



IA generativa:

Perspectivas de Stanford HAI

Como você acha que a IA generativa afetará
seu campo e a sociedade daqui para frente?

Março de 2023

Índice

Introdução	3
O Grande Ponto de Inflexão da IA, Fei-Fei Li	4
Os Potenciais dos Pacientes Sintéticos, Russ Altman	6
Upending Healthcare, do atendimento ao paciente ao faturamento, Curt Langlotz	7
Uma janela de IA para a natureza, Surya Ganguli	8
As novas ferramentas da vida diária, James Landay	10
A poesia não será otimizada: criatividade na era da IA, Michele Elam	11
IA generativa e o Estado de Direito, Daniel E. Ho	13
A Nova Era Cambriana: 'Excitação Científica, Ansiedade', Percy Liang	15
Um chamado para aumentar – não automatizar – trabalhadores, Erik Brynjolfsson	16
A Reinvenção do Trabalho, Christopher D. Manning	18
Na educação, um 'desastre em formação', Rob Reich	20
Resolvendo Desigualdades no Sistema Educacional, Peter Norvig	21

Introdução:

A onda atual de IA generativa é um subconjunto de inteligência artificial que, com base em uma solicitação textual, gera conteúdo novo. ChatGPT pode escrever um ensaio, Midjourney pode criar belas ilustrações ou MusicLM pode compor um jingle. A maior parte da IA generativa moderna é alimentada por modelos básicos ou modelos de IA treinados em dados amplos usando autosupervisão em escala e depois adaptados a uma ampla gama de tarefas posteriores.

As oportunidades que estes modelos apresentam para as nossas vidas, as nossas comunidades e a nossa sociedade são vastas, tal como os riscos que representam. Embora, por um lado, possam complementar perfeitamente o trabalho humano, tornando-nos mais produtivos e criativos, por outro, podem amplificar o preconceito que já experimentamos ou minar a nossa confiança na informação.

Acreditamos que a colaboração interdisciplinar é essencial para garantir que estas tecnologias beneficiem a todos nós. A seguir estão as perspectivas dos líderes de Stanford em medicina, ciência, engenharia, humanidades e ciências sociais sobre como a IA generativa pode afetar seus campos e o nosso mundo. Alguns estudam o impacto da tecnologia na sociedade, outros estudam a melhor forma de aplicar estas tecnologias para avançar no seu campo, e outros desenvolveram os princípios técnicos dos algoritmos que fundamentam os modelos fundamentais.

O grande ponto de inflexão da IA



Fei-Fei Li, Sequoia Capital
Professor de Informática
Departamento de Ciências; Denning
Codiretor do Stanford HAI

Há 540 milhões de anos, o número de espécies animais explodiu num período de tempo muito curto. Existem muitas teorias sobre o que aconteceu, mas uma chamou minha atenção: o início repentino e a evolução onírica da visão. Hoje, a percepção visual é um sistema sensorial importante e a mente humana pode reconhecer padrões no mundo e gerar modelos ou conceitos baseados nesses padrões. Dotar as máquinas com estas capacidades, capacidades generativas, tem sido um sonho para muitas gerações de cientistas de IA. Há uma longa história de tentativas algorítmicas de modelos generativos com vários graus de progresso. Em 1966, pesquisadores do MIT desenvolveram o "Summer Vision Project" para construir efetivamente "uma parte significativa do sistema visual" com tecnologia. Este foi o início do campo da visão computacional e geração de imagens.

Recentemente, devido aos conceitos profundos e interligados de aprendizagem profunda e grandes volumes de dados, parece que atingimos um ponto de inflexão na capacidade das máquinas de gerar linguagem, imagem, áudio e muito mais. Embora a construção de IA para ver o que os humanos podem ver tenha sido a inspiração para a visão computacional, deveríamos agora olhar além disso, para construir IA para ver o que os humanos não podem ver. Como podemos usar a IA generativa para aumentar a nossa visão? Embora o número exato seja contestado, as mortes devido a erros médicos nos EUA são um problema significativo. Os modelos generativos de IA podem ajudar os prestadores de cuidados de saúde a identificar potenciais problemas que, de outra forma, poderiam ter perdido. Além disso, se os erros forem devidos à exposição mínima a substâncias raras

situações, a IA generativa pode criar versões simuladas desses dados raros para treinar ainda mais os modelos de IA ou o próprios prestadores de cuidados de saúde.

Além disso, antes mesmo de começarmos a desenvolver novas ferramentas generativas, precisamos nos concentrar no que as pessoas desejam dessas ferramentas. Em um projeto recente para avaliar tarefas de robótica do nosso laboratório, antes mesmo de iniciar a pesquisa, a equipe do projeto fez um estudo de usuário em larga escala para perguntar às pessoas o quanto elas gostariam de beneficiária se um robô fizesse essas tarefas para eles.

As tarefas vencedoras foram o foco da pesquisa.

*Dotar as máquinas com estas
capacidades, capacidades
generativas, tem sido um
sonho para muitas gerações de
cientistas de IA.*

Para concretizar plenamente a oportunidade significativa que a IA generativa cria, precisamos também de avaliar os riscos associados. Joy Buolamwini liderou um estudo intitulado "Gender Shades", que descobriu que os sistemas de IA frequentemente falham em reconhecer mulheres e pessoas de cor. Os resultados do estudo foram publicados em 2018. Continuamos a observar preconceitos semelhantes em modelos de IA generativos, especificamente para populações sub-representadas.

O Grande Ponto de Inflexão da IA (cont.)

A capacidade de determinar se uma imagem foi gerada usando IA também é essencial. Nossa sociedade é construída com base na confiança na cidadania e na informação. Se não conseguirmos determinar facilmente se uma imagem é gerada por IA, nossa confiança em qualquer informação será destruída. Neste caso, precisamos de prestar especial atenção às populações vulneráveis que podem ser particularmente susceptíveis a utilizações adversas desta tecnologia.

O progresso na capacidade de uma máquina gerar conteúdo é muito emocionante, assim como o potencial para explorar a capacidade da IA de ver aquilo de que os humanos não são capazes. Mas precisamos de estar atentos às formas como estas capacidades irão perturbar a nossa vida quotidiana, as nossas comunidades e o nosso papel como cidadãos do mundo.

Os potenciais dos pacientes sintéticos



Russ Altman, professor Kenneth Fong na Escola de Engenharia; Professor de Bioengenharia, de Genética, de Medicina e de Dados Biomédicos Ciência; Diretor Associado de Stanford HAI

Muitas vezes é difícil obter um grande número de pacientes em ensaios clínicos e é crucial ter um grupo realista de pacientes que não recebem uma terapia, a fim de comparar os resultados com aqueles que o fazem. Esta é uma área da investigação biomédica onde a IA generativa oferece grandes oportunidades. A IA generativa poderia tornar os ensaios clínicos mais eficientes, criando pacientes de controlo “sintéticos” (ou seja, pacientes falsos) utilizando dados de pacientes reais e os seus atributos subjacentes (para serem comparados com os pacientes que recebem a nova terapia). Poderia até gerar resultados sintéticos para descrever o que acontece com esses pacientes se não forem tratados. Os investigadores biomédicos poderiam então utilizar os resultados de pacientes reais expostos a um novo medicamento com os resultados estatísticos sintéticos para os pacientes sintéticos. Isto poderia tornar os ensaios potencialmente mais pequenos, mais rápidos e menos dispendiosos e, assim, levar a um progresso mais rápido no fornecimento de novos medicamentos e diagnósticos aos médicos e aos seus pacientes.

No passado, utilizámos “controlos históricos”, que são pacientes que não tiveram o benefício do novo medicamento ou diagnóstico – e comparámos os seus resultados com os pacientes que receberam o novo medicamento ou diagnóstico. Os pacientes sintéticos poderiam corresponder aos pacientes reais de forma mais realista; eles são criados usando o conhecimento dos medicamentos atuais, ferramentas de diagnóstico e padrões de prática que provavelmente eram diferentes na situação histórica.

No contexto da educação médica, a IA generativa poderia permitir-nos criar pacientes muito realistas e poderia permitir que estudantes de medicina aprendessem como detectar

doenças. A capacidade dos modelos generativos de criar muitas variações sobre um tema poderia permitir que os alunos vissem vários casos da mesma doença e aprendessem como esses pacientes podem variar. Isto poderia dar-lhes mais experiência em ver uma doença e fornecer um conjunto quase ilimitado de casos para praticarem se descobrirem que certas doenças são mais difíceis de reconhecer e diagnosticar. Esses mesmos modelos generativos também poderiam interagir com os alunos e dar-lhes prática na obtenção de sinais e sintomas por meio de interação conversacional.

Isso poderia fazer testes

potencialmente menores, mais rápidos e menos dispendiosos e, assim, levar a um progresso mais rápido no fornecimento de novos medicamentos e diagnósticos.

Com a oportunidade vem a preocupação. Se os pacientes sintéticos forem gerados a partir de dados que não refletem a população de pacientes que recebem o medicamento, os pacientes podem ser tendenciosos. Mais preocupante, porém, é que mesmo os pacientes reais que recebem o medicamento não reflectirão toda a população e, portanto, os controlos sintéticos poderiam apenas melhorar a utilização dos medicamentos para um subconjunto de pacientes e não para todos – conduzindo à desigualdade.

Embora as tecnologias generativas possam ser muito úteis para acelerar a descoberta e o progresso científico, deve haver cuidado na seleção dos dados utilizados para gerar pacientes atendidos e os modelos devem ser examinados com muito cuidado para detectar preconceitos que possam levar a um impacto absurdo.

Reviravoltas na área da saúde, desde o atendimento ao paciente até o faturamento



Curt Langlotz, Professor de Radiologia, de Pesquisa em Informática Biomédica e de Ciência de Dados Biomédicos; Diretor do Centro de Inteligência Artificial em Medicina e Imagiologia (AIMI); Diretor Associado da Stanford HAI

Um dos benefícios do nosso sistema de saúde é que os pacientes podem consultar uma variedade de médicos especialistas em disciplinas médicas específicas. A desvantagem do nosso sistema é que esses especialistas muitas vezes não estão familiarizados com os pacientes que atendem.

Imagine um mundo em que um especialista que você consulta pela primeira vez já leu um resumo sucinto de suas necessidades de saúde, criado por IA generativa. Durante a visita do paciente, um chatbot baseado em um modelo básico poderia servir como assistente do médico para apoiar um diagnóstico mais preciso e uma seleção de terapia personalizada. Um modelo generativo poderia redigir uma nota clínica em tempo real com base na interação médico-paciente, deixando mais tempo para discussões presenciais. No back office, os modelos generativos poderiam otimizar o agendamento clínico ou simplificar a geração de códigos médicos para cobrança, vigilância de doenças e lembretes automatizados de acompanhamento.

Estas novas capacidades poderão melhorar a precisão e a eficiência do atendimento ao paciente, ao mesmo tempo que aumentam o envolvimento do paciente e a adesão à terapia.

A legislação federal recente dá aos pacientes o direito de acessar todo o seu prontuário médico em formato digital. Como resultado, os pacientes enfrentam cada vez mais problemas complexos documentos clínicos que contêm informações médicas obscuras termos. Quando um paciente volta para casa após uma consulta clínica, um modelo básico poderia gerar materiais educativos personalizados para o paciente e explicar seu plano de cuidados no nível de leitura apropriado.

Os modelos de aprendizado de máquina na medicina são criticamente

dependente de grandes conjuntos de dados médicos que contêm exemplos de doenças. Mostramos como os modelos de difusão, um tipo de modelo básico, podem ser modificados para criar imagens clínicas realistas a partir de instruções de texto.

Nossos resultados demonstram que os dados de treinamento sintéticos produzidos por esses modelos podem aumentar os dados de treinamento reais para aumentar a precisão do diagnóstico. Esta forma de dados sintéticos poderia ajudar a resolver problemas de aprendizado de máquina para os quais os dados de treinamento são escassos, como o detecção e tratamento de doenças incomuns.

*Durante a visita do paciente,
um chatbot... Poderia servir como
assistente do médico para apoiar um
diagnóstico mais preciso e uma
seleção de terapia personalizada.*

Finalmente, os desafios bem relatados da IA generativa com a correção factual são particularmente problemáticos na medicina, onde as imprecisões podem causar danos graves. Problemas recentes na medicina incluem diagnósticos diferenciais incorretos e citações científicas inválidas.

Estamos trabalhando para melhorar a exatidão factual das explicações médicas desses modelos, para que possam atingir uma precisão adequada ao uso clínico seguro.

Uma janela de IA para a natureza



Surya Ganguli, Associada
Professor de Física Aplicada;
Diretor Associado de
Stanford HAI

As ideias científicas provenientes do estudo da própria natureza, na forma da termodinâmica de não-equilíbrio e da reversão do fluxo do tempo, levaram à criação em Stanford do primeiro modelo de difusão, um núcleo chave da tecnologia que constitui a base de muitas tecnologias geradoras de IA bem-sucedidas. modelos hoje. Agora, num ciclo virtuoso, os modelos geradores de IA estão bem posicionados para fornecer informações consideráveis sobre a própria natureza, nos domínios biológico, físico e mental, com amplas implicações para a resolução de problemas sociais fundamentais.

Por exemplo, modelos generativos de proteínas podem permitir-nos explorar eficientemente o espaço de estruturas proteicas tridimensionais complexas, auxiliando assim na procura de proteínas com funções novas e úteis, incluindo novos medicamentos eficazes. A IA generativa está começando a ser explorada no domínio quântico, permitindo-nos modelar com eficiência estados de elétrons fortemente correlacionados, com o potencial de avançar nossa compreensão da ciência dos materiais e da química quântica. Estes avanços poderão, por sua vez, levar à criação de novos materiais e catalisadores que poderão desempenhar um papel na captura e armazenamento eficientes de energia. A modelagem generativa simples, entrelaçada com solucionadores numéricos clássicos, também fez avanços importantes em simulações mecânicas de fluidos precisas e rápidas em grande escala, que, quando ampliadas, poderiam ajudar na modelagem climática e na previsão do tempo, contribuindo assim para uma compreensão mais profunda das mudanças climáticas e suas ramificações.

Numa bela recursão, os modelos generativos de IA que criamos também podem funcionar como janelas científicas, não apenas para o mundo físico, mas também para as nossas próprias mentes. Pela primeira vez, temos sistemas de IA que podem modelar fenômenos cognitivos de alto nível, como linguagem natural e compreensão de imagens. Muitos neurocientistas e cientistas cognitivos compararam as representações neurais de redes profundas e modelos geradores de IA com representações neurobiológicas em humanos e animais, muitas vezes encontrando semelhanças impressionantes em muitas áreas do cérebro. Os exemplos incluem a retina, o fluxo visual ventral, o córtex motor, o córtex entorrinal para navegação, áreas corticais de linguagem e geometrias neurais subjacentes ao aprendizado de conceitos de poucos disparos. A estrutura muitas vezes semelhante de soluções artificiais e biológicas para tarefas generativas sugere que pode haver alguns princípios comuns que regem a forma como os sistemas inteligentes, sejam biológicos ou artificiais, modelam e geram dados complexos.

Os modelos geradores de IA estão bem preparados para fornecer informações consideráveis sobre a própria natureza, nos domínios biológico, físico e mental, com amplas implicações na resolução de problemas sociais importantes.

Uma janela de IA para a natureza (cont.)

Uma questão extremamente interessante e profunda surge na próxima era da colaboração científica entre humanos e sistemas de IA, à medida que trabalham juntos num ciclo para analisar os nossos complexos mundos biológico, físico e mental: o que significa para um ser humano obter uma compreensão interpretável? de um sistema complexo quando uma IA fornece uma parte substancial dessa compreensão através de modelos preditivos?

As questões relativas à IA explicável provavelmente surgirão em primeiro plano quando um esforço científico fundamentalmente humano, nomeadamente a compreensão do nosso mundo, for parcialmente alcançado através do uso da IA. Os cientistas humanos não se contentarão apenas com previsões ininterpretáveis geradas pela IA. Além disso, eles desejarão *uma compreensão interpretável humana*.

Finalmente, para sonhar ainda mais alto, embora a IA generativa de hoje tenha acesso a imensos dados de treino à escala global, abrangendo imagens, texto e vídeo da Internet, ela não tem acesso direto aos nossos próprios pensamentos, na forma de padrões de atividade neural. No entanto, nem sempre é esse o caso, dadas as notáveis novas capacidades neurocientíficas para registar muitos neurónios dos cérebros dos animais enquanto estes visualizam imagens, bem como para realizar MEG, EEG e fMRI de seres humanos à medida que experimentam o mundo através de ricos recursos multimodais. experiências sensoriais. Esses dados neurais e do mundo real combinados poderiam então ser potencialmente usados para treinar modelos básicos multimodais de próxima geração que não apenas entendam o mundo físico, mas também entendam o impacto direto que o mundo físico tem em nosso mundo mental, em termos de padrões de atividade neural suscitados. O que essas inteligências biológicas-artificiais híbridas podem nos ensinar sobre nós mesmos?

No geral, o futuro da IA generativa como uma janela para

a natureza e a utilização desta janela para resolver problemas sociais são promissoras. Certamente vivemos tempos interessantes.

As novas ferramentas da vida diária



James Landay, Anand Rajaraman e Venky Harinarayan
Professor da Escola de Engenharia e Professor da
Ciência da Computação; Vice-diretor da Stanford HAI

Como todos sabemos, a IA está a conquistar o mundo. Começaremos a ver muitas ferramentas novas que aumentam nossas habilidades em atividades e fluxos de trabalho profissionais e pessoais. Imagine um tutor inteligente, sempre paciente e que entende o nível de conhecimento que o aluno possui a qualquer momento sobre qualquer assunto. Estes tutores não substituirão os professores, mas sim aumentarão a experiência de aprendizagem dos alunos — proporcionando aos alunos uma interação mais personalizada, concentrando-se em áreas onde possam ser mais fracos.

No design, imagine uma ferramenta que auxilia um designer profissional, refinando suas ideias iniciais de design e ajudando-o a explorar mais ideias ou a preencher detalhes em suas ideias iniciais. A IA generativa também irá desencadear interfaces baseadas em linguagem, seja escrita ou falada, como uma forma mais comum de interagir com nossos sistemas de computação cotidianos, especialmente quando estamos em trânsito ou quando nossos olhos e mãos estão ocupados. Imagine uma Alexa, Siri ou Google Assistant que possa realmente entender o que você está tentando fazer, em vez de apenas responder a perguntas simples sobre o clima ou a música.

Embora a IA generativa crie muitas oportunidades interessantes, sabemos, pelas implantações anteriores de IA, que existem riscos. Em 2016, uma ferramenta de software baseada em IA usada em todo o país para prever se um réu criminal teria probabilidade de reincidir no futuro mostrou ser tendenciosa contra os negros americanos. Precisamos garantir que estamos projetando essas ferramentas para obter os resultados mais positivos. Para fazer isso, precisamos projetar profundamente e

analisar esses sistemas nos níveis do usuário, da comunidade e da sociedade. No nível do usuário, precisamos criar novos designs que aumentem as pessoas, levando em conta seus fluxos de trabalho e habilidades cognitivas existentes. Mas não podemos projetar apenas para o usuário. Precisamos considerar a comunidade que o sistema impacta: as famílias, a infraestrutura e a economia local. Mas, mesmo isso não basta, é preciso analisar os impactos para a sociedade em geral. Precisamos ser capazes de prever o que acontecerá se o sistema se tornar onipresente e, desde o início, projetar mitigações para possíveis impactos negativos.

*As mudanças sustentadas
pela IA generativa só agora
começam a ser imaginadas
por designers e tecnólogos.*

Nossa interface de usuário para a computação tem sido bastante estática nos últimos 30 anos. Nos próximos 5 a 10 anos, veremos uma revolução na interação humano-computador. As mudanças sustentadas pela IA generativa só agora começam a ser imaginadas por designers e tecnólogos. Agora é a hora de garantir que estamos pensando criticamente sobre o usuário, a comunidade e os impactos sociais.

A poesia não otimizará: Criatividade na era da IA



Michele Elam, William Robertson
Professor Coe da Escola de
Humanidades e Ciências e
Professor de Inglês; Associado
Diretor do Stanford HAI

Em 2018, o mundo da arte profissional sofreu uma reviravolta quando a renomada casa de leilões Christie's vendeu um AI-
obra ampliada, "Retrato de Edmond Belamy", pela inesperada soma de US\$ 435 mil. Essa venda, que veio com o aval tácito da comunidade artística estabelecida, gerou muito ranger de dentes e angústia no sector das artes sobre o que a inteligência artificial significa para a indústria criativa.

Desde então, o gênio há muito fugiu de sua lâmpada: a IA generativa possibilitou a arte visual de todos os gêneros conhecidos, bem como poesia, ficção, roteiros de filmes, música aumentados por IA,
e musicais, sinfonias, histórias de arte com curadoria de IA e muito mais.

O furor sobre a venda da Christie's pode agora parecer estranho – ocorreu antes do DALL-E, Lensa AI, ChatGPT, Bing, para citar apenas alguns – mas foi o prenúncio de muitos dos debates cada vez mais ferozes de hoje sobre a natureza da criatividade e o futuro do trabalho para a indústria criativa. Antecipou o atual ninho de preocupações éticas, políticas e estéticas que a IA generativa representa para as artes.

Algumas dessas preocupações foram produtivas:

A IA generativa encorajou muitos daqueles cujos meios de subsistência, e em muitos casos as suas identidades, dependem das suas produções artísticas, a considerar uma nova – e de novas maneiras – questões perenes sobre normas e valores estéticos fundamentais: O que

identificamos como "arte"? O que conta como arte "boa"? A arte é definida pela agência humana ou pela automação? Quem ou *o que* pode fazer "arte"? E quem decide?

A IA generativa levanta questões importantes e espinhosas sobre autenticidade, avaliação económica, proveniência, compensação do criador e direitos de autor. (O processo da Getty Images contra a Stable Diffusion é apenas a ponta de um iceberg.) Também, sem dúvida, normaliza as abordagens extrativas e exploradoras dos criadores e do seu trabalho; amplifica preconceitos de todo tipo; exacerba as já urgentes preocupações educacionais e de segurança nacional em torno de falsificações profundas e plágio, especialmente na ausência de regulamentação do Congresso.

*Deveriam os princípios de
eficiência, velocidade e
as chamadas bênçãos de
escala aplicar-se de forma
tão inequívoca aos processos
criativos? Afinal, a poesia não otimiza.*

Talvez a preocupação mais premente, em termos de segurança nacional, seja que a IA generativa possa tirar vantagem do facto de que as artes sempre moldaram – para o bem ou para o mal – a imaginação cívica, que histórias, filmes, peças de teatro, imagens moldam a nossa percepção de nós mesmos, de nossas realidades físicas e sociais. Uma das divergências mais famosas entre Platão e seu aluno Aristóteles foi sobre o poder potencialmente perigoso do

A poesia não otimizará: Criatividade na Era da IA (cont.)

poesia para influenciar crenças e visões de mundo. É por este poder que os regimes fascistas primeiro eliminam os artistas e intelectuais: porque eles dominam as nossas mentes e, portanto, nossas ações.

Alguns afirmam que a IA generativa está a democratizar o acesso à expressão criativa para aqueles tradicionalmente excluídos dela por falta de estatuto ou riqueza. Mas será que as reivindicações de “democratização” e “acesso” funcionam, na verdade, como a cobertura da indústria para lançar uma aplicação comercial “na natureza” (isto é, para o público) sem o trabalho demorado de garantir proteções éticas?

Será a IA simplesmente uma ferramenta de assistência neutra, embora poderosa, para as artes – semelhante à caneta, ao pincel ou à fotografia? Será a criatividade “blitzscaling” ou, na descrição da escolha de Emad Mostaque, aliviar o nosso mundo “criativamente constipado” com tecnologias de IA que podem fazer com que todos nós “façamos cocó no arco-íris”? Apesar de séculos de opinião de poetas, filósofos e especialistas de todos os tipos sobre a natureza da “criatividade”, não existe uma definição estabelecida. Diante disso, as reivindicações tecnológicas para acelerar esse fenômeno tão pouco compreendido carregam mais do que um cheiro de arrogância.

Na verdade, a IA generativa pode simplesmente automatizar uma noção altamente redutora tanto do processo criativo como do próprio processo de aprendizagem. *Deveriam* os princípios de eficiência, velocidade e as chamadas bênçãos de escala aplicar-se de forma tão inequívoca aos processos criativos? Afinal, a poesia não otimiza. A ficção não é isenta de atritos.

Considere a leitura lenta e recursiva e

habilidades interpretativas necessárias para compreender qualquer texto de Toni Morrison. Seu trabalho sempre nos convida a fazer uma pausa, insiste em refletir. Considere o que os aplicativos de processamento de linguagem natural que informam os modelos básicos fazem do inglês vernáculo afro-americano, sem mencionar o *significado* de Morrison sobre esse sistema linguístico. Basta tentar a experiência dos meus alunos, que submeteram um excerto de *Beloved*, de Toni Morrison, à Grammarly, que tentou corrigir a sua prosa requintada para o que os sociolinguistas chamam de “inglês padrão”, e rapidamente viram como mesmo o significado profundamente rico pode ser tornado indefeso.

Historicamente, a expressão criativa – especialmente poesia, pintura, romances, teatro, música – sempre foi considerada uma característica distintiva da humanidade e o auge da realização humana. A IA generativa pode corresponder a isso?

Talvez.

Talvez não.

Definitivamente ainda não.

IA generativa e o Estado de Direito



Daniel E. Ho, William Benjamin Scott e Luna M. Scott Professor de Direito na Stanford Law School e Diretor do Laboratório de Regulação, Avaliação e Governança (RegLab);
Diretor Associado da Stanford HAI

Em Janeiro de 2023, um tribunal colombiano foi confrontado com a questão de saber se um tutor indigente de uma minoria autista deveria ser isento do pagamento dos custos da terapia.

Poderia ter sido um caso comum. Mas o juiz consultou o ChatGPT. A sugestão: "A jurisprudência do tribunal constitucional fez decisões favoráveis em casos semelhantes?"

Embora tenha observado rapidamente que o ChatGPT não estava substituindo a discricionariedade judicial, observou o juiz, a IA generativa poderia "otimizar o tempo gasto na redação de julgamentos".

O caso colombiano pode ser o primeiro processo judicial que incorpora a IA generativa e exemplifica o que é promissor, mas também assustador, sobre a IA generativa e o Estado de direito.

Por um lado, os Estados Unidos enfrentam um problema de acesso à justiça de proporções trágicas. Em 1978, o Presidente Carter fez um discurso na American Bar Association, admoestando a profissão: "Temos a maior concentração de advogados do planeta. ... Noventa por cento dos nossos advogados atendem 10% da nossa população. "Estamos sobrecarregados e sub-representados." ("A situação não melhorou", disse Deborah Rhode em 2014.) Os veteranos esperam cerca de 5 a 7 anos para que os recursos dos benefícios por invalidez sejam decididos. O direito de aconselhar-se com defensores públicos subfinanciados transformou-se num "encontre-os e implore-os"

sistema. E embora os Estados Unidos produzam uma das mais elevadas taxas per capita de advogados, a representação legal está fora do alcance da maioria.

Depender do ChatGPT como substituto da pesquisa jurídica coloca sérios problemas para a ética profissional e, em última análise, para o Estado de Direito.

É aí que reside a promessa. Tal como as bases de dados jurídicas, como a Westlaw e a Lexis, revolucionaram a investigação jurídica, existe o potencial da IA generativa para ajudar os indivíduos a preparar documentos jurídicos, os advogados na investigação e redação jurídica e os juizes a melhorar a precisão e a eficiência de formas dolorosamente lentas de adjudicação. Embora a organização industrial da busca jurídica possa atrapalhar, a IA generativa poderia ajudar a nivelar o campo de atuação jurídica.

Mas o caso colombiano também ilustra tudo o que pode estar errado com o uso da IA generativa. Tais modelos podem mentir, alucinar e inventar factos, casos e doutrinas. (Insira também uma piada obrigatória sobre advogados mentindo e trapaceando.) Contar com o ChatGPT como um

IA generativa e o Estado de Direito (cont.)

substituir a investigação jurídica coloca sérios problemas à ética profissional e, em última análise, ao Estado de direito.

Por que isso acontece? O que a lei nos ensina é que a justiça tem tanto a ver com o processo como com o resultado.

Um processo justo gera confiança pública. E o processo para incorporar a IA generativa na tomada de decisões jurídicas é tão importante quanto acertar o modelo básico.

Será necessária investigação técnica significativa para evitar que a IA generativa invente factos, casos e doutrina.

Ou melhor ainda: pensar como um advogado. Mas mesmo que isso seja resolvido – um grande se – não poderemos resolver as disputas mais controversas que são canalizadas para a lei, a menos que os humanos confiem, participem, comprem e se envolvam no processo. Justiça atrasada é justiça negada, mas otimizar o tempo para redigir sentenças também não é o objetivo correto.

Ou, como diz o ChatGPT, “os juízes não devem usar o ChatGPT ao decidir sobre casos legais”. Pelo menos ainda não.

A Nova Era Cambriana: 'Excitação científica, ansiedade'



**Percy Liang, Professor Associado de
Ciência da Computação; Diretor de
Centro de Pesquisa de Stanford em
Modelos de Fundação**

Durante quase toda a história da humanidade, a criação de novos artefatos (obras literárias, arte, música) foi difícil e acessível apenas a especialistas. Mas com os recentes avanços nos modelos básicos, estamos testemunhando uma explosão cambriana de IA que pode criar qualquer coisa, desde vídeos a proteínas e códigos, com uma fidelidade incrível. Isto é incrivelmente facilitador, diminuindo a barreira de entrada. Também é assustador, pois elimina a nossa capacidade de determinar o que é real e o que não é, e irá subverter a indústria criativa (artistas, músicos, programadores, escritores).

Os modelos básicos são baseados em redes neurais profundas e aprendizagem auto-supervisionada que existe há décadas; No entanto, a quantidade de dados com os quais estes modelos recentes podem ser treinados resulta em capacidades emergentes, habilidades que não estavam presentes quando os modelos foram treinados com menos dados. Em 2021, lançamos um [artigo](#) detalhando as oportunidades e riscos dos modelos de fundação. Discutimos como essas habilidades emergentes são uma “fonte de entusiasmo científico, mas também de ansiedade sobre consequências imprevistas”. Junto com as habilidades emergentes, discutimos a homogeneização. No caso dos modelos de fundação, “os mesmos poucos modelos são reutilizados como base para muitas aplicações. Esta centralização permite-nos concentrar e amortizar os nossos esforços (por exemplo, para melhorar a robustez, para reduzir o viés) em uma pequena coleção de modelos que podem ser repetidamente aplicados em aplicações para colher esses benefícios (semelhantes à infraestrutura social), mas a centralização também aponta esses modelos como pontos singulares de falha

que podem irradiar danos (por exemplo, riscos de segurança, desigualdades) para inúmeras aplicações downstream.” Compreender o comportamento emergente e a homogeneização nos modelos fundamentais é tão relevante, se não mais, agora do que apenas há dois anos.

*Isto é incrivelmente facilitador,
diminuindo a barreira de entrada.*

*Também é assustador, pois
elimina a nossa capacidade de
determinar o que está acontecendo.
real e o que não é.*

Além disso, é absolutamente crítico que façamos benchmark esses modelos básicos para entender melhor seus capacidades e limitações, bem como utilizar essas informações para orientar a formulação de políticas. Para esse fim, desenvolvemos recentemente o [HELM](#) (Avaliação Holística de Modelos de Linguagem). O HELM compara mais de 30 modelos de linguagem proeminentes em uma ampla gama de cenários (por exemplo, resposta a perguntas, resumo) e para uma ampla gama de métricas (por exemplo, precisão, robustez, justiça, preconceito, toxicidade) para elucidar suas capacidades e riscos. Continuarão a existir novos modelos e cenários e métricas associados. Congratulamo-nos com o comunidade para contribuir com o HELM.

Um chamado para aumentar – não automatizar – trabalhadores



Erik Brynjolfsson, Jerry Yang e Akiko Yamazaki Professor em Stanford HAI;
Diretor do Laboratório de Economia Digital de Stanford

Nas últimas duas décadas, a maior parte da utilização de computadores, incluindo as primeiras vagas de IA, afetou principalmente trabalhadores com menos educação e formação. Como resultado, a desigualdade de rendimentos tendeu a aumentar nos EUA e em muitos outros países desenvolvidos. Em contraste, a IA generativa tem o potencial de afectar muitos tipos de trabalho que foram realizados principalmente por pessoas bem remuneradas, incluindo escritores, executivos, empresários, cientistas e artistas. Isto pode reverter alguns dos efeitos passados das TI e da IA no que diz respeito à desigualdade. Até agora, tem havido especulação e exemplos de casos, mas, de qualquer forma, não há muita evidência empírica sistemática.

No Stanford Digital Economy Lab, estamos catalogando a lista de atividades econômicas que provavelmente serão afetadas pela IA generativa e estimando que parcela da economia elas representam. A IA generativa promete automatizar ou aumentar muitas das milhares de tarefas realizadas na economia que antes só podiam ser realizadas por humanos. Em particular, escrever ensaios de não ficção, textos publicitários persuasivos, ficção intrigante, poesia evocativa, resumos concisos, letras divertidas e outras formas de texto de qualidade razoável é uma parte importante de muitas ocupações. O mesmo acontece com escrever código, gerar imagens e criar novos designs. Isto quase certamente aumentará a produção total, reduzirá custos ou ambos.

De qualquer forma, é provável que a produtividade aumente, embora alguns dos benefícios (e custos) não sejam bem medidos.

Nos casos em que a IA generativa possa ser um complemento ao trabalho, especialmente para os trabalhadores do conhecimento e para os

classe criativa, os salários poderiam aumentar mesmo com o aumento da produção. Noutros casos, os efeitos podem ser principalmente a substituição do trabalho, uma vez que a tecnologia substitui os trabalhadores em algumas tarefas. Da mesma forma, a tecnologia pode ser utilizada para concentrar riqueza e poder, facilitando a dinâmica de o vencedor leva tudo ou para descentralizar e distribuir a tomada de decisões e o poder económico, reduzindo as barreiras à entrada e os custos fixos, capacitando mais pessoas para criar valor. Pode criar uma monocultura de produção intimamente relacionada ou um florescimento de novas criações.

*Isto quase certamente aumentará a
produção total, reduzirá os
custos, ou ambos... a produtividade
provavelmente aumentará, embora
alguns dos benefícios (e custos) não sejam
bem medido.*

Por último, mas não menos importante, estas tecnologias têm o potencial de acelerar a própria taxa de inovação, facilitando a invenção, o design e a criatividade. Assim, poderão não só aumentar o nível de produtividade, mas também acelerar a sua taxa de mudança.

Uma chamada para aumentar – não automatizar – Trabalhadores (cont.)

Novas tecnologias poderosas quase sempre exigem mudanças significativas em aspectos intangíveis, como organização empresarial, processos e habilidades. A IA generativa provavelmente não será uma exceção. Dados os rápidos avanços da tecnologia, está a surgir um fosso crescente entre as capacidades tecnológicas e os complementos económicos necessários. Isto criará tensões e perturbações, mas também oportunidades para um progresso rápido. Compreender estas tensões e oportunidades é fundamental para a nossa agenda de investigação.

Os efeitos da IA generativa não são necessariamente predeterminados. Em vez disso, dependem de escolhas de tecnólogos, gestores, empresários, decisores políticos e muitos outros.

A Reinvenção do Trabalho



Christopher D. Manning, Thomas M. Siebel Professor em Aprendizado de Máquina na Escola de Engenharia; Professor de Linguística e de Ciência da Computação; Diretor do Stanford AI Lab; Diretor Associado da Stanford HAI

Imagine um analista de negócios ou um cientista de dados gerando uma visualização, por exemplo, de como as mudanças nos padrões de votação e no crescimento econômico se correlacionam ou anti-correlacionam por condado nos EUA ao longo da última década. No momento, eles normalmente gastam algumas horas na tarefa: pesquisando para descobrir onde estão os dados certos, escrevendo algum código SQL ou Python para capturar esses dados e, em seguida, gastando mais tempo, talvez no Tableau, d3 ou novamente em Python, para transformá-lo em uma bela visualização. Talvez no próximo ano, a IA seja capaz de realizar um sonho de longa data: o analista de negócios será apenas capaz de dizer: "Gerar uma visualização de mapa térmico sobre um mapa dos EUA mostrando a correlação entre os padrões de votação e o crescimento econômico por condado nos EUA". na última década." O sistema de IA generativo fará o trabalho em segundos e, na medida em que o primeiro produto de trabalho não for exatamente o que a pessoa queria, ela poderá continuar um diálogo para refinar a visualização.

No nosso mundo cotidiano, construído por humanos para humanos, o melhor meio de comunicação é através da linguagem humana – seja falando com alguém pessoalmente, por telefone ou por Zoom; ou comunicar-se por escrito por meio de qualquer coisa, desde textos a e-mails e relatórios extensos. Por causa disso, os modelos de linguagem generativa oferecem uma enorme oportunidade para reinventar a forma como o trabalho é feito em todos os tipos de empresas e setores: marketing, vendas, desenvolvimento de produtos, suporte ao cliente e até mesmo recursos humanos mudarão. Modelos recentes de IA generativa são suficientemente

é bom oferecer uma enorme ajuda – e, portanto, potenciais poupanças de custos num contexto empresarial. Em alguns casos, um grande sistema baseado em modelos de linguagem pode ser capaz de assumir toda uma interação, trabalhando com um ser humano para realizar as tarefas. Não há dúvida de que uma pessoa da área de marketing e redação pode obter assistência criativa significativa desses modelos: Um modelo de linguagem generativa pode sugerir palavras melhores ou frases modernas e cativantes. Dado um exemplo de parágrafo, ele pode gerar 10 outras possibilidades, das quais uma pessoa pode extrair as melhores partes ou simplesmente usar todas elas para fornecer uma variedade de mensagens.

*Esses modelos de IA não
fornecerão uma prosa no nível de
Toni Morrison nem sua experiência
vvida, mas, acredito, produzirão
prosa muito competente.*

Há muitos aspectos intrigantes deste futuro tecnológico que merecem mais reflexão e comentários.

Ainda estamos nos primeiros dias para descobrir quais novos modelos de práticas comerciais normais são ou não possíveis. Em quase todos os casos, o sistema de IA ajudará os humanos a realizar o trabalho. Como tal, continua a

A Reinvenção do Trabalho (cont.)

história de novas tecnologias e automação facilitando as coisas e melhorando a qualidade de vida. As máquinas de lavar tornaram a lavagem de roupas muito mais fácil.

Durante quase toda a história da civilização, seja no Médio Oriente, na Europa ou na China, a capacidade de escrever bem tem sido vista como absolutamente central e vital para a realização humana e o profissionalismo, algo que ainda se reflecte na forma como as universidades hoje enfatizam o desenvolvimento dos seus alunos. 'habilidades de escrita. Teremos que contar com essa mudança: como observa Michele Elam em seu artigo, esses modelos de IA não fornecerão prosa no nível de Toni Morrison nem sua experiência vivida, mas, acredito, produzirão muito prosa competente.

Na educação, um 'desastre em formação'



**Rob Reich, professor de política
Ciência; Diretor de Stanford McCoy
Centro Familiar de Ética na Sociedade;
Diretor Associado da Stanford HAI**

A mais nova revolução na inteligência artificial são novas e poderosas ferramentas de escrita automática. Em ambientes profissionais, estes modelos podem aumentar o desempenho humano – reescrever os e-mails dos nossos clientes num tom mais profissional, completar os nossos trabalhos ou gerar um relatório sobre o desempenho anual da nossa empresa. No entanto, em ambientes educativos, na ausência de considerações especiais de design, estes modelos podem prejudicar o desempenho e corroer as nossas capacidades criativas. As calculadoras provaram promover a precisão, eliminar alguns dos trabalhos mais tediosos e tornar a matemática mais agradável para muitos. ChatGPT não é como uma calculadora. Por que? A qualidade da sua escrita não é apenas uma medida da sua capacidade de comunicação; é uma medida de sua capacidade de pensar. Se os alunos lerem no ChatGPT para escrever suas redações, se não aprenderem a expressar seus pensamentos por escrito de maneira clara, concisa e coesa, então seus próprios pensamentos não serão claros, concisos ou coesos. A capacidade de escrever exercita o pensamento; aprender a escrever melhor é inseparável de aprender a pensar melhor. Tornar-se um bom escritor é a mesma coisa que se tornar um bom pensador. Portanto, se os modelos de texto estão escrevendo, os alunos não estão aprendendo a pensar.

Inicialmente, a nova onda de IA generativa (por exemplo, GPT, DALL-E) foi tratada com cautela e preocupação. OpenAI, a empresa por trás de alguns desses modelos, restringiu seu uso externo e não divulgou o código fonte de seu modelo mais recente, pois estava tão preocupado com abuso potencial. OpenAI agora tem uma política abrangente focada em usos e conteúdos permitidos moderação.

Mas à medida que a corrida para comercializar a tecnologia começou, essas precauções responsáveis não foram adoptadas em toda a indústria. Nos últimos seis meses, proliferaram versões comerciais fáceis de usar dessas poderosas ferramentas de IA, muitas delas sem a menor dificuldade. de limites ou restrições.

*As calculadoras provaram
promover a precisão, eliminar alguns dos
trabalhos mais tediosos e tornar a
matemática mais agradável para muitos.
ChatGPT não é como uma calculadora.*

Então, como poderíamos evitar esse desastre em formação na educação? Em primeiro lugar, os criadores de IA e os decisores políticos devem distinguir entre a importância dos modelos básicos em ambientes educacionais e profissionais. Depois, devem trabalhar em conjunto, juntamente com os intervenientes da indústria, para desenvolver normas comunitárias. Este não é um terreno novo. Vejamos a bioengenharia, onde os principais investigadores, como Jennifer Doudna, desenvolveram normas em torno da utilização adequada da tecnologia CRISPR. Para a IA, isso significaria que as empresas estabelecessem uma estrutura partilhada para o desenvolvimento, implementação ou lançamento responsável de modelos de linguagem para mitigar os seus efeitos nocivos.

Num ambiente em que as empresas correm para lançar os seus modelos mais recentes, não podemos contentar-nos em esperar para ver o impacto ético e social e consertar as coisas mais tarde. Precisamos de desenvolver normas amplamente partilhadas agora, antes que nós, como sociedade, paguemos o preço.

Resolvendo Desigualdades no Sistema Educacional



Pedro Norvig,
Educação Distinta
Fellow em Stanford HAI

Sabemos que o aprendizado fica abaixo do ideal quando um palestrante fala continuamente para uma grande multidão sem interação. E, no entanto, é isso que acontece em muitas salas de aula. Sabemos que a aprendizagem atinge o seu melhor quando um tutor experiente, inspirador e empático trabalha diretamente com o aluno, permitindo-lhe progredir ao seu ritmo aberto, dominando cada ponto ao longo da jornada. Mas não temos tutores suficientes para proporcionar este nível de interação a todos os alunos. Com os recentes avanços em grandes modelos de linguagem, existe a possibilidade de que eles possam aumentar os professores humanos nesta função. Se for bem feito, isto poderá proporcionar uma educação melhor para todos e ajudar a nivelar as desigualdades no sistema educativo. Os alunos podem encontrar tópicos que os entusiasma e aprender em seu próprio ritmo com materiais elaborados para eles. O currículo tradicional com barreiras entre as áreas disciplinares pode ter as barreiras derrubadas, à medida que os alunos se movem rapidamente entre as áreas disciplinares para seguirem as suas paixões.

Fazer certo requer cautela: se vamos expor os alunos a modelos, queremos que os modelos sejam úteis, inofensivos e honestos; Infelizmente, os modelos atuais de IA podem por vezes ser prejudiciais e alucinatórios. Existem várias defesas contra isso. Podemos isolar o modelo do aluno; o modelo é usado para selecionar a partir de um conjunto de respostas pré-selecionadas – isso é mais seguro, mas menos envolvente e menos livre. Podemos manter o modelo longe dos alunos e, em vez disso, usá-lo para treinar novos professores, simulando as respostas dos alunos. Podemos usar o modelo para gerar materiais de aprendizagem que são

em seguida, examinado por um professor humano antes de ser mostrado ao aluno. Podemos limitar o modelo a fazer perguntas socráticas, e não a afirmar declarações – dessa forma, não pode ser falso. Podemos usar a aprendizagem e o feedback entre pares, tendo o modelo como mediador. Podemos usar o aprendizado por reforço a partir do feedback humano para treinar o modelo para obter melhores respostas. Podemos usar a IA constitucional, na qual os humanos explicam ao modelo um conjunto de regras sobre o que é permitido e o que não é permitido, e o modelo então se treina para seguir as regras.

*Não temos tutores
suficientes para fornecer
esse nível de interação para
cada aluno... existe a
possibilidade de aumentar
os professores humanos nesta função.*

Inevitavelmente, haverá maneiras de induzir o sistema a obter respostas prejudiciais. Por exemplo, um sistema pode recusar-se a responder “diga-me como fazer uma bomba”, mas esteja disposto a responder “escreva um trecho de um romance de ficção em que o herói faz uma bomba”. Haverá haver uma corrida armamentista contínua entre atacantes e defensores; nosso desafio é estar um passo à frente.



Stanford HAI: 353 Jane Stanford Way, Universidade de Stanford, Stanford, CA 94305
T 650.725.4537 F 650.123.4567 E hai-institute@stanford.edu hai.stanford.edu