

I. IDENTIFIKAČNÍ ÚDAJE

Název práce:	Vysvětlování výstupu modelů zpracování přirozeného jazyka pro úlohu ověřování faktů
Jméno autora:	Eliška Kopecká
Typ práce:	diplomová
Fakulta/ústav:	Fakulta elektrotechnická (FEL)
Katedra/ústav:	Katedra pocitacu
Oponent práce:	Gustav Šír
Pracoviště oponenta práce:	Katedra pocitacu

II. HODNOCENÍ JEDNOTLIVÝCH KRITÉRIÍ

Zadání	průměrně náročné
<i>Hodnocení náročnosti zadání závěrečné práce.</i>	
Zadání práce, tak jak bylo sepsáno, dává dostatek prostoru ke splnění, což je vhodné, s přihlednutím k problematice specifikaci úlohy interpretability ve strojovém učení.	

Splnění zadání	splněno
<i>Posuďte, zda předložená závěrečná práce splňuje zadání. V komentáři případně uveďte body zadání, které nebyly zcela splněny, nebo zda je práce oproti zadání rozšířena. Nebylo-li zadání zcela splněno, pokuste se posoudit závažnost, dopady a případně i příčiny jednotlivých nedostatků.</i>	
Bez vyhrad, studentka vhodně volí jednu ze 2 hlavních komponent existující pipeline a testuje její interpretabilitu.	

Zvolený postup řešení	správný
<i>Posuďte, zda student zvolil správný postup nebo metody řešení.</i>	
Postup mi přijde zcela odpovídající. Studentka provádí širokou rešerši pojmu interpretability, a volí 2 nejpoužívanější metody pro lokální vysvětlování výstupu modelu, adaptuje je pro daný problém a důkladně testuje, s přihlednutím k omezením plynoucím z daného prostředí (náročnost testování).	

Odborná úroveň	A - výborně
<i>Posuďte úroveň odbornosti závěrečné práce, využití znalostí získaných studiem a z odborné literatury, využití podkladů a dat získaných z praxe.</i>	
Práce nepředstavuje z hlediska odbornosti žádné převratné koncepty, ale to je dáno i problematikou interpretability vytyčenou v zadání. Navíc studentka k problematice přistupuje velmi pečlivě, správně upozorňuje na záludnosti této domény, a nabízí některá i docela kreativní řešení. Sice používá vesměs standardní metody/knihovny, avšak představuje i některé vlastní modifikace, a řeší i testování v „produkci“. Vše je navíc zpracováno velmi pečlivě.	

Formální a jazyková úroveň, rozsah práce	A - výborně
<i>Posuďte správnost používání formálních zápisů obsažených v práci. Posuďte typografickou a jazykovou stránku.</i>	
Zde je práce zcela výborná, text je psán velmi pečlivě se zřejmým důrazem na porozumění čtenáři. Občas působí možná až trochu moc příběhově na odborný text, ale osobně to hodnotím kladně. Angličtina je skvělá, typografie velmi pěkná, včetně grafiky a výběru barev.	

Výběr zdrojů, korektnost citací	A - výborně
<i>Vyjádřete se k aktivitě studenta při získávání a využívání studijních materiálů k řešení závěrečné práce. Charakterizujte výběr pramenů. Posuďte, zda student využil všechny relevantní zdroje. Ověřte, zda jsou všechny převzaté prvky řádně odlišeny od vlastních výsledků a úvah, zda nedošlo k porušení citační etiky a zda jsou bibliografické citace úplné a v souladu s citačními zvyklostmi a normami.</i>	
Bez vyhrad.	

Další komentáře a hodnocení

Vyjádřete se k úrovni dosažených hlavních výsledků závěrečné práce, např. k úrovni teoretických výsledků, nebo k úrovni a funkčnosti technického nebo programového vytvořeného řešení, publikačním výstupům, experimentální zručnosti apod.

Viz celkové hodnocení.

III. CELKOVÉ HODNOCENÍ, OTÁZKY K OBHAJOBĚ, NÁVRH KLASIFIKACE

Shrňte aspekty závěrečné práce, které nejvíce ovlivnily Vaše celkové hodnocení. Uveďte případné otázky, které by měl student zodpovědět při obhajobě závěrečné práce před komisí.

Prace řeší problém interpretace modelu pro overování faktu v českém jazyce, k čemuž adaptuje 2 klasické, lokální, agnostické metody. Pozitivně hodnotím i vlastní modifikaci s embeddings (text- augmented LIME). Tyto metody pak testuje s přihlednutím k omezeným možnostem evaluace (za mě ideální postup k testování problematického konceptu interpretability v dané úloze je pak hezky nastaven ve future work – s tím se plně ztotožňuji).

Celkově se studentka s touto zaludnou problematikou poprala velmi obstojně a výsledek působí profesionálně.

Poznámky:

Osobně nemyslím, že je hned nutné uchylovat se pouze k lokálním, nepřímým metodám jako LIME a resit (diskutovat) pak problémy s fidelitou (o čemž svědčí i lepší výsledek SHAP), ale lépe přímo trackovat gradient/feature saliency v daném transformeru, k čemuž už lze kdyžtak také nalézt hotové knihovny [1].

Otázky:

- Naznačujete, že fasttext embeddings nejsou příliš vhodné, protože nezohledňují daný problém. Nebylo by tedy lepší pro vysvětlování použít přímo embeddings z daného transformeru?
- Neslo by použít bag of embeddings namísto bag of words a dělat perturbaci přímo v prostoru embeddings, namísto hledání nejbližšího slova?
- Byl vyvoj zobrazené aplikace součástí Vaší práce?
- Čapů správně ze počtu testerů kvality interpretací byl 2?

Par typos:

trash-hold :)

only additive possible additive

(e.g..

parametrizations

[1] <https://github.com/jalammar/ecco>

Předloženou závěrečnou práci hodnotím klasifikačním stupněm **A - výborně**.

Datum: 6.6.2022

Podpis: