# Northeastern University

College of Engineering

*INFO 7290: Data Warehousing and Business Intelligence*

## Sales Data Warehouse

## Group 8

JayaChandra Chimakurthi (001065324)

Vyshnavi Bhyravajosula (001085050)

Kirtikumar Waykos (001029199)

Pradeep Poonati (001316491)

# Index

## Revision History:

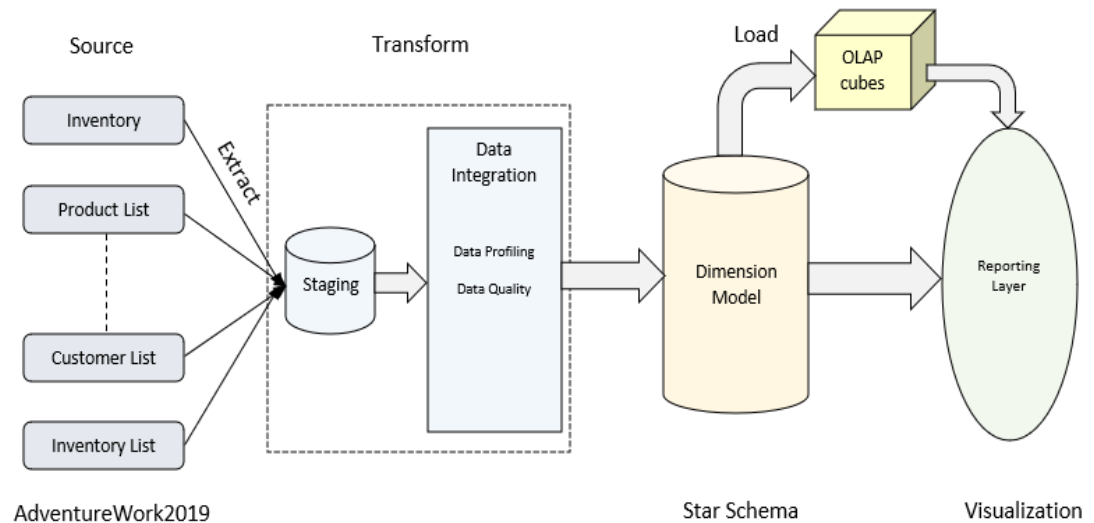| Date | Version | Changes | Author |
|---|---|---|---|
| 03/27/2021 | 1.0 | Initial Document | All members |
| 04/06/2021 | 1.1 | Added pipelines for managing SCD in the tables. Created facts and dimensions for all the tables. Built Dimension Model Chart. Added Professor's comment as appendix | All members |
| 04/18/2021 | 1.2 | Added how the data is getting populated to fact tables, expanded the business scenarios. Explained how the data is analyzed and insights were drawn to business users using OLAP cubes. | All members |
| 04/20/2021 | 1.3 | Created and populated data into dimensions. Used required transformations in SSIS process. | All members |
| 04/21/2021 | 1.4 | Added SCD's and populated fact tables with data. | All members |
| 04/22/2021 | 1.5 | Created OLAP cubes and extracted data into Excel. | All members |
| 04/26/2021 | 1.6 | Analysis in PowerBI | All members |
| 04/27/2021 | 1.7 | Documentation and Presentation | All members |

# 1. Objective

We are using AdventureWorks19 database to understand and analyze the sales by monitoring costs, discounts and selling prices. We are identifying factors that influence sales and analyze the Sales based on different quarters and periods. We understand the sales of products based on their geographical location and analyze the list price and standard cost based on the product category.

# 2. Data Description

Adventure Work is a fictional company that manufactures the bicycles and other related accessories products such as Helmet, Water bottle, Clothing, Frames, and other premium services. It distributes its products primarily through retail stores located across multiple countries and directly to individual customers through the internet.

AdventureWorks database is meant to be a one-stop solution for all data-related needs of the organization, and with time, more data will be added to the Database. Our process for warehousing and analysis will concern only with the Products, Sales, Customers, along with geographical data.
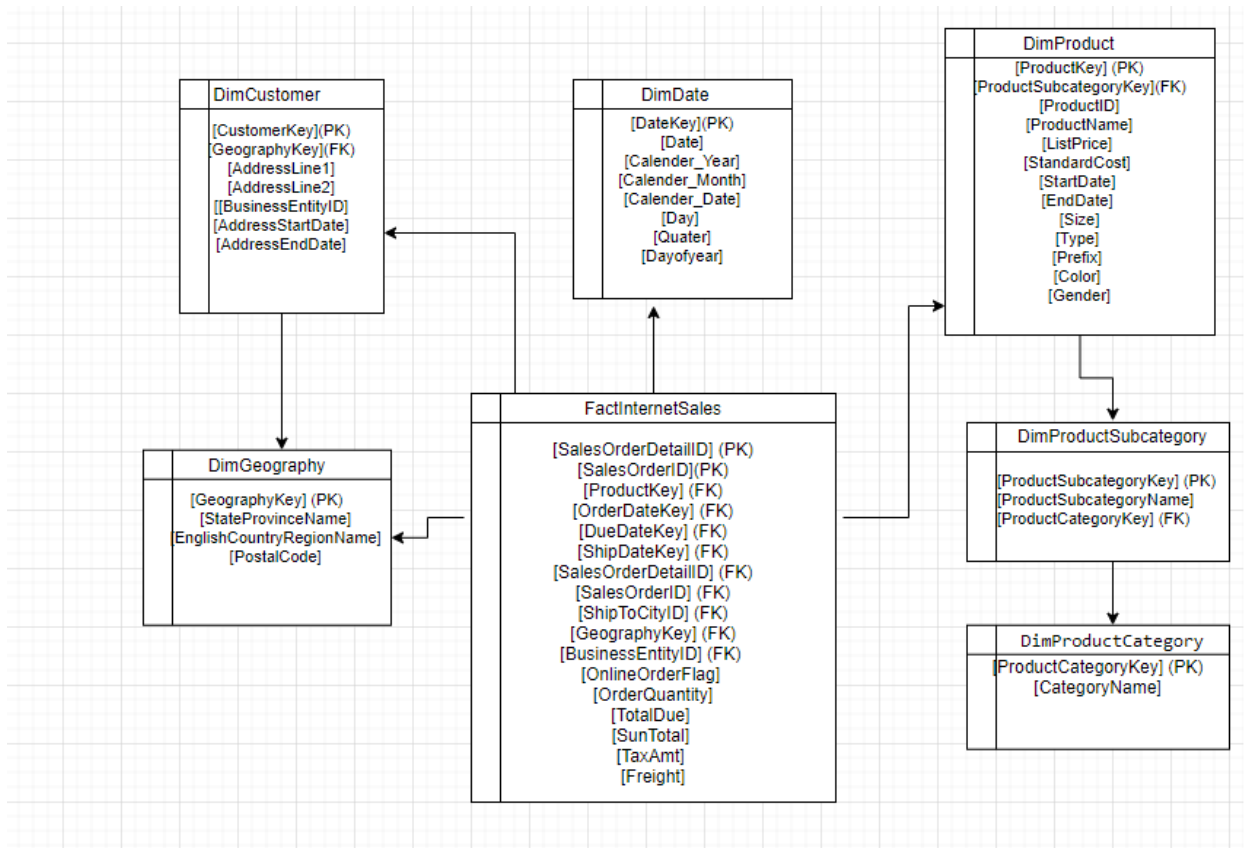
## 3. Business Intelligence Dataflow Architecture



## 4. Workflow:

- Data is extracted from the AdventureWork2019 database. We extracted tables according to business purpose and load the data into the centralized platform, i.e., staging area. We do not make any changes to data while loading into the Staging area.

- We apply transformations where required; we also perform integration tasks by applying business rules.

- After the integration process, data is then loaded into the dimensions and facts.

- With SQL Server Analysis Services, we build OLAP cubes and create hierarchies. We deploy the cubes and extract the data into excel files by creating offline cubes to analyze the relations between Product, Sales, Customers by different geographic territories and other factors.

- The final dimensional model is visualized using PowerBI to justify the use of DataWarehouse and answer various questions through visual techniques.

# 5. Dimension Model:



We built a snowflake schema, which is used to represent multidimensional data.

## 5.1 Dimensional Model Description

We have worked on one Fact table, i.e, FactInternetSales for analysis and the other fact table (FactInternetSalesReason) did not meet our business scenarios.

- FactInternetSales
  Fact internetsales table consists of 15 columns with one primary key: SalesOrderDetailID. The FactInternetSales table has many to one relation with the dimension tables.
- FactInternetSalesReason:
  FactInternetSalesReason table consists of 3 columns with three primary keys columns; SalesOrderNumber, SalesOrderLineNumber and SalesReasonKey.

In our dimensional model, the dimensions include:

DimProduct, DimDate, DimCustomer, DimGeography, DimProductSubCategory, DimProductCategory

- DimProduct:
  This has 13 columns with one primary key (ProductKey), and DimProductSubCategory reference the DimProduct with the help of ProductSubCategoryKey.
- DimDate:
  This has 8 columns with primary key (DateKey), it deals with the different dates that are present in various tables i.e, Ship Date, Order Date and Due Date.
- DimCustomer: This has 7 columns with one primary key (CustomerKey) and this table has complete details of the customer and has the reference to DimGeography using GeographyKey
- DimGeography has six columns with Geographykey as primary key this section tells about the geographical categories like city ,state, region zip code etc.
- DimProductCategory has two rows with productcategorykey as primary key, it gives information on product categories like Bikes,Components etc.
- Dimproductsubcategory has three columns with productsubcategorykey as primary key and productCategory key as foreign key, this section gives details about subcategories of the product.

## 5.2 First time load

Until loading files, it's critical to understand what kind of cleaning is needed. We need to double-check the rows after they've been loaded. After the data has been loaded, we would search to see how the rows in the original data source match the rows in the destination.
Staging:
• Raw data is loaded into the staging area from different tables, errors are handled, and we set referential integrity constraints as required.
• Each data flow task undergoes a truncation stage to avoid data loading duplication, even though the data flow task is run several times.
• Current data is loaded followed by an incremental load data from the AdventureWorks Database.

## 5.3 How do Dimensions and Facts get populated with data?

Dimension tables are used as a collection of reference information about a business. In a snowflake schema model, dimension tables offer descriptive characteristics of the fact table with the help of their attributes.

Fact tables are loaded from the transaction tables such as Production.Product by mapping the fact columns.

The columns of the transaction tables become measures and the primary keys on the transaction table such as Productkey, Customerkey becomes the dimensional key columns of the fact table when loading the data from OLTP transaction table to Fact table.

For our analysis we are storing the details about the Products (Product, ProductCategory, ProductSubcategory) along with the Sales details with respect to the products. We are also creating and modifying the Date Columns according to our business scenarios.

The dimensions are populated with the following source tables from AdventureWork2019

a. DimProduct
[Production].[Product]
[Production].[ProductCategory],
[Production].[ProductSubcategory],
[Production].[ProductModel]
b. DimCustomer
[Person].[Person]
[Person].[PersonPhone]
[Person].[EmailAddress]
[Person].[Address]
c. DimGeography
[Person].[Address]
[Person].[StateProvince]
[Person].[CountryRegion]
[person].[AddressType]
d. DimDate
[sales].[SalesOrderHeader]
e. DimProductCategory
[production].[ProductCategory]
f. DimProductSubCategory
[production].[ProductSubcategory]

# 6. Slowly Changing Dimensions (SCD)

Slowly Changing Dimensions are used to capture the changing data within the dimension over time.  In this project we are implementing Type 2 SCD for Customer and Product dimension. We are adding a new record when the change occurs on the dimension and updating the old record as inactive with the help of enddate attribute.

This helps the business to do the historical reporting purposes as the change is captured in the new record.

- AddressLine1 attribute on the DimCustomer dimension.

If Customer updates Address, then we track changes with AddressStartDate and AddressEndDate. (SCD Type2).

- StandardCost attribute on DimProduct
  If StandardCost is updated, then we track changes with CostStartDate and CostEndDate. (SCD Type2).

## 7. Fact Tables and Dimension Tables Design- SSIS

**Control Flow:**



**Execute SQL Task:**

We are cleaning our staging area for every process.

**Staging Area:**

We are extracting data from Source table from AdventureWorks2019 and loading into the staging area



Following workflows shows tables being populated after performing cleansing and transformations:

1. **DimProduct Category and DimDate**

   We are populating DimProductCategory and DimDate(for ShipDate) by using a multicast, checking the nulls and removing nulls. We then further remove the duplicates and do data conversions as required and add lookups to get no match outputs only.

## 2. DimSubCategory

We are populating DimSubCategory from the staging area. We are performing the required transformations and adding a lookup to get the foreign keys of 'ProductCategoryKey' from the DimProductCategory table.

### 3. DimCustomer and DimDate:

We have populated DimCustomer and DimDate(for DueDate) from the staging. We used a multicast, checking the nulls and removing nulls. We then further remove the duplicates and do data conversions as required and add lookups to get no match outputs only. We have used SCD on AddressLine1 to track changes when a update occurs.

4. **Geography and DimDate**

We are populating DimProductCategory and DimDate(for OrderDate) by using a multicast, checking the nulls and removing nulls. We then further remove the duplicates and do data conversions as required and add lookups to get no match outputs only.

### 5. DimProduct:

We are populating DimProduct from the staging. We are performing the required transformations and adding a lookup to get the foreign keys of 'ProductSubCategoryKey' from the DimProductSubCategory table. We are adding an SCD to monitor any changes occurring in StandardCost.

## 1. Fact table:

## 8. Cube Design-SSAS:

Online analytical processing cubes are used to perform multidimensional analysis at high speeds on large volumes of data from a data warehouse. We created two hierarchies for Dimgeography and Dimdate to perform analysis and extract insights at a higher granular level.

**OLAP CUBE:**

## Hierarchy of DimDate:

## Hierarchy of DimGeography

Creating Offline files using Excel:

# 9. Analysis and Visualizations Using PowerBI:

1. Tax Amount by Country:



The above pie chart Company has paid highest amount of tax in USA and Australia. Both sharing total 62% of company's total tax payment in between. While other countries account for less than 10% of tax.

2. Sales by Year and Country:



The highest sales are seen by United States and Australia. Company sales are less in European countries. Hence, it is suggested that company should maintain market in these countries. However, Company also needs to focus on increasing it's sales in European Countries.

3. Standard cost versus List price by Product



While some Product category has higher than average Standard cost, Mountain Bikes can be sold for higher Margin. Road bike - 250 and Road bike-150 have higher Margin as Mountain Bike.

4. Decomposition of Sales



Decomposition of Sales

In USA, Q4 of fiscal year 2013 has highest company had highest sales in November. Possibly because of the Holiday season occurring in November and December.

5. Top 10 States by Total due

| EnglishCountryRegionName | StateProvinceName | TotalDue ▼ | Sales Total | TaxAmt | Freight | OrderQuantity |
|---|---|---|---|---|---|---|
| United States | California | 6,174,983.32 | 5,588,220.11 | 447,057.62 | 139,705.58 | 2995 |
| Australia | New South Wales | 4,295,807.30 | 3,887,608.36 | 311,008.68 | 97,190.27 | 1944 |
| United Kingdom | England | 3,685,839.63 | 3,335,601.43 | 266,848.12 | 83,390.08 | 1858 |
| United States | Washington | 2,657,470.12 | 2,404,950.30 | 192,396.03 | 60,123.79 | 1310 |
| Australia | Victoria | 2,488,573.62 | 2,252,102.78 | 180,168.23 | 56,302.60 | 1117 |
| Australia | Queensland | 2,171,457.79 | 1,965,120.14 | 157,209.62 | 49,128.03 | 974 |
| Canada | British Columbia | 2,015,845.46 | 1,824,294.51 | 145,943.56 | 45,607.39 | 914 |
| United States | Oregon | 1,253,017.74 | 1,133,952.69 | 90,716.22 | 28,348.83 | 594 |
| Germany | Saarland | 785,722.99 | 711,061.52 | 56,884.92 | 17,776.55 | 395 |
| Germany | Hamburg | 765,614.28 | 692,863.59 | 55,429.09 | 17,321.60 | 388 |
| **Total** | | **26,294,332.24** | **23,795,775.43** | **1,903,662.08** | **594,894.73** | **12489** |

Company must receive the above amount from the respective Countries. It also shows the Tax amount and Order Quantity in those countries. We can see highest selling regions in each country.

# 10. Software and Tools Used:

- Microsoft SQL Server Management Studio
- ER/Studio
- Microsoft Visual Studio 2017(SSDT)
- PowerBI
- Excel

*References*

- https://www.mssqltips.com/
- https://www.sqlshack.com/

Professors Comments (03/28/2021):

- Intro paragraph telling me what the project / document is for. You have some of it later in the document but you s
- As the document gets large you would want to consider a table of contents.
- I don't understand what the Timeline section is for
- Data pipeline looks good.  I'd like a full list of tables…
- You only show me one fact table and you are strictly reverse engineering the DW dbase, you can make your life easier but trimming it back.
- You will need a real model not just a reversed engineered one. They have a number of columns.  I need to know the full list -- Analyse the business queries for our model and find redundant columns.
- The business scenarios are good to list
- Softwares typically isn't used as a plural-done
- I see no mention of SCDs

-

Professors Comments (04/11/2021):

- Pipeline is a good level of detail
- Explain how the sources will also populate the facts
- Telling me the number of columns doesn't tell me anything… tell me what you are storing
- SCD descriptions are good
- Add a modification history to the document
- Expand the business scenarios some more
- Explain your olap model
- Describe the visualizations