

ANALISI

La fonte dell'intelligenza artificiale incorpora un serio problema etico

PAOLO BENANTI

Un'inchiesta giornalistica illustra le possibili derive dovute alla composizione dei «dataset» di informazioni. Si chiama «C4» ed è il colossale campione di dati e pagine web all'origine dei sempre più evoluti sistemi di AI. Rischi per la democrazia ma anche di una colonizzazione culturale. Il termine C4 farà pensare a molti, vista la notorietà acquisita tra film e videogiochi, all'esplosivo al plastico. In realtà c'è un altro C4, non meno esplosivo nei contenuti, che in questi giorni sta facendo notevolmente discutere. Uno dei processi fondamentali che permette a motori di ricerca come Google, Bing e Yahoo di indicizzare un contenuto su internet è il cosiddetto crawling, ovvero – semplificando un po' – un software che analizza i contenuti di una rete (o di un database) in un modo metodico e automatizzato acquisendo una copia testuale di tutti i documenti presenti e creando un indice che ne permetta, successivamente, la ricerca e la visualizzazione. Esiste un corpus, detto Common Crawl, che contiene petabyte di dati raccolti in 12 anni di web crawling. Il corpus contiene dati grezzi di pagine web, estratti di metadati ed estratti di testo.

Di tutto questo è stata fatta una versione definitiva colossale e ripulita (cleaned). Si ottengono così componenti – le quattro C – di questo C4: Colossal, Cleaned version of Common Crawl.

Il quotidiano statunitense "Washington Post" ha reso pubblico di recente in un lungo articolo investigativo che il dataset C4 contiene dati provenienti da fonti quali Stormfront, Kiwi Farms, 4chan e altri siti web tutti considerati potenzialmente problematici. Tra le fonti emerse dall'indagine ci sono almeno 27 siti web identificati dal governo statunitense come connessi ai mercati della contraffazione e della pirateria o dati provenienti da siti come "Vdare", un sito statunitense di estrema destra che promuove l'opposizione all'immigrazione ed è associato alla supremazia bianca, al nazionalismo bianco e all'alt-right, e "Breitbart", un sito di notizie di estrema destra considerato da accademici e giornalisti misogino, xenofobo, ma anche l'emittente russa "RT", il primo tra i canali della Russia completamente in digitale direttamente finanziata dal Cremlino.

Questo testo che compone il C4, di fatto, è la base che costituisce la principale fonte di addestramento e di acquisizione di informazioni che le intelligenze artificiali (AI) mostrano di possedere sul mondo, e inevitabilmente influenza il modo in cui ogni AI risponde alle richieste e alle interazioni degli utenti. Se costruiamo intelligenze artificiali come Gpt-4 che sono in grado di superare alcuni dei più severi test di ammissione alle facoltà universitarie, ad esempio, è molto probabile che questa capacità emergente del sistema sia connessa ai dati di addestramento che hanno incluso migliaia di siti con test di esercitazione per questi esami.

Le aziende tecnologiche, però, hanno innalzato una ferrea cortina di mistero su ciò che hanno dato in



Avvenire

pasto in fase di addestramento all'intelligenza artificiale. E se, nell'utilizzarle, ci sorprende quanto sembrano in grado di fare, di fatto rimaniamo ciechi sulle fonti e sulle origini di questo sapere. La cosa sembra costituire una vera e propria inversione rispetto alla modalità scientifica della conoscenza che ha fatto della trasparenza sui dati e sulle fonti una delle cifre della stessa scientificità. Per guardare all'interno di questa scatola nera il "Post" ha analizzato il set di dati C4 di Google collaborando con i ricercatori dell'Allen Institute for AI e hanno classificato i siti web utilizzando i dati di Similarweb, una società di analisi web. Circa un terzo dei siti web da cui sono estratti i dati in origine non ha potuto essere classificato, soprattutto perché non sono più presenti su Internet.

Stando a quanto dichiara Google, C4 è stato inizialmente sviluppato come "versione ripulita" dei dati di Common Crawl ed è stato utilizzato per addestrare alcune AI di alto profilo in lingua inglese, chiamate modelli linguistici di grandi dimensioni, o Llm, tra cui il T5 di Google e LLaMA di Facebook. OpenAI, di contro, non rivela quali set di dati utilizza per addestrare i modelli che supportano il suo popolare chatbot, ChatGpt, appena tornato fruibile in Italia.

Quello che ci interessa sottolineare, prima di fare ulteriori analisi, è il fatto che un sito web viene indicizzato in C4 solo se è in inglese e che il dataset, non contenendo dati in altre lingue, è anglofono. Grazie al tool fornito dal quotidiano statunitense abbiamo fatto ulteriori indagini cercando alcune fonti. Un dato interessante, per esempio, è che il dominio vatican.va è al 4.967 posto avendo fornito quasi 2 milioni di token (i piccoli frammenti di testo che costituiscono la base delle informazioni con cui è addestrato il sistema). Non bisogna farsi impressionare dal numero di classifica perché le fonti sono oltre 15,7 milioni e di fatto un risultato sotto i primi 5.000 è altissimo, soprattutto se si considera che si prendono in esame solo le pagine in inglese e non tutte le pagine del sito. Per fare un confronto, la Cia, che contiene il Cia World Factbook, una pubblicazione annuale che riporta i dati statistici fondamentali e una sintesi di informazioni riguardanti tutti i Paesi del mondo, è dietro il sito vaticano di quasi 600 posizioni. Parlando di questo con Denis "Jaromil" Roio, il famoso programmatore, "hacker etico", artista digitale e attivista, ci è venuto in mente di cercare anche il sito della Nasa: ebbene, Nasa.gov è staccato di 100 posizioni. La battuta è sorta spontanea: per sapere del cielo C4 si affida più alla Chiesa che non all'astrofisica.

Il "Washington Post" riporta un'interessante analisi aggregata del dataset: «I siti web commerciali e industriali costituiscono la categoria più grande (16% dei token categorizzati), guidata da fool.com al n. 13, che fornisce consigli sugli investimenti. Poco distante kickstarter.com, al n. 25, che consente agli utenti di finanziare in crowdfunding progetti creativi, e più in basso patreon.com, n. 2.398, che aiuta i creatori a raccogliere quote mensili dagli abbonati per contenuti esclusivi. Kickstarter e Patreon potrebbero dare all'intelligenza artificiale l'accesso alle idee e alle copie di marketing degli artisti, sollevando il timore che la tecnologia possa copiare questo lavoro per suggerirlo agli utenti». L'iniziativa del "Washington Post", che ha per la prima volta permesso di analizzare uno di questi set di dati per rivelare completamente i tipi di siti web proprietari, personali e spesso offensivi che entrano nei dati di addestramento di un'intelligenza artificiale, è

di grande interesse, utilissima per iniziare a pensare se e come questi sistemi siano adeguati e rispettosi delle identità culturali e democratiche dei Paesi occidentali e in particolare di quell'area sempre più isolata nella difesa della rule of law che è l'Europa.

L'indagine ha rilevato che C4 è in maggior parte dominato da siti web legati al giornalismo, alla creazione di contenuti, all'intrattenimento e allo sviluppo di software, con patents. google.com, wikipedia. org e scribd.com elencati come primi tre siti. Tuttavia i dati di addestramento provenienti da siti più discutibili potrebbero potenzialmente indurre i modelli di intelligenza artificiale a generare testi indesiderati, razzisti, pornografici, inaffidabili e in generale dannosi. L'Algoritmica chiede di riflettere su questa sorta di "materia prima" per le AI: la qualità dei dati influenza la qualità e l'affidabilità dei sistemi su cui avviene l'addestramento. Dobbiamo chiederci se e come una scelta di cosa includere in C4 non sia di fatto anche un'opzione di natura politica e con severe conseguenze geopolitiche. Di fatto la scelta dei dati è una scelta – una tokenizzazione, per usare un termine tecnico – di una cultura. Questa scelta potrebbe, estremizzando un po', trasformare il tradizionale softpower culturale esercitato dall'industria dei media e da Hollywood al livello di un vero e proprio colonialismo culturale.

Forse anche per questo di recente il governo britannico ha stanziato 100 milioni di sterline per la creazione di una task force sull'AI incaricata di creare modelli di base o modelli di IA pre-addestrati, come Gpt di OpenAI. Creare un set di dati pubblico è la prima forma di difesa. Dobbiamo chiederci se non sia arrivato il momento di pensare alla creazione di un nostro dataset culturalmente pesato ed eticamente bilanciato per permettere al Paese e ai servizi pubblici di beneficiare dell'impatto trasformativo dell'AI. RIPRODUZIONE RISERVATA L'addestramento proveniente da siti discutibili potrebbe potenzialmente indurre a generare testi indesiderati, razzisti, pornografici, oltre che inaffidabili. Dobbiamo chiederci se e come una scelta di cosa includere non sia anche un'opzione politica e con severe conseguenze geopolitiche. Il Washington Post ha reso noto che la fonte di dati «C4» contiene informazioni provenienti da siti considerati pericolosi.

Nella foto: un'immagine tratta dal film "Matrix".