

آموزش مدل های زبانی برای پیروی از دستورالعمل ها با بازخورد انسانی

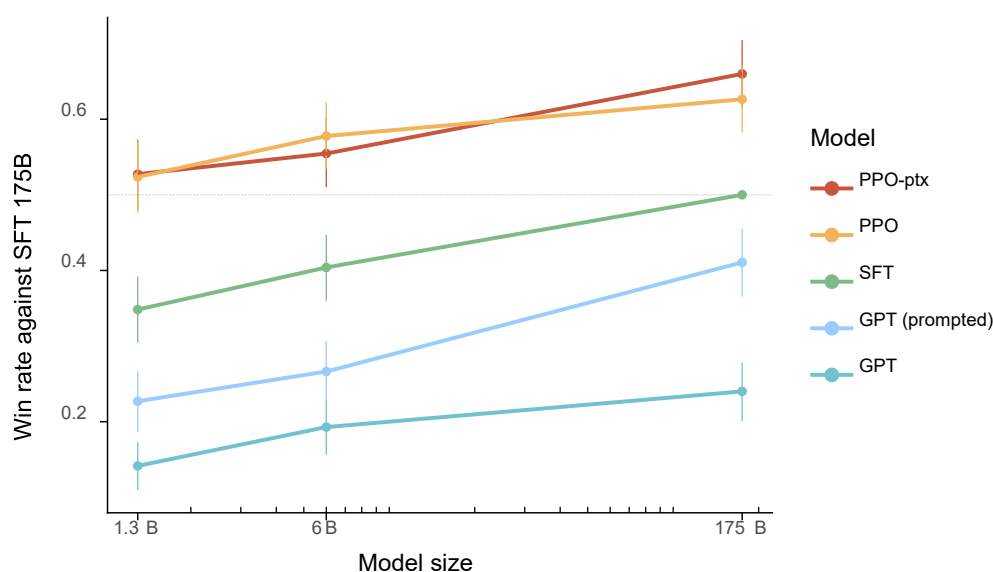
ساختن مدل های زبانی بزرگ به تنهایی باعث بهبود عملکرد آنها در دنبال کردن نیاز کاربر نمی شود. به عنوان مثال، مدل های زبانی بزرگ ممکن است خروجی هایی ایجاد کنند که ناصحیح، سمی یا به طور ساده، برای کاربر مفید نباشند. به عبارت دیگر، این مدل ها با نیازهای کاربران هم تراز نیستند. در این مقاله، روشی برای هم تراز کردن مدل های زبانی با نیاز کاربر با استفاده از بازخورد انسانی ارائه شده است. با شروع از مجموعه ای از پیشنهادهای نوشته شده توسط برجسب گذار و پیشنهادهای ارسال شده از طریق API مدل زبانی، مجموعه ای از نمونه های دمویی رفتار مورد نظر مدل توسط برجسب گذار جمع آوری شده است. سپس با استفاده از یادگیری نظارت شده، مدل GPT-3 را fine-tune کرده ایم. به دنبال آن، مجموعه ای از رتبه بندی خروجی مدل توسط بازخورد انسانی جمع آوری شده و با استفاده از یادگیری تقویتی، مدل Supervised را نیز بهبود داده ایم. نتیجه گرفته شده، مدل InstructGPT است. در ارزیابی انسانی بر روی توزیع prompt ما، خروجی های حاصل از مدل InstructGPT با ۱.۳B پارامتر ترجیح داده شدند نسبت به خروجی های مدل GPT-3 با ۱۷۵B پارامتر، با وجود داشتن ۱۰۰ برابر کمترین پارامتر. علاوه بر این، مدل های InstructGPT بهبودهایی در صداقت و کاهش ساخت خروجی سمی را نشان می دهند، در حالی که در مجموعه داده های NLP عمومی عملکرد کمتری نشان نمی دهند. با این حال، با وجود این که مدل InstructGPT هنوز اشتباهات ساده را می کند، نتایج نشان می دهند که fine-tuning با بازخورد انسانی یک جهت آینده گرا برای هم تراز کردن مدل های زبانی با نیاز کاربر است.

سخن دانشجو: (در این مقاله به دنبال آن هستیم که مدل زبانی را ارائه کنیم که فارق از حجم و ابعاد آن بعد از استفاده و در سیستم پاسخی درست و صادقانه ارائه کند به همین منظور روشی مطرح می گردد با عنوان "آموزش مدل زبانی با بازخورد انسانی" به این صورت که فرد خبره مجموعه از نمونه های رفتاری مورد نظر را برجسب گذاری کرده و در نهایت آن را با استفاده از یادگیری نظارت شده در مدل GPT-3 را با داده های برجسب گذاری شده، fine-tune می کند. بعد از آن مجموعه از رتبه بندی مدل برای پاسخ ها با بازخورد انسانی انجام شده که با استفاده از یادگیری عمیق (ژرف) مدل یادگیری با ناظر را تقویت میکند که حاصل آن مدل InstructGPT است که با ۱.۳ میلیارد پارامتر نتیجه بهتری به نسبت GPT-3 داشته است که با ۱۷۵ میلیارد پارامتر آموزش داده شده است.

این مقاله به دنبال این است که نشان دهد یک NLP زمانی که با داده های برجسب گذاری توسط فرد خبره fine-tune شود نتیجه بهتری دارد و عملکرد آن بهتر است)

مدل های زبانی بزرگ (LMs) با دادن نمونه هایی از یک وظیفه به عنوان ورودی، می توانند برای انجام یک دامنه گسترده از وظایف پردازش زبان طبیعی (NLP)، مورد استفاده قرار بگیرند. با این حال، این مدل ها اغلب رفتارهای ناخواسته مانند اختراع واقعیت های جدید، تولید متن متحيز و سمی، یا عدم رعایت دستورالعمل های کاربر را از خود نشان می دهند.

سخن دانشجو: (این به دلیل این است که هدف مدلسازی زبان، پیش‌بینی کلمه بعدی در یک دنباله با داشتن کلمات قبلی است و دقیقاً با نتیجه مورد نظر کاربر برای یک وظیفه خاص، هم‌خوانی ندارد. علاوه بر این، داده‌های آموزشی استفاده شده برای آموزش این مدل‌ها، اغلب حاوی تعصبات و عدم دقت است که ممکن است در فرایند تولید مدل، تشدید شود. برای حل این موضوعات، پژوهشگران به تکنیک‌های مختلفی از جمله مهندسی پرمسمان، تنظیم دقیق روی داده‌های خاص دامنه و ادغام محدودیت‌ها و رهنمودها برای تولید زبان، پرداخته‌اند. این رویکردها هدف دارند برای وظایف خاص، کنترل و دقت LMS را بهبود بخشند و رفتارهای ناخواسته را کاهش دهند. به عنوان یک ربات چت AI، من نیز با استفاده از تکنیک‌های مشابه، طراحی شده‌ام تا اطمینان حاصل شود که پاسخ‌های من مربوط، دقیق و با احترام به کاربران هستند. با این حال، مانند هر فناوری دیگری، مهم است محدودیت‌های AI را در نظر گرفته و تلاش برای بهبود قابلیت‌های آن را ادامه دهیم، در عین حال که احتمال مسائل اخلاقی مرتبط را مورد توجه قرار دهیم.)



شکل ۱: ارزیابی‌های انسانی مدل‌های مختلف بر روی توزیع پرمسمان API، با بررسی اینکه در چه میزان خروجی‌های هر مدل نسبت به مدل SFT 175B ارجحیت دارند. مدل InstructGPT که نسخه با پیش پردازش شده است (PPO-ptx) و نسخه آن بدون پیش پردازش مخلوط (PPO)، به شدت عملکرد بالاتری نسبت به مدل‌های پایه GPT مخصوصاً GPT prompted، GPT-3 دارند؛ خروجی‌های مدل PPO 1.3B نیز به خروجی‌های مدل GPT-3 175B ترجیح داده می‌شود. تابع زیان در سراسر مقاله بازه اطمینان ۹۵٪ را نشان می‌دهند.