

فناوری پشت ChatGPT، (مدل زبان بزرگ) چیست؟

تاریخ خبر: ۲۰۲۳/۰۳/۱۵

مبانی داده‌ها، موارد استفاده و پروژه‌ها Kurt Muehmel

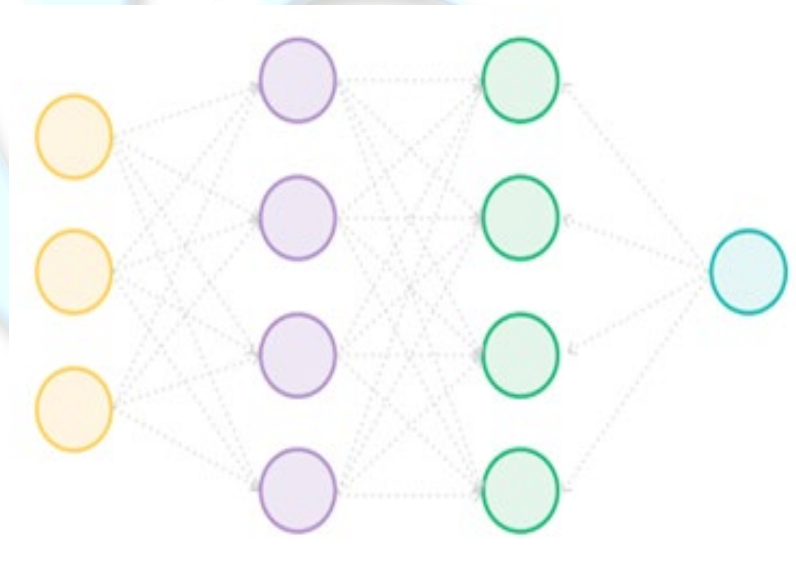


ارائه و معرفی ChatGPT توسط آزمایشگاه تحقیقاتی هوش مصنوعی آمریکایی OpenAI در دسامبر ۲۰۲۲ توجه فوق‌العاده‌ای را به خود جلب کرده است. این کنجکاوی درمورد هوش مصنوعی به‌طورکلی تا کلاس فناوری‌هایی که به‌طورخاص زیربنای چت ربات هوش مصنوعی هستند گسترش می‌یابد. این مدل‌ها که مدل‌های زبان بزرگ (LLM) نامیده می‌شوند، قادر به تولید متن درطیف به‌ظاهر بی‌پایانی از موضوعات هستند. درک LLM برای درک نحوه کار ChatGPT کلیدی است.

چیزی که LLMها را چشمگیر می‌کند توانایی آنها در تولید متنی شبیه متن نوشته شده توسط انسان در تقریباً هرزبانی (از جمله زبان‌های کدنویسی) است. این مدل‌ها یک نوآوری واقعی هستند. هیچ چیز مشابه آنها در گذشته وجود نداشته است.

این مقاله توضیح می‌دهد که این مدل‌ها چیستند، چگونه توسعه یافته‌اند و چگونه کار می‌کنند. و برای اینکه کاملاً بفهمیم چطور کار می‌کنند. همانطور که معلوم است، درک ما از اینکه چرا آنها کار می‌کنند بسیار ناچیز است.

مدل زبان بزرگ (LLM) نوعی شبکه عصبی است



شبکه عصبی نوعی مدل یادگیری ماشینی است که براساس تعدادی توابع کوچک ریاضی به نام نورون‌ها ساخته شده است. مانند نورون‌های مغز انسان، آنها پایین‌ترین سطح محاسباتی را دارند.

هر نورون یک تابع ریاضی ساده است که یک خروجی را براساس مقداری ورودی محاسبه می کند. با این حال، قدرت شبکه عصبی از اتصالات بین نورون‌ها ناشی می شود.

هر نورون به برخی از همتایان خود متصل است و قدرت هر اتصال از طریق یک وزن عددی تعیین می شود. آنها درجه‌ای را تعیین می کنند که خروجی یک نورون بعنوان ورودی به نورون بعدی در نظر گرفته می شود.

یک شبکه عصبی می تواند بسیار کوچک باشد. بعنوان مثال، یک نورون پایه می تواند شش نورون با مجموع هشت اتصال بین آنها داشته باشد. با این حال، یک شبکه عصبی نیز می تواند بسیار بزرگ باشد، همانطور که در مورد LLM ها صدق می کند. اینها ممکن است میلیون‌ها نورون با صدها میلیارد اتصال بین آنها داشته باشند که هر اتصال وزن خاص خود را دارد.

یک LLM از معماری ترانسفورماتور استفاده می کند

ما تا کنون می دانستیم که LLM نوعی شبکه عصبی است. به طور خاص، LLM ها از معماری شبکه عصبی خاصی به نام ترانسفورماتور یا مبدل استفاده می کنند که برای پردازش و تولید داده های متوالی مانند متن طراحی شده است.

معماری در این زمینه نحوه اتصال نورون‌ها به یکدیگر را توصیف می کند. همه شبکه‌های عصبی نورون‌های خود را در چندین لایه مختلف گروه بندی می کنند. اگر لایه‌های زیادی وجود داشته باشد، شبکه بعنوان «عمیق» توصیف می شود، که اصطلاح «یادگیری عمیق» از آنجا آمده است.

در یک معماری شبکه عصبی بسیار ساده، هر نورون ممکن است به هریک از نورون‌ها در لایه بالای خود متصل شود. در برخی دیگر، یک نورون ممکن است فقط به برخی از نورون‌های دیگر که در نزدیکی آن در یک شبکه قرار دارند متصل شود.

مورد دوم در شبکه‌های عصبی کانولوشنال (CNN) وجود دارد. سی‌ان‌ان‌ها پایه‌و‌اساس تشخیص تصویر مدرن را در دهه گذشته تشکیل داده‌اند. این واقعیت که CNN در یک شبکه (مانند پیکسل‌های یک تصویر) ساختار یافته است تصادفی نیست. در واقع، این دلیل مهمی است برای اینکه چرا آن معماری برای داده‌های تصویری به‌خوبی کار می‌کند.

با این حال، ترانسفورماتور تاحدودی متفاوت است. یک ترانسفورماتور که در سال ۲۰۱۷ توسط محققان گوگل ساخته شد، ایده «توجه» را معرفی می‌کند، به موجب آن نورون‌های خاصی که قوی‌تر هستند به نورون‌های دیگر در یک توالی متصل می‌شوند یا «به آنها توجه بیشتری می‌کنند».

از آنجایی که متن در یک دنباله خوانده می‌شود، یکی‌پس‌ازدیگری، با بخش‌های مختلف یک جمله که به دیگران اشاره می‌کند یا آن را تغییر می‌دهد (مانند صفتی که اسم را تغییر می‌دهد اما فعل را تغییر نمی‌دهد) همچنین تصادفی نیست که معماری ای که برای کار متوالی، با نقاط قوت ارتباط متفاوت بین بخش‌های مختلف آن دنباله ساخته شده‌است، باید روی داده‌های متنی به‌خوبی کار کند.

یک LLM خودش را می‌سازد



بعبارت ساده‌تر، مدل LLM یک برنامه کامپیوتری است. مجموعه‌ای از دستورالعمل‌ها است که محاسبات مختلفی را روی داده‌های ورودی خود انجام می‌دهد و یک خروجی ارائه می‌دهد. با این حال، چیزی که در مورد یادگیری ماشین یا مدل هوش مصنوعی مهم است، این است که به جای نوشتن آن دستورالعمل‌ها به‌طور صریح، در عوض برنامه‌نویسان انسانی مجموعه‌ای از دستورالعمل‌ها (یک الگوریتم) را می‌نویسند که سپس حجم زیادی از داده‌های موجود را برای تعریف خود مدل بررسی می‌کند. به این ترتیب، برنامه‌نویسان انسانی مدل را نمی‌سازند، بلکه الگوریتمی را می‌سازند که مدل را می‌سازد.

در مورد LLM، این بدان معناست که برنامه‌نویسان معماری مدل و قوانینی را که براساس آن ساخته می‌شود، تعریف می‌کنند. اما آنها نوروها یا وزنه‌های بین نوروها را ایجاد نمی‌کنند. این در فرآیندی به نام "آموزش" انجام می‌شود که در طی آن مدل، به دنبال دستورالعمل‌های الگوریتم، خود آن متغیرها را تعریف می‌کند.

در مورد LLM، داده‌ای که بررسی می‌شود متن است. دربرخی موارد، ممکن است تخصصی‌تر یا عمومی‌تر باشد. در بزرگ‌ترین مدل‌ها، هدف، ارائه هرچه بیشتر متن دستوری به مدل برای یادگیری است.

در ابتدا، خروجی نامفهوم است، اما از طریق یک فرآیند عظیم آزمون و خطا و با مقایسه مداوم خروجی آن با ورودی آن کیفیت خروجی به تدریج بهبود می‌یابد و متن قابل فهم‌تر می‌شود.

باتوجه به زمان کافی، منابع محاسباتی کافی و داده‌های آموزشی کافی، مدل یاد می‌گیرد که متنی را تولید کند که برای خواننده انسانی، از متن نوشته شده توسط انسان قابل تشخیص نیست. در برخی موارد، خوانندگان انسانی ممکن است بازخوردی را به شکل نوعی مدل پاداش ارائه دهند و زمانی که متن به خوبی خوانده می‌شود یا زمانی که خوانده نمی‌شود به آن بگویند (به این می‌گویند «یادگیری تقویتی از بازخورد انسانی» یا RLHF). مدل این را در نظر می‌گیرد و به طور مداوم خود را براساس آن بازخورد بهبود می‌بخشد.

یک LLM پیش‌بینی می‌کند که کدام کلمه باید کلمه قبلی را دنبال کند

یک توصیف بسیار ساده از LLMها این است که آنها «به سادگی کلمه بعدی را در یک دنباله پیش‌بینی می‌کنند». این درست است، اما این واقعیت را نادیده می‌گیرد که این فرآیند ساده می‌تواند به این معنی باشد که ابزارهایی مانند ChatGPT متن با کیفیت بسیار بالایی تولید می‌کنند.

به همین سادگی می توان گفت که «مدل به سادگی محاسبه ریاضی انجام می دهد»، که این نیز درست است، اما برای کمک به درک نحوه عملکرد مدل یا درک قدرت آن چندان مفید نیست.

نتیجه فرآیند آموزشی که در بالا توضیح داده شد یک شبکه عصبی با صدها میلیارد اتصال بین میلیون ها نورون است که هر کدام توسط خود مدل تعریف شده اند. بزرگترین مدل ها حجم زیادی از داده ها را نشان می دهند، شاید چند صد گیگابایت فقط برای ذخیره تمام وزن ها.

هریک از وزن ها و هریک از نورون ها یک فرمول ریاضی است که باید برای هر کلمه (یا در برخی موارد، بخشی از یک کلمه) که برای ورودی آن در اختیار مدل قرار می گیرد و برای هر کلمه (یا بخشی از یک کلمه) محاسبه شود که به عنوان خروجی خود تولید می کند.

این جزئیات فنی است، اما به این «کلمات کوچک یا بخش هایی از کلمات» «نشان ها» یا «توکن ها» گفته می شود، که معمولاً وقتی استفاده از این مدل ها بعنوان یک سرویس ارائه می شوند، قیمت گذاری می شود. در ادامه در مورد آن بیشتر توضیح خواهیم داد.

کاربر در حال تعامل با یکی از این مدل ها، ورودی را در قالب متن ارائه می دهد. برای مثال، می توانیم دستور زیر را به ChatGPT ارائه کنیم:

سلام ChatGPT ، لطفاً یک توضیح ۱۰۰ کلمه ای از Dataiku به من ارائه دهید.

شرحی از نرم افزار و ارزش اصلی آن را درج کنید

سپس مدل‌های پشت ChatGPT این درخواست را به توکن تبدیل می‌کنند. به طور متوسط، یک نشانه $\frac{4}{5}$ از یک کلمه است، بنابراین دستور بالا و ۲۳ کلمه آن ممکن است منجر به حدود ۳۰ نشانه شود. مدل GPT-3 که مدل gpt-3.5-turbo مبتنی بر آن است، ۱۷۵ میلیارد وزن دارد، به این معنی که ۳۰ توکن متن ورودی به $۱۷۵ \times ۳۰ = ۵,۲۵$ تریلیون محاسبات منجر می‌شود. مدل GPT-4 که در ChatGPT نیز موجود است، دارای وزن نامشخصی است.

سپس، مدل شروع به تولید پاسخی می‌کند که براساس حجم متنی که در طول آموزش مصرف کرده، درست به نظر می‌رسد. نکته مهم این است که چیزی در مورد سوال جستجو نمی‌کند. هیچ حافظه‌ای ندارد که بتواند «dataiku»، «value proposition»، «software» یا هر عبارت مرتبط دیگری را جستجو کند. در عوض، تولید هر نشانه متن خروجی را آغاز می‌کند، ۱۷۵ میلیارد محاسبات را دوباره انجام می‌دهد، و رمزی را تولید می‌کند که به احتمال قوی‌تر درست به نظر می‌رسد.

LLMها متنی را تولید می‌کنند که درست به نظر می‌رسد اما نمی‌توانند تضمین کنند که درست باشد.

ChatGPT نمی‌تواند تضمینی برای درست بودن خروجی‌ش ارائه دهد، آن فقط درست به نظر می‌رسد. پاسخ‌های آن در حافظه‌اش جستجو نمی‌شوند آنها بر اساس ۱۷۵ میلیارد وزنی که قبلاً توضیح داده شد، ایجاد می‌شوند.

این نقص مختص ChatGPT نیست، بلکه مربوط به وضعیت فعلی همه LLMها است. آنها مهارتی در یادآوری واقعیات ندارند. ساده‌ترین پایگاه‌های داده این کار را به خوبی انجام می‌دهند. در عوض، نقطه قوت آنها در تولید متنی است که مانند متن نوشته شده توسط انسان خوانده می‌شود و خوب به نظر می‌رسد. در بسیاری از موارد، متنی که درست به نظر می‌رسد نیز در واقع درست خواهد بود، اما نه همیشه.

در آینده، این احتمال وجود دارد که LLMها در سیستم‌هایی ادغام شوند که قدرت تولید متن LLM را با یک موتور محاسباتی یا پایگاه دانش ترکیب می‌کنند تا پاسخ‌های واقعی را در متن زبان طبیعی بصورت قانع کننده‌ای ارائه دهند. آن سیستم‌ها امروزه وجود ندارند، اما به راحتی می‌توان تخمین زد که چقدر طول می‌کشد تا آنها را ببینیم.

امکان دیگر این است که اگر می‌خواهید اطلاعاتی را که قبلاً دارید در قالب پاسخ زبان طبیعی به کاربران ارائه دهید، می‌توانید آن پاسخ‌ها را به ابزارهایی مانند ChatGPT ارائه دهید و از آنها بخواهید براساس آن پاسخ‌ها جوابی بسازند. Dataiku یک نسخه نمایشی با استفاده از GPT-3 برای ارائه پاسخ از اسناد Dataiku ایجاد کرده است که دقیقاً این کار را انجام می‌دهد.

آیا GPT-4 یک LLM است؟

در ۱۴ مارس ۲۰۲۳، OpenAI، آخرین نسخه از مدل‌های خود را در خانواده GPT به نام GPT-4 منتشر کرد. علاوه بر تولید متن با کیفیت بالاتر در مقایسه با GPT-3.5، GPT-4 توانایی تشخیص تصاویر را نیز ارائه می‌دهند. ممکن است قادر به تولید تصاویر نیز باشد. با این حال، این قابلیت، اگر وجود داشته باشد، هنوز در دسترس نیست. توانایی مدیریت داده‌های ورودی و خروجی از انواع مختلف (متن، تصاویر، ویدئو، صدا و غیره) به این معنی است که GPT-4 چندوجهی است.

اصطلاحات مربوط به این مدل‌های آخر به سرعت در حال تکامل است، مطابق با برخی از بحث‌ها در جامعه متخصص استدلال می‌شود که "مدل زبان" بسیار محدودکننده است. اصطلاح "مدل بنیاد" توسط محققان در استنفورد رایج شده است، اما همچنین منبع بحث‌هایی است. مانند خود فناوری، زبان مورد استفاده برای توصیف فناوری به سرعت به تکامل خود ادامه خواهد داد.

استفاده از ChatGPT، GPT-4 و مدل‌های زبان بزرگ (LLM) در سازمان

ما از ChatGPT و یکی از مدل‌های آن، gpt-3.5-turbo، بعنوان مثال در سراسر این مقاله استفاده کرده‌ایم، اما این تنها یک مدل و یک محصول درمیان بسیاری از آنها است. برخی از LLMها اختصاصی هستند و از طریق یک رابط وب یا یک API مانند ChatGPT قابل دسترسی هستند. سایر LLMها منبع باز (open source) هستند و اگر توان محاسباتی و مهارت لازم برای انجام این کار را داشته باشند، می‌توانند توسط یک دانشمند یا مهندس داده باهوش دانلود و اجرا شوند. برای هر رویکرد جانشین‌هایی وجود دارد.

#چت جی پی تی #مدل زبان بزرگ #ChatGPT #GPT-4 #gpt-3.5-turbo #LLM

تماس با ما:



شرکت عصر ارتباطات بین الملل پارس کار (ایکاست)

ارائه دهنده خدمات ارتباطات داده ماهواره ای با مجوز SAP، از سازمان تنظیم مقررات و ارتباطات رادیویی

آدرس: تهران، سعادت آباد، میدان بهرود، خیابان عابدی، پلاک ۱۵ ساختمان صبا، طبقه سوم واحد ۸ - کد پستی: ۱۹۸۱۸۶۳۶۹۵

تلفن: +۹۸۲۱۷۵۲۲۹۲۲۹

فکس: +۹۸۲۱۷۵۲۲۹۲۳۹

وبگاه: www.icasat.org

پست الکترونیک: cmo@icasat.net