

INTRODUCTION

System Background:

The Medical Clinic Management System is a comprehensive database system put in place by QualityLife Healthcare with the aim of streamlining and enhancing the operations of a busy medical clinic. The system was created in response to the demand for data-driven decision-making in the healthcare sector, the increased complexity of patient care, and the requirement for effective provider management.

Features and Objectives:

1. Patient Management: The system allows clinic staff to efficiently manage patient information, including demographics, medical history, and appointment scheduling.
2. Provider Management: It enables healthcare providers to maintain their profiles, including specialization, licenses, and contact information, making it easy for patients to choose the right provider for their needs.
3. Appointment Scheduling: Patients can conveniently schedule appointments online or through the clinic's front desk. The system optimizes appointment slots, minimizes wait times, and ensures efficient allocation of resources.
4. Tracking of procedures: Names, descriptions, and related costs of procedures carried out throughout appointments are documented to help with invoicing and reporting.

Intended Report:

Provider Performance Report. This report evaluates the performance of healthcare providers based on metrics such as appointment counts, average satisfaction ratings, appointment durations, and wait times. It assists in recognizing top-performing providers and areas that may need improvement.

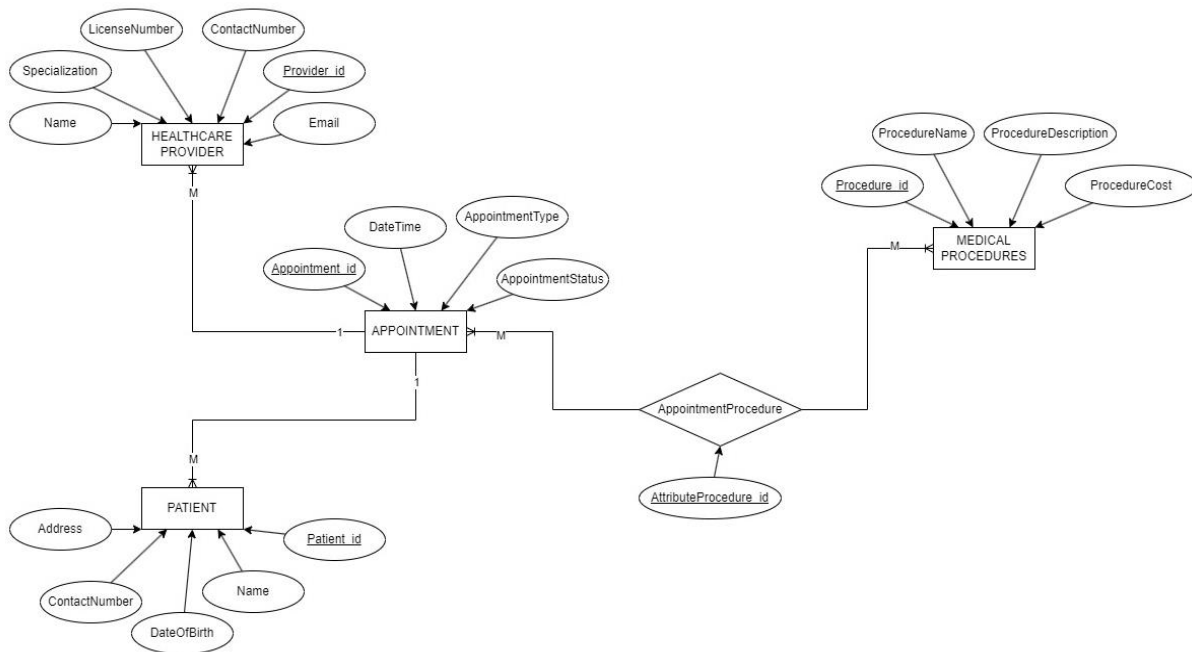


Figure 1. Entity-Relational Diagram

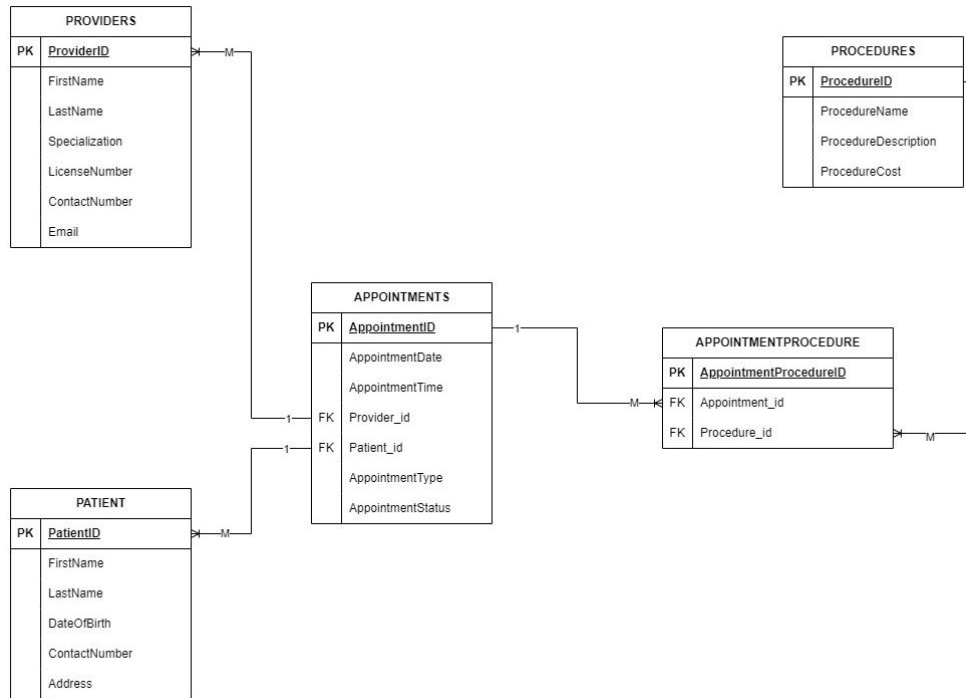


Figure 2. Normalized Relational Model

After creating ERD Figure 1 on page 1, we proceeded to convert it into a Normalized Relational Model (NRM), as depicted in Figure 2. Utilizing the NRM, we can derive the dimension tables necessary for our data warehouse. These dimension tables, in conjunction with the Fact table, will enable us to generate the desired reports mentioned on page 1. Please refer to the table below for the derived dimension tables.

Methodology of Dimensional Normal Form:

DIMENSION TABLE	ATTRIBUTES	DIMENSION TABLE	ATTRIBUTES
Time	TimeID (Primary Key) Date Day of the Week Month Quarter Year	Provider	ProviderID (Primary Key) FirstName LastName Specialization LicenseNumber ContactNumber Email
Specialization	SpecializationID (Primary Key) Specialization		
FACT TABLE		ATTRIBUTES	
Provider Performance		PerformanceID (Prime Key) TimeID (Foreign Key) ProviderID (Foreign Key) SpecializationID (Foreign Key) AppointmentCount AverageSatisfactionRating AverageAppointmentDuration AverageWaitTime	

In the Dimensional Normal Form (DNF), each attribute is stored in its own dimension table, facilitating efficient storage and retrieval of time-related data. In this design, there are three Dimension Tables: Time, Provider, and Specialization. The Time and Provider Dimensions were derived from the NRM displayed in Figure 2 on page 2, while the Specialization Dimension, and the attributes found in the Fact Table named "Provider Performance" in Figure 3 below were sourced from another database.

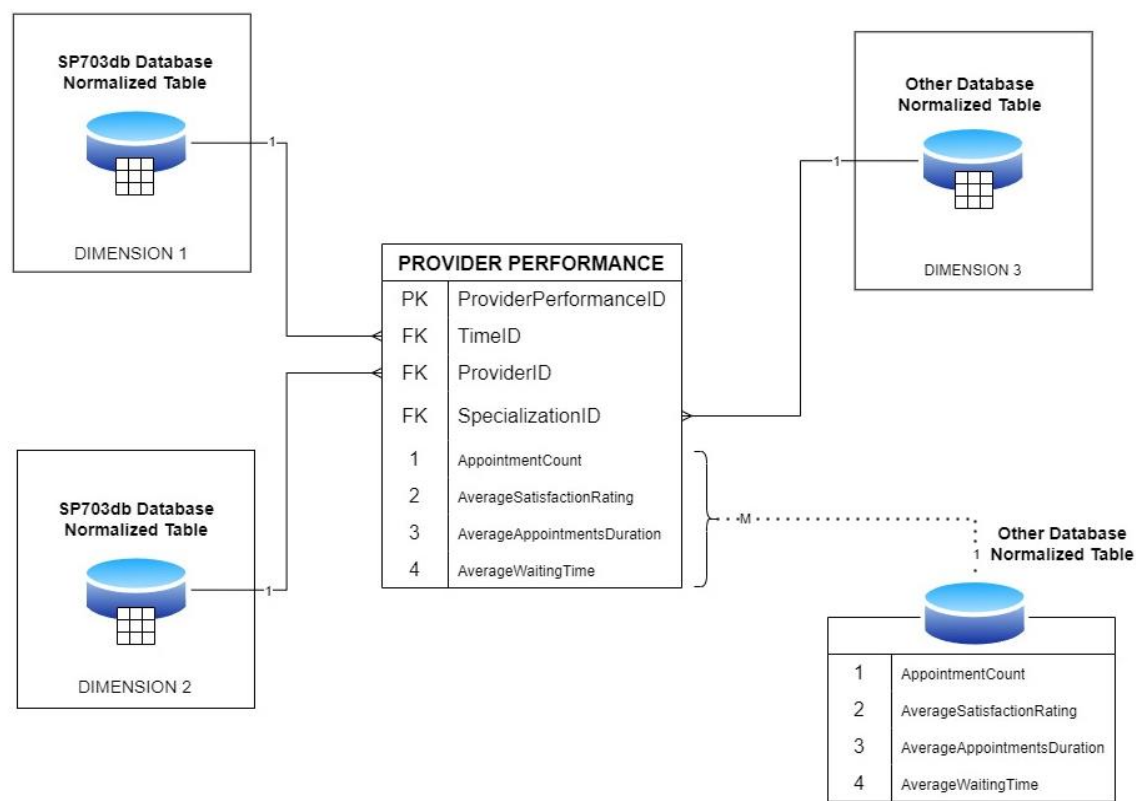


Figure 3. SP703dw Data Warehouse source databases

In Figure 3, we can see that Dimension 1 (Time) and Dimension 2 (Provider) were derived from the same database, named SP703db. On the other hand, Dimension 3 (Specialization) and the attributes within the Fact Table (Provider Performance) were sourced from a different database, as mentioned previously.

The Fact Table plays a critical role in a data warehousing system, as it provides a structured and efficient means of storing and analyzing data. By structuring dimension tables in this manner, the design enhances query performance for the Provider Performance fact table. See Figure 4 on page 4 for details.

For clarification, we will only create the database SP703db and use the Normalized Relational Model shown in Figure 2 on page 2 as our guide for the structure.

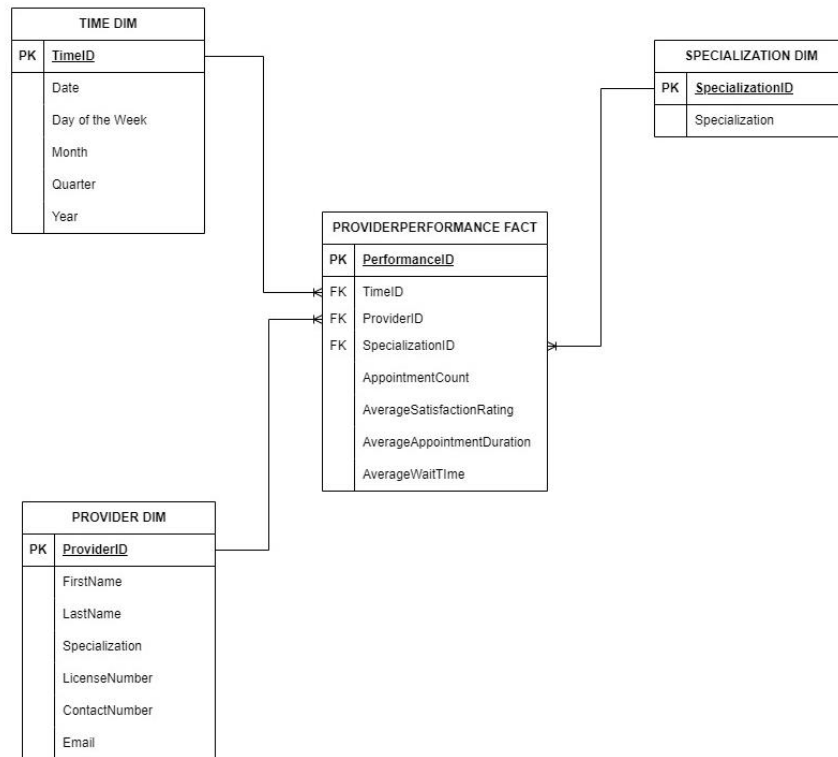


Figure 4. Dimensional Model (Fact & Dimensions table)

CREATING DATABASE

SQL script for Database in PostgreSQL (based on the NRM Figure 2 on page 2)

Query	Query History
1	-- Create the Patients table
2	CREATE TABLE Patients (
3	PatientID INT PRIMARY KEY,
4	FirstName VARCHAR(255),
5	LastName VARCHAR(255),
6	DateOfBirth DATE,
7	ContactNumber VARCHAR(20),
8	Address VARCHAR(255)
9);
11	-- Create the Providers table
12	CREATE TABLE Providers (
13	ProviderID INT PRIMARY KEY,
14	FirstName VARCHAR(255),
15	LastName VARCHAR(255),
16	Specialization VARCHAR(255),
17	LicenseNumber VARCHAR(255),
18	ContactNumber VARCHAR(20),
19	Email VARCHAR(255)
20);
22	-- Create the Appointments table
23	CREATE TABLE Appointments (
24	AppointmentID INT PRIMARY KEY,
25	AppointmentDate DATE,
26	AppointmentTime TIME,
27	PatientID INT,
28	ProviderID INT,
29	AppointmentType VARCHAR(50),
30	AppointmentStatus VARCHAR(50),
31	FOREIGN KEY (PatientID) REFERENCES Patients(PatientID),
32	FOREIGN KEY (ProviderID) REFERENCES Providers(ProviderID)
33);

```

35 -- Create the Procedures table
36 CREATE TABLE Procedures (
37     ProcedureID INT PRIMARY KEY,
38     ProcedureName VARCHAR(255),
39     ProcedureDescription TEXT,
40     ProcedureCost DECIMAL(10, 2)
41 );
42
43 -- Create the AppointmentProcedure junction table
44 CREATE TABLE AppointmentProcedure (
45     AppointmentProcedureID INT PRIMARY KEY,
46     AppointmentID INT,
47     ProcedureID INT,
48     FOREIGN KEY (AppointmentID) REFERENCES Appointments(AppointmentID),
49     FOREIGN KEY (ProcedureID) REFERENCES Procedures(ProcedureID)
50 );

```

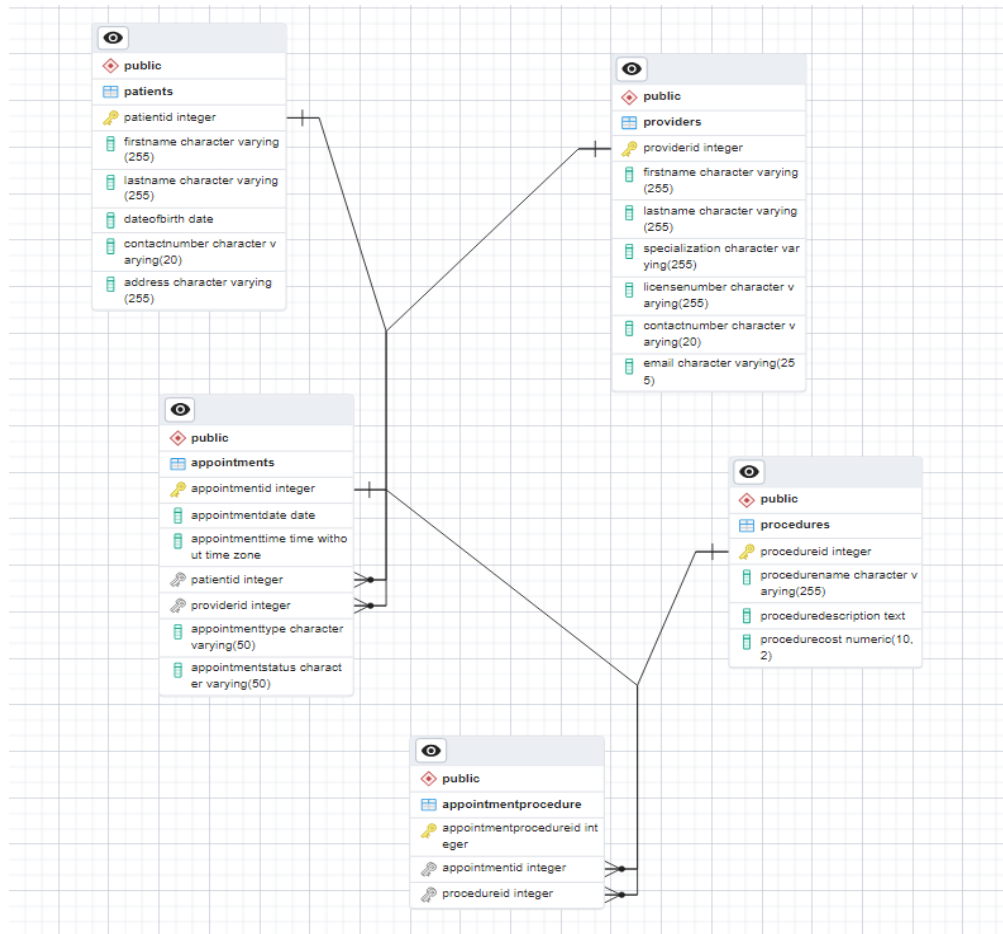


Figure 5. Database SP703db (PostgreSQL)

After creating the **database in PostgreSQL** named “**SP703db**” and creating the tables using the NRM, we will use **Python + SQL to load the dummy data** instead of using SQL only.

LOADING DUMMY DATA TO SP703db

SQL + Python (Jupyter Notebook)

Setting credentials to access database to establish a database connection.

```
1 PGHOST = 'localhost'
2 PGDATABASE = 'SP703db'
3 PGUSER = '*****' # hiding sensitive info
4 PGPASSWORD = '*****' # hiding sensitive info
```

```
1 def create_connect():
2     dbconn = psycopg2.connect(
3         dbname=PGDATABASE,
4         user=PGUSER,
5         password=PGPASSWORD,
6         host=PGHOST,
7         port='5432'
8     )
9     return dbconn
```

Loading Patients data into patients table in SP703db

```
1 #creating connection
2 dbconn = create_connect()
3
4 # Create a cursor
5 cur = dbconn.cursor()
6
7 # Insert data from csv file to Patient table in SP703 db
8 for item, row in patientdf.iterrows():
9     cur.execute(
10         '''INSERT INTO patients (PatientID, FirstName, LastName, DateOfBirth, ContactNumber, Address)
11           VALUES (%s, %s, %s, %s, %s, %s) ON CONFLICT (PatientID) DO NOTHING''',
12         (row['PatientID'], row['FirstName'], row['LastName'],
13          row['DateOfBirth'], row['ContactNumber'], row['Address']))
14
15
16 # commit the changes
17 dbconn.commit()
18
19 # Execute SQL queries to check if we successfully inserted the data
20 cur.execute('SELECT * FROM patients')
21 patientdb_tbl = pd.DataFrame(cur.fetchall(), columns=patientdf.columns)
22
23 # Close the cursor and connection
24 cur.close()
25 dbconn.close()
```

SQL query in PostgreSQL

```
1 SELECT patientid, firstname, lastname, dateofbirth, contactnumber, address
2 FROM public.patients;
```

	patientid [PK] integer	firstname character varying (255)	lastname character varying (255)	dateofbirth date	contactnumber character varying (20)	address character varying (255)
1	1	John	Doe	1985-05-15	123-456-7890	123 Main St
2	2	Jane	Smith	1990-08-20	987-654-3210	456 Elm St
3	3	David	Johnson	1998-03-10	555-555-5555	789 Oak Ave
4	4	Emily	Williams	1980-11-28	111-222-3333	321 Pine Rd
5	5	Michael	Brown	1975-07-12	444-777-8888	555 Cedar Ln
6	6	Susan	Anderson	1982-09-05	222-333-4444	789 Maple St
7	7	Robert	White	1993-04-30	666-999-1111	234 Birch Ave
8	8	Linda	Miller	1977-12-15	888-333-2222	567 Oak St
9	9	William	Moore	1989-06-22	444-666-7777	890 Pine Rd
10	10	Mary	Johnson	1987-03-08	777-888-9999	456 Elm St

Loading Providers data into Providers table in SP703db

```
1 #creating connection
2 dbconn = create_connect()
3
4 # Create a cursor
5 cur = dbconn.cursor()
6
7 # Insert data from csv file to Providers table in SP703 db
8 for item, row in providerdf.iterrows():
9     cur.execute(
10         '''INSERT INTO providers (ProviderID, FirstName, LastName, Specialization,
11             LicenseNumber, ContactNumber, Email)
12             VALUES (%s, %s, %s, %s, %s, %s, %s) ON CONFLICT (ProviderID) DO NOTHING''',
13         (row['ProviderID'], row['FirstName'], row['LastName'],
14         row['Specialization'], row['LicenseNumber'], row['ContactNumber'],
15         row['Email']))
16
17
18 # commit the changes
19 dbconn.commit()
20
21 # Execute SQL queries to check if we successfully inserted the data
22 cur.execute('SELECT * FROM providers')
23 providerdb_tbl = pd.DataFrame(cur.fetchall(), columns=providerdf.columns)
24
25 # Close the cursor and connection
26 cur.close()
27 dbconn.close()
```

SQL query in PostgreSQL

1	SELECT providerid, firstname, lastname, specialization, licensenumber, contactnumber, email						
2	FROM public.providers;						

Data Output Messages Notifications

	providerid [PK] integer	firstname character varying (255)	lastname character varying (255)	specialization character varying (255)	licensenumber character varying (255)	contactnumber character varying (20)	email character varying (255)
1	1	Dr. Sarah	Anderson	Cardiology	MD123456	555-123-7890	sarah.anderson@example.com
2	2	Dr. Robert	Clark	Orthopedics	MD789012	444-987-6543	robert.clark@example.com
3	3	Dr. Lisa	Smith	Internal Medicine	MD345678	777-333-2222	lisa.smith@example.com
4	4	Dr. John	Wilson	Neurology	MD234567	888-555-9999	john.wilson@example.com
5	5	Dr. Emily	Jones	Dermatology	MD567890	999-111-2222	emily.jones@example.com

Loading Appointments data into Appointments table in SP703db

```
1 #creating connection
2 dbconn = create_connect()
3
4 # Create a cursor
5 cur = dbconn.cursor()
6
7 # Insert data from csv file to Appointments table in SP703 db
8 for item, row in appointmentdf.iterrows():
9     cur.execute(
10         '''INSERT INTO appointments (AppointmentID, AppointmentDate, AppointmentTime, PatientID,
11                                     ProviderID, AppointmentType, AppointmentStatus)
12         VALUES (%s, %s, %s, %s, %s, %s, %s) ON CONFLICT (AppointmentID) DO NOTHING''',
13         (row['AppointmentID'], row['AppointmentDate'], row['AppointmentTime'],
14          row['PatientID'], row['ProviderID'], row['AppointmentType'],
15          row['AppointmentStatus']))
16
17
18 # commit the changes
19 dbconn.commit()
20
21 # Execute SQL queries to check if we successfully inserted the data
22 cur.execute('SELECT * FROM appointments')
23 appointmentdb_tbl = pd.DataFrame(cur.fetchall(), columns=appointmentdf.columns)
24
25 # Close the cursor and connection
26 cur.close()
27 dbconn.close()
```


SQL query in PostgreSQL

Query

Query History

1

2

SELECT

appointmentid, appointmentdate, appointmenttime, patientid, providerid, appointmenttype, appointmentstatus

FROM public.appointments;

Data Output

Messages

Notifications

	appointmentid [PK] integer	appointmentdate date	appointmenttime time without time zone	patientid integer	providerid integer	appointmenttype character varying (50)	appointmentstatus character varying (50)
1	1	2023-10-10	09:00:00	1	1	General Cardiology Check-up	Scheduled
2	2	2023-10-11	10:30:00	2	2	Orthopedic Assessment	Confirmed
3	3	2023-10-12	14:15:00	3	1	Cardiology Check-up	Scheduled
4	4	2023-10-13	11:45:00	4	4	Neurological Evaluation	Confirmed
5	5	2023-10-14	08:30:00	5	5	Dermatological Examination	Scheduled
6	6	2023-10-15	13:30:00	6	2	Orthopedic Assessment	Confirmed
7	7	2023-10-16	10:15:00	7	5	Dermatological Examination	Scheduled
8	8	2023-10-17	15:45:00	8	1	Cardiology Check-up	Confirmed
9	9	2023-10-18	09:30:00	9	4	Neurological Evaluation	Scheduled
10	10	2023-10-19	11:00:00	10	3	Internal Medicine Check-up	Confirmed

Loading Procedure data into Procedure table in SP703db

```
1 #creating connection
2 dbconn = create_connect()
3
4 # Create a cursor
5 cur = dbconn.cursor()
6
7 # Insert data from csv file to Procedures table in SP703 db
8 for item, row in proceddf.iterrows():
9     cur.execute(
10         '''INSERT INTO procedures (ProcedureID, ProcedureName, ProcedureDescription,
11                                     ProcedureCost)
12         VALUES (%s, %s, %s, %s) ON CONFLICT (ProcedureID) DO NOTHING''',
13         (row['ProcedureID'], row['ProcedureName'], row['ProcedureDescription'],
14           row['ProcedureCost']))
15
16
17 # commit the changes
18 dbconn.commit()
19
20 # Execute SQL queries to check if we successfully inserted the data
21 cur.execute('SELECT * FROM procedures')
22 proceddb_tbl = pd.DataFrame(cur.fetchall(), columns=proceddf.columns)
23
24 # Close the cursor and connection
25 cur.close()
26 dbconn.close()
```

SQL query in PostgreSQL

```
1 SELECT procedureid, procedurename, proceduredescription, procedurecost
2 FROM public.procedures;
```

Data Output Messages Notifications

	procedureid [PK] integer	procedurename character varying (255)	proceduredescription text	procedurecost numeric (10,2)
1	1	Cardiology Check-up	Comprehensive cardiology check-up including ECG and echocardiogram.	250.00
2	2	Orthopedics Consultation	Orthopedic assessment and consultation, including X-ray examination if needed.	150.00
3	3	Internal Medicine Check-up	Comprehensive internal medicine check-up and health assessment.	200.00
4	4	Neurology Consultation	Neurological evaluation, consultation, and diagnostic tests as required.	180.00
5	5	Dermatology Check-up	Dermatological examination, consultation, and minor skin procedures if necessa...	120.00

Loading AppointmentProcedure data into AppointmentProcedure table in SP703db

```
1 #creating connection
2 dbconn = create_connect()
3
4 # Create a cursor
5 cur = dbconn.cursor()
6
7 # Insert data from csv file to AppointmentProcedure table in SP703 db
8 for item, row in apptproceddf.iterrows():
9     cur.execute(
10         '''INSERT INTO appointmentprocedure (AppointmentProcedureID, AppointmentID, ProcedureID)
11            VALUES (%s, %s, %s) ON CONFLICT (AppointmentProcedureID) DO NOTHING''',
12         (row['AppointmentProcedureID'], row['AppointmentID'], row['ProcedureID']))
13
14
15 # commit the changes
16 dbconn.commit()
17
18 # Execute SQL queries to check if we successfully inserted the data
19 cur.execute('SELECT * FROM appointmentprocedure')
20 apptproceddb_tbl = pd.DataFrame(cur.fetchall(), columns=apptproceddf.columns)
21
22 # Close the cursor and connection
23 cur.close()
24 dbconn.close()
```

SQL query in PostgreSQL

```

1 SELECT appointmentprocedureid, appointmentid, procedureid
2 FROM public.appointmentprocedure;

```

	appointmentprocedureid [PK] integer	appointmentid integer	procedureid integer
1	1	1	1
2	2	2	2
3	3	3	1
4	4	4	4
5	5	5	5
6	6	6	2
7	7	7	5
8	8	8	1
9	9	9	4
10	10	10	3

CREATING DATA WAREHOUSE

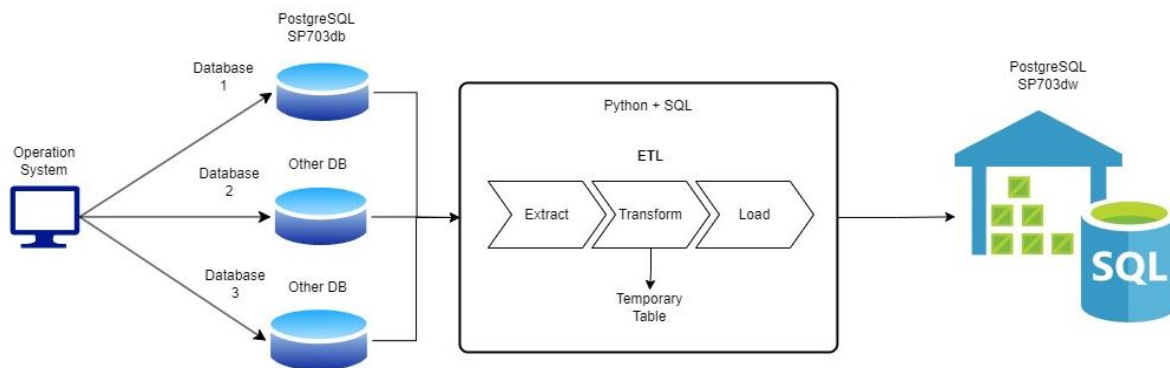


Figure 6. Data warehouse architecture

This is the diagram of the data warehouse architecture for the proposed system. As discussed on page 2 - 3, the Fact table (Provider Performance) which is in the SP703dw data warehouse will get its dimensions and data from 3 databases.

SQL script for Data warehouse in PostgreSQL (Based on the Dimensional Model Figure 4 on page 4)

Query	Query History
<pre> 1 -- Create Dimension Tables 2 CREATE TABLE Time_dim (3 TimeID INT NOT NULL PRIMARY KEY, 4 Date DATE, 5 DayOfWeek INT, 6 Month INT, 7 Quarter INT, 8 Year INT 9); </pre>	<pre> 11 CREATE TABLE Provider_dim (12 ProviderID INT NOT NULL PRIMARY KEY, 13 FirstName VARCHAR(255), 14 LastName VARCHAR(255), 15 Specialization VARCHAR(255), 16 LicenseNumber VARCHAR(255), 17 ContactNumber VARCHAR(20), 18 Email VARCHAR(255) 19); </pre>

```

21 CREATE TABLE Specialization_dim (
22     SpecializationID INT NOT NULL PRIMARY KEY,
23     Specialization VARCHAR(255)
24 );
25
26 -- Create Fact Table
27 CREATE TABLE ProviderPerformance_fact (
28     PerformanceID INT PRIMARY KEY,
29     TimeID INT,
30     ProviderID INT,
31     SpecializationID INT,
32     AppointmentCount INT,
33     AverageSatisfactionRating DECIMAL(5, 2),
34     AverageAppointmentDuration DECIMAL(5, 2),
35     AverageWaitTime DECIMAL(5, 2),
36     FOREIGN KEY (TimeID) REFERENCES Time(TimeID),
37     FOREIGN KEY (ProviderID) REFERENCES Provider(ProviderID),
38     FOREIGN KEY (SpecializationID) REFERENCES Specialization(SpecializationID)
39 );

```

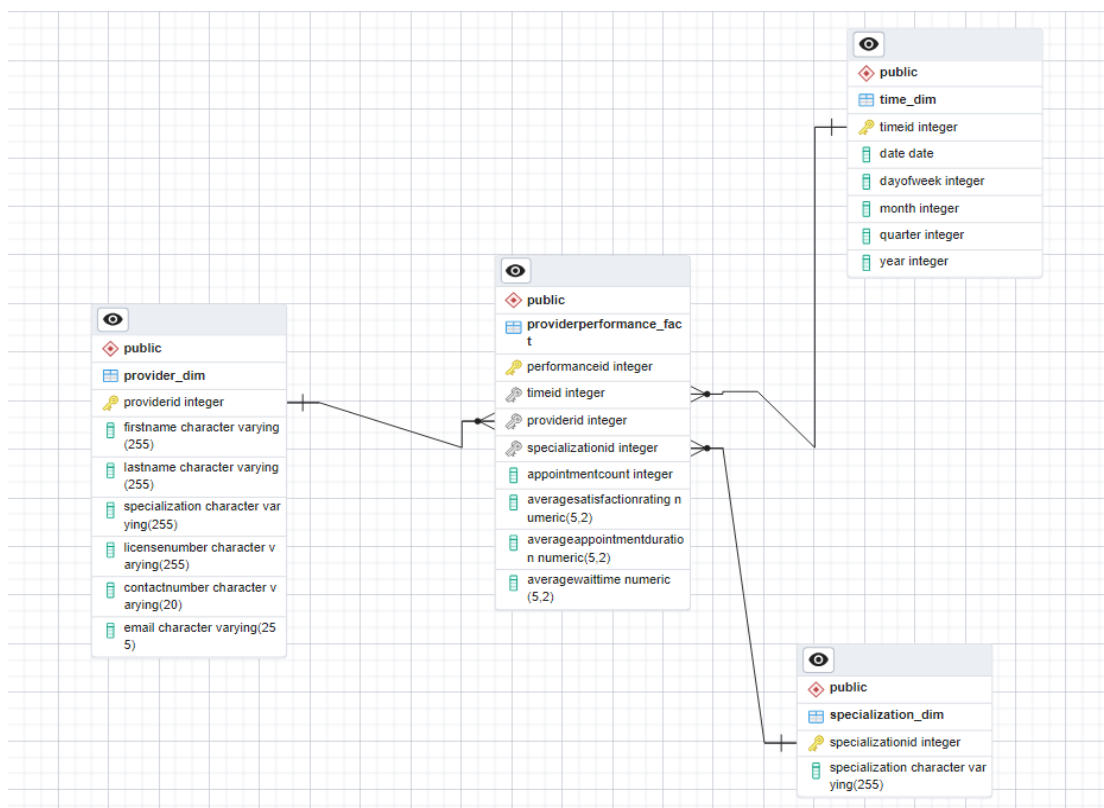


Figure 7. Data warehouse SP703dw (PostgreSQL)

Next is we need to extract, then transform the data from our database SP703db to ensure it aligns with the format and structure required before loading it to the SP703dw data warehouse. See Figure 8 on page 11.

EXTRACTING DATA FROM SP703db DATABASE

Python + SQL script

Providers data from SP703db database

```
1 #creating connection
2 dbconn = create_connect()
3
4 # Create a cursor
5 cur = dbconn.cursor()
6
7 # Execute SQL queries to retrieved data
8 cur.execute('SELECT * FROM providers')
9 providerdf = pd.DataFrame(cur.fetchall(), columns=providerdf.columns)
10
11 # Close the cursor and connection
12 cur.close()
13 dbconn.close()
```

Result from fetching data from database SP703db

This data doesn't need transformation as it has the same structure in the SP703dw data warehouse.

1	providerdf						
	ProviderID	FirstName	LastName	Specialization	LicenseNumber	ContactNumber	Email
0	1	Dr. Sarah	Anderson	Cardiology	MD123456	555-123-7890	sarah.anderson@example.com
1	2	Dr. Robert	Clark	Orthopedics	MD789012	444-987-6543	robert.clark@example.com
2	3	Dr. Lisa	Smith	Internal Medicine	MD345678	777-333-2222	lisa.smith@example.com
3	4	Dr. John	Wilson	Neurology	MD234567	888-555-9999	john.wilson@example.com
4	5	Dr. Emily	Jones	Dermatology	MD567890	999-111-2222	emily.jones@example.com

Appointments data from SP703db database

```
1 #creating connection
2 dbconn = create_connect()
3
4 # Create a cursor
5 cur = dbconn.cursor()
6
7 # Execute SQL queries to retrieved data
8 cur.execute('SELECT appointmentid, appointmentdate FROM appointments')
9 appointdf = pd.DataFrame(cur.fetchall(), columns=appointmentdb_tbl.columns[:2])
10
11 # Close the cursor and connection
12 cur.close()
13 dbconn.close()
```

Result from fetching data from database SP703db.

Next is we need to manipulate this data during transformation stage to meet the required structure before loading it to SP703dw data warehouse

```
1 # result from fetching data from database SP703db
2 appointdf
```

	AppointmentID	AppointmentDate	Date	DayOfWeek	Month	Quarter	Year
0	1	2023-10-10	10	Tuesday	October	4	2023
1	2	2023-10-11	11	Wednesday	October	4	2023
2	3	2023-10-12	12	Thursday	October	4	2023
3	4	2023-10-13	13	Friday	October	4	2023
4	5	2023-10-14	14	Saturday	October	4	2023
5	6	2023-10-15	15	Sunday	October	4	2023
6	7	2023-10-16	16	Monday	October	4	2023
7	8	2023-10-17	17	Tuesday	October	4	2023
8	9	2023-10-18	18	Wednesday	October	4	2023
9	10	2023-10-19	19	Thursday	October	4	2023

Specialization data from "other database"

Note: We will just assume that this data comes from "other database" within the medical clinic, this has already been discussed in the documentation figure 8 page 11

```
1 # this is a List of specialization by the providers, but it will add more if the provider does have two specialization.
2 # but for now, we will just assume a 1:1 ratio.
3 spec = pd.read_excel('Specializationdata.xlsx')
4 spec
```

	SpecializationID	Specialization
0	1	Cardiology
1	2	Orthopedics
2	3	Internal Medicine
3	4	Neurology
4	5	Dermatology

Fact Table data from "other database" and SP703db

Note: We will just assume that this data comes from "other database" within the medical clinic, this has already been discussed in the documentation figure 8 page 11

```
1 spec_dict = {}
2
3 # storing the providers id and specialization into a dictionary
4 for i, r in providerdf.iterrows():
5     spec_dict[r['Specialization']] = r['ProviderID']
6
7 spec_dict
```

```
{'Cardiology': 1,
 'Orthopedics': 2,
 'Internal Medicine': 3,
 'Neurology': 4,
 'Dermatology': 5}
```

```
1 # data for AppointmentCount, AverageSatisfactionRating, AverageAppointmentDuration, and AverageWaitingTime
2 mixdata = pd.read_excel('Ratingsetcdata.xlsx')
3 mixdata
```

	AppointmentCount	AverageSatisfactionRating	AverageAppointmentDuration	AverageWaitingTime
0	1	4.6	30	5
1	1	4.4	40	6
2	1	4.2	25	4
3	1	4.1	45	3
4	1	4.8	35	5
5	1	4.5	38	6
6	1	4.7	42	1
7	1	4.3	28	10
8	1	4.0	50	5
9	1	4.9	33	6

Note: Code block below will extract data comes from the SP703db database not from the "other database"

```
1 #creating connection
2 dbconn = create_connect()
3
4 # Create a cursor
5 cur = dbconn.cursor()
6
7 # new column names for extracted data
8 cols = ['TimeID', 'ProviderID']
9
10 # Execute SQL queries to retrieved data
11 cur.execute('SELECT appointmentid, providerid FROM appointments')
12 patxprovd = pd.DataFrame(cur.fetchall(), columns=cols)
13
14 # Close the cursor and connection
15 cur.close()
16 dbconn.close()
```

1	patxprovd	
	TimeID	ProviderID
0	1	1
1	2	2
2	3	1
3	4	4
4	5	5
5	6	2
6	7	5
7	8	1
8	9	4
9	10	3

TRANSFORMING DATA FROM SP703db DATABASE

Next is we need to transform some of the data based on the structure of Dimension and Fact tables in SP703dw Datawarehouse. Transforming the Appointment table.

```
1 # Convert 'AppointmentDate' column to datetime
2 appointdf['AppointmentDate'] = pd.to_datetime(appointdf['AppointmentDate'])
3
4 # Split the 'AppointmentDate' column into separate columns
5 appointdf['Date'] = appointdf['AppointmentDate'].dt.date
6 appointdf['DayOfWeek'] = appointdf['AppointmentDate'].dt.day
7 appointdf['Month'] = appointdf['AppointmentDate'].dt.month
8 appointdf['Quarter'] = appointdf['AppointmentDate'].dt.quarter
9 appointdf['Year'] = appointdf['AppointmentDate'].dt.year
10
11 # dropping unnecessary column
12 trfndappointdf = appointdf.drop(columns='AppointmentDate')
13 trfndappointdf
```

	AppointmentID	Date	DayOfWeek	Month	Quarter	Year
0	1	2023-10-10	10	10	4	2023
1	2	2023-10-11	11	10	4	2023
2	3	2023-10-12	12	10	4	2023
3	4	2023-10-13	13	10	4	2023
4	5	2023-10-14	14	10	4	2023
5	6	2023-10-15	15	10	4	2023
6	7	2023-10-16	16	10	4	2023
7	8	2023-10-17	17	10	4	2023
8	9	2023-10-18	18	10	4	2023
9	10	2023-10-19	19	10	4	2023

Merging the data from SP703 database and "other database" for the Provider Performance Fact Table

```
1 # Add a new column 'Specialization' based on 'ProviderID'
2 patxprovdf['SpecializationID'] = patxprovdf['ProviderID'].replace(spec_dict)
```

```
1 factdw_tbl = pd.concat([patxprovdf, mixdata], axis=1)
2 factdw_tbl
```

TimeID	ProviderID	SpecializationID	AppointmentCount	AverageSatisfactionRating	AverageAppointmentDuration	AverageWaitingTime
0	1	1	1	4.6	30	5
1	2	2	1	4.4	40	6
2	3	1	1	4.2	25	4
3	4	4	1	4.1	45	3
4	5	5	1	4.8	35	5
5	6	2	1	4.5	38	6
6	7	5	1	4.7	42	1
7	8	1	1	4.3	28	10
8	9	4	1	4.0	50	5
9	10	3	1	4.9	33	6

LOADING DATA TO SP70dw DATA WAREHOUSE

We will load the transformed data into SP703dw Datawarehouse but first we need to connect to PostgreSQL to access the data warehouse.

Setting credentials to access data warehouse to establish a data warehouse connection.

```
1 PGHOST = 'localhost'
2 PGDATAWAREHOUSE = 'SP703dw'
3 PGUSER = '*****' # need to hide sensitive info
4 PGPASSWORD = '*****' # need to hide sensitive info
```

```
1 def create_connectdw():
2     dwconn = psycopg2.connect(
3         dbname=PGDATAWAREHOUSE,
4         user=PGUSER,
5         password=PGPASSWORD,
6         host=PGHOST,
7         port='5432'
8     )
9     return dwconn
```

Load the Providers data from SP702 database to Providers Dimension in SP703dw data warehouse.

```
1 #creating connection
2 dwconn = create_connectdw()
3
4 # Create a cursor
5 cur = dwconn.cursor()
6
7 # Insert data from providers table into SP703 dw provider_dim
8 for item, row in providerdf.iterrows():
9     cur.execute(
10         '''INSERT INTO provider_dim (ProviderID, FirstName, LastName,
11                                     Specialization, LicenseNumber, ContactNumber, Email)
12         VALUES (%s, %s, %s, %s, %s, %s, %s) ON CONFLICT (ProviderID) DO NOTHING''',
13         (row['ProviderID'], row['FirstName'], row['LastName'],
14          row['Specialization'], row['LicenseNumber'], row['ContactNumber'],
15          row['Email']))
16
17
18 # commit the changes
19 dwconn.commit()
20
21 # Execute SQL queries to check if we successfully inserted the data
22 cur.execute('SELECT * FROM provider_dim')
23 providerdw_tbl = pd.DataFrame(cur.fetchall(), columns=providerdf.columns)
24
25 # Close the cursor and connection
26 cur.close()
27 dwconn.close()
```

SQL query in PostgreSQL

Query

Query History

1

2

SELECT

providerid, firstname, lastname, specialization, licensenumber, contactnumber, email

FROM public.provider_dim;

Data Output

Messages

Notifications

	providerid [PK] integer	firstname character varying (255)	lastname character varying (255)	specialization character varying (255)	licensenumber character varying (255)	contactnumber character varying (20)	email character varying (255)
1	1	Dr. Sarah	Anderson	Cardiology	MD123456	555-123-7890	sarah.anderson@example.com
2	2	Dr. Robert	Clark	Orthopedics	MD789012	444-987-6543	robert.clark@example.com
3	3	Dr. Lisa	Smith	Internal Medicine	MD345678	777-333-2222	lisa.smith@example.com
4	4	Dr. John	Wilson	Neurology	MD234567	888-555-9999	john.wilson@example.com
5	5	Dr. Emily	Jones	Dermatology	MD567890	999-111-2222	emily.jones@example.com

Load the transformed Appointment data from SP702 database to Time Dimension in SP703dw data warehouse

```
1 #creating connection
2 dwconn = create_connectdw()
3
4 # Create a cursor
5 cur = dwconn.cursor()
6
7 # Insert data from transformed Appointment table into SP703 dw time_dim
8 for item, row in trfndappointdf.iterrows():
9     cur.execute(
10         '''INSERT INTO time_dim (TimeID, Date, DayOfWeek, Month, Quarter, Year)
11           VALUES (%s, %s, %s, %s, %s, %s) ON CONFLICT (TimeID) DO NOTHING''',
12         (row['AppointmentID'], row['Date'], row['DayOfWeek'],
13          row['Month'], row['Quarter'], row['Year']))
14
15
16 # commit the changes
17 dwconn.commit()
18
19 #we need to declare new columns
20 cols = ['TimeID', 'Date', 'DayOfWeek', 'Month', 'Quarter', 'Year']
21
22 # Execute SQL queries to check if we successfully inserted the data
23 cur.execute('SELECT * FROM time_dim')
24 timedw_tbl = pd.DataFrame(cur.fetchall(), columns=cols)
25
26 # Close the cursor and connection
27 cur.close()
28 dwconn.close()
```

SQL query in PostgreSQL

Query

Query History

1

2

SELECT

timeid,

date,

dayofweek,

month,

quarter,

year

FROM

public.time_dim;

Data Output

Messages

Notifications

≡

+

📄

▼

📋

▼

🗑️

📦

⬇️

📈

	timeid [PK] integer	date date	dayofweek integer	month integer	quarter integer	year integer
1	1	2023-10-10	10	10	4	2023
2	2	2023-10-11	11	10	4	2023
3	3	2023-10-12	12	10	4	2023
4	4	2023-10-13	13	10	4	2023
5	5	2023-10-14	14	10	4	2023
6	6	2023-10-15	15	10	4	2023
7	7	2023-10-16	16	10	4	2023
8	8	2023-10-17	17	10	4	2023
9	9	2023-10-18	18	10	4	2023
10	10	2023-10-19	19	10	4	2023

Load the Specialization data from "other database" to Specialization Dimension in SP703dw data warehouse

```

1 #creating connection
2 dwconn = create_connectdw()
3
4 # Create a cursor
5 cur = dwconn.cursor()
6
7 # Insert data from "other database" table in SP703 dw
8 for item, row in spec.iterrows():
9     cur.execute(
10         '''INSERT INTO specialization_dim (SpecializationID, Specialization)
11         VALUES (%s, %s) ON CONFLICT (SpecializationID) DO NOTHING''',
12         (row['SpecializationID'], row['Specialization']))
13
14
15 # commit the changes
16 dwconn.commit()
17
18 # Execute SQL queries to check if we successfully inserted the data
19 cur.execute('SELECT * FROM specialization_dim')
20 specdw_tbl = pd.DataFrame(cur.fetchall(), columns=spec.columns)
21
22 # Close the cursor and connection
23 cur.close()
24 dwconn.close()

```

SQL query in PostgreSQL

Query	Query History
1	SELECT specializationid, specialization
2	FROM public.specialization_dim;

Data Output	Messages	Notifications
<div> <div>+</div> <div>📄</div> <div>▼</div> <div>📋</div> <div>▼</div> <div>🗑️</div> <div>📦</div> <div>📥</div> <div>📈</div> </div>		
specializationid [PK] integer	specialization character varying (255)	
1	Cardiology	
2	Orthopedics	
3	Internal Medicine	
4	Neurology	
5	Dermatology	

Load the data from SP703db database to Provider Performance Fact in SP703dw data warehouse.

```

1 #creating connection
2 dwconn = create_connectdw()
3
4 # Create a cursor
5 cur = dwconn.cursor()
6
7 # getting the latest key
8 latest_key = 0
9 cur.execute('SELECT MAX(performanceid) AS key FROM providerperformance_fact')
10 get_key = pd.DataFrame(cur.fetchall(), columns=['key'])
11 if get_key['key'].item() == None:
12     latest_key = 0
13 else:
14     latest_key = get_key['key'].item()
15
16
17 # Insert data from other database + SP703db table into providerperformance fact table in SP703dw
18 for item, row in factdw_tbl.iterrows():
19     latest_key+=1
20     cur.execute(
21         '''INSERT INTO providerperformance_fact (PerformanceID, TimeID, ProviderID, SpecializationID,
22             AppointmentCount, AverageSatisfactionRating, AverageAppointmentDuration, AverageWaitTime)
23             VALUES (%s, %s, %s, %s, %s, %s, %s, %s) ON CONFLICT (PerformanceID) DO NOTHING''',
24             (latest_key, row['TimeID'], row['ProviderID'], row['SpecializationID'],
25             row['AppointmentCount'], row['AverageSatisfactionRating'],
26             row['AverageAppointmentDuration'], row['AverageWaitingTime']))
27
28 # commit the changes
29 dwconn.commit()
30
31 cols = ['PerformanceID', 'TimeID', 'ProviderID', 'SpecializationID',
32         'AppointmentCount', 'AverageSatisfactionRating', 'AverageAppointmentDuration', 'AverageWaitTime']
33
34 # Execute SQL queries to check if we successfully inserted the data
35 cur.execute('SELECT * FROM providerperformance_fact')
36 specdw_tbl = pd.DataFrame(cur.fetchall(), columns=cols)
37
38 # Close the cursor and connection
39 cur.close()
40 dwconn.close()

```

SQL query in PostgreSQL

Query

Query History

```
1 SELECT performanceid, timeid, providerid, specializationid, appointmentcount, averagesatisfactionrating,
2     averageappointmentduration,
3     averagewaittime
4 FROM public.providerperformance_fact;
```

Data Output

Messages

Notifications

	performanceid [PK] integer	timeid integer	providerid integer	specializationid integer	appointmentcount integer	averagesatisfactionrating numeric (5,2)	averageappointmentduration numeric (5,2)	averagewaittime numeric (5,2)
1	1	1	1	1	1	4.60	30.00	5.00
2	2	2	2	2	1	4.40	40.00	6.00
3	3	3	1	1	1	4.20	25.00	4.00
4	4	4	4	4	1	4.10	45.00	3.00
5	5	5	5	5	1	4.80	35.00	5.00
6	6	6	2	2	1	4.50	38.00	6.00
7	7	7	5	5	1	4.70	42.00	1.00
8	8	8	1	1	1	4.30	28.00	10.00
9	9	9	4	4	1	4.00	50.00	5.00
10	10	10	3	3	1	4.90	33.00	6.00

After loading all of the data to datawarehouse SP703dw, we will do some test queries.

Query

Query History

```
1 -- sample query in the datawarehouse SP703dw
2
3 SELECT  prov.firstname,
4         ROUND(AVG(averagesatisfactionrating), 2) AS Ratings,
5         spec.specialization
6 FROM public.providerperformance_fact per
7     LEFT JOIN public.provider_dim prov ON per.providerid = prov.providerid
8     LEFT JOIN public.specialization_dim spec ON per.specializationid = spec.specializationid
9 GROUP BY prov.firstname, spec.specialization;
```

Data Output

Messages

Notifications

	firstname character varying (255)	ratings numeric	specialization character varying (255)
1	Dr. Emily	4.75	Dermatology
2	Dr. John	4.05	Neurology
3	Dr. Lisa	4.90	Internal Medicine
4	Dr. Robert	4.45	Orthopedics
5	Dr. Sarah	4.37	Cardiology

The result table above is achieved by joining the 2 dimension table with the fact table