

# Posudek oponenta diplomové práce

Jméno a příjmení autora posudku: doc. RNDr. Iveta Mrázová, CSc.

Jméno a příjmení autora práce: Bc. Matěj Sochor

Název práce: **Semi-supervised Learning from Unfavorably Distributed Data**

Vlastní text (sem prosím napište text posudku, délka textu posudku není omezena):

Diplomová práce Bc. Matěje Sochora se věnuje problematice částečně řízeného učení (semi-supervised learning). Tato oblast patří k mimořádně aktuálním. Moderní paradigma tzv. hlubokého učení totiž vyžaduje apriorní znalost obrovského množství značkových (labeled) dat. Důvodem je nutnost nastavit správně velké množství parametrů učeného systému tak, aby mohl správně zobecňovat extrahované znalosti. Získání velkého množství značkových dat ovšem bývá extrémně náročné, a to jak vzhledem k finančním nárokům, tak také vzhledem k nárokům kladeným na lidské zdroje. Trendem je tedy možnost využít při učení kromě dostupných označkových dat i data neoznačovaná, která je ovšem nutné vhodným způsobem (před)zpracovat. Neoznačkových dat bývá na druhou stranu k dispozici dostatečné množství a jejich správné využití může i výrazně zvýšit přesnost příslušného klasifikátoru.

Cílem předkládané práce proto bylo prostudovat dostupnou literaturu o částečně řízeném učení. Na základě provedené rešerše měl student vytipovat možné problémy spojené s charakterem použitých dat, které by bránily širšímu využití částečně řízeného učení v praxi. Velmi problematické by např. mohlo být použití dat odlišného typu, resp. třídy v rámci označkových a neoznačkových vzorků. Takový případ by dokonce mohl vést i k horším výsledkům, než kdybychom neoznačovaná data nepoužili vůbec. Diplomant měl proto navrhnout nové metody nebo vylepšení stávajících metod částečně řízeného učení tak, aby byly robustnější vzhledem k rozdílným distribucím označkových a neoznačkových dat. Navržené metody, resp. vylepšení měl uchazeč následně porovnat s existujícími metodami částečně řízeného učení.

Těžiště práce spočívá v návrhu původní metody nazvané Unfavorable Data Filtering (UDF) pro extrakci vzorků vhodných pro částečně řízené učení. Technika elegantním způsobem využívá model autoenkodérů. Neoznačovaná data je tak možné využít i k odhadu vnitřní struktury zpracovávaných dat v příznakovém prostoru. Vnitřní struktura (vhodně transformovaných) dat totiž může mít mnohem nižší dimenzi než původní příznakový prostor. Pravděpodobnostní rozložení příznaků neoznačkových dat zároveň může (ale také nemusí) mít vztah k třídě, která odpovídá vzorkům s příslušnými hodnotami příznaků. Aby mělo částečně řízené učení smysl, měla by struktura dat zhruba odpovídat výsledku jejich klastrování. Součástí navrhované strategie je proto i klastrování (před)zpracovaných dat. V rámci testování nového přístupu se studentovi podařilo jak nastínit jeho limity, tak také ukázat, že pro jistá nastavení standardně dosahované výsledky částečně řízeného učení zlepšuje.

Práce sama je napsaná v angličtině a jisté rezervy má zejména z formálního hlediska. K ověření vlastností navrhované metody autor provedl mnoho experimentů, některé z nich by ovšem bylo vhodné podrobněji vyhodnotit. Této skutečnosti si je však autor vědom a ve své práci ji zmiňuje. Lépe zdůvodnit by bylo vhodné volbu některých kroků a nastavení, např. počet neuronů pro hrdlo použitého autoenkodéru. Autor by měl lépe citovat použité prameny. Z textu uvedeném ve 2. odstavci na str. 4 např. není jasné, která pozorování budou autorova a která převzatá, navíc z jakých pramenů (pohromadě jsou tu uvedeny tři různé). Reference A. Krizhevsky na str. 45 a některé další nejsou úplné. Většina z pojmů a použitých postupů je v práci popsána slovně, a to i tam, kde by jejich formální zavedení (definice, vztah apod.) mohlo výrazně přispět ke srozumitelnosti textu. Srozumitelnosti práce by ovšem prospělo i členění textu do vět a odstavců adekvátní délky (výjimkou nejsou věty přes 5 řádků a odstavce o více než 20 řádcích, např. na str. 23 anebo na str. 25, resp. na str. 20). Příliš dlouhý je i popis k obrázku 3.2 uvedený v seznamu na str. 48. Postup znázorněný v diagramu 3.3 na str. 29 by bylo vhodnější formulovat spíše formou pseudokódu. V českém abstraktu práce jsou hrubé gramatické chyby a práce obsahuje i velké množství typografických chyb, viz např. vdovy a sirotci na str. 6, 19, 20, 42, 44.

Hlavní přínos práce tedy vidím v návrhu původní metody UDF pro eliminaci irelevantních trénovacích vzorků, její implementaci a ověření jejích vlastností v rámci částečně řízeného učení. I přes výše zmíněné nedostatky, které by byl diplomant jistě schopen snadno odstranit, předkládaná práce splňuje svůj původní cíl. Diplomant pronikl dostatečně hluboko do problematiky umělých neuronových sítí a strojového učení. Přítom prokázal schopnost samostatně řešit i poměrně náročné úlohy z oblasti analýzy dat. Podrobná diskuse ohledně volby navrhovaného řešení a možného testování jejích vlastností svědčí o orientaci autora v řešené problematice. Vlastní výsledky a zkušenosti dokázal uchazeč objektivně vyhodnotit.

## Doporučení k obhajobě:

Z výše uvedených důvodů práci doporučuji k obhajobě.

Vynikající práce vhodná pro soutěž studentských prací	NE [x]
---	--------

Seznam soutěží studentských prací, viz <http://www.mff.cuni.cz/studium/bcmgr/prace/>

Pokud jste výše zaškrtnli ANO, zdůvodněte prosím svůj návrh, případně uveďte konkrétní soutěž, pro kterou je práce vhodná (rámeček lze nechat prázdný, pokud za dostatečné zdůvodnění považujete text posudku):

V Praze dne: 29. 6. 2020

Podpis:\*\*

\* *nehodící se škrtněte (vymažte)*

\*\* *do SISu vkládejte formulář nepodepsaný (ve formátu PDF), podpis je potřeba doplnit až na vytištěný posudek.*