons, 1971.
- [RSO94] Vwani Roychowdhury, Kai-Yeung Siu, and Alon Orlitsky. *Theoretical Advances in Neural Computation and Learning*. Kluwer Academic Publishers, 1994.

- [Sau72] N. Sauer. On the Density of Families of Sets. *Journal of the Association for Computing Machinery*, 13:145–147, feb 1972.
- [ŠJ96] Miroslav Šnorek and Marcel Jiřina. *Neuronové sítě a neuropočítače*, chapter Kapitola 2. ČVUT, 1996.
- [Smi65] David R. Smith. Bounds on the number of threshold functions. *IEEE Transactions on Electronic Computers*, pages 368–369, 1965.
- [ŠN96] Jiří Šíma and Roman Neruda. *Teoretické otázky neuronových sítí*. MATFY-ZPRESS, MFF UK Praha, 1996.
- [Spr93] David A. Sprecher. A universal mapping for Kolmogorov’s superposition theorem. *Neural Networks*, 6:1089–1094, 1993.
- [YI65] S. Yajima and T. Ibaraki. A lower bound of the number of threshold functions. *IEEE Transaction on Electronic Computers*, pages 926–929, 1965.

Rejstřík

- $\widetilde{\mathbf{LT}}_{d,\tilde{p}}$, 35
 \mathbf{LT}_d , 35
 $\mathbf{PT}_{1,\tilde{p}}$, 36
 \mathbf{PT}_1 , 36
 $\widetilde{\mathbf{LT}}_d$, 36
 $\widetilde{\mathbf{PT}}_1$, 36
 $\widetilde{G}_\phi(\vec{x})$, 61
 $\vec{p}_i^{(n)}$, 32
 α -symetrický vektor, 36
 \bar{I}_k , 56
 \bar{J}_ϵ^{2m} , 77
 $\bar{Q}_{m,\epsilon}$, 77
 $\mathbf{EF}_{A^*,m,s}$, 85
 \mathbf{l}_C , 87
 $\mathbf{l}_{k,C}$, 87
 $\mathbf{R}_{\tilde{P},\epsilon}$, 76
 \mathbf{U}_C , 87
 $\mathbf{U}_{k,C}$, 87
 $\mathbf{Y}_{m,s}$, 85
 \mathbf{Z}_s , 85
 ϵ -transversála, 76
 (ϵ, δ) -učicí algoritmus, 75
 \bar{P}^m , 77
 $\mathbf{S}_{l,\bar{V}}$, 71
 d -polynomiálně naučitelný systém konceptů, 81
 l -rozlišitelná hranová ohodnocení, 71
 s, m -efektivní prostor, 85
 $\mathbf{HALFSPACE}_n$, 66
 \mathbf{SYM}_2 , 36
 $\mathbf{e}_{\tilde{P}}(\bar{c}, \bar{h})$, 74
 $\mathbf{I}^{(n)}$, 33
 $\mathbf{E}^{(n)}$, 33
 $\omega(\tilde{f}, \delta)$, 48
 $\Pi_C(\bar{S})$, 68
 $\mathbf{M}_{\tilde{y}}$, 32
 $\mathbf{B}^{(n)}$, 24
 $\mathbf{VC}_{dim}(C)$, 66
C rozděljuje \bar{Y} , 66
 \tilde{P} -učicí algoritmus složitosti $\tilde{m}(\epsilon, \delta)$, 75, 15
řady matice, 14
back-propagation, 17, 89
booleovský obvod, 34
Borelovská množina, 76
celočíslně nazávislá posloupnost, 56
chyba hypotézy \bar{h} vzhledem ke konceptu \bar{c} , 74
délka aproximace optimálního pokrytí, 83
delta pravidlo, 91
dimenze booleovského obvodu, 35
diskrétní delta pravidlo, 92
dobře utvořená třída konceptů, 77
dopředný obvod, 70
efektivní prostor hypotéz, 85
funkce parity, 24
Hadamardova matice, 24
hloubka booleovského obvodu, 35
hloubka vrcholu, 35
jádro prahového vektoru, 41
koncept nad množinou \bar{X} , 65
konceptuální míra složitosti, 85
konceptuální složitost učicího algoritmu, 85
konvoluční jádro, 61
konzistentní množina, 74
lineárně uspořádaná třída konceptů, 70
lineární sigmoidální obvod, 72
lokální třída konceptů, 71
matice parity, 24
maximum spektra, 30
minimální pokrytí množiny, 84

- modul spojitosti, 48
- modulus of continuity, 48
- nedeterministický polynomiální generátor hypotéz, 81
- Occamův algoritmus, 86
- PAC learning, 15
- polynomiálně definovaná třída hypotéz, 81
- polynomiálně definovaná třída konceptů, 80
- polynomiálně definovaný systém hypotéz, 81
- polynomiálně definovaný systém konceptů, 80
- polynomiální generátor hypotéz, 87
- práh ostražitosti, 21
- práh vektorů, 28
- Probably Approximately Correct, 75
- problém konzistence, 87
- prostor vzorů systému konceptů, 80
- prostor vzorů třídy konceptů, 74
- Radial Basis Function, 52
- RBF funkce, 52
- reprezentace třídy konceptů, 66
- s-polynomiálně naučitelná třída konceptů, 85
- sdružená matice, 32
- semilokální funkce, 53
- sigmoidální funkce, 17
- singular value decomposition, 94
- skalární součin matic, 27
- systém hypotéz, 80
- systém konceptů, 80
- třída hypotéz, 74
- třída konceptů, 65
- třída konceptů dopředného obvodu, 71
- tenzorový součin matic, 23
- totálně naučitelný systém konceptů, 81
- triviální třída konceptů, 79
- uniformní naučitelnost, 75
- univerzálně separabilní třída hypotéz, 77
- váha hrany, 34, 70
- výstupní vrchol, 34, 70
- Vapnik-Chervonenkova dimenze, 66
- VC-dimenze dopředného obvodu, 71
- velikost booleovského obvodu, 35
- vlastní očíslování grafu, 70
- vnitřní vrchol, 34, 70
- vstupní vrchol, 34, 70
- vzorek konceptu \bar{c} délky m , 74
- vzorová složitost, 75, 81
- well-behaved, 77
- základní vektory parity, 32
- zobecněné spektrum vektoru, 30