

ТЕМА	ЭЛЕМЕНТ	БАЗОВЫЙ	ПРОДВИНУТЫЙ	ЭКСПЕРТНЫЙ	МАТЕРИАЛЫ
Понятие искусственного интеллекта	Знания	Общее понятие искусственного интеллекта. Базовые принципы работы и обучения нейронных сетей. История термина "Искусственный интеллект": история возникновения в странах мира и в России. - Современное понимание ИИ в различных юрисдикциях и в различных видах документов.	Особенности технологии искусственного интеллекта, связанные с возможностью возникновения этических и социальных эрозий. Контекст ИИ: неподобие или антиподность, социальный уровень, автономия системы ИИ, зависимость качества работы ИИ-сервисов от качества, пропрезентивность, балансированности и разметки данных, используемых для обучения.	Принципы работы и обучения нейронных сетей. Глубокое обучение, обучение с подкреплением: обобщимость и интерпретируемость моделей. Понимание, какие задачи решает ИИ, каковы основные области применения ИИ, включая распознавание речи и образов, обработка естественного языка, рекомендации и пр..	<p>ПРИКАЗ МИНИСТЕРСТВА ЭКОНОМИЧЕСКОГО РАЗВИТИЯ РОССИЙСКОЙ ФЕДЕРАЦИИ от 15.02.2024 г. № 392-о УТВЕРЖДЕН КРИТЕРИИ ПРИНАДЛЕЖНОСТИ ПРОЕКТОВ ИСКУССТВЕННОГО ИНТЕЛЛЕКТА [http://publications.mre.gov.ru/document/VIEW/120210729008] НАЧАЛЬНИКА СЛУЖБЫ ГИДРОГЕОЛОГИИ УКАЗ ПРЕЗИДЕНТА РОССИЙСКОЙ ФЕДЕРАЦИИ от 15.02.2024 г. № 124-о ВНЕСЕНИИ ИЗМЕНЕНИЙ В КАЗЫМ-ХАСЫНСКИЙ УЧЕБНИК ПО ПРАВОМ ИНТЕЛЛЕКТА В РОССИЙСКОЙ ФЕДЕРАЦИИ [http://publications.mre.gov.ru/document/VIEW/120210729008] МАРКОВ С. ОХОТА НА ЗЕМЛЮ: ПРАВОМОСТЬ КИИ ИСКУССТВЕННОГО ИНТЕЛЛЕКТА. МОСКВА, 2024. КОДЕКС ЭТИКИ И РЕГУЛИРОВАНИЯ #BOOK + ДОКУМЕНТЫ ПО СТАНДАРТИЗАЦИИ КОНВЕНЦИИ ЮНЕСКО ПО ЭТИКЕ</p>
	Практическое задание	Сформулируйте дефиницию термина ИИ в различных юрисдикциях. Ответьте на вопрос – существует ли отличие в понимании указанного термина?	Представьте, что Вы участвуете в разработке системы ИИ для решения задачи по подбору персонала, выдачи кредитов и т.п.). Какие данные требуется для обучения этой системы ИИ? Какими могут быть социальные последствия недостаточной пропрезентивности или несбалансированности данных?	Выберите в списке технологий те, что не относятся к технологиям ИИ и поясните выбор: 1. Автоматизация. 2. Глубокое обучение. 3. Компьютерное программирование. 4. Экспертные системы.	
	Вопрос для контрольной/экзамена	Дайте определение ИИ и приведите различия в их понимании в различных странах.	Чем ИИ отличается от автоматизации или обычной ИИ-системы?	Объясните основные принципы работы и обучения нейронных сетей, выделяя различия между нейронными и рекуррентными нейросетями. Опишите, в чем заключаются преимущества и недостатки обучения с подкреплением, а также обсудите важность обобщимости и интерпретируемости моделей. Приведите примеры задач, которые требуют применения и изучение областей применения этих технологий, такие как распознавание речи, обработка естественного языка и рекомендательные системы.	
Механизмы регулирования в сфере искусственного интеллекта	Знания	I. Особенности различных типов регулирования ИИ. Принцип взаимодополнения трех типов регулирования в сфере ИИ: 1. Кодекс этики в сфере ИИ как инструмент саморегулирования источников будущих правовых норм в сфере ИИ. 2. Правовое регулирование ИИ: возможности и ограничения. 3. Нормативно-техническое регулирование: стандартизация и сертификация в сфере ИИ.  II. Вертикальный и горизонтальный подходы к регулированию ИИ: определения и цели.	Специфика регулирования ИИ в РФ. Ключевые стратегические документы: 1. Общие положения Национальной стратегии РФ в области ИИ; 2. Концепция регулирования ИИ – общее направление развития регулирования технологий; 3. ФП "Искусственный интеллект" нац. программы "Цифровая экономика РФ" – общий вектор развития;	Индексы развития ИИ в разных странах – регуляторная политика и этика ИИ: 1. AI Index Report Stanford 2. The Global AI Index (Tortoise) 3. Этические индексы Канады и Китая	<p>НАЦИОНАЛЬНАЯ СТРАТЕГИЯ: УДА ГЛАВЫ РЕСПУБЛИКИ РОССИЙСКОЙ ФЕДЕРАЦИИ от 15.02.2024 г. № 124-о ВНЕСЕНИИ ИЗМЕНЕНИЙ В КАЗЫМ-ХАСЫНСКИЙ УЧЕБНИК ПО ПРАВОМ ИНТЕЛЛЕКТА В РОССИЙСКОЙ ФЕДЕРАЦИИ от 10 ОКТЯБРЯ 2019 г. № 490-о РАЗВИТИЕ ИСКУССТВЕННОГО ИНТЕЛЛЕКТА В РОССИЙСКОЙ ФЕДЕРАЦИИ И НАЦИОНАЛЬНАЯ СТРАТЕГИЯ, УТВЕРЖДЕННУЮ ЭТИМ КАЗОМ [http://publications.mre.gov.ru/document/VIEW/120210729008]</p>
	Практическое задание	1. Формулируйте преимущества и ограничения каждого из 3 типов регулирования для развития технологий ИИ.	Проанализируйте конкретный работающий сервис, основанный на ИИ, на предмет его соответствия Кодексу этики ИИ. Установите, почему проверка некоторых положений Кодекса оказалась неосуществимой из-за недостатка информации из открытых источников, оцените, насколько это повлияет на Ваш пользовательский опыт.	1. Глобальный опыт регулирования: создайте таблицу, обобщающую мировой опыт регулирования ИИ, с выделением встроенных в него этических подходов и перспектив развития технологий.	<p>Декларация об ответственном регулировании ИИ в зарубежных странах лидеров на разработку и / или внедрение ИИ: ЕС, Китай, США, Сингапур.</p>
	Вопрос для контрольной/экзамена	Каковы отличия правового регулирования этического? В чем заключается преимущества этического регулирования в сфере ИИ по сравнению с правовым и нормативно-техническим?	Опишите основные компоненты отечественной системы регулирования в сфере ИИ. Как компоненты связаны между собой?	Какой ключевой фактор, по вашему мнению, необходимо учитывать при оценке соответствия этических и правовых норм этическому развитию технологий ИИ? Объясните, почему этот фактор важен и как он может влиять на разработку более эффективного регулирования.	<p>Декларация об ответственном регулировании ИИ на основе генеративного ИИ: Кодекс этики и фундаментальная стратегия в сфере ИИ и РЕГУЛИРОВАНИЕ ИСКУССТВЕННОГО ИНТЕЛЛЕКТА ГЕНЕРАТИВНЫМ ИИ [http://publications.mre.gov.ru/document/VIEW/120210729008]</p>
Ключевые этические принципы, ценности и риски в сфере ИИ	Знания	Основные этические принципы в сфере ИИ: справедливость, прозрачность, конфиденциальность, обобщимость, надежность, ответственность, человекоцентрированность, автономия человека, непредвзятость, открытость, прозрачность, дисциплинированность, неискованность частной жизни, защита от опасного использования ИИ. Умени анализировать социальные, экономические и экологические последствия применения ИИ-технологий с использованием методов оценки возможностей, таких как анализ сценариев или социально-экономическое моделирование.	Концепция риско-ориентированного подхода в рамках этики в сфере ИИ. Понимание влияния этических принципов на принятие решений в работе ИИ и на разработки и использование ИИ на примере конкретных кейсов. Знание подходов для поиска компромиссных решений при столкновении различных этических требований в контексте ИИ (например, между конфиденциальностью и безопасностью) и умение смоделировать ситуацию.	Особенности разработки этических стандартов и политик в сфере ИИ с учетом особенностей ИИ-сервисов, а также положений национального регулирования сферы ИИ; специфика разработки нормативных рекомендаций и стандартов, применяемых в различных отраслевых и специфических областях ИИ-сервисов и международных аспектах регулирования ИИ.	<p>КОДЕКС ЭТИКИ В СФЕРЕ ИИ РФ: СТАНДАРТЫ-АЛЮР ДЕКЛАРАЦИЯ ОБ ОТВЕТСТВЕННОМ РАЗРАБОТКЕ И ИСПОЛЬЗОВАНИИ ИСКУССТВЕННОГО ИНТЕЛЛЕКТА ГЕНЕРАТИВНЫМ ИИ [http://publications.mre.gov.ru/document/VIEW/120210729008]</p>
	Практическое задание	Приведите пример существующего (существовавшего ранее) сервиса на основе ИИ, анализирующего, какие этические принципы утверждаются или нарушаются при разработке и использовании данного сервиса (примерный перечень ценностей включает в себя: человекоцентрированность, автономия человека, непредвзятость, открытость, прозрачность, дисциплинированность, неискованность частной жизни, защита от опасного использования ИИ).	Как ответственность за ошибки или негативные последствия работы ИИ должна быть распределена между разработчиками, пользователями и владельцами системы? Какие этические принципы помогут правильно определить границы этой ответственности?	Крупная компания разрабатывает ИИ-сервис (можно взять, например, образовательный сервис, сервис по подбору персонала, сервис генерации изображений и т.д.), который планируется запустить на международный рынок. Подготовьте стратегию управления этическими рисками и нормативные рекомендации, учитывающие как национальные, так и международные аспекты регулирования в сфере ИИ.	<p>ДЕКЛАРАЦИЯ ОБ ОТВЕТСТВЕННОМ РАЗРАБОТКЕ И ИСПОЛЬЗОВАНИИ ИСКУССТВЕННОГО ИНТЕЛЛЕКТА ГЕНЕРАТИВНЫМ ИИ ETHICS-AI.RU</p>
	Вопрос для контрольной/экзамена	2. Перечислите потенциальные риски, которые могут быть связаны с использованием сервиса. Поясните и обоснуйте свой выбор.	Что подразумевается под риском-ориентированным подходом в этике в сфере ИИ? Как отдельные этические принципы могут повлиять на возникновение различных негативных последствий? Приведите пример того, как такой подход можно применить на практике при разработке и внедрении системы ИИ в социальной сфере.	Опишите ключевые аспекты разработки этических стандартов для ИИ-сервисов, предназначенного для использования в международном контексте.	<p>ДЕКЛАРАЦИЯ ОБ ОТВЕТСТВЕННОМ РАЗРАБОТКЕ И ИСПОЛЬЗОВАНИИ ИСКУССТВЕННОГО ИНТЕЛЛЕКТА ГЕНЕРАТИВНЫМ ИИ AI.RISK.REPOSITORY AIRISK.MIT.EDU [http://publications.mre.gov.ru/document/VIEW/120210729008]</p>
Способы контроля возникновения этических рисков в сфере ИИ	Знания	Объясните, как принципы справедливости, прозрачности и обобщимости применяются при разработке и внедрении систем искусственного интеллекта.	Подход комплексной имитации реальных атак (red teaming) для выявления этически значимых уязвимостей систем ИИ включает следующие методы:	Концепция организации внешнего независимого этического комитета, консультирующего по этически значимым аспектам проектирования, разработки, тестирования, внедрения и использования искусственного интеллекта. Это формирование внутренней политики по повышению уровня ответственности ИИ-разработчиков и создателей для разработки ИИ-решений.	<p>КОДЕКС ЭТИКИ В СФЕРЕ ИИ РФ: СТАНДАРТЫ-АЛЮР ДЕКЛАРАЦИЯ ОБ ОТВЕТСТВЕННОМ РАЗРАБОТКЕ И ИСПОЛЬЗОВАНИИ ИСКУССТВЕННОГО ИНТЕЛЛЕКТА ГЕНЕРАТИВНЫМ ИИ ISO/IEC JTC1/SC2/WG1/AI/ARTIFICIAL INTELLIGENCE ISO.ORG [http://publications.mre.gov.ru/document/VIEW/120210729008]</p>
	Практическое задание	Приведите пример существующего (существовавшего ранее) сервиса на основе ИИ, анализирующего, какие этические принципы утверждаются или нарушаются при разработке и использовании данного сервиса (примерный перечень ценностей включает в себя: человекоцентрированность, автономия человека, непредвзятость, открытость, прозрачность, дисциплинированность, неискованность частной жизни, защита от опасного использования ИИ).	Проведите тестирование одной из моделей и найдите этические уязвимости. Опишите возможные способы устранения выявленных уязвимостей.	Разработайте пример пункта дополнительного соглашения к договору между заказчиком и разработчиком, который регулируется бы этическими аспектами (например, прозрачность, механизмы мониторинга и ответственности).	<p>ДЕФИНИЦИЯ СТАНДАРТА AI RISK REPOSITORY AIRISK.MIT.EDU [http://publications.mre.gov.ru/document/VIEW/120210729008]</p>
	Вопрос для контрольной/экзамена	Объясните, как принципы справедливости, прозрачности и обобщимости применяются при разработке и внедрении систем искусственного интеллекта.	Составьте перечень задач, которые решаются с помощью метода red teaming	Какие преимущества связаны с созданием внешнего независимого этического комитета для компании-разработчика ИИ?	<p>ДЕФИНИЦИЯ СТАНДАРТА AI RISK REPOSITORY AIRISK.MIT.EDU [http://publications.mre.gov.ru/document/VIEW/120210729008]</p>
Основы отраслевого этического регулирования	Знания	Технологические методы интеграции этических ценностей и принципов в технологиях ИИ, обучение с помощью обратной связи ( обратной связью от людей (RLFH)), работа ИИ-тренеров, RAG-решения, оценка возможностей моделей по решению этических задач (бенчмаркинг), методы формирования тестовых наборов для создания ИИ-решений, соответствующих ожиданиям общества.	Целенаправленный джейлбрейкинг. Испытания, направленные на взлом и обнаружение механизмов ИИ-сервисов, чтобы выявить уязвимости, которые могут быть злоупотреблены и выведены из строя. Сценарии, в которых алгоритмы могут быть использованы во вред. Такой подход помогает разработчикам понять, как ИИ может быть манипулирован и использован с нарушением этических норм.	Концепция организации внешнего независимого этического комитета, консультирующего по этически значимым аспектам проектирования, разработки, тестирования, внедрения и использования искусственного интеллекта. Это формирование внутренней политики по повышению уровня ответственности ИИ-разработчиков и создателей для разработки ИИ-решений.	<p>ДЕФИНИЦИЯ СТАНДАРТА AI RISK REPOSITORY AIRISK.MIT.EDU [http://publications.mre.gov.ru/document/VIEW/120210729008]</p>
	Практическое задание	Составьте перечень существующих бенчмарков, позволяющих оценить ответы модели ИИ на предмет соответствия нормам этики ответы модели ИИ.	Проведите тестирование одной из моделей и найдите этические уязвимости. Опишите возможные способы устранения выявленных уязвимостей.	Разработайте пример пункта дополнительного соглашения к договору между заказчиком и разработчиком, который регулируется бы этическими аспектами (например, прозрачность, механизмы мониторинга и ответственности).	<p>ДЕФИНИЦИЯ СТАНДАРТА AI RISK REPOSITORY AIRISK.MIT.EDU [http://publications.mre.gov.ru/document/VIEW/120210729008]</p>
	Вопрос для контрольной/экзамена	Какие этические ценности помогают внедрить в ИИ обучение с подкреплением на основе обратной связи от людей (RLFH)?	Составьте перечень задач, которые решаются с помощью метода red teaming	Какие преимущества связаны с созданием внешнего независимого этического комитета для компании-разработчика ИИ?	<p>ДЕФИНИЦИЯ СТАНДАРТА AI RISK REPOSITORY AIRISK.MIT.EDU [http://publications.mre.gov.ru/document/VIEW/120210729008]</p>
Основы отраслевого этического регулирования	Знания	Понимание ключевых этических принципов и стандартов, применимых к конкретным отраслям (например, здравоохранение, транспорт, финансы). Способность идентифицировать основные этические вызовы, связанные с применением ИИ в отраслевых приложениях. Умение объяснять важность соблюдения этических норм и стандартов в профессиональной практике.	Знание основных нормативных актов и кодексов регулирующих ИИ в этих областях. Глубокое понимание отраслевых требований (образование, медицина, право, судоисправление, правоохранительные органы и т.д.) и международных и национальных стандартов, а также этических нормативов в различных юрисдикциях. Знание механизмов контроля и обеспечения соблюдения этических норм. Способность анализировать и оценивать этические риски, связанные с отраслью в конкретной сфере, а также в других отраслях, а также предпринимать меры по их минимизации. Умение разрабатывать и обосновывать политику этического использования ИИ для конкретных проектов.	Особенности существующего нормативно-технического, правового и этического регулирования социально значимых отраслей ИИ. Специфика разработки этических и правовых стандартов, а также правового регулирования для социально значимых отраслей ИИ. Национальных и международных законодательных актов, регулирующих этические нормы и стандарты в различных юрисдикциях. Способность адаптировать и применять эти нормы в конкретной отрасли различий в регулировании ИИ в разных странах и отраслях. Умение адаптировать эти стандарты к специфическим условиям разных юрисдикций и отраслей.	<p>ДЕФИНИЦИЯ СТАНДАРТА AI RISK REPOSITORY AIRISK.MIT.EDU [http://publications.mre.gov.ru/document/VIEW/120210729008]</p>
	Практическое задание	Составьте перечень наиболее важных и применимых этических принципов применения ИИ в одной из социально значимых отраслей.	Провести анализ кейса, связанного с нарушением этических норм при применении ИИ в отрасли (например, в медицинских данных), и предложить способы улучшения этического регулирования.	Представьте, что вы являетесь консультантом по этике и правовому регулированию ИИ. Вам поручили подготовить рекомендации для внедрения системы искусственного интеллекта в социальную сферу, например, в образование или здравоохранение.	<p>ДЕФИНИЦИЯ СТАНДАРТА AI RISK REPOSITORY AIRISK.MIT.EDU [http://publications.mre.gov.ru/document/VIEW/120210729008]</p>
	Вопрос для контрольной/экзамена	Привести пример наиболее важных и применимых этических принципов применения ИИ в одной из социально значимых отраслей.	Какие основные этические риски связаны с использованием технологий искусственного интеллекта в отрасли (например, в здравоохранении, образовании, транспорте и т.д.)? К каким социальным последствиям могут привести ошибки или смещения данных, используемых для обучения систем ИИ в транспорте или системах изображений в медицине? Какие меры могут быть приняты для их минимизации?	Представьте, что вы являетесь консультантом по этике и правовому регулированию ИИ. Вам поручили подготовить рекомендации для внедрения системы искусственного интеллекта в социальную сферу, например, в образование или здравоохранение.	<p>ДЕФИНИЦИЯ СТАНДАРТА AI RISK REPOSITORY AIRISK.MIT.EDU [http://publications.mre.gov.ru/document/VIEW/120210729008]</p>