

نقشه‌برداری مرتبط جمع آوری داده ها مکانیسم هایی برای آموزش هوش مصنوعی

مصنوعی OECD
اسناد اطلاعاتی
اکتبر ۲۰۲۵ شماره ۴۸

مقالات هوش مصنوعی OECD

نقشه برداری از مجموعه داده های مرتبط مکانیسم های آموزش هوش مصنوعی



پیشگفتار

این مقاله توسط سرگی گالوز دوران، به سرپرستی کلاریس ژیروت، با پشتیبانی کریستین ریمزباخ-کوناتزه نوشته شده است. بررسی سرمقاله و کمک برای انتشار توسط آندریا فورتادوارائه شده است. این مقاله شامل بازخورد نمایندگان کمیته سیاست دیجیتال (DPC) OECD و همچنین نمایندگان گروه کاری آن در زمینه حاکمیت داده و حریم خصوصی (DGP) است. مشارکت جهانی در زمینه هوش مصنوعی (GPAI) نیز در طول جلسه عمومی پاییزی 2024 خود در مورد این کار بحث کرد. نویسنده با کمال تشکر از نظرات ارزشمند لیمور شمزلینگ ماگازانیک، سارا بروبه، آنتونیا فون بورن-فالویس، سلین کایرا و کاسومی سوگیموتو قدردانی می کند. این مقاله به طور قابل توجهی از مشارکت های خردمندان اعضای گروه متخصص OECD در زمینه هوش مصنوعی، داده و حریم خصوصی بهره برده است. این مقاله در 18 ژوئیه 2025 توسط کمیته سیاست دیجیتال (DPC) تصویب و از طبقه بندی خارج شد و توسط دبیرخانه OECD برای انتشار آماده شد.

یادداشت برای هیئت های نمایندگی:

این سند همچنین در ONE Members & Partners با کد مرجع زیر موجود است:

نهایی/4/2025(DGP/DPC/DSTI)

این سند، و همچنین هرگونه داده و نقشه ای که در اینجا آمده است، هیچ گونه خدشه ای به وضعیت یا حاکمیت بر هیچ قلمرویی، تعیین حدود و مرزهای بین المللی و نام هیچ قلمرو، شهر یا منطقه ای وارد نمی کند.

داده های آماری مربوط به اسرائیل توسط و تحت مسئولیت مقامات ذیصلاح اسرائیلی تهیه شده است. استفاده از چنین داده هایی توسط سازمان همکاری و توسعه اقتصادی (OECD) هیچ گونه خدشه ای به وضعیت بلندی های جولان، بیت المقدس شرقی و شهرک های اسرائیلی در کرانه باختری طبق شرایط حقوق بین الملل وارد نمی کند.

اصلاحیه های نشریات OECD را می توانید در اینجا بیابید. <https://www.oecd.org/en/publications/support/corrigenda.html>.

تصویر جلد: © Kjpgargetter/Shutterstock

© سازمان همکاری و توسعه اقتصادی ۲۰۲۵



انتساب ۴.۰ بین المللی (CC BY 4.0)

این اثر تحت مجوز Creative Commons Attribution 4.0 International در دسترس است. با استفاده از این اثر، شما موافقت می کنید که به شرایط این مجوز پایبند باشید (<https://creativecommons.org/licenses/by/4.0/>).

انتساب - شما باید به اثر استناد کنید.

ترجمه - شما باید به اثر اصلی استناد کنید، تغییرات در اثر اصلی را مشخص کنید و متن زیر را اضافه کنید: در صورت وجود هرگونه اختلاف بین اثر اصلی و ترجمه، فقط متن اثر اصلی باید معتبر تلقی شود.

سازگاری ها - شما باید به اثر اصلی ارجاع دهید و متن زیر را اضافه کنید: این اقتباسی از یک اثر اصلی از سازمان همکاری و توسعه اقتصادی (OECD) است. نظرات بیان شده و استدلال های به کار رفته در این اقتباس نباید به عنوان نماینده دیدگاه های رسمی سازمان همکاری و توسعه اقتصادی یا کشورهای عضو آن گزارش شود.

مطالب شخص ثالث - این مجوز شامل مطالب شخص ثالث در اثر نمی شود. در صورت استفاده از چنین مطالبی، شما مسئول اخذ اجازه از شخص ثالث و هرگونه ادعای نقض حق نشر هستید.

شما مجاز به استفاده از لوگو، هویت بصری یا تصویر روی جلد OECD بدون اجازه صریح یا بدون اشاره به اینکه OECD استفاده شما از اثر را تأیید می کند، نیستید.

هرگونه اختلاف ناشی از این مجوز، از طریق داوری مطابق با قوانین داوری دیوان دائمی داوری (PCA) مصوب ۲۰۱۲ حل و فصل خواهد شد. محل داوری پاریس (فرانسه) خواهد بود. تعداد داوران یک نفر خواهد بود.

چکیده

هنگام توسعه سیستم های هوش مصنوعی، متخصصان اغلب بر ساخت مدل تمرکز می کنند، در حالی که گاهی اوقات اهمیت تجزیه و تحلیل مکانیسم های متنوع جمع آوری داده ها را دست کم می گیرند. با این حال، تنوع مکانیسم های مورد استفاده برای جمع آوری داده ها شایسته توجه دقیق تری است زیرا هر یک از آنها پیامدهای متفاوتی برای توسعه دهندگان هوش مصنوعی، صاحبان داده ها و سایر دارندگان حقوقی که داده های آنها جمع آوری شده است، دارد. این مقاله سیاستی، مکانیسم های اصلی مورد استفاده در حال حاضر برای تهیه داده ها برای آموزش سیستم های هوش مصنوعی را ترسیم می کند و یک طبقه بندی برای پشتیبانی از مباحث سیاستی پیرامون حریم خصوصی، مدیریت داده ها و توسعه مسئولانه هوش مصنوعی پیشنهاد می دهد.

فهرست مطالب

۲	پیشگفتار
۳	چکیده
۵	خلاصه اجرایی
۷	امقدمه
۸	۲ مکانیسم های جمع آوری داده ها برای آموزش هوش مصنوعی
۱۲	۳ داده های جمع آوری شده به طور مستقیم از افراد و سازمان ها
۱۲	۳.۱. داده های ارائه شده و داده های مشاهده شده
۱۳	۳.۲. کمک های داوطلبانه ی داده ها
۱۵	۴ داده های جمع آوری شده از ارائه دهندگان شخص ثالث
۱۵	۴.۱. داده های جمع آوری شده از اشخاص ثالث بر اساس توافقات تجاری
۱۶	۴.۲. داده های جمع آوری شده از اشخاص ثالث بر اساس رویه های غیرتجاری
۱۹	۵ نتیجه گیری
۲۰	منابع
۲۵	یادداشت ها
	ارقام
۱۰	شکل ۲.۱. چرخه حیات توسعه مدل هوش مصنوعی
۱۰	شکل ۲.۲. حوزه های شخصی، اختصاصی و عمومی داده ها شکل ۲.۳.
۱۱	سازوکارهای کلیدی جمع آوری داده ها برای آموزش هوش مصنوعی شکل ۴.۱.
۱۶	درجه باز بودن داده ها
	جعبه ها
۱۷	کادر ۴.۱. درجات باز بودن داده ها

خلاصه اجرایی

عملکرد و قابلیت اطمینان مدل های هوش مصنوعی (AI) ارتباط نزدیکی با کیفیت و تنوع داده های مورد استفاده در آموزش آنها دارد. در حالی که جنبه های فنی توسعه مدل اغلب در مرکز توجه قرار می گیرند، **روش های اساسی برای تهیه و جمع آوری داده های آموزشی نیز به همان اندازه مرتبط هستند.** مکانیسم های مختلف جمع آوری داده ها، مزایا و چالش های متمایزی را نه تنها برای توسعه دهندگان هوش مصنوعی که به دنبال مدل های قوی و نماینده تر هستند، بلکه برای افرادی که داده هایشان ممکن است شامل شود، به همراه دارد. در عمل، توسعه دهندگان هوش مصنوعی اغلب از چندین مکانیسم جمع آوری داده به طور همزمان برای ساخت مجموعه داده های آموزشی جامع استفاده می کنند. **درک این مکانیسم ها برای پیشرفت سیستم های هوش مصنوعی قابل اعتماد و برای پرداختن به ملاحظات مربوط به حریم خصوصی و مدیریت داده ها در فرآیند توسعه ضروری است.**

براین اساس، این مقاله نقشه برداری و طبقه بندی از سازوکارهای اصلی که در حال حاضر برای تهیه داده ها برای آموزش سیستم های هوش مصنوعی استفاده می شوند را ارائه می دهد. هدف این طبقه بندی، ارائه مبنایی برای تجزیه و تحلیل های آینده در مورد پیامدهای حریم خصوصی و مدیریت داده ها در هر سازوکار است.

این طبقه بندی، این مکانیسم های کلیدی جمع آوری داده ها را در ساختار زیر سازماندهی می کند:

۱. داده های جمع آوری شده به طور مستقیم از افراد و سازمان ها

- **داده های ارائه شده و مشاهده شده:** حجم فزاینده ای از داده های آموزشی از داده هایی که توسط افراد ارسال می شود یا به صورت غیرفعال در طول تعاملات آنها با سیستم های هوش مصنوعی، به ویژه در محیط های کسب و کار به مصرف کننده (B2C) مانند چت بات ها، دستیاران مجازی و میزهای کمک خودکار، جمع آوری می شود، سرچشمه می گیرد. علاوه بر این، برخی از توسعه دهندگان هوش مصنوعی، مانند پلتفرم های رسانه های اجتماعی، ممکن است از داده های ارائه شده یا مشاهده شده از افراد در سراسر مجموعه گسترده تر خود برای پشتیبانی از آموزش مدل هوش مصنوعی استفاده کنند.
- **کمک های داوطلبانه ی داده ها:** اگرچه هنوز در حال ظهور است، اما مشارکت های داوطلبانه افراد یا سازمان ها در جمع آوری داده ها، پتانسیل غنی سازی مجموعه داده های آموزشی با اطلاعات متنوع و واقعی را فراهم می کند که در غیر این صورت دسترسی به آنها دشوار خواهد بود.

۲. داده های جمع آوری شده از ارائه دهندگان شخص ثالث

- **صادر مجوز داده های تجاری:** توافق نامه های صدور مجوز داده ها با سازمان ها، راه دیگری را برای توسعه دهندگان هوش مصنوعی جهت دسترسی به مجموعه داده ها فراهم می کند. بازارهای داده و دلالان داده نقش مهمی را به عنوان واسطه های داده در این اکوسیستم ایفا می کنند و دسترسی به طیف گسترده ای از داده های شخص ثالث را فراهم می کنند.
- **رویه های غیرتجاری:** توسعه دهندگان هوش مصنوعی همچنین ممکن است مجموعه داده ها را از طریق روش های غیرتجاری به دست آورند. ابتکارات داده باز، شامل داده های بخش دولتی و خصوصی که تحت مجوزهای باز منتشر می شوند، منابع کلیدی برای توسعه مدل های هوش مصنوعی هستند. مشارکت کنندگان مهم در این زمینه، ناشران مجموعه داده هستند که مجموعه داده ها را از منابع مختلف گردآوری و سازماندهی می کنند و آنها را به صورت رایگان و آشکار در دسترس قرار می دهند. با توجه به نیاز به مجموعه داده های بزرگ و متنوع برای پشتیبانی از فرآیندهای آموزش هوش مصنوعی، جمع آوری داده ها به عنوان یک مکانیسم جمع آوری داده که به طور گسترده برای رسیدگی به این نیازها پذیرفته شده است، ظهور کرده است.

با توسعه این طبقه بندی، این مقاله به سیاست گذاران و ذینفعان رویکردی ساختاریافته برای بحث های سیاسی در مورد **حریم خصوصی، مدیریت داده ها و توسعه قابل اعتماد هوش مصنوعی ارائه می دهد.** این خروجی، پیچیدگی و تنوع مکانیسم های جمع آوری داده ای را که توسعه دهندگان هوش مصنوعی به آنها متکی هستند، برجسته می کند و خاطرنشان می سازد که رویکردهای نوظهور شامل محیط های پردازش امن و ابزارهایی مانند ابزارهای افزایش حریم خصوصی

فناوری‌ها (PETs) روش‌هایی را برای بهبود قابلیت استفاده از این سازوکارهای جمع‌آوری داده‌ها ارائه می‌دهند، ضمن اینکه از حریم خصوصی و سایر حقوق و منافع مانند مالکیت معنوی محافظت می‌کنند. این طبقه بندی، زمینه را برای تجزیه و تحلیل بیشتر در مورد چگونگی ایجاد تعادل بین تقاضای رو به رشد برای داده‌های آموزشی هوش مصنوعی (از نظر حجم و تنوع) فراهم می‌کند، در حالی که جنبه‌های حریم خصوصی و مدیریت داده‌ها مانند کیفیت داده‌ها و قابلیت ردیابی را نیز در نظر می‌گیرد.

۱ مقدمه

ساخت مدل های یادگیری ماشینی هوش مصنوعی اغلب (البته نه به صورت سیستماتیک) به حجم زیادی از داده ها نیاز دارد که ممکن است عمده یا سهواً شامل داده های شخصی باشد. در حالی که نیازهای خاص داده ها بسته به نوع و هدف مدل متفاوت است، دسترسی بیشتر به داده ها عموماً مدل های هوش مصنوعی را قادر می سازد تا عملکرد بهتری داشته باشند، زیرا توانایی یادگیری از نقاط داده و الگوهای استنباط شده در یک فرآیند تکراری را دارند (OECD، 2022)^[1] داشتن داده های متنوع و باکیفیت (مثلاً دقت، کامل بودن، ثبات، قابلیت اطمینان، اعتبار، به موقع بودن) به همان اندازه در پشتیبانی از توسعه سیستم های هوش مصنوعی قوی و قابل اعتماد اهمیت دارد، زیرا داده های بهتر می توانند به رفع سوگیری ها، کاهش خطاها و محدود کردن نتایج ناخواسته کمک کنند.

هنگام توسعه سیستم های هوش مصنوعی، متخصصان اغلب بر ساخت مدل - از جمله وزن ها و پارامترها - تمرکز می کنند، در حالی که گاهی اوقات اهمیت تجزیه و تحلیل مکانیسم های متنوع جمع آوری داده ها را دست کم می گیرند. با این حال، تنوع مکانیسم های مورد استفاده برای جمع آوری داده ها شایسته توجه دقیق تری است زیرا هر یک از آنها پیامدهای متفاوتی برای توسعه دهندگان هوش مصنوعی، صاحبان داده ها و سایر دارندگان حق که داده های آنها جمع آوری شده است، دارد (2022، OECD)^[1] در همین راستا، توصیه نامه OECD در مورد هوش مصنوعی (از این پس، «اصول هوش مصنوعی OECD») (۲۰۱۹، OECD)^[2] بر اهمیت احترام به حریم خصوصی و اولویت بندی سازوکارهای خاص جمع آوری داده ها که مجموعه داده های با کیفیت بهتری ارائه می دهند و در مقایسه با سایر سازوکارها، به صاحبان داده ها کنترل بیشتری بر داده هایشان می دهند، تأکید دارند. همانطور که در ادامه بررسی می شود، روش هایی مانند جمع آوری داده ها اغلب بدون اطلاع افراد رخ می دهد و اعمال حقوق قابل اجرا، مانند حقوق قانونی برای اطلاع رسانی، دسترسی، حذف و اصلاح داده هایشان را دشوار می کند.

کارهای قبلی OECD، به ویژه «حفاظت از حریم خصوصی در اقتصاد داده محور» (OECD، ۲۰۱۴)^[3] و «افزایش دسترسی و اشتراک گذاری داده ها» (OECD، ۲۰۱۹)^[4]، داده ها را بر اساس نحوه جمع آوری آنها طبقه بندی می کند و بین داده های داوطلبانه، داده های مشاهده شده، داده های مشتق شده و داده های اکتسابی تمایز قائل می شود. طبقه بندی منشأ داده ها نیز در چارچوب OECD برای طبقه بندی سیستم های هوش مصنوعی گنجانده شده است (OECD، 2022)^[11]. با پیشرفت سیستم های هوش مصنوعی، آنها نه تنها به مقادیر زیادی داده، بلکه به مجموعه داده های متنوع و با کیفیت بالا نیز نیاز دارند. این نیاز دوگانه یک سوال کلیدی را مطرح می کند: چگونه می توانیم دسترسی و اشتراک گذاری داده ها را برای آموزش هوش مصنوعی - هم از نظر حجم و هم از نظر کیفیت - تسهیل کنیم و در عین حال از منافع داده های افراد محافظت کنیم؟

در این زمینه، سازمان همکاری و توسعه اقتصادی (OECD) در سال ۲۰۲۴ تصمیم گرفت یک تحلیل ساختاری در مورد مکانیسم های مختلف مورد استفاده برای جمع آوری داده های شخصی برای آموزش سیستم های هوش مصنوعی انجام دهد. چنین تحلیلی به ویژه امروزه اهمیت دارد، زیرا چارچوب های حفظ حریم خصوصی و مدیریت داده ها به طور فزاینده ای روش های جمع آوری داده ها برای ساخت سیستم های هوش مصنوعی توسط متخصصان هوش مصنوعی را شکل می دهند. این امر همچنین در آینده نیز مهم خواهد بود، زیرا با وجود عدم قطعیت پیرامون اشکال دقیق هوش مصنوعی آینده، آنها همچنان به شدت به داده ها متکی خواهند بود.

این مقاله با تکیه بر کارهای قبلی OECD و بینش های کلیدی چندرشته ای از گروه متخصصان OECD.AI در زمینه هوش مصنوعی، داده و حریم خصوصی، گزیده ای از سازوکارهای رایج جمع آوری داده برای آموزش هوش مصنوعی را ارائه می دهد. این مقاله دو هدف را دنبال می کند. اول، سازوکارهای کلیدی جمع آوری داده مورد استفاده در آموزش هوش مصنوعی را تشریح می کند و ویژگی های اصلی آنها را برجسته می سازد. هدف این است که پیشینه لازم را برای درک بهتر مزایا و چالش های بالقوه هر سازوکار در اختیار سیاست گذاران، تنظیم کنندگان و متخصصان قرار دهد. دوم، این مقاله پایه و اساس تحقیقات آینده در مورد پیامدهای حفظ حریم خصوصی و مدیریت داده ها در این سازوکارها، و همچنین ابتکارات سیاستی را که به دنبال ایجاد تعادل بین نیازهای حفظ حریم خصوصی و مدیریت داده ها با تقاضای رو به رشد برای داده ها در آموزش هوش مصنوعی هستند، بنامی نهد.

۲

مکانیسم های جمع آوری داده ها برای هوش مصنوعی آموزش

دولت هادر سراسر جهان در حال تدوین سیاست های متنوعی برای حمایت و هدایت توسعه سیستم های هوش مصنوعی قابل اعتماد و انسان محور هستند. در ژوئن ۲۰۲۴، اتحادیه اروپا قانون هوش مصنوعی را تصویب کرد (اتحادیه اروپا، ۲۰۲۴) [۵] و کمیسیون اروپا اخیراً پیش نویس دوم آیین نامه اجرایی هوش مصنوعی همه منظوره (کمیسیون اروپا، ۲۰۲۴) را منتشر کرده است. [۶] مجلس ملی کره همچنین «قانون توسعه هوش مصنوعی و ایجاد اعتماد» را تصویب کرده است (مجلس ملی کره، ۲۰۲۴) [۷] که قرار است از ژانویه ۲۰۲۶ لازم الاجرا شود. در دسامبر ۲۰۲۴، سنای فدرال برزیل لایحه شماره ۲۳۳۸/۲۰۲۳ را تصویب کرد (مجلس سنای برزیل، ۲۰۲۴) [۸] که هدف آن تنظیم توسعه و استفاده از هوش مصنوعی نیز هست. در ماه مه ۲۰۲۵، ژاپن قانون ارتقای تحقیق، توسعه و استفاده از فناوری های مرتبط با هوش مصنوعی را تصویب کرد (مجلس ملی ژاپن، ۲۰۲۵) [۹] که بر ترویج نوآوری هوش مصنوعی در عین مدیریت ریسک ها از طریق هماهنگی دولت و انطباق داوطلبانه صنعت تأکید دارد. یک نکته مشترک در این چارچوب های سیاستی، به رسمیت شناختن حاکمیت داده ها به عنوان امری مرتبط برای تضمین توسعه هوش مصنوعی قابل اعتماد است. در این زمینه، ابتکارات سیاستی بین المللی مانند «آیین نامه رفتار داوطلبانه گروه ۷» (۲۰۲۳، G7) [۱۰] بر اهمیت اجرای اقداماتی توسط توسعه دهندگان هوش مصنوعی تأکید می کنند که (۱) قابلیت ردیابی مجموعه داده ها را فراهم می کند؛ (۲) داده های آموزشی با کیفیت بالا را حفظ می کند؛ (۳) از داده های شخصی و حقوق مالکیت معنوی محافظت می کند؛ و (۴) با انتشار ارزیابی هایی از تأثیر مدل هوش مصنوعی بر حریم خصوصی و حفاظت از داده های شخصی، شفافیت و پاسخگویی را افزایش می دهد. OECD نقش محوری در پیشبرد آیین نامه رفتاری G7 و چارچوب گزارش دهی ایفا کرد که مورد دوم به عنوان یک پلتفرم آنلاین برای نظارت بر پذیرش داوطلبانه آیین نامه رفتاری توسط سازمان ها در فوریه ۲۰۲۵ راه اندازی شد (OECD، ۲۰۲۵) [۱۱].

یکی از جنبه های کلیدی مدیریت داده های هوش مصنوعی، فرآیند جمع آوری داده ها برای ساخت مجموعه داده های آموزشی هوش مصنوعی است. اندازه و پیچیدگی روزافزون مدل های مدرن هوش مصنوعی در مقیاس بزرگ، به ویژه مدل های زبان بزرگ (LLM)، ابرای آموزش مؤثر به مجموعه داده های گسترده ای نیاز است. از همه مهم تر، حوزه عمومی یا مجموعه داده های تحقیقاتی دارای مجوز ممکن است همیشه برای پشتیبانی از مقیاس مورد نیاز برای آموزش اکثر LLM های امروزی کافی نباشند (لی، ۲۰۲۳) [۱۲] در نتیجه، اندازه LLM ها چالش هایی در زمینه منبع یابی داده ها ایجاد می کند. به عنوان مثال، نگرانی عمده این است که LLM های مدرن در درجه اول بر اساس داده های استخراج شده از وب آموزش می بینند (فدر کوپر، ۲۰۲۳) [۱۳]. دعاوی و تحقیقات متعددی در سراسر جهان آغاز شده است که ادعا می کنند این مجموعه داده ها با نقض قوانین مربوط به حق چاپ، حریم خصوصی و حفاظت از داده ها جمع آوری شده اند که هنوز حل نشده اند.

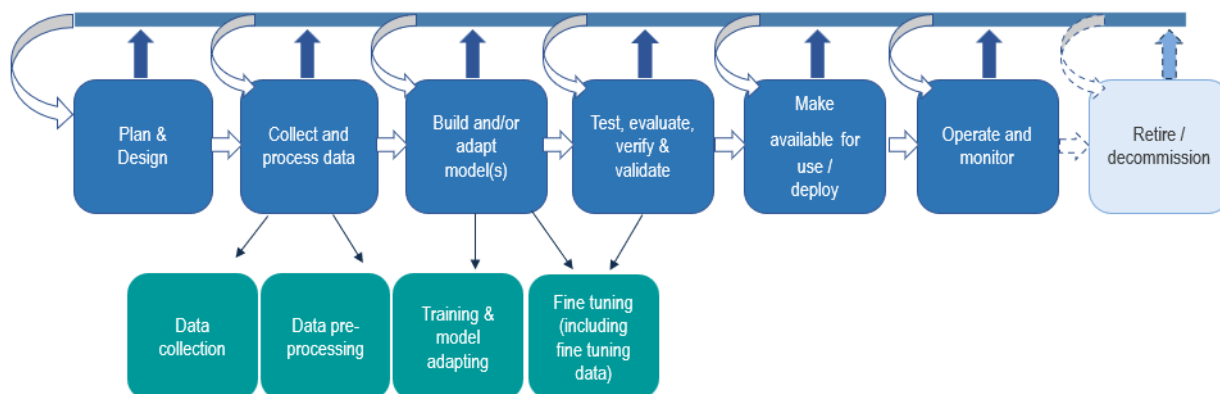
در حالی که توجه سیاست گذاران تا حد زیادی بر شیوه های جمع آوری داده ها برای LLM ها متمرکز شده است، امروزه جمع آوری داده ها تنها راه جمع آوری حجم زیادی از داده ها نیست. از آنجایی که هیچ بازار یا مخزن متمرکزی برای به دست آوردن داده های آموزشی وجود ندارد، شرکت هایی که مدل های هوش مصنوعی را توسعه می دهند به منابع داده متعددی متکی هستند. مجموعه داده های آموزشی دارای مجوز آزاد، یکی دیگر از مکانیسم های محبوب برای به دست آوردن داده های آموزشی هوش مصنوعی است. به عنوان مثال، پروژه گوتنبرگ، یک کتابخانه دیجیتال است که بیش از ۷۰،۰۰۰ کتاب الکترونیکی رایگان ارائه می دهد که عمدتاً شامل آثار ادبی قدیمی تر است که در حوزه عمومی قرار دارند (Gerlach, 2018) [۱۴]. هر مکانیسم جمع آوری داده ها با سطوح مختلفی از دسترسی، قابلیت همکاری، قابلیت اطمینان و کیفیت همراه است - ویژگی هایی که لزوماً منحصر به فرد نیستند. به عنوان مثال، دسترسی بالای مجموعه داده های دارای مجوز آزاد تضمین نمی کند که منشأ داده های زیربنایی به خوبی مستند شده باشد. جنبه دیگری از این وضعیت

این است که اغلب چندین بازیگر در جمع آوری مجموعه داده های آموزشی هوش مصنوعی مشارکت می کنند. به عنوان مثال، مجموعه داده های ImageNet در ابتدا توسط محققان دانشگاه استنفورد ایجاد شد که میلیون ها تصویر متنوع را از سراسر جهان جمع آوری کردند (دنگ، ۲۰۰۹) [15]. این محققان همچنین تصاویر را حاشیه نویسی کردند، فرآیندی که می تواند توسط سازمان هایی که در برچسب گذاری و حاشیه نویسی داده ها تخصص دارند نیز انجام شود. پس از حاشیه نویسی، مجموعه داده ها توسط سازمان غیرانتفاعی (ImageNet, nd) ImageNet گردآوری و در دسترس عموم قرار گرفت. [16] به عنوان منبعی برای آموزش مدل های مختلف هوش مصنوعی در کارهایی مانند طبقه بندی تصویر، تشخیص اشیا و تشخیص چهره عمل می کند. یک توسعه دهنده ممکن است از مجموعه داده های ارائه شده توسط ImageNet استفاده کند یا در طول مرحله آموزش، تغییرات بیشتری در آن ایجاد کند.

اغلب، دید محدودی به مجموعه داده های مورد استفاده برای ساخت مدل های هوش مصنوعی وجود دارد که ارزیابی کیفیت آنها رادشوار می کند. به عنوان مثال، شاخص شفافیت مدل بنیاد ۲۰۲۴، ۲۳ زیردامنه را مشخص می کند که جنبه های مختلف توسعه و استقرار مدل های پایه، از جمله دسترسی به داده ها را ارزیابی می کنند (بوماسانی، 2024) [17]. از بین تمام توسعه دهندگان مدل پایه که در این شاخص ارزیابی شده اند، تنها یکی از آنها داده های آموزشی خود را آشکارا به اشتراک گذاشته است (بوماسانی، 2024) [17]. نکته قابل توجه این است که بین انتشار اولین گزارش در اکتبر 2023 و پیگیری آن در مه 2024، توسعه دهندگان امتیاز شفافیت خود را در تمام زیردامنه ها به جز افشای داده ها بهبود بخشیدند (بوماسانی، 2024) [17]. پیشرفت کند در نمرات شفافیت پیرامون داده ها را می توان تا حدودی با خطرات قانونی که کسب و کارها هنگام افشای مجموعه داده های مورد استفاده برای ساخت مدل های هوش مصنوعی متحمل می شوند، و همچنین نگرانی ها در مورد اطلاعات اختصاصی توضیح داد، زیرا کیفیت یک مدل اغلب با ارزش داده هایی که بر اساس آن آموزش داده می شود، ارتباط نزدیکی دارد. فراتر از پیچیدگی ها و عدم قطعیت های قانونی مربوط به مقررات حفظ حریم خصوصی و حقوق مالکیت معنوی، نگرانی هایی نیز در مورد موانع ورود به بازار برای جمع آوری داده ها وجود دارد، زیرا برخی از شرکت ها ممکن است دسترسی ممتازی به مجموعه داده های به دست آمده از سایر فعالیت ها در بازارهای دیجیتال داشته باشند (Competition & Markets Authority، 2023) [18]. همچنین بحث های فزاینده ای در مورد محدودیت های موجودی داده های عمومی که از طریق اینترنت به عنوان منبعی برای توسعه هوش مصنوعی قابل دسترسی هستند، وجود دارد. مطالعات نشان می دهد که اگر روند فعلی گسترش اندازه مجموعه داده های آموزشی برای LLM ها به طور خاص ادامه یابد، مدل های هوش مصنوعی ممکن است بین سال های 2026 تا 2032 عرضه داده های متنی عمومی تولید شده توسط انسان را به پایان برسانند (پابلو ویلاوبوس، 2022) [19]. این کمبود بالقوه داده ها، سوالاتی را در مورد قابلیت اتکای بلندمدت عمدتاً به داده های عمومی موجود از اینترنت برای آموزش هوش مصنوعی مطرح می کند و اهمیت در نظر گرفتن طیف گسترده ای از مکانیسم های جمع آوری داده ها را افزایش می دهد. با این حال، روند روبه رشدی در صنعت به سمت توسعه مدل های بنیادی کوچک تر و کارآمدتر وجود دارد که به داده های کمتر یا مجموعه داده های دقیق تر متکی هستند. در نتیجه، اگرچه پتانسیل «تمام شدن داده ها» یک نگرانی قابل توجه است، اما ممکن است یک مانع غیرقابل مدیریت نباشد.

تصمیمات مربوط به داده های مورد استفاده برای ساخت مجموعه داده های آموزشی در مرحله اولیه توسعه مدل گرفته می شود. توسعه مدل های یادگیری ماشین هوش مصنوعی معمولاً شامل چهار مرحله است: جمع آوری داده ها، پیش پردازش، آموزش و تنظیم دقیق مدل (OECD، 2024) [20]. در طول فرآیند جمع آوری داده ها، توسعه دهندگان هوش مصنوعی در مورد کمیت و نوع داده هایی که باید در مجموعه داده های آموزشی گنجانده شوند و همچنین محل جمع آوری داده ها تصمیم می گیرند. انتخاب های انجام شده در طول توسعه مجموعه داده ها بر خروجی مدل هوش مصنوعی تأثیر خواهد گذاشت. در مرحله بعد، مرحله پیش پردازش داده ها، هر نقطه داده را به یک قالب قابل استفاده برای مدل یادگیری ماشین تبدیل می کند. این نقاط داده با هم ترکیب شده و در یک مجموعه داده برای آموزش جمع آوری می شوند. مرحله آموزش (یا پیش آموزش) مرحله اولیه ای است که در آن یک «مدل پایه» با استفاده از یک مجموعه داده گسترده ساخته می شود. این مدل پایه به عنوان پایه ای عمل می کند که تنظیمات و بهینه سازی های بیشتر در طول مرحله تنظیم دقیق، که معمولاً نیاز به استفاده از مجموعه داده های خاص دامنه دارد، بر روی آن انجام می شود. ۴.

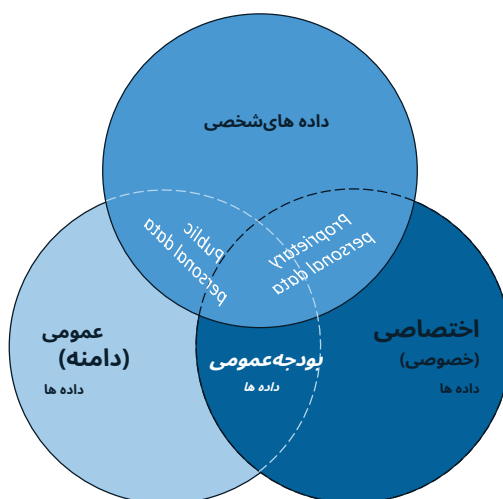
شکل ۲.۱. چرخه حیات توسعه مدل هوش مصنوعی



منبع: تصویر OECD بر اساس چرخه عمر سیستم هوش مصنوعی (OECD, 2024) (20) (سازمان همکاری و توسعه اقتصادی، ۲۰۲۵) (21)

هر دو مجموعه داده های آموزشی و تنظیم دقیق ممکن است شامل داده هایی باشند که دامنه های همپوشانی را در بر می گیرند. (OECD, ۲۰۱۹) [4] بین سه حوزه داده زیر تمایز قائل می شود: (i) دامنه شخصی که شامل تمام داده های شخصی «مربوط به یک فرد مشخص یا قابل شناسایی» است که صاحبان داده ها در مورد آنها منافع حریم خصوصی دارند؛ (ii) دامنه اختصاصی که شامل تمام داده های اختصاصی می شود که معمولاً منافع اقتصادی برای حذف دیگران وجود دارد؛ و (iii) دامنه عمومی که شامل تمام داده هایی است که دسترسی و استفاده مجدد از آنها رایگان است. این سه حوزه متقابلاً منحصر به فرد نیستند و اغلب همانطور که در شکل 2.2 نشان داده شده است، همپوشانی دارند که منعکس کننده منافع مختلف ذینفعان و چارچوب های حاکمیت داده قابل اجرا است (OECD, 2019) [22] (23). برای مثال، داده های عمومی می توانند شامل داده هایی باشند که توسط حقوق مالکیت فکری یا هرگونه حقوق دیگری با اثرات مشابه محافظت نمی شوند، و همچنین مجموعه داده های اختصاصی با بودجه عمومی که به طور مجاز مجوز گرفته اند.

شکل ۲.۲. دامنه های شخصی، اختصاصی و عمومی داده ها



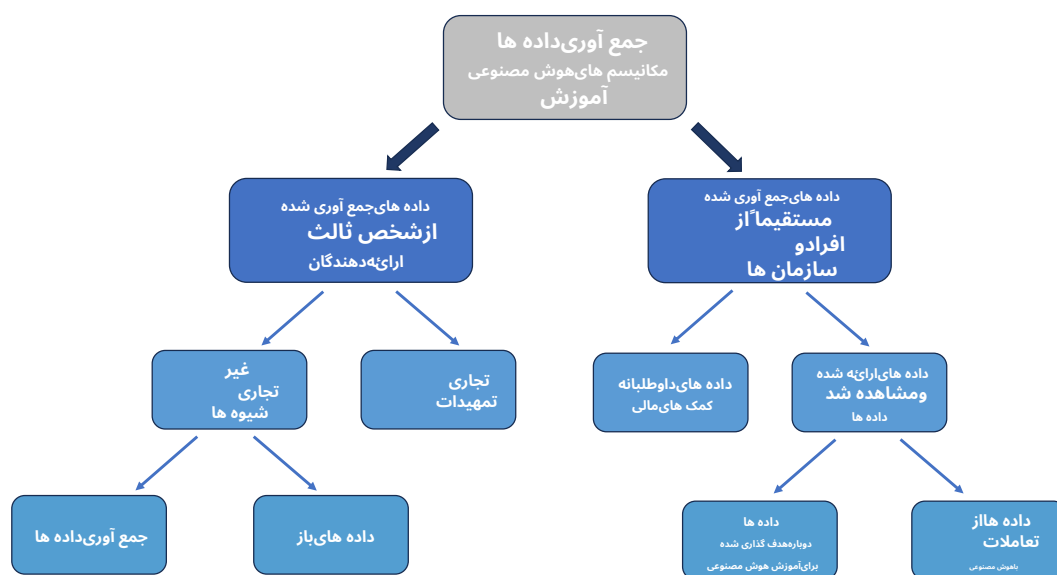
منبع: (OECD, ۲۰۱۹) [4] افزایش دسترسی و اشتراک گذاری داده ها: تطبیق خطرات و مزایای استفاده مجدد از داده ها در جوامع، انتشارات OECD، پاریس، <https://doi.org/10.1787/b4d546a9-en>

مرحله آموزش بر اندازه و تنوع مجموعه داده های آموزشی تمرکز دارد زیرا مدل های یادگیری ماشین با قابلیت های عمومی به داده های بزرگ و متنوع نیاز دارند. مجموعه داده های در مقیاس بزرگ، مدل هایی را با

پایه محکمی از دانش عمومی، که به آنها اجازه می دهد طیف متنوعی از خروجی ها را تولید کنند و بسیار مقیاس پذیر باشند. مدل های از پیش آموزش دیده را می توان به صورت اختیاری با استفاده از مجموعه داده های تخصصی برای به دست آوردن قابلیت های خاص، تنظیم دقیق کرد. در طول مرحله تنظیم دقیق، حجم داده های مورد نیاز معمولاً کمتر از مرحله آموزش است، زیرا هدف آن تطبیق یک مدل از پیش آموزش دیده با وظایف و دامنه های خاص است.

بخش های بعدی بر مکانیسم های جمع آوری داده مورد استفاده برای هر دو مرحله آموزش و تنظیم دقیق توسعه مدل یادگیری ماشین تمرکز دارند. آن ها مکانیسم های جمع آوری داده مرتبط برای آموزش هوش مصنوعی را با تمرکز ویژه بر ویژگی های هر روش، تشریح می کنند. این طبقه بندی ابتدا بر اساس منابع، موارد زیر را متمایز می کند: توسعه دهندگان هوش مصنوعی از طریق کدام منابع داده، داده ها را دریافت می کنند: (۱) مستقیماً از افراد و سازمان ها؛ و (۲) از ارائه دهندگان شخص ثالث. در این دو منبع داده اصلی، مکانیسم های کلیدی و مختلف جمع آوری داده ها را شناسایی می کنیم که در شکل زیر مشخص شده اند و در بخش های بعدی با جزئیات بیشتری مورد بحث قرار می گیرند:

شکل ۲.۳. مکانیسم های کلیدی جمع آوری داده ها برای آموزش هوش مصنوعی



توجه: این شکل، مروری بر سازوکارهای کلیدی جمع آوری داده ها برای آموزش هوش مصنوعی که در سند به آنها اشاره شده است، ارائه می دهد. این شکل لزوماً تمام روش های مختلف جمع آوری داده ها برای آموزش هوش مصنوعی را در بر نمی گیرد - اکوسیستم هوش مصنوعی پویا است و احتمالاً با ورود بازیگران و مدل های جدید کسب و کار به حوزه هوش مصنوعی، سازوکارهای جدید به تکامل خود ادامه خواهند داد.

لازم به ذکر است که PET ها نقش مهمی در فراهم کردن امکان جمع آوری داده ها برای توسعه هوش مصنوعی، کاهش خطرات مربوط به مدیریت داده ها و حریم خصوصی ایفا می کنند. PET ها با افزایش محرمانگی جمع آوری و استفاده از داده ها، از عملکرد مؤثر سازوکارهای ذکر شده در این طبقه بندی پشتیبانی می کنند. به عنوان مثال، داده های مصنوعی به طور فزاینده ای به عنوان وسیله ای برای جمع آوری محرمانه داده ها در نظر گرفته می شوند (OECD، 2025). [24] این طبقه بندی به طور خاص بر مکانیسم های جمع آوری داده های دنیای واقعی تمرکز دارد که اغلب مبنای تولید مجموعه داده های مصنوعی را تشکیل می دهند.

۳

داده‌های جمع‌آوری شده به طور مستقیم از افراد و سازمان‌ها

۳.۱ داده‌های ارائه شده و داده‌های مشاهده شده

داده‌های ارائه شده به اطلاعاتی اشاره دارد که از اقدامات مستقیم انجام شده توسط یک فرد ناشی می‌شود، به طوری که فرد کاملاً از اقداماتی که منجر به تولید داده‌ها می‌شود، آگاه است (OECD، ۲۰۱۴) [۳]. این داده‌ها معمولاً از طریق اقدامات مستقیم کاربرانند پر کردن فرم‌ها، ثبت نام در خدمات، ارائه بازخورد یا آپلود محتوا به صورت آنلاین جمع‌آوری می‌شوند. این داده‌ها ممکن است شامل انواع مختلفی از داده‌ها، از اطلاعات شخصی اولیه گرفته تا ورودی‌های پیچیده‌تر مانند تنظیمات برگزیده، رتبه‌بندی‌ها یا محتوای تولید شده توسط کاربر (مثلاً عکس، متن و ویدیو) باشد. در حالی که افراد مربوطه ممکن است از پیامدهای ارائه این داده‌ها بی‌اطلاع باشند، این واقعیت که این داده‌ها ایجاد می‌شوند باید واضح - یا حداقل شهودی - باشد (OECD، ۲۰۱۴) [۳].

داده‌های مشاهده شده داده‌هایی هستند که توسط دیگران مشاهده شده و در قالب دیجیتال ثبت شده‌اند. این داده‌ها می‌توانند در لحظه ایجادشان ثبت شوند یا پس از مشاهده به یک حامل دیجیتال منتقل شوند (OECD، ۲۰۱۴) [۳]. این نوع جمع‌آوری داده‌ها معمولاً از طریق ردیابی دیجیتال (مثلاً کوکی‌های آنلاین، برنامه‌های ردیابی موقعیت مکانی) یا نظارت فیزیکی (مثلاً دوربین‌های مداربسته) انجام می‌شود. در حالی که ممکن است افراد از ایجاد داده‌های مشاهده شده آگاه شوند (مثلاً به دلیل تعامل فعال)، بخش زیادی از ایجاد داده‌های مشاهده شده ممکن است مورد توجه قرار نگیرد.

در زمینه سیستم‌های هوش مصنوعی، تمایز بین داده‌های ارائه شده و مشاهده شده اغلب مبهم است. هنگامی که افراد با سیستم‌های هوش مصنوعی تعامل می‌کنند، معمولاً هر دو نوع داده به طور همزمان جمع‌آوری و پردازش می‌شوند. به عنوان مثال، هنگامی که کاربر در یک چت بات (داده‌های ارائه شده) درخواستی را وارد می‌کند، سیستم ممکن است در همان جلسه، فراداده‌هایی مانند مهرهای زمانی، مدت زمان تعامل، اطلاعات دستگاه یا الگوهای تعامل کاربر (داده‌های مشاهده شده) را نیز ثبت کند. این همپوشانی عملی به این معنی است که در برنامه‌های دنیای واقعی، ترسیم مرز مشخصی بین این دو دسته می‌تواند دشوار باشد. بنابراین، اغلب تجزیه و تحلیل دقیق‌تر در سطح زیرگروه ضروری است.

براین اساس، در بخش بعدی بین دو سناریو تمایز قائل می‌شویم: اول، داده‌هایی که توسط افراد در طول تعامل مستقیم با خدمات یا محصولات هوش مصنوعی ارائه یا مشاهده می‌شوند؛ و دوم، داده‌هایی که در ابتدا در زمینه سایر فعالیت‌های دیجیتال جمع‌آوری شده‌اند اما متعاقباً برای آموزش هوش مصنوعی استفاده می‌شوند.

۳.۱.۱ داده‌های ارائه شده توسط یا مشاهده شده از افراد هنگام تعامل مستقیم با سیستم‌های هوش مصنوعی

یک مکانیسم فزاینده مرتبط برای آموزش سیستم‌های هوش مصنوعی، استفاده از اطلاعات تولید شده از پیام‌ها در ابزارهای هوش مصنوعی کسب و کار به مصرف‌کننده (B2C) مانند چت بات‌ها، دستیاران مجازی و سیستم‌های خدمات مشتری خودکار است. این چت بات‌های هوش مصنوعی به طور فزاینده‌ای "داده‌های تعاملی" تولید می‌کنند - محتوایی که هنگام تعامل یک سیستم هوش مصنوعی با یک کاربر تولید می‌شود. با ادغام بیشتر سیستم‌های هوش مصنوعی در محیط‌های مکالمه، حجم و اهمیت این داده‌های تعاملی همچنان رو به افزایش است. داده‌های تعاملی همچنین می‌توانند از طریق ادغام‌های کسب و کار به کسب و کار (B2B) با سرویس‌های میزبانی شده، که در آن قابلیت‌های هوش مصنوعی

از طریق مشارکت مستقیم یا دسترسی به API در محصولات تعبیه شده اند. به عنوان مثال، ChatGPT از طریق یک مشارکت استراتژیک بین مایکروسافت و OpenAI در جستجوی مایکروسافت بینگ ادغام شده است.

دسترسی به داده های تعاملی محدود به شرکت هایی که تحت این مدل های کسب و کار B2C و B2B فعالیت می کنند، نیست. این نوع داده ها، که به ویژه برای تنظیم دقیق مدل های هوش مصنوعی مفید هستند، از طریق بازارهایی مانند promptbase.com در دسترس توسعه دهندگان هوش مصنوعی قرار دارند.

علاوه بر داده های تعاملی، شرکت هایی که خدمات و محصولات هوش مصنوعی ارائه می دهند، به داده های ارائه شده توسط کاربر یا قابل مشاهده، مانند اطلاعات سطح حساب یا فراداده نیز دسترسی دارند. این شرکت ها ممکن است استفاده از این نوع اطلاعات را برای آموزش مدل های هوش مصنوعی خود مفید ببینند.

۳.۱.۲ داده های (ارائه شده و/یا مشاهده شده) که برای آموزش هوش مصنوعی مورد استفاده مجدد قرار گرفته اند

علاوه بر داده های جمع آوری شده در طول استفاده از ابزارهای هوش مصنوعی، برخی از توسعه دهندگان هوش مصنوعی ممکن است به داده هایی (چه توسط کاربران ارائه شده باشد و چه از تعاملات آنها مشاهده شده باشد) از سایر سرویس های دیجیتال خود نیز دسترسی داشته باشند و ممکن است به دنبال استفاده از این داده ها برای آموزش مدل های هوش مصنوعی باشند. این مورد در مورد صاحبان پلتفرم های بزرگی که میزبان محتوای تولید شده توسط کاربر هستند، صدق می کند (FTC، 2024).^[25] این پلتفرم ها انگیزه های تجاری بسیار قدرتمندی برای استفاده از محتوایی مانند پست های رسانه های اجتماعی، نظرات، بررسی ها، سابقه خرید یا بازخورد کاربر برای بهبود مدل های هوش مصنوعی خود دارند.

با اعمال نفوذ داده های ارائه شده توسط کاربر توسعه دهندگان هوش مصنوعی می توانند سیستم ها و فرآیندهایی را طراحی کنند که احتمالاً به نیازها و ترجیحات دنیای واقعی پاسخ می دهند. این داده ها به بهبود عملکرد برنامه های هوش مصنوعی در صناعی مانند مراقبت های بهداشتی، آموزش، امور مالی یا سفر، که در آن ها شخصی سازی و پاسخگویی به نیازهای فردی مرتبط است، کمک می کند. به عنوان مثال، ویژگی های مبتنی بر هوش مصنوعی که توسط Booking.com برای شخصی سازی برنامه ریزی سفر ارائه می شوند، با استفاده از داده های ناشناس از رزرو هتل و جزئیات رزرو آموزش داده می شوند (گلدنبرگ، 2021).^[26]

داده های مشاهده شده داده های جمع آوری شده، به ویژه هنگامی که از طریق فناوری های ردیابی آنلاین جمع آوری می شوند، نقش مهمی در آموزش مدل های هوش مصنوعی طراحی شده برای تجربیات شخصی سازی شده ایفا می کنند. سیستم های هوش مصنوعی از این داده ها یاد می گیرند تا ترجیحات فردی را درک کرده و رفتار آینده را پیش بینی کنند. به عنوان مثال، پلتفرم های تجارت الکترونیک مانند آمازون، داده های خرید مشتری - مانند ترجیحات، جستجوها و رفتار مرور - را به LLM های خود ارائه می دهند تا توصیه های محصول شخصی سازی شده تری ارائه دهند (لوی، 2024).^[27]

در برخی موارد، این اطلاعات ممکن است در ابتدا برای سایر خدمات دیجیتال جمع آوری شده باشد، نه آموزش مدل هوش مصنوعی. اینکه آیا چنین استفاده ثانویه ای مجاز است یا خیر، به قوانین حفظ حریم خصوصی و حفاظت از داده ها در حوزه قضایی مربوطه بستگی دارد و معمولاً نیاز به ارزیابی مبنای قانونی مربوطه برای پردازش دارد. به عنوان مثال، برزیل *Pessoais* (در تاریخ 2 ژوئیه 2024 یک اقدام پیشگیرانه صادر کرد و از متا خواست که فوراً استفاده از داده های شخصی برزیلی ها را از پلتفرم رسانه اجتماعی خود برای آموزش مدل هوش مصنوعی مولد خود متوقف کند و نگرانی هایی را در مورد مبنای قانونی چنین پردازشی تحت قانون حفاظت از داده های برزیل مطرح کرد (بادیلو، 2024) ANPD^[28]). مقاله آئی OECD تحلیل بیشتری در مورد این ملاحظات ارائه خواهد داد.

۳.۲ اهدافی داوطلبانه داده ها

داده های شخصی یا غیرشخصی ممکن است به صورت داوطلبانه توسط صاحبان داده ها یا دارندگان داده ها برای اهداف آموزشی هوش مصنوعی، بدون دریافت غرامت، و به نفع جامعه ارائه شوند. این عمل گاهی اوقات به عنوان «اهدای داده» یا «نوع دوستی داده» شناخته می شود، همانطور که در قوانینی مانند قانون حاکمیت داده اتحادیه اروپا (DGA) (اتحادیه اروپا، 2022) آمده است.^[29] این شامل افراد (یا سازمان ها) می شود که برای به اشتراک گذاشتن داده های خود با محققان، سازمان های غیرانتفاعی و سایر نهادها برای اهدافی مانند بهبود مراقبت های بهداشتی یا افزایش خدمات عمومی، رضایت می دهند (Kirstein، 2023).^[30]

اگرچه این رویه ها ممکن است سؤالاتی را ایجاد کنند،^۸ اهدای داده ها همچنین می تواند با فراهم کردن دسترسی به داده های متنوع و واقعی که در غیر این صورت دشوار یا پرهزینه خواهد بود، نقش مهمی در آموزش سیستم های هوش مصنوعی ایفا کند.

به عنوان مثال، در زمینه مراقبت های بهداشتی، نظارت بر محیط زیست و آموزش، اهدای داوطلبانه داده ها می تواند به ایجاد مجموعه داده های غنی تر و جامع تر کمک کند (Hirsch MC, 2020) [31]. این با اصول هوش مصنوعی OECD همسو است که دولت ها را تشویق می کند تا اعتماد داده ها و سایر سازوکارها را برای پشتیبانی از اشتراک گذاری ایمن، منصفانه، قانونی و اخلاقی داده ها ترویج دهند. همچنین با توصیه OECD در مورد افزایش دسترسی و اشتراک گذاری داده ها (OECD, 2021) همسواست. [32] که از انواع مختلف همکاری های داده ای و بهره برداری از منابع داده ای جدید و موجود برای تقویت اکتشافات و نوآوری های علمی مبتنی بر داده در بخش های خصوصی و دولتی پشتیبانی می کند.

در زمینه هایی مانند مراقبت های بهداشتی و ژنومیک، دسترسی به مجموعه داده های بزرگ برای توسعه راه حل های مبتنی بر هوش مصنوعی بسیار مهم است. در این زمینه، اهدای داده ها می تواند برای کمک به افراد جهت ارائه داده های شخصی خود (مثلاً اطلاعات ژنتیکی یا داده های سلامت) به طرح های تحقیقاتی با منافع عمومی، تسريع اکتشافات در توسعه دارو، پزشکی شخصی سازی شده و اپیدمیولوژی مورد استفاده قرار گیرد. به عنوان مثال، می توان به UK Biobank، یک پایگاه داده و منبع تحقیقاتی زیست پزشکی در مقیاس بزرگ حاوی اطلاعات ژنتیکی، سبک زندگی و سلامت و نمونه های بیولوژیکی از نیم میلیون شرکت کننده در بریتانیا (UK Biobank, nd) اشاره کرد. [33] جالب توجه است که مطالعه ای که تمایل به اشتراک گذاری داده های سلامت دیجیتال برای تحقیقات در آلمان و اسرائیل را اندازه گیری می کرد، نشان داد که اکثریت قابل توجهی از شهروندان در هر دو کشور (۸۲٪ از آلمانی ها و ۸۱٪ از اسرائیلی ها) نگرش مثبتی نسبت به ایجاد یک پایگاه داده متمرکز برای تحقیقات پزشکی ابراز کرده اند. البته تا زمانی که داده ها ناشناس باشند (Weisband, ۲۰۲۳) [34] (سازمان همکاری و توسعه اقتصادی، ۲۰۲۵) [35] در این زمینه، شایان ذکر است که قانون استفاده و دسترسی به داده ها (DUA) که در ۱۹ ژوئن ۲۰۲۵ به تصویب دارایی سلطنتی رسید، مورد توجه قرار گیرد. DUA توسط دولت بریتانیا برای ارتقای استراتژی دیجیتال بریتانیا و آزادسازی ارزش داده ها، بهره برداری از قابلیت های آن برای تقویت خدمات عمومی و کمک به اقتصاد بریتانیا در نظر گرفته شده است. در میان مفاد آن، DUA تصریح می کند که افراد می توانند «رضایت گسترده» خود را به یک حوزه تحقیقات علمی بدهند (ICO, 2025) [36].

به همین ترتیب، دولت ها و مؤسسات عمومی اغلب فاقد داده های لازم برای بهبود خدمات عمومی، توسعه شهرهای هوشمندتر یا مقابله با چالش های اجتماعی هستند. استدلال می شود که نوع دوستی داده ها می تواند با ارائه داده های لازم برای آموزش مدل های هوش مصنوعی که مسائلی مانند زیرساخت های عمومی، حمل و نقل یا آموزش را بهبود می بخشد، به حمایت از این تلاش ها کمک کند. به عنوان مثال، شهروندان داوطلبانه داده های مربوط به مصرف انرژی یا عادات حمل و نقل خود را به اشتراک می گذارند، می توانند به آموزش سیستم های هوش مصنوعی مورد استفاده برای بهینه سازی برنامه ریزی شهری، بهبود مدیریت ترافیک یا کاهش اتلاف انرژی در شهرهای هوشمند کمک کنند. یک مورد قابل توجه از نوع دوستی داده ها، مربوط به پروژه DECODE است. در این ابتکار، شهروندان بارسلونا با استفاده از حسگرهای محیطی که در داخل و خارج از خانه هایشان قرار داده شده بودند، داده هایی در مورد سر و صدا، آلودگی هوا، دما و رطوبت جمع آوری کردند. این داده های ارزشمند متعاقباً از طریق آنچه به عنوان "داده های مشترک" شناخته می شود، در دسترس عموم قرار گرفت (ساگارا، ۲۰۱۹) [37].

داده‌های جمع‌آوری شده از ارائه‌دهندگان شخص ثالث

وقتی داده‌ها مستقیماً از شخص مربوط به داده‌ها تهیه نمی‌شوند، توسعه‌دهندگان هوش مصنوعی که به دنبال داده‌ها هستند، معمولاً آن‌ها را از طریق شیوه‌های تجاری و غیرتجاری به دست می‌آورند. این بخش این دو مکانیسم را به عنوان وسیله‌ای برای تأمین داده‌ها برای اقتصاد هوش مصنوعی بررسی می‌کند.

۴.۱ داده‌های جمع‌آوری شده از اشخاص ثالث بر اساس توافقات تجاری

صدور مجوز داده‌های تجاری یکی از راه‌هایی است که توسعه‌دهندگان هوش مصنوعی به داده‌ها دسترسی پیدا می‌کنند. به عنوان مثال، OpenAI قراردادهای صدور مجوز را با رسانه‌های خبری، از جمله آسوشیتدپرس (O'Brien, 2023) امضا کرده است. [38] و اکسل اسپرینگر (OpenAI, ۲۰۲۳) [39] و همچنین با پلتفرم‌های رسانه‌های اجتماعی مانند ردیت (OpenAI, 2024) [40] بسیاری از توسعه‌دهندگان هوش مصنوعی، از جمله متا، قراردادهای مجوز داده را نیز با پلتفرم توزیع تصویر، شاتر استوک، امضا کرده‌اند (شاتر استوک، ۲۰۲۳) [41].

توسعه‌دهندگان هوش مصنوعی همچنین ممکن است برای دسترسی به مجموعه داده‌های شخص ثالث به واسطه‌هایی مانند بازارهای داده متکی باشند. این بازارهای داده، تبادل داده‌ها بین تأمین‌کنندگان و توسعه‌دهندگان هوش مصنوعی را تسهیل می‌کنند. برای روشن شدن موضوع، بازارهای داده از این نوع شامل پلتفرم‌هایی مانند [appen.com](https://www.appen.com)، [ai](https://www.ai.com)، یا [AWS](https://aws.com) [تبادل داده‌ها](#).

قراردادهای تجاری داده برای پشتیبانی از توسعه هوش مصنوعی همچنین می‌توانند شامل مدل‌های تجاری باشند که از دسترسی و اشتراک‌گذاری داده‌ها، مانند دلان داده، درآمد کسب می‌کنند. دلان داده، داده‌ها را از منابع مختلف جمع‌آوری و تجمیع می‌کنند و این داده‌ها را به اشخاص ثالث می‌فروشند. یکی از عوامل مهم تمایز بین دلان داده و ارائه‌دهندگان بازار داده این است که دلان داده به طور فعال در جمع‌آوری داده‌های اضافی و تجمیع آنها مشارکت دارند، در حالی که ارائه‌دهندگان بازار داده واسطه‌های غیرفعالی هستند که از طریق آنها کنترل‌کنندگان داده، از جمله دلان، می‌توانند مجموعه داده‌های خود را ارائه دهند (OECD, 2019) [44]. دلان داده مدت‌هاست که بخشی از اقتصاد داده محور بوده‌اند و در برخی حوزه‌های قضایی، شیوه‌های آنها توسط سازمان‌های نظارتی مورد تجزیه و تحلیل قرار گرفته است. یکی از فعالیت‌های فعلی آنها تهیه مجموعه داده‌های آماده هوش مصنوعی برای آموزش مدل هوش مصنوعی است (PR Newswire, 2024) [42].

مجموعه داده‌های توسعه هوش مصنوعی همچنین ممکن است توسط کسب و کارها به عنوان بخشی از فعالیت‌های تجاری ثانویه آنها به صورت تجاری ارائه شود. به عنوان مثال، یک تولیدکننده تجهیزات کشاورزی ممکن است از داده‌های حسگر جمع‌آوری شده از تراکتورها استفاده نکند، اما می‌تواند آن را به خریداران علاقه‌مند برای توسعه‌دهندگان هوش مصنوعی که روی راه‌حل‌های کشاورزی دقیق کار می‌کنند، بفروشد. اخیراً، بایر با مایکروسافت برای توسعه یک مدل هوش مصنوعی که با داده‌های بایر تنظیم شده است، همکاری امضا کرده است. این مدل برای ارائه بینش در مورد زراعت و حفاظت از محصولات در نظر گرفته شده است و برای توزیع کنندگان بایر، استارت‌آپ‌های فناوری کشاورزی و احتمالاً رقبا قابل مجوز خواهد بود. از جمله موارد دیگر، مدل هوش مصنوعی قادر به پاسخگویی به سؤالاتی در مورد موضوعاتی مانند اجزای حشره کش و مناسب بودن محصول برای پنبه خواهد بود (2024, Bousquette) [43].

۴.۲ داده های جمع آوری شده از اشخاص ثالث بر اساس رویه های غیرتجاری

۴.۲.۱ ترتیبات داده باز

داده های باز به سرعت به یک رویکرد برجسته برای افزایش دسترسی به داده ها تبدیل شده است (OECD، 2019) [44]. همچنین، در دسترس قرار دادن داده ها از بخش های دولتی و خصوصی تحت مجوزهای داده باز، برای رشد اقتصاد هوش مصنوعی به طور فزاینده ای اهمیت پیدا می کند. در بخش دولتی، نمونه هایی از آن شامل دستورالعمل داده باز اتحادیه اروپا (اتحادیه اروپا، ۲۰۱۹) است. [44] یا ابتکار ملت هوشمند سنگاپور (Smart Nation Singapore، nd) [45] دولت کره همچنین نقش فعالی در ایجاد و انتشار مجموعه داده های آموزشی هوش مصنوعی برای تقویت نوآوری داخلی داشته است، با هدف ایجاد پایه ای قوی برای توسعه مدل های هوش مصنوعی در بخش خصوصی. به عنوان مثال، این کشور پلتفرم AI Hub را راه اندازی کرده است که دسترسی آزاد به طیف گسترده ای از مجموعه داده ها را فراهم می کند (AI Hub، 2025) [46]. یک نمونه قابل توجه از بخش خصوصی، مجموعه داده COCO (اشیاء مشترک در متن) مایکروسافت است که شامل بیش از ۳۳۰۰۰۰ تصویر با حاشیه نویسی های دقیق برای ۸۰ دسته شیء است که تحت مجوز CC-BY 4.0 برای حاشیه نویسی ها قابل دسترسی است (اشیاء مشترک در متن، ۲۰۲۵) [47].

پیشرفت هادر توسعه هوش مصنوعی همچنین توسط ناشران مجموعه داده تسهیل می شود که داده ها را از منابع مختلف گردآوری و سازماندهی می کنند و دسترسی رایگان و آزاد به آنها را فراهم می کنند. در این راستا، ImageNet، یک مجموعه داده با مجوز آزاد از تصاویر در مقیاس بزرگ (J. Deng، ۲۰۰۹) [48] برای آموزش بسیاری از مدل های بینایی کامپیوتر اساسی بوده است. به همین ترتیب، پیشرفت های اخیر در تauxردگی پروتئین که توسط سیستم یادگیری عمیق AlphaFold حاصل شده است، تنها به لطف دسترسی آزاد به بانک داده های پروتئین، که تقریباً نیم قرن پیش تأسیس شد، امکان پذیر شد (لیزنفلد، 2023) [49]. مخازن مبتنی بر متن مانند Common Crawl، nd (Common Crawl) [50] همچنین مقادیر عظیمی از داده های وب را برای آموزش مدل های پردازش زبان طبیعی (NLP) ارائه می دهند. همچنین مخازنی از مجموعه داده های تنظیم دقیق وجود دارد که به صورت آزاد در دسترس هستند، مانند کتابخانه مجموعه داده های Hugging Face، nd (Hugging Face) [51].

اگرچه داده های باز (بدون محدودیت) ممکن است مطلوب باشند، شرایط حاکم بر اشتراک گذاری داده ها منجر به درجات مختلفی از باز بودن می شود، همانطور که در شکل ۴.۱ نشان داده شده و در کادر ۴.۱ در زیر توضیح داده شده است.

شکل ۴.۱. درجات باز بودن داده ها



منبع: بر اساس (OECD 2015) [52]، نوآوری مبتنی بر داده: کلان داده برای رشد و رفاه، <https://doi.org/10.1787/9789264229358-en>.

کادر ۴.۱. درجات باز بودن داده ها

باز بودن داده ها یک «مفهوم دوتایی» نیست. بلکه پیوستاری از درجات مختلف باز بودن است که از دسترسی بسته و استفاده فقط توسط دارنده داده (سطح ۰ در شکل ۳) تا موارد زیر را شامل می شود:

- سطح ۱: ترتیبات دسترسی و اشتراک گذاری مشروط که در آن دسترسی و اشتراک گذاری داده ها تابع «شرایطی است که شامل محدودیت هایی برای کاربران مجاز به دسترسی به داده ها (ترتیبات تبعیض آمیز)، شرایط استفاده از داده ها از جمله اهدافی که داده ها می توانند برای آنها استفاده شوند و الزامات مربوط به مکانیسم های کنترل دسترسی به داده ها که از طریق آنها دسترسی به داده ها اعطا می شود» می شود. به عنوان مثال، توافق نامه های دوجانبه صدور مجوز داده ها.
- سطح ۲: ترتیبات دسترسی و اشتراک گذاری بدون تبعیض داده ها، که در آن می توان به داده ها دسترسی پیدا کرد و با پرداخت هزینه به اشتراک گذاشت، اما «بر اساس شرایطی که مستقل از هویت کاربران داده ها هستند». به عنوان مثال، محتوای دارای حق اشتراک، مانند مقالاتی که در وب سایت های خبری با حق اشتراک منتشر می شوند یا مجموعه داده هایی که از طریق بازارهای داده آنلاین فروخته می شوند، معمولاً در این دسته قرار می گیرند. هر کسی که مایل به پرداخت هزینه باشد، می تواند صرف نظر از اینکه چه کسی است، به داده ها دسترسی پیدا کند. ابتکاراتی مانند فضای داده های سلامت اروپا و مراکز ملی داده های سلامت، به عنوان مثال دیگر، به نهادهای دولتی اجازه می دهند تا دسترسی به داده ها را تحت شرایط تعریف شده و در محیط های پردازش امن فراهم کنند. این ترتیبات همچنین ممکن است شامل هزینه هایی برای پوشش هزینه های عملیاتی باشد.
- سطح ۳: داده های باز (ترتیبات) به عنوان شکل افراطی باز بودن داده ها، که عبارتند از «ترتیبات دسترسی و اشتراک گذاری بدون تبعیض داده ها، که در آن داده ها قابل خواندن توسط ماشین هستند و می توانند به صورت رایگان به آنها دسترسی پیدا کرده و به اشتراک گذاشته شوند و توسط هر کسی برای هر هدفی مورد استفاده قرار گیرند، حداکثر با رعایت الزاماتی که یکپارچگی، منشأ، انتساب و باز بودن را حفظ می کنند» (OECD، ۲۰۲۵) [۵۳] برای مثال، مجموعه داده های علمی منتشر شده تحت مجوزهای Creative Commons که به صورت رایگان در دسترس همه قرار دارند.

منبع: بر اساس (OECD، ۲۰۲۵) [۵۳]، افزایش دسترسی و اشتراک گذاری داده ها در عصر هوش مصنوعی، و (OECD، ۲۰۱۵) [۵۲] نوآوری مبتنی بر داده: کلان داده برای رشد و رفاه <https://doi.org/10.1787/9789264229358-en>

میزان باز بودن داده ها به عوامل مختلفی بستگی دارد: محدودیت های صدور مجوز و استفاده، نگرانی های قانونی در مورد حریم خصوصی، حفاظت از داده ها، امنیت ملی و حقوق مالکیت معنوی. علاوه بر این، فرآیند آماده سازی، گردآوری و میزبانی مجموعه داده های باز در مقیاس بزرگ پرهزینه است. حفظ داده های باز با کیفیت بالا نیازمند سرمایه گذاری مداوم در تخصص انسانی و زیرساخت های فناوری است و همین امر، پشتیبانی مالی بلندمدت را به یک چالش تبدیل می کند (باک، ۲۰۲۵) [۵۴].

لازم به ذکر است که چارچوب باز بودن داده ها از یک سو، برای منابع داده اصلی و از سوی دیگر، برای مجموعه داده ایجاد شده از ترکیبی از این منابع داده اعمال می شود. این بدان معناست که یک مجموعه داده آموزشی ممکن است شامل داده هایی با سطوح مختلف باز بودن باشد، اما مجموعه داده همچنان بسته باقی می ماند. بسیاری از LLM های محبوب بر روی داده هایی با درجات مختلف باز بودن (مانند داده های وب اسکریپینگ و دارای مجوز) آموزش دیده اند، اما مجموعه داده های آموزشی آنها همچنان اختصاصی است. شرح مجموعه داده ممکن است در دسترس باشد، اما دسترسی به مجموعه داده کامل یا استفاده از آن امکان پذیر نیست. به عنوان مثال، گزارش فنی GPT-4 OpenAI بیان می کند که GPT-4 از قبل آموزش دیده است. *با استفاده از داده های عمومی (مانند داده های اینترنتی) و داده های دارای مجوز از ارائه دهندگان شخص ثالث* (OpenAI، ۲۰۲۳) [۵۵]. با این حال، خود مجموعه داده های آموزشی بسته باقی مانده است. این ممکن است به دلیل ترکیبی از عوامل، مانند منافع تجاری در حفظ مالکیت معنوی و اختصاصی مجموعه داده ها یا ملاحظات صدور مجوز مربوط به نحوه تهیه، سازماندهی و اشتراک گذاری داده ها باشد. پاسخ های سیاستی بالقوه برای پشتیبانی از در دسترس بودن داده ها برای توسعه هوش مصنوعی در مقاله ای که به زودی توسط OECD منتشر خواهد شد، بررسی خواهد شد.

۴.۲.۲. استخراج داده های اینترنتی عمومی

استخراج داده ها، که به آن استخراج وب نیز گفته می شود، استخراج خودکار داده های وب قابل دسترس عموم با استفاده از یک عامل نرم افزاری ("ربات") است. استفاده از عامل های نرم افزاری برای مرور خودکار وب چیز جدیدی نیست و سال هاست که در بسیاری از زمینه ها به کار گرفته می شود. به عنوان مثال، سایت های مقایسه پروازهای هواپیمایی از "اسکرپینگ صفحه" برای اسکن وب سایت های هواپیمایی برای قیمت ها استفاده می کنند (ویتاکر، 2024) [56] به طور مشابه، «خزنده های وب» توسط موتورهای جستجو مانند گوگل برای فهرست بندی محتوای آنلاین و پیوند دادن کاربران به صفحات مرتبط استفاده می شوند (گوگل، و [57]).

از منظر نحوه ی ثبت داده ها، می توان خراشیدن صفحه نمایش را از خزیدن در وب متمایز کرد. خراشیدن صفحه نمایش شامل دانلود داده های نمایش داده شده در وب سایت ها است، در حالی که خزیدن در وب بر پیوند دادن وب سایت ها از طریق تجزیه و تحلیل کلمات کلیدی و فراداده (بدون نیاز به دانلود محتوای آن وب سایت ها) تمرکز دارد. تفاوت های تحلیلی بیشتر بین این فرآیندها در مقاله (OECD، 2025) مورد بحث قرار گرفته است. [21] «مسائل مربوط به مالکیت معنوی در هوش مصنوعی آموزش دیده بر روی داده های خراشیده».

جذابیت جمع آوری داده ها به عنوان ابزاری برای جمع آوری داده ها، منجر به استفاده گسترده از ربات های جمع آوری داده ها برای اهداف تجاری و غیرتجاری، از جمله آموزش مدل هوش مصنوعی، شده است. جالب توجه است که مدل های سنتی یادگیری ماشینی در درجه اول با استفاده از مجموعه داده های ساختاریافته ای که برای وظایف خاص انتخاب شده بودند، آموزش داده می شدند. با این حال، LLM ها برای درک پیچیدگی های زبان طبیعی به مجموعه داده های بزرگتر و متنوع تری نیاز دارند. بنابراین، آنها معمولاً بر اساس داده های جمع آوری شده از وب آموزش داده می شوند (لی، 2023) [12].

شرکت های فناوری و اپراتورهای پلتفرم، مانند پلتفرم های رسانه های اجتماعی، موتورهای جستجو و سایت های تجارت الکترونیک، هم منبع داده برای استخراج داده ها و هم مشارکت کنندگان فعال در اکوسیستم استخراج داده ها هستند (OECD، 2025) [21]. تقاضای روزافزون برای داده های وب اسکرپ شده، باعث ایجاد سازمان هایی مانند Common Crawl یا LAION، به نام تجمیع کنندگان داده های هوش مصنوعی، شده است که داده ها را جمع آوری کرده و با اشخاص ثالث به اشتراک می گذارند (OECD، 2025) [21] در واقع، تجمیع کننده های هوش مصنوعی نقش قابل توجهی در تسهیل داده های جمع آوری شده برای توسعه هوش مصنوعی ایفا می کنند؛ مطالعه ای که ۴۷ LLM منتشر شده بین سال های ۲۰۱۹ تا اکتبر ۲۰۲۳ را تجزیه و تحلیل کرد، نشان داد که حداقل ۶۴٪ از آنها با استفاده از داده های Common Crawl آموزش دیده اند (Baack، ۲۰۲۴) [58].

در رابطه با نوع داده های به دست آمده از طریق جمع آوری اطلاعات، لازم به ذکر است که محتوای جمع آوری شده ممکن است نه تنها شامل داده های شخصی باشد که مستقیماً از افراد گرفته شده یا در مورد آنها مشاهده شده است، بلکه شامل داده های شخصی نیز باشد که به طور عمومی به صورت آنلاین توسط اشخاص ثالث به اشتراک گذاشته می شود. به عنوان مثال، تصاویر، نام ها یا نظرات ممکن است توسط کاربران رسانه های اجتماعی در مورد افراد دیگر ارسال شود - مانند تگ کردن کسی در یک عکس یا ذکر نام او در یک نظر. در این موارد، داده های جمع آوری شده به افرادی اشاره دارد که خودشان محتوا را آپلود نکرده اند، که نشان می دهد مجموعه داده های جمع آوری شده ممکن است شامل داده های شخصی باشد که هم از شخص موضوع داده و هم از اشخاص ثالث سرچشمه می گیرد.

اگرچه داده ها ممکن است در یک وب سایت قابل دسترسی باشند، اما این به طور خودکار به این معنی نیست که واجد شرایط داده های باز هستند که می توانند آزادانه دوباره مورد استفاده قرار گیرند (مجمع جهانی حریم خصوصی، ۲۰۲۴) [59]. قانونی بودن جمع آوری داده ها برای آموزش مدل های هوش مصنوعی در حال حاضر در حوزه های قضایی مختلف، با تمرکز ویژه بر مالکیت معنوی، امنیت سایبری، حریم خصوصی و مسائل مربوط به مدیریت داده ها، مورد بحث است. در همین حال، شرکت های میزبانی وب به طور فزاینده ای در حال اجرای اقدامات فنی و قانونی، مانند پروتکل های حذف ربات (robots.txt) و ممنوعیت های صریح در شرایط استفاده وب سایت های خود، برای کاهش فعالیت های جمع آوری داده ها از وب هستند.

۵ نتیجه گیری

این مقاله، سازوکارهای کلیدی جمع آوری داده ها که توسط شرکت ها در توسعه مدل های هوش مصنوعی استفاده می شوند را برجسته می کند و بر منابعی که داده ها معمولاً از آنها به دست می آیند تمرکز دارد: (۱) مستقیماً از افراد و سازمان ها، و (۲) از ارائه دهندگان شخص ثالث. در این دو دسته کلی، سازوکارهای متمایزی را شناسایی می کنیم که نقش مهمی در ایجاد مجموعه داده های آموزشی هوش مصنوعی ایفا می کنند.

یکی از سازوکارهای مهم و رو به رشد برای آموزش هوش مصنوعی، استفاده از داده های تولید شده از تعاملات با ابزارهای هوش مصنوعی B2C، مانند چت بات ها، دستیاران مجازی و سیستم های خدمات مشتری خودکار است. داده های تعاملی همچنین از طریق ادغام های B2B جمع آوری می شوند، جایی که قابلیت های هوش مصنوعی از طریق مشارکت مستقیم یا دسترسی به API در محصولات تعبیه می شوند.

علاوه بر داده های جمع آوری شده در طول استفاده از ابزارهای هوش مصنوعی، برخی از توسعه دهندگان ممکن است به داده های سایر سرویس های دیجیتال خود نیز دسترسی داشته باشند و ممکن است به دنبال بهره برداری از این داده ها برای آموزش مدل هوش مصنوعی باشند.

صدور مجوز داده های تجاری یکی دیگر از رویه های رایج است که در آن توسعه دهندگان هوش مصنوعی از طریق توافق نامه های صدور مجوز با سازمان هایی که مجموعه داده های شخص ثالث را ارائه می دهند، به داده ها دسترسی پیدا می کنند. بازارهای داده و دلالان داده نیز منابع مهمی برای توسعه دهندگان هوش مصنوعی هستند که دسترسی به طیف گسترده ای از داده های شخص ثالث را ارائه می دهند.

داده های باز برای آموزش مدل های هوش مصنوعی به طور فزاینده ای اهمیت پیدا می کنند. مفهوم باز بودن داده ها در دو سطح عمل می کند: منابع داده اصلی و مجموعه داده های ایجاد شده از این منابع. این بدان معناست که در حالی که یک مجموعه داده آموزشی ممکن است شامل داده هایی با سطوح مختلف باز بودن باشد، خود مجموعه داده حاصل می تواند بسته و اختصاصی باقی بماند.

محبوبیت روزافزون جمع آوری داده ها به عنوان ابزاری برای جمع آوری داده ها، منجر به استفاده گسترده از آن برای آموزش مدل هوش مصنوعی شده است. با این حال، قانونی بودن جمع آوری داده ها در حال حاضر موضوع بحث در حوزه های قضایی مختلف است و نگرانی های اصلی در مورد مالکیت معنوی، امنیت سایبری، حریم خصوصی و مدیریت داده ها وجود دارد. به موازات آن، شرکت های میزبانی وب به طور فزاینده ای اقدامات فنی و قانونی را برای محدود کردن یا جلوگیری از فعالیت های جمع آوری داده ها اجرا می کنند.

اگرچه اهدای داوطلبانه داده ها هنوز منبع اصلی برای آموزش مدل هوش مصنوعی نیستند و ماهیت داوطلبانه واقعی آنها می تواند در برخی شرایط به چالش کشیده شود، اما در نهایت می توانند نقش مهمی در آموزش مدل هوش مصنوعی ایفا کنند و دسترسی به مجموعه داده های متنوع و واقعی را که در غیر این صورت دستیابی به آنها دشوار و پرهزینه است، فراهم کنند.

همچنین لازم به ذکر است که ابزارهای نوظهوری مانند PET ها نقش مهمی در فراهم کردن امکان جمع آوری داده ها برای توسعه هوش مصنوعی ایفا می کنند. PET ها با ایمن تر کردن جمع آوری و استفاده از داده ها، خطرات مربوط به حریم خصوصی و حاکمیتی را کاهش داده و از عملکرد مؤثر سازوکارهای ذکر شده در این طبقه بندی پشتیبانی می کنند.

این یافته ها، پیچیدگی و تنوع مکانیسم های جمع آوری داده ها که توسعه دهندگان هوش مصنوعی به آن ها متکی هستند، و همچنین چالش های قانونی در حال تکاملی که ممکن است با آن ها همراه باشد را برجسته می کند. چشم انداز دایماً در حال تغییر است و بررسی مداوم این مکانیسم ها برای شکل دهی سیاست هایی که توسعه هوش مصنوعی را با نیازهای حفظ حریم خصوصی و مدیریت داده ها متعادل می کنند، ضروری است.

منابع

- [46] مرکز هوش مصنوعی (۲۰۲۵). <https://aihub.or.kr/>.
- [54] باک، ب. (۲۰۲۵)، «به سوی بهترین شیوه ها برای مجموعه داده های باز برای آموزش LLM»، <https://arxiv.org/abs/2501.08365>, arXiv:2501.08365.
- [58] باک، س. (۲۰۲۴)، داده های آموزشی برای قیمت یک ساندویچ، <https://foundation.mozilla.org/es/research/library/generative-ai-training-data/common-crawl/> (دسترسی در ۱۹ آگوست ۲۰۲۵).
- [28] بدیلوم، م. (۲۰۲۴)، پردازش داده های شخصی برای آموزش هوش مصنوعی در برزیل: نکات کلیدی از تصمیمات اولیه ANPD در پرونده متا، <https://fpf.org/blog/processing-of-personaldata-for-ai-training-in-brazil-takeaways> (دسترسی در ۱۹ آگوست ۲۰۲۵).
- [17] بوماسانی، ک. (۲۰۲۴)، شاخص شفافیت مدل بنیاد نسخه ۱.۱.
- [43] بوسکت، آی. (۲۰۲۴)، این یک شرکت کشاورزی قدیمی است—و جدیدترین فروشنده هوش مصنوعی شماست، <https://www.wsj.com/articles/bayer-microsoft-generative-ai-90754f54> (دسترسی در ۱۹ آگوست ۲۰۲۵).
- [8] مجلس سنای فدرال برزیل (۲۰۲۴)، چارچوب قانونی هوش مصنوعی برزیل.
- [50] خزیدن مشترک (دوم)، <https://commoncrawl.org/>.
- [47] اشیاء رایج در متن (۲۰۲۵)، <https://cocodataset.org/>.
- [18] سازمان رقابت و بازارها (۲۰۲۳)، مدل های بنیاد هوش مصنوعی: گزارش اولیه.
- [15] دنگ، جی. (۲۰۰۹)، «ImageNet: یک پایگاه داده تصویر سلسله مراتبی در مقیاس بزرگ. در سال ۲۰۰۹ IEEE کنفرانس بینایی کامپیوتر و تشخیص الگوکنفرانس IEEE در زمینه بینایی کامپیوتر و تشخیص الگو، صفحات ۲۴۸-۲۵۵.
- [6] کمیسیون اروپا (۲۰۲۴)، پیش نویس دوم هوش مصنوعی همه منظوره، <https://digital-strategy.ec.europa.eu/en/library/second-draft-general-purpose-ai-code-practice-published-written-independent-experts>.
- [5] اتحادیه اروپا (۲۰۲۴)، آیین نامه (EU) ۲۰۲۴/۱۶۸۹ پارلمان اروپا و شورای ۱۳ ژوئن ۲۰۲۴، وضع قوانین هماهنگ در مورد هوش مصنوعی و اصلاح مقررات (EC) شماره ۲۰۰۸/۳۰۰، (EU) شماره ۲۰۱۳/۱۶۷، (EU) شماره ۲۰۱۳/۱۶۸، (EU) ۲۰۱۸/۱۱۳۹، (EU) ۲۰۱۸/۱۵۸۰ و.
- [29] اتحادیه اروپا (۲۰۲۲)، آیین نامه (EU) ۲۰۲۲/۸۶۸ پارلمان اروپا و شورای ۳۰ مه ۲۰۲۲ در مورد حاکمیت داده اروپا و اصلاح مقررات (EU) ۲۰۱۸/۱۷۲۴ (قانون حاکمیت داده)، <https://eur-lex.europa.eu/legalcontent/EN/TXT/PDF/?uri=CELEX:32022R0868>.

- [44] اتحادیه اروپا (۲۰۱۹)، دستورالعمل (EU) 2019/1024 پارلمان اروپا و شورای 20 ژوئن 2019 در مورد داده های باز و استفاده مجدد از اطلاعات بخش عمومی (بازنگری)، [/oj/eng](https://eur-lex.europa.eu/eli/dir/2019/1024/oj/eng)، (دسترسى در ۱۹ آگوست ۲۰۲۵).
- [13] فدرکوپر، ل. (۲۰۲۳)، «گزارش اولین کارگاه آموزشی هوش مصنوعی مولد و قانون»، حقوق و قانون ییل، مقاله پژوهشی اقتصاد، https://papers.ssrn.com/sol3/papers.cfm?abstract_id=4634513.
- [25] کمیسیون تجارت فدرال (2024)، شرکت های هوش مصنوعی (و سایر): تغییر بی سروصدای شرایط خدمات شما می تواند ... ناعادلانه یا فریبنده، <https://www.ftc.gov/policy/advocacy-research/tech-at-ftc/2024/02/aiother-companies-quietly-changing-your-terms-service-could-be-unfair-or-deceptive>.
- [10] گروه ۷ (۲۰۲۳)، آیین نامه بین المللی رفتار سازمانی فرآیند هیروشیما.
- [14] گراخ، ف. (۲۰۱۸)، «یک پیکره استاندارد پروژه گوتنبرگ برای تحلیل آماری داده های طبیعی» زبان و زبان شناسی کمی، <https://arxiv.org/abs/1812.08092v1>، *arXiv:1812.08092*.
- [59] مجمع جهانی حریم خصوصی (۲۰۲۴)، بیانیه مشترک پایانی در مورد جمع آوری داده ها و حفاظت از آنها از حریم خصوصی، https://www.priv.gc.ca/en/opc-news/speeches-and-statements/2024/jsdc_20241028/، (دسترسى در ۱۵ مه ۲۰۲۵).
- [26] گلدنبرگ، ل. (۲۰۲۱)، «مجموعه داده های سفرهای چند مقصدی Booking.com»، مجموعه مقالات '21 SIGIR چهل و چهارمین کنفرانس بین المللی ACM SIGIR در زمینه تحقیق و توسعه در بازیابی اطلاعات، صفحات ۲۴۵۷-۲۴۶۲، <https://doi.org/10.1145/3404835.3463240>.
- [62] گوگل (دوم)، .
- [57] گوگل (دوم)، جستجوی گوگل مرکزی، <https://developers.google.com/search/docs/crawling-indexing/robots-user-agent>، (دسترسى در ۱۹ آگوست ۲۰۲۵).
- [31] هیرشام سی، آر. (۲۰۲۰)، «بیماری های نادر ۲۰۳۰: چگونه هوش مصنوعی تقویت شده از تشخیص و درمان بیماری های نادر در آینده»، صفحات ۷۴۰-۷۴۳، <https://doi.org/10.1136/annrheumdis-2020-217125>.
- [51] چهره در آغوش گرفته (دوم)، .
- [36] (ICO) 2025، قانون استفاده و دسترسى به داده ها ۲۰۲۵ (DUA) - برای سازمان ها چه معنایی دارد؟، <https://ico.org.uk/about-the-ico/what-we-do/legislation-we-cover/data-use-and-access-act-2025/the-data-use-and-access-act-2025-what-does-it-mean-for-organisations/#innovate> (دسترسى در ۲۶ ژوئن ۲۰۲۵).
- [16] ایمیج نت (دوم)، <http://www.image-net.org/about>.
- [48] جی. دنگ، دبلیو. (۲۰۰۹)، «ImageNet: یک پایگاه داده تصویر سلسله مراتبی در مقیاس بزرگ»، کنفرانس IEEE بینایی کامپیوتر و تشخیص الگو، فلوریدا، ایالات متحده آمریکا، صفحات ۲۴۸-۲۵۵، <https://doi.org/10.1109/CVPR.2009.5206848>.
- [9] رژیم ملی ژاپن (۲۰۲۵)، قانون ترویج تحقیق و توسعه و بهره برداری فناوری های مرتبط با هوش مصنوعی.
- [30] کیستین، ل. (۲۰۲۳)، داده های تحرک به عنوان یک دارایی مشترک - به سوی یک زیرساخت داده های تحرک مشترک، https://fsr.eui.eu/mobility-data-as-a-commons-towards-a-common-mobility-datainfrastructure/#_ftnref6، (دسترسى در ۱۹ آگوست ۲۰۲۵).
- [7] مجلس ملی کره (۲۰۲۴)، قانون توسعه هوش مصنوعی و ایجاد اعتماد.

- [12] لی، سی. (۲۰۲۳)، «هوش مصنوعی و قانون: نسل بعدی»، <https://blog.genlaw.org/explainers/>.
- [27] لوین (۲۰۲۴)، چگونه آمازون از هوش مصنوعی مولد برای بهبود توصیه های محصول استفاده می کند و توضیحات، <https://www.aboutamazon.com/news/retail/amazon-generative-ai-productsearch-results-and-descriptions> (دسترسی در ۱۹ آگوست ۲۰۲۵).
- [60] لیبرمن، ب. (۲۰۲۰)، ردیابی انتشار کربن جهانی در لحظه با استفاده از *Climate TRACE*، <https://yaleclimateconnections.org/2020/08/climate-trace-to-track-real-time-global-carbonemissions/>.
- [49] لیزنفلد، آ. (2023)، «باز کردن ChatGPT: ردیابی باز بودن، شفافیت و پاسخگویی در مولدهای متن تنظیم شده با دستورالعمل»، مجموعه مقالات پنجمین کنفرانس بین المللی رابط های کاربری محاوره ای، *CUI 2023*، <https://dl.acm.org/doi/10.1145/3571884.3604316>.
- [66] لورنز، پ. (۲۰۲۳)، ملاحظات اولیه سیاست گذاری برای هوش مصنوعی مولد، <https://doi.org/10.1787/fae2d1e6-en>.
- [38] اُریان، م. (۲۰۲۳)، شرکت *OpenAI*، سازنده ی *ChatGPT*، برای صدور مجوز انتشار اخبار با آسوشیتدپرس قرارداد امضا کرد...، <https://apnews.com/article/openai-chatgpt-associated-press-apf86f84c5bcc2f3b98074b38521f5f75a> (دسترسی در ۱۹ آگوست ۲۰۲۵).
- [53] سازمان همکاری و توسعه اقتصادی (2025)، افزایش دسترسی و اشتراک گذاری داده ها در عصر هوش مصنوعی، انتشارات OECD، پاریس.
- [35] سازمان همکاری و توسعه اقتصادی (2025)، تسهیل استفاده ثانویه از داده های سلامت برای اهداف عمومی در سراسر مرزها انتشارات OECD، پاریس، <https://doi.org/10.1787/d7b90d15-en>.
- [21] سازمان همکاری و توسعه اقتصادی (2025)، مسائلی مربوط به مالکیت معنوی در هوش مصنوعی آموزش دیده بر اساس داده های خراشیده شده، سازمان همکاری و توسعه اقتصادی انتشارات، پاریس، <https://doi.org/10.1787/d5241a23-en>.
- [11] سازمان همکاری و توسعه اقتصادی (2025)، راه اندازی چارچوب گزارش دهی فرآیند هوش مصنوعی هیروشیما (*HAIP*)، <https://www.oecd.org/en/events/2025/02/launch-of-the-hiroshima-ai-process-reportingframework.html> (دسترسی در ۷ فوریه).
- [24] سازمان همکاری و توسعه اقتصادی (2025)، اشتراک گذاری مدل های هوش مصنوعی قابل اعتماد با فناوری های تقویت کننده حریم خصوصی، سازمان همکاری و توسعه اقتصادی انتشارات، پاریس، <https://doi.org/10.1787/a266160b-en>.
- [20] سازمان همکاری و توسعه اقتصادی (2024)، یادداشت توضیحی در مورد تعریف به روز شده OECD از سیستم هوش مصنوعی، انتشارات OECD، <https://doi.org/10.1787/623da898-en>.
- [65] سازمان همکاری و توسعه اقتصادی (2023)، مدل های زبانی هوش مصنوعی: ملاحظات فناوری، اجتماعی-اقتصادی و سیاست گذاری، انتشارات OECD، پاریس، <https://doi.org/10.1787/13d38f92-en>.
- [23] سازمان همکاری و توسعه اقتصادی (2022)، راهنمای سیاست گذاری در حوزه حاکمیت داده انتشارات OECD، پاریس، <https://doi.org/10.1787/40d53904-en>.
- [1] «برای طبقه بندی سیستم های هوش مصنوعی OECD چارچوب»، (2022) OECD اقتصاد دیجیتال *OECD* مقالات، شماره ۳۲۳، انتشارات OECD، پاریس، <https://doi.org/10.1787/cb6d9eca-en>.
- [64] سازمان همکاری و توسعه اقتصادی (2021)، ترسیم ابتکارات، فرصت ها و چالش ها در زمینه قابلیت انتقال داده ها، سازمان همکاری و توسعه اقتصادی انتشارات، پاریس، <https://doi.org/10.1787/a6edfab2-en>.

- [32] سازمان همکاری و توسعه اقتصادی (2021)، توصیه شورا در مورد افزایش دسترسی و اشتراک گذاری داده ها، <https://legalinstruments.oecd.org/fr/instruments/OECD-LEGAL-0463> (دسترسی در ۱۹ اگوست ۲۰۲۵).
- [4] سازمان همکاری و توسعه اقتصادی (۲۰۱۹)، افزایش دسترسی و اشتراک گذاری داده ها: تطبیق خطرات و مزایا برای استفاده مجدد از داده ها در جوامع مختلف انتشارات OECD، پاریس، <https://doi.org/10.1787/b4d546a9-en>.
- [22] سازمان همکاری و توسعه اقتصادی (۲۰۱۹)، افزایش دسترسی و اشتراک گذاری داده ها: تطبیق خطرات و مزایا برای استفاده مجدد از داده ها در جوامع مختلف انتشارات OECD، پاریس، <https://doi.org/10.1787/b4d546a9-en>.
- [2] سازمان همکاری و توسعه اقتصادی (۲۰۱۹)، توصیه شورای هوش مصنوعی، <https://legalinstruments.oecd.org/en/instruments/OECD-LEGAL-0449> (دسترسی در ۱۹ اگوست ۲۰۲۵).
- [52] سازمان همکاری و توسعه اقتصادی (۲۰۱۵)، نوآوری مبتنی بر داده: کلان داده برای رشد و رفاه انتشارات OECD، پاریس، <https://doi.org/10.1787/9789264229358-en>.
- [3] سازمان همکاری و توسعه اقتصادی (۲۰۱۴)، حفاظت از حریم خصوصی در اقتصاد داده محور: بررسی تفکرات فعلی، <https://one.oecd.org/document/DSTI/ICCP/REG/2014/3/en/pdf>.
- [40] اوپن ای آی (۲۰۲۴)، همکاری OpenAI و Reddit، <https://openai.com/index/openai-and-reddit-> مشارکت/ (دسترسی در ۱۹ اگوست ۲۰۲۵).
- [39] اوپن ای آی (۲۰۲۳)، اکسل اسپرینگر استفاده مفید از هوش مصنوعی در روزنامه نگاری را تعمیق می بخشد، <https://openai.com/index/axel-springer-partnership/> (دسترسی در ۱۹ اگوست ۲۰۲۵).
- [55] اوپن ای آی (۲۰۲۳)، گزارش فنی GPT-4، <https://arxiv.org/pdf/2303.08774>.
- [19] پابلویلاوبوس، آ. (2022)، «آیا داده ها تمام خواهند شد؟ محدودیت های مقیاس پذیری LLM بر اساس داده های انسانی» داده های تولید شده، <https://arxiv.org/abs/2211.04325>.
- [42] (2024)، PR Newswire نکسیس از نکسیس دیتا پلاس، یک پلتفرم تک API برای کمک های مالی رونمایی کرد دسترسی بی سابقه سازمان ها به محتوای خبری دارای مجوز و داده های باکیفیت شرکتی که توسط هوش مصنوعی عمومی تأیید شده اند، <https://www.prnewswire.com/news-releases/lexisnexis-unveils-nexisdata-a-single-api-platform-giving-organizations-unprecedented-access-to-gen-ai-approvedlicensed-news-content-and-high-quality-company-data-302325319.html>.
- [37] ساگارا، اچ. (۲۰۱۹)، گزارش نهایی در مورد طرح های آزمایشی بارسلونا، ارزیابی های BarcelonaNow و طرح های پایداری.
- [41] شاتراستوک (۲۰۲۳)، شاتر استوک رابطه طولانی مدت خود با متا را گسترش می دهد، <https://www.prnewswire.com/news-releases/shutterstock-expands-long-standing-relationship-with-meta-301719769.html> (دسترسی در ۱۹ اگوست ۲۰۲۵).
- [45] ملت هوشمند سنگاپور (دوم)،.
- [33] بانک زیستی بریتانیا (دوم)، <https://www.ukbiobank.ac.uk/learn-more-about-uk-biobank>.
- [61] وزارت بازرگانی ایالات متحده (۲۰۲۴)، آماده سازی داده های باز برای عصر هوش مصنوعی، <https://www.commerce.gov/news/blog/2024/01/preparing-open-data-age-ai>.
- [63] مجلس نمایندگان ایالات متحده (۲۰۲۴)، گزارش کارگروه دو حزبی مجلس نمایندگان در مورد هوش مصنوعی: راهنما اصول، توصیه های آینده نگر و پیشنهادهای سیاستی برای اطمینان از اینکه آمریکا همچنان رهبری جهان را در نوآوری مسیولانه هوش مصنوعی بر عهده دارد..

[34] وایزبند، س. (۲۰۲۳)، مطالعه *GIHF-AI 2023*، اعتماد به استفاده از داده های سلامت - یک مقایسه روابط آلمان و اسرائیل.

[56] ویتاکر، جی. (۲۰۲۴)، رایان ایر در پرونده ی «اسکرین اسکرین» *Booking.com* پیروز شد، [/industry-news/ryanair-wins-booking-com-screenscraper-case/](https://www.aviationbusinessnews.com/industry-news/ryanair-wins-booking-com-screenscraper-case/) (دسترسی در ۱۹ آگوست ۲۰۲۵).

یادداشت‌ها

۱) OECD، 2023) ها الگوریتم‌های پیشرفته یادگیری ماشین هستند که از طریق فرآیندی که اطلاعات را از داده‌های ورودی به پارامترهای مدل خود کدگذاری می‌کند، توسعه یافته و اصلاح می‌شوند. این امر به آنها امکان یادگیری و انجام وظایفی مانند ترجمه، تولید محتوا یا تشخیص تصویر را می‌دهد (LLM [65] (لورنز، 2023 [66]).

۲) دامنه عمومی شامل تمام داده‌هایی می‌شود که اشتراک‌گذاری آنها توسط حقوق مالکیت فکری یا هرگونه حقوق دیگری با اثرات مشابه محدود نشده است. بنابراین، در اینجا به طور گسترده‌تری از صرفاً عاری بودن از حمایت حق چاپ درک می‌شود (OECD، 2019 [14]). حوزه عمومی می‌تواند با حوزه‌های شخصی و اختصاصی همپوشانی داشته باشد (شکل ۲ را ببینید).

۳) شاخص شفافیت مدل بنیادی ۲۰۲۴ شامل ۲۳ زیردامنه است که جنبه‌های مختلف مدل‌های بنیادی را از توسعه تا استقرار ارزیابی می‌کند: داده‌ها، تأثیر نیروی کار، دسترسی به داده‌ها، محاسبات، کد، تأثیر زیست محیطی، ویژگی‌های مدل، قابلیت‌ها، محدودیت‌ها، خطرات، قابلیت اعتماد، کاهش اثرات مدل، ارزیابی‌ها، به روزرسانی‌های مدل، فرایند انتشار، کانال‌های توزیع، رابط کاربری، حفاظت از داده‌ها، مکانیسم‌های بازخورد، مستندسازی، سیاست‌های استفاده، جغرافیای تحت تأثیر و تأثیر (یوماسانی، ۲۰۲۴ [17]).

۴) در برخی موارد، مانند مدل‌های استدلال، تنظیم دقیق ممکن است شامل یادگیری تقویتی از بازخورد انسانی باشد، نه اینکه صرفاً بر اساس مجموعه داده‌های خاص دامنه آموزش داده شود.

دیه یادداشت ۲ مراجعه کنید.

۵) اگرچه این طبقه بندی بر منابع جمع‌آوری داده‌ها تمرکز دارد، اما توجه به این نکته ضروری است که روش‌های دسترسی به داده‌های می‌تواند در این مکانیسم‌ها به طور قابل توجهی متفاوت باشد. در برخی موارد، داده‌ها به طور کامل به دست می‌آیند - به عنوان مثال، از طریق قراردادهای مجوز تجاری، که در آن دارنده مجوز ممکن است مالکیت یا حقوق استفاده مجدد گسترده‌ای را به دست آورد. در موارد دیگر، داده‌ها فقط تحت شرایط تعریف شده‌ای مانند محیط‌های پردازش امن قابل دسترسی هستند، جایی که می‌توان آن‌ها را تجزیه و تحلیل کرد اما نمی‌توان آن‌ها را دانلود یا نگهداری کرد.

۷) طبق قانون حاکمیت داده، نهادها می‌توانند به عنوان «سازمان‌های نوع دوستی داده» ثبت شوند که به اشتراک‌گذاری داده‌ها را برای منافع عمومی مدیریت و تسهیل می‌کنند و تضمین می‌کنند که داده‌ها به صورت مسئولانه و مطابق با مقررات حفظ حریم خصوصی استفاده می‌شوند.

۸) ماهیت داوطلبانه‌ی اهدای داده‌ها ممکن است تحت تأثیر عواملی مانند سواد دیجیتال، آگاهی از حقوق داده‌ها و درک خطرات مرتبط باشد. الگوهای مشارکت ممکن است در گروه‌های جمعیتی مختلف متفاوت باشد و زمینه‌ی درخواست‌های اهدای داده‌های می‌تواند بر ماهیت واقعاً داوطلبانه‌ی رضایت تأثیر بگذارد.

۹) «ترتیبات داده‌های باز» به ترتیبات دسترسی و اشتراک‌گذاری غیرتبعیض‌آمیز داده‌ها اشاره دارد، که در آن داده‌ها قابل خواندن توسط ماشین هستند و می‌توانند به صورت رایگان به آنها دسترسی پیدا کرده و به اشتراک گذاشته شوند و توسط هر کسی برای هر هدفی مورد استفاده قرار گیرند، حداکثر با رعایت الزاماتی که یکپارچگی، منشأ، انتساب و باز بودن را حفظ می‌کنند (OECD، 2021 [32]).