

Posudek oponenta na bakalářskou práci
Klasifikace dat popsaných stromovou strukturou

Autor práce: Lukáš Kulička

Oponent práce: Kateřina Henclová

Bakalářská práce se zabývá aktuálním tématem z oboru strojového učení: využití stromových struktur pro reprezentaci dat ve více-istančním učení. Teprve na takto zpracovaná data je možno nasadit nástroje klasického strojového učení, např. v klasifikačních úlohách.

Při zpracování náročného tématu si student osvojil rozsáhlé znalosti z různých oblastí matematiky a strojového učení: teorie pravděpodobnosti, teorie informace, optimalizace, teorie grafů, neuronové sítě. V první kapitole jsou pečlivě a obsáhle vysvětleny základní pojmy, věty a jiné nástroje potřebné pro pochopení daného tématu. Druhá kapitola je věnována obtížnému odvození nástroje ELBO (Evidence Lower Bound) a vysvětlení logistické regrese. Ve třetí kapitole autor nejprve představuje myšlenku a důležitost vysvětlitelnosti modelu. Dále se soustředí na souvislost vysvětlitelnosti a řídkosti parametrizace. Nakonec přechází k hlavnímu cíli práce: více-istančnímu učení. Zde ukazuje, jak se dá problém více-istančního učení pomocí představených nástrojů transformovat do klasické úlohy strojového učení a dále řešit. Příklady byly naprogramovány v jazyce Julia.

I přes delší seznam nalezených nedostatků (viz dále) je nutné zdůraznit, že se nejedná o závažné chyby a jejich množství je dáno především délkou práce. Většina chyb je soustředěna v méně důležitých či přímo okrajových pasážích textu, zatímco ty podstatné jsou odladěné lépe. Student by měl zlepšit zejména přesnost formulací v definicích a větách, v textu se lépe odkazovat na své zdroje a v závěru přidat lepší ilustrační příklad. Je škoda, že u prakticky motivované práce chybí popis cílové aplikace.

Tato práce splňuje podmínky kladené na bakalářskou práci. Svým rozsahem i obsahem navíc vysoce převyšuje obvyklý standard bakalářské práce. Zpracování je pečlivé, přehledné a s minimálním počtem překlepů (a i ty jsou pouze kosmetické jako velikost závorek apod.) Autor si dal záležet na vybudování potřebné teorie od základních kamenů až po značně pokročilé prvky. Výpočty provádí precizně a i složitými odvozovacími pasážemi provádí čtenáře velmi srozumitelně.

I přes místy kolísavou úroveň se vzhledem k vysoké obtížnosti i rozsáhlosti zpracovávaného tématu a délku práce přikláním k hodnocení A za předpokladu zodpovězení níže uvedených otázek 1 a 2 při obhajobě.

Připomínky:

1. (strana 24) Pojmy “pravděpodobnost” a “věrohodnost” nejsou v našem kontextu záměnné, obzvlášť jedná-li se o tzv. maximálně věrohodný odhad parametru. Tato připomínka se týká jen slovního popisu, v matematické vzorce jsou správně.
2. (strana 25, 44 aj.) Pojem “metoda” je v práci opakovaně nevhodně používán. Např. ELBO - evidence lower bound - není sama o sobě optimalizační metoda ani metoda strojového učení (kromě toho z uvedeného není zřejmé, k čemu tento nástroj slouží). Také feature selection je označeno jako metoda, ačkoli se jedná o typ úlohy a/nebo celou třídu metod k jejímu řešení.
3. (strana 27) V obrázku 1.4 je označení $f(x)$ zavádějící, protože x v téže obrázku již označuje něco jiného. Vhodnější by bylo např. $f(\cdot)$.
4. (strana 30 a dále) V sekci Teorie grafů autor neustále zaměňuje pojmy “graf” a “grafická reprezentace grafu”. I kvůli tomu výsledné formulace definic a vět mnohdy připomínají spíše neformální poznámky než korektně formulovanou matematiku. Důvodem může být použití zdroje [5], což je diplomová práce s názvem “Teorie grafů ve výuce na střední škole”.
5. (strana 30) Definice 1.5.2 je nevyhovující. Formulace “Pokud jsou hrany grafu doplněny o tzv. šipky, které určují jejich směr...” klade důraz na šipky, které ovšem jsou pouhým grafickým označením faktu, že hrany mají orientaci.
6. (strana 30) Definice 1.5.7: “nejkratší délku cesty z V_i do V_j ” by bylo lepší nahradit formulací “délka nejkratší cesty z V_i do V_j ”.
7. (strana 31) Ve formulaci “strom [je] pouze konkrétní realizace grafu” by bylo vhodnější říci např. “strom je pouze jedním typem grafu”. Slovo realizace s sebou totiž nese jiný význam.
8. (strana 31) V bodě 4 teoremu 1.5.9 chybí předpoklad souvislosti grafu G .
9. (strana 31) Definice 1.5.11 je nejednoznačně formulovaná, a tedy zmatečná (navíc v ní opět chybí předpoklad souvislosti). Lepší formulace by byla např. “polystrom je orientovaný strom”. Dále, pokud už nepříliš častý pojem “polystrom” zavádíme, tak by bylo vhodné jej dále používat (nebo jej vůbec nezavádět a vystačit si s pojmem “orientovaný strom”).
10. (strana 31) Věta “Pro nás je ovšem důležité se v grafu orientovat, respektive znát směr.” dle mého subjektivního dojmu nevhodně zaměňuje pojmy “orientace grafu” a “orientovat se”.

Také formulace “tímto strom zakořeníme” není per se špatně, ale působí poněkud neformálně.

11. (strana 32) V definici 1.5.12 chybí předpoklad orientovanosti (poly)stromu G . Pořadí definovaných pojmů by bylo vhodnější uspořádat tak, aby definice nenásledovala až poté, co je pojem použit. Chybí definice pojmu “kořen” (odstaveček na str. 31 tento pojem motivuje, ale nedefinuje; poznámka na str. 32 také definice nenahradí).
12. (strana 41) Část “Reparametrizační trik” by zasloužila podstatných změn a upřesnění. Např. funkce $f(m)$ je zavedena bez jakýchkoliv předpokladů (a navíc je označení f v předcházející i následující části již použito). Formulace “Toto je přesné pro velká n .” a “Budeme-li výpočet mnohokrát opakovat, docílíme rovnosti [(2.26)]” jsou nepravdivé, protože se v obou případech jedná o rovnost v limitě (možná i za uvedení dalších předpokladů, např. konečnost střední hodnoty ve (2.26)). Ve větě “Ovšem pokud budeme brát $n = 1$, získáme nestranný (nevychýlený) odhad gradientu.” je zmíněn gradient, ale čeho, jaký? Kromě toho z uvedeného není zcela zřejmé, co je cílem tohoto triku. Tuto část doporučuji do obhajoby důkladněji prostudovat. Vhodným zdrojem může být např. kapitola 2 ve článku *Implicit Reparameterization Gradients* (Figurnov et al.).
13. (strana 43) Definice 3.0.1 je formulována nesrozumitelně. Jedná se o důsledek nevhodného překladu ze zdroje [3] (“Explainable Artificial Intelligence is a system that produces details or reasons to make its functioning clear or easy to understand.”)
14. (strana 43) V úvodním odstavci pro část 3.1 dává autor motivaci, proč se zabývat řídkými parametrizacemi. Nicméně přitom nešťastně míchá dohromady vysvětlitelnost a náročnost trénování modelu. Vzhledem ke své ideové důležitosti by tato část zasloužila lepší vysvětlení.
15. (strana 44) L_0 regularizace - kvůli velkému praktickému významu by bylo vhodné zmínit také L_1 regularizaci (ať už jako samostatný nástroj nebo relaxaci L_0 regularizace).
16. (strana 49) Správný překlad anglického termínu “map, mapping” je “zobrazit, zobrazení”, nikoliv “mapa”.
17. (strana 52-53) Jako jediný příklad ilustrující použití stromové struktury pro klasifikaci je uveden příklad, jehož řešení (z dobrých důvodů) nevyjde tak, jak je očekáváno. Navíc modelová stromová struktura je extrémně jednoduchá. Chybí mi zde tedy:

- 1) příklad, který vyjde “správně”, 2) příklad, u kterého nebude stromová struktura takto triviální.
18. (Obecná připomínka.) Kvůli přehlednosti není vhodné začínat větu matematickým symbolem, např. “ α zde má význam přesnosti...”
 19. (Práce se zdroji.) Bylo by vhodnější mít v textu více citací, zejména u odvozovacích pasáží (zdroje jsou řádně uvedeny v seznamu literatury, ale odkazů v textu by mělo být více). Zejména kapitola 2 postrádá odkazy na literaturu. Např. podstatná část odvozování byla převzata ze zdroje [1], aniž by na něj bylo (znovu) odkázáno. Počet referencí je pro bakalářskou práci relativně vysoký, i když netriviální podíl tvoří méně standardní zdroje: nerecenzované online zdroje, učební texty, diplomové práce apod.

Otázky k obhajobě:

1. Můžete uvést konkrétní praktický příklad-motivaci víceinstančního učení? Co je zde bag, instance?
2. Uveďte netriviální ilustrační příklad (včetně jeho řešení) klasifikace dat popsaných stromovou strukturou.
3. Co je reparametrizační trik a proč je důležitý?

Návrh na hodnocení bakalářské práce: A (výborně).

V Praze dne 30. 7. 2020

Kateřina Henclová