

HW #10: Data Layout

1. Описание задания	2
2. Критерии оценивания	2
3. Описание данных	2
4. Задания	3
4.1. Задание #1, Task ID: hive.optimize_storage	4
4.2. Задание #2, Task ID: hive.speedup_query	4
4.3. Задание #3, Task ID: hive.skew	5
4.4. Задание #4, Task ID: hive.optimize_aggregate	5
5. Правила оформления задания	6

автор задания:

- Драль Алексей, aadral@bigdatateam.org



1. Описание задания

В этом задании будем оптимизировать производительность хранилища и скорость выполнения аналитических запросов с помощью правильного выбора Data Layout. Нужно решить **4 задачи**. Для решения используем Hive.

Сами задания несложные, но на выходе вы получите полезные скрипты, которые сможете применять для оптимизации работы с вашими данными на работе.

Полезные материалы:

- [stackoverflow: использование конструкции --hivevar;](#)

2. Критерии оценивания

Веса задач:

1. 50%
2. 50%
3. 0% (для самостоятельного изучения)
4. 0% (для самостоятельного изучения)

Discounts (скидки и другие акции):

- **100%** за плагиат в решениях (всем участникам процесса)
- **100%** за посылку решения после deadline
- **5%** за каждую новую посылку (одна дополнительная посылка бесплатно)

Формула подсчета финальной оценки¹:

$\max(0, 0.95^{\max(0, \# \text{доп.посылок} - 1)} * (1 - \text{штраф. за дедлайн. и списывание})) * \text{последняя. оценка. из. grader}$

3. Описание данных

3.1. Логи запросов пользователей новостных сайтов.

logs_raw:

- Путь на кластере: полный датасет - /data/user_logs/user_logs_M
- Семпл (для тестирования): /data/user_logs/user_logs_S
- Формат: текст

¹ результат округляется до целого



- В каждой строке находятся следующие поля, разделенные знаком табуляции (иногда не одним):
 1. ip STRING - ip-адрес, с которого пришел запрос,
 2. date STRING - время запроса,
 3. request STRING - пришедший с ip-адреса http-запрос,
 4. page_size INT - размер переданной клиенту страницы в байтах,
 5. http_status INT - http-статус запроса.
 6. user_agent STRING - User Agent, информация о клиентском приложении, с которого осуществлялся запрос на сервер, в том числе, информация о браузере.

Пример:

```
135.124.143.193          20150601013300
http://newsru.com/4712386 235 412 Firefox/5.0 (compatible; MSIE
9.0; Windows NT 6.1; Win64; x64; Trident/5.0)n
```

Важно:

- разделитель между IP и временем запроса состоит из 3 символов табуляции;
- Будем считать, что информация о браузере содержится в начале 6-ого поля лога - символы с нулевой позиции до позиции первого пробельного символа.
 - пример User Agent:
 - Chrome/5.0 (compatible; MSIE 9.0; Windows NT 8.0; WOW64; Trident/5.0; .NET CLR 2.7.40781; .NET4.0E; en-SG)
 - тогда браузером будет: Chrome/5.0

Подсказка:

- поскольку нас не интересует оставшаяся часть User Agent, то получить тип браузера пользователя можно с помощью правильного регулярного выражения в период чтения logs_raw.

4. Задания

В рамках решения ДЗ по Hive, у вас появилась таблица с логами пользователей новостных сайтов `logs`. Вам предлагается решить следующие задачи (отработать задачи на семплах _S, _M и получить решение или оценки роста производительности для полного датасета). Рекомендуется использовать Managed таблицы и перезаписывать logs_ с помощью INSERT OVERWRITE запроса.

4.1. Задание #1, Task ID: hive.optimize_storage

Переложите данные logs_raw в таблицу logs_orc, где будет использоваться формат хранения данных ORC. С помощью параметров TBLPROPERTIES найдите оптимальный набор параметров, чтобы получить максимальное сжатие данных.

Проверка будет производиться на датасете _M с помощью следующего кода:

```
CREATE TABLE logs_orc
STORED AS orc
TBLPROPERTIES (
    <content of your HQL is here>2
)
AS SELECT *
FROM logs_raw;
```

Балл за задачу складывается из:

- **0%** - правильное решение задачи
- **0%** - поддерживаемость и читаемость кода
 - в общем случае см. Clean Code и [Google Python Style Guide](#)
- **100%** - эффективность решения, для ориентира - размер данных в HDFS в эталонном решении на порядок меньше, чем объем данных в датасете _M (проверяем размер с помощью `hdfs dfs -du -s /path/to/table`).

Вопрос для самостоятельной проработки: какой оптимизации пространства удалось добиться для датасетов _S, _M и _full? Сохраняется ли динамика между _M и _full?

4.2. Задание #2, Task ID: hive.speedup_query

Придумайте аналитические запросы, которые должны работать быстрее за счет использования ORC. Проверьте скорость выполнения таких запросов на таблицах logs_raw и logs_orc. Какая оптимизация по скорости выполнения получена в зависимости от типа запроса? Сделайте релевантные таблицы для датасетов _S, _M и _full и сравните наблюдения. Производительность решения будет проверяться на датасете _full.

Балл за задачу складывается из:

- **0%** - правильное решение задачи

² Таким образом вам нужно сохранить в HiveQL файл только свойства ORC файла для DDL



- **0%** - поддерживаемость и читаемость кода
 - в общем случае см. Clean Code и [Google Python Style Guide](#)
- **100%** - эффективность решения:
 - **80%** - использование MapReduce CPU Time должно быть в разы меньше в случае использования logs_orc (эталонное решение работает эффективнее более чем 4.7 раза).
 - **20%** - wall time выполнения задачи (эталонное решение работает в 1.25 раза быстрее)

Сохраните ваш запрос в HiveQL файле, где название таблицы `${table_name}` для работы будет передаваться через `hivevar`. Оптимизированная таблица `logs` - это данные в формате ORC со значениями TBLPROPERTIES по умолчанию.

4.3. Задание #3, Task ID: hive.skew

Для самостоятельного изучения

Попробуйте заменить в логах информацию про браузер таким образом, чтобы 90% данных содержало одинаковый браузер (или браузер "unknown"), запишите результат в таблицу `logs_broken`. Попробуйте посчитать запрос в задаче "identify browser sex". Оцените время на выполнение запроса. Для того, чтобы пофиксить проблему:

1. В реальной жизни рекомендуется сделать запрос в формате TABLESAMPLE, чтобы увидеть по каким параметрам происходит перекос;
2. Теперь вы знаете по каким данным происходит перекос, дайте эту информацию в формате SKEWED TABLE для Hive.

Оцените скорость выполнения запроса для датасетов `_S`, `_M` и `_full`. Не забывайте отслеживать параметр числа редьюсеров, если их недостаточно для выполнения запроса.

4.4. Задание #4, Task ID: hive.optimize_aggregate

Для самостоятельного изучения

Придумайте запрос, содержащий конструкцию GROUP BY или JOIN, который можно выполнить на стадии Map с помощью правильной укладки данных. Под правильной укладкой данных подразумевается бакетирование и сортировка данных. Сколько времени тратится на переукладку данных? Какая оптимизация по скорости выполнения запроса получена?

5. Правила оформления задания

Оформление задания:

- Код задания (Short name): **HW10:DataLayout**
- Выполненное ДЗ запакуйте в архив **MADEBD2021Q1_<Surname>_<Name>_HW#.zip**, например, для Алексея Драля -- **MADEBD2021Q1_Dral_Alexey_HW10.zip**. (Проверяйте отсутствие пробелов и невидимых символов после копирования имени отсюда.³) Если ваше решение лежит в папке `my_solution_folder`, то для создания архива `hw.zip` на Linux и Mac OS, выполните команду⁴:
 - `zip -r hw.zip my_solution_folder/*`
- На Windows 7/8/10: необходимо выделить все содержимое директории `my_solution_folder/` нажать правую кнопку мыши на одном из выделенных объектов, выбрать в открывшемся меню "Отправить >", затем "Сжатая ZIP-папка". Теперь можно переименовать архив.
- По результатам решения ожидается отчет в формате PDF с описанием результатов оптимизации (ответов на поставленные исследовательские вопросы).
- HQL-скрипты для запуска решений следует называть по суффиксу Task ID задачи **task_<Surname>_<Name>_<#task_ID_suffix>.hql**:
 - например решение задачи 2 должно называться `task_<Surname>_<Name>_speedup_query.hql` и его можно запустить с помощью команды:

```
$ hive -v --database=${DB_NAME}5 --hivevar  
table_name=${table_name} -f task_*_speedup_query.hql
```
- Перед проверкой убедитесь, что дерево вашего архива выглядит так:
 - | **MADEBD2021Q1_<Surname>_<Name>_HW10.zip**
 - | ---- **task_<Surname>_<Name>_optimize_storage.hql**
 - | ---- **task_<Surname>_<Name>_speedup_query.hql**
 - | ---- **task_<Surname>_<Name>_skew.hql (optional)**
 - | ---- **task_<Surname>_<Name>_optimize_aggregate.hql (optional)**
 - При несовпадении дерева вашего архива с представленным деревом, ваше решение будет невозможно автоматически проверить, а значит, и оценить его.
- Для того, чтобы сдать задание необходимо:
 - Зарегистрироваться и залогиниться в сервисе [Everest](#)
 - Перейти на страницу приложения: [BDT-grader-MADE-BD](#)

³ Онлайн инструмент для проверки: <https://www.soscisurvey.de/tools/view-chars.php>

⁴ Флаг `-r` значит, что будет совершен рекурсивный обход по структуре директории

⁵ Это означает, что Вы не должны использовать "use <database_name>" внутри скриптов



- Выбрать вкладку Submit Job (если отображается иная).
- Выбрать в качестве "Task" значение: **HW10:DataLayout**⁶
- Загрузить в качестве "Task solution" файл с решением
- В качестве Sender ID указать тот, который был выслан по почте
- Если Вы видите надпись "You are not allowed to run this application" во вкладке Submit Job в Everest, то на данный момент сдача закрыта (нет доступных для сдачи домашних заданий, по техническим причинам или другое). Попробуйте, пожалуйста, еще раз через некоторое время. Если Вы еще ни разу не сдавали, у коллег сдача работает, но Вы видите такое сообщение, сообщите нам об этом.
- Ситуации:
 - * система оценивания показывает оценку (Grade) < 0, а отчет (Grading report) не помогает решить проблему (пример помощи: в случае неправильно указанного Sender ID система вернет -2 и информацию о том, что его нужно поправить);
 - * показывает 0 и в отчете (Grading report) не указано, какие тесты не пройдены. Если Вы столкнулись с какой-то из них присылайте ссылку на выполненное задание (Job) на почту с темой письма "Short name. ФИО.". Например: **"HW10:DataLayout. Иванов Иван Иванович."**Пример ссылки: <https://everest.distcomp.org/jobs/67893456230000abc0123def>
Внимание: Если до дедлайна остается меньше суток, и Вы знаете (сами проверили или коллеги сообщили), что сдача решений сломана, обязательно сдайте свое решение и напишите письмо, как написано выше, чтобы мы видели, какое решение Вы имели до дедлайна и смогли его оценить.
- Перед отправкой задания, оставьте, пожалуйста, отзыв о домашнем задании по ссылке: http://rebrand.ly/mailbd2021q1_feedback_hw. Это позволит нам скорректировать учебную нагрузку по следующим заданиям (в зависимости от того, сколько часов уходит на решение ДЗ), а также ответить на интересующие вопросы.

Любые вопросы / комментарии / предложения можно писать в телеграм-канал курса или на почту bigdata_made2021q1@bigdatateam.org.

Peace, love, обнимашки, интересности скидываем в общий чат курса :)

⁶ Сервисный ID: `hive.layout_hw`