

## Задание 3 "Качество данных"

### Проверки на пригодность данных

Проверка качества	Название	Описание	Источник данных	Требование	Метод проверки
Полнота	Наличие <b>client_id</b>	Каждая запись должна содержать идентификатор клиента	IB_SYS TEM, CRM_BA NK	client_id IS NOT NULL	<code>SELECT COUNT(*) FROM raw_ib WHERE client_id IS NULL</code>
Полнота	Наличие реквизитов клиента	Если нет <b>client_id</b> , должны быть ФИО, паспорт, дата рождения	PARTNER_API_*, PARTNER_FILE_*	При client_id IS NULL: все поля заполнены	<code>SELECT COUNT(*) FROM raw_partner WHERE client_id IS NULL AND (first_name IS NULL OR passport IS NULL OR birth_date IS NULL)</code>
Полнота	Наличие <b>product_id</b>	Каждое подключение должно быть связано с продуктом	Все источники	product_id IS NOT NULL	<code>SELECT COUNT(*) FROM raw_source WHERE product_id IS NULL</code>
Точность	Корректность значения <b>channel_cd</b>	Поле <b>channel_cd</b> должно быть одним из: IB, BRANCH, PARTNER	Все источники	channel_cd IN ('IB', 'BRANCH', 'PARTNER')	<code>SELECT COUNT(*) FROM raw_source WHERE channel_cd NOT IN ('IB', 'BRANCH', 'PARTNER')</code>
Точность	Формат и валидность <b>activation_dttm</b>	Дата активации должна быть в правильном формате и $\geq$ 01.01.2024	Все источники	activation_dttm >= '2024-01-01' и IS TIMESTAMP	<code>SELECT COUNT(*) FROM raw_source WHERE activation_dttm &lt; '2024-01-01' OR activation_dttm IS NOT A VALID TIMESTAMP</code>
Точность	Статус подключения = ' <b>active</b> '	Только активные подключения учитываются	Все источники	status = 'active'	<code>SELECT COUNT(*) FROM raw_source WHERE status != 'active'</code>

Уникальность	Отсутствие дубликатов	Не должно быть одинаковых записей (по <code>client_id, product_id, activation_dttm</code> )	Все источники	Нет дублей	<code>SELECT client_id, product_id, COUNT(*) FROM raw_source GROUP BY client_id, product_id, activation_dttm HAVING COUNT(*) &gt; 1</code>
Полнота	Наличие всех обязательных полей	Проверка, что все ключевые поля присутствуют в схеме	Все источники	Должны быть: <code>client_id</code> (или реквизиты), <code>product_id</code> , <code>channel_cd</code> , <code>activation_dttm</code> , <code>status</code>	<code>DESCRIBE raw_source</code> или проверка при загрузке

### Проверки на регулярной основе (ежедневно)

Наличие <code>client_id</code> (для <code>IB_SYSTEM, CRM_BANK</code> )	Гарантия идентификации клиента
Наличие реквизитов клиента (для <code>PARTNER_*</code> )	Возможность мэтчинга без <code>client_id</code>
Корректность <code>channel_cd</code>	Предотвращение искажения каналов
Проверка диапазона дат ( <code>activation_dttm ≥ 2024-01-01</code> )	Соответствие периоду анализа
Проверка статуса ( <code>status = 'active'</code> )	Учёт только активных подключений
Проверка на дубликаты	Исключение завышения метрик

### Документ описания метрик

Код метрики	Описание метрики	Слой DWH	Единица измерения	Нижний порог хорошего качества	Верхний порог хорошего качества	Система источник	Бизнес владелец	SQL запрос
missing_client_id_count	Количество записей без client_id	raw (IB, CRM)	count	0	0.1	IB_SYSTEM, CRM_BANK	Владелец данных "Клиент"	SELECT COUNT(*) FROM raw_ib WHERE client_id IS NULL
missing_requisites_count	Количество записей без реквизитов для мэтчинга	raw (PARTNER)	count	0	0.1	PARTNER_API_*, PARTNER_FILE_*	Владелец данных "Клиент"	SELECT COUNT(*) FROM raw_partner WHERE client_id IS NULL AND (first_name IS NULL OR passport IS NULL OR birth_date IS NULL)
invalid_channel_cd_count	Количество записей с некорректным channel_cd	raw	count	0	0.1	Все источники	Владелец данных "Канал"	SELECT COUNT(*) FROM raw_source WHERE channel_cd NOT IN ('IB', 'BRANCH', 'PARTNER')
out_of_range_date_count	Количество записей с activation_dttm < 2024-01-01	raw	count	0	0.1	Все источники	Владелец данных "Продукт"	SELECT COUNT(*) FROM raw_source WHERE activation_dttm < '2024-01-01'
inactive_status_count	Количество подключений с status != 'active'	raw	count	0	0.1	Все источники	Владелец данных "Продукт"	SELECT COUNT(*) FROM raw_source WHERE status != 'active'
duplicate_count	Количество дублирующихся	raw	count	0	0.1	Все источники	ETL-разработчик	SELECT COUNT(*) FROM (SELECT client_id, product_id,

	подключен ий							<code>activation_dttm FROM raw_source GROUP BY client_id, product_id, activation_dttm HAVING COUNT(*) &gt; 1) AS dups</code>
<code>total_re cords_lo aded</code>	Общее количество загружен ых записей	raw	coun t	100	$\infty$	Все источники	ETL-разр аботчик	<code>SELECT COUNT(*) FROM raw_source</code>

## Шаблон для информирования об инциденте качества данных

(этим цветом выделены примеры заполнения шаблона)

Дата и время инцидента:[ГГГГ-ММ-ДД НН:ММ]

Источник данных: PARTNER\_API\_X

Тип инцидента: Пропущенные поля

Описание проблемы: В выгрузке от партнёра за 15.04.2024 отсутствует `client\_id` в 18% записей. Также обнаружены значения `channel\_cd = 'partner\_site'`, что не соответствует стандарту.

Последствия:

- Невозможно идентифицировать 18% клиентов
- Некорректное распределение подключений по каналам
- Риск искажения доли партнёрских каналов в отчёте

Решение/Рекомендации:

1. Немедленно уведомить владельца источника ('Представитель партнёрской интеграции').
2. Временно использовать мэтчинг по ФИО+паспорт+дата рождения.
3. Заблокировать загрузку данных с некорректным `channel\_cd`.
4. Запросить исправление формата выгрузки от партнёра.

Ответственный: [ФИО, должность]

Статус: [Новый / В работе / Решено / Закрыто]

Срок исправления:\*\* [ДД.ММ.ГГГГ]