# Leveraging SAP HANA & Hortonworks Data Platform to analyze Wikipedia Page Hit Data

## 1 Introduction

SAP HANA is the leading OLTP and OLAP platform delivering instant access and critical business insight using a single in-memory columnar data store. Hortonworks Data Platform is the leading 100% open source, fully tested and certified Apache Hadoop platform enabling businesses to store, enrich and analyze the ever-increasing volumes of observational data such as web-clicks, sensor readings and social sentiments. Together SAP HANA and Hortonworks Data platform (HDP) offer unique opportunity for enterprises to cost effectively capture all enterprise data (ERP & non-ERP) and analyze them in real time to drive efficiency and enable new business opportunities.

This demonstration showcases a common integration scenario of SAP HANA and HDP, where-in raw observational data (web-clicks) is collected in Hortonworks Data Platform. The raw data is enriched and filtered in Hortonworks Data Platform to create data subsets critical for further analysis. Filtered dataset is subsequently transferred to HANA for real-time aggregation and predictive analytics.

## 2 Audience

The demonstration is targeted for solution architects and field engineers interested in understanding and showcasing the joint capabilities of SAP HANA and Hortonworks Data Platform.

Demonstration setup instructions leverage prior work and refer to them where appropriate.

## 3 Prerequisites

To understand and enable a working demonstration, following conditions must be met:
1. User must have basic knowledge of Linux
2. User must have an Amazon Web Services (AWS) account
   - Users sign up for an AWS account at http://aws.amazon.com
3. User must have a basic understanding of AWS

# 4 Cost

The demonstration leverages SAP HANA Developer edition and Hortonworks Data Platform, both available at no software charge. However, the user is responsible for the Amazon EC2 instance charges that serve as the hardware infrastructure for the demonstration scenario. Based on recommended hardware configurations, the following charges are anticipated:

| Server | Specification | Projected Cost/Per Hour |
| --- | --- | --- |
| EC2 server for SAP HANA Developer Edition | EBS Optimized M1.xlarge 250 GB EBS Block Storage | $0.65* |
| EC2 server for Hortonworks Data Platform | M2.xlarge 130 GB EBS Block Storage | 0.550* |
| Total Cost | | $1.20/per hour |

* Cost per hour is an approximation based on EC2 rates as of December 2013. Above server specifications are recommended for balanced cost and performance.

# 5 Demonstration Scenario

The demonstration leverages HDP and HANA to analyze Wikipedia web click data. Wikipedia maintains webclick statistics for it's pages and makes them available for download as compressed files. The demonstration illustrates the following scenario:

1) Download Wikipedia pagehit data and load it to HDP
2) Leverage HDP to filter raw data, retaining pages with hit count of over 99
   o Filtering focuses the analysis on high value pages
3) Retrieve filtered data in HDP into SAP HANA using Smart Data Access
4) Develop analytical view and visualize the data using SAP Hana Studio

# 6 Setup Instructions

Following steps must be followed to setup the demonstration:
1) Setup SAP HANA, developer edition on Amazon Web Services (AWS) and HANA Studio
   a. Note: SAP HANA on AWS is not mandatory -- if you have access to your own instances of SAP HANA, you can leverage your current instance.

2) Setup Hortonworks Data Platform on AWS
3) Install Hortonworks Hive ODBC Driver on SAP HANA AWS server
4) Access HDP from SAP HANA using Smart Data Access

## 6.1 Setup SAP HANA Developer Edition on AWS & SAP HANA Studio

Detailed step-by-step instruction to create SAP HANA Developer Edition on AWS and setting HANA Studio has been published by SAP. Follow the instructions as outlined in the document http://scn.sap.com/docs/DOC-28294.

## 6.2 Setup Hortonworks Data Platform on AWS

➢ Log into your AWS console and click on the AMI menu in the left navigation pane

➢ Search for the AMI "HDP2_HANA_DEMO_AMI" in the **US East Region**
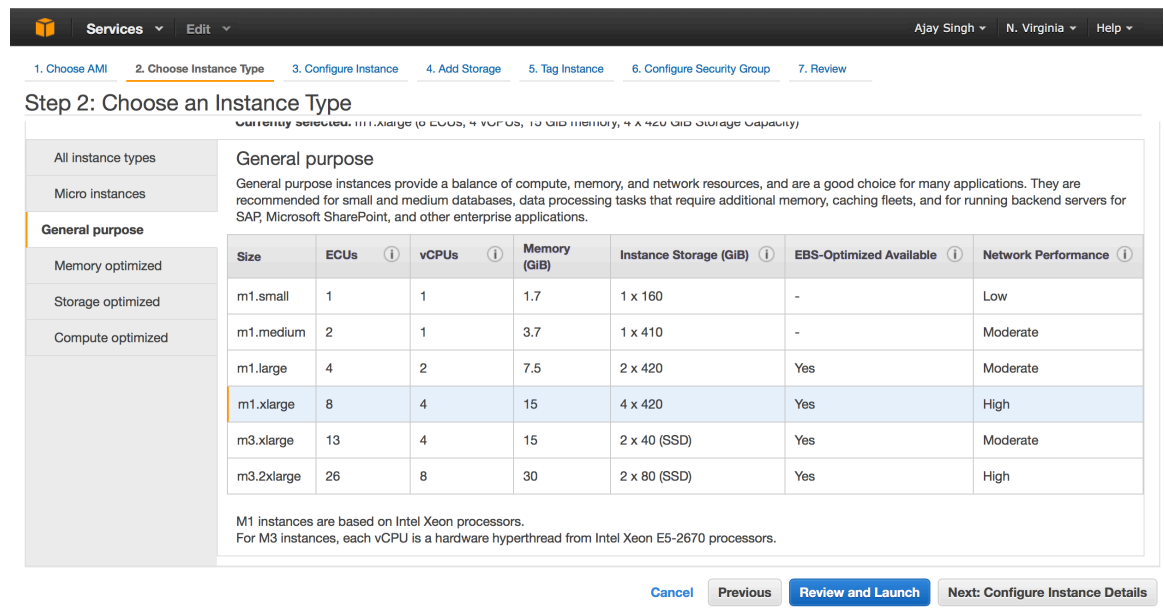


➢ Select the HDP2_HANA_DEMO_AMI and start it using the Launch button at the top.



➢ In the resulting wizard, select m1.xlarge for the Instance Type by selecting General purpose menu on the left.  Then, click "Next Configure Instance Details".



➢ In the step 3 "Configure Instance Details", select EBS-optimized instance and leave the remaining values as their defaults.  Click Next:Add Storage when done.

1. Choose AMI    2. Choose Instance Type    3. Configure Instance    4. Add Storage    5. Tag Instance    6. Configure Security Group    7. Re

## Step 3: Configure Instance Details

Configure the instance to suit your requirements. You can launch multiple instances from the same AMI, request Spot Instances to take adv management role to the instance, and more.

| | |
|---|---|
| Number of instances ⓘ | 1 |
| Purchasing option ⓘ | ☐ Request Spot Instances |
| Network ⓘ | Launch into EC2–Classic ⬍   C Create new VPC |
| Availability Zone ⓘ | No preference ⬍ |
| IAM role ⓘ | None ⬍ |
| Shutdown behavior ⓘ | Stop ⬍ |
| Enable termination protection ⓘ | ☐ Protect against accidental termination |
| Monitoring ⓘ | ☐ Enable CloudWatch detailed monitoring<br>Additional charges apply. |
| EBS-optimized instance ⓘ | ☑ Launch as EBS-optimized instance<br>Additional charges apply. |

➢ Step 4 "Add Storage" leave the defaults and select Next: Tag Instance

1. Choose AMI    2. Choose Instance Type    3. Configure Instance    4. Add Storage    5. Tag Instance    6. Configure Security Group    7. Rev

## Step 4: Add Storage

Your instance will be launched with the following storage device settings. You can attach additional EBS volumes and instance store volume edit the settings of the root volume. You can also attach additional EBS volumes after launching an instance, but not instance store volumes storage options in Amazon EC2.

| Type ⓘ | Device ⓘ | Snapshot ⓘ | Size (GB) ⓘ | Volume Type ⓘ | IOPS ⓘ | Delete o |
|---|---|---|---|---|---|---|
| Root | /dev/sda1 | snap-c2c0b2d2 | 20 | Standard ⬍ | N/A | ☑ |
| EBS ⬍ | /dev/sdc ⬍ | snap-c5c0b2d5 | 50 | Standard ⬍ | N/A | ☑ |
| EBS ⬍ | /dev/sdd ⬍ | snap-c8c0b2d8 | 50 | Standard ⬍ | N/A | ☑ |
| EBS ⬍ | /dev/sde ⬍ | snap-cec0b2de | 50 | Standard ⬍ | N/A | ☑ |
| EBS ⬍ | /dev/sdf ⬍ | snap-f4c0b2e4 | 50 | Standard ⬍ | N/A | ☑ |
| EBS ⬍ | /dev/sdb ⬍ | snap-fdc0b2ed | 50 | Standard ⬍ | N/A | ☑ |

**Add New Volume**

➢ Step 5 "Tag Instance", add the tag "Name = HDP".  Click Next: Configure Security Group when done.

> Step 6 "Configure Security Group", select "Create a new security group" with the following rules

| Protocol | Type | Port Range | Source |
|---|---|---|---|
| • SSH | • TCP | • 22 | • Anywhere |
| • Custom TCP Rule | • TCP | • 10000 | • Anywhere |



Once done select "Review and Launch"

➢ In the review window, validate the following
  o Instance Type matches the settings below:
    ▪ Instance Type: m1.xlarge
    ▪ EBS-Optimized: Yes
    ▪ Network Performance: High

  o Security Groups matches the settings below

| Protocol | Type | Port Range | Source |
|---|---|---|---|
| • SSH | • TCP | • 22 | • 0.0.0.0/0 |
| • Custom TCP Rule | • TCP | • 10000 | • 0.0.0.0/0 |

  o Click Launch when done

➢ Select a key in the resulting window and launch the instance.

Launch Status

The following instance launch has been initiated: i-fc1bf0d2    View launch log

💬 **Get notified of estimated charges**
Create billing alerts to get an email notification when estimated charges on your AWS bill exceed $0.0 (in other words, when you have exceeded the free usage tier).

How to connect to your instance

Your instance is launching, and it may take a few minutes until it is in the **running** state, when it will be ready for you to use. Usage hours on your new instance will start immediately and continue to accrue until you stop or terminate your instance.

Click **View Instances** to monitor your instance's status. Once your instance is in the **running** state, you can **connect** to it from the Instances screen. Find out how to connect to your instance.

▼ Here are some helpful resources to get you started

• How to connect to your Linux instance        • Amazon EC2: User Guide
• Learn about AWS Free Usage Tier               • Amazon EC2: Discussion Forum

While your instances are launching you can also

  Create status check alarms to be notified when these instances fail status checks. (Additional charges may apply)
  Create and attach additional EBS volumes (Additional charges may apply)
  Manage security groups

**View Instances**

➢ Navigate to the Instances view and note the public DNS & IP of the launched instance. We will use the public DNS & IP to establish a connection from SAP HANA.



The instance will take a few minutes to initialize. Upon initialization, the instance will contain a running instance of Hortonworks Data Platform with following key services enabled:

- HDFS
- YARN
- MapReduce HistoryServer
- Amabri
- Hive
- Zookeeper

## Sample Data Set

The Hortonworks Data Platform comes pre-loaded with Wikipedia page hit data for the Months of March and April, 2014.  Following filtered datasets are available as hive tables.

| Table Name | Description | Num of Records |
|---|---|---|
| pagecountfilter100 | Pages that have received over 99 hits in an hour | 22 Million |
| pagecountfilter1000 | Pages that have received over 999 hits in an hour | 9 Million |
| pagecountfilter10000 | Pages that have received over 9999 hits in an hour | 1 Million |

Different size datasets have been made available to facilitate experimentation using variably sized SAP HANA & HDP AWS server instances.  Note the demonstration leverages modest sized servers for SAP HANA and HDP.  As such the usage of table pagecountfilter10000 is recommended.

All the tables share a common schema, as noted below:

| Column Name | Column Type |
|---|---|
| project code | String |
| Pagename | String |
| Hitsperhour | Bigint |
| Bytesperhour | Bigint |
| Year | String |
| Month | String |
| Day | String |
| Hour | String |

## 6.3 Install Hortonworks Hive ODBC Driver on SAP HANA AWS Instance

Now that we have successfully launched the SAP HANA Developer and HDP instances on AWS, we will next install the Hortonworks Hive ODBC driver on SAP HANA instance, so one can access data in Hortonworks Data Platform from within SAP HANA.  To install the ODBC driver, perform the following steps:

### Setup

➢ Log in SAP HANA Developer Edition as root
➢ In the shell window, run the following commands:

mkdir /tmp/HDP

cd /tmp/HDP

wget http://public-repo-1.hortonworks.com/HDP/hive-odbc/1.2.13.1018/suse11/hive-odbc-native-1.2.13.1018.tar.gz

tar -xvzf hive-odbc-native-1.2.13.1018.tar.gz

rpm -i **hive-odbc-native-1.2.13.1018**/hive-odbc-native-1.2.13.1018-1.x86_64.rpm

create file /usr/sap/HDB/home/.hortonworks.hiveodbc.ini with following content

```
[Driver]
DriverManagerEncoding=UTF-16
ErrorMessagesPath=/usr/lib/hive/lib/native/hiveodbc
/ErrorMessages/
LogLevel=0
LogPath=
ODBCInstLib=libodbcinst.so
```

create the file /usr/sap/HDB/home/.odbc.ini and copy the contents below.
Note that the HOST & HS2HostFQDN (identified in red) should be updated

```
[ODBC]
# Specify any global ODBC configuration here such as ODBC tracing.

[ODBC Data Sources]
HDP=Hortonworks Hive ODBC Driver 64-bit
[HDP]

Description=Hortonworks Hive ODBC Driver (64-bit) DSN
Driver=/usr/lib/hive/lib/native/Linux–amd64–
64/libhortonworkshiveodbc64.so
HOST= ec2-54-225-6-165.compute-1.amazonaws.com
HS2HostFQDN= ec2-54-225-6-165.compute-1.amazonaws.com
PORT=10000
Schema=default
FastSQLPrepare=0
UseNativeQuery=0
HiveServerType=2
HS2AuthMech=2
UserName=hive
```

to the public DNS of Hortonworks Data Platform launched in step 6.2


chmod 777 /usr/sap/HDB/home/.hortonworks.hiveodbc.ini

chmod 777 /usr/sap/HDB/home/.odbc.ini

export
LD_LIBRARY_PATH=$LD_LIBRARY_PATH:/usr/lib/hive/lib/native/Linux-amd64-64

echo export
LD_LIBRARY_PATH=\$LD_LIBRARY_PATH:/usr/lib/hive/lib/native/Linux-amd64-64 >> /usr/sap/HDB/home/.sapenv.sh

echo export

LD_LIBRARY_PATH=\$LD_LIBRARY_PATH:/usr/lib/hive/lib/native/Linux-

amd64-64 >> /usr/sap/HDB/home/.sapenv.csh

export ODBCINI=/usr/sap/HDB/home/.odbc.ini

echo export ODBCINI=/usr/sap/HDB/home/.odbc.ini >>

/usr/sap/HDB/home/.sapenv.sh

echo export ODBCINI=/usr/sap/HDB/home/.odbc.ini >>

/usr/sap/HDB/home/.sapenv.csh

su - hdbadm

HDB stop

HDB start

## Test the Setup

Run the following command to validate the setup
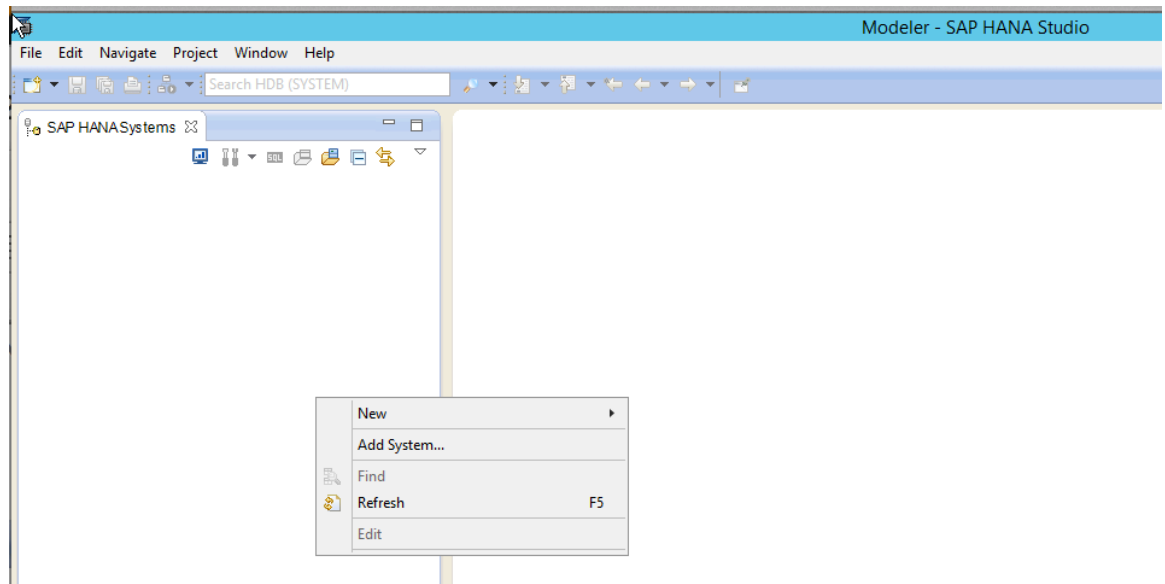
isql –v HDP

isql should successfully connect and show the sql prompt.  In the event a connection is not established, please review to make sure all the shell commands above were executed successfully.

# 6.4 Access HDP from SAP HANA using Smart Data Access

As of this point, we have a working SAP HANA instance with Hortonworks Hive ODBC driver installed and a working Hortonworks Data Platform with sample data. We will now use SAP HANA Studio to establish a connection between HANA and HDP.

## Establish connection between SAP Studio & SAP HANA on AWS

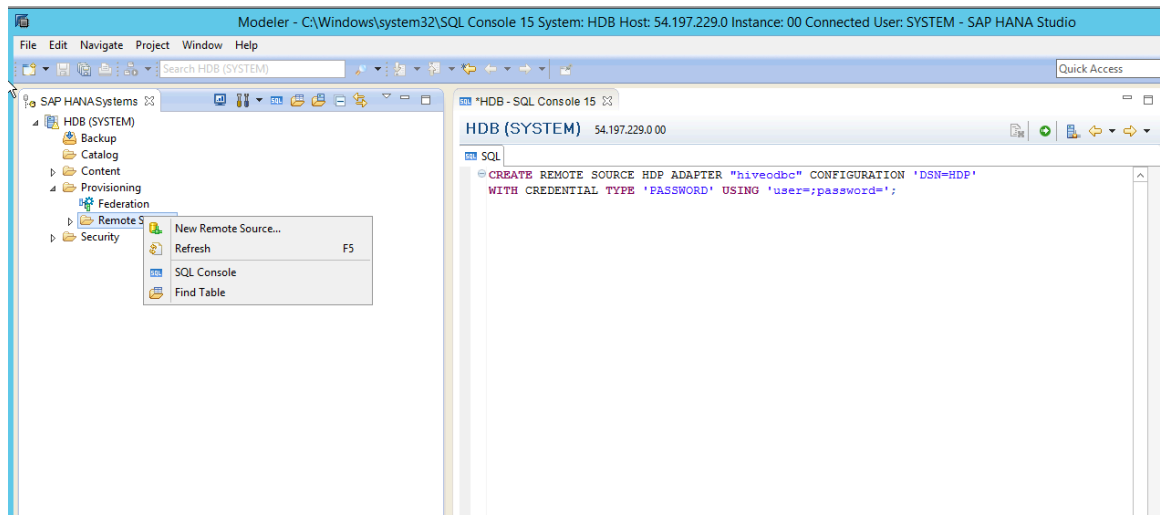➢ Log into SAP HANA Studio and open the Modeler Perspective
➢ Add a new system

> In the resulting wizard, provide the following information:

- Host Name: Public IP of your SAP HANA on AWS

- Instance Number: 00

- User Name: SYSTEM

- Password: manager

Accept the default value for all the other fields and click Finish.

## Create Virtual Table to access HDP from SAP HANA

> Navigate to Provisioning -> Remote Source and right click to select SQL Console

> In the SQL console, enter the following statement to create a remote data source

> CREATE REMOTE SOURCE HDP ADAPTER "hiveodbc"
> CONFIGURATION 'DSN=HDP' WITH CREDENTIAL TYPE
> 'PASSWORD' USING 'user=;password=';

> Create virtual table for HDP table pagehitfilter10000 using the following SQL command:

> create virtual table "pagecountfilter10000"
> ("projectcode" varchar(50),
>  "pagename" varchar(2000),
>  "hitsperhour" BIGINT,
>  "bytesperhour" BIGINT,
>  "year" VARCHAR(4),
>  "month" VARCHAR(2),
>  "day" VARCHAR(2),
>  "hour" VARCHAR(2)
>  ) AT "HDP"."default"."default"."pagecountfilter10000";

➢ Review the data in the virtual table by selecting the table and opening it for Data preview



## Transfer data to SAP HANA for Instant Access and Interactive Analysis

➢ Create a column table using the following SQL

```
create column table "pagehitfacttable"
("pagename" varchar(2000),
 "year" varchar(4),
 "month" varchar(2),
 "day" varchar(2),
 "hour" varchar(2),
 "pagehitperhour" BIGINT,
```
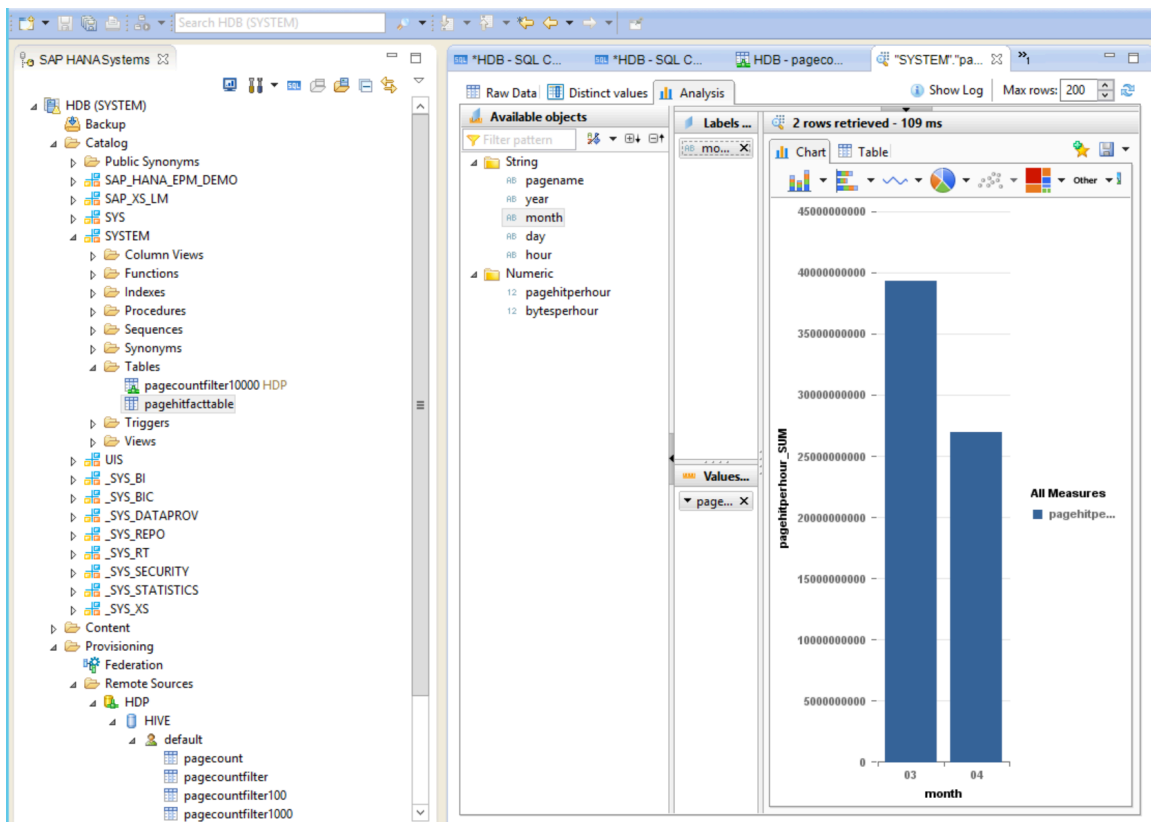
```
                    "bytesperhour" BIGINT);
```

➢ Transfer data from virtual table pagecountfilter10000 to column table
   pagehitfacttable using the SQL

```
select
"pagename","year","month","day","hour","hitsperhour","bytesper
hour"
from "SYSTEM"."pagecountfilter10000"
where "hitsperhour" > 100000
into "SYSTEM"."pagehitfacttable";
```

## Analyze the Column Table in SAP HANA

Select the table pagehitfacttable and open it in Data Preview mode for further
analysis.

# 7 Summary

The demonstration illustrates the unique capability of SAP HANA and Hortonworks Data Platform to deliver instant access at infinite scale. Organizations can leverage the joint capabilities of HANA and HDP to meet the needs of ever growing data volumes and a business imperative to act in real-time.

# 8 References

A number of blog posts were referenced to create the demonstration scenario. Of special note is the blog post by Bill Ramos on Importing Wikipedia Hive data into SAP HANA One. The original work by Bill was extended to include a live connection between Hortonworks Data Platform and SAP HANA.

# 9 Questions

If you have any question please reach us at sap@hortonworks.com