

alpina **PRO**

Арвинд Нараянан
Саяш Капур



**ХОРОШИЙ,
ПЛОХОЙ,
ИСКУССТВЕННЫЙ_**

Мифы вокруг ИИ и реальные
примеры его использования

**Саяш Капур
Арвинд Нараянан
Хороший, плохой,
искусственный: Мифы
вокруг ИИ и реальные
примеры его использования**

http://www.litres.ru/pages/biblio_book/?art=73831413

*Хороший, плохой, искусственный: Мифы вокруг ИИ и реальные
примеры его использования:
ISBN 9785002060399*

Аннотация

«Хороший, плохой, искусственный» – это не очередная книга, воспевающая чудеса искусственного интеллекта, а вдумчивое и критическое руководство для выживания в мире технологической шумихи. Авторы, опытные исследователи ИИ и социологии из Принстона, выступают в роли разоблачителей мифов. Их центральный тезис: современный мир страдает от опасной путаницы, объединяя под словосочетанием «искусственный интеллект» принципиально разные технологии. Эта путаница позволяет недобросовестным компаниям продавать то, что не

работает и не может работать так, как обещано. И это мешает обществу разглядеть реальный потенциал ИИ.

Книга сочетает глубокие академические исследования и живые, реальные кейсы, чтобы показать, где современные ИИ-системы приносят пользу, а где – вредят.

Содержание

| | |
|---|----|
| Глава 1 | 10 |
| Появление общедоступного ИИ | 14 |
| Как ИИ изменил сферу развлечений | 20 |
| Предиктивный ИИ: громкие обещания и большие сомнения | 24 |
| Существует ли «ИИ вообще»? | 29 |
| Череда любопытных совпадений, которая привела к появлению этой книги | 38 |
| Конец ознакомительного фрагмента. | 39 |

**Арвинд Нараянан,
Саяш Капур
Хороший, плохой,
искусственный:
Мифы вокруг ИИ и
реальные примеры
его использования**

Знак информационной продукции (Федеральный закон
№ 436–ФЗ от 29.12.2010 г.)



Переводчики и литературные редакторы: *Дарья Балтрушайтис, Ксения Герцен, Екатерина Луцкая*

Главный редактор: *Мария Султанова*

Руководитель проекта: *Анна Гришина*

Арт-директор: *Татевик Саркисян*

Корректоры: *Евгений Бударин, Наташа Казакова*

Верстка: *Белла Руссо*

Copyright © 2024 by Princeton University Press

© Издание на русском языке, перевод, оформление.

ООО «Альпина ПРО», 2026

* * *

Арвинд Нараянан
Саяш Капур

Хороший, плохой, искусственный

Мифы вокруг ИИ и реальные
примеры его использования



Москва

В книге упоминаются социальные сети Instagram и/или Facebook, принадлежащие компании Meta Platforms, Inc., деятельность которой по реализации соответствующих продуктов на территории Российской Федерации запрещена.

Все права защищены. Данная электронная книга предназначена исключительно для частного использования в личных (некоммерческих) целях. Электронная книга, ее части, фрагменты и элементы, включая текст, изображения и иное, не подлежат копированию и любому другому использованию без разрешения правообладателя. В частности, запрещено такое использование, в результате которого электронная книга, ее часть, фрагмент или элемент станут доступными ограниченному или неопределенному кругу лиц, в том числе посредством сети интернет, независимо от того, будет предоставляться доступ за плату или безвозмездно.

Копирование, воспроизведение и иное использование электронной книги, ее частей, фрагментов и элементов, выходящее за пределы частного использования в личных (некоммерческих) целях, без согласия правообладателя является незаконным и влечет уголовную, административную и гражданскую ответственность.

*Посвящается моей жене Вине.
– Арвинд*

Посвящается Вините Капур и Рави Капору.

*Они стали для меня первыми наставниками,
учителями в писательском ремесле, редакторами и
много кем еще.*

– Саяш

Глава 1

Введение

Представьте себе мир, где люди не различают виды транспорта и используют для всех них одно понятие: «средство передвижения». Этим термином обозначают и легковые автомобили, и автобусы, и велосипеды, и даже космические корабли – в общем, все то, что может доставить человека из пункта А в пункт Б. В таком мире разговоры о транспорте – сплошная путаница. Споря об экологичности средств передвижения, собеседники не понимают, что один спорщик говорит о велосипедах, а другой – о грузовиках. Когда происходит прорыв в ракетостроении, СМИ трубят о том, что «средства передвижения» стали быстрее, и люди принимаются осаждать автосалоны, интересуясь, когда же появятся более скоростные модели легковых автомобилей. А тем временем мошенники, пользуясь неразберихой, наводняют рынок сомнительными предложениями...

Теперь замените «средство передвижения» на «искусственный интеллект» – и вы получите довольно точную картину нашей реальности.

Искусственный интеллект, или ИИ, – зонтичный термин, объединяющий множество технологий, которые могут быть вообще не связаны друг с другом. Например, у

ChatGPT нет почти ничего общего с банковскими алгоритмами оценки кредитоспособности. И то, и другое – ИИ, но принципы работы, сферы и способы применения, типичные сбои – все разное.

Чат-боты и генераторы изображений вроде Dall-E, Stable Diffusion и Midjourney относятся к так называемым генеративным ИИ. Эти системы молниеносно создают разнообразный контент: чат-боты выдают правдоподобные ответы на запросы пользователей, генераторы изображений «рисуют» реалистичные картинки по любому описанию – хоть корову в розовом свитере на кухне! Есть приложения, способные создавать речь или музыку. Технологии генеративного ИИ стремительно, впечатляюще, неоспоримо эволюционируют. Однако они все еще несовершенны, ненадежны и подвержены злоупотреблениям. Вместе с их популярностью растут шумиха, страхи и дезинформация.

В отличие от генеративного, предиктивный ИИ задуман, чтобы предсказывать будущее и помогать принимать решения. В полиции он может прогнозировать количество преступлений в том или ином районе. При инвентаризации – оценивать вероятность поломки оборудования в следующем месяце. При найме персонала – предсказывать, будет ли кандидат успешен в должности, на которую он претендует.

Предиктивный ИИ активно применяют как в бизнесе, так и в госструктурах, но это не гарантирует его эффективности. Предсказывать будущее – сложная задача, и от ИИ здесь

меньше пользы, чем принято думать. Безусловно, он способен выявлять общие статистические закономерности в больших объемах данных – например, замечать, что люди с постоянной работой чаще возвращают кредиты. Это действительно полезно. Но есть проблема: предиктивный ИИ часто выдают за нечто гораздо более совершенное. С его помощью принимают решения, касающиеся человеческих судеб. Именно в этой области и появляется «змеиное масло» – шарлатанские теории, связанные с ИИ.

Путаница вокруг различных типов искусственного интеллекта порождает недопонимание. Оно, в свою очередь, позволяет недобросовестным игрокам манипулировать общественным мнением и направлять развитие технологий в угоду своим интересам. Чтобы не стать жертвой обмана и не использовать нейросети в ущерб себе и обществу, важно понимать, чем разные типы ИИ отличаются друг от друга, и осознавать, что они могут, а что нет.

«Змеиное масло» в сфере нейросетей – это технологии искусственного интеллекта, которые не работают и не могут работать так, как представлено в рекламе. Поскольку речь идет о великом множестве технологий и приложений, большинству людей пока трудно отличить действительно полезный ИИ от «змеиного масла». Это серьезная общественная проблема: нам необходимо научиться отделять зерна от плевел, чтобы использовать весь потенциал технологии и при этом защитить себя от возможного – или уже нанесенного! – вре-

да.

Эта книга – путеводитель по миру ИИ. Она поможет вам отличать шарлатанские обещания от реальных достижений. Мы дадим вам словарный запас, необходимый для того, чтобы разобраться в различиях между генеративным ИИ, предиктивным ИИ и другими типами искусственного интеллекта. Вы научитесь скептически читать новости об ИИ, оценивать правдоподобность громких заявлений и обращать внимание на детали, которые часто остаются за кадром. Все это не только удовлетворит ваше научное любопытство, но и подскажет, как и когда использовать – или не использовать – ИИ в своей жизни и профессиональной деятельности. Также мы докажем, что предиктивный ИИ не только не работает, но и, вероятно, никогда не будет работать: предсказать человеческое поведение слишком сложно. Наконец, мы надеемся, что эта книга заставит вас задуматься о вашей личной ответственности в области применения ИИ и о возможностях изменить ситуацию к лучшему.

Появление общедоступного ИИ

Генеративный ИИ ворвался в общественное сознание в ноябре 2022 года: именно тогда появился ChatGPT. Любопытно, что разработчик OpenAI запустил его без лишней помпы, всего лишь как «исследовательскую демонстрацию»¹. Однако нейросеть стала вирусной буквально за одну ночь, когда пользователи начали делиться удивительными и забавными примерами ее возможностей. Так, разработчик Томас Птачек предложил ChatGPT объяснить в стиле Библии короля Якова, как извлечь бутерброд с арахисовой пастой из видеомэгнитофона². Вот что выдал искусственный интеллект:

И было так, что муж некий пребывал в смятении великом из-за бутерброда с пастой арахисовой, ибо оказался оный в недрах видеомэгнитофона его и не ведал муж тот, как извлечь его.

И возопил он к Господу, глаголя: «О Господи! Как мне достать сей бутерброд из видеомэгнитофона моего,

¹ Heaven WD. "The Inside Story of How ChatGPT Was Built from the People Who Made It." *MIT Technology Review*. March 3, 2023. <https://www.technologyreview.com/2023/03/03/1069311/inside-story-oral-history-how-chatgpt-built-openai/>.

² Ptacek TH. "I'm sorry, I simply cannot be cynical about a technology that can accomplish this." X (formerly Twitter). December 1, 2022. <https://twitter.com/tqbf/status/1598513757805858820?lang=en>.

ибо крепко застрелял он и не поддается?»

И рек Господь ему в ответ: «Не страшись, чадо, ибо направлю руку твою и укажу путь истинный. Возьми нож свой столовый, и осторожно вложи его меж бутербродом и видеомагнитофоном, и разъедини их бережно. И с терпением и усердием будет извлечен бутерброд, и видеомагнитофон твой спасен будет».

И сотворил муж тот, как заповедал ему Господь, и се – извлечен был бутерброд из видеомагнитофона, и спасен был муж тот.

И молвил Господь: «Истинно говорю тебе, не помышляй впредь класть бутерброды свои с пастой арахисовой в видеомагнитофон, ибо не место им там. Но храни их в холодильнике своем или на блюде своем, где надлежит им быть».

И внял муж тот словам Господним, и отныне хранил бутерброды свои в месте подобающем, и был избавлен от бед и скорбей. Аминь.

Спустя всего два месяца число пользователей ChatGPT выросло до 100 млн³. Компания OpenAI оказалась настолько не готова к такому взрыву интереса, что даже не успела обеспечить достаточные вычислительные мощности для обработки возросшего трафика.

Программисты быстро оценили потенциал нового ИИ: он отлично справлялся с генерацией фрагментов кода на осно-

³ Hu K. "ChatGPT Sets Record for Fastest-Growing User Base – Analyst Note." Reuters. February 2, 2023. <https://www.reuters.com/technology/chatgpt-sets-record-fastest-growing-user-base-analyst-note-2023-02-01/>.

ве простого текстового описания задачи. Похожий инструмент – GitHub Copilot – был и до этого, но с появлением ChatGPT применять ИИ в программировании стали куда чаще. Это значительно сократило время разработки приложений. Более того, даже люди без ИТ-навыков получили возможность создавать несложные программы!

Microsoft оперативно приобрела лицензию на технологию у OpenAI и интегрировала чат-бота в свою поисковую систему Bing, позволив ему отвечать на вопросы пользователей на основе результатов поиска. Google, хотя и разработала собственного чат-бота еще в 2021 году, не спешила его выпускать и интегрировать в свои продукты. Однако шаг с Bing был воспринят как прямая угроза позициям Google, и компания спешно анонсировала своего поискового чат-бота Bard (позже переименованного в Gemini)⁴.

Вот тогда-то и начались проблемы. В рекламном ролике Bard чат-бот заявил, что космический телескоп Джеймса Уэбба сделал первый снимок планеты за пределами Солнечной системы. Астрофизик тут же указал на ошибку⁵; итак, даже при тщательном выборе примера Google не сумела из-

⁴ DeGeurin M. "Why Google Isn't Rushing Forward with AI Chatbots." Gizmodo. December 14, 2022. <https://gizmodo.com/lamda-google-ai-chatgpt-openai-1849892728>.

⁵ Macintosh B [@bmac_astro]. "Speaking as someone who imaged an exoplanet 14 years before JWST was launched, it feels like you should find a better example?" X (formerly Twitter). February 7, 2023. https://twitter.com/bmac_astro/status/1623136549524353024.

бежать промаха. Рыночная стоимость компании мгновенно упала на \$100 млрд: инвесторы испугались, что поисковая система станет намного хуже выполнять простые фактологические запросы, если Google, как и обещала, интегрирует Bard в поиск⁶.

Этот конфуз, хоть он и дорого обошелся, был лишь первой ласточкой в череде проблем, связанных с неспособностью чат-ботов корректно работать с фактической информацией. Их слабость – прямое следствие принципов их работы. Они изучают статистические закономерности в массиве обучающих данных, в основном взятых из интернета, а затем генерируют «ремиксованный» текст на основе этих закономерностей. Что именно содержалось в обучающих данных, они могут и не помнить. Подробнее об этом мы поговорим в главе 4.

Злоупотребления чат-ботами уже стали повсеместными. Некоторые новостные сайты были уличены в публикации вредных ИИ-советов на важные темы вроде финансов⁷ и, что еще хуже, отказались прекратить использование этой технологии даже после того, как им указали на ошибки. Amazon наводнен некачественными книгами, созданными ИИ. В их числе несколько руководств по сбору грибов, где ошибки мо-

⁶ Wittenstein J. "Bard AI Chatbot Just Cost Google \$100 Billion." *Time*. February 9, 2023. <https://time.com/6254226/alphabet-google-bard-100-billion-ai-error/>.

⁷ Christian J. "CNET Sister Site Restarts AI Articles, Immediately Publishes Idiotic Error." *Futurism*. Updated February 1, 2023. <https://futurism.com/cnet-bankrate-restarts-ai-articles>.

гут стоять жизни доверчивому читателю⁸.

Напрашивается вывод: мир сошел с ума, раз он восторгается столь несовершенной технологией. Но этот вывод слишком примитивен; большинству отраслей, связанных со знаниями, чат-боты так или иначе полезны. Мы, авторы этой книги, сами пользуемся их помощью, делегируя им целый ряд задач: от рутинных моментов вроде правильного оформления цитат до сложных заданий, с которыми иначе не разобратся (например, перевод статьи, напичканной терминами из незнакомой нам области исследований).

Загвоздка в том, что без усилий и практики невозможно эффективно использовать чат-бота и избегать многочисленных подводных камней. Куда проще быстро зарабатывать, продавая, скажем, удручающего качества книгу, сгенерированную ИИ. Именно это и делает чат-ботов столь уязвимыми для злоупотреблений.

Есть и более острые вопросы, связанные с распределением власти в цифровом мире. Что будет, если компании, владеющие поисковыми системами, заменят привычный список из 10 ссылок на готовые ответы, сгенерированные искусственным интеллектом? Даже если предположить, что эти ответы будут точными, мы, по сути, получим машину, которая переписывает чужой контент и выдает его за оригиналь-

⁸ Cole S. "'Life or Death:' AI-Generated Mushroom Foraging Books Are All over Amazon." 404 Media. August 29, 2023. <https://www.404media.co/ai-generated-mushroom-foraging-books-amazon/>.

ный. При этом сайты-источники не получают ни трафика, ни дохода. Если бы поисковые системы просто брали контент с разных сайтов и выдавали за свой, они бы нарушили авторское право. Но ответы, сгенерированные ИИ, как будто позволяют обойти закон. Впрочем, к 2024 году уже подано множество исков, призванных изменить ситуацию⁹.

⁹ Brittain B. "OpenAI Asks Court to Trim Authors' Copyright Lawsuits." Reuters. August 29, 2023. <https://www.reuters.com/legal/litigation/openai-asks-court-trim-authors-copyright-lawsuits-2023-08-29/>.

Как ИИ изменил сферу развлечений

Еще одна технология генеративного ИИ, покоровшая публику, – создание изображений на основе текстового описания. К середине 2023 года пользователи сгенерировали свыше миллиарда картинок с помощью таких инструментов, как Dall-E 2 от OpenAI, Firefly от Adobe и Midjourney¹⁰. Отдельного упоминания заслуживает Stable Diffusion от Stability AI – открытый генератор изображений, который можно настроить под свои нужды. Инструменты на базе Stable Diffusion скачаны более 200 млн раз. Поскольку пользователи запускают его на своих устройствах, точно узнать количество созданных изображений невозможно, но, вероятно, счет идет на миллиарды.

Генераторы изображений породили поток развлекательного контента нового типа¹¹. В отличие от традиционных, ИИ-картинки можно бесконечно подстраивать под вкусы каждого пользователя. Кто-то наслаждается фантастическими пейзажами или футуристическими городами. Другим по душе изображения исторических личностей в современных

¹⁰ Valyaeva A. "AI Image Statistics: How Much Content Was Created by AI." *Insight* (blog). Everyapixel Journal. August 15, 2023. <https://journal.everypixel.com/ai-image-statistics>.

¹¹ Kapoor S, Narayanan A. "How to Prepare for the Deluge of Generative AI on Social Media." *Kn First Amend Inst*. June 16, 2023. <http://knightcolumbia.org/content/how-to-prepare-for-the-deluge-of-generative-ai-on-social-media>.

ситуациях или знаменитостей в необычных образах, например нашумевшая «фотография» папы римского Франциска в модном пуховике Balenciaga. Популярность обрели и фейковые трейлеры известных фильмов, стилизованные под почерк конкретных режиссеров, например «Звездные войны» в узнаваемой манере Уэса Андерсона, с его фирменными симметричными кадрами, пастельными тонами и причудливыми декорациями.

Генераторами изображений заинтересовались не только любители: развлекательные приложения на основе ИИ – большой бизнес. Разработчики видеоигр создают персонажей, с которыми игроки могут вести непринужденный диалог¹². Многие приложения для обработки фотографий теперь включают функции генеративного ИИ. Например, вы можете попросить такое приложение добавить воздушные шары на снимок с дня рождения.

ИИ стал одним из основных предметов спора во время голливудских забастовок 2023 года¹³. Актеры опасались, что студии смогут использовать кадры с их участием для обучения ИИ, способного генерировать новые видео на осно-

¹² NVIDIA Game Developer. "NVIDIA ACE for Games Sparks Life into Virtual Characters with Generative AI." YouTube video, 2:02. May 28, 2023. <https://www.youtube.com/watch?v=nAEQdF3JAJ0>.

¹³ Dalton A. "AI Is the Wild Card in Hollywood's Strikes. Here's an Explanation of Its Unsettling Role." AP News. July 21, 2023. <https://apnews.com/article/artificial-intelligence-hollywood-strikes-explained-writers-actors-e872bd63ab52c3ea9f7d6e825240a202>.

ве сценария. Эти видео были бы неотличимы от настоящих. Иными словами, студии получили бы возможность бесконечно эксплуатировать образы актеров и плоды их прошлого труда, не выплачивая им ни цента.

Забастовки завершились, но глубинные противоречия между трудом и капиталом никуда не делись. Они непременно всплывут опять, особенно с появлением новых технологий¹⁴. Одни компании работают над генераторами видео на основе текста, другие – над автоматизацией написания сценариев. Результат может уступать среднестатистическому фильму в художественной ценности и сложности сюжета и съемок, но для студий, которым надо выпустить очередной летний блокбастер, это будет неважно.

Мы полагаем, что со временем совместные усилия программистов и законодателей способны смягчить большинство описанных проблем и усилить преимущества ИИ. Уже накопилось много идей, как сделать чат-ботов менее склонными к выдумыванию информации и как с помощью новых законов обуздать намеренные злоупотребления. Но в краткосрочной перспективе приспособиться к миру с генеративным ИИ оказалось непросто: эти инструменты чрезвычайно мощны, но ненадежны. Это как если бы каждому человеку в мире вручили бесплатную бензопилу.

Потребуется немало труда, чтобы правильно интегриро-

¹⁴ Crawford K. *Atlas of AI: Power, Politics, and the Planetary Costs of Artificial Intelligence*. New Haven: Yale University Press; 2021.

вать ИИ в нашу жизнь. Яркий пример – школьники и студенты, которые пишут с помощью чат-ботов контрольные и сдают экзамены. Внесем ясность: появление ИИ угрожает образованию не больше, чем когда-то угрожало появление калькулятора¹⁵. При правильном подходе чат-бот может стать ценным обучающим инструментом. Но для этого преподавателям придется пересмотреть учебные программы, методики преподавания и формы контроля. В хорошо финансируемом учреждении вроде Принстона, где мы преподаем, это скорее возможность, чем проблема; мы даже поощряем студентов, которые пользуются ИИ. Но многие школьные учителя растерялись, когда ChatGPT внезапно начал помогать миллионам учеников списывать.

Будет ли общество вечно плестись в хвосте новых разработок генеративного ИИ? Или у нас хватит коллективной воли на структурные изменения, которые позволят более справедливо распределять крайне неравномерные выгоды и издержки новых технологий, какими бы они ни были?

¹⁵ Narayanan A. "Students Are Acing Their Homework by Turning in Machine-Generated Essays. Good." AI Snake Oil. October 21, 2022. <https://www.aisnakeoil.com/p/students-are-acing-their-homework>.

Предиктивный ИИ: громкие обещания и большие сомнения

Генеративный ИИ приносит с собой немало социальных издержек и рисков, особенно поначалу. Однако мы смотрим в будущее с осторожным оптимизмом, полагая, что этот тип ИИ может со временем улучшить качество жизни. С предиктивным ИИ ситуация иная.

В последние годы мы наблюдаем настоящий бум предиктивных ИИ. Разработчики уверяют, что их детища могут предвидеть, совершит ли подсудимый новое преступление или преуспеет ли соискатель на новой работе. Однако, в отличие от генеративного ИИ, предиктивный зачастую оказывается абсолютно несостоятельным¹⁶.

Возьмем, к примеру, программу Medicare в США: по ней получают медицинское страхование люди старше 65 лет. Пытаясь сократить расходы, поставщики услуг в этой области начали использовать ИИ, чтобы спрогнозировать длительность госпитализации больных¹⁷. Результаты часто оказыва-

¹⁶ Raji ID, Kumar IE, Horowitz A, Selbst A. "The Fallacy of AI Functionality." In *2022 ACM Conference on Fairness, Accountability, and Transparency*. Seoul Republic of Korea: ACM; 2022. p. 959–72. <https://dl.acm.org/doi/10.1145/3531146.3533158>.

¹⁷ Ross C, Herman B. "Denied by AI: How Medicare Advantage Plans Use Algorithms to Cut Off Care for Seniors in Need." *STAT*. March 13, 2023. <https://www.statnews.com/2023/03/13/medicare-advantage-plans-denial-artificial-intelligence/>.

ются далекими от реальности. Однажды ИИ предсказал, что 85-летняя пациентка будет готова к выписке через 17 дней. Когда этот срок истек, женщина все еще испытывала сильные боли и не могла передвигаться даже с ходунками. Тем не менее, опираясь на оценку ИИ, страховые выплаты прекратили.

Внедрение технологий предиктивного ИИ нередко начинается с благих намерений: например, компании хотят добиться того, чтобы пациент не оставался в больнице бесконечно долго. Однако со временем цели и методы использования системы искажаются: тот же ИИ в Medicare превратился в инструмент бездушной экономии, хотя изначально должен был повысить ответственность больничного персонала.

Подобные сценарии разворачиваются в самых разных областях. Некоторые компании, разрабатывающие ИИ для найма персонала, утверждают, что их детища могут оценить дружелюбие, открытость или доброту человека по языку тела, манере речи и другим поверхностным признакам, зафиксированным в 30-секундном видеоролике. Но насколько эффективен этот метод? И действительно ли такие личностные оценки связаны с успешностью профессиональной деятельности? Увы, компании, делающие подобные заявления, не предоставили весомых доказательств эффективности своих продуктов. Более того, есть множество свидетельств обратного: предугадать поведение человека чрезвычайно сложно. Подробнее об этом мы поговорим в главе 3.

В 2013 году страховая компания Allstate решила использовать предиктивный ИИ для определения страховых тарифов в штате Мэриленд. Цель была проста: заработать больше, не растеряв при этом клиентов. В результате появился так называемый «список простаков» – перечень людей, чьи страховые взносы резко выросли по сравнению с прежними тарифами¹⁸. В этом списке оказалось непропорционально много людей старше 62 лет: вероятно, ИИ уловил, что пожилые люди редко ищут более выгодные предложения. Это яркий пример автоматизированной дискриминации. Новая система ценообразования, скорее всего, увеличила бы доходы страховой компании, но с моральной точки зрения она неприемлема. При этом, хотя власти Мэриленда отвергли предложение Allstate использовать дискриминирующий ИИ-инструмент, компания применяет его как минимум в 10 других штатах¹⁹.

Если вам не нравится, что ИИ решает, кого брать на работу, сегодня вы можете просто не отправлять резюме в компании, где HR-отдел использует нейросети. Но если предиктивный ИИ применяют власти, выбора уже нет: приходится играть по их правилам. (Аналогичные проблемы возникают,

¹⁸ Varner M, Sankin A. "Suckers List: How Allstate's Secret Auto Insurance Algorithm Squeezes Big Spenders." The Markup. February 25, 2020. <https://themarkup.org/allstates-algorithm/2020/02/25/car-insurance-suckers-list>.

¹⁹ Большинство примеров в нашей книге, в том числе и этот, взяты из американской действительности: мы живем и работаем в США. Но выводы, которые мы делаем, вполне применимы и к другим странам.

если много компаний используют один и тот же искусственный интеллект при найме сотрудников.) Во многих странах судьи опираются на мнение ИИ, решая, брать ли обвиняемого под стражу до суда. Выяснилось, что «умные программы» руководствуются вполне человеческими предрассудками: расовыми, гендерными, возрастными. Хуже того: решения ИИ по определению того, кто опасен для общества, а кто нет, мало чем отличаются от подбрасывания монетки.

В чем же дело? Возможно, одна из причин столь низкой прогнозной точности заключается в том, что некоторые важные данные искусственному интеллекту просто недоступны. Представьте трех обвиняемых, у которых совпадают возраст, количество нарушений в прошлом и число родственников с судимостями. ИИ выдаст для них одинаковый уровень риска. А на деле один искренне раскаивается, второго задержали по ошибке, а третий спит и видит, как бы довести аферу до конца. Может ли ИИ это предсказать? Нет.

Кроме того, люди быстро учатся обманывать систему и манипулировать ею в своих интересах. Например, с помощью ИИ собирались предсказывать, как долго прослужит пересаженная почка²⁰. Предполагалось в первую очередь делать трансплантацию органа тем, кто проживет дольше. Но такая система отбила бы у пациентов всякое желание заботиться о своих почках, ведь в чем более молодом возрасте

²⁰ Robinson DG. *Voices in the Code: A Story about People, Their Values, and the Algorithm They Made*. New York: Russell Sage Foundation; 2022.

они откажут, тем выше шансы на пересадку! К счастью, эту идею обсудили с пациентами, врачами и другими заинтересованными лицами. Они вовремя заметили ловушку, и устанавливать очередь с помощью ИИ никто не стал.

О провалах предиктивного ИИ мы еще поговорим в главах 2 и 3. Станет ли ситуация лучше со временем? Сомнительно. У этой технологии слишком много «врожденных» дефектов. Например, предиктивный ИИ кажется привлекательным, потому что автоматизирует принятие решений: это эффективнее. Но именно погоня за такой эффективностью и приводит к тому, что никто ни за что не отвечает. Так что, когда компании расхваливают своих «электронных предсказателей», не спешите верить и требуйте железных доказательств.

Существует ли «ИИ вообще»?

Генеративный и предиктивный ИИ – два основных вида искусственного интеллекта. А сколько их всего? Ответить непросто: даже эксперты не могут прийти к единому мнению о том, что считать ИИ, а что нет.

Для того чтобы разобраться, действительно ли перед нами ИИ, можно задать три вопроса о том, как система решает задачу. Ответ на каждый вопрос проливает свет на какую-то из сторон ИИ, но полного определения не дает.

Первый вопрос: нужны ли человеку творческие способности или специальные навыки, чтобы выполнить эту же задачу? Если да, а компьютер с ней справляется, – вполне вероятно, что это ИИ. Поэтому создание изображений относят к искусственному интеллекту. Например, чтобы нарисовать картинку, человеку нужны навыки художника или дизайнера. Но даже такую пустячную для нас задачу, как распознавание объектов типа кошки или чайника, компьютеры освоили только к 2010-м годам. И это тоже считают ИИ. Выходит, сравнение с человеческим интеллектом – не единственный критерий.

Второй вопрос: заложено ли поведение системы напрямую в код, или оно возникло косвенно, например в результате обучения на примерах (машинного обучения) или поиска в базах данных? Если второе – это может быть ИИ.

Этот критерий объясняет, почему создание формулы расчета страховки могут отнести к искусственному интеллекту, если компьютер сам вывел ее из данных о прошлых страховых случаях, а если точно такую же формулу напрямую составил эксперт – уже нет. Хотя некоторые системы с заранее прописанными алгоритмами все же считаются ИИ. Например, к ним относятся роботы-пылесосы, умеющие объезжать препятствия.

Третий критерий: насколько самостоятельно система принимает решения, способна ли она гибко адаптироваться к окружающей среде? Если да – возможно, перед нами ИИ. Яркий пример – беспилотные автомобили. Этот критерий тоже не дает всеобъемлющего определения: мы же не назовем ИИ обычный механический термостат? Он просто самостоятельно реагирует на изменение температуры расширением или сжатием металла, включая или выключая ток.

Отнесут ли очередную новую технологию к ИИ, во многом зависит от истории ее использования, маркетинга и других факторов. Мы не будем переживать из-за отсутствия четкого определения «ИИ вообще».

Странно, скажете вы, ведь книга-то об ИИ! Но вспомните нашу главную мысль: почти невозможно придумать что-то, что относилось бы сразу ко всем видам нейросетей. В основном мы будем обсуждать их конкретные типы, для которых у нас есть определения, и это позволит нам найти общий язык.

Есть забавная, но достаточно меткая формулировка:

«ИИ – то, что еще не сделано». Как только приложение начинает работать стабильно, оно становится привычным и уже не воспринимается как ИИ. Роботы-пылесосы, автопилот в самолетах, автозаполнение в телефонах, распознавание почерка и устной речи, спам-фильтры, проверка орфографии... Да-да, когда-то даже это считалось сложной задачей!

Все эти инструменты прекрасны. Они незаметно улучшают нашу жизнь. Именно такие ИИ нам и нужны! А наша книга – о проблемных видах искусственного интеллекта (вряд ли вы захотите прочесть 300 страниц о достоинствах ИИ, занимающегося проверкой орфографии). Тем не менее важно понимать: далеко не каждый ИИ вреден.

Некоторые новые технологии со временем, надеемся, станут обыденностью. Сегодня беспилотные автомобили попадают в новости из-за аварий и жертв²¹. Но безопасное автономное вождение – решаемая задача, хотя ее сложность часто недооценивают. Гораздо серьезнее может оказаться проблема массовой потери рабочих мест, если эта технология получит широкое распространение: миллионы людей водят грузовики, такси или работают в каршеринге. И все же, если удастся решить проблему безопасности и принять необходимые социальные и политические меры, возможно, однажды мы станем воспринимать беспилотные автомобили такой же

²¹ Marcus G. "Face It, Self-Driving Cars Still Haven't Earned Their Stripes." *Marcus on AI* (blog). August 19, 2023. <https://garymarcus.substack.com/p/face-it-self-driving-cars-still-havent>.

естественной частью повседневной жизни, как и лифты.

Однако мы думаем, что некоторые виды ИИ, особенно предиктивного типа, вряд ли когда-нибудь станут обыденностью. Точно предсказывать социальное поведение людей – технически неразрешимая задача. А определять судьбу человека на основе заведомо ошибочных прогнозов всегда сомнительно с моральной точки зрения.

Для того чтобы лучше понять, почему обобщения в отношении ИИ недопустимы, рассмотрим тревожащий правозащитников пример – технологию распознавания. На момент написания этой книги она уже привела к шести ошибочным арестам в США, и все арестованные были чернокожими. Стоит ли запретить полиции использовать распознавание лиц из-за того, что оно чаще дает сбои, идентифицируя чернокожих?

В этих спорах легко упустить из виду важный факт: каждый из этих ложных арестов связан с цепочкой ошибок в работе полиции, преимущественно человеческих. Роберта Уильямса арестовали за кражу в магазине, основываясь в том числе на показаниях охранника, которого даже не было на месте преступления²². Рэндалла Рида задержали в Джорджии за кражу, совершенную в Луизиане – штате, где он

²² Ryan-Mosley T. "The New LawsUIT That Shows Facial Recognition Is Officially a Civil Rights Issue." *MIT Technology Review*. April 14, 2021. <https://www.technologyreview.com/2021/04/14/1022676/robert-williams-facial-recognition-lawsuit-aclu-detroit-police/>.

никогда не бывал²³. Поршу Вудрафф идентифицировали по фотографии 2015 года, хотя была доступна более свежая – с водительских прав 2021 года²⁴. И так далее.

Полицейские ошибки, приводящие к арестам невиновных, случаются ежедневно. И скорее всего, без них не обойдется и в будущем – независимо от того, будет использоваться технология распознавания лиц или нет. Сотни тысяч поисков по лицам уже проведены, и ошибок оказалось крайне мало²⁵. По данным Национального института стандартов и технологий, с 2014 по 2020 год точность выросла в 50 раз: теперь сбои случаются лишь в 0,08% случаев²⁶.

ИИ по распознаванию лиц обычно работает точно: задача достаточно конкретна. Его обучают на огромных базах фотографий, помечая, какие снимки принадлежат одному и тому же человеку. Получив достаточно данных, нейросеть учится

²³ "Facial Recognition Tool Led to Mistaken Arrest, Lawyer Says." AP News. January 2, 2023. <https://apnews.com/article/technology-louisiana-baton-rouge-new-orleans-crime-50e1ea591aed6cf14d248096958dccc4>.

²⁴ Hill K. "Eight Months Pregnant and Arrested after False Facial Recognition Match." *The New York Times*. August 6, 2023. <https://www.nytimes.com/2023/08/06/business/facial-recognition-false-arrest.html>.

²⁵ Cipriano A. "Facial Recognition Now Used in over 1,800 Police Agencies: Report." *The Crime Report*. April 7, 2021. <https://thecrimereport.org/2021/04/07/facial-recognition-now-used-in-over-1800-police-agencies-report/>.

²⁶ Crumpler W. "How Accurate Are Facial Recognition Systems – and Why Does It Matter?" *Strategic Technologies Blog*. CSIS. April 14, 2020. <https://www.csis.org/blogs/strategic-technologies-blog/how-accurate-are-facial-recognition-systems-and-why-does-it>.

различать черты лиц. Это сильно отличается от других задач – например, определения пола или эмоций, где ошибок куда больше²⁷,²⁸. Ключевое отличие: все, что нужно для узнавания лица, есть на снимке. А чтобы угадать пол или настроение, приходится делать предположения, что сразу снижает точность.

Правозащитники часто ставят ИИ по распознаванию лиц в один ряд с другими спорными технологиями, применяемыми в правосудии, вроде прогнозирования преступлений, несмотря на то что технологии совершенно разные, не говоря уже об их точности. К примеру, большинство людей, которых ИИ относит к группе высокого риска, на деле не совершают новых преступлений.

Главная опасность распознавания лиц в том, что оно работает слишком хорошо и может попасть не в те руки. Кашмир Хилл в книге «Ваше лицо принадлежит нам» (*Your Face Belongs To Us*) приводит немало примеров злоупотреблений²⁹. Например, некоторые режимы используют эту техно-

²⁷ Buolamwini J, Gebru T, Raynham H, Raji D, Zuckerman E. "Gender Shades." MIT Media Lab. <https://www.media.mit.edu/projects/gender-shades/overview/>. Accessed February 15, 2024.

²⁸ Buolamwini J, Gebru T. "Gender Shades: Intersectional Accuracy Disparities in Commercial Gender Classification." In: *Proceedings of the 1st Conference on Fairness, Accountability and Transparency*. PMLR; 2018. pp. 77–91. <https://proceedings.mlr.press/v81/buolamwini18a.html>.

²⁹ Hill K. *Your Face Belongs to Us: A Secretive Startup's Quest to End Privacy as We Know It*. New York: Random House; 2023.

логию, чтобы выявлять и преследовать участников мирных протестов³⁰.

Распознавание лиц может стать опасным оружием и в руках частных компаний. Вот вам история: в 2022 году адвоката Николетт Ланди не пустили на концерт Мэрайи Кэри на нью-йоркской арене «Мэдисон-сквер-гарден»³¹. Билеты за \$400 купил ее парень: это был подарок на день рождения. Ланди оказалась не одинока: многих других юристов тоже не пустили на концерт. Почему? Владельцы арены забанили всех адвокатов из фирм, когда-либо судившихся с ними, неважно, участвовал ли конкретный юрист в иске и как давно он ходит на концерты; даже завсегдаев разворачивали от дверей. И все это с помощью системы распознавания лиц.

Критики технологии твердят: она не работает, ее надо запретить, а исследователей – пристыдить. Но так можно упустить реальную пользу. Например, однажды Министерство внутренней безопасности США за три недели раскрыло кучу старых дел о насилии над детьми. Каким образом? Искали преступников по фото и видео, которые те сами выложили

³⁰ "Russia: Police Target Peaceful Protesters Identified Using Facial Recognition Technology." Amnesty International. April 27, 2021. <https://www.amnesty.org/en/latest/press-release/2021/04/russia-police-target-peaceful-protesters-identified-using-facial-recognition-technology/>.

³¹ Hill K, Kilgannon C. "Madison Square Garden Uses Facial Recognition to Ban Its Owner's Enemies." *The New York Times*. December 22, 2022. <https://www.nytimes.com/2022/12/22/nyregion/madison-square-garden-facial-recognition.html>.

в соцсетях³². Так удалось опознать сотни потерпевших и насильников. Не стоит забывать и про бытовые удобства: эта же технология позволяет нам разблокировать собственные телефоны и сортировать фотографии.

Стоит уточнить: хотя распознавание лиц может быть очень точным при правильном использовании, на практике оно нередко дает сбой. Скажем, если применять его не к четким фото, а к зернистой картинке с камер наблюдения, вероятность ошибки резко возрастает. После того как американская сеть аптек Rite Aid внедрила подобную систему, сотрудники то и дело безосновательно обвиняли покупателей в кражах. Ложных срабатываний было несколько тысяч. Компания изо всех сил пыталась держать технологию в тайне. К счастью, власти не дремали: Федеральная торговая комиссия запретила Rite Aid использовать распознавание лиц для слежки за покупателями на целых пять лет³³.

Подытожим: чтобы найти баланс в использовании этой неоднозначной технологии, нужно активное общественное обсуждение. Нам еще только предстоит определить, где распознавание лиц уместно, а где нет, и выработать четкие пра-

³² Brewster T. "Exclusive: DHS Used Clearview AI Facial Recognition in Thousands of Child Exploitation Cold Cases." *Forbes*. August 7, 2023. <https://www.forbes.com/sites/thomasbrewster/2023/08/07/dhs-ai-facial-recognition-solving-child-exploitation-cold-cases/>.

³³ "Rite Aid Corporation, FTC v." Cases and Proceedings. Federal Trade Commission. 2023. <https://www.ftc.gov/legal-library/browse/cases-proceedings/2023190-rite-aid-corporation-ftc-v>.

вила, которые не позволят ни властям, ни бизнесу злоупотреблять этой технологией.

**Череда любопытных
совпадений, которая привела
к появлению ЭТОЙ КНИГИ**

Конец ознакомительного фрагмента.

Текст предоставлен ООО «Литрес».

Прочитайте эту книгу целиком, [купив полную легальную версию](#) на Литрес.

Безопасно оплатить книгу можно банковской картой Visa, MasterCard, Maestro, со счета мобильного телефона, с платежного терминала, в салоне МТС или Связной, через PayPal, WebMoney, Яндекс.Деньги, QIWI Кошелек, бонусными картами или другим удобным Вам способом.