



## داده‌های جهان برای تغذیه هوش مصنوعی در حال اتمام است

کارشناسان هشدار می‌دهند که داده‌های جهان برای تغذیه هوش مصنوعی در حال اتمام است.

کارشناسان هشدار می‌دهند که داده‌های جهان برای تغذیه هوش مصنوعی در حال اتمام است. به گزارش ایسنا و به نقل از اس‌ای، در حالی که هوش مصنوعی این روزها به اوج محبوبیت خود می‌رسد، پژوهشگران هشدار داده‌اند که این صنعت ممکن است با کمبود داده‌های آموزشی روبرو شود.

داده‌ها به مانند سوختی هستند که سیستم‌های هوش مصنوعی قدرتمند را نیرو می‌دهند و این مشکل می‌تواند رشد مدل‌های هوش مصنوعی، به ویژه مدل‌های زبانی بزرگ را کاهش دهد و حتی ممکن است مسیر انقلاب هوش مصنوعی را تغییر دهد.

اما چرا کمبود بالقوه داده با توجه به میزان موجود در شبکه جهانی وب یک مشکل است و اینکه آیا راهی برای مقابله با این خطر وجود دارد؟

### اهمیت داده‌های با کیفیت بالا برای هوش مصنوعی

ما برای آموزش الگوریتم‌های هوش مصنوعی قدرتمند، دقیق و با کیفیت به داده‌های زیادی نیاز داریم. به عنوان مثال، چت بات ChatGPT بر روی ۵۷۰ گیگابایت داده متنی یا حدود ۳۰۰ میلیارد کلمه آموزش داده شده است.

به طور مشابه، الگوریتم انتشار پایدار (که پشت بسیاری از برنامه‌های تولید تصویر هوش مصنوعی مانند DALL-E و Midjourney قرار دارد) بر روی مجموعه داده LIAON-5B متشکل از ۵.۸ میلیارد جفت تصویر-متن آموزش داده شده است. اگر الگوریتمی روی مقدار ناکافی داده آموزش داده شود، خروجی‌های نادرست یا با کیفیت پایین تولید می‌کند.

کیفیت داده‌های آموزشی نیز مهم است. دسترسی به داده‌های با کیفیت پایین مانند پست‌های رسانه‌های اجتماعی یا عکس‌های تار آسان است، اما برای آموزش مدل‌های هوش مصنوعی با کارایی بالا کافی نیستند.

متنی که از سکوها رسانه‌های اجتماعی گرفته می‌شود ممکن است مغرضانه یا تعصب آمیز باشد یا ممکن است حاوی اطلاعات نادرست یا محتوای غیرقانونی باشد که می‌تواند توسط مدل هوش مصنوعی تکرار شود. به عنوان مثال، زمانی که مایکروسافت سعی کرد ربات هوش مصنوعی خود را با استفاده از محتوای توییتر آموزش دهد، دریافت که خروجی‌های نژادپرستانه و زن ستیز تولید می‌کند.

به همین دلیل است که توسعه دهندگان هوش مصنوعی به دنبال محتوای باکیفیت مانند متن از کتاب‌ها، مقالات آنلاین، مقالات علمی، ویکی‌پدیا و محتوای وب فیلتر شده خاص هستند. مثلاً Google Assistant بر روی ۱۱ هزار رمان عاشقانه برگرفته از سایت خودانتشار Smashwords آموزش دیده است تا اطلاعات بیشتری در چنته داشته باشد.

### آیا ما داده‌های کافی داریم؟

صنعت هوش مصنوعی، سیستم‌های هوش مصنوعی را بر روی مجموعه داده‌های بزرگ‌تر آموزش می‌دهد، به همین دلیل است که ما اکنون مدل‌هایی با عملکرد بالا مانند ChatGPT یا DALL-E ۲ داریم. در عین حال، تحقیقات نشان می‌دهد که ذخیره داده‌های آنلاین بسیار کندتر از مجموعه داده‌های مورد استفاده در حال رشد برای آموزش هوش مصنوعی هستند.

در مقاله‌ای که سال گذشته منتشر شد، گروهی از پژوهشگران پیش‌بینی کردند که اگر روند آموزشی فعلی هوش مصنوعی ادامه یابد، قبل از سال ۲۰۲۶ داده‌های متنی با کیفیت بالا تمام خواهد شد.

آنها همچنین تخمین زدند که داده‌های زبانی با کیفیت پایین بین سال‌های ۲۰۳۰ تا ۲۰۵۰ و داده‌های تصویری با کیفیت پایین بین سال‌های ۲۰۳۰ تا ۲۰۶۰ به پایان می‌رسد.

به گفته گروه مشاوره و حسابداری PWC، هوش مصنوعی تا سال ۲۰۳۰ می‌تواند تا ۱۵.۷ تریلیون دلار به اقتصاد جهان کمک کند. اما تمام شدن داده‌های قابل استفاده می‌تواند توسعه آن را آهسته کند.

### آیا باید نگران باشیم؟

در حالی که نکات بالا ممکن است برخی از طرفداران هوش مصنوعی را نگران کند، این وضعیت ممکن است آنقدرها هم که به نظر می‌رسد بد نباشد. ناشناخته‌های زیادی در مورد چگونگی توسعه مدل‌های هوش مصنوعی در آینده و همچنین چند راه برای مقابله با خطر کمبود داده وجود دارد.

یکی از فرصت‌ها برای توسعه دهندگان هوش مصنوعی است که الگوریتم‌ها را بهبود ببخشند تا از داده‌هایی که در حال حاضر در اختیار دارند به طور کارآمدتر استفاده کنند.

این احتمال وجود دارد که در سال‌های آینده آنها بتوانند سیستم‌های هوش مصنوعی با عملکرد بالا را با استفاده از داده‌های کمتر و احتمالاً قدرت محاسباتی کمتر آموزش دهند. این همچنین به کاهش ردپای کربن هوش مصنوعی کمک می‌کند.

گزینه دیگر استفاده از هوش مصنوعی برای ایجاد داده‌های مصنوعی برای آموزش سیستم‌هاست. به عبارت دیگر، توسعه دهندگان به سادگی می‌توانند داده‌های مورد نیاز خود را تولید کنند که متناسب با مدل هوش مصنوعی خاص آنها باشد.

چندین پروژه در حال حاضر از محتوای مصنوعی استفاده می کنند که اغلب از سرویس های تولید داده مانند Mostly AI تهیه می شود. این امر در آینده رایج تر خواهد شد.

توسعه دهندگان همچنین به دنبال محتوایی خارج از فضای آنلاین رایگان هستند، مانند محتوایی که توسط ناشران بزرگ و مخازن آفلاین نگهداری می شود. به میلیون ها متن منتشر شده قبل از دوران فراگیری اینترنت فکر کنید. آنهایی که به صورت دیجیتالی درنیامده اند و در دسترس قرار بگیرند، می توانند منبع جدیدی از داده ها را برای پروژه های هوش مصنوعی فراهم کنند.

بنیاد News Corp یکی از بزرگترین دارندگان محتوای خبری در جهان اخیراً اعلام کرده است که در حال مذاکره با توسعه دهندگان هوش مصنوعی برای قراردادهای محتواست. چنین معاملاتی شرکت های هوش مصنوعی را مجبور می کند برای داده های آموزشی پول بپردازند، در حالی که تاکنون بیشتر آنها را به صورت رایگان از اینترنت حذف کرده اند.

سازندگان محتوا نسبت به استفاده غیرمجاز از محتوای خود برای آموزش مدل های هوش مصنوعی اعتراض کرده اند و برخی از شرکت هایی مانند مایکروسافت، OpenAI و Stability AI شکایت کرده اند. دریافت پاداش برای کار آنها ممکن است به بازگرداندن برخی از عدم تعادل قدرت بین خلاقان و شرکت های هوش مصنوعی کمک کند.