Модельная система

Рассмотрим простую модельную систему: у нас есть фронты (веб-сайт, приложение и т.д.) и бэк (серверное приложение). Фронтов может быть много, потому что одновременно может быть несколько версий сайта, приложения и пр. Кроме того, может быть сервис по приемке логов за балансером с несколькими инстансами. Все пользовательские действия записываются на фронте и на бэке. Для простоты будем считать, что любому событию на фронте соответствует некоторый лог на сервере.

Мы хотим считать какуюнибудь пользовательскую метрику, например, время, которое пользователь чтонибудь ждет (крутится splash окошечко, прогресс бар и т.д.). Для этого, мы сливаем логи из браузера и из сервера в потоковую систему, которая считает эту метрику. Далее, потоковая система отдает значения метрики для

Важно, что считать метрику мы хотим по пользовательским сессиям, т.е. по промежуткам активности. Чтобы определять, что сессия закончилась, мы используем ML модель, которая по специальным

метрики отображения в дашборд. Потоковая система фичам определяет, Фронт (браузер) Бэкенд (сервер) Фронт (мобильное

приложение)

Фронт (Телеграмбот)

Дашборд для

отображения

Naive graph

операция для

состояния

согласованности

шардируем по user

Источник браузерных событий

Это источник данных, через который идет доставка логов из браузера. Каждый элемент потока содержит структуру лога, описанную выше.

Это искусственная операция для того, чтобы из двух потоков сделать один. Ее можно представить просто как операцию, которая принимает любые элементы на вход и их отдает на выход. Вводим ее для наглядности, на самом деле, например во Flink, любая операция может иметь несколько входов.

Мержим

потоки

В этой операции мы соединяем в один элемент лог из браузера и соответствующий ему лог из сервера. На выходе получается элемент, который содержит в себе оба входных

Join

ЛОГОВ

события. Чтобы правильно соединить два лога, нам нужно, чтобы логи от одного и того же пользователя шли на одну машинку физическую. Поэтому, данные шардируются по user id. Bo Flink

По этой ветке мы должны получить ответ от модели: закончилась сессия или нет. Сначала добавляем фичи для фронта, например, как часто идут события от этого фронта и т.д. Соответственно, также меняем и состояние этой операции, например, ту же частоту. Чтобы состояние было полным, нужно

Фичи для

фронта

Структура логов

user id

front id

payload

ts

получать все события от одного фронта на одну

Фичи для Применяем юзера модель Аналогичная Применяем модель к

полученным фичам для пользователей. Для пользователя и фронта, опредяляем, закончилась ли сессия для пользователя. Т.к. модель может иметь разные параметры для каждого пользователя и изменяться во времени (online learning), то шардировать нужно по user id.

Выпускаем Фильтрация метрику

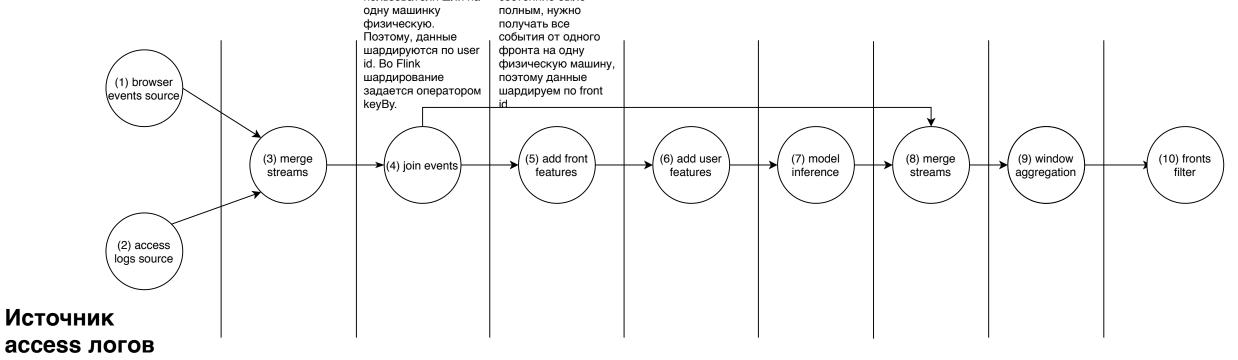
Здесь нам приходят два типа элементов:

1. Соединенные логи фронт-юзер. Их мы аггрегируем в метрику.

2. Событие конца окна от модели. Выпускаем метрику наружу. Чистим состояние для последней сессии.

Опциональная операция. Можем захотеть пофильтровать самые редкие фронты, или старые версии фронтов. Важно, что фичи для этих фронтов посчитаны и учтены в модели. Шардируем по front id, так как считаем, что здесь есть состояние.

Шардируем по user id.



Источник данных, через который идет доставка логов с сервера. Структура логов аналогична браузерным.

partition by user id

partition by front id

partition by user id

partition by user id

partition by user id

partition by front id

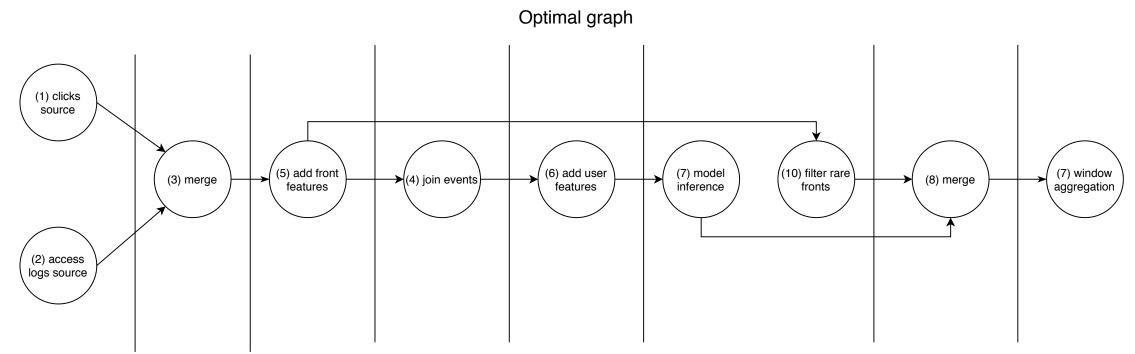
partition by user id

Принцип построения оптимального графа

Заметим, что в исходном графе:

- 1. Много ре-шардирований
- 2. Добавление фичей и соединение логов можно менять местами, т.к. добавление фичей только добавляет поля в элемент
- 3. Фильтр можно делать до агрегации, но после обновления модели

Таким образом, мы можем переставить операции так, чтобы они "склеились" по шардированям и мы бы не гоняли данные по



partition by front id

partition by user id

partition by user id

partition by user id