

Explaining Explanations: An Overview of Interpretability of Machine Learning

In order for humans to trust black-box methods, we need explainability – models that are able to summarize the reasons for neural network behavior, gain the trust of users, or produce insights about the causes of their decisions.

Interpretovateľnosť - chceme chápať čo model robí

Vysvetľiteľnosť - chceme pochopiť prečo sa model správa tak ako sa správa,

- interpretovateľnosť sama o sebe nestačí
- Aby ľudia dôverovali metódam black box, potrebujú vysvetliteľné modely, ktoré sú schopné zhrnúť dôvody pre správanie neurónovej siete, tieto získajú dôveru používateľov alebo produkujú informácie o príčinách ich rozhodnutí
- Hoci pojmy interpretovateľnosť a vysvetľiteľnosť sa zvykne často zamieňať je veľmi dôležité rozlišovať medzi nimi
- vysvetliteľné metódy sú štandardne interpretovateľné ale naopak to nemusí platiť

Zaujímavosť:

- Niektoré existujúce zavedené systémy a nariadenia spôsobujú, že je potrebné, aby boli systémy vysvetliteľné, srozumiteľné a včasné
- S blížiacimi sa nariadeniami, ako je „právo Európskej únie na vysvetlenie“ [12], ktoré požaduje rozmanitosť a začlenenie do systémov AI [13], zistenia, že niektoré automatizované systémy môžu posilniť nerovnosť a zaujatosť [14], spôsobili že v poslednej dobe extrémne narastol záujem o vysvetliteľné systémy
- Tieto black box modely sú implementované na veľa miestach a v rôznych disciplínach kde je ich vysvetľiteľnosť potrebná
- Príklady všeobecných „vysvetliteľných systémov“ : interpretovateľné AI, vysvetliteľné ML, kauzalita, bezpečná AI, výpočtová sociálna veda a automatický vedecký objav. Ďalší výskum vo vysvetľeniach a ich vyhodnotení sa nachádzajú v strojovom učení, interakcii človeka s počítačom (HCI) a mnoho ďalších iných disciplín.

sekcia 2: základné pojmy - explanation, interpretability, explainability

- explanation: možno vyhodnotiť dvoma spôsobmi: podľa jeho interpretovateľnosti a podľa jej úplnosti.
- cieľom interpretovateľnosti je opísať vnútro systému tak aby bol zrozumiteľný pre ľudí
- Úspech tohto cieľa je spojený s poznaním, znalosťami a zaujatosťou používateľa: na to, aby bol systém interpretovateľný, musí poskytnúť opisy, ktoré sú dostatočne jednoduché na to, aby ich človek pochopil pomocou

slovnej zásoby, ktorá je pre používateľa pochopiteľná.

- Cieľom úplnosti je presne popísať fungovanie systému.
- vysvetlenie je úplnejšie keď umožňuje predpovedať správanie sa systému vo viacerých situáciách
- Pri vysvetľovaní programu počítača ako napríklad DNN, dokonale kompletne vysvetlenie možno vždy získať odhalením všetkých matematických operácií a parametrov v systéme.
- Výzva, ktorej čelí vysvetliteľná AI, spočíva vo vytváraní vysvetlení, ktoré sú úplné a interpretovateľné: je ťažké dosiahnuť súbežnú interpretovateľnosť a úplnosť.
- Najpresnejšie vysvetlenia nie sú ľahko interpretovateľné ľuďom; a naopak najviac interpretovateľné popisy často neposkytujú prediktívnu silu.
- Herman [18] poznamenáva, že by sme mali byť opatrní pri vyhodnocovaní interpretovateľných systémov iba pomocou ľudských hodnotení interpretovateľnosti, pretože hodnotenia ľudí implikujú silne skreslenie k jednoduchším opisom. Varuje že spoliehanie sa na hodnotenia ľudí môže viesť vedcov k vytváraniu presvedčivých systémov a nie transparentných
- Dve dôležité etické dilemy

Kedy je neetické manipulovať s vysvetlením na to aby sme lepšie presvedčili používateľov?

Ako vyvážíme naše obavy týkajúce sa transparentnosti a etiky s našou túžbou po interpretovateľnosti?

- veria že je neetické prezentovať zjednodušený opis komplexného systému len preto aby sme zvýšili dôveryhodnosť, najmä ak rôzne limitácie a obmedzenia zjednodušeného opisu používateľa nie sú schopní pochopiť
 - a čo je horšie, ak je optimalizované vysvetlenie skryva nežiaduce atribúty systému
- na to aby sme tomu predišli vysvetlenia by mali predstavovať akýsi kompromis medzi interpretovateľnosťou a úplnosťou

Explainability of Deep Networks processing

- vysvetliteľnosť neuronových sietí sa zameriava buď na vysvetlenie spracovania údajov sieťou alebo na vyseveleínie reprezentácie dát vo vnútri siete
- Bežne používané hlboké siete odvodzujú svoje rozhodnutia pomocou veľkého množstva základných operácií - veľa parametrov a operácií
- Teda základným problémom, ktorému čelia vysvetlenia takéhoto spracovania, je nájsť spôsoby, ako znížiť zložitosť všetkých týchto operácií.
- To sa dá dosiahnuť vytvorením proxy modelu, ktorý sa správa podobne ako pôvodný model

Linear Proxy Models:

- LIME, Riberio
- Pri LIME sa systém čiernych skriniek vysvetľuje skúšaním správania poruchy vstupu a potom sa tieto údaje použijú na zostavenie lokálneho lineárneho modelu, ktorý slúži ako zjednodušený úplný model v blízkosti k

vstupu.

- Ribeiro ukazuje metódu, ktoru možno použiť na identifikáciu oblastí vstupu, ktoré majú najväčší vplyv na rozhodovanie v rôznych druhoch modelov a problémových doménach
-

Decision trees

- obmedzena skalovateľnosť
- Úsilie rozložiť neurónové siete na rozhodovacie stromy nedávno rozšírili prácu od 90. rokov, zameraná na plytké siete, na zovšeobecnenie procesu hlbokých neurónových sietí.
- jednou z takýchto metod je DeepRed
- ANN DT
-

Automatic-Rule Extraction:

- Automatická extrakcia pravidiel je ďalším dobre preštudovaným prístupom k sumarizácii rozhodnutí.
- Dekompozičné prístupy pracujú na úrovni neurónov, extrahujú pravidlá napodobňujúce správanie jednotlivých jednotiek.

Saliency Mapping

- sieť je opakovane testovaná so zakrytými časťami vstupu aby sa vytvorila mapa ukazujúca, ktoré časti údajov skutočne majú vplyv na sieťový výstup.
 - Príkladmi sú LRP [40], DeepLIFT [41], CAM [42], GradCAM [43], integrované prechody [44] a SmoothGrad [45]. Každá technika vytvára rovnováhu medzi zobrazením oblastí, vysokou aktiváciou siete, kde neuróny sú najsilnejšie a областami s vysokou citlivosťou na sieť, kde by zmeny mohli najviac ovplyvniť výstup.
 - Porovnanie niektorých z týchto metód Ancona[46].

explainability of deep learning representation

Role of layers

Role of individual units

Role of representation vectors

Zaujímavé zdroje: 46