

Итоговый проект по программе

«Дата-инженер»

Разработал:

Олейников Михаил Николаевич

Цели проекта



Разработка MVP системы, которая собирает, обрабатывает и анализирует данные по использованию услуги интерактивного телевидения.

Дополнительная цель – в рамках основной цели использование всех компонентов кластера.

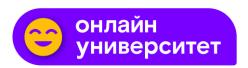
Поставленные задачи



Реализовать следующие пункты технического задания:

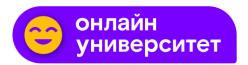
- 1. Сбор данных с использованием RT.Streaming:
 - Создание продюсера на Python, который будет симулировать данные о поведении пользователей интерактивного телевидения и отправлять их в топик Kafka.
- 2. Хранение сырых данных в RT.DataLake:
 - Создание потребителя на Python для RT.Streaming, который будет читать данные и сохранять их в HDFS в формате CSV.
- 3. Обработка и агрегация данных с использованием Apache Hive в <u>RT.DataLake</u>: Создание таблиц Hive для хранения данных из HDFS. Написание запросов для агрегации данных, таких как:
 - общее время просмотра по дням
 - популярность различного контента
 - активность пользователей по времени суток и т.д.

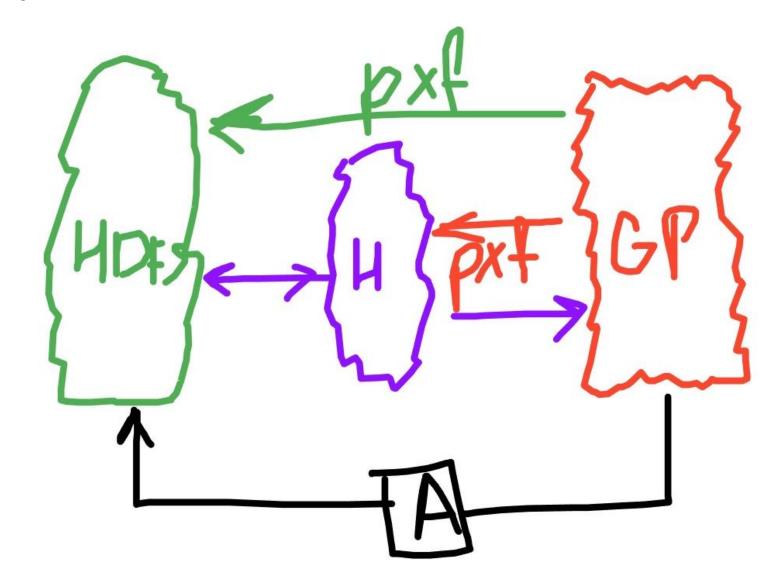
Поставленные задачи



- 4. Перенос данных в GreenPlum (<u>RT.Warehouse</u>) и / или ClickHouse (<u>RT.WideStore</u>) Настройка процесса ETL на основании Apache Airflow / Nifi продукта <u>RT.Streaming</u>, чтобы перенести обработанные данные из HDFS/Hive в GreenPlum/ClickHouse для сложных аналитических запросов.
- 5. Аналитика с использованием Python (Apache Zeppelin + Apache Spark = RT.DataLake) Использование библиотек Python для глубокого анализа данных, выявления инсайтов по данным, прогнозирования поведения пользователей
- 6. Визуализация данных с использованием Apache Superset (продукт RT.DataVision)
 Создание интерактивных дашбордов на основе данных из GreenPlum и ClickHouse.
 Дашборды могут включать в себя графики такие как
 - активности пользователей
 - рейтинг просмотра каналов
 - гистограммы длительности просмотров

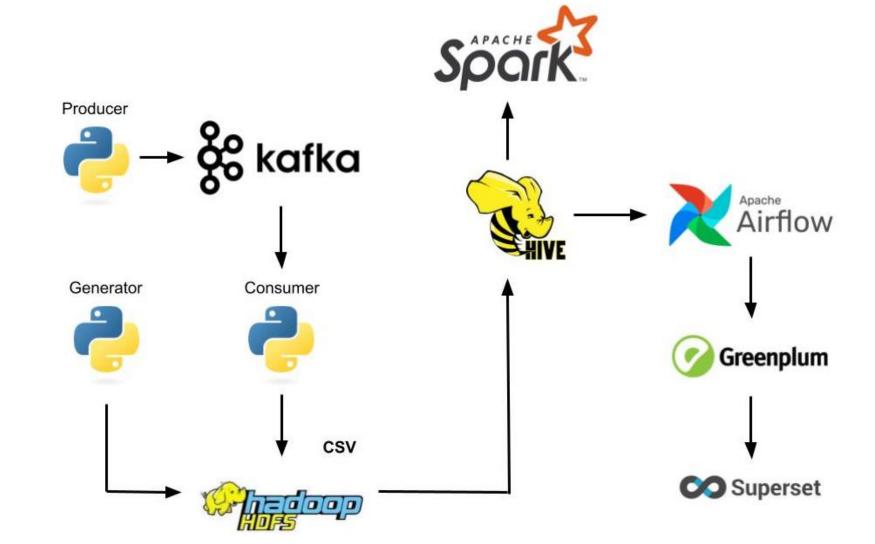
Концептуальная схема



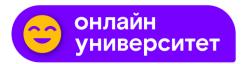


Архитектура MVP



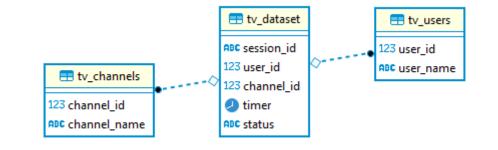


Генерация датасета за предыдущую неделю





Место	Телеканал	Рейтинг по времени пр	Рейтинг по количеству зр
1	Россия 1	7.98	4.04
2	Первый канал	6.79	4.29
3	PEH TB	5.83	3.63
4	НТВ	5.7	3.23
5	THT	4.84	2.85
6	СТС	4.79	3.65
7	Домашний	2.91	1.81
8	Россия 24	2.68	2.38
9	5 Канал	2.66	1.82
10	Пятница	2.46	1.82

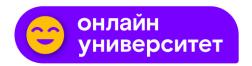


session_id	user_id	channel_id	timer	status
c877f3d1fdfaceecabf95706b2e1ac82	1	1	2023-09-16 10:25:20	enabled
c877f3d1fdfaceecabf95706b2e1ac82	1	1	2023-09-16 17:25:45	disabled
8e489eb90bc7ca76c6b36858e3317a81	2	1	2023-09-16 16:56:03	enabled
8e489eb90bc7ca76c6b36858e3317a81	2	1	2023-09-16 23:57:01	disabled

channel_id	channel_name
1	Россия 1
2	Первый канал
3	PEH TB
4	HTB
5	THT

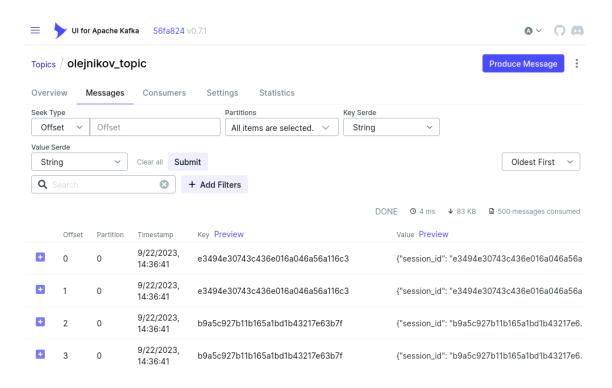
user_id	user_name
1	Рыбаков Егор Иосипович
2	Денисов Эмиль Дмитриевич
3	Копылов Карп Трофимович
4	Никифор Харлампьевич Константинов
5	Большакова Наина Кузьминична

Сбор данных с использованием RT.Streaming



Сбор данных с использованием распределенного программного брокера сообщений Apache Kafka, являющегося компонентом RT.Streaming. продукта данных реализован языке Python. Для этого был продюсер, который симулирует 10000 поведении данные пользователей интерактивного телевидения и отправляет топик Kafka.

Генерируем активность 10000 пользователей за текущий день.



Хранение сырых данных в RT.DataLake



Хранение собранных данных в распределенной файловой системе HDFS, являющейся элементом корпоративного хранилища данных RT.DataLake. Хранение осуществляется с помощью потребителя на Python, который читает данные с Kafka и сохраняет их в HDFS в формате CSV.

□ ↓	Permission	↓↑ ↓ Owner	↑ ↓ Group	Size	Last Modified	Replication	† Block Size	Name	ŢŢ
	drwxrw-rw-	olejnikov	hadoop	0 B	Aug 04 08:00	0	0 B	.Trash	â
	drwxrwxrwx	olejnikov	hadoop	0 B	Aug 02 08:12	0	0 B	.sparkStaging	
	-rw-rr	dr.who	hadoop	3.14 KB	Sep 23 06:23	3	128 MB	tv_channels.csv	â
	-rw-rr	dr.who	hadoop	8.13 MB	Sep 23 06:23	3	128 MB	tv_dataset.csv	
	-rw-rr	olejnikov	hadoop	1.33 MB	Sep 23 08:01	3	128 MB	tv_stream.csv	â
	-rw-rr	dr.who	hadoop	551.02 KB	Sep 23 06:23	3.	128 MB	tv_users.csv	â

Обработка данных с использованием Apache Hive



в RT.DataLake

```
create external table olejnikov.tv_users (
    user_id int,
    user_name varchar(50)
)
row format delimited fields terminated by
','
lines terminated by '\n'
tblproperties("skip.header.line.count"="1");

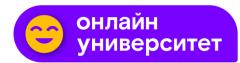
load data inpath
'/user/olejnikov/tv_users.csv' overwrite
into table olejnikov.tv_users;

select * from olejnikov.tv_users limit 5
```

Обработка агрегации данных с использованием СУБД Apache Hive, являющейся компонентом корпоративного хранилища данных RT.DataLake. При ЭТОМ HDFS расположенные на сырые данные, полученные В результате генерации и стриминга загружаются и преобразовываются в таблицы Hive.

₽BC session_id ▼	123 user_id 🔻	123 channel_id 🔻	timer	ABC status 🔻
c877f3d1fdfaceecabf95706b2e1ac82	1	1	2023-09-16 10:25:20.000	enabled
c877f3d1fdfaceecabf95706b2e1ac82	1	1	2023-09-16 17:25:45.000	disabled
8e489eb90bc7ca76c6b36858e3317a81	2	1	2023-09-16 16:56:03.000	enabled
8e489eb90bc7ca76c6b36858e3317a81	2	1	2023-09-16 23:57:01.000	disabled
42e3b49eabda9fd99e2cbff5b5719685	3	1	2023-09-16 16:02:53.000	enabled

Агрегация данных с использованием Apache Hive

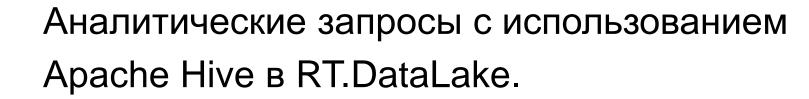


в RT.DataLake

```
create view tv_dataset_v as
with ts as (
    select session id, user id, channel id, timer as time start
    from (select * from tv dataset where status = 'enabled') t1
    union
    select session id, user id, channel id, timer as time start
    from (select * from tv stream where status = 'enabled') t2
), te as (
    select session id, timer as time end
    from (select * from tv dataset where status = 'disabled') t3
    union
    select session id, timer as time end
    from (select * from tv stream where status = 'disabled') t4
select ts.session id, user id, channel name, time start, time end
from ts
join te on ts.session id = te.session id
join tv channels on ts.channel id = tv channels.channel id
select * from tv dataset v limit 5
```

Для агрегации нескольких таблиц в одну используется виртуальная таблица. В ней уже нет колонки со статусом сессии, а непосредственно указаны ее время начали и конца. Вместо id канала стоит его название.

ABC session_id ▼	123 user_id 🔻	ABC channel_name ▼	time_start	② time_end
000718b18ee41c39a73c7aa40eb125cd	1,171	Русский бестселлер	2023-09-23 19:43:10.000	2023-09-23 21:42:10.000
00081155034cba88e8d8afb7aded6f81	66,006	Кинопоказ	2023-09-22 08:32:58.000	2023-09-22 08:33:10.000
000c0ff5e40fa1cc34b42624f67e6dbd	57,893	Иллюзион плюс	2023-09-21 20:01:07.000	2023-09-21 20:01:07.000
00113b378de9ca5d2ae0aa36ee9dd32e	9,375	Уникум	2023-09-16 20:38:52.000	2023-09-16 20:39:43.000
0016e6584372aaf8ae993c7cfb171034	32,208	5 Канал	2023-09-19 01:33:01.000	2023-09-19 03:33:55.000





```
-- Рейтинг каналов
select channel_name, round((sum(UNIX_TIMESTAMP(time_end)
- UNIX_TIMESTAMP(time_start))) / 3600) as sum_time,
count(channel_name) as count_views
from tv_dataset_v
group by channel_name
order by sum_time desc

-- Количество просмотров по дням
select to_date(time_start) as time_dates, round((sum(
UNIX_TIMESTAMP(time_end) - UNIX_TIMESTAMP(time_start)))
/ 3600) as sum_time, count(*) as count_views
from tv_dataset_v
group by to_date(time_start)
order by to_date(time_start)
```

asc channel_name 🔻	123 sum_time 🔻	123 count_views 🔻
Россия 1	15,839	2,313
Первый канал	14,515	2,460
PEH TB	10,703	2,188
НТВ	9,587	1,948
THT	8,526	1,752

1 time_dates	123 sum_time 🔻	123 count_views 🔻
2023-09-16	16,169	9,990
2023-09-17	16,168	9,990
2023-09-18	12,935	8,002
2023-09-19	12,940	7,973
2023-09-20	12,925	8,017
2023-09-21	12,864	7,968
2023-09-22	12,866	7,964
2023-09-23	19,882	10,000

Перенос данных в GreenPlum (RT.Warehouse)

Перенос обработанных данных из Hive в массивнопараллельную СУБД GreenPlum (компонент RT.Warehouse) при помощи ETL процесса – ориентированного ациклического графа Airflow DAG. B Greenplum SQL-диалект свежее и можно построить более сложный запрос – например распределение количества просмотров в разрезе каждого часа.

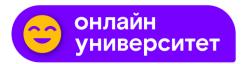


coloct hours count/hours) as
select hours, count(hours) as
count_views
from (
<pre>select generate_series(date_part(</pre>
'hour', time_start::time)::int,
date_part('hour', time_end::time
)::int) as hours
<pre>from olejnikov_tv_mv) h</pre>
group by hours
order by hours



123 hours	•	123 count_views -
	0	504
	1	812
	2	998
	3	1,354
	4	1,674
	5	1,908
	6	2,077
	7	2,099
	8	2,066
	9	2,063
	10	2,051
	11	2,065
	12	1,995
	13	2,045
	14	2,004
	15	1,992
	16	3,853
	17	5,774
	18	12,693
	19	25,103
	20	55,598
	21	32,097
	22	13,509
	23	9,660

Аналитика (Apache Zeppelin + Apache Spark)

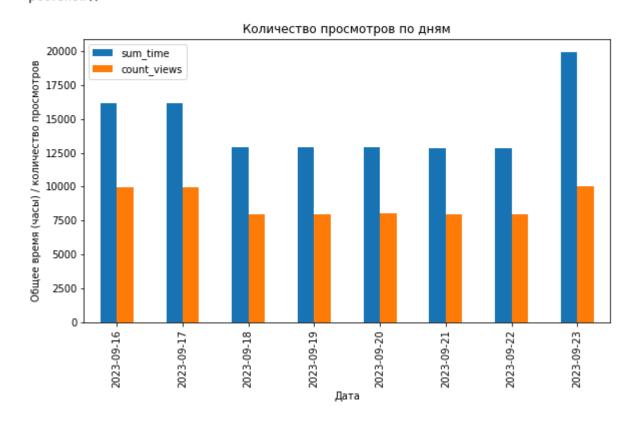


Загружаем виртуальную таблицу из Hive



```
%spark.pyspark
from pyspark.sql import SparkSession
spark = SparkSession.builder\
       .master("local[*]")\
       .appName('olejnikov tv')\
       .enableHiveSupport()\
       .get0rCreate()
df = spark.sql("select * from olejnikov.tv dataset v")
df.show()
          session id|user id|
                               channel name
                                                  time start
  -----+-
                              Первый канал 2023-09-16 17:30:45 2
|52ac12716e958a22c...|
                      503
10312364efb5c75fc3...l
                      682 l
                              Первый канал 2023-09-16 17:33:26 2
                                    PEH TB|2023-09-16 16:03:54|2
|8ee7986b4a5866c42...|
                      926
|392d735affcf7a1fa...|
                      987
                                    PEH TB|2023-09-16 18:54:29|2
|fd80418e594409f98...|
                                    PEH TB|2023-09-16 18:26:41|2
                     1037
[67c66ea5ea6adc22c...]
                                       CTC|2023-09-16 18:09:43|2
                     1806
                                   5 Канал | 2023-09-16 19:17:58 | 2
l f4c73750039200999...l
                     2251
```

```
df2.plot(x='time_dates', y=['sum_time', 'count_views'], kind='bar')
plt.ylabel('Общее время (часы) / количество просмотров')
plt.xlabel("Дата")
plt.title("Количество просмотров по дням")
plt.tight_layout()
plt.show()
```



Аналитика (Apache Zeppelin + Apache Spark)

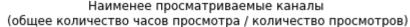


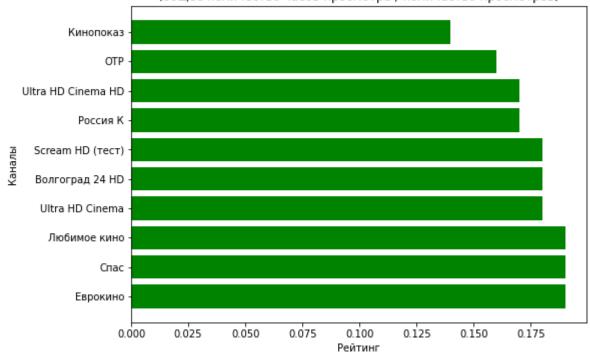
```
hours = df1['sum_time'].head(6)
hours.loc[len(hours.index)] = df1['sum_time'].tail(144).sum()
channels = df1['channel_name'].head(6)
channels.loc[len(channels.index)] = 'Прочие'
myexplode = [0, 0, 0, 0, 0, 0.1,]
fig, ax = plt.subplots()
ax.pie(hours, labels=channels, autopct='%1.2f%%', explode = \
myexplode, shadow = True)
plt.title("Топ просматриваемых каналов")
plt.show()
```

Топ просматриваемых каналов



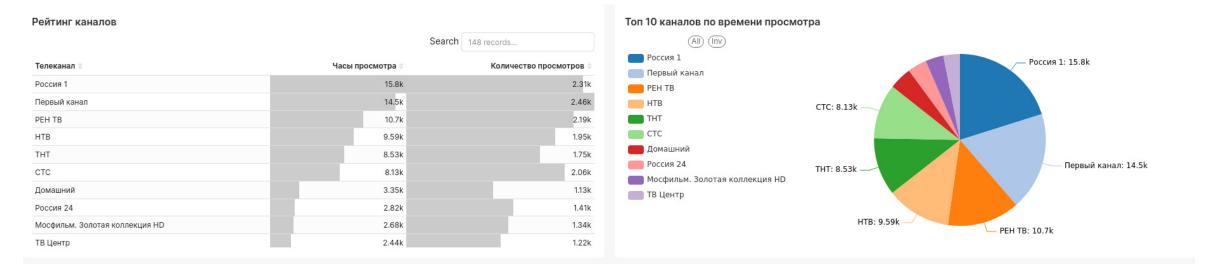






Визуализация данных в Apache Superset



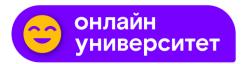


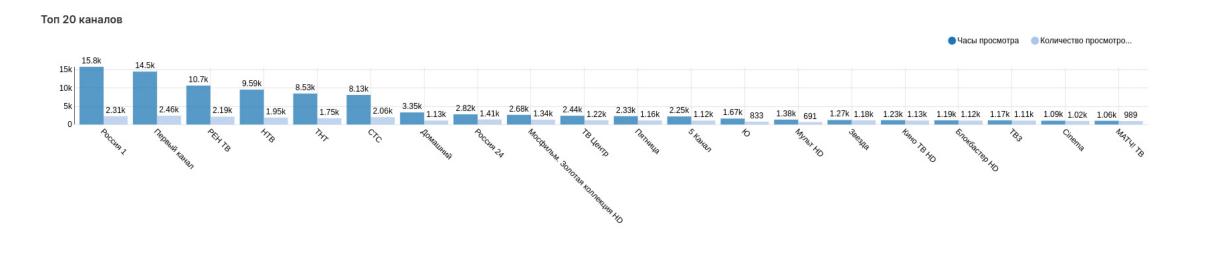
Дата 🗘	Часы просмотра 🛊	Количество просмотров
2023-09-23	19.9k	10k
2023-09-16	16.2k	9.99k
2023-09-17	16.2k	9.99k
2023-09-19	12.9k	7.97k
2023-09-18	12.9k	8k
2023-09-20	12.9k	8.02k
2023-09-22	12.9k	7.96k
2023-09-21	12.9k	7.97k

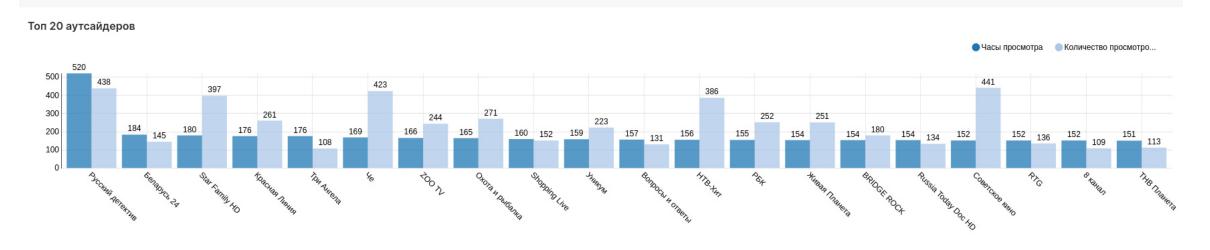




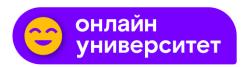
Визуализация данных в Apache Superset







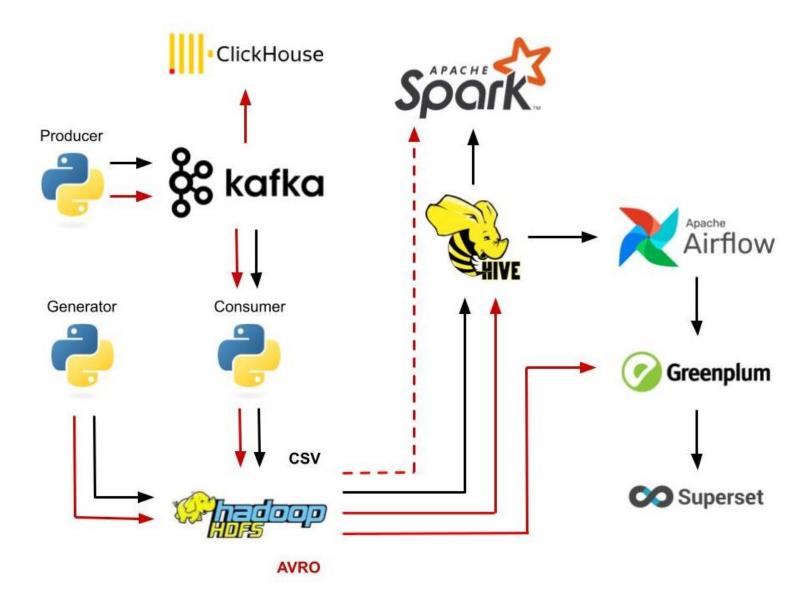
Расширенные задачи



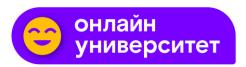
- 1. Реализовать долгосрочное хранение датасета в AVRO-формате.
- 2. Реализовать стриминг так, чтоб данные одновременно попадали на HDFS и ClickHouse.
- 3. Реализовать загрузку датасета непосредственно из HDFS в Hive, Spark, GreenPlum.
- 4. Реализация расширенной части задания не должна влиять на выполнение основной части.

Архитектура MVP - расширенная



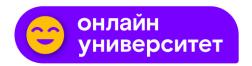


Расширенная часть задания



```
%spark.pyspark
olejnikov
                         1002.35 KB tv stream.avro
            hadoop
                                                                              avro df = spark.read.format("avro").load('/user/olejnikov/tv stream.avro')
olejnikov
            hadoop
                         1.33 MB
                                   tv stream.csv
                          CREATE TABLE olejnikov.tv avro
                              ROW FORMAT SERDE
                                   'org.apache.hadoop.hive.serde2.avro.AvroSerDe'
                              STORED AS INPUTFORMAT
                                   'org.apache.hadoop.hive.gl.io.avro.AvroContainerInputFormat'
                              OUTPUTFORMAT
                                   'org.apache.hadoop.hive.ql.io.avro.AvroContainerOutputFormat'
                              TBLPROPERTIES ('avro.schema.literal'='
                              "doc": "olejnikov avro schema",
                              "name": "tv stream",
                              "namespace": "tv stream",
                              "type": "record",
                              "fields": [
                                   {"name": "session id", "type": "string"},
                                   {"name": "user id", "type": "int"},
                                   {"name": "channel id", "type": "int"},
                                   {"name": "timer", "type": {"type": "long", "logicalType": "timestamp-millis"}},
                                  {"name": "status", "type": "string"}
                          }')
                          LOAD DATA INPATH '/user/olejnikov/tv stream.avro' OVERWRITE INTO TABLE olejnikov.tv avro;
```

Итоги и инсайты

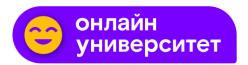


По итогам работы удалось реализовать систему анализа и сбора данных услуги интерактивного телевидения с использованием продуктов платформы управления данными ПАО «Ростелеком».

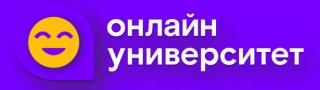
При помощи аналитических инструментов выявлены наиболее популярные каналы у пользователе как по количеству, так и по длительности просмотра. А так же каналы аутсайдеры. Выявлено использование услуги интерактивного телевидения в разрезе часов, дней недели.

На основании запросов были построены аналитические графики, витрина данных.

Список ссылок



- 1. GitHub проекта: https://github.com/Olmeor/Data-engineer_Rostelecom_programming_school
- 2. Рейтинг популярности телеканалов: https://www.powernet.com.ru/channels-stat
- 3. Аналитика просмотров: https://journal.tinkoff.ru/television-stat/
- 4. Spark notebook http://vm-dlake2-m-2:8180/#/notebook/2JBN33XNC (только внутри кластера)
- 5. Superset http://vm-datavision.test.local:8090/superset/dashboard/p/rYymq0Yxn3R/ (только внутри кластера)
- 6. Платформа управления данными ПАО «Ростелеком» https://data.rt.ru/products



Спасибо за внимание!

Готов ответить на ваши вопросы

Разработал:

Олейников Михаил Николаевич