



**DEPARTMENT OF MANAGEMENT SCIENCE AND TECHNOLOGY**

**MASTER OF SCIENCE IN BUSINESS ANALYTICS**

**Lesson:** Data Management and Business Intelligence

**Professor:** D. Chatziantoniou

**Assignment:** 2

**Editors:** Vogiatzis George, Zourou Myrsini

**AM:** P2821827, P2821828

**Date:** Monday 17 December 2018

## Table of Contents

1. Data Description .....	3
2. Description of the fact table and the dimensions.....	9
2.1 Deletions and alterations of our dataset.....	10
2.2 Creation of fact and dimension tables in sql-server management studio. ....	13
2.3 Visual Studio – Cube Creation .....	18
2.4 Olap Reports .....	27
3. TABLEAU.....	29

## 1. Data Description

Road traffic accidents (RTAs) have emerged as an important public health issue, which needs to be tackled by a multi-disciplinary approach. The trend in RTA injuries and death is becoming alarming in many countries around the world. The number of fatal and disabling road accident is increasing day by day. Therefore, it is a real public health challenge for all the concerned agencies to prevent it. The approach to implement the rules and regulations available to prevent road accidents is often ineffective and half-hearted. Awareness creation, strict implementation of traffic rules, and scientific engineering measures can prevent this public health catastrophe. The dataset, which selected as part for this assignment, is relative to this serious social issue. Because contains informations about the accidents which take place in New York City about the last 6 six years (2012-2018). Specifically contains informations about the date, the time and the location of each accident, the number of the vehicles, which was collided, the numbers of the injured or death people per accident and the type of vehicle that take part in the accident. Also gives information about the reasons that led in the accident. The original form of the dataset was in NYOpenData and it is shown on the table below:

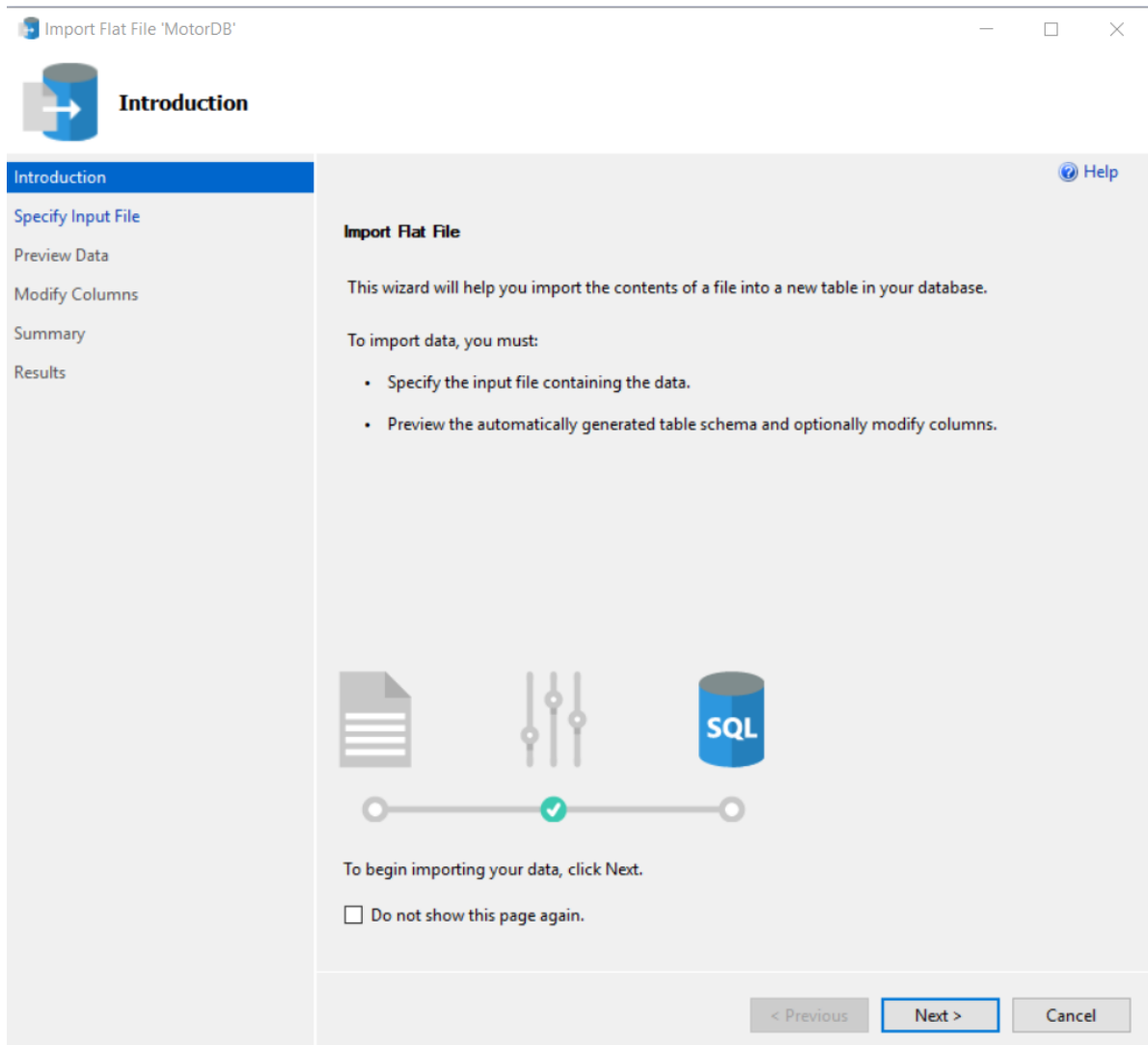
DATE ↓	TIME	BOROUGH	ZIP CODE	LATITUDE	LONGITU...	LOCATION	ON STREET NAME	CROSS STREET NAME
11/16/2018	0:10	MANHATTAN	10010	40.742275	-73.988914	(40.742275°, -73.988...	5 AVENUE	BROADWAY
11/16/2018	0:45	BROOKLYN	11211	40.710197	-73.95843	(40.710197°, -73.958...	BORINQUEN PLACE	HAVEMEYER STREET
11/16/2018	0:40			40.76272	-73.72817	(40.76272°, -73.7281...	LONG ISLAND EXPRESSWAY	
11/16/2018	1:00	BRONX	10454	40.803555	-73.91184	(40.803555°, -73.911...	EAST 137 STREET	WILLOW AVENUE
11/16/2018	1:00	BROOKLYN	11221	40.694923	-73.915565	(40.694923°, -73.915...	WILSON AVENUE	PALMETTO STREET
11/16/2018	1:20			40.8047	-73.91243	(40.8047°, -73.91243°)	EAST 138 STREET	BRUCKNER BOULEVARD
11/16/2018	1:41			40.737682	-73.85206	(40.737682°, -73.852...	108 STREET	HORACE HARDING EXPRESSWAY
11/16/2018	2:50	BRONX	10451	40.819057	-73.92923	(40.819057°, -73.929...	EAST 149 STREET	GERARD AVENUE
11/16/2018	4:22			40.81175	-73.93144	(40.81175°, -73.9314...	MAJOR DEEGAN EXPRESSWAY	
11/16/2018	5:30	QUEENS	11417	40.67887	-73.83419	(40.67887°, -73.8341...	ROCKAWAY BOULEVARD	CENTREVILLE AVENUE

Through this analysis, we want to identify the most often reason which leads to an accident, the type of the cars which take part in an accident and to show the variation of the number of deaths or injuries by car accidents change through the years. Furthermore, we will specify which area of the following Bronx, Brooklyn, Manhattan, Queens or Staten Island has the more accidents. The results of our analysis could be used from the local authorities in order to define the tax for the violations of the road traffic code and the days, which happens the most accidents in order to be on hand.

The initial entries of our dataset were 1,385,920 lines and the dataset consisted of 29 columns in total, but after the ETL processes, the dataset ended up having 918,846 entries. We load the dataset into a local database ('ny\_acc') through the Import Flat File option from the Microsoft SQL Server Management Studio. We chose datatype nvarchar(100) null for all the columns except for the metric ones which were imported as int, to avoid any insertion errors and we proceeded with the alteration of some types later. The steps, which follow to insert the data, were the following:

## Step 1

The first page of the wizard is the welcome page



## Step 2

We click the browse button to select the file of our dataset and after we give the new table name, which will contain our dataset.

Import Flat File 'ny\_acc'

**Specify Input File**

Introduction  
**Specify Input File**  
Preview Data  
Modify Columns  
Summary  
Results

**Specify Input File**  
This operation will create a table from your input file.

Location of file to be imported  
C:\Users\George\Desktop\NYPD\_Motor\_Vehicle\_Collisions.csv Browse...

New table name:  
ny\_collisions


Table schema:  
dbo

< Previous Next > Cancel

### Step 3

The wizard generates a preview where you can view for the first 50 rows of the dataset.

Import Flat File 'ny\_acc'

**Preview Data**

Introduction

Specify Input File

**Preview Data**

Modify Columns

Summary

Results

Preview Data

This operation analyzed the input file structure to generate the preview below for up to the first 50 rows.

DATE	TIME	BOROUGH	ZIP_CODE	LATITUDE	LONGIT
11/16/2018	0:10	MANHATTAN	10010	40.742275	-73.9889
11/16/2018	0:40			40.76272	-73.7281
11/16/2018	0:45	BROOKLYN	11211	40.710197	-73.9584
11/16/2018	1:00	BRONX	10454	40.803555	-73.9118
11/16/2018	1:00	BROOKLYN	11221	40.694923	-73.9155
11/16/2018	1:20			40.8047	-73.9124
11/16/2018	1:41			40.737682	-73.8520
11/16/2018	2:50	BRONX	10451	40.819057	-73.9292
11/16/2018	4:22			40.81175	-73.9314
11/16/2018	5:30	QUEENS	11417	40.67887	-73.8341
11/16/2018	5:45	MANHATTAN	10128	40.77918	-73.9508
11/16/2018	6:00			40.74315	-73.8319
11/16/2018	6:02	BROOKLYN	11206	40.7	-73.9406
11/16/2018	6:10	MANHATTAN	10128	40.78169	-73.9489
11/15/2018	0:00	BROOKLYN	11228	40.620872	-74.0022
11/15/2018	0:00	QUEENS	11004	40.752037	-73.7211
11/15/2018	0:00	QUEENS	11368	40.74989	-73.8625
11/15/2018	0:00	QUEENS	11375	40.71614	-73.8335
11/15/2018	0:00			40.64927	-74.0096
11/15/2018	0:00			40.610256	-74.0073

< Previous

Next >

Cancel

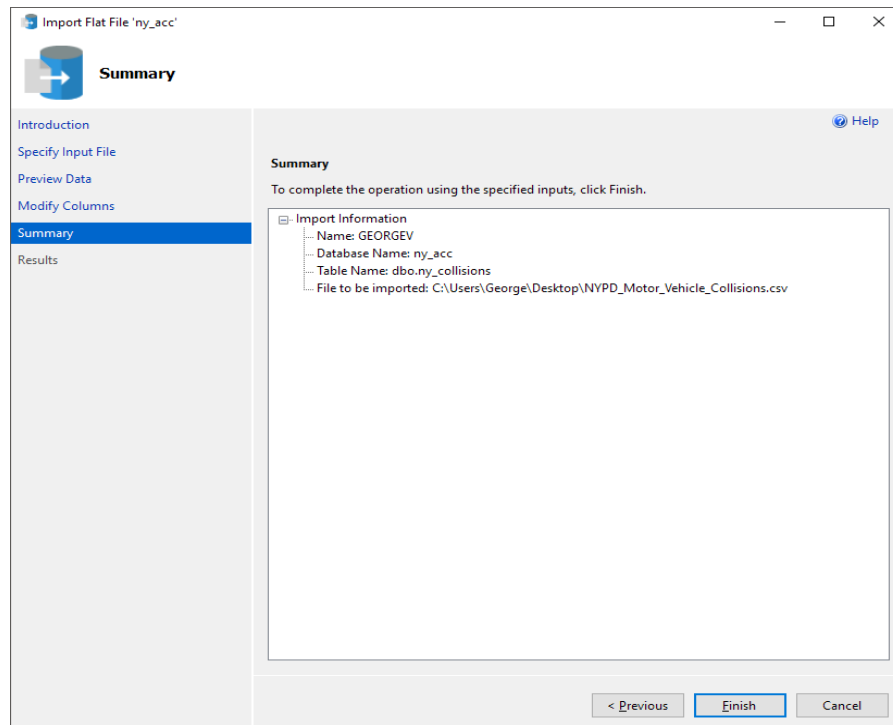
## Step 4

In this step, the wizard identifies what data type believes that describe better our variables. However, we proceeded with some alterations in the proposed data types in order to insert the data with the right variables types, names, data types, etc

Column Name	Data Type	Primary Key	Allow Nulls
DATE	nvarchar(100)	<input type="checkbox"/>	<input checked="" type="checkbox"/>
TIME	nvarchar(100)	<input type="checkbox"/>	<input checked="" type="checkbox"/>
BOROUGH	nvarchar(50)	<input type="checkbox"/>	<input checked="" type="checkbox"/>
ZIP_CODE	nvarchar(100)	<input type="checkbox"/>	<input checked="" type="checkbox"/>
LATITUDE	float	<input type="checkbox"/>	<input checked="" type="checkbox"/>
LONGITUDE	float	<input type="checkbox"/>	<input checked="" type="checkbox"/>
LOCATION	nvarchar(100)	<input type="checkbox"/>	<input checked="" type="checkbox"/>
ON_STREET_NAME	nvarchar(100)	<input type="checkbox"/>	<input checked="" type="checkbox"/>
CROSS_STREET_NAME	nvarchar(100)	<input type="checkbox"/>	<input checked="" type="checkbox"/>
OFF_STREET_NAME	nvarchar(100)	<input type="checkbox"/>	<input checked="" type="checkbox"/>
NUMBER_OF_PERSONS_INJURED	int	<input type="checkbox"/>	<input checked="" type="checkbox"/>
NUMBER_OF_PERSONS_KILLED	int	<input type="checkbox"/>	<input checked="" type="checkbox"/>
NUMBER_OF_PEDESTRIANS_INJURED	int	<input type="checkbox"/>	<input checked="" type="checkbox"/>
NUMBER_OF_PEDESTRIANS_KILLED	int	<input type="checkbox"/>	<input checked="" type="checkbox"/>
NUMBER_OF_CYCLIST_INJURED	int	<input type="checkbox"/>	<input checked="" type="checkbox"/>
NUMBER_OF_CYCLIST_KILLED	int	<input type="checkbox"/>	<input checked="" type="checkbox"/>
NUMBER_OF_MOTORIST_INJURED	int	<input type="checkbox"/>	<input checked="" type="checkbox"/>
NUMBER_OF_MOTORIST_KILLED	int	<input type="checkbox"/>	<input checked="" type="checkbox"/>
CONTRIBUTING_FACTOR_VEHICLE_1	nvarchar(100)	<input type="checkbox"/>	<input checked="" type="checkbox"/>
CONTRIBUTING_FACTOR_VEHICLE_2	nvarchar(100)	<input type="checkbox"/>	<input checked="" type="checkbox"/>
CONTRIBUTING_FACTOR_VEHICLE_3	nvarchar(100)	<input type="checkbox"/>	<input checked="" type="checkbox"/>
CONTRIBUTING_FACTOR_VEHICLE_4	nvarchar(100)	<input type="checkbox"/>	<input checked="" type="checkbox"/>
CONTRIBUTING_FACTOR_VEHICLE_5	nvarchar(100)	<input type="checkbox"/>	<input checked="" type="checkbox"/>
UNIQUE_KEY	ntext	<input type="checkbox"/>	<input checked="" type="checkbox"/>
VEHICLE_TYPE_CODE_1	nvarchar(50)	<input type="checkbox"/>	<input checked="" type="checkbox"/>
VEHICLE_TYPE_CODE_2	nvarchar(50)	<input type="checkbox"/>	<input checked="" type="checkbox"/>
VEHICLE_TYPE_CODE_3	nvarchar(50)	<input type="checkbox"/>	<input checked="" type="checkbox"/>
VEHICLE_TYPE_CODE_4	nvarchar(50)	<input type="checkbox"/>	<input checked="" type="checkbox"/>
VEHICLE_TYPE_CODE_5	nvarchar(50)	<input type="checkbox"/>	<input checked="" type="checkbox"/>

## Step 5

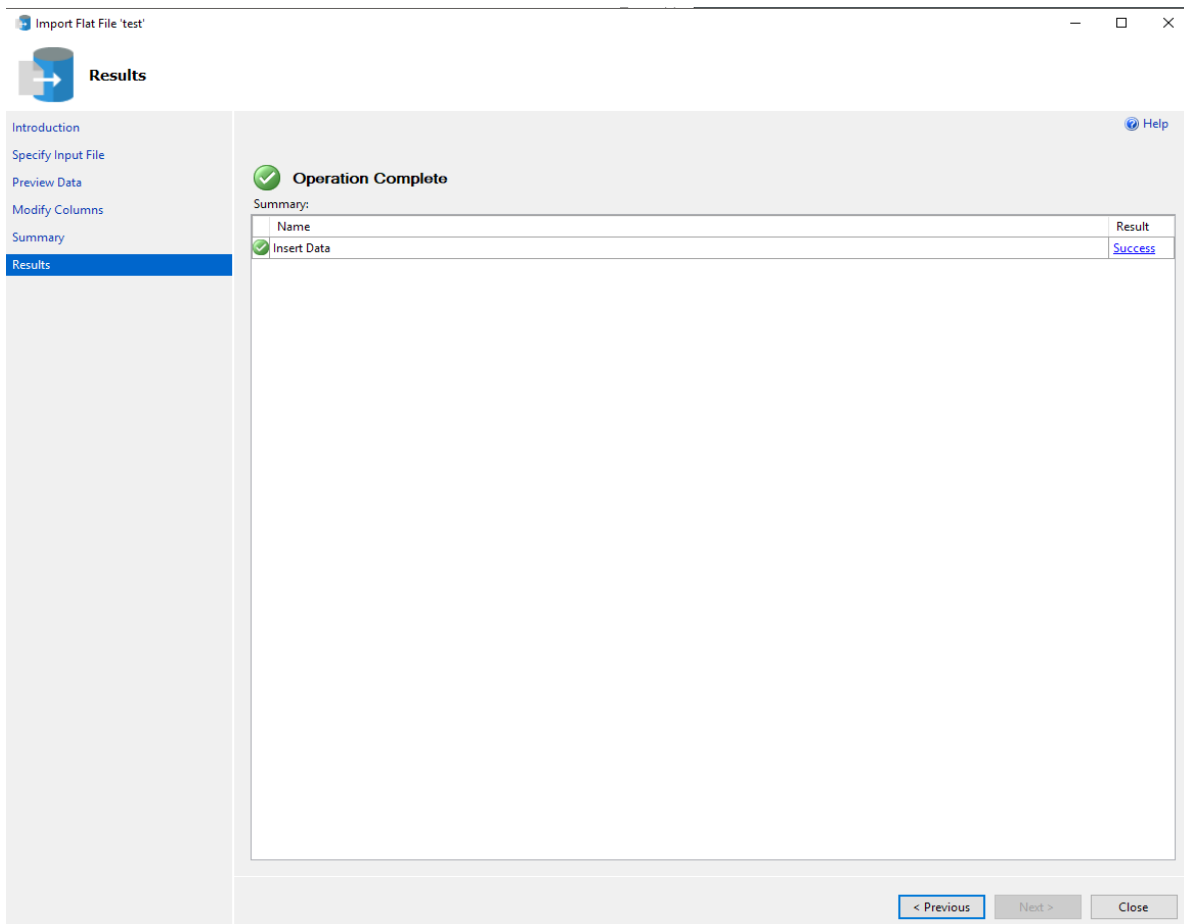
This is simply a summary page displaying our current configuration. If there are issues, we should go back to previous sections to alter the problems, which appears. Otherwise, we click the finish bottom to attempt the import process.



## Step 6

This page indicates whether the import was successful. If a green check mark appears, it was a success, otherwise you may need to review our configuration or input file for any errors.





## 2. Description of the fact table and the dimensions

After the successful import of the dataset, we continued with the determination of the fact table and the dimensions tables, which are connected. Our dataset as we already mentioned, refers to traffic accidents occurred in New York City the last 6 years. Fact table will represent all the information about each accident and the dimensions table will give further information about some variables. Our dimensions are the following tables:

- Time of accident
- Date of accident
- Location of accident
- Factors of accident
- Type of vehicles

In the next section, we will analyze in a more descriptive way the configuration of the fact table and of the dimensions. Just before, we created the schema of our database we had to clean the data in order to be in a desired form and organize them properly. In the next section, we explain how we do the cleansing of our data and their configuration.

## 2.1 Deletions and alterations of our dataset.

When we inserted the flat file on the server we had to define almost all datatypes as nvarchar or integer (if they were measures) to avoid getting insertion errors. Therefore, after this procedure we had to change some datatypes in order to represent better their fields.

We have successfully load the dataset in the database with the name Collisions. After that, we searched if our data had some unreasonable values. Specifically every line of our dataset should describe one vehicle collision.

Firstly, we change the empty cells of the variables with NULL and we continue with the deletion of NULLS for the following variable values: Longitude, Latitude, Zip Code, contributing factor vehicle and type of the vehicle because we want to keep the lines, which have all the information. Secondly, we delete some observations that do not make sense and maybe are wrong entries. For example, latitude and longitude equal to zero should be deleted because we know that our data refers to New York City and this coordination represent a position in Atlantic Ocean. In addition, we delete the observations where the contributing factor vehicle was equal to 1 and 80 because these observations did not reveal any reason why the accident happened. Below we give the query, which we use to do the aforementioned deletions.

### Query1:

```
alter table [ny_acc].[dbo].[ny_collisions] alter column UNIQUE_KEY nvarchar(50);

update [ny_acc].[dbo].[ny_collisions] set ZIP_CODE = null where ZIP_CODE = ' ';
update [ny_acc].[dbo].[ny_collisions] set ON_STREET_NAME = null where ON_STREET_NAME = ' ';

delete from [ny_acc].[dbo].[ny_collisions] where CONTRIBUTING_FACTOR_VEHICLE_1 is Null;

delete from [ny_acc].[dbo].[ny_collisions] where VEHICLE_TYPE_CODE_1 is Null
and VEHICLE_TYPE_CODE_2 is Null
and VEHICLE_TYPE_CODE_3 is Null
and VEHICLE_TYPE_CODE_4 is Null
and VEHICLE_TYPE_CODE_5 is Null

delete from [ny_acc].[dbo].[ny_collisions] where LOCATION is Null;
delete from [ny_acc].[dbo].[ny_collisions] where ZIP_CODE is null;
delete from [ny_acc].[dbo].[ny_collisions] where LATITUDE = 0 and LONGITUDE = 0;

delete from [ny_acc].[dbo].[ny_collisions]
where CONTRIBUTING_FACTOR_VEHICLE_1 = '80'
or CONTRIBUTING_FACTOR_VEHICLE_2 = '80'
or CONTRIBUTING_FACTOR_VEHICLE_3 = '80'
or CONTRIBUTING_FACTOR_VEHICLE_4 = '80'
or CONTRIBUTING_FACTOR_VEHICLE_5 = '80'
or CONTRIBUTING_FACTOR_VEHICLE_1 = '1'
or CONTRIBUTING_FACTOR_VEHICLE_2 = '1'
or CONTRIBUTING_FACTOR_VEHICLE_3 = '1'
or CONTRIBUTING_FACTOR_VEHICLE_4 = '1'
or CONTRIBUTING_FACTOR_VEHICLE_5 = '1'

alter table [ny_acc].[dbo].[ny_collisions] drop column OFF_STREET_NAME;
alter table [ny_acc].[dbo].[ny_collisions] drop column LOCATION;
alter table [ny_acc].[dbo].[ny_collisions] drop column CROSS_STREET_NAME;
```

Secondly, we delete the column Location of our dataset because it contains the same information with the columns Longitude and Latitude. Furthermore, we delete the columns with the name 'OFF\_STREET\_NAME' and 'CROSS\_STREET\_NAME', because most of the rows were empty and generally the information provided by

these fields is unnecessary for our analysis. Below we can see the query we used to do the aforementioned deletions.

## Query 2:

```
alter table [collisions].[dbo].[Collisions] drop column OFF_STREET_NAME;
alter table [collisions].[dbo].[Collisions] drop column LOCATION;
alter table [collisions].[dbo].[Collisions] drop column CROSS_STREET_NAME;
```

Furthermore, we altered the type of the 'TIME' variable, which we inserted as varchar. We transform it in time type because it will suit better our analysis tools. The last step of the deletion process was to alter the entries in the columns CONTRIBUTING\_FACTOR\_VEHICLE and VEHICLE\_TYPE\_CODE because there were categories which described the same entity but with different literal so we altered the names in order to group in the same category. For example, the VEHICLE\_TYPE\_CODE in the beginning had 380 unique categories and after the process, we end up with 46 unique categories. We alter the names which described the same car type but with the abbreviation and not with the total name. In addition, we added many vehicle types in the category unknown because it was not clear in which vehicle type they refer to. Below we can see the queries, we executed concerning the CONTRIBUTING\_FACTOR\_VEHICLE records and the first column of vehicle type from a total of five columns. However, we executed the same query for the rest vehicles' contributing factor records.

## Query 4:

### CONTRIBUTING\_FACTOR\_VEHICLE

```
update [collisions].[dbo].[Collisions] set CONTRIBUTING_FACTOR_VEHICLE_1 = 'Illness'
where CONTRIBUTING_FACTOR_VEHICLE_1 = 'Illnes'
update [collisions].[dbo].[Collisions] set CONTRIBUTING_FACTOR_VEHICLE_2 = 'Illness'
where CONTRIBUTING_FACTOR_VEHICLE_2 = 'Illnes'

UPDATE [collisions].[dbo].[Collisions] SET CONTRIBUTING_FACTOR_VEHICLE_1 = 'Reaction to Uninvolved Vehicle'
where CONTRIBUTING_FACTOR_VEHICLE_1 in ('Reaction to Other Uninvolved Vehicle');
UPDATE [collisions].[dbo].[Collisions] SET CONTRIBUTING_FACTOR_VEHICLE_2 = 'Reaction to Uninvolved Vehicle'
where CONTRIBUTING_FACTOR_VEHICLE_2 in ('Reaction to Other Uninvolved Vehicle');
UPDATE [collisions].[dbo].[Collisions] SET CONTRIBUTING_FACTOR_VEHICLE_3 = 'Reaction to Uninvolved Vehicle'
where CONTRIBUTING_FACTOR_VEHICLE_3 in ('Reaction to Other Uninvolved Vehicle');
UPDATE [collisions].[dbo].[Collisions] SET CONTRIBUTING_FACTOR_VEHICLE_4 = 'Reaction to Uninvolved Vehicle'
where CONTRIBUTING_FACTOR_VEHICLE_4 in ('Reaction to Other Uninvolved Vehicle');
UPDATE [collisions].[dbo].[Collisions] SET CONTRIBUTING_FACTOR_VEHICLE_5 = 'Reaction to Uninvolved Vehicle'
where CONTRIBUTING_FACTOR_VEHICLE_5 in ('Reaction to Other Uninvolved Vehicle');
```

### VEHICLE\_TYPE\_CODE

```
UPDATE [collisions].[dbo].[Collisions] SET VEHICLE_TYPE_CODE_1 = 'AMBULANCE'
where VEHICLE_TYPE_CODE_1 in ('Ambulance', 'AMBUL', 'AM', 'Ambul', 'Ambul', 'ambul', 'AMB', 'AMBU', 'ambu', 'AMBULANCEANCE', 'ABULA');
UPDATE [collisions].[dbo].[Collisions] SET VEHICLE_TYPE_CODE_1 = 'DELIV'
where VEHICLE_TYPE_CODE_1 in ('deliv', 'Deliv', 'delv', 'DELV', 'DELV');
UPDATE [collisions].[dbo].[Collisions] SET VEHICLE_TYPE_CODE_1 = 'ARMORED TRUCK'
where VEHICLE_TYPE_CODE_1 in ('armor', 'ARMOR', 'AR', 'Armored Truck');
UPDATE [collisions].[dbo].[Collisions] SET VEHICLE_TYPE_CODE_1 = 'EBIKE'
where VEHICLE_TYPE_CODE_1 in ('E BIK', 'E-BIK', 'EBIKE', 'e-bik', 'E SCO', 'e sco', 'elect', 'ELECT');
UPDATE [collisions].[dbo].[Collisions] SET VEHICLE_TYPE_CODE_1 = 'VAN' where VEHICLE_TYPE_CODE_1 in
('VAN', 'van', 'REFRI', 'Refrigerated Van', 'RF', 'MINIV', 'GMC V', 'Van Camper', 'VAN T', 'VAN/B', 'VAN/T', 'Vanete', 'VN', 'Vanette');
UPDATE [collisions].[dbo].[Collisions] SET VEHICLE_TYPE_CODE_1 = 'UTILITY VEHICLES'
where VEHICLE_TYPE_CODE_1 in ('UTILI', 'Utili', 'utili', 'UTIL', 'UT', 'Util', 'BOB C', 'bobca', 'uliti', 'BR', 'BROOM', 'Bucke', 'bulld', 'BULLD',
'CRANE', 'UTLL', 'cemen', 'CMIX', 'Concrete Mixer', 'CONST', 'Hopper', 'JCB40', 'JLG B', 'Lift', 'Lift Boom', 'Lunch Wagon', 'CAT', 'cat 3', 'CATER',
'Pallet', 'scaff', 'utlit', 'FORK', 'Forkl', 'FORK-', 'forkl', 'FORKL', 'FORKLIFT', 'Well Driller');
UPDATE [collisions].[dbo].[Collisions] SET VEHICLE_TYPE_CODE_1 = 'UNKNOWN' where VEHICLE_TYPE_CODE_1 in ('Unkno', 'unkno',
'unk', 'UKN', 'OTHER', 'N/A', 'NA', 'ND', 'Enclosed Body - Removable Enclosure', 'stree', '(ceme', '00', '013', '1',
'15 Pa', '18 WH', '1S', '315 e', '994', 'ACCES', 'B5-44', 'BA', 'BACK', 'BACKH', 'BS', 'CASE', 'CB', 'CHERR', 'DUNBA', 'e com', 'E ONE', 'E PAS', 'E-MOT', 'EMRGN', 'EMS',
'EMS H', 'EN', 'ENGIN', 'f550', 'FORTL', 'FREE', 'FREIG', 'frieg', 'FRONT', 'G Sem', 'Glass Rack', 'GN', 'GR', 'H/WH', 'Hand', 'HELP', 'HIGHL', 'HI-LO', 'HINO', 'HO',
'HOTDO', 'HWY C', 'Inter', 'IP', 'JOHN', 'LIGHT', 'LL', 'LP', 'RYDER', 'ROAD', 'RGS', 'REP', 'Unk', 'Unkn', 'UNKON', 'Comix', 'CHART', 'Carri', 'east', 'GAS S',
'GE/SC', 'LF', 'LW', 'MACK', 'MAN L', 'MARK', 'Marke', 'MB', 'mcy', 'MD', 'Mecha', 'MH', 'MILLI', 'MK', 'NEH Y', 'NS AM', 'NYC A', 'nyc d', 'NYC F',
'OBJEC', 'QWR', 'OP', 'Open Body', 'PL', 'PM', 'Porta', 'POWER', 'PSD', 'PUMP', 'SC', 'RD/S', 'red', 'refg', 'RENTA', 'RESCU', 'seagr', 'SELF', 'SELF-',
'small', 'SP', 'SPC', 'spc p', 'spec', 'Sprin', 'ST', 'ST150', 'STAK', 'Subn', 'SUBN', 'SUBUR', 'SWEEP', 'SWT', 'TL TR', 'Tln', 'TRAFF', 'TRL', 'TRLPM',
'TRLR', 'TT', 'U.S', 'Uber', 'UD', 'U-HAU', 'UHUAL', 'VC', 'VERIZ', 'WC', 'WD', 'WHEEL', 'WORK', 'WORKH', 'work', 'p/sh', 'P/SH', 'nyc a', 'NYC a', 'NYC A', 'EB',
'TN', 'red', 'U.S.');
```

```

UPDATE [collisions].[dbo].[Collisions] SET VEHICLE_TYPE_CODE_1 = 'PASSENGER VEHICLE' where
VEHICLE_TYPE_CODE_1 in ('PAS','pas','passa','Passe','Pavin','2 DOO','2 dr sedan','3D','3DC-','3-Door','4 dr sedan','4DSD','4wheel','Chevrn','City','CONV',
'Convertible','ford','GOLF','GOVER','jeep','nissa','Sedan','Smart','WHITE','YELLOW','SPORT UTILITY / STATION WAGON','Station Wagon/Sport Utility Vehicle',
'SPORT UTILITY / STATION WAGON','SPORT UTILITY/STATION WAGON','COUPE','GRAY','Hume','KENNO','MINI','Mini');
UPDATE [collisions].[dbo].[Collisions] SET VEHICLE_TYPE_CODE_1 = 'TRUCK'
where VEHICLE_TYPE_CODE_1 in ('Truck','truck','uhaul','BOOM','BOOML','BOX','Beverage Truck','trk','tk','Pick-up TRUCK','Flat Rack','Flat Bed',
'Pickup with mounted Camper','pick','PK','PICKU','Pick-','PICK-UP TRUCK','FLAT','flat','FLATB','ice c','icecr','box t','BOX T','boxtr','COM',
'COMER','COMME','COMM','Food','COM T','COM','COMME','CO','CH','HEAVY','LARGE COM VEH(6 OR MORE TIRES)','WAGON','TOW T','Tow T','Tow t','tow t',
'TOW','tow','SEMI-','SEMI','Semi','FB','CARGO','Cargo','Tow Truck','Tow Truck / Wrecker','TOW TRUCK','TR','2- to','2 TON','Box Truck','OIL T',
'Tract','TRACT','tract','TRAC','Tractor Truck Diesel','Tractor Truck Gasoline','Tow Truck','TRANS');
UPDATE [collisions].[dbo].[Collisions] SET VEHICLE_TYPE_CODE_1 = 'TRAILER' where VEHICLE_TYPE_CODE_1 in ('TRAI','trail','rv','Trail','trail','trail');
UPDATE [collisions].[dbo].[Collisions] SET VEHICLE_TYPE_CODE_1 = 'BICYCLE' where VEHICLE_TYPE_CODE_1 in ('mta b','MTA B','Bike','BK','Minibike','Minicycle');
UPDATE [collisions].[dbo].[Collisions] SET VEHICLE_TYPE_CODE_1 = 'TAXI' where VEHICLE_TYPE_CODE_1 in ('Taxi','CAB');
UPDATE [collisions].[dbo].[Collisions] SET VEHICLE_TYPE_CODE_1 = 'Tanker' where VEHICLE_TYPE_CODE_1 in ('TANK');
UPDATE [collisions].[dbo].[Collisions] SET VEHICLE_TYPE_CODE_1 = 'Subn' where VEHICLE_TYPE_CODE_1 in ('SUBN','subn');
UPDATE [collisions].[dbo].[Collisions] SET VEHICLE_TYPE_CODE_1 = 'SPC' where VEHICLE_TYPE_CODE_1 in ('spc p|Spc');
UPDATE [collisions].[dbo].[Collisions] SET VEHICLE_TYPE_CODE_1 = 'SNOW PLOW' where VEHICLE_TYPE_CODE_1 in ('SNOW','Snow Plow');
UPDATE [collisions].[dbo].[Collisions] SET VEHICLE_TYPE_CODE_1 = 'SCOOTER'
where VEHICLE_TYPE_CODE_1 in ('SCOOT','Scoot','scoot','SCOOTERER','Motorscooter','MS','mot s');
UPDATE [collisions].[dbo].[Collisions] SET VEHICLE_TYPE_CODE_1 = 'BUS'
where VEHICLE_TYPE_CODE_1 in ('bus','BUS','SCHOO','Schoo','schoo','School Busl Bus','School Busl Bus','BU','School Bus','omni','omnib');
UPDATE [collisions].[dbo].[Collisions] SET VEHICLE_TYPE_CODE_1 = 'GARBAGE TRUCK'
where VEHICLE_TYPE_CODE_1 in ('garba','GARBA','Garba','Garbage or Refuse','SANIT','sanit','Sanit','DS','dsny','Dump','DUMPT','dump','DUMP','Dumps');
UPDATE [collisions].[dbo].[Collisions] SET VEHICLE_TYPE_CODE_1 = 'FIRE TRUCK'
where VEHICLE_TYPE_CODE_1 in ('FIRE','Firet','fire','Fd fi','FIRET','Fire','FIRE TRUCK TRUCK','fdny','FDNY');
UPDATE [collisions].[dbo].[Collisions] SET VEHICLE_TYPE_CODE_1 = 'MINI' where VEHICLE_TYPE_CODE_1 in ('Mini');
UPDATE [collisions].[dbo].[Collisions] SET VEHICLE_TYPE_CODE_1 = 'FORD' where VEHICLE_TYPE_CODE_1 in ('ford');
UPDATE [collisions].[dbo].[Collisions] SET VEHICLE_TYPE_CODE_1 = 'FORKLIFT' where VEHICLE_TYPE_CODE_1 in ('FORK','Forkl','FORK-','forkl','FORKL');
UPDATE [collisions].[dbo].[Collisions] SET VEHICLE_TYPE_CODE_1 = 'MOTORCYCLE' where VEHICLE_TYPE_CODE_1 in ('MOTORCYCLE','motor','MOTOR','Motorbike');
UPDATE [collisions].[dbo].[Collisions] SET VEHICLE_TYPE_CODE_1 = 'FIRE TRUCK'
where VEHICLE_TYPE_CODE_1 in ('FIRE TRUCK TRUCK','fd tr','FD tr','FIRE','FIRET','FR','FD LA');
UPDATE [collisions].[dbo].[Collisions] SET VEHICLE_TYPE_CODE_1 = 'LIMO' where VEHICLE_TYPE_CODE_1 in ('LIMO','Limou');

```

We also altered the data type of column Time from nvarchar(100) to time(7) because the SSAS tool could handle this column in its native form, which is time.

```

alter table [ny_acc].[dbo].[ny_collisions] add coltime time;
UPDATE [ny_acc].[dbo].[ny_collisions] SET coltime = convert (time, time,8);
alter table [ny_acc].[dbo].[ny_collisions] alter column time nvarchar(100) NULL;
UPDATE [ny_acc].[dbo].[ny_collisions] SET TIME = null;
UPDATE [ny_acc].[dbo].[ny_collisions] SET time = coltime;
alter table [ny_acc].[dbo].[ny_collisions] alter column time time(7) NOT NULL;
alter table [ny_acc].[dbo].[ny_collisions] drop column coltime;

```

Apart from Time column we altered the type of 'Latitude' and 'Longitude' columns from float to Decimal(9,6) so that the bi tool could understand these two columns as geographical.

```

UPDATE [ny_acc].[dbo].[ny_collisions]
SET LATITUDE = CAST(LATITUDE as decimal(9, 6))
ALTER TABLE [dbo].[ny_collisions] ALTER COLUMN LATITUDE DECIMAL(9,6);

UPDATE [ny_acc].[dbo].[ny_collisions]
SET LONGITUDE = CAST(LONGITUDE as decimal(9, 6))
ALTER TABLE [dbo].[ny_collisions] ALTER COLUMN LONGITUDE DECIMAL(9,6);

```

## 2.2 Creation of fact and dimension tables in sql-server management studio.

We decided to use the star schema for the creation of our tables. In this section, we describe the creation of the dimension tables along with the steps followed for the creation of the fact table.

### Location:

We used the below queries to create the table with all the necessary features of Location dimension. The columns of our original dataset that could be included in the Location are the borough of the accident, the street name, the zip code and the latitude and longitude variables. But we observe that there is no unique key that can be set as primary for this dimension, so we used a group-by for all these characteristics that together make up a complete address and we transferred the results of this group-by in 'loc\_table' and created the address\_id. Finally, we assign the content of 'loc\_table' table in our Location dimension table.

```
--LOCATION

create table loc_table( borough nvarchar(100), on_street_name nvarchar(100), cross_street_name nvarchar(100),
zip_code nvarchar(100),longitude decimal(9,6),latitude decimal(9,6));

insert into loc_table (borough,on_street_name,zip_code,longitude,latitude)
select BOROUGH, ON_STREET_NAME, ZIP_CODE, LONGITUDE, LATITUDE FROM [ny_acc].[dbo].[ny_collisions]
group by BOROUGH, ON_STREET_NAME, ZIP_CODE, LONGITUDE, LATITUDE;

alter table loc_table
add location_id int identity(1,1);

select * from loc_table order by location_id;

create table Location(
location_id nvarchar(100),
borough nvarchar(100),
on_street_name nvarchar(100),
zip_code nvarchar(100),
longitude decimal(9,6) ,
latitude decimal(9,6)
PRIMARY KEY(location_id));

insert into Location(location_id,borough,on_street_name,zip_code,longitude,latitude)
select location_id, borough,on_street_name,zip_code,longitude,latitude from loc_table;
```

### Date:

For the Date dimension we used the date column of our original dataset. To achieve this we created a temp table where we inserted the distinct values of dates present in the dataset and afterwards we add the Date\_id to represent the primary key. Finally, we assigned the content of the temp table in our Date dimension table where we set the Date\_id as a primary key.



```
--DATE

create table Accident_Date1(Acc_Date nvarchar(100));

insert into Accident_Date1(Acc_Date) select distinct(DATE) from [ny_acc].[dbo].[ny_collisions];
alter table Accident_Date1 add Date_id int identity(1,1);

create table Accident_Date (
Date_id nvarchar(100),
Acc_Date nvarchar(100)
primary key(Date_id)
);

insert into Accident_Date(Date_id,Acc_Date)select Date_id, Acc_Date from Accident_Date1;
```

### **Time:**

Regarding the time dimension table, we created it in the same manner with Date table, simply by inserting the distinct values of time we found in the original dataset, and then assigning a time id as primary key.

```
--TIME

create table time1(Acc_Time time(7));

insert into time1(Acc_Time) select distinct(TIME) from [ny_acc].[dbo].[ny_collisions];
alter table time1 add Time_id int identity(1,1);

create table Accident_Time (
Time_id nvarchar(100),
Acc_Time time(7)
primary key(Time_id)
);

insert into Accident_Time(Time_id,Acc_Time)select Time_id, Acc_Time from time1;
```

### **Accident factor:**

This table is about the factor that caused the accident. To populate this table we had to select the distinct values from the 5 different contributing factor columns of our original dataset.

```

--Accident Factor

create table factor( factor nvarchar(100));

insert into factor(factor)
select factor from (
    select CONTRIBUTING_FACTOR_VEHICLE_1 as factor from [ny_acc].[dbo].[ny_collisions] where CONTRIBUTING_FACTOR_VEHICLE_1 is not null
    union all select CONTRIBUTING_FACTOR_VEHICLE_2 from [ny_acc].[dbo].[ny_collisions] where CONTRIBUTING_FACTOR_VEHICLE_2 is not null
    union all select CONTRIBUTING_FACTOR_VEHICLE_3 from [ny_acc].[dbo].[ny_collisions] where CONTRIBUTING_FACTOR_VEHICLE_3 is not null
    union all select CONTRIBUTING_FACTOR_VEHICLE_4 from [ny_acc].[dbo].[ny_collisions] where CONTRIBUTING_FACTOR_VEHICLE_4 is not null
    union all select CONTRIBUTING_FACTOR_VEHICLE_5 from [ny_acc].[dbo].[ny_collisions] where CONTRIBUTING_FACTOR_VEHICLE_5 is not null
) t1 group by factor order by factor;

alter table factor add factor_id int identity(1,1);

create table Accident_factor(
    factor_id nvarchar(100),
    factor nvarchar(100)
    primary key(factor_id));

insert into Accident_factor (factor_id, factor)
select factor_id, factor from factor;

```

### **Vehicle type:**

The specific table is about the vehicle types observed in all these accidents. It was created in the same manner as the accident\_factor. Therefore, we had to perform a union on all vehicle type codes to extract the distinct values, which will be in the vehicle\_type dimension table.

```

--vehicle type code
create table vtype( vtype nvarchar(100));

insert into vtype(vtype)
select vtype from (
    select VEHICLE_TYPE_CODE_1 as vtype from [ny_acc].[dbo].[ny_collisions] where VEHICLE_TYPE_CODE_1 is not null
    union all select VEHICLE_TYPE_CODE_2 from [ny_acc].[dbo].[ny_collisions] where VEHICLE_TYPE_CODE_2 is not null
    union all select VEHICLE_TYPE_CODE_3 from [ny_acc].[dbo].[ny_collisions] where VEHICLE_TYPE_CODE_3 is not null
    union all select VEHICLE_TYPE_CODE_4 from [ny_acc].[dbo].[ny_collisions] where VEHICLE_TYPE_CODE_4 is not null
    union all select VEHICLE_TYPE_CODE_5 from [ny_acc].[dbo].[ny_collisions] where VEHICLE_TYPE_CODE_5 is not null
) t1 group by vtype order by vtype;

alter table vtype add vtype_id int identity(1,1);

create table Vehicle_type(
    vtype_id nvarchar(100),
    vtype nvarchar(100)
    primary key(vtype_id));

insert into Vehicle_type (vtype_id, vtype)
select vtype_id, vtype from vtype;

```

### **Fact table:**

Finally, to complete the star schema, we have to create our facttable, which is about accidents and has as foreign keys the primary keys of our dimension tables. As metrics, we defined the persons\_injured, persons\_killed, cyclists\_injured, cyclists\_killed, pedestrians\_injured, pedestrians\_killed, motorists\_injured and motorists\_killed, which represent the populations of victims in these accidents. To create our fact table we had to alter our original table containing the dataset (ny\_collisions) and add the id columns of our dimensions.

```

[+] alter table [ny_acc].[dbo].[ny_collisions]
    add Date_id nvarchar(100);

[+] alter table [ny_acc].[dbo].[ny_collisions]
    add Time_id nvarchar(100);

[+] alter table [ny_acc].[dbo].[ny_collisions]
    add Factor_id1 nvarchar(100);
[+] alter table [ny_acc].[dbo].[ny_collisions]
    add Factor_id2 nvarchar(100);
[+] alter table [ny_acc].[dbo].[ny_collisions]
    add Factor_id3 nvarchar(100);
[+] alter table [ny_acc].[dbo].[ny_collisions]
    add Factor_id4 nvarchar(100);
[+] alter table [ny_acc].[dbo].[ny_collisions]
    add Factor_id5 nvarchar(100);

```

```

[+] alter table [ny_acc].[dbo].[ny_collisions]
    add Vtype_id1 nvarchar(100);
[+] alter table [ny_acc].[dbo].[ny_collisions]
    add Vtype_id2 nvarchar(100);
[+] alter table [ny_acc].[dbo].[ny_collisions]
    add Vtype_id3 nvarchar(100);
[+] alter table [ny_acc].[dbo].[ny_collisions]
    add Vtype_id4 nvarchar(100);
[+] alter table [ny_acc].[dbo].[ny_collisions]
    add Vtype_id5 nvarchar(100);

[+] alter table [ny_acc].[dbo].[ny_collisions]
    add Location_id nvarchar(100);

```

The next step was to insert the corresponding values in the id columns:



```

update [ny_acc].[dbo].[ny_collisions]
set [ny_acc].[dbo].[ny_collisions].Date_id = [ny_acc].[dbo].[Accident_Date].Date_id
from [ny_acc].[dbo].[Accident_Date]
where [ny_acc].[dbo].[ny_collisions].DATE = [ny_acc].[dbo].[Accident_Date].Acc_Date

update [ny_acc].[dbo].[ny_collisions]
set [ny_acc].[dbo].[ny_collisions].Time_id = [ny_acc].[dbo].[Accident_Time].Time_id
from [ny_acc].[dbo].[Accident_Time]
where [ny_acc].[dbo].[ny_collisions].TIME = [ny_acc].[dbo].[Accident_Time].Acc_Time

update [ny_acc].[dbo].[ny_collisions]
set [ny_acc].[dbo].[ny_collisions].Factor_id1 = [ny_acc].[dbo].[Accident_factor].factor_id
from [ny_acc].[dbo].[Accident_factor]
where [ny_acc].[dbo].[ny_collisions].CONTRIBUTING_FACTOR_VEHICLE_1 = [ny_acc].[dbo].[Accident_factor].factor

update [ny_acc].[dbo].[ny_collisions]
set [ny_acc].[dbo].[ny_collisions].Factor_id2 = [ny_acc].[dbo].[Accident_factor].factor_id
from [ny_acc].[dbo].[Accident_factor]
where [ny_acc].[dbo].[ny_collisions].CONTRIBUTING_FACTOR_VEHICLE_2 = [ny_acc].[dbo].[Accident_factor].factor

update [ny_acc].[dbo].[ny_collisions]
set [ny_acc].[dbo].[ny_collisions].Factor_id3 = [ny_acc].[dbo].[Accident_factor].factor_id
from [ny_acc].[dbo].[Accident_factor]
where [ny_acc].[dbo].[ny_collisions].CONTRIBUTING_FACTOR_VEHICLE_3 = [ny_acc].[dbo].[Accident_factor].factor

update [ny_acc].[dbo].[ny_collisions]
set [ny_acc].[dbo].[ny_collisions].Factor_id4 = [ny_acc].[dbo].[Accident_factor].factor_id
from [ny_acc].[dbo].[Accident_factor]
where [ny_acc].[dbo].[ny_collisions].CONTRIBUTING_FACTOR_VEHICLE_4 = [ny_acc].[dbo].[Accident_factor].factor

update [ny_acc].[dbo].[ny_collisions]
set [ny_acc].[dbo].[ny_collisions].Factor_id5 = [ny_acc].[dbo].[Accident_factor].factor_id
from [ny_acc].[dbo].[Accident_factor]
where [ny_acc].[dbo].[ny_collisions].CONTRIBUTING_FACTOR_VEHICLE_5 = [ny_acc].[dbo].[Accident_factor].factor

update [ny_acc].[dbo].[ny_collisions]
set [ny_acc].[dbo].[ny_collisions].Vtype_id1 = [ny_acc].[dbo].[Vehicle_type].vtype_id
from [ny_acc].[dbo].[Vehicle_type]
where [ny_acc].[dbo].[ny_collisions].VEHICLE_TYPE_CODE_1 = Vehicle_type.vtype

update [ny_acc].[dbo].[ny_collisions]
set [ny_acc].[dbo].[ny_collisions].Vtype_id2 = Vehicle_type.vtype_id
from [ny_acc].[dbo].[Vehicle_type]
where [ny_acc].[dbo].[ny_collisions].VEHICLE_TYPE_CODE_2 = Vehicle_type.vtype

update [ny_acc].[dbo].[ny_collisions]
set [ny_acc].[dbo].[ny_collisions].Vtype_id3 = [ny_acc].[dbo].[Vehicle_type].vtype_id
from [ny_acc].[dbo].[Vehicle_type]
where [ny_acc].[dbo].[ny_collisions].VEHICLE_TYPE_CODE_3 = Vehicle_type.vtype

update [ny_acc].[dbo].[ny_collisions]
set [ny_acc].[dbo].[ny_collisions].Vtype_id4 = [ny_acc].[dbo].[Vehicle_type].vtype_id
from [ny_acc].[dbo].[Vehicle_type]
where [ny_acc].[dbo].[ny_collisions].VEHICLE_TYPE_CODE_4 = Vehicle_type.vtype

update [ny_acc].[dbo].[ny_collisions]
set [ny_acc].[dbo].[ny_collisions].Vtype_id5 = Vehicle_type.vtype_id
from [ny_acc].[dbo].[Vehicle_type]
where [ny_acc].[dbo].[ny_collisions].VEHICLE_TYPE_CODE_5 = Vehicle_type.vtype

update [ny_acc].[dbo].[ny_collisions]
set [ny_acc].[dbo].[ny_collisions].Location_id = Location.location_id
from [ny_acc].[dbo].[Location]
where [ny_acc].[dbo].[ny_collisions].BOROUGH = Location.borough and
[ny_acc].[dbo].[ny_collisions].ZIP_CODE = Location.zip_code and
[ny_acc].[dbo].[ny_collisions].LONGITUDE = Location.longitude and
[ny_acc].[dbo].[ny_collisions].LATITUDE = Location.latitude;

```

Final we conclude the procedure of the creation of the fact table and the insertion of the data we run the following queries.

```

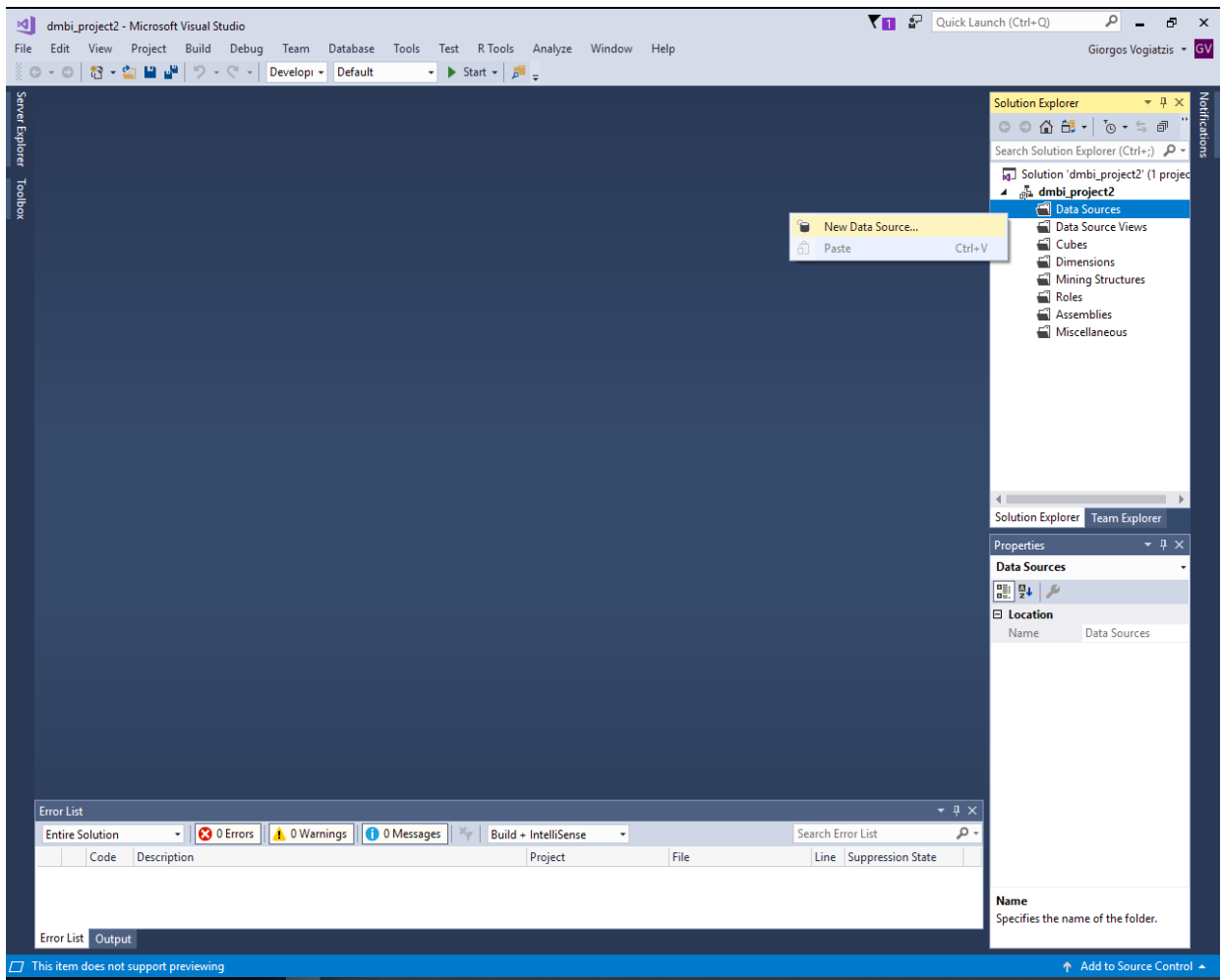
create table fact_table
(
    unique_key nvarchar(50),
    time_id nvarchar(100),
    date_id nvarchar(100),
    location_id nvarchar(100),
    persons_injured int,
    persons_killed int,
    pedestrians_injured int,
    pedestrians_killed int,
    cyclists_injured int,
    cyclists_killed int,
    motorists_injured int,
    motorists_killed int,
    factor_vehicle1_id nvarchar(100),
    factor_vehicle2_id nvarchar(100),
    factor_vehicle3_id nvarchar(100),
    factor_vehicle4_id nvarchar(100),
    factor_vehicle5_id nvarchar(100),
    vtype1_id nvarchar(100),
    vtype2_id nvarchar(100),
    vtype3_id nvarchar(100),
    vtype4_id nvarchar(100),
    vtype5_id nvarchar(100),
    foreign key(time_id) references [ny_acc].[dbo].[Accident_Time](Time_id),
    foreign key(date_id) references [ny_acc].[dbo].[Accident_Date](Date_id),
    foreign key(location_id) references [ny_acc].[dbo].[Location](location_id),
    foreign key(factor_vehicle1_id) references [ny_acc].[dbo].[Accident_factor](factor_id),
    foreign key(factor_vehicle2_id) references [ny_acc].[dbo].[Accident_factor](factor_id),
    foreign key(factor_vehicle3_id) references [ny_acc].[dbo].[Accident_factor](factor_id),
    foreign key(factor_vehicle4_id) references [ny_acc].[dbo].[Accident_factor](factor_id),
    foreign key(factor_vehicle5_id) references [ny_acc].[dbo].[Accident_factor](factor_id),
    foreign key(vtype1_id) references [ny_acc].[dbo].[Vehicle_type](vtype_id),
    foreign key(vtype2_id) references [ny_acc].[dbo].[Vehicle_type](vtype_id),
    foreign key(vtype3_id) references [ny_acc].[dbo].[Vehicle_type](vtype_id),
    foreign key(vtype4_id) references [ny_acc].[dbo].[Vehicle_type](vtype_id),
    foreign key(vtype5_id) references [ny_acc].[dbo].[Vehicle_type](vtype_id)
)

insert into fact_table (unique_key,time_id,date_id,location_id,persons_injured,persons_killed,pedestrians_injured,pedestrians_killed,
cyclists_injured,cyclists_killed,motorists_injured,motorists_killed,factor_vehicle1_id,factor_vehicle2_id,factor_vehicle3_id,factor_vehicle4_id,factor_vehicle5_id,
vtype1_id,vtype2_id,vtype3_id,vtype4_id,vtype5_id )
select UNIQUE_KEY,Time_id,Date_id,Location_id,NUMBER_OF_PERSONS_INJURED,NUMBER_OF_PERSONS_KILLED,NUMBER_OF_PEDESTRIANS_INJURED,NUMBER_OF_PEDESTRIANS_KILLED,
NUMBER_OF_CYCLIST_INJURED,NUMBER_OF_CYCLIST_KILLED,NUMBER_OF_MOTORIST_INJURED,NUMBER_OF_MOTORIST_KILLED,
Factor_id1,Factor_id2,Factor_id3,Factor_id4,Factor_id5,Vtype_id1,Vtype_id2,Vtype_id3,Vtype_id4,Vtype_id5
from [ny_acc].[dbo].[ny_collisions]

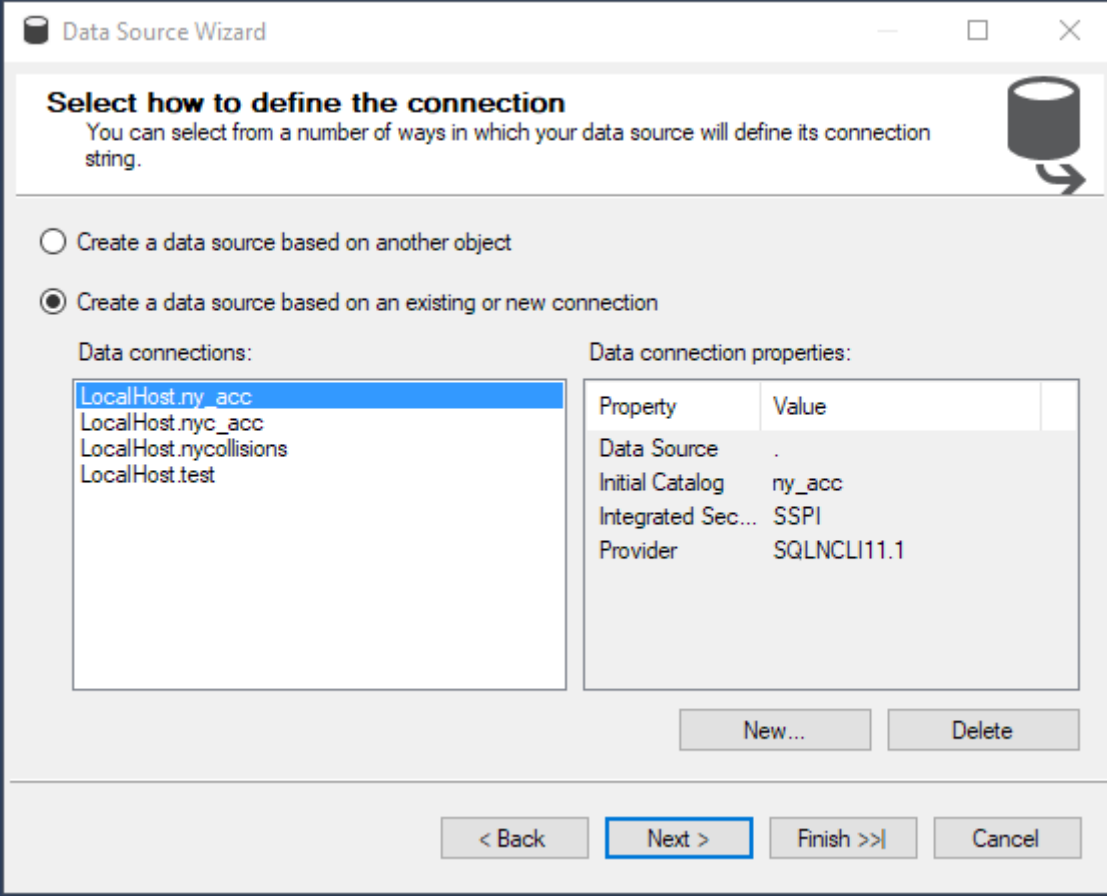
```

## 2.3 Visual Studio – Cube Creation

Initially, we created an analysis services multidimensional project. Screenshots provided show the steps followed to create the cube. Firstly, we choose the data source:



Then, we prompted to choose our data connection



The Data Source Wizard dialog box is shown. It has a title bar with a cylinder icon and the text "Data Source Wizard". The main area has a header "Select how to define the connection" with a sub-header "You can select from a number of ways in which your data source will define its connection string." and a cylinder icon with an arrow. There are two radio buttons: "Create a data source based on another object" (unselected) and "Create a data source based on an existing or new connection" (selected). Below the radio buttons are two sections: "Data connections:" and "Data connection properties:". The "Data connections:" section contains a list box with four items: "LocalHost.ny\_acc" (selected), "LocalHost.nyc\_acc", "LocalHost.nycollisions", and "LocalHost.test". The "Data connection properties:" section contains a table with two columns: "Property" and "Value". The table has four rows: "Data Source" with value ".", "Initial Catalog" with value "ny\_acc", "Integrated Sec..." with value "SSPI", and "Provider" with value "SQLNCLI11.1". Below the table are two buttons: "New..." and "Delete". At the bottom of the dialog are four buttons: "< Back", "Next >" (highlighted with a blue border), "Finish >>|", and "Cancel".

**Select how to define the connection**  
You can select from a number of ways in which your data source will define its connection string.

☐ Create a data source based on another object

☒ Create a data source based on an existing or new connection

Data connections:

- LocalHost.ny\_acc
- LocalHost.nyc\_acc
- LocalHost.nycollisions
- LocalHost.test

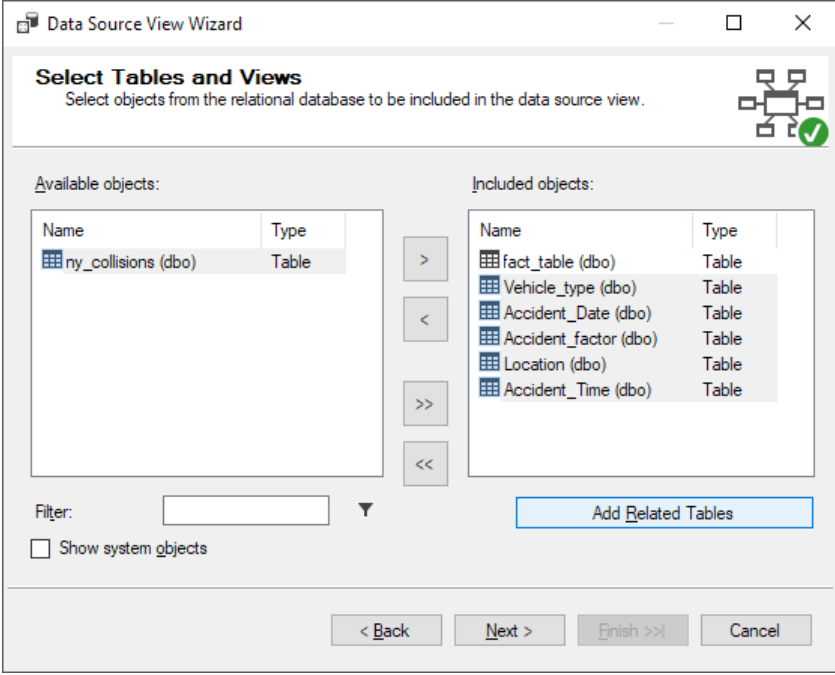
Data connection properties:

Property	Value
Data Source	.
Initial Catalog	ny_acc
Integrated Sec...	SSPI
Provider	SQLNCLI11.1

New... Delete

< Back Next > Finish >>| Cancel

After, he had to create a new data source view and transfer our fact table and its related dimensions in the data source views.



The Data Source View Wizard dialog box is shown. It has a title bar with a cylinder icon and the text "Data Source View Wizard". The main area has a header "Select Tables and Views" with a sub-header "Select objects from the relational database to be included in the data source view." and a cylinder icon with a green checkmark. There are two sections: "Available objects:" and "Included objects:". The "Available objects:" section contains a table with two columns: "Name" and "Type". It has one row: "ny\_collisions (dbo)" with type "Table". The "Included objects:" section contains a table with two columns: "Name" and "Type". It has six rows: "fact\_table (dbo)" with type "Table", "Vehicle\_type (dbo)" with type "Table", "Accident\_Date (dbo)" with type "Table", "Accident\_factor (dbo)" with type "Table", "Location (dbo)" with type "Table", and "Accident\_Time (dbo)" with type "Table". Between the two tables are four buttons: ">", "<", ">>", and "<<". Below the "Available objects:" table is a "Filter:" text box and a "Show system objects" checkbox. Below the "Included objects:" table is an "Add Related Tables" button. At the bottom of the dialog are four buttons: "< Back", "Next >" (highlighted with a blue border), "Finish >>|", and "Cancel".

**Select Tables and Views**  
Select objects from the relational database to be included in the data source view.

Available objects:

Name	Type
ny_collisions (dbo)	Table

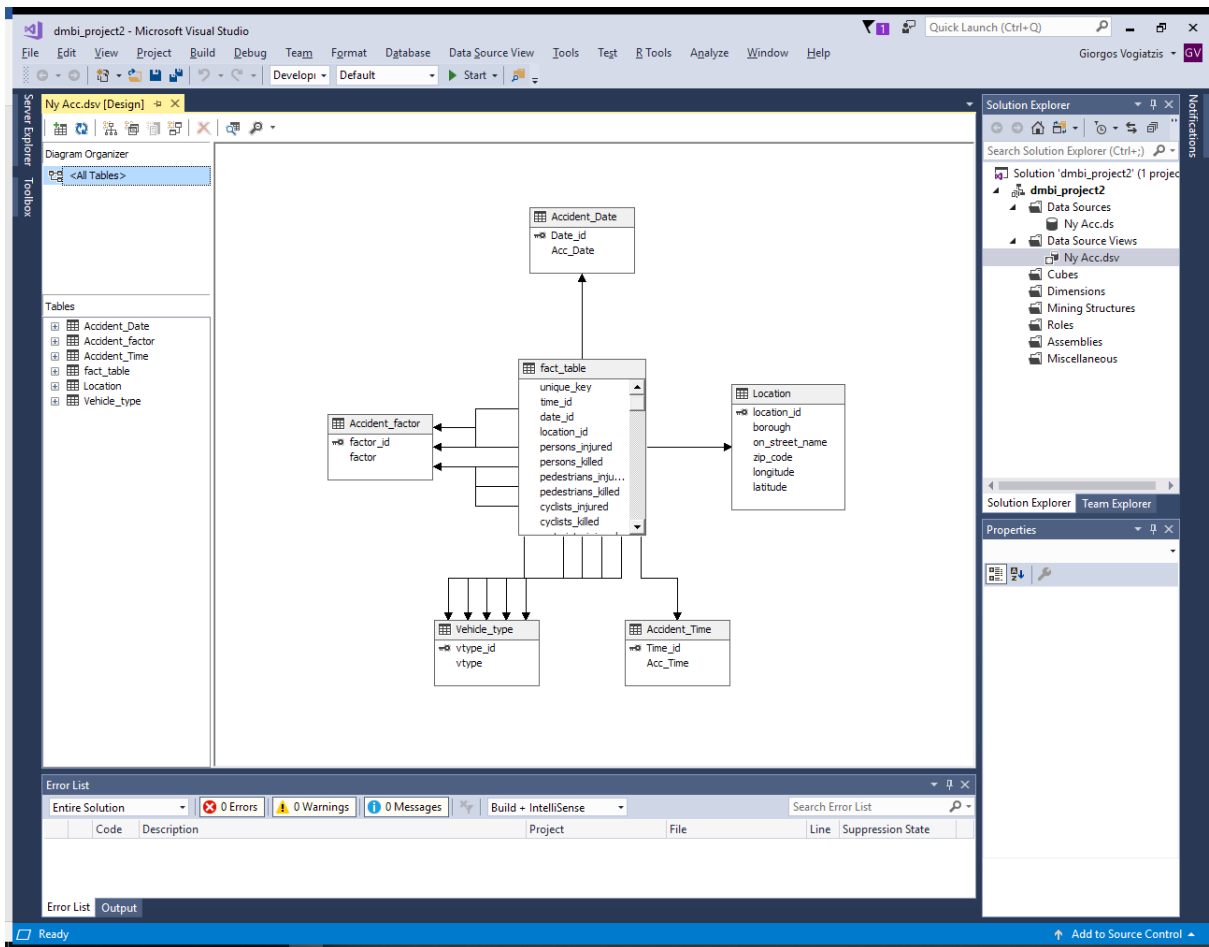
Included objects:

Name	Type
fact_table (dbo)	Table
Vehicle_type (dbo)	Table
Accident_Date (dbo)	Table
Accident_factor (dbo)	Table
Location (dbo)	Table
Accident_Time (dbo)	Table

Filter: Show system objects

Add Related Tables

< Back Next > Finish >>| Cancel



Moreover, we continue with the creation of the cube and its dimensions with the following steps:

### Select Creation Method

Cubes can be created by using existing tables, creating an empty cube, or generating tables in the data source.

How would you like to create the cube?

☒ Use existing tables  
☐ Create an empty cube  
☐ Generate tables in the data source

Template:

(None)

Description:

Create a cube based on one or more tables in a data source.

Cube Wizard

### Select Measure Group Tables

Select a data source view or diagram and then select the tables that will be used for measure groups.

Data source view:

Ny Acc

Measure group tables:

Suggest

☒

fact_table
------------

☐

Vehicle_type
--------------

☐

Accident_Date
---------------

☐

Accident_factor
-----------------

☐

Location
----------

☐

Accident_Time
---------------

< Back

Next >

Finish >>

Cancel

Cube Wizard

### Select Measures

Select measures that you want to include in the cube.

☐ Measure

☒

Fact Table
------------

☒

Persons Injured
-----------------

☒

Persons Killed
----------------

☒

Pedestrians Injured
---------------------

☒

Pedestrians Killed
--------------------

☒

Cyclists Injured
------------------

☒

Cyclists Killed
-----------------

☒

Motorists Injured
-------------------

☒

Motorists Killed
------------------

☐

Fact Table Count
------------------

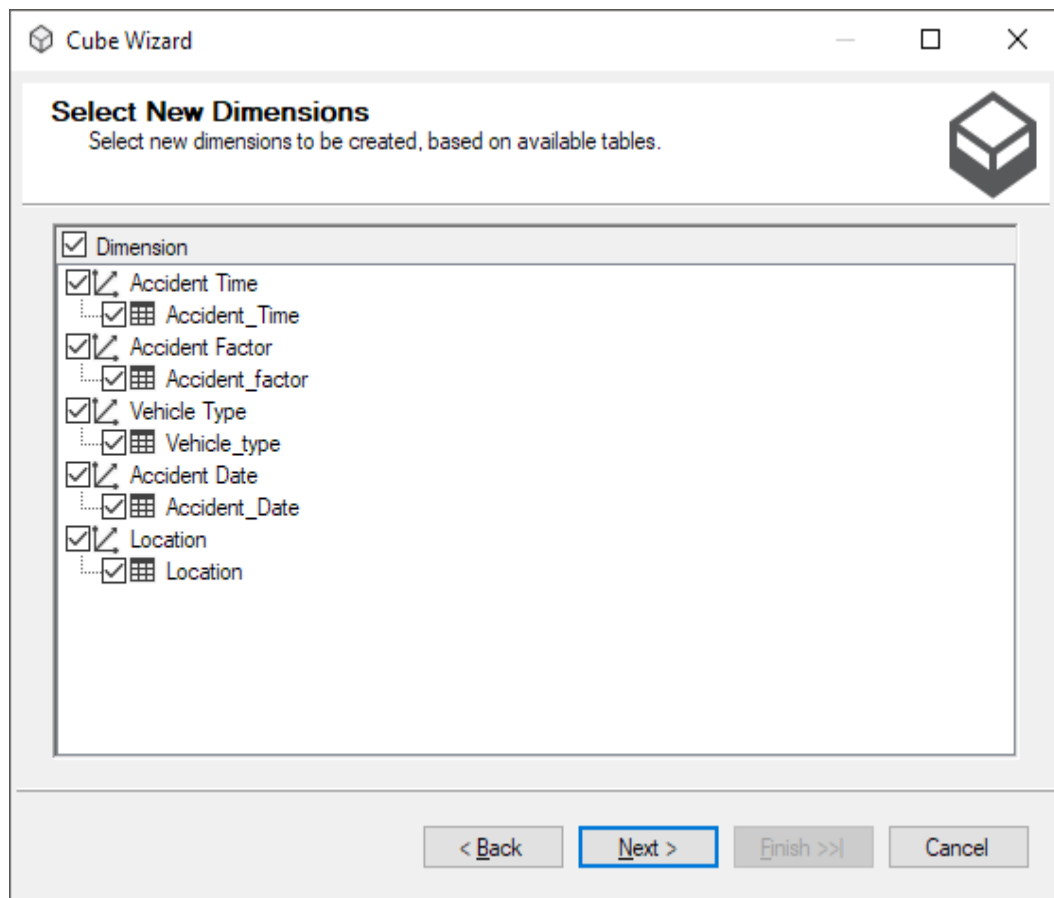
< Back

Next >

Finish >>

Cancel

After the choice of the fact table, we continue with the dimensions:

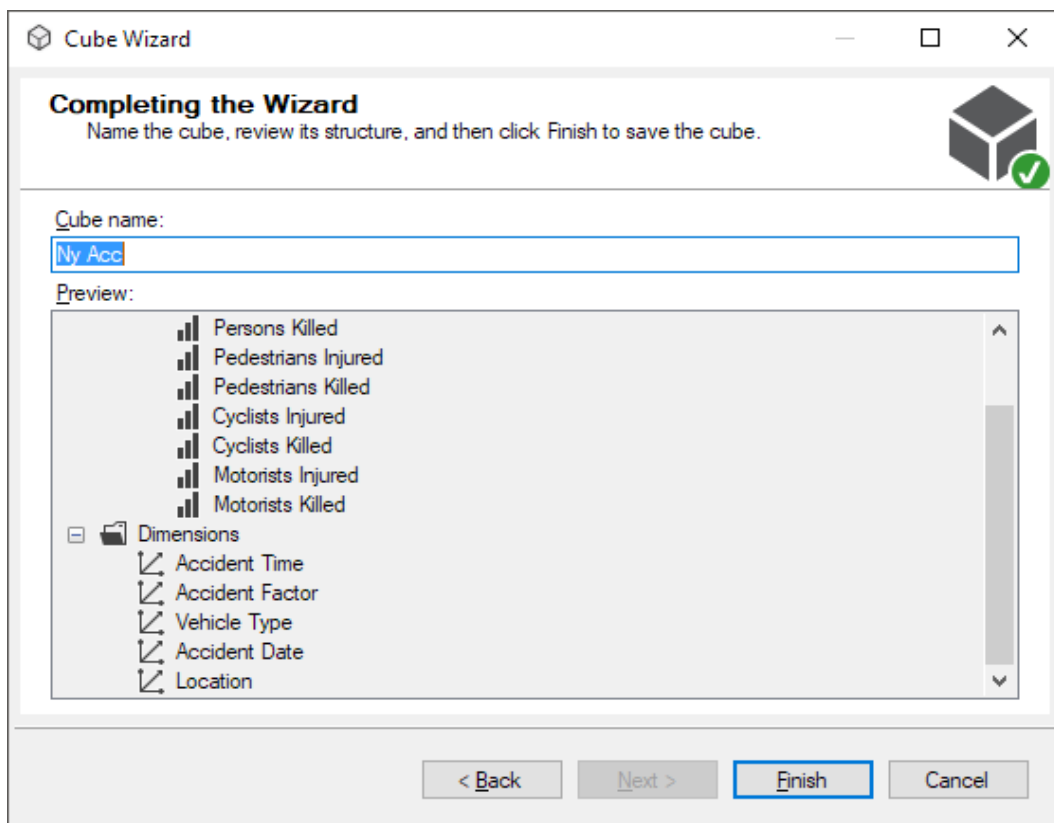


The screenshot shows the 'Select New Dimensions' window of the Cube Wizard. The title bar reads 'Cube Wizard'. The main heading is 'Select New Dimensions' with a subtitle 'Select new dimensions to be created, based on available tables.' and a cube icon. A list of dimensions is shown, each with a checked checkbox and a small grid icon: 'Accident Time', 'Accident Factor', 'Vehicle Type', 'Accident Date', and 'Location'. Below each dimension name is a sub-item with a checked checkbox and a grid icon: 'Accident\_Time', 'Accident\_factor', 'Vehicle\_type', 'Accident\_Date', and 'Location'. At the bottom are four buttons: '< Back', 'Next >', 'Finish >>', and 'Cancel'.

**Select New Dimensions**  
Select new dimensions to be created, based on available tables.

- ☒ Dimension
  - ☒ Accident Time
    - ☒ Accident\_Time
  - ☒ Accident Factor
    - ☒ Accident\_factor
  - ☒ Vehicle Type
    - ☒ Vehicle\_type
  - ☒ Accident Date
    - ☒ Accident\_Date
  - ☒ Location
    - ☒ Location

< Back   Next >   Finish >>   Cancel



The screenshot shows the 'Completing the Wizard' window of the Cube Wizard. The title bar reads 'Cube Wizard'. The main heading is 'Completing the Wizard' with a subtitle 'Name the cube, review its structure, and then click Finish to save the cube.' and a cube icon with a green checkmark. A text field for 'Cube name:' contains 'N.Y. Acc'. Below it is a 'Preview:' section showing a tree structure of the cube's contents: 'Persons Killed', 'Pedestrians Injured', 'Pedestrians Killed', 'Cyclists Injured', 'Cyclists Killed', 'Motorists Injured', and 'Motorists Killed' under a 'Dimensions' folder. At the bottom are four buttons: '< Back', 'Next >', 'Finish', and 'Cancel'.

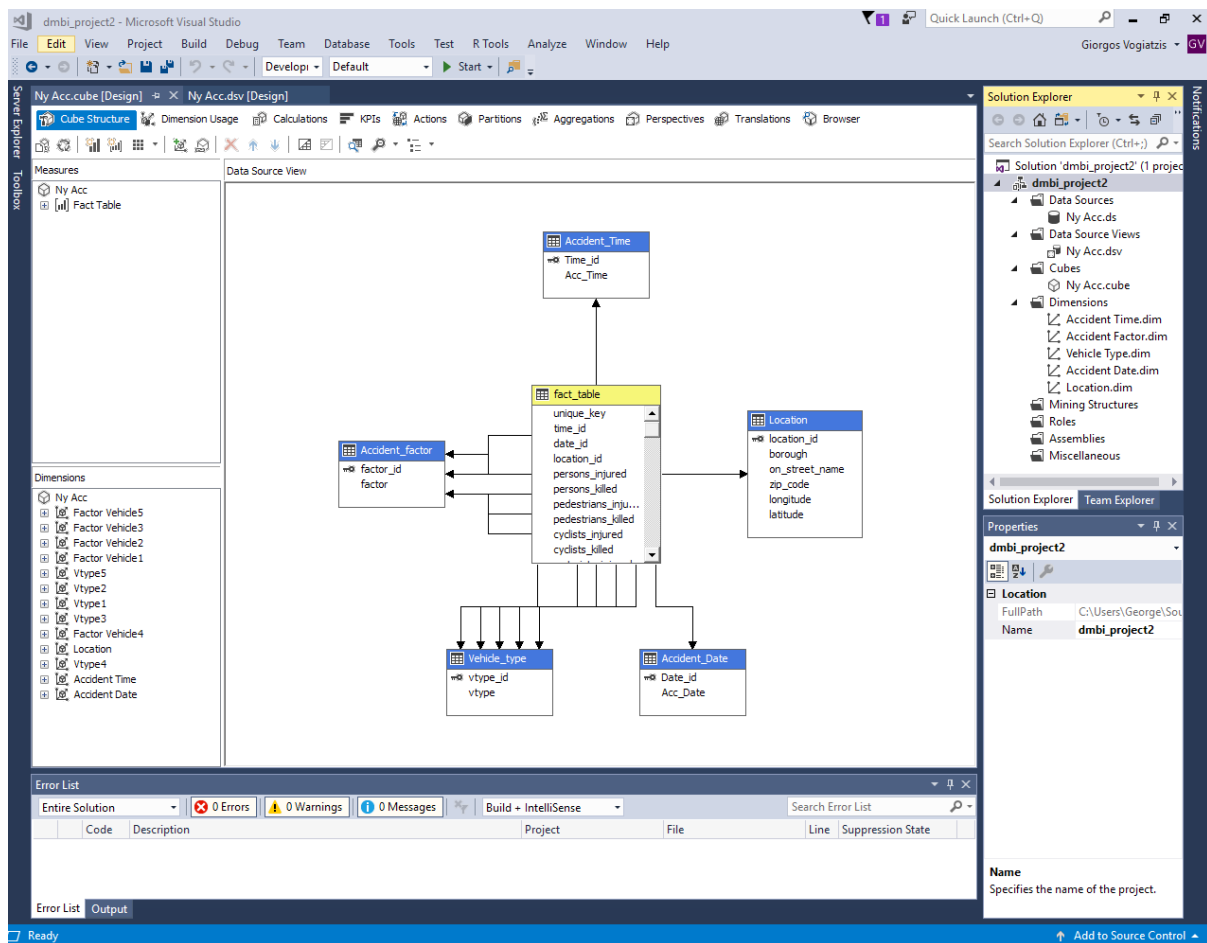
**Completing the Wizard**  
Name the cube, review its structure, and then click Finish to save the cube.

Cube name: N.Y. Acc

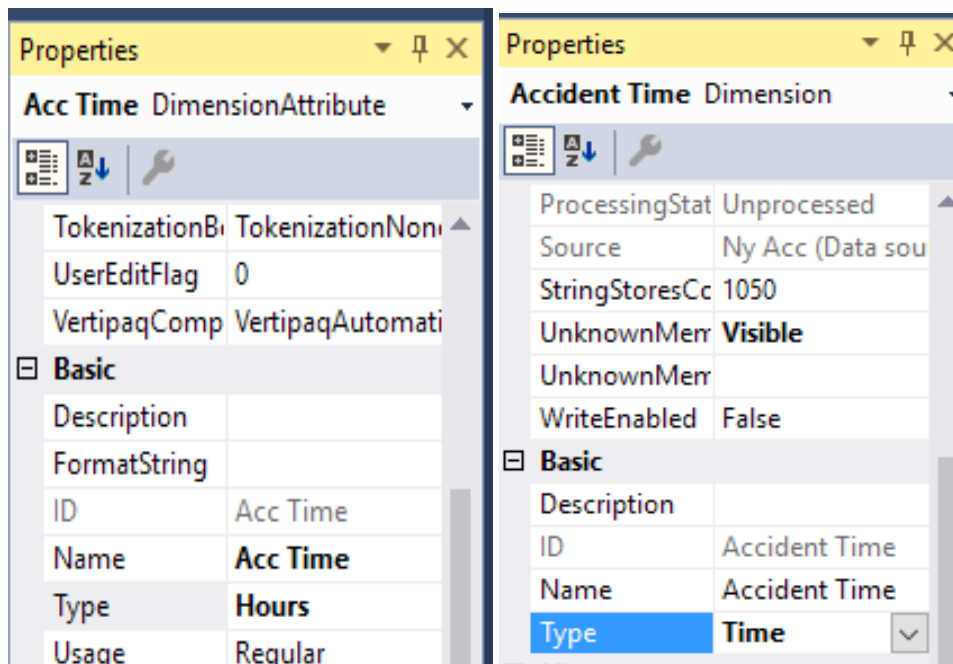
Preview:

- Persons Killed
- Pedestrians Injured
- Pedestrians Killed
- Cyclists Injured
- Cyclists Killed
- Motorists Injured
- Motorists Killed
- Dimensions
  - Accident Time
  - Accident Factor
  - Vehicle Type
  - Accident Date
  - Location

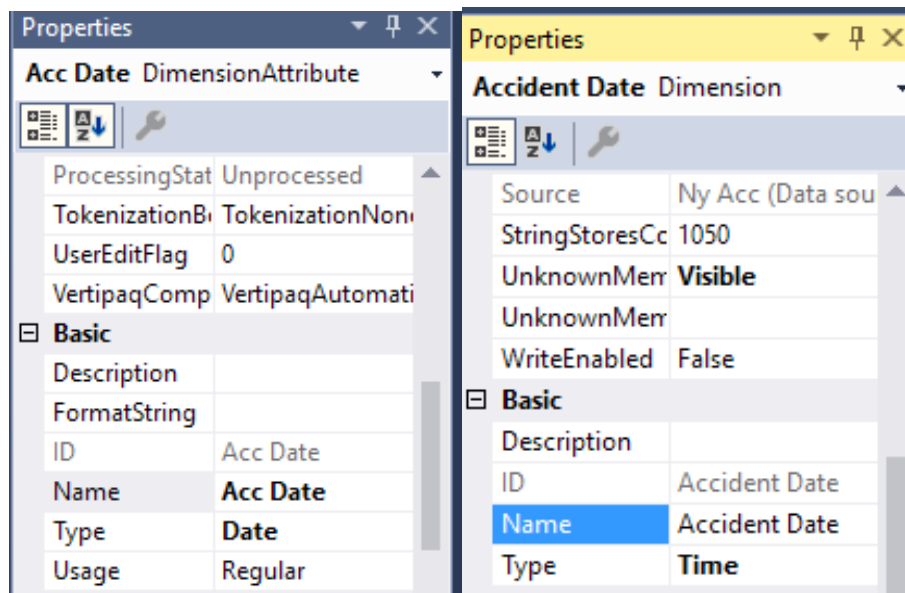
< Back   Next >   Finish   Cancel



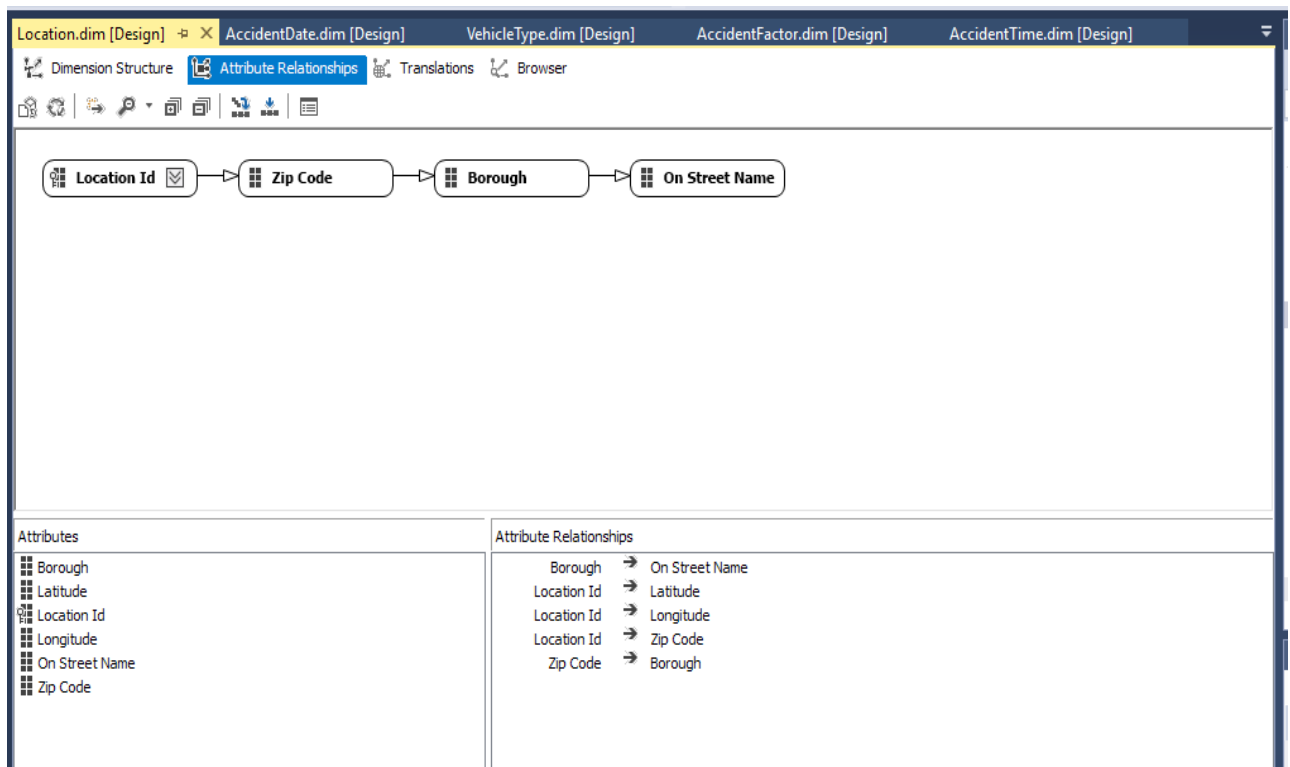
Then we had to change the data type of 'acc\_time' and 'acc\_date' attributes dimensions to Hours and Date respectively, and change the dimensions time and date to Time type.

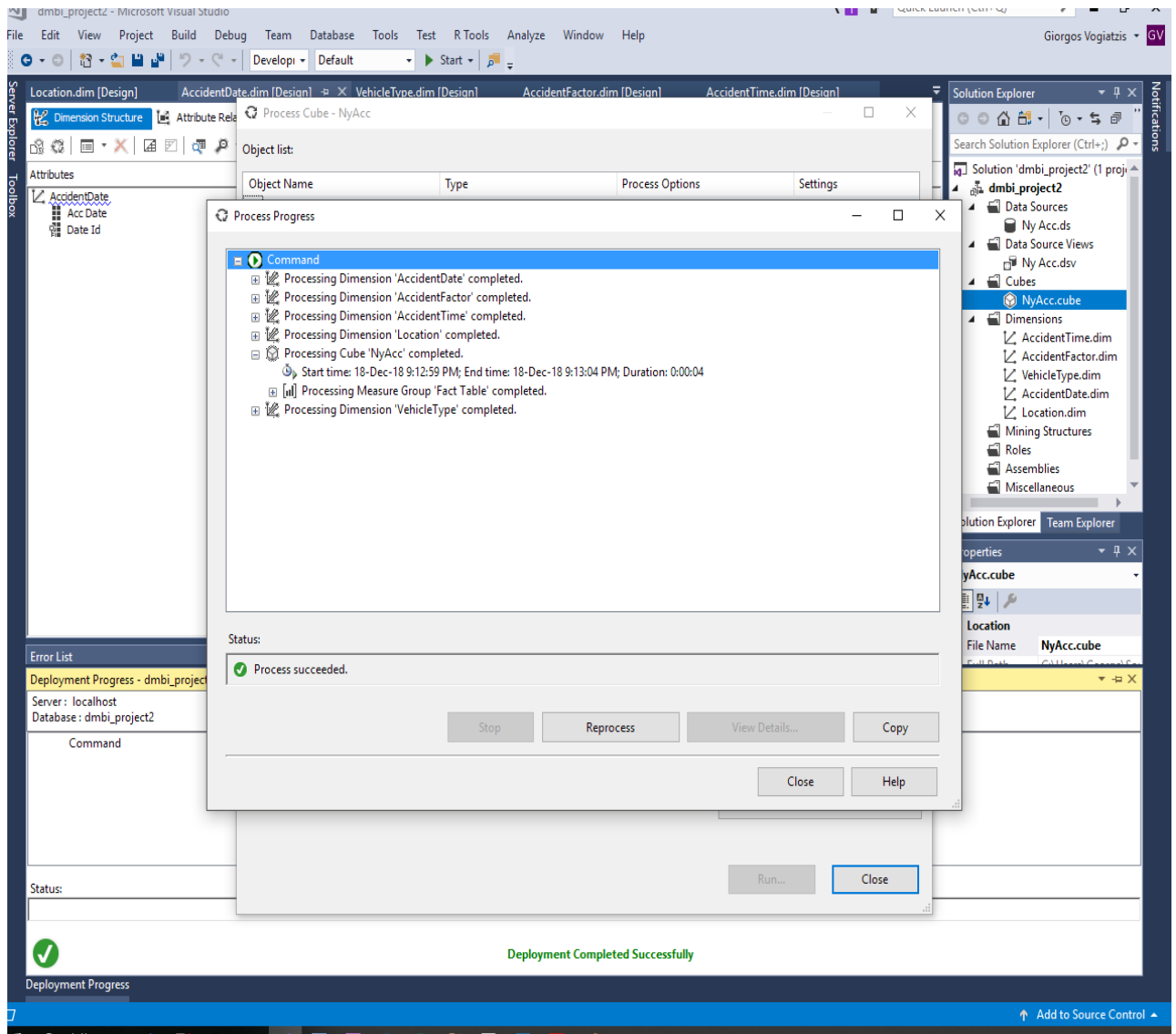






After, we defined the hierarchies in the dimension location table. So, we processed and deployed the cube successfully.





## 2.4 Olap Reports

In the following report we can see the number of pedestrians and motorists that got injured in the borough of Brooklyn on the date 14<sup>th</sup> of may 2015

The screenshot displays the SAP BW OLAP report interface. The left pane shows the 'NyAcc' cube structure with a tree view of dimensions and measures. The right pane shows the report data.


**Dimension Usage Table:**











Dimension	Hierarchy	Operator	Filter Expression	Parameter
Location	Borough	Equal	{BROOKLYN}	<input type="checkbox"/>
Accident Date	Acc Date	Equal	{05/14/2015}	<input type="checkbox"/>
<Select dimension>				<input type="checkbox"/>


**Report Data Table:**










Pedestrians Injured	Motorists Injured
10	28

Here we can see the result of our rollup for all the injuries occurred in Brooklyn and Manhattan at 13:47 o'clock for all the days in our calendar.

Location.dim [Design]    AccidentDate.dim [Design]    VehicleType.dim [Design]    AccidentTime.dim [Design]    NyAcc.cube [Design]    

 Cube Structure    Dimension Usage    Calculations    KPIs    Actions    Partitions    Aggregations    Perspectives    Translations    Browser

Language: Default 

 Edit as Text    Import...   MDX                   


NyAcc


Metadata


Search Model


Measure Group:

<All>


 Pedestrians Injured

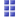
 Pedestrians Killed


 Persons Injured


 Persons Killed


KPIs


 Accident Date


 Acc Date


 Date Id


 Accident Time


 Factor Vehicle1


 Factor Vehicle2


 Factor Vehicle3


 Factor Vehicle4


 Factor Vehicle5


 Location


 Borough


 Latitude


 Location Id

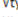
 Longitude


 On Street Name


 Members


 On Street Name


 Zip Code


 Hierarchy

 Vtype1

 Vtype2

 Vtype3

 Vtype4

 Vtype5

Calculated Members

Dimension	Hierarchy	Operator	Filter Expression	Param...
Location	Borough	Equal	{ BROOKLYN, MANHATTAN }	<input type="checkbox"/>
Accident Time	Acc Time	Equal	{ 13:47:00.0000000 }	<input type="checkbox"/>
<Select dimension>				<input type="checkbox"/>

Acc Date	Motorists Injured	Pedestrians Injured	Persons Injured
01/02/2...	1	0	1
01/04/2...	0	0	0
01/05/2...	1	0	1
01/07/2...	0	1	1
01/08/2...	0	0	0
01/12/2...	0	0	1
01/25/2...	0	0	0
01/29/2...	0	0	0
01/31/2...	0	0	1
02/08/2...	0	0	0
02/09/2...	1	0	1
02/12/2...	5	1	6
02/14/2...	0	0	0
02/23/2...	0	0	0
02/26/2...	0	0	0
02/27/2...	0	0	0
02/28/2...	0	0	0
03/03/2...	0	0	0
03/03/2...	0	0	0
03/05/2...	0	0	0
03/10/2...	0	0	0
03/10/2...	0	0	0
03/12/2...	0	0	0
03/12/2...	0	0	0
03/16/2...	1	0	1
03/16/2...	0	0	0
03/16/2...	0	0	0
03/20/2...	0	0	0
03/24/2...	0	0	0
03/26/2...	0	0	0












Calculated Members

And finally a slice and dice example which gives us the number of pedestrians and motorists killed in Brooklyn between the 4<sup>th</sup> of January 2014 and the 2<sup>nd</sup> of January 2018.

Edit as Text

Import...

MDX



NyAcc

Metadata

Search Model

Measure Group:

<All>

Cyclists Killed

Motorists Injured

Motorists Killed

Pedestrians Injured

Pedestrians Killed

Persons Injured

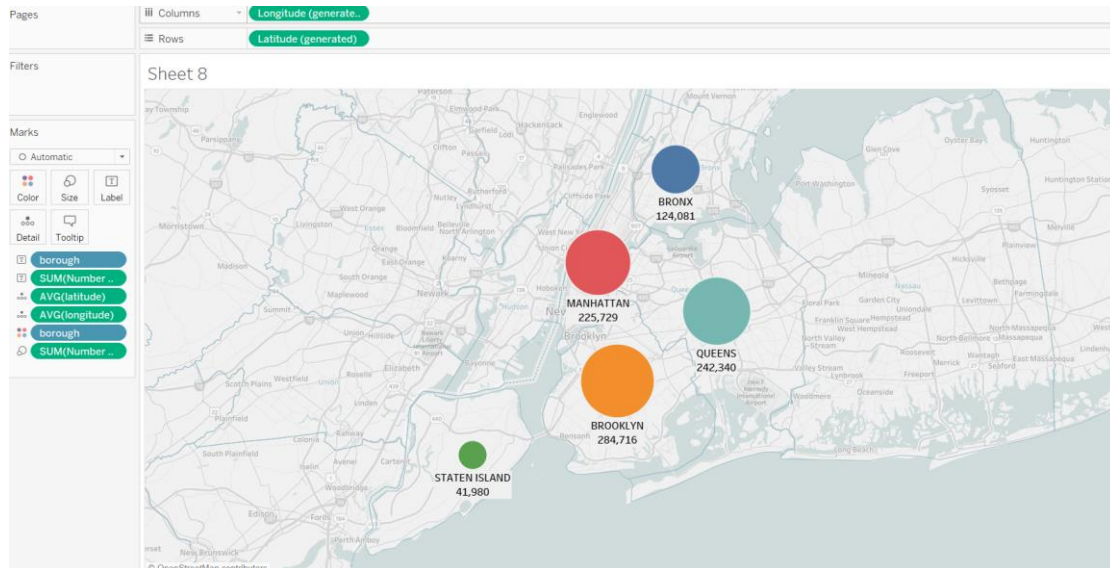
Persons Killed

Dimension	Hierarchy	Operator	Filter Expression	Param...
Location	Borough	Equal	{ BROOKLYN }	<input type="checkbox"/>
Accident Date	Acc Date	Range (Inclusive)	01/04/2014 : 01/02/2018	<input type="checkbox"/>
<Select dimension>				<input type="checkbox"/>

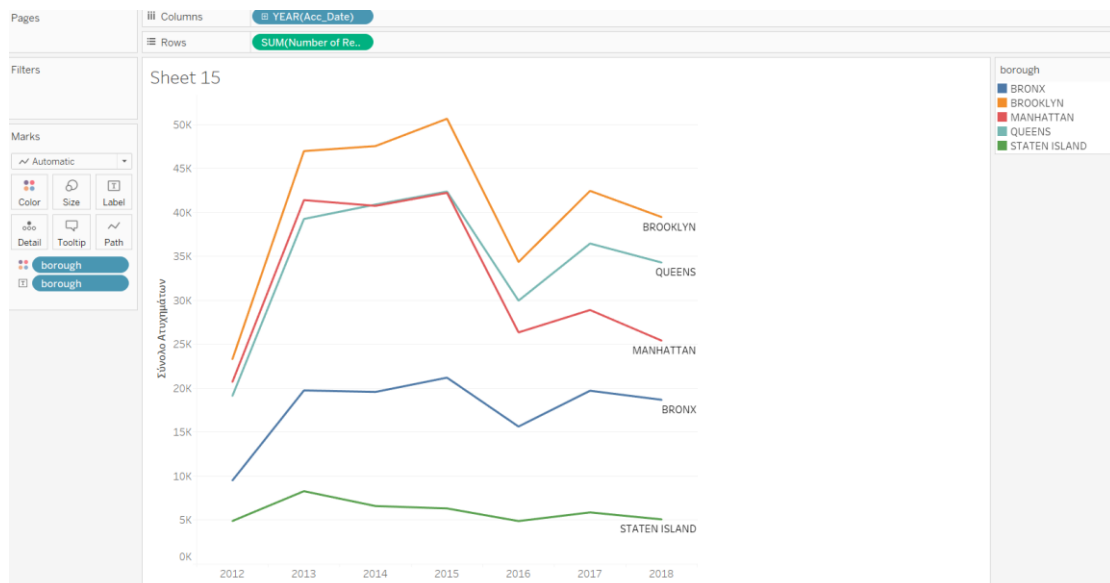
Motorists Injured	Pedestrians Injured
171	92

### 3. TABLEAU

We use the tableau program in order to make the visualization analysis of our data. In the first diagram, we observe that the Brooklyn have the more accidents and the Staten Islands the lowest collisions in New Work City.



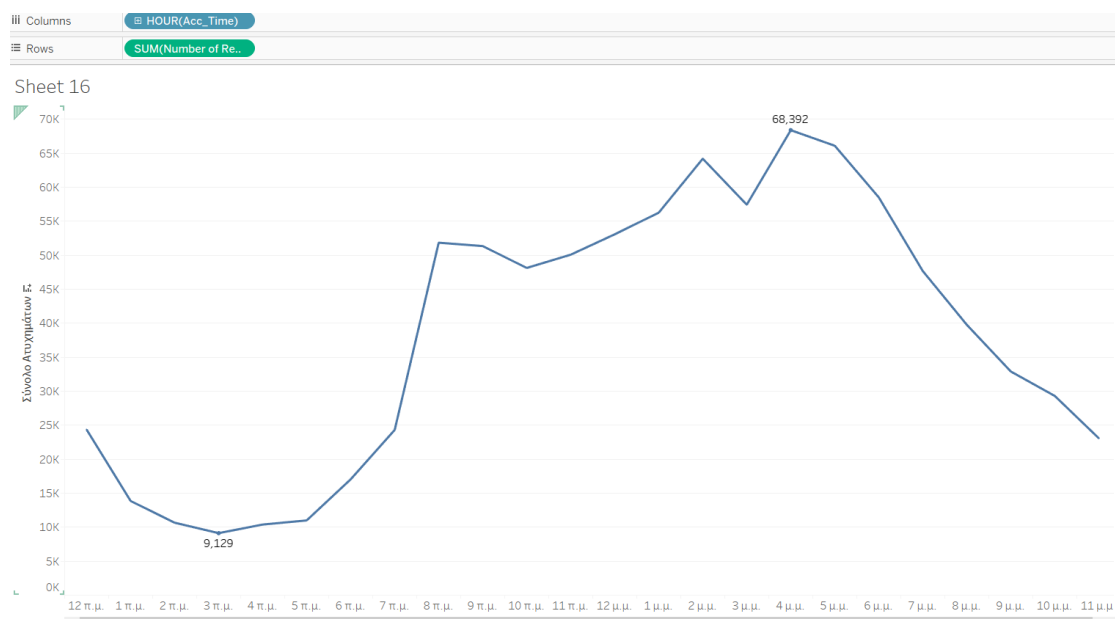
In the second diagram, we observe that the Brooklyn have the more accident in all the examed years and the Staten Islands has the less collisions in all the examed years. Moreover, we observed that from 2012 until 2014 Manhattan have more accidents than Queens has but that change in 2016 when Queens has more accidents. That pattern continues until today.



From the below diagram we see that the most collisions happened in 2015 and the less in the 2012.



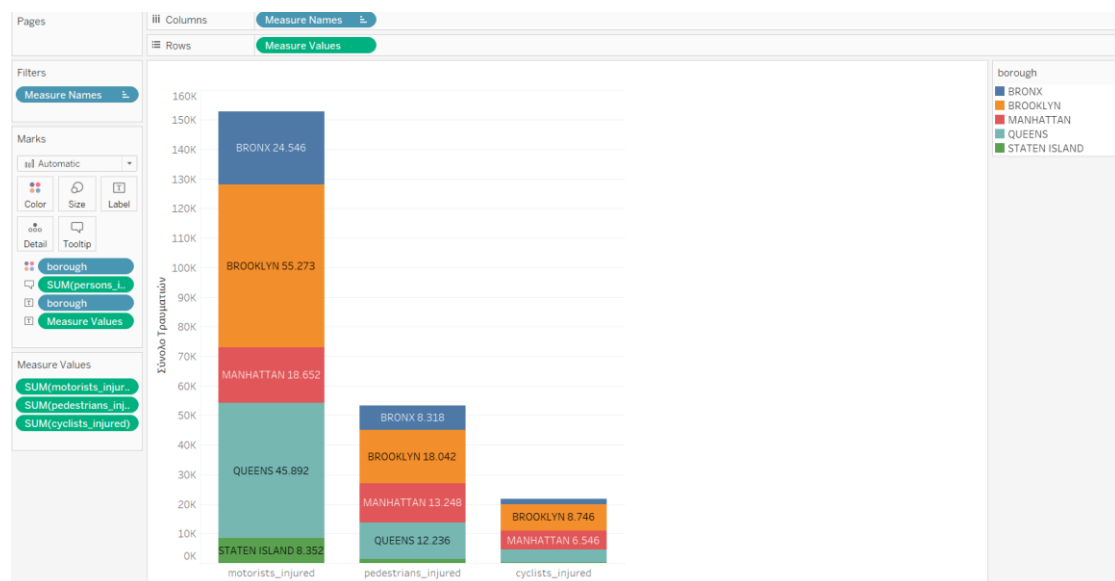
Except from the year, we want to observe how the number of collisions alter through the time. We observe that the most collisions happened in 4 a.m. and the less in 3 p.m...



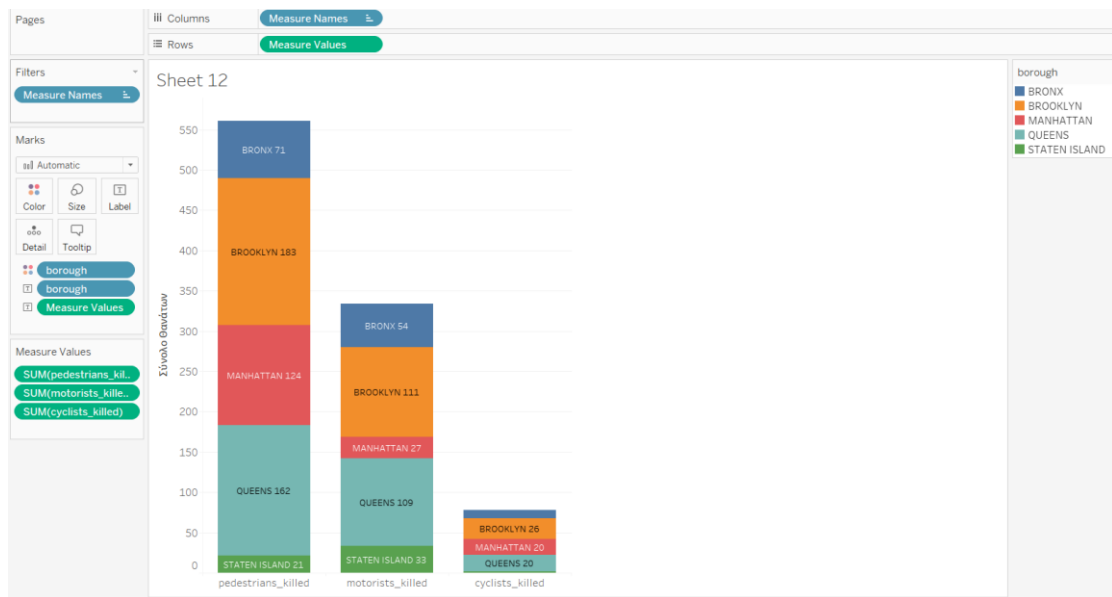
Furthermore, with the above diagram we see that the most collisions happened in October and the less in April.



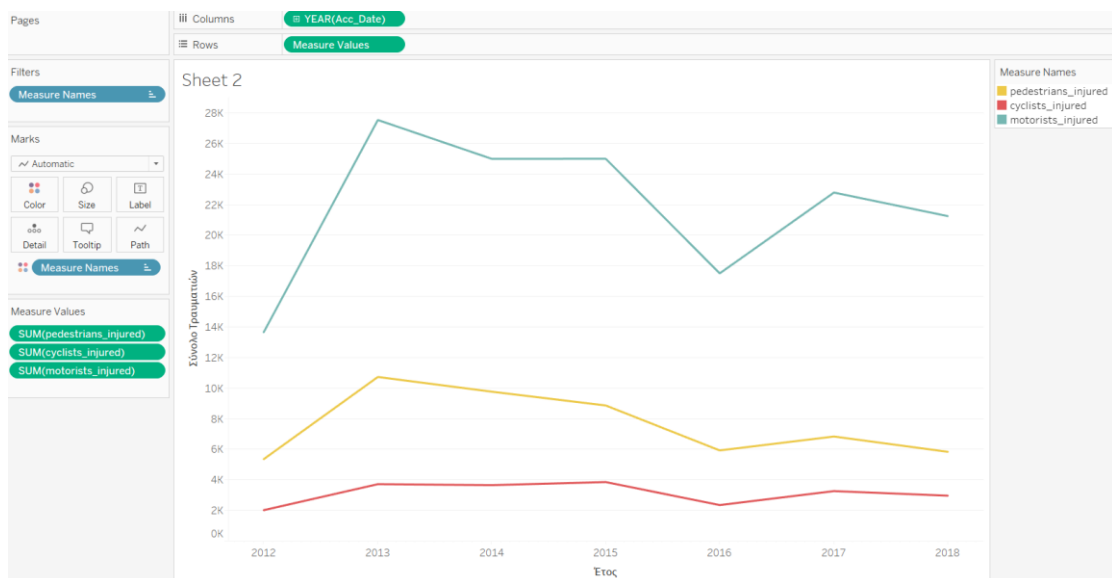
From the above diagram we see that from the total of the injured people after a collision are the motorists are the most and the cyclist the less. In addition, we see how this number alter in towns. We observe that the Brooklyn have the most injures.



From the above diagram we see that the Brooklyn have the most deaths which affect most the pedestrians.

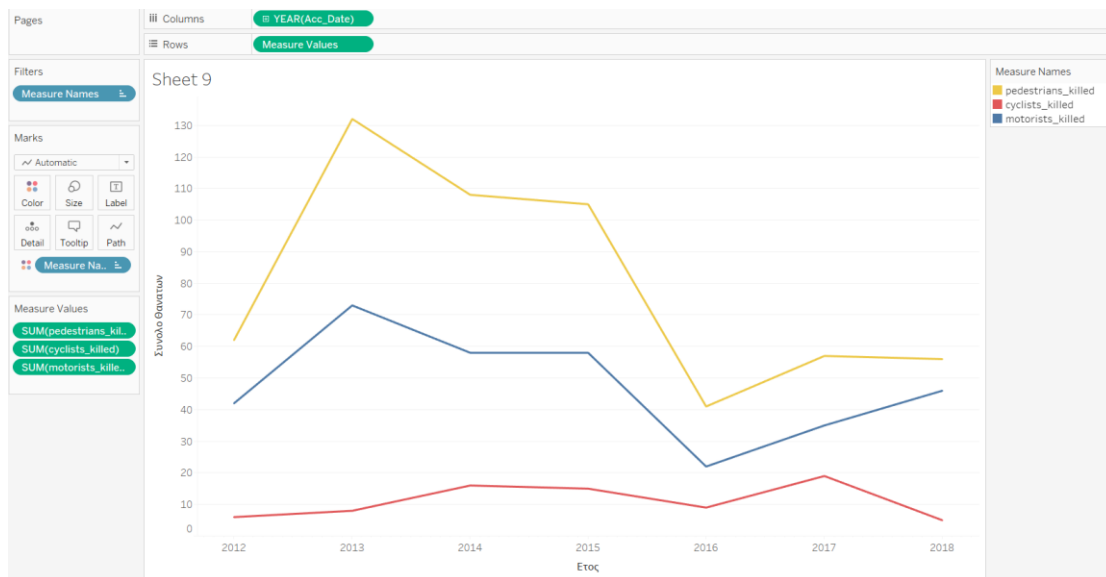


From the above plot, we see that the most injuries have the motorist in 2013. In 2016 we observe a sharp reduce in all the categories but the in the motorists was more observed.

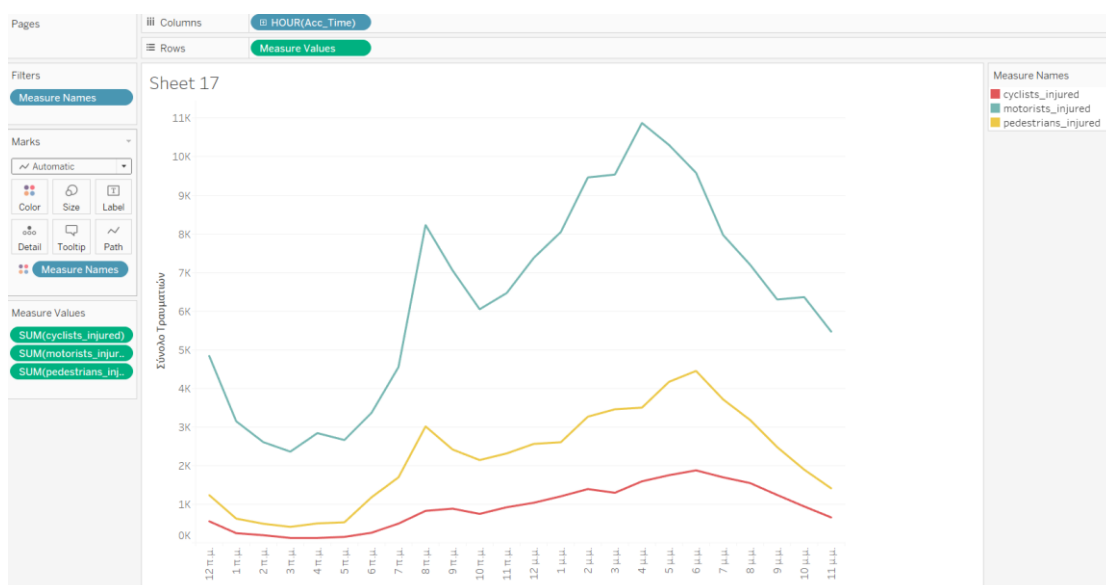


From the above diagram, we see that the most deaths affect the pedestrians and this is stable through the years. In addition, the numbers of deaths in cyclist are the lowest and the only one who reduced in 2018 in contrast to other categories, which increased.

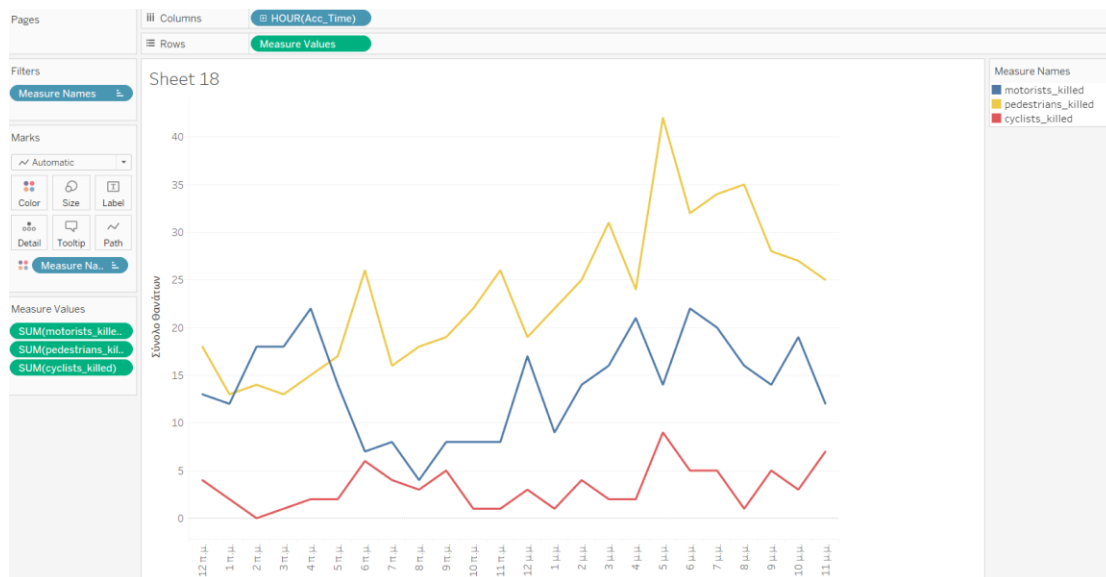




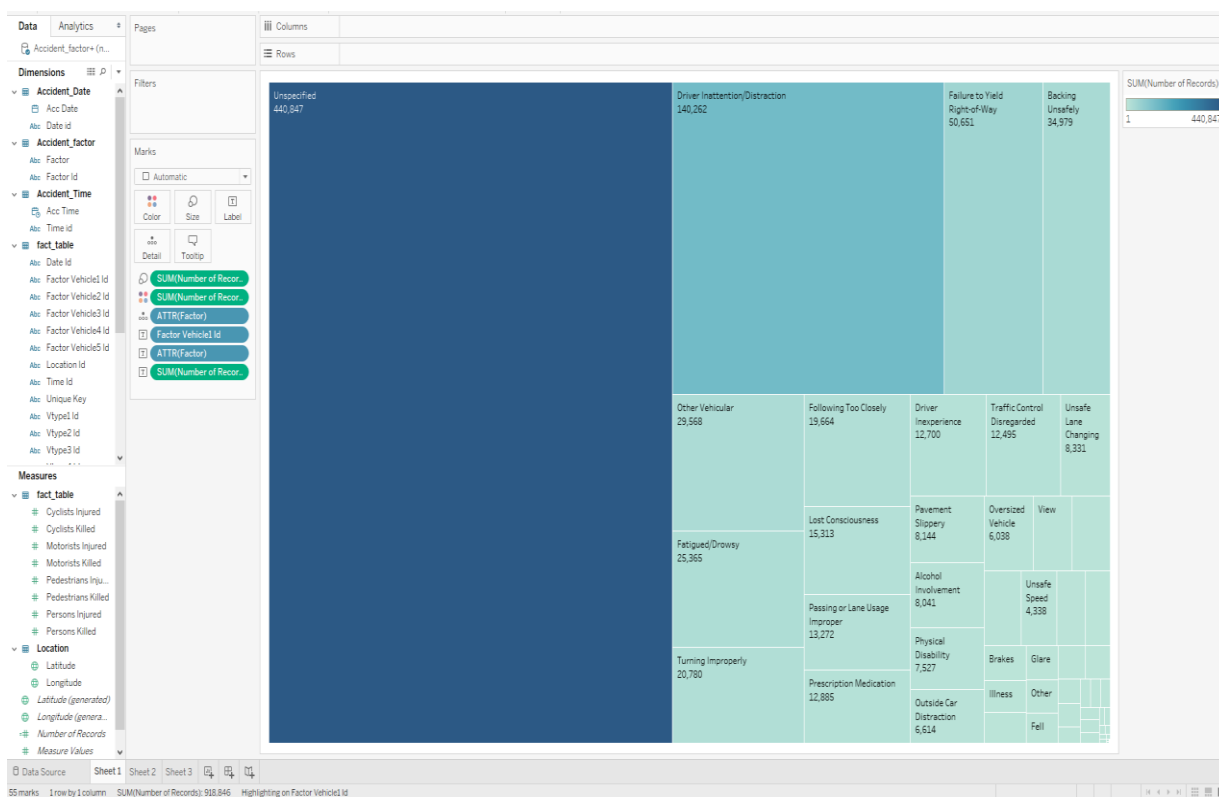
The most injuries happen among 8 a.m. and 6 p.m. seems all the categories follow the same pattern with time.



However the number of people who dead in the collision do not seem to follow the same pattern through time. The most deaths, which refers in pedestrians and cyclist, happen in 6 p.m. and the most deaths of motorists happened in 3a.m.



The reason, which affects a collision in the most cases, are unspecified. The second most often reason is the distraction of the driver and the third more often reason is the failure to yield right of way. For this diagram, we use the variable of our dataset 'factor vehicle\_1' because we think that the first vehicle is responsible for the accident.



In corrspondence that we use the 'factor vehicle\_1' we use the variable about the 'vehicle type\_1'. In order to see what type of vehicle causes a collision. Final we see that the type of vehicle which take part in the most collisions are the passenger vehicle, the second category are the taxi and third the category with unknown vehicle type.

