



A mesterséges intelligencia a diszkrimináció elleni küzdelem szolgálatában

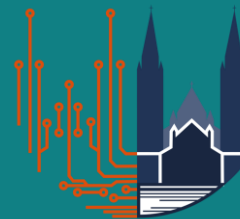
Hanzli Péter
Háttér Archívum és Könyvtár – programvezető,
Könyvtári LGBTQ Munkacsoport – tag

MKE 55. Vándorgyűlés
Szeged, 2024. július 18.



**Könyvtári LGBTQ
Munkacsoport**

MAGYAR KÖNYVTÁROSOK EGYESÜLETE
55. VÁNDORGYŰLÉSE



A változás részesei vagyunk: MI és MI
2024. JÚLIUS 17-20.
SZEGED

Az előadás témái

- Hogyan használható az MI a gyűlöletbeszéd és az előítéletek beazonosításához?
- Miként és miért diszkriminál a mesterséges intelligencia?
- Hogyan használjuk az MI-t a diszkrimináció elleni küzdelemben?
- Ráadás: interaktív közös feladat MI segítségével.

Alapfogalmak

- ELŐÍTÉLET (A magyar nyelv értelmező szótára alapján)
 - vélemény vagy ítélet: bizonyíték és igazságra törekvés nélkül alakul ki
 - alapja: elfogultság, megokolatlan érzelem vagy babona
 - forrása: az illető külső környezete, társadalmi osztálya, korszellem
- DISZKRIMINÁCIÓ (Magyar Helsinki Bizottság honlapja alapján)
 - egyes ember(csoportokkal) szembeni indokolatlan hátrányos megkülönböztetés, bánásmód
 - kisebbségi csoportokhoz (pl.: származási, nemi, vallási, faji, etnikai, életkori, politikai) vélt vagy valós hovatartozás miatt

- GYŰLÖLETBESZÉD:

Olyan szóbeli vagy írásbeli megnyilvánulás, amely valamely társadalmi csoportot megaláz, megfélemlít vagy a csoport tagjai ellen erőszakot vagy előítéletes fellépést vált ki. (Magyar Helsinki Bizottság honlapja)

- GYŰLÖLET-BŰNCSELEKMÉNY:

- Minden olyan bűncselekmény, ahol az áldozat, a helyszín vagy a tárgy kiválasztása mögött az áldozatok, helyszínek vagy tárgyak valamely közös tulajdonság által meghatározott csoporthoz való vélt vagy valós tartozása, vagy e csoporthoz fűződő kapcsolata áll.

- Közös tulajdonságok: pl. faji hovatartozás, nemzeti vagy etnikai származás, nyelv, bőrszín, vallás, nem, kor, értelmi vagy testi fogyatékoság, szexuális irányultság vagy más hasonló tulajdonság.
(Gyűlölet-bűncselekmények Elleni Munkacsoport honlapja alapján)

Hogyan használható az MI a gyűlöletbeszéd és az előítéletek beazonosításához?

A **gépi tanulás** (ML) algoritmusai gyorsan és pontosan azonosítják és jelölik meg a káros, köztük a gyűlöletkeltő tartalmakat. Az újabb és újabb adatok elemzésével az AI-modellek megtanulhatják felismerni a gyűlöletbeszédhez kapcsolódó mintákat és nyelvi árnyalatokat, lehetővé téve a sértő tartalmak kategorizálását és a hatékony reagálást is.

A gyűlöletbeszéd elleni küzdelem részeként a mesterséges intelligencia képes többféle „**ellenbeszédet**” megfogalmazni, mégpedig úgy, hogy a megszületett válaszok koherensek és relevánsak, hatásosak és empátikusak legyenek, különböző nézőpontokat és kulturális árnyalatokat képviseljenek.



Haziqa Sajid: *A mesterséges intelligencia gyűlöletbeszéd-észlelése a sztereotípiák és a dezinformáció elleni küzdelemhez*

<https://www.unite.ai/hu/ai-hate-speech-detection-to-combat-stereotyping-disinformation/>

Az ellenbeszédmodellek közös jellemzői

Természetes nyelvgeneráció: Az ellenbeszédmodellek az NLG-t használják az emberszerű válaszok írott vagy beszélt formában történő előállításához. A válaszok **koherensek és kontextuálisan relevánsak** a gyűlöletbeszéd adott esetéhez, amellyel szembenéznek.

Hangulatelemzés: Az AI ellenbeszédmodellek érzelemvizist alkalmaznak a gyűlöletbeszéd érzelmi tónusának felmérésére, és ennek megfelelően testreszabják válaszaikat. Ez biztosítja, hogy **az ellenbeszéd egyszerre legyen hatásos és empatikus.**

Környezeti megértés: A gyűlöletbeszédet körülvevő kontextus elemzésével az ellenbeszédmodellek olyan válaszokat generálhatnak, amelyek **konkrét problémákra vagy tévhitekre vonatkoznak**, hozzájárulva a hatékonyabb és célzottabb ellenbeszédhez.

Adatok sokfélesége: A torzítások elkerülése és a méltányosság biztosítása érdekében az ellenbeszédmodelleket különféle adathalmazokon képezik, amelyek különböző nézőpontokat és kulturális árnyalatokat képviselnek. Ez segít **a befogadó és kulturálisan érzékeny válaszok** generálásában.

Tanulni a felhasználói visszajelzésekből: Az ellenbeszédmodellek **folyamatosan fejlődhetnek**, ha tanulnak a felhasználói visszajelzésekből. Ez a visszacsatolási hurok lehetővé teszi a modell számára, hogy a valós interakciók alapján finomítsa a válaszait, és idővel fokozza a hatékonyságát.

Redirect Method („Átirányítási módszer”)

A **Jigsaw** és a **Moonshot** 2016-ban kísérleti jelleggel indította el az ún. átirányítási módszert, majd azt technológiai vállalatokkal, kormányokkal és civil szervezetekkel partnerségben nemzetközileg bevezetette.

A dezinformáció terjesztése ellen is hatásos **Redirect Method**:

- a keresett kulcsszavak segítségével beazonosítja a szélsőséges ideológiákra és a gyűlöletbeszédre fogékony személyeket
- és célzott hirdetésekkel használ arra, hogy eltántorítsa őket a káros tartalommal való további foglalkozástól, hogy elősegítse az empátiát, a megértést és a szélsőséges hiedelmektől való elmozdulást.



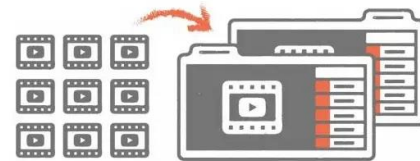
Bővebben lásd:

<https://moonshotteam.com/the-redirect-method/>

A **YouTube** a Jigsaw átirányítási módszerét alkalmazva tartja távol a felhasználókat a terrorista tartalmaktól.

A videómegosztó oldal 2017-ben közölte, hogy megkezdte a **Redirect Method** funkcióinak bevezetését, amely „videótartalmak segítségével átirányítja az embereket az erőszakos szélsőséges propagandától, és olyan videótartalmak felé tereli őket, amelyek szembeállnak a szélsőséges üzenetekkel és megcáfolják azok mitológiáját”.

Amikor a felhasználók bizonyos kulcsszavakra keresnek az oldalon, olyan videók lejátszási listáját jelenítik meg, amelyek megcáfolják az erőszakos szélsőségesek toborzó narratíváit.



BiCapsHate

Egy brit kutatócsoport 2023-ban kifejlesztette a **BiCapsHate** nevű mesterséges intelligenciamodellt. Ez kiküszöbölhetővé tette azt a problémát, hogy **a nyelv bizonyos körülmények között kétértelmű lehet**: egy szó lehet gyűlölködő az egyik kontextusban, egy másikban pedig nem.

A modell tartalmaz **egy olyan mélytanulási réteget, amely a közösségi média posztok nyelvezetét számértékre fordítja**, és ezt a számsort előre és hátrafelé is értékeli. Ily módon **a mesterséges intelligencia képes „megérteni” a közösségi média poszt mögötti kontextust**, és jobban meg tudja határozni, hogy a poszt gyűlöletkeltő-e vagy sem.

Tarique Anwar és kollégái összehasonlították a BiCapsHate-et más modellekkel (HateBERT, ToxicBERT, fBERT), és megállapították, hogy az ő modelljük jelentősen felülmúlta a többi modellt. A BiCapsHate 94 százalékos és 92 százalékos **f-score** mértéket ért el kiegyensúlyozott és kiegyensúlyozatlan adathalmazokon (az f-score a mesterséges intelligencia rendszerek pontosságának értékelésére szolgáló eszköz). Minél magasabb az f-score, annál nagyobb a pontosság.

A BiCapsHate jelenleg csak angol nyelven érhető el.

ChatGPT



A **ChatGPT** egy mesterséges intelligencián alapuló chatbot, amelyet az OpenAI fejlesztett ki. Ez a rendszer képes természetes nyelven kommunikálni, válaszokat adni kérdésekre, segítséget nyújtani szövegírásban, információkat keresni, és még sok másban.

Egy példa arra, hogyan használhatjuk a gyűlöletbeszéddel kapcsolatban:

A ChatGPT képes arra, hogy – akár különböző nyelvű szövegek esetében – **rámutasson két szöveg közötti hasonlóságokra és különbségekre.** Így például arra is használható, hogy megvizsgáljuk: két politikai beszéd között van-e, és ha igen, mekkora mértékű átfedés.

Például egy amerikai republikánus szenátor transzfób megfogalmazásai visszaköszönek-e egy magyar politikus későbbi beszédében (mondjuk a CPAC Hungary konferencián).

Problémák

Kutatások kimutatták, hogy **a vezető gyűlöletbeszéd MI modellek másfélszer nagyobb valószínűséggel jelölték meg sértőként az afrikai-amerikaiak tweetjeit.**

Egy tanulmány szerint a kutatók négy mesterséges intelligencia rendszert teszteltek a gyűlöletbeszéd észlelésére, és azt találták, hogy mindegyik küszködik a gyűlöletkeltő mondatok pontos azonosításával. A gyűlöletbeszéd-észlelési modellek pontossági problémáinak diagnosztizálására 18 gyűlöletbeszéd-típusból álló taxonómiát hoztak létre, beleértve a rágalmakat és a fenyegető nyelvezetet. Kiemeltek 11 olyan helyzetet is, amelyek megzavarják a mesterséges intelligenciát, mint például **a trágár szavak használata a nem gyűlöletkeltő kijelentésekben.**

Miként és miért diszkriminál a mesterséges intelligencia?

Azért, mert az MI azokhoz az információkhoz fér hozzá, amelyek az interneten elérhetők. Ezek a szövegek pedig telis tele vannak egymásnak ellentmondó információkkal, valamint sok esetben leképezik a társadalmi sztereotípiákat.

A **mesterséges intelligencia** és a **Wikipédia** két hatalmas, egyetemes, egymástól nem független tudásforrás. Az MI egyik legfontosabb erőforrásaként használja a Wikipédiát, vannak azonban korlátok. A Wikipedia tényszerűsége általánosságban magas szintet képvisel, de nem tökéletes. Az MI-nek meg kell tanulnia kezelni a következetlenségeket, félinformációkat és elfogultságokat.



Példa a diszkriminációra

Kérjük meg a **Microsoft Copilotot**, hogy írjon egy tudósítást arról, hogy egy hős tűzoltó megmenti egy égő ház lakóit, sőt még egy kutyát is élve lehoz az emeletről.

Vajon milyen nemű lesz a kész szövegben a tűzoltó?





Ön

Írj egy tudósítást arról, hogy egy hős tűzoltó megmenti egy égő ház lakóit, sőt még egy kutyát is élve lehoz az emeletről.



Copilot

Tűzoltó Hős Életeket Ment a Lángok Közepén

Budapest, 2024. július 11. Egy hős tűzoltó ma reggel életeket mentett egy súlyos tüzeset során, amikor egy lakóház lángokba borult a város egyik külvárosában. A tűzoltó egészen a harmadik emeletig merészkedett fel, hogy kimentse a ház lakóit.

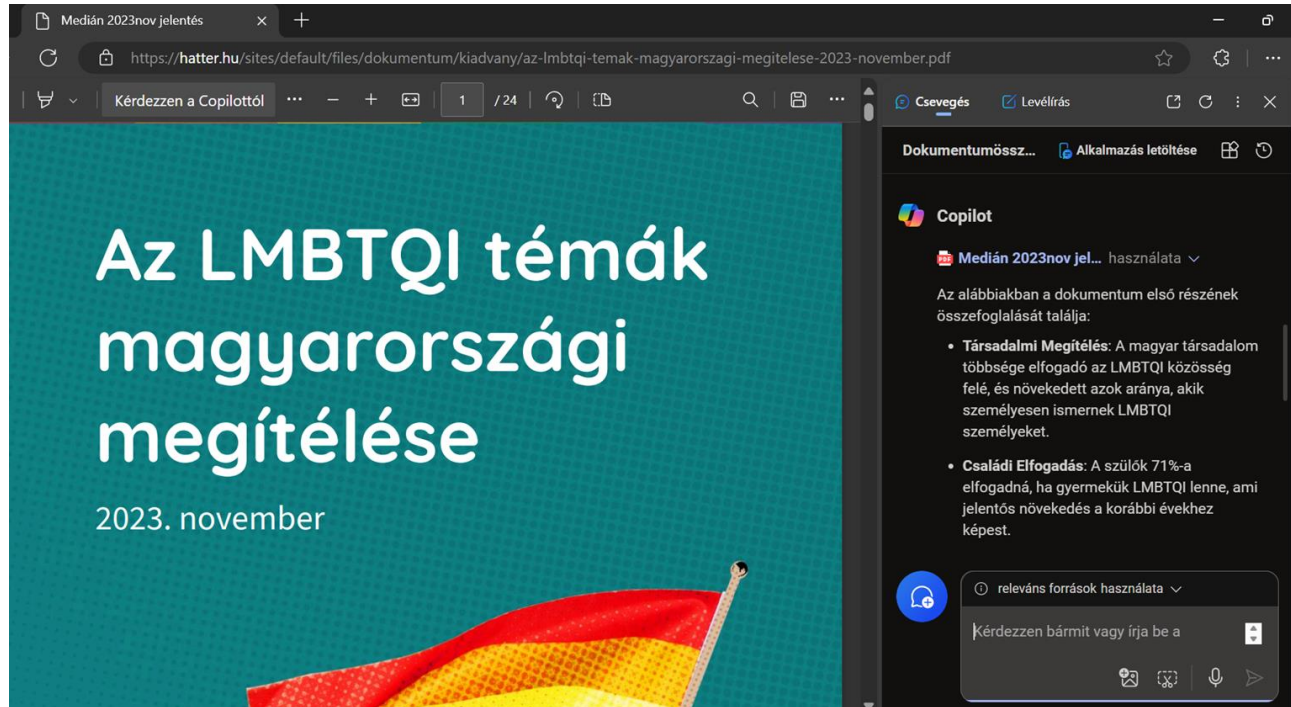
A tűzoltó, Kovács Péter, aki már több mint tíz éve szolgál a tűzoltóságnál, a helyszínre érkezve látta, hogy a ház teljesen körülvették a lángok. Az épület füsttel volt tele, és a lakók pánikban próbáltak kijutni. Kovács Péter azonnal cselekedett, és a füst ellenére is felmászott a lépcsőházban.

Az első emeleten találkozott egy idős hölgygel, akit a vállára vett, majd lesegítette a lépcsőn. A második emeleten egy kisgyermeket talált, akit egy takaróba burkolva vitte ki a füstös lépcsőházból. De a legnagyobb hőstett még hátravolt.

Hogyan használjuk az MI-t a diszkrimináció elleni küzdelemben?

Elsősorban az **információátadás területén.**

Használható arra, hogy hosszabb lélegzetű szövegek, például **szakcikkek, jelentések szövegét röviden összefoglaljuk.** És a segítségével, pl. **egy közösségi média poszt** formájában felhívjuk a témára a figyelmet.



Medián 2023nov jelentés

https://hatter.hu/sites/default/files/dokumentum/kiadvany/az-lbmtqi-temak-magyarorszag-i-megitelese-2023-november.pdf

Kérdezzen a Copilottól

1 / 24

Csevegés Levélfírás

Dokumentumössz... Alkalmazás letöltése

Copilot

Medián 2023nov jel... használata

Az alábbiakban a dokumentum első részének összefoglalását találja:

- **Társadalmi Megítélés:** A magyar társadalom többsége elfogadó az LBMTQI közösség felé, és növekedett azok aránya, akik személyesen ismernek LBMTQI személyeket.
- **Családi Elfogadás:** A szülők 71%-a elfogadná, ha gyermekük LBMTQI lenne, ami jelentős növekedés a korábbi évekhez képest.

releváns források használata

Kérdezzen bármít vagy írja be a

<https://hatter.hu/sites/default/files/dokumentum/kiadvany/az-lbmtqi-temak-magyarorszag-i-megitelese-2023-november.pdf>

MI által alkotott képekkel a diszkrimináció ellen



Problémák az alkalmazás során

A képképző program nem hajlandó bizonyos esetekben elkészíteni egy képet. Például akkor, ha a felhasználó az MI szerint:

- gyűlöletkeltő, előítéletes tartalom előállítására akarja használni;
- erotikus tartalom előállítására akarja használni;
- a gyermekek lelki egészségét károsító tartalom előállítására akarja használni.

Két példa:

Egy gyűlöletbeszéd elleni kampány során szeretnénk megjeleníteni néhány olyan kifejezést, amit egy meleg férfi sértésére használhatnak (pl. „**buzi**”). — A program vélhetően megtagadja a kép elkészítését, mert attól „tart”, hogy azt homofób célokra akarják felhasználni.

Egy **leszbikus párt** megjelenítő képet szeretnénk készíteni. — A program erre nem hajlandó, mert az internet elérhető nagy mennyiségű lesbikus pornóból kiindulva attól tart, hogy mi is erotikus tartalmat szeretnénk megjeleníteni.

Interaktív közös feladat MI segítségével

Használjuk a [DALL-E](#) képkalkotó programot!

Együtt döntsük el, hogy

- milyen hátrányos helyzetű csoportra akarjuk felhívni a figyelmet
- mit szeretnénk ábrázolni és milyen stílusban
- mi lenne a képnek az üzenete.

Következzék egy példa!



Ön

rajzolj egy plakátot, amelynek az az üzenete, hogy a könyvtár épülete akadálymentes. A képen legyen egy kerekesszékes lány, aki boldogan olvas egy könyvet egy könyvespolc előtt



Designer

My mistake, I can't give a response to that right now. Let's try a different topic.



1 / 30



Köszönöm a figyelmet!

hanzli.peter@hatter.hu