

О возможности косвенных утечек данных и концепция противодействия им

УДК 519.724.6; 004.89

Доверенность систем искусственного интеллекта (ИИ) обусловлена, в частности, выполнением комплекса мероприятий по технической защите информации, которые должны быть усилены защитой от новых видов атак, направленных на извлечение защищаемых данных из совокупности отчетов (результатов). Утечки, связанные с таким видом атак, идентифицируются как «косвенные утечки». Для систем искусственного интеллекта определение возможности косвенной утечки данных является актуальной научной задачей. В статье показана ограниченность методов дифференциальной конфиденциальности при их применении к системам ИИ. Предложен новый метод анализа совокупности запросов, основанный на применении матриц Якоби.

Ключевые слова: косвенные утечки, доверенность систем искусственного интеллекта, матрицы Якоби, дифференциальная конфиденциальность, статистическая база данных, защита информации в системах искусственного интеллекта, слепая обработка данных

Валерий Аркадьевич Конявский,
доктор технических наук
konyavskiy@gospochta.ru

Данила Дмитриевич Агапитов
agapitov.dd@phystech.edu

Московский физико-технический институт
(национальный исследовательский университет)

Введение

Традиционные подходы к защите информации ориентированы преимущественно на управление доступом в соответствии с обоснованной

Интуиция подсказывает, где искать истину, но лишь интеллект способен отличить возможное от невозможного.
Макс Борн, «My Life and Views»

политикой информационной безопасности [1] и обеспечивают предотвращение утечек защищаемых данных посредством предотвращения несанкционированного доступа к ним. Такие утечки можно считать прямыми, и меры по защите данных от этого вида утечек хорошо изучены и активно регулируются¹. Нужно от-

En On the Possibility of Indirect Data Leaks and the Concept of Counteracting Them

V. A. Konyavskiy,
PhD (Eng., Grand Doctor)
konyavskiy@gospochta.ru

D. D. Agapitov
agapitov.dd@phystech.edu

Moscow Institute of Physics and Technology.

The trustworthiness of artificial intelligence systems is conditioned, in particular, by the implementation of a set of measures for the technical protection of information, which must be strengthened by protection against new types of attacks aimed at extracting protected data from a set of reports (results). Leaks associated with this type of attack are identified as «indirect leaks». For artificial intelligence systems, determining the possibility of indirect data leakage is a pressing scientific challenge. The limitations of differential privacy methods when applied to artificial intelligence systems have been demonstrated. To identify the possibility of indirect data leaks, a new method for analyzing a set of queries based on the use of Jacobian matrices is proposed.

Keywords: Indirect leaks, trustworthiness of artificial intelligence systems, Jacobian matrices, differential privacy, statistical database, information security in artificial intelligence systems, blind data processing

¹ См., например, приказы ФСТЭК России: от 11 февраля 2013 года № 17 «Об утверждении требований о защите информации, не составляющей государственную тайну, содержащейся в государственных информационных системах»; от 18 февраля 2013 года № 21 «Об утверждении состава и содержания организационных и технических мер по обеспечению безопасности персональных данных при их обработке в информационных системах персональных данных»; от 14 марта 2014 года № 31 «Об утверждении требований к обеспечению защиты информации в автоматизированных системах управления производственными и технологическими процессами на критически важных объектах, потенциально опасных объектах, а также объектах, представляющих повышенную опасность для жизни и здоровья людей и для окружающей природной среды».

метить, что подход управления доступом полностью оправдывает себя при защите данных в корпоративных системах, но требует значительного развития в системах открытого типа [2], а также в системах совместной обработки привлеченных данных, то есть в системах, в которых для улучшения качества моделей обрабатываются данные из различных источников, в том числе данные, накапливаемые разными операторами персональных данных.

При развитии «экономики знаний» особое значение приобретают системы искусственного интеллекта (СИИ), позволяющие на основе обработки большого числа статистических данных строить прогнозные модели приемлемого качества [3, 4]. Источником данных для СИИ зачастую служат статистические базы данных (БД) – то есть базы, предназначенные для статистической обработки данных. Для предотвращения утечек конфиденциальной информации из баз такого типа запросы к таковым ограничивают, предоставляя доступ только к агрегированным данным, и запрещая доступ к отдельным записям, предотвращая утечку данных об отдельном субъекте.

Пусть запросы к статистической БД носят, как следует из названия, статистический характер: то есть узнать статистические параметры можно, а аналитические, касающиеся одного конкретного субъекта, – нельзя. В этом случае говорят, что обеспечивается дифференциальная конфиденциальность субъектов данных.

Свойство дифференциальной конфиденциальности обеспечивает сохранение конфиденциальности одного отдельного субъекта данных при неконтролируемых запросах к статистической БД.

Однако работа с системами обработки больших данных продемонстрировала [5], что на основе анализа значительного количества результатов обработки зачастую можно восстановить начальные наборы данных, в том числе – защищаемые данные. При этом фактически будет зафиксирована утечка данных, хотя нарушения правил доступа не будет. Такой вид утечки

данных можно назвать косвенной утечкой.

Косвенная утечка возникает при успешном восстановлении значений из исходных наборов данных по совокупности вычисленных параметров.

В результате косвенных утечек вполне могут быть нарушены права субъекта на неприкосновенность частной жизни. Поскольку косвенная утечка, нарушающая права субъекта, возникает при изучении отличий результатов вычислений при различных наборах исходных данных, говорят, что в этом случае нарушается дифференциальная конфиденциальность [6]. Накопление данных, достаточных для идентификации отдельного субъекта, может носить случайный характер, а может быть целенаправленным. В этом случае говорят о дифференциальной атаке.

Простейший пример косвенной утечки – использование близких измерений: из средней зарплаты n сотрудников и $n + 1$ сотрудников не составит труда вычислить зарплату сотрудника с номером $n + 1$.

Появление и широкое распространение СИИ, статистических баз данных, переход к связанным данным резко увеличивают риски косвенных утечек. При этом возможности предупреждения и противодействия им еще не разработаны в достаточной мере.

Возникает противоречие между потребностью в обеспечении противодействия косвенным утечкам данных и отсутствием обоснованных научных методов. Разрешение данного противоречия обеспечивается разработкой методов анализа зависимости функций в информационных технологиях обработки данных.

Статистическая база данных и дифференциальная конфиденциальность

Будем называть двумерную таблицу (размерностью) статистической базой данных (СБД), если каждая ее строка описывает некую сущность (физическое лицо, юридическое лицо, момент времени и т. д.), а каждый столбец, кроме первого, являю-

щегося индексным, – некоторое свойство сущности. Каждая сущность должна иметь одинаковую однозначную природу, как и каждое свойство сущности.

Пусть x^i – i -я сущность в статистической базе данных;
 x_j^i – j -е свойство i -й сущности в статистической базе данных;

$$\forall i \in \{1, 2, \dots, N\}, \forall j \in \{1, 2, \dots, n\}.$$

В принципе, индексным может быть любой столбец и даже несколько столбцов, но это не имеет особого значения для идентификации понятий.

Свойства сущностей могут быть количественными и качественными, в зависимости от множества значений в соответствующих столбцах.

Свойство под номером j будем называть количественным, если множество значений включено в множество действительных чисел, иными словами:

j -е свойство – количественное по определению, если:

$$\forall i \in \{1, 2, \dots, N\}: x_j^i \in \mathbb{R}.$$

В случае

$$\exists i \in \{1, 2, \dots, N\}: x_j^i \notin \mathbb{R}.$$

j -е свойство будем называть качественным.

Отметим еще раз: природа сущностей должна быть однозначно определена. Если это условие не выполняется, таблица не является статистической базой данных.

Если сущность – физическое лицо, а свойства сущностей – некоторые статистические, временные или другие параметры – как количественные, так и качественные, то актуальной становится защита данных от утечки – распространения, обеспечения и предоставления доступа, не санкционированного защитными механизмами, политикой информационной безопасности (ПИБ), нормативными правовыми актами (НПА).

Для контроля запросов к СБД обычно предоставляется доступ только к агрегированным данным, а не к отдельным записям. С этой целью устанавливается правило применения агрегированных запросов: на-

пример, SUM, COUNT, AVG и т. д. Можно также возвращать неточные относительно запроса подсчеты (например, вместо параметров 101 записи, удовлетворяющей запросу, указываются параметры 90–120 записей). Иногда стараются возвращать не точные значения конфиденциальных данных, а только диапазоны, к которым они принадлежат: например, при величине дохода условно в 37 тыс. руб. можно указать диапазон 35–40 тыс.

Одним из методов обеспечения дифференциальной конфиденциальности является добавление шума (случайной величины) ко всем значениям свойств сущности или к результатам выполнения запросов. Этот метод широко применяется, хотя и снижает качество информации, содержащейся в изначальном наборе данных [7].

Есть мнение [8], что метод дифференциальной конфиденциальности обеспечивает защиту от выделения предикатом. Действительно, его применение позволяет во многих случаях предоставлять статистику (агрегированные данные) достаточной точности из БД, сохраняя при этом высокий уровень конфиденциальности данных. Однако вывод о достаточности методов дифференциальной конфиденциальности касается лишь статистических данных (да и то не в полной мере), и представляется некорректным распространить его на использование аналитических данных, что совершенно необходимо в СИИ.

Даже ухудшая данные добавлением шума, нельзя предотвратить утечку при неконтролируемых механизмах атаки. Если модификация механизма конфиденциальности известна злоумышленнику, то что может помешать ему изменить параметры выделяющего предиката? Например, при добавлении шума изменить диапазон выделения на аналогичную величину.

Эти и подобные механизмы защиты не решают задачи обеспечения конфиденциальности защищаемых данных, так как при незначительных усилиях злоумышленники могут использовать комбинацию агрегированных запросов для по-

лучения информации об отдельном субъекте.

Утечка защищаемых данных осуществляется, если можно сформулировать набор критериев, сужающих группу до одного лица. В этом случае говорят, что произошло выделение (или идентификация) физического лица. Возможность выделения предикатом данных отдельного физического лица, по сути, создает предпосылки к распространению персональных данных. Возможность идентификации, и, следовательно, утечки персональных данных физического лица означает существование предиката, позволяющего однозначно выделить физическое лицо из набора данных вне зависимости от природы предиката.

Идентификация является зависимой от контекста и набора данных.

Предикат «Зарплата Кузнецова больше 25 000 руб. в месяц» вряд ли идентифицирует конкретного субъекта в масштабах страны, где фамилия «Кузнецов» – самая распространенная. Чего не скажешь о предикате «Зарплата Кузнецова больше 5 000 000 руб. в месяц». В этом случае предикату будет удовлетворять едва ли несколько человек, а скорее всего, лишь один. Тогда вполне можно будет считать, что идентификация состоялась.

Методы дифференциальной конфиденциальности обеспечивают защиту от выделения предикатом при предоставлении статистики, но не обеспечивают достаточного уровня конфиденциальности при использовании аналитических данных. Представляется, что следует искать другой механизм, который учитывал бы применяемые информационные технологии обработки данных.

Технология слепой обработки данных

В [9] предложена технология слепой обработки данных (ТСОД), позволяющая разрешить существенную часть проблем, возникающих в области защиты информации при обработке больших данных в системах искусственного интеллекта. ТСОД основана на сочетании традиционных механизмов эшелонированной

защиты информации, способов создания и поддержки доверенной среды функционирования криптографии, разделения функций участников обработки данных, ограничения информационных технологий обработки данных только доверенными технологиями (ИТ-конвейерами), с учетом публикуемых отчетов по доверенным аналитикам и произвольным пользователям. При этом учитывается риск утечек, возникающих при объединении данных – как прямых, так и косвенных. Вводятся в научный оборот понятия слепой обработки данных, неизвлекаемых данных и моделей, привлекаемых данных.

Представим теперь СИИ или систему машинного обучения как совокупность информационных технологий.

Обработка данных в информационной системе выполняется в соответствии с зафиксированной информационной технологией, понимаемой как последовательность информационных операций преобразования исходных данных в данные результатов [10]. Тогда на входе СИИ мы обнаружим статистическую БД, затем последовательность операций трансформации данных из входных наборов данных в модель и далее – отчеты по прогнозу. При этом на входе может быть не одна база, а несколько, если мы обрабатываем привлекаемые данные. Именно поэтому и возникает потребность в слепой обработке данных, так как любой оператор персональных данных может ознакомиться только с теми данными, относительно которых именно у него есть информированное согласие субъекта этих данных. ТСОД позволяет разрешить эту коллизию при правильно определенной совокупности ИТ-конвейеров, так как в системах слепой обработки данных известны набор ИТ-конвейеров для подготовки отчетов по наборам входных данных и совокупность публикуемых отчетов.

Если не существует механизма, выделяющего признаки субъекта в статистической БД по значениям отчетов, то дифференциальная конфиденциальность не будет наруше-

на. Если же косвенная утечка возможна, то возникающие риски нужно оценить, и принять меры по их блокированию (например, изменить состав отчетов в процессе разработки модели) и/или компенсации методами страхования информационных рисков [11]. Таким образом, актуальной научной задачей является **определение возможности косвенной утечки данных из входного набора в СИИ при применении конкретных ИТ-конвейеров.**

Постановка задачи

Рассматриваем работу с привлекаемыми данными. Здесь для повышения качества модели используются данные, которые собираются разными операторами: банками, ритейлом, операторами связи, страховыми компаниями и др. Таким образом, данные вначале обогащаются (в нашем случае – за счет слияния наборов данных (НД) разных операторов), а затем используются в процессах разработки модели.

На каждом из этих этапов есть особенности работы с наборами данных. Так, обогащенный набор данных (ОНД) формируется объединением имеющихся НД различных операторов. В процессе разработки модели доступ к данным для разработчика может быть ограничен классическими механизмами защиты информации. Однако при этом не исключена возможность вычисления защищаемых данных (и/или способа доступа к ним) с помощью других механизмов. Таким образом, **необходимо определить, возможна ли косвенная утечка данных из ОНД при разработке модели.**

Применение матриц Якоби и их расширение для выявления возможных косвенных утечек данных

Предполагается использовать привлекаемые данные, накапливаем

ые k операторами данных (ОД). Обозначим их как $ОД_i, i = 1, k$. $ОД_i$ накапливает набор данных $НД_i, i = 1, k$. Очевидно, что со временем $НД_i$ постоянно растет в соответствии с появлением отношений с новыми субъектами данных, то есть объем набора данных не уменьшается. Если НД касается m субъектов и по каждому субъекту содержит p признаков, то очевидно, что в $НД_i$ содержится $d_i = m_i \cdot p_i$ данных.

Вначале необходимо сформировать ОНД. Пусть для определенности $НД_i$ называется обогащаемым набором данных при $i = 1$, то есть $НД_1$ является обогащаемым НД. В результате объединения $НД_1$ с $НД_i, i = 2, k$ формируется ОНД, который содержит $d \leq \sum_{i=1}^k m_i \cdot p_i$ данных². Очевидно также, что $m \leq \sum_{i=1}^k m_i$, а $p \leq \sum_{i=1}^k p_i$, и $m \ll p$. Здесь m – число строк (субъектов) в ОНД, p – число столбцов (признаков) в ОНД.

Заметим, что обогащение называют горизонтальным при $p > p_1$, и вертикальным при $m > m_1$. При этом ОНД содержит $d = m \cdot p$ данных. Обозначим всю совокупность этих данных: $X = \{x^i\}, i = 1, d$.

Далее происходит выбор типа модели и ее обучение. Здесь разработчик модели выбирает из множества типовых моделей лучшую с использованием ОНД. Технологии машинного обучения хорошо описаны в литературе³. Отметим лишь, что оценка качества модели выполняется на основе вычисления значений ограниченного набора показателей⁴. Обозначим совокупность значений показателей: $Y = \{y^i\}, i = 1, M$.

Для вычисления этих значений используются Минформационных технологий, понимаемых как зафиксированная последовательность операций над данными – ИТ-конвейер, в соответствии с [9]. По сути, такое преобразование может быть описано некоторой функцией. Введем обозначение этих функций: $F = \{f^i\}, i = 1, M$.

Теперь, с использованием введенных обозначений, можно описать деятельность разработчика модели так: оценить качество модели по оценкам:

$$Y = F(x). \quad *$$

Это формализация стандартной деятельности. С точки зрения технической защиты информации в части косвенных утечек формулируем постановку:

Возможно ли узнать $x_j^i, x^i \in X$ для каких-либо $i \in \{1, d\}, j \in \{1, p\}$, зная F и Y ?

В виде * задача сводится [12] к вычислению ранга матрицы Якоби.

Метод обнаружения косвенных утечек на основе матриц Якоби

Задача обнаружения косвенной утечки для системы ИТ-конвейеров схожа с задачей определения функциональной зависимости системы функций:

1) для обнаружения косвенной утечки требуется определить, можно ли на основании результатов выполнения ИТ-конвейеров (зная вычисленные оценки из *) вычислить исходные данные

2) для определения функциональной зависимости требуется определить, может ли одна функция из набора быть выражена в качестве преобразования над остальными функциями этого же набора.

Следовательно, построив изоморфизм между этими задачами, получим возможность решить первую, решив вторую. Для этого представим набор конкретных ИТ-конвейеров, для которых требуется проверить наличие косвенных утечек, в виде набора функций:

$$\begin{aligned} y_1 &= f_1(x_1, x_2, \dots, x_n) \\ y_2 &= f_2(x_1, x_2, \dots, x_n) \\ &\dots \\ y_m &= f_m(x_1, x_2, \dots, x_n) \end{aligned} \quad (1)$$

При этом $x_i, i = 1, n$ – значения параметров из набора данных.

² Не большие, так как некоторые признаки в разных наборах могут совпадать, и дублировать их нет смысла.

² Воронцов К. В. Математические методы обучения по прецедентам (теория обучения машин) [Электронный ресурс]. – URL: https://mathprofi.com/uploads/files/4210_f_41_lekcii-voronova-k.v.-mashinnoe-obuchenie.pdf?key=77680f456097da8ff038c97ac64842b8/ (дата обращения: 05.01.2024).

³ Разработчик модели не может получить доступ к данным, так как доступ к персональным данным и ознакомление с ними без информированного согласия субъекта является нарушением законодательства.

Пусть значение y_i однозначно определяется значениями остальных функций, то есть справедливо

$$y_i = \phi(y_1, y_2, \dots, y_{i-1}, y_{i+1}, \dots, y_m). \quad (2)$$

В этом случае говорят, что функции y_i зависят от остальных, а функции из (1) называются зависимыми. В другом случае функции из (1) называют независимыми в рассматриваемой области.

Полагая эти функции непрерывно дифференцируемыми на всей области определения, рассмотрим матрицу Якоби, составленную из частных производных этих функций по всем независимым переменным:

$$J = \begin{pmatrix} \frac{dy_1}{dx_1} & \frac{dy_1}{dx_2} & \dots & \frac{dy_1}{dx_n} \\ \frac{dy_2}{dx_1} & \frac{dy_2}{dx_2} & \dots & \frac{dy_2}{dx_n} \\ \dots & \dots & \dots & \dots \\ \frac{dy_m}{dx_1} & \frac{dy_m}{dx_2} & \dots & \frac{dy_m}{dx_n} \end{pmatrix}. \quad (3)$$

Ранг $rg(J)$ матрицы Якоби J определяет характеристики независимости функций, а именно: если $rg(J) = m$, то функции из (1) независимы [12].

Теперь можно установить, возможна ли косвенная утечка значения параметра x_i при применении ИТ-конвейеров, описанных (1). Для этого добавим к (1) функцию тривиального конвейера: обозначим его $y_{m+1} = x_i$.

Составим матрицу (3) для расширенного набора ИТ-конвейеров. Вычислим ранг расширенной матрицы Якоби J^+ . Если при этом $rg(J^+) = m + 1$, то все функции независимы, и косвенная утечка x_i исключена. Если же $rg(J^+) = m$, то в число зависимых ИТ-конвейеров попадет добавленный тривиальный конвейер $y_{m+1} = x_i$ и мы должны заключить, что косвенная утечка x_i возможна. Иными словами, существует такая функция ϕ , что $x_i = \phi(y_1, y_2, \dots, y_m)$, то есть значение x_i вычислимо.

Эти действия проделаем для каждого значения $i = 1, n$, проверив таким образом наличие косвенной утечки для всех x_i .

Следовательно, с помощью этого метода для любого набора ИТ-конвейеров, которым соответствуют непрерывно дифференцируемые функции, можно определить возможность косвенной утечки. Тем самым задача обнаружения косвенных утечек для таких ИТ-конвейеров решена.

После обнаружения вероятности косвенной утечки можно либо применить механизм зашумления к результатам вычисления ИТ-конвейеров, допуская обнаруженную утечку, либо заблокировать выполнение данного набора ИТ-конвейеров, либо сократить объемы вычислений, обеспечив этим защиту от эксплуатации косвенных утечек.

Если же по какой-либо причине устранить зависимость функций, входящих в состав ИТ-конвейеров, невозможно, следует использовать механизмы дифференциальной конфиденциальности в применении к данным отчетов.

Заключение

Введенное понятие «косвенная утечка» раскрывает особенности СИИ, которые необходимо учитывать при создании доверенных систем. На основе общей теории зависимости функций для анализа возможности косвенных утечек предложен новый метод анализа с использованием расширяемых матриц Якоби. Предложенный метод можно расширить для работы с метриками, соответствующие функции которых в общем случае не являются непрерывно дифференцируемыми на всей области определения. Это решает поставленную задачу выявления возможности косвенных утечек данных в СИИ. ■

ЛИТЕРАТУРА

1. Язов, Ю. К. Организация защиты информации в информационных системах от несанкционированного доступа: монография / Ю. К. Язов, С. В. Соловьев. – Воронеж: Кварт, 2018. – 558 с.
2. Конявский, В. А. Защищенные информационные технологии в цифровой экономике / В. А. Конявский, В. В. Медведев, Г. В. Росс // Вопросы защиты информации. – 2022. – № 2 (137). – С. 34–44.

3. Талалаев, А. А. Анализ эффективности применения искусственных нейронных сетей для решения задач распознавания, сжатия и прогнозирования / А. А. Талалаев, И. П. Тищенко, В. П. Фраленко [и др.] // Искусственный интеллект и принятие решений. – 2008. – № 2 (2). – С. 24–33.
4. Смирнов, А. В. Паттерны человеко-машинного сотрудничества в системах поддержки принятия решений / А. В. Смирнов, Т. В. Левашова // Искусственный интеллект и принятие решений. – 2024. – № 2 (66). – С. 3–17.
5. Dinur, I. Revealing information while preserving privacy / I. Dinur, K. Nissim // Proc. of the Twenty-Second ACM SIGMOD-SIGACT-SIGART Symposium on Principles of Database Systems (PODS '03). – ACM, New York, NY, USA, 2003. – DOI: 10.1145/773153.773173.
6. Dwork, C. Differential Privacy / C. Dwork // International Colloquium on Automata, Languages and Programming (ICALP), 2006. – DOI: 10.1007/11787006_1.
7. Wood, A. Differential Privacy: A Primer for a NonTechnical Audience / A. Wood, M. Altman, A. Bembek [et al.] // 21 Vanderbilt Journal of Entertainment and Technology Law. – 2020. – V. 21. Iss. 1. – P. 209–276 [Электронный ресурс]. – URL: <http://scholarship.law.vanderbilt.edu/jetlaw/vol21/iss1/4/> (дата обращения: 01.07.2025).
8. Kenthapadi, K. Releasing Search Queries and Clicks Privately / K. Kenthapadi, A. Korolova, N. Mishra [et al.] // WWW '09: The 18th International World Wide Web Conference. Madrid, Spain. April 20–24. – 2009. – P. 171–180. – DOI: 10.1145/1526709.1526733.
9. Конявский, В. А. Технология «слепой» обработки привлекаемых данных в системах машинного обучения / В. А. Конявский, С. В. Конявская-Счастливая, Г. В. Росс [и др.] // Вопросы защиты информации. – 2024. – № 2. – С. 17–32.
10. Конявский, В. А. Доверенные информационные технологии: от архитектуры к системам и средствам / В. А. Конявский, С. В. Конявская. – М.: URSS, 2019. – 264 с. (то же: Конявский В. А., Конявская С. В. Доверенные информационные технологии: от архитектуры к системам и средствам : 2-е изд. – М.: URSS, 2021. – 264 с.).
11. Вусс, Г. Система страхования информационных рисков / Г. Вусс, В. Конявский, В. Хованов // Финансовый бизнес. – 1998. – № 3. – С. 34.
12. Фихтенгольц, Г. М. Курс дифференциального и интегрального исчисления. Том 1 / Г. М. Фихтенгольц. – СПб.: Лань, 1997. – 608 с.