

Вас пригласили на работу в качестве специалиста по данным в коммерческую компанию, занимающуюся предоставлением банковских услуг населению. Одной из услуг компании, пользующейся большой популярностью у населения, является проведение онлайн-платежей с банковских счетов граждан. К сожалению, в последнее время участились случаи мошенничества при проведении таких операций. Руководство компании встревожено такой ситуацией. На ближайшем совещании было принято решение внедрить в IT-инфраструктуру компании модуль на базе машинного обучения, способный в реальном времени оценивать, является ли проводимая транзакция мошеннической. Это, в свою очередь, позволило бы системе онлайн-платежей отклонить подозрительную операцию. Вам предлагают принять участие в проектировании и построении такой системы. В ходе детального обсуждения проекта с представителями заказчика выяснились следующие моменты:

1. Руководство готово выделить на реализацию данного проекта не более 10 млн. руб. (не считая зарплат специалистам). Через три месяца оно хотело бы увидеть первые результаты, что позволило бы понять, стоит ли дальше развивать проект. В случае положительного решения, проект нужно завершить за полгода.
2. Из открытых источников известно, что у ближайших конкурентов на каждую сотню проводимых их клиентами транзакций фиксируется не более двух мошеннических, приводящих к потере денежных средств. При этом общий ущерб клиентов за месяц не превышает 500 тыс. руб. Разрабатываемая система должна выдавать результаты не хуже, иначе компания станет неконкурентоспособной.
3. В среднем, компания обрабатывает около 50 транзакций в секунду, однако, перед праздниками это число может достигать 400.
4. Если система определит корректную транзакцию как мошенническую, эта транзакция будет отклонена, а пользователь будет недоволен. Опыт бизнес-аналитиков подсказывает, что если доля таких транзакций превысит 5 %, то начнется отток клиентов из компании.
5. Представители компании не готовы размещать разрабатываемый модуль на собственных вычислительных ресурсах.
6. При проведении транзакции системой вся информация о ней сохраняется в виде csv-файлов, каждая строка которых соответствует одной транзакции. Информация о транзакциях некоторого фиксированного периода помещается в отдельный файл. 7. Файлы с данными содержат всю информацию о транзакциях, включая данные клиента, который ее проводил. Данная информация является конфиденциальной, ее утечка недопустима. Вам предлагается на основе представленной информации:

**1. Сформулировать цели проектируемой антифрод-системы в соответствии с требованиями заказчика.**

Целью проекта является разработка сервиса оценки транзакций, проводимых с банковских счетов через системы онлайн-платежей. Сервис должен работать в режиме, близком к реальному времени и построен на базе машинного обучения. Проект должен быть реализован не более чем за 6 месяцев при бюджете в 10 млн. рублей. Первый этап (MVP) должен быть завершен не позднее чем через 3 месяца после начала работ.

**2. Аргументированно выбрать метрику машинного обучения, оптимизация которой, на Ваш взгляд, позволит достичь поставленных целей.**

Предлагается использовать метрику Precision (точность). Данная метрика покажет, какая доля из всех предсказанных сервисом мошеннических транзакций на самом деле были таковыми. Это позволит

управлять показателем False Positive, так как при большом количестве ложных срабатываний ожидается рост недовольства и отток клиентов.

Для оценки качества работы системы дополнительно следует применить метрику F1-score, которая позволит учесть, что сервис должен выявлять максимальное количество мошеннических транзакций.

Применение метрики False Positive Rate позволит оценить количество транзакций, ошибочно оцененных как мошеннические.

Чем выше Precision, тем лучше работает сервис.

Чем ниже False Positive Rate, тем лучше работает сервис.

Чем выше F1-score, тем лучше работает сервис.

## 3. Проанализировать особенности проекта с использованием, например, MISSION Canvas.

<b>EDITION TASK</b> <i>Задача:</i> выявление мошеннических транзакций из потока транзакций, обрабатываемых предприятием  <b>Вход:</b> транзакция  <b>Выход:</b> прогноз оценки мошенничества при выполнении транзакции, от 0 до 1, где 1 - транзакция мошенническая, 0 - мошенничество не определено.	<b>DECISIONS</b> Внедрение сервиса на базе машинного обучения, определяющего является ли транзакция мошеннической	<b>VALUE PROPOSITION</b> Конечным пользователем продукта являются потребители услуг заказчика, взаимодействующие непосредственно с транзакционной системой, а также заказчик, осуществляющий анализ деятельности предприятия и оценку влияния мошеннических транзакций на показатели деятельности предприятия. Целями потребителей являются: сокращение издержек от отсекаемых ошибочно транзакций, не являющихся мошенническими. Целями заказчика являются: сокращение издержек за счёт сокращения числа исполненных мошеннических транзакций, повышение репутации предприятия за счёт минимизации невыполненных транзакций, ошибочно оценённых как мошеннические.	<b>DATA COLLECTION</b> Для сбора обучающего набора данных допускается использование данных, собранных в процессе операционной деятельности заказчика, а также данные из открытых источников в случае возможности их адаптации к процессу обучения модели. Минимально требуется сбор собственных данных за 6 месяцев. В процессе обучения набор данных должен обогащаться новыми данными из внутренних и внешних источников не реже одного раза в две недели.	<b>DATA SOURCES</b> Внутренние данные заказчика в формате csv должны быть преобразованы к унифицированному для обучения модели
<b>OFFLINE EVALUATION</b>  Средний поток обрабатываемых транзакций: 50 транзакций в секунду	<b>MAKING PREDICTIONS</b>  Выводы о мошеннических транзакциях формируются на основе внутренних и внешних данных, сведений о	Бизнес-система предприятия запрашивает у БД-системы оценку транзакции, передавая информацию по защищённым каналам связи. На основе оценки бизнес-система исполняет транзакцию, либо отказывает в исполнении.	<b>BUILDING MODELS</b>  Требуется одна модель. В процессе разработки требуется выполнять обновления каждые 2 недели в течение первых 3 месяцев и не реже 1 раза в	<b>FEATURES</b>  Сумма транзакции Страна, откуда была отправлена транзакция

<p>Пиковый поток обрабатываемых транзакций: 400 транзакций в секунду (в предпраздничные дни)</p>	<p>мошеннических схемах, об утечках конфиденциальных данных.</p>		<p>месяц во вторые 3 месяца. Для обновления модели, в т.ч. выявление признаков и анализ доступно две недели.</p>	<p>Частота и время проведения транзакций</p> <p>Информация об устройстве, с которого была отправлена транзакция (например, IP-адрес, браузер и т.д.)</p> <p>История платежей клиента</p> <p>Поведенческие данные клиента (например, скорость ввода данных, поведение мыши и т.д.)</p> <p>Информация о том, были ли ранее обнаружены мошеннические транзакции с этого банковского счета</p>
	<p><b>LIVE MONITORING</b></p> <ol style="list-style-type: none"> <li>1. Непрерывный мониторинг всех транзакций и быстрое определение мошеннических транзакций с помощью применения алгоритмов машинного обучения, обученных на данных о ранее проводимых транзакциях</li> <li>2. Оповещение ответственных лиц при обнаружении мошеннических транзакций, оповещение корпоративных систем для автоматической блокировки исполнения транзакции</li> <li>3. Для оценки производительности системы требуется настроить контроль и мониторинг ее работы путем использования средств, предоставляемых поставщиками облачных услуг</li> <li>4. Так как система должна обеспечивать в предпраздничные дни 400 транзакций в минуту, тогда как в штатном режиме только 50, требуется предусмотреть горизонтальное масштабирование вычислительных компонентов системы, интегрированное с системами мониторинга и оркестрации, предоставляемыми провайдером облачных услуг</li> <li>5. Для обеспечения безопасности данных требуется постоянный контроль доступа, шифрование, резервное копирование, антивирусная защита</li> </ol>			

**4. Попытаться декомпозировать планируемую систему, определить ее основные функциональные части.**

Основные функциональные компоненты:

1. Интерфейс пользователя. Это часть системы, которая позволяет пользователю взаимодействовать с ней. Она должна быть интуитивно понятной, удобной и легко настраиваемой. В ее состав могут входить такие элементы, как окна настроек, кнопки управления, всплывающие подсказки и т.д.
2. Обработка данных. Эта часть системы отвечает за обработку и анализ данных, полученных от пользователя или из других источников. Здесь могут использоваться алгоритмы машинного обучения и статистические методы для анализа данных.
3. Хранение данных. Эта часть системы отвечает за хранение и управление данными. Здесь могут использоваться базы данных, облачные хранилища и другие технологии для хранения и обработки данных.
4. Визуализация данных. Эта часть системы отвечает за представление данных в удобном для пользователя формате. Здесь могут использоваться графики, таблицы, диаграммы и другие элементы визуализации данных.
5. Управление системой. Эта часть системы отвечает за управление всеми ее компонентами. Здесь могут использоваться автоматические системы управления и мониторинга, которые позволяют обнаруживать и устранять проблемы в работе системы.
6. Подсистема обучения. Включает в себя имплементацию компонентов сбора данных, предварительной обработки, обучения и валидации модели машинного обучения, ее упаковки для развертывания в продуктивных средах.

**5. Определить задачи, решение которых необходимо для достижения поставленных целей с учетом проведенного анализа. Задачи рекомендуется формулировать по принципу S.M.A.R.T. На текущий момент, пока не конкретизированы детали антифрод-системы, они могут быть представлены в достаточно общем виде. Но они позволят сформировать некоторый roadmap по реализации проекта.**

1. Создание базы данных для хранения информации о пользователях и их транзакциях.
  - Specific (конкретная): создание базы данных для хранения информации о пользователях и их транзакциях.
  - Measurable (измеримая): количество созданных таблиц и столбцов в базе данных.
  - Achievable (достижимая): разработка и использование баз данных - стандартный подход в современной разработке программных продуктов.
  - Relevant (релевантная): база данных является необходимой частью системы для хранения и обработки информации о пользователях и их транзакциях.
  - Time-bound (ограниченная по времени): создание базы данных должно быть завершено в течение первых двух месяцев проекта.
2. Разработка и обучение модели машинного обучения для обнаружения мошеннических транзакций.
  - Specific (конкретная): разработка и обучение модели машинного обучения для обнаружения мошеннических транзакций.
  - Measurable (измеримая): точность работы модели машинного обучения в процентах.
  - Achievable (достижимая): на сегодняшний день существует множество алгоритмов и фреймворков для разработки и обучения моделей машинного обучения, которые можно применить для решения этой задачи.
  - Relevant (релевантная): модель машинного обучения является ключевой частью системы антифрод-защиты и позволит автоматически определять мошеннические транзакции.
  - Time-bound (ограниченная по времени): разработка и обучение модели машинного обучения должно быть завершено в течение первых четырех месяцев проекта.
3. Создание системы алертинга для оперативного уведомления о мошеннических транзакциях.

- Specific (конкретная): создание системы алертинга для оперативного уведомления о мошеннических транзакциях.
  - Measurable (измеримая): время реакции системы на обнаружение мошеннических транзакций.
  - Achievable (достижимая): существуют готовые решения для создания систем алертинга, которые можно адаптировать
  - Relevant (релевантная):
  - Time-bound (ограниченная по времени): завершить разработку и внедрить систему алертинга в рамках платформы в первые четыре месяца работы над проектом.
4. Разработка системы мониторинга аномальной активности на платформе с целью своевременного обнаружения и предотвращения фрода.
- Specific (конкретная): конкретизировать требования к системе мониторинга и установить ее ключевые метрики (например, частота обновления данных, время ответа, точность обнаружения фрода).
  - Measurable (измеримая): определить необходимые ресурсы (например, специалисты по анализу данных, инфраструктуру для сбора и обработки данных) и распределить их для реализации системы мониторинга.
  - Achievable (достижимая): разработать систему мониторинга с использованием алгоритмов машинного обучения, обеспечивающих высокую точность обнаружения фрода.
  - Relevant (релевантная): оценить эффективность системы мониторинга с помощью метрик качества и улучшить ее, если это необходимо.
  - Time-bound (ограниченная по времени): завершить разработку и внедрить систему мониторинга в рамках платформы в первые четыре месяца работы над проектом.
5. Создание базы знаний для обучения моделей машинного обучения с целью улучшения точности определения фрода.
- Specific (конкретная): определить требования к базе знаний (например, размер, содержание, формат данных) и установить ключевые метрики (например, точность определения фрода на основе данных из базы знаний).
  - Measurable (измеримая): собрать и обработать данные для создания базы знаний.
  - Achievable (достижимая): создать базу знаний и обучить модели машинного обучения на ее основе, используя различные алгоритмы и техники обработки данных.
  - Relevant (релевантная): оценить точность определения фрода на основе данных из базы знаний и улучшить базу знаний, если это необходимо.
  - Time-bound (ограниченная по времени): завершить создание базы знаний и использовать ее для улучшения точности определения фрода на платформе в первые три месяца работы над проектом.