

Index

1. Introduction
2. Business queries
3. Tables
4. Starnet schema and the diagram
5. Concept hierarchy
6. Starnet foot prints for business queries
7. Starnet database diagram (ER diagram)
8. ETL process
9. Process of creating the data cube
10. Data cube hierarchy in Visual Studio
11. SQL Analysis service in Power BI
12. Analysis Report screenshot using live connection
13. Final Analysis Report

1. Introduction

The US Adult Income data gives us the information about the income of people depending on various factors such as education, occupation, age, location, race etc. The dataset has the data of 32561 people. This consists the anonymous information of people without their names. The dataset has people from different countries however most of the people has their native as United States.

Tasks to do

Major task:

Create a Data Warehouse for this project

1. Business questions that the data warehouse could help answer.
2. Draw Starnet with the aim to identify the dimensions and concept hierarchies for each dimension. This should be based on the lowest level information you have access to.
3. Use the Starnet footprints to illustrate how the business queries can be answered with the design.
4. Implement a star or snowflake schema using SQL Server Management Studio (SSMS).
5. Load the data from the csv files to populate the tables.
6. Use SQL Server Data Tools to build a multi-dimensional analysis service solution, with a cube designed to answer your business queries. Make sure the concept hierarchies match your Starnet design. Paste the cube diagram to your Power BI Dashboard.
7. Use Power BI to visualize the data returned from your business queries.

2. Business Queries

1. Finding number of people having income maximum 50000 USD per annum also the number of people who have income more than 50000 USD a year from US, India, China, Germany and Mexico of all age group.
2. Finding the number of young (18-30) people having income more than 50000 USD a year from United States.
3. Finding total young (age 18-30) private employees having earnings more than USD 50000 a year from Mexico and Canada.
4. Finding the number of people having income more than 50000 USD per year and who works 10-30 hours a week, from United States.
5. Finding the preferred education of people who are earning more than 50000 USD a year and along with that people who are getting maximum 50000 USD a year working for Federal-govt from United-States.
6. Getting all the males and females working as a sales representative, Exce-managerial and a Machine-op-inspector having maximum income USD 50000 a year from US, Canada, Cuba having race white.

3. Tables

a. Fact Table :

- i. **Fact_user_info:** This table here in this project is our fact table.

This table consists measures such as user_count, age, working hour, final weight. This table will have the foreign keys of the dimension table in order make the relationship with the dimension tables for this project .

	Results												
	Messages												
	user_id_fact	income	age	finalWeight	totalworkHrs	location_dim_id	race_dim_id	occupation_dim_id	sector_dim_id	gender_dim_id	relation_dim_id	education_dim_id	
1	1	<=50K	39	77516	40	1	3	11	4	1	1	1	
2	2	<=50K	50	83311	13	1	3	2	5	1	2	1	
3	3	<=50K	38	215646	40	1	3	12	9	1	3	2	
4	4	<=50K	53	234721	40	1	5	12	9	1	2	3	
5	5	<=50K	28	338409	40	2	5	13	9	2	2	1	
6	6	<=50K	37	284582	40	1	3	2	9	2	2	4	
7	7	<=50K	49	160187	16	3	5	7	9	2	4	5	
8	8	>50K	52	209642	45	1	3	2	5	1	2	2	

b. Dimension Tables:

- i. **LocationDim:** This table will consist the countries of the users having unique ID for each country.

A screenshot of a database query results window titled "Results". The table has two columns: "locationDim" and "native_country". The data shows 10 rows of user native countries with their corresponding IDs.

	locationDim	native_country
1	1	United-States
2	2	Cuba
3	3	Jamaica
4	4	India
5	5	Other
6	6	Mexico
7	7	South
8	8	Puerto-Rico
9	9	Honduras
10	10	England

- ii. **RaceDim:** This dimension table will have the race such as Black, White ,Asia-Pacific etc.

A screenshot of a database query results window titled "Results". The table has two columns: "dim_race_id" and "race". The data shows 5 rows of user races with their corresponding IDs.

	dim_race_id	race
1	1	Asian-Pac-Islander
2	2	Other
3	3	White
4	4	Amer-Indian-Eskimo
5	5	Black

- iii. **OccupationDim:** This table will have the occupation for all the user with a unique primary key. Example : Machine-op-incpt

A screenshot of a database query results window titled "Results". The table has two columns: "dim_occupation_id" and "occupation". The data shows 4 rows of user occupations with their corresponding IDs.

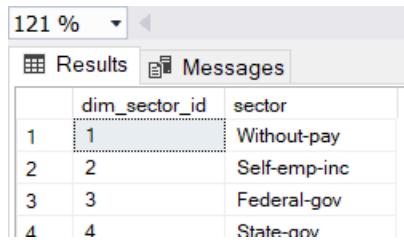
	dim_occupation_id	occupation
1	1	Machine-op-inspct
2	2	Exec-managerial
3	3	Transport-moving
4	4	Farmmn-fishing

- iv. **EducationDim:** This dimension table education will have the education qualification along with the education number of a particular user.

A screenshot of a database query results window titled "Results". The table has three columns: "educationDim_id", "edu_level", and "edu_num". The data shows 4 rows of user education levels with their corresponding IDs and numbers.

	educationDim_id	edu_level	edu_num
1	1	Bachelors	13
2	2	HS-grad	9
3	3	11th	7
4	4	Masters	14

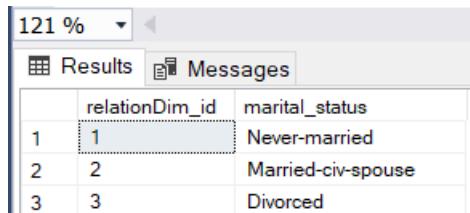
- v. **SectorDim:** This table holds the information of different sectors where people work or have a business of their own.



A screenshot of a Microsoft SQL Server Management Studio (SSMS) results window titled "Results". The table has two columns: "dim_sector_id" and "sector". The data is as follows:

	dim_sector_id	sector
1	1	Without-pay
2	2	Self-emp-inc
3	3	Federal-gov
4	4	State-nov

- vi. **RelationDim:** Relation dimensional tables give us the data of marital status of people.
Example : Married, Never-married, Divorced.



A screenshot of a Microsoft SQL Server Management Studio (SSMS) results window titled "Results". The table has two columns: "relationDim_id" and "marital_status". The data is as follows:

	relationDim_id	marital_status
1	1	Never-married
2	2	Married-civ-spouse
3	3	Divorced

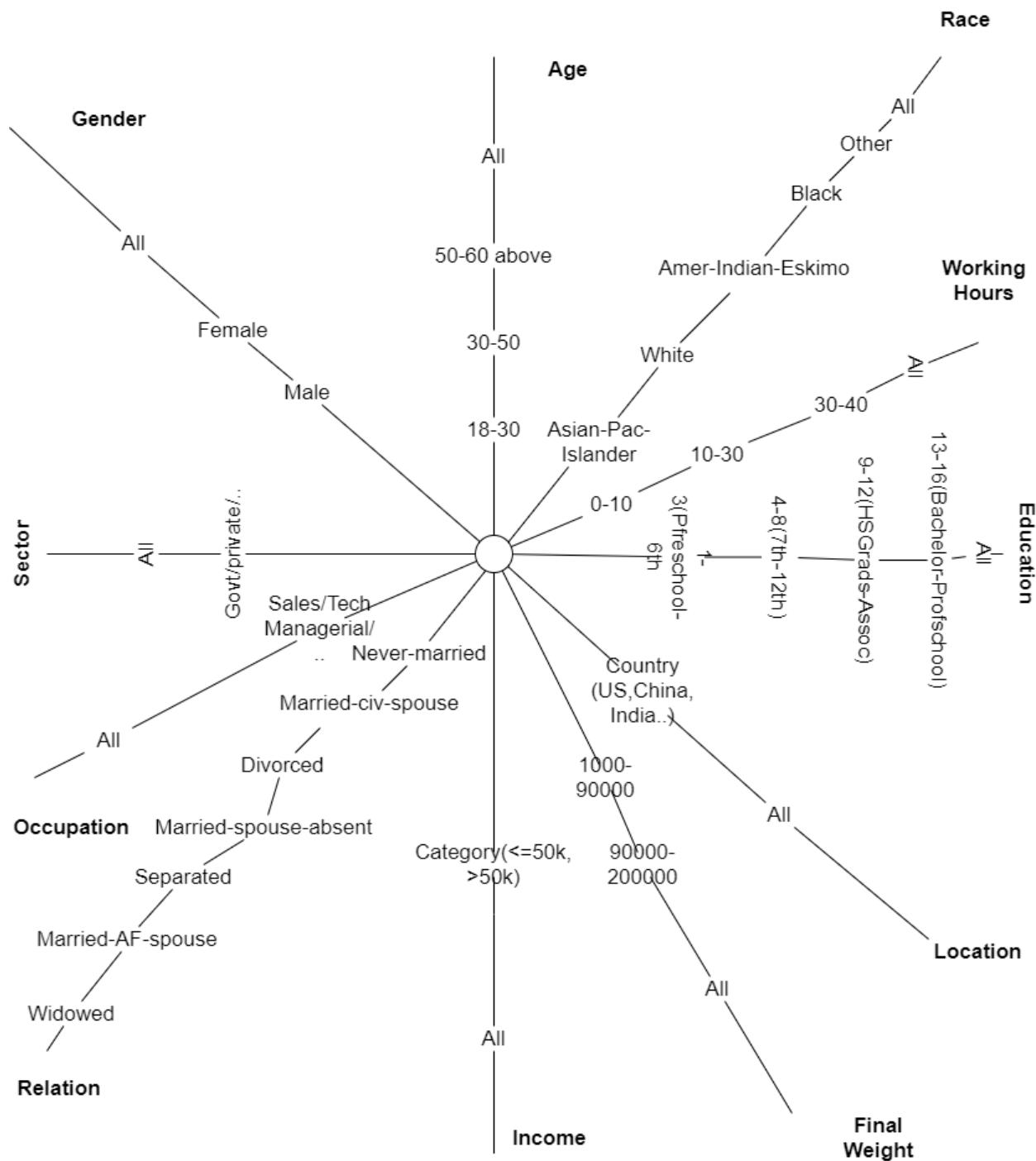
- vii. **GenderDim:** This dimension table holds the gender of people. Example : Male, female.



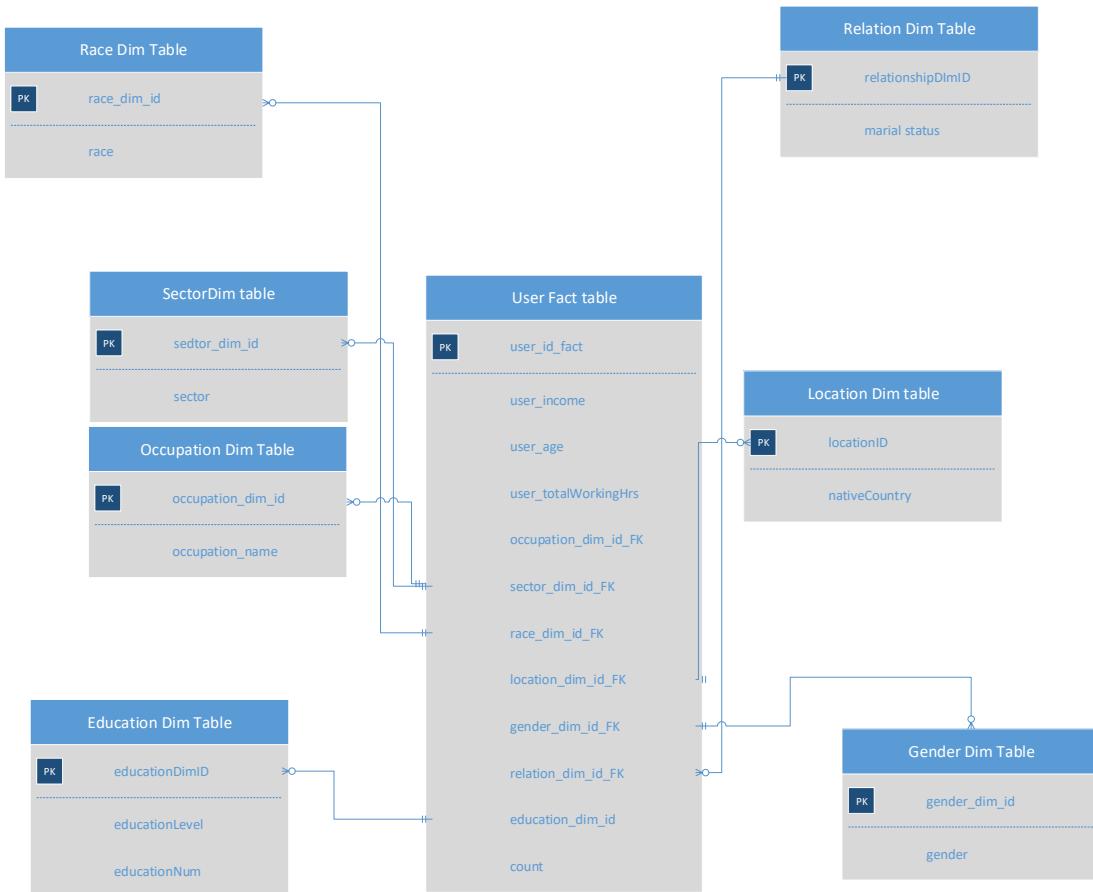
A screenshot of a Microsoft SQL Server Management Studio (SSMS) results window titled "Results". The table has two columns: "dim_gender_id" and "gender". The data is as follows:

	dim_gender_id	gender
1	1	Male
2	2	Female

4. Starnet schema and the diagram



Startnet Diagram

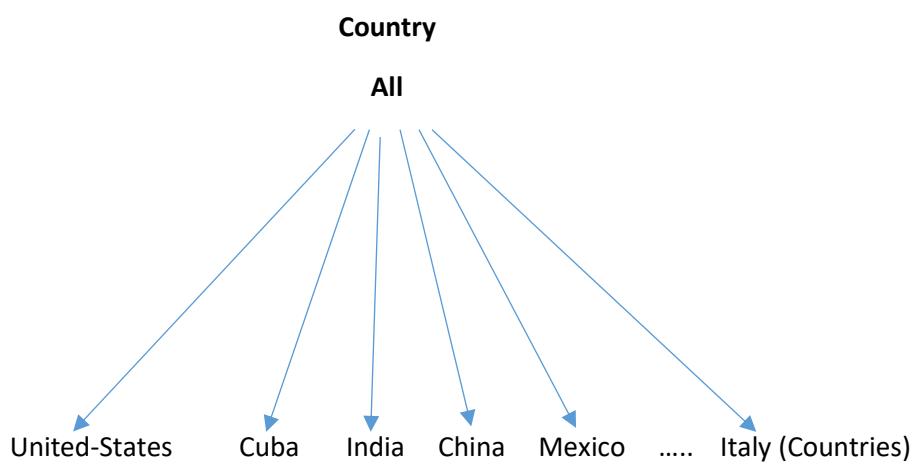
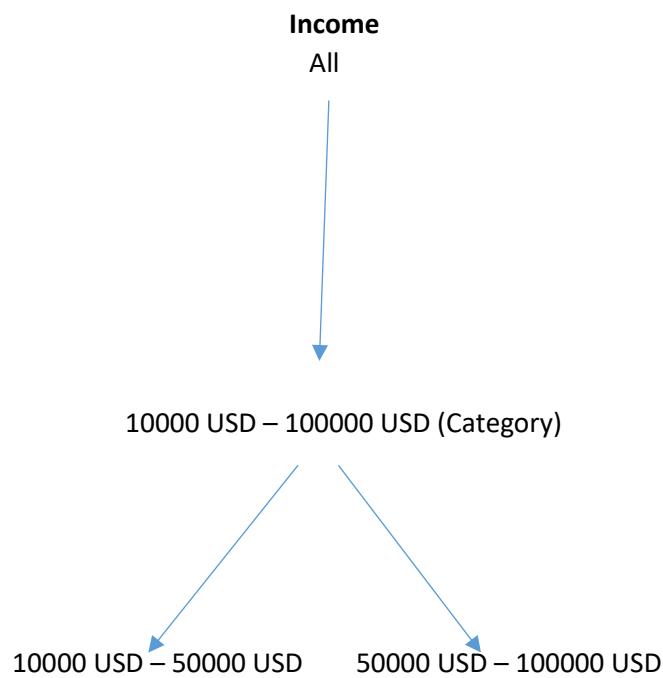


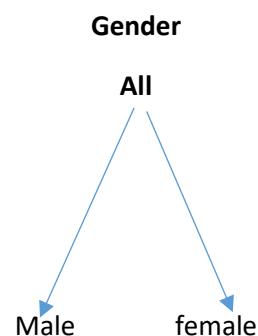
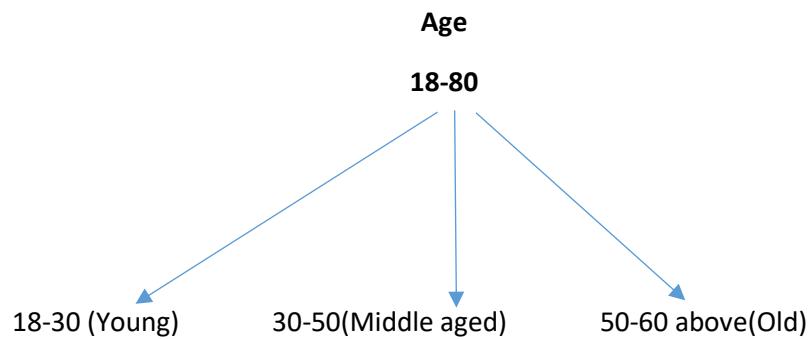
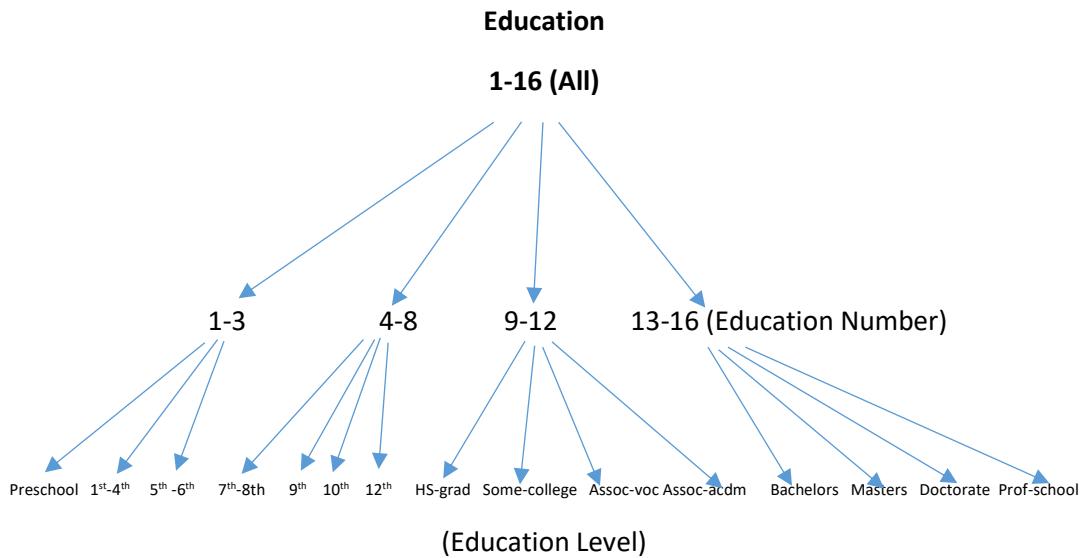
Numeric Measures of fact tables:

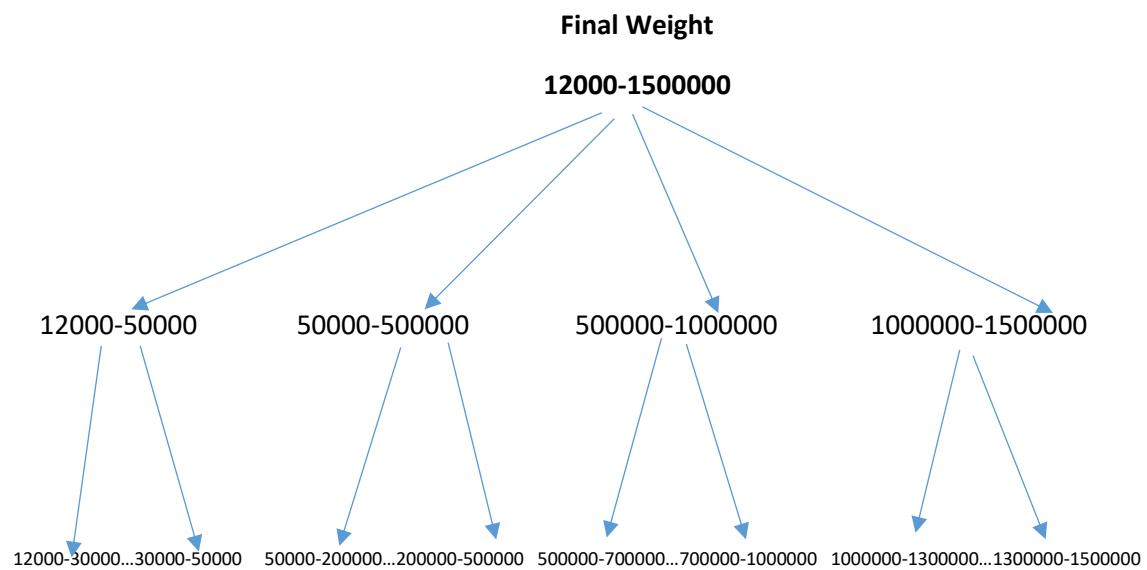
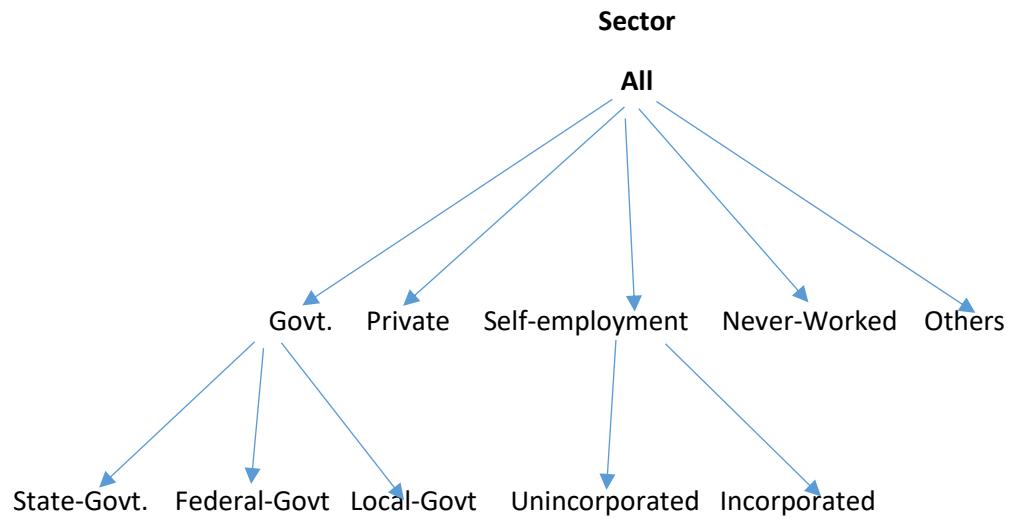
2. Age
3. Final Weight
4. Total working hours
5. User fact ID count

5. Concept hierarchy

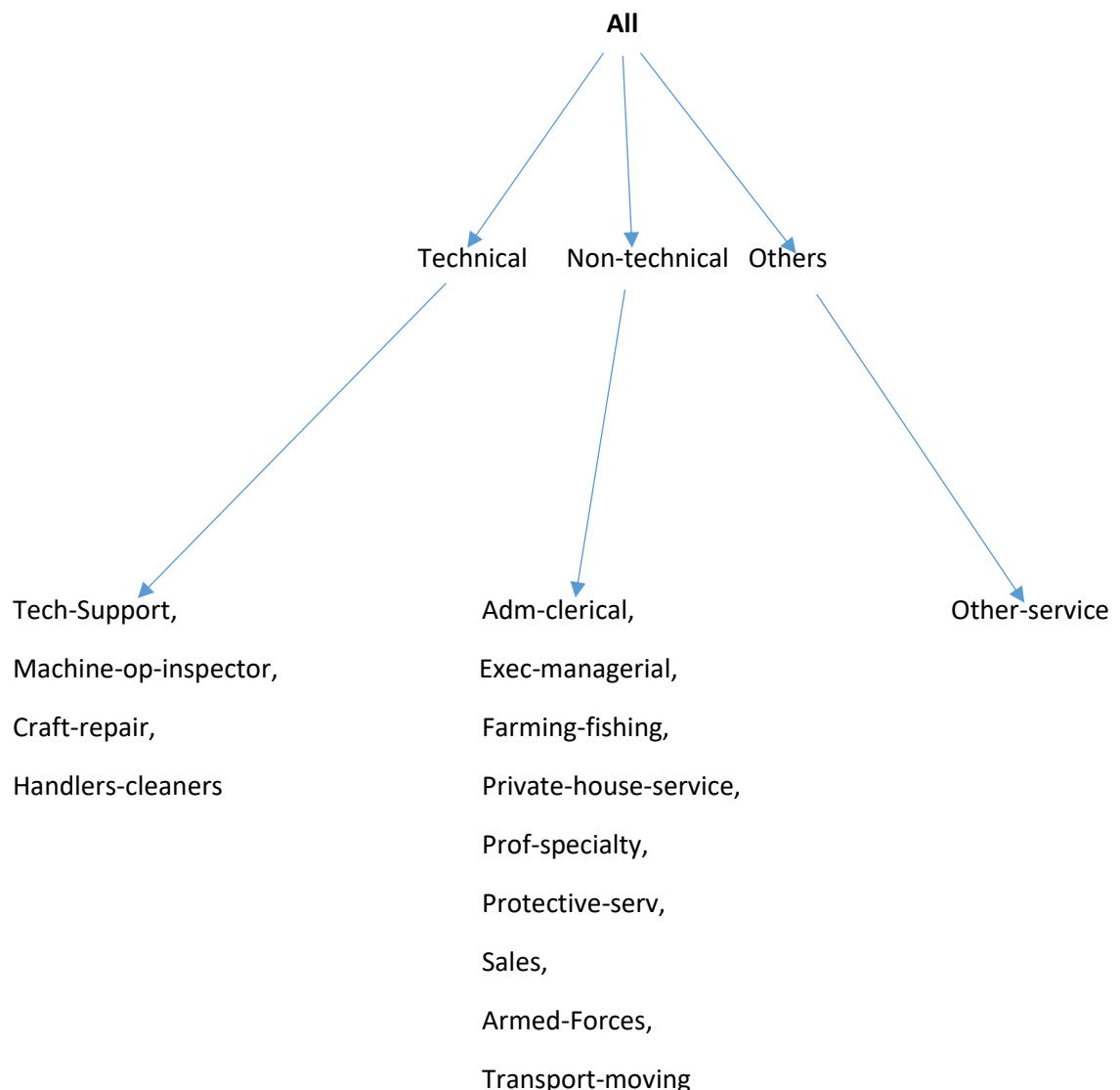
Here, the Drill-down process has been used to find the concept hierarchy for the US-Adult-Income. We will get information by using the lowest level hierarchy.



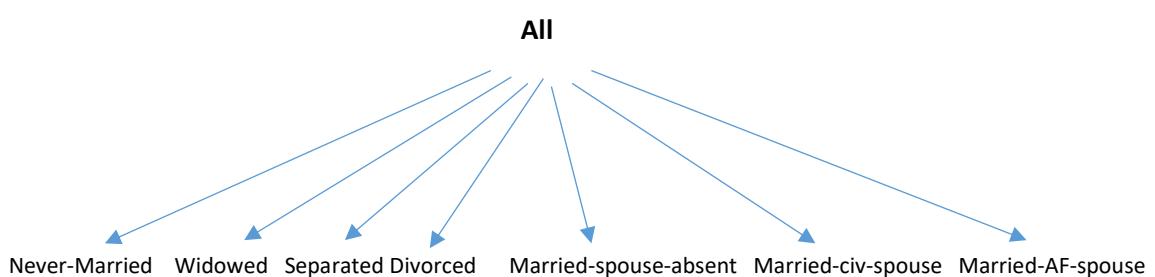




Occupation

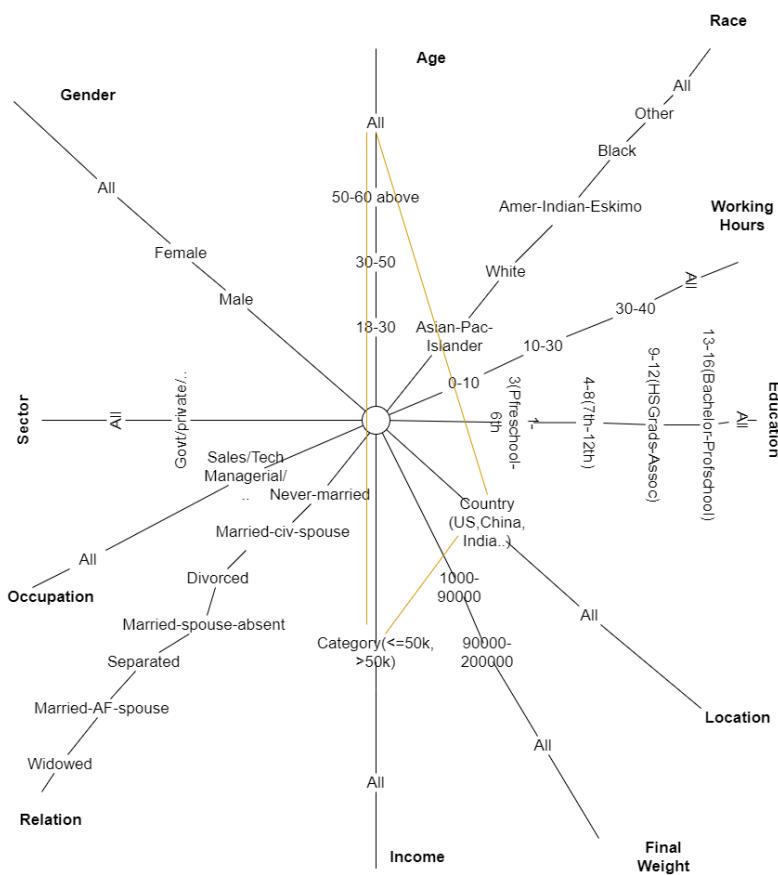


Marital status

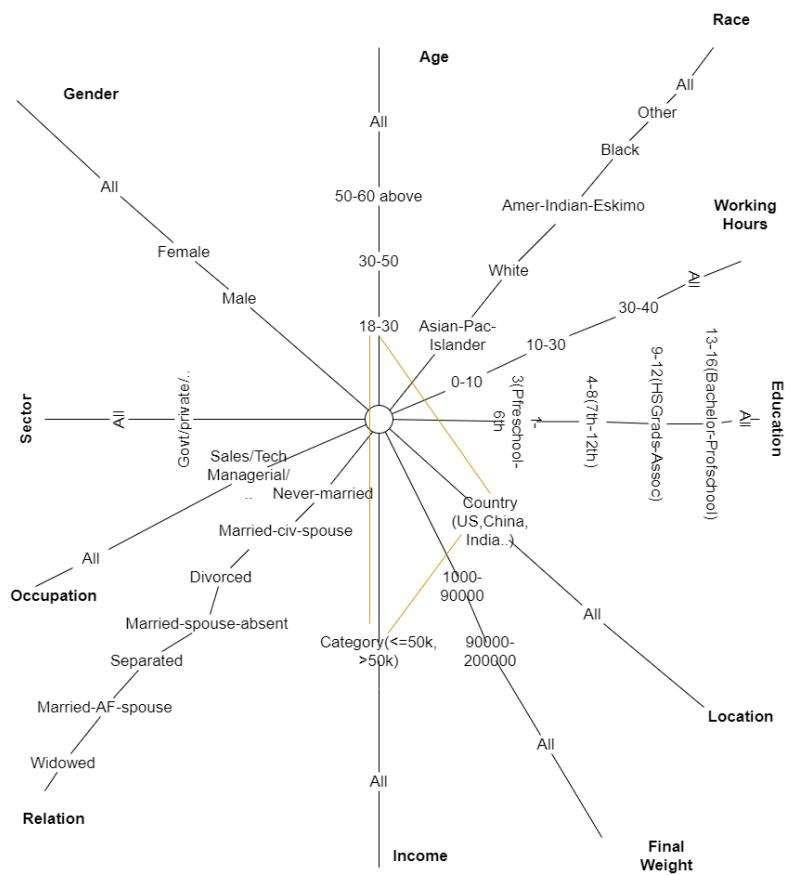


6. Starnet footprints for business queries

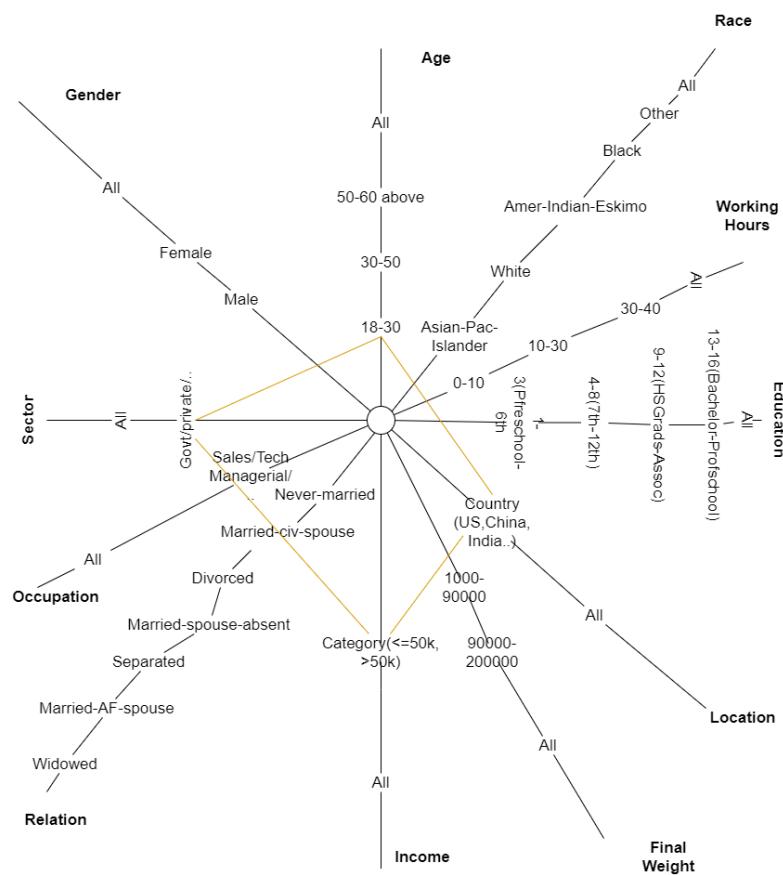
- We will go with the mentioned footprint in order to solve the first business query. This footprint helps in finding the number of people having income maximum 50000 USD per annum also we are finding the number of people who have income more than 50000 USD a year in the same business query from US, China, India, Germany and Mexico of all age group.



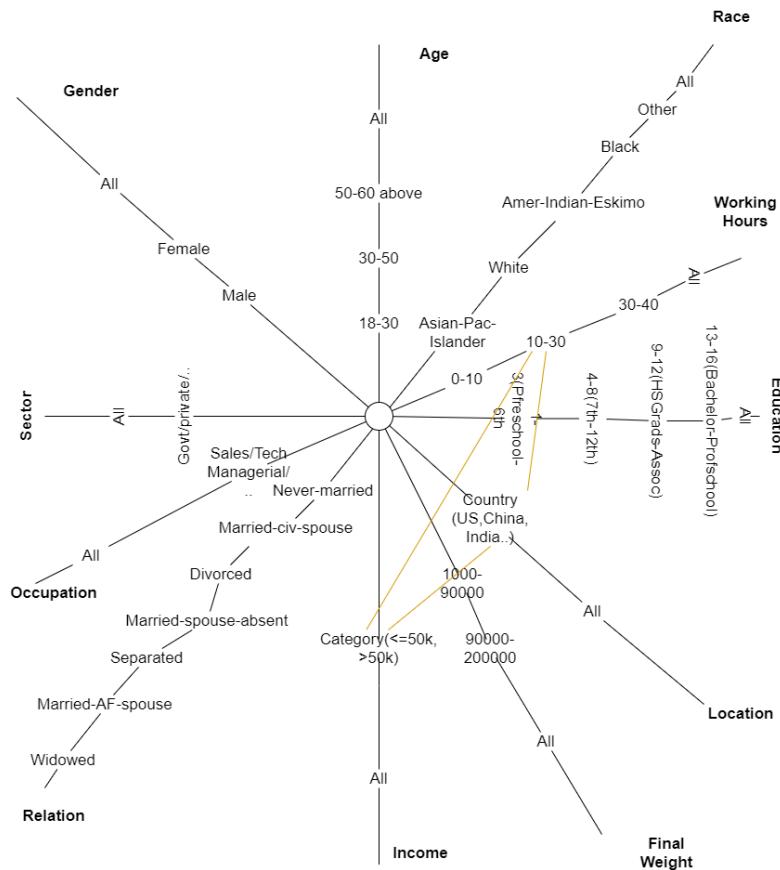
2. This footprint focuses on second business query which finds the number of young people having income more than 50000 USD in a year in United States.



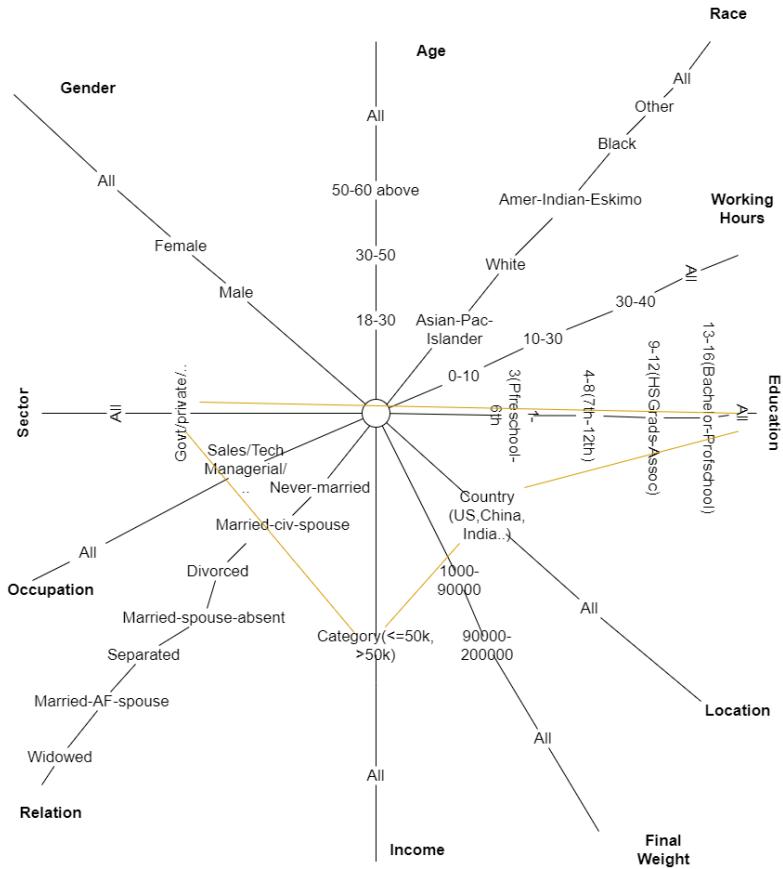
3. Finding the number of young employees of a private sector companies in Mexico and Canada who is having income over 50000 USD per annum.



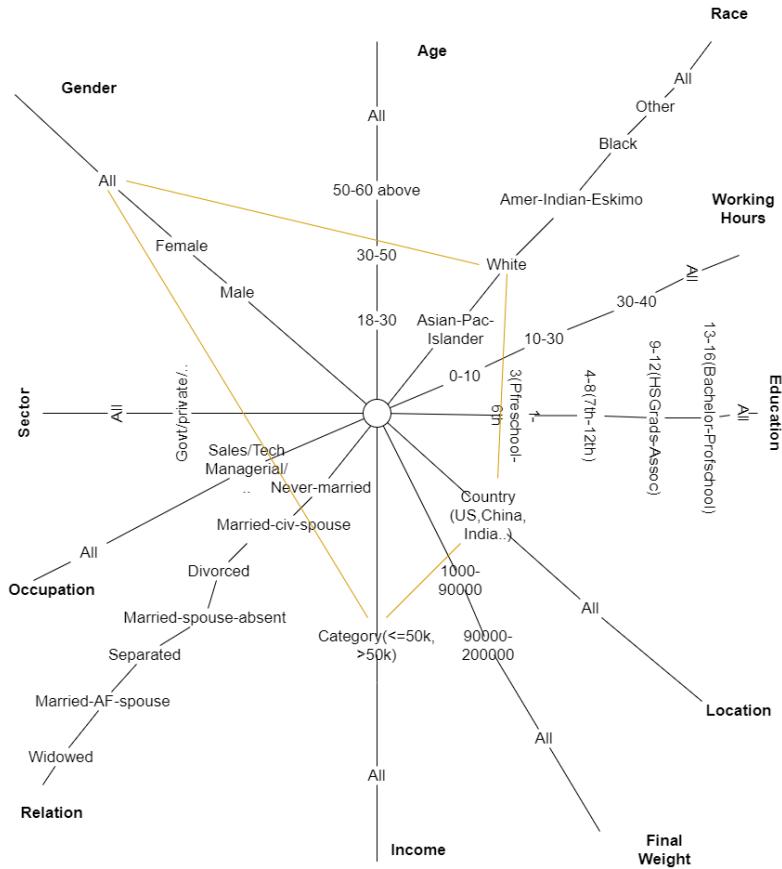
4. 4th business query also helps in find the people having highest wages rate per week because in this query the working hour is less and income is high. So, we are finding the people having income 50000 USD per annum having shifts from 10- 30 hours a week for people from US.



5. Finding the preferred education people who is earning more than 50000 USD a year and maximum 50000 USD a year from United-States.

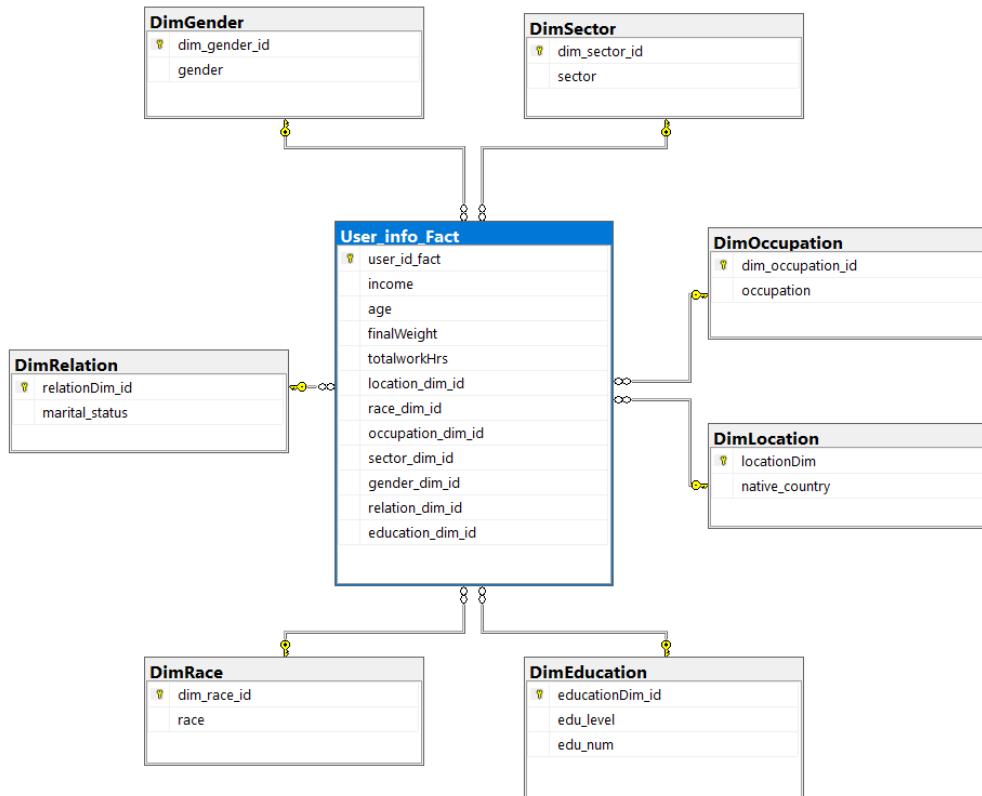


7. This starnet foot print helps in getting all the males and females working as a sales representative, executive manager and a machine-op-inspector(Occupation: sales,exec-managerial, machine-op-inspector) having maximum income USD 50000 in a year from US, Canada and Cuba having race white.



7. Starnet database diagram (ER diagram) :

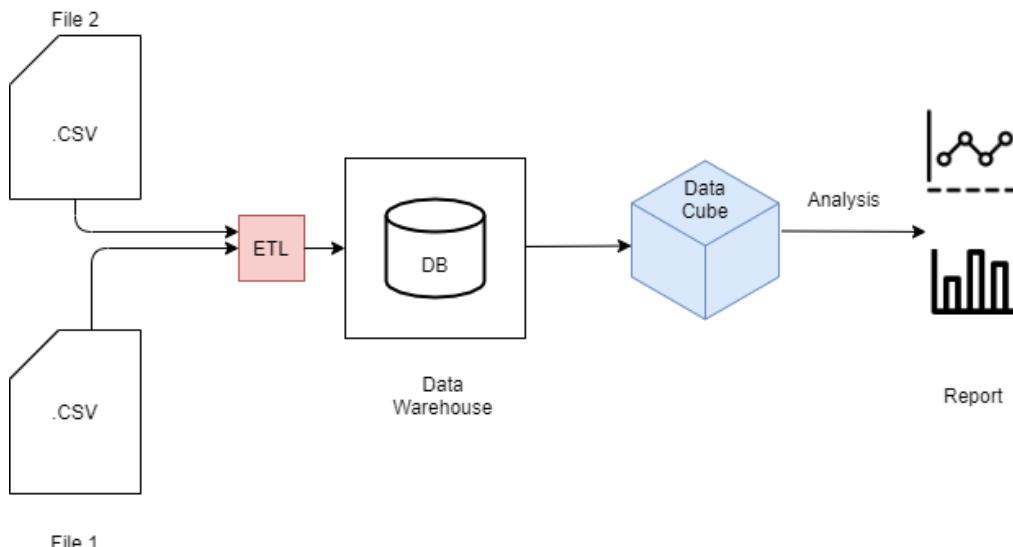
Fact tables	Dimension tables
User	Gender
	Location
	Education
	Occupation
	Relation
	Race
	Sector



8.ETL process

Extract Transform Load gets the data from different sources to a destination. Here, in this case our source is the .csv files. The .csv files will contain the `User_info_fact.csv` as fact table the rest of the file will be the dimension tables, dimensional tables are `DimSector.csv`, `DimRelation.csv`, `DimRace.csv`, `DimSector.csv`, `DimEducation.csv`, `DimGender.csv` and the `DimLocation.csv`.

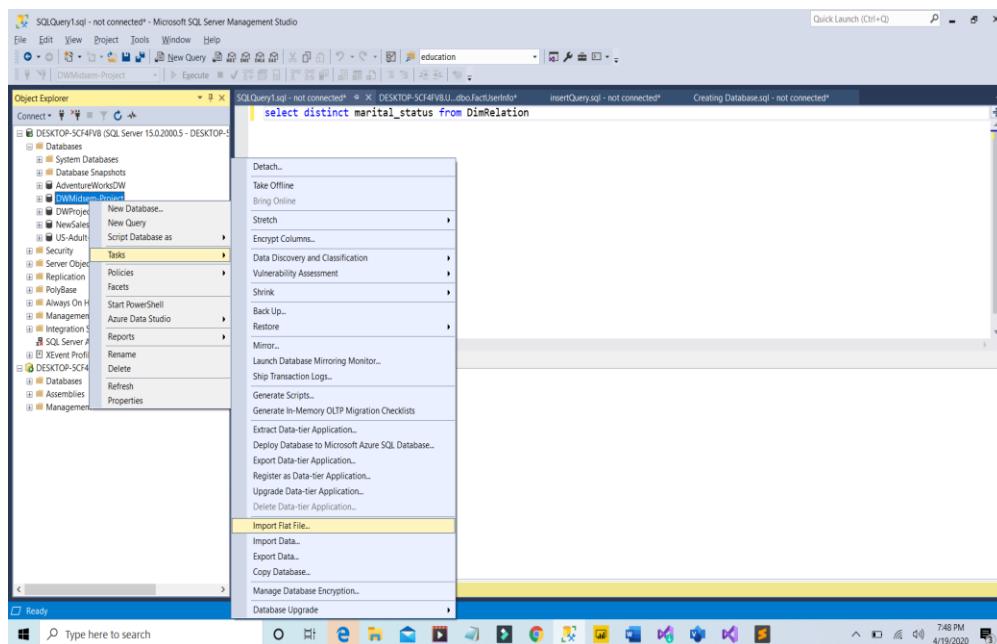
A diagram to show the ETL process for this project



How to import the .csv files and populate the database:

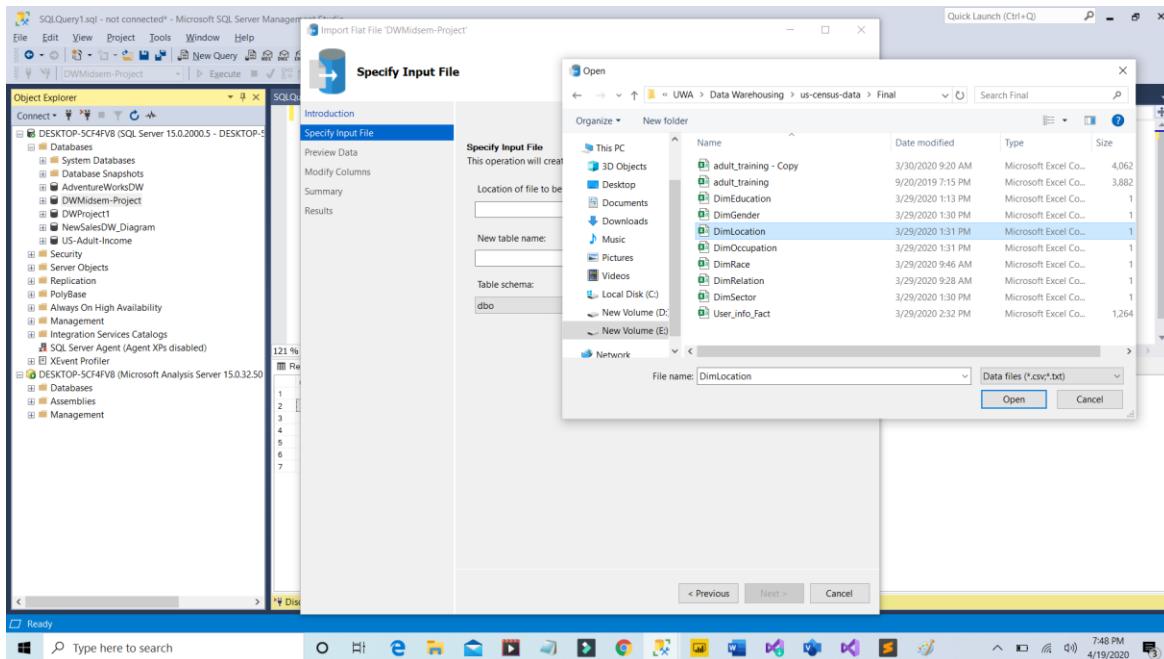
Step 1 :

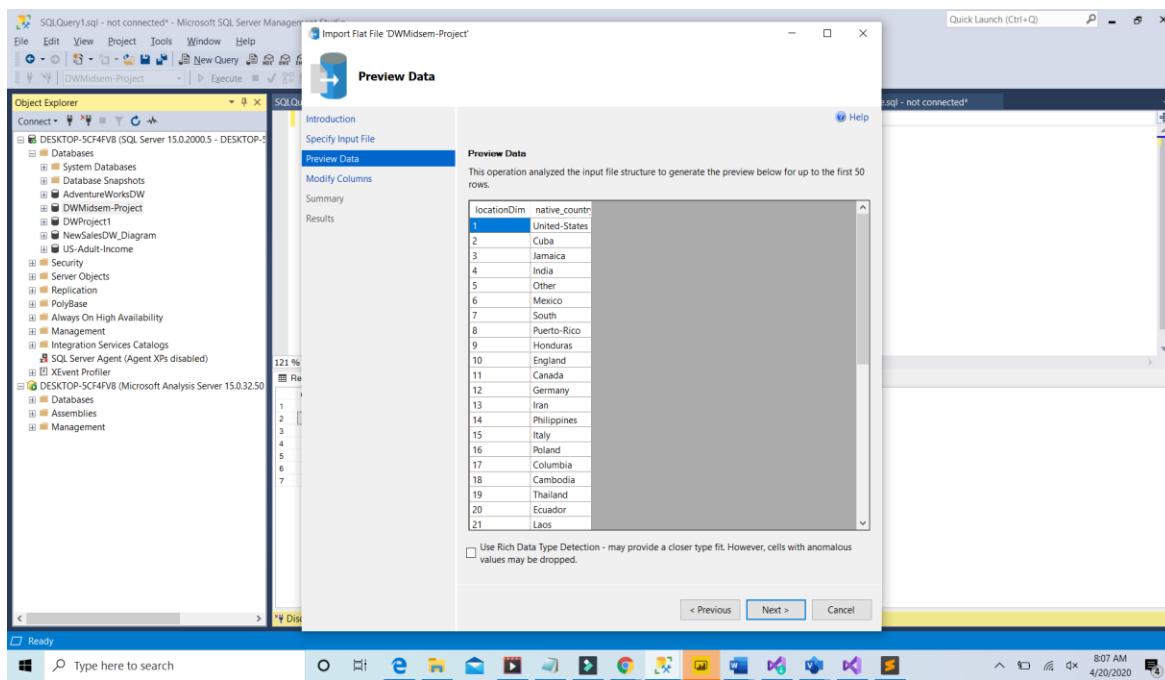
1. Go to object explorer select the name of the database where will populate the files.
2. Right click on to the Database go to Tasks.
3. Then go to import flat file at the bottom.



Step 2:

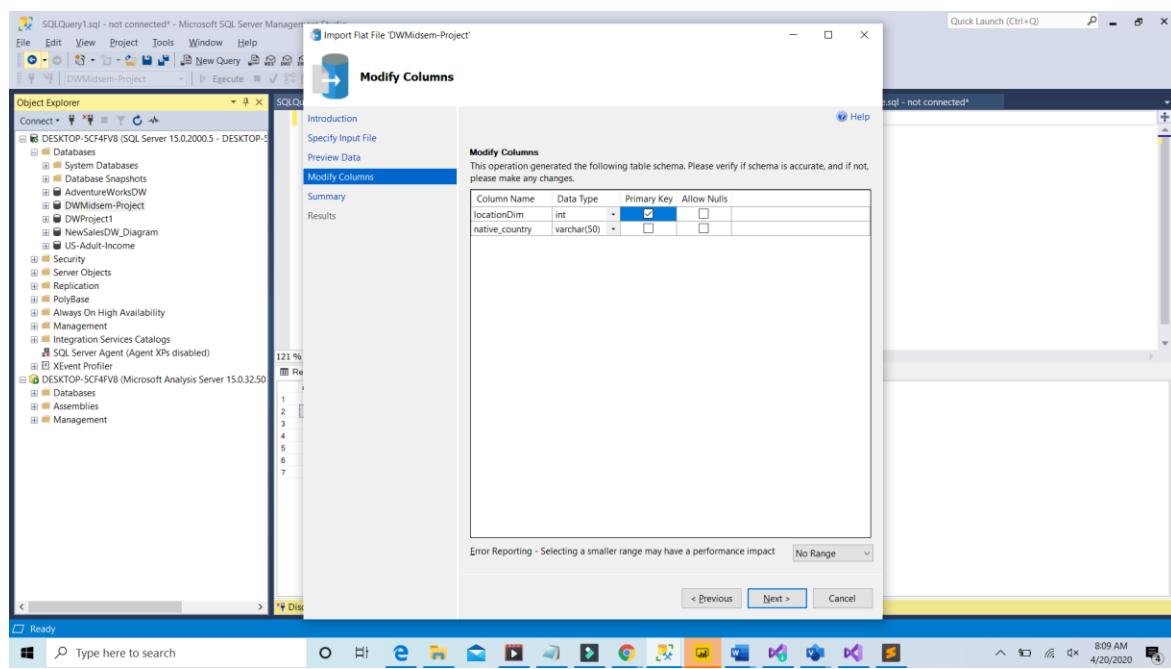
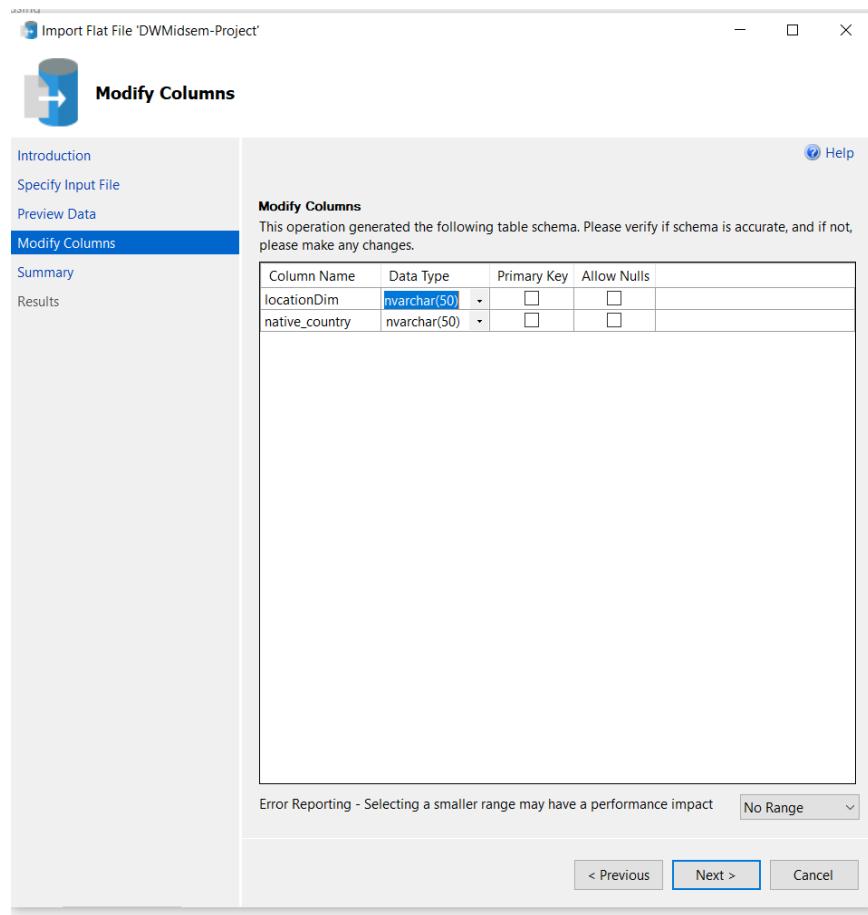
1. Then select the location where the csv file located. In my case the location is "E:\UWA\Data Warehousing\us-census-data\Final"
2. Load the .CSV





Step 3:

1. Change the datatypes. The IDs should be integer and the name column will be varchar because all the IDs of dimension table will work as a foreign key in our fact table and the data type of surrogate keys in the fact table is integer so both the datatype should match in order to create the relationship between two tables.



Step 4:

1. Now, click next and finished. The .csv file will now populate the database.
2. Do the same process from all the csv file, including all the dimension and fact tables.

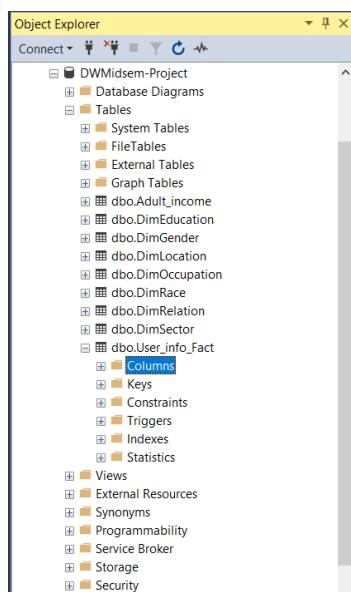
Step 5:

1. Now creating the relationship between the fact and the dimension tables.

Step 6:

1. Select the fact table Go to columns, right click and select the new column.
2. Once the new window opens for the columns of the fact table which is User_fact_info in this case. Then select the columns of the dimension table Right Click and choose relationships. (Do the same process for all the dimension table one by one).
3. Now, the new window will open then click "Add". This will add a relation to the particular column.
4. Now, open Tables and column Space as shown in the screenshot.
5. Now, in the Primary key table section choose your dimension table in the dropdown list and choose the column which is the **primary key** of that particular table.
6. Go to the Foreign key table section and select the column which makes the foreign key in the fact table of that particular dimension table.

The whole process to generate a foreign key



DESKTOP-5CF4FV8.DWMidsem-Project - dbo.User_info_Fact - Microsoft SQL Server Management Studio

File Edit View Project Table Designer Tools Window Help

Quick Launch (Ctrl+Q) x

Object Explorer

- DWMidsem-Project
 - Tables
 - System Tables
 - FileTables
 - External Tables
 - Graph Tables
 - dbo
 - dbo.Adult_income
 - dbo.DimEducation
 - dbo.DimGender
 - dbo.DimLocation
 - dbo.DimOccupation
 - dbo.DimRace
 - dbo.DimRelation
 - dbo.DimSector
 - dbo.DimSector
 - dbo.User_info_Fact
 - Columns
 - Keys
 - Constraints
 - Triggers
 - Indexes
 - Statistics
 - Views
 - External Resources
 - Synonyms
 - Programmability
 - Service Broker
 - Storage
 - Security
- DWProject1
- NewSalesDW_Diagram
- US-Adult-Income
- Security
- Server Objects
- Replication

DESKTOP-5CF4FV8.D..dbo.User_info_Fact SQLQuery1.sql - not connected* DESKTOP-5CF4FV8U..dbo.FactUserInfo* insertQuery.sql - not connected*

Column Name	Data Type	Allow Nulls
user_id_fact	int	<input type="checkbox"/>
income	varchar(50)	<input type="checkbox"/>
age	int	<input type="checkbox"/>
finalWeight	int	<input type="checkbox"/>
totalworkhrs	int	<input type="checkbox"/>
location_dim_id	int	<input type="checkbox"/>
race_dim_id	int	<input type="checkbox"/>
occupation_dim_id	int	<input type="checkbox"/>
sector_dim_id	int	<input type="checkbox"/>
gender_dim_id	int	<input type="checkbox"/>

Set Primary Key Insert Column Delete Column Relationships... Indexes/Keys... Full Text Index XML Indexes... Check Constraints... Spatial Indexes... Generate Change Script... Properties Alt+Enter

Allow Nulls Data Type Default Value or Binding

Table Designer (General)

location_dim_id
No int

Ready Type here to search 8:31 AM 4/20/2020

DESKTOP-5CF4FV8.DWMidsem-Project - dbo.User_info_Fact - Microsoft SQL Server Management Studio

File Edit View Project Table Designer Tools Window Help

Quick Launch (Ctrl+Q) x

Object Explorer

- DWMidsem-Project
 - Tables
 - System Tables
 - FileTables
 - External Tables
 - Graph Tables
 - dbo
 - dbo.Adult_income
 - dbo.DimEducation
 - dbo.DimGender
 - dbo.DimLocation
 - dbo.DimOccupation
 - dbo.DimRace
 - dbo.DimRelation
 - dbo.DimSector
 - dbo.DimSector
 - dbo.User_info_Fact
 - Columns
 - Keys
 - Constraints
 - Triggers
 - Indexes
 - Statistics
 - Views
 - External Resources
 - Synonyms
 - Programmability
 - Service Broker
 - Storage
 - Security
- DWProject1
- NewSalesDW_Diagram
- US-Adult-Income
- Security
- Server Objects
- Replication

DESKTOP-5CF4FV8.D..dbo.User_info_Fact SQLQuery1.sql - not connected* DESKTOP-5CF4FV8U..dbo.FactUserInfo* insertQuery.sql - not connected*

Column Name	Data Type	Allow Nulls
user_id_fact	int	<input type="checkbox"/>
income	varchar(50)	<input type="checkbox"/>
age	int	<input type="checkbox"/>
finalWeight	int	<input type="checkbox"/>
totalworkhrs	int	<input type="checkbox"/>
location_dim_id	int	<input type="checkbox"/>
race_dim_id	int	<input type="checkbox"/>
occupation_dim_id	int	<input type="checkbox"/>
sector_dim_id	int	<input type="checkbox"/>
gender_dim_id	int	<input type="checkbox"/>
relation_dim_id	int	<input type="checkbox"/>
education_dim_id	int	<input type="checkbox"/>

Foreign Key Relationships Selected Relationship: FK_User_info_Fact_DimGender

Editing properties for existing relationship.

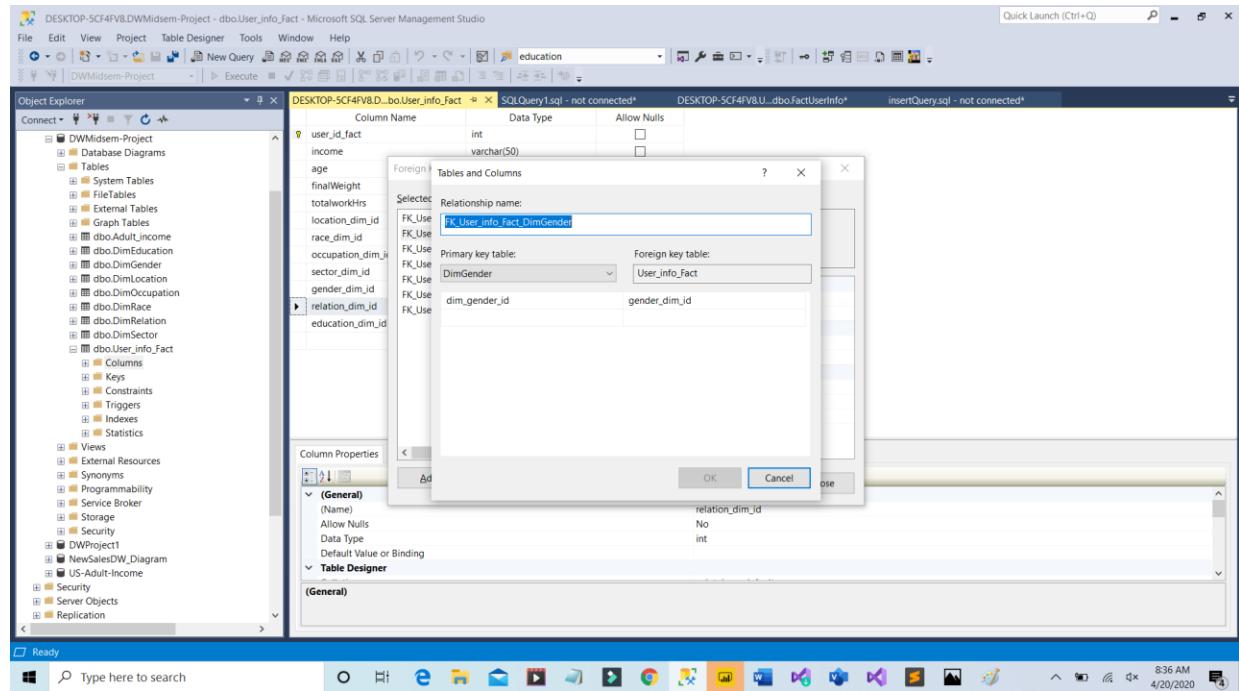
(General) (Name) FK_User_info_Fact_DimGender
Check Existing Data On Go Yes
Tables And Columns Spec:
Identity (Name) Description
Table Designer Enforce For Replication Yes
Enforce Foreign Key Consistency Yes
INSERT And UPDATE Spec

Column Properties Add Delete Close

relation_dim_id
No int

Table Designer (General)

Ready Type here to search 8:34 AM 4/20/2020

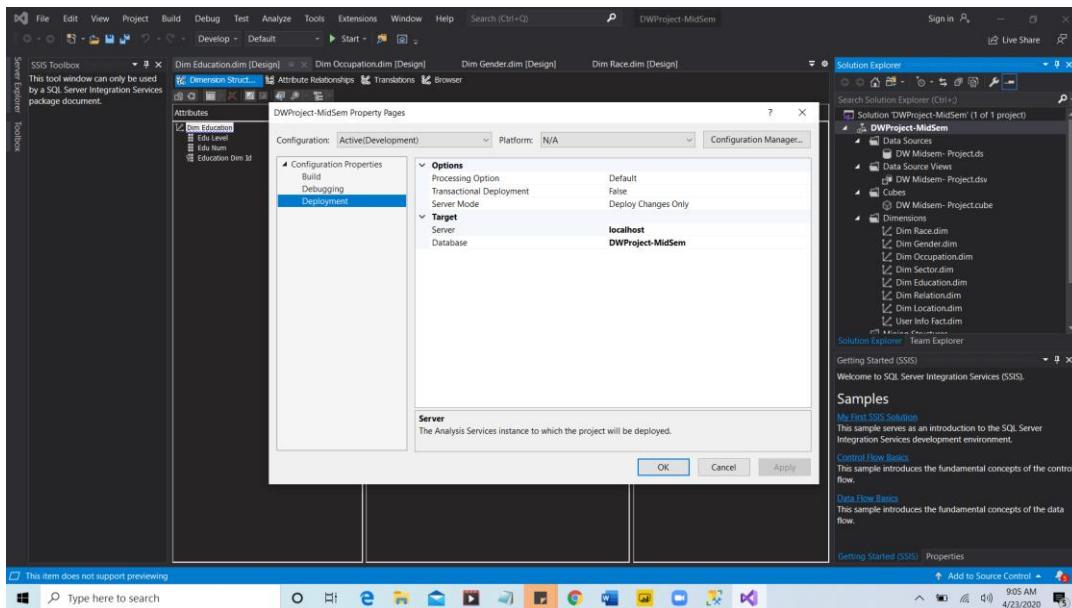


Do the same for all the dimension tables.

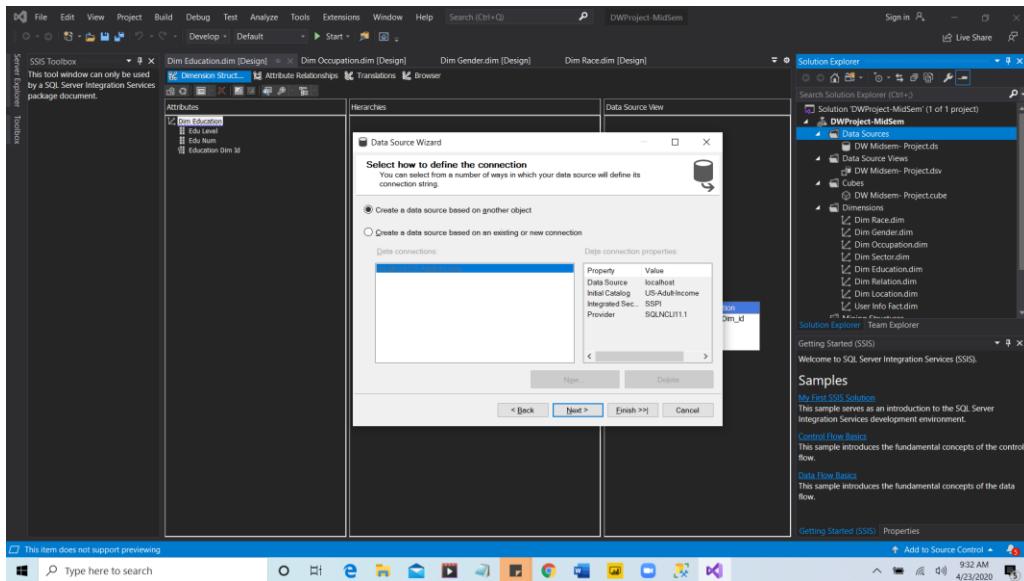
9. Process of creating Data Cube:

In order to design a data cube we use visual studio 2019 Enterprise. Visual should have the extension for Analysis Service Multidimensional and Data Mining project (ASMD). Then we create a project in ASMD in visual studio.

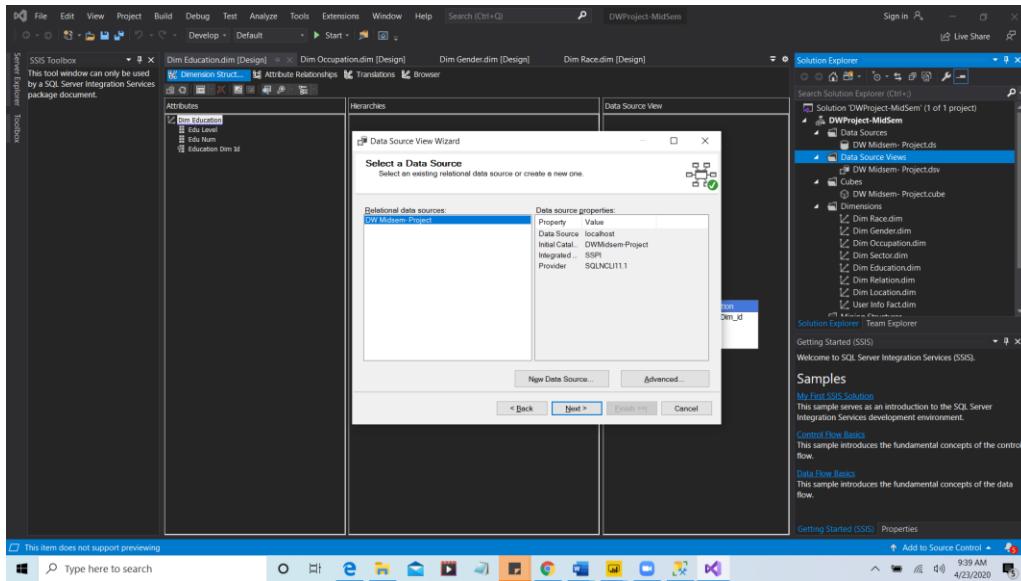
Step 1. Make sure in the properties of the project “DW Midsem-Project” the sever should be connected with the “localhost” .



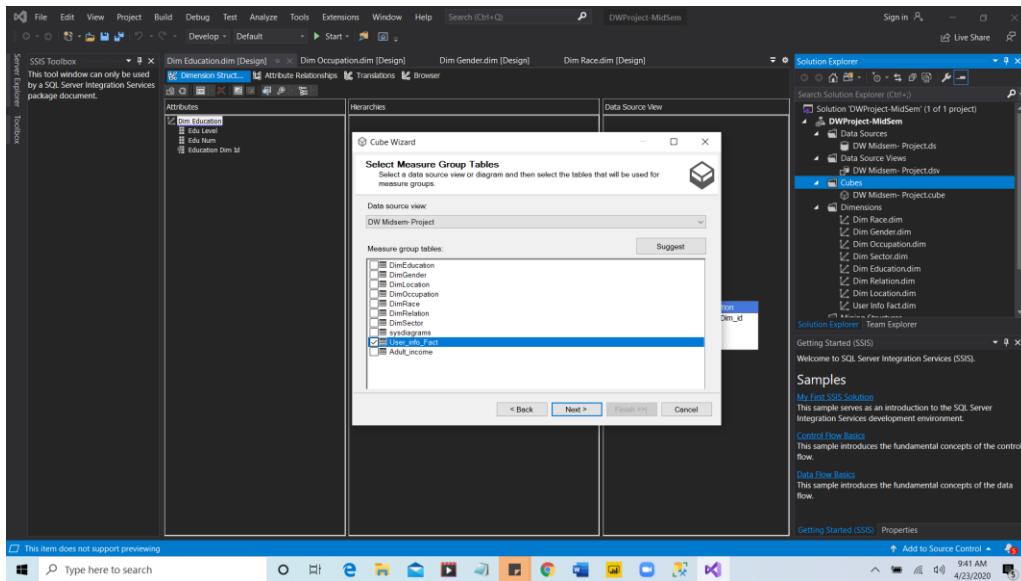
Step 2: Right click on the project solution (.ds file) and connect the Datasource.

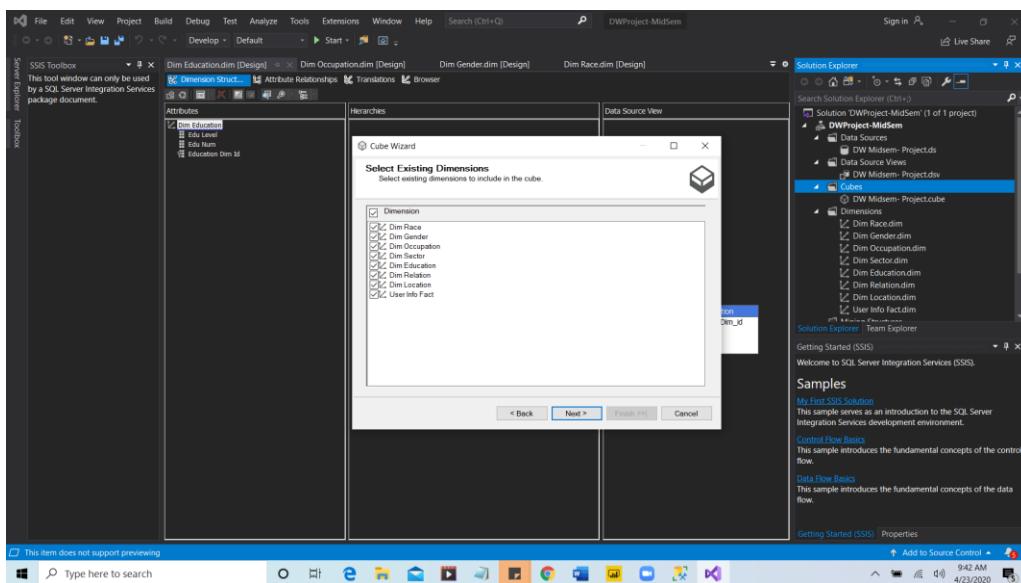
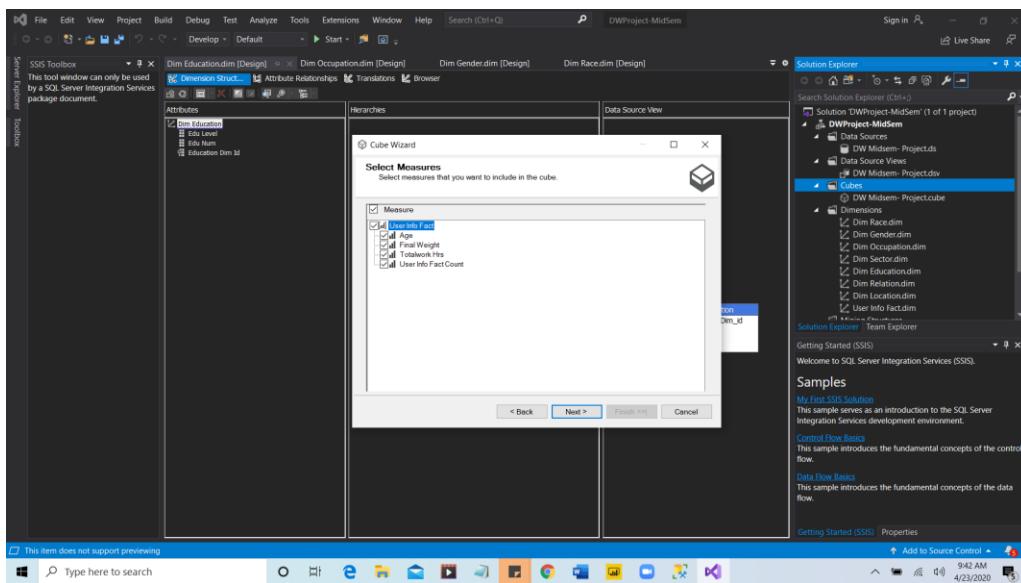


Step 3: Then go to the Data Source View folder and create a data source view.

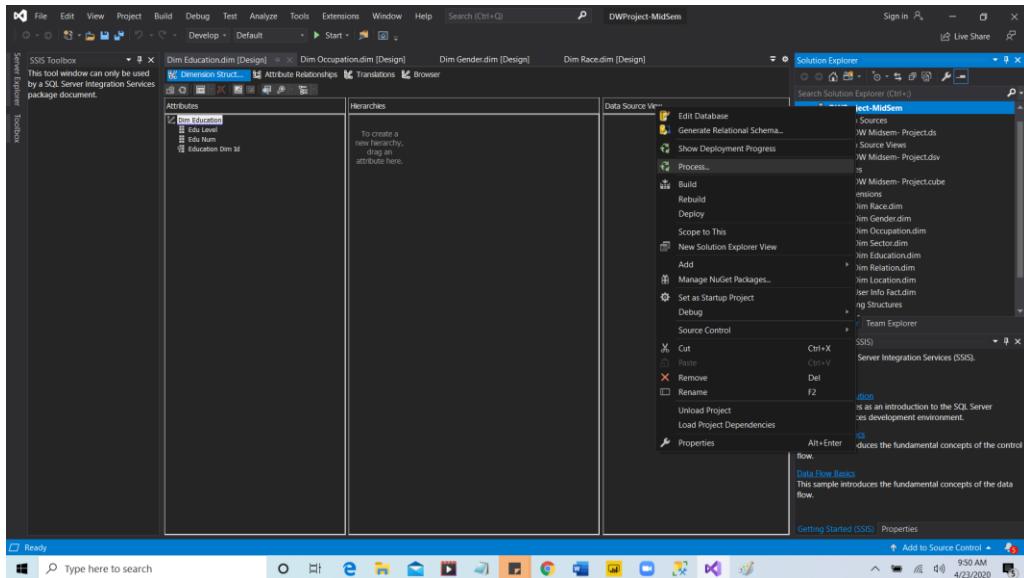
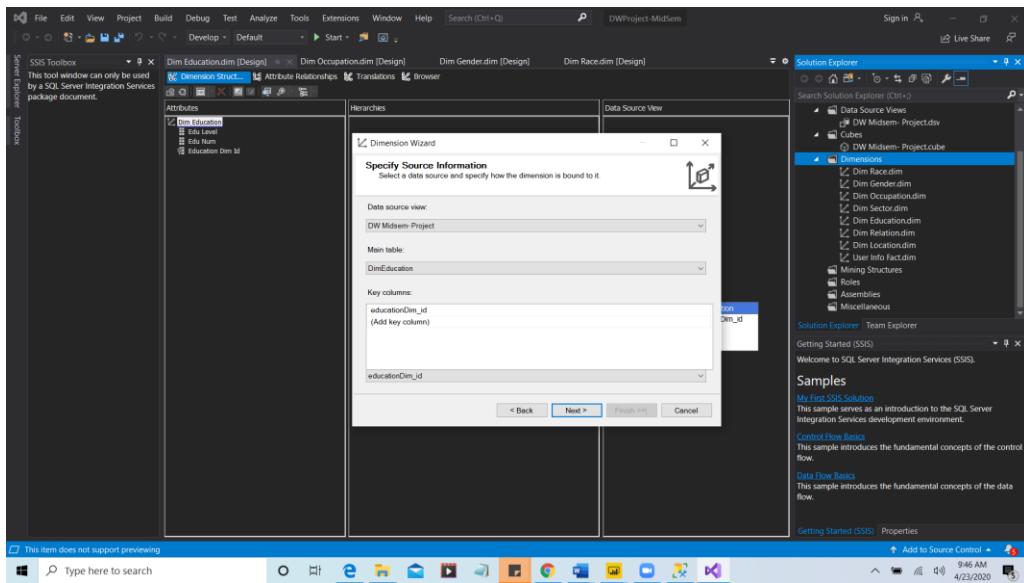


Step 4: Now, create the cube by right clicking on “Cubes”. Select only the fact table which is User_fact_info, once the fact table is selected the measures and the dimension tables linked with the fact table will automatically be selected because creates a relationship using surrogate keys. Now, load the cube.



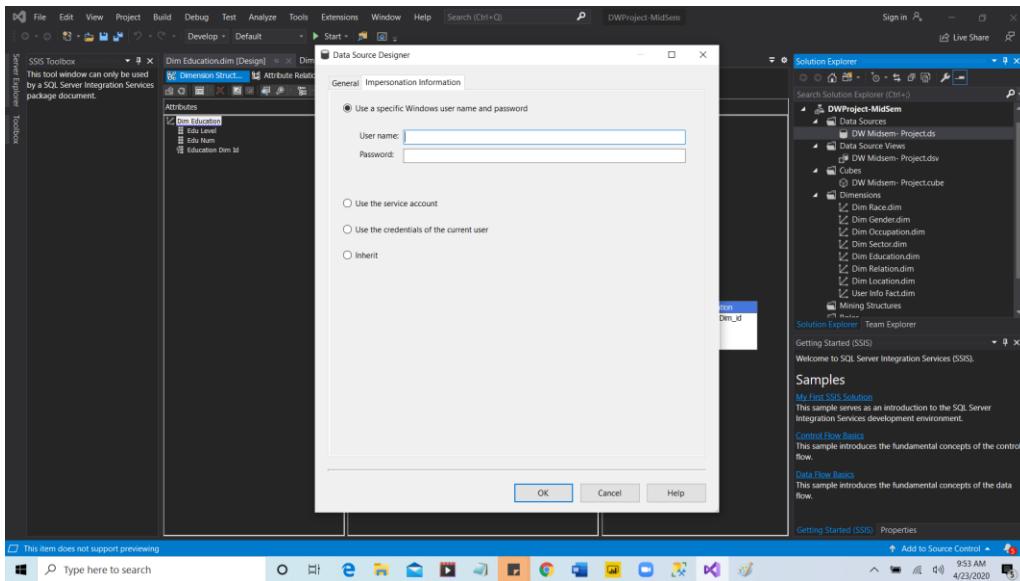


Step 5: Now, create the dimension by right clicking on “Dimensions” and then process it.

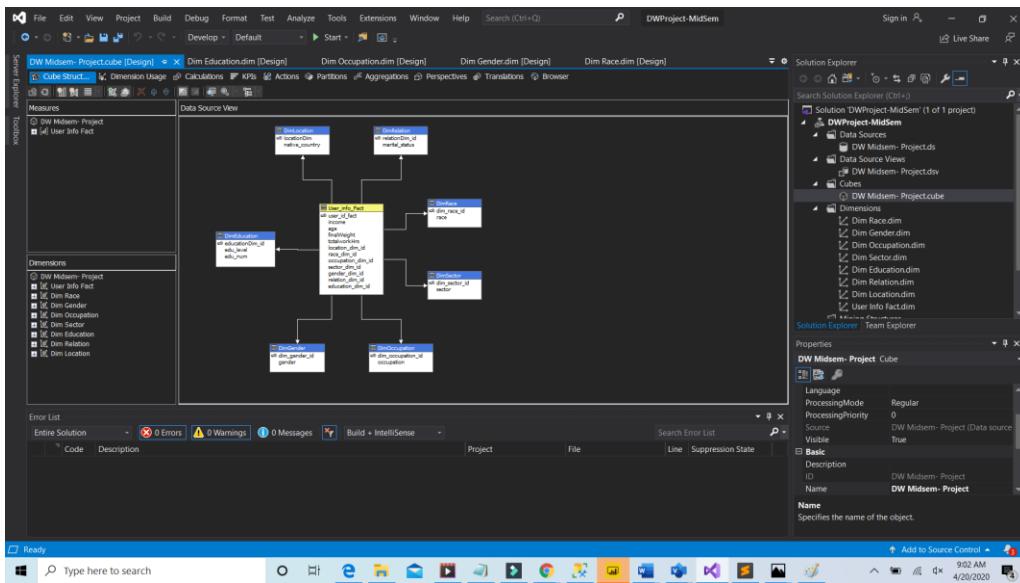
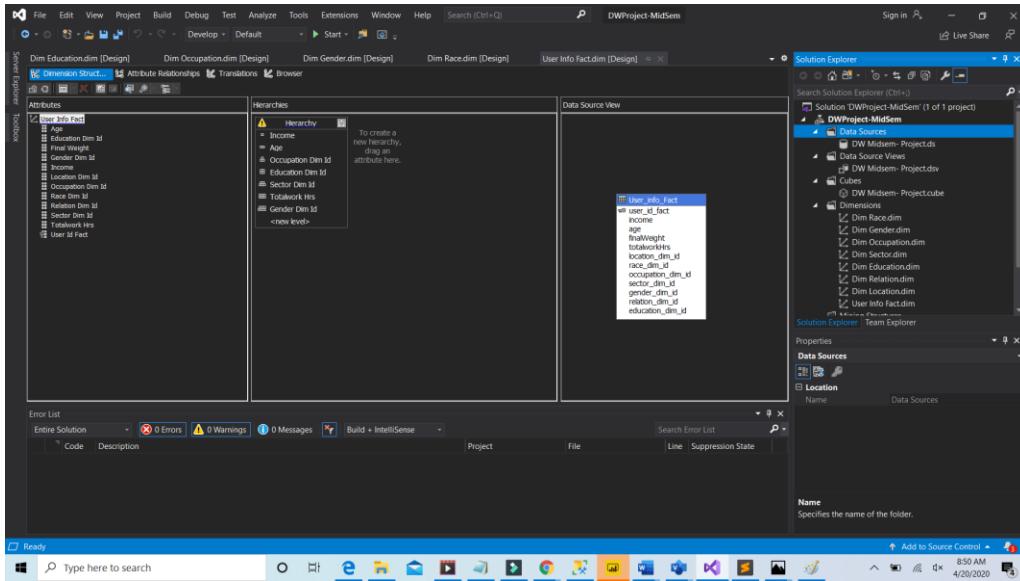


If any error is occurred while processing the data then we have to follow the 6th step.

Step 6: Make sure you have authentication of your system that means before processing the data double click on the .ds file in Data Sources then go to **Impersonation information** enter the username of the system and the password then click OK and then process again.



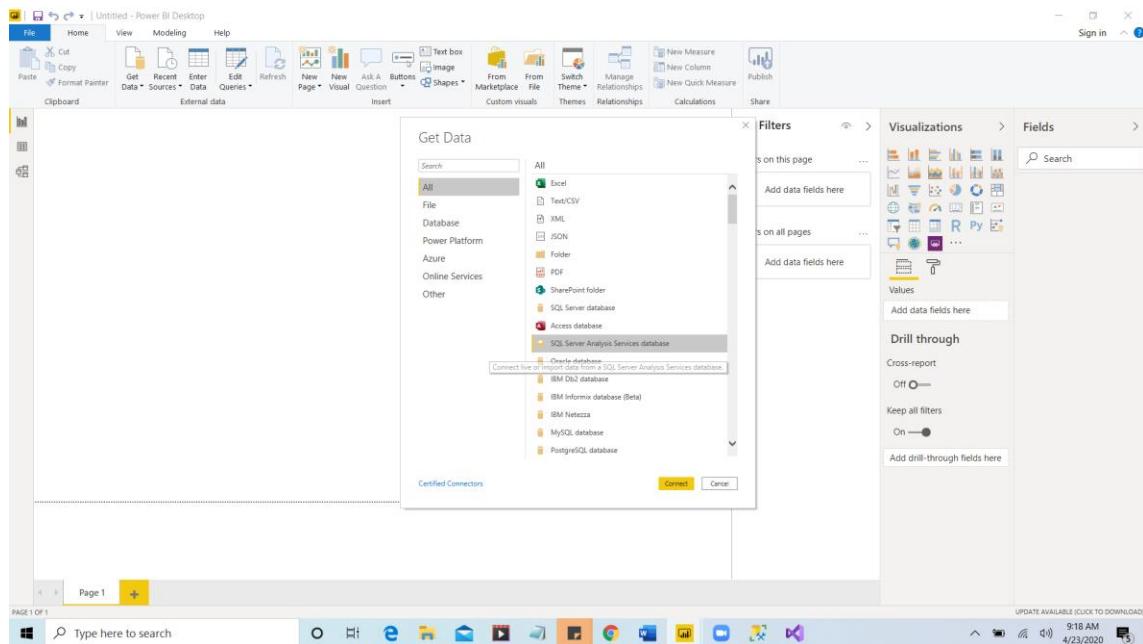
10. Data cube hierarchy in Visual Studio



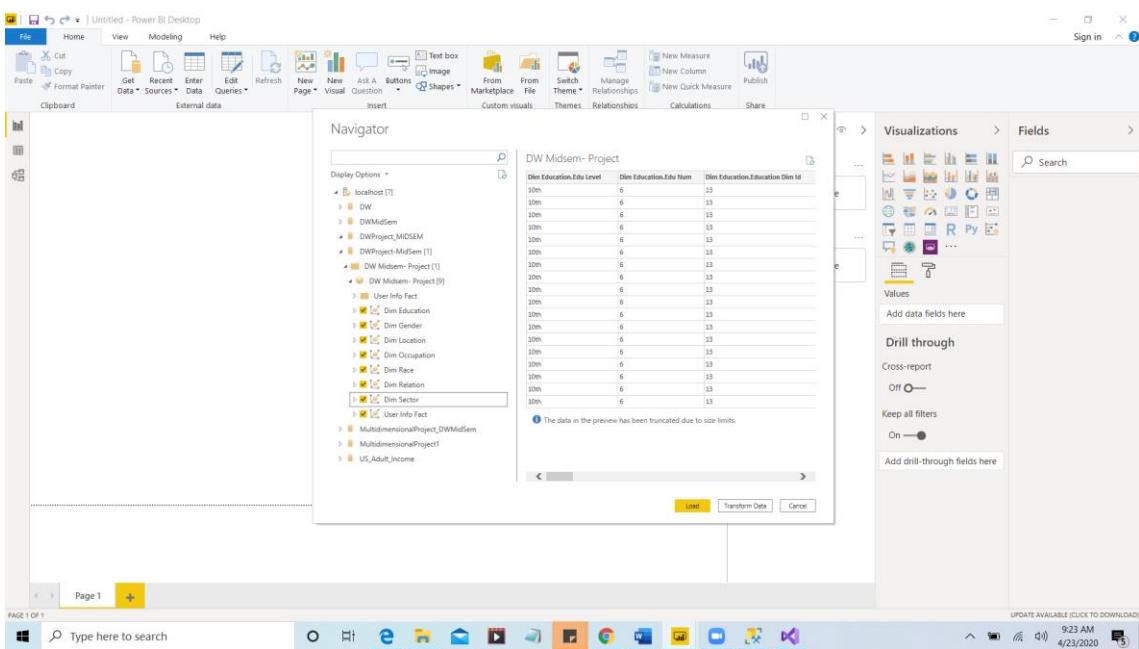
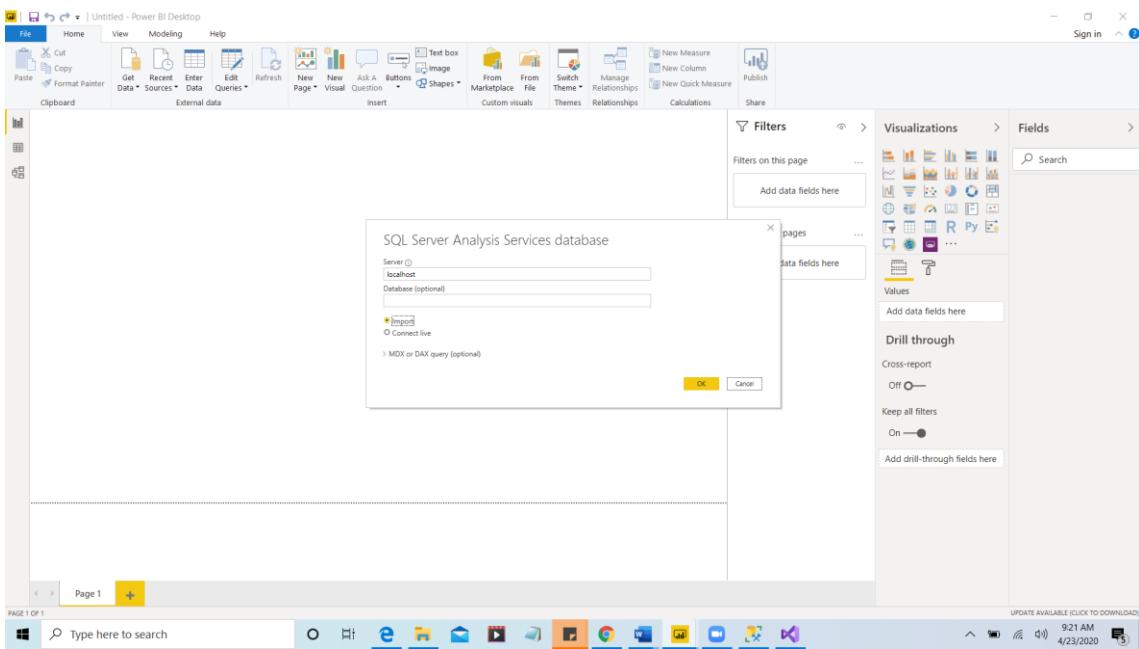
11. SQL Analysis service in Power BI

Process of importing the data cube in the power BI is shown below.

1. Open Microsoft PowerBI and click on **Get Data** select **Server Analysis Services database** and then **connect**.

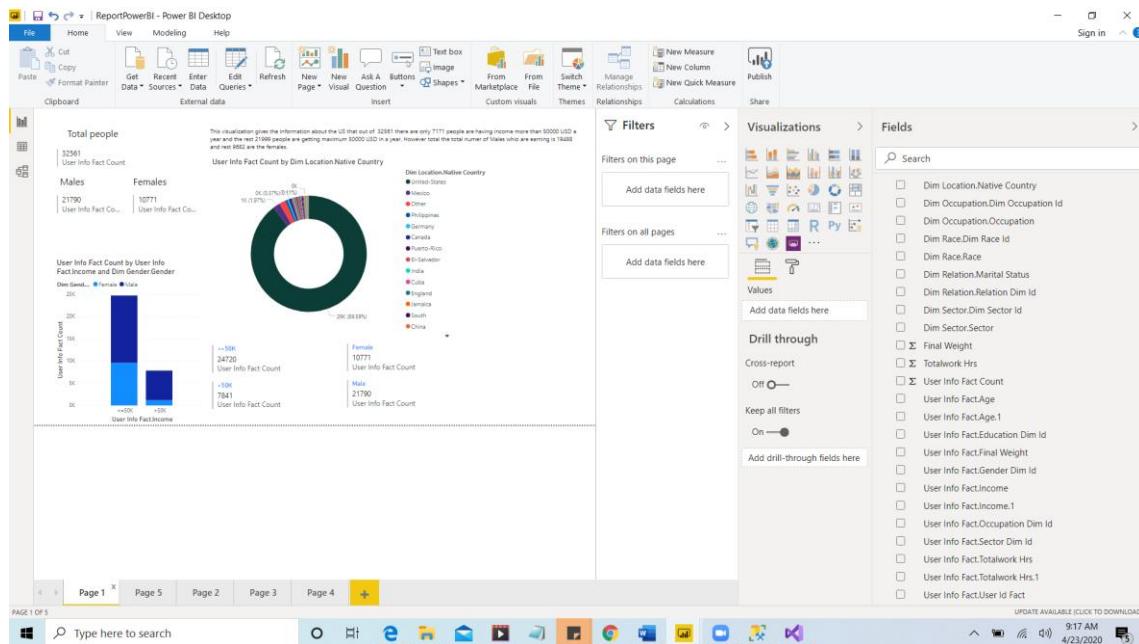


2. Enter your server name “localhost” and choose import then OK.
3. Select the relevant data cube which is created by the visual studio from the server and select the tables (Fact and the dimension) and then load.



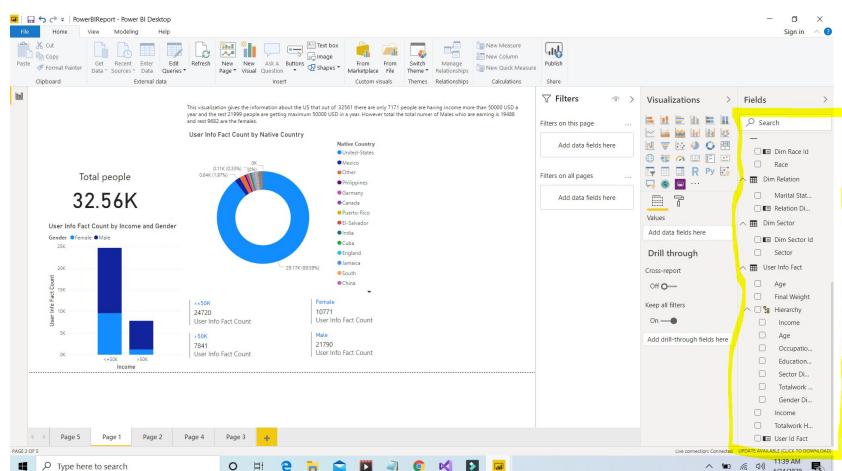
4. After the loading the data cube in Microsoft PowerBI it looks something like this. All the tables (Fact and dimension tables) and measures are situated at the top right side and we can do the analysis and visualization using the functions dragging from the visualization section.

Using Import Connection



12. Analysis Report screenshot using live connection :

Here, the screenshots are added to show the hierarchy of the cube using live connection with the database. But the final report has been made using import method in PowerBI. The purpose this is to show the hierarchy of the data cube as is not shown in PowerBI using the live connection. The proper analysis report has been added below using import method.



The hierarchy using live connection

This screenshot shows the Power BI Desktop interface with two charts displayed on the left and a detailed Fields pane on the right.

Left Side (Visualizations):

- User Info Fact Count by Totalwork Hrs:** A horizontal bar chart showing the count of users based on their total working hours per week. The x-axis ranges from 0 to 30 hours, with the highest count at 20 hours (75 users).
- User Info Fact Count by Edu Level and Sector:** Two stacked horizontal bar charts. The top chart, under the "Sector" filter (Federal-gov), shows counts for various education levels like HS-grad, Some-college, Bachelor's, etc. The bottom chart, also under "Sector" (Federal-gov), shows counts for different job sectors like Prof-school, Assoc-voc, etc.

Right Side (Fields pane):

A large yellow box highlights the "Fields" pane, which lists all available data fields categorized into dimensions and measures. The dimensions include:

- Dim Education:** Age, Final Weight, Totalwork Hrs, User Info Fact, Edu Level, Edu Num, Education ...
- Dim Gender:** Dim Gender, Gender
- Dim Location:** Dim Location Dim, Native Cou...
- Dim Occupation:** Dim Occup..., Occupation
- Dim Race:** Dim Race Id, Race
- Dim Relation:** Marital Stat..., Relation Di...
- Dim Sector:** Dim Sector Id, Sector

The measures listed are:

- User Info Fact:** Age, Final Weight, Hierarchy, Income, Totalwork Hrs, User Id Fact

This screenshot shows the same Power BI Desktop interface as the previous one, but with a specific focus on the "Hierarchy" field in the Fields pane.

Fields pane details:

A yellow box highlights the "Hierarchy" field under the "User Info Fact" measure. This indicates that the "Totalwork Hrs" dimension has been assigned a hierarchy, likely corresponding to the "Sector" dimension in the visualizations.

13. Final analysis report :

Total people

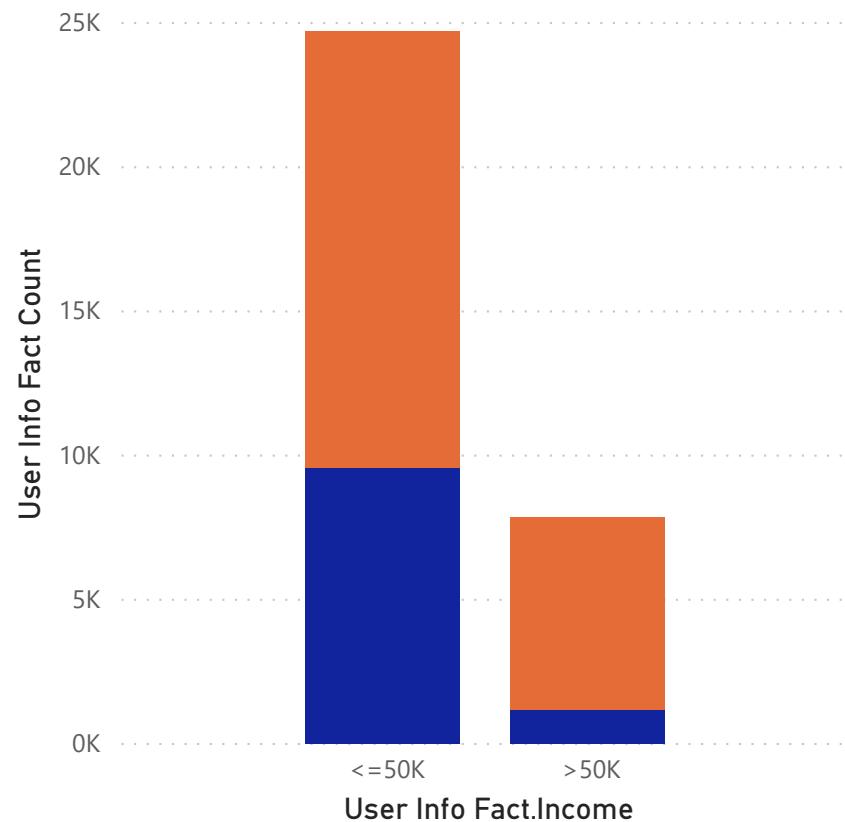
32561
User Info Fact Count

Males

21790
User Info Fact Co...

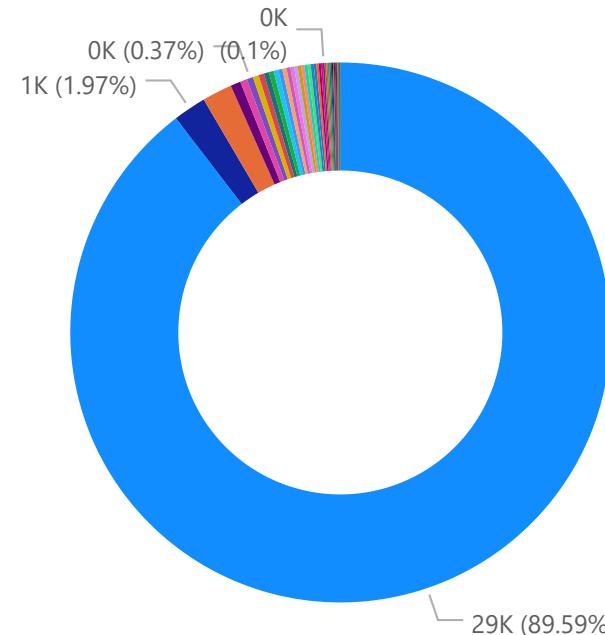
User Info Fact Count by User Info Fact.Income and Dim Gender.Gender

Dim Gender.G... ● Female ● Male



US-Adult-Income Dataset Over all description

User Info Fact Count by Dim Location.Native Country



Dim Location.Native Country

- United-States
- Mexico
- Other
- Philippines
- Germany
- Canada
- Puerto-Rico
- El-Salvador
- India
- Cuba
- England
- Jamaica

Female

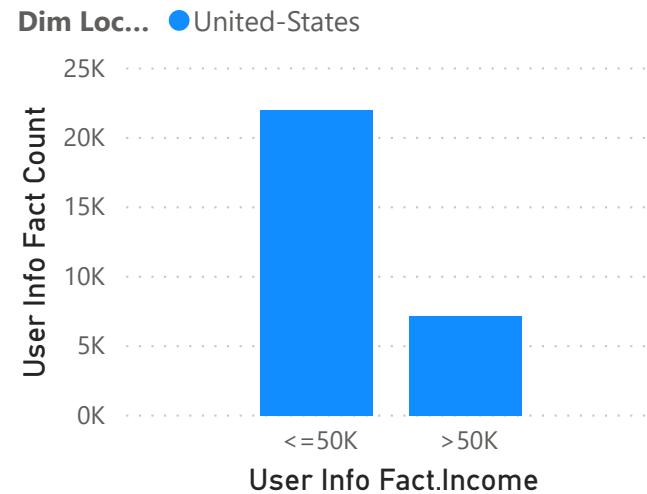
10771
User Info Fact Count

Male

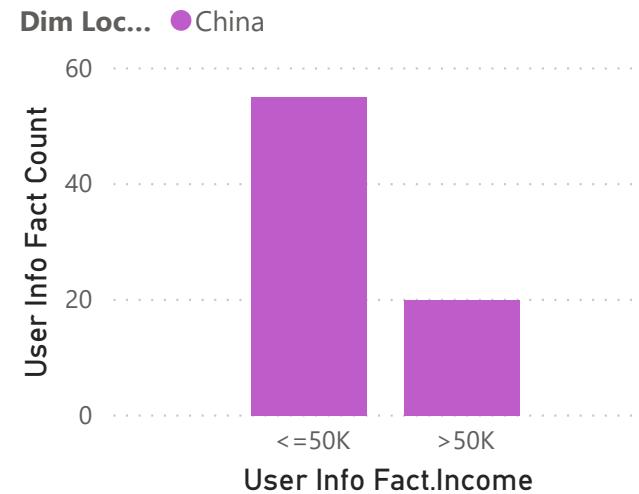
21790
User Info Fact Count

Income of people who earn more than USD 50000 also maximum USD 50000 a year from US, China, Germany, India and Mexico. (Solution to business Query 1)

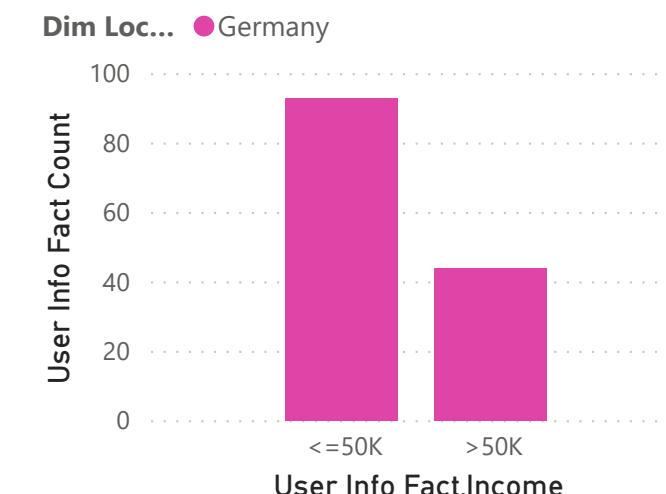
User Info Fact Count by User Info
Fact.Income and Dim Location.Native
Country



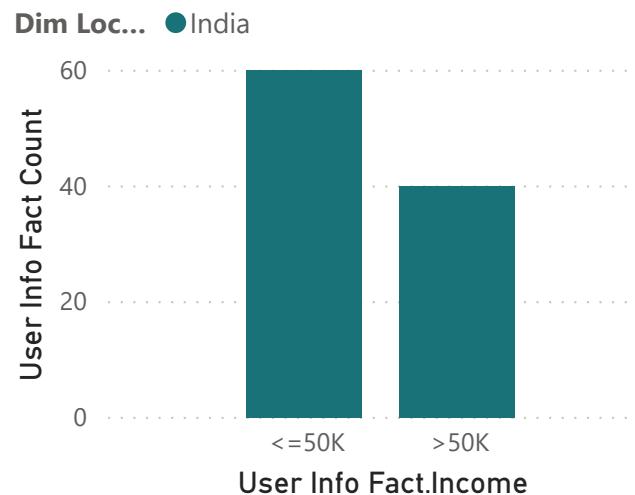
User Info Fact Count by User Info
Fact.Income and Dim Location.Native
Country



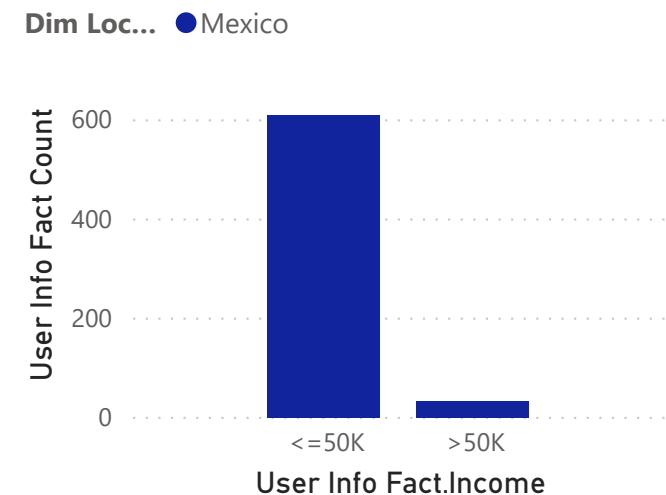
User Info Fact Count by User Info
Fact.Income and Dim Location.Native
Country



User Info Fact Count by User Info
Fact.Income and Dim Location.Native
Country



User Info Fact Count by User Info
Fact.Income and Dim Location.Native
Country

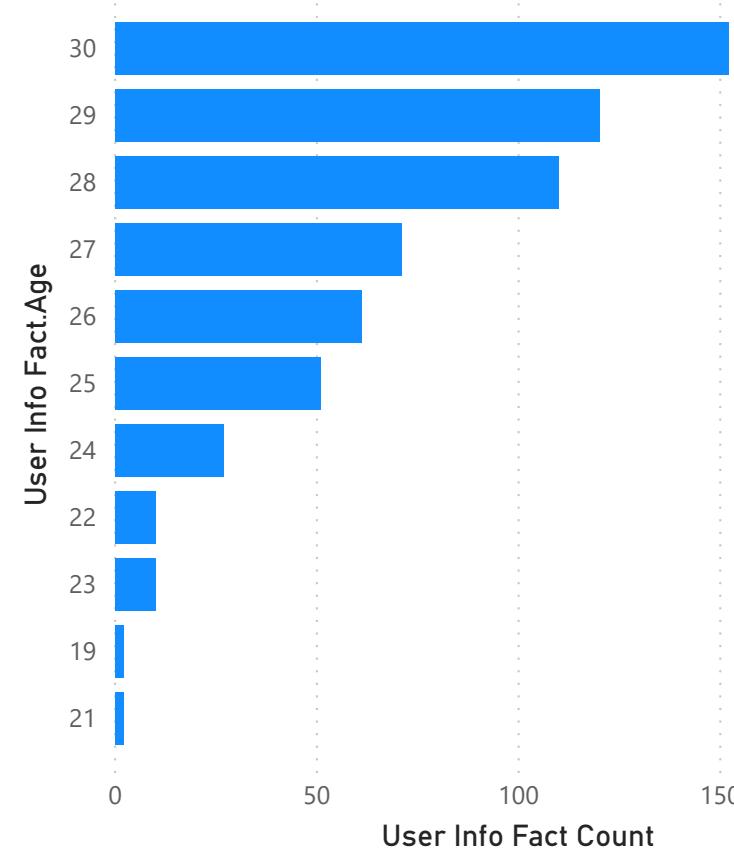


A very few people have got income
more than USD 50000 a year in
Mexico

The number of young people of age group 18-30 having income at least 50000 USD a year in the US (Solution to business Query 2)

Fact.Income

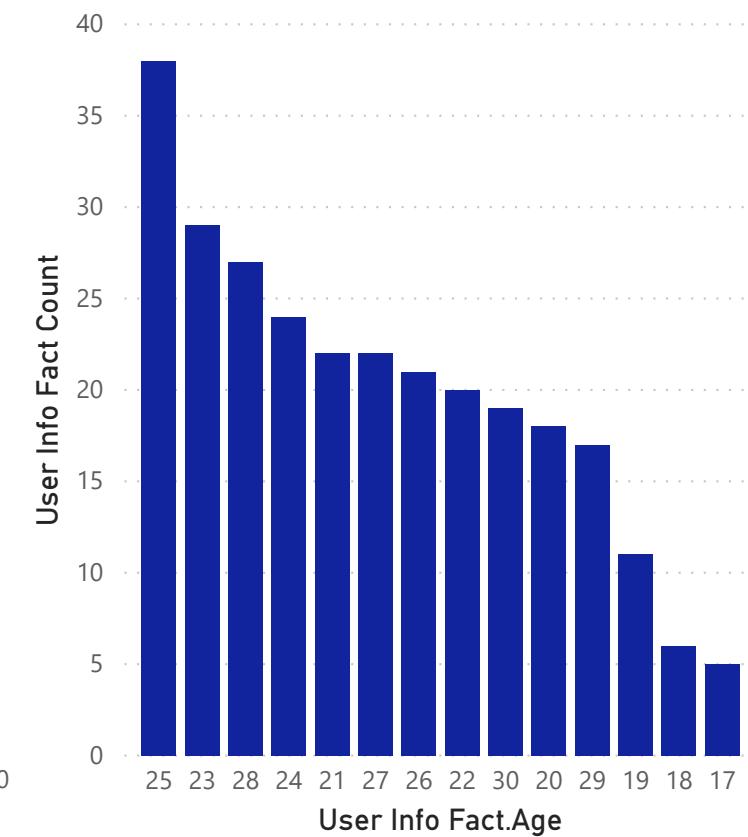
User Info Fact.In... ● >50K



Young people working for private companies in Mexico and Canada having income more than USD 50000 per annum. (Solution to business Query 3)

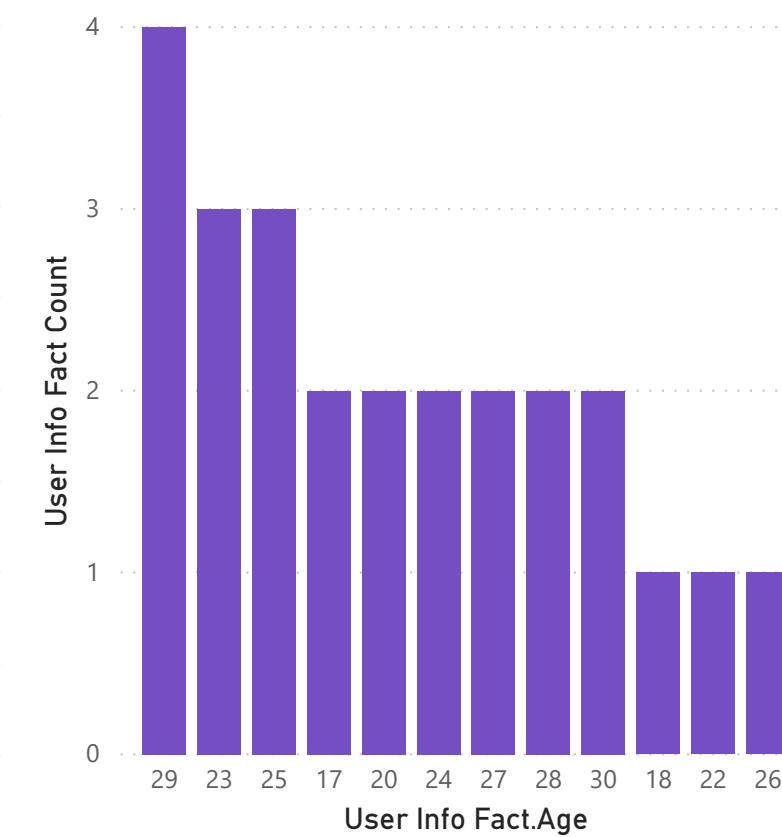
User Info Fact Count by User Info Fact.Age and Dim Location.Native Country

Dim Locati... ● Mexico

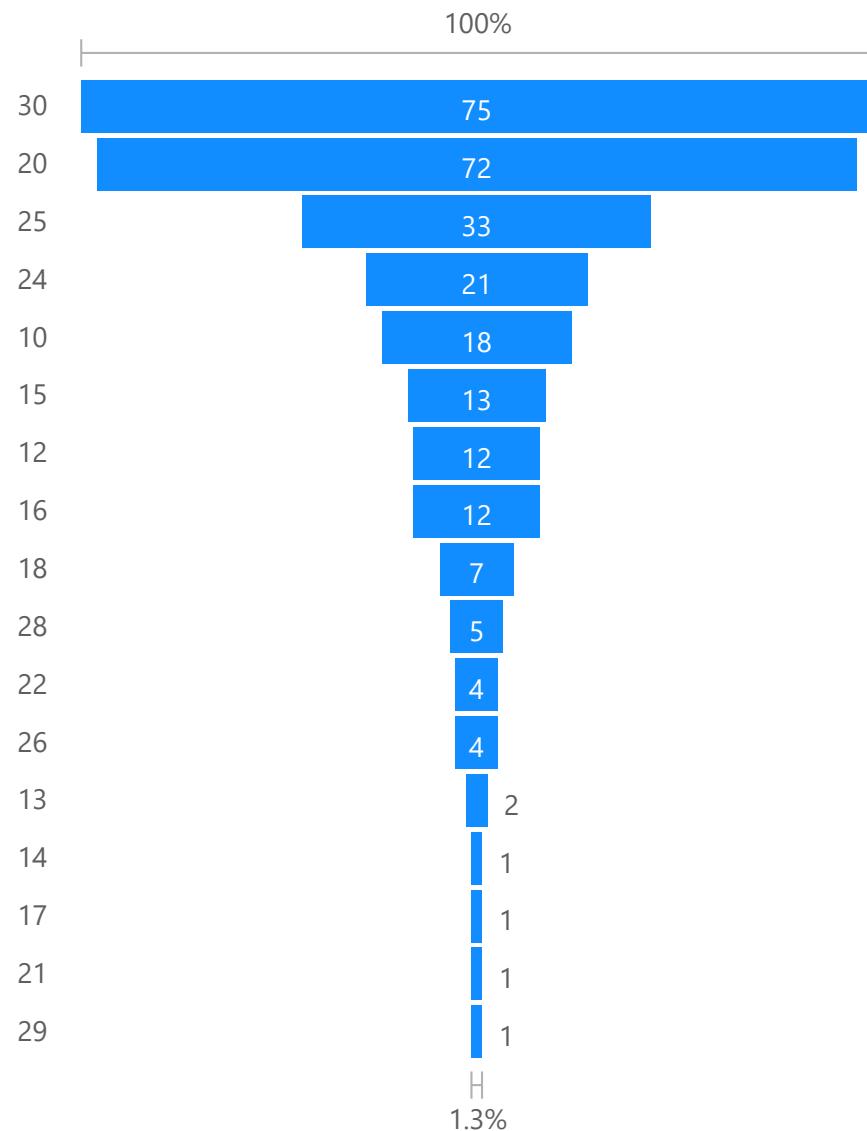


User Info Fact Count by User Info Fact.Age and Dim Location.Native Country

Dim Locatio... ● Canada



User Info Fact Count, First Dim Location.Native Country and First User Info Fact.Income by User Info Fact.Totalwork Hrs

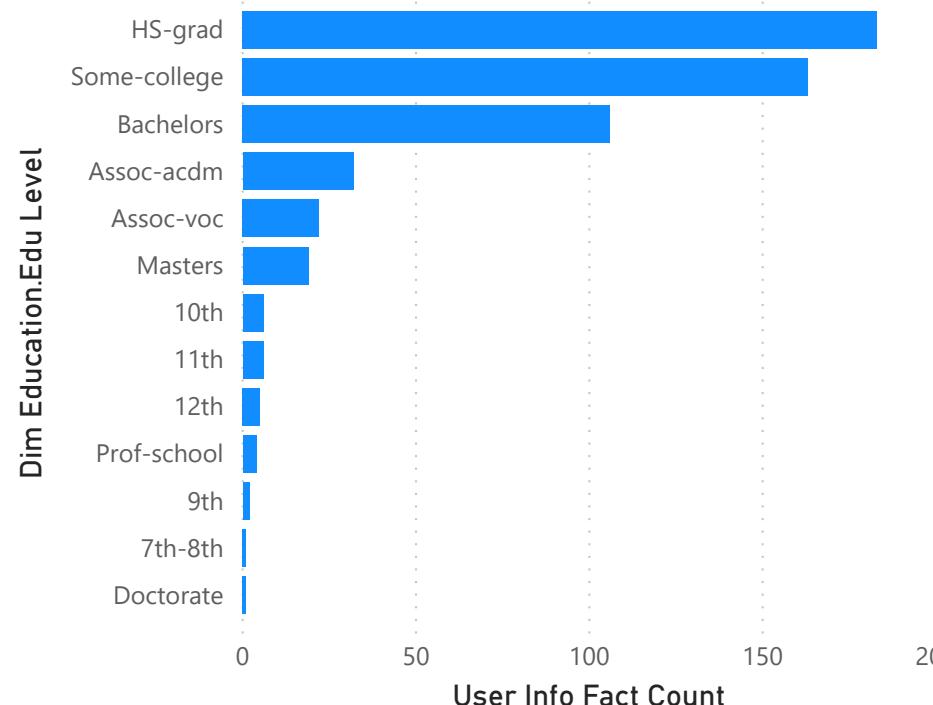


The people who have good wages rate in US her hour.
The presentation is made depending on variables such as location and the working hour of users.

This presentation shows an interesting insight. There is a good number of people having income more than 50000 USD an annum who are woking 20 hours a week and the number of people are for the same income woking from 21-29 hours are less , especially US has a very few number of people having the same income who is woking for 29 hours a week.
(Solution to business Query 4)

User Info Fact Count by Dim Education.Edu Level and Dim Sector.Sector

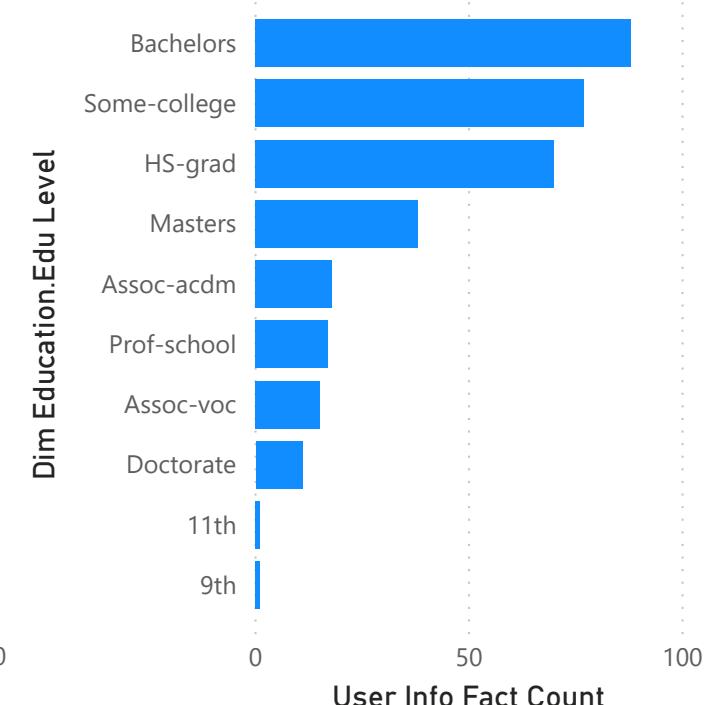
Dim Sector.Sector ● Federal-gov



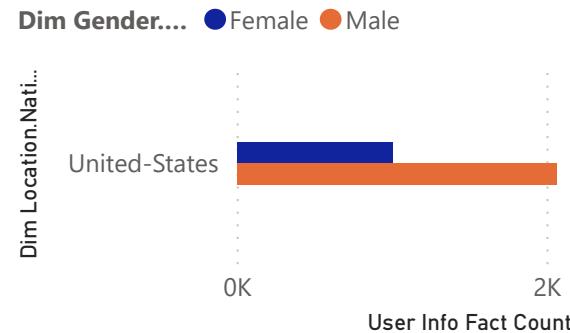
The two charts of US for education level shows,the US is having the maximum number of high school graduates working under federal govt who earn 50000 USD a year(Maximum) and the people who are woking for federal govt and getting at least 50000 USD a year most of them are Bachelor degree holders. (Solution to business query 5)

User Info Fact Count by Dim Education.Edu Level and Dim Sector.Sector

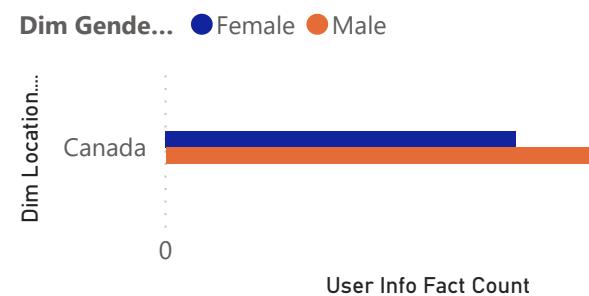
Dim Sect... ● Federal-gov



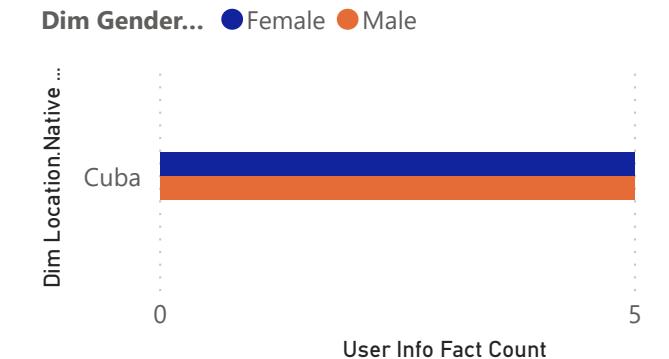
User Info Fact Count by Dim Location.Native Country and Dim Gender.Gender



User Info Fact Count by Dim Location.Native Country and Dim Gender.Gender

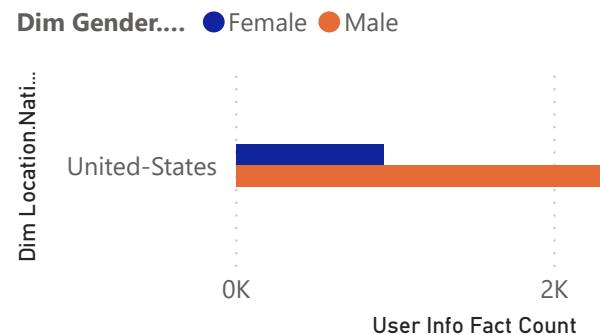


User Info Fact Count by Dim Location.Native Country and Dim Gender.Gender

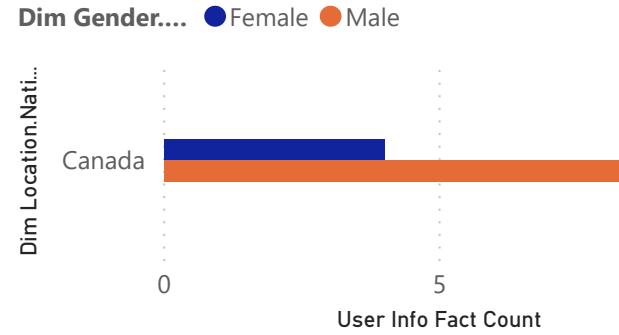


(Solution to business Query
6)

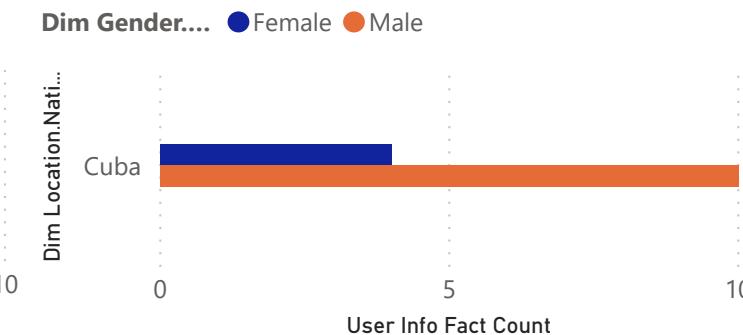
User Info Fact Count by Dim Location.Native Country and Dim Gender.Gender



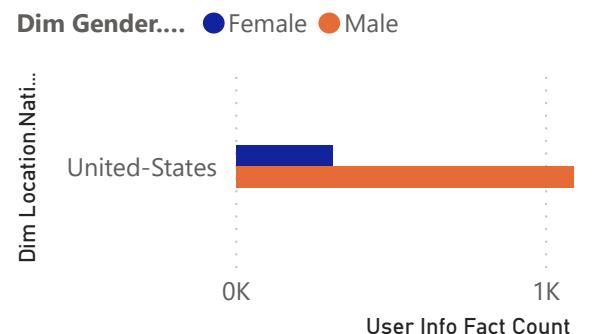
User Info Fact Count by Dim Location.Native Country and Dim Gender.Gender



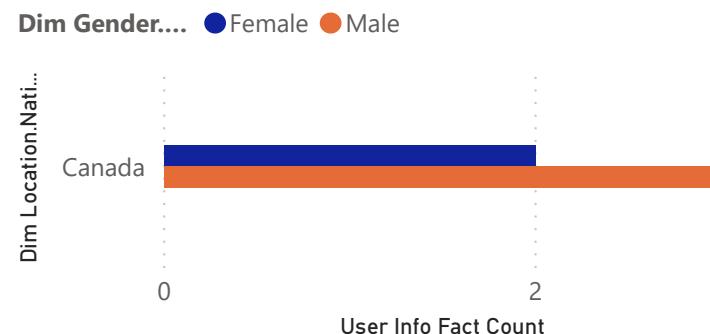
User Info Fact Count by Dim Location.Native Country and Dim Gender.Gender



User Info Fact Count by Dim Location.Native Country and Dim Gender.Gender



User Info Fact Count by Dim Location.Native Country and Dim Gender.Gender



User Info Fact Count by Dim Location.Native Country and Dim Gender.Gender

