

Введение в искусственный интеллект. Современное компьютерное зрение

Семинар 10. Проблемы и вызовы современного искусственного интеллекта

Бабин Д.Н., Иванов И.Е., Петюшко А.А.

кафедра Математической Теории Интеллектуальных Систем

27 апреля 2021 г.



Что же такое искусственный интеллект?

Естественный интеллект (человек)

- Может воспринимать информацию, ее анализировать, принимать решения на основе анализа



Что же такое искусственный интеллект?

Естественный интеллект (человек)

- Может воспринимать информацию, ее анализировать, принимать решения на основе анализа

Искусственный интеллект

- (Сильный) то же самое, что и естественный, только на месте человека — компьютер



Что же такое искусственный интеллект?

Естественный интеллект (человек)

- Может воспринимать информацию, ее анализировать, принимать решения на основе анализа

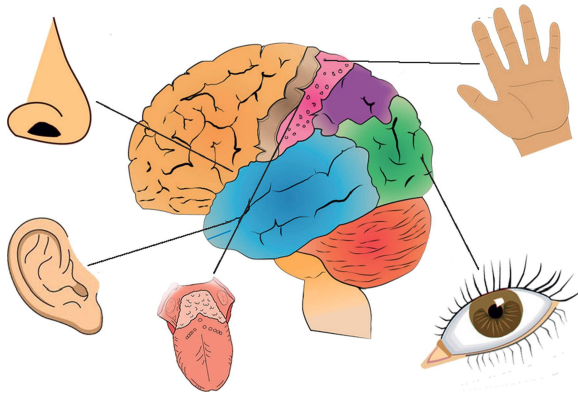
Искусственный интеллект

- (Сильный) то же самое, что и естественный, только на месте человека — компьютер
- (**Слабый**) алгоритм, способный обучиться на основе массива входных данных, чтобы затем выполнять задачу вместо человека



Взаимодействие со средой

- Около 90 % информации поступает через зрение¹
- Около 9 % информации поступает через слух



¹https://www.rlsnet.ru/books_book_id_2_page_40.htm

ИСКУССТВЕННЫЙ ИНТЕЛЛЕКТ ПРЕВЗОШЕЛ ПРОФЕССИОНАЛЬНЫХ ЮРИСТОВ ПО СКОРОСТИ И ЭФФЕКТИВНОСТИ ОЦЕНКИ ДОГОВОРОВ



- Самоуправляемые автомобили
- Магазины без продавцов
- Голосовые помощники
- Рекомендательные системы
- Интернет вещей
- Умный дом
- Умный город



- Самоуправляемые автомобили
- Магазины без продавцов
- Голосовые помощники
- Рекомендательные системы
- Интернет вещей
- Умный дом
- Умный город

Вывод

Влияние искусственного интеллекта на жизнь человека огромно и оно продолжает расти



Вы могли не замечать, но уже

- Вы смотрите видео, которое рекомендует вам youtube
- Вы передвигаетесь по маршруту, который вам построил навигатор
- В банке робот решает давать вам кредит или нет
- Штрафы за нарушение ПДД выписываются автоматически
- ИИ гугла ранжирует вам доступ к информации
- ...



Вы могли не замечать, но уже

- Вы смотрите видео, которое рекомендует вам youtube
- Вы передвигаетесь по маршруту, который вам построил навигатор
- В банке робот решает давать вам кредит или нет
- Штрафы за нарушение ПДД выписываются автоматически
- ИИ гугла ранжирует вам доступ к информации
- ...

Вывод

Искусственный интеллект уже принимает активное участие в вашей жизни!



Вопросы

- Что будет если ИИ начнет себя вести не совсем этично по отношению к некоторым группам населения?



Вопросы

- Что будет если ИИ начнет себя вести не совсем этично по отношению к некоторым группам населения?
- Кто будет за это отвечать?



Вопросы

- Что будет если ИИ начнет себя вести не совсем этично по отношению к некоторым группам населения?
- Кто будет за это отвечать?
- Как вообще определить, что такое происходит и происходит ли?



Вопросы

- Что будет если ИИ начнет себя вести не совсем этично по отношению к некоторым группам населения?
- Кто будет за это отвечать?
- Как вообще определить, что такое происходит и происходит ли?
- Как такое предотвратить?



Пример дискриминации

Пример

Объявление о высокооплачиваемой работе чаще показывается белым мужчинам среднего возраста



Пример дискриминации

Пример

Объявление о высокооплачиваемой работе чаще показывают белым мужчинам среднего возраста

Причины

- Потому что компания скорее всего наймет белого мужчину среднего возраста



Пример дискриминации

Пример

Объявление о высокооплачиваемой работе чаще показывается белым мужчинам среднего возраста

Причины

- Потому что компания скорее всего наймет белого мужчину среднего возраста
- Есть функция потерь, и ИИ ее оптимизирует. В данном случае дискриминирование эквивалентно эффективному показу рекламы



Пример дискриминации

Пример

Объявление о высокооплачиваемой работе чаще показывается белым мужчинам среднего возраста

Причины

- Потому что компания скорее всего наймет белого мужчину среднего возраста
- Есть функция потерь, и ИИ ее оптимизирует. В данном случае дискриминирование эквивалентно эффективному показу рекламы

Возможное решение

- Если в нанимающей компании есть квоты на дискриминируемые меньшинства, то такой алгоритм выдачи перестаёт быть эффективным

Вопрос

Кто отвественен за подобное поведение ИИ?



Отвественность перед обществом

Вопрос

Кто отвественен за подобное поведение ИИ?

Возможный ответ

Компания, которая производит такой продукт



Отвественность перед обществом

Вопрос

Кто отвественен за подобное поведение ИИ?

Возможный ответ

Компания, которая производит такой продукт

Возможное возражение

Компания, которая производит такой продукт, настраивается на данные, за которые она не отвечает



Ответственность перед обществом

Вопрос

Кто ответственен за подобное поведение ИИ?

Возможный ответ

Компания, которая производит такой продукт

Возможное возражение

Компания, которая производит такой продукт, настраивается на данные, за которые она не отвечает

Вопрос

Можем ли требовать от ИИ быть более этичным, чем само общество в целом?



Второй пример дискриминации

Пример

Банк не был готов дать кредит супругам вместе. Вместо этого предлагал заключить брачный договор. Хотя кредитный рейтинг жены был выше, чем у мужа. На 100% нельзя утверждать, что тут был случай дискриминации, но банк отказался прокомментировать своё решение.



Второй пример дискриминации

Пример

Банк не был готов дать кредит супругам вместе. Вместо этого предлагал заключить брачный договор. Хотя кредитный рейтинг жены был выше, чем у мужа. На 100% нельзя утверждать, что тут был случай дискриминации, но банк отказался прокомментировать своё решение.

Непрозрачность

В данном случае непрозрачность системы принятия решений позволяет ИИ дискриминировать.



Проблема

Во многих случаях ИИ работает как черный ящик. В некоторых областях такое поведение недопустимо

Проблема

Во многих случаях ИИ работает как черный ящик. В некоторых областях такое поведение недопустимо

Пример

Медицина

Третий пример дискриминации

Распознавание лиц

Известный факт, что на африканцах распознавание лиц работает хуже, чем на европейцах

Причины

Датасеты для европейцев более представительные.



- Этот тип ИИ вызывает наиболее оживленную дискуссию



- Этот тип ИИ вызывает наиболее оживленную дискуссию
- Самоуправляемые автомобили уже на дорогах Москвы



- Этот тип ИИ вызывает наиболее оживленную дискуссию
- Самоуправляемые автомобили уже на дорогах Москвы
- Известны случаи, когда люди гибли по вине самоуправляемых автомобилей



Самоуправляемые автомобили

- Этот тип ИИ вызывает наиболее оживленную дискуссию
- Самоуправляемые автомобили уже на дорогах Москвы
- Известны случаи, когда люди гибли по вине самоуправляемых автомобилей
- Известны случаи, когда люди были спасены ими



- Этот тип ИИ вызывает наиболее оживленную дискуссию
- Самоуправляемые автомобили уже на дорогах Москвы
- Известны случаи, когда люди гибли по вине самоуправляемых автомобилей
- Известны случаи, когда люди были спасены ими
- Перед ИИ стоит множество этических вопросов: как вести в той или иной ситуации



- Этот тип ИИ вызывает наиболее оживленную дискуссию
- Самоуправляемые автомобили уже на дорогах Москвы
- Известны случаи, когда люди гибли по вине самоуправляемых автомобилей
- Известны случаи, когда люди были спасены ими
- Перед ИИ стоит множество этических вопросов: как вести в той или иной ситуации

Проблема

Рынок труда изменится при всеобщем внедрении этой технологии.



Проблема

Добавление малозаметного шума к изображению ставит нейронную сеть в тупик

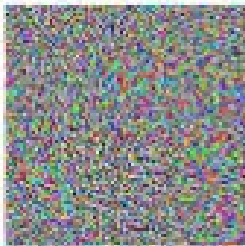


x

"panda"

57.7% confidence

+ .007 ×



$\text{sign}(\nabla_x J(\theta, x, y))$

"nematode"

8.2% confidence

=



$x + \epsilon \text{sign}(\nabla_x J(\theta, x, y))$

"gibbon"

99.3 % confidence

Атаки на ИИ в реальной жизни

Проблема

Первые три знака распознаются как «Ограничение скорости 45», а последний — как знак «STOP»



Атаки на ИИ в реальной жизни

Проблема

Специальные элементы одежды (очки, шапки и т.д.) позволяют обмануть систему распознавания лиц



Некоторые факты

- Сотовые операторы продают данные о местоположении пользователей



Некоторые факты

- Сотовые операторы продают данные о местоположении пользователей
- Говорят, что Амазон по вводу запроса может определять беременность



Некоторые факты

- Сотовые операторы продают данные о местоположении пользователей
- Говорят, что Амазон по вводу запроса может определять беременность
- Ваши поисковые запросы влияют на ваши рекомендации



Некоторые факты

- Сотовые операторы продают данные о местоположении пользователей
- Говорят, что Амазон по вводу запроса может определять беременность
- Ваши поисковые запросы влияют на ваши рекомендации

Вопрос

Кому должны принадлежать данные, производимые людьми? Какие данные являются приватными, а какие нет?



Обеспечение безопасности

Для обеспечения безопасности государство устанавливает камеры видеонаблюдения

Обеспечение безопасности

Для обеспечения безопасности государство устанавливает камеры видеонаблюдения

Вопрос

Где границы того, что должно быть доступно государству, а что нет?

- Secure AI. Контролируемость и управляемость систем ИИ
- Explainable AI. Прозрачность и предсказуемость функционирования
- Reliable AI. Стабильность и надежность систем ИИ
- Responsible AI. Ответственное применение ИИ
- Fair AI. Непредвзятый ИИ



- Жизнь современного человека невозможно представить без ИИ
- Сейчас нет общего понимания и согласия, каким должен быть этический ИИ
- Но уже многие компании задумываются об этике ИИ
- Исследователи и разработчики должны больше обращать внимание на проблемы этики ИИ



Спасибо за внимание!

