

Москва — 2021

DWH как продукт

Евгений Николаев

Руководитель юнита DWH

Евгений Николаев

Руководитель юнита DWH

- выпускник ВМК МГУ
- 6+ лет опыта разработки
- > 3+ года опыта руководства
- фанат классных продуктов
- 🕨 капитан ФК Авито
- КМС по шахматам
- говорю по-испански



Авито

40,8м месячная аудитория

Больше четверти населения России

Каждый день создаётся почти 600К новых объявлений

11

объявлений

182м сделок в 2020

Это 60% от всех сделок в России

DWH в Авито



План выступления

- Зарождение аналитики в Авито
- Платформизация
- DWH как продукт
- Масштабирование
- Выводы и планы

Зарождение аналитики в Авито

Предпосылки

- Монетизация данных.
- А/В-тестирование.
- Автоматическая модерация.
- Финансовая отчётность.
 - 🗢 расширяемая ВІ-платформа.

С чего мы начинали?

- Vertica + Tableau.
- ▶ ~50 Tb данных.
- ▶ ~6 дата-инженеров.
- ~25 пользователей-аналитиков.

кейс Vertica в Авито

Пример задачи: витрины

Analytical Processing / OLAP-2904

Подготовить расширенную витрину по активным пользователям

Description

К текущей витрине current_user добавить следующие поля:

- Item_Location_id самый частый локейшн среди 5 последних размещенных айтемов
- Item_Category_id самая частая категория среди 5 последних размещенных айтемов
- Item_Active количество активных объявлений у юзера (External_id статус<10)
- Item_Total общее количество объявлений у юзера
- Item_Blocked_Total общее количество заблокированных объявлений у пользователя
- Item_Blocked_1W количество заблокированных объявлений за последние 7 дней
- Item_New_1M количество новых объявлений за последние 30 дней
- Item_New_1W количество новых объявлений за последние 7 дней
- VAS_Amount_Total объем потраченных средств на VAS за всю историю
- VAS_Amount_1M объем потраченных средств на VAS за последний месяц
- VAS_Amount_1W объем потраченных средств на VAS за последнюю неделю
- VAS_Count количество VAS транзакций за всю историю
- Last_visit_time датавремя последнего действия юзера на сайте

Пример задачи: ETL

Analytical Processing / OLAP-67 Необходимо создать свойство multi_vas для события 273 на step = 1

Description

Нужен сателит, который позволял бы понимать, что данное применение vas включает в себя несколько vasservice (dds.S_WebLog_VASService), то есть является мультиплатежом. Свойство нужно видеть на первом шаге воронки (dds.S_WebLog_Step), когда платежу еще не присвоен идентификатор трансакции (а значит нет идентификатора мультиплатежа – vasgroup), но по факту пользователь уже добавил в корзину несколько услуг.

Проблемы

Как справиться с ростом задач?

- Как не погрязнуть в рутине?
- Как масштабировать команду?
- Как ставить цели?

Платформизация

В чём идея платформизации?

- Мы создаём инструменты.
- Пользователи сами решают задачи.
- Мы помогаем им, где не получается.

Цель — 100% задач решаются пользователями.

Витрины создают аналитики!

OLAP-6457 → master MERCED OLAP-6457 new dma Overview Diff Commits Builds ↓ All changes in this pul ↓ ↓ All changes in this pul ↓ ↓ Eliher file	sql/data s.sql At	Analytical Processing / OLAF Витрина dma.subs marts / dma_subscription_type	P-6457 scription_types Blame 2 other comme
S Filler III 75 Searc	1 +	/**	
🚔 sql/datamarts	2 +	<pre>* @datamart DMA.subscription</pre>	n_types
dma_subscription_types.sql	3 +	* @key event_date	
	4 +	* @key LastVisit_Month	
	5 +	* @key UserType	
	0 +	* @key UserStatus_10	
	/ + 8 +	* of TNCREMENTAL REEDESH	
	9 +	* @naram first date	
	10 +	* Oparam last date	
	11 +	* @param launch id	
	12 +	*/	
	13 +		
	14 +	CREATE TABLE IF NOT EXISTS de	ev_DMA.subscription_types
	15 +	(
	16 +	event_date	date NOT NULL,
	17 +	LastVisit_Month	date,
	18 +	UserType	varchar(80),
	19 +	UserStatus_id	int,
	20 +	email_avito_news	int,
	21 +	email_action	int,
	22 +	email_poll	int,

Build		Statu	IS
n	DWH / DailyScheduleCycleDetection #13562 LATEST Svetlana Trushina § 5946abf3c6a REAL-8596	0	Passed
n	DWH / TestDatamart #19937 LATEST Svetlana Trushina ∳ 5946abf3c6a REAL-8596	0	Passed
n	DWH / ValidateInputs #20616 LATEST Svetlana Trushina 🕴 5946abf3c6a REAL-8596	0	Passed
n	DWH / TestDatamartLimit0 #21439 LATEST Svetlana Trushina 🕴 5946abf3c6a REAL-8596	0	Passed
n	DWH / CheckCodeStandard #19824 LATEST Svetlana Trushina 🕴 5946abf3c6a REAL-8596	•	Passed
n	DWH / TestDatamartViews #30596 LATEST Svetlana Trushina 🕴 5946abf3c6a REAL-8596	•	Passed
n	DWH / TestUploaders #19430 LATEST Svetlana Trushina § 5946abf3c6a REAL-8596	•	Passed
n	DWH / PyTest #35303 LATEST Svetlana Trushina ∮ 5946abf3c6a REAL-8596	9	Passed

Datamart Framework

- SQL-based синтаксис.
- Инструмент оркестрации.
- Итог: 800+ ежедневных расчётов.

Queue States			L	Launches		Name
dma_calcul	ator_new	 ✓ States 	~	7486335 🛞	\sim	Nan
id \$	launch 🌲	name 🗢		updated 🗢	st	ate ≑
4782469	7486335	DMA.current_cthulhu_te	estplan_testcase	2021-08-24 13:44	:50	done

- /**
- * @datamart DMA.subscription_types
- * @key event_date
- * @key LastVisit_Month
- * @key UserType
- * @key UserStatus_id
- *
- * @fn INCREMENTAL_REFRESH
- * @param first_date
- * @param last_date
- * @param launch_id
- */

ETL делают аналитики!

Analytical Processing / OLAP-11555 Разложить в DDS данные из realty-ownership

Description

Что за данные:

Наименование сервиса:

realty-ownership

Бизнес-смысл данных:

Мы в ближайшее время запускаем MVP верификации собственника в RE. Соответственно из сервиса забираем данные и статусы для аналитики продукта.

Здесь подробно про верификацию описано - https://cf.avito.ru/pages/viewpage.action?pageId=175181829

Сценарии использования:

Данные в DWH нужны, чтобы настроить аналитику по продукту верификации собственника в RE.

PR с описанием раскладки тут:

http://stash.msk.avito.ru/projects/BI/repos/avito-dwh/pull-requests/21556/overview

ETL Framework



Масштабирование команды

dwh-growth

Цель — растить скорость решения аналитических задач:

поиск данных, построение отчетов, проверка гипотез, ETL.

dwh-infra

Цель — развитие инфраструктуры:

запросы, витрины, отчёты, антибот, realtime-аналитика.

Проблемы

- dwh-growth: много пользовательских контекстов;
- dwh-infra: мало взаимодействия с пользователями;
- платформа медленно развивается;
- инфраструктура усложняется (+ClickHouse);
- больше бизнесовых пользователей;

⇒ юнит растёт.

DWH как продукт

DWH — продукт?



План

Полезные практики:

- User Story Mapping;
- Customer Development.
- Примеры проблем:
 - о поиск данных;
 - высокий порог входа в аналитику.

User Story Mapping

- Визуализировать и лучше представлять продукт.
- ▶ Что делаем?
- Для кого?
- ▶ Зачем?

Пример: проверка гипотезы



<u>плагин в miro</u>

Как составить?

- Определяем активности пользователя.
- Разбиваем активности на задачи.
- Нарезаем карту на куски по важности.
- Свадебный торт -> капкейки.

Jeff Patton лекция на ютубе

Jeff Patton «User Story Mapping» книга

Customer Development

- Как понять потребности пользователей?
- Как проверять гипотезы по развитию продукта?
- Как собирать обратную связь?

CustDev: примеры вопросов

- В чём для вас главная ценность продукта?
- З главных преимущества продукта.
- Что изменится, если продукта не будет?
- Какой продукт служит альтернативой?
- Вы рекомендовали знакомым, как?
- Есть ли факторы, которые мешают пользоваться?

CustDev: советы

- Все вопросы открытые.
- Правильно ли я понял, что «…»?
- Сценарий канва разговора.
- Вопросы про прошлое, а не будущее.
- Интересуют факты, а не мнения и оценки.
- ▶ 5 почему.

книга Роберта Фитцпатрика

Проблема поиска данных

"Данные искать сложно. Чаще всего, приходится искать человека, который знает, где что лежит и как работает."

"Часто приходится обращаться в чаты в слаке, чтобы понять, есть ли та или иная информация в вертике. Мне кажется, очень не хватает описания dds таблиц."

Система поиска dwh-docs



agolovanev 🐡 Aug 3rd at 12:11 PM in #bi-and-dwh

всем привет, а где-то ведь наверняка хранятся кадастровые номера указанные в объявлениях? Кто , то может подсказать - где?)

2 replies



mbnekrasov 1 month ago

select * from tables where table_name
ilike '%cadastral%'

1



agolovanev 🐡 1 month ago

спасибо, нашел

кадастровые номера				
Искать везде	• Все кластеры	v	Все схемы	
Все юниты	Все типы табл	иц 👻	Все витрины	
Название		Юнит	Описание	
DDS.H_CadastralNumber		Common not Allocatable	Кадастровый номер	
DDS.S_ItemCadastralNumberLi	nk_IsDuplicate	-	Повторяется ли кадастрового номера на сайте в другом объявлении	
DDS.S_ItemCadastralNumberLi	nk_CheckBitMask	-	Битовая маска всех проверок, что происходили по кадастровому номеру	
DDS.L_ItemCadastralNumberLin	nk_CadastralNumber	-	Связка cadastralnumber с уникальной связкой item и кадастрового номера	
DDS.L_ItemCadastralNumberLin	nk_Item	-	Связка item с уникальной связкой item и кадастрового номера	
DDS.S_CadastralNumber_VerifiedRosreestr		-	Подтверждение существование кадастрового номера от Росреестра	
DDS.H_ItemCadastralNumberLi	nk	Common not Allocatable	Ускуственный уникальный ключ между item и кадастровым номером	
DDS.S_ItemCadastralNumberLink_IsShown		-	Будет ли показано на сайте, что кадастровый номер проверен	

dwh-docs

- Поиск таблиц и отчётов.
- Документирование
 таблиц и колонок.
- 🕨 Связи с Jira и Stash.
- Отображение аномалий в данных.
 - Подписки на объекты.

Тип Таблица Последний расчет 30 сентября 2021 г. 7:48 🗸 DuplicateKey 🗸 EmptyLaunch 🗸 Описание Витрина с текущим состоянием отзывов Исходный код sgl/datamarts/dma_current_reviews.sgl megoncharov, avshovko, aarvsmyatova Подписчики Подписаться Юнит T&S

DMA.current reviews

Проблема DBeaver

«Он стремный, убогий интерфейс, сложные настройки, после обновления периодически ломается (например, работа с параметрами), сложно было настраивать самому (аналитик помогал выбрать драйвер и ещё что-то), не так много документации/инфы в интернете/на stackoverflow.»

SQL-клиент в браузере

- Поддержка Vertica, CH, Postgres.
- Шаринг запросов через URL.
- Создание витрин.
- Работа с Tableau.
 - Проверка планов запросов.

Обратите внимание на план запроса

explain select count(*) from dma.current_item join dma.current_user

| +---> JOIN HASH [Cost: 509K, Rows: 2B] (PATH ID: 2) Outer (RESEGM | | +-- Inner -> STORAGE ACCESS for current_user [Cost: 15K, Rows:

$\leftarrow \rightarrow \mathbf{C}$ ($\hat{\mathbf{a}}$ dwh.avito.ru/queries)	inew 🖈 🙂 🖸 🛅	0 🔕
🗄 Приложения 🎽 Рабочий стол - А.,	📌 System Dashboar 🗞 Meragawaan DWH 🧧 B rayda saaka Pyt 🍼 Pythen Cookbook 🕸 Pythen-coodiaject 🚺 Vertica® 8.1 x Doc » 🛅 Другие з	акладки
master(idap proxy:5433) •	Run Cancel 🕮 🕭 🔁 🕼 🏋 master	I
Search table C	1 select dotar:dute, swr(cost), lawnh_id 2 fram 540_MARCITHK_UMARLONS, swade.direct.ye.guten_report 3 where dute:disclist between "2004-641" and "2024-65-20"	
 #60 L.Beldog, Loskie L.Beldog, Loskie L.Beldog, Longalescenger L.Beldog, Dragalescenger L.Beldog, J.Bendog, Software L.Beldog, Software L.Beldog, Userdynet L.Beldog, Userdynet L.Beldog, Userdynet S.Beldog, Userdynet	S order by 1 dec	
 cookie_day_olap_5416_bckp cookie_day_olap_5416_bckp_vali_ 	1.059 seconds 148 rows 🛓 .csv ½ .xisx ½ .tsb 🕍 select date:rdate, sum(cost), launch_id _	0
<pre>export_click_stream_1000000_va_</pre>	dete aun laureh (d	
search_stream	1 2020-05-26 385 491-9 5 528 134	
* Search_stream_valiable	2 2020-05-26 387 490.12 5 539 030	
+ API ALL ACTIVE USERS	3 2020-05-26 1 916.77 5 566 000	
<pre>> actual_title</pre>	4 2020-05-26 387 207.6 5 571 390	
<pre>> octual_title1</pre>	5 2020-05-26 387 207.6 5 634 769	
adm_host	6 2020-05-26 320 093.03 5 540 276	
all_us	7 2020-05-26 345 715.14 5 556 319	
h oll ur?		

Queries		>	K master I v ± ± T	4
My queries All connections search	Order by last	saved •	definition d	
office_jps vertica-dwh-proxy:5433/DWH defivery.ymong_orders vertica-dwh-enroy-5633/DWH	agverkhovtseva 2020-06-18719-06:12 agverkhovtseva 2020-06-18720:38:12		No. 1 attle (max)(10, max)(10, max) as (max)(10, max)(10, max)	
delivery_orders vertica-dwh-proxy:5433/DWH	agverkhovtseva 2020-06-16T20:37:09	2	9 unios all select "46.43.455.44", 29340 office power 10 unios all select "35.76.347.244", 22359 NRC ip 11), 22 office_ips as (
utm_campaigns_matching vertica-dwh-proxy:5433/DWH	agverkhovtseva 2020-05-21T20:14:56		1 partic pris (now-ip.d 14 INELANDED) Dow-ip.i 15 gene also mark is not nill then power(2, 32 - mask)::int else 0 end ip.size 17 mg ing sith off(n hander	I
forecast_bands_checker_sellers vertica-dwh-proxy:5435/DWH	agverkhovtseva 2020-04-29T11:40:50	2	18), 19 ips os (20 select	
repaired vertica-dwh-proxy:5435/DWH	agverkhovtseva 2020-02-24T14-00:40.172562		22 entermail.dis (p 23 for 64.8,1) ent 1 for 64.9,1) 24 enter land,14 for 566550 25 effect 25 effect 27 effect 27 effect 28 effect 29 effect 20	
s://dwh.avito.ra/gueries/db/640c6b8/289a3	7		20 from las	

Итоги внедрения sqlpad

- Доля пользователей
 DWH среди
 инженеров выросла на 10%.
 Новые сотрудники
 - используют только sqlpad.

0.8 function Analyst dwh_users_share Engineer Product 0,6 Product 0,4 0,2 +10% 0,0 декабря 2019 декабря 2020 июня 2020

dwh_users_share_by_function

Полезные продуктовые практики

User Story Mapping:

- Jeff Patton лекция на ютубе
- Jeff Patton «User Story Mapping» книга

- Customer Development:
 - книга Роберта Фитцпатрика

Масштабирование

Масштабирование

- Как сформировать цели команд?
- Как поделиться на команды?
- Как приоритизировать задачи?
- Как вовлечь бизнес-пользователей?

Цели команд

Activities из User Story Мар (высокоуровневые задачи пользователя):

- Integration (интегрироваться с DWH).
- Datamart (создать регулярный отчёт).
- Usage (проверить гипотезу).

Integration



Datamart



Usage



Итоги разделения

- **Разнообразие задач**: рутина/вдохновляющая цель.
- Фокус команд: экспертиза, стимул автоматизации.
- Прозрачность: зоны ответственности.
- Масштабирование: пошарили взаимодействие с пользователями.

Как делиться?

- Определяем граничные условия (лиды, цели, число людей, роли).
- Предлагаем выбрать команду и роль самостоятельно.
- Делаем итеративно:
 - выбор команды;
 - о оценка.

Если оценки >= X — завершаем.



Советы

- Подробно объяснить процесс участникам и предпосылки. Ответственность — на команде.
- Кто-то не выберет команды: оставить время пообщаться и снять сомнения.
- Снять страх, что переход навсегда.
- Как вариант снятия ступора предложить всем выписать приоритеты команд.



Как приоритизировать задачи?



Структура бэклога

Квартальные OKR:

- сбор кандидатов, оценки;
- о отбираем TOP, влезающий в 60% ёмкости.

Спринтовые задачи:

- 60%: OKR-задачи;
- 20%: помощь аналитикам;
- 20%: ретро AI + техдолг.

Как ранжировать задачи?

Reach × Impact × Confidence Effort

Использование RICE: советы



Reach оценивать сложно.

Холивар по ценности разных пользователей.

- Обсудить голосом перед голосованием.
- Откалибровать Impact/Effort (или отказаться?)

81 способ приоритизации фичей

Демократизация аналитики



Упрощение задач



Масштабирование: выводы

Как сформировать цели команд?

о Activities из User Story Map.

Как поделиться на команды?

- Ребята сами решают.
- Как приоритизировать задачи?
 - RICE голосование.
- Как вовлечь бизнес-пользователей?
 - Упрощение задач, обучение.

Выводы

DWH как продукт

Производство данных

Гранитик-контрибьютор DWH Аналитик-контрибьютор Продакт Загрузка данных integration Построение отчётов datamart Usage Usage

avito.tech

Использование данных

Цифры

1 Рb 3000 объём данных в CH активных & Vertica пользователей

Доступны для ежедневной аналитики Ежедневно решают задачи в хранилище данных **3300** отчетов в Tableau

Ежедневно обновляются по утрам

Планы

демократизация аналитики



доверие к данным

realtime-аналитика

Q1 2022

аналитика под ключ

opensource

DQ фреймворк

Q2 2022

IDE (SQL + Python + BI)

холодные данные

облачные СУБД

avito.tech

Москва — 2021

Евгений Николаев

Руководитель DWH





<u>@eanikolaev</u>



Если вы хотите построить корабль, Не надо будоражить народ, отправлять собирать древесину, делить работу и отдавать приказы. Лучше научите людей тосковать По обширному и бесконечному морю.