

Schrettner Attila

AI Ethics - 2. rész

* *Data governance* szerepe a mesterséges intelligencia etikus használatában *

Ez a cikk az AI Ethics téma feldolgozását folytatja, a sorozat előző, bevezető cikke elérhető a következő linken: <https://ia.hu/images/hirlevel/ethics.pdf>. A sorozat egyes részei a kapcsolódó EU direktíva három nagyobb elvárás csoportja köré szerveződnek: *data governance* (jelen cikk), *kockázatkezelés* és *AI specifikus kontroll tevékenységek*

I. Adatok és mesterséges intelligencia

A mesterséges intelligencia megoldások rendszerint nagy mennyiségű adathalmazra épülnek¹. Így az etikus AI megoldások megteremtésének első lépése az őket „tápláló” adatok etikus kezelése. Ez a lépés jól - de nem teljesen - elkülöníthető a már a rendszer által generált outputokból származó kockázatoktól.

Az adatok etikus kezelésére két különböző szempontból is tekinthetünk:

- Az egyik nézőpont folyamatok szemléletű, amely az üzleti igény definiálásának, adatok beszerzésének, tárolásának, felhasználásának, archiválásának és törlésének folyamatát, azaz az adatok teljes életciklusát vizsgálja
- A másik nézőpont pedig mesterséges intelligencia „tanítása” közben zajló három fázis szerint - a tréning, validációs és tesztelési („*training*”, „*validation*”, „*testing*”) adatokat vizsgálja

A fő megfontolás mindkét szempont esetén a teljeskörűség, azaz a *data governance*-nek az adatok teljes életciklusára és mindhárom elkülönült adathalmazra ki kell terjednie. Fontos hangsúlyozni, hogy itt a *data governance* nem csupán a GDPR-nak² való megfelelést jelenti. Az AI rendszerek használata során adatvédelmi kockázatokon túl számos más kockázat is felmerülhet azáltal, hogy az adatok rossz inputot jelentenek az AI rendszer és a kapcsolódó algoritmus(ok) számára³:

¹ The European Commission’s High-Level Expert Group on Artificial Intelligence (2018): A Definition of AI: Main Capabilities and Scientific Disciplines

² Az Európai Parlament és a Tanács (EU) 2016/679 Rendelete (2016. április 27.) a természetes személyeknek a személyes adatok kezelése tekintetében történő védelméről és az ilyen adatok szabad áramlásáról, valamint a 95/46/EK rendelet hatályon kívül helyezéséről (általános adatvédelmi rendelet)

³ University of Pennsylvania, Wharton (2022): Artificial Intelligence/Machine Learning Risk & Security Working Group (AIRS): Artificial Intelligence Risk & Governance

- **Nem megfelelő adatminőség** - technikai szempontból, pl.: az adatok nem megfelelő formátumban vannak
- **Nem megfelelő adatmennyiség** – a rendelkezésre álló adatmennyiség kevés az AI rendszer megfelelő működéséhez
- **„Prekonceptió („bias”) az adatokban”** – az adatok technikai és mennyiségi szempontból is megfelelők, de tartalmilag valamilyen hiba van bennük (lásd: jelen cikk következő fejezete), például a rendelkezésre álló adatok egy adott társadalmi csoportra koncentrálnak, így belőlük előállított modell és a modellt használó AI rendszer kevésbé lesz alkalmas más társadalmi csoportokra vonatkozó döntési helyzetek kezelésére
- **Tiltott adatkategóriák** – valamilyen törvény vagy más szabály által tiltott adat található az adathalmazban

A fenti kockázatok általánosak, az egyes vállalatok iparágának, működési sajátosságainak vagy egyéb paramétereinek megfelelően további kockázatok is azonosíthatók és a *data governance* keretrendszernek ezekre is ki kell terjednie.

II. „Prekonceptió az adatokban”

A fenti kockázatok közül érdemes kiemelni az „prekonceptió az adatokban” kockázatot. Ez két szempontból is érdekes, egyrészt ez a kockázat leginkább AI rendszerek kontextusában jelentkezik, azok tanulásához kapcsolódóan, másrészt ez egy olyan eset, ami látens módon is jelen lehet, nehéz első ránézésre, vagy felületesebb vizsgálatokkal kiszűrni.

Ráadásul az adatokban lévő prekonceptiók akár az adott vállalat ügyfeleinek alapvető jogainak sérüléséhez is vezethetnek, mint méltósághoz való jog, egyenlőség és diszkriminációmentességhez való jog, fogyasztóvédelemhez való jog, alapvető szolgáltatásokhoz való hozzáférés, stb⁴. Ezért nagyon fontos, hogy a *data governance* kiemelt figyelmet fordítson erre az aspektusra.

Ezt a kockázatot leginkább példákon keresztül lehet szemléletesen bemutatni⁵:

- **Egy globális technológiai cég AI rendszeren alapuló toborzási rendszere hátrányosan megkülönböztette a női jelentkezőket**⁶ - az eltérés oka az volt, hogy az AI rendszer mögött álló algoritmust nagy arányban férfi alkalmazottak önéletrajzaiból álló adatokkal tanították, aminek többek között az is volt az oka, hogy a cég alkalmazottai között eleve magasabb arányban találhatóak férfiak. A példa amiatt is érdekes, mert rámutat, hogy még egy, a digitális technológiák terén élenjáró vállalat is eshet ilyen hibába.

⁴ European Union Agency for Fundamental Rights (2020): Getting The Future Right Artificial Intelligence and Fundamental Rights

⁵ James Manyika, Jake Silberg, and Brittany Presten (2019): What Do We Do About the Biases in AI?, Harvard Business Review

⁶ Jeffrey Dastin (2018): Amazon scraps secret AI recruiting tool that showed bias against women, Reuters

- **Két szakszervezet is támadta az egyik személyszállítással foglalkozó vállalatot a sötétebb bőrtónusú sofőröket nem felismerő AI rendszere miatt**⁷ - A cég biztonsági okok miatt bevezetett egy AI rendszert, ami bizonyos időközönként ellenőrizte a sofőrök személyazonosságát úgy, hogy az okostelefonjuk élő kamera képét összehasonlította a sofőr arcáról előzetesen eltárolt fotókkal. Két sofőr szakszervezet is azt állítja az AI rendszer alacsonyabb hatékonysággal működik sötétebb bőrtónusú sofőrök esetén, amely szélsőséges esetben ki is zárhatja a sofőröket a rendszerből, így gyakorlatilag elvesztik a munkájukat.
- **Irányítószám szerinti megkülönböztetés hitelképesség elbírálásakor**⁸ - Számos pénzügyi szolgáltatónál felmerülő AI rendszer rendellenesség, hogy más paramétereiktől függetlenül, a szegényebb irányítószámú településeken vagy kerületekben élő ügyfelek alacsonyabb „hitelképességi” pontszámot kapnak, így magas kamat mellett jutnak hitelhez, vagy el is eshetnek a hitelhez jutás lehetőségétől.

Ezek a prekoncepciók ráadásul üzleti értelemben is hátrányt jelenthetnek egy vállalat számára, a fenti példákra reflektálva:

- A toborzási rendszer hibája miatt tovább lehetnek betöltetlenek a nyitott pozíciók, ami miatt projektek késhetnek, ami bevételkieséshez vezethet
- A belső folyamatok indokolatlan megakasztása költségtöbbletet okozhat
- A hitelkérelmek elbírálásakor kiszűrésre kerülhetnek olyan ügyfelek, akiket egyébként ki tudna szolgálni a pénzügyi szolgáltató, ami bevétel kieséshez vezethet

III. Data governance jó gyakorlatok

A fenti problémák kezelésére azonban léteznek jógyakorlatok. A megfelelő data governance teremti meg az AI rendszerhez kapcsolódó adatok kezelésének szervezeti kereteit, kijelölve azokat a felelősöket és döntési pontokat, valamint felállítva azokat a vállalaton belüli intézményeket, amelyek segítenek kiküszöbölni a fent említett esetek bekövetkezését.

⁷ Chris Vallance (2021): Legal action over alleged Uber facial verification bias

⁸ Sian Townson (2020): AI Can Make Bank Loans More Fair, Harvard Business Review

Az EY is rendelkezik saját *data governance* módszertannal⁹. Eszerint egy jól átgondolt *data governance* keretrendszer kialakítása során figyelembe kell venni legalább:

- **Stratégiai kérdések:**
 - A vállalat stratégiáját, technológia/IT stratégiáját és mesterséges intelligencia stratégiáját (ha léteznek ezek a dokumentumok)
 - Az AI rendszerekhez kapcsolódó érintettek elvárásait
 - A vállalat által használt AI rendszerek és kapcsolódó felhasználási területek mennyiségét és milyenségét
 - A vállalat által kezelt adatvagyon méretét és milyenségét
 - A vállalat által kezelt adatvagyon monetizációs modelljét
- **Szervezési kérdések:**
 - A vállalat szervezeti struktúráját, bizottságait és más *governance* rendszerek működését
 - A szervezeti modellt, amelyben meg kívánják valósítani a *data governance* keretrendszert
 - A *data governance*-ért felelős szervezeti egység felállítását vagy kijelölését és átalakítását
 - A *data governance* keretrendszer egyéb szereplőinek kijelölését
 - A *data governance* keretrendszerhez kapcsolódó folyamatok kialakítását vagy a meglévő folyamatok átalakítását
- **Működési kérdések:**
 - A *data governance* keretrendszer működtetéséhez szükséges erőforrásokat
 - A *data governance* rendszer működését leíró mutatószámokat (KPI, KRI)
 - A *data governance* keretrendszer működtetéséhez szükséges informatikai támogatást
 - A kapcsolódó szabályzatok, kódexek, kézikönyvek és egyéb dokumentumok elkészítését vagy frissítését
 - Folyamatos fejlesztést és változásvezetés beépítését

⁹ EY Data Governance Playbook

A *data governance* egy remek terület lehet az *ESG rating* javítására is főként azoknak a vállalatoknak, amelyeknél az üzleti modell középpontjában állnak az adatok, pl.: IT, pénzügyi vagy telco szektorokban.

Illetve a *data governance* témakörébe tartozik számos informatika, informatikai biztonsági és adattudomány („*data science*”) területére vonatkozó kérdés is, amiket jelen cikknek nem célja tárgyalni.

IV. EU direktíva tervezet a *data governance*-ről

Ezen jó gyakorlatok egy részére reflektál az EU direktíva tervezete is kötelező elvárásként a magas kockázatú AI rendszerek esetén¹⁰. A direktíva a *data governance* kifejezésre magyarul az „adatkormányzás” kifejezést használja.

Az EU direktíva is kiköti, hogy a *data governance* keretrendszer kialakításánál legalább az alábbi döntési pontok, kompetenciák kijelöléséről és szervezeti keretek kialakításáról rendelkeznie kell a vállalatnak¹¹:

- Az AI rendszer megtervezésének főbb döntési pontjai;
- Adatgyűjtés mikéntje;
- Releváns adat-előkészítési műveletek, mint például annotáció, címkézés, tisztítás, gazdagítás és összesítés;
- Vonatkozó hipotézisek, előfeltevések megfogalmazása, különös tekintettel azokra az információkra, amelyeket az adatoknak mérniük kell és meg kell jeleníteniük;
- Szükséges adathalmaz elérhetőségének, mennyiségének és alkalmasságának előzetes értékelése;
- Esetleges torzítások – beleértve a prekonceptiók - meglétének vizsgálata (itt a direktíva eredeti, angol nyelvű szövege a „*bias*” kifejezést használja);
- Esetleges adathiányok vagy hiányosságok azonosítása, valamint e hiányok és hiányosságok kezelésének módja.

Érdemes itt kiemelni, hogy a direktíva tervezet is alkalmazza az adathalmaz hármass bontását a tréning, validációs és tesztelési adatokra. Így a fenti feltételeket mindhárom adattípus esetén biztosítani kell.

¹⁰ European Commission: Regulatory framework proposal on artificial intelligence, Link: <https://digital-strategy.ec.europa.eu/en/policies/regulatory-framework-ai>

¹¹ European Commission: Regulatory framework proposal on artificial intelligence, Article 10.

V. Data governance és belső ellenőrzés

A *data governance* keretrendszer megfelelőségének felmérése a belső ellenőrzés számára is jelenthet kihívásokat. Ilyen területek lehetnek^{12, 13}:

- Törzsadatok („*master data*”) ellenőrzése
- Metaadatok ellenőrzése
- Back-up adatok ellenőrzése
- Üzleti és IT terület közti megállapodások, SLA-k ellenőrzése
- Adatminőség ellenőrzése

A kihívások áthidalásával azonban jelentős értéket teremthet a belső ellenőrzés, hiszen garanciát tud nyújtani az üzleti vezetőknek, hogy a döntéseket alátámasztó adatok megfelelőek, így bátran használhatók új *insight*-ok előállítására.

VI. Összefoglalás

A *data governance* keretrendszer tehát az egész *AI Ethics* téma alapját, gerincét adja. Ez biztosítja, hogy az AI rendszerek a megfelelő adatokkal legyenek ellátva és így segíteni tudják azt az üzleti célt, amelyre létrehozták őket. Mivel a *data governance* keretrendszer outputja – a megfelelően előkészített adat és az optimalizált algoritmus – az AI rendszer inputját jelenti, így az itteni nemmegfelelőségeknek a teljes AI rendszer szempontjából jelentős tovaryűrűző hatása lehet. Érdemes tehát az adott vállalatnak az *AI Ethics* felkészülését a *data governance* keretrendszer kiépítésével kezdeni.

¹² IIA Belgium (2020): Data Ethics, Where does internal audit fit?, Link: <https://iia.be/nl/wp-content/uploads/2020/08/GKB-Data-Ethics.pdf>

¹³ IIA Global (2017): Understanding and Auditing Big Data, Link: <https://www.iiainet.org/SiteFiles/Publications/GTAG-Understanding-and-Auditing-Big-Data.pdf>