



# Цель модуля “SQL поверх больших данных (Hive)”

**Драль Алексей**, [study@bigdatateam.org](mailto:study@bigdatateam.org)

CEO at BigData Team, <https://bigdatateam.org>

<https://www.facebook.com/bigdatateam>



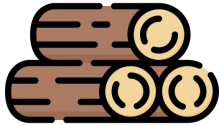
# 1. Наиболее популярные регионы





**BIGDATA  
TEAM**

# 1. Наиболее популярные регионы



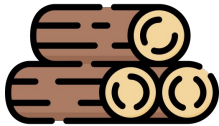
**Web-service  
access logs**



# 1. Наиболее популярные регионы



**Web-service  
access logs**



**Geobase**

**IPv4: 109.188.67.224**



**area: Moscow City Center**



**city: Moscow**



**country: Russia**

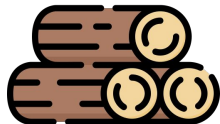


**Earth**

**109.188.67.224, Moscow City Centre, Moscow, Russia, Earth**  
**109.188.67.221, Moscow City Centre, Moscow, Russia, Earth**  
...



# 1. Наиболее популярные регионы



**Web-service access logs**

format: ip, request, status\_code,...



**Geobase**

format: ip, region<sub>city</sub>, region<sub>country</sub>,...

Join (ip)

+

WordCount(region)

+

TOP(100)  
Sort + LIMIT



format: region, hit\_count



# 1. Наиболее популярные регионы



**Web-service access logs**

format: ip, request, status\_code,...



**Geobase**

format: ip, region<sub>city</sub>, region<sub>country</sub>,...

Join (ip)



+

WordCount(region)



+

TOP(100)  
Sort + LIMIT



format: region, hit\_count

**WordCount**

Hive QL



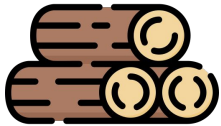
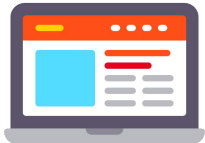
```
SELECT regioncity, COUNT(1) AS hit_count  
FROM access_log JOIN geo_base  
ON (access_log.host = geo_base.host)  
GROUP BY regioncity ORDER BY hit_count LIMIT 100
```

**join**

**TOP(100)**

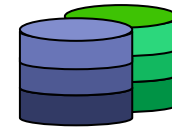


## 2. Доля роботных запросов



**Web-service  
access logs**

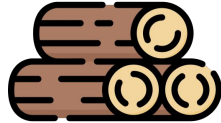
bot name + user\_agent<sub>1</sub> + ip<sub>1</sub>  
user\_agent<sub>2</sub> + ip<sub>2</sub>  
...



**[ro]bot database**

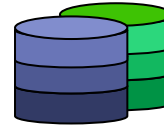


## 2. Доля роботных запросов



### Web-service access logs

format: ip, request, user\_agent,...



### [ro]bot database

format: bot\_name, user\_agents, ips

Join(ip, user\_agent) + WordCount(request / user, bot)

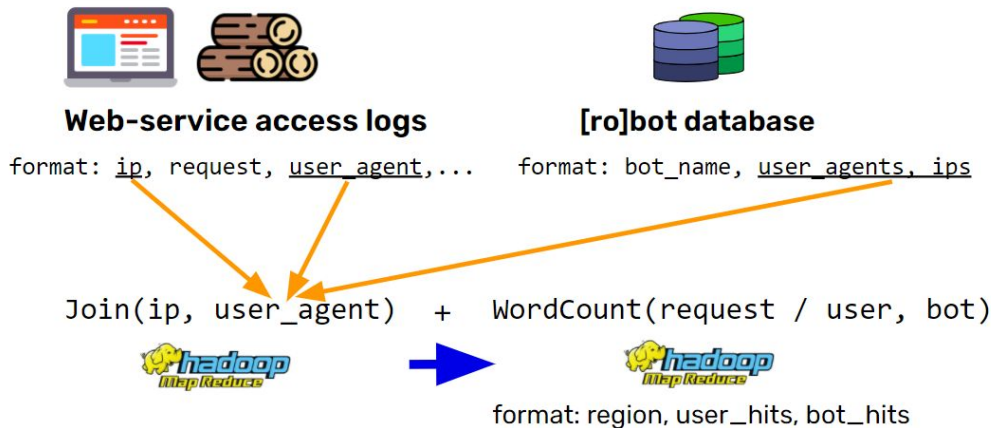


format: region, user\_hits, bot\_hits





## 2. Доля роботных запросов



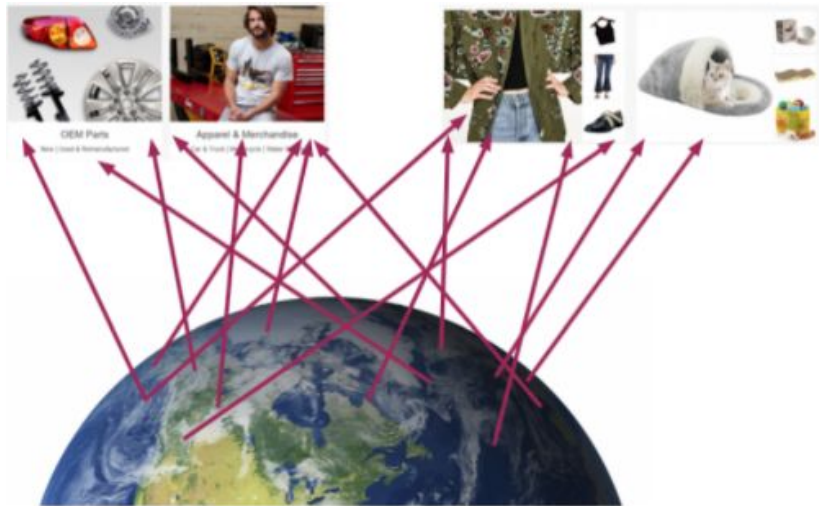
```
SELECT request,  
       SUM(IF(robot.bot_name IS NULL, 1, 0)) as user_hit_count,  
       SUM(IF(robot.bot_name IS NOT NULL, 1, 0)) as bot_hit_count  
FROM access_log LEFT OUTER JOIN robot ON (  
    access_log.host = robot.host  
    AND access_log.user_agent = robot.user_agent  
)  
GROUP BY request
```

Hive QL

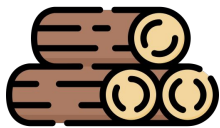




### 3. Гендерное распределение



**Web-service  
access logs**



**Geobase**

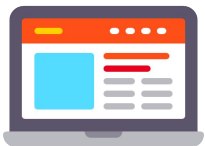


**User personal data**

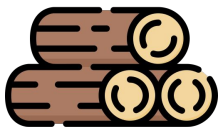
user_id	age	gender	occupation	zipcode
1	24	M	technician	85711
2	53	F	other	94043
3	23	M	writer	32067
4	24	M	technician	43537
5	33	F	other	15213
6	42	M	executive	98101
7	57	M	administrator	91344
8	36	M	administrator	05201
9	29	M	student	01002
10	53	M	lawyer	90703



### 3. Гендерное распределение в регионах



**Web-service access logs**



**User personal data**



**Geobase**

format: user\_id, gender, age,...

format: ip, request, user\_agent,...

format: ip, region, city,  
region, country',...

Join(user\_id)

+

Join(ip)

+

WordCount(region/gender)



user\_id = ip + user\_agent

format:  
region, male\_hits, female\_hits



# 3. Гендерное распределение в регионах



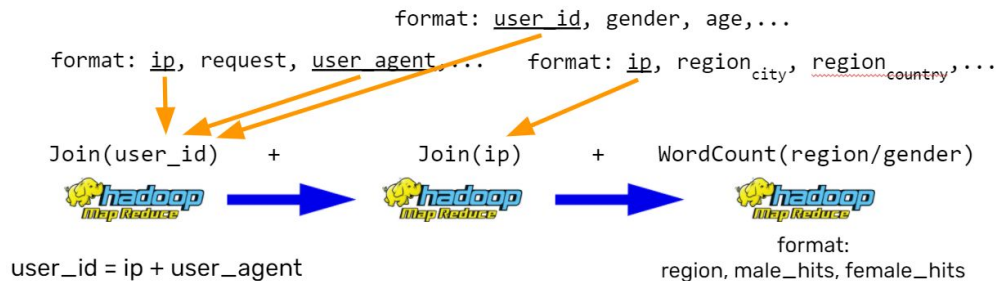
Web-service access logs



User personal data



Geobase



Hive QL



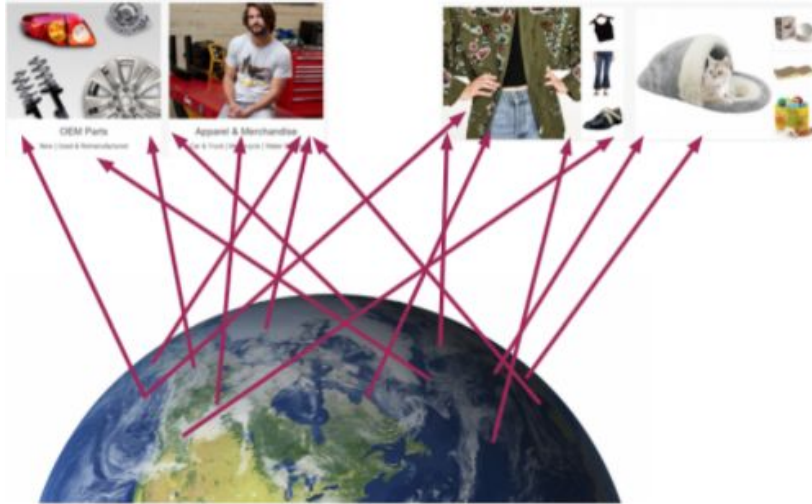
```
SELECT regioncity,  
       SUM(IF(user.gender = "M",1,0)) as male_hit_count,  
       SUM(IF(user.gender = "F",1,0)) as female_hit_count  
FROM access_log  
JOIN geo_base ON (access_log.host = geo_base.host)  
JOIN user ON (access_log.host = user.host  
              AND access_log.user_agent = user.user_agent  
             )  
GROUP BY regioncity
```

two joins

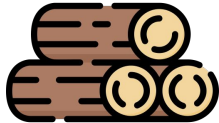


**BIGDATA  
TEAM**

## 4. Средний возраст клиента



**Web-service  
access logs**



**Geobase**



**User personal data**

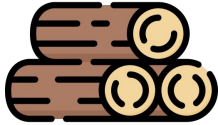
user_id	age	gender	occupation	zipcode
1	24	M	technician	85711
2	53	F	other	94043
3	23	M	writer	32067
4	24	M	technician	43537
5	33	F	other	15213
6	42	M	executive	98101
7	57	M	administrator	91344
8	36	M	administrator	05201
9	29	M	student	01002
10	53	M	lawyer	90703



## 4. Средний возраст клиента



**Web-service access logs**



**User personal data**



**Geobase**

format: user\_id, gender, age,...

format: ip, request, user\_agent,...

format: ip, region, city,  
region, country',...

Join(user\_id) +

Join(ip) +

**Average**(age)



format: region, average\_age

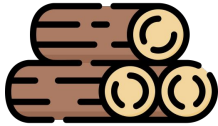
user\_id = ip + user\_agent



## 4. Средний возраст клиента



**Web-service access logs**



**User personal data**



**Geobase**

format: user\_id, gender, age,...

format: ip, request, user\_agent,...

format: ip, region, city,  
region, country',...

Join(user\_id) +

Join(ip) +

**Average**(age)



format: region, average\_age

**use Combiner for optimisation**

user\_id = ip + user\_agent



## 4. Средний возраст клиента



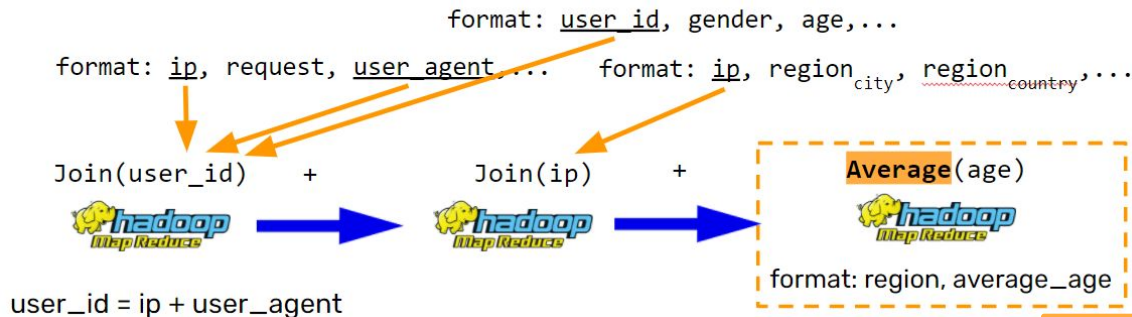
Web-service access logs



User personal data



Geobase



GenericUDAFAverage  
Hive Java source code:  
<https://rebrand.ly/hql-avg>

Hive QL



```
SELECT regioncity, AVG(user.age)
FROM access_log
  JOIN geo_base ON (access_log.host = geo_base.host)
  JOIN user ON (access_log.host = user.host
               AND access_log.user_agent = user.user_agent
              )
GROUP BY region2
```





**BIGDATA  
TEAM**

План



- ▶ Map-Side Join
- ▶ Reduce-Side Join
- ▶ Bucket Map-Side и SMB Join



- ▶ Map-Side Join
- ▶ Reduce-Side Join
- ▶ Bucket Map-Side и SMB Join

А также:

- ▶ Как теряют production данные
- ▶ Как правильно пользоваться RegExpSerDe