

Искусственный интеллект в кибербезопасности:

защита будущего или новые угрозы?

чл.-корр. РАН, д.т.н., профессор
Зегжда Д.П.

д.т.н., доцент
Полтавцева М.А.



ПОЛИТЕХ

Институт компьютерных
наук и кибербезопасности

ИИ и Кибербезопасность



Поиск уязвимостей в логике приложений

Автоматизированные атаки / тестирование

Анализ поведения

Генерация / обнаружение спама

Дипфейки и их обнаружение
ИИ внутри ВПО...



Использование
ИИ
в атаках и защите

Развитие известных атак на ИИ

Появление первых атак, нацеленных на
большие языковые модели (LLM)

Появление атак, нацеленных на ИИ,
которые уже используются в бизнес-
процессах компаний

Защита цепочек поставок

Проверка
целостности

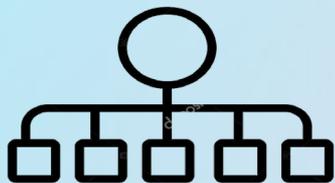
Защита
систем ИИ

Атаки на ИИ-
системы

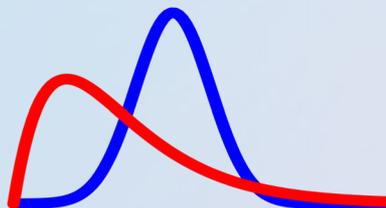
Доверие
к системам ИИ

Развитие атак на конфиденциальные
ИИ системы (федеративное обучение)

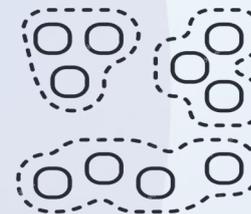
Типовые задачи ИИ в Кибербезопасности



Классификация



Регрессия



Кластеризация

ИИ в задачах кибербезопасности:

- обнаружение атак
- обнаружение фишинга
- антиспам-системы
- антивирусы
- биометрическая аутентификация
- обнаружение аномалий, ботнетов
- DLP-системы

Северная Америка
крупнейший рынок (2019г)



2019
8.6 млрд

Объем рынка ИИ-решений для
кибербезопасности

Азиатско-Тихоокеанский
регион

самый быстрорастущий
рынок (2020-2030 гг)

2030
101.8
млрд

Прогнозируемый
объем рынка ИИ-
решений для
кибербезопасности

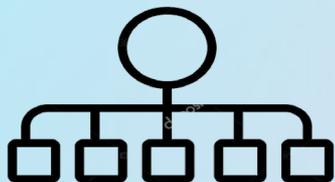
Что мы знаем об атаках на ИИ?

Модель угроз



СВОЙСТВО	ADVERSARIAL ML THREAT MATRIX (MITRE)
ОБЪЕКТ ИССЛЕДОВАНИЯ	Система машинного обучения, подверженная состязательным атакам
СООТВЕТСТВИЕ НОРМАТИВНЫМ ТРЕБОВАНИЯМ	Нет (представляет собой реестр типовых атак, построенный по сводкам, полученным от 24 ИТ-компаний)
Полнота модели	Нет (построена на основе коллекции зафиксированных примеров реализации состязательных атак на системы ИИ)
Анализ природы угроз, способов их реализации	Нет (фиксация случаев атак, типизация по базовым классам)
Анализ возможных механизмов защиты	Нет (декларированы лишь сами угрозы, техники защиты требуют дополнительного изучения)
РАСШИРЯЕМОСТЬ МОДЕЛИ	Нет (фиксированная структура матрицы по этапам атак на ИИ)

Уязвимости типовых решений

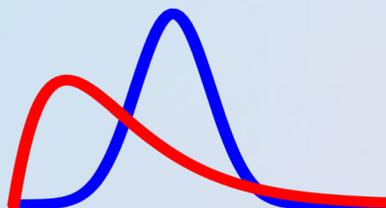


Классификация

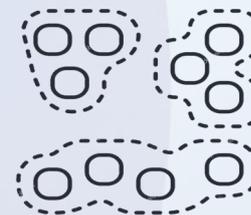
Через добавление дополнительно сгенерированных значений изменяется область распределения значений (линия предсказаний);
⇒ изменяется область допустимых значений параметра для определения класса
⇒ классификатор будет давать неверный ответ
системы классификации вирусов, вредоносной активности и т.п.

Hint: Злоумышленник может заранее произвести исследование по выявлению классов для создания дальнейших «отравленных» данных

Hint: Если злоумышленнику известны параметры классификатора – процесс корректировки еще больше упрощается



Регрессия



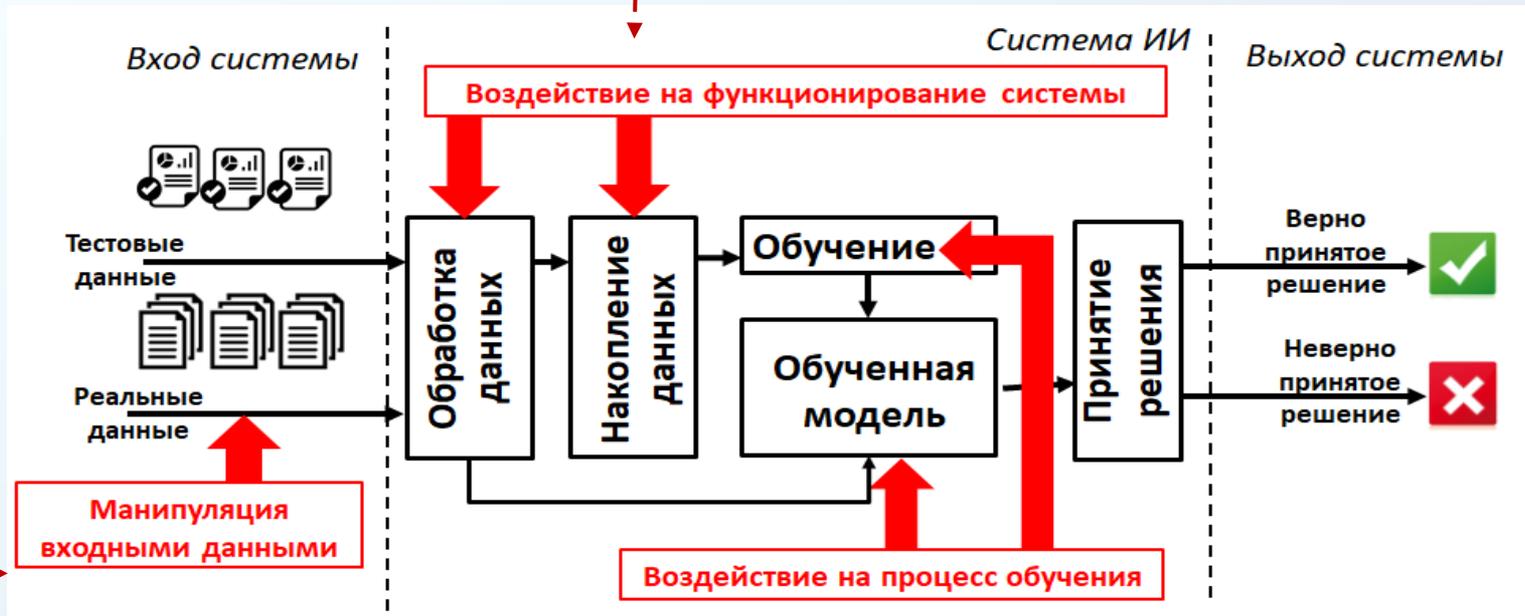
Кластеризация

Расширить границы кластеров – аналогично.

Также можно создать дополнительные кластеры путем создания группы объектов со схожими характеристиками.

Чувствительные компоненты систем ИИ

- **модули системы** – нарушение корректной работы системы, сбой в работе одного из ее компонентов или модулей

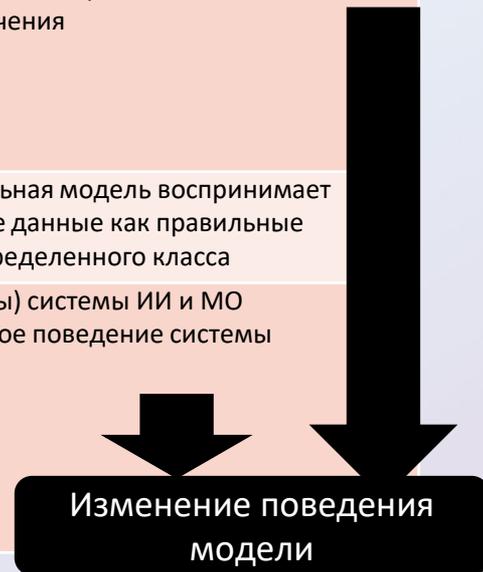


- **наборы данных** – воздействие на реальные входные данные

- **механизм обучения** (недостатки самого принципа «обучения модели») – как воздействие на сам процесс обучения модели, так и на уже обученную модель

Как заставляют ИИ «изменять» себе

	Локализация	Способ эксплуатации	Результат эксплуатации
Уязвимости входа модели	Обучающие данные , поступающие в вычислительную модель Обученная модель ИИ и МО Обучающий модуль	<ul style="list-style-type: none"> - Воздействие на обучающие данные до того, как они попадут в систему ИИ и МО (отравление данных) - Поиск «слепых зон» в обученной вычислительной модели, генерация данных, попадающих в эти слепые зоны - Генерация большого объема специальных входных данных, искажающих поле экспертизы системы ИИ и МО 	Неправильное обучение вычислительной модели Наложение конкретного решения на вычислительную модель Изменение логики работы модуля обучения
	Реальные (входные) данные , поступающие в уже обученную вычислительную модель ИИ и МО	Реализация сопоставительных примеров : искажение	Вычислительная модель воспринимает искаженные данные как правильные данные определенного класса
Уязвимости модели	Ввод данных в вычислительную модель Предварительная обработка и обработка данных Хранилище данных	<ul style="list-style-type: none"> - Перехват и модификация ввода данных в модель ИИ - Поиск и эксплуатация уязвимостей в реализации обработки/обработки данных, изменение параметров модели (метод оптимизации, гранулярность, масштаб и т.д.) - Удаление фрагментов данных из хранилища 	Сбои (отказы) системы ИИ и МО Некорректное поведение системы ИИ и МО



Классы актуальных угроз

направленных на снижение достоверности результатов

По характеру
воздействия на
модель

угроза отравления
угроза искажения
исследовательские угрозы

По стратегии возможной
реализации угрозы по
отношению к модели

угроза «белого ящика»
угроза «серого ящика»
угроза «черного ящика»

По способу
воздействия

через вмешательство в
обучающую выборку
через ошибки модели

По специфике
воздействия

нецеленаправленная
целенаправленная

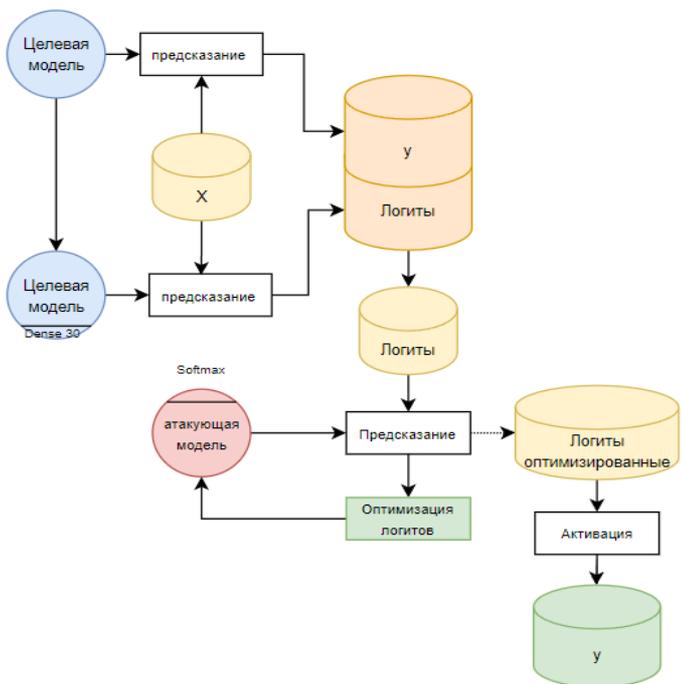
По нарушаемому
аспекту безопасности

нарушение целостности модели
нарушение доступности модели
нарушение конфиденциальности
модели

↓
Модель угроз
в соответствии методикой ФСТЭК

Центр компетенций НТИ
Технологии доверенного
взаимодействия

Конфиденциальность данных и модели



Метод «Умный шум»:

- гибкое задание количества шума, налагаемого на выходе защищаемой модели;
- градиентный шаг для оптимизации шума;
- делает невозможным для нарушителя вывод о принадлежности отдельных образцов набору данных, на котором обучалась атакуемая модель (атака «вопрос о вхождении»)

Метод на базе контроля расстояний в запросах

- выявляет поток запросов, поступающих от нарушителя при зондировании атакуемой им модели, по отклонению от нормального распределения данных, характерного для поведения обычных пользователей.
- обеспечивает выявление и затем блокировку нарушителя, исследующего атакуемую модель
- не воздействует на механизмы защищаемой модели

Результат

Центр компетенций НТИ
Технологии доверенного
взаимодействия

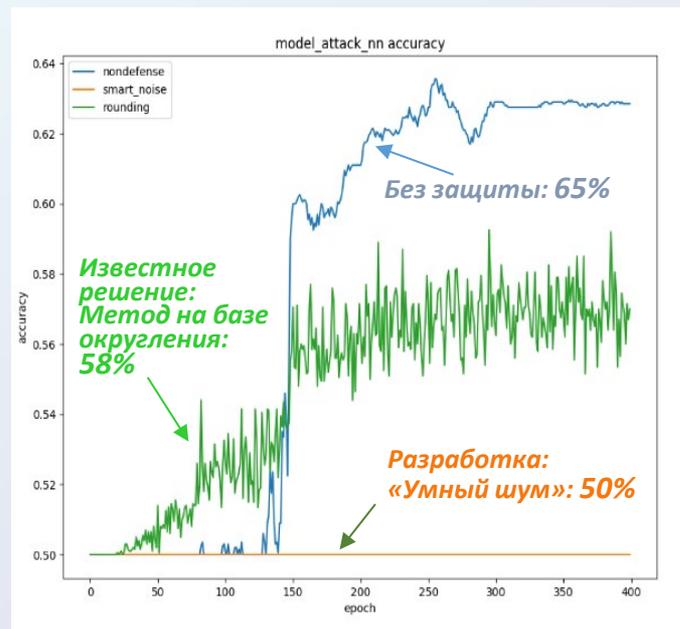
Сравнение с альтернативными решениями

Название	Воздействие на обучающие данные	Воздействие на параметры модели	Воздействие на архитектуру модели	Воздействие на функцию потерь	Воздействие на выходные данные	Влияние на точность модели	Сложность	Уровень безопасности
Отсев		+				Повышает	Низкая	Низкая
Ансамблевое обучение ML			+			Повышает	Средний	Средний
Ансамблевое обучение DL			+			Повышает	Высокая	Высокий
Усложнение архитектуры модели			+			-	Высокая	Высокий
Регуляризация L_1 и L_2				+		Повышает	Средний	Низкая
Эластичная сеть				+		Повышает	Средний	Высокий
Стохастическое округление					+	Снижает	Низкая	Средний
Округление вывода					+	Снижает	Низкая	Средний
Отсечение					+	Снижает	Низкая	Высокий
Округление вывода Class labels					+	Снижает	Низкая	Высокий
Умный шум					+	Сохраняет	Средний	Высокий

Оценка точности метода защиты на базе контроля расстояния в запросах

Набор данных	Частота ложных срабатываний (метрика FPR)	Количество запросов пользователя до обнаружения атаки	
		Атака на модели с совместным обучением	Атака с использован ием GAN
MNIST	0,0	5 560	120
GTSRB	0,1	5 020	430

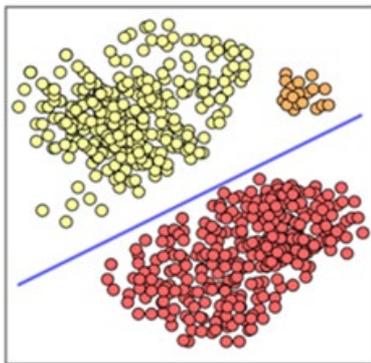
Оценка точности атакующей модели-классификатора нарушителя, который определяет принадлежность данных набору данных целевой модели



Примеры угроз системам ИИ

направленных на снижение достоверности результатов

Угроза отравления:



Отравленные данные

Данные, которые должны были бы быть помечены как данные класса 2, но были помечены злоумышленником как данные класса 1



Класс 1

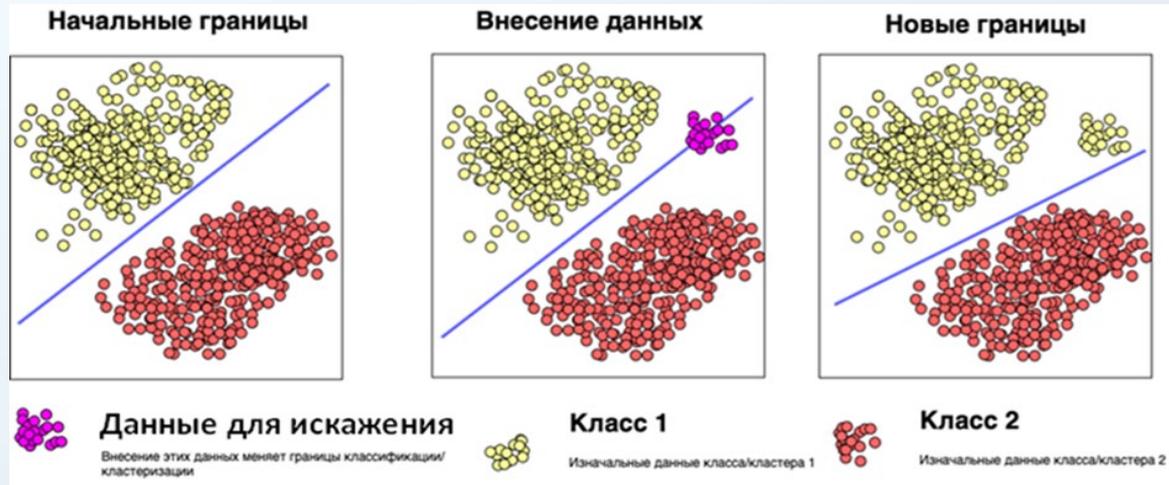
Изначальные данные класса/кластера 1



Класс 2

Изначальные данные класса/кластера 2

Угроза искажения:



Данные для искажения

Внесение этих данных меняет границы классификации/кластеризации



Класс 1

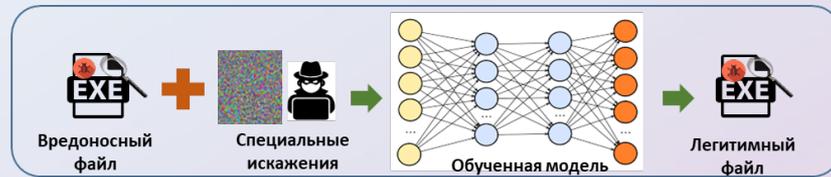
Изначальные данные класса/кластера 1



Класс 2

Изначальные данные класса/кластера 2

Пример
деструктивного
воздействия
искажений
на классификатор
компьютерных атак



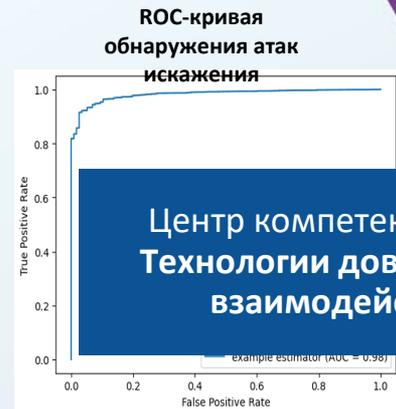
ПОЛИТЕХ
Санкт-Петербургский
политехнический университет
Петра Великого



Институт компьютерных
наук и кибербезопасности

Обнаружение искажений

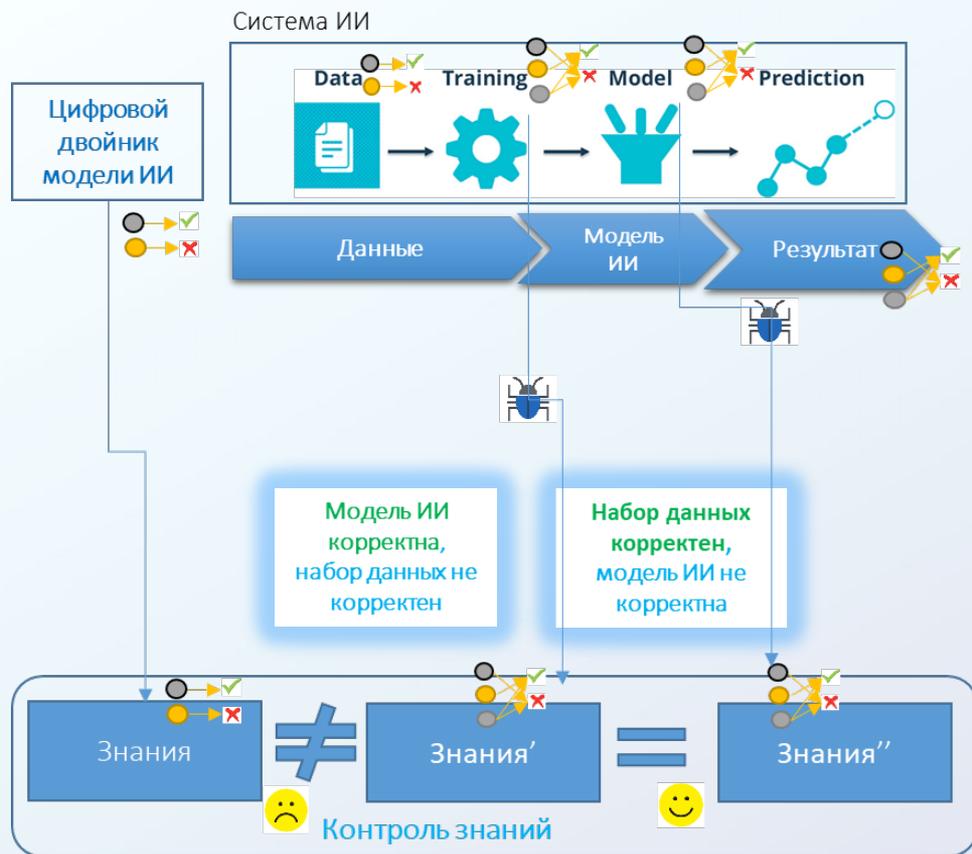
- ✓ незащищенная модель (точность 97%) при воздействии искажений демонстрирует падение точности до 49%. Точность обнаружения атак искажения – **98%**, точность защищенной модели при воздействии атак искажения – **>91%** (вместо 49%);
- ✓ метод защиты **не зависит от вычислительной модели ИИ**, безопасность которой он обеспечивает;



Центр компетенций НТИ
Технологии доверенного
взаимодействия

КОНКУРИРУЮЩИЕ РЕШЕНИЯ	РЕСУРСОПОТРЕБЛЕНИЕ	ТОЧНОСТЬ ОБНАРУЖЕНИЯ / КЛАССИФИКАЦИИ
Обнаружение искажений с помощью метода k-ближайших соседей (kNN) с маркировкой состязательных образцов [1]	Высокое (обучение алгоритма, работа классификатора)	97% / Не реализует
Поддержание устойчивости модели с помощью генеративной нейросети GAN для генерации состязательных образцов [2]	Высокое (работа нейросети-генератора состязательных образцов)	Не реализует / 70%
Разработанный в рамках НТИ метод, агрегирующий обнаружение искажений на основе анализа тестовых данных и защитную дистилляцию	Низкое (нет необходимости формировать обучающий набор данных, обучать и переобучать нейросети)	98% / >91%

Доверие к результату



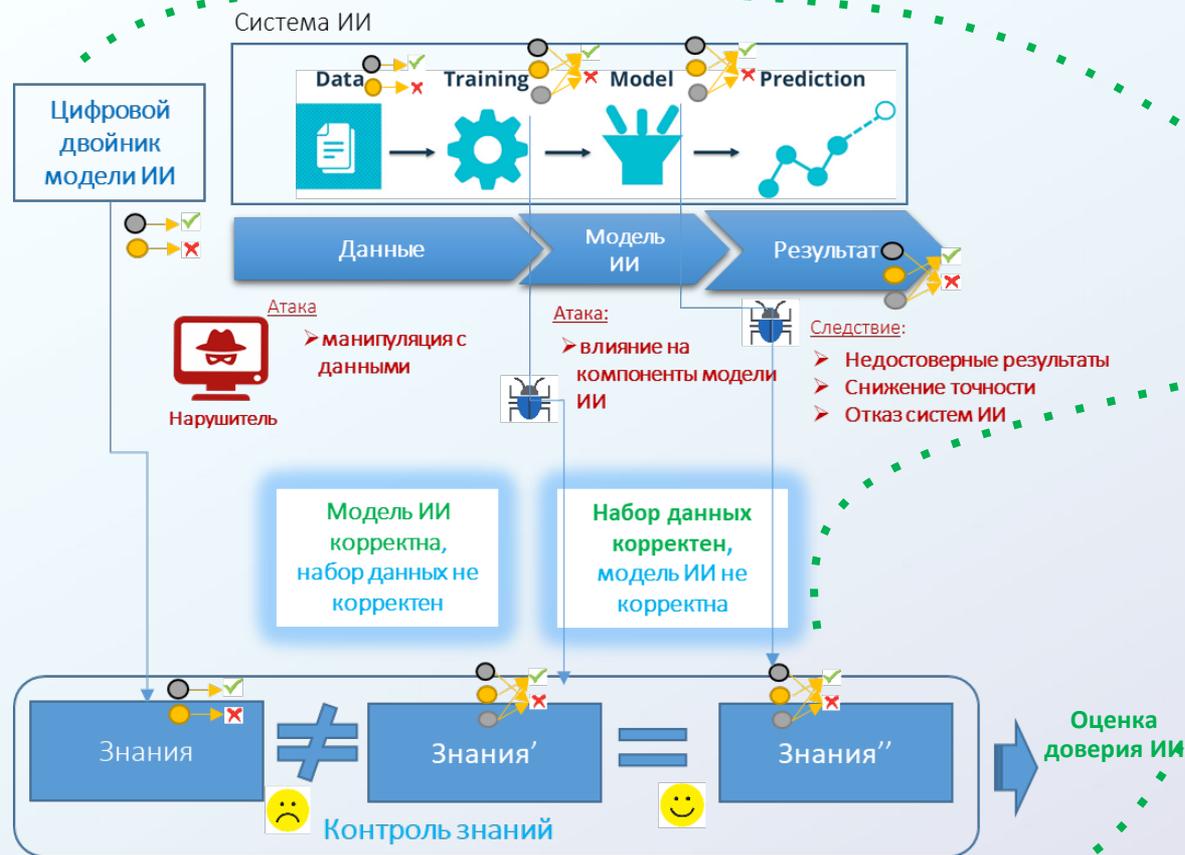
- Самотестирование
- Кросс-валидация

Гибридизация взаимодополняющих методов и алгоритмов

- Сокращение размера «враждебного пространства данных» для систем ИИ.
- Предотвращение использования мобильности (переноса) атак на вычислительные модели систем ИИ.
- Затруднение действий нарушителей ИИ.

- Контроль цепочек поставок данных в системах ИИ.
- Выявление искажающих данных в схемах обучения систем ИИ.
- Защитная дистилляция и состязательная тренировка вычислительной модели.
- Быстрая кросс-валидация решений, полученных системами ИИ и машинного обучения.
- Аудит и трассировка событий на уровне алгоритмов ИИ и МО, с помощью которых возможно быстро проверить состояние классификаторов, которое привело к недостоверному решению

Безопасность ИИ в киберсреде



(1) **Использование «двойника» ИИ:** мониторинг поведения исходной модели ИИ и контроль результатов по двойнику

(1) **«Песочница» ИИ:** проверка результатов на копии данных ДО внесения их в незагрязненный набор данных

(1) **Обучение с переносом:** обучение на выборках из другой области и сверка результатов с уязвимой моделью ИИ

ФГАОУ ВПО «СПБПУ»

**ИНСТИТУТ КОМПЬЮТЕРНЫХ НАУК
И КИБЕРБЕЗОПАСНОСТИ**

Зегжда Дмитрий Петрович

чл.-корр. РАН, д.т.н., профессор

Полтавцева Мария Анатольевна

д.т.н., доцент

poltavtseva@ibks.spbstu.ru

Главный учебный корпус, к. 173

Политехническая ул., 29, Санкт-Петербург 195251

Тел: +7 (812) 552-76-32