

Kontextus alapu érzelemdetekci6

Keszitette: Szaniszló Csongor Adam

Konzulens: Dr. Hullám Gabor, Revy Gabor

Kib6vitett Absztrakt

Az rzelemdetektalásnak számos felhasználási területe van, ilyen például az egészségügy, ahol mentális betegségeket tudnak detektálni ezáltal, a vevőszolgálat, ahol ennek segítségével személyre szabott ajánlásokat tudnak adni a cégek vagy akár az oktatás is, ahol például egy tanárral kaphat bővebb visszajelzést az órájáról. Így tehát az érzelemdetektalásnak, mint feladatnak számos motivációja van.

Ugyanakkor míg ez emberek számára egy egyszerűbb feladat, hiszen emberként ezt születésünk óta tanuljuk, addig a számítógépek számára ez már közel sem egy egyértelmű megoldással rendelkező probléma. Ennek nagyon sok oka lehet, például, hogy a bemeneti képek vagy videofelvétel minősége nem feltétlenül megfelel, amely megnehezíti a gépek feladatát. A neurális hálók pontosabban a Convolution Neural Network (CNN) megjelenése jelentősen elősegítette a jobb eredmények elérését ezen a területen.

Az utóbbi években számos kutatás foglalkozott mélyebben ezzel a témával, melyek során a fő cél annak megállapítása volt, hogy különböző jellemzők/bemeneti információk, mint például egy kép az arcáról vagy egy hangfelvétel hogyan segítik el az érzelmek helyes detekcióját. A tematikus laboratóriumi munkám során is például egy ilyen kutatásnak olvastam utána, illetve próbáltam ki a cikk által javasolt megoldást. Ennek során a FER2013 (Facial Emotion Recognition) adathalmazzal foglalkoztam, amely az érzelemdetektálásra használt adathalmazok közül kiemelkedően fontos, gyakorlatilag egy mércéjeve vált az arc alapú érzelemdetektálásnak.

Ugyanakkor pszichológiailag is alátámasztható, hogy az érzelemdetektálás során több aspektust is figyelni kell ahhoz, hogy helyesen tudjuk behatározni a vizsgálandó személy érzelmi állapotát. Így például az előzőekben említett FER2013-ra visszautalva nem elég az, ha csupán egy személy arc kifejezését vizsgáljuk

1. es 2. abra: Kivagott arckép az EMOTIC(2] adathalmazból, illetve a hozzá tartozó kép, kontextussal

Peldaül a fenti kép esetén is látható, hogy ha csupán az arcát vizsgáljuk a személynek, akkor valamilyen negatív érzélemmel asszociálnak az érzelmi állapotát, míg ha a kontextust is figyelembe vesszük, akkor már evidens, hogy a személy valószínűleg a sikerének örövend.

Ezt figyelembe véve tehát az érzelemdetektálás során érdemes az arcon kívül más is elemezni. Ilyen például a kontextus, amelyben az alany részt vesz, hiszen ebből megállapítható az, hogy a személy boldog, ha például egy koncerten van. Egy másik aspektus, amit érdemes fontolara venni, az az alanyról kivagott kép, amit a kontextust tartalmazó képből nyerünk ki. Ezen kép alapján elemezhető az alany testtartása, ami ugyanúgy jelentősen hozzájárulhat az érzelemdetektáláshoz.

Az önálló laboratóriumi munkám során tehát a kontextus alapú érzelemdetektálással foglalkoztam. A motiváció az volt, hogy megvizsgáljam, hogy a kontextus ténylegesen mennyire befolyasolja az érzelemdetektálást.

A munkám során először egy publikációval, illetve az írók által létrehozott adathalmazzal, az EMOTIC-kal foglalkoztam. Az EMOTIC egy olyan adathalmaz, amely minden egyes vizsgálandó személyhez két képet társít. Az első egy olyan kép, amely a személyt egy adott kontextusban jeleníti meg. A második pedig az első képből egy kivagott kép, ez magát a személyt tartalmazza.

A publikáció írói egy modellt is bemutatottak, amin megvizsgáltam, hogy az eszköz milyen eredményeket képes elérni az adathalmazon. A modell az előbb említett két bemeneti kép alapján hozott két döntést. Az egyik döntés 26 érzelmi kategóriából statisztikailag mennyiségűbe sorolta be a képet, míg a második döntés 3 folytonos érték predikálására vonatkozott, amely értékek a VAD/PAD[3] modell szerint értékeltek a személy érzelmi állapotát. A kiértékelés során a kapott

eredményeket_összehasonlítottam a publikációban leírt eredményekkel. Bár ez nem sokban különbözött (az alkalmazott metrika szerint 2%-os volt az eltérés), mégis arra a megfontolásra jutottam, hogy a modell által elért eredmények viszonylag alacsonyak, így megpróbáltam javítani az eszköz teljesítményen.

Ennek során számos megközelítéssel próbalkoztam. Először megpróbáltam image padding-et használni. Ennek motivációja az volt, hogy az adat előfeldolgozás során úgy ítéltém meg, hogy a csupán csak a személyt tartalmazó kép esetében a kép átmeretezésének hatására túl sok információ veszik el, amely megnehezíti a modell feladatát a döntésben. Ezenkívül megvizsgáltam az adathalmazt és azt figyeltem meg, hogy a képek nagy részénél a vizsgaland személy arca részben, vagy teljes egészében nem látszik. Úgy ítéltém meg, hogy bár nem elegendő feltétele az érzelemdetektálásnak az arc elemzése, de ugyanakkor szükséges, így az adathalmazból kiszűrtem az arcot, nem tartalmazó képeket. A harmadik megközelítés pedig egy j modellel integrálása volt, amely a FER adathalmazon lett előtanítva. Ezzel a célom az volt, hogy a kontextuson, illetve a testtartáson kívül a személyek arcát is elemezze a modell, ezzel segítve azt a döntésben.

adathalmazt is annak reményében, hogy a modell azon jobban teljesít majd. Ez az adathalmaz a HECO[4] volt. A HECO nagyon hasonlít az EMOTIC-hoz, az egyetlen különbség az, hogy 26 érzelme helyett 8 érzelmi címke van, illetve minden képhez egy darab érzelmet társítottak az alkotók.

A HECO-n való kiértékelés során sajnos ugyanúgy teljesítménnyel kapcsolatos problémákba ütköztem, így az előzőekben bemutatott megközelítésekkel próbáltam a modell teljesítmény javítani. Ebben az esetben sem tapasztaltam semmilyen javulást.

Az eredmények hatására úgy döntöttem, hogy célszerű az alkalmazott modell, illetve a használt adathalmazok mélyebb vizsgálata. A modellek vizsgálatánál arra jutottam, hogy a vizsgálandó személyt elemző modell, illetve az arcot elemző modell egyik adathalmaz esetében sem volt képes releváns jellemzőket kiemelni. Ennek az előbbi esetben az volt az oka, hogy a modell egy olyan adathalmazon lett előtanítva, amely egyáltalán nem kapcsolódik az érzelemdetektáláshoz, míg utóbbi esetben a probléma a modell komplexitásával (EMOTIC) volt vagy az adathalmaz inkonzisztens címkézéssel (HECO). A kontextust elemző modell vizsgálatával így arra jutottam, hogy a modell feladata érzelme-kontextus párok összekapcsolása volt, viszont megvizsgáltam azt

is, hogy a modell milyen kepre milyen kontextus kimenetet adja, es arra jutottam, hogy semmilyen korrelacio nincs az erzelmek es a kontextusok kozott

Az adathalmazok vizsgalata soran felfedeztem, hogy mindket adathalmazzal tobb problema volt. Peldaul hibas bounding box-ok (egy dobozon belül tobb személy is van), inkonzisztens cimkezes, illetve az annotalt személyek arcanak nagy resze nem latszodik bar utobbi javithato, ugyanakkor a kepek kiszürese nem cel, ha azt szeretnenk tesztelni, hogy egy modell hogyan teljesit egy adott adathalmazon

Forrasok

[1] FER2013 - <https://www.kaggle.com/datasets/msambare/fer2013>

[2] EMOTIC website - <https://s3.sunai.uoc.edu/emotic/index.html>

[3] VAD/PAD model - https://en.wikipedia.org/wiki/PAD_emotional_state_model

[4] HECO website - <https://heco2022.github.io/>