

نگرانی دانشمندان از تعصب «چت جی پی تی»

مطالعات تأیید می‌کنند که هوش مصنوعی چت جی پی تی که محبوب‌ترین مدل زبانی هوش مصنوعی در جهان است در پاسخ‌های خود دارای تعصب است.



مطالعات تأیید می‌کنند که هوش مصنوعی چت جی پی تی که محبوب‌ترین مدل زبانی هوش مصنوعی در جهان است در پاسخ‌های خود دارای تعصب است.

به گزارش ایسنا و به نقل از آی‌ای، براساس مطالعه محققان دانشگاه کالیفرنیا، برکلی، چت جی پی تی متعلق به شرکت اپن‌ای‌آی (OpenAI) که در حال حاضر محبوب‌ترین مدل زبان مبتنی بر هوش مصنوعی در جهان به حساب می‌آید، در اطلاعاتی که ارائه می‌کند بی‌طرف نیست. مقاله آنها نشان می‌دهد که از آنجایی که چت جی پی تی از حجم عظیمی از مقالات دارای حق نشر به عنوان مجموعه داده‌ای که بر اساس آن آموزش یافته و پاسخ‌های خود را پایه‌گذاری می‌کند، استفاده می‌کند، سوگیری ذاتی در آن داده‌ها به نتایج چت جی پی تی نیز راه پیدا می‌کند.

محققان در مطالعه خود نوشتند که مدل‌های اپن‌ای‌آی تعداد زیادی از انواع مختلف مقالات دارای حق نشر را در خود جای داده‌اند و این فرآیند تحت تأثیر تعداد دفعاتی است که قسمت‌های آن مقالات در فضای وب ظاهر می‌شوند. به این ترتیب، اطلاعات مرتبط با مقالاتی که در فرآیند آموزش بیشتر مورد استفاده قرار گرفته‌اند، نسبت به اطلاعاتی که به طور مکرر مورد استفاده قرار نگرفته‌اند، بیشتر در پاسخ‌ها گنجانده می‌شوند.

همانطور که محققان در مورد این مطالعه توضیح دادند، دقت چنین مدل‌هایی به شدت به فرکانس مشاهده اطلاعات در داده‌های آموزشی وابسته است و این، توانایی آنها در تعمیم را زیر سؤال می‌برد.

یکی از نمونه‌های آن این است که کتاب‌های علمی تخیلی و فانتزی ظاهراً در فهرست کتاب‌هایی که برای آموزش استفاده می‌شوند بیش از حد استفاده شده‌اند. این یک «سوگیری درونی» ایجاد می‌کند تا آنجا که می‌توان گفت از چت جی پی تی چه نوع پاسخ‌هایی را می‌توان انتظار داشت.

این مطالعه اشاره می‌کند که کتاب‌هایی که در فهرست داده‌های آموزشی چت جی پی تی به طور غالب حضور داشته‌اند، عناوین محبوبی مانند هری پاتر، ۱۹۸۴، ارباب حلقه‌ها، بازی‌های گرسنگی، راهنمای مسافران مجانی کهکشانشان، فارنهایت ۴۵۱، بازی تاج و تخت، و تپه‌های شنی را شامل می‌شوند.

آنچه در داده‌های آموزشی پیدا نمی‌کنید

در حالی که چت جی پی تی موارد زیادی در مورد مقالات موجود در حوزه عمومی می‌داند، آنچه در داده‌های آموزشی نمی‌بینید، عناوین کمتر شناخته شده‌ای مانند به اصطلاح ادبیات جهانی آنگلوفون هستند. این‌ها کتاب‌هایی هستند که برای و توسط مخاطبانی به جز کشورهای اصلی انگلیسی زبان نوشته شده‌اند. چنین مناطقی شامل آفریقا، آسیا و کارائیب می‌شود.

برای اینکه نتایج تولید شده توسط چت جی پی تی معتبرتر باشد، محققان از مدل‌های هوش مصنوعی می‌خواهند که شرکت اپن‌ای‌آی در مورد داده‌های مورد استفاده در آموزش مدل‌هایش شفاف‌تر باشد. به گفته آنها دانستن اینکه یک مدل از روی چه کتاب‌هایی آموزش دیده است، برای ارزیابی چنین منابع سوگیرانه‌ای بسیار مهم است.

یکی از محققان دانشگاه برکلی در توییتری به این سوگیری بالقوه پرداخت و نوشت که با توجه به ترجیح آنها نسبت به منابع علمی تخیلی و فانتزی، برنامه‌های محبوب هوش مصنوعی احتمالاً معیارهای اندازه‌گیری کاملی برای عملکرد مدل‌ها نیستند و ما باید به این فکر کنیم که تجربیات روایت چه کسی در این مدل‌ها کدگذاری شده است و چگونه بر سایر رفتارها تأثیر می‌گذارد. استفاده از مطالب دارای حق نشر در آموزش مدل‌ها نیز انبوهی از سوالات حقوقی را ایجاد می‌کند. چه کسی دارای حق نشر متنی را که توسط چت جی پی تی ایجاد شده دارد. چت جی پی تی که خود پیش از این از روی کتاب‌های دارای حق نشر آموزش داده شده است؟ آیا مفهوم «استفاده منصفانه» در چنین موردی می‌تواند کاربرد داشته باشد؟

اگر تعدادی از افراد سعی کنند برای خروجی‌های مشابه یا یکسان توسط مدلی مانند چت جی پی تی حق نشر بگیرند، چه خواهد شد؟ از آنجایی که ماشین، انسان نیست، آیا هر چیزی که توسط آن تولید می‌شود در وهله اول دارای حق نشر است؟

سوگیری که توسط مطالعات دیگر آشکار شد

علاوه بر تحقیقات برکلی، مطالعات دیگری نیز مواردی از سوگیری را در پاسخ‌های چت جی پی تی پیدا کرده‌اند. مطالعه‌ای که توسط گروهی در دانشگاه انگلیا شرقی در بریتانیا انجام شد، سوگیری سیاسی را در برخی از پاسخ‌های این مدل نشان داد. هنگامی که صدها سؤال در مورد اعتقادات سیاسی از این هوش مصنوعی پرسیده شد، به نظر می‌رسید چت جی پی تی بیشتر به دموکرات‌ها در ایالات متحده، حزب کارگر در بریتانیا و رئیس‌جمهور لولا داسیلوا (Lula da Silva) از حزب کارگران در برزیل متمایل باشد.

نویسنده اصلی آن مطالعه، دکتر فابیو موتوکی (Fabio Motoki)، از دانشکده تجارت نورویچ در دانشگاه انگلیا شرقی، می‌گوید: وجود سوگیری سیاسی می‌تواند بر دیدگاه‌های کاربران نیز تأثیر بگذارد و پیامدهای بالقوه‌ای برای فرآیندهای سیاسی و

انتخاباتی داشته باشد.

یافته های ما این نگرانی را تقویت می کند که سیستم های هوش مصنوعی می توانند چالش های موجود ناشی از اینترنت و رسانه های اجتماعی را تکرار یا حتی تقویت کنند.

در مطالعه دیگری، محققان دانشگاه واشنگتن، دانشگاه کارنگی ملون و دانشگاه شیان جیائوتنگ، ۱۴ مدل هوش مصنوعی را تحت یک آزمون سوگیری سیاسی قرار دادند و به پاسخ های هر مدل به انواع مختلف ۶۲ بیانیه سیاسی پرداختند. چیزی که آن ها دریافتند این بود که پاسخ های تولید شده توسط چت جی پی تی و نسخه جدیدتر آن موسوم به جی پی تی-۴ (GPT-4) چپ گرایانه و آزادیخواهانه بودند.

شرکت اِپن ای آی رویکرد خود را در یک پست شرکتی با عنوان چگونه سیستم ها باید رفتار کنند به اشتراک گذاشته است. براساس این پست، در حالی که بسیاری نگران سوگیری ها در طراحی سیستم های هوش مصنوعی هستند، ما متعهد هستیم که به طور جدی به این موضوع رسیدگی کنیم و در مورد اهداف و پیشرفت خود شفاف باشیم. این شرکت همچنین خاطرنشان می کند که دستورالعمل های ما صریح است که بازبینان نباید از هیچ گروه سیاسی حمایت کنند. آن ها اظهار داشتند که ممکن است سوگیری ها همچنان در این فرآیند ظاهر شوند، اما ادعا کردند که این ها اشکالات به حساب می آیند و جزو ویژگی های هوش مصنوعی آن ها نیستند.