

BLUE SKY ONLINE RETAILER - CASE STUDY

BLUE SKY ONLINE RETAILER -CASE STUDY

INTRODUCTION:

According to the given case study, Blue Sky Online is a consumer electronics retailer that offers more than five hundred brands. The company conducts online business across Europe and the UK. They have multiple selling channels like Amazon, eBay, Tesco, and their website. Blue Sky has a huge customer base. Due to the competitive nature of the electronics retail industry, businesses in this industry are facing multiple challenges like pricing, targeting of customers, customer satisfaction, and so on. Blue Sky Company is facing an issue with accessing the customer base as some valuable information is missing due to minimal maintenance. Therefore, the company requires the analyst to analyze and manage the customer data to improve the business and relationship with its customers, to promote customer retention, and to improve their sales. The following business requirements have been identified:

1. Maintain customer database.
2. Analyze purchase or transaction history based on demographic data such as age group, income group, marital status, and postcode/City/Country.
3. Identify customer preferences in terms of selling channels/ payment payments/.
4. Identify the best and most profitable customers at different periods (month, quarter, and year).

According to the brief, this project aims to:

- 1.1 Identify the grain detail of the fact table and explain in brief the hierarchies.
- 1.2 Identify the dimensions of the central fact table.
- 1.3 Populate the fact table and dimension table.
- 1.4 Create three graphical representations of simple star schemas.

1. Analysis and design of dimensional data model:

In this analysis and design phase, we are required to understand the given data and produce simple star schemas.

1.1 Identify and define the grain of a central fact table, to represent the primary business activity. Explain the hierarchies if there are any.

To effectively design a dimensional data model for Blue Sky's Online data warehouse, the team has collaboratively identified and defined critical elements relating to the central fact table(s) and hierarchies.

Central Fact Table(s) / Grain Detail:

We first identified that we needed three fact tables to meet the three defined business requirements of Blue Sky Online. Each fact table would need to have its granularity – the level of detail in which the data is stored. The granularity decides the grain level and a grain refers to the specific unit of analysis for a fact table within the dimensional data model.

We first looked at the second business requirement. To meet this, we decided to focus on customer transactions, specifically individual customer sales transactions named. We named this the 'CustomerSalesTransactions' fact. This central fact table captures the most granular details at the atomic level, allowing for subsequent aggregation or grouping.

The agreed-upon grain is set to the most granular level, each row captures an individual transaction made by a customer. This ensures a comprehensive representation of the primary business activity.

For the third business requirement, we identified the need for a fact table that captures customer preferences for selling channels and payment – we called this the 'CustomerPreferences' fact. Again, the grain level was set on the customer level. Each row represents a customer transaction, including information on the selling channel and payment type. The individual rows can then be aggregated by grouping and counting to identify the most popular selling channels and payment methods.

Lastly, for the fourth business requirement, we identified the usefulness of a period snapshot fact table that can capture the profit per customer aggregated over a set period. This will allow for quicker and easier management reports instead of having to write repetitive queries every (by month, quarter, year). We call this the 'CustomerProfitability' fact.

For this table, the grain level is the period (by customer), and each row represents a periodic snapshot of the profitability per customer. For the requirement, we could pick 'Month' as the grain. However, it could also be useful to consider 'daily' as the lowest grain, should business requirements change, and it becomes necessary to find the most profitable customer on a given day or week. It is then possible to aggregate the report into quarters and years if needed.

Hierarchies:

Hierarchies are of paramount importance in the context of dimensional data modeling. They play a crucial role in structuring and organizing data for efficient querying and reporting within data warehousing systems. These hierarchies serve as the structural framework that represents

interrelationships and pathways to dive deep into various attributes and data levels. They empower users to effectively explore and analyze data, enhancing user understanding and overall usability and utilization of the data model.

There are a few hierarchies that we should be aware of when deciding the grain of the fact table. Dates are a fixed depth hierarchy – a series of many-to-one relationships – with agreedupon hierarchy level names, such as Year -> Quarter -> Month -> Day. The level we chose needs to fit the grain of the fact(s). in the case of ‘CustomerSalesTransactions’ and ‘CustomerPreferences’ the grain is each customer purchase, and this is captured on the daily hierarchy. We can later use the different hierarchies to ‘drill up’ to higher levels. Note that if a customer submits more than one order per day, it might be useful to consider including more identifiable attributes such as invoice number or timestamp (the data source does not include a timestamp).

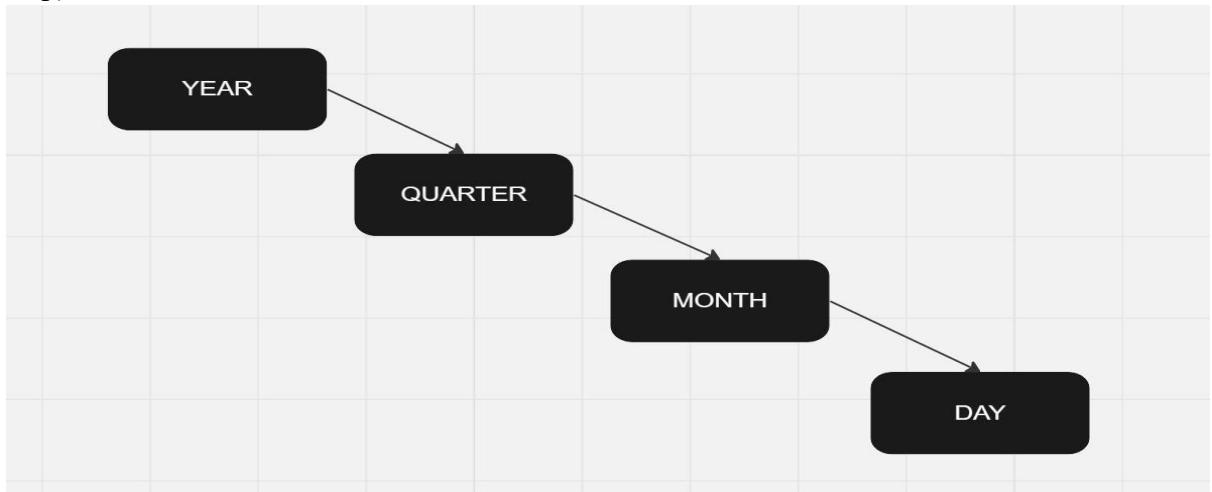


Figure 1: Time Hierarchy

The source data also includes information on the geographic location of the customer. Geographic hierarchies are considered *slighter ragged* hierarchies because they have variable depth. Our data source only includes geographic information as part of customer data, and the hierarchies are city ->postcode. As Blue Sky Online’s transactions happen online, the customer’s geographic location will not be a part of the fact table but instead used as a description of the customer for filtering and grouping.

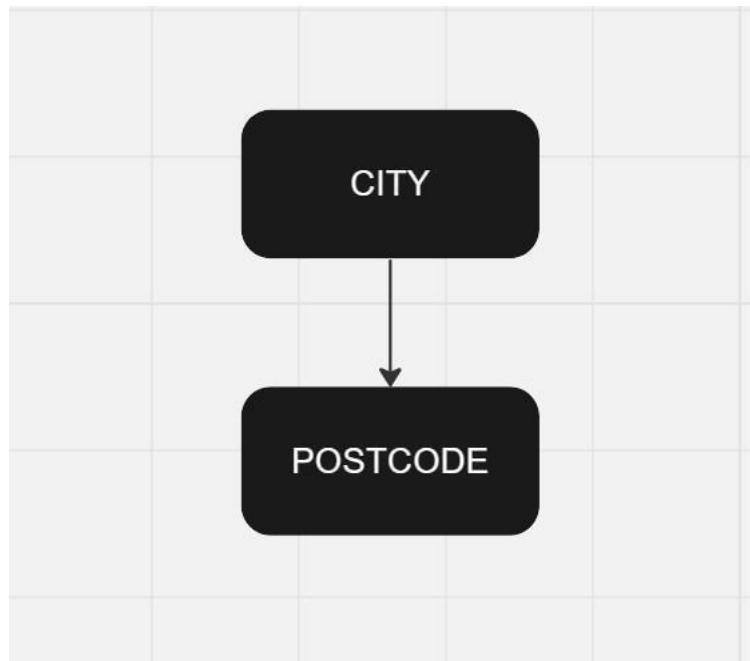


Figure 2: Location Hierarchy

1.2 Identify dimensions of central fact table(s) based on the available data source that will meet the reporting and analysis requirements mentioned in the business case scenario.

Before identifying the dimensions of our fact table, it can be useful to understand the given data sources. They are:

1. **CustomerDetails-1:** A CSV file containing the customer's full name, customer codes, date of birth, marital status, gender, postcode, and city.
2. **CustomerDetails-2:** A text file containing the customer's first name, last name, and income.
3. **SellingChannels:** A text file containing the channel name, channel code, and commission rate.
4. **GeneratedDateTime:** A text file containing 21 columns; Date Key, Full Date, Date Name, Day of Week, Day of Month, Day of Year, WeekdayWeekend, Week of Year, Month Name, Month of Year, Last Day of Month, Calendar Quarter, Calendar Year, Calendar Year Month, Calendar Year Quarter, Fiscal Month of Year, Fiscal Quarter, Fiscal Year, Fiscal Year Mont, Fiscal Year Quarter.
5. **PaymentDate:** A CSV file containing the payment type name, and Retailer Payment Type ID.
6. **CustomerSaleTransactions:** A CSV file containing the Date, Customer Code, Invoice Number, Total Cost, Toral Retail Price, Payment Type ID, and Selling Channel.

The dimensions of the fact table were identified by asking ‘when,’ ‘who,’ ‘where,’ ‘what,’ ‘how,’ and ‘why’ related to the event. They represent the descriptions of the attributes for filtering and grouping the facts.

For the first fact table – CustomerTransactions – the event is capturing transactions daily (when), by customers (who). The ‘what’ – product – is not provided in the source information.

Based on the above, the following dimension can be identified based on the available data source and requirements:

- Date dimension
- Customer dimension

The second fact table – CustomerPreferences – also captures events daily (when), by customers (who), on online sales channels (where), and by payment methods (how).

- Date Dimension
- Customer Dimension
- Selling channel Dimension
- Payment Dimension

For the third fact table – CustomerProfitability – the answer to the above question is shorter; The event is captured monthly (when) and by (who). This gives the following dimensions:

- Date dimension
- Customer dimension

1.3 Identify attributes of dimensions and measures of the fact table(s) based on the available data source descriptions that will meet the reporting and analysis requirements mentioned in the business case scenario.

To meet the business requirements, it is necessary to populate the dimension models with attributes that can be used for filtering and grouping. We also need to add numeric values as measurements to our fact table.

Date dimension:

Despite the source data containing more than twenty columns/attributes for each date, it is not necessary to add all of these to meet Blue Sky Online’s business requirements. The minimum number of attributes for this dimension are:

- FullDate (PK)

- Month of Year
- Calendar Quarter
- Calendar Year

Though this is the minimum number of attributes to meet the business requirements, it can be useful to consider adding more when building the data warehouse, in case the business requirements change.

Customer Dimension:

The customer dimension has combined attributes from two sources – CustomerDetails-1 and CustomerDetails-2 – to create a comprehensive customer dimension. It includes the following:

- Customer Key (PK)
- Customer Code
- First Name
- Last Name
- Birth Date
- Marital Status
- Gender
- Post Code
- City
- Income

Note that commission is a numeric value. However, it is a ratio rather than a fully additive numeric value. It has been included in the Selling Channel dimension because it is used for filtering and grouping.

Selling channel dimension:

- Selling Channel Key (PK)
- Selling Channel Name
- Selling Channel Code
- Commission

Note: Commission is a numeric value. However, it is a ratio rather than a fully additive numeric value. It has been included in the Selling Channel dimension because, in addition to being a measurement (see fact measurements below), it can also be used for filtering and grouping.

Payment Type Dimension:

- Payment Type Key (PK)
- Payment Type Name
- Payment Type ID

1.4 Use simple star schema to define the structure that shows how the fact tables are related to their dimensions. Create a graphical representation of the three simple star schemas on one page.

After careful consideration, we decided on a structure that combines the three business requirements into one fact table. The justifications for the design of the simple star schema below are:

First, as both ‘CustomerSalesTransactions’ and ‘CustomerPreferences’ have the same granularity – customer transaction level – it is useful to combine these two fact tables into one single fact table. The fact table will be connected to all four dimensions as we will need date, customer details, selling channel details, and payment type details to meet business requirements 2 and 3. Business requirement two needs transaction information, date, and details from the customer dimension to analyze purchase history by customer demographics. To meet the third business requirement, we will need information on the selling channel, payment method, customer, and date. We would then have to aggregate the individual transactions to identify the most used selling channel or payment type.

We were then left with two fact tables, ‘CustomerSalesTransactions’ and ‘CustomerProfitability.’ The initial thought was to design two simple star schemas based on these, as the preferred granularity would be on transaction level for the former, and period and customer for the latter. The idea was the table would show the total profit per customer for that given period. We could then meet the fourth business requirement by filtering the customers by profit to get the most profitable customer.

However, this design produced a challenge. Profit is calculated by taking the ‘Total Retail Price’ and subtracting ‘Total Costs’ and ‘Commission.’ Both retail prices and costs can easily be aggregated across time and customer. Commission, on the other hand, is a ratio, so it cannot be aggregated in the same way across dimensions. In addition, the commission is based on what type of selling channel has been used for a specific transaction. So, unless the customer has used the same selling channel across the period it will also be harder to calculate the exact commission rate for the month without calculating this on the transaction level.

To get around this issue most simply, we decided that it would be better to combine all three fact tables into one. The grain would be the sales transaction level for all; each row corresponds to one customer sales transaction. We can then calculate the profitability for each transaction, and then aggregate this by customer and time to get the most profitable customer in a month, quarter, or year. This design requires a little more aggregation than separating profitability as an individual fact table, but it allows us to meet business requirement four without having to make more complicated commission calculations.

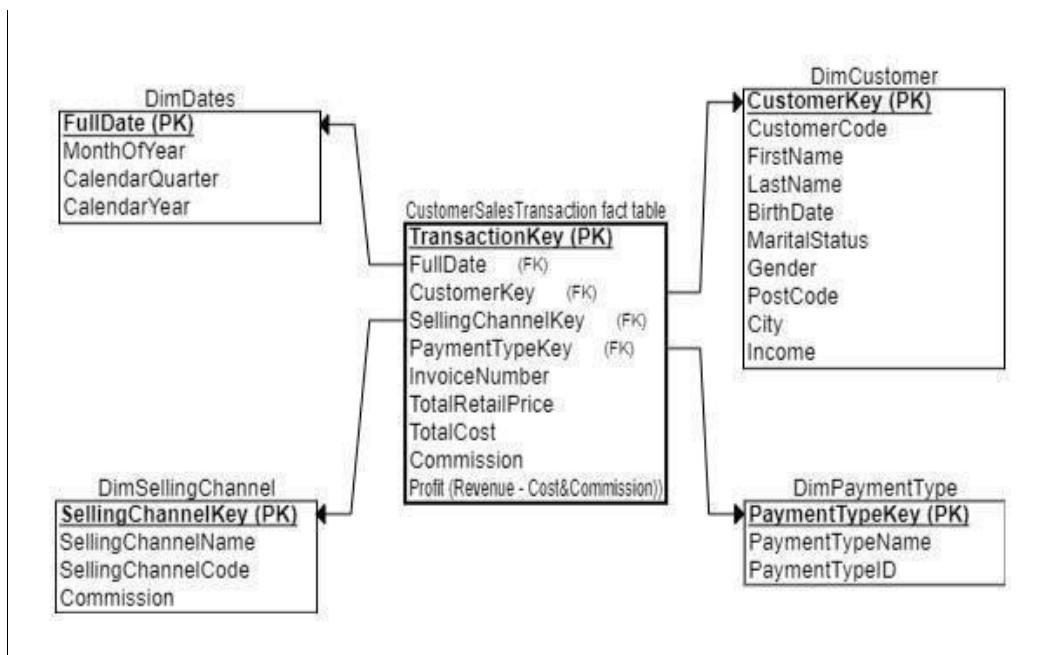


Figure 3: Simple Star Schema

We have also added the following SQL queries that show that our design can meet the business requirements.

Business Requirement: Identify customer preferences in terms of selling channels/ payment payments.

```
SELECT pt.PaymentTypeName as most_popular_payment_type, sc.
```

```
SellingChannelName as most_popular_selling_channel,
```

```
COUNT (*) as frequency
```

```
FROM CustomerSalesTransaction cst
```

```
JOIN DimPaymentType pt ON cst.PaymentTypeKey = pt.PaymentTypeKey
```

```
JOIN DimSellingChannel sc ON cst. SellingChannelKey= sc. SellingChannelKey  
GROUP BY pt. PaymentTypeName, sc. SellingChannelName  
ORDER BY frequency DESC.
```

This query counts the frequency of each combination of payment type and selling channel together. To find the most popular payment type and selling channel individually, the query needs to be two separate subqueries for each dimension.

Business Requirement: Identify the best and most profitable customers at different periods (month, quarter, and year).

```
SELECT dc. LastName, SUM  
(cst. Profit) as TotalProfit  
  
FROM CustomerSalesTransaction  
  
JOIN DimCustomer dc ON cst. CustomerKey = dc. CustomerKey  
  
JOIN DimFullDate dd ON cst. FullDate = dd. FullDate  
  
WHERE EXTRACT (month FROM dd. FullDate) = 1  
  
GROUP BY DC. LastName  
  
ORDER BY TotalProfit DESC.
```

This query will return based on the customer's last name. We might want to consider either joining first name and last name together to a name attribute or writing a more complex query to make sure that the results show exact customers in cases where there are customers with the same last name.

Business Requirement: Analyse purchase or transaction history based on demographic data such as age group, income group, marital status, and postcode/City/Country.

One example could be to find the TotalRetailPrice for customers who are married.

```
SELECT SUM (cst. TotalRetailPrice) as totalretailsales_married.  
  
FROM CustomerSalestTransactions cst  
  
JOIN DimCustomer dc ON cst. CustomerKey = dc. CustomerKey  
  
WHERE dc. MaritalStatus = 'Married.'
```

From the above Conceptual Design, a physical design needs to be implemented as individual work.

2. Create a Relational Database to store dimensions and fact tables using Microsoft SQL Server Management Studio.

To create a relational database for Blue Sky Online, the database must first be created in Microsoft SQL Server Management Studio (SSMS). SSMS offers a graphical interface and a SQL query for performing this task.

2.1 Create a database.

Before starting to create a database there few rules that need to be followed, SQL Server can create up to 32,767 databases. It is necessary to execute the **CREATE DATABASE** command in auto-commit mode rather than in explicit or implicit transactions. The master database should be backed up when we create, modify, or drop the database.

To create a database Open Object Explorer, connect to a SQL Server Database Engine instance, and then expand it.

- To Select New Database by right-clicking on **Databases**.
- Give a database name as “**BlueSkyOnlineDB**” in the New Database. • Click **OK** to create the database using all of the default options.

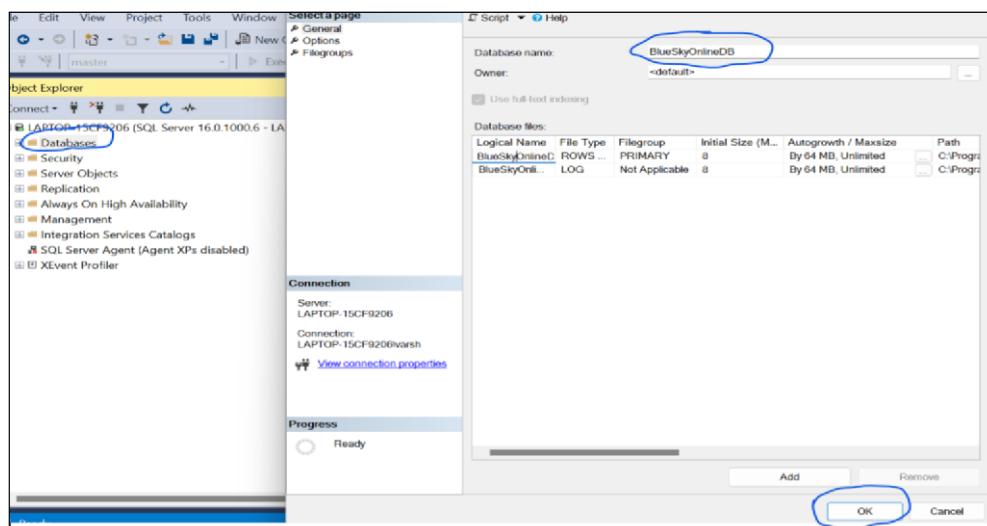


Fig: Creating Database “BlueSkyOnlineDB”

In a similar way to create a staging area, follow the previously mentioned steps for “BlueSkyOnlineDB” and name as “**BS_Staging**”

It acts as a temporary storage place for raw data extraction from source systems, and a staging area is an essential intermediate stage in the ETL process. It facilitates data cleansing, transformation, and validation, ensuring accuracy and consistency. To prevent direct impacts during the data extraction and transformation phases, the staging area separates the extraction process from the data warehouse. Additionally, before the data is loaded into the final destination, it provides verification and data validation, making it easier to find inconsistencies or abnormalities.

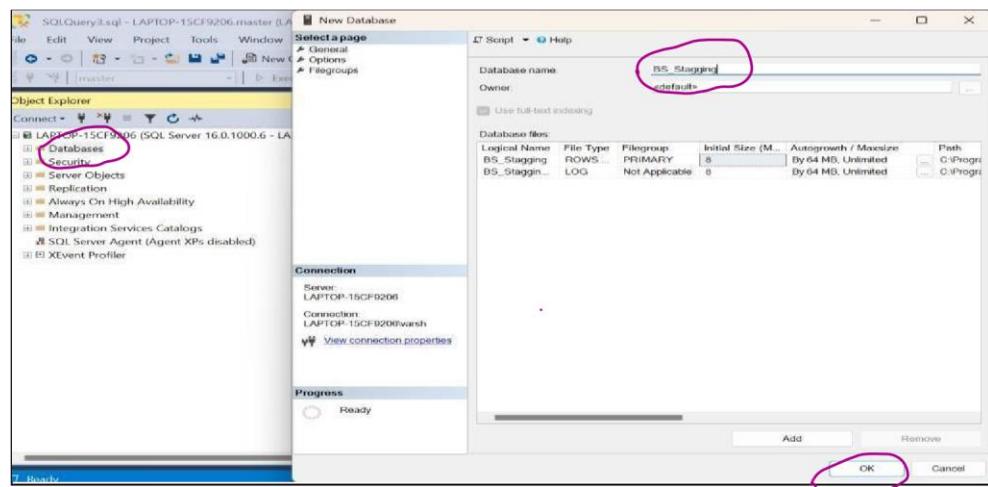


Fig: Creating staging area “BS_Staging”

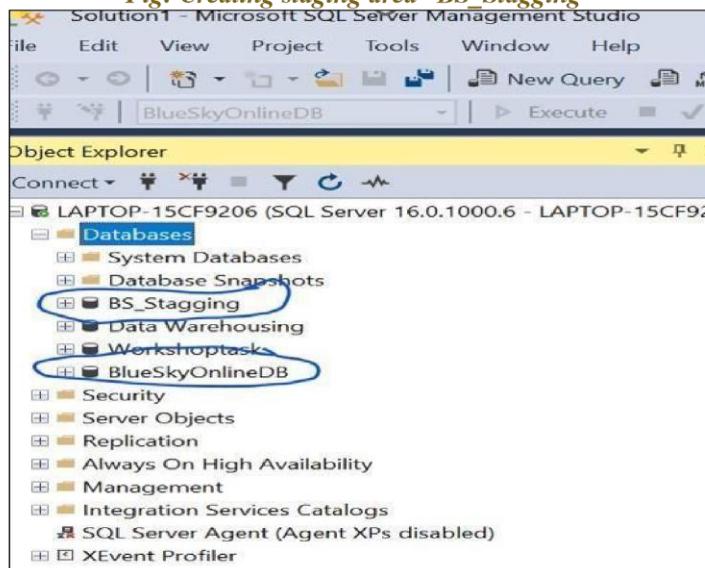


Fig: Databases “BlueSkyOnlineDB” and “BS_Staging”

2.2 Create dimension tables.

To create a new table in a given database, SQL provides the create TABLE statement. The structure of a table must be specified in an SQL query before one can be created. The structure consists of the table name, the column names, and the data type for each column in the table.

Note: We have to make sure each table in a database needs a unique name.

After the database has been created, the dimension tables must be designed and implemented. These dimensional tables represent the main entities like Customers, Dates, Payments, and Selling Channels. Each dimension table has been created to include specific attributes related to its entity. Primary keys are assigned for uniquely identifying records inside each dimension table.

Creating the Customer Dimension table in SSMS. These are the customer dimensions that have been created by combining attributes from two sources.

Customer Key: It is the primary key for uniquely identifying for each customer.

Customer Details: FirstName, LastName, Gender and BirthDate.

Demographic Data: Age Group, Income, and Marital Status.

Location Details: Postcode and City

These are the steps for creating the Customer Dimension that have been followed.

- Connect to the database server by launching SQL Server Management Studio.
- Click on “New Query” in the toolbar to open a new query window.
- Make sure we are using the right database when creating the Customer Dimension “**BlueSkyOnline**.”
- To run the query, use the “Execute” button in the toolbar.

The screenshot shows the SQL Server Management Studio interface. The Object Explorer on the left shows the database structure, including the 'BlueSkyOnlineDB' database and its tables. The 'CustomerDimension' table is selected and highlighted with a red oval. The central pane displays the SQL code for creating the table:

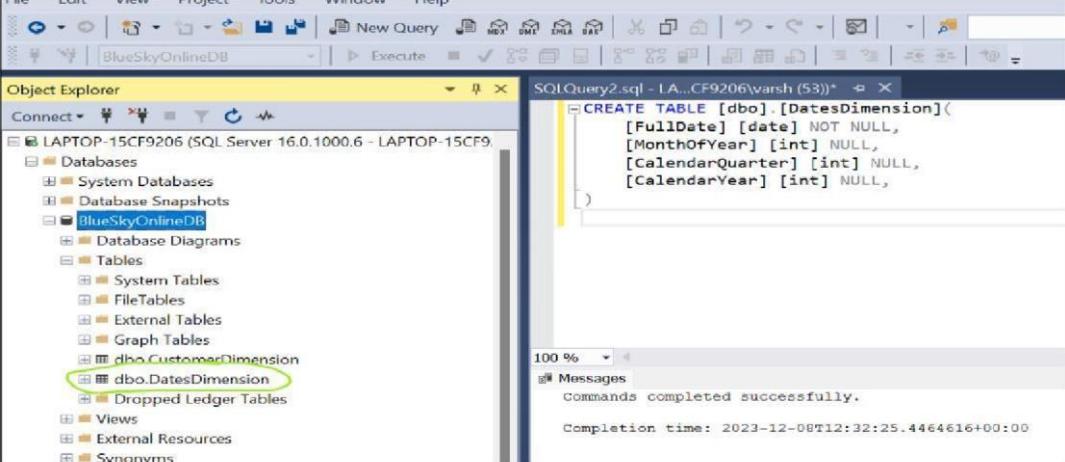
```
CREATE TABLE [dbo].[CustomerDimension](
    [CustomerKey] [int] IDENTITY(1,1) NOT NULL,
    [CustomerCode] [nvarchar](25) NOT NULL,
    [FirstName] [varchar](50) NOT NULL,
    [LastName] [varchar](50) NOT NULL,
    [BirthDate] [date] NOT NULL,
    [MaritalStatus] [varchar](20) NOT NULL,
    [Gender] [varchar](10) NOT NULL,
    [PostCode] [varchar](20) NULL,
    [City] [varchar](50) NULL,
    [Income] [int] NULL
)
```

The status bar at the bottom indicates "Commands completed successfully." and the completion time: 2023-12-08T13:24:59.314Z.

Fig: Creating CustomerDimension Table

To create a customer dimension table, use this SQL Query.

Follow the above steps for the other dimension tables and the fact table. These lines of code have been used for creating a table with date-related attributes. The query creates a “**DateDimension**” table with FullDate, MonthOfYear, CalendarQuarter, and CalendarYear.



The screenshot shows the SQL Server Management Studio interface. The Object Explorer on the left lists the database structure, including the BlueSkyOnlineDB database with its tables. The 'dbo.DatesDimension' table is highlighted with a green oval. The central pane displays the T-SQL code for creating the table:

```

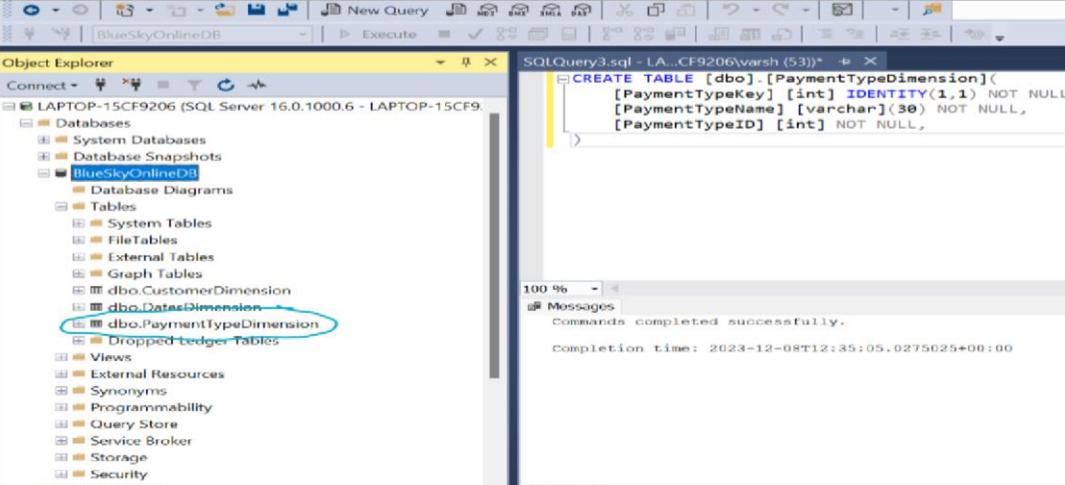
CREATE TABLE [dbo].[DatesDimension](
    [FullDate] [date] NOT NULL,
    [MonthOfYear] [int] NULL,
    [CalendarQuarter] [int] NULL,
    [CalendarYear] [int] NULL,
)

```

The status bar at the bottom right indicates the command completed successfully with a completion time of 2023-12-08T12:32:25.4464616+00:00.

Fig: Creating DatesDimension Table

To create the “**PaymentTypeDimension**” table in the database. These lines of SQL query have been used for payment-related attributes like PaymentTypeKey, PaymentTypeName, and PaymentTypeID.



The screenshot shows the SQL Server Management Studio interface. The Object Explorer on the left lists the database structure, including the BlueSkyOnlineDB database with its tables. The 'dbo.PaymentTypeDimension' table is highlighted with a blue oval. The central pane displays the T-SQL code for creating the table:

```

CREATE TABLE [dbo].[PaymentTypeDimension](
    [PaymentTypeKey] [int] IDENTITY(1,1) NOT NULL,
    [PaymentTypeName] [varchar](30) NOT NULL,
    [PaymentTypeID] [int] NOT NULL,
)

```

The status bar at the bottom right indicates the command completed successfully with a completion time of 2023-12-08T12:35:05.0275025+00:00.

Fig: Creating PaymentDimension Table

For creating specific a table for selling channels as “**SellingChannelDimension**” attributes like sellingchannelkey, sellingchannelname, sellingchannelcode, and commission in the database. This dimension will provide insights into the various channels through which Blue Sky Online conducts its business.

The screenshot shows the SQL Server Management Studio (SSMS) interface. In the Object Explorer on the left, under the 'Tables' node of the 'BlueSkyOnlineDB' database, the 'dbo.SellingChannelDimension' table is highlighted with a red oval. In the center pane, a query window titled 'SQLQuery3.sql' contains the following SQL code:

```

CREATE TABLE [dbo].[SellingChannelDimension](
    [SellingChannelKey] [int] IDENTITY(1,1) NOT NULL,
    [SellingChannelName] [varchar](30) NOT NULL,
    [SellingChannelCode] [varchar](10) NULL,
    [Commission] [int] NOT NULL,
)

```

The status bar at the bottom right of the query window shows the completion time as 'Completion time: 2023-12-08T12:40:23.2202246+00:00'.

Fig: Creating SellingChannelDimension Table

2.3 Create fact tables.

In the Simple star schema of a data warehouse, the fact table is the central table. It stores quantitative information for analysis and is often denormalized. Dimension tables are working with this fact table. In addition, it contains the data that needs to be analyzed and a dimension table stores data about how the data in the fact table can be analyzed. As a result, the fact table has two distinct types of columns. The foreign keys column helps connect with dimension tables, while the measures columns contain the data that is being analyzed.

To create a Fact table in SSMS to store key details about customer sales transactions. The details including Customer details, Invoice Number, purchase Total retail price, and related data will be stored in this database as “**CustomerSalesTransactionFactTable**.” Open a new SSMS query window and execute the following SQL code.

The screenshot shows the SQL Server Management Studio (SSMS) interface. In the Object Explorer on the left, under the 'Tables' node of the 'BlueSkyOnlineDB' database, the 'dbo.CustomerSalesTransactionFactTable' table is highlighted with a yellow oval. In the center pane, a query window titled 'SQLQuery6.sql' contains the following SQL code:

```

CREATE TABLE [dbo].[CustomerSalesTransactionFactTable](
    [TransactionKey] [int] IDENTITY(1,1) NOT NULL,
    [CustomerKey] [int] NOT NULL,
    [SellingChannelKey] [int] NOT NULL,
    [PaymentTypeKey] [int] NOT NULL,
    [FullDate] [date] NOT NULL,
    [InvoiceNumber] [varchar](20) NULL,
    [TotalRetailPrice] [decimal](15, 2) NULL,
    [TotalCost] [decimal](15, 2) NULL,
    [Commission] [int] NOT NULL,
    [Profit] [decimal](15, 2) NULL
)

```

The status bar at the bottom right of the query window shows the completion time as 'Completion time: 2023-12-08T13:01:31.8121643+00:00'.

Fig: Creating CustomerSalesTransactionFact Table

2.4 Add appropriate primary keys and foreign key constraints.

To Build Primary Keys: We can use the *ALTER TABLE* command to change a primary key in a dimensional table after it has been created. Suppose if we want to modify the primary key of the *CustomerDimension* table from the column *CustomerKey* as *PrimaryKey*.

In this SQL code, the *ALTER TABLE* statement is used to modify as primary key constraint (*PK_CustomerDimension*) from the *CustomerDimension* table along with other dimensional tables like *DateDimension*, *PaymentTypeDimension*, and *SellingChannelDimension*.

```
Object Explorer
Connect > BlueSkyOnlineDB > Execute > SQLQuery7.sql - LA...CF9206varsh (58)*
SQLQuery6.sql - LA...CF9206varsh (53)* > X

ALTER TABLE [dbo].[DatesDimension]
ADD CONSTRAINT PK_DatesDimension PRIMARY KEY (FullDate) ON [PRIMARY];

ALTER TABLE [dbo].[CustomerDimension]
ADD CONSTRAINT PK_CustomerDimension PRIMARY KEY (CustomerKey) ON [PRIMARY];

ALTER TABLE [dbo].[PaymentTypeDimension]
ADD CONSTRAINT PK_PaymentTypeDimension PRIMARY KEY (PaymentTypeKey) ON [PRIMARY];

ALTER TABLE [dbo].[SellingChannelDimension]
ADD CONSTRAINT PK_SellingChannelDimension PRIMARY KEY (SellingChannelKey) ON [PRIMARY];

ALTER TABLE [dbo].[CustomerSalesTransactionFactTable]
ADD CONSTRAINT PK_CustomerSalesTransactionFactTable PRIMARY KEY (TransactionKey) ON [PRIMARY];

Messages
Commands completed successfully.

Completion time: 2023-12-08T16:41:19.4754015+00:00
```

Fig: Altering CustomerDimension table to add PrimaryKey

By following the above code, the same as adding constraints for the customer dimension table other dimensional table constraints have been modified.

```

File Edit View Project Tools Window Help
Object Explorer
Connect SQL Server Object Explorer
Tables
System Tables
FileTables
External Tables
Graph Tables
dbo.CustomerDimension
dbo.CustomerSalesTransactionFactTable
dbo.DatesDimension
Columns
FullDate (PK, date, not null)
MonthYear (int, not null)
CalendarQuarter (int, not null)
CalendarYear (int, not null)
Keys
PK_DatesDimension
Constraints
Triggers
ALTER TABLE [dbo].[DatesDimension]
ADD CONSTRAINT PK_DatesDimension PRIMARY KEY (FullDate) ON [PRIMARY];
ALTER TABLE [dbo].[CustomerDimension]
ADD CONSTRAINT PK_CustomerDimension PRIMARY KEY (CustomerKey) ON [PRIMARY];
ALTER TABLE [dbo].[PaymentTypeDimension]
ADD CONSTRAINT PK_PaymentTypeDimension PRIMARY KEY (PaymentTypeKey) ON [PRIMARY];
ALTER TABLE [dbo].[SellingChannelDimension]
ADD CONSTRAINT PK_SellingChannelDimension PRIMARY KEY (SellingChannelKey) ON [PRIMARY];
ALTER TABLE [dbo].[CustomerSalesTransactionFactTable]
ADD CONSTRAINT PK_CustomerSalesTransactionFactTable PRIMARY KEY (TransactionKey) ON [PRIMARY];
100% - 4
Messages
Commands completed successfully.
Completion time: 2023-12-08T16:41:19.4754015+00:00

```

Fig: Altering DateDimension table to add PrimaryKey

```

File Edit View Project Tools Window Help
Object Explorer
Connect SQL Server Object Explorer
Tables
System Tables
FileTables
External Tables
Graph Tables
dbo.CustomerDimension
dbo.CustomerSalesTransactionFactTable
dbo.DatesDimension
dbo.PaymentTypeDimension
Columns
PaymentTypeKey (PK, int, not null)
PaymentTypeName (varchar(30), not null)
PaymentTypeID (int, not null)
Keys
PK_PaymentTypeDimension
Constraints
Triggers
ALTER TABLE [dbo].[DatesDimension]
ADD CONSTRAINT PK_DatesDimension PRIMARY KEY (FullDate) ON [PRIMARY];
ALTER TABLE [dbo].[CustomerDimension]
ADD CONSTRAINT PK_CustomerDimension PRIMARY KEY (CustomerKey) ON [PRIMARY];
ALTER TABLE [dbo].[PaymentTypeDimension]
ADD CONSTRAINT PK_PaymentTypeDimension PRIMARY KEY (PaymentTypeKey) ON [PRIMARY];
ALTER TABLE [dbo].[SellingChannelDimension]
ADD CONSTRAINT PK_SellingChannelDimension PRIMARY KEY (SellingChannelKey) ON [PRIMARY];
ALTER TABLE [dbo].[CustomerSalesTransactionFactTable]
ADD CONSTRAINT PK_CustomerSalesTransactionFactTable PRIMARY KEY (TransactionKey) ON [PRIMARY];
100% - 4
Messages
Commands completed successfully.
Completion time: 2023-12-08T16:41:19.4754015+00:00

```

Fig: Altering PaymentTypeDimension table to add PrimaryKey

```

File Edit View Project Tools Window Help
Object Explorer
Connect SQL Server Object Explorer
Tables
System Tables
FileTables
External Tables
Graph Tables
dbo.CustomerDimension
dbo.CustomerSalesTransactionFactTable
dbo.DatesDimension
dbo.PaymentTypeDimension
dbo.SellingChannelDimension
Columns
SellingChannelKey (PK, int, not null)
SellingChannelName (varchar(30), not null)
SellingChannelCode (varchar(10), null)
Commission (int, not null)
Keys
PK_SellingChannelDimension
Constraints
ALTER TABLE [dbo].[DatesDimension]
ADD CONSTRAINT PK_DatesDimension PRIMARY KEY (FullDate) ON [PRIMARY];
ALTER TABLE [dbo].[CustomerDimension]
ADD CONSTRAINT PK_CustomerDimension PRIMARY KEY (CustomerKey) ON [PRIMARY];
ALTER TABLE [dbo].[PaymentTypeDimension]
ADD CONSTRAINT PK_PaymentTypeDimension PRIMARY KEY (PaymentTypeKey) ON [PRIMARY];
ALTER TABLE [dbo].[SellingChannelDimension]
ADD CONSTRAINT PK_SellingChannelDimension PRIMARY KEY (SellingChannelKey) ON [PRIMARY];
ALTER TABLE [dbo].[CustomerSalesTransactionFactTable]
ADD CONSTRAINT PK_CustomerSalesTransactionFactTable PRIMARY KEY (TransactionKey) ON [PRIMARY];
100% - 4
Messages
Commands completed successfully.
Completion time: 2023-12-08T16:41:19.4754015+00:00

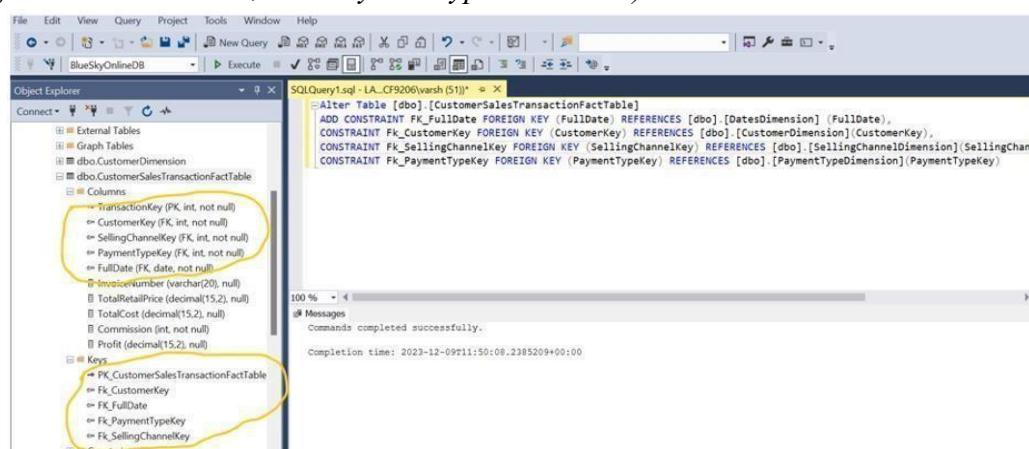
```

Fig: Altering SellingChannelDimension table to add PrimaryKey

To Build Foreign Keys: When we create foreign keys in the fact table, we are making connections to data stored in other dimension tables. Here, we have a *CustomerSalesTransactionFactTable*, which captures details about the customer transactions, such as which channel and when.

The below SQL code refers to the primary keys of relevant dimension tables through the use of foreign keys (*CustomerKey*, *FullDate*, *SellingChannelKey*, and *PaymentTypeKey*) in the *CustomerSalesTransactionFactTable*.

We have added the PRIMARY KEY constraint to the “**TransactionKey**” column to make sure that each transaction has a unique identification. Furthermore, each foreign key constraint refers to the associated primary key in the dimension tables (*CustomerDimension*, *DateDimension*, *SellingChannelDimension*, and *PaymentTypeDimension*).



```

File Edit View Query Project Tools Window Help
BlueSkyOnlineDB Execute
Object Explorer
Connect* External Tables Graph Tables dbo.CustomerDimension dbo.CustomerSalesTransactionFactTable Columns
    TransactionKey (PK, int, not null)
    > CustomerKey (FK, int, not null)
    > SellingChannelKey (FK, int, not null)
    > PaymentTypeKey (FK, int, not null)
    > FullDate (FK, date, not null)
    > InvoiceNumber (varchar(20), null)
    > TotalRetailPrice (decimal(15,2), null)
    > TotalCost (decimal(15,2), null)
    > Commission (int, not null)
    > Profit (decimal(15,2), null)
Keys
    PK_CustomerSalesTransactionFactTable
    > Fk_CustomerKey
    > Fk_FullDate
    > Fk_PaymentTypeKey
    > Fk_SellingChannelKey
SQLQuery1.sql - LA_C9206\varsh (5)* 0 X
ALTER Table [dbo].[CustomerSalesTransactionFactTable]
ADD CONSTRAINT FK_FullDate FOREIGN KEY (FullDate) REFERENCES [dbo].[DatesDimension] (FullDate),
CONSTRAINT Fk_CustomerKey FOREIGN KEY (CustomerKey) REFERENCES [dbo].[CustomerDimension](CustomerKey),
CONSTRAINT Fk_SellingChannelKey FOREIGN KEY (SellingChannelKey) REFERENCES [dbo].[SellingChannelDimension](SellingChannelKey),
CONSTRAINT Fk_PaymentTypeKey FOREIGN KEY (PaymentTypeKey) REFERENCES [dbo].[PaymentTypeDimension](PaymentTypeKey)
100 % 4 Messages Commands completed successfully.
Completion time: 2023-12-09 11:50:08.23805209+00:00

```

Fig: Altering CustomerSalesTransactionFactTable to add PrimaryKey and ForeignKey

3. Populate the implemented data warehouse based on the provided data sources using SQL Server Management Studio.

A sequence of steps must be followed before loading data into dimension tables to perform ETL processes in SQL Server Management Studio. A staging area has been created and data has been inserted in the staging area as “**BS_Staging**.”

3.1 Implement ETL (extract, transform, and load) processes to populate dimension tables.

Once the staging is created, the available data will be temporarily stored in the staging area in the database. Before inserting the data into the final dimension tables, this staging area serves as an intermediate stage where we can clean, transform, and validate the data. In these cases, we could manage missing values, standardize date formats, or carry out other data-cleaning tasks.

Importing customer details:

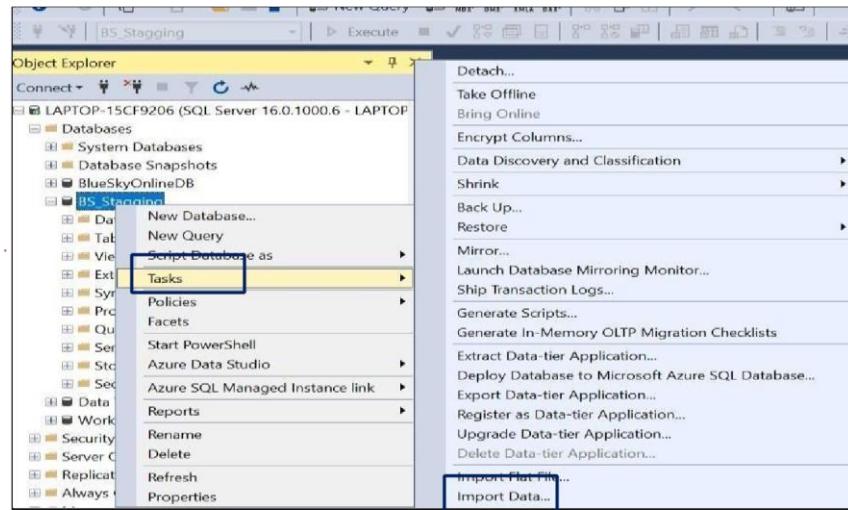


Fig: Importing customer details by selecting Task

- Select the Task to start the SQL Server import and export wizard, and then import.

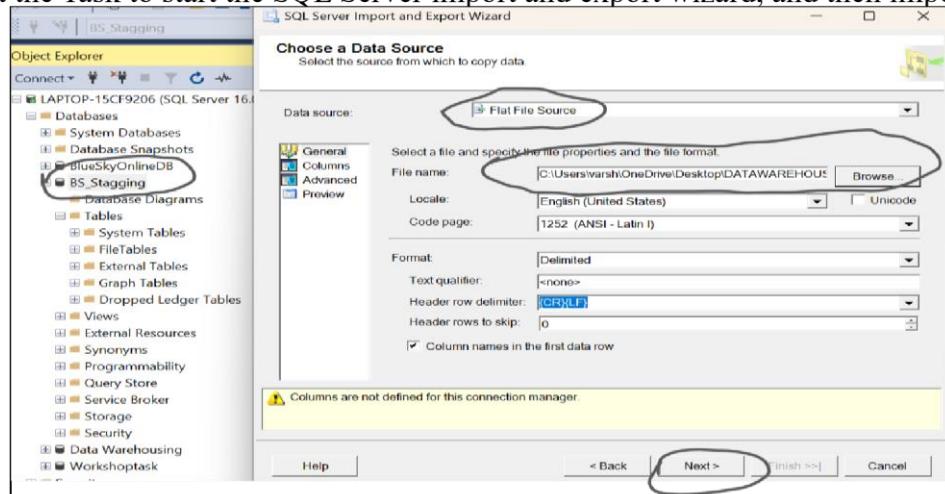


Fig: Importing data files into BS_Staging

- Select the customer details data source “Flat File Source.”
- Select a flat file source because the customer data that is generated has been stored in a CSV file with comma separated.

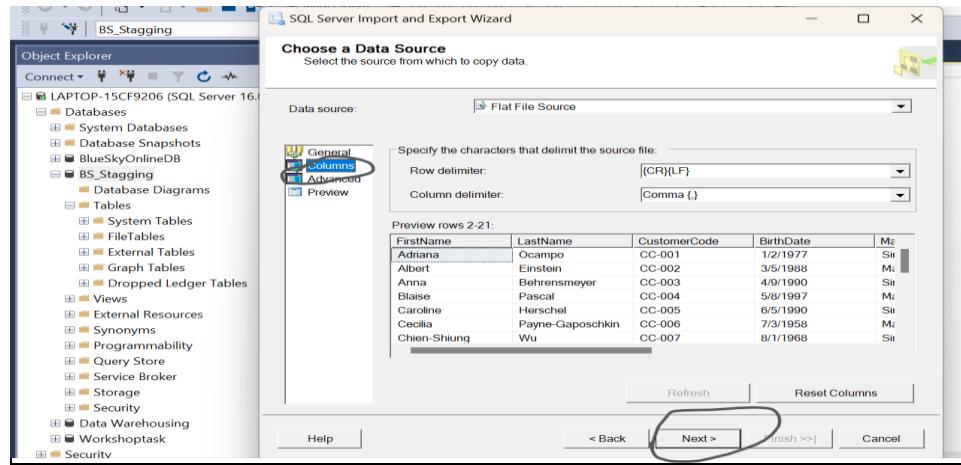


Fig: Displaying the columns and preview of the data source

- When we select a file, information about that file is displayed. We can also see the columns and preview the data source. It is suggested to preview the data source to quickly know the data type of each column.

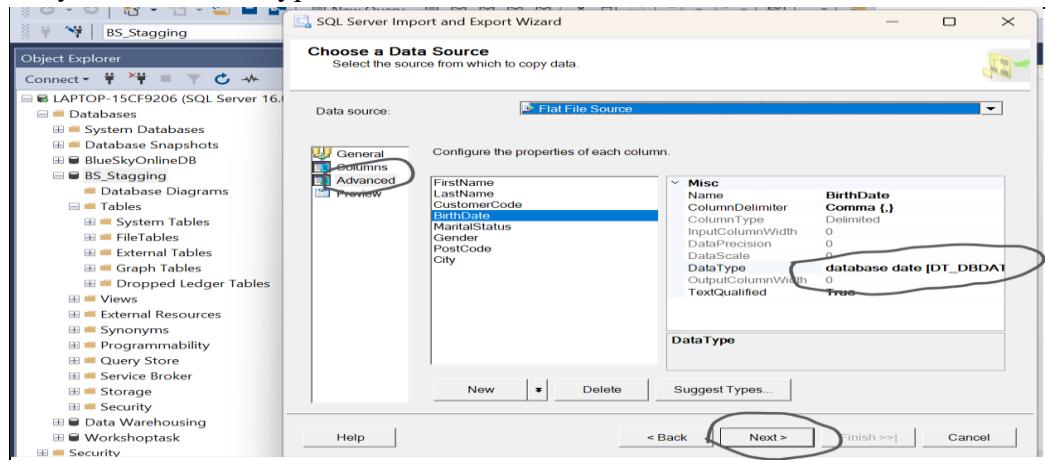


Fig: Showing an advanced option to change datatype

- To specify the data type and length, use the advanced feature. When importing or extracting data from a flat file source.

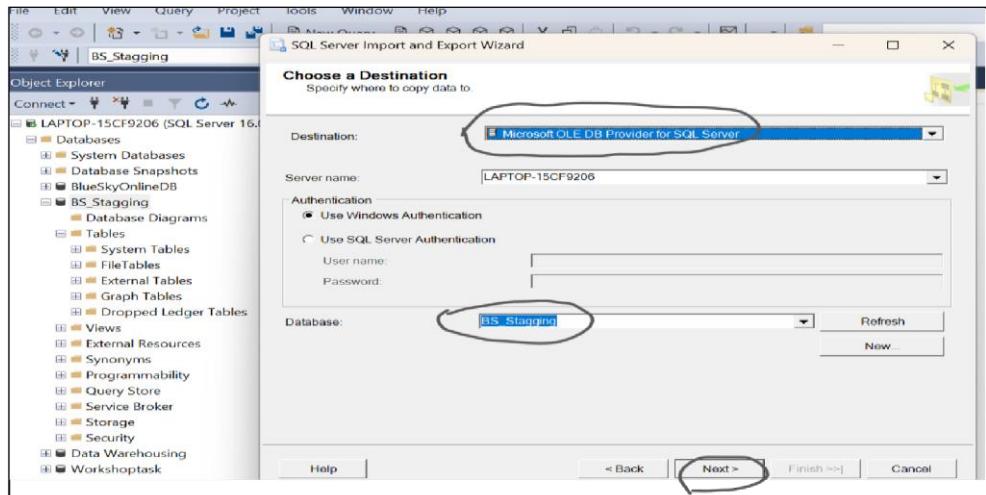


Fig: Choosing a destination for the ETL process

- The final destination is the Microsoft OLE DB provider for SQL Server for the Customer Details
- Select a destination for the ETL process. Although the purpose is to populate the customer details in Staging, the final destination is the data warehouse, which is managed by an SQL server.
- Customer Details 1 is the destination dimension table to be specified. If the ETL process is successful, a popup will appear indicating that the execution was successful and how many records were extracted and loaded to the Customer Details.

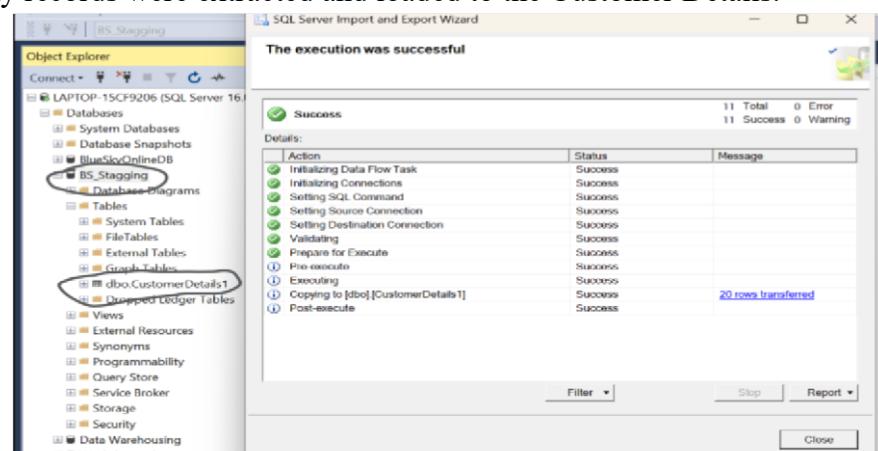


Fig: Execution window showing successful

Similarly for importing other details available in these case studies, to populate the other details available for Blue Sky online using flat file sources, repeat the procedure used for the previous customer details.

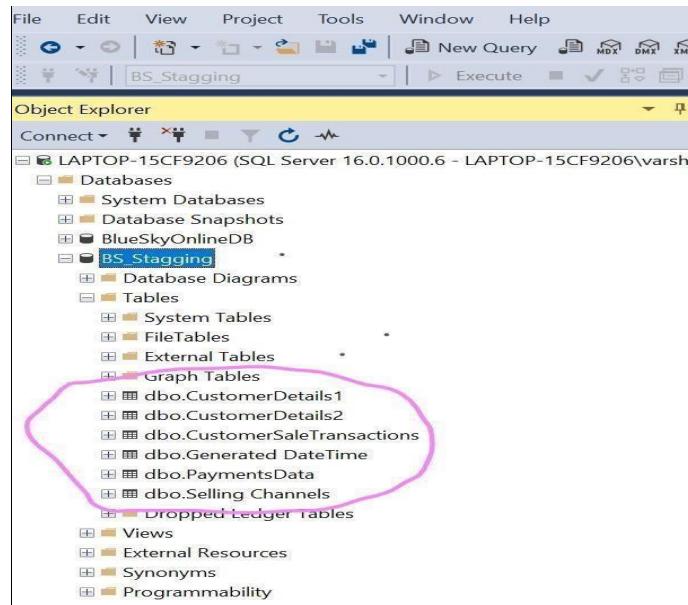


Fig: Imported data sources in BS_Staging

Once the data files are imported into the final dimension tables, it has been validated and loaded into the staging area. With this approach, we can ensure that thoroughly cleaned and structured data is utilized in data warehouses by controlling the quality of the data before inserting it into dimension tables. We may transfer data from the staging to the original dimension tables using `INSERT INTO...SELECT` query.

3.2 Implement ETL (extract, transform, and load) processes to populate a fact table.

The staging area is a temporary storage space, where raw data enters before moving it to the dimension tables. The required data from this staging area is subsequently extracted and transformed into a format that is suitable for distinct dimension tables.

We have customer data in the staging area, so we extract it, apply any necessary transformations (such as formatting addresses or providing unique customer IDs), and then place it into the “**CustomerDimension**” table.

```

SQLQuery2.sql - LA..CF9206\varsh [53]* < X SQLQuery4.sql - LA..CF9206\varsh [70]
--Insert into [dbo].[DatesDimension] select FullDate,MonthOfYear,CalendarQuarter,CalendarYear
From [BS_Staging].[dbo].[Generated Datetime]

--Insert into dbo.CustomerDimension select CustomerCode, cust1.FirstName, cust1.LastName,BirthDate,
MaritalStatus,Gender,Postcode,City,income
From BS_Staging.[dbo].[CustomerDetails1] cust1
, BS_Staging.[dbo].[CustomerDetails2] cust2
where cust1.firstname= cust2.FirstName and cust1.lastname = cust2.LastName

```

100 % < Messages
(19 rows affected)
Completion time: 2023-12-09T16:21:20.78446176+00:00

Fig: Inserting data from BS_Staging to CustomerDimiension

In this SQL code, we have used a *SELECT* query to select relevant columns from the staging area “BS_Staging” and *INSERT* them into the “CustomerDimension” table. *WHERE* clause filters the data to include only records where the customer's *firstname* and *lastname* are same.

```

SQLQuery4.sql - LA..CF9206\varsh [70] SQLQuery3.sql - LA..CF9206\varsh [69] SQLQuery2.sql - LA..CF9206\varsh [53]* < X
--Insert into [dbo].[PaymentTypeDimension] select Name, RetailerPaymentTypeID
From [BS_Staging].[dbo].[PaymentsData]

--Insert into SellingChannelDimension select [Name],[Code],[CommissionRate]
From [BS_Staging].[dbo].[Selling Channels]

```

100 % < Messages
(5 rows affected)
Completion time: 2023-12-09T15:51:40.7006335+00:00

Fig: Inserting data from BS_Staging to PaymentTpeDimiension and SellingChannelDimiension

This procedure is done for each dimension table and is a key step in keeping a data warehouse up-to-date and dependable.

In this SQL script, we are inserting particular data from the BS_Staging to the *DateDimension* table, such as *FullDate*, *CalendarQuarter*, *MonthOfYear*, and *CalenderYear*.

```
SQLQuery4.sql : [A..CF9206\varsh(70)] SQLQuery3.sql : [A..CF9206\varsh(69)] SQLQuery2.sql : [A..CF9206\varsh(53)]  
Insert into [dbo].[PaymentTypeDimension] select Name, RetailerPaymentTypeID  
From [B5_Staging].[dbo].[PaymentsData]  
  
Insert into [dbo].[SellingChannelDimension] select [Name],[Code],[CommissionRate]  
From [B5_Staging].[dbo].[Selling Channels]  
  
Insert into [dbo].[DateDimension] select FullDate,MonthOfYear,CalendarQuarter,CalendarYear  
From [B5_Staging].[dbo].[Generated DateTime]  
  
100 % - 4  
# Messages  
  
(7671 rows affected)  
Completion time: 2023-12-09 15:57:20.5952265+00:00
```

Fig: Inserting data from BS_Staging to DateDimiension

To import data from the Blue-Sky Staging area (BS_Staging) into the *CustomerSalesTransactionFactTable*, use a similar procedure. This process involves extracting relevant data from the staging area, possibly transforming it if needed, and then loading it into the fact table.

We can get transactional data from the staging area, alter it as needed, and load it into our fact table to gain insights into customer transactions. The “*CustomerSalesTransactionFactTable*” stores data such as *TransactionKey*, *CustomerKey*, *SellingChannelKey*, *FullDate*, *InvoiceNumber*, *Commission* and *TotalCost*.

In this SQL code, we use a **SELECT** statement to extract data from the staging area (BS_Staging) and **INSERT** them into the *CustomerSalesTransactionFactTable* along with that we used a **WHERE** statement to filter the records from the BS_Staging

The screenshot shows the SQL Server Management Studio interface. The Object Explorer on the left lists database objects for 'BlueSkyOnlineDB'. A red oval highlights the 'CustomerSalesTransactionFactTable' under 'Tables'. The central pane contains a query window with two T-SQL scripts. The top script inserts data into the 'CustomerSalesTransactionFactTable' from 'CustomerDimension', 'PaymentTypeDimension', and 'SellingChannelDimension'. The bottom script selects all columns from the same fact table. Below the queries is a results grid showing transaction data from January 2017. The status bar at the bottom indicates the query was executed successfully.

```
File Edit View Project Tools Window Help
[BlueSkyOnlineDB] New Query Execute
Object Explorer
Connect Database Snapshots
[BlueSkyOnlineDB]
  □ Database Diagrams
  □ Tables
    □ System Tables
    □ FileTables
    □ External Tables
    □ Graph Tables
  □ dbo.CustomerDimension
  □ dbo.CustomerSalesTransactionFactTable
  □ dbo.ObjectDimension
  □ dbo.PaymentTypeIDimension
  □ dbo.SellingChannelDimension
  □ Dropped Ledger Tables
  □ Views
  □ External Resources
  □ Synonyms
  □ Programmability
  □ Query Store
  □ Service Broker
  □ Storage
  □ Security
BS_Staging
Data Warehousing

SQLQuery2:sql [LA_CSF9206\varsh (53)] + SQLQuery1.sql LA_CSF9206\varsh (59)
Insert into [dbo].[CustomerSalesTransactionFactTable] select cd.CustomerKey,
scd.SellingChannelKey, ptd.PaymentTypeKey, dd.FullDate, InvoiceNumber, TotalRetailPrice, TotalCost,
scd.Commission, Profit, (TotalRetailPrice - TotalCost)
from BS_Staging.[dbo].[CustomerSaleTransactions] cst, [dbo].[CustomerDimension] cd, [dbo].[DatesDimension] dd,
[dbo].[PaymentTypeIDimension] ptd, [dbo].[SellingChannelDimension] scd
where cst.PaymentTypeID = ptd.PaymentTypeID and cst.SellingChannel = scd.SellingChannelCode and
cst.Date = dd.FullDate and cst.CustomerCode = cd.CustomerCode

Select * from [dbo].[CustomerSalesTransactionFactTable]

100 %
Results Messages
1 TransactionKey CustomerKey SellingChannelKey PaymentTypeKey FullDate InvoiceNumber TotalRetailPrice TotalCost Commission Profit
2 1 2 1 3 2017-01-02 A02129649 56.00 34.00 0 22.00
2 2 6 3 3 2017-01-02 A02129649 56.00 34.00 0 22.00
3 3 2 1 14 2017-01-02 A02129649 56.00 34.00 0 22.00
4 4 2 6 14 2017-01-02 A02129649 56.00 34.00 0 22.00
5 5 21 1 3 2017-01-02 A02129649 56.00 34.00 0 22.00
6 6 21 0 3 2017-01-02 A02129649 56.00 34.00 0 22.00
7 7 21 1 14 2017-01-02 A02129649 56.00 34.00 0 22.00
8 8 21 6 14 2017-01-02 A02129649 56.00 34.00 0 22.00
9 9 3 2 4 2017-02-14 A02129648 297.00 250.00 12 37.00
10 10 3 7 4 2017-02-14 A02129648 297.00 250.00 12 37.00
11 11 3 2 15 2017-02-14 A02129648 297.00 250.00 12 37.00
-- -- -- -- -- -- -- -- -- --
```

Fig: Inserting data from BS Staging to CustomerSalesTransactionFactTable output

4. Migrate test data from the data warehouse to an Apache Hadoop platform for further analysis of Big Data using Hortonworks Data Platform (HDP)

An individual physical machine's storage capacity becomes exhausted when the amount of data increases. This expansion requires the need to split our data across multiple devices. Distributed File Systems are File Systems that handle data storage over a network of machines. HDFS is an integral part of Apache Hadoop that is designed to store huge data files with streaming data access patterns on commodity hardware clusters. Hortonworks Data Platform (HDP) has enhanced HDFS to allow heterogeneous storage media within the HDFS cluster.

4.1 Export the data warehouse database data into an external data file.

SQL Server Import and Export Wizard is a graphical tool for transferring data between SQL Server databases and other data sources, as well as exporting data to CSV files. To start with, launch SQL Server Management Studio (SSMS). Connect to SQL Server machine.

To get the customer dimension table data in a CSV file, In the Object Explorer, right-click on the database containing the table that want to export, go to "Task," and then select "Export Data."

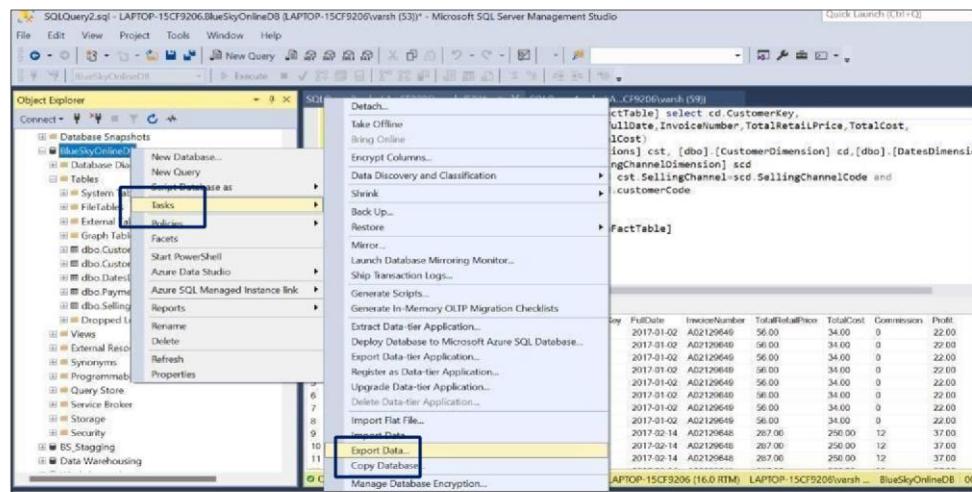


Fig: Exporting customer's data by selecting Task

The data source is the Microsoft OLE DB provider for SQL Server for the Customer Details

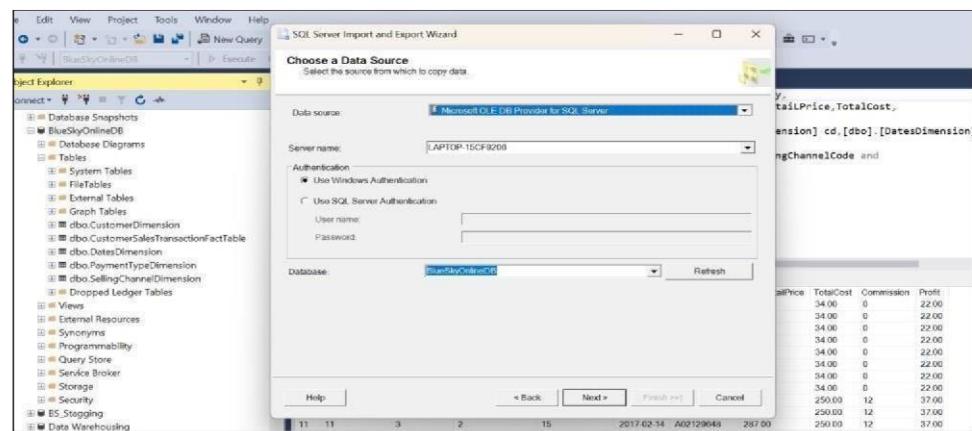


Fig: Selecting data source for data files

For the Customer dimension, choose a data source. Select a flat file source because the date data that is created remains stored as tab-separated text.

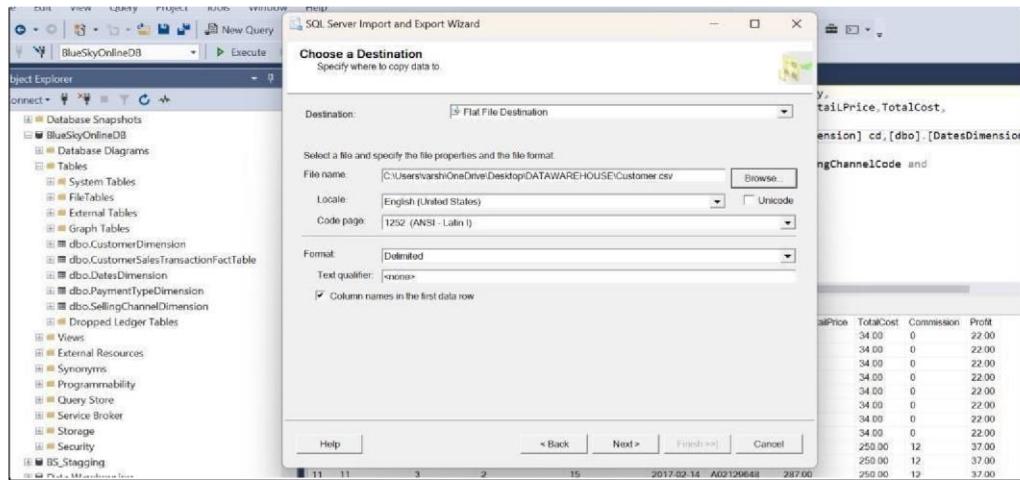


Fig: Selecting a destination for Customer details

Specifying the customer dimension as the target dimension table

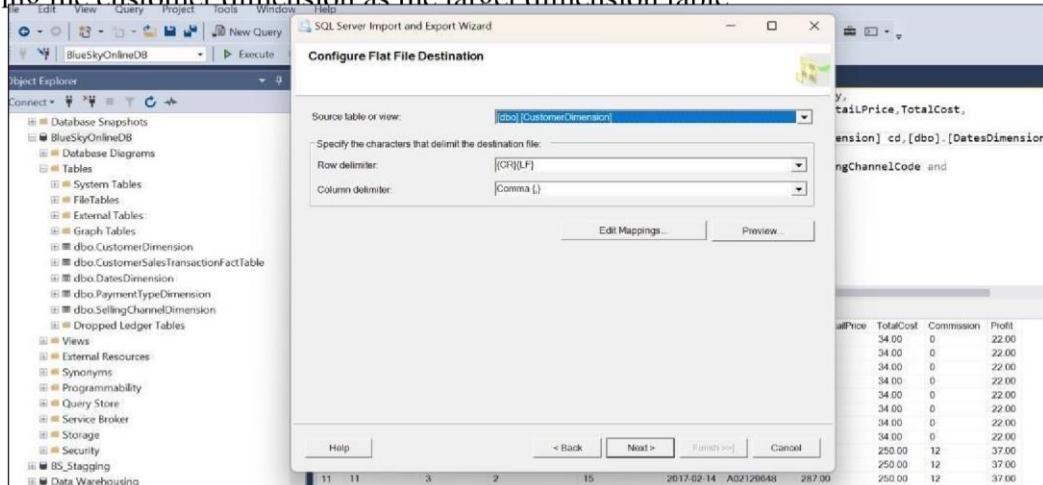


Fig: Selecting Flat File Destination for Customer dimension

If the exported data process is successful, a dialog box will appear showing that the execution was successful as well as the number of records generated and stored in the Customer dimension.csv.

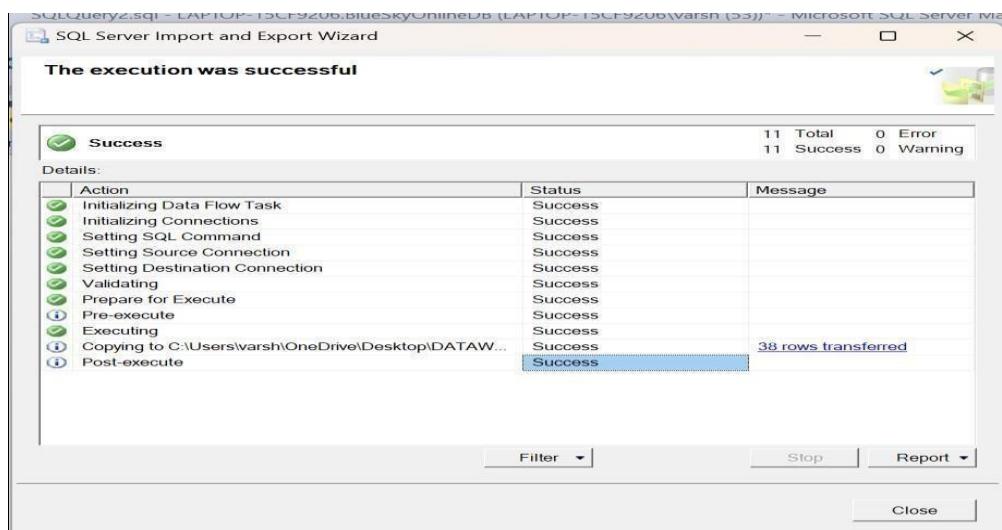


Fig: Execution window showing successful for customer dimension

In a similar way to exporting data to the other data available for Blue Sky online utilizing flat file sources, repeat the process used for the previous customer dimension.

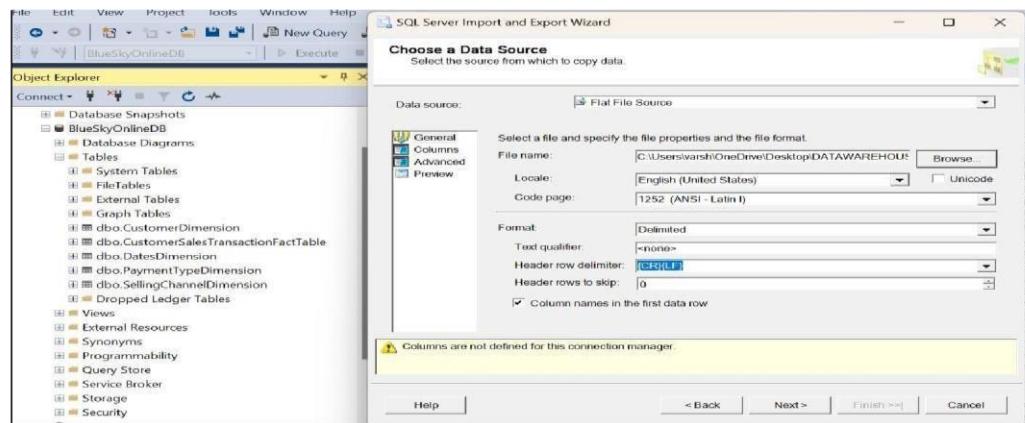


Fig: Choosing data source for customer sales transaction fact table

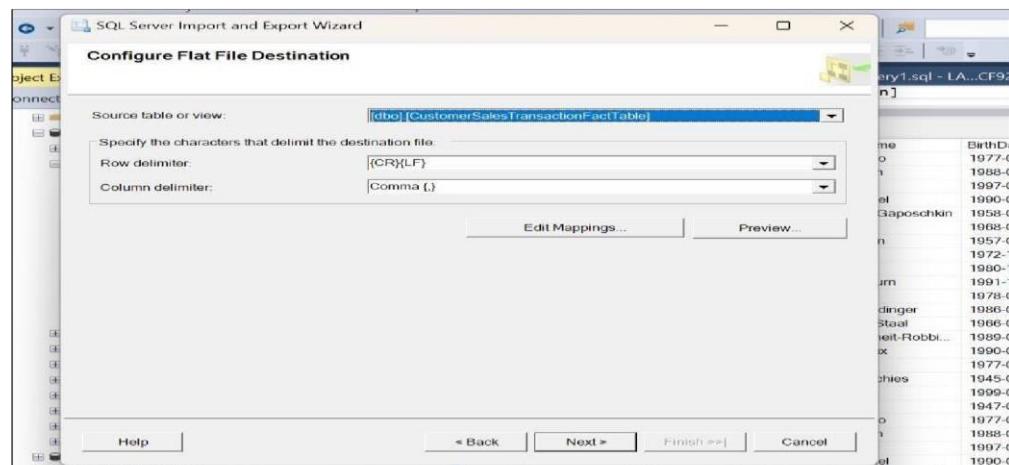


Fig: Selecting destination file for customer sales transaction fact

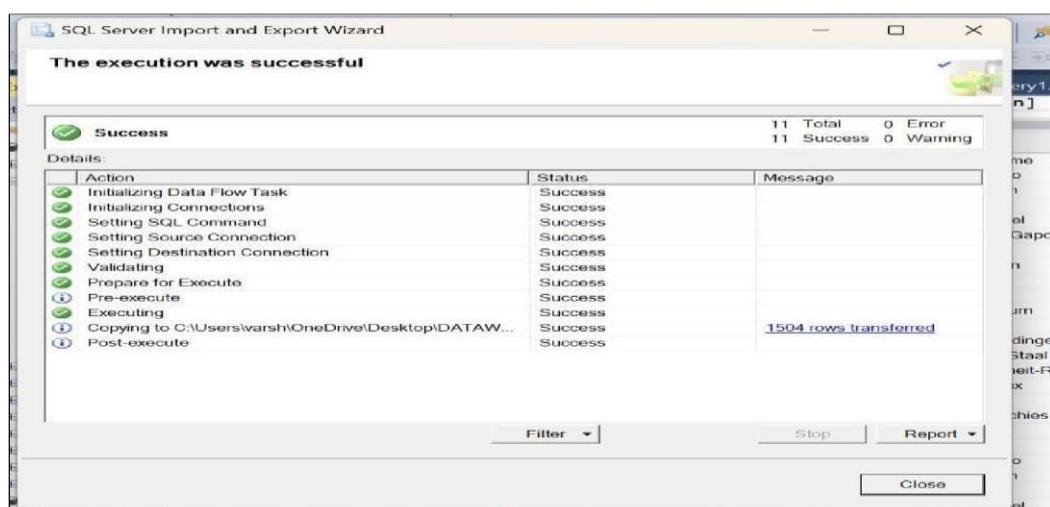


Fig: Execution window for Customer Sales Transaction Fact

4.2 Migrate the data file from the file system to Apache HDFS

The Hadoop Distributed File System (HDFS) is a critical component in the field of big data. Consider it a virtual space that spans multiple machines, allowing huge amounts of data to be

stored and processed efficiently. In simple words, instead of storing all of our data on a single computer, HDFS breaks it into smaller portions and distributes them over a cluster of machines. This technique has two major advantages:

- It allows parallel processing, which means allowing many computers to work on various sections of the data at the same time, speeding up tasks.
- It provides fault tolerance, which means that even if one computer fails, the data is still accessible and may be retrieved from other nodes in the cluster.

HDFS serves as the foundation for maintaining and analyzing large datasets in big data applications.

Load Data in HDFS: To begin with, select the HDFS Files to view from the Side panel or Offcanvas menu at the top. We can see the Hortonworks Data Platform (HDP) file store using the HDFS Files view. This is a standalone of the local file system. It will be included as part of the file system in the Hortonworks Sandbox VM.

After getting CSV files from the SQL server, we have to load the data into HDFS. To load data files in HDFS we have to log in to Ambari.

Name	Size	Last Modified	Owner	Group	Permission	Erasure Coding	Encrypted
app-logs	—	2018-11-29 17:56	yarn	hadoop	drwxrwxrwx	No	
apps	—	2018-11-29 19:01	hdfs	hdfs	drwxr-xr-x	No	
ats	—	2018-11-29 17:25	yarn	hadoop	drwxr-xr-x	No	
atsv2	—	2018-11-29 17:26	hdfs	hdfs	drwxr-xr-x	No	
hdp	—	2018-11-29 17:26	hdfs	hdfs	drwxr-xr-x	No	
key2-recovery	—	2018-11-29 17:55	livy	hdfs	drwx-----	No	
mapred	—	2018-11-29 17:26	mapred	hdfs	drwxr-xr-x	No	
mr-testery	—	2018-11-29 17:26	mapred	hadoop	drwxrwxrwx	No	
ranger	—	2018-11-29 18:54	hdfs	hdfs	drwxr-xr-x	No	
spark2-history	—	2020-09-21 16:55	spark	hadoop	drwxrwxrwx	No	
tmp	—	2023-12-18 19:52	hdfs	hdfs	drwxrwxrwx	No	

Fig: Ambari Dashboard

- Log in to the Ambari web interface using your web browser <http://localhost:8080>
- Enter the login credentials Maria_dev for both username and password.
- Go to the burger icon and select file by clicking on “File View.”
- Here we can see all of the files to which the logged-in user maria_dev has access.
- Switch to the “tmp” directory by clicking on the directory links.

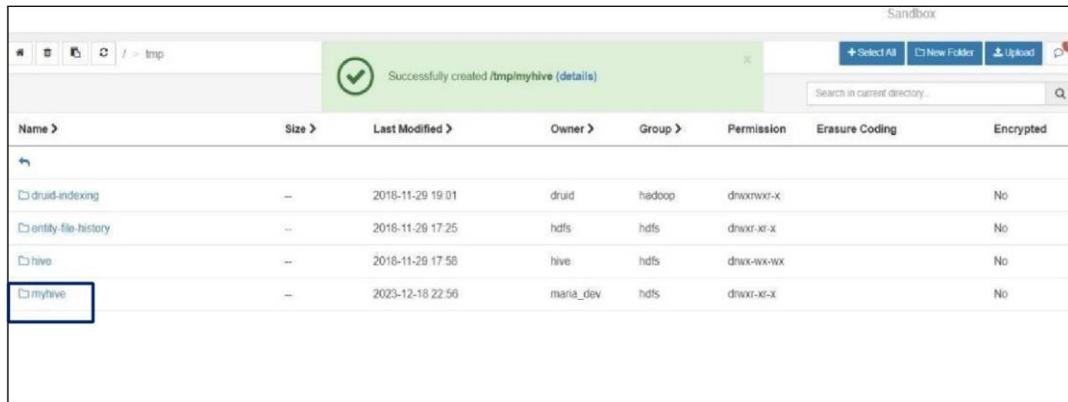


Fig: Creating a new folder for data files

- Create a separate folder for data files as my hive. To add that directory, click the New Folder option, then name the folder and create it, along with giving it permission to read, write, and execute. The folder's path we should see /tmp/myhive.

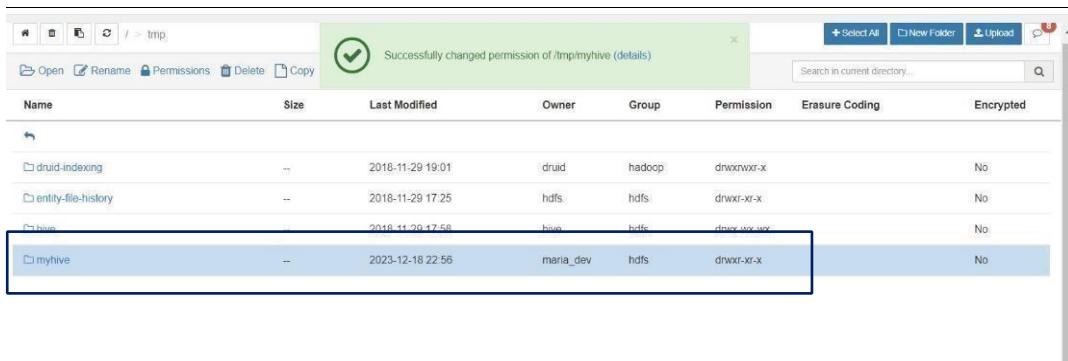


Fig: Giving permissions to the folder

- We are on the newly created directory path /tmp/myhive, so navigate to the myhive folder. Then, click the upload button to upload into the Hortonworks Sandbox environment with suitable data files **CustomerDetails.csv** and **PaymentDetails.csv** files.

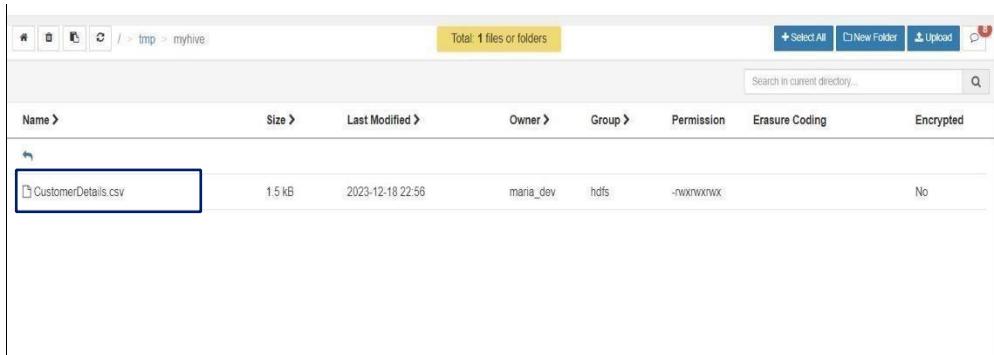
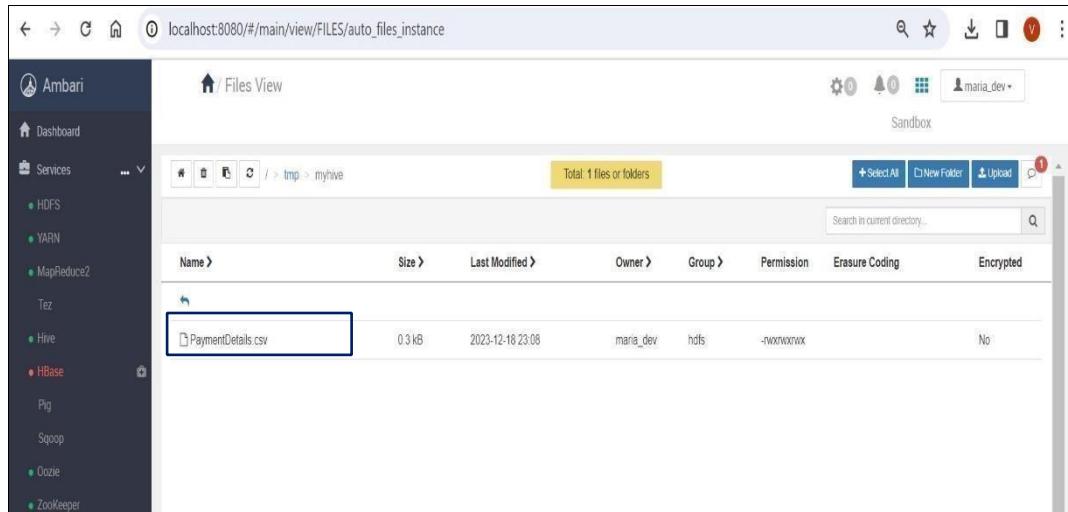


Fig: Uploading Customer Details file in HDFS

By clicking on the entity's row, we can execute the following operations such as *open*, *rename*, *permissions*, *delete*, *copy*, *move*, *download*, and *concatenate*.



The screenshot shows the Ambari File View interface at the URL `localhost:8080/#/main/view/FILES/auto_files_instance`. The left sidebar lists various services: Dashboard, Services (HDFS, YARN, MapReduce2, Tez, Hive, HBase, Pig, Sqoop, Oozie, ZooKeeper). The main area is titled 'Files View' and shows a file tree: '/tmp > myhive'. A single file, 'PaymentDetails.csv', is listed in the table below. The table columns are: Name, Size, Last Modified, Owner, Group, Permission, Erasure Coding, and Encrypted. The 'PaymentDetails.csv' file has a size of 0.3 kB, was last modified on 2023-12-18 23:08, belongs to owner 'maria_dev' and group 'hdfs', has permissions '-rwxrwxrwx', and is not encrypted.

Name	Size	Last Modified	Owner	Group	Permission	Erasure Coding	Encrypted
PaymentDetails.csv	0.3 kB	2023-12-18 23:08	maria_dev	hdfs	-rwxrwxrwx		No

Fig: Uploading Payment Details file in HDFS

4.3 Create a suitable data structure for loading the data file into the HIVE.

Apache Hive is a data warehouse infrastructure built on Hadoop Environment, enabling data summarization, query, and analysis using the Hive Query Language (HiveQL). It facilitates reading, writing, and managing large datasets residing in distributed storage using a SQL-like language.

This SQL-like language allows users to express queries using a SQL-like syntax. Hive translates HiveQL queries into MapReduce tasks, allowing them to be executed on a Hadoop cluster for efficient data management.

To perform HIVE, Hortonworks provides Data Analytics Studio (DAS) as a platform for doing a variety of analytics activities, such as interfacing with Hive to query and analyze data. DAS is an easy-to-use web-based interface for browsing data, writing, and executing Hive queries, and carrying out analysis on the Hadoop environment.

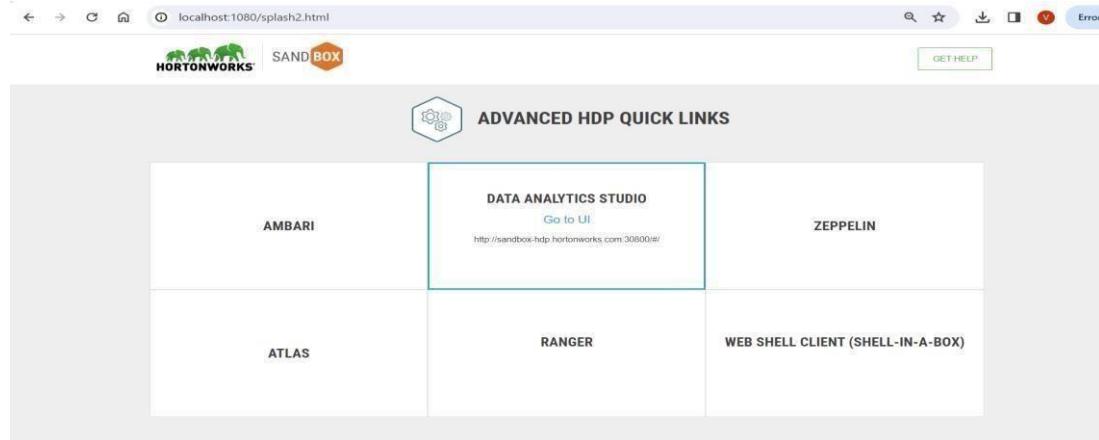


Fig: HortonWorks Dashboard

For accessing the Data Analytics Studio user interface. Open the internet browser and type DAS's URL into the address bar. <http://localhost:1080/sandbox-hdp>. when DAS running hortonworks.com at port 30800. enter our login information to sign in.

Hive is suitable for data warehousing tasks using a language similar to SQL and uses HiveQL. Using the Apache Hadoop data store as its operating platform, Hive offers a data workspace for analyzing, modifying, and altering data. It performs data processing operations one query or line at a time, it gives the ability to a user to solve problems and rearrange priorities in response to results. This adaptability makes it possible to manipulate data quickly and efficiently.

At first, we have created a table to store the data. Using the Query Editor, we will write the query. Click the **Execute** button located at the bottom of the query once you have completed it.

The screenshot shows the Data Analytics Studio Query Editor at localhost:30800/#/queries/worksheet4/results. The query is:

```

1. SELECT CREATE DATABASE IF NOT EXISTS MyBlueSkyOnline; use MyBlueSkyOnline; CREATE EXTERNAL TABLE IF NOT EXISTS CustomerDetails (
2. CUSTOMERID INT,
3. CustomerName VARCHAR(30),
4. FirstName VARCHAR(30),
5. LastName VARCHAR(30),
6. BirthDate DATE,
7. MaritalStatus VARCHAR(10),
8. Gender CHAR(1),
9. PostCode VARCHAR(10),
10. City VARCHAR(30),
11. Income INT,
12. ) ROW FORMAT DELIMITED
13. FIELDS TERMINATED BY ','
14. STORED AS TEXTFILE
15. LOCATION '/tmp/myhive/CustomerDetails';LOAD DATA INPATH '/tmp/myhive/CustomerDetails.csv' OVERWRITE INTO TABLE CustomerDetails;
  
```

The 'RESULTS' tab is selected at the bottom. The status bar at the bottom right says 'Query completed.'

Fig: Created a Query to Populate Hive Table CustomerDetails with CustomerDetails.csv Data

After the execution of LOAD DATA, we can see that the table CustomerDetails has been filled with data from the CSV file, and Hive consumed the data file CustomerDetails.csv into these stages. CustomerDetails.csv is no longer available in the File Browser. Here the LOAD command in HIVE provides the structure for the data after it has been stored in Hadoop. This is the result we got after executing the select * from CustomerDetails command which holds CustomerCode, Firstname, Lastname, Birthdate, MaritalStatus, Gender and few more.

The screenshot shows the Apache Tez Data Analytics Studio interface. The left sidebar has 'Compose' selected under 'Queries'. The main area shows a query window with the following code:

```

    1 MaritalStatus VARCHAR(10),
    2 FirstName VARCHAR(50),
    3 LastName VARCHAR(50),
    4 City VARCHAR(50),
    5 DOB DATE,
    6 Gender CHAR(1),
    7 STORE AS TEXTFILE
    8 LOCATION '/tmp/myhive/CustomerDetails';
    9 LOAD DATA INPATH '/tmp/myhive/CustomerDetails.csv' OVERWRITE INTO TABLE CustomerDetails;
    10 Select * from CustomerDetails;
  
```

Below the code, there are buttons for 'EXECUTE', 'SAVE AS', 'VISUAL EXPLAIN', 'Show Results' (which is checked), and 'Download Results'. The results section displays a table with 9 rows of data:

	CUSTOMERDETAILS.CUSTOMERCODE	CUSTOMERDETAILS.FIRSTNAME	CUSTOMERDETAILS.LASTNAME	CUSTOMERDETAILS.BIRTHDATE	CUSTOMERDETAILS.MARITALSTATUS	CUSTOMERDETAILS.GENDER
1	CC-001	Adriana	Ocampo	1977-02-01	Single	M
2	CC-002	Albert	Eriksson	1980-09-03	Married	M
3	CC-003	Anita	Henneberry	1990-09-04	Single	F
4	CC-004	Bilal	Pascal	1997-08-05	Married	M
5	CC-005	Carmen	Henschel	1980-05-06	Single	F
6	CC-006	Cecilia	Payne-Gaspardien	1995-03-07	Married	F
7	CC-007	Cheri-Shawn	Wu	1968-01-08	Single	M
8	CC-008	Dorothy	Hughes	1997-09-09	Married	F
9	CC-009	Edmond	Haley	1972-03-10	Single	M

Fig: Output of Customer Details query in Hive

The screenshot shows the Apache Tez Data Analytics Studio interface. The left sidebar has 'Compose' selected under 'Queries'. The main area shows a query window with the following code:

```

    1 CREATE DATABASE IF NOT EXISTS myBlueSkyOnline; use myBlueSkyOnline; CREATE EXTERNAL TABLE IF NOT EXISTS PaymentDetails (
    2 PaymentTypeID INT,
    3 PaymentTypeLabel VARCHAR(25),
    4 PaymentTypeID2 INT,
    5 ) ROW FORMAT DELIMITED
    6 FIELDS TERMINATED BY ","
    7 STORED AS TEXTFILE
    8 LOCATION '/tmp/myhive/PaymentDetails';
    9 LOAD DATA INPATH '/tmp/myhive/PaymentDetails.csv' OVERWRITE INTO TABLE PaymentDetails; select * from PaymentDetails;
  
```

Below the code, there are buttons for 'EXECUTE', 'SAVE AS', 'VISUAL EXPLAIN', 'Show Results' (which is checked), and 'Download Results'. The results section displays a table with 9 rows of data:

	PaymentTypeID	PaymentTypeLabel
1	1	Credit Card
2	2	Debit Card
3	3	Check
4	4	Cash
5	5	Gift Card
6	6	Mobile Payment
7	7	Bank Transfer
8	8	Digital Wallet
9	9	Other

Fig: Created a Query to Populate Hive Table PaymentDetails with PaymentDetails.csv Data

Similarly for Payment Details, we created an additional table and then loaded the PaymentDetails.csv file in HDFS. By executing the above query, we have seen our results using the select * from payment details command. The results which stored in the table like PaymenttypeID and Name with paymentkeys. As previously explained, we managed the task using Hive in step by step.

The screenshot shows a browser window for 'localhost:30800/#/queries/worksheet4/results'. The left sidebar has 'Compose' selected. The main area shows a 'Compose' tab with a query editor containing the following code:

```

1 CREATE DATABASE IF NOT EXISTS MyBlueSkyOnline;
2 USE MyBlueSkyOnline;
3 CREATE EXTERNAL TABLE IF NOT EXISTS PaymentDetails (
4   PaymentTypeKey INT,
5   PaymentTypeValue VARCHAR(20),
6   PaymentTypeDesc STRING,
7   RowFormat STRING,
8   FIELDS TERMINATED BY ','
9   GROUP BY TEXTFILE,
10 LOCATED '/tmp/myhive/PaymentDetails');
11 LOAD DATA INPATH '/tmp/myhive/PaymentDetails.csv' OVERWRITE INTO TABLE PaymentDetails;
12 select * from PaymentDetails;

```

Below the code, there are tabs for 'EXECUTE', 'SAVE AS...', 'VISUAL EXPLAIN', 'Show Details', and 'Download Results'. The 'RESULTS' tab is selected, displaying a table titled 'PaymentDetails PAYMENTTYPEKEY PAYMENTTYPEVALUE PAYMENTTYPEDESC'. The data in the table is:

Row	PaymentTypeKey	PaymentTypeValue	PaymentTypeDesc
1		Unknown	-1
2		Amazon	1
3		Gift Card	2
4		Finance	3
5		Postal	4
6		Digital	5

Fig: Output of PaymentDetails query in Hive

4.4 Demonstrate the use of Apache Pig for manipulating the loaded data.

Apache Pig is a high-level platform for Hadoop processing and analysis of large data sets. It abstracts the MapReduce programming design, allowing developers to express data processing tasks using Pig Latin, as a scripting language. Pig Latin is a scripting language for data flows that is used in Pig. It is similar to SQL and is intended to be easy to use. Pig scripts are used to define a set of data transformations that are carried out on Hadoop clusters. Pig Latin is a scripting language that is excellent for data processing tasks.

In the previous task, we have used Hive, a dataflow-focused scripting language. Pig and Hive are both part of the Hadoop ecosystem and provide a higher-level abstraction for data processing tasks, making it more accessible to users who may not have significant programming or MapReduce skills.

Pig Platform is developed by Yahoo, to make programming with MapReduce, Hadoop's assembler task easier. Pig Latin commands can be executed in one of two ways, either in Shell Commands (Grunt shell) to enter commands one at a time, and in Pig scripts we can create as a text file with a Pig script and run it with the Pig program. Here, the Grunt shell can be used to run Pig Latin scripts as well as HIVE and HDFS commands. Pig's logical plan does not run until a DUMP or STORE command is executed.

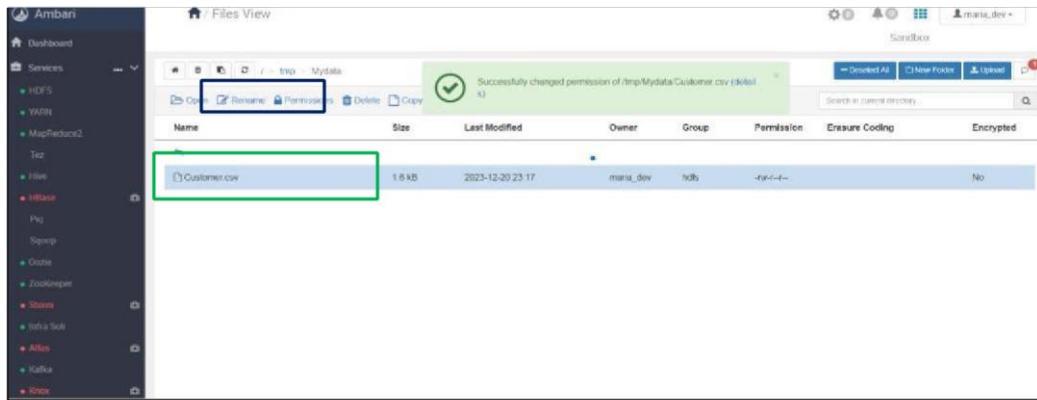


Fig: Uploaded Customer file in HDFS for Pig Task

To launch Apache Pig on Hortonworks Data Platform (HDP), browse to localhost:4200 (Shellon-A-Box) and log in with root as the login and password, then change the user name to maria_dev and change the directory to home (cd). Create a file for pig script as *vi Customer* and execute it.

Fig: Creating a script to load data and define the relation

In this step, we created a pig script to load data and construct a relationship with the data file. In the first line, we created a relation named Customer, which represents customer details, and in the following line, we used the DESCRIBE command to review customer details and to

```
 sandbox hdp login: root
root@sandbox-hdp.hortonworks.com's password:
root@sandbox-hdp.hortonworks.com's password:
root@ sandbox-hdp ~]# curl -s http://131.144.10.223 from 172.18.0.2
[root@sandbox hdp ~]# su mariadb dev
mariadb_devs@ sandbox-hdp root ]$ cd
[mariadb_devs@ sandbox-hdp root ]$ Customer
[mariadb_devs@ sandbox-hdp ~]# █
```

Fig: Web shell in pig with login details

To create a script for loading data and defining a relation with the data file `Customer.csv`

```
customer = Load '/tmp/Mydata/Customer.csv' USING PigStorage(',');
DESCRIBE customer;
-- INSERT --
```

execute the script pig -f Customers then Enter before that we have to save changes by clicking esc :x.

Fig: Script creates map-reduce tasks as unknown

By using the following code in the shell with the schema to give pig datatypes to the Customer and save it. Once the execution done, enter vi script.

```

customer = Load '/tmp/hydata/customer.csv' Using PigStorage(',') As (customerkey: int,customercode : chararray,firstname : chararray,lastname : chararray,birthdate : chararray,maritalstatus : chararray,income : long);
DESCRIBE customer;

```

Fig: Defining relation with `customer`

```

root@sandbox-hdp.hortonworks.com's password:
Last login: Wed Dec 20 23:13:34 2028 from 172.18.0.2
[root@maria-dev sandbox-hdp] ~
[root@maria-dev sandbox-hdp] ~
[root@maria-dev sandbox-hdp] ~$ cd /tmp
[root@maria-dev sandbox-hdp] ~$ ls
[root@maria-dev sandbox-hdp] ~$ pig -f customer
WARNING: HADOOP_PREFIX has been replaced by HADOOP_HOME. Using value of HADOOP_PREFIX.
2023-12-20 23:24:48 INFO org.apache.pig.ExecTypeProvider: Trying ExecType : LOCAL
2023-12-20 23:24:48 INFO org.apache.pig.ExecTypeProvider: Trying ExecType : MAPREDUCE
2023-12-20 23:24:48 INFO org.apache.pig.ExecTypeProvider: Trying ExecType : TEE_LOCAL
2023-12-20 23:24:48 INFO org.apache.pig.ExecTypeProvider: Trying ExecType : TEE
2023-12-20 23:24:48 INFO org.apache.pig.ExecTypeProvider: Picked TEE as the ExecType
2023-12-20 23:24:48,999 [main] INFO org.apache.pig.Main - Apache Pig version 0.16.0-3.9.1-0-187 (unversioned directory) compiled Sep 19 2018, 10:13:33
2023-12-20 23:24:49,000 [main] INFO org.apache.pig.Main - Logging error messages to /home/maria/dev/pig_170352810243-1.log
2023-12-20 23:24:49,505 [main] INFO org.apache.pig.impl.util.Utils: Default bootstrap file /home/maria/dev/pigbootup not found
2023-12-20 23:24:49,583 [main] INFO org.apache.pig.backend.hadoop.executionengine.HExecutionEngine - Connecting to hadoop file system at: hdfs://sandbox-hdp.hortonworks.com:8020
2023-12-20 23:24:49,583 [main] INFO org.apache.pig.PigServer - Pig Script ID for the session: PIG-CUSTOMER-d958a22-6e2d-4571-be12-83fbad54ccc
2023-12-20 23:24:49,583 [main] WARN org.apache.pig.PigServer - ATS is disabled since yarn.timeline-service.enabled set to false
2023-12-20 23:24:51,129 [main] INFO org.apache.pig.Main - Pig script completed in 2 seconds and 734 milliseconds (2734 ms)
[maria_dev@sandbox-hdp ~]$ vi customer
[maria_dev@sandbox-hdp ~]$ pig -f customer
WARNING: HADOOP_PREFIX has been replaced by HADOOP_HOME. Using value of HADOOP_PREFIX.
2023-12-20 23:43:39 INFO org.apache.pig.ExecTypeProvider: Trying ExecType : LOCAL
2023-12-20 23:43:39 INFO org.apache.pig.ExecTypeProvider: Trying ExecType : MAPREDUCE
2023-12-20 23:43:39 INFO org.apache.pig.ExecTypeProvider: Trying ExecType : TEE_LOCAL
2023-12-20 23:43:39 INFO org.apache.pig.ExecTypeProvider: Trying ExecType : TEE
2023-12-20 23:43:39 INFO org.apache.pig.ExecTypeProvider: Picked TEE as the ExecType
2023-12-20 23:43:39,245 [main] INFO org.apache.pig.Main - Apache Pig version 0.16.0-3.9.1-0-187 (unversioned directory) compiled Sep 19 2018, 10:13:33
2023-12-20 23:43:39,245 [main] INFO org.apache.pig.Main - Logging error messages to /home/maria/dev/pig_170352810243-1.log
2023-12-20 23:43:39,245 [main] INFO org.apache.pig.impl.util.Utils: Default bootstrap file /home/maria/dev/pigbootup not found
2023-12-20 23:43:39,818 [main] INFO org.apache.pig.backend.hadoop.executionengine.HExecutionEngine - Connecting to hadoop file system at: hdfs://sandbox-hdp.hortonworks.com:8020
2023-12-20 23:43:39,818 [main] INFO org.apache.pig.PigServer - Pig Script ID for the session: PIG-CUSTOMER-daa28be-6ab6-48c6-862f-735351209bde
[maria_dev@sandbox-hdp ~]$
```

Fig: Successfully Defined scheme for the Customer Relation

In the next step, we created a new relation with the available Customer data, which is a collection of 10 entries from the Customer relation. Now save and execute the Pig query.

```

Customer = Load '/tmp/Metadata/Customer.csv' USING PigStorage(',') AS (CustomerKey:int,CustomerCode:chararray,FirstName:chararray,LastName:chararray,BirthDate:chararray,MaritalStatus:chararray,Gender:chararray,PostCode:chararray,City:chararray,Income:long);
DESCRIBE Customer;
Customer_subset = LIMIT Customer 10;
DESCRIBE Customer_subset;

```

```

2013-12-20 23:13:24 [51,629] [main] INFO org.apache.pig.Main - Pig script completed in 2 seconds and 754 milliseconds (2734 ms)
[maria1a_devisandbox-hdp ~]$ pig -f Customer
WARNING: HADOOP_PREFIX has been replaced by HADOOP_HOME. Using value of HADOOP_PREFIX.
23/12/20 23:14:33: INFO pig.ExecTypeProvider: Trying ExecType : LOCAL
23/12/20 23:14:36: INFO pig.ExecTypeProvider: Trying ExecType : MAPREDUCE
23/12/20 23:14:36: INFO pig.ExecTypeProvider: Trying ExecType : LOCAL
23/12/20 23:14:36: INFO pig.ExecTypeProvider: Trying ExecType : TEZ
23/12/20 23:14:36: INFO pig.ExecTypeProvider: Picked TEZ as the ExecType
2013-12-20 23:14:37:0,245 [main] INFO org.apache.pig.Main - Apache Pig version 0.16.0-0.1.0-187 (r11734) directory compiled Sep 19 2010, 10:13:33
2013-12-20 23:14:37,0,245 [main] INFO org.apache.pig.Main - Logging error messages to: /home/maria/dev/pig_1703165059424.log
2013-12-20 23:14:37,0,245 [main] INFO org.apache.pig.impl.util.Utils - Default bootstrap file /home/maria/dev/pigbootup not found
2013-12-20 23:14:37,0,245 [main] INFO org.apache.pig.backend.hadoop.executionengine.ExecutingEngine - Connecting to hadoop file system at: hdfs://sandbox-hdp.hortonworks.com:8020
2013-12-20 23:14:37,0,245 [main] INFO org.apache.pig.PigServer - ATS is disabled since yarn.timeline-service.enabled set to false
Customer = (CustomerKey:int,CustomerCode:chararray,FirstName:chararray,LastName:chararray,BirthDate:chararray,MaritalStatus:chararray,Gender:chararray,PostCode:chararray,City:chararray,Income:long);
2013-12-20 23:14:37,2,255 [main] INFO org.apache.pig.Main - Pig script completed in 2 seconds and 88 milliseconds (2088 ms)
[maria1a_devisandbox-hdp ~]$ vi Customer
WARNING: HADOOP_PREFIX has been replaced by HADOOP_HOME. Using value of HADOOP_PREFIX.
23/12/20 23:15:00: INFO pig.ExecTypeProvider: Trying ExecType : LOCAL
23/12/20 23:15:00: INFO pig.ExecTypeProvider: Trying ExecType : MAPREDUCE
23/12/20 23:15:00: INFO pig.ExecTypeProvider: Trying ExecType : LOCAL
23/12/20 23:15:00: INFO pig.ExecTypeProvider: Trying ExecType : TEZ
23/12/20 23:15:00: INFO pig.ExecTypeProvider: Picked TEZ as the ExecType
2013-12-20 23:15:00,506 [main] INFO org.apache.pig.Main - Apache Pig version 0.16.0-0.1.0-187 (r11734) directory compiled Sep 19 2010, 10:13:33
2013-12-20 23:15:00,506 [main] INFO org.apache.pig.Main - Logging error messages to: /home/maria/dev/pig_1703165059425.log
2013-12-20 23:15:01,277 [main] INFO org.apache.pig.impl.util.Utils - Default bootstrap file /home/maria/dev/pigbootup not found
2013-12-20 23:15:01,277 [main] INFO org.apache.pig.backend.hadoop.executionengine.ExecutingEngine - Connecting to hadoop file system at: hdfs://sandbox-hdp.hortonworks.com:8020
2013-12-20 23:15:01,277 [main] INFO org.apache.pig.PigServer - Pig Script ID for the session: DFO-Customer_d4a28be-66b6-49c-b62f-73535a12b0d6
2013-12-20 23:15:01,277 [main] INFO org.apache.pig.PigServer - ATS is disabled since yarn.timeline-service.enabled set to false
Customer = (CustomerKey:int,CustomerCode:chararray,FirstName:chararray,LastName:chararray,BirthDate:chararray,MaritalStatus:chararray,Gender:chararray,PostCode:chararray,City:chararray,Income:long);
Customer_subset = (CustomerKey:int,CustomerCode:chararray,FirstName:chararray,LastName:chararray,BirthDate:chararray,MaritalStatus:chararray,Gender:chararray,PostCode:chararray,City:chararray,Income:long);
2013-12-20 23:15:03,425 [main] INFO org.apache.pig.Main - Pig script completed in 3 seconds and 16 milliseconds (3016 ms)
[maria1a_devisandbox-hdp ~]$ 
```

Fig: Successfully executed above query by defining subset Customer

To view the data Customer relation by using the DUMP command as (DUMP Customer_subset), inserting these commands in pig script, and again save and execute the code.

```

Customer_subset = CustomerKey: int,CustomerCode: chararray,FirstName: chararray,LastName: chararray,BirthDate: chararray,MaritalStatus: chararray,Gender: chararray,PostCode: chararray,City: chararray,Income: long;
2013-12-20 23:15:09,57,37,444 [main] INFO org.apache.pig.tools.pigstats.ScriptState - Pig features used in the script: LIMIT
2013-12-20 23:15:09,57,444 [main] INFO org.apache.pig.data.SchemaTupleBackend - Key [pig.schematuple] was not set... will not generate code.
2013-12-20 23:15:09,57,444 [main] INFO org.apache.pig.newplan.logical.optimizer.LogicalPlanOptimizer - {RULES_ENABLED=[AddOrReplace, columnMapKeyPrune, constantCalculator, GroupByConstParallelSetter, LimitOptimizer, LoadTypeCastInserter, MergeFilter, MergeForEach, PartitionFilterOptimizer, PredicatePushdownOptimizer, PushUpFilter, SplitFilter, StreamTypeCastInserter]}
2013-12-20 23:15:09,57,444 [main] INFO org.apache.pig.impl.util.SplittableMemoryManager - Selected heap (PS Old Gen) of size 699400192 to monitor. collectionsUsageThreshold = 489580128, usSystemMemory = 489580128
2013-12-20 23:15:09,58,371 [main] INFO org.apache.hadoop.mapreduce.lib.output.FileOutputCommitter - File Output Committer APIVersion is 2
2013-12-20 23:15:09,58,371 [main] INFO org.apache.hadoop.mapreduce.lib.output.FileOutputCommitter - FileOutputCommitter: skip cleanup _temporary folders under output directory:false, ignore cleanup failures: false
2013-12-20 23:15:09,58,453 [main] INFO org.apache.pig.data.SchemaTupleBackend - Key [pig.schematuple] was not set... will not generate code.
2013-12-20 23:15:09,58,526 [main] WARN org.apache.pig.data.SchemaTupleBackend - SchemaTupleBackend has already been initialized
2013-12-20 23:15:09,58,526 [main] INFO org.apache.pig.backend.PigInputFormat - Input paths to process : 1
2013-12-20 23:15:09,58,526 [main] INFO org.apache.pig.backend.hadoop.executionengine.util.NeededUtil - Total input paths to process : 1
2013-12-20 23:15:09,58,526 [main] INFO org.apache.pig.backend.hadoop.executionengine.PigHadoopMapper - org.apache.pig.builtin.Utf8StorageConverter[FIELD_DISCARDED_TYPE_ON_VERSION_FAILED]: Unable to interpret value [07, 117, 115, 116, 111, 109, 101, 114, 79, 101, 121] in field being converted to int, caught NumberFormatException <for input string: "CustomerKey"> Field discarded
2013-12-20 23:15:09,58,526 [main] WARN org.apache.pig.backend.hadoop.executionengine.wapireducerlayer.PigHadoopMapper - org.apache.pig.builtin.Utf8StorageConverter[FIELD_DISCARDED_TYPE_ON_VERSION_FAILED]: <for input string: "CustomerCode"> Field discarded
2013-12-20 23:15:09,58,526 [main] INFO org.apache.hadoop.mapreduce.lib.output.FileOutputCommitter - Saved output of task 'attempt_0001_m_000001' to hdfs://sandbox-hdp.hortonworks.com:8020/tmp/_teep990977210/tm109175019
2013-12-20 23:15:09,59,117 [main] WARN org.apache.pig.data.SchemaTupleBackend - SchemaTupleBackend has already been initialized
2013-12-20 23:15:09,59,143 [main] INFO org.apache.hadoop.mapreduce.lib.input.FileInputFormat - Total input files to process : 1
2013-12-20 23:15:09,59,143 [main] INFO org.apache.pig.backend.hadoop.executionengine.util.NeededUtil - Total input paths to process : 1
Customer,FirstName,LastName,BirthDate,MaritalStatus,Gender,PostCode,City
1,CC-001,Albert,Pminster,01/05/1968,Married,Male,NE1 8PN,London,45000
2,CC-002,Anna,Behrensmeyer,04/09/1990,Single,Female,NB 9GW,London,23000
3,CC-003,Elaise,Pascal,09/05/1997,Married,Male,SE8 9GW,London,24000
4,CC-004,Caroline,Herschel,05/06/1990,Single,Female,NB 9YH,London,25000
5,CC-005,Cecilia,Payne-Gapschkin,03/07/1958,Married,Female,E3 6YW,London,23000
6,CC-006,Chien-Shiung,Wu,01/08/1968,Single,Male,E8 3PW,London,45000
7,CC-007,Albert,Shiung,Wu,01/08/1968,Single,Male,E8 3PW,London,45000
8,CC-008,Albert,Shiung,Wu,01/08/1968,Single,Male,E8 3PW,London,45000
9,CC-009,Albert,Shiung,Wu,01/08/1968,Single,Male,E8 3PW,London,45000
10,CC-010,Albert,Shiung,Wu,01/08/1968,Single,Male,E8 3PW,London,45000
11,CC-011,Albert,Shiung,Wu,01/08/1968,Single,Male,E8 3PW,London,45000
12,CC-012,Albert,Shiung,Wu,01/08/1968,Single,Male,E8 3PW,London,45000
13,CC-013,Albert,Shiung,Wu,01/08/1968,Single,Male,E8 3PW,London,45000
14,CC-014,Albert,Shiung,Wu,01/08/1968,Single,Male,E8 3PW,London,45000
15,CC-015,Albert,Shiung,Wu,01/08/1968,Single,Male,E8 3PW,London,45000
16,CC-016,Albert,Shiung,Wu,01/08/1968,Single,Male,E8 3PW,London,45000
17,CC-017,Albert,Shiung,Wu,01/08/1968,Single,Male,E8 3PW,London,45000
18,CC-018,Albert,Shiung,Wu,01/08/1968,Single,Male,E8 3PW,London,45000
19,CC-019,Albert,Shiung,Wu,01/08/1968,Single,Male,E8 3PW,London,45000
20,CC-020,Albert,Shiung,Wu,01/08/1968,Single,Male,E8 3PW,London,45000
21,CC-021,Albert,Shiung,Wu,01/08/1968,Single,Male,E8 3PW,London,45000
22,CC-022,Albert,Shiung,Wu,01/08/1968,Single,Male,E8 3PW,London,45000
23,CC-023,Albert,Shiung,Wu,01/08/1968,Single,Male,E8 3PW,London,45000
24,CC-024,Albert,Shiung,Wu,01/08/1968,Single,Male,E8 3PW,London,45000
25,CC-025,Albert,Shiung,Wu,01/08/1968,Single,Male,E8 3PW,London,45000
26,CC-026,Albert,Shiung,Wu,01/08/1968,Single,Male,E8 3PW,London,45000
27,CC-027,Albert,Shiung,Wu,01/08/1968,Single,Male,E8 3PW,London,45000
28,CC-028,Albert,Shiung,Wu,01/08/1968,Single,Male,E8 3PW,London,45000
29,CC-029,Albert,Shiung,Wu,01/08/1968,Single,Male,E8 3PW,London,45000
30,CC-030,Albert,Shiung,Wu,01/08/1968,Single,Male,E8 3PW,London,45000
31,CC-031,Albert,Shiung,Wu,01/08/1968,Single,Male,E8 3PW,London,45000
32,CC-032,Albert,Shiung,Wu,01/08/1968,Single,Male,E8 3PW,London,45000
33,CC-033,Albert,Shiung,Wu,01/08/1968,Single,Male,E8 3PW,London,45000
34,CC-034,Albert,Shiung,Wu,01/08/1968,Single,Male,E8 3PW,London,45000
35,CC-035,Albert,Shiung,Wu,01/08/1968,Single,Male,E8 3PW,London,45000
36,CC-036,Albert,Shiung,Wu,01/08/1968,Single,Male,E8 3PW,London,45000
37,CC-037,Albert,Shiung,Wu,01/08/1968,Single,Male,E8 3PW,London,45000
38,CC-038,Albert,Shiung,Wu,01/08/1968,Single,Male,E8 3PW,London,45000
39,CC-039,Albert,Shiung,Wu,01/08/1968,Single,Male,E8 3PW,London,45000
40,CC-040,Albert,Shiung,Wu,01/08/1968,Single,Male,E8 3PW,London,45000
41,CC-041,Albert,Shiung,Wu,01/08/1968,Single,Male,E8 3PW,London,45000
42,CC-042,Albert,Shiung,Wu,01/08/1968,Single,Male,E8 3PW,London,45000
43,CC-043,Albert,Shiung,Wu,01/08/1968,Single,Male,E8 3PW,London,45000
44,CC-044,Albert,Shiung,Wu,01/08/1968,Single,Male,E8 3PW,London,45000
45,CC-045,Albert,Shiung,Wu,01/08/1968,Single,Male,E8 3PW,London,45000
46,CC-046,Albert,Shiung,Wu,01/08/1968,Single,Male,E8 3PW,London,45000
47,CC-047,Albert,Shiung,Wu,01/08/1968,Single,Male,E8 3PW,London,45000
48,CC-048,Albert,Shiung,Wu,01/08/1968,Single,Male,E8 3PW,London,45000
49,CC-049,Albert,Shiung,Wu,01/08/1968,Single,Male,E8 3PW,London,45000
50,CC-050,Albert,Shiung,Wu,01/08/1968,Single,Male,E8 3PW,London,45000
51,CC-051,Albert,Shiung,Wu,01/08/1968,Single,Male,E8 3PW,London,45000
52,CC-052,Albert,Shiung,Wu,01/08/1968,Single,Male,E8 3PW,London,45000
53,CC-053,Albert,Shiung,Wu,01/08/1968,Single,Male,E8 3PW,London,45000
54,CC-054,Albert,Shiung,Wu,01/08/1968,Single,Male,E8 3PW,London,45000
55,CC-055,Albert,Shiung,Wu,01/08/1968,Single,Male,E8 3PW,London,45000
56,CC-056,Albert,Shiung,Wu,01/08/1968,Single,Male,E8 3PW,London,45000
57,CC-057,Albert,Shiung,Wu,01/08/1968,Single,Male,E8 3PW,London,45000
58,CC-058,Albert,Shiung,Wu,01/08/1968,Single,Male,E8 3PW,London,45000
59,CC-059,Albert,Shiung,Wu,01/08/1968,Single,Male,E8 3PW,London,45000
60,CC-060,Albert,Shiung,Wu,01/08/1968,Single,Male,E8 3PW,London,45000
61,CC-061,Albert,Shiung,Wu,01/08/1968,Single,Male,E8 3PW,London,45000
62,CC-062,Albert,Shiung,Wu,01/08/1968,Single,Male,E8 3PW,London,45000
63,CC-063,Albert,Shiung,Wu,01/08/1968,Single,Male,E8 3PW,London,45000
64,CC-064,Albert,Shiung,Wu,01/08/1968,Single,Male,E8 3PW,London,45000
65,CC-065,Albert,Shiung,Wu,01/08/1968,Single,Male,E8 3PW,London,45000
66,CC-066,Albert,Shiung,Wu,01/08/1968,Single,Male,E8 3PW,London,45000
67,CC-067,Albert,Shiung,Wu,01/08/1968,Single,Male,E8 3PW,London,45000
68,CC-068,Albert,Shiung,Wu,01/08/1968,Single,Male,E8 3PW,London,45000
69,CC-069,Albert,Shiung,Wu,01/08/1968,Single,Male,E8 3PW,London,45000
70,CC-070,Albert,Shiung,Wu,01/08/1968,Single,Male,E8 3PW,London,45000
71,CC-071,Albert,Shiung,Wu,01/08/1968,Single,Male,E8 3PW,London,45000
72,CC-072,Albert,Shiung,Wu,01/08/1968,Single,Male,E8 3PW,London,45000
73,CC-073,Albert,Shiung,Wu,01/08/1968,Single,Male,E8 3PW,London,45000
74,CC-074,Albert,Shiung,Wu,01/08/1968,Single,Male,E8 3PW,London,45000
75,CC-075,Albert,Shiung,Wu,01/08/1968,Single,Male,E8 3PW,London,45000
76,CC-076,Albert,Shiung,Wu,01/08/1968,Single,Male,E8 3PW,London,45000
77,CC-077,Albert,Shiung,Wu,01/08/1968,Single,Male,E8 3PW,London,45000
78,CC-078,Albert,Shiung,Wu,01/08/1968,Single,Male,E8 3PW,London,45000
79,CC-079,Albert,Shiung,Wu,01/08/1968,Single,Male,E8 3PW,London,45000
80,CC-080,Albert,Shiung,Wu,01/08/1968,Single,Male,E8 3PW,London,45000
81,CC-081,Albert,Shiung,Wu,01/08/1968,Single,Male,E8 3PW,London,45000
82,CC-082,Albert,Shiung,Wu,01/08/1968,Single,Male,E8 3PW,London,45000
83,CC-083,Albert,Shiung,Wu,01/08/1968,Single,Male,E8 3PW,London,45000
84,CC-084,Albert,Shiung,Wu,01/08/1968,Single,Male,E8 3PW,London,45000
85,CC-085,Albert,Shiung,Wu,01/08/1968,Single,Male,E8 3PW,London,45000
86,CC-086,Albert,Shiung,Wu,01/08/1968,Single,Male,E8 3PW,London,45000
87,CC-087,Albert,Shiung,Wu,01/08/1968,Single,Male,E8 3PW,London,45000
88,CC-088,Albert,Shiung,Wu,01/08/1968,Single,Male,E8 3PW,London,45000
89,CC-089,Albert,Shiung,Wu,01/08/1968,Single,Male,E8 3PW,London,45000
90,CC-090,Albert,Shiung,Wu,01/08/1968,Single,Male,E8 3PW,London,45000
91,CC-091,Albert,Shiung,Wu,01/08/1968,Single,Male,E8 3PW,London,45000
92,CC-092,Albert,Shiung,Wu,01/08/1968,Single,Male,E8 3PW,London,45000
93,CC-093,Albert,Shiung,Wu,01/08/1968,Single,Male,E8 3PW,London,45000
94,CC-094,Albert,Shiung,Wu,01/08/1968,Single,Male,E8 3PW,London,45000
95,CC-095,Albert,Shiung,Wu,01/08/1968,Single,Male,E8 3PW,London,45000
96,CC-096,Albert,Shiung,Wu,01/08/1968,Single,Male,E8 3PW,London,45000
97,CC-097,Albert,Shiung,Wu,01/08/1968,Single,Male,E8 3PW,London,45000
98,CC-098,Albert,Shiung,Wu,01/08/1968,Single,Male,E8 3PW,London,45000
99,CC-099,Albert,Shiung,Wu,01/08/1968,Single,Male,E8 3PW,London,45000
100,CC-100,Albert,Shiung,Wu,01/08/1968,Single,Male,E8 3PW,London,45000
101,CC-101,Albert,Shiung,Wu,01/08/1968,Single,Male,E8 3PW,London,45000
102,CC-102,Albert,Shiung,Wu,01/08/1968,Single,Male,E8 3PW,London,45000
103,CC-103,Albert,Shiung,Wu,01/08/1968,Single,Male,E8 3PW,London,45000
104,CC-104,Albert,Shiung,Wu,01/08/1968,Single,Male,E8 3PW,London,45000
105,CC-105,Albert,Shiung,Wu,01/08/1968,Single,Male,E8 3PW,London,45000
106,CC-106,Albert,Shiung,Wu,01/08/1968,Single,Male,E8 3PW,London,45000
107,CC-107,Albert,Shiung,Wu,01/08/1968,Single,Male,E8 3PW,London,45000
108,CC-108,Albert,Shiung,Wu,01/08/1968,Single,Male,E8 3PW,London,45000
109,CC-109,Albert,Shiung,Wu,01/08/1968,Single,Male,E8 3PW,London,45000
110,CC-110,Albert,Shiung,Wu,01/08/1968,Single,Male,E8 3PW,London,45000
111,CC-111,Albert,Shiung,Wu,01/08/1968,Single,Male,E8 3PW,London,45000
112,CC-112,Albert,Shiung,Wu,01/08/1968,Single,Male,E8 3PW,London,45000
113,CC-113,Albert,Shiung,Wu,01/08/1968,Single,Male,E8 3PW,London,45000
114,CC-114,Albert,Shiung,Wu,01/08/1968,Single,Male,E8 3PW,London,45000
115,CC-115,Albert,Shiung,Wu,01/08/1968,Single,Male,E8 3PW,London,45000
116,CC-116,Albert,Shiung,Wu,01/08/1968,Single,Male,E8 3PW,London,45000
117,CC-117,Albert,Shiung,Wu,01/08/1968,Single,Male,E8 3PW,London,45000
118,CC-118,Albert,Shiung,Wu,01/08/1968,Single,Male,E8 3PW,London,45000
119,CC-119,Albert,Shiung,Wu,01/08/1968,Single,Male,E8 3PW,London,45000
120,CC-120,Albert,Shiung,Wu,01/08/1968,Single,Male,E8 3PW,London,45000
121,CC-121,Albert,Shiung,Wu,01/08/1968,Single,Male,E8 3PW,London,45000
122,CC-122,Albert,Shiung,Wu,01/08/1968,Single,Male,E8 3PW,London,45000
123,CC-123,Albert,Shiung,Wu,01/08/1968,Single,Male,E8 3PW,London,45000
124,CC-124,Albert,Shiung,Wu,01/08/1968,Single,Male,E8 3PW,London,45000
125,CC-125,Albert,Shiung,Wu,01/08/1968,Single,Male,E8 3PW,London,45000
126,CC-126,Albert,Shiung,Wu,01/08/1968,Single,Male,E8 3PW,London,45000
127,CC-127,Albert,Shiung,Wu,01/08/1968,Single,Male,E8 3PW,London,45000
128,CC-128,Albert,Shiung,Wu,01/08/1968,Single,Male,E8 3PW,London,45000
129,CC-129,Albert,Shiung,Wu,01/08/1968,Single,Male,E8 3PW,London,45000
130,CC-130,Albert,Shiung,Wu,01/08/1968,Single,Male,E8 3PW,London,45000
131,CC-131,Albert,Shiung,Wu,01/08/1968,Single,Male,E8 3PW,London,45000
132,CC-132,Albert,Shiung,Wu,01/08/1968,Single,Male,E8 3PW,London,45000
133,CC-133,Albert,Shiung,Wu,01/08/1968,Single,Male,E8 3PW,London,45000
134,CC-134,Albert,Shiung,Wu,01/08/1968,Single,Male,E8 3PW,London,45000
135,CC-135,Albert,Shiung,Wu,01/08/1968,Single,Male,E8 3PW,London,45000
136,CC-136,Albert,Shiung,Wu,01/08/1968,Single,Male,E8 3PW,London,45000
137,CC-137,Albert,Shiung,Wu,01/08/1968,Single,Male,E8 3PW,London,45000
138,CC-138,Albert,Shiung,Wu,01/08/1968,Single,Male,E8 3PW,London,45000
139,CC-139,Albert,Shiung,Wu,01/08/1968,Single,Male,E8 3PW,London,45000
140,CC-140,Albert,Shiung,Wu,01/08/1968,Single,Male,E8 3PW,London,45000
141,CC-141,Albert,Shiung,Wu,01/08/1968,Single,Male,E8 3PW,London,45000
142,CC-142,Albert,Shiung,Wu,01/08/1968,Single,Male,E8 3PW,London,45000
143,CC-143,Albert,Shiung,Wu,01/08/1968,Single,Male,E8 3PW,London,45000
144,CC-144,Albert,Shiung,Wu,01/08/1968,Single,Male,E8 3PW,London,45000
145,CC-145,Albert,Shiung,Wu,01/08/1968,Single,Male,E8 3PW,London,45000
146,CC-146,Albert,Shiung,Wu,01/08/1968,Single,Male,E8 3PW,London,45000
147,CC-147,Albert,Shiung,Wu,01/08/1968,Single,Male,E8 3PW,London,45000
148,CC-148,Albert,Shiung,Wu,01/08/1968,Single,Male,E8 3PW,London,45000
149,CC-149,Albert,Shiung,Wu,01/08/1968,Single,Male,E8 3PW,London,45000
150,CC-150,Albert,Shiung,Wu,01/08/1968,Single,Male,E8 3PW,London,45000
151,CC-151,Albert,Shiung,Wu,01/08/1968,Single,Male,E8 3PW,London,45000
152,CC-152,Albert,Shiung,Wu,01/08/1968,Single,Male,E8 3PW,London,45000
153,CC-153,Albert,Shiung,Wu,01/08/1968,Single,Male,E8 3PW,London,45000
154,CC-154,Albert,Shiung,Wu,01/08/1968,Single,Male,E8 3PW,London,45000
155,CC-155,Albert,Shiung,Wu,01/08/1968,Single,Male,E8 3PW,London,45000
156,CC-156,Albert,Shiung,Wu,01/08/1968,Single,Male,E8 3PW,London,45000
157,CC-157,Albert,Shiung,Wu,01/08/1968,Single,Male,E8 3PW,London,45000
158,CC-158,Albert,Shiung,Wu,01/08/1968,Single,Male,E8 3PW,London,45000
159,CC-159,Albert,Shiung,Wu,01/08/1968,Single,Male,E8 3PW,London,45000
160,CC-160,Albert,Shiung,Wu,01/08/1968,Single,Male,E8 3PW,London,45000
161,CC-161,Albert,Shiung,Wu,01/08/1968,Single,Male,E8 3PW,London,45000
162,CC-162,Albert,Shiung,Wu,01/08/1968,Single,Male,E8 3PW,London,45000
163,CC-163,Albert,Shiung,Wu,01/08/1968,Single,Male,E8 3PW,London,45000
164,CC-164,Albert,Shiung,Wu,01/08/1968,Single,Male,E8 3PW,London,45000
165,CC-165,Albert,Shiung,Wu,01/08/1968,Single,Male,E8 3PW,London,45000
166,CC-166,Albert,Shiung,Wu,01/08/1968,Single,Male,E8 3PW,London,45000
167,CC-167,Albert,Shiung,Wu,01/08/1968,Single,Male,E8 3PW,London,45000
168,CC-168,Albert,Shiung,Wu,01/08/1968,Single,Male,E8 3PW,London,45000
169,CC-169,Albert,Shiung,Wu,01/08/1968,Single,Male,E8 3PW,London,45000
170,CC-170,Albert,Shiung,Wu,01/08/1968,Single,Male,E8 3PW,London,45000
171,CC-171,Albert,Shiung,Wu,01/08/1968,Single,Male,E8 3PW,London,45000
172,CC-172,Albert,Shiung,Wu,01/08/1968,Single,Male,E8 3PW,London,45000
173,CC-173,Albert,Shiung,Wu,01/08/1968,Single,Male,E8 3PW,London,45000
174,CC-174,Albert,Shiung,Wu,01/08/1968,Single,Male,E8 3PW,London,45000
175,CC-175,Albert,Shiung,Wu,01/08/1968,Single,Male,E8 3PW,London,45000
176,CC-176,Albert,Shiung,Wu,01/08/1968,Single,Male,E8 3PW,London,45000
177,CC-178,Albert,Shiung,Wu,01/08/1968,Single,Male,E8 3PW,London,45000
178,CC-179,Albert,Shiung,Wu,01/08/1968,Single,Male,E8 3PW,London,45000
179,CC-180,Albert,Shiung,Wu,01/08/1968,Single,Male,E8 3PW,London,45000
180,CC-181,Albert,Shiung,Wu,01/08/1968,Single,Male,E8 3PW,London,45000
181,CC-182,Albert,Shiung,Wu,01/08/1968,Single,Male,E8 3PW,London,45000
182,CC-183,Albert,Shiung,Wu,01/08/1968,Single,Male,E8 3PW,London,4500
```

This process takes a few minutes to finish, and the result should be ten fields from the Customer.csv file.

By adding the FOREACH command in the existing pig script, it defines new relations to the specific columns and remove the dump command from the query, it is no longer needed once it's executed.

```
Specific_columns = FOREACH Customer_subset GENERATE CustomerCode,PostCode
,Income;DESCRIBE specific_columns;
```

```
2023-12-20 23:59:38,899 [main] WARN org.apache.pig.backend.hadoop.executionengine.mapReduceLayer.PigHadoopLogger - org.apache.pig.builtin.Utf8StorageConverter{FIELD_DISCARDED_TYPE_CODE=VERSION_FAILED}: Unable to interpret value [73, 110, 99, 111, 109, 101] in field being converted to long, caught NumberFormatException <for input string: "Income"> field discarded
2023-12-20 23:59:39,081 [main] INFO org.apache.hadoop.mapreduce.lib.output.FileOutputCommitter - Saved output of task 'attempt_0001_m_000001_1' to hdfs://sandbox-hdp.hortonworks.com:8020/tmp/_temp/_0000957218/tmp1991750198
2023-12-20 23:59:39,117 [main] WARN org.apache.pig.data.SchemaTupleBackend - SchemaTupleBackend has already been initialized
2023-12-20 23:59:39,143 [main] INFO org.apache.hadoop.mapreduce.lib.input.FileInputFormat - Total input files to process : 1
2023-12-20 23:59:39,143 [main] INFO org.apache.pig.backend.hadoop.executionengine.util.MapRedUtil - Total input paths to process : 1
(,CustomerCode,FirstName,LastName,BirthDate,MaritalStatus,Gender,PostCode,City,)
(1,CC-001,Adriana,Ocampo ,02/01/1977,Single,Male,NB 80V,London,30000)
(2,CC-002,Albert,Einstein ,03/05/1988,Married,Male,NE1 8GH,London,45000)
(3,CC-003,Anna,Behrensmeyer ,04/09/1990,Single,Female,NB 9GV,London,23000)
(4,CC-004,Blaise,Pascal ,08/05/1997,Married,Male,SE8 9GV,London,34000)
(5,CC-005,Caroline,Herschel ,05/06/1990,Single,Female,N6 9YH,London,35000)
(6,CC-006,Cecilia,Payne-Gaposchkin ,03/07/1958,Married,Female,EC3 0VU,London,23000)
(7,CC-007,Chien-Shiung,Wu ,01/08/1968,Single,Male,E8 9YU,London,45000)
(8,CC-008,Dorothy,Hodgkin ,09/09/1957,Married,Female,M1 9UH,Manchester,50000)
(9,CC-009,Edmond,Halley ,03/10/1973,Single,Male,M8 9J3,Manchester,23000)
2023-12-20 23:59:39,263 [main] INFO org.apache.pig.Main - Pig script completed in 4 seconds and 626 milliseconds (4626 ms)
[maria_dev@sandbox-hdp ~]$ vi Customer
[maria_dev@sandbox-hdp ~]$ pig -f Customer
WARNING: HADOOP_HOME has been replaced by HADOOP_HOME. Using value of HADOOP_PREFIX.
23/12/21 00:15:20 INFO pig.ExecTypeProvider: Trying ExecType : LOCAL
23/12/21 00:15:20 INFO pig.ExecTypeProvider: Trying ExecType : MAPREDUCE
23/12/21 00:15:20 INFO pig.ExecTypeProvider: Trying ExecType : TEZ_LOCAL
23/12/21 00:15:20 INFO pig.ExecTypeProvider: Trying ExecType : TEZ
23/12/21 00:15:20 INFO pig.ExecTypeProvider: Picked TEZ as the ExecType
2023-12-21 00:15:20,118 [main] INFO org.apache.pig.Main - Apache Pig version 0.16.0.3.0.1.0-187 (rUnVersioned directory) compiled Sep 19 2018, 10:13:33
2023-12-21 00:15:20,118 [main] INFO org.apache.pig.Main - Logging error messages to: /home/maria_dev/pig_1703117720116.log
2023-12-21 00:15:20,613 [main] INFO org.apache.pig.impl.util.Utils - Default bootstrap file /home/maria_dev/.pigbootup not found
2023-12-21 00:15:20,613 [main] INFO org.apache.pig.backend.hadoop.executionengine.HExecutionEngine - Connecting to hadoop file system at: hdfs://sandbox-hdp.hortonworks.com:8020
2023-12-21 00:15:21,291 [main] INFO org.apache.pig.PigServer - Pig Script ID for the session: PIG-Customer-47942551-f004-abf6-ad88-f2d2c47c3709
2023-12-21 00:15:21,291 [main] WARN org.apache.pig.PigServer - ATS is disabled since yarn.timeline-service.enabled set to false
Customer: {CustomerKey: int,CustomerCode: chararray,FirstName: chararray,LastName: chararray,Birthdate: chararray,MaritalStatus: chararray,Gender: chararray,PostCode: chararray,City: chararray,Income: long}
specific_columns: {CustomerCode: chararray,PostCode: chararray,Income: long}
2023-12-21 00:15:22,172 [main] INFO org.apache.pig.Main - Pig script completed in 2 seconds and 180 milliseconds (2180 ms)
[maria_dev@sandbox-hdp ~]$
```

Fig: Successfully Executed Specific Columns from Customer Relation

To fulfill the requirements outlined below, store the data relation into HDFS as a file by using the STORE specific_columns INTO 'output/specific_columns using PigStorage(',')' command in pig script, then save and execute it.

```

← → ⌂ ⌂ localhost:4200
File Name: Customer
Started At: 2023-12-21 00:40:05
Finished At: 2023-12-21 00:41:01
Features: LIMIT

Success!

DAG 0:
    Name: PigLatin:Customer-0_scope-0
    ApplicationId: job_1703103117930_0002
    Total Launched Tasks: 2
        File Bytes Read: 250
        File Bytes Written: 250
        HDFS Bytes Read: 1590
        HDFS Bytes Written: 207
    Spillable Memory Manager spill count: 0
        Bags proactively spilled: 0
        Records proactively spilled: 0

DAG Plan:
Tez vertex scope-20 -> Tez vertex scope-22,
Tez vertex scope-22

Vertex Stats:
Vertex Parallelism Total Tasks Input Records Reduce Input Records Output Records File Bytes Read File Bytes Written HDFS Bytes Read HDFS Bytes Written Alias Feature Outputs
scope-20 1 1 10 0 10 0 250 1590 0 Customer,Customer_subset
scope-22 1 1 10 0 10 0 0 0 207 Customer,Customer_subset,specific_col
ums LIMIT hdfs://sandbox-hdp.hortonworks.com:8020/user/maria_dev/output/specify_columns

Input(s):
Successfully read 10 records (1590 bytes) from: "/tmp/Mydata/customer.csv"

Output(s):
Successfully stored 10 records (207 bytes) in: "hdfs://sandbox-hdp.hortonworks.com:8020/user/maria_dev/output/specify_columns"

2023-12-21 00:41:01,062 [main] INFO org.apache.pig.Main - Pig script completed in 58 seconds and 30 milliseconds (58030 ms)
2023-12-21 00:41:01,062 [main] INFO org.apache.pig.backend.hadoop.executionengine.tez.TezLauncher - Shutting down thread pool
2023-12-21 00:41:01,079 [shutdown-hook-0] INFO org.apache.pig.backend.hadoop.executionengine.tez.TezSessionManager - Shutting down Tez session org.apache.tez.client@735d6d8e

```

Fig: Storing relationship data in HDFS

Once it has been executed, we have to go to HDFS Files View to view the output, which is saved under the user as maria_dev. We need to check for a newly created folder named "Output"

Name	Size	Last Modified	Owner	Group	Permission	Erasure Coding	Encrypted
app-logs	--	2023-12-21 00:40	yarn	hadoop	drwxrwxrwx	No	
apps	--	2018-11-29 19:01	hdfs	hdfs	drwxr-xr-x	No	
ats	--	2018-11-29 17:25	yarn	hadoop	drwxr-xr-x	No	
atsv2	--	2018-11-29 17:26	hdfs	hdfs	drwxr-xr-x	No	
hdp	--	2018-11-29 17:26	hdfs	hdfs	drwxr-xr-x	No	
ivy2-recovery	--	2018-11-29 17:55	ivy	hdfs	drwx-----	No	
mapred	--	2018-11-29 17:26	mapred	hdfs	drwxr-xr-x	No	
mr-history	--	2018-11-29 17:26	mapred	hadoop	drwxrwxrwx	No	
ranger	--	2018-11-29 18:54	hdfs	hdfs	drwxr-xr-x	No	
spark2-history	--	2023-12-21 00:44	spark	hadoop	drwxrwxrwx	No	
tmp	--	2023-12-21 00:41	hdfs	hdfs	drwxrwxrwx	No	
user	--	2018-11-29 19:21	hdfs	hdfs	drwxr-xr-x	No	
warehouse	--	2018-11-29 17:51	hdfs	hdfs	drwxr-xr-x	No	

Fig: Navigating the “User” folder in the HDFS file view

Now click on the **Output** folder to see the sub folder named as Specific_columns.

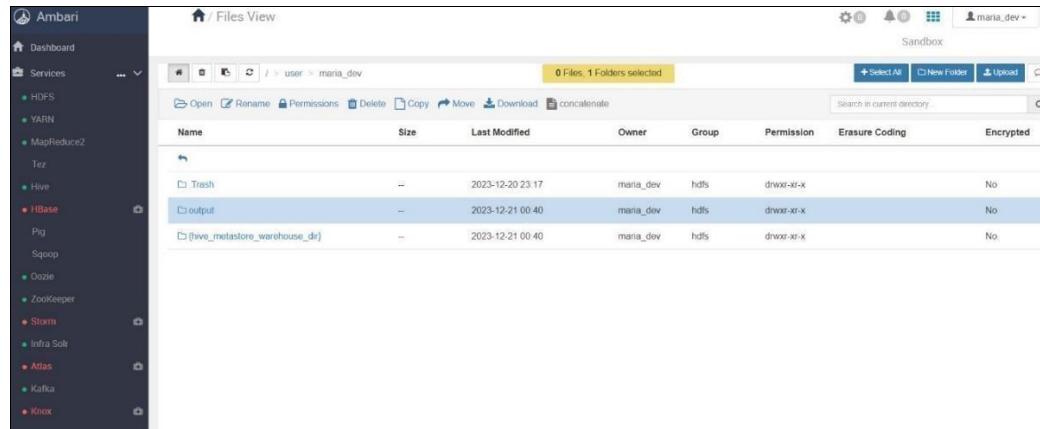


Fig: Output for Customer folder created in HDFS

By clicking the Output folder, we can see the specific _columns folder. The output is named “part-v001-o000-r-00000” in that folder.

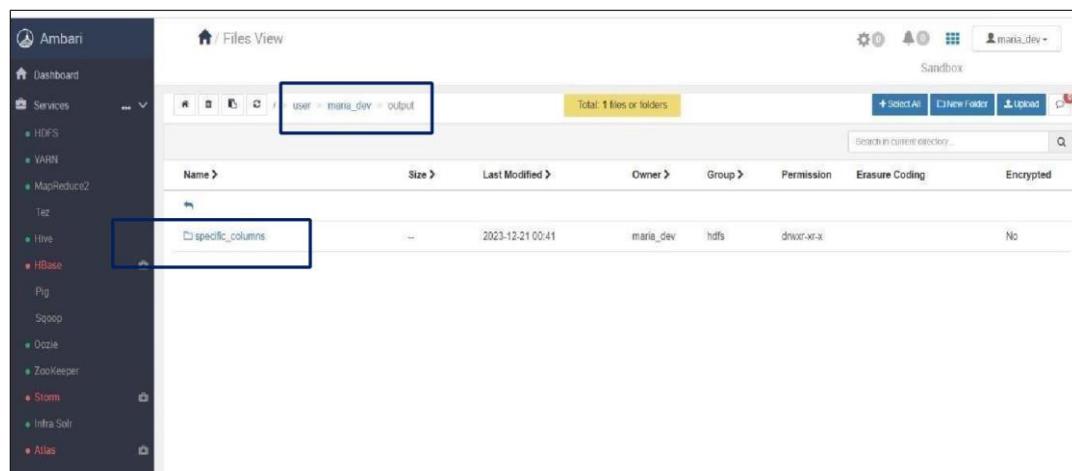
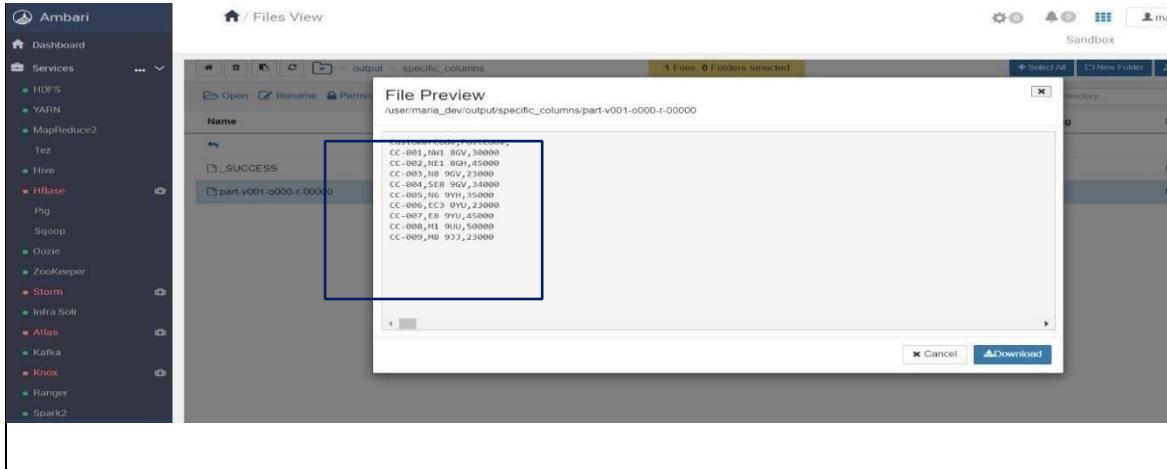


Fig: Specific Column created in HDFS

Once we click on the “part-v001-o000-r-00000” file and then Open we can see output by clicking on Open.



Conclusion:

In conclusion, the Blue-Sky Online Retailer project provided me with great insights into the field of big-data processing using technologies such as SSMS, HDFS, Pig, and Hive in the context of Apache Hadoop. From a technical point of view, I have developed an in-depth understanding of Pig and Hive's features, knowing that they are feature-rich and capable of completing tasks in the big data ecosystem.

One important takeaway from this case study is the flexibility and choice provided by these technologies. The fact that organizations may choose from Pig and Hive based on their individual needs and preferences demonstrates Apache Hadoop's versatility in fulfilling a wide range of objectives. Having the ability to choose between tools improves the simplicity and adaptability of big data processing by allowing businesses to modify their strategy to their specific data processing challenges.

The hands-on experience that I got from this project in uploading data into HDFS files and reviewing data using various Hive queries has helped me to develop my abilities. The usage of queries like CREATE, LOAD, INSERT, SELECT, FROM, and GROUP has given me a good foundation in handling and querying large data sets. The actual implementation of these queries in creating and handling tables for customer and payment data has deepened my understanding of real-world scenarios in big data processing.

Furthermore, working on tasks such as creating tables for customers as well as extracting data from databases, and sorting outcomes based on specified criteria has provided me with a practical understanding of how these technologies can be used to extract useful insights from enormous datasets. The project not only taught me technical skills, but also highlighted the value of proper data management and extraction for making accurate decisions.

Reflecting on my experience working in an agile setting, I have learned to appreciate the iterative and collaborative aspects of agile approach phases. The courses required adaptation, rapid problem-solving, and effective team communication. These abilities are not only important in the team, but they are also highly transferable to real-world professional circumstances where flexibility and teamwork are essential for success.

Overall, by doing this project I have gained a rewarding experience, which taught me technical skills, practical insights, and a deeper understanding of the flexible work environment. The information learned will surely help me manage and contribute effectively to the rapidly changing world of big data and analytics.