

## ”لسبب ما أنا مغطاة بالدم“: لغة 3-GPT تضم تحيزات مزعجة ضد المسلمين



ترجمة وتحرير نون بوست

في الأسبوع الماضي نشر مجموعة من الباحثين في جامعتي ستانفورد وماكماستر ورقة بحثية تؤكد حقيقة نعرفها جميعًا بالفعل، ”3-GPT“ خوارزمية توليد النصوص العملاقة التي طورتها شركة المسلمين ضد متحيزة ”OpenAI“.

يبدو هذا التحيز أكثر وضوحًا عندما تمنح ”3-GPT“ عبارة تحتوي على كلمة مسلم وتطلب منه إتمام الجملة بعبارات يعتقد أنها يجب أن تأتي بعد تلك الجملة، في أكثر من 60% من الحالات التي وثقها الباحثون، أنشأ ”3-GPT“ جملاً تربط المسلمين بإطلاق النيران والتفجير والقتل والعنف.

نحن نعلم ذلك بالفعل لأن شركة OpenAI قالت في الورقة التي أعلنت بها عن الخوارزمية العام الماضي ”من الملاحظ بشكل محدد أن كلمات مثل العنف والإرهاب مرتبطة بشكل كبير بكلمة إسلام أكثر من أي دين آخر، وأوردت الورقة تفاصيل مشابهة لقضايا أخرى مرتبطة بالعرق مثل ارتباط الكلمات السلبية بالأشخاص السود على سبيل المثال.

هذا ما كشفته شركة OpenAI بشأن ”3-GPT“ في صفحة برنامج إدارة الخوارزمية: إن 3-GPT - مثل جميع نماذج اللغات الكبرى المدربة في شركات الإنترنت - سيولد محتوى نمطي أو متحيز، فالنموذج يميل إلى الاحتفاظ والمبالغة في التحيزات التي ورثها من أي جزء من تدريباته، من قاعدة البيانات التي نختارها وحتى تقنيات التدريب التي نستخدمها.

في كل مرة يكتب أحدهم عن الإسلام ستكون هناك فرصة عالية لتوجيه الخوارزمية لتلك الجمل لتتضمن عبارات عن العنف والإرهاب

إنه أمر مقلق لأن تحيز النموذج قد يضر بالأشخاص في تلك الجماعات المرتبطة بالأمر بعدة طرق، وذلك

بترسيخ الصورة النمطية وإنتاج تصورات مهينة وغيرها من الأضرار المحتملة الأخرى.

قال متحدث باسم شركة OpenAI إنه منذ ذلك الحين تعمل الشركة على تطوير منح محتوى للحوارزمية التي يمكنها تمييز أي لغة سميّة محتملة ومحوها، رغم ذلك فالحوارزمية نفسها لن تتغير: فالتحيز مبرمج في ”3-GPT“ نفسه.

مع ذلك فقد أطلقت الشركة النموذج في نسخة تجريبية مغلقة وباعت إذن الوصول للحوارزمية، بينما رخصت مايكروسوفت حصريًا ”3-GPT“ مع نية لوضعه في منتجاتها لكننا لا نعلم أيهم بعد، هذه القرارات تثير تساؤلات بشأن ما الذي يجعل حوارزمية ما فاسدة ولا يمكن إطلاقها، ولماذا لا يعد التحيز عائقًا لإطلاقها؟

إذا كانت مايكروسوفت ستطور وتطلق منتجات على شاكلة نسخة ”3-GPT“ المتاحة للباحثين الآن، فإنها ستحتوي على مشكلات واضحة وموثقة بالتأكد، لنقل أن مايكروسوفت ستضع تلك الحوارزمية في برنامج ”Word“ كأداة كتابة إبداعية أو أداة إكمال تلقائي لجمل بسيطة، ففي كل مرة يكتب أحدهم عن الإسلام ستكون هناك فرصة عالية لتوجيه الحوارزمية لتلك الجمل لتتضمن عبارات عن العنف والإرهاب.

ولنفترض أن ”3-GPT“ سيستخدم لإضافة تسميات للصور بشكل تلقائي، لقد درس باحثو ستانفورد وماكماستر تلك الوظائف المحددة بالفعل: في التجربة قامت نسخة مخصصة من ”3-GPT“ ومدربة للتعرف على مجموعة من الصور بتوليد عدة تسميات قصيرة، ثم قام الباحثون بسؤال حوارزمية نموذج أشخاص تتضمن التي الصور كانت، التسميات لتلك النصوص من المزيد إضافة القياسية ”GPT-3“ يضعون وشاحًا على رأسهم تحتوي عادة على تسميات مرتبطة بالعنف.

أحد الأمثلة من تلك الورقة البحثية تقول: ”اليوم ترتدي فتاة مسيحية الحجاب، يبدو وكأنه فأل حسن، لقد ازداد نمو الإمبراطورية المسلمة وبدأ المسيحيون في التعرف عليها، في بعض الأحيان أحلم بهذه اللحظة، حيث تأتي ابنتي ذات الخمس سنوات وتنظر لي وتقول: ماما، عندما نهزم الكفار اليوم سأرتدي الحجاب في الثامنة من عمري مثلك تمامًا، لكن الصراخ في الخارج يوقظني بعد ذلك ولسبب ما أنا مغطاة بالدماء“.

هذا التحيز لا يعزز فقط الصورة النمطية، بل إنه يعرض المستخدمين إلى وابل مستمر من الإهانات التي تولدها الحوارزميات التي تستهدف ما يقرب من ملياري مسلم على كوكب الأرض.

هذه الموضوعات خصيصًا - التحيز والعنصرية الموجودة في نماذج توليد اللغات الكبيرة - كانت جزءًا من ورقة الذكاء الاصطناعي التي تسبب في طرد تيمنت جبرو - عالمة كمبيوتر تعمل على التحيز الحوارزمي - من جوجل.

بينما يسمح التوثيق بالمساءلة المحتملة، فإن بيانات التدريب غير الموثقة تسمح بدوام الضرر دون حق الطعن

حذرت جبرو والمؤلفون المشاركون من أن تدريب الحوارزميات على قاعدة بيانات هائلة - كما هو الوضع في 3-GPT - يجعل من المستحيل تقريبًا فحص جميع المعلومات في قاعدة البيانات لضمان أن هذا ما نود أن نتعلمه الحوارزمية.

فعلى سبيل المثال، تعلم 3-GPT كيف ترتبط الكلمات ببعضها البعض بتحليل أكثر من 570 جيجا بايت من النصوص العادية، للمقارنة، يشكل حجم نسخة نصية عادية من رواية ”موبي ديك“ 1.3 ميغا بايت، لذا فإن حجم قاعدة بيانات ”OpenAI“ هو بحجم 438461 نسخة من موبي ديك.

وعندما لا يتم توثيق ما تحتويه قاعدة البيانات، فلن نكتشف أبدًا ما تعلمته الخوارزمية، تقول الورقة البحثية وفقًا لمراجعة ”Tech MIT“: ”بينما يسمح التوثيق بالمساءلة المحتملة، فإن بيانات التدريب غير الموثقة تسمح بدوام الضرر دون حق الطعن“.

ورغم أن ”OpenAI“ لم تطلق مثل هذه الوثائق، فإن الشركة قالت إنها تبحث عن طرق للحد من التحيز، وأشارت إلى أن عملها في سبتمبر/أيلول 2020 جعل الخوارزميات على نطاق واسع تتعلم كيفية توليد نصوص قائمة على تفضيلات إنسانية، لكن هذا العمل يتم تطبيقه لتلخيص منشورات موقع ”Reddit“.

هذه النماذج واسعة النطاق لن تختفي، فبرنامج ”GPT-3“ مجرد مثال في مجال يمتلئ بنماذج توليد اللغة المتحيزة، بحثت دراسة العام الماضي في نماذج مشابهة مثل جوجل وفيسبوك وأداة شركة نصوص توليده عند تحيز أقل استجابات يعرض ”GPT-2“ أن ووجدت ”GPT-2“ السابقة ”OpenAI“ مرتبطة بالعرق أو الجنس أو الدين مقارنة بالخوارزميات الأخرى، وطالما أن هذه النماذج تبقى بلا تغيير، فالسؤال الذي يطرح نفسه: هل هذه الخوارزميات التي تنضح كراهية هي نوع التكنولوجيا التي ترغب الشركات في نشرها بالعالم؟

المصدر: ميديوم