

Final Project

Disease database

Group member: Ruoyu Chen, Wei Lo, Yu Wang

Content

Business Problem Definition:	2
Model Description:	2
Use Cases and DML Operations:	3
Analysis of Dashboard	7
Dimensional Database	8
ELT	11
Commentary on Snowflake	14
Transitioning from Relational to NoSQL Database Structure	16
Structure & AWS Implementation	错误!未定义书签。

Business Problem Definition:

A Disease Data management system for maintaining records for hospitals, their staff (doctors and nurses), patients, diseases, diagnoses, medications, and related information. The system aims to organize and manage the following entities:

1. Address Information:
Contains details about different addresses such as city, state, and country. Hospitals are associated with specific addresses.
2. Hospital Information:
Includes hospital names and their respective addresses. This table stores unique identifiers for hospitals.
3. Doctor and Nurse Information:
Stores records of doctors and nurses associated with hospitals. It includes their names, genders, qualifications, departments, and references the hospitals they work in.
4. Disease Information:
Manages records of various diseases, including their names and types.
5. Diagnosis Information:
Contains data about diagnoses, including issue dates, treatments, and references to nurses and doctors responsible for the diagnosis.
6. Medicine Information:
Maintains details about medicines, including names, prices, companies producing them, active ingredients, and standard industry numbers.
7. Diagnosis-Medicine Relationships:
Represents the relationship between diagnoses and medicines, including quantities of medicines prescribed for specific diagnoses.
8. Patient Information:
Stores patient details such as names, ages, genders, races, and references to their diseases, diagnoses, home addresses, and associated hospitals.

Model Description:

In a disease-centered data model, the core entity is the "Disease." This entity typically holds a one-to-many relationship with the "Diagnosis" entity, as multiple diagnoses can be associated with one disease. Each "Diagnosis" is linked to a "Patient," creating a many-to-one relationship, as multiple patients can be diagnosed with the same disease. Patients, in turn, have a one-to-many relationship with "Medication," as a patient might receive various medications for a single diagnosis. Additionally, patients may be connected to "Doctor" entities through the "Diagnosis" entity, indicating the doctor responsible for the diagnosis, thus creating a many-to-one relationship. Furthermore, "Doctor" entities can be linked to "Hospital" entities, showcasing the association of doctors with particular healthcare facilities. Overall, this interconnected structure captures the relationships between diseases, patient diagnoses, medications, healthcare providers, and healthcare facilities in the context of disease management.

Use Cases and DML Operations:

Use case 1: Get Patient Information with their Diagnoses information.

Objective: Pull all the information of the patients with their diagnoses, medication, address, hospital. And calculate the total medicine price of each diagnose. This can show almost all the data in the database.

```
SELECT
    CONCAT(p.First_Name, ' ', p.Last_Name) AS Patient_Name,
    p.Age AS Patient_Age,
    p.Gender AS Patient_Gender,
    p.Race AS Patient_Race,
    dis.Disease_Type AS Disease_Type,
    hos.Hosp_Name AS Hospital_Name,
    CONCAT(a.City_Name, ', ', a.State_Name, ', ', a.Country_Name) AS
Address,
    d.Issue_Date,
    d.Treatment,
    dis.Disease_Name,
    m.Med_Name
FROM
    Patient p
JOIN
    Diagnosis d ON p.Diagnos_ID = d.Diagnos_ID
JOIN
    Disease dis ON p.Disease_ID = dis.Disease_ID
LEFT JOIN
    DiagnosisMedicine dm ON d.Diagnos_ID = dm.Diagnos_ID
LEFT JOIN
    Medicine m ON dm.Med_ID = m.Med_ID
LEFT JOIN
    Hospital hos ON p.Hosp_ID = hos.Hosp_ID
LEFT JOIN
    "Address" a ON p.Address_ID = a.Address_ID;
```

Result:

	patient_name	patient_age	patient_gender	patient_race	disease_type	hospital_name	address	issue_date	treatment	disease_name	med_name
	text	integer	character	character varying (20)	character varying (200)	character varying (100)	text	date	character varying (30)	character varying (200)	character varying (20)
1	John Doe	35	M	Caucasian	Viral Infection	New York General Hospital	San Francisco, California, USA	2023-01-15	Prescribed medication	Influenza	Aspirin
2	Jane Smith	50	F	African American	Metabolic Disorder	Los Angeles Medical Center	San Francisco, California, USA	2023-02-20	Physical therapy	Diabetes	Amoxicillin
3	David Johnson	28	M	Asian	Cardiovascular	London Health Care	Miami, Florida, USA	2023-03-25	Surgery recommended	Hypertension	Lipitor
4	Emily Davis	28	F	Caucasian	Neurological Disorder	Chicago General Hospital	Miami, Florida, USA	2023-04-10	Prescribed medication	Migraine	Amoxicillin
5	Michael Wilson	45	M	African American	Viral Infection	Houston Medical Center	Miami, Florida, USA	2023-05-20	Physical therapy	Asthma	Metformin
6	Sophia Martinez	38	F	Hispanic	Metabolic Disorder	Philadelphia Health Care	Atlanta, Georgia, USA	2023-06-15	Surgery recommended	Arthritis	Atorvastatin
7	Olivia Anderson	32	F	Asian	Viral Infection	New York General Hospital	Boston, Massachusetts, USA	2023-07-10	Prescribed medication	Influenza	Paracetamol
8	Olivia Anderson	32	F	Asian	Viral Infection	New York General Hospital	Boston, Massachusetts, USA	2023-07-10	Prescribed medication	Influenza	Ibuprofen
9	Daniel Thompson	50	M	Caucasian	Metabolic Disorder	Los Angeles Medical Center	Denver, Colorado, USA	2023-08-20	Physical therapy	Diabetes	Ibuprofen
10	Isabella Hernandez	42	F	Hispanic	Cardiovascular	London Health Care	London, England, UK	2023-09-15	Surgery recommended	Hypertension	Omeprazole
11	William Brown	40	M	Caucasian	Neurological Disorder	Chicago General Hospital	Chicago, Illinois, USA	2023-10-10	Physical therapy	Migraine	Aspirin
12	Emma Wilson	28	F	African American	Viral Infection	Houston Medical Center	Houston, Texas, USA	2023-11-15	Prescribed medication	Influenza	Amoxicillin
13	James Garcia	55	M	Hispanic	Cardiovascular	Philadelphia Health Care	Portland, Oregon, USA	2023-12-20	Surgery recommended	Hypertension	Lipitor
14	Ara Rodriguez	33	F	Asian	Metabolic Disorder	New York General Hospital	Boston, Massachusetts, USA	2024-01-05	Physical therapy	Diabetes	Amoxicillin
15	Sophie Lopez	60	F	Hispanic	Viral Infection	London Health Care	Seattle, Washington, USA	2024-02-10	Prescribed medication	Asthma	Metformin
16	Emma Wilson	25	F	Caucasian	Viral Infection	New York General Hospital	Austin, Texas, USA	2023-10-05	Prescribed medication	Influenza	Amoxicillin

Use Case 2: Patients Count and Average Age per Disease Type

Objective: To retrieve the count of patients and their average age grouped by disease type.

```

SELECT d.Disease_Type,
       COUNT(*) AS Total_Patients,
       ROUND(AVG(p.Age), 2) AS Average_Age
FROM Patient p
JOIN Disease d ON p.Disease_ID = d.Disease_ID
GROUP BY d.Disease_Type;

```

Result:

	disease_type character varying (200) 🔒	total_patients bigint 🔒	average_age numeric 🔒
1	Cardiovascular	3	41.67
2	Neurological Disorder	2	34.00
3	Metabolic Disorder	4	42.75
4	Viral Infection	6	37.50

Use Case 3: Analyzing Disease Severity and Treatment Counts

Objective: To examine the average severity of diseases and the most frequent of treatments prescribed for each disease.

DiseaseSeverity CTE: Calculates the average severity level for each disease.

TreatmentCounts CTE: Determines the count of treatments given for each diagnosis.

Main Query: Combines the CTEs to display disease names, their average severity, and treatment counts. And sorts the output by disease name and treatment count in descending order.

```

WITH DiseaseSeverity AS (
    SELECT p.Disease_ID, dis.Disease_Name, ROUND(AVG(p.Severity_lv), 2)
    AS Avg_Severity
    FROM Patient p
    JOIN Disease dis ON p.Disease_ID = dis.Disease_ID
    GROUP BY p.Disease_ID, dis.Disease_Name
),
TreatmentCounts AS (
    SELECT d.Diagnos_ID, d.Treatment, COUNT(*) AS Treatment_Count
    FROM Diagnosis d
    GROUP BY d.Diagnos_ID, d.Treatment
)
SELECT ds.Disease_Name, ds.Avg_Severity, tc.Treatment
FROM DiseaseSeverity ds
LEFT JOIN TreatmentCounts tc ON ds.Disease_ID = tc.Diagnos_ID
ORDER BY ds.Disease_Name, tc.Treatment_Count DESC;

```

Result:

	disease_name character varying (200) 🔒	avg_severity numeric 🔒	treatment character varying (30) 🔒
1	Arthritis	6.00	Surgery recommended
2	Asthma	4.50	Physical therapy
3	Diabetes	4.33	Physical therapy
4	Hypertension	4.67	Surgery recommended
5	Influenza	2.50	Prescribed medication
6	Migraine	3.50	Prescribed medication

DML operations:

In order to help maintain data integrity, we have a trigger-function pair to remove related diagnosis and medication records when a patient is deleted from the database. The function `delete_patient_related_records` contains two `DELETE` statements: First one deletes records from the `DiagnosisMedicine` table where the `Diagnos_ID` matches the `OLD.Diagnos_ID` (related to the patient being deleted). The second one deletes records from the `Diagnosis` table where the `Diagnos_ID` matches the `OLD.Diagnos_ID` (related to the patient being deleted). The trigger `delete_patient_trigger` is set to execute `delete_patient_related_records` before the deletion of a record from the `Patient` table.

```
-- Create a function that deletes related diagnoses and medication
records
CREATE OR REPLACE FUNCTION delete_patient_related_records()
RETURNS TRIGGER AS $$
BEGIN
    -- Delete DiagnosisMedicine records related to the patient being
deleted
    DELETE FROM DiagnosisMedicine WHERE Diagnos_ID = OLD.Diagnos_ID;
    -- Delete Diagnosis records related to the patient being deleted
    DELETE FROM Diagnosis WHERE Diagnos_ID = OLD.Diagnos_ID;
    RETURN OLD;
END;
$$ LANGUAGE plpgsql;

-- Create a trigger that activates the function before deleting a
patient
CREATE TRIGGER delete_patient_trigger
BEFORE DELETE ON Patient
FOR EACH ROW
EXECUTE FUNCTION delete_patient_related_records();
```

First add a new patient that has Asthma. Patient has a reference key from diagnosis, so first create the diagnosis for the patient.

```
INSERT INTO Diagnosis (Diagnos_ID, Issue_Date, Treatment, Nurse_ID,
Doc_ID)
```

VALUES

```
(20, '2023-12-18', 'Prescribed medication', 1, 3);
```

Now create the patient information.

```
INSERT INTO Patient (SSN, Disease_ID, Diagnos_ID, First_Name,
Last_Name, Age, Gender, Race, Severity_lv, Address_ID, Hosp_ID)
```

VALUES

```
(222222222, 1, 20, 'Jacky', 'Chan', 25, 'M', 'Asian', 1, 1, 1);
```

Also add the information on DiagnosisMedicine.

```
INSERT INTO DiagnosisMedicine (Diagnos_ID, Med_ID, Med_Qty)
```

VALUES

```
(20, 1, 5);
```

Check the insert:

SELECT

```
    p.First_Name || ' ' || p.Last_Name AS Patient_Name,
    p.Age AS Patient_Age,
    p.Gender AS Patient_Gender,
    p.Race AS Patient_Race,
    d.Issue_Date,
    d.Treatment,
    dis.Disease_Name AS Disease_Type,
    h.Hosp_Name AS Hospital_Name,
    a.City_Name || ', ' || a.State_Name || ', ' || a.Country_Name AS
```

Address

FROM Patient p

JOIN Diagnosis d ON p.Diagnos_ID = d.Diagnos_ID

JOIN Disease dis ON p.Disease_ID = dis.Disease_ID

LEFT JOIN DiagnosisMedicine dm ON d.Diagnos_ID = dm.Diagnos_ID

LEFT JOIN Medicine m ON dm.Med_ID = m.Med_ID

JOIN Hospital h ON p.Hosp_ID = h.Hosp_ID

JOIN "Address" a ON p.Address_ID = a.Address_ID

WHERE p.SSN = 222222222;

	patient_name text	patient_age integer	patient_gender character	patient_race character var	issue_date date	treatment character varying (30)	disease_type character var	hospital_name character varying (100)	address text
1	Jacky Chan	25	M	Asian	2023-12-18	Prescribed medication	Influenza	New York General Hospital	New York City, New York, USA

Now let's remove it and check if the diagnosis still exists.

```
DELETE FROM Patient WHERE SSN = 222222222;
```

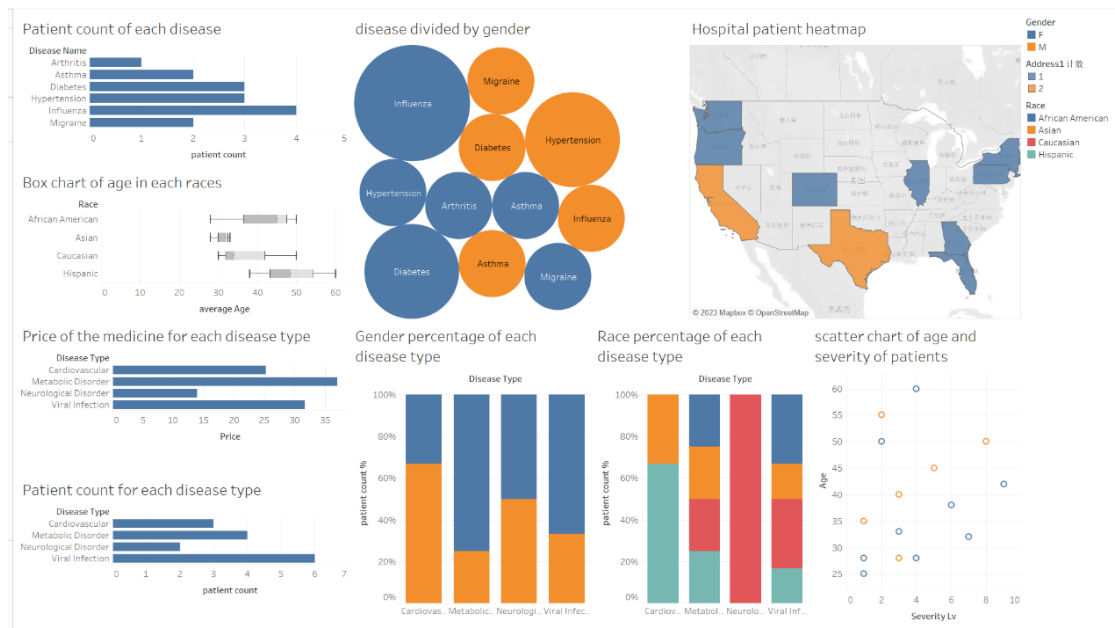
```
SELECT * FROM Diagnosis;
```

18	18	2023-11-02	Prescribed medication	2	3
19	19	2023-10-18	Physical therapy	1	1

Diagnosis 20 has been removed.

Analysis of Dashboard

We create a dashboard using tableau.



- Patient count for each disease type: This chart shows the number of patients with each disease type. The most common diseases are cardiovascular diseases, metabolic disorders, and neurological disorders.
- Hospital patient heat map: This chart shows the number of patients with each disease by state. The states with the most patients are California, Texas, and Florida.
- Patient count of each disease: This chart shows the number of patients with each disease by gender. There are more female patients than male patients for all diseases except hypertension, which is more common in men.
- Box chart of age in each race: This chart shows the box chart of patients with each disease, broken down by race. African American patients have the highest average age for all diseases except influenza, for which Hispanic patients have the highest average age.
- Price of the medicine for each disease type: This table shows the average price of the medicine for each disease type. The most expensive diseases to treat are cardiovascular diseases and neurological disorders.
- Gender percentage of each disease type: This chart shows the percentage of patients with each disease who are male or female.
- Race percentage of each disease type: This chart shows the percentage of patients with each disease who are of each race.
- Scatter chart of age and severity of patients: This chart shows the relationship between the age of patients and the severity of their disease. There is a general trend that older patients have more severe diseases.

Dimensional Database

This code is the design of a medical data warehouse, which includes the creation of dimension tables and fact tables. The dimension table describes various dimensions in the business, such as address, hospital, doctor, nurse, disease and medicine, etc. The fact table contains business facts such as diagnosis, diagnosis and drug correlation, and patient information.

Dimension Table:

Dim_Address: Stores address information, including city, state, and country.

Dim_Hospital: Describes hospital information, including hospital name and address ID associated with the address table.

Dim_Doctor: Doctor information table, recording the doctor's name, gender, qualifications, department and hospital ID.

Dim_Nurse: Nurse information table, including the nurse's name, gender and hospital ID.

Dim_Disease: Stores disease information, including disease name and type.

Dim_Medicine: Describes drug information, including drug name, price, production company, active ingredients, and industry standard number.

```
-- Create dimension table
-- Dim_Address
CREATE TABLE Dim_Address (
    Address_ID SERIAL PRIMARY KEY,
    City_Name VARCHAR(100),
    State_Name VARCHAR(100),
    Country_Name VARCHAR(100)
);

-- Dim_Hospital
CREATE TABLE Dim_Hospital (
    Hosp_ID SERIAL PRIMARY KEY,
    Hosp_Name VARCHAR(100),
    Address_ID INTEGER REFERENCES Dim_Address(Address_ID)
);

-- Dim_Doctor
CREATE TABLE Dim_Doctor (
    Doc_ID SERIAL PRIMARY KEY,
    Doc_Name VARCHAR(30),
    Gender CHAR(1),
    Qualification VARCHAR(30),
    Department VARCHAR(30),
    Hosp_ID INTEGER REFERENCES Dim_Hospital(Hosp_ID)
);

-- Dim_Nurse
CREATE TABLE Dim_Nurse (
```



```

    Nurse_ID SERIAL PRIMARY KEY,
    Nurse_Name VARCHAR(30),
    Gender CHAR(1),
    Hosp_ID INTEGER REFERENCES Dim_Hospital(Hosp_ID)
);

-- Dim_Disease
CREATE TABLE Dim_Disease (
    Disease_ID SERIAL PRIMARY KEY,
    Disease_Name VARCHAR(200),
    Disease_Type VARCHAR(200)
);

-- Dim_Medicine
CREATE TABLE Dim_Medicine (
    Med_ID SERIAL PRIMARY KEY,
    Med_Name VARCHAR(20),
    Price DOUBLE PRECISION,
    Company VARCHAR(30),
    Active_Ingredient VARCHAR(30),
    Standard_Industry_Number VARCHAR(25)
);

```

Fact Table:

Fact_Diagnosis: Record diagnosis information, including diagnosis date, treatment method, and IDs associated with nurses and doctors.

Fact_DiagnosisMedicine: Information related to diagnosis and medicine, including diagnosis ID, medicine ID and medicine quantity.

Fact_Patient: Stores patient information including social security number, disease ID, name, age, gender, race, diagnosis ID, disease severity, address ID, and hospital ID.

```

-- Create fact table
-- Fact_Diagnosis
CREATE TABLE Fact_Diagnosis (
    Diagnos_ID SERIAL PRIMARY KEY,
    Issue_Date DATE,
    Treatment VARCHAR(30),
    Nurse_ID INTEGER REFERENCES Dim_Nurse(Nurse_ID),
    Doc_ID INTEGER REFERENCES Dim_Doctor(Doc_ID)
);

-- Fact_DiagnosisMedicine
CREATE TABLE Fact_DiagnosisMedicine (
    Diagnos_ID INTEGER REFERENCES Fact_Diagnosis(Diagnos_ID),
    Med_ID INTEGER REFERENCES Dim_Medicine(Med_ID),

```

```

    Med_Qty INTEGER,
    PRIMARY KEY (Diagnos_ID, Med_ID)
);

-- Fact_Patient
CREATE TABLE Fact_Patient (
    SSN INTEGER PRIMARY KEY,
    Disease_ID INTEGER REFERENCES Dim_Disease(Disease_ID),
    First_Name VARCHAR(100),
    Last_Name VARCHAR(100),
    Age INTEGER,
    Gender CHAR(1),
    Race VARCHAR(20),
    Diagnos_ID INTEGER REFERENCES Fact_Diagnosis(Diagnos_ID),
    Severity_lv INTEGER,
    Address_ID INTEGER REFERENCES Dim_Address(Address_ID),
    Hosp_ID INTEGER REFERENCES Dim_Hospital(Hosp_ID)
);

```

This database schema can be used to record and analyze patient diagnoses, drug usage, and related information about doctors, nurses, and hospitals in medical institutions. This database model aims to provide a clear way to analyze patient diagnosis, drug use, and distribution of medical teams in medical institutions, providing data support for medical decision-making and process optimization.

ELT

This code defines a function that executes the ELT (Extract, Load, Transform) process:

ELT function (perform_elt_process()):

Extract: Extract data from the table named "public". Patient into the temporary table temp_patient.

Load: Create a target table named target_patient and load data from the temporary table into the target table.

Transform: No specific data conversion is performed here, just copy the data from one table to another, and clean up the temporary table at the end.

```
-- Create a function to execute the ELT process
CREATE OR REPLACE FUNCTION perform_elt_process()
RETURNS VOID AS $$
BEGIN
    -- Extract data to a temporary table
    CREATE TEMP TABLE temp_patient AS
    SELECT *
    FROM "public".Patient;

    -- Create the target table
    CREATE TABLE IF NOT EXISTS "dimension".Fact_Patient (
        SSN INTEGER PRIMARY KEY,
        Disease_ID INTEGER,
        First_Name VARCHAR(100),
        Last_Name VARCHAR(100),
        Age INTEGER,
        Gender CHAR(1)
        -- Add other columns...
    );

    -- Load data into the target table
    INSERT INTO "dimension".Fact_Patient (SSN, Disease_ID, First_Name,
Last_Name, Age, Gender /*, ...*/)
    SELECT SSN, Disease_ID, First_Name, Last_Name, Age, Gender /*, ...*/
    FROM temp_patient;

    -- Clean up the temporary table
    DROP TABLE IF EXISTS temp_patient;

    -- Output message
    RAISE NOTICE 'ELT process completed.';
END;
```

```

$$ LANGUAGE plpgsql;

-- Execute the ELT process
SELECT perform_elt_process();

```

Analytical query:



1. Count the number of patients by age group: Group patients by age and count the number of patients in different age groups.

```

-- 1.Number of patients by age group
SELECT
    CASE
        WHEN Age BETWEEN 0 AND 10 THEN '0-10'
        WHEN Age BETWEEN 11 AND 20 THEN '11-20'
        WHEN Age BETWEEN 21 AND 30 THEN '21-30'
        WHEN Age BETWEEN 31 AND 40 THEN '31-40'
        ELSE 'Over 40'
    END AS age_group,
    COUNT(*) AS num_patients
FROM "dimension".Fact_Patient
GROUP BY age_group
ORDER BY age_group;

```

Result:

	age_group  text	num_patients  bigint
1	21-30	4
2	31-40	5
3	Over 40	6

2. Count the number of patients by disease type: Group patients according to disease type and count the number of patients for each type.

```

-- 2.Number of patients by disease type
SELECT Disease_ID, COUNT(*) AS num_patients
FROM "dimension".Fact_Patient
GROUP BY Disease_ID
ORDER BY num_patients DESC;

```

Result:

	disease_id integer	num_patients bigint
1	1	4
2	3	3
3	2	3
4	5	2
5	4	2
6	6	1

3. Gender ratio: Count the number of male and female patients to obtain gender ratio information.

```
-- 3.Sex ratio
SELECT Gender, COUNT(*) AS num_patients
FROM "dimension".Fact_Patient
GROUP BY Gender;
```

Result:

	gender character	num_patients bigint
1	M	6
2	F	9

Commentary on Snowflake

Introduction to Snowflake

Snowflake is a cloud-based data warehousing platform that stands out for its unique architecture and capabilities. It separates compute and storage, allowing for scalable and efficient data processing. Given the context of the disease management system, let's explore how Snowflake's features could enhance the application.

Advantages of Snowflake for the Disease Management System

1. Scalability and Performance

Elastic Scalability: Snowflake can dynamically scale compute resources, adapting to the workload demands of the disease management system. This means during times of high query load or data processing requirements; Snowflake can scale up to ensure performance is not compromised.

Concurrent Processing: Snowflake's architecture allows multiple users and applications to run queries simultaneously without performance degradation. This is essential for a disease management system where multiple stakeholders (doctors, administrators, researchers) may need access to the data concurrently.

2. Data Sharing Capabilities

Secure Data Sharing: Snowflake facilitates secure sharing of data across different departments or even with external entities like research organizations. This can be highly beneficial for collaborative research or analysis without the need to copy or transfer data, maintaining data governance and security.

3. Storage and Support for Diverse Data Formats

Diverse Data Support: Snowflake supports structured and semi-structured data (like JSON, XML, Parquet). This is advantageous for integrating data from various sources in the healthcare sector, which often includes a mix of formats.

Automatic Compression and Optimization: Snowflake automatically compresses data and optimizes it for storage and retrieval, leading to cost savings and improved performance.

4. Data Governance and Security

Robust Security Features: Snowflake offers comprehensive security features like end-to-end encryption, role-based access control, and compliance with various standards (HIPAA for healthcare data). This is crucial for sensitive patient data in the disease management system.

Data Cloning and Time Travel: Features like data cloning and time travel allow for easy recovery of data and facilitate complex analyses over historical data, which can be pivotal in disease trend analysis and research.

5. Integration and Maintenance

Ease of Integration: Snowflake integrates well with various data integration tools and BI tools, allowing for a seamless connection with existing systems in the disease management framework.

Low Maintenance: Being a fully managed service, Snowflake reduces the overhead of maintenance, allowing the healthcare organization to focus on analytical and operational aspects rather than database management.

Conclusion

Incorporating Snowflake into the disease management system could significantly enhance the system's scalability, performance, and data analysis capabilities. Its ability to handle large-scale data workloads efficiently, combined with robust security features and support for diverse data types, makes it a suitable choice for managing complex and sensitive healthcare data. The flexibility and reduced maintenance burden further underscore its suitability for dynamic and data-intensive environments like healthcare and disease management.

Transitioning from Relational to NoSQL

Introduction

The purpose of this report is to explore the transition of a disease management system database from a traditional relational model to a NoSQL database structure. The focus will be on utilizing MongoDB, a widely-used NoSQL database, to illustrate the changes in database design and the implications of this transition.

Relational vs. NoSQL Database Structure

Relational Database Characteristics

Structured Schema: Relational databases like PostgreSQL use a fixed schema with predefined tables and columns.

Relationships: They maintain relationships using foreign keys and data from multiple tables is often combined using JOIN operations.

Normalization: Data is normalized to reduce redundancy, leading to multiple interlinked tables.

NoSQL Database Characteristics (MongoDB)

Flexible Schema: MongoDB stores data in JSON-like documents which can have varied structures.

Document Relationships: Relationships are represented through embedded documents or references.

Denormalization: Data is often denormalized, allowing for comprehensive documents that aggregate data from what would be multiple relational tables.

Translating Relational Model to MongoDB

Patient Document Example

In MongoDB, a Patient document could embed data from Disease, Diagnosis, and Address tables:

```
{
  "_id": "123456789",
  "first_name": "John",
  "last_name": "Doe",
  "age": 35,
  "gender": "M",
  "race": "Caucasian",
  "severity_lv": 1,
  "address": {
    "city_name": "San Francisco",
    "state_name": "California",
    "country_name": "USA"
  },
  "hospital_id": "1",
  "disease": {
```



```

        "name": "Influenza",
        "type": "Viral Infection"
    },
    "diagnosis": {
        "issue_date": "2023-01-15",
        "treatment": "Prescribed medication",
        "medicine": [
            {
                "name": "Aspirin",
                "quantity": 1
            }
        ]
    }
}

```

Hospital Collection Example

Separate Hospital collection with reference in Patient documents:

```

{
    "_id": "1",
    "hospital_name": "New York General Hospital",
    "address_id": "1"
}

```

Advantages and Disadvantages of Transition

Advantages in NoSQL

Flexibility: MongoDB's schema flexibility accommodates varying data structures.

Scalability: Efficient in handling large volumes of diverse data.

Performance: Can offer improved performance for certain queries due to data denormalization.

Disadvantages in NoSQL

Complex Aggregations: More complex to perform queries involving multiple levels of data aggregation.

Data Redundancy: Potential increase in data redundancy and storage requirements.

Consistency: Managing data consistency can be more challenging in a denormalized schema.

Conclusion

Transitioning to a NoSQL database like MongoDB offers flexibility and scalability, beneficial for handling large, unstructured datasets and evolving data structures. However, it requires a different approach to data management and may pose challenges in complex data aggregations and consistency. The decision to use a relational or NoSQL database should be driven by specific application requirements, considering the nature of the data, query types, and scalability needs.

Transition to NoSQL Database (MongoDB/Neo4J)

MongoDB

- **Document-Oriented Structure:** MongoDB would store data in JSON-like

documents, allowing for a more flexible schema. Each patient's record, including embedded details like diseases and treatments, would be in one document.

- **No Joins Needed:** Relations between entities like patients and diseases would be handled through embedded documents, eliminating the need for JOIN operations.
- **Scalability and Flexibility:** MongoDB excels in scalability and managing unstructured data, making it a good fit for diverse and evolving healthcare datasets.

Neo4J

- **Graph-Based Model:** Neo4J, a graph database, represents data as nodes and edges. This could be advantageous for mapping complex relationships, like patient interactions or disease pathways.
- **Data Relationships:** Relationships are first-class entities in Neo4J, which would help in visualizing connections between different entities in the healthcare ecosystem.

AWS Architecture for the Application

Core Components

1. **Amazon RDS/Amazon DynamoDB:** For relational data, Amazon RDS can be used; for NoSQL, DynamoDB is suitable.
2. **Amazon S3:** For storing large datasets, backups, and logs.
3. **AWS Lambda:** To handle real-time data processing and batch processing tasks.
4. **Amazon Redshift/Snowflake:** For data warehousing solutions. Redshift if using AWS's native solution or Snowflake as an alternative.
5. **Amazon EC2:** To host application servers and other components that require a virtual server environment.
6. **Amazon API Gateway:** To create, publish, maintain, monitor, and secure APIs.

Data Loading (Batch and Real-Time)

- **Batch Loading:** Use AWS Glue or Lambda functions to periodically transfer data from RDS/DynamoDB to Redshift/Snowflake.
- **Real-Time Loading:** Implement Kinesis Streams for capturing and loading real-time data into the database. Lambda can be used for light processing and transferring data to the appropriate services.

Resilience, Performance, and Security

- **Auto-Scaling and Load Balancing:** Use EC2 Auto Scaling and Elastic Load Balancing to manage and distribute incoming application traffic, ensuring high performance.
- **Data Redundancy and Backup:** Implement S3 versioning and RDS snapshots for data backup and redundancy.
- **Security:** Utilize AWS Identity and Access Management (IAM) for secure access control. Employ encryption for data at rest (using RDS/DynamoDB encryption features) and in transit (using SSL/TLS).

Advantages of Snowflake Over PostgreSQL in AWS

Snowflake Advantages

1. **Separation of Compute and Storage:** Snowflake's architecture allows independent scaling of compute and storage, leading to cost-efficiency and performance optimization.
2. **Zero-Copy Cloning and Time Travel:** Enables easy data recovery and historical data analysis without additional storage costs.

3. **Data Sharing:** Facilitates secure and easy data sharing, beneficial for collaborative healthcare environments.
4. **Automatic Performance Tuning:** Reduces the need for manual performance tuning and maintenance.
5. **Native Support for Semi-Structured Data:** Efficiently handles JSON, Avro, XML directly, which is advantageous for varied healthcare data formats.

Conclusion

Transitioning to a NoSQL database architecture, like MongoDB or Neo4J, offers flexibility, scalability, and efficient handling of complex relationships, essential for a dynamic and data-intensive field like healthcare. Implementing the application in AWS provides a robust, scalable, and secure environment, leveraging services like Lambda, S3, and RDS/DynamoDB. Snowflake, as an alternative to PostgreSQL, brings additional benefits in terms of scalability, data recovery, and performance optimization, making it a compelling choice for modern healthcare data warehousing needs.